

# Bacteriophage replication modules

## supplement: compendium of origins and proteins

Christoph Weigel & Harald Seitz

Max-Planck-Institut für molekulare Genetik, Ihnestr. 73, Berlin, Germany

**Correspondence:** Christoph Weigel,  
Technical University Berlin, Faculty III, Inst. for  
Biotechnology, Sekr. MA 5-11, Straße des 17.  
Juni 136, D-10623 Berlin, Germany.  
Tel. +49 30 84131614,  
e-mail: weigel@molgen.mpg.de

Received 22 March 2005, revised 6 September  
2005, accepted 8 November 2005

doi: 10.1111/j.1574-6976.2006.00015.x

Editor: Ramón Díaz Orejas

**Abbreviations:** bp base pair(s); cp $\phi$  cryptic prophage, in most cases defect, i.e. non-inducible; cSDR constitutive stable DNA replication; D-loop DNA-based displacement loop; DSB double-strand break; dsDNA double-stranded DNA; EMSA electrophoretic mobility shift assay; F4 family 4 (of helicases), e.g. DnaB<sub>F4</sub>; HGT horizontal gene transfer; hmUra 5-(hydroxymethyl)-2'-deoxyuridine; HTH helix-turn-helix DNA-binding motif; iSDR inducible stable DNA replication; NMR nuclear magnetic resonance; NTP nucleoside 5'-triphosphate; orf open reading frame; PCNA proliferating cell nuclear antigen (sliding clamp);  $\phi$  (phi) prefix to mark bacteriophage (replicon); p prefix to mark plasmid (replicon); p $\phi$  prophage; ppDR protein-primed DNA replication; RCR 'rolling circle' DNA replication; RDR recombination-dependent DNA replication; RF replicative form (of ssDNA phages); R-loop displacement loop, based on a RNA-DNA hybrid; RFC replication factor C (clamp loader); SDR stable DNA replication;  $\sigma$ DR sigma( $\sigma$ )-type DNA replication; SF2 superfamily 2 (of helicases), e.g.  $\phi$ T4 UvsW; SSB single-strand (DNA) binding protein; SAP single-strand (DNA) annealing protein; ssDNA single-stranded DNA; TBD thioredoxin binding domain; tDR transcript-driven initiation of DNA replication; (ts) temperature-sensitive mutation;  $\theta$ DR theta( $\theta$ )-type DNA replication; TP terminal protein; UDG uracil-DNA glycosylase.

### Contents

C1.	Introduction	2
C2.	Bacteriophage replication origins	2
C2.1.	Origins for 'rolling circle' DNA replication (RCR)	4
	The replication origin region of the filamentous phage fd of <i>E. coli</i>	6
	The <i>dso</i> and <i>sso</i> replication origins of the isometric phage $\phi$ X174 of <i>E. coli</i>	7
	The <i>dso</i> of <i>V. cholerae</i> phage CTX $\phi$	8
	The <i>dso</i> of <i>E. coli</i> phage P2	8
	Conclusions	10
C2.2.	Origins for DNA replication in the 'theta-mode' ( $\theta$ DR)	10
	The replication origin <i>ori</i> $\lambda$ of phage $\lambda$	12
	The replication origin <i>ori</i> L of phage $\phi$ SPP1	13
	The replication origin <i>ori</i> J of the <i>E. coli</i> Rac prophage	13
	Replication origins of lambdoid (pro)phages	15
	Other bacteriophage origins for $\theta$ DR	19
C2.3.	Other bacteriophage replication origins	20
C3.	Bacteriophage replication proteins	21
C3.1.	Initiator proteins	21
C3.1.1.	Initiator proteins for 'rolling circle' DNA replication (RCR)	22
	Initiator proteins of phages with ssDNA genomes	22
	<i>Gene II protein (gpII) of phage fd - Protein A of phage <math>\phi</math>X174 - RCR initiators of ssDNA phages</i>	
	Initiator proteins of phages with dsDNA genomes: P2 A protein	25
C3.1.2.	Initiator proteins for DNA replication in the 'theta-mode' ( $\theta$ DR)	26
	Bacteriophage Lambda O protein	27
	Gene 38 protein (G38P) of <i>B. subtilis</i> phage SPP1	28
	Initiator proteins of the Lambda O / SPP1 G38P-type	28
	ReplL, the initiator for bacteriophage P1 replication in the lytic cycle	33
	The 'initiator domain' of phage P4 alpha ( $\alpha$ ) protein and related helicases	33
C3.2.	Helicase loaders	35
	Gene 59 protein (gp59) of phage T4	36
	Bacteriophage Lambda P protein	37
	Gene 39 protein (G39P) of <i>B. subtilis</i> phage SPP1	38
	B protein of phage P2	38
	<i>E. coli</i> DnaC and phage-encoded homologues	38
	The helicase loading mechanism in <i>B. subtilis</i>	41
C3.3.	Helicases	42
	Family 4 helicases (5'→3' helicases)	42
	<i>Phage T7 gene 4 helicase - Gene 41 (gp41) helicase of phage T4 - Gp65 helicase of phage D29 - E. coli DnaB and phage-encoded homologues - Miscellaneous phage-encoded family 4 helicases</i>	
	Helicases with similarity to phage P4 alpha ( $\alpha$ ) primase-helicase	47
	Phage-encoded superfamily 2 helicases	50
	Phage-encoded superfamily 1 helicases	52
C3.4.	Primases	52
	The gene 4A primase domain of phage T7	54
	Gene 61 (gp61) primase of phage T4	54
	Phage-encoded <i>E. coli</i> DnaG-type primases	54
C3.5.	DNA polymerases and accessory proteins	59
C3.5.1.	DNA polymerases	59
	<i>E. coli</i> Pol I ( <i>polA</i> )-type DNA polymerases	60
	<i>Phage T7 gene 5-type subfamily - Phage T5 DNA polymerase-type subfamily - Phage D29 gp44-type subfamily - Phage 12 p12-type subfamily</i>	

## Contents (continued)

<i>E. coli</i> Pol II ( <i>polB</i> )-type DNA polymerases	62
<i>E. coli</i> Pol III $\alpha$ ( <i>dnaE</i> )-type DNA polymerases	64
<i>E. coli</i> Pol IV ( <i>dinB</i> )-type and <i>E. coli</i> Pol V ( <i>umuC</i> )-type DNA polymerases	66
5'→3' exonucleases	66
3'→5' exonucleases	67
C3.5.2. Sliding clamps, clamp loaders and DNA polymerase accessory proteins	68
C3.6. Single-strand DNA binding and recombination proteins	69
C3.6.1. Single-strand DNA binding proteins (SSBs)	69
Group 1: phage T4 gene 32 ( <i>gp32</i> ) SSB - Group 2: phage T7 gene 2.5 SSB - Group 3: phage A2 <i>orf34</i> SSB - Group 4: phage M13 gene V SSB - Group 5: phage $\lambda$ . <i>Ea10</i> SSB - Group 6: miscellaneous SSBs - Group 7: <i>E. coli</i> SSB and phage-encoded homologues	
C3.6.2. Recombination proteins	72
The $\lambda$ Red $\alpha$ /Red $\beta$ (Exo/Bet) recombination pathway	74
The <i>E. coli</i> $\phi$ P1 Rac RecE/RecT recombination pathway	74
The $\phi$ SPP1 G34. 1P/G35P recombination pathway	74
The $\phi$ P22 Erf recombination protein	74
Evolutionary relationship of the RecE/RecT, G34. 1P/G35P, and Red $\alpha$ /Red $\beta$ protein pairs	74
Other phage-encoded (putative) exonuclease/SAP gene pairs	76
Holliday junction resolvases	80
Gene 49 ( <i>gp49</i> ) endonuclease VII of phage T4 - Gene 3 endonuclease I of phage T7 - <i>Rap</i> protein of phage Lambda - Phage-encoded resolvases with similarity to <i>E. coli</i> RuvC - Phage-encoded resolvases with similarity to <i>E. coli</i> RusA	
Other phage-encoded recombination proteins	84
C4. Concluding remarks	86
References	87
Note added in proof	101

## C1. Introduction

This supplement of the review 'bacteriophage replication modules' provides a comprehensive basis of knowledge and data – for the discussion of replication mechanisms as well as for the discussion of phage replication modules. It is structured in a way that allows rapid access to particular replication components.

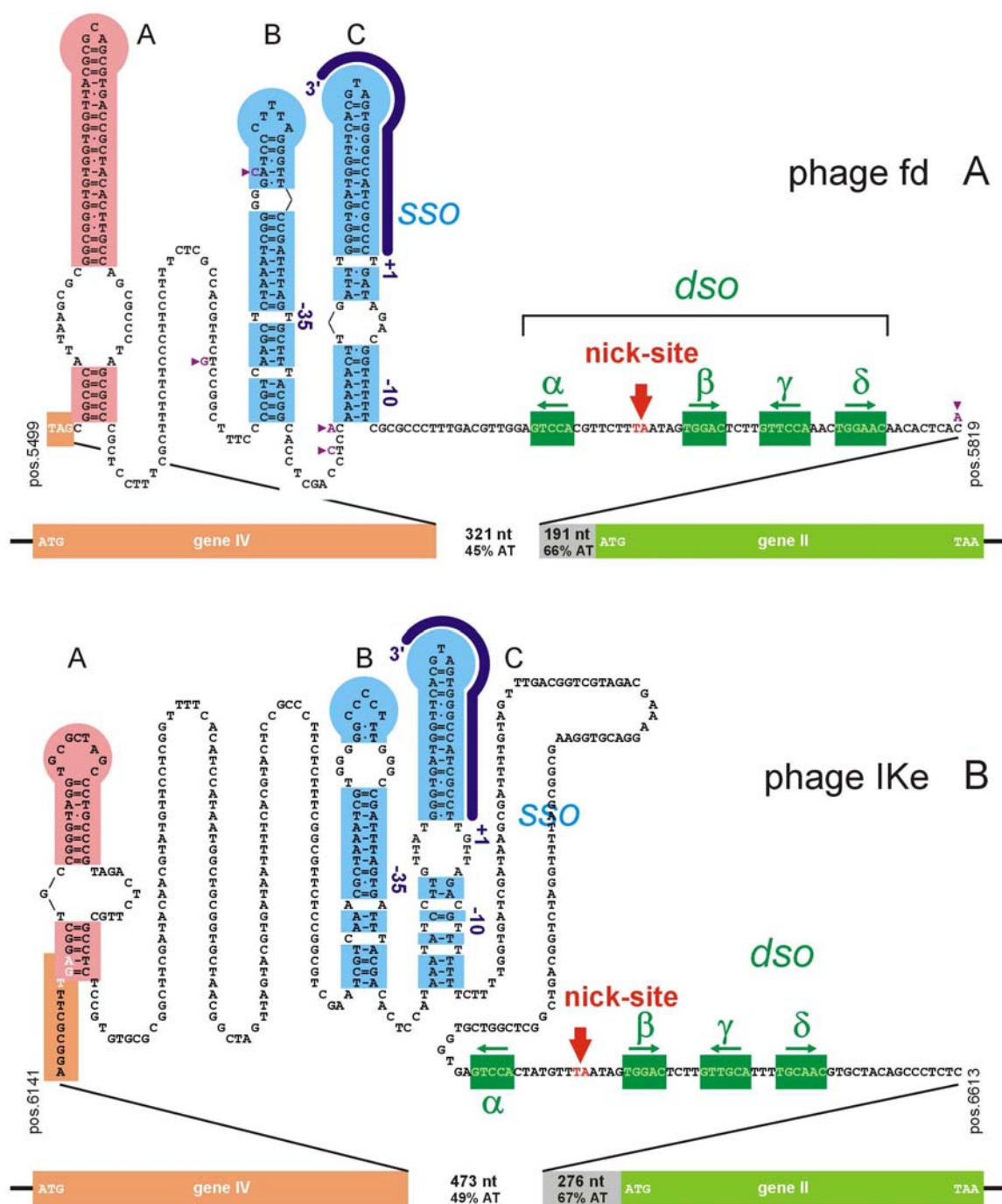
Within the individual chapters of each section, the role of the respective component for bacteriophage replication is discussed first in the context of analogous components from plasmids and chromosomes. We then review the genetics and biochemistry of well-known examples for phage replication origins (Section C2.) and replication proteins (Section C3.). Each chapter concludes with a discussion of homologous or similar origins/proteins that were – in their majority – detected by (DNA/protein) BLAST searches and have not yet been studied experimentally. This latter part within each chapter will appear less 'readable' due to numerous Figures and Tables designed to present the available data, but it should be easily 'searchable', and thus provide a compendium.

*Note:* all sections of the main body of the review 'bacteriophage replication modules' referred to in the

following are marked with 'BRM' (see also last page). All numbers of Sections, Tables, and Figures of this compendium are 'tagged' with the prefix 'C' in order to facilitate unambiguous navigation between both parts of the review. The table of contents and all tables, figures, citations and *internal* cross-references are hyperlinked in the PDF version (not visible).

## C2. Bacteriophage replication origins

A replication origin is defined as a specific segment of a replicon where the synthesis of daughter DNA molecules commences. Prokaryotic replicons contain one single replication origin. There are two exceptions from this rule: i. phage or plasmid replicons that propagate by the rolling circle mechanism of DNA replication (RCR) contain two replication origins, one for (+)-strand synthesis and one for (–)-strand synthesis (BRM Section 2.1.), and ii. joint replicons that use different replication origins for different stages of their life-cycles. E.g. phage P1 replication is initiated in the  $\theta$ -mode at *oriL* (within the *repL* gene) during the lytic cycle, and also in the  $\theta$ -mode at *oriR* during propagation as (plasmid) prophage of its host. Conjugative plasmids use dif-



**Fig. C1** Structure of the  $\phi$ fd and  $\phi$ Ike replication origins.

**A.** Sequence and structural elements of the  $\phi$ fd [J02451] origin region. The sequence of the viral (+)-strand is shown for pos. 5499 to 5819. Stemloops (A = pink, B+C = light blue), the promoter structural motifs (-35 region, -10 region, and +1 transcript start) with the priming transcript (dark blue line) of the *SSO*, and the structural motifs of the *dso* (gpII binding-sites  $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\delta$  = dark green; nick-site = red + red arrowhead) were assigned according to Horiuchi [181]. Genes IV, II, and the intergenic region are shown in a different scale. Sequence and structural elements of the origin region are identical in  $\phi$ f1 and  $\phi$ M13 [NC\_003287] except for the residues indicated and marked by a dark violet triangle. **B.** Sequence and structural elements of the  $\phi$ Ike [NC\_002014] origin region. The sequence of the viral (+)-strand is shown for pos. 6141 to 6613. Numbering and colour code for structural elements as in **A**.

ferent origins for DNA transfer during conjugation and during propagation in their hosts. Formally, also most bacterial chromosomes are 'joint replicons' because they contain varying numbers of prophages. But since initiation from prophage origins is normally tightly repressed, the replication origin *oriC* is the only 'active' replication origin of bacterial chromosomes under most growth conditions (replication restart sites are not considered 'replication origins'). Within this theoretical framework, different experimental approaches result in slightly differing origin definitions – besides their specific resolution.

Electron microscopy was frequently used in early studies for the analysis of bacteriophage replication intermediates. Origins can be mapped by this technique to the smallest detectable 'replication bubbles', or by extrapolating a common 'origin' from a series of bubbles of different size. By this technique, the positions of replication origins can be pinpointed to regions of approximately 100 bp in length (mostly expressed as a position corresponding to % of the genome). Despite its somewhat limited resolution, electron microscopy can contribute valuable details, e.g. indications for unidirectional and/or bidirectional replication, and also indications for coupled leading- and lagging-strand DNA synthesis and/or displacement synthesis. Replication intermediates of RCR and  $\theta$ DR can be distinguished by electron microscopy but a discrimination between replication intermediates of  $\theta$ DR and tDR is not possible.

Genetics defines the replication origin of a prokaryotic replicon as the only element that is strictly required *in cis*. The experimental strategy for the isolation of the *E. coli* chromosomal replication origin, *oriC*, based on this definition [437,286]. Also, this definition has been widely exploited to detect the replication origins of various phages – the history of the  $\lambda$ *dv* plasmids may serve as example (BRM Section 3.). In combination with mutational analysis and DNA sequencing, the genetic approach can define a replication origin more exactly than any other method – at least with respect to size. However, the raw DNA sequence does not allow for a prediction of the replication mechanism operating on this origin unless it is reasonably similar to a known replication origin. Even if the replication mechanism is known, the set of *trans*-acting factors has to be determined in subsequent experiments – mostly, but not always a trivial task.

The major contribution of protein biochemistry to the definition of the replication origin(s) of a given replicon has always been the elucidation of the complete set of factors – cognate and host-encoded – that are required to obtain *de novo* DNA synthesis *in vitro*. But in addition, for most experimental systems the first indications that the conformation of the origin-containing DNA substrate – single/double-stranded, linear/circular, relaxed/negatively supercoiled – determines its proficiency

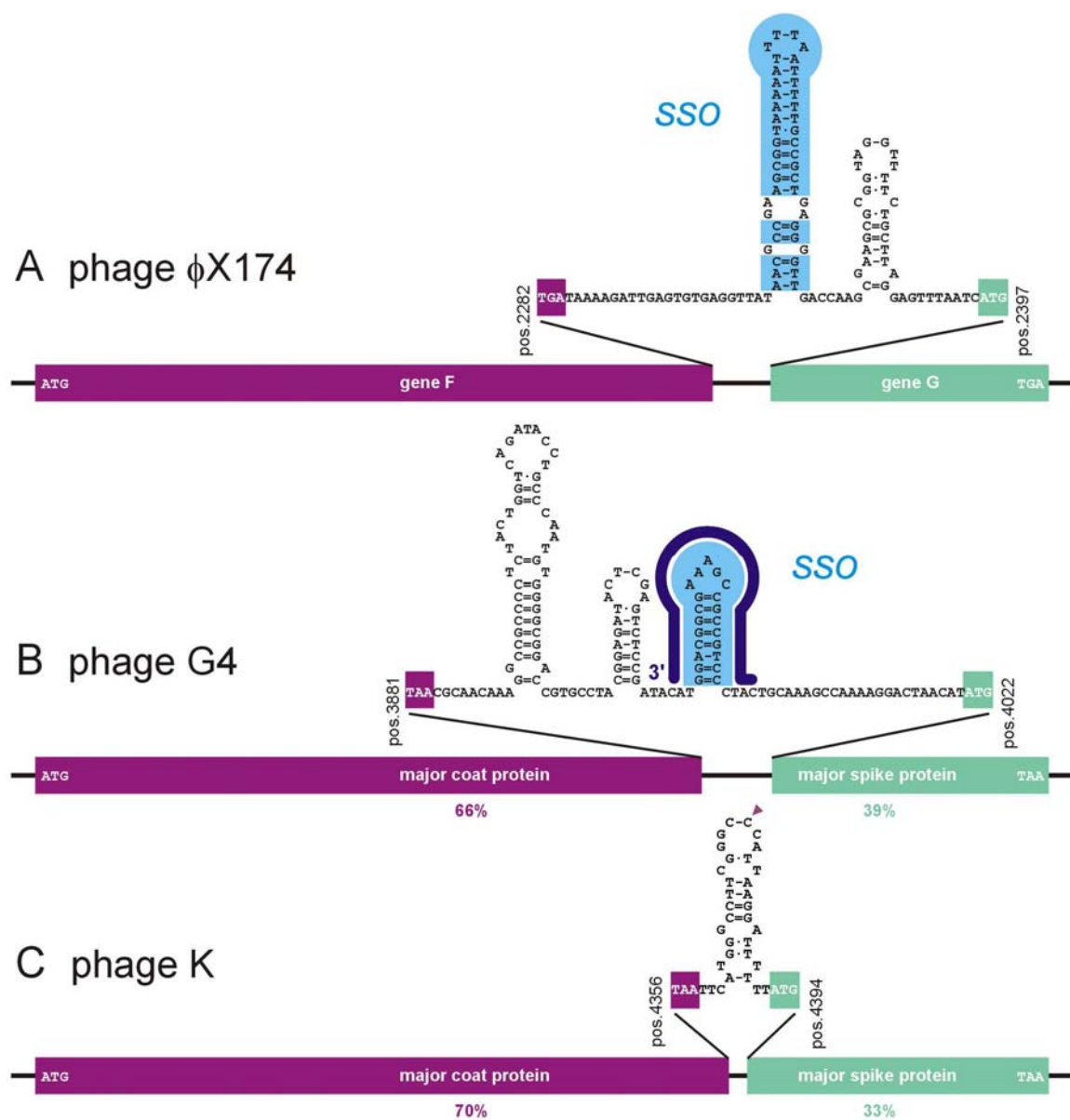
to direct DNA synthesis came from biochemical studies. Biochemistry can determine the first nucleotides synthesised for the leading- and the lagging-strand, but this may lead to a biased definition of the origin if the primase acts distributively rather than at fixed positions (compare e.g. [406] and [169]).

State-of-the-art analyses of replication origins include i. the genetic definition, ii. the DNA sequence, iii. the biochemically defined 'catalogue' of factors required for initiation, and iv. footprinting analyses which determine the nucleic acid-protein interaction(s) and the DNA distortions that characterise the 'DNA unwinding' – either for 'nicking' or for 'melting'. We wish to emphasise that the interaction of a replication origin with its cognate initiation protein(s) is a highly dynamic process and determined by the properties of both, the nucleic acid and the protein(s). Horiuchi and co-workers found a point mutation in the  $\phi$ f1 gpII initiator that renders the protein proficient for unwinding at reduced levels of origin superhelicity [171]. On the other side, Leonard and co-workers describe point mutations in *E. coli oriC* that allow initiation by ADP-DnaA, which is inactive under 'normal conditions' [250] (Section C3.1.2.). The dynamic nature of the origin-initiator interaction has until now frustrated attempts to describe the origin-unwinding reaction in thermodynamically sound terms. We can expect considerable progress by forthcoming analyses employing the novel 'single molecule analysis' technique (reviewed in [433,434]).

The replication mechanism driving the propagation of a replicon determines the structural elements that constitute its replication origin. We will therefore discuss origins for 'rolling circle' DNA replication (RCR) and theta ( $\theta$ )-type DNA replication ( $\theta$ DR) in separate chapters. The  $\phi$ 29 'replication origin' for protein-primed DNA replication (ppDR) is discussed in BRM Section 2.3. . Also, we will not discuss in this section the replication origins for transcription-driven DNA replication (tDR) of the phages T7 and T4. The promoters that serve as replication origins for these phages are discussed in BRM Sections 2.4. & 2.5., respectively. We will discuss instead the replication origins of  $\phi$ c2 and related phages, which probably also replicate by tDR.

## C2.1. Origins for 'rolling circle' DNA replication (RCR)

We characterise in BRM Section 2.1. 'rolling circle' DNA replication (RCR) as a mechanism that avoids topological problems during replication of circular dsDNA molecules by directing the synthesis of both DNA daughter-strands in consecutive steps. Both synthesis steps employ different initiation mechanisms – and initiation proteins – and different replication origin



**Fig. C2** Structure of the  $\phi$ X174,  $\phi$ G4, and  $\phi$ K single-strand origins (*ssos*).

**A.** Sequence and structural elements of the  $\phi$ X174 [NC\_001422] *ssos*. The sequence of the viral (+)-strand is shown for pos. 2282 to 2397. The light blue stemloop corresponds to the PAS site detected by Shlomai and Kornberg [413]. Genes F, G, and the intergenic region are shown in a different scale. Sequence and structural elements of the *ssos* region are identical in  $\phi$ S13 [NC\_001424]. **B.** Sequence and structural elements of the  $\phi$ G4 [NC\_001420] *ssos*. The sequence of the viral (+)-strand is shown for pos. 3881 to 4022. The light blue stemloop corresponds to the structure, and the dark blue line to the DnaG-dependent transcript detected by Bouché *et al.* [41]. The genes for the major coat protein, the major spike protein, and the intergenic region are shown in a different scale; %-values indicate similarity to the corresponding proteins of  $\phi$ X174 (ident. res.). **C.** Sequence and structural elements of the putative  $\phi$ K [NC\_001730] *ssos*. The sequence of the viral (+)-strand is shown for pos. 4356 to 4394. Sequence and structural elements of the *ssos* region are identical in  $\phi\alpha$ 3 [NC\_001330] except for the missing C residue in the stemloop marked with a lilac triangle. The genes for the major coat protein, the major spike protein, and the intergenic region are shown in a different scale; %-values indicate similarity to the corresponding proteins of  $\phi$ X174 (ident. res.).

structures. Therefore, phage replicons that propagate by RCR contain a 'double-strand origin' (*dso*) for the synthesis of single-stranded DNA circles – the viral (+)-strand – and a 'single-strand origin' (*sso*) for the synthesis of the complementary strand – the viral (–)-strand. The *dso* contains the binding site(s) for the cognate initiator protein responsible for the 'nicking reaction' that precedes RCR. Replication initiation driven by a *sso* on single-stranded DNA involves host RNA polymerase ( $\phi$ fd), primase ( $\phi$ G4), or a restart primosome ( $\phi$ X174). A *sso* is characterised by DNA secondary structure elements that are specifically recognised by the initiation protein(s).

**The replication origin region of the filamentous phage fd of *E. coli*.** The structural and functional elements of the almost identical replication origin regions of the filamentous phages fd, f1 and M13 have been extensively studied by Horiuchi and co-workers [181].

The *dso* of  $\phi$ fd spans a "core region" of ~40 bp, which is absolutely required for viral (+)-strand synthesis. The *dso* contains the four binding-sites  $\alpha$ – $\delta$  for the phage-encoded initiator gpII (Section C3.1.1.) arranged as inverted repeats, and the nick-site between  $\alpha$  and  $\beta$  (Fig. C1; A). While the specific sequence of the gpII binding-site  $\alpha$  is not absolutely required for initiation, the sequence around and including  $\alpha$  is necessary for the termination reaction. The AT-rich region downstream of the *dso* is important for origin function but its exact role unknown. A deletion of the AT-rich region can be compensated for by mutations elsewhere in the genome, and also insertion in this region – as in cloning vectors derived from  $\phi$ f1 or  $\phi$ M13 – do not abolish origin function [181].

In the 'left' part of the origin region, stemloop A is thought to serve as 'morphogenetic signal' for packaging of viral (+)-strands into phage heads and therefore not part of the replication origin in the strict sense [181], stemloops B and C constitute the *sso*. Both stemloops B and C together create a partially single-stranded promoter structure with the –35 region residing within stemloop B, and the –10 region and the transcription start within stemloop C. Except for the direct neighbourhood of the –35 and –10 regions, the base-pairing in the stems is more important than a specific nucleotide sequence, and the entire structure is recognised by *E. coli* RNA polymerase holoenzyme ( $\alpha\beta\beta'\sigma^{70}$ ). The impressive promoter strength – approximately 15-fold stronger than the *lacUV5* promoter – is probably due to the partial single-strandedness of the otherwise non-canonical promoter [174,173,172].

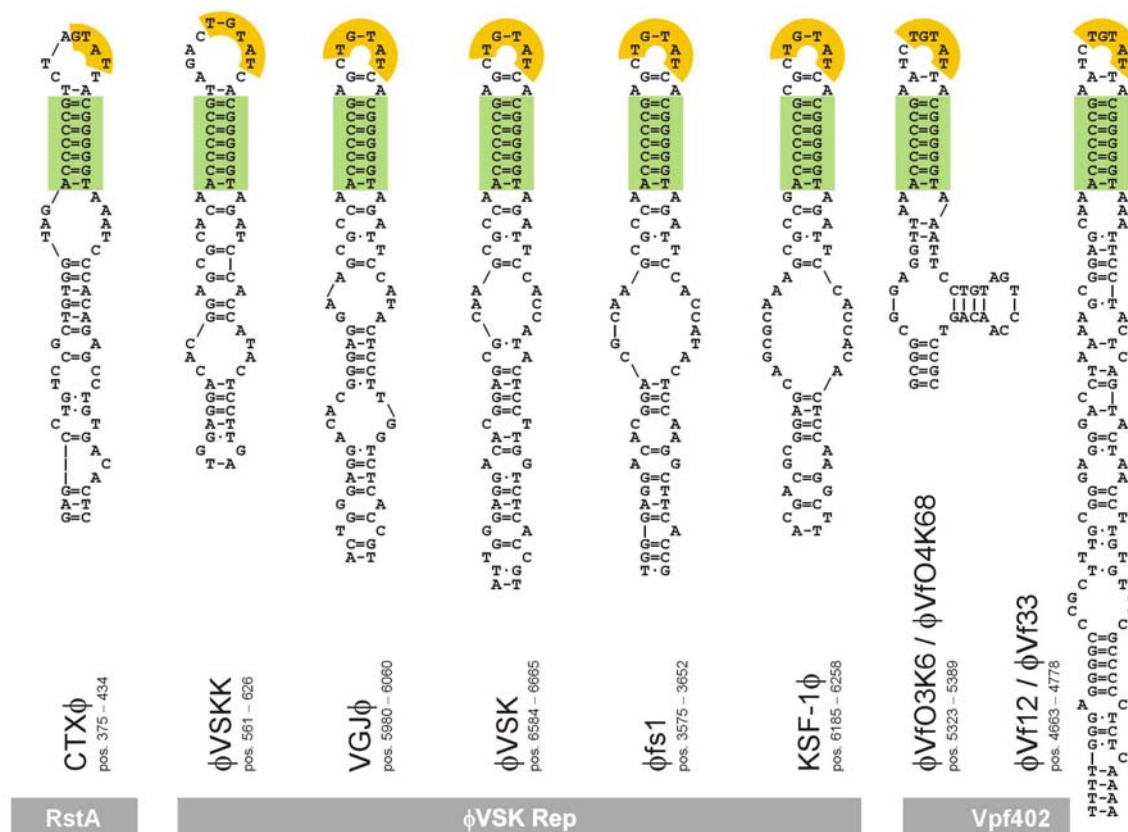
In Figure C1, we show the *dso* regions of  $\phi$ fd and  $\phi$ IKe without the possible additional stemloop structures D and E – encompassing four gpII-binding sites  $\alpha$  through  $\delta$  – because they are irrelevant for initiation by

gpII on the doublestranded RF (compare to Fig. 1 in [181]). Initiation by nicking occurs after the sequential binding of gpII to all 4 binding-sites, which would be impossible if  $\alpha + \beta$  form part of stemloop D, and  $\gamma + \delta$  part of stemloop E, respectively, on the viral (+)-strand DNA [151,170]. The (putative) stemloops D + E might nevertheless have a biological role, e.g. preventing untimely binding of gpII and nicking of the viral ssDNA. We note that also in a number of *dsos* of RCR plasmids the nick-sites are embedded in (putative) stemloop structures [208].

According to Konings and co-workers, the filamentous phage Ike – infecting host cells carrying IncI2 plasmids – is a distant relative of the F-specific phages fd, f1 and M13 [343]. The genome organisation of  $\phi$ IKe is strikingly similar to that of  $\phi$ fd in the gene IV/gene II region. In addition, two short stretches of (almost perfect) sequence homology between  $\phi$ IKe and  $\phi$ fd allowed the straightforward identification of stemloops B + C (*sso*) and the nick-site region (*dso*) (Fig. C1; B). Noteworthy, the individual motifs are separated by intervening sequences of different lengths and varying sequences. This observation emphasises that *sso* and *dso* have to be considered independent structural motifs. The colocalisation of these motifs in  $\phi$ fd ( $\phi$ f1,  $\phi$ M13) seems to be rather the exception than the rule: *sso* and *dso* do not co-localise in  $\phi$ X174 and related phages (see below), and also not in various RCR-plasmids [208].

Together with gpII of  $\phi$ fd and  $\phi$ IKe, the (putative) initiator protein gpII of phage I2-2 belongs to group 3 of RCR-initiators (Section C3.1.1.; Table C9). Except for a comparable size and the conserved 'active tyrosine'-motif 3 however,  $\phi$ I2-2 gpII is unrelated to  $\phi$ fd gpII (<20% ident. res.) but shows instead significant homology to the orf1(Rep) protein of *Xanthomonas campestris* RCR-plasmid pXV64 (36% ident. res.) [476], and to the RepC protein of *Moraxella sp.* RCR-plasmid pTA144up (39% ident. res.). Konings and co-workers assumed therefore that  $\phi$ I2-2 inherited its replication module from a RCR-plasmid [427]. However, the genome organisation of  $\phi$ I2-2 is comparable to that of  $\phi$ fd and  $\phi$ IKe: the genes IV and II are separated by an intergenic region of 835 nt. The part proximal to gene II can be classified as 'AT-rich' (65% A+T) in contrast to the distal part (51% A+T), and a putative stemloop structure analogous to stemloop A could be readily identified directly downstream of gene IV (motifs as in Fig. C1). The few short and patchy sequence homologies to either  $\phi$ fd or  $\phi$ IKe did not allow pinpointing the putative *sso* and *dso* structures (not shown). This example demonstrates that RCR-origin structures cannot always be simply deduced by comparing nucleotide sequences of (more distantly) related phages.





**Fig. C3** (Putative) replication origin (*dso*) of CTX $\phi$  and related vibriophages.

Secondary structure of the region containing the conserved *dso*-motif (light green) of CTX $\phi$  taken from Moyer *et al.* [306] with slight modifications; the putative nick-site is shown in orange. The positions of the regions shown were taken from the entries for CTX $\phi$  [(partial sequence) VCU83796],  $\phi$  VSKK [NC\_003311], VGJ $\phi$  [NC\_004736],  $\phi$  VSK [NC\_003327],  $\phi$  fs1 [NC\_004306], KSF-1 $\phi$  [NC\_006294],  $\phi$  Vf12 [NC\_005949],  $\phi$  Vf33 [NC\_005948],  $\phi$  Vfo4K68 [NC\_002363], and  $\phi$  Vfo3K6 [NC\_002362]. Protein names given in the grey bar in the lower part refer to (putative) initiator subgroups as discussed in Section C3.1.1. .

**The *dso* and *ssso* replication origins of the isometric phage  $\phi$ X174 of *E. coli*.** The *dso* of  $\phi$ X174 is located at a distance of  $\sim$ 1900 nt from the *ssso*. The *dso* of  $\phi$ X174 was mapped to a discrete 30 bp stretch within the A gene (see Table C1; corresponding to amino acid res. 107–116) and the position of the nick-site has been determined as A<sub>4307</sub> [14,187,115]. The *dsos* of phages  $\phi$ X174,  $\phi$ S13,  $\phi$ G4,  $\phi$ U3,  $\phi$ G14,  $\phi$ K,  $\phi\alpha$ 3, and  $\phi$ St-1 are virtually identical, and all are recognised by  $\phi$ X174 A protein (Table C1) [166]. Although the A proteins of the isometric phages are close homologues (similarity ranging from 43% to 96%; see group 2A in Table C8), the nucleotide sequence similarity drops significantly to either side of the 30 bp *dso* core sequence. The  $\phi$ X174 *dso* lacks detectable secondary structures or sequence repeats, and the binding-site(s) of A protein is not known exactly.

Despite their uniform *dsos*, the isometric phages possess different *ssos*, and also the mechanism for comple-

mentary strand synthesis differs from that of the filamentous phages (see above and BRM Section 2.1.). Shlomai and Kornberg could show that efficient *in vitro* complementary strand synthesis of  $\phi$ X174 depends on the host PriA protein and assembly of the restart primosome on a primosome-assembly site (PAS) located between the genes F and G (Fig. C2; A) [413]. The sequence containing the PAS site is identical in  $\phi$ X174 and in  $\phi$ S13, and  $\phi$ S13 complementary strand synthesis most likely follows the  $\phi$ X174 pathway. In contrast, *in vitro* complementary strand synthesis of  $\phi$ G4 depends on the synthesis of a specific  $\sim$ 26–28 nt long primer by *E. coli* DnaG [41]. Bouché *et al.* could locate the template of the RNA primer to a stemloop structure in the intergenic region between the genes coding for the major coat protein and major spike protein (Fig. C2; B). Despite considerable protein sequence similarity between  $\phi$ X174 and  $\phi$ G4 for F protein/ major coat protein and G protein/ major spike protein, respectively, the nucleotide sequen-

**Table C1** Double-strand origin (*dso*) of  $\phi$ X174 and related phages.

$\phi$	gene	acc.	5'-flanking	pos.	<i>dso</i>	pos.	3'-flanking
$\phi$ X174	A	NC_001422	AGTGCTCCCC	4299	CAACTTG↓ATATTAATAACACTATAGACCAC	4328	CGCCCCGAAG
$\phi$ S13	A	NC_001424	.....	382	.....↓.....	411	.CG.....
$\phi$ G4	A	NC_001420	C.....GGA	500	.....↓.....	529	AAA...CT..
$\phi$ U3	-	M10630	C.....GGA	125	.....↓.....	154	AC...AA...
$\phi$ G14	-	M10632	.....GGA	55	.....↓.....	84	AC...TA...
$\phi$ K	A	NC_001730	GTGTGCAG..	1001	.....↓...A..G.....	1030	...A..CC.T
$\phi$ $\alpha$ 3	A	NC_001330	GTGTGCTG..	1002	.....↓...A..G.....	1031	..TA..CC.T
$\phi$ St-1	-	1)	GTGTGCTG..	-	.....↓...A..G.....	-	..TA

Identity with the nucleotide sequence of the  $\phi$ X174 gene A sequence is indicated by '!'. The nick-site is indicated by an arrow. 1) partial sequence taken from Heidekamp *et al.* [166].

ce and (putative) secondary structure elements in the intergenic regions are widely different. The *ssos* of the closely related phages  $\phi$ K and  $\phi\alpha$ 3 are not known, but – by analogy to  $\phi$ X174 and  $\phi$ G4 – a (putative) secondary structure in the intergenic region between the genes coding for the major coat protein and major spike protein could serve this purpose (Fig. C2; C).

The  $\phi$ X174 A protein is distantly related to the (putative) replication protein encoded by *orf4* of  $\phi$ chp1 (Table C8; groups 2A + 2B). The replication of the small chlamydial phages has yet to be studied, and their replication origins remain to be determined.

**The *dso* of *V. cholerae* phage CTX $\phi$ .** The small ssDNA phages were instrumental in developing the present models for prokaryotic replication and have been widely used as smart cloning vehicles. However, their biological importance was generally underestimated until it became known that they can be responsible for the spread of pathogenicity factors: e.g. the cholera toxin genes are disseminated among *Vibrio cholerae* strains by a filamentous phage, CTX $\phi$  [95]. Waldor and co-workers analysed the replication of CTX $\phi$ , which may reside in (mostly multiple) tandem copies as prophage in the genomes of *V. cholerae* strains. The mechanism of CTX $\phi$  prophage formation in *V. cholerae* is not known, but lysogenic host cells release mature phages [379]. Using an elegant genetic technique developed by Horiuchi [179], Waldor and co-workers could show that a (presumed) *dso* of CTX $\phi$  present twice on a test plasmid is reduced to a single chimeric origin after replication [306]. This result offered a (mechanistically) simple explanation for the puzzling observation that unit-length closed-circular viral (+)-strands were produced from tandemly arranged chromosomal copies: initiation of RCR at the first *dso* is followed by ring closure at a second *dso* in the 'prophage array' (BRM Section

2.1.). Supporting these results, the phage-encoded RstA protein shares the 'active tyrosine' motif 3 with numerous initiators for RCR of phages and plasmids (Section C3.1.1.). We note, however, that neither the function of RstA as initiator for RCR has yet been proven biochemically, nor the position and functionality of the presumed nick-site confirmed.

A number of filamentous vibriophages encode (putative) initiators for RCR with significant similarity to the CTX $\phi$  RstA protein (see Section C3.1.1. and Table C8; subgroup 1A). We were interested whether also the *dso* of CTX $\phi$  is conserved in these phages. Nucleotide sequence comparison with the 'stemloop 2' defined by Moyer *et al.* [306] as query led to the detection of corresponding structures in all phage genomes (Fig. C3). Although we show the conserved structures as stem-loops in Figure C3, we should cautiously speak of 'stretches with dyad symmetry' rather, until their stem-loop structures are experimentally proven. The conserved 'stemloop structures' are embedded in an only partially conserved context but are invariably located in intergenic regions in the respective genomes. Although not a stringent proof, the presence of a highly conserved 'stemloop 2'-like structure in all compared phage genomes encoding similar initiator proteins supports the hypothesis that they constitute the *dsos* of these phages. A *ssso* has yet to be identified for CTX $\phi$  and the other vibriophages.

**The *dso* of *E. coli* phage P2.** Besides the group of phages with small single-stranded DNA genomes (<10 kb for the RF) that replicate via RCR, a second group of phages with mid-sized dsDNA genomes (~30 kb) is known to replicate by this mechanism. The best known phage of the latter group is *E. coli* phage P2, for which the nick-site of the *dso* could be located within the 3'-end of gene A encoding the initiator protein



**Table C2** Double-strand origin (*dso*) of  $\phi$ P2 and related phages.

$\phi$ / $p\phi$	gene	5'-flanking	pos.	<i>dso</i>	pos.	3'-flanking	ident.	translation <sup>1)</sup>
$\phi$ P2	A	CGGCATCGCC	29885	GCGC <b>CTCG</b> ↓GAGTCCTGTCA <b>ATAACTGT</b>	29911	GGAAAGCTCA	x	AAPRSPVNNC
$\phi$ L-413C	gpA	.....	28691	.....↓.....	28717	.....	97%	.....
$\phi$ W $\phi$	gpA	.....	28077	.....↓.....	28103	.....CTG	96%	.....
<i>Y. pseudo.</i> $p\phi$	yptb1766	..CATCT.G.	2119424	.....↓.....	2119398	..GTTAGGTG	56%	.....
$\phi$ PSP3	gp36	..CTCCT.G	28618	C.CT.....↓.....T.....	28644	ACGGGAAG.G	39%	..PS..S....
$\phi$ 186	A	..CTTCT.G	28568	C.CT.....↓.....T.....	28594	ACGGGAAG.G	37%	APS..S....
$\phi$ Fels-2	stm2729	G...GCTTG.	2870666	..C...T.↓..C..G..G.....	2870640	CCCCTTGCTG	34%	C..WTRG...
$p\phi$ CP-933T	z2978	...CCT...	2676644	..C...T.↓..C..G..G.....	2676670	CCCCGTG.AC	31%	A..WTRG...
$\phi$ HP2 / $\phi$ HP1	Rep	..CACG.AGT	7444	.....T.↓..C.TG.....G.....	7470	AACCGCTCA.	31%	VERS.TA..S
$\phi$ PM2	p12	..A..G..A.	3504	..C...T.↓...CA...AA.....	3530	AACC.C.C.G	33%	D..W.TE...
$\phi$ fs2	orf716	ACA..ATTGG	2357	..T...T.↓..C..G...T.....	2383	C..CC.TA.TT	33%	IGCLGL.SLT
$\phi$ V86	-	C...TCT...	5922	..CT.....↓.....T.....	5948	ACG <sup>2)</sup>	36%	..S..S....
$\phi$ K139	Rep	C...TCT...	10212	..CT.....↓.....T.....	10238	ACGG.T.C.T	31%	..S..S....

The nick-site in the  $\phi$ P2 *dso* (within the coding region for A protein) determined by Liu and Haggård-Ljungquist is indicated by an arrow [257]. Identity with the nucleotide sequence of the  $\phi$ P2 gene A sequence is indicated by '.'; bold letters (in the  $\phi$ P2 A lane) indicate positions conserved in all 13 sequences. ident. = BLAST similarity (% ident. res.) as in Table C9; group 4. 1) reading frame in the *dso* region; identity with  $\phi$ P2 A protein is indicated by '1. 2) end of sequenced region. Sequences were taken from:  $\phi$  P2 [NC\_001895],  $\phi$  L-413C [NC\_004745], W $\phi$  [NC\_005056], *Y. pseudotuberculosis*  $p\phi$  [NC\_006155],  $\phi$  PSP3 [NC\_005340],  $\phi$  186 [NC\_001317],  $\phi$  K139 [NC\_003313],  $\phi$  V86 (partial sequence) [AF008938], Fels-2  $p\phi$  [NC\_003197],  $\phi$  fs2 [NC\_001956],  $p\phi$  CP-933T of *E. coli* O157:H7 EDL933 [NC\_002655],  $\phi$  PM2 [NC\_000867], and  $\phi$  HP2 [NC\_003315] (the DNA sequence of  $\phi$  HP1 [NC\_001697] is identical in this region). For each (putative) *dso*, the relative position within the gene is shown as projection onto the protein sequence in Fig. C6.

[256,257] (BRM Section 2.1.). In Section C3.1.1., we discuss the similarity of 12 known phage-encoded homologues of  $\phi$ P2 A protein. Here we can show that a stretch of 27 bp around the  $\phi$ P2 nick-site is well conserved in all 13 initiator genes, suggesting an identical nick-site in all sequences (Table C2). This 27 bp 'core' is contained within the region of 35 bp suggested by Liu and Haggård-Ljungquist as *dso* by comparison of the  $\phi$ P2 and  $\phi$ 186 sequences with other phage or plasmid *dsos* [257]. Except for the  $\phi$ L413-C and W $\phi$  genes that are (almost) identical with the entire  $\phi$ P2 A sequence, the similarities drop immediately on either side of the 27 bp 'core' sequence. This is reminiscent of the *dsos* of  $\phi$ X174 and related phages, where the conservation of the nucleotide sequence was also confined to the 30 bp 'core' (see above). The ' $\phi$ P2-type' *dso* shares more features with the  $\phi$ X174-type *dso*: i. the 3'-end of the nick contains the conserved CTXG motif also present in a number of plasmid *dsos*, ii. all sequences lack detectable direct or inverted repeats within the *dso* and the adjacent regions, and iii. they also lack a detectable AT-rich region. In contrast to the  $\phi$ X174 A gene however, the (putative) *dsos* are located close to the 3'-end of the genes. Except for  $\phi$ f2,  $\phi$ V86, and  $\phi$ K139, the putative *dsos* are located at exactly the same position in the genes as the  $\phi$ P2 *dso* (the position of the putative *dsos* within the initiator genes is shown in Fig. C6).

Given the reasonable protein sequence similarity of the  $\phi$ P2 A homologues and the conserved position of

the *dsos*, the observed nucleotide sequence conservation of the putative *dsos* may simply reflect the requirement for conservation of structurally/functionally important amino acid residues. The degree of conservation of the amino acid residues in the *dso* region reflects the overall similarity of the proteins, and does not allow for a clear decision (Table C2; rightmost column). However, alternative 'translation' of the  $\phi$ f2 *dso* in the -1 frame results in a sequence WVPWTRVI-NC, which is closer to the  $\phi$ P2 A sequence. This indicates that the conservation of the nucleotide sequence is important, supporting the tentative assignment of the double-strand origin to this part of initiator genes of the 13 phages.

Vibriophage  $\phi$ fs2 is a typical filamentous phage – only distantly related to  $\phi$ fs1, though (see above) – with a ssDNA genome size of 8.7 kb (RF). It is the first example for a small phage replicon with a  $\phi$ P2-type replication module for RCR. Also *Alteromonas* phage MP2 has a relatively small genome (10.1 kb dsDNA).  $\phi$ MP2 is unique among the known bacteriophages because its capsid contains lipids [272]. However, its initiator protein p12 (Section C3.1.1.) and its (putative) *dso* are clearly related to the ' $\phi$ P2-type'. Both examples demonstrate that the  $\phi$ P2-type replication module for RCR is not confined to the closer relatives of *E. coli* phage P2. However, plasmid replicons with a  $\phi$ P2-type module are not (yet) known, which suggests that this module is not promiscuously distributed among different genetic elements. We note that neither for  $\phi$ P2 nor for any other of

the mentioned phages the *ssos* have yet been determined.

**Conclusions.** Despite accomplishing identical functions and containing the nick-site for initiation of RCR, the double-strand origins discussed above differ considerably with respect to their specific structural elements. The  $\phi$ d *dso* contains repeats as binding sites for the initiator, flanking the nick-site. The *dso* of CTX $\phi$  contains a region of dyad symmetry that may form a stemloop, with the nick-site in the loop. Neither of the above-mentioned structural elements could be detected in the ~30 bp long *dsos* of  $\phi$ X174 and  $\phi$ P2. There are thus three clearly distinct types of phage *dsos* for RCR, but the more subtle differences in the initiation and termination reactions that can be expected thereof remain to be studied.

The *dsos* of the phages also differ with respect to their localisation within the replicon. The *dsos* of the filamentous *E. coli* phages are located directly upstream of gene II encoding the initiator in an intergenic region. Also the *dsos* of the vibriophages are located in intergenic regions, but unlinked to the initiator gene. In contrast, the *dsos* of the isometric phages and of the  $\phi$ P2-group are located within the coding region of the initiator: the  $\phi$ X174-type *dsos* in the 5'-region, and the  $\phi$ P2-type *dsos* in the 3'-region, respectively. A localisation of the *dso* within the cognate initiator gene is a feature also found in plasmids replicating by RCR, e.g. in pT181 and related plasmids (Table C8; group 1B) [210]. Origin-containing initiator genes are known, in addition, for phage replicons propagating in the  $\theta$ -mode, e.g. the lambdoid phages (see below). Although this 'double layer' of genetic information imposes strong selective pressure on the initiator gene – simultaneous selection for a functional origin and for a functional protein – we have to realise that this intricate organisation is apparently evolutionary successful.

## C2.2. Origins for DNA replication in the 'theta-mode' ( $\theta$ DR)

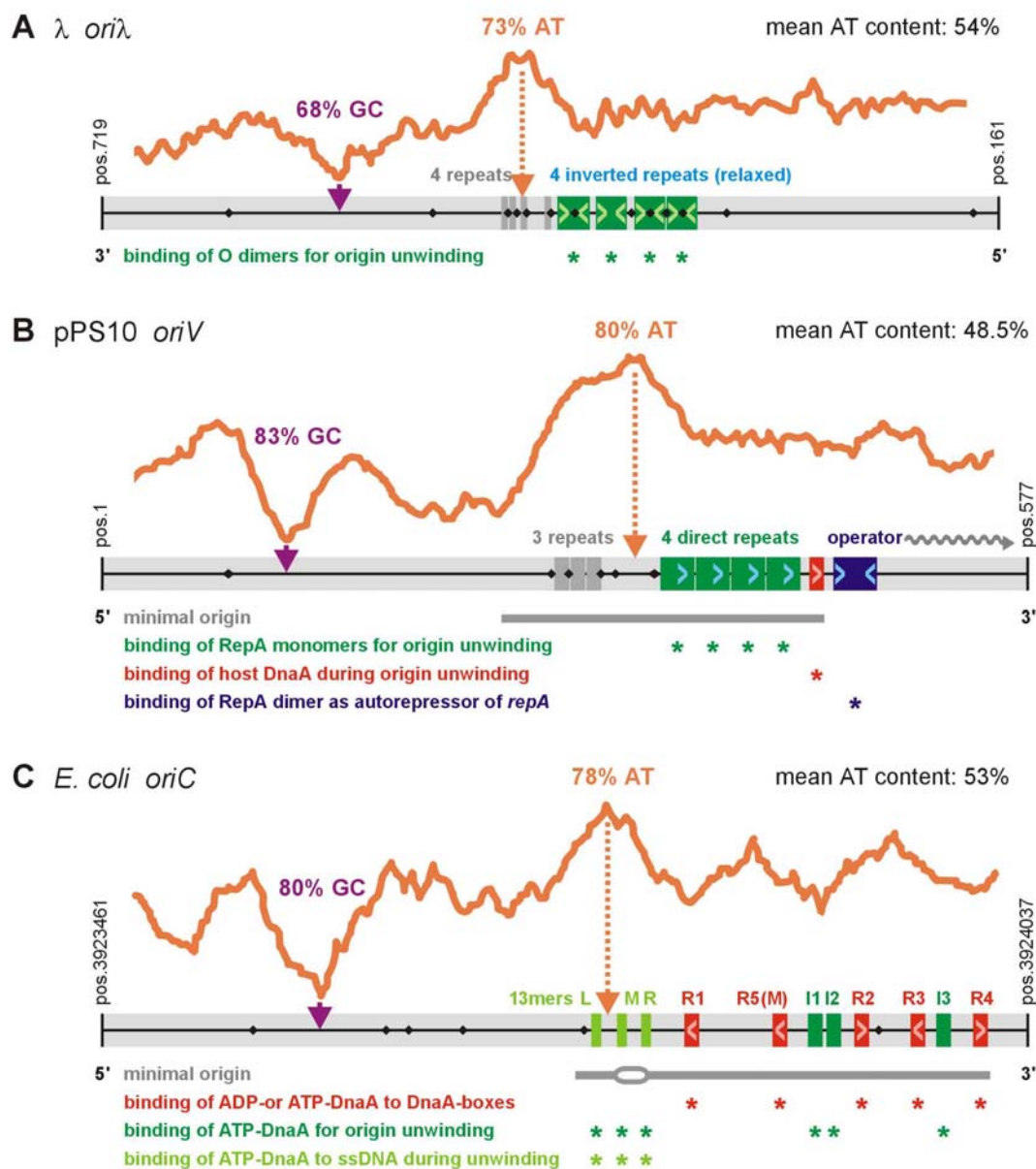
All known prokaryotic origins for DNA replication in the theta( $\theta$ )-mode have a strikingly similar architecture irrespective of whether they drive the replication of plasmids, phages, or chromosomes (Fig. C4). In the generalised 'origin design', a stretch of 100–200 bp containing a number of short repeated sequence motifs (iterons) is located on one side of an AT-rich region of ~50 bp in length. Initiation occurs after binding of the cognate initiator protein to the iterons on a negatively supercoiled template: concomitant with oligomerisation of the iteron-bound initiators, a long-range distortion of the DNA occurs that is sufficiently strong to promote

the 'melting' of a short stretch of neighbouring bases within the AT-rich region. The partially single-stranded (unwound) region serves as entry site for primosomal and replisomal proteins, culminating in the establishment of a replication fork(s).

Two types of iterons can be discriminated in origins for  $\theta$ DR: inverted repeats (*ori $\lambda$* ) or direct repeats (pPS10 *oriV*, *E. coli oriC*) (Fig. C4). Direct repeats can be perfect (*oriV*) or degenerate (*oriC*). They can be arranged head-to-tail (*oriV*) or with irregularly alternating orientation (*oriC*). Regularly alternating orientations of direct repeats would appear as array of larger (imperfect) palindromes. Iterons of the inverted repeat type can be either perfect or imperfect palindromes, as in *ori $\lambda$*  (Fig. C4; A). We note that although any of these motifs are usually detectable in putative origin sequences at first glance their correct classification requires knowledge of the binding-mode of the cognate initiator protein. As a rule of thumb, initiators bind as monomers to direct repeats (RepA→*oriV*, DnaA→*oriC*), while dimeric initiators bind to inverted repeats (O→*ori $\lambda$* ). Iterons of either type can be regularly spaced (*oriV*, *ori $\lambda$* ) or irregularly spaced (*oriC*). A regular iteron spacing can coincide with helical phasing (*oriV*: 4 iterons with a length of 22 bp, directly adjacent to each other = 8 helical turns), but also may deviate from it. It could be shown for phage  $\lambda$  *ori $\lambda$*  and *E. coli oriC* that disruption of the phasing of the iterons had a more pronounced (negative) effect on origin function than had their mutation [299,291,484,239]. This latter observation emphasises that the spacing of the iterons is dictated by the architecture of the nucleoprotein complex that is responsible for origin unwinding. Even in the best-studied model systems however, the architecture of this nucleoprotein complex is not known to the necessary detail – present efforts concentrate on the elucidation of the initiator-initiator interactions and on describing the path of the origin DNA within the complex (see e.g. [250]).

In addition to the iterons that serve as binding-sites for the initiator protein, many origins contain additional repeats for which a function is not yet known (*ori $\lambda$* , *oriV*). They are mostly direct repeats and their consensus sequence differs from that of the iterons. These 'secondary repeats' are often located more closely to (or within) the AT-rich region than the initiator-binding sites. The observation that *E. coli* DnaA binds to such repeats in the AT-rich region of *oriC* when they become single-stranded during unwinding has so far not been repeated in other systems and may be a unique property of the DnaA-type of initiators (Section C3.1.2.) [423].

The regulation of replication initiation at origins for  $\theta$ DR is as strikingly different as is their overall architecture similar. In the pPS10 system, two types of binding-sites for the RepA initiator exist that allow for the regulation of initiation in response to the oligomerisation



**Fig. C4** Structural elements of three origins for replication in the  $\theta$ -mode.

Origin regions are shown as grey bar, structural elements as coloured boxes; all sequence motifs are shown exactly to scale. Red boxes = 'canonical' DnaA binding-sites (DnaA boxes [389, 239]); dark green boxes = binding sites of the cognate initiator protein that are crucial for origin opening (I-boxes in the case of *ori*C [250]); light green boxes = 'ATP-DnaA boxes' bound as ssDNA by ATP-DnaA during origin unwinding [423]; blue boxes = initiator binding sites not required for origin unwinding; grey boxes: repeats of unknown function. '<' and '>' indicate directionality/dyad symmetry within the protein-binding sites. The function of a DNA-binding protein at a particular binding-site is indicated by '\*' and additional explanatory text. '♦' indicate stretches of  $\geq 4$  A (or T bases, respectively). Note that in **A** the conventional 5'→3' orientation of the sequence was reversed to allow for a better comparison with **B** and **C**. AT-plots (orange) were computed by the spin v1.1 software from the Staden Package [425], using a window of 41 nt (therefore, the plots are shown for -20 bp from either end only). The AT-plots in **A**, **B**, and **C** are shown in the same scale with the upper edge of the sequence bar as baseline. The positions of the (partial) sequences analysed for AT-content are shown. The genetically defined 'minimal origin' is indicated by a dark grey bar in **B** and **C**. **A**: partial sequence of the phage  $\lambda$  *o* gene [J02459] (see text for details). **B**: partial sequence of plasmid pPS10 [X58896] (for details see [143]), the genetically defined 'minimal origin' is indicated by a dark grey bar. **C**: partial sequence of the *E. coli* K12 chromosome [U00096] (for details see [290, 250]), The genetically defined 'minimal origin' is indicated by a dark grey bar and the unwound region by an open bubble [232].

state of RepA (reviewed in [143]). Both binding-sites share the motif GGACAG, which is present once in each *oriV* iteron, and – as inverted repeat – twice in the operator of the *repA* promoter (Fig. C4; B).

Monomeric RepA binds to the *oriV* iterons. Dimeric RepA acts as autorepressor: *repA* transcription is repressed by binding of a RepA dimer to the operator site. In contrast to the *oriV*/RepA paradigm, binding of DnaA to *oriC* is not regulated by the type of the iterons. The DNA-binding and the oligomerisation properties of *E. coli* DnaA depend crucially on the bound nucleotide: ADP-DnaA and ATP-DnaA bind as monomers – and with comparable affinity – to DnaA boxes that match the consensus TTATNCACA (Fig. C4; C: R1, R2, R4) [389]. Binding to 'low affinity' DnaA boxes that deviate from the consensus (Fig. C4; C: R3, R5) occurs only when a close-by 'high affinity' DnaA box promotes cooperative binding, and requires ATP-DnaA. Also the non-canonical binding sites (Fig. C4; C: I-boxes, ATP-DnaA boxes in the L, M, R 13mer region) are only bound by ATP-DnaA and require neighbouring 'high affinity sites' for cooperative binding [424,423,250]. Different from the above systems,  $\lambda$  O protein binds as dimer to single iterons, and further oligomerisation occurs among iteron-bound dimers (see below). A specific role for the O monomer and iteron half-sites as targets for binding of O monomers are not known. Unlike when bound to DNA, O is highly unstable in solution. Therefore the initiation of  $\lambda$  replication may be simply regulated by the availability of O, and may not involve more intricate regulatory loops as known for the initiators of plasmids and chromosomes.

The AT-rich region is a prominent feature in the three origins shown in Figure C4 and also readily detected in many replication origins for  $\theta$ DR. In some origins that have been studied experimentally (*Saccharomyces cerevisiae* ARS, Simian Virus SV40), spontaneous unwinding of negatively supercoiled DNA occurs close to AT<sub>max</sub> in the absence of divalent cations – not necessarily at *exactly* the same position where initiator-induced unwinding occurs [460,252]. This region has been termed 'DNA unwinding element' (DUE) therefore, and it could be shown for *E. coli oriC* that it is the helical 'instability' that is required, and not the exact sequence [231]. In all three origin sequences shown in Figure C4, the AT-rich region is flanked on one side by the iteron-region, and on the opposite side by a GC-rich region. Although this 'GC-peak' adds to the overall similarity of the three origin regions, it does not play a (known) role for origin functioning and is not part of the genetically defined 'minimal' origins of pPS10 and *E. coli*.

A feature of origin sequences that is important for the DNA distortion required for unwinding is sequence-induced (intrinsic) curvature. Intrinsic curvature is detectable by an abnormal electrophoretic behaviour of DNA

in gels, and was first described for the minicircle DNA of *Leishmania tarentolae* kinetoplasts. *oriλ* DNA was the first prokaryotic DNA for which curvature could be demonstrated [499,500]. Intrinsic curvature of DNA can in many instances be predicted for a sequence with periodically spaced A runs of ~5-6 bp in length. This is the case for *oriλ*, but not for *oriV* and *oriC* (see Fig. C4; '♦' in the sequence bar). Minimal *oriC* DNA has a close-to normal electrophoretic behaviour that varies however slightly depending on the degree of *dam* methylation (H.Seitz and C.Weigel, unpublished). A functional role has not been demonstrated for the region with significant intrinsic curvature immediately downstream of minimal *oriC* (not shown in Fig. C4; C) [215]. Intrinsic curvature is thus not a sufficient criterion for the detection/definition of replication origins. Initiator-induced unwinding requires negative superhelicity of the origin DNA. We refer the reader to BRM Section 2. for a detailed discussion of this topic.

The three examples given in Figure C4 show origins for  $\theta$ DR which differ with respect to localisation in their cognate replicon. The phage  $\lambda$  replication origin, *oriλ*, is located within the *O* initiator gene, *oriV* of pPS10 upstream of the *repA* initiator gene, and *E. coli oriC* at a distance of ~40 kb from the *dnaA* gene – in contrast to *B. subtilis* where *oriC* and *dnaA* co-locate. As already pointed out in the previous chapter, there is no preference for origin localisation in prokaryotic replicons.

**The replication origin *oriλ* of phage  $\lambda$ .** The  $\lambda$  initiator protein O binds to the relaxed palindromic iteron ATCCCT<sub>c</sub>AAAA<sub>TCTG</sub><sub>A</sub>G<sub>0-1</sub>RGGGAT<sub>AC</sub> present four times in *oriλ* as part of the *O* gene (Fig. C4; Table C3). The iterons are arranged as an array with exact helical phasing, and 5'- upstream of the AT-rich region at a distance of approximately two helical turns (23 bp). A considerable variation exists among the number of iterons and their lengths even among phage initiators with highly conserved protein sequence: there are e.g. three iterons in gene 15 of  $\phi$ 80, but 6 in the *O* gene of  $\phi$ H19-B. In both cases, the 'backbone' of the relaxed palindromic iteron is highly similar to the *oriλ* iterons and iteron spacing corresponds exactly to helical phasing, i.e. one helical turn (Table C3). It could be shown for *oriλ* that only three of the four iterons are absolutely required for origin function [299]. The phenomenon that initiator-binding sites defined by footprinting may be disabled without grossly impairing origin function is also known for *E. coli oriC* [23,239,472]. As a consequence, the number of iterons detected in *oriλ* and other phage origins for  $\theta$ DR (see below) is probably higher in most cases than the minimal number required for origin functioning. The interaction of  $\lambda$  O with a single iteron and with the entire *oriλ* is discussed in more detail in Section C3.1.2. .

**Table C3** (Putative) replication origins of Gram(-)-specific phages; part A.

(pro)phage	gene	nt	AT	consensus sequence of iteron(s)	nt	#	groups	type	spacing	pos.	d
φ D3	orf73	990	424	TGAGACCAAACWCCKCTCACTC	21	2	1+2+1	rIR <sup>1</sup> )	-	5'	13
				TGAGC-N <sub>8/13</sub> -TCTC	17/22	2		rIR	-		
φ P22	gp18	816	451	GGTAAAA-N <sub>1-6</sub> TWACC	13-18	3	3	rIR	-	5'	25
φ 80	15	909	466	TTCCCGAAAWC-N <sub>0-1</sub> GGWA	15-16	3	3	rIR	5 nt / 2 t	5'	24
λ / φ 21	O	900	459	ATCCCT <sub>/c</sub> AAAA <sub>/Tc</sub> /T <sub>G</sub> /A <sub>G</sub> 0-1RGGGAT <sub>/AC</sub>	19-20	4	4	rIR	1 t	5'	23
φ 4795	O	936	496	ATCCCT <sub>/c</sub> AAAA <sub>/Tc</sub> /T <sub>R</sub> GGGA <sub>/GW</sub>	19	4	4	rIR	1 t	5'	22
φ H-19B	O	939	493	CCCTC <sub>0-1</sub> AAAA <sub>/T</sub> YG <sub>0-1</sub> AGGG	14-15	6	6	rIR	1 t	5'	22
φ 933W	O	939	490	CCCT <sub>/c</sub> AAAA <sub>/Tc</sub> /T <sub>(A/G)</sub> 0-1A <sub>/G</sub> GGG	14-15	6	6	rIR	1 t	5'	19
cpφ CP-933V	z3356	939	496	CCCT <sub>/Tc</sub> (C <sub>/T</sub> ) <sub>0-1</sub> A <sub>/c</sub> AAAC <sub>/T<sub>G</sub>/T<sub>R</sub>A<sub>/G</sub></sub> GGG	14-15	6	6	rIR	1 t	5'	24
cpφ CP-933X	z1868	744	376	TTACCC <sub>/T<sub>G</sub>/A<sub>A</sub>T<sub>/A</sub>C<sub>/G</sub></sub> AGGTAA	17	5	5	rIR	1 t	5'	17
pφ LambdaSo	O	1017	385	-							
				478	ATTGTCRRAACCGAYAG	17	2	2	rep	2 t	5'
<i>P. putida</i>	pp3894	831	409	GTGT <sub>/G</sub> A <sub>/T</sub> GAAA <sub>/T<sub>A</sub>/C<sub>T</sub>B<sub>0-1</sub>W<sub>0-1</sub>A<sub>0-1</sub></sub> CCAC <sub>/G</sub>	13-16	3	3	rIR	-	5'	5
φ Nil2	o	822	442	GCT <sub>/c</sub> N <sub>7-11</sub> AGC	13-17	4	1+3	rIR	2 t	5'	29
φ HK97	gp54	822	446	GCT <sub>/c</sub> N <sub>7-11</sub> AGC	13-17	4	1+3	rIR	2 t	5'	33
cpφ Gifsy-2	stm1014	984	628	ATCCCVGAAAA <sub>/G</sub> RGGGA <sub>/W</sub>	18	6	4+3	rIR	2 t	5'	22
				GTACCTGAAAC <sub>/G</sub> AGGCA <sub>/T</sub>	18	1					
cpφ CP-933P	z6070	978	433	TC <sub>/A</sub> AA <sub>/T</sub> MDG <sub>/T</sub> TTGA	11	4	8	rIR	6 nt	5'	14
				C <sub>/G</sub> A <sub>/G</sub> AA <sub>/T</sub> T <sub>/A</sub> AT <sub>/A</sub> TG <sub>/A</sub>	9	4					
				GMA <sub>/S</sub> AK <sub>/TK</sub> G	9	2	3	rIR	6 nt	5'+3'	3/2
<i>P. putida</i>	pp1551	855	415	AAACMG	6 <sup>2</sup> )	4	2+2	rep	1 t/4 t/1 t	5'	19
φ Aaφ23	O	783	409	AC <sub>/T</sub> AC <sub>/AT</sub> T <sub>/A</sub> VG <sub>/c</sub> NT <sub>/AG</sub> T <sub>/AG</sub> A <sub>/G</sub>	11	7	-	rep	-	5'+3'	1/0
				AC <sub>/T</sub> AT <sub>/A</sub> VG <sub>/c</sub> NT <sub>/AG</sub> T <sub>/AG</sub> A <sub>/G</sub>	10	1					
φ P27	L17	927	493	T <sub>/G</sub> CA <sub>/T</sub> G <sub>/A</sub> AT <sub>/A</sub> T <sub>/c</sub> C <sub>/AT</sub> T <sub>/A</sub>	9	8	1+7+2	rep	1.5 t	5'	18
				KCAK <sub>2</sub> WSKW	9	2					
φ SfV	orf39	819	400	A <sub>/c</sub> TTCAN <sub>7-10</sub> GTG <sub>0-1</sub> CA	17-20	5	6	rep	~2 t	5'	16
				TTGCACCAGGGGGTAGTGC	19	1		rep			
φ ST64B	sb42	819	403	A <sub>/c</sub> TTCAN <sub>7-10</sub> GTG <sub>0-1</sub> CA	17-20	5	6	rep	~2 t	5'	19
				TTGCACCAGGGGGTAGTGA	19	1		rep			

**The replication origin *oriL* of phage φSPP1.** The approximate localisation of the φSPP1 replication origin could be inferred from analysis of phage replication mutants with an 'immediate stop' phenotype and mapped to a region containing the genes 38, 38.1, 39, and 40 [284,55,341,475]. Alonso and co-workers were the first to show that the gene order in the replication origin region of φSPP1 closely resembles that of other lambdoid phages [341]. They could demonstrate specific binding of G38P protein to a 393 bp DNA segment from within the coding sequence of gene 38. A prominent AT-rich region approximately in the middle of gene 38 is flanked 5'-upstream by an array of five direct repeats (AB boxes) with the consensus sequence A<sub>G</sub>G<sub>A</sub>AAA-C<sub>A</sub>G<sub>A</sub>C<sub>A</sub>A<sub>T</sub> (Table C5). Together with the single intervening base pair, the 10 bp long iterons are aligned in a phasing of exactly one helical turn per iteron. Footprint analysis confirmed that the iterons are the binding-

sites for the G38P initiator protein (Section C3.1.2.) [284,297]. The importance for origin function of the two more degenerate 'ab boxes' located at a distance of 47 bp upstream of the major iteron cluster is presently not known (Table C5).

A second origin-like structure, *oriR*, is present in the φSPP1 genome at a distance of ~32 kb from *oriL*. Compared to *oriL*, *oriR* lacks one AB box in the major iteron cluster, and although G38P binds efficiently to *oriR*, unwinding does not occur [341,297]. Alonso and co-workers could show that *oriR* is instrumental for the switch from θDR to σDR during replication of φSPP1. In their model, DNA-bound G38P 'initiates' this switch by blocking replication fork progression through *oriR* (BRM Section 2.2.) [12].

**The replication origin *oriJ* of the *E. coli* Rac phage.** Díaz and Pritchard identified the replication ori-

**Table C4** (Putative) replication origins of Gram(-)-specific phages; part B.

(pro)phage	gene	nt	AT	consensus sequence of iteron(s)	nt	#	groups	type	spacing	pos.	d	
cpϕ CP-933M	z1337	1203	487	AAG <sub>3-4</sub> CTT <sub>/G</sub>	8-9	3	3	rIR	1 t	5'	161	
				643	AAG <sub>4</sub> CTT	9	3	3	rIR	1 t	5'	60
					AAGAACAGGAACAGGA	16	2	2	rep	2 t	5'	14
ϕ Bcep22	gp26	774	391	C <sub>/G</sub> CC <sub>/R</sub> C <sub>/T</sub> T <sub>/C</sub> T <sub>/G</sub> C <sub>/G</sub> T <sub>/G</sub> GG <sub>/R</sub> A <sub>/G</sub> A <sub>1-2</sub> M-G <sub>/C</sub> A <sub>R</sub> HCG	16-17	5	1+4	rep	3 t + 3x 2 t	5'	24	
ϕ E125	gp60	993	435	AAAG <sub>/C</sub> GGTGCA	10	4	9	rep	-	5'	11	
				GGGTTC A	7	3						
				TGCACCC (= inverted GGGTGCA)	7	2		rep <sup>5)</sup>				
<i>P. luminescens</i>	plu3473	753	250	YTGA <sub>/C</sub> AC <sub>/T</sub> C <sub>/R</sub>	7	6	-	rep	-	5'+3'	10/5	
<i>B. bronchis.</i>	bb1683	957	415	CACCCCGCA	10	7	7	rep	-	5'	27	
cpϕ CP-933O	z2048	1041	346	GGCAAGAACTAAGCAAAATGCTGGTGTAAAC, AGTGTGCGGAGCCTGAATACACC, GGTCTGTTAATAAAACAGTACC, GGTATGTTTCAGCACAAACAAGACC	24-28	4	4	rIR	-	5'+3'	-	
<i>S. enterica</i>	sty2063	789	421	TTGTAACGCTCGCGCGCTTACAA, GTTACAAAAGGCGAGTTAC, GTAACGCTGCCAGCGTTAC, AGCGTTAC	19-20	3,5	3+1/2	IR	-	5'	11	
cpϕ CP-933R	z2397	789	421	TGTAACGCTCGCGCGCTTACA, GTAACGCTGCAGCGTTAC, AGCGTTAC, GTTACAAAAGGAGAGTAAC	19-20	3,5	3+1/2	IR	-	5'	11	
cpϕ Rac	<i>ydaU</i>	858	490	C <sub>/G</sub> A <sub>/T</sub> G <sub>/T</sub> VTGTT <sub>/C</sub>	8	6	1+1+4	rep	1 t	5'	24	
<i>S. flexneri</i>	sf0873	633	292	not detected	-	-	-	-	-	-	-	
ϕ ε15	-	678	397	G <sub>1</sub> TCCGCN <sub>5</sub> GTCCGCN <sub>7</sub> GCGGACN <sub>5</sub> -GCGGAY	41	2	1+1	IR	-	5'	22	
<i>B. bronchis.</i>	bb2207	825	517	CGT <sub>N</sub> ACC   CGT <sub>N</sub> ACC	14	5.5	0.5+1+1+3	IR <sup>4)</sup>	-1 / -2 nt	5'	29	
<i>X. fastidiosa</i>	xf1560	906	400	AAACCCA	7	5	2+7	rep	[1 t]	5'+3'	-	
				A <sub>/TG</sub> AA <sub>/G</sub> C <sub>/G</sub> CC <sub>/T</sub> A <sub>/G</sub> (1 mismatch)	6	1+1						
<i>H. somnus</i>	hsom1669	714	334	TCCGGAWA <sub>/CT</sub>   AA <sub>/G</sub> A <sub>/G</sub> TCCGA	16	4	4	IR	2 t	5'	20	
cpϕ CP-933N	z1773	1038	475	GTCACG <sub>/TR</sub> A <sub>/GTC</sub>	7	11	6+5	rep	1 t	5'	17	
cpϕ CP-933U	z3124	1038	466	G <sub>/R</sub> TCACG <sub>/T</sub> A <sub>/GTC</sub>	7	11	6+5	rep	1 t	5'	8	
ϕ ST64T	gp18	900511		TGTCCAA <sub>/T</sub> C <sub>/R</sub> G <sub>/A</sub>	9	6	3+3	rep	2 nt / 1 t	5'	17	
ϕ HK022	O	900	502	TACCAGCA	8	3	1+6+3	rep	[3 nt / 1 t]	5'	11	
				T <sub>/R</sub> A <sub>/CG</sub> C <sub>/T</sub> CAG <sub>/C</sub> A <sub>R</sub> HN	8 <sup>3)</sup>	7						
ϕ VT2-Sa	O	900	502	TACCAGCA	8 <sup>3)</sup>	4	1+6+3	rep	[3 nt / 1 t]	5'	11	
				T <sub>/R</sub> A <sub>/G</sub> C <sub>/T</sub> CAG <sub>/C</sub> A <sub>R</sub> HN	8	6						
ϕ HK620	O	900	505	TACCAGCA	8 <sup>3)</sup>	4	1+9	rep	[3 nt / 1 t]	5'	14	
				T <sub>/R</sub> A <sub>/G</sub> C <sub>/T</sub> C <sub>/R</sub> A <sub>/T</sub> G <sub>/C</sub> A <sub>/T</sub> A <sub>R</sub> D	8	6						

Legend to Tables C3 + C4 . nt = length of gene; AT = AT<sub>max</sub>; nt = length of iteron; number of iterons; groups = clustering of iterons; type = iteron type (rep: repeat, IR: inverted repeat, rIR: relaxed inverted repeat); spacing = spacing between iterons, indicated as distance (nt) between outside bases or helical turns (t) calculated from the iteron middle; pos. = position of iterons/major iteron cluster relative to AT<sub>max</sub> ('AT-peak'); d = distance (nt) between the proximal iteron and AT<sub>max</sub>. Inverted repeats in the iterons are indicated by boxes. Abbreviations used in the 'consensus' column: R = A or G; Y = C or T; M = A or C; K = G or T; S = C or G; W = A or T; H = A or T or C; D = A or G or T; B = G or C or T; V = A or C or G; N = any. 1) relaxed inverted repeat. 2) repeats embedded in 2 rIR: GA(GT)(T/C)CGAAACAGACCAAAC(A/C)GCCGACTC. 3) 2 perfect repeats: TGTACCAG-CAGATTACCAGCAAATTACCAC. 4) nested palindromes: cgtcaccgCGTAAcgtTCACGCCTAAcgttacccgGTGACC. 5) 2 iterons arranged with the inverted iterons in two 19 and 21 bp palindromes, respectively. DNA sequences for Table C3 were taken from: ϕD3 [NC\_002484], ϕP22 [NC\_002371], ϕ80 [X13065], λ / ϕ21 [J02459], ϕ4795 [NC\_004813], ϕH-19B [AF034975], ϕ933W [NC\_000924], cpϕ CP-933V [NC\_002655], cpϕ CP-933X [NC\_002655], ϕϕ LambdaSo [NC\_004347], *P. putida* [NP\_746024], ϕNil2 [ECO413274], ϕHK97 [NC\_002167], cpϕ Gifsy-2 [NC\_003197], cpϕ CP-933P [NC\_002655], *P. putida* [NP\_743708], ϕAaϕ23 [NC\_004827], ϕP27 [NP\_543069], ϕSV [NC\_0034441], and ϕST64B [NC\_004313]. DNA sequences for Table C4 were taken from: cpϕ CP-933M [NC\_002655], ϕBcep22 [AY349011], ϕE125 [NC\_003309], *P. luminescens* [BX571870], *B. bronchiseptica* [NC\_002927], cpϕ CP-933O [NP\_287511], *S. enterica* [NP\_456424], cpϕ CP-933R [NC\_002655], cpϕ Rac [NC\_000913], *S. flexneri* [NC\_004337], ϕε15 [NC\_004775], *B. bronchiseptica* [NC\_002927], *X. fastidiosa* [NC\_002488], *H. somnus* [NZ\_AABO02000023], cpϕ CP-933N [NC\_002655], cpϕ CP-933U [NC\_002655], ϕST64T [NC\_004348], ϕHK022 [NC\_002166], ϕVT2-Sa [NC\_000902], and ϕHK620 [NC\_002730].



**Table C5** (Putative) replication origins of Gram(+)-specific phages; part A.

(pro)phage	gene	nt	AT	consensus sequence of iteron(s)	nt	#	groups	type	spacing	pos.	d
φ SPP1	gp38	771	439	A <sub>/G</sub> G <sub>/A</sub> AAAC <sub>/A</sub> GA <sub>/G</sub> C <sub>/A</sub> A <sub>/T</sub> (AB boxes)	10	5	2+5	rep	1 t	5'	25
				A <sub>/TG</sub> A <sub>/T</sub> A <sub>/T</sub> AAA <sub>/C</sub> GAC <sub>/T</sub> A (ab boxes)	10	2					
φ mv4	orf292	876	565	CAGAA <sup>GT</sup> T <sup>G</sup> AAC <sup>AA</sup>	14	2	2	rIR	2 t	5'	118
		670	-								
φ r1t	orf11	735	382	TGGAACC	7	4	4	rep	1 t	5'	27
φ 3626	gp34	753	268	-							
				451	CC <sup>TG</sup> TAG <sup>CA</sup> AAAT	14	3	3	IR	2 t	5'
φ 13	phi13_15	861	493	TGG <sup>AAAA</sup> ACC <sup>GTT</sup> AAC <sup>GTT</sup> AAAA <sup>CCA</sup>	26	2	2	rIR	4 nt	5'	29
				TG <sup>GA</sup> AAAACC <sup>GTT</sup> AAT <sup>GTT</sup> AAA <sup>TC</sup> GC							
φ SLT	orf256	771	400	TGG <sup>AAAA</sup> ACC <sup>GTT</sup> AAC <sup>GTT</sup> AAAA <sup>CCA</sup>	26	2	2	rIR	4 nt	5'	35
				TG <sup>GA</sup> AAAACC <sup>GTT</sup> AAT <sup>GTT</sup> AAA <sup>TC</sup> GC							
φ ETA	orf22	771	409	TGG <sup>AAAA</sup> ACC <sup>GTT</sup> AAC <sup>GTT</sup> AAAA <sup>CCA</sup>	26	2	2	rIR	4 nt	5'	35
				TG <sup>GA</sup> AAAACC <sup>GTT</sup> AAT <sup>GTT</sup> AAA <sup>TC</sup> GC							
φ LL-H	orf299	900	469	G <sub>/A</sub> G <sub>/A</sub> AA <sub>/T</sub> V <sup>EC</sup> <sub>/A</sub>	8	6	6	rIR	1.5 t	5'	44
φ 315.1	spyM3_0691	762	448		44	2	2	rIR	5 t	5'	27
				ATGGA <sup>ACTGTA</sup> AAAA <sup>TACAGT</sup> ATCGGAA <sup>CTG</sup> GAA <sup>TTT</sup> A <sup>CAG</sup>							
<i>S. pyogenes</i>	sps1159	836	520		44	2	2	rIR	5 t	5'	23
				ATGGA <sup>ACTGTA</sup> AAAA <sup>TACAGT</sup> ATCGGAA <sup>CTG</sup> GAA <sup>TTT</sup> A <sup>CAG</sup>							
<i>S. pyogenes</i>	spyM18_1803	939	487	AGAG <sup>TTA</sup> G <sup>A</sup> TAA <sup>AGAGA</sup>	17	4	4	rIR	1 nt	5'+3'	-
<i>L. gasseri</i>	lgas 0588	924	459	AAGA <sup>TAA</sup> G <sup>T</sup> TAA <sup>G</sup>	13	2	1+4	rIR	2 nt		
				AWGATAA	7	3					
pφ LambdaBa04	ba3819	759	448	C <sup>AG</sup> AG <sup>A</sup> Y <sup>T</sup> A	9	3	3	rIR	1 t	5'	8
		520	-								
φ TP901-1	REP	819	382	TTTTCCAGATGTGG	14	2	4	rep	1 nt	5'	45
				YCWGAT	6	2				9 nt	
φ PV83	orf20	804	469	2x <sup>GTCAC</sup> <sup>GTGAC</sup> ; 1x <sup>GT</sup> <sup>CAC</sup> <sup>G</sup> <sup>CAAC</sup>	10	3	3	rIR	1 t	5'	40
pφ LambdaBa02	ba4121	717	382	CCAC <sup>AAAA</sup> T <sup>CA</sup> CCCA <sup>GTGG</sup>	19/22	2	2	rIR	8 nt	5'	44
				CC <sup>ACCA</sup> GT <sup>GGA</sup> AAA <sup>AC</sup> CACC <sup>GG</sup>							
φ BC6A52	bc2563	747	400	R <sup>TGGARRW</sup> WY <sup>CAB</sup>	13	4	6	rIR	2 nt	5'	33
				A <sup>TGRARRW</sup> W <sup>CMC</sup>	12	2					
φ PVL	orf46	894	442	TCAA <sub>/G</sub> HMC <sub>/A</sub> A <sub>/G</sub> AC	10	8	1+7	rep	1 t	5'	18
φ N315	sa1791	894	451	TCAA <sub>/G</sub> HMC <sub>/A</sub> A <sub>/G</sub> AC	10	8	1+7	rep	1 t	5'	27
φ A118	gp49	933	448	AC <sub>/T</sub> AA <sub>/T</sub> CG <sub>/A</sub> DHHD	10	10	1+2+6+1	rep	1 t	5'+3' 1)	36

gin *oriJ* of the Rac prophage of *E. coli* K12 by its property to drive replication of a selectable plasmid [105]. The location of *oriJ* within in the ~13 kb insert of the original plasmid pLG-2 could be narrowed down to ~2.5 kb, encompassing the *ydaSTUV* genes (H.Seitz and C.Weigel, unpublished). In analogy to the arrangement of the λ replication genes and origin, Wegrzyn and co-workers proposed for *oriJ* four inverted repeats (consensus sequence GCATCTN<sub>11-13</sub>AGATG-C) upstream of an AT-rich region within the *ydaU* gene as putative binding sites of the YdaU initiator [354]. These 4 and 2 additional more degenerate inverted repeats are rather widely spaced, and located preferentially in regions enriched for GC within the *ydaU* 5'-half. We could de-

tect, in addition, six non-palindromic iterons (Table C4). When compared with the arrangement of the iterons in other phage replication origins (see below), the size of these latter iterons, their number, their regular spacing, and their position closely upstream of the 'AT-peak' in the *ydaU* gene correspond better to a generalised design of a phage replication origin for θDR than the 4 inverted repeats proposed by Wegrzyn and co-workers.

**Replication origins of lambdoid (pro)phages.** Experimental data on the replication origin structure of lambdoid phages other than those discussed above are scarce. By sequence comparison and genetic analysis of phage hybrids, the replication origin of φP22 was map-

ped within gene 18 encoding the initiator (analogous to  $\lambda$  *O*) [346,15].

The conserved arrangement of structural elements in the replication origins of  $\lambda$  and  $\phi$ SPP1 led us to assume that the replication origins should be easily detectable also for other lambdoid (pro)phages: i. by the presence of an AT-rich region located approximately in the middle of the (putative) initiator gene, and ii. by the presence of iterons 5'-adjacent to the AT-rich region. For the analysis, we chose the 40 (putative) initiator genes of Gram(-)-specific phages and the 40 (putative) initiator genes of Gram(+)-specific phages discussed in Section C3.1.2.

A+T plots revealed a prominent 'AT-peak' in ~90% of the analysed genes. The 'AT-peak' located approximately to the middle of most genes (Fig. C5; see also Tables C3 – C6). The AT-rich region 3'-adjacent to the iterons in the  $\lambda$  *O* gene and  $\phi$ SPP1 gene 38 appeared as 'AT-peaks' in the A+T plots, suggesting that this peak is indicative of an origin-specific AT-rich region also in the other phage initiator genes. We found one single prominent 'AT-peak' in 73% of the (putative) initiator genes, and two centrally located 'AT-peaks' in 16%. In the remaining sequences (11%), either multiple peaks were found, or no significant AT-clustering. Interestingly, no 'AT-peak' could be detected in the  $\lambda$  *cI* gene encoding the  $\lambda$  repressor;  $\lambda$  CI protein is comparable in size and architecture to  $\lambda$  *O* (N-terminal DNA-binding domain, C-terminal oligomerization domain) [58].

All but 3 of the 80 initiator gene sequences contain easily detectable repeats, in >90% of the cases clustered 5'-upstream of an 'AT-peak'. Although only assigned theoretically, the presence of such short repeated sequences in virtually all (putative) initiator genes suggests that these repeats represent iterons, i.e. binding sites for the phage replication initiator, and are thus part of the replication origins of these phages (Tables C3 – C6). However, with respect to type (simple or inverted repeat), number, spacing, distance to the 'AT-peak', and orientation to each other these (putative) iterons present a colourful bouquet. Only in those cases where the initiator proteins show a high degree of similarity in their N-termini (>60% ident. res.; Section C3.1.2.) the iterons found in the respective genes are also very similar or identical. With one exception discussed below, we did not detect any features of the iterons that could discriminate replication origins of phages of Gram(+) hosts from those of Gram(-) hosts.

We found iterons of the 'simple repeat'-type (rep) as in *oriL* of  $\phi$ SPP1 in roughly one half of the initiator genes. Iterons of the 'relaxed inverted repeat'-type (rIR) as in *ori $\lambda$*  or perfect palindromic iterons (IR) were found in the remaining genes. In some cases a decision whether the iterons were of the 'rep'- or the 'rIR'-type was impossible because simple repeats were found embedded in

larger (relaxed) inverted repeats. E.g. in the prophage gene pp1551 of *P. putida* KT2440 two AAACMG repeats are embedded twice in a larger (relaxed) inverted repeat gagtcgAAACAGaccaAAACCGccgactc, which is found twice 5'-upstream of a prominent 'AT-peak'. There is thus no clear preference for a particular iteron type in the (putative) phage replication origins. The length of 'rep'-type iterons varies between 6 and 34 nt, with a mean of ~10 nt. The length of 'rIR/IR'-type iterons varies between 8 and 41 nt, with a mean of ~20 nt. In most cases, 'rIR'-type iterons with a length of ~20 nt are characterised by short inverted repeats flanking a less conserved central part. These motifs are reminiscent of the DNA-binding sites of transcription factors of the  $\lambda$  CI-type, which bind to DNA as dimers with each monomer binding to one half-side of the binding-site, prior to forming higher-order nucleoprotein complexes by protein-protein oligomerisation [58].

The number of iterons in the individual initiator genes varies between 2 and 11, with a mean of ~5. The *Streptococcus agalactiae* prophage gene sag0559 seems unique in that it contains one single 31 nt long iteron of the 'rIR'-type, which is located 36 nt upstream of a prominent 'AT-peak'. As a rule of thumb, we found higher numbers of iterons for short iteron sequences and lower numbers for longer ones. Invariably, however, we found a clustering of iterons at a distance of ~20 nt from the calculated maximum (ATmax) of the 'AT-peak' (measured from the 3'-end of the proximal iteron) for initiator genes of phages and prophages of Gram(-) hosts. For the genes of phages from Gram(+) hosts, the mean distance is more close to ~30 nt; the significance of this difference with respect to phage host range is unclear. Note that the values for ATmax in AT-plots is a calculated value and does not indicate the position where DNA unwinding could possibly occur. Therefore the distance between ATmax of the 'AT-peak' and the proximal iteron given in Tables C3 – C6 cannot be more but a crude estimate of the spacing between the iteron region and the AT-rich region of a putative phage replication origin. Correspondingly, however, the distances between the 3'-end of the iteron proximal to the AT-rich region and the position where unwinding commences are generally in the order of 10-50 nt in chromosomal and plasmid replication origins [98,290].

There also exists a great variation with respect to iteron spacing. In close to one half of the initiator genes we found a spacing of exactly one helical turn between individual iterons (centre-to-centre) when they are part of the iteron cluster 5'-adjacent to the 'AT-peak'. Such a regular spacing of 1 helical turn is particularly observed for short iteron sequences. In the other half of the (putative) initiator genes we found an iteron spacing of exactly 2, 3, and 4 helical turns, less frequently a regular spacing deviating from complete helical turns, and only in a

**Table C6** (Putative) replication origins of Gram(+)-specific phages; part B.

(pro)phage	gene	nt	AT	consensus sequence of iteron(s)	nt	#	groups	type	spacing	pos.	d
<i>L. innocua</i>	lin2412	933	445	AC <sub>/T</sub> AA <sub>/T</sub> CG <sub>/A</sub> DHHD	10	11.5	1+1+2 +6.5+1	rep	1 t	5'+3' 1)	33
<i>L. innocua</i>	lin0086	912	445	THABAA <sub>/cA</sub> / <sub>c</sub> BDD	10	8	2+6	rep	1 t	5' 1)	31
				634	THABAA <sub>/cA</sub> / <sub>c</sub> BDD	10	2	1+1	rep	-	3'
<i>L. monocyt.</i>	lmo2317	975	493	A <sub>/T</sub> G <sub>/A</sub> TDA <sub>/G</sub> HG	7	7	7	rep	≤1 t	5'	23
φ SM1	gp19	798	475	GTAAC	10	2	2+0.5+0.5	IR	1 t	5'	39
				GTTAC	5	2					
<i>B. cereus</i>	bc0955	864	343	not detected	-	-	-	-	-	-	-
<i>B. anthracis</i>	ba0933	861	343	not detected	-	-	-	-	-	-	-
				523	-	-	-	-	-	-	-
<i>B. anthracis</i>	ba3331	741	339	GAA	24/22	2	2	rIR	-	3'	-
				520							
pφ Lp1	orf20	933	459	ATTAAC	18	3	1+2	rIR	3 t	5'	35
				ATTAAT							
pφ Lp2	orf20	978	337	AGTAGGYCARCTGGTT	16	2	2	rep	32 nt	3'	30
				487							
<i>L. innocua</i>	lin1244	918	411	CA	16/25	2	2	IR	16 nt	3'	16
<i>S. agalactiae</i>	sag0559	861	442	TGGTGTACAAA	31	1	1	rIR	-	5'	36
<i>C. acetobu.</i>	cac1934	840	-	GTAAC	10	4	4	rIR	≤1 t	-	-
				GTTAC							
φ P335	orf11	813	556	ATTGGACA	15	1	1+2x0.5	IR	1 t	5'	137
				GTCCAA							
φ bIL309	orf14	777	520	ATTGGACA	15	1	1+2x0.5	IR	1 t	5'	101
				GTCCAA							
φ Tuc2009	orf16	783	259	CAACCA	6	4	4	rep	1 t	3'	82
				535							
φ bIL285	orf16	783	274	CAACCA	6	4	4	rep	1 t	3'	67
				538							
φ bIL286	bIL286p16	840	505	ATTAA	15/18	2	1+1	IR	-	5'	33
φ 11	phi11_15	807	478	GTAAC	10/5	2/2	0.5+2.5	IR	1 t	5'	53
				GTTAC							
φ 7201	orf4	807	430	ACG	34	2	2	rep <sup>2)</sup>	4 t	5'	34
φ BK5-T	orf49	810	385	T	16	4	4	rIR	6 nt / 2 t	5'	24

Legend to Tables C5 + C6 . nt = length of gene; AT = AT<sub>max</sub>; nt = length of iteron; number of iterons; groups = clustering of iterons; type = iteron type (rep: repeat, IR: inverted repeat, rIR: relaxed inverted repeat); spacing = spacing between iterons, indicated as distance (nt) between outside bases or helical turns (t) calculated from the iteron middle; pos. = position of iterons/major iteron cluster relative to AT<sub>max</sub> ('AT-peak'); d = distance (nt) between the proximal iteron and AT<sub>max</sub>. Inverted repeats in the iterons are indicated by boxes. Abbreviations used in the 'consensus' column: R = A or G; Y = C or T; M = A or C; K = G or T; S = C or G; W = A or T; H = A or T or C; D = A or G or T; B = G or C or T; V = A or C or G; N = any. 1) major cluster of iterons 5' to AT<sub>max</sub>. 2) 2 repeats, each containing 2 16 nt long palindromes with the order > < < >. DNA sequences for Table C5 were taken from: φSPP1 [NP\_690731], φmv4 [AAG31328], φr1t [NP\_695039], φ3626 [NP\_612863], φ13 [NP\_803370], φSLT [NP\_075484], φETA [BAA97608], φLL-H [AAL77546], φ315.1 [NP\_664495], *S. pyogenes* [NP\_802421], *S. pyogenes* [NP\_607824], *L. gasseri* [ZP\_0004-6421], pφ LambdaBa04 [AAP27557], φTP901-1 [NP\_112676], φPV83 [NP\_061610], pφ LambdaBa02 [NP\_846360], φBC6A52 [NP\_832321], φPVL [NP\_058485], φN315 [NP\_835538], and φA118 [NP\_463514]. DNA sequences for Table C6 were taken from: *L. innocua* [NP\_471742], *L. innocua* [NP\_469432], *L. monocytogenes* [NP\_465841], φSM1 [NP\_862858], *B. cereus* [NP\_830741], *B. anthracis* [NP\_843439], *B. anthracis* [NP\_845619], pφ Lp1 [NP\_784408], pφ Lp2 [NP\_785894], *L. innocua* [NP\_470581], *S. agalactiae* [NP\_687588], *C. acetobutylicum* [NP\_348556], φP335 [NP\_839902], φbIL309 [NP\_076709], φTuc2009 [NP\_108693], φbIL285 [NP\_076588], φbIL286 [NP\_076650], φ11 [NP\_803268], φ7201 [NP\_038304], and φBK5-T [NP\_116541].

few cases a completely irregular spacing. Despite many attempts to resolve the stoichiometry and the three-dimensional structure of 'open complexes' [228] formed by initiator proteins with their cognate replication origins, their exact architecture(s) still remain enigmatic. It is even not known whether there exists a generalised architecture for 'open complexes' but it is assumed that a regular spacing of iterons within origins reflects a highly ordered structure of the 'open complex'. Based on electron microscopic studies, the Kornberg lab proposed a model for the 'open complex' formed by *E. coli* DnaA and *oriC* in which ~200 bp *oriC* DNA are wrapped around an ellipsoid DnaA core [86]. The great variation in iteron length and spacing found here in the set of 80 initiator genes leads us to speculate that it is more likely that 'open complexes' formed by phage initiators at their cognate origins have a three-dimensional structure resembling the RecA-filament [112], as was recently proposed also for the *E. coli* DnaA-*oriC* complex by Messer [290].

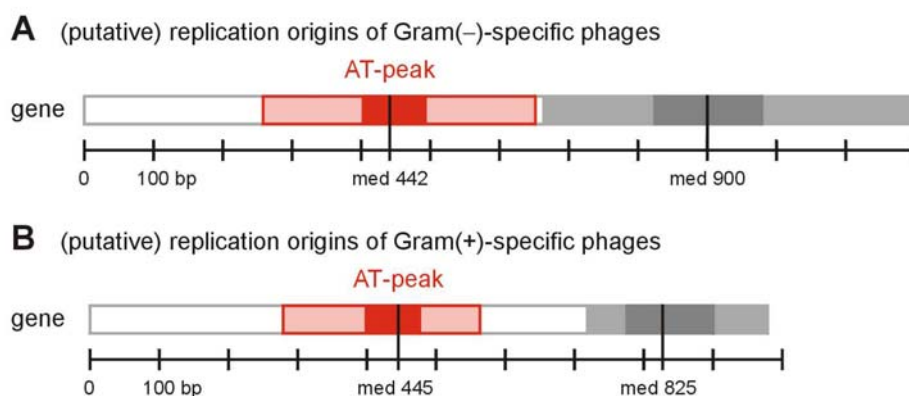
Alonso and co-workers detected a threefold repeated sequence AACAAATGA within the AT-rich region of the  $\phi$ SPP1 *oriL* for which no function is known [12]. During our search for iterons in the set of 80 initiator genes we found that they contain – in addition to iterons – gene-specific repeats at a positions corresponding to those in  $\phi$ SPP1 *oriL*. In 69 of the 80 genes we found several small repeats with the sequence AAGAA, AAGA, or AGAA in close neighbourhood of the 'AT-peak', in many cases in addition to the gene-specific repeats mentioned above. These small AAGAA, AAGA, and AGAA repeats were also found on the complementary strand, albeit less frequently. Because in 22 of the 69 genes containing such small repeats these are found as part of larger, gene-specific repeats their presence might be meaningful in all 69 genes despite their short length. Invariably, the gene-specific repeats and the small repeats – we tentatively like to coin both 'secondary repeats' – are located closer to AT<sub>max</sub> of the 'AT-peak' than the iterons, and mostly on both sides of AT<sub>max</sub>. A straightforward (trivial) explanation for the existence of these AT-rich 'secondary repeats' in all but one of the 80 phage and prophage initiator genes is that they have evolved by micro-duplications and conserved simply because they contribute positively to origin unwinding. However, the 13mer repeats present within the AT-rich region of the chromosomal replication origin *oriC* of *E. coli* and related species overlap with recently detected ssDNA-binding sites for the initiator protein DnaA [423], which was long before only known to bind to the dsDNA *oriC* iterons, the DnaA boxes [290]. Therefore, the possibility cannot be excluded that the 'secondary repeats' within the phage replication origins represent preferential entry sites – on dsDNA, or ssDNA in the unwound origin – for one or more of the proteins

involved in the initiation of replication, i.e. initiator, SSB, helicase loader(s), helicase, or primase. Clearly, the preferential though not exclusive arrangement of these 'secondary repeats' on one strand adds asymmetry to the putative replication origins, which might be instrumental for the choice between unidirectional and/or bidirectional replication (a detailed listing of all 'secondary repeats' is available from the authors upon request).

Our search for structural elements of a replication origin in the set of 80 (putative) phage initiator genes revealed that the origin structures determined experimentally for the phages  $\lambda$  and  $\phi$ SPP1 are well conserved in the majority of the genes analysed. Therefore,  $\lambda$  and  $\phi$ SPP1 can be safely considered valid model systems for phage replicons that replicate by  $\theta$ DR.

Although the localisation of the replication origin within the initiator gene seems to be a common feature among phage replicons that replicate by  $\theta$ DR, also some plasmid replicons with the same characteristic have been described. Cohen and co-workers identified the replication origin of the 17 kb linear plasmid pSLA2 of *Streptomyces rochei* within its *rep1* initiator gene [71]. With 864 bp, the *rep1* gene has the typical length of a phage initiator gene, it contains a centrally located 'AT-peak', three (+ one degenerate) 21 bp long iterons of the 'rep'-type with regular spacing (2 helical turns) 5'-upstream of the 'AT-peak', and secondary repeats on either side of the 'AT-peak' (data not shown). Despite weak sequence similarity, the replication origin of the 12 kb *Streptomyces clavuligerus* plasmid pSCL shares the arrangement of structural elements with pSLA2. Interestingly, both replicons can drive bidirectional  $\theta$ DR of linear and circular forms of their plasmids [411,70], and encode in addition a DnaB<sub>Eco</sub>-type helicase [NP\_0442-86]. Although the *rep1* initiator protein of pSLA2 lacks similarity with any known or putative phage initiator protein, the weak similarity (~30% ident. res.) of the (putative) initiator of pSCL with the N-termini of L17 of  $\phi$ P27 and orf39 of  $\phi$ V (32% ident. res.) suggests that these plasmid replicons and phage replicons have a common evolutionary origin. The identification of the *rep2* gene located downstream of *rep1* in pSLA2 as DnaB<sub>Eco</sub>-type helicase supports this hypothesis (BRM Section 3.) [71].

From our analysis one can derive a 'set of rules' which could – in the case of initiator genes – help to assign prophage gene functions in bacterial genomes with higher reliability than can be achieved by mere protein comparison. A candidate gene of ~900 bp / 300 res. in length (Fig. C5) with N- or C-terminal similarity to at least one of the 80 (putative) initiator protein sequences (Section C3.1.2.) should be considered as phage replication initiator containing the replication origin of this phage if: i. a prominent 'AT-peak' is located in its middle, ii. it contains several repeated sequences 5'-adjacent



**Fig. C5** (Putative) replication origins of lambdoid phages.

**A** Gram(-)-specific phages: summary of the 40 sequences from Tables C3 + C4 . **B** Gram(+)-specific phages: summary of the 40 sequences from Tables C5 + C6 . The length variation of the genes is indicated: light grey box: upper and lower quartil, respectively; dark grey box: 50% region; the median (med) is given as discrete value. AT-peaks were determined using the spin v1.1 software from the Staden Package [425], using a window of 41 nt. The variation of the position of the AT-peak in the genes is indicated: light red box: upper and lower quartil, respectively; red box: 50% region; the median (med) is given as discrete value.

to the 'AT-peak', and iii. it contains 'secondary repeats' to both sides of  $AT_{max}$ . These rules might even be translatable into robust algorithms for an automated search for prophage replication genes and origins in genomic sequences. But we expect from the great variability already apparent in our limited set of sequences that a closer inspection of candidate sequences by eye would still be necessary. Computer-aided searches for (pro)phage replication origins could profit from GC-skew analyses, in addition [152,375].

As examples, we discuss the *Shigella flexneri* prophage gene sf0873 and the *Bordetella bronchiseptica* prophage gene bb1683. The sf0873 protein sequence is highly similar to the *E. coli* cp $\phi$  CP-933R putative initiator protein encoded by gene z2397 (93% ident. res.) and to the YdaU initiator protein of the *E. coli* Rac prophage (65% ident. res. in the C-terminus). With only 633 bp, the sf0873 gene is considerably shorter than most other origin-containing phage initiator genes. It contains 2 prominent 'AT-peaks', one located approximately in its middle. However, no iterons or 'secondary repeats' could be detected in the sequence. In addition, the gene is not found in a context of other replication genes (BRM Section 3.1.). Therefore, the indications for a replication origin within this gene are weak. In contrast, the 957 bp long gene bb1683 of *B. bronchiseptica* contains a very prominent 'AT-peak', preceded by a cluster of 7 iterons, and 'secondary repeats' to both sides of  $AT_{max}$ . Although the translated sequence has only moderate similarity to the  $\phi$ P27 ini-

tiator L17 (37% ident. res. in the N-terminal 75 res.), this gene should be considered a phage replication initiator gene containing the replication origin.

**Other bacteriophage origins for  $\theta$ DR.** Very recently, the  $\phi$ PY54 replication origin for  $\theta$ DR could be located in a segment of 212 bp from the 3'-part of the *repA* gene encoding the multifunctional helicase-primase protein [507]. A sub-segment of 54 bp contains four *dam* methylation sites – 3 as parts of short 6 bp repeats – and an AT-rich region with additional (different) repeats. These structural elements are conserved with respect to sequence and location in the related phages  $\phi$ N15 and  $\phi$ KO2 [365]. The coding capacity of the origin region is not important for RepA function, and the origin is functional if RepA is provided *in trans* [507].

The replication origin *oriI* of *E. coli* phage P4 resembles plasmid origins more than any known phage origin with respect to complexity.  $\phi$ P4 *oriI* has been defined as a bi-partite structure, consisting of *oriI* approximately 4.5 kb upstream of the  $\alpha$  gene and the *crr* region downstream of  $\alpha$  ([123]; reviewed in [47]). Replication is initiated at *oriI*, but *oriI* and *crr* are both required for origin (*oriI*) function. The relative orientation of *oriI* to *crr* is important but the distance between both sites is not. *OriI* contains 6 'type I' iterons GGTGAACAGA/T to which  $\alpha$  protein binds for initiation; the 'type I' iterons are arranged in irregularly alternating orientations but with helical phasing to both sides of an AT-rich region in a DNA segment of 123 bp in length. *Crr* con-

**Table C7** (Putative) replication origins for  $\theta$ DR of other (non-lambdoid) phages.

phage	gene	location	AT	consensus sequence of iteron(s)	nt	#	groups	type	pos.	d	
$\phi$ adh	orf771	3'	ig	+	AGTGTAG <sub>/TG</sub> GTT	11	4	2+ <u>2</u>	rep	5'	45
$\phi$ A2	orf35	3'	ig	+	CGGGAG <sub>/TG</sub> AT <sub>/A</sub>	9	4	3+ <u>1</u>	rep	5'	18
$\phi$ PSA / $\phi$ 2389	pri	3'	ig	+	HRRATAGTT <sub>/AG</sub>	9	4	1+1+ <u>2</u>	rep	5'	0
$\phi$ DT1	orf36	3'	ig	+	GTTA <sub>/GT</sub> C <sub>/TG</sub>	5	10	1+3+ <u>3</u> +3	rep	5'+3'	22
$\phi$ O1205	orf13	3'	ig	+	RGTTA <sub>/TG</sub> C	6	16	4+4+ <u>4</u> +4	rep	5'+3'	23
$\phi$ Sfi11	orf504	3'	ig	+	RGTT <sub>/AA</sub> A <sub>/TG</sub> C <sub>/TG</sub>	6	11	4+4+3	rep	5'+3'	24
$\phi$ 31	primase	3'	ig	+	T <sub>/C</sub> GTTCCA	7	7	3+ <u>2</u> +2	rep	5'+3'	9
$\phi$ bIL310	orf24	3'	ig	+	GT <sub>/GT</sub> A <sub>/A</sub> C <sub>/TA</sub>	5	9	3+4+ <u>2</u>	rep	5'	12
$\phi$ 105	orf11	3'	ig	+	GTT <sub>/CA</sub> A <sub>/CG</sub> C <sub>/GTA</sub>	5	9	<u>8</u> +1	rep	5'+3'	30
$\phi$ P1	repL	int.	-	+	<u>GTCA</u> R <sub>3</sub> <u>TGAC</u> <sub>/T</sub>	11	2.5	2+0.5	rIR	5'	3

gene = gene upstream of (putative) origin; location = '3' indicates location 3'-downstream of 'gene', 'int.' indicates location within coding region of 'gene'; 'ig' indicates location in intergenic region. In the 'groups' column, underlining indicates the iteron group closest to AT<sub>max</sub>. Other abbreviations and criteria as in legend to Table C6. Sequences were taken from:  $\phi$ adh [NC\_000896],  $\phi$ A2 [NC\_004112],  $\phi$ PSA /  $\phi$ 2389 [NC\_003291],  $\phi$ DT1 [NC\_002072],  $\phi$ O1205 [NC\_004303],  $\phi$ Sfi11 [NC\_002214],  $\phi$ 31 [LBA292531],  $\phi$ bIL310 [NC\_002669],  $\phi$ 105 [NC\_004167], and  $\phi$ P1 [NC\_005856].

tains two direct repeats each 120 bp long. A number of (somewhat more degenerate) 'type 1' iterons are found in the *crr* repeats, less regularly spaced as in *ori1*, but also with alternating orientations. Origin function is diminished but not completely suspended upon deletion of one of the *crr* repeats.  $\phi$ P4 contains a second, cryptic replication origin, *ori2*, located in the 5'-half of the  $\alpha$  gene [451]. The 36 bp long *ori2* sequence together with *crr* has been shown to confer replication proficiency to a test plasmid. *Ori2* lacks type 1 iterons and  $\alpha$  protein did not bind to it *in vitro* [451].

In Section C3.1.2., we will discuss the 'initiator domain' of *Lactobacillus gasseri* phage adh and other phage-encoded helicase-primase enzymes that are related to the  $\phi$ P4  $\alpha$  protein. The replication origin of  $\phi$ adh was mapped to the intergenic region downstream of orf771 (encoding the  $\phi$ P4 $\alpha$ -type primase-helicase) [4]. Applying the criteria developed for the lambdoid phages (see above), we could readily detect (putative) origin structures downstream of the  $\phi$ P4 $\alpha$ -type helicase gene in all phage genomes of this group (Table C7) [124,426,303,268]. The detected iterons are invariably of the 'direct repeat' type and in several cases found on both sides of the AT-rich region, in this respect resembling the  $\phi$ P4 *ori1* configuration.

The phage P1 origin for theta( $\theta$ )-mode replication in the lytic cycle is located within the *repL* initiator gene as in the lambdoid phages (Table C7) [82].

### C2.3. Other bacteriophage replication origins

The genome of the *Lactococcus lactis* phage c2 (22.2 kb) is organised in two divergently transcribed blocks of – early and late – genes separated by an intergenic region. This intergenic region, designated *ori*, was shown to support the replication of a selectable plasmid in *L. lactis* in the absence of phage-encoded proteins [469]. By analysing  $\phi$ c2 replication intermediates by two-dimensional gel electrophoresis, it could be shown that replication is initiated in the *ori* region *in vivo* and occurs in the  $\theta$ -mode [60]. Deletion of the early promoter P<sub>E1</sub> in the *ori* region abolished origin function. Introduction of mutations into P<sub>E1</sub> or replacement of P<sub>E1</sub> with an unrelated but functional promoter did not abolish replication. Replacement of the P<sub>E1</sub> transcript template sequence by an unrelated sequence with a similar G+C content abolished replication, showing that the sequence encoding the transcript is essential for origin function. [394]. Phage c2 may thus represent the first example for a phage replicon that – like the ColE1-type plasmids – entirely depends on host functions for replication. Although it is presently not known to which extent possible secondary structures of the P<sub>E1</sub> transcript are important – and whether they are processed in the lactococcal host – we may assume that  $\phi$ c2 replicates by the tDR mechanism. Unlike ColE1 however, replication of  $\phi$ c2 occurs in a *polA* host and therefore does not require host PolA for transcript elongation [359]. We note that although the  $\phi$ c2 genome lacks replication genes, several (putative) recombination genes



could be readily identified (Section C3.6.2.). Also lactococcal phages of the  $\phi$ bIL67 and  $\phi$ 923-type replicate by tDR like  $\phi$ c2. All three phage groups contain a non-homologous  $P_{E1}$  transcript template embedded in highly homologous PE1 and PE2 regions [359]. Also the origin regions of  $\phi$ bIL67 and  $\phi$ 923 can sustain plasmid replication in the absence of phage proteins but chimeric origins were inactive. The latter observation suggests that the particular secondary structure of each  $P_{E1}$  transcript is important.

Hillier and co-workers have identified the replication origin of *L. lactis* phage sk1 by its property to drive replication of a selectable plasmid; the origin is located in the intergenic region between orfs 46 and 47 and includes the first 179 residues of orf 47 [68]. Although the intergenic region contains an AT-rich stretch and also short direct repeats, a closer analysis of this origin region is warranted: at present, neither the replication mechanism nor a putative initiator protein are known.  $\phi$ sk1 orf43 encodes a putative replication protein, but it may rather be a variant of the plasmid-encoded RepA-type or  $\phi$ PG1 gp59 F4-type helicases than a polymerase subunit as proposed by Schouler *et al.* [398].

A similar unsatisfactory situation exists for the replication origin of the *Leptospira biflexa* (Spirochaetes) phage LE1. Girons and co-workers could show that the  $\phi$ LE1 prophage exists as plasmid in its host, and the replication origin could be located on a ~2 kb fragment encompassing orfs 3–5. An AT-rich region and a number of direct repeats could be identified, but the replication mechanism and a putative initiator are not known [144].

### C3. Bacteriophage replication proteins

Having set the stage with a discussion of the different molecular mechanisms driving phage replicon propagation (BRM Section 2.) and origin structures in the preceding section (Section C2.) we will proceed here with the discussion of the bewildering variety of phage-encoded replication proteins. A discussion of the primosomal proteins – initiators, helicase loaders, helicases, and primases – will be followed by a discussion of replisomal proteins – DNA polymerases and their accessory proteins. We will complete this section with a survey of phage-encoded single-strand DNA binding proteins (SSBs) and recombination proteins. Although not usually considered in the context of replication, there is compelling evidence that recombination is closely interlocked with the replication of different types of phages replicons (see BRM Section 2.).

As outlined in the introduction, we will first discuss within each chapter of this section the genetic and biochemical characterisation of well-known proteins – with special emphasis on experimentally determined

interactions among them, and the domains responsible for such interaction(s). This part will be accompanied by a discussion of homologous proteins or proteins with significant similarity that were detected by BLAST searches. Unless indicated otherwise, we performed BLAST searches with default settings: blastp; expect = 10.0; word size = 3; BLOSUM62 matrix; existence = 11; extension = 1; low complexity filter = off [5]. BL2seq default settings: blastn; reward = 1; penalty = -2; gap\_x\_dropoff = 50; expect = 10.0; word size = 11; [445]. All searches were performed using the NCBI BLAST- server(s) (<http://www.ncbi.nlm.nih.gov/BLAST/>). To increase the stringency in the discussion of 'protein similarity', we will exclusively refer to '% identical residues (ident. res.)'; for full-length comparisons, we define 'orthologous' or 'identical' for proteins which share >80% ident. res., 'significant similarity' for proteins with 40–80% ident. res., and 'similar' for proteins with >20% ident. res. . In a few cases, it was necessary – and mentioned explicitly in each case – to confine the 'region of similarity' to parts of the proteins, but usually not less than ~ 50. res. . More refined 'data mining' procedures can certainly reveal relationships among proteins that cannot be pinpointed by the 'crude' BLAST approach used here. However, such procedures require thorough discussions and are therefore not appropriate as tools for a review article. In several cases, BLAST results required further *in silico* analyses to gain reliability. In these cases, secondary structure predictions were performed using either the Jpred method or the PHD method, which gave comparable results; both methods claim to reach performances better than 70% [374,88]. The determination of protein tertiary structures is crucial for an understanding of the mechanisms driving replication on the molecular level. However, discussing replication *in toto* and replication modules in particular does not allow us to discuss each single step, and each single protein in minute details. Therefore, we will mention tertiary structures – where known – and refer the reader to the appropriate literature.

#### C3.1. Initiator proteins

Initiation of DNA replication in the 'rolling circle'-mode (RCR) and initiation in the theta-mode ( $\theta$ DR) depend on the positive action of a protein that is usually called 'initiator' – following the nomenclature proposed by Jacob, Brenner and Cuzin in the 'replicon model' [195]. The molecular mechanisms for both initiation reactions are fundamentally different (BRM Sections 2.1. + 2.2.). Reflecting these differences, initiators for RCR and  $\theta$ DR are functionally and structurally unrelated: initiators for RCR can be classified as site-specific topoisomerases, initiators for  $\theta$ DR are related to transcrip-

tion factors. Both types of initiators have one highly specific property in common: they direct primosomal proteins – helicases and primases – to the replication origin, i.e. they function as 'replisome organiser'.

### C3.1.1. Initiator proteins for 'rolling circle' DNA replication (RCR)

Proteins that initiate replication in the 'rolling circle'-mode perform four functions: i. they bind to the replication origin of their cognate replicon and promote duplex unwinding in a small region (nick-site) adjacent to the initiator-binding site(s), ii. they catalyse the breakage of the phosphodiester backbone of one strand in the single-stranded region at a specific site and transfer the 5'-phosphate end to a tyrosine residue within their polypeptide chain – creating a covalent protein-DNA bond – and leave the free 3'-OH end exposed as primer for DNA synthesis, iii. they direct a helicase to the unwound region, thus functioning as 'replisome organiser', and iv. they participate in the termination reaction by catalysing the transfer of the covalently bound 5'-phosphate end of the displaced strand to its 3'-OH end, liberating a covalently closed circular ssDNA molecule. In the known replication systems, these initiators are not involved in the conversion of single-stranded phage genomes into the double-stranded replicative form (RF), or do not participate in the synthesis of the lagging-strand.

Genes encoding RCR-initiators are found in prokaryotic plasmids and phages, but also in archaeal plasmids [276] and plant gemini viruses [429,226]. The RCR-initiators share sequence motifs containing 'active tyrosine' residue(s) involved in the (transient) covalent linkage to DNA with topoisomerases and relaxases (Tra, Mob), proteins that initiate the DNA transfer of conjugative plasmids [334,470]. (Super)families and motifs were defined for the different types of Rep and Mob proteins on the basis of an in-depth sequence comparison [188,227]. From these motifs, we will only consider 'motif 3' containing the catalytically active tyrosine residue(s) in the following; the functional importance of the other motifs remains to be established. Specific features of plasmid-encoded RCR-initiators have been discussed in recent reviews [119,208,98,209]. The molecular mechanisms driving the replication of ssDNA phages have been concisely summarised by Kornberg and Baker [228], and later by Horiuchi [181].

#### Initiator proteins of phages with ssDNA genomes.

Upon entry into the bacterial cell, the closed-circular single-stranded phage genome (+-strand or viral strand) is converted into double-stranded DNA, the replicative form (RF). Subsequently, the RF serves as template for

the synthesis of progeny viral genomes by the 'rolling circle' mechanism [141]. The first stage is entirely dependent on host proteins. The second stage involves a phage-encoded function – the initiator for RCR – and a different subset of host proteins. The development of *in vitro* replication systems in the labs of Jansz, Hurwitz, Denhardt, and Kornberg was instrumental in dissecting the timely order of the individual enzyme-catalysed steps and in identifying the proteins involved [477,478,390,118,117,401,116,110]. These studies culminated in the definition of 'primosomal' and 'replisomal' proteins in the terms of enzymology – making replication of the *E. coli* phage  $\phi$ X174 the best-understood system in the 1970s. Also during these years, the determination of the nucleotide sequence of the  $\phi$ X174 genome by Sanger and co-workers in 1977 marks the onset of the era of genome sequencing [386].

*Gene II protein (gpII) of phage  $\phi$ d.* Meyer and Geider characterised  $\phi$ d gpII biochemically as protein with endonuclease and topoisomerase activity. The endonuclease activity was found to be specific for  $\phi$ d RF DNA (supercoiled dsDNA) because other supercoiled substrates were not cleaved. Also, neither the viral DNA (circular ssDNA) nor relaxed (doubly-closed circular) RF DNA were appropriate substrates. With a maximum of activity at pH 8.5 and 80 mM KCl, gpII introduced a single nick in supercoiled  $\phi$ d RF at low  $Mg^{2+}$  concentrations at a specific site (Section C2.1.) [296]. At higher concentrations of  $Mg^{2+}$ , the topoisomerase activity of gpII became apparent: roughly one half of the supercoiled RF substrate was converted to the relaxed (doubly-closed circular) form [295].

Geider, Meyer and co-workers established an *in vitro* replication system that depended on gpII, and the *E. coli* proteins Rep, SSB, and DNA Pol III holoenzyme for the conversion of  $\phi$ d RF DNA into viral (+)-strand DNA. [138]. In another, partially heterologous *in vitro* system, viral DNA could be generated from RF DNA by the action of gpII,  $\phi$ T7 gene 4 helicase, and  $\phi$ T7 gene 5 DNA polymerase. The generation of closed circular ssDNA by both systems added support to the hypothesis that the cleaving/joining activities of gpII are required at successive steps during  $\phi$ d replication. Although the *in vitro* system using  $\phi$ T7 components also produced closed circular  $\phi$ d ssDNA, the yield was quenched by the strand-switching activity of the  $\phi$ T7 replisome – detected as double-stranded tails of the 'rolling circle' structures [164]. This observation may lead to an answer to the puzzling question why – in all known cases – RCR depends on the *E. coli* Rep helicase *in vivo* and *in vitro*, rather than on the replicative helicase DnaB. DnaB is involved in the coupling of leading- and lagging-strand synthesis as part of the *E. coli* replisome, and the same

Table C8 RCR initiators, part A.

$\phi$ / p $\phi$ / plasmid	gene	res.	ident.	motif 3		$\phi$ / p $\phi$ / plasmid	gene	res.	ident.	motif 3	
<b>1A CTX<math>\phi</math>-group</b>						<b>1B pT181-group</b>					
			cons.	SRI A	WRI NKAALQ KAQL				cons.	SxR F K Q V L	IRI NKKxER R Q
CTX $\phi$	RstA	276	x	SRI	WRI	pT181	RepC	314	x	SNR F	IRI
$\phi$ VGJ $\phi$	orf359	359	43%	SAI	WRI	pCW7	RepN	314	86%	SDR F	IRI
$\phi$ KSF-1 $\phi$	orf1	367	38%	SAI	WRI	pTZ12	Rep	314	83%	SDR F	IRI
$\phi$ VSK	Rep	368	37%	SRI	WRI	pC223	Rep	314	77%	SNR F	IRI
$\phi$ VSKK	VSKKp3	367	38%	SAI	WRI	<i>S. agalactiae</i>	sag0224	332	29%	SEK Q	VRL
$\phi$ Vf12	Vpf402	402	36%	SRI	WRI	pRS2	Rep	319	27%	SEK Q	IRL
$\phi$ Vf33	Vpf402	402	36%	SRI	WRI						
$\phi$ VfO4K68	Vpf402	402	35%	SRI	WRI						
$\phi$ VfO3K6	Vpf402	402	35%	SRI	WRI						
$\phi$ fs1 <sup>1</sup> )	fs1p12	166	33%								
	fs1p13	208	46%	SAI	WRI						
<b>2A <math>\phi</math>X174-group</b>						<b>2B <math>\phi</math>chp1-group</b>					
			cons.	Vx F N W	VAK T	Vx Q F V			cons.	xhx T K	VAR T K
$\phi$ X174	A	513	x	VG F	VAK	VNKKSD	$\phi$ chp1	orf4	399	x	NIF
$\phi$ S13	A	513	96%	VG F	VAK	VNKKSD	$\phi$ MH2K	Vp4	315	30%	SAS
$\phi$ G4	A	554	61%	VG F	VAK	VNKKSD	$\phi$ 3	orf4	315	30%	SAG
$\phi$ K ( <i>E. coli</i> )	A	494	44%	VAW	VTK	VAKQSD	$\phi$ 2	orf4	336	29%	SAG
$\phi$ $\alpha$ 3	A	494	43%	VAW	VTK	VAKQSD	$\phi$ CPAR39	p7	327	27%	SAG
$\phi$ chp1	orf4	399	25%	NIF	VAR	VQKKFV	$\phi$ CPG1	p9	263	28%	SAG
							$\phi$ 4	gene 2	320	28%	SAN
							$\phi$ X174	A	513	25%	VG F

holds for the gene 4 helicase in the  $\phi$ T7 replisome. It would be highly interesting to learn whether the *E. coli* Rep helicase – or PcrA, its orthologue in *B. subtilis* – is inefficient for the coupling of leading- and lagging-strand DNA synthesis. A reason for such a hypothetical inefficiency could be the inability of Rep helicase to recruit the DnaG primase.

The binding of gpII to the  $\phi$ fd replication origin was studied to great detail in the lab of Horiuchi. Initially, it was found by filter-binding assays that gpII binds to  $\phi$ fd DNA in the superhelical or linear form. The gpII binding-site(s) could be mapped close to the site of initiation of DNA synthesis, overlapping the nick-site. The nick-site itself was found to be dispensable for gpII binding [180]. Using more advanced techniques – e.g. EMSA, DNase I footprinting and methylation interference – two distinct complexes between  $\phi$ fd origin DNA and gpII could be defined: complex I forms at low gpII concentrations at two inverted repeats (~25 bp) to the right of the nick-site (binding-sites  $\beta + \gamma$ ; Fig. C1), complex II requires approximately twice the amount of gpII, and the bound region comprises of ~40 bp, including the nick-site (binding-sites  $\alpha + \delta$ ). Binding of gpII

to the origin in complex II occurs to the inverted repeat ( $\beta + \gamma$ ) as in complex I, to a third repeat at the right of the inverted repeat ( $\delta$ ) and to the nick-site to the left without sequence specificity ( $\alpha$ ). This latter observation indicates that protein-protein interactions among gpII protomers stabilise the entire complex, allowing one protomer to cleave the nick-site [151]. Also it was found that complex I induces bending of the origin DNA, an effect that was even more pronounced for complex II [170]. Analysis of gpII complexes formed on supercoiled  $\phi$ fd RF DNA by the *in vitro* KMnO<sub>4</sub>-footprinting technique revealed that complex II induces a local unwinding in the nick-site region as a prerequisite for efficient nicking [170].

In 1999, Horiuchi and co-workers could isolate the covalent complex between gpII and an artificial substrate mimicking nicked-origin DNA. By peptide sequencing, the covalent link could be traced down to Y<sub>197</sub>. Substituting Y<sub>197</sub> by phenylalanine resulted in a protein with wildtype DNA-binding properties *in vitro*. Also, the mutant protein induced origin bending like the wildtype protein but did not show any origin-nicking or topoisomerase activity [10]. The 'active' tyrosine of  $\phi$ fd

Table C9 RCR initiators, part B.

$\phi$ / p $\phi$ / plasmid	gene	res.	ident.	motif 3	$\phi$ / p $\phi$ / plasmid	gene	res.	ident.	motif 3
<b>3 pTLC-group</b>					<b>4 <math>\phi</math>P2-group</b>				
			cons.	Rxx x hxh Q S				cons.	Axx P
				xKxxEx G S					IAK V T
				TKHDEI LKHVEV LKHVEL LKHVEL LKHVEL WGKGST SKYDEV					IxKNID ISKNID ISKNID ISKNID ISKNID ISKNID ISKNID ISKNID IGKNLD IAKNID ISKNIN LSKNID ISKNID ISKNID
pTLC	Cri	522	x	RSF E LCI	$\phi$ P2	A	761	x	AAG IAK ISKNID
$\phi$ lKe	gpII	421	34%	RHR T LVA	$\phi$ L-413C	gpA	761	97%	AAG IAK ISKNID
$\phi$ M13	gpII	410	32%	RHR T LVA	$\phi$ W $\phi$	gpA	761	96%	AAG IAK ISKNID
$\phi$ f1	gpII	410	32%	RHR T LVA	<i>Y. pseudo.</i> p $\phi$	yptb1766	760	56%	AAG IAK IAKNID
$\phi$ fd	gpII	410	32%	RHR T LVA	$\phi$ PSP3	gp36	708	39%	ATG VAK ISKNID
pTA144up	RepC	336	24%	QFK G STL	$\phi$ 186	A	694	37%	ATG VAK ISKNID
$\phi$ l2-2	gpII	344	24%	RRW S IKA	$\phi$ Fels-2	stm2729	809	34%	PTS IAK ISKNID
plasmid RCR initiator-groups lacking phage homologues (5+6 taken from [98])					p $\phi$ CP-933T	z2978	940	31%	PTX IAT IGKNLD
<b>5</b> pLS1 / pMV158	RepB	210	x	NVE N MYL	$\phi$ HP2 / $\phi$ HP1	Rep	775	31%	ATA IAK IAKNID
<b>6</b> pC194	Rep	232	x	ELY E MAK	$\phi$ PM2	p12	634	33%	AVG IAK ISKNIN
<b>7</b> pJV1	Rep	528	x	DDV A LIE	$\phi$ fs2	orf716	716	33%	AVG VAK LSKNID
				LTKNQD	$\phi$ V86 <sup>2</sup> )	-	[532]	36%	ATG IAK ISKNID
					$\phi$ K139	Rep	800	31%	ATG IAK ISKNID

The position of the signature motif 3 [188] containing the conserved active tyrosine are: in pT181 res. 184-198,  $\phi$ f1 res. 190-204, in pLS1 res. 92-106, in pJV1 res. 268-282, in pC194 res. 207-221, in  $\phi$ fs2 res. 437-453. The active tyrosine residues are boxed; bold letters indicate conservation of residues in motif 3 among subgroups. ident.: % identical residues. cons.: consensus sequence. Sequences were taken from: CTX $\phi$  [AAN06952], VGJ  $\phi$  [NP\_835472]; KSF-1  $\phi$  [YP\_087070];  $\phi$ VSK [NP\_542355];  $\phi$ VSKK [NP\_536619];  $\phi$ Vf12 [YP\_031686];  $\phi$ Vf33 [YP\_031679];  $\phi$ VfO4K68 [NP\_059541];  $\phi$ VfO3K6 [NP\_059531];  $\phi$ fs1 [NP\_695201, NP\_695202],  $\phi$ X174 [NP\_040703],  $\phi$ S13 [AAG29971],  $\phi$ G4 [AAL51008],  $\phi$ K (*E. coli*) [NP\_043942],  $\phi$  $\alpha$ 3 [NP\_039590],  $\phi$ chp1 [NP\_044320, NP\_044321],  $\phi$ MH2K [NP\_073537],  $\phi$ 3 [YP\_022485],  $\phi$ 2 [NP\_054653],  $\phi$ CPAR39 [NP\_063900],  $\phi$ CPG1 [NP\_510879],  $\phi$ 4 [NP\_598335],  $\phi$ Ike [NP\_040570],  $\phi$ M13 [NP\_510885],  $\phi$ f1 [CAA23876],  $\phi$ fd [AAA32303],  $\phi$ I2-2 [NP\_039615],  $\phi$ P2 [NP\_046795],  $\phi$ L-413C [NP\_839887], W $\phi$  [NP\_878237]; *Y. pseudotuberculosis* p $\phi$  [YP\_070292],  $\phi$ PSP3 [NP\_958093],  $\phi$ 186 [NP\_052289],  $\phi$ PM2 [NP\_049896],  $\phi$  Fels-2 [NP\_461656], *E. coli* p $\phi$  CP-933T [NP\_288356],  $\phi$ HP2 [NP\_536816],  $\phi$ HP1 [NP\_043478],  $\phi$ K139 [NP\_536641],  $\phi$ V86 [AAB66811],  $\phi$ fs2 [NP\_047364], pLS1 [NP\_040421], pJV1 [NP\_044352], pC194 [CAA24585], pTA144up [NP\_052248], pTLC [NP\_862693], pT181 [NP\_040469], pCW7 [AAA26669], pTZ12 [AAA72571], pC223 [CAA30291], *S. agalactiae* [NP\_687259], pRS2 [NP\_443751]. Table C8: 1) Overlapping orfs encode the proteins fs1p12 and fs1p23. Both predicted proteins show significant similarity to CTX $\phi$  RstA in the N- and C-terminus, respectively. The splitting may be a mere sequencing artefact (reading-frame shift), therefore. Table C9: 2) the length of the unnamed  $\phi$ V86 protein is given in brackets because the sequence may be only partly correct (see Legend to Fig. C6 for details).

gpII is embedded in the motif 3-variant RHRTLVA<sub>197</sub>LKHVEL (res. 190-204) which is conserved in other initiators (see below). Neither the domain responsible for DNA- binding, nor the interaction domain with the host Rep helicase have yet been determined for gpII, and also a crystal structure is not available.

*Protein A of phage  $\phi$ X174.* Synthesis of  $\phi$ X174 viral-strand circles *in vitro* requires gene A protein and in addition the *E. coli* proteins Rep, SSB, and DNA Pol III holoenzyme [118]. During replication, A protein nicks  $\phi$ X174 RF DNA and binds covalently to it; a complex with 1:1 stoichiometry could be isolated [116,110,50]. It was shown that the DNA bound covalently to A corresponds to (+)-strand DNA directly adjacent to the replication origin, and that the bond involved tyrosine residues of A protein [377]. By peptide sequencing it could be shown that either of the two tyrosine residues in the motif YVALKYVNK can bind covalently to DNA [462]. The 'active' tyrosines of  $\phi$ X174 A protein are embedded in the motif 3-variant VGFY<sub>343</sub>VAKY<sub>347</sub>VNKKSD (res. 340-353) which is conserved in other

initiators (see below). Neither the domain responsible for DNA- binding, nor the interaction domain with the host Rep helicase have yet been determined for A protein, and also a crystal structure is not available.

*RCR initiators of ssDNA phages.* The results of BLAST searches suggests the definition of three distinct groups of RCR-initiators of ssDNA phages:

- the  $\phi$ X174/ $\phi$ chp1-group. The orf14 protein of *Chlamydia* sp. phage  $\phi$ chp1 is only weakly similar to  $\phi$ X174 A protein, but all proteins detected by BLAST searches with either 'prototype' as query contain a well conserved motif 3 of the 'two-tyrosine'-type (Table C8; groups 2A + 2B). Historically, studies of  $\phi$ S13 replication preceded the extended studies of  $\phi$ X174 [448].
- the CTX $\phi$ /pT181-group. The proteins of the CTX $\phi$ - and pT181-subgroups are of comparable size but there is no detectable sequence similarity among them (Table C8; groups 1A + 1B). The proteins of the CTX $\phi$ -subgroup contain a conserved motif 3 of

the 'two-tyrosine'-type, while the proteins of the pT181-subgroup contain a single tyrosine in their motif 3. Despite this difference which certainly reflects subtle functional differences of the proteins, the motif 3-variants present in both subgroups are more closely related to each other than to any other known motif 3-variant. The CTX $\phi$ - subgroup can be further divided in three 'branches': i. the CTX $\phi$  RstA-type, ii. the  $\phi$ VSK Rep-type ( $\phi$ VGJ, KSF-1 $\phi$ ,  $\phi$ VSKK,  $\phi$ fs1), and iii. the  $\phi$ Vf12 Vpf402-type ( $\phi$ Vf33,  $\phi$ VfO3K6,  $\phi$ VfO4K68). Within each branch, the proteins are virtually identical, but the similarities to proteins of the other two branches does not exceed 45% (ident. res.) and excludes the N-termini.

iii. the pTLC-group. The Cri protein of plasmid pTLC is similar to the gene II proteins of phages M13, f1, and fd, and all proteins of this group contain a motif 3 variant of the 'one tyrosine'-type.

We have listed in Tables C8 + C9 groups of plasmid- and phage-encoded RCR-initiators by sequence similarity which in all cases corresponded with a conserved variant of motif 3. Groups 1A + 1B and 3 contain proteins encoded by phages and plasmids. Plasmid-encoded initiators belonging to groups 2 or 4 ( $\phi$ P2-group, discussed below) could not be found. Also *vice versa*, phage-encoded initiators with similarity to pLS1/pMV158 (group 5), pC194 (group 6), and pJV1 (group 7) were not detected. This shows that there is no promiscuous exchange of genes for these functionally analogous proteins among the different replicon types. However, a report of Waldor and co-workers suggests that an intricate molecular interplay exists between the cholera-toxin encoding CTX $\phi$  phage (Table C8; group 1A) of *V. cholerae* and the plasmid pTLC (Table C9; group 3) [379].

Recently, Chopin and co-workers have characterised  $\phi$ B5, a filamentous phage with ssDNA genome of a Gram-positive host [77]. However, the putative replication protein orf9 does not bear any resemblance with any of the known initiators discussed above, and the replication modus  $\phi$ B5 remains unknown therefore.

#### Initiator proteins of phages with dsDNA genomes:

**P2 A protein.** Haggård-Ljungquist and co-workers could show that the 76.7 to 91.6% region of bacteriophage  $\phi$ P2 (genome size 33.6 kb, dsDNA) is necessary and sufficient to drive plasmid replication [258]. This region contains – among several orfs with unknown function – the replication genes A and B, and the replication origin located within gene A (pos. 89 $\pm$ 1%). The helicase loader encoded by gene B is discussed in the following chapter, the  $\phi$ P2 replication origin in Section C2.1. . These authors could show that the isolated A gene is sufficient to drive  $\phi$ P2 mini-chromosome repli-

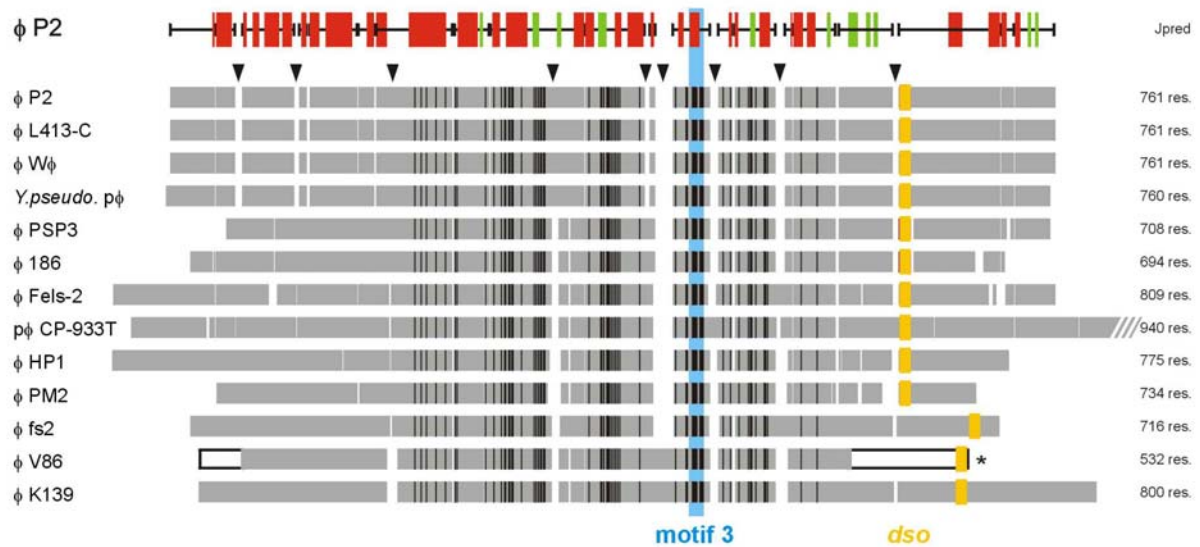
cation. Egan and co-workers could show that gene A protein (CP87) of  $\phi$ 186 is also sufficient to drive plasmid replication [417].

$\phi$ P2 A protein contains an AAGY<sub>450</sub>IAKY<sub>454</sub>ISKNIID motif that corresponds to motif 3 defined by Koonin and Ilyina [227]. Both tyrosine residues could be shown to be essential for RCR of  $\phi$ P2 [325]. Y<sub>454</sub> is apparently more proficient in the initial nicking reaction while the nicking activity of Y<sub>450</sub> can be observed when Y<sub>454</sub> is covalently linked to a specific oligonucleotide. These observations led Haggård-Ljungquist and co-workers to suggest a coupling of termination and re- initiation of  $\phi$ P2 replication by A protein, mediated by the two tyrosine residues. It is known that ~100 res. of the C-terminus are dispensable for  $\phi$ P2 A functioning during replication *in vivo* and *in vitro* [325], but domains involved in DNA-binding and the interaction with B protein or the host helicases – DnaB and Rep [59,43] – have not yet been identified.

(Putative) initiators for RCR with similarity to  $\phi$ P2 A were detected in the genomes of several phages of  $\gamma$ -proteobacterial hosts, and also encoded by prophage genes in the genomes of many  $\gamma$ -proteobacteria (Table C9; group 4). The similarities to  $\phi$ P2 A and the sizes of these proteins vary considerably, and a sequence alignment could only reveal several patches of homology in a 'core region' of ~300 res. in length (Fig. C6) but no true consensus sequence. Interestingly, the alignment suggests the existence of several distinct sub-families, which share a (predicted) comparable secondary structure in the 'core region'. The motif containing the two catalytically important tyrosine residues is well conserved (Table C9; group 4), and located at a corresponding position in the C-terminal part of the 'core region' in all sequences (Fig. C6; light-blue vertical bar). The A proteins of  $\phi$ P2 and  $\phi$ 186, respectively, have both been shown to function as initiators although their similarity is limited. From the above we conclude that all proteins of the  $\phi$ P2-group are *bona fide* initiators (Table C9; group 4).

A distantly related member of the  $\phi$ P2-group of RCR-initiators is the Arp protein encoded by the 1282 bp long phasyl (phage-plasmid hybrid) [140]. With only 285 res. Arp is considerably shorter than  $\phi$ P2 A and a reasonable alignment of Arp with the  $\phi$ P2-group could not be obtained (not shown). However, similarity (27% ident. res.) could be detected for a stretch of ~150 residues within the 'core region' between Arp and the p $\phi$ Fels-2 stm2729 protein, and Arp contains the motif 3-variant IGRYVVGKYISKGIE at a corresponding position.

Plasmid-encoded initiator proteins for RCR with significant similarity to  $\phi$ P2 A could not be detected by BLAST searches, but there is a puzzling similarity of  $\phi$ P2 A protein to the orf2 protein [SWISS YR72\_ECO-



**Fig. C6**  $\phi$ P2 A-type (pro)phage-encoded initiators for RCR.

First line: secondary structure prediction for  $\phi$ P2 A by the Jpred method [88]. Colour code: red = predicted  $\alpha$ -helical region; green = predicted  $\beta$ -stranded region; black line = unstructured/no prediction; the sequence was split into segments to allow comparison with the sequence alignment below. Following lines: sequence alignment of proteins with similarity to  $\phi$ P2 A by 'MultAlin' [84]. For gene names and accession numbers see Table C9. Aligned sequences are shown as grey blocks. Gaps are shown in white. Identical residues in all sequences are indicated by black bars. The vertical light-blue bar indicates the position of motif 3; orange squares mark the position of the (putative) replication origin (*dso*) in the respective gene (see Table C2; Section C2.1. for discussion). The sequence of p $\phi$  CP-933T z2978 protein was shortened at the C-terminus to fit to the figure. The sequences of the homologous  $\phi$ K139 Rep and (unnamed)  $\phi$ V86 proteins are identical on the nucleotide level (99.5%). We therefore tentatively corrected the  $\phi$ V86 sequence at two positions (indicated by open boxes) i. choosing the alternate GUG start codon 40 triplets upstream of the AUG assigned as start in entry AF008938 allows the perfect alignment with the  $\phi$ K139 Rep protein in the N-terminus (the Rep gene also starts with GUG). ii. 2 additional bases inserted shortly after codon 527 (original sequence) result in a reading-frame shift and premature stop, *in silico* correction restores the  $\phi$ K139 Rep reading frame. The sequenced part of the  $\phi$ V86 gene ends after codon 678 (counted from the corrected N-terminus), i.e. one triplet downstream of the putative *dso* (indicated by '\*').

LI] encoded by the *E. coli* retro-element Ec67 [256]. However, the similarity does not exceed 27% ident. res. for the N-terminal ~350 res & does not include motif 3. Therefore, a role for this protein for retron propagation remains speculative.

### C3.1.2. Initiator proteins for DNA replication in the 'theta-mode' ( $\theta$ DR)

Proteins that initiate replication in the  $\theta$ -mode perform two major functions: i. they bind to the replication origin and promote duplex unwinding, and ii. they direct other replication proteins – the helicase loader(s) and/or the replicative helicase – to the exposed single-stranded regions, thus functioning as a 'replisome organiser'.

Initiators bind to the iterons present in the replication origins for  $\theta$ DR adjacent to an AT-rich region (Section C2.2.), and binding occurs – depending on the iteron type – in the monomeric form to non-palindromic iterons and in the dimeric form to palindromic iterons. In addition to binding of initiator protomers to the iterons, unwinding requires protein-protein interactions between two or more protomers. In most cases, the properties for DNA-binding and for oligomerisation are located in

distinct regions of the proteins, which retain their partial function upon separation. Therefore, prokaryotic initiators are composed of – at least – two distinct domains: a DNA-binding domain, and an oligomerisation domain. A third domain contiguous with – or overlapping – the oligomerisation domain is required for the physical interaction with the replicative helicase, or the helicase loader(s).

(Almost) all sequenced bacterial genomes contain *dnaA* genes, and DnaA may be the 'universal' initiator for the replication of chromosomal replicons in bacteria. *E. coli* DnaA binds to dsDNA with its C-terminal domain 4 [376]. The N-terminal domain 1 of DnaA and – in addition – a recently identified structural motif in the central NTP-binding domain 3 are responsible for oligomerisation [473,121]. The physical interaction between DnaA and the replicative helicase DnaB could be shown to involve two domains of either protein: res. 24-86 (domain 1) and 130-148 (domain 3) of DnaA and res. 154-210 ( $\beta$ -fragment) and 1-156 ( $\alpha$ -fragment) of DnaB, respectively [278,404]. It is not known whether *E. coli* DnaA interacts specifically also with the helicase loader, DnaC, but experimental results supporting this possibility are not available (see following chapter). In contrast to the known initiators of bacterial plasmids



and phages, DnaA is a NTP-binding protein, and the proficiency of DnaA for origin unwinding is strictly dependent on the ATP-bound form [423]. The properties and functions of DnaA have been reviewed by Skarstad and Boye [418], and more recently by Messer [290].

The molecular architecture of various initiators of plasmids replicating by the  $\theta$ -mode – including e.g. pPS10 RepA, R6K  $\pi$ , R1 RepA – is comparable to that of DnaA: a N-terminal dimerisation/oligomerisation domain, and a C-terminal DNA binding domain. However, there is no apparent sequence similarity between the plasmid initiators and DnaA, and the structural motifs performing analogous functions are clearly different. A physical interaction between the plasmid initiator and the host replicative helicase was shown for R6K  $\pi$  and pSC101 RepA, and involves similar regions in DnaB as found for DnaA [362,91]. The structures and functions of plasmid initiators have been reviewed by Giraldo [142] and del Solar and colleagues [98].

Besides being essential for replication initiation, DnaA serves as transcription factor in *E. coli* [293]. Most importantly, DnaA acts as repressor of its own gene in a complex regulatory feedback loop involving other factors in addition. Autoregulation is also known for plasmid initiators and serves as part of the copy-number control mechanisms operating in plasmid replicons [98, 97]. A role as transcription factor has not been found for the initiator proteins of bacteriophages –  $\lambda$  O and  $\phi$ SPP1 G38P – and a role of these proteins as auto-repressors can be excluded.

We discuss in detail bacteriophage Lambda O protein, and Gene 38 protein (G38P) of *B. subtilis* phage SPP1. Few other phage-encoded replication initiators of the  $\lambda$  O'-type have attained attention. We could detect by BLAST searches 140 putative  $\lambda$  O'-type initiator sequences from phages or prophages in the sequenced bacterial genomes, and we present (theoretical) evidence that most if not all proteins are *bona fide* initiators for the replication of their cognate (pro)phages. In addition, the RepL initiator for bacteriophage P1 replication in the lytic cycle and the C-terminal 'initiator domain' of phage P4 alpha ( $\alpha$ ) protein and related multifunctional helicases are discussed.

**Bacteriophage Lambda O protein.** Genetic experiments with hybrid phages obtained by *in vivo* crosses between  $\lambda$  and  $\phi$ 80 suggested a dual role for  $\lambda$  O protein as initiator: recognition of the  $\lambda$  replication origin by the N-terminus, and interaction with the  $\lambda$  P helicase loader by the C-terminus [440,133, 134]. The crossover points in these hybrids were mapped to the middle of the O and 15 genes, respectively, and it was demonstrated by complementation assays that only hybrid O proteins containing the  $\lambda$  O N-terminus showed specificity for *ori $\lambda$* . Similar complementation assays showed that a

hybrid O protein carrying the C-terminus of  $\phi$ 80 gene 15 protein could only function in replication with the gene 14 product of  $\phi$ 80. It was concluded that the C-terminus of O protein interacts specifically with P. The physical interaction of O and P was later confirmed biochemically by Zylicz and co-workers [509] (see next chapter).

Matsubara and Tsurimoto analysed purified  $\lambda$  O for its binding specificity to the four *ori $\lambda$*  iterons *in vitro* by exonuclease III protection, DNase I footprinting, and filter-binding assays [458,459]. At low protein concentration, only the inner two iterons were bound, while binding to all four iterons occurred at higher protein concentrations, covering up to ~95 bp of linear *ori $\lambda$*  DNA fragments. The authors could not detect DNA-binding for a mutant O protein lacking 20 res. in the N-terminus, corroborating the results from genetic studies that the O N-terminus is responsible for specific binding to *ori $\lambda$*  (see above). In addition, their results suggested that dimers of O bind to single iterons, i.e. eight protomers per *ori $\lambda$* . Wickner and Zahn confirmed the above results and could show, in addition, that the isolated  $\lambda$  O N-terminus is necessary and sufficient for specific binding to *ori $\lambda$* . The O N-terminal fragment cannot drive  $\lambda$  ( $\lambda$ dv) replication *in vitro* but inhibits this reaction by competing with wildtype O. They could also show that *ori $\lambda$* -bound O protein – but not the isolated O N-terminus – forms a stable complex with P [481].

Detailed structural studies of the  $\lambda$  O-*ori $\lambda$*  complex and of the binding of O to isolated iterons were performed by Blattner and co-workers, employing a battery of *in vitro* techniques [499,501,396,395]. The emergent picture for the interaction of O with its cognate DNA binding-site(s) is: binding of O augments the intrinsic curvature of *ori $\lambda$*  DNA (Section C2.2.), and the binding of O (dimers) already induces bending of the isolated iteron DNA. The 19 bp (imperfect) inverted repeat of the iteron is the minimal sequence required for efficient binding by O, and iteron binding + bending is exclusively dependent on the N-terminal domain of O suggesting that this domain has the property to dimerise. Asymmetrically located protected regions found by footprinting experiments for the O-bound iteron suggest, in addition, its partial distortion. The DNA distortions induced by O-binding in supercoiled *ori $\lambda$*  DNA can be traced well into the AT-rich region. Up to ~84 bp of *ori $\lambda$*  DNA are wrapped around – or partially buried within – a core of oligomerised O protomers with a diameter of ~50 Å.

The four 19 bp iterons of the  $\lambda$  replication origin, *ori $\lambda$* , possess dyad symmetry, are regularly spaced to each other (repeat: 1 helical turn), and located adjacent to the AT-rich region (Section C2.2.). Three iterons are sufficient for efficient replication, as suggested by the normal growth and burst size of  $\lambda$ r93hot5 [299]. For *in*

*in vitro* replication of *oriλ*-based plasmids, however, only the two iterons proximal to the AT-rich region seem to be required [480]. Similar observations were made with the DnaA-binding sites (DnaA boxes) in the *E. coli* replication origin, *oriC*: alterations of the helical phasing of the DnaA boxes had more severe negative effects than inactivation, and these effects were most pronounced for DnaA box R1, directly adjacent to the AT-rich region [239,290]. It should be kept in mind that the initiators bound to the iterons proximal to the AT-rich region are probably not only important for the unwinding reaction but also for the correct positioning of the replicative helicase [474].

When  $\lambda$  O protein remains bound to DNA after departure of the replication forks from *oriλ* it becomes resistant to degradation by ClpP/ClpX protease in the so-called 'inherited replication complex' [149,485,471,446]. Initiation of new replication rounds from the inherited replication complex *in vivo* does not require *de novo* synthesis of O but still depends on host DnaA-dependent transcriptional activation of  $p_R$ , which triggers chaperone action for the liberation of the DnaB helicase from the  $\lambda$ O- $\lambda$ P-DnaB complex (see below). It is not entirely clear at present whether and how the switch from  $\theta$ DR to  $\sigma$ DR during  $\lambda$  replication is triggered by the inherited replication complex.

### Gene 38 protein (G38P) of *B. subtilis* phage SPP1.

Studies of the replication of phage SPP1 date back to the 1970s, and could reveal over the years the essential nature of the products of genes 38, 39 and 40.  $\phi$ SPP1 replication is independent of the host DnaA, DnaC helicase, the DnaB, DnaD and DnaI helicase loaders, but requires the host DnaG primase, host gyrase, and DNA polymerase III holoenzyme [55,341]. G40P was identified as the cognate helicase for  $\phi$ SPP1 replication (Section C3.3.), and G39P as the cognate helicase loader (Section C3.2.). G38P binds specifically two discrete regions in the SPP1 genome: *oriL*, the replication origin located within gene 38, and *oriR*, an origin-like structure at a distance of ~32 kb from *oriL* (Section C2.2.) [341]. Alonso and co-workers could measure specific, cooperative, and tight binding of G38P to *oriL* ( $K_{app} = 1$  nm) and to *oriR* ( $K_{app} = 4$  nm) *in vitro*, and protection of the AB iterons present in *oriL* and *oriR* – but not of the AT-rich region – by DNase I footprinting experiments [297]. As discussed in the following chapter, G38P does by itself not interact with its cognate helicase: when bound to the replication origin *oriL*, G38P interacts specifically with G39P in the G39P-G40P-ATP complex, which results in G40P binding to the origin [13]. The interaction domains of G38P and G39P have not yet been determined.

The generally accepted model for the complex between the initiator and its cognate origin DNA is that of

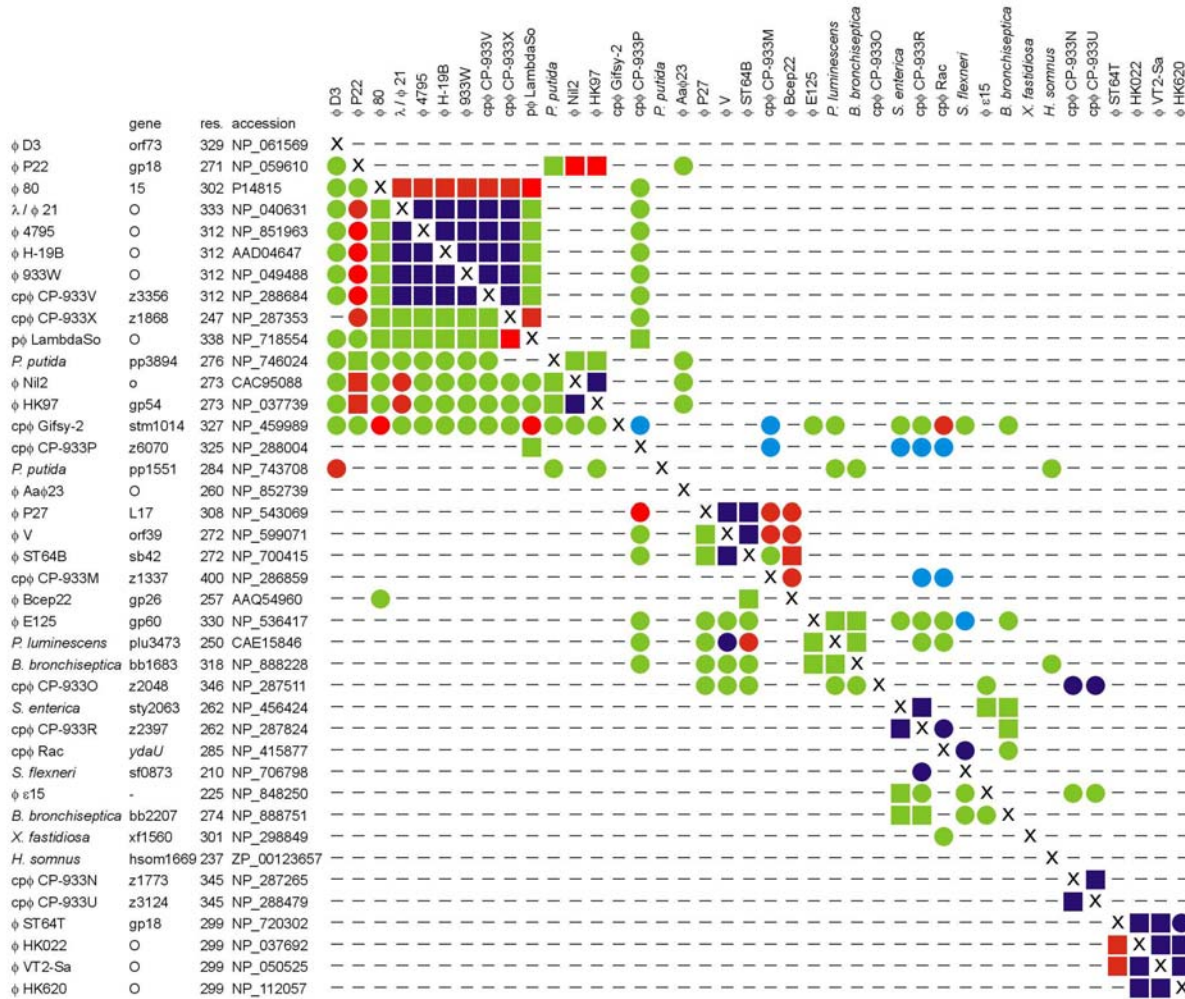
a nucleosome-like structure, i.e. the DNA wrapped around a protein core – the  $\lambda$  O-*oriλ* complex ('O-some') may serve as example (see above) [446]. The core of the complexes would be held together by protein-protein interactions between iteron-bound initiator protomers, and the wrapping would create sufficient stress on the DNA helix to promote its unwinding. However, this model is not supported by electron microscopic analysis of the G38P-*oriL* complex [297]. This analysis suggested that the binding of G38P to the *ori* iterons results in DNA bending – as was observed for  $\lambda$  O (see above) –, stabilised by protein-protein interactions between G38P protomers. However, a 'shortening' of the DNA fragment – indicative of DNA wrapping – was not observed, which makes a nucleosome-like structure of the G38P-*oriL* complex unlikely. Until more experimental studies of initiator-origin complexes are presented the question of their higher-order structure(s) remains an intriguing issue.

Few other putative phage replication initiators have attained attention. Purified orf16(Rep) protein of  $\phi$ Tuc-2009 was found to bind specifically to an internal fragment of the orf16 gene containing the putative replication origin, which qualifies the protein as initiator for the replication of this *L. lactis* phage [282]. The orf11 initiator of  $\phi$ r1t was identified by its similarity to  $\phi$ SPP1 G38P; (partially) purified orf11 protein binds to its own coding region, and the nucleoprotein complexes analysed *in vitro* support the notion of extensive oligomerisation of orf11 bound to the (putative) replication origin [508]. A preliminary genetic analysis of  $\phi$ TP901-1 replication identified REP(orf13) as replication protein containing the replication origin [333].

### Initiator proteins of the Lambda O / SPP1 G38P-type.

BLAST similarity searches were performed with the  $\lambda$  O and  $\phi$ SPP1 G38P protein sequences as query. A final set of 140 putative initiator sequences was accumulated in an iterative manner, i.e. the searches were repeated with the least similar (putative) phage initiator from the results lists as new query sequence until no further new matching sequences were found. In most cases, the results of BLAST searches for a given query sequence showed four types of matching sequences:

- i. orthologues with similarity values close to 100% (ident. res.) and full-length homologues with values >60%;
- ii. proteins with full-length similarity in the range from 40% to 60%;
- iii. distantly related sequences with values in the range from 20% to 40% over the entire length, for the N-terminus, or for the C-terminus;
- iii. distantly related sequences with values in the range from 20% to 40% over the entire length, for the N-terminus, or for the C-terminus;



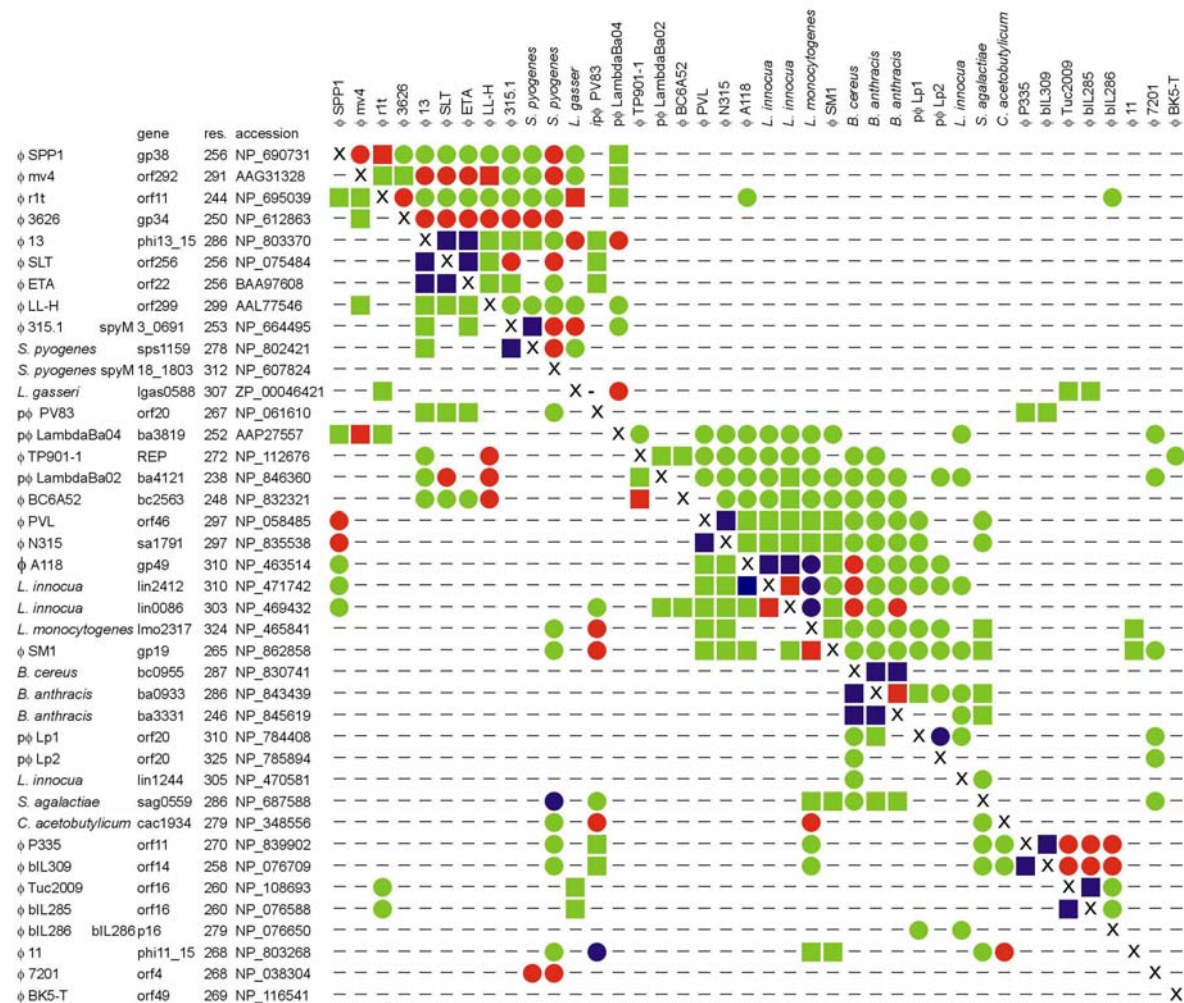
**Fig. C8** Protein sequence similarity among initiator proteins of Gram(-)-specific phages.

BLAST searches were performed separately for the N-termini (res. N~150) and C-termini (res. ~150-C) of the query-proteins. Similarities found for the N-termini are given below the diagonal; similarities found for the C-termini are given above the diagonal. Only regions of similarity with a length ~40 res. or more were considered significant except for a small region (~20 res.) with particularly high similarity (~70% ident. res.) at the extreme C-terminus in a subset of the proteins compared here (light blue filled circles). Similarities are shown in green (≥20-<40% ident.res.), red (≥40-<60% ident. res.), and blue (≥60% ident. res.). '-' no detectable similarity. Filled squares indicate similarity of two proteins in the N-terminus and in the C-terminus; filled circles indicate similarity in the N-terminus, or the C-terminus, respectively.

iv. proteins with highly significant similarity – even close to identity – exclusively to the N-terminus or to the C-terminus of the query sequence.

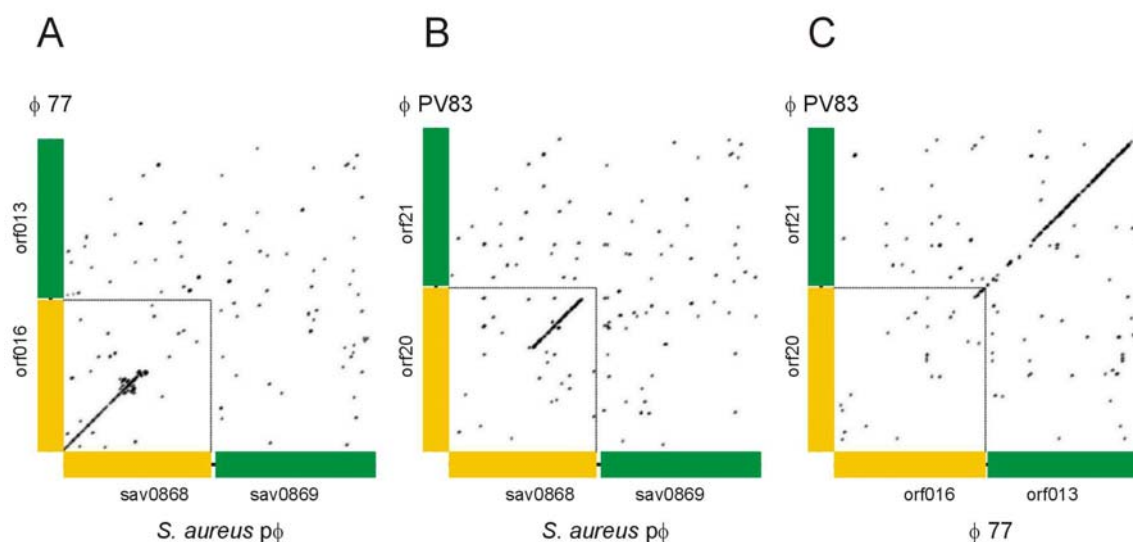
The majority of the matches were phage proteins or 'hypothetical phage proteins' from prophages in the sequenced bacterial genomes. Not surprisingly, a number of matches to putative transcription factors were obtained for the N-termini (putative DNA-binding domains). Genetic drift can easily explain the similarity types i.–iii.. For type iii., detectable low similarity in either

the N- or the C-terminus could have coincided with a degree of similarity in the second part of the proteins that was below the threshold level set by the BLAST parameters (~20% ident. res.). However, we observed for a number of proteins with full-length similarity in the ~40% to 60% range that the similarity for one domain was considerably higher than for the other, pointing to a two-domain architecture of all these proteins (Figs. C7 + C8). Similarity type iv. cannot be explained by genetic drift, and we have to assume instead an exchange of initiator domains among phages



**Fig. C7** Protein sequence similarity among initiator proteins of Gram(+) specific ph





**Fig. C9** DNA sequence similarity among (putative) replication genes:  $\phi 77$ orf016/orf013,  $\phi$ PV83 orf20/orf21, and *S. aureus* prophage sav0868/sav0869.

Dot-plot matrix analysis of the DNA sequences of the (putative) initiator and helicase loader genes of phages  $\phi 77$  (pos. 32867-34399, 1532 bp [NC\_005356]),  $\phi$ PV83 (pos. 9476-11062, 1586 bp [NC\_002486]), and an *S. aureus* p $\phi$  (pos. 927419-928948, 1529 bp [NC\_002758]). Dot-plots were obtained using 'method 2' (K\_tuple value = 8) of the dot matrix subprogram of the DNAMAN<sup>®</sup> software (version 4.0; Lynnon Inc.). **A** : *S. aureus* p $\phi$  X  $\phi 77$ ; **B** : *S. aureus* p $\phi$  X  $\phi$ PV83; **C** :  $\phi 77$  X  $\phi$ PV83. Yellow blocks: putative initiator; green blocks: (putative) DnaE<sub>eco</sub>-type helicase loader; black line: intergenic region. Gene sizes are shown exactly to scale. Dotted lines indicate the 3'-ends of the initiator genes.

**Table C10** Sequence similarity among (putative) replication proteins:  $\phi 77$ orf016/orf013,  $\phi$ PV83 orf20/orf21, and *S. aureus* prophage sav0868/sav0869.

$\phi$ / p $\phi$	gene	res. accession	<i>S. aureus</i> p $\phi$		$\phi$ PV83		$\phi 77$		gene	res. accession	<i>S. aureus</i> p $\phi$	$\phi$ PV83	$\phi 77$
			N	C	N	C	N	C					
<b>initiator</b>													
<i>S. aureus</i> p $\phi$	sav0868	243 NP_371392	x	x									
$\phi$ PV83	orf20	267 NP_061610	–	95	x	x							
$\phi 77$	orf016	247 NP_958644	98	–	–	32	x	x					
$\phi$ Bcep22	gp26	257 NP_944256	28	–	–	–	26	–					
<b>helicase loader</b>													
	sav0869	261 NP_371393	x										
	orf21	257 NP_061611	43		x								
	orf013	259 NP_958645	44		87			x					
	gp27	265 NP_944257	26		30			30					

BLAST (bl2seq; [445]) similarity searches were performed for the N- and C-termini of the initiator proteins separately, and for the (putative) helicase loaders with the complete sequences. Values are % ident. res.; '–' no detectable similarity.

exchanged by recombination among ancestors of these (pro)phages (Fig. C9; A+B). As was also observed for most of the (putative) initiator genes with high N-terminal similarity, the similarity on the DNA level extends from the 5'-end to a region downstream of the iteron region defining the putative replication origin (Fig. C9; A: note small strings parallel to the diagonal indicative of repeats close to the middle of the (putative) initiator

genes sav0868 and orf016). Contrarily,  $\phi 77$  orf016 and  $\phi$ PV83 orf20 show only moderate similarity to each other, confined to the C-termini, but with complete identity for the C-terminal 20 res. (Fig. C9; C: see DNA sequence identity in the extreme 3'-ends of the initiator genes). It is obvious that conventional methods used to determine evolutionary distances in a set of analogous proteins on the basis of mutation rates in the full-length

sequences would produce highly biased results for mosaic proteins like the phage initiators.

From the 140 (putative) initiator sequences found by 'saturation' BLAST searches 40 were chosen as examples from Gram(+)-specific (pro)phages (Fig. C7) and another 40 from Gram(-)-specific (pro)phages (Fig. C8) for a matrix-type comparison. The number of 40 was chosen so that the matrix would still fit to a page-size figure; the selected sequences represent closely related as well as seemingly unrelated sequences. Similarities between (putative) initiators of Gram(+)-specific and Gram(-)-specific (pro)phages were rarely found, and mostly confined to the N-terminal DNA-binding domains with similarity values in the range of ~30% or below. Although we cannot exclude interbreeding between phages of both groups – including exchange of parts of their (putative) initiator genes by recombination – such rare events are hardly detectable by crude methods like BLAST comparisons.

The matrices shown in Figs. C7 + C8 give an simultaneous synopsis of similarities calculated separately for the N- and C- termini because it was clear from the above that only this procedure would allow for a meaningful comparison. The sequences were arranged manually in a way that nearly identical sequences were placed as closely as possible to each other. Around these 'foci', the remaining sequences were placed according to their decreasing similarity. It is immediately apparent from the matrices that several subgroups exist within both sets. Tentatively, we would call them the  $\lambda$ , the  $\phi$ P27/  $\phi$ Rac, and the  $\phi$ VT2-Sa subgroups for the Gram(-)-specific phages. For the Gram(+)-specific phages, we would tentatively define the  $\phi$ ETA, the  $\phi$ A118, and the  $\phi$ P335 subgroups. However, the 'borders' of these subgroups are fuzzy, and there are examples in both matrices for sequences that fit into different subgroups with their N- and C- termini, respectively, e.g. the initiators of the Gram(-)- specific prophages Gifsy-2 and CP-933P, and the initiators of the Gram(+)-specific phages  $\phi$ SPP1,  $\phi$ 11 and  $\phi$ PV83. These protein sequences belonging to different subgroups support the notion of recombinatorial crosstalk between the subgroups, and create a 'link' – on the 'BLAST level' – of unrelated sequences, e.g.  $\phi$ SPP1 G38P and  $\phi$ BK5-T orf49. For both sets, however, the sample size is not sufficient for a prediction whether further sequences would consolidate the existing subgroups, create novel subgroups, or add more examples for sequences belonging to different subgroups with each of their two domains.

For most of the initiator proteins from both sets with significant similarity to each other, a correlation exists with the type of replication module driving the replication of these phages (BRM Section 3.1.). However, there are exceptions: the initiators of the Gram(-)-specific phages  $\phi$ Bcep22 gp26,  $\phi$ V orf39, and  $\phi$ P27 L17 have

significantly similar C-termini, but  $\phi$ Bcep22 has an IL-type replication module,  $\phi$ V an 'Initiator-solo' (I-solo type) replication module, and  $\phi$ P27 an 'Initiator-helicase loader-helicase' (ILH-type) replication module. Also, the initiators of the Gram(+)-specific phages  $\phi$ SPP1 G38P and  $\phi$ r1t orf11 are similar, but  $\phi$ SPP1 has an ILH-type replication module, and  $\phi$ r1t an IL-type module. Therefore, reliable predictions about the replication mechanism of a given (pro)phage can only be made after an analysis of the complete set of replication genes in its genome, particularly with respect to the point of entry of the host replication machinery (see below).

In conclusion, there are four lines of (theoretical) evidence that most if not all proteins analysed in the 2 matrices are *bona fide* initiators for the replication of their cognate (pro)phages: i. although a direct one-to-one comparison of two proteins fails to reveal similarity for most of them, it is possible to find 'connecting' protein(s) with similarity to both, ii. all proteins are of similar size and show the two-domain architecture of  $\lambda$  O with N-terminal DNA-binding and C-terminal oligomerisation/interaction domains, iii. for most of the initiator genes a co-localisation with other replication genes – i.e. helicase loader and/or helicase genes – can be detected in a defined part of the respective phage genomes (BRM Section 3.1.), and iv. in >90% of the 80 initiator genes, (putative) replication origin structures could be detected located approximately in the middle of the genes (Section C2.2.).

Several (putative) initiator proteins of Gram(+)-specific (pro)phages contain in their C-terminal domains regions of significant similarity to the DnaB and/or DnaD helicase loaders of their hosts. This topic is discussed in more detail in BRM Section 3.1.1. In a few instances, we found similarities of the N-termini of phage initiators to plasmid-encoded proteins. The *Streptomyces* sp. plasmid pSCL contains a replication module strongly resembling the IH-type phage replication module (BRM Section 3.1.3.). The N-terminus of the putative initiator of pSCL shows similarity to the N-terminus of *E. coli* phage P27 initiator L17 (27% ident. res.).  $\phi$ P27 L17 also shows similarity (31% ident. res.) in its N-terminus to the N-terminus of the RepC protein of pTAV320 (RSF1010-type; host: *Paracoccus pantotrophus*,  $\alpha$ -proteobacteria). Taken together, there are only few indications for a closer relationship of plasmid and phage initiator proteins. An intriguing finding is the N-terminal similarity (25% ident. res.) of the *Bacteroides thetaiotaomicron* (Chlorobium/Bacteroides group) gene bt1515 to  $\phi$ SPP1 G38P and  $\phi$ PVL orf46 because this gene is followed downstream by a (putative) DnaB<sub>Eco</sub>-type helicase, reminiscent of the IH-type replication module (BRM Section 3.1.3.). No phages specific for species of this phylon have yet been described at



the sequence or molecular level. These two genes may indicate that prophages similar to Gram(+)-specific phages are present in the genomes of species from this phylon.

**RepL, the initiator for bacteriophage P1 replication in the lytic cycle.** *E. coli* cells lysogenic for phage P1 carry the prophage either integrated into their chromosomes or as low-copy number plasmid [76,259]. Plasmids containing the  $\phi$ P1 R (plasmid) replicon have been among the favourite models for studies of plasmid replication and plasmid copy number control [98,74]. The L (lytic) replicon – responsible for  $\phi$ P1 propagation in the lytic cycle – has attracted considerably less attention.

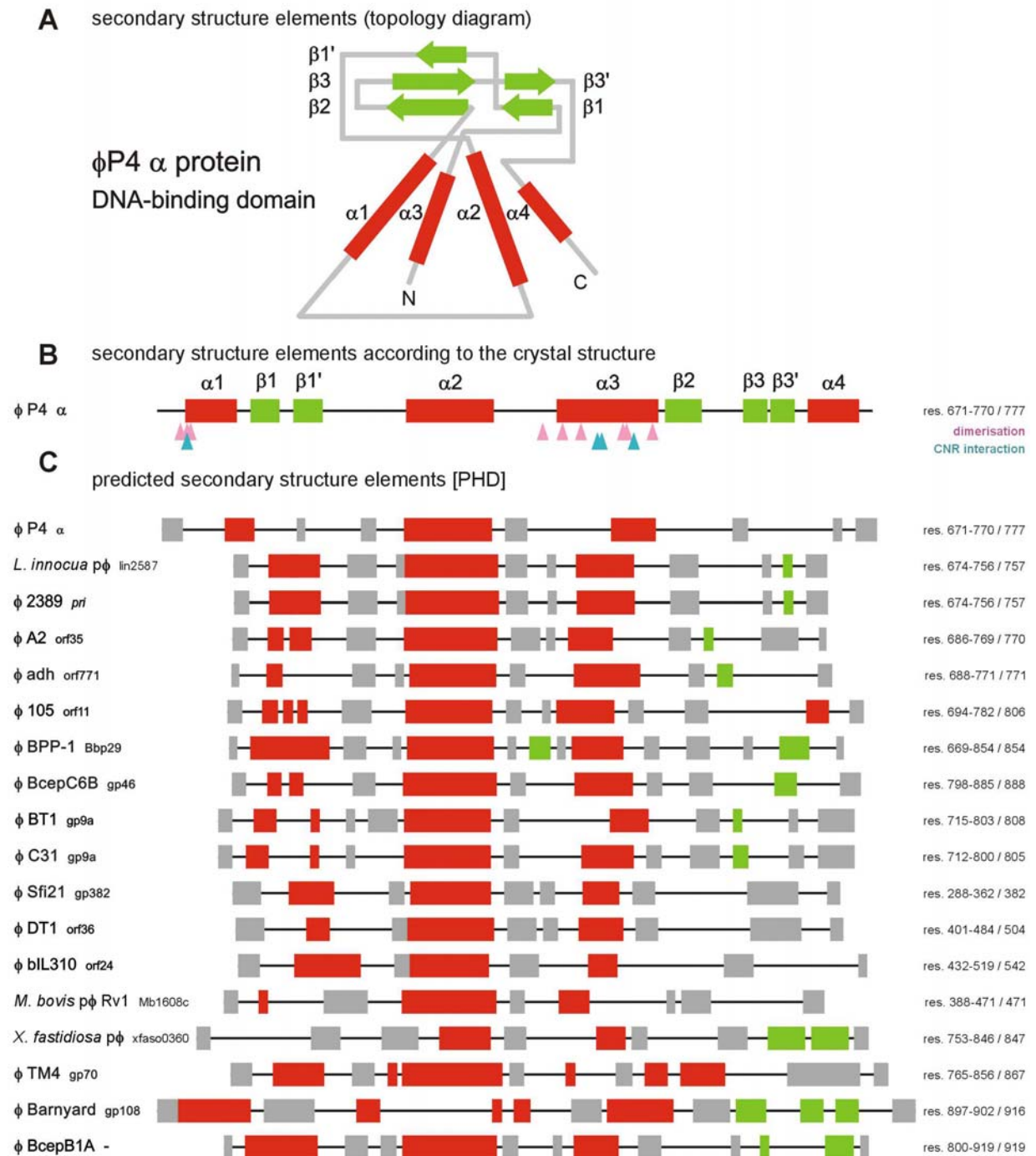
Early genetic analyses suggested that different replicons might be responsible for the propagation of  $\phi$ P1 as plasmid or as phage. Experimentally, the phage burst size and specific DNA synthesis were measured in a set of *E. coli* strains carrying temperature-sensitive mutations in replication genes. Phage propagation could be shown to depend on host DnaC, DnaG, and Pol III holoenzyme, but did not require DnaA and DnaB [165]. In contrast,  $\phi$ P1 plasmid replication crucially depended on host DnaA and DnaB, and – in addition – on DnaC, DnaG, and Pol III holoenzyme [163]. Analysis of  $\phi$ P1 replication intermediates by electron microscopy revealed theta-shaped molecules early after infection of *E. coli* and predominantly sigma-shaped molecules later on. Such  $\sigma$ -shaped molecules were rarely found in *recA* host cells. Cohen concluded from these observations that initiation of  $\phi$ P1 replication occurs in the  $\theta$ -mode and switches to the  $\sigma$ -mode later depending on host recombination functions [81]. A switch from  $\theta$ DR to  $\sigma$ DR is also known for  $\lambda$  and  $\phi$ SPP1 replication but is dependent only in the case of  $\phi$ SPP1 on the cognate recombination functions (Section C3.6.2.). In independent studies, Hansen [162] and Cohen + Sternberg [82] could show that the  $\phi$ P1 L replicon is contained within a DNA segment that carries the *c4*, *ant*, *kilA* – all three non-essential for replication – and *repL* genes. Transcription of the *kilA* and *repL* genes is driven by an upstream promoter, P53, which is negatively regulated by the  $\phi$ P1 cI repressor. Insertional inactivation of the *repL* gene abolishes L replicon activity. A putative replication origin structure is located within the *repL* gene, comparable to the localisation of the replication origin in ' $\lambda$  O-type' initiators (Section C2.2.). P53-driven transcription is required for L replicon activity but the promoter can be replaced by e.g. the *lac* promoter. Schuster and co-workers could show that *repL* mRNA transcription is also subject to negative control by an 180 nt long anti-sense transcript [167]. The  $\phi$ P1 RepL protein has not yet been studied biochemically.

The  $\phi$ P1 *ban* (DnaB-analogue) gene encodes the cognate replicative helicase, an orthologue of *E. coli* DnaB (Section C3.3.) Because  $\phi$ P1 replication from *oriL* requires host DnaC but not DnaB, we have to assume that RepL directs a  $Ban_6DnaC_6$  double-hexamer(s) to the unwound region, but experimental evidence is lacking. Phage proteins with detectable similarity to  $\phi$ P1 RepL are not known.

**The 'initiator domain' of phage P4 alpha ( $\alpha$ ) protein and related helicases.** Like P1 but for different reasons, *E. coli* phage P4 belongs to the group of replicons that withstand a straightforward classification: it is a natural occurring phage-plasmid hybrid. All protein components of the P4 virus capsid stem from *E. coli* phage  $\phi$ P2, and progeny can only be obtained from *E. coli* hosts lysogenic for  $\phi$ P2 – or from mixed infections. Therefore,  $\phi$ P4 is usually described as satellite phage of  $\phi$ P2, although the 11.6 kb-long dsDNA genome of  $\phi$ P4 is completely unrelated to the  $\phi$ P2 genome with respect to replication functions. Following infection of a host that is lysogenic for  $\phi$ P2, the circular  $\phi$ P4 genome can either enter the lytic cycle or integrate into the host chromosome as prophage, mediated by its cognate integration system. In a  $\phi$ P2<sup>-</sup> host, the  $\phi$ P4 prophage can either integrate or replicate as a free plasmid. Integration is preferred, however, and maintaining  $\phi$ P4 as plasmid under laboratory conditions requires continuous counter-selection. The extraordinary life-style of  $\phi$ P4 has been reviewed in detail by Dehò, Ghisotti, and colleagues [47].

Replication of  $\phi$ P4 in the lytic cycle or as free plasmid depends on the alpha ( $\alpha$ ) protein, host SSB, and DNA Pol III holoenzyme, but is independent of DnaA, DnaB, DnaC, and DnaG [253,22,43,233,103].  $\phi$ P4  $\alpha$  is a multifunctional polypeptide that exhibits primase activity (Section C3.4.), helicase activity (Section C3.3.), and binds specifically to the  $\phi$ P4 replication origin, *oriI* (Section C2.2.).  $\phi$ P4  $\alpha$  binds to the decameric iterons GGTGAACAGA/T in *oriI* and *crr*, which both constitute the P4 *oriI* replication origin [123, 506].  $\phi$ P4 replication proceeds bidirectionally from *oriI* in the  $\theta$ -mode [234]. Because it is responsible for origin unwinding,  $\phi$ P4  $\alpha$  has the properties of an initiator and is discussed in this context.

The  $\phi$ P4  $\alpha$  C-terminus has DNA binding activity, and retains this activity upon isolation [505]. In addition, mutations were found in the  $\alpha$  C-terminus that render the protein insensitive to negative regulation by the Cnr protein (see below). This suggested to Lanka and co-workers that the C-terminal DNA-binding and Cnr-interaction domains overlap [503]. Dehò, Ghisotti and co-workers could show by a yeast two-hybrid approach that  $\alpha$  protein can oligomerise, and this interaction domain was also located in the  $\alpha$  C-terminus [452]. The



**Fig. C10** (Putative)  $\phi$ P4  $\alpha$ -type initiator domains.

**A** topology diagram of the  $\phi$ P4  $\alpha$  DNA-binding domain adapted from Fig. 3B from the paper by Yeo *et al.* [492];. **B** linear projection of the secondary structure elements according to the crystal structure of the  $\phi$ P4  $\alpha$  DNA-binding domain [PDB 1KA8]; pink arrowheads point to residues involved in dimerisation, light blue arrowheads point to residues involved in the interaction with  $\phi$ P4 Cnr according to Yeo *et al.* [492]. **C** secondary structure prediction by the PHD method for the listed proteins. For accession numbers see Table C16 . Colour code: red = predicted  $\alpha$ -helical region; green = predicted  $\beta$ -stranded region; grey = predicted loop; black line = unstructured/no prediction. The sequences were aligned along  $\alpha$ -helix 2 of the  $\phi$ P4  $\alpha$  DNA-binding domain. Right column gives the regions of the proteins analysed by PHD and the length of the protein sequences.

three-dimensional structure of the  $\phi$ P4  $\alpha$  DNA-binding domain (or OBD, for origin-binding domain) is known at a resolution of 2.95 Å [PDB 1KA8]. Waksman and co-workers describe the  $\alpha$  DNA-binding domain as double winged-helix with pseudo-twofold symmetry (see Fig. C10; A); they could pinpoint the residues responsible for dimerisation and for interaction with  $\phi$ P4 Cnr (see Fig. C10; B) [492].

In the lytic cycle of  $\phi$ P4, ~100 progeny copies accumulate in the host cell, in contrast to a copy number of ~40 for the plasmid state [2]. Therefore, mechanism(s) must operate that control  $\phi$ P4 plasmid copy number. One mechanism involves the phage encoded Cnr (copy number regulation) protein which down-regulates  $\alpha$  activity through interaction (see above) [447,503]. It is presently not known exactly at which step inhibition of  $\alpha$  by Cnr occurs: although it does apparently not inhibit DNA-binding of  $\alpha$ , it may inhibit dimerisation of  $\alpha$  bound to iterons in *ori1* or *crr*. Alternatively, it may prevent the additional oligomerisation of  $\alpha$  protomers bound to these two sites, which is necessary for origin unwinding. Modulation of the oligomerisation state of the initiator proteins is one out of several mechanisms for copy-number control operating in plasmid replicons [98]. The other known mechanism regulating  $\phi$ P4 copy number is a transcriptional network involving a highly complex and timely ordered interplay of activities of cognate transcription factors acting on several  $\phi$ P4 promoters (for details see Briani et al. [47]). The *cnr* gene is located immediately upstream of the  $\alpha$  gene in the same transcription unit and the expression of both genes is co-regulated [96,350]. The available data suggest that  $\phi$ P4 replication strictly requires a balanced expression of Cnr and  $\alpha$  [447].

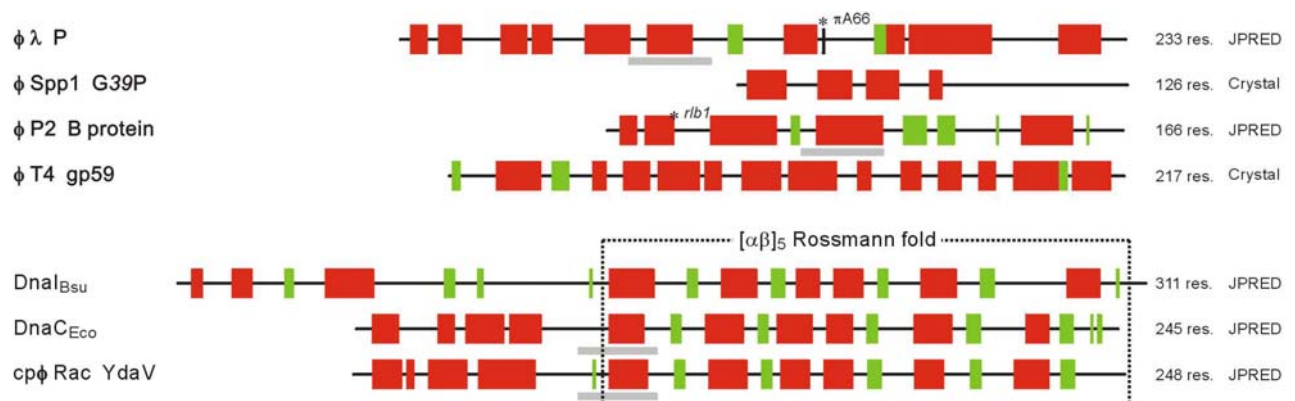
BLAST searches with  $\phi$ P4  $\alpha$  as query led to the detection of numerous (pro)phage-encoded proteins with varying degrees of similarity, mostly confined to the helicase domain but in several cases also including the N-terminal primase domain (Sections C3.3. + C3.4.). In contrast, BLAST searches with the  $\phi$ P4  $\alpha$  DNA-binding domain (res. 671-770) led only to the detection of few orthologous prophage sequences (>80% ident. res.). In all proteins with similarity to the  $\phi$ P4  $\alpha$  helicase domain, the homologous region is located either at a position corresponding to that of the  $\phi$ P4  $\alpha$  helicase domain, or at the N-terminus, suggesting that – in the latter case – these proteins lack a primase domain. However, all proteins have a C-terminal domain of approximately 100 res. in length, similar in length but lacking detectable BLAST-similarity to the  $\phi$ P4  $\alpha$  DNA-binding domain. We speculated that this domain may be the DNA-binding domain of these proteins and performed a 'MultAlin' [84] analysis of a subset of these proteins. We obtained a reasonable alignment with gaps only introduced between  $\beta 1+\beta 1'$  and to both sides of  $\beta 3'$ , but

no significant consensus sequence (not shown; for numbering of structural elements see Fig. C10; B). This subset of proteins was subjected to secondary structure prediction analysis by the PHD method [374]. The results were projected onto the secondary structure of  $\phi$ P4  $\alpha$ , which is known from the crystal structure (Fig. C10; C). Except for  $\alpha$ -helix 2, the PHD method failed to predict the known secondary structure elements of  $\phi$ P4  $\alpha$  precisely: neither the  $\beta$ -stranded regions nor  $\alpha$ -helix 4 were predicted, nor the length and positions of  $\alpha$ -helices 1 and 3. However, there was no false prediction, and the 'loop' predictions are approximately correct. In spite of these limitations, the secondary structure prediction for  $\phi$ P4  $\alpha$  seems reliable enough to discuss the predictions for the other proteins in a meaningful way. Except for  $\phi$ Barnyard gp108, an equivalent of  $\alpha$ -helix 2 is predicted at a corresponding position for all sequences. Also for  $\alpha$ -helices 1 + 3, predictions could be obtained but less stringent than for  $\alpha$ -helix 2. Except for  $\phi$ 105 orf11, no predictions for  $\alpha$ -helices were obtained at the position of  $\alpha$ -helix 4. The most surprising result was the prediction of  $\beta$ -stranded regions at positions roughly corresponding to the  $\beta$ -strands 2, 3 + 3'. Given the low similarity on the protein sequence level, the highly similar secondary structure prediction for all proteins indicates that they may be structural homologues. Together with the position C-terminal to a helicase domain and a rather uniform length, the results of the secondary structure predictions support the hypothesis that the C-terminal domains of the (putative) helicases represent 'initiator domains' alike the  $\phi$ P4  $\alpha$  DNA-binding domain. However, experimental results supporting this hypothesis are not yet available – and homology modelling on the basis of the known crystal structure of  $\phi$ P4  $\alpha$  beyond the scope of this review.

### C3.2. Helicase loaders

Helicase loaders perform one or more of the following functions during primosome formation: i. they interact with the replicative helicase and promote its binding to unwound or forked DNA, ii. they interact with the DNA-bound pre-primosome, iii. they compete with SSB for binding to ssDNA, and iv. they link the replication of their cognate replicon to the host replication machinery in the case of numerous phage replicons. Genes encoding helicase loaders have so far not been described for genuine plasmid replicons.

A comparison of the primary sequence of several helicase loaders, or – more informative – of their secondary structure reveals that they are poorly related, if at all (Fig. C11). The similarity between the DnaC<sub>Eco</sub> and DnaI<sub>Bsu</sub> helicase loaders is entirely confined to the C-terminal nucleotide binding domain (Rossmann-fold



**Fig. C11** (Predicted) secondary structure of helicase loaders.

For  $\phi$ T4 gp59 [PDB 1C1K] and  $\phi$ SPP1 G39P [PDB 1N01] the secondary structure elements were taken from the crystal structures. Jpred secondary structure predictions were obtained for  $\lambda$  P [NP\_040632],  $\phi$ SPP1 G39P [Q38151],  $\phi$ P2 B [P07696],  $\phi$ T4 gp59, DnaI<sub>Bsu</sub> [P06567], DnaC<sub>Eco</sub> [P07905], and cp $\phi$  Rac YdaV [AAC74442] [88]. Colour code: red = predicted  $\alpha$ -helical region; green = predicted  $\beta$ -stranded region; black line = unstructured. The length of the protein sequences and the sizes of the structural elements are shown to scale. The (approximate) positions of the *rlb1* mutation in  $\phi$ P2 B and the  $\pi$ A66 mutation in  $\lambda$  P are indicated by an asterisk. The dotted line indicates structural elements of the Rossmann fold [373]. The grey bar underlining an  $\alpha$ -helical region in  $\lambda$  P, DnaC<sub>Eco</sub>, and cp $\phi$  Rac YdaV show the positions of the hypothetical interaction sites with DnaB<sub>Eco</sub> (see text and Fig. C12 for details).

type) [373], as has already been noted previously [421]. In contrast, with 50% identity over their entire length DnaC<sub>Eco</sub> and the protein encoded by the *ydaV* (b1360) gene of the *E. coli* K12 Rac prophage are certainly orthologues [489], which is reflected in the similarity of their secondary structure predictions (see below).

The discussion focusses on the three phage-encoded helicase loaders that have been studied extensively: gp59 protein of phage T4 [495], P protein of phage  $\lambda$ , and G39P of *B. subtilis* phage SPP1 [341]. In these three model systems the participation of the cognate helicase loader in primosome formation is obligatory. B protein, the helicase loader of coliphage P2, seems to be required for lagging-strand synthesis [326]. DnaC<sub>Eco</sub> and DnaI<sub>Bsu</sub> are the helicase loaders for chromosome replication in *E. coli* and *B. subtilis*, respectively, and are discussed together with their phage-encoded homologues in detail at the end of the chapter.

**Gene 59 protein (gp59) of phage T4.** Gp59 is required for both types of  $\phi$ T4 replication initiation: origin-dependent replication (R-loops) early after infection and RDR (D-loops) later on [198]. Also, gp59 efficiently co-ordinates leading and lagging-strand synthesis by blocking replication until the gp41 helicase is loaded [198]. The crystal structure of gp59 protein has been determined, and the protein appears to be composed of two, largely  $\alpha$ -helical domains [PDB 1C1K] [307]. Mueser and co-workers proposed a model that predicts N-terminal DNA-binding sites in gp59 as well as putative domains for the interactions with gp32 SSB and gp41 helicase, respectively [307]. It had been shown earlier

that purified gp59 is a monomer in solution, and binds to both gp41 helicase and gp32 SSB [495]. Gp59 was also found to bind to single and double-stranded DNA, most tightly however to forked DNA substrates or DNA/RNA hybrids mimicking R-loops [495,198]. Gp59 binding to artificial substrates *in vitro* occurred without apparent sequence specificity, and binding sites for each fork arm required more than 6 but less than 12 single-stranded nucleotides for efficient helicase loading to the lagging-strand [495,199,198]. The N-terminus of gp59 interacts with the C-terminus of gp32 SSB [191]. Because gp32 SSB binds to ssDNA through its C-terminal DNA-binding domain it was proposed that the interaction with gp59 might weaken the association with ssDNA. Oligomerisation of gp59 requires its contact with gp32, and it was concluded from kinetic studies that gp59 oligomerises as a hexamer which forms stable complexes with the hexameric gp41 helicase in the presence of nucleotide [191]. In the gp59-gp41 double-hexamer, the six gp59 monomers are arranged in a head-to-head orientation, i.e. as trimers of dimers, according to Ishmael and co-workers [192]. The gp59 C-terminus interacts with the gp41-gp41 subunit interface, involving N- and C-terminus of gp41 [192]. Gp 41 helicase lacking 20 res. from the C-terminus does not form a complex with 59 protein on fork DNA [199]. Upon hexamerisation, a conformational switch of the gp41-gp59 complex is assumed to occur, which finally displaces gp32 from fork DNA [192].

Homologues of  $\phi$ T4 gene59 (gp59) are encoded by phages  $\phi$ T2 (99%; res. 26-217),  $\phi$ RB69 (85%; res. 1-217),  $\phi$ RB49 (58%; res. 1-216),  $\phi$ Aeh1 (38%; res. 1-

126), and  $\phi$ KVP40 (33%; res. 13-216). BLAST searches could not identify other proteins in the data collections showing significant similarity with  $\phi$ T4 gp59. However, the identification of gp59 homologues in the genomes of phages  $\phi$ Aeh1 and  $\phi$ KVP40 that infect  $\gamma$ -proteobacterial hosts other than *E. coli* shows that this type of helicase loader is not uniquely found in the T-even phages of *E. coli*.

**Bacteriophage Lambda P protein.** P protein is the prototype phage-encoded helicase loader that links the replication of its replicon to the replication machinery of the host: it interacts tightly with the initiator protein O [509] and also with the host replicative helicase DnaB<sub>Eco</sub>, with which it forms a stable 3:6 complex [271]. Like its (functional) analogue DnaC<sub>Eco</sub>,  $\lambda$  P exhibits ssDNA-binding activity in the complex with DnaB<sub>Eco</sub>, which is further augmented by the interaction of the P·DnaB complex with ori $\lambda$ -bound O protein, the so-called 'O-some' [106,244]. The interaction of P with DnaB<sub>Eco</sub> is significantly more stable *in vitro* than the DnaB·DnaC<sub>Eco</sub> interaction, which allows the re-direction of the host replicative helicase to the phage replicon [271]. The ATPase-, ssDNA-binding, and helicase properties of DnaB<sub>Eco</sub> remain suppressed in the P·DnaB complex and in the O·P·DnaB complex, and the action of chaperones is required for the release of the helicase from the complex as the last step of pre-primosome formation at ori $\lambda$  [106,271].

Despite the wealth of data concerning the function of P during  $\lambda$  replication its molecular anatomy is poorly understood, and a crystal structure is not available. The ssDNA-binding domain within P has not yet been mapped. Wickner and Zahn could show that – in agreement with genetic observations – the O C-terminus interacts with P, but the O-interaction domain within P has not yet been determined (see above; [481]). There exist conflicting hypotheses concerning the interaction domain with DnaB<sub>Eco</sub>. Ogawa and co-workers found that the N-termini of the homologous phage helicase loaders  $\lambda$  P and  $\phi$ 80 gene14 protein show a higher similarity in their N-termini, which suggested to them that this part of the proteins interacts with a host protein, most likely DnaB<sub>Eco</sub> [329]. On the other hand it was proposed that the P C-terminus interacts with DnaB<sub>Eco</sub> because the  $\pi$ (pi) mutations of the P gene map exclusively to the C-terminus of P and result in a lowered affinity of the mutant proteins for DnaB<sub>Eco</sub> [368, 224] (see Fig. C11).

A small stretch (~ 30 res.) of similarity exists between  $\lambda$  P, DnaC<sub>Eco</sub>, and  $\phi$ P2 B, which we were unable to detect by BLAST analysis employing the complete protein sequences [326]. This region corresponds to  $\alpha$ -helix 1 of the Rossmann fold of DnaC<sub>Eco</sub> with few neighbouring residues, and is also present in the  $\phi$ 80 gene14 protein,  $\phi$ 27 L18, and cp $\phi$  Rac YdaV pro-

teins (see Figs. C11 + C12). The identity between the  $\lambda$  P and  $\phi$ 80 gene 14 proteins raises significantly to 77% in this region, in contrast to 47% over-all identity, and 65% in the N-terminal 110 res. [329]. Roughly one third of the residues are conserved among several of the helicase loaders, but also in the *E. coli* DnaA protein. Interestingly, one of the two interaction domains of DnaA<sub>Eco</sub> with the DnaB<sub>Eco</sub> helicase has been mapped to this region [279,404]. Presently, experimental data in support of the hypothesis that this motif might be involved in an interaction of  $\lambda$  P-type or DnaC<sub>Eco</sub>-type helicase loaders with the DnaB<sub>Eco</sub> helicase are not available.

DnaA <sub>Eco</sub>	137	F . DNF . VEGKSNQLA . RAAARQVADNPG	162
cp $\phi$ Rac YdaV	73	F . SNYQVQNEGQRYA . LSQAKSIADLM	99
$\phi$ P27 L18	73	F . DNY . LEVNPDAARNLAACRRYAENWP	98
DnaC <sub>Eco</sub>	70	F . ENYRVECEGQMNA . LSKARQYVEEFD	96
$\lambda$ P	74	FRENG . ITTMEQVNAGMRVARRQNRPF	100
$\phi$ 80 14	75	FQENG . IHTMAQVDAGMRIARRQERPFL	101
$\phi$ P2 B	68	VERERLVCAIDELRGAFSKRRQVGASEY	95
consensus		F . DNY . I . . . . Q . . A . L . . . AR . . . . .	
		E G V D M C	

**Fig. C12** Hypothetical interaction site of initiators and helicase loaders with DnaB<sub>Eco</sub>.

Sequence alignment for DnaA<sub>Eco</sub> [P03004], p $\phi$ Rac YdaV [AAC74442],  $\phi$ P27 L18 [NP\_543070], DnaC<sub>Eco</sub> [P07905],  $\lambda$  P [NP\_040632],  $\phi$ 80 14 [CAA31475],  $\phi$ P2 B [P07696]. The hypothetical site of interaction is shown as pink box, the regions in the individual sequences are indicated by residue numbering.

Full-length homologues of the  $\lambda$  P helicase loader are encoded in the genomes of phages  $\phi$ 21 (99% ident.),  $\phi$ H-19B (97% ident.),  $\phi$ 4795 (96% ident.),  $\phi$ 933W (96% ident.), and – as mentioned above –  $\phi$ 80 (47% identity with res. 1-201 of  $\lambda$  P). Various homologues of the  $\lambda$  p gene are also present in prophages within the genomes of several *E. coli*, *Salmonella* and *Shigella* strains, e.g. in *E. coli* O157:H7 EDL933 the genes z3355 of cp $\phi$  CP-933V (96% ident. res.), and z1869 of cp $\phi$  CP-933X (97% ident. res.). Two partial homologues probably represent truncated, non-functional genes: the hypothetical 77 residues long ybcJ gene product of *E. coli* K12 cp $\phi$  DLP12 (96% identity with res.131-233 of  $\lambda$  P), and the hypothetical 103 residues long z0313 gene product of *E. coli* O157:H7 EDL933 cp $\phi$  CP-933H. Three more prophage genes in *E. coli* O157:H7 EDL933 (z1774 of cp $\phi$  CP-933N, z2094 of cp $\phi$  CP-933O, and z3123 cp $\phi$  CP-933U) show partial similarity with the  $\lambda$  P C-terminus but none with its N-terminus (~20% identity with res. 94-202 of  $\lambda$  P). This



observation points to the possibility that  $\lambda$  P-type helicase loaders are composed of two distinct domains, which might be found in various combinations also in yet unidentified (pro)phage-encoded helicase loaders.

#### Gene 39 protein (G39P) of *B. subtilis* phage SPP1.

A crystal structure for the essential helicase loader G39P of  $\phi$ SPP1 phage has been obtained that shows a largely unstructured C-terminus of the protein [16] [PDB 1NO1]. G39P does not bind to DNA, but interacts with the initiator, G38P, and with the G40P helicase with the C-terminal 14 res. being essential for this interaction *in vitro* [13,17]. In the complex with G40P helicase and ATP, G39P inactivates the single-stranded DNA binding, ATPase and unwinding activities of G40P. G38P does by itself not interact with its cognate helicase in solution but with the G39P·G40P·ATP complex. When bound to the replication origin *oriL*, G38P interacts specifically with the G39P·G40P·ATP complex, which results in G40P binding to the origin. The dissociation of G38P·G39P heterodimers from the unstable intermediate complex is supposed to be required for subsequent helicase action [13]. No homologues of the  $\phi$ SPP1 G39P are presently found in the databases by BLAST searches.

**B protein of phage P2.** During RCR of  $\phi$ P2 *in vivo*, the displaced strand remains single-stranded in the absence of B protein [132]. It was therefore assumed that B protein is required to direct the host replicative helicase DnaB<sub>Eco</sub> to the displaced strand. Odegrip *et al.* could show that  $\phi$ P2 B protein interacts with DnaB<sub>Eco</sub> *in vitro*, [326]. They could show in addition that B protein competes with DnaC<sub>Eco</sub> for binding to DnaB<sub>Eco</sub>, comparable to  $\lambda$  P. From the observation of partial suppression of an *E. coli dnaB(ts)* mutation by the *rib1* mutation in the  $\phi$ P2 B gene it was speculated that the B N-terminus interacts with DnaB<sub>Eco</sub> (see Fig. C11). DNA-binding activity of  $\phi$ P2 B protein has not been shown, and its exact role during primosome formation for lagging-strand synthesis remains to be elucidated.

BLAST searches could identify full-length  $\phi$ P2 B homologues in phages L-413C (98%; res. 1-166) and W $\phi$  (98%; res. 1-166). Despite a comparable size, the limited similarity of  $\phi$ P2 B with gp77 of  $\phi$ Rosebush (31%; res. 41-109) and with orf3 of  $\phi$ K139 (28%; res. 36-80) does not support a straight-forward qualification of these proteins as helicase loaders; this applies also to the *E. coli* K12 *ybgA* gene product (32%; res. 18-81).

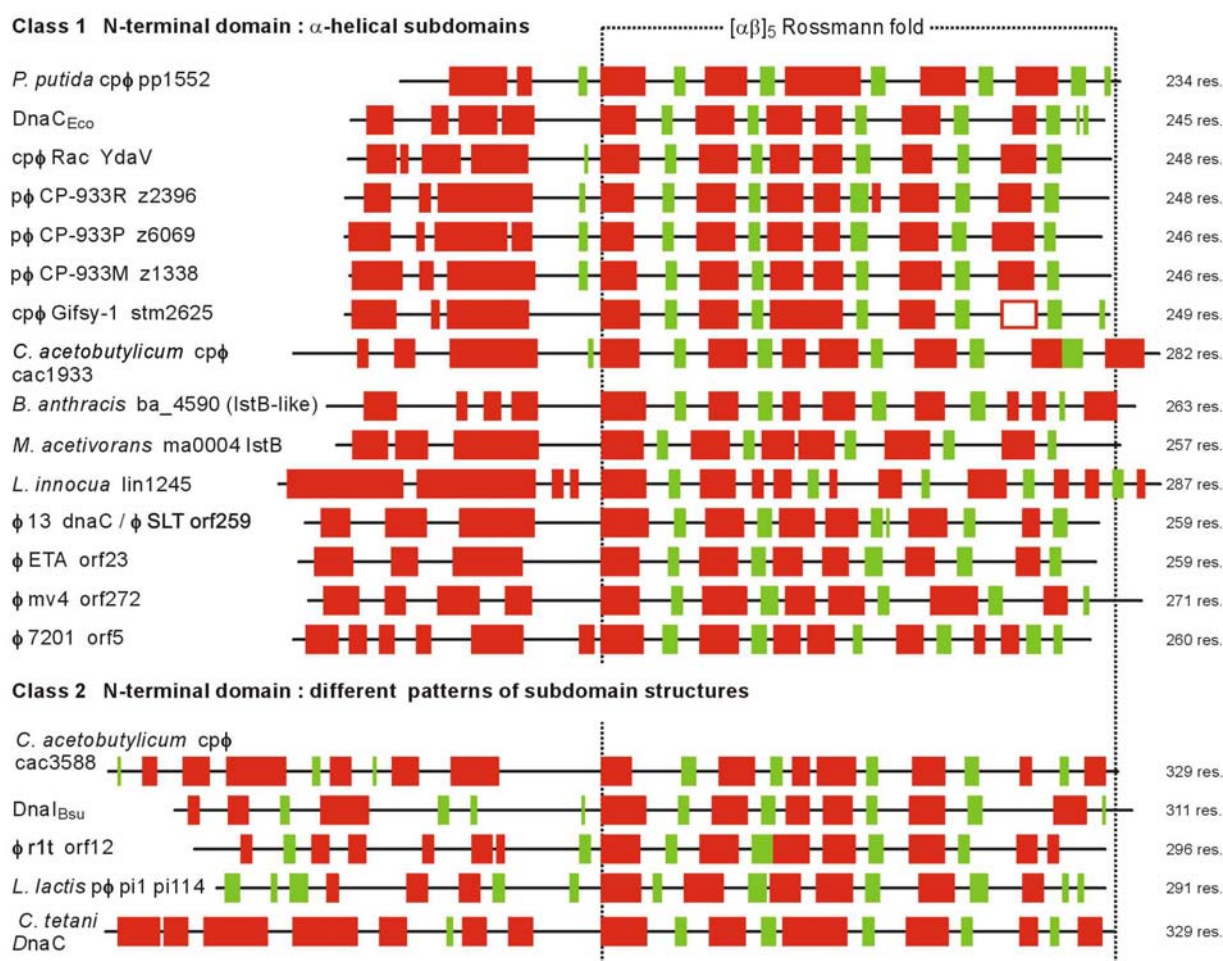
#### *E. coli* DnaC and phage-encoded homologues.

DnaC was originally detected as host factor required for  $\phi$ X174 complementary strand synthesis [479,390]. Primosome formation on  $\phi$ X174 *in vitro* requires the *E. coli* primosomal proteins PriA, PriB, PriC, DnaT, the

replicative helicase DnaB in complex with DnaC, and the DnaG primase, in addition to  $\phi$ X174 ssDNA and SSB [319]. The ' $\phi$ X174-primo-some' was later renamed to 're-start primosome' when it became apparent that its formation is crucial for replication re-start and SDR (BRM Section 2.5.). Primosome formation at the *E. coli oriC* during initiation of chromosome replication *in vitro* requires DnaA, DnaB, DnaC, and DnaG (reviewed in [292]). DnaC binds and hydrolyses ATP and contains a nucleotide-binding domain of the Rossmann fold-type with the characteristic 'Walker A+B motifs' [373]. Together with DnaA proteins DnaC belongs to a superfamily of ATPases whose members are found throughout all phylogenetic kingdoms [225]. DnaC does not belong to the AAA+ family of chaperone-like ATPases [318] – as erroneously stated by Davey *et al.* [92] – because it lacks the AAA+-specific RFC-Box VII', RFC-Box VII", and Sensor II motifs [318].

Although probably evolutionary related, DnaC and DnaA are not functionally analogous proteins, and DnaC lacks the DNA-binding domain of DnaA [376]. However, both proteins use ATP-hydrolysis for a conformational switch, which dramatically alters their interaction properties with other proteins and DNA [424,92]. In the presence of ATP DnaC co-purifies with DnaB, and forms a stable double-hexameric complex [243]. The helicase action of DnaB requires the ATP-driven release of DnaC from the complex following its binding to DNA [131,467,468,120]. The double-hexameric DnaB<sub>6</sub>·DnaC<sub>6</sub>·ATP complex binds to ssDNA mediated by the cryptic DNA-binding property of DnaC [244]. Davey and co-workers could show that the ATP-form of DnaC in the DnaB<sub>6</sub>·DnaC<sub>6</sub>·ATP complex is required for its efficient binding to unwound *E. coli oriC* concomitant with an enlargement of the single stranded region, while ATP-hydrolysis is driven by the combined action of DnaB and ssDNA on DnaC [92]. Hydrolysis of bound ATP by DnaC is required to liberate the inactive helicase from the complex, in analogy to the action of chaperones required to separate DnaB from  $\lambda$  P during primosome formation at *oriL* (see above).

By cryoelectron microscopy of DnaB<sub>6</sub>·DnaC<sub>6</sub> complexes and three-dimensional image reconstruction Bárcena, San Martín and co-workers obtained structures in which six DnaB monomers were assembled as three asymmetric dimers in a hexamer with threefold and sixfold symmetries on its both sides, respectively. The six DnaC protomers bound tightly to the sixfold face of the DnaB hexamer [382, 20]. Despite these efforts, the molecular anatomy of DnaC<sub>Eco</sub> is poorly understood, and a crystal structure is not available. The ssDNA-binding domain within DnaC has not yet been mapped. The *dnaC810* mutations was shown to confer to the mutant protein the property of helicase loading to SSB-covered ssDNA in the absence of the re-start primosomal pro-



**Fig. C13** (Putative) helicase loaders; part A.

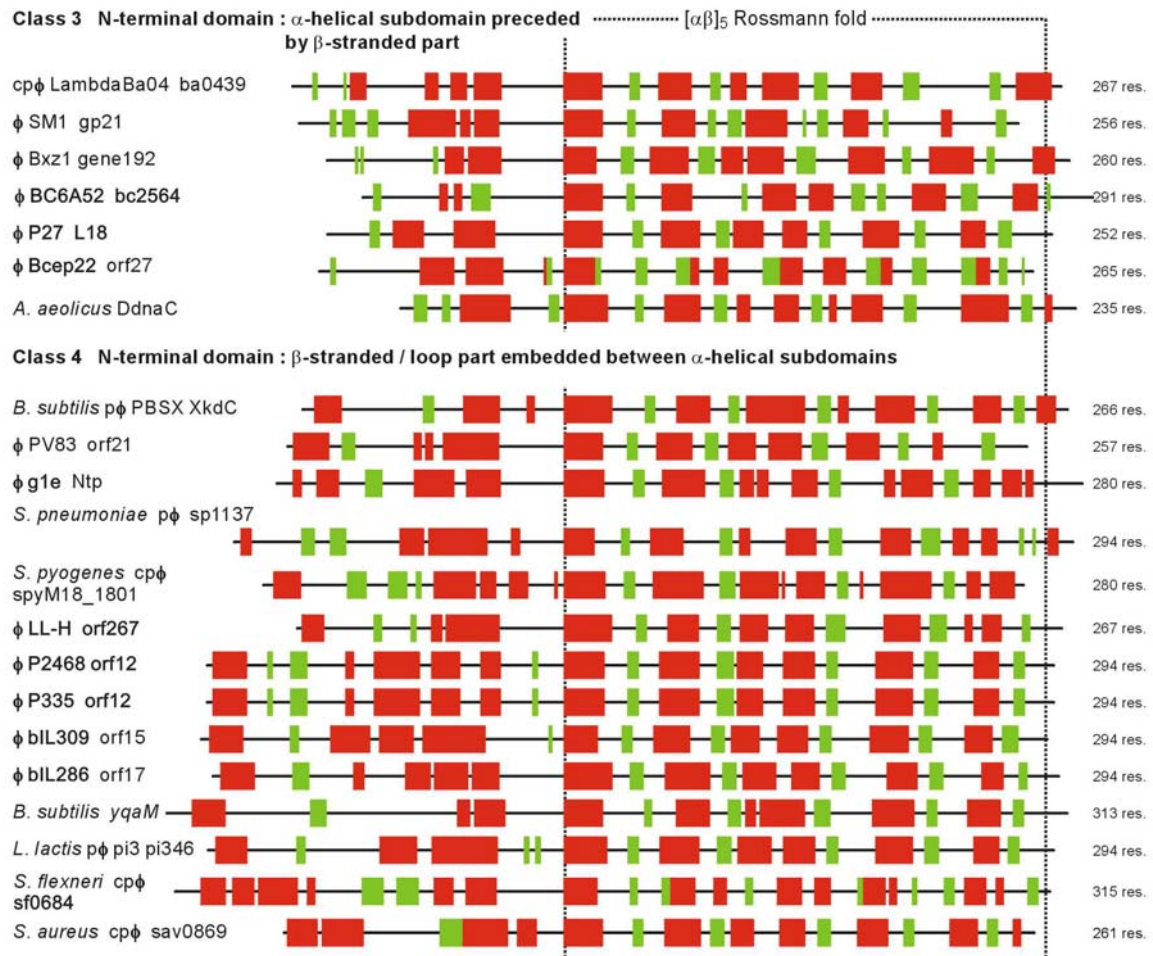
Jpred secondary structure predictions [88] were obtained for *P. putida* cp $\phi$  pp1552 [NP\_743709], DnaC<sub>Eco</sub> [P07905], cp $\phi$ Rac YdaV [AAC74442], p $\phi$  CP-933R z2396 [NP\_287823], p $\phi$  CP-933P z6069 [NP\_288003], p $\phi$  CP-933M z1338 [NP\_286860], cp $\phi$  Gifsy-1 stm2625 [NP\_461560], *C. acetobutylicum* cp $\phi$  cac1933 [NP\_348555], *B. anthracis* ba\_4590 [NP\_657950], *M. acetivorans* ma0004 [NP\_614978], *L. innocua* p $\phi$  lin1245 [NP\_470582],  $\phi$ 13 dnaC [NP\_803371] /  $\phi$ SLT orf259 [NP\_075485],  $\phi$ ETA orf23 [NP\_510917],  $\phi$ mv4 orf272 [AAG31329],  $\phi$ 7201 orf5 [NP\_038305], *C. acetobutylicum* cp $\phi$  cac3588 [NP\_350171], DnaI<sub>Bsu</sub> [P06567],  $\phi$ r1t orf12 [NP\_695040], *L. lactis* p $\phi$  pi1 pi114 [NP\_266605], and *C. tetani* DnaC [NP\_780838]. Colour code: red = predicted  $\alpha$ -helical region; green = predicted  $\beta$ -stranded region; black line = unstructured. The length of the protein sequences and the sizes of the structural elements are shown to scale. The dotted line indicates structural elements of the Rossmann fold [373].

teins and may therefore have an increased affinity for DNA [385]. However, the E176G alteration in the *dnaC810* allele is located close to the Walker B-motif within the Rossmann fold and may therefore influence the reaction of the mutant protein with its bound nucleotide rather than affecting its DNA-binding property directly. Also DnaA<sub>Eco</sub> binds to ssDNA [423,474] but mediated by its C-terminal domain 4, which is not present in DnaC (F.Blaesing, and H.Seitz; unpublished). Unlike  $\lambda$  P, which interacts with the O initiator, DnaC<sub>Eco</sub> apparently does not interact with DnaA<sub>Eco</sub> in solution (F.Blaesing, and H.Seitz; unpublished); it is not known whether DnaC interacts with DnaA in the DnaB<sub>6</sub>

DnaC<sub>6</sub>ATP double-hexameric complex. During re-start primosome formation, DnaC in the DnaB<sub>6</sub>·DnaC<sub>6</sub>ATP complex interacts with the PriABC·DnaT pre-primosome, most likely with the DnaT moiety, but the DnaC subdomains involved in this interaction are not known [319]. Kaguni and co-workers presented evidence from mutational analyses that the N-terminus of DnaC<sub>Eco</sub> is responsible for the interaction with DnaB<sub>Eco</sub>; All 6 mutations characterised in this study map to the three predicted N-terminal  $\alpha$ -helices (see Fig. C13) [266].

Homologues of the *E. coli* K12 *dnaC* gene are present in the genomes of all sequenced *E. coli* strains (100% ident.), *Salmonella* (and sub-species) strains





**Fig. C14** (Putative) helicase loaders; part B.

Jpred secondary structure predictions [88] were obtained for *B. anthracis* cp $\phi$  LambdaBa04 ba0439 [NP\_842981],  $\phi$ SM1 gp21 [NP\_862860],  $\phi$ Bxz1 gene192 [NP\_818243],  $\phi$ BC6A52 bc2564 [NP\_852569],  $\phi$ 27 L18 [NP\_543070],  $\phi$ Bcep22 orf27 [NP\_944257], *A. aeolicus* DnaC [NP\_213618], *B. subtilis* p $\phi$  PBSX XkdC [NP\_389135],  $\phi$ PV83 orf21 [NP\_061611],  $\phi$ g1e Ntp [NP\_695183], *S. pneumoniae* cp $\phi$  sp1137 [NP\_345607], *S. pyogenes* cp $\phi$  spyM18\_1801 [NP\_607823],  $\phi$ LL-H orf267 [AAL77547],  $\phi$ P335 orf12 [NP\_839903],  $\phi$ bIL309 orf15 [NP\_076710],  $\phi$ bIL286 orf17 [NP\_076651], *B. subtilis* YqaM [P45910], *L. lactis* p $\phi$  pi3 pi346 [NP\_267575], *S. flexneri* cp $\phi$  sf0684 [NP\_706613], and *S. aureus* cp $\phi$  sav0869 [NP\_371393]. Colour code: red = predicted  $\alpha$ -helical region; green = predicted  $\beta$ -stranded region; black line = unstructured. The length of the protein sequences and the sizes of the structural elements are shown to scale. The dotted line indicates structural elements of the Rossmann fold [373].

(~93% ident.), *Shigella flexneri* (100% ident.), *Klebsiella pneumoniae* (93% ident.), and *Buchnera* strains (65% ident.). *Y. pestis* is a close relative of *E. coli* with which it shares most relevant replication genes (*dnaA* 87%, *dnaB* 84%, *dnaG* 74%, *priA* 73%, *dnaE* 88% ident. res.) but a *dnaC* gene is lacking from the genomes of both sequenced *Y. pestis* strains [450]. It is likely, therefore, that the 'classical' helicase loader DnaC<sub>Eco</sub> is a very specific solution for helicase loading during primosome formation in few Enterobacteria. In fact, the distribution of *dnaC* genes among  $\gamma$ -proteobacteria seems to be even more restricted than the distribution of the *B. subtilis*

*dnaI* gene among Bacillaceae (see also BRM Section 4.3.).

Annotations of 'DnaC-like' genes are found for various sequenced (pro)phage genomes in the databases. BLAST searches with DnaC<sub>Eco</sub> as query produce numerous hits (>250) due to the presence of the Rossmann fold-type nucleotide binding domain. We excluded all the hits produced by enterobacteriaceal DnaC proteins (discussed in the preceding paragraph) and by DnaA proteins from the further analysis. As mentioned above, *dnaA* and *dnaC* genes are structurally but not functionally related. Also we excluded – for simplicity – all

hits produced by proteins which are considerably smaller or larger than DnaC. For the majority of the remaining proteins, the similarity with DnaC<sub>Eco</sub> was confined to the C-terminal 2/3 of the protein (res. ~70-234), i.e. the Rossmann fold. We focussed our interest on those proteins that showed similarity with DnaC<sub>Eco</sub> also in the N-terminal portion of the protein. Among the proteins which showed identities ranging between >20% and ~40% were numerous (pro)phage-encoded 'putative replication' proteins, and – indistinguishable from the former – a significant number of *istB* genes encoding the small subunit of two-protein transposases. Because the BLAST-approach failed to detect 'DnaC-type' helicase loaders reliably we performed a JPred-based secondary structure prediction analysis for 40 of the proteins showing similarity with DnaC<sub>Eco</sub> [88,87] (Figs. C13 + C14). With few exceptions, the Jpred method could identify the ( $\alpha\beta$ )<sub>5</sub>-structure of the Rossmann fold, which made up between 1/2 and 2/3 of the entire structure of the analysed proteins. Note that in almost all cases the third  $\alpha$ -helix is predicted to be either longer than the others, or bi-partite – as was found in the known crystal structures of other AAA- or AAA+-type NTPases (see e.g.[154]). For the majority of the proteins analysed this C-terminal nucleotide-binding domain appears to be separated from the N-terminus by an unstructured stretch of ~25 res. in length. Tentatively, we consider the proteins as two-domain proteins with a flexible connecting linker. Although we find four formally distinct classes of N-terminal structures as depicted in Figs. C13 + C14 we cannot relate known or possible functions to any of these classes unequivocally. Within each class there are examples of structurally diverse but also structurally closely related proteins. In class 4, the best similarity is found between the proteins of  $\phi$ PV83 and  $\phi$ g1e. In class 3, the proteins of *Aquifex aeolicus* and  $\phi$ P27 are certainly structurally very similar. In class 1, the proteins of  $\phi$ 13,  $\phi$ SLT,  $\phi$ ETA, and  $\phi$ mv4 are structurally more related among each other than with the small group of truly DnaC-like proteins of the *E. coli* Rac, Gifsy-1 (and Gifsy-2, not shown), CP-933M, CP-933P, and CP-933R prophages. Note that the predicted structures of two *istB*-like proteins closely resemble that of DnaC<sub>Eco</sub>. We conclude therefore that neither BLAST searches nor secondary structure predictions allow for a clear discrimination between putative DnaC homologues and IstB proteins, for example. Thus, the identification of phage encoded helicase loader genes certainly requires additional criteria. We explore one such possible criterion – modular gene arrangement – in BRM Section 3. .

Neither the YdaV protein of the *E. coli* Rac prophage, nor the putative helicase loader of the *S. typhimurium* Gifsy-1 prophage were able to suppress temperature-sensitivity of an *E. coli* *dnaC28* mutant strain

(B.Wróbel, G.Wegrzyn, H.Seitz, and C.Weigel, unpublished). On the other hand, autonomously replicating *oriJ*-plasmids [105,104] lacking the *ydaV* gene can be obtained but they are lost rapidly even with selection (R.Díaz O., H.Seitz, and C.Weigel; unpublished). This may indicate that the host DnaC protein can at least partially substitute for YdaV in the replication of the *oriJ* plasmid.

**The helicase loading mechanism in *B. subtilis*.** Genetic and biochemical analyses have shown that helicase loading during primosome formation follows different routes in *E. coli* and *B. subtilis* chromosome replication, for initiation at *oriC* as well as for replication re-start. Primosome formation in *B. subtilis* involves the DnaB<sub>Bsu</sub>, DnaD<sub>Bsu</sub>, and DnaI<sub>Bsu</sub> proteins, which have no obvious homologues in *E. coli*. DnaB<sub>Bsu</sub> and DnaD<sub>Bsu</sub> bind to DNA, and both proteins oligomerise [277, 53]. DnaB<sub>Bsu</sub> and DnaD<sub>Bsu</sub> have similar C-Termini (~100 res.), which is not detectable by BLAST analysis but becomes apparent by cross-wise comparison of DnaB and DnaI proteins from other Bacillaceae (BRM Section 3.1. 1.). The DnaD<sub>Bsu</sub>-DnaB<sub>Bsu</sub> hetero-oligomer or DnaD<sub>Bsu</sub> alone interact with DNA-bound PriA<sub>Bsu</sub> *in vitro* [277], and also the interaction of DnaD<sub>Bsu</sub> with DnaA<sub>Bsu</sub> has been demonstrated by yeast two-hybrid analysis [190]. Marsin and co-workers propose a sequence of events for re-start primosome formation that might also be applicable for DnaA-primosome formation: DNA-bound PriA<sub>Bsu</sub> attracts DnaD<sub>Bsu</sub>, which in turn interacts with DnaB<sub>Bsu</sub> [277]. This pre-primosome subsequently binds DnaI<sub>Bsu</sub> through a DnaB<sub>Bsu</sub>-DnaI<sub>Bsu</sub> interaction [421,464]. DnaB<sub>Bsu</sub> and DnaI<sub>Bsu</sub> jointly recruit monomers of DnaC<sub>Bsu</sub> (the DnaC<sub>Eco</sub> orthologue) resulting in hexamerisation of the replicative helicase around single-stranded DNA [189,464, 52].

BLAST searches readily identify *dnaI*, *dnaD*, and *dnaB* genes among the members of the Bacillales and Lactobacillales sub-groups of the Firmicutes (Gram-positives). The detection of *bona fide* homologues of the *B. subtilis* *dnaI* gene in the genomes of other bacteria or in phage genomes is difficult by this approach due to conservation of the Rossmann fold in functionally analogous (*dnaC*), related (*istB*), or unrelated genes (*dnaA*). The *orf12* gene product (296 rs.) of *L. lactis*  $\phi$ r1t may represent a true DnaI<sub>Bsu</sub> homologue because it shows 23% identity over a stretch of 270 res.. Other phage-encoded helicase loaders of the DnaI-type (similarity extending beyond the Rossmann fold towards the N-terminus) could not be detected by BLAST searches. Because  $\phi$ r1t does not encode genes similar to the *L. lactis* *dnaB* and *dnaD* genes the *orf12* gene product might be used by this phage replicon to link its own replication to that of the host machinery (see Fig. C14). This hypothesis is supported by the observation that *orf11* preceding

orf12 in the  $\phi$ r1t genome shows partial similarity with the gene38 initiator of *B. subtilis*  $\phi$ SPP1 (Section C3.1.2. and BRM Section 3.1.). In addition, the latter observation supports the notion that DnaC in *E. coli* and DnaI in *B. subtilis* perform analogous functions during helicase loading. *dnaD* genes and the *dnaB-dnaI* gene pairs are exclusively found in the genomes of the Bacillales and Lactobacillales sub-groups of the Firmicutes, and not present in the genomes of Clostridia and Mollicutes (analysis restricted to the fully sequenced genomes). This makes it very likely that the helicase loading mechanism operating in *B. subtilis* chromosome replication is specific for these two bacterial subfamilies. The partial similarity of several phage-encoded (putative) initiators with *dnaD* genes, and the partial similarity of DnaB<sub>Bsu</sub> with the  $\phi$ A118 gene49 (gp49) protein is discussed in BRM Section 3.1.1. .

### C3.3. Helicases

Helicases separate the complementary strands of duplex nucleic acids in a processive manner, depending on NTP hydrolysis as energy supply. All known genomes of prokaryotes, eukaryotes, and also the genomes of many viruses and several plasmids encode helicases of several different types, underscoring the fundamental importance of helicases for all metabolic processes involving duplex nucleic acids (DNA, RNA, DNA-RNA hybrids): replication, recombination and repair, transcription, translation, and RNA splicing.

Helicases can be monofunctional proteins that interact with proteins involved in upstream and/or downstream processes, or alternatively domains of multifunctional proteins [466]. Helicases are usually classified with respect to the polarity of the strand to which they remain bound during duplex unwinding: 5'→3' or 3'→5' helicases. Helicases can function as monomers but many helicases – in particular replicative helicases – form oligomers, often homo-hexamers. The structure and function of hexameric helicases – including phage-encoded helicases – as well as present models for helicase movement have been extensively reviewed by Patel and Picha [339]. These authors also propose a very useful 'pathway' to guide studies of oligomeric helicases: i. the nature of the bound nucleotide dictates the oligomerisation state and the conformation required for DNA binding, ii. DNA binding stimulates the NTPase-activity, which in turn supplies the energy required for conformational changes that result in translocation along the bound DNA/RNA and duplex unwinding (see Fig. 1 in ref.[339]).

The replicative helicases  $\phi$ T7 gene 4A,  $\phi$ T4 gp41, and  $\phi$ P1 belong to the family 4 helicases with *E. coli* DnaB as the prototype, and will be discussed

first together with their phage-encoded homologues. As we will show then, the importance of the helicase type represented by the  $\phi$ P4  $\alpha$  primase-helicase has been somewhat underestimated in the past. Next we will discuss the phage-encoded superfamily 2 helicases (SF2), which are encoded too frequently in phage genomes to be neglected. Their function as replicative helicases has yet to be demonstrated, but the  $\phi$ N15 RepA protein is the best candidate for such a role. Finally we shall briefly discuss phage-encoded superfamily 1 (SF1) helicases, represented by the Dda helicases of the  $\phi$ T4 group of phages. For a thorough discussion of the classification of the various types of helicases by their conserved structurally and functionally important motifs we refer to Hall and Matson [158].

#### Family 4 helicases (5'→3' helicases)

*Phage T7 gene 4 helicase.* The gene A4 primase-helicase protein is involved in the initiation of replication, in DNA synthesis, and in recombination of *E. coli* phage T7 [436,222, 220,221]. The  $\phi$ T7 replisome is formed by three phage-encoded proteins: gene 2.5 SSB, gene4A primase-helicase, and gene 5 DNA polymerase together with host thioredoxin [111,369].

Gene 4 of  $\phi$ T7 encodes two proteins: the full-length gene 4A or 63 kDa protein, and – starting from codon 64 – a shorter gene 4B or 56 kDa protein that lacks the N-terminal Zn- finger motif. The gene 4B protein is active as helicase but lacks primase activity [111,29]. The gene 4A protein is absolutely required for phage propagation but optimal rates of DNA synthesis require both translation products, pointing to a physiological role of the smaller gene 4B protein [288]. Both gene 4 proteins form hexamers in the presence of nucleotide, with or without Mg<sup>2+</sup> [338,323,347]. 50-fold higher protein concentrations were required for hexamer formation in the absence of nucleotide, while hexamers were already observed at a concentration of 0.2  $\mu$ M with dTMP [338]. Furthermore, gene 4A hexamers were shown to be more compact with dTTP than with dTDP [496]. Analysis of  $\phi$ T7 gene4A protein by electron microscopy revealed a two-tiered ring with six-fold symmetry, and with one DNA strand passing through the central cavity [114,497]. At present, RepA helicase of plasmid RSF1010 and the  $\phi$ T7 gene 4 helicase domain are the only F4-type helicases for which a crystal structure is known [PDB 1Q57] [387,415]. Acidic residues in the C terminus of the  $\phi$ T7 helicase have been implicated in binding the  $\phi$ T7 DNA polymerase [315]. Deletion of the 17 C-terminal residues of the gene 4A protein abolishes this interaction, but the mutant protein retains primase and helicase activities [322].

Homologues of the  $\phi$ T7 gene4A helicase are found in all phages that also encode DNA polymerases with high

**Table C11** Family 4 helicases: subfamily of  $\phi$ T7 gene 4A-type helicases.

F4 motifs:				<b>1<sub>F4</sub></b>	<b>1a<sub>F4</sub></b>	<b>2<sub>F4</sub></b>	<b>3*<sub>F4</sub></b>	<b>4*<sub>F4</sub></b>
consensus				xhxxhhxARxxhGKS SG T	VLxhSLEM G M M E	hIhhDY L H I	LKAhAxoLxxPhxxh xQ RG T V C	DLR x SGxIxQxADxIh L S
$\phi$ T7 gene 4 subfamily consensus				EhIhhxSGSGhGKS Q V A T T	VGxhxxEE D	xIVLDxh II	LxxhxxxxxxhxxhhI I V S	DhR SGALxQhSxxhI K GI A L
replicon	accession	gene	res.					
$\phi$ T7	P03692	gene 4A	566	EVIMVTS <b>SGSGMGKS</b>	<b>VGLAMLEE</b>	VIILDHI	LKGFAKSTGVVLVVI <b>CH</b>	DLR <b>SG</b> SGALRQLSDTII
$\phi$ A1122	NP_848279	p17	566	EVIMVTS <b>SGSGMGKS</b>	<b>VGLAMLEE</b>	VIILDHI	LKGFAKSTGVVLVVI <b>CH</b>	DLR <b>SG</b> SGALRQLSDTII
$\phi$ T3	P20315	gene 4	566	EVVMVTS <b>SGSGMGKS</b>	<b>VGLAMLEE</b>	VIVLDHI	LKGFAKSTGVVLVVI <b>CH</b>	DLR <b>SG</b> SGALRQLSDTII
$\phi$ YeO3-12	NP_052087	gene 4A	566	EVIMVTS <b>SGSGMGKS</b>	<b>VGLAMLEE</b>	VIVLDHI	LKGFAKSTGVVLVVI <b>CH</b>	DLR <b>SG</b> SGALRQLSDTII
$\phi$ gh-1	NP_813761	p15	562	EVILVTS <b>SGSGGKS</b>	<b>CGVAMLEE</b>	VIVLDHI	IKAFAKTKNVAVFVI <b>CH</b>	DLR <b>SG</b> SGGLRQLSDTII
$\phi$ P60	NP_570328	P60_18	531	ELVTIT <b>AGSGTGKS</b>	<b>VGTYALEE</b>	WIVLDHL	LRSFVEETGIGMILI <b>SH</b>	QLR <b>SG</b> SHSIAQLSDLVI
$\phi$ SIO1	NP_064752	p15	522	ELITFV <b>AGTGVGKT</b>	<b>VGTFLEEE</b>	FIILDHI	LKTLTVELDICLLMV <b>SH</b>	DIR <b>SG</b> TAGIGQLSNIII
$\phi$ PaP3	NP_775217	p40	569	EIVGVGGT <b>GIGKT</b>	<b>VGTFLEEE</b>	TVILDNM	LAGMADELGIRVFIF <b>SH</b>	QFT <b>SG</b> SRAMQRWCQLMI
$\phi$ SP6	NP_853570	gp10	661	QLIIV <b>GAGSGVGKT</b>	<b>VGIISTED</b>	NIIIDNL	IGTIKDRHPVTIFLV <b>SH</b>	DFR <b>SG</b> SGAIGFWASYAL
$\phi$ K1-5	AAL84819	gp9	662	QLIIV <b>GAGSGVGKT</b>	<b>VGIISTED</b>	NIIIDNL	IGTIKDRHPVTIFLV <b>SH</b>	DFR <b>SG</b> SGAIGFWASYAL
$\phi$ Felix01	NP_944967	p188	663	EITTLA <b>APSSVGKS</b>	<b>IGVIPVED</b>	IIILDPI	LLRRCKRYQYQVNV <b>CH</b>	DIK <b>SG</b> SGAYFQISMNNI

Motif numbering is according to Hall + Matson [158]. Shorter versions of the motifs 3 and 4 were used here, indicated by their re-naming as 3\*<sub>F4</sub> and 4\*<sub>F4</sub>, respectively. Bold letters indicate matches with the conserved residues of the family 4 consensus sequences. Residues within the individual motifs that are specific for the subfamily are boxed.

similarity to  $\phi$ T7 gene 5 (Section C3.5.), and SSBs with high similarity to  $\phi$ T7 gene 2.5 (Section C3.6.1.). A common characteristic of the helicases of this subfamily is the replacement of the conserved glutamine residue in motif 3\* by histidine (Table C11).

*Gene 41 (gp41) helicase of phage T4.* Bacteriophage T4 replicates its genome using a complex of seven phage-encoded proteins: DNA polymerase (gp43), three polymerase accessory proteins (gp44, gp62, and gp45), SSB (gp32), primase (gp61), and helicase (gp41). The 54 kDa gp41 helicase is essential for DNA replication and recombination [494,381], and loaded to DNA by the gp59 helicase loader (see previous chapter).

$\phi$ T4 gp41 forms ring-shaped hexamers in the presence of nucleotide, preferably ATP, and bound nucleotide-triphosphate is required for hexamer formation [108]. The half-life for  $\phi$ T4 gp41 dissociation from DNA was measured to be around 11 min [399]. Deletion of ~20 res. from the gp41 C-terminus strongly reduced the *in vitro* interaction with the  $\phi$ T4 clamp loader and clamp proteins gp44, gp62, and gp45, respectively [371] (Section C3.5.2.).  $\phi$ T4 gp41 and gp61 are believed to form the primosome that catalyses coupled leading- and lagging-strand DNA synthesis. When assayed *in vitro*, the reciprocal stimulation of the helicase and primase activities by gp41 and gp61, respectively, could be shown, emphasising their cooperation [465,176,370].

*E. coli* DnaB and  $\phi$ T4 gp41 both belong to the F4-family of helicases, and possess highly similar biochemical properties. The sequence similarity between the two proteins is in the range of 25% (ident. res.) over the C-terminal nucleotide- and DNA-binding domains, but their N-terminal domains seem completely unrelated: there is no detectable primary sequence similarity, and also the secondary structure prediction for both sequences does not reveal any conservation of structural features (Jpred, not shown). Full-length homologues of the  $\phi$ T4 gp41 helicase sharing specific variations of the F4-helicase family signature motifs are found in all genomes of phages that also encode homologues of the  $\phi$ T4 gp59 helicase loader, of the  $\phi$ T4 gp43 DNA polymerase (Section C3.5.), of the DNA polymerase accessory proteins  $\phi$ T4 *rmh*, gp44,45,62, and of the  $\phi$ T4gp32 SSB (Section C3.6.1.). The similarity values for the proteins of this F4-type-helicase range from 45% for *Aeromonas* sp.  $\phi$ 65 to 85% for  $\phi$ RB69 (Table C12).

*Gp65 helicase of phage D29.* Four mycobacterial phages encoded F4-type helicases show specific variants of the F4 family consensus motifs, therefore we chose to place them in a distinct subfamily (Table C13). A prominent feature of these helicases is their small size (~270 res.). The direct comparison with DnaB<sub>Eco</sub> reveals that they lack the entire N-terminal domain of DnaB<sub>Eco</sub>, which was shown to be absolute re-

**Table C12** Family 4 helicases: subfamily of  $\phi$ T4 gp41-type helicases.

Family 4 motifs:					<b>1<sub>F4</sub></b>	<b>1a<sub>F4</sub></b>	<b>2<sub>F4</sub></b>	<b>3*<sub>F4</sub></b>	<b>4*<sub>F4</sub></b>
consensus					x h x h h x A R x x h G K S SG T	V L x h S L E M G M M E	h I h h D Y L H I	L K A h A x o L x x P h x x h x Q R G T V C	D L R x S G x I x Q x A D x I h L S
$\phi$ T4 gp41 subfamily consensus					<b>L</b> <b>N</b> h h h A G x N V G K S	V h Y h S M E M	I V I V D Y L V I M	h R G h A V x x x x V x W T A A Q	D h <b>A E</b> S <b>A G</b> L x <b>A T</b> A D F h L <b>D</b> <b>H G</b>
replicon	accession	gene	res.						
$\phi$ T4	NP_049654	gp41	475	<b>L</b> <b>N</b> V L M A G V N V G K S	V L Y I S M E M	I I I V D Y L	L R A L A V E T E T V L W T A A Q	D I <b>A E</b> S <b>A G</b> L P <b>A T</b> A D F M L	
$\phi$ RB69	NP_861732	gp41	458	<b>L</b> <b>N</b> V L M A G V N V G K S	V L Y I S M E M	V I M V D Y L	L R A L A V E S E T V L W T A A Q	D I <b>A E</b> S <b>A G</b> L P <b>A T</b> A D F M L	
$\phi$ RB49	NP_891593	gp41	470	<b>L</b> <b>N</b> V L L A G V N V G K S	V L Y I S M E M	V V I V D Y L	L R G F F V K W D V V G W T A A Q	D T <b>A E</b> S <b>A G</b> L P <b>A T</b> A D F M L	
$\phi$ KVP40	NP_899258	gene 41	466	<b>L</b> <b>N</b> L L L A G S N V G K S	V L Y I S M E M	I V I V D Y L	I R G F A V E H N V A V W S A A Q	D I <b>A E</b> S <b>A G</b> L A <b>H T</b> A D L I L	
$\phi$ Aeh1	NP_943891	gp41	485	<b>L</b> <b>N</b> L I L A G V N V G K S	V L Y V S M E M	I V I V D Y L	L R G F A V E H N V V W S A A Q	D I <b>A E</b> S <b>A G</b> L A <b>A T</b> A D F I L	
$\phi$ 44RR2.8t	NP_932387	gp41	469	<b>L</b> <b>N</b> V I M A G V N V G K S	V V Y F S M E M	I V I V D Y L	L R G L A V Q H Q V V L W T G A Q	D V <b>A E</b> S <b>A G</b> L P <b>A T</b> A D F M L	
$\phi$ 65	AAR90925	gp41	520	<b>L</b> <b>N</b> V I L A G T G V G K S	V L Y V S F E M	I I M V D Y L	V R G L A I E K K L V A W S G A Q	D I <b>A D</b> S <b>Y A</b> L L <b>H G</b> C D F V L	

Motif numbering is according to Hall + Matson [158]. Shorter versions of the motifs 3 and 4 were used here, indicated by their re-naming as 3\*<sub>F4</sub> and 4\*<sub>F4</sub>, respectively. Bold letters indicate matches with the conserved residues of the family 4 consensus sequences. Residues within the individual motifs that are specific for the subfamily are boxed.

**Table C13** Family 4 helicases: subfamily of  $\phi$ D29 gp65-type (putative) helicases.

F4 motifs:					<b>1<sub>F4</sub></b>	<b>1a<sub>F4</sub></b>	<b>2<sub>F4</sub></b>	<b>3*<sub>F4</sub></b>	<b>4*<sub>F4</sub></b>
consensus					x h x h h x A R x x h G K S SG T	V L x h S L E M G M M E	h I h h D Y L H I	L K A h A x o L x x P h x x h x Q R G T V C	D L R x S G x I x Q x A D x I h L S
$\phi$ D29 gp65 subfamily consensus					Q L h L h C A G <b>G</b> T G K S A G	S T R A <b>V R</b> E Q A A G	L I V h D <b>N</b> h	L <b>H E</b> x x R E T G S <b>C</b> h h G L <b>H H</b>	<b>L S G</b> h K G A L G R V <b>P</b> E h h h
replicon	accession	gene	res.						
$\phi$ D29 <sup>1)</sup>	NP_046882	gp65	[232]	Q L V L V C A G <b>G</b> T G K S	S T R A <b>V R</b> E Q	L I V V D <b>N</b> I	L <b>H E</b> M A R E T G S <b>C</b> V I G L <b>H H</b> <b>L G G</b> I K G Q I G R V <b>P</b> E M V L		
$\phi$ L5	NP_039729	L5p62	268	Q L V L V C A G <b>G</b> T G K S	A T R A <b>V R</b> E Q	L I V V D <b>N</b> I	L <b>H E</b> M G R E T G S <b>C</b> V V G L <b>H H</b> <b>L S G</b> I K G Q I G R V <b>P</b> E M I L		
$\phi$ Bxz2	NP_817653	gp64	269	Q L A L V C A G <b>G</b> T G K S	S A R A <b>V R</b> E G	L V V I D <b>N</b> V	L <b>H D</b> M A R S T S A <b>C</b> V V G L <b>H H</b> <b>L S G</b> V K G Q I S R V <b>P</b> E M I L		
$\phi$ Bxb1	NP_075324	gp57	269	Q L A L I A A A <b>P</b> G G A K S	S A R A <b>V R</b> E G	L I V I D <b>N</b> I	L <b>H E</b> K A R E T G A <b>C</b> I I G L <b>H H</b> <b>L S G</b> I K G Q I G R V <b>P</b> E L V A		
pRSF1010	NP_044309	RepA	279	T V G A L V S P G G A G K S	V I Y L P A E D	Q A V A D G L	M E A I A A D T G C S I V F L H H	G V S K A N Y G A P F A D R W F	

Motif numbering is according to Hall + Matson [158]. Shorter versions of the motifs 3 and 4 were used here, indicated by their re-naming as 3\*<sub>F4</sub> and 4\*<sub>F4</sub>, respectively. Bold letters indicate matches with the conserved residues of the family 4 consensus sequences. Residues within the individual motifs that are specific for the subfamily are boxed. 1) Tentative translation of the sequence upstream of the rare start codon UUG used to define the N-terminal Met residue of  $\phi$ D29 gp 65 of this NCBI entry results in a complete motif 1<sub>F4</sub> (underlining) and full-length similarity in the N-terminus to e.g. the  $\phi$ L5 L5p62 protein. Hence the length of the protein is shown in brackets. The sequence of the RSF1010 RepA helicase is included to allow for an easy comparison.

quired for *in vitro* helicase activity of the latter [37]. Small size is also a feature of the RepA helicase (F4-type) of plasmid RSF1010 with which the  $\phi$ D29 gp65-type helicases share no detectable sequence similarity, though (<20% ident. res.; see Table C13 for comparison). Purified RepA<sub>RSF1010</sub> binds nucleotide, shows GTP-dependent 5'→3' helicase activity in the presence

of divalent cations, and forms a hexameric ring structure like DnaB<sub>Eco</sub> [392]. The crystal structure of RepA<sub>RSF1010</sub> protein is known [PDB 1G8Y] [320]. In contrast to DnaB<sub>Eco</sub>, hexamer formation by RepA does not require bound nucleotide, and nucleotide binding is co-ordinated between two adjacent protomers in the hexamer, while DnaB<sub>Eco</sub> protomers bind nucleotide individually

**Table C14** Family 4 helicases: subfamily of *E. coli* DnaB-type helicases.

F4 motifs:				<b>1<sub>F4</sub></b>	<b>1a<sub>F4</sub></b>	<b>2<sub>F4</sub></b>	<b>3*<sub>F4</sub></b>	<b>4*<sub>F4</sub></b>
consensus				xhxhhxAR x xhGKS SG T	VLxhSLEM G MM E	hIhhDYL HI	LKAhA x oLxxPhxxh x Q RG T V C	DLRxSGxIxQ x ADxIh L S
DnaB <sub>Eco</sub> consensus				xLhhhAAR [P]xMGKT GG V S	hhhhSLEM M	hIhhDYL H	LKxhA [K]ELxxPVhxL [S]Q [R]D V [I]	DLRESGxIEQ [D]ADhhh S DT C
replicon	accession	gene	res.					
φ P1	YP_006543	ban	454	DLIIVAAR [P]SMGKT	VLVFSLEM	MIMIDYL	LKALA [K]ELQVPVVAL [S]Q	DLRESGATEQ [D]ADLIM
φ D3	NP_061570	orf74	446	QMIVIAGR [P]AMGKT	VLVISLEM	LIVIDYL	MKLLA [R]EIGCPVLP L C Q	DLRESGATEQ [D]ADIVM
φ SPP1	NP_690733	gene 40	442	NFVLIAGR [P]SMGKT	VNLHSLEM	IVMIDYL	LKKMA [R]ELDVVVIAL [S]Q	DLRESGQLEQ [D]ADIIE
φ HK97	NP_037740	gp55	458	DLVIAAR [P]GMGKT	VLIFSLEM	LIMADYL	LKAMA [K]DLKTPVISL [S]Q	DLRDSGSIEQ [D]ADSII
φ P27	NP_543071	L19	463	DLVFIAAR [P]SMGKT	VLLFTMEM	LVVVDYL	LKGLA [K]SGGFPLIAL [S]Q	DLKNSGIEA [D]ADIIL
φ Nil2	CAC95089	P	458	DLVIAAR [P]GMGKT	VLIFSLEM	LIMADYL	LKAMA [K]DLRTPVFSL [S]Q	DLRDSGSIEQ [D]ADSII
φ P22 <sup>1)</sup>	NP_059611	gene 12	458	DLVIAAR [P]GMGKT	VLIFSLEM	LIMADYL	LKAMA [K]DLKTPVISL [S]Q	DLRDSGSIEQ [D]ADSII
φ ST104	YP_006384	gene 12	458	DLVIAAR [P]GMGKT	VLIFSLEM	LIMADYL	LKAMA [K]DLKTPVISL [S]Q	DLRDSGSIEQ [D]ADSII
φ 3626	NP_612864	orf35	427	EFIVLGAR [P]SMGKT	VLYIQ LDM	VIVDHI	LKSIA [K]ELNVAMVGL C Q	DLRDTGSIEE [D]ADVIG
φ 11	NP_803269	16	413	QLIVIAAR [P]SVGKT	TSFFSLET	VIFIDYL	LKIIANETGAIIVLL [S]Q	DMKESGGIEA [D]ASLAM
φ CJW1	NP_817531	gp82	420	NLVVVGGR [P]GAGKS	SLVFSLEM	VVMVDYA	LKVAA [K]QLHMVVVLA [S]Q	NLAGGDSLGR [D]ADVVL
φ Lc-Nu	AAG43560	orfA	416	RLLTIGAR [P]GVGKS	VDMFSLEM	LAIVDYL	FKVLTNELGIPVLL [S]Q	DLRESGSIEQ [D]SNAV
φ MAV1	NP_047257	RepB	276	QLYIFAGR [P]ATGKT	LVFFSLEM	AIFIDHL	LKKLA [K]KLNVNVFAL [S]Q	DLRESGSIEQ [D]ADVVL
<i>B. cereus</i> φ	NP_830742	bc0956	431	DFVVLGAR [P]SMGKT	VGLFSLEM	LIIVDYL	LKLLA [R]ELNVCVVAL [S]Q	DLRETGQIEQ [D]ADVIM
φ 315.1	NP_664494	SpyM3_0690	446	QLIILAAR [P]AMGKT	VAVFSLEM	LIVIDYL	LKILA [K]ELKVPVIAL [S]Q	DLRESGSIEQ [D]ADIVA
φ HK022 <sup>2)</sup>	NP_037693	P	478	SLFVIGAR [P]KMGKT	ALMFSLEM	MILVDYL	LKNLA [K]ELDCVVVLL [I]Q	DSRDTGQIEQ [D]CDYVW
φ Aaφ23	NP_852741	p	485	SLVALGAR [P]KCGKT	VLLFSLEM	FIGIDYL	LKNLA [R]EMDCVVVLL [I]Q	DSRDTGQIEQ [D]CDYWL
<i>E. coli</i>	P03005	<i>dnaB</i>	471	DLIIVAAR [P]SMGKT	VLIFSLEM	LIMIDYL	LKALA [K]ELNVPVVAL [S]Q	DLRESQSIEQ [D]ADLIM
<i>B. subtilis</i>	P37469	<i>dnaC</i>	454	DLIIVAAR [P]SVGKT	VAIFSLEM	MILIDYL	LKSIA [R]ELQVPVIAL [S]Q	DIRESGSIEQ [D]ADIVA

Motif numbering is according to Hall + Matson [158]. Shorter versions of the motifs 3 and 4 were used here, indicated by their re-naming as 3\*<sub>F4</sub> and 4\*<sub>F4</sub>, respectively. Bold letters indicate matches with the conserved residues of the family 4 consensus sequences. Residues within the individual motifs that are specific for the subfamily are boxed. 1) a virtually identical helicase (98% ident. res.) is encoded by phage φST104 (YP\_006384). 2) virtually identical helicases (>98% ident. res.) are encoded by phages φSf6 (NP\_958219), φVT2-Sa (NP\_050526), φHK620 (NP\_112058), and φST64T (NP\_720303). The sequences of the *E. coli* DnaB and the homologous *B. subtilis* DnaC helicase are included to allow for an easy comparison.

[320]. It has been reported that RepA performs a C3–C6 conformational switch despite the lack of an equivalent of the DnaB<sub>Eco</sub> N-terminal domain (cited in [320]). Since none of the φD29 gp65-type helicases has been characterised structurally or biochemically so far, the results obtained with RepA<sub>RSF1010</sub> may provide a rough guideline for their understanding. Also in this subfamily the conserved glutamine residue in motif 3\* is replaced by histidine (Table C11).

*E. coli* DnaB and phage-encoded homologues. The DnaB protein (52 kDa) is the replicative helicase of *E. coli*, and required for both initiation and elongation of DNA replication [246,228]. All completely sequenced bacterial genomes have highly similar *dnaB* genes, and in addition a number of plasmid and phage replicons.

Homologues of DnaB are also encoded by several mitochondria and chloroplast replicons. Although readily detected by sequence comparison few DnaB proteins have been characterised.

Biochemical methods including gel-filtration and cross-linking, and also studies by electron microscopy could establish that *E. coli* DnaB forms hexamers in solution [367,240, 241,8,54,383,498,382]. In the absence of Mg<sup>2+</sup> ions and nucleotides DnaB forms mostly trimers, which can be converted to hexamers by the addition of Mg<sup>2+</sup> or high concentration of salt [54]. When bound to nucleotide, DnaB oligomerises to hexamers with threefold symmetry (C3) and sixfold symmetry (C6) in addition to intermediates [498]. Also, the preferred conformation is greatly influenced by the pH of the sample buffer [107]. The significance of the

two symmetry forms and their interconversion for helicase function are not well understood. The C6 symmetry was found for the DnaB<sub>6</sub>C<sub>6</sub> double-hexamer, and it may represent the inactive conformation because the ATPase activity of DnaB is inactivated in the DnaB<sub>6</sub>C<sub>6</sub> complex [20].

The dimensions of the *E. coli* DnaB hexamer (Ø: 12.5–14 nm, cavity: 3–4 nm, height: 5.7 nm; taken from [339]) change in the presence of different nucleotides. Hexamers formed in the presence of ADP are somewhat more compact than those formed in the presence of AMP-PNP [498]. When bound to DNA, the DnaB hexamers appear even more compact [197].

As discussed in Section C3.1.2., the physical interaction between *E. coli* DnaA and DnaB was shown to involve two domains of either protein: res. 24–86 (domain 1) and 130–148 (domain 3) of DnaA and res. 154–210 (βγ-fragment) and 1–156 (α-fragment) of DnaB, respectively [278,404]. Kaguni and co-workers presented evidence from mutational analyses that the N-terminus of DnaB<sub>Eco</sub> is responsible for the interaction with DnaB<sub>Eco</sub>, probably with a site within the γ-fragment [266] (Section C3.2.). From the observation that the similarity between the λ P helicase loader and its φ80 homologue, gene 14 protein, is considerably higher in the N-terminus than in the C-terminus, Ogawa and co-workers deduced that the C-termini of the helicase loaders could interact with their cognate initiators, and the conserved N-termini with a host protein, possibly DnaB [329]. These results support the hypothesis of Chatteraj and co-workers that the known helicase loaders contact DnaB via their N-termini [326]. In contrast however, λ π mutations with the phenotype of reduced affinity to DnaB map exclusively in the C-termini of mutant P proteins [224]. Certainly however, the interaction of DnaB with the different types of helicase loaders is even more complex: image reconstruction for electron micrographs of DnaB<sub>6</sub>C<sub>6</sub> double-hexamers suggests multiple (putative) sites of mutual interactions, with one DnaC protomer contacting two neighbouring DnaB protomers of the hexamer [20].

*E. coli* DnaB can unwind duplex DNA *in vitro* without any accessory proteins but this activity is stimulated by DnaG primase [265]. The interaction between the *E. coli* DnaB and DnaG has always been difficult to analyse *in vitro* because of instability of the complex. Tougu and Marians could show that the DnaG C-terminus is responsible for interaction with DnaB [455]. By analysing mutant DnaB proteins in priming reactions, Chang and Marians found that residues in the DnaB N-terminus (α-fragment) are responsible for a productive DnaB-DnaG interaction [69]. This is at variance with the observation of Bastia and co-workers who found that a region partially overlapping with the interaction sites for plasmid initiators and DnaA in the

N-terminus of the βγ-fragment is important [1]. A recent study by Sultanas and co-workers suggests that the flexible linker – connecting the DnaB α- and βγ-fragments – is involved in the interaction with DnaG in addition to several residues in the N-terminal α-fragment, and it seems possible that up to three DnaG monomers bind to one DnaB hexamer [449].

In the nomenclature of Hall and Matson [158], the signature motifs 1<sub>F4</sub> and 1a<sub>F4</sub> of Family 4-helicases are located in the nucleotide-binding β domain of DnaB<sub>Eco</sub> (res. 157–302), and represent the 'Walker A' and 'Walker B' motifs, respectively; motifs 2<sub>F4</sub> to 4<sub>F4</sub> are located in the DnaB<sub>Eco</sub> γ domain (res. 303–471), responsible for DNA binding [37] (See Tables C11 – C14, C16).

With 78% identical residues, the φP1 ban (DnaB analogue) protein shows extremely high overall similarity to *E. coli* DnaB, in the 3'-half of the gene also on the DNA level. This is in fact the highest similarity found so far between a (pro)phage-encoded protein and a host-encoded orthologue (Table C14). Not surprisingly therefore, ban greatly influenced the development of models for horizontal gene transfer (HGT) among prokaryotic replicons. Following its detection by genetic experiments [89,328], biochemical analysis established that ban and DnaB are orthologues with very similar properties, which form functional hetero-hexamers [242,400,454]. Only recently it was shown that ban supports growth of an *E. coli* dnaB<sup>o</sup>(null) mutant strain albeit not at all physiological conditions [249]. The role of ban for φP1 replication is not clear: in the (plasmid) prophage state, ban is repressed, and suppression of a host dnaB mutation was first observed in a φP1 mutant expressing ban constitutively [168]. Since the availability of *E. coli* DnaB is restricted to ~4–8 hexamers/cell under normal growth conditions it seems possible that de-repression of the gene encoding the cognate helicase allows for phage propagation under certain conditions – in analogy to the situation found for the φP1 SSB [26].

As discussed in Section C3.1.2 & C3.2., the initiation of replication of *B. subtilis* phage SPP1 requires the joint action of the phage-encoded G38P initiator, the G39P helicase loader, and the G40P helicase. G40P bound to the unwound *oriL* recruits the host DnaG primase, and both proteins together the host DNA polymerase III holoenzyme [341,297]. The DNA unwinding activity of G40P was shown *in vitro*, and – similar to φT4 gp41 and *E. coli* DnaB – the helicase-activity of G40P is stimulated by the DnaG primase, suggesting that the primase and helicase interact during DNA replication [11,13]. G40P was shown to form a ring-shaped hexamer with a two-tiered structure by electron microscopy [19]. Comparable to *E. coli* DnaB, two types of hexamers were observed, one with a C6- and the other with a C3-symmetry.



**Table C15** Family 4 helicases: miscellaneous phage-encoded (putative) helicases.

F4 motifs:				<b>1<sub>F4</sub></b>	<b>1a<sub>F4</sub></b>	<b>2<sub>F4</sub></b>	<b>3*<sub>F4</sub></b>	<b>4*<sub>F4</sub></b>
consensus				xhxhhxARxxhGKS SG T	VLxhSLEM G MM E	hIhhDYL HI	LKAhAxoLxxPhxxhxQ RG T V C	DLRxSGxIxQxADxIh L S
replicon	accession	gene	res.					
φ Rosebush	NP_817815	gp54	920	GLS <b>C</b> IIGPPGV <b>GKS</b>	VLYMP <b>G</b> EG	EVGN <b>D</b> LL	FDQIRRH <b>T</b> NAGVLI <b>V</b> HH	TGPTGTVDPMQ <b>G</b> DIVL
φ PG1	NP_943837	gp59	915	GLT <b>S</b> IIGPPG <b>I</b> GKS	VLYLP <b>G</b> EG	DLGN <b>D</b> LM	YDKLREL <b>T</b> GAGVCV <b>V</b> HH	TGPN <b>S</b> VDPMQ <b>G</b> EIVL
φ RM378	NP_835606	p019	416	ELGIVMLPSGW <b>GKS</b>	V <b>I</b> YFT <b>L</b> EL	VVLID <b>Y</b> A	<b>L</b> R <b>L</b> IAKVYNTAV <b>S</b> AS <b>Q</b>	YIAD <b>S</b> FAKV <b>V</b> EID <b>F</b> GM
pφ SPBc2	NP_046682	yorl	504	KFYLRSSITGG <b>GKT</b>	STV <b>I</b> ST <b>E</b> M	YVY <b>F</b> D <b>Y</b> I	ILL <b>L</b> MSDKL <b>G</b> LCNK <b>Y</b> D	Y <b>L</b> R <b>G</b> S <b>K</b> AIAD <b>K</b> T <b>D</b> AAM
φ T5	YP_006947	D6	502	ETLL <b>L</b> GGWR <b>G</b> T <b>GKS</b>	APY <b>F</b> S <b>I</b> <b>E</b> M	AK <b>M</b> SDF <b>Y</b>	AR <b>Y</b> GDK <b>V</b> T <b>V</b> ALLD <b>Y</b> IN <b>Q</b>	D <b>G</b> R <b>T</b> R <b>M</b> SK <b>G</b> IL <b>D</b> S <b>A</b> D <b>M</b>
φ VpV262	NP_640278	p17	408	DH <b>V</b> I <b>I</b> FGPV <b>N</b> AG <b>S</b>	VLY <b>V</b> GN <b>E</b> D	IC <b>I</b> VD <b>Q</b> A	<b>L</b> R <b>M</b> LY <b>K</b> R <b>M</b> K <b>I</b> V <b>G</b> SV <b>T</b> Q	V <b>F</b> GS <b>R</b> RE <b>V</b> AA <b>Q</b> AD <b>V</b> M <b>I</b>
φ Bxz1	NP_818244	gp193	418	ELAA <b>I</b> A <b>A</b> Y <b>T</b> K <b>V</b> G <b>S</b>	VM <b>.</b> FT <b>L</b> <b>E</b> M	FVLID <b>Q</b> L	ET <b>S</b> RS <b>S</b> V <b>G</b> EL <b>P</b> T <b>F</b> L <b>A</b> I <b>Q</b>	N <b>F</b> AN <b>S</b> S <b>M</b> IE <b>Q</b> T <b>V</b> D <b>I</b> A <b>L</b>
φ K	YP_024501	orf71	480	EV <b>G</b> L <b>I</b> A <b>P</b> T <b>G</b> R <b>G</b> S	VLY <b>I</b> A <b>L</b> <b>E</b> <b>E</b>	V <b>V</b> I <b>I</b> D <b>Y</b> P	IR <b>R</b> LS <b>Q</b> Q <b>Y</b> GF <b>V</b> CW <b>T</b> L <b>A</b> Q	H <b>V</b> E <b>G</b> S <b>R</b> K <b>I</b> V <b>N</b> A <b>V</b> E <b>V</b> S <b>L</b>
φ KMV	NP_877454	orf15	397	AS <b>V</b> LV <b>A</b> AP <b>P</b> D <b>A</b> G <b>K</b> T	ILAL <b>D</b> <b>P</b> <b>E</b> <b>E</b>	V <b>V</b> F <b>W</b> D <b>M</b> M	V <b>R</b> EM <b>A</b> VR <b>H</b> D <b>F</b> IS <b>F</b> M <b>T</b> W <b>Q</b>	CL <b>K</b> D <b>S</b> K <b>T</b> AV <b>Q</b> GA <b>V</b> D <b>V</b> Q
φ Xp10 <sup>1</sup> )	NP_858988	xp10p41	261	DS <b>I</b> I <b>W</b> A <b>A</b> R <b>P</b> D <b>Q</b> G <b>K</b> T	HE <b>Q</b> G <b>V</b> L <b>E</b> T	-	-	-
	NP_858987	xp10p40	139	-	-	AN <b>A</b> P <b>D</b> E <b>V</b>	VR <b>D</b> Y <b>A</b> S <b>I</b> E <b>G</b> Y <b>A</b> A <b>I</b> N <b>T</b> S <b>Q</b>	ML <b>K</b> D <b>S</b> K <b>T</b> A <b>K</b> A <b>G</b> A <b>C</b> E <b>A</b> I

Motif numbering is according to Hall + Matson [158]. Shorter versions of the motifs 3 and 4 were used here, indicated by their re-naming as 3\*<sub>F4</sub> and 4\*<sub>F4</sub> respectively. Bold letters indicate matches with the conserved residues of the family 4 consensus sequences. 1) The two orfs encoding the N- and C-terminus of the putative φXp10 helicase, respectively, are separated by an untranslated stretch of 82 bp; the comparison to the highly similar φKMV orf14 protein suggests that this gene split is not a sequencing artefact.

The function of the small RepB helicase (276 res.) for the replication of the mycobacterial phage φMAV1 (15.6 kb) is not known. However, RepBMAV1 adds another example to the group of smaller versions of DnaB-type helicases, but this one is more related to DnaB<sub>Eco</sub> than RepBRSF1010, or φD29 gp65 (see above, and Table C14).

*Miscellaneous phage-encoded family 4 helicases.* BLAST searches with various DnaB<sub>Eco</sub>-type helicases as query yielded matches with several phage-encoded (putative) helicases that – considering the reasonable conservation of the signature motifs of the F4-type helicases – clearly belong to this family (Table C15). Because they lack the specific 'signature residues' of the DnaB-type subfamily within the motifs we placed them in a somewhat heterogeneous subfamily. The ~500 C-terminal residues of the φRosebush gp54 and φPG1 gp59 proteins contain the signature motifs and show similarity to the other helicases, the function of the extended N-termini of these proteins is not known (Section C3.4.).

**Helicases with similarity to phage P4 alpha (α) primase-helicase.** Replication of the circular φP4 prophage *in vitro* requires the phage-encoded α protein, host (*E. coli*) DNA PolIII holoenzyme, SSB, and gyrase; the host replication proteins DnaA, DnaB, DnaC, and DnaG are not required [103]. Thus, the initiator, helicase and

primase functions are performed by the α protein (reviewed in [47]). φP4 Replication *in vitro* proceeds (mostly) bidirectionally from *ori*, and elongation occurs probably as coupled leading- and lagging-strand synthesis, without a requirement for host DnaG primase and DnaB helicase [103]. Therefore, α is not only responsible for priming of leading-strand synthesis, but also for priming of lagging-strand synthesis. This suggests, in turn, an association of α with *E. coli* DNA PolIII holoenzyme. In short: α is the replicative helicase for φP4 replication. However, and unlike *E. coli* DnaB, φP4 α has been shown to unwind DNA in the 3'→5' direction *in vitro* (with respect to the polarity of the strand to which it binds), and its oligomeric status during helicase action is not known [506].

When expressed as isolated domain, the N-terminus of φP4 α could be shown to retain its primase function [505], discussed in the following chapter (Section C3.4.). Lanka and co-workers could show that the isolated C-terminus of α still binds specifically to DNA, but helicase activity could only be detected for constructs expressing the nucleotide-binding and the DNA-binding domains in a contiguous polypeptide [504]. This situation is reminiscent of RNA polymerases where the polymerising activity also cannot be uncoupled from the helicase activity [466]. The specific binding of α to the octamer repeats in the φP4 *ori* qualifies this protein as initiator, as discussed in Section C3.1.2. . We will refer

**Table C16** (Putative) helicases of the  $\phi P4\alpha$ -type.

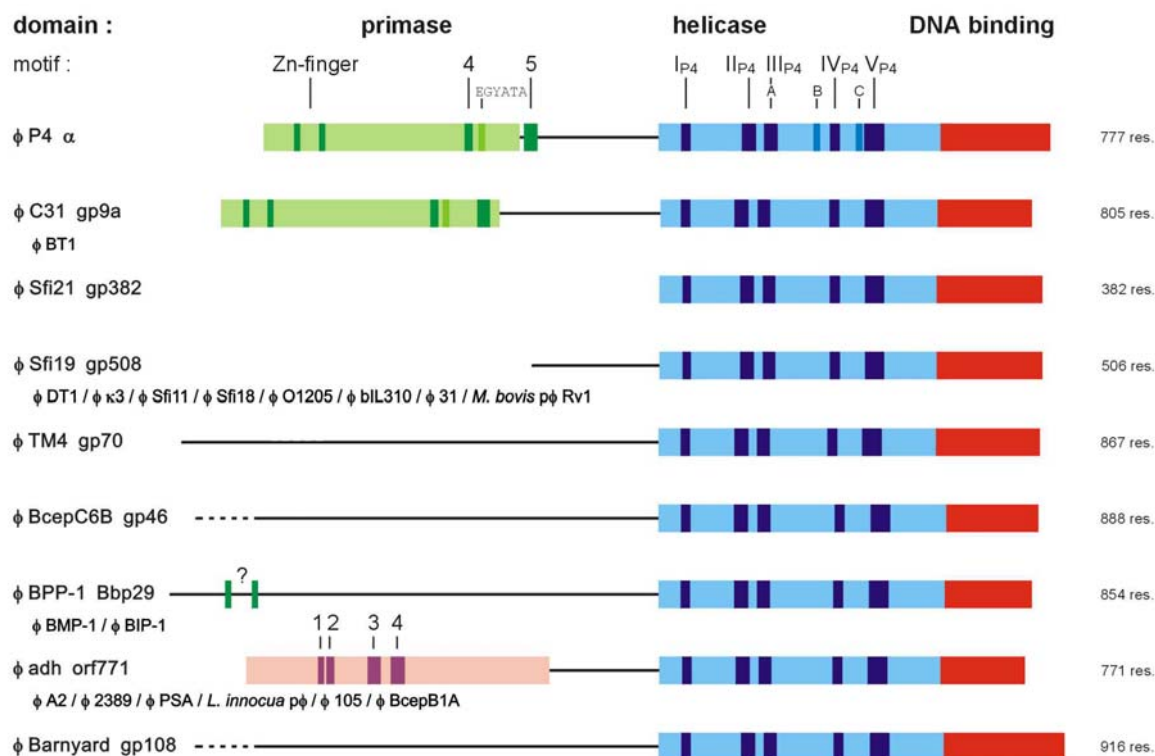
$\phi P4$ $\alpha$ motifs:				I <sub>P4</sub>	II <sub>P4</sub>	III <sub>P4</sub>	IV <sub>P4</sub>	V <sub>P4</sub>
consensus				hxxxNGhhD S N	ELxxxhQxhhGYSh W AA	hhhhGxGxNGKS D G	hKxhTGGDxh N VT E	Px.hRxxxxGhWRRhxxhhPF K Y
replicon	accession	gene	res.					
$\phi P4 / \phi R73$	P10277	alpha	777	IGFRNGVLD	EKRVDILAALFMVL	LEVTPGGSGKS	LKAITGGDAV	PMRFTDRSGGVSRRRVLIHF
$\phi BPP-1$ <sup>1)</sup>	NP_958698	Bbp29	854	LGVNGVVD	DMIGFFQRLIGYSL	AIPYGSNGSNGKS	IKSMTGGEPL	PI.VKGDDHAIWRRLLLPVPF
$\phi BcepC6B$	AAT38405	gp46	888	LGVNGAVD	EQVEFFQRLVGYAL	IIPHGTGSNGKS	IKAMTGGDPI	PI.VKGDDHAIWRRLLMLVPF
$\phi BT1$	NP_813724	gp9a	808	LSFRNGVVD	ELAEYMQRLTGYGI	AVLWKGKGSNGKS	LKRVTGKDKV	PK.FKSQDEGLWRRVVKLIPF
$\phi C31$	NP_047955	gp9a	805	LSFANGVVD	DLVGYMRRLVGYGI	AVLWKGKGSNGKS	LKRVTGKDKV	PK.FKSQDEGLWRRVVKLIPF
$\phi TM4$	NP_569803	gp70	867	LNTPSGVVN	SMVEYVQRLAGYAA	PFLFGAGSNGKS	VKLLTGGDVL	PD.VKAGGTSFFRRFRLIPF
$\phi BcepB1A$	AAT37747	-	919	IGAPNCIID	ELIVFFKRLMGYAL	VIPYGHGANGKS	VKSLTGGDDTI	PV.IKGSNDGIWRRIMMIPF
<i>X. fastidiosa</i> $\phi P$	ZP_00040564	xfaso0360	847	LNCTNGTVD	PLSDFLQRWFGYCA	AVMYGMGRNGKS	VKQATGGDSL	PV.IKQDQVGIWSRLMLIPF
$\phi Barnyard$	NP_818646	gp108	916	LAVANGVVE	SYRRDVQIILGHAL	IIFKGAANTGKS	MKTATGNDYI	PE.IEGHDKALRELRVISF
$\phi A2$	NP_680517	orf35	770	LNTINGYVD	ELIDYLQKAIIGYSL	FILYGNRNGKS	VKELTGGDMV	PI.IRGTDGIWRRLLMLIPF
$\phi 2389 / \phi PSA$	NP_511033	pri	757	LNTQNGYIN	ELINYIQKAVGYSL	FILFNGRNGKS	VKQLTGGDKV	PI.IRGRDDGIWRRLLHLVPF
$\phi adh$	NP_050131	orf771	771	LNTPSGYID	ELIHVQKLIIGYSL	FILYGNRNGKS	VKQMTGGDTL	PK.IYGTDEGIWRRLLVLPF
<i>L. innocua</i> $\phi P$	NP_471917	lin2587	757	LNTQNGYIN	ELINYMOKAVGYSL	FILFNGRNGKS	VKQLTGGDKV	PI.IRGRDDGIWRRLLHLVPF
$\phi 105$	NP_690795	orf11	806	FNCENGVID	EIIIEFLQKAIIGYSL	FFLFNGRNGKS	VKQITGGDKM	PI.VKGSDEGIWRRIRLVPF
$\phi DT1 / \phi \kappa 3$	NP_049424	orf36	504	L.VKNGIYD	ELVELLWQVIAASL	IWLVGNGNDGKG	FNSVVTGEPV	PV.FKNKSNGTYRRRIVLPF
$\phi Sfi11$ <sup>2)</sup>	NP_056712	orf504	504	L.VKNGIYD	ELVELLWQVIAASL	IWLVGNGNDGKG	FNSVVTGEPV	PV.FKNKSNGTYRRRIVLPF
$\phi Sfi19$	NP_049955	orf508	506	L.VKNGIYD	ELVELLWQVIAASL	IWLVGNGNDGKG	FNSVVTGEPV	PV.FKNKSNGTYRRRIVLPF
$\phi Sfi21$	NP_050002	orf382	382	L.VKNGIYD	ELVELLWQVIAASL	IWLVGNGNDGKG	FNSVVTGEPV	PV.FKNKSNGTYRRRIVLPF
$\phi bIL310$	NP_076775	orf24	542	IPVANGIFN	ELVSLWQIISAST	VWLVGKNGNDGKG	YFSVVTGDPV	PK.FRNKSNGTYRRLLIVPF
$\phi 31$	CAC04163	primase	492	EVESTASRD	ELVKLLWQVISASL	IWFVGEKNDGKG	FNSVVTGEPV	PK.VRNKNTGTYRRFLIIPF
<i>M. bovis</i> $\phi P$ Rv1	NP_855261	Mb1608c	471	LNVANGTLD	GVRGFVQRLAGVGL	AILIGVGGANGKS	IKRLTGGDDTI	PR.VPGDDTAIWRIRRVVPF

Motifs I - V were detected by BLAST alignment of the sequences as the regions with the highest degree of conservation. Bold typing indicates particularly well conserved residues. 1) identical sequences are present in *Bordetella* sp. phages  $\phi BPP-1$ , and  $\phi BIP-1$ . 2) identical sequences are present in *S. thermophilus* phages  $\phi Sfi18$ , and  $\phi O1205$ .

to the central nucleotide-binding domain as 'helicase domain' in the following, and to the C-terminus as 'DNA-binding domain'.

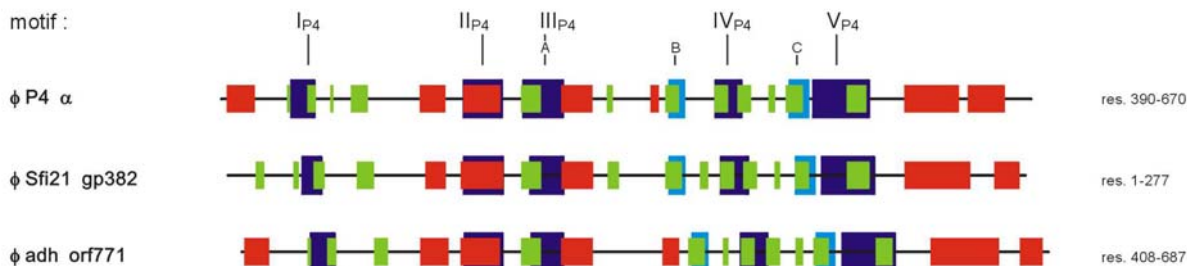
BLAST searches with  $\phi P4 \alpha$  as query (entire sequence or res. 390-777) led to the detection of >20 phage-encoded proteins with significant similarity, and many more (putative) prophage sequences in the genomes of bacteria from different phyla. Within the helicase domain ( $\phi P4 \alpha$ : res. 390-670) we could detect five particularly well conserved short stretches, which we propose as signature motifs I<sub>P4</sub>-V<sub>P4</sub> although it is presently not possible to correlate them – except for motif III<sub>P4</sub> – with known structural or functional elements (Table C16). Mutations of G506 or K507 in motif III<sub>P4</sub> were found to result in a reduction or complete loss of nucleotide hydrolysis proficiency and helicase activity of  $\alpha$  protein [504]. Motif III<sub>P4</sub> 'hhhhGxGxNGKS' corresponds to the 'A' (Walker A) motif 'hhhxGPxGTGKS' proposed ear-

lier by Koonin and co-workers for a group of helicases that are now placed into the helicase superfamily 1 (SF1) [148,158]. Although a rather formal approach, it should be noted that among the variants of the 'Walker A' or 'P-loop' motifs found in nucleotide-binding proteins – e.g. the AAA+ family – the specific positioning of the glycine residues is highly conserved within the subfamilies. We tend to place the  $\phi P4 \alpha$ -type helicases in a separate family rather than grouping them into SF1, therefore. The signature motifs 'B' and 'C' suggested by Koonin and co-workers were also detected but considered less useful due to their degeneracy (see Figs. C15 + C16) [148]. The choice of signature motifs I<sub>P4</sub> – V<sub>P4</sub> is supported by comparing the secondary structure predictions for three  $\phi P4 \alpha$ -type helicases that showed <25% identity to each other: all motifs are positioned within a fairly similar (predicted) secondary structure context (Fig. C16).



**Fig. C15** Anatomy of  $\phi$ P4  $\alpha$ -type (putative) helicases.

All (putative) helicase sequences, and the size/position of the signature motifs are drawn to scale. For accession numbers and signature motif numbering see Table C16 for the helicase motifs, and Tables C23 + C25 for the primase motifs; the A, B + C motifs proposed by Gorbalenya *et al.* for  $\phi$ P4  $\alpha$  are indicated at their respective positions for this protein only [148]. In addition, the 'EGYATA' motif conserved in some primases is indicated (see text for details [431, 335]). Line length indicates the size of the protein (dotted line: protein sequence longer than shown), domains are shown as coloured blocks. Colour code: red = DNA-binding domain of  $\phi$ P4  $\alpha$  [PDB 1KA8] and homologous regions detected by BLAST comparison + secondary structure prediction Jpred (not shown; [88]); blue = helicase domain, the  $\phi$ P4  $\alpha$ -type helicase signature motifs are shown in dark blue (size of the domain tentatively adjusted to the smallest detected helicase domain of  $\phi$ Sfi21 gp382); green = primase domain, the  $\phi$ P4  $\alpha$ -type primase signature motifs are shown in dark green (domain size in  $\phi$ P4  $\alpha$  is shown according to the minimal primase domain determined by Ziegelin *et al.* [505]); pink = conserved domain of unknown function, the signature motifs are shown in dark pink. Phages encoding (putative) helicases with identical or highly similar anatomy are indicated below the prototype sequence by smaller font size. '?' weak similarity to Zn-finger motif.



**Fig. C16** Secondary structure prediction for the  $\phi$ P4  $\alpha$  helicase domain.

The (putative) helicase domain sequences, and the size/position of the signature motifs are drawn to scale. For signature motif numbering see Table C16. The A, B + C motifs proposed by Gorbalenya *et al.* for  $\phi$ P4  $\alpha$  are indicated at their respective positions [148]. Secondary structure predictions were obtained from the Jpred server [88]. Colour code: red = predicted  $\alpha$ -helical region; green = predicted  $\beta$ -stranded region. Light blue = A, B + C motifs according to Gorbalenya *et al.*; dark blue = signature motifs proposed in this study (see Text for details).

**Table C17** Superfamily II (SF2) helicases:  $\phi$ T4 UvsW-type helicases.

SF2 motifs:				I <sub>SF2</sub>	Ia <sub>SF2</sub>	II <sub>SF2</sub>	III <sub>SF2</sub>	VI <sub>SF2</sub>
consensus				hhxxxoGxGKT S	xhnhxPoo	hhhDExH	hxhSATxxx TGS	QxxGRxxR
replicon	accession	gene	res.					
$\phi$ T4	NP_049796	UvsW	587	ILNLPT <b>S</b> AGKS	ILIIVPTT	MMNDECH	FGLSGSLRD	QTIGRVLR
$\phi$ RB69	NP_861886	UvsW	504	ILNLPT <b>S</b> AGKS	ILIIVPTT	MMNDECH	FGLSGSLRD	QTVGRVLR
$\phi$ RB49	NP_891741	UvsW	500	MLVLPT <b>S</b> AGKS	VLIIVPTT	VIVDECH	IGMTGSPRD	QSIGRALR
$\phi$ 44RR2.8t	NP_932529	UvsW	493	ILNLPT <b>S</b> AGKS	VLVIVPTT	VMNDECH	YGLSGSLRD	QTIGRVLR
$\phi$ KVP40	NP_899623	UvsW	507	LLNLPT <b>S</b> AGKS	VLIIVPTT	LLVDECH	IGLTGSRLD	QSIGRILR
$\phi$ Aeh1	NP_944130	UvsW	503	IGNLPT <b>S</b> AGKS	VLVLVPTT	IIVDECH	IGMSGSLRD	QSIGRVLR

Motif numbering is according to Hall + Matson [158]. Motifs IV and V were omitted due to their high degree of degeneration. Bold letters indicate matches with the conserved residues of the SF2 consensus sequences. Residues within the individual motifs that are specific for the subfamily are boxed.

The proteins with similarity to  $\phi$ P4  $\alpha$  show a considerable variation in size indicating that they are probably also multi-domain proteins (see Table C16). A tentative anatomy of these proteins is shown in Figure C15, using  $\phi$ P4  $\alpha$  as 'template'. The length of the helicase domain is very similar in all proteins (~270 res.), and the signature motifs are evenly spaced.  $\phi$ Sfi21 gp382 was the smallest protein detected in the databases with full-length similarity to  $\phi$ P4  $\alpha$  and might therefore represent the 'minimal version' of this type of helicases. All proteins shown in Figure C15 have a C-terminal domain with significant similarity in length and sequence to the  $\phi$ P4  $\alpha$  DNA-binding domain.  $\phi$ Sfi21 gp382 and  $\phi$ Sfi19 gp508 lack any N-terminal extension and are unlikely to function as 'primase-helicase' proteins. In analogy to the function of  $\alpha$  for  $\phi$ P4 replication, we speculate that these proteins function as 'initiator-helicases' (Section C3.1.2.).  $\phi$ C31 gp9a is a full-length homologue to  $\phi$ P4  $\alpha$ , containing a clearly detectable N-terminal primase domain (see next chapter). The function of the long N-terminal extensions – responsible for the size variation – of the remaining proteins is not known, but they might represent yet unknown primase domains (see next chapter).

In the past, the  $\alpha$  protein of satellite phage P4 has been studied mainly as a rare example of a multifunctional replication protein, comparable to the large T antigen of simian virus 40 (SV40). Here we can show that  $\alpha$ -type helicases are found as widespread among phage replicons as the DnaB<sub>Eco</sub>-type helicases. In contrast to DnaB however, homologues of  $\phi$ P4  $\alpha$  are not found in the sequenced bacterial genomes, except for the many  $\alpha$ -type proteins encoded by (putative) prophages. It should be interesting to learn why chromosomal replicons rely on DnaB as their replicative helicase, and (with few exceptions) on DnaA as replication initiator,

for which a phage-encoded homologue has not yet been found – contrary to the expectation of Campbell and Botstein [62].

**Phage-encoded superfamily 2 helicases.** Next to genes encoding F4-type helicases, genes encoding superfamily 2 helicases (SF2) are found most frequently in phage genomes. Clear evidence for the participation of SF2-type helicases in the initiation and elongation steps of replication has not yet been presented, but the SF2-type helicase UvsW of  $\phi$ T4 may participate in the later (RDR) stages of  $\phi$ T4 replication. SF2-type helicases homologous to  $\phi$ T4 UvsW are present in the other completely sequenced genomes of the  $\phi$ T4 group phages (Table C17), and could be detected in further >30 (pro)phage genomes (Table C18). The listing in Table C18 is complete for phage-encoded SF2 helicases, but numerous (putative) prophage sequences detected in the fully sequenced genomes of species from all phyla have been omitted. We note that for several phage genome entries in the NCBI database the SF2-type genes were annotated 'putative helicase', and the genes showing similarity to  $\phi$ P4- $\alpha$  annotated 'putative primase' although the region of similarity did not include the N-terminal primase domain of  $\phi$ P4 $\alpha$  (e.g.  $\phi$ Sfi21,  $\phi$ Sfi11) [102]. This imprecision calls for more expert assistance during gene function annotation procedures.

In addition to the SF2-type helicases with a size of ~450-650 res., four large phage-encoded proteins (~1300 res.) with significant homology to the RepA protein of  $\phi$ N15 were detected, which show partial similarity to the  $\phi$ P4  $\alpha$  primase-helicase in their N-terminal (putative) primase domain, but have a central nucleotide binding domain that – as judged from the signature motifs – suggests their grouping in the superfamily 2 (Table C18). Very recent studies suggest that these pro-

**Table C18** Phage-encoded superfamily II (SF2) helicases.

SF2 motifs:				I <sub>SF2</sub>	Ia <sub>SF2</sub>	II <sub>SF2</sub>	III <sub>SF2</sub>	VI <sub>SF2</sub>
consensus				hhxxxoGxGKT S	xhhhxPoo	hhhDExH	hxxSATxxx TGS	QxxGRxxR
replicon	accession	gene	res.					
φ K	AAO47518	orf69	582	IAHLATNGGKT	VAFFTGST	MIVDEAH	IALTGSIDK	QRIGRALR
φ Bcep1	NP_944368	gp60	614	IAAIVTGGGKS	VMSLAPSM	VLVDEAH	VGMTATDFI	QIVGRGLR
φ Bcep781	NP_705681	gp58	614	IAAIVTGGGKS	VMSLAPSM	VLVDEAH	VGMTATDFI	QIVGRGLR
φ Bcep43	NP_958163	gp57	614	IAAIVTGGGKS	VMSLAPSM	VLVDEAH	VGMTATDFI	QIVGRGLR
φ T5	P11107	D10	450	IINGKPGFGKT	TLVICNT	VIVDEVH	IGLSGTLKR	QLAGRVQR
φ PG1	NP_943831	gp53	559	GVVLPTGTGKS	VLMVAHRG	ILWDEVH	AGFTATMRR	QMIGRALR
φ Rosebush	NP_817811	gp50	571	GVVLPTGAGKS	VVALAHRA	ILWDEFH	CGFTATMYR	QMVGRALR
φ VP16T	AAQ96511	orf44T	550	LAVLPTGAGKT	HRVIAPAK	SVYDEGH	LFVTATPER	QMVGRALR
φ VP16C	AAQ96576	orf45C	718	LAVLPTGAGKT	HRIIAPAK	GVYDEGH	LFVTATPER	QMVGRALR
φ BcepB1A	AAT37757	-	614	MAVSATGSGKT	ARAGIPHN	LIGDEGH	VLFTATPER	QQWGRVLR
φ Fels-1	NP_459877	STM0900	527	MVYAPTGGGKT	VMFCVPYQ	LIIDEAH	IGLSGTFFS	QAIGRGLR
φ Asp2	AAT36801	Pas53	567	LVLVPTGTGKT	IGMVAPEL	VWVDEAH	VGFTATPER	QMIGRGTR
φ 105	NP_690792	ORF8	311	GFIGPSGSGKT	MQEAYPDA	GVVDTEH	VIDSLSHNW	nd
<i>L. innocua</i> pφ	NP_471919	lin2589	418	CCSLWLGGGKS	NDVDMKFV	IIIDESH	VGFTATPVR	QQSMRGMR
φ adh	NP_050128	orf455	455	LIQSPAGSGKS	WGVNMLC	ILVDEAH	LGFATPVR	QQSMRCMR
φ 2389 / φ PSA	NP_511031	hel	418	RCSLWVGAGKS	NEVDMSYV	IIIDESH	VGFTATPVR	QQSMRGMR
φ 31	CAC04160	helicase	448	IVQSPPE <sup>1</sup> SGKT	NQVNMEYV	ILVDEAH	LFFTGTPIR	QFAMRPLN
φ O1205	NP_695088	orf10	443	MVQSPPE <sup>1</sup> SGKT	NGVDMNLV	ILIDEAH	LMFTGTTPAR	QFAMRALN
φ Sfi19	NP_049952	orf443	443	MVQSPPE <sup>1</sup> SGKT	NGVDMNLV	ILIDEAH	LMFTGTTPAR	QFAMRALN
φ Sfi21 <sup>1</sup> )	NP_049999	orf443	443	MVQSPPE <sup>1</sup> SGKT	NGVDMNLV	ILIDEAH	LMFTGTTPAR	QFAMRALN
φ κ3	AAL74161	hel3.1	443	MVQSPPE <sup>1</sup> SGKT	NGVDMNLV	ILIDEAH	LMFTGTTPAR	QFAMRALN
φ DT1	NP_049421	orf33	443	MVQSPPE <sup>1</sup> SGKT	NGVDMNLV	ILIDEAH	LMFTGTTPAR	QFAMRALN
φ A2	NP_680515	orf33	455	LLVSPAGSGKS	DEIDLQAT	IITDESH	LGFASAPWR	QQSMRGMR
φ LE1	CAE14775	LE1-0071	501	LIIMATGTGKT	SGLERAEB	IIIDECH	VGLTATPDG	QMVGRGTR
φ T1	YP_003923	22	672	FKASVSAGKT	SIFSASLG	IGIDECH	FGMTGSEFR	QLLGRGMR
φ mi7-9	P24125	-	452	MIQSPPGSGKS	KVHGVP LN	IITDEGH	LGFATPWR	QQSMRSMR
φ RM378	NP_835691	p104	670	LIADEMGTGKT	VLVCPAS	VIVDECH	LFLTGTPIV	QAEDRLHR
φ Bamyard	NP_818603	gp65	653	VIADQPLGKT	ILVVAPKT	VIVDESQ	IAISGTPFR	QVENRAYG
pφ BC6A51	NP_831634	bc1861	397	IGVAPAGSGKT	PGIGRIGF	VVVEAH	IGLTATPSR	nd
φ 12 <sup>2</sup> )	NP_803332	p26	423	<u>GLFLD</u> <sup>2</sup> MGLGKT	MLVIAPKQ	VVIDEL <sup>2</sup> S	IGLTGTSPSP	QANARLYR
φ BPP-1	NP_958716	Bbp47	464	AVWAGMGMGKT	SLVLAPLR	VVADES <sup>1</sup> T	IELTGTSPSP	QIIERIGP
φ APSE-1	NP_051002	P41	460	NVWAGMGMGKT	TLVLAPLR	IIADES <sup>1</sup> T	VNLTGTSPSP	QIIERIGP
φ BcepNazgul	NP_919000	gp47	513	ALFIDLGMGKT	CLVIAPLR	VFIDES <sup>1</sup> S	HLLTATPAA	QLIGRLAR
<i>X. fastidiosa</i> pφ	NP_299806	xf2528	472	NLFVPMGLGKT	MLVIAPLR	VVADECS <sup>1</sup> S	IGLTGTSPSP	QIIERIGP
<i>S. pyogenes</i> pφ	AAR83209	-	458	AVILDMMGMGKT	VLVIAPLR	VVIDEL <sup>1</sup> S	VGLTGTSPSS	QTNARLWR
<i>S. pyogenes</i> pφ	NP_801687	sps0425	442	GLLLDMGLGKT	ILIVAPKK	VVIDEL <sup>1</sup> S	VGLTGTSPAP	QANARLDR
φ VP5 / φ VP2	AAR92075	-	489	ALFMEMGTGKT	VLVIVPTP	VVVEDS <sup>1</sup> S	MILTGTPIIT	QSEDRNHR
φ PY54	NP_892081	RepA	1330	IVRAGMGSGKS	AKDIAPYA	FGLDEA <sup>1</sup> T	ESIAMTDDH	QMLRRDRT
φ KO2	YP_006615	gp35	1328	IVRAGMGSGKS	YQEMAPYA	FGFDEA <sup>1</sup> T	DALARTEDH	QMLRRDRT
φ VP882	AAS38505	RepA	1239	ILRAPMGSGKT	MAVEMPYT	LCIDEA <sup>1</sup> S	AAMRSRRV	QMMRRDRT
φ N15	AAC48876	RepA(gp37)	1228	IVRAGMGSGKS	YQEMAPYA	FGFDEA <sup>1</sup> T	DALARTEEH	QMLRRDRT

Motif numbering is according to Hall + Matson [158]. Motifs IV and V were omitted due to their high degree of degeneration. Bold letters indicate matches with the conserved residues of the SF2 consensus sequences. Residues within the individual motifs that are specific for the subfamily are boxed. nd: not detected. 1) Homologues of the φSfi21 orf443 protein are encoded by phages φSfi18 and φSfi11 (including identical signature motifs). 2) the M<sub>AUG</sub> in motif I<sub>SF2</sub> residue (boxed) was assigned as start codon for this NCBI entry; alternative translation starting from a rare start codon located 28 triplets upstream results in a complete motif I<sub>SF2</sub> (indicated by underlining) and full-length similarity in the N-terminus to the Bbp47 protein of φBPP-1.

**Table C19** Superfamily I (SF1) helicases:  $\phi$ T4 Dda-type helicases.

SF1 motifs:				I <sub>S1</sub>	II <sub>S1</sub>	III <sub>S1</sub>	V <sub>S1</sub>	VI <sub>S1</sub>
consensus				hhxGxAGoGKS A P T	hhhDExo	hhhhhGDxoQ	xxThxxxQGhohooV S K	VAhTRxoo G S
replicon	accession	gene	res.					
$\phi$ T4	NP_049632	Dda	439	TINGPAGT <b>GKT</b>	LICDEVS	TIIGIGDN <b>KQ</b>	ASTFHKA <b>QG</b> MSVDRA	<b>VG</b> VTRGRY
$\phi$ RB69	NP_861705	Dda	437	TINGPAGT <b>GKT</b>	LLFDEAS	TVIGIGDR <b>CQ</b>	VSTFHKA <b>QG</b> MSVDTA	<b>VG</b> TTRGRF
$\phi$ RB49	NP_891582	Dda	463	TVRGAAGT <b>GKT</b>	IIVDEAS	VIIAVGD <b>LYQ</b>	ACTIHKS <b>QG</b> ISVDRV	<b>VG</b> VTRARY
$\phi$ KVP40	NP_899402	Dda	411	TISGPAGS <b>GKT</b>	LIVVEAS	RILAVGD <b>KYQ</b>	ASTVHKS <b>QG</b> TTVKGV	<b>VGL</b> TRPTD
$\phi$ Aeh1	NP_944136	Dda	454	TISGPAGS <b>GKS</b>	LIVDEAS	QIIAIGD <b>KHQ</b>	ASTFHKS <b>QG</b> TTVIGV	<b>VG</b> CTRAQK
$\phi$ 44RR2.8t	NP_932377	Dda	439	TIRGPAGT <b>GKT</b>	LICDEVS	VIIGIGD <b>KAQ</b>	AMTFHKS <b>QG</b> STFKNA	<b>VG</b> NTRARE
$\phi$ PaP2	YP_024782	p54	683	ALTGAAGT <b>GKT</b>	LIIEEAS	QIIYLGD <b>INQ</b>	CLTVHKS <b>QG</b> SEWRKV	<b>TAL</b> TRAKE

Motif numbering is according to Hall + Matson [158]. Motifs I<sub>SF1</sub> + IV<sub>SF1</sub> were omitted due to their degeneracy. Bold letters indicate matches with the conserved residues of the SF1 consensus sequences.

teins perform multiple functions for the replication of their cognate replicons comparable to  $\phi$ P4  $\alpha$ , including a role of RepA as replicative helicase [507,274].

**Phage-encoded superfamily 1 helicases.** The Dda helicase of  $\phi$ T4 has been characterised as a monomeric helicase with 5'→3' DNA duplex unwinding activity [360,317]. Together with the  $\phi$ T4 UvsX SAP, Dda can rescue stalled replication forks by re-activating DNA synthesis by  $\phi$ T4 gp43 DNA polymerase when the polymerase encounters a block in the template strand. This reaction may require two sequential template-switches, and Dda and UvsW may provide redundant functions [202]. Homologues of the  $\phi$ T4 Dda helicase are encoded by all (completely sequenced) phage genomes from the  $\phi$ T4-group (Table C19). Except for the  $\phi$ PaP2 gp54 protein, other phage-encoded proteins with significant similarity to the  $\phi$ T4 Dda helicase were not detected. The  $\phi$ PaP2 gp54 protein contains all signature motifs of SF1-type helicases, but has a N-terminal extension of ~150 res. of unknown function, and may be more closely related to the RecD-type of exodeoxyribonuclease V subunits.

The PcrA helicase of Gram(+) bacteria, and the Rep [PDB 1UAA] and UvrD helicases of Gram(-) bacteria are three closely related bacterial helicases belonging to the helicase superfamily 1. PcrA and Rep have been shown to be essential for the propagation of plasmids and phages that replicate by RCR (BRM Section 2.1.), and this has now also been shown for several *E. coli* plasmids in the case of UvrD [51]. None of these helicases plays a direct role during chromosome replication and their involvement in plasmid and phage replication is a long-standing enigma. Ehrlich and co-workers propose that these helicases may counter-act recombination

processes of the RecFOR pathway that would upset chromosome replication [270,345]. Phage-encoded homologues of Rep or PcrA are not yet known.

### C3.4. Primases

Primases are DNA-dependent RNA polymerases that synthesise short oligo-ribonucleotides on ssDNA as template, which are then elongated by DNA polymerases. In contrast to RNA polymerases, all known replicative DNA polymerases are incapable of *de novo* synthesis, and therefore require priming (Section C3.5.). Given the sterically correct positioning of a 3'-hydroxyl group, priming of leading-strand DNA synthesis can occur in principle by any one of five different ways:

- by the 3'-end of DNA in a recombination intermediate (D-loop) or in a partially single-stranded region;
- by the 3'-end of an aborted RNA transcript which persists annealed to its template (R-loop);
- by the CCA-OH-3' end of a tRNA annealed to a unwound stretch of DNA (cauliflower mosaic virus (CAMV));
- by the hydroxyl group of an amino acid side chain (serine, threonine, tyrosine) of a DNA-bound protein ( $\phi$ 29);
- by the 3'-end of a short RNA transcript synthesised in a template-dependent manner on unwound DNA by a primase.

In all studied cases, however, priming of lagging-strand synthesis during coupled leading- and lagging-strand synthesis is performed by a primase. Accordingly, all phage replicons that employ coupled leading- and lag-

**Table C20** Phage-encoded primases of the  $\phi$ T7 gene 4A primase-type.

motifs:					Zn finger		4	5	
consensus					xHxxCx	1.6	xCxxCx	xEGEhDxh	FDxDxxGxxAxx
replicon	accession	gene	res.	ident. res.					
$\phi$ T7	P03692	gene 4A	566	[1-245]	YHIPC <b>D</b>	1.6	FCYV <b>C</b> E	TEGEIDML	FDMDEAGRKA <b>V</b> E
$\phi$ A1122	NP_848279	p17	566	100%	YHIPC <b>D</b>	1.6	FCYV <b>C</b> E	TEGEIDML	FDMDEAGRKA <b>V</b> E
$\phi$ T3	P20315	gene 4	566	82%	FHAP <b>C</b> E	1.6	WCFV <b>C</b> E	TEGEIDAL	FDMDDAGRKA <b>V</b> E
$\phi$ YeO3-12	NP_052087	gene 4A	566	85%	YHIPC <b>E</b>	1.6	YCYA <b>C</b> E	TEGEIDAL	FDMDDAGRKA <b>V</b> E
$\phi$ gh-1	NP_813761	p15	562	56%	KHVP <b>C</b> E	1.6	FCFAC <b>D</b>	TEGEIDCL	FDMDDAGRAAS <b>Q</b>
$\phi$ P60	NP_570328	P60_18	531	36%	RHEPC <b>P</b>	1.6	YCFSC <b>G</b>	TEGEFDAL	FDNDDAGIQAA <b>E</b>
$\phi$ SIO1	NP_064752	p15	522	33%	KHQPC <b>Q</b>	1.5	YCHV <b>C</b> K	TEGEFDAM	FDNDKAGQDA <b>A</b> M
$\phi$ Felix01	NP_944967	p188	663	27%	ACPRC <b>G</b>	1.7	TCFSC <b>N</b>	WEGEMECA	FDNDEAGAKA <b>T</b> K
$\phi$ PaP3	NP_775217	p40	569	25%	ACPGC <b>R</b>	2.4	KCNRC <b>G</b>	SEDELSAM	HDADEEGRKS <b>V</b> E
$\phi$ SP6	NP_853570	gp10	661	28%	PCPAC <b>Q</b>	1.9	YCNRG <b>H</b>	VGGE <b>L</b> DAL	FDGDEIGQK <b>L</b> NQ
$\phi$ K1-5	AAL84819	gp9	662	30%	PCPAC <b>Q</b>	1.9	YCNRG <b>H</b>	VGGE <b>L</b> DAL	FDGDEVGQK <b>Q</b> NQ

The presence of the N-terminal zinc finger motif of  $\phi$ T7 gene 4 primase is shown together with the spacing of the cysteine residues probably involved in zinc ion binding (indicated by gray shading). The signature motifs 4 + 5 correspond to the conserved functional important motifs of the TOPRIM domain [CDD 8119] defined by Aravind *et al.* [9]. Highly conserved residues are indicated by bold type.

**Table C21** Phage-encoded primases of the  $\phi$ T4 gp61 primase-type.

motifs:					Zn finger		4	5*	
consensus					xCxxCx	1.6	xCxxCx	hEGPhDSh	KDhNDhh
replicon	accession	gene	res.	ident. res.					
$\phi$ T4	P04520	gp61	342	[1-342]	RCPV <b>C</b> G	2.2	HCY <b>N</b> CN	LEGPIDSL	KDVNDMI
$\phi$ RB69	NP_861720	gp61	340	73%	RCPV <b>C</b> G	2.1	HCY <b>N</b> CQ	MEGPIDSL	KDVNDMV
$\phi$ RB49	NP_891587	gp61	342	51%	SCP <b>I</b> CG	2.2	GCF <b>N</b> CN	LEGPIDSL	KDVNAMI
$\phi$ KVP40	NP_899268	gp61	352	40%	RCP <b>I</b> CG	2.2	LHL <b>R</b> CS	VEGPIDSL	KDINDFI
$\phi$ Aeh1	NP_943889	gp61	344	46%	RCH <b>I</b> CG	2.2	GCF <b>N</b> CG	FEGPIDSV	KDINDII
$\phi$ 44RR2.8t	NP_932385	gp61	334	53%	RCP <b>I</b> CG	2.2	HCF <b>N</b> CG	MEGPLDSL	KDINDMI
$\phi$ 65	AAR90927	gp61	348	43%	KCP <b>L</b> CG	2.1	GCF <b>N</b> CG	LEGPIDSV	KDINDMV

The presence of the N-terminal zinc finger motif of  $\phi$ T4 gp61 primase is shown together with the spacing of the cysteine residues probably involved in zinc ion binding (indicated by gray shading). The signature motif 4 corresponds to the conserved functional important motifs of the TOPRIM domain [CDD 8119] defined by Aravind *et al.* [9], and was detected by alignment of the  $\phi$ T4 gp1 and  $\phi$ T5 *pri* sequences. Motif 5\* was tentatively assigned to a conserved motif with the approximate position of motif 5 in a comparable secondary structure context. Highly conserved residues are indicated by bold type.

ging-strand DNA synthesis for their propagation either depend on the host primase or encode their cognate primase.

Because primer synthesis by primases depends on a single-stranded DNA template they are often closely associated with replicative helicases – phages provide examples for multifunctional proteins consisting of primase and helicase domains (see previous chapter). Although a physical interaction was also demonstrated for some primases with their cognate replicative DNA polymerase, a primase domains within a multi-domain DNA

polymerase has not yet been detected. A closer association with DNA polymerases exists for eukaryotic primases, which can be purified as two distinct subunits of the four-protein Pol  $\alpha$  complex (for references see [129]). The N-terminal primase domain of the  $\phi$ T7 gene 4A primase-helicase is among the best understood primases at the molecular level, has been studied extensively by various biochemical approaches, and is the only primase for which a complete crystal structure is available [204].



We discuss the different types of phage-encoded primases, including homologues of  $\phi$ T7 gene 4, of  $\phi$ T4 gene 61 (gp61), phage-encoded DnaG<sub>Eco</sub>-type primases, and  $\phi$ P4  $\alpha$ -type primase domains. The structure and function of bacterial and phage-encoded primases has also been reviewed recently by Frick and Richardson [129].

**The gene 4A primase domain of phage T7.** As early as 1969, it became known that  $\phi$ T7 gene 4 is essential for replication [435], and that the protein shows primase activity [432]. It was shown that gene 4 mutants synthesise only small amounts of DNA hybridising exclusively to one parental strand – the leading-strand template as we know today. As primase,  $\phi$ T7 gene 4A protein has the preferred recognition sequence 5'-GTC on the single-stranded template, and the primers synthesised *in vitro* or *in vivo* and attached to Okazaki fragments were found to be mostly tetra-ribonucleotides with the sequence 5'AC(A/C)(AC) [432,412,391,405,330].

Deletion analyses could show that the primase activity of gene 4A protein resides in the N-terminus, and the helicase activity in the C-terminus. These results were in accordance with earlier observations, that the two translation products of gene 4 occurring *in vivo* differ with respect to these two activities depending on the start codon used. The truncated N-terminus (271 res.) retained the primase activity but lacked helicase activity [128]. The crystal structure was determined for a polypeptide containing the first 182 res. of the gene 4A primase domain, showing its organisation in two separate domains: the N-terminal Zn-binding domain, and the RNA polymerase domain [PDB 1NUI] [204]. Interestingly, primer extension by  $\phi$ T7 DNA gene 5 polymerase is already stimulated by the isolated Zn-binding domain [204].

A physical interaction between the  $\phi$ T7 gene 4A protein and the acidic C-terminus of the gene 2.5 SSB has been demonstrated *in vitro* [212] (Section C3.6.1.). The  $\phi$ T7 gene 2.5 SSB bound to ssDNA stimulates the rate of primer synthesis by gene 4A *in vitro*, and also stimulates primer extension by gene 5 DNA polymerase in a coupled system [316,289].

As mentioned in the preceding chapter, gene 4 of  $\phi$ T7 encodes two proteins: the full-length gene 4A or 63 kD protein, and – starting from codon 64 – a shorter gene 4B or 56 kDa protein that lacks the N-terminal Zn-finger motif. The gene 4B protein is active as helicase but lacks primase activity [111,29]. BLAST searches revealed the presence of genes encoding homologues of the  $\phi$ T7 gene 4 primase-helicase in all phage genomes that also encode DNA polymerases of the  $\phi$ T7-subfamily of the Pol I-type DNA polymerases (Table C11; see also BRM Section 3.4.). The  $\phi$ SP6

gp10 primase has been characterised biochemically and found to be largely comparable to  $\phi$ T7 gene 4A primase except that primer synthesis occurs preferentially at 5'-GCA trinucleotide template sequences [457]. A Met codon at pos. 64 as in  $\phi$ T7 gene 4 is present in the corresponding genes of phages  $\phi$ A1122,  $\phi$ T3, and  $\phi$ YeO3-12, and the translation pattern of these genes might be comparable to that of  $\phi$ T7 gene 4 (see above). Met codons close to the position of the secondary start codon in  $\phi$ T7 gene 4 were found in the  $\phi$ SIO1 p15 gene, and the  $\phi$ P60 gene 18, but not in the other primase-helicase genes of this subfamily (Table C20).

**Gene 61 (gp 61) of phage T4.** Genetic analyses of  $\phi$ T4 replication mutants and inspection of their progeny DNA molecules by electron microscopy led to the detection of gene 61, responsible for the phenotype of a reduced rate of DNA synthesis and accumulation of ssDNA [48,309,160]. The product of  $\phi$ T4 gene 61 (gp61) was then shown to be responsible for the requirement for NTPs in a T4 DNA *in vitro* replication system, qualifying gp61 as primase [302].

$\phi$ T4 gp61 primase is a monofunctional protein that requires gp41 helicase for optimal activity [321,465,176,370] reminiscent of the  $\phi$ T7 gene 4A protein containing both activities as separate domains in a contiguous polypeptide. The analysis of the 5'-ends of Okazaki fragments synthesised *in vivo* after infection of *E. coli* cells with  $\phi$ T4 or in an *in vitro* system revealed that  $\phi$ T4 gp61 primase synthesizes mostly pentanucleotide primers with the sequence 5'-ACNNN [236, 332,321,255]. Similar to the interaction of  $\phi$ T7 gene 4A primase with gene 2.5 SSB, also  $\phi$ T4 gp61 primase and gp32 SSB protein could be shown to interact *in vitro*: both proteins form a stable complex in the absence of DNA, and the negatively charged C-terminus of gp32 is responsible for this interaction (Section C3.6.1.) [56,57]. There is a complex interplay of gp61 with gp32 on one hand, and of both with the accessory proteins gp44, gp45, and gp62 (clamp-loader and clamp proteins; see Section C3.5.2) on the other hand. This interplay is thought to be instrumental in synchronising lagging-strand primer synthesis with subsequent primer elongation by gp43 DNA polymerase [66].

Except for *Aeromonas* sp.  $\phi$ 25 (sequence not yet complete), all phages of the  $\phi$ T4 group encode homologues of the gp 61 primase, and all proteins of this subfamily contain the N-terminal Zn-finger motif and the 'TOPRIM domain' motifs (Table C21).

**Phage-encoded *E. coli* DnaG-type primases.** *E. coli* DnaG was the first bacterial protein shown to catalyse the synthesis of oligo-ribonucleotides required for *in vitro* replication [42]. Two standard assays for the analysis of primase activity *in vitro* were established by

**Table C22** Phage-encoded (putative) primases of the *E. coli* DnaG-type.

motifs:				Zn finger		4	5
consensus				x <b>C</b> xx <b>C</b> x	1.6 x <b>C</b> xx <b>C</b> x	h <b>EG</b> Ph <b>D</b> Sh	FDx <b>D</b> xx <b>G</b> xx <b>A</b> xx
replicon	accession	gene	res.				
<i>E. coli</i>	NP_417538	<i>dnaG</i>	581	<b>CC</b> PF <b>HN</b> 15	<b>HC</b> FG <b>CG</b>	VEGYMDVV	YDGDRAGRDAAW
φ P4 / φ R73	P10277	alpha	777	R <b>H</b> Q <b>P</b> CP 19	Y <b>C</b> N <b>Q</b> CG	LEGQNQAG	ADRDLSGDGQKK
φ C31	NP_047955	gp9a	805	L <b>C</b> PA <b>HS</b> 18	T <b>C</b> RAGC	TEGPGDAL	GDNDTAGVGFLL
φ BT1 <sup>1</sup> )	NP_813724	gp9a	808	V <b>C</b> PA <b>HA</b> 18	T <b>C</b> RAGC	TEGPGDAL	GDNDTAGTGFTT
<i>B. subtilis</i> pφ SPBc2	NP_046683	<i>yorJ</i>	378	YRTV <b>CH</b> 18	<b>HC</b> Y <b>TE</b> C	nd	LDYKDSPADKGG
φ T5	YP_006949	<i>pri</i>	296	LNP <b>D</b> HD 15	<b>HC</b> L <b>S</b> CG	VEGIFDML	LDNDASGNKAAQ
φ Bxz1	NP_818246	gp195	327	IVHS <b>CL</b> 26	V <b>C</b> WAFW	GEGRIDRV	ADDDAGRFMER
φ K	AAO47525	orf76	355	<b>CC</b> PF <b>CG</b> 19	<b>HC</b> KK <b>CD</b>	TEGVFDAL	LDTDALDNNIDL
φ Xp10	NP_858989	xp10p42	280	RHA <b>CG</b> 18	Y <b>CH</b> RCG	TEDAISAY	LDNDTGHSSGSN
φ KMV	NP_877453	orf14	274	HV <b>L</b> GC <b>Q</b> 20	Y <b>C</b> Y <b>S</b> C <b>Q</b>	TEDYLSAL	LDGDPAGVRGSA
φ VpV262	NP_640276	p15	287	V <b>C</b> PC <b>CR</b> 17	LAY <b>RC</b> Y	nd	LDADATSKAASM
φ T1	YP_003921	24	306	P <b>C</b> PN <b>CG</b> 19	I <b>C</b> NS <b>CG</b>	TESFDVGI	SDKDTLYMADDR
pφ Fels-1	NP_459878	stm0901	322	E <b>C</b> PV <b>CG</b> 16	I <b>C</b> -V <b>CG</b>	LQEDNYLD	ADRDENASATGLA
φ RM378	NP_835688	p101	310	L <b>C</b> PK <b>CA</b> 19	I <b>C</b> FR <b>CG</b>	FEGMFDML	FSDSVSVEEISK
φ PY54	NP_892081	RepA	1330	nd	nd	IIGKLRGA	LDNDQKSAAEGK
φ KO2	YP_006615	gp35	1328	nd	nd	VIGDLKGA	LDNDRKSSAEGK
φ VP882	AAS38505	RepA	1239	nd	nd	IIGSLTDA	ADNDAWKPHVGN
φ N15	AAC48876	RepA(gp37)	1228	nd	nd	LIGDLQGA	LDNDRKSSAEGK

The presence of the N-terminal zinc finger motif is shown together with the spacing of the cysteine residues probably involved in zinc ion binding (indicated by gray shading). The signature motifs 4 + 5 correspond to the conserved functional important motifs of the TOPRIM domain [CDD 8119] defined by Aravind *et al.* [9]. Highly conserved residues are indicated by bold type. 'nd' not detected.

the Kornberg lab: in the first assay, primer synthesis by DnaG starts at a hairpin region in φG4 or φX174 DNA, which is left uncovered by SSB; the oligo-ribonucleotides synthesised in this assay are 26–29-nt in length [41,183,378]. The other assay – termed 'general priming reaction' – uses ssDNA, DnaB helicase and primase [6]; in this assay, the priming reaction depends on the ability of DnaB to load DnaG to the ssDNA. Primers synthesised by DnaG *in vitro* start with 5'GA but are considerably longer than those synthesised by the phage-encoded primases: 10–60 nt, with the majority around 11 nt [7,502]. The interaction of *E. coli* DnaG with the replicative helicase DnaB is weak *in vitro*, however, the orthologues from *Bacillus stearothermophilus* were shown to form a stable complex [34]. Proteolytic removal of <50 res. from the DnaG<sub>Eco</sub> C-terminus resulted in polypeptides that were still able to direct the synthesis of RNA primers although they were inactive in the helicase-dependent priming assay [265,455,438]. The biochemical analysis of the DnaG Q576A mutation – in the extreme C-terminus – was particularly rewar-

ding: the interaction of the mutant protein with DnaB is significantly weakened [265], the initiation of bi-directional replication is less efficient than with the wildtype protein [169], and Okazaki fragments are significantly longer [455]. Intriguingly, the interaction of DnaB with DnaG is apparently not only necessary to direct the primase to its template but, in addition, modulates the DnaG conformation such that a reduction of sequence specificity allows more efficient lagging-strand priming [285,493,32].

As already addressed in the previous chapter, the initiator, helicase and primase functions required for φP4 replication are performed by the α protein: φP4 replication *in vitro* proceeds (mostly) bidirectionally from *ori*, and elongation occurs probably as coupled leading- and lagging-strand synthesis, without a requirement for host DnaG primase and DnaB helicase [103]. Primers synthesised by the φP4 α protein on ssDNA are 2–5-nucleotides long and begin with 5'-AG [506]. Host DnaG can only poorly substitute for the α E214Q primase(null) mutation in φP4 replication, and – *vice versa* –

**Table C23** Split genes encoding (putative) primases of the *E. coli* DnaG-type.

motifs:				Zn finger		4	5	
consensus		N	C	x <b>C</b> xx <b>H</b> x	15	x <b>C</b> xx <b>C</b> x	x <b>EGYMD</b> xh	x <b>DGD</b> xx <b>G</b> xxx <b>A</b> x
replicon	accession	gene	res.	ident.	ident.			
<i>E. coli</i>	NP_417538	<i>dnaG</i>	581	32%	28%	<b>CCPFH</b> N 15 <b>HCFGC</b> G	<b>VEGYMD</b> VV	<b>YDGD</b> RAGRDAAW
φ D29	NP_046874	gp58	129	[1-129]	-	<b>LCPFH</b> G 15 <b>NCLAC</b> G	-	-
	NP_046873	gp57	152	-	[1-152]	-	<b>TEGE</b> IDAI	<b>ADGDE</b> PGMEFAK
φ L5	NP_039722	L5p55	130	84%	-	<b>LCPFH</b> G 15 <b>NCMAC</b> G	-	-
	NP_039721	L5p54	152	-	80%	-	<b>CEGEL</b> DTI	<b>ADGDD</b> DAGMEFAR
φ Bxz2	NP_817647	gp58	135	69%	-	<b>RCPFH</b> G 15 <b>RCFAC</b> P	-	-
	NP_817646	gp57	157	-	51%	-	<b>TEGE</b> IDAI	<b>ADGDE</b> PGLQFAN
φ Bxb1	NP_075318	gp51	171	62%	-	<b>LCWHE</b> E 15 <b>NCLAC</b> S	-	-
	NP_075317	gp50	157	-	39%	-	<b>AEGE</b> IDAL	<b>ADGDD</b> DAGMQFAE

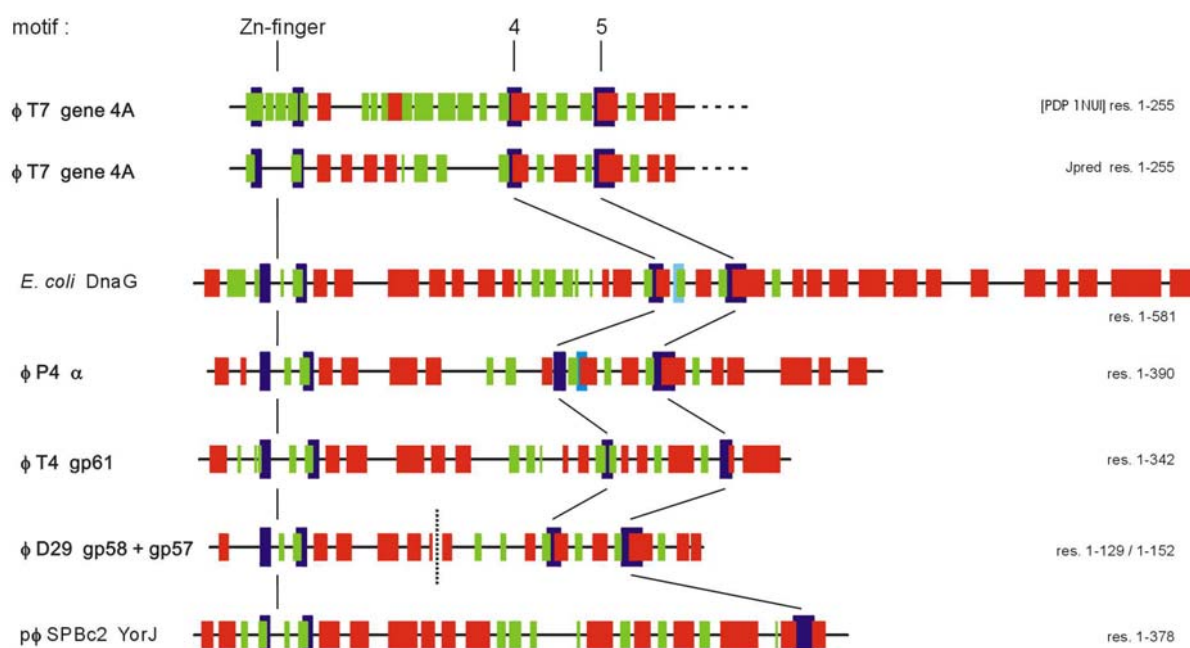
The presence of the N-terminal zinc finger motif of *E. coli* DnaG is shown together with the spacing of the cysteine residues probably involved in zinc ion binding (indicated by gray shading). The signature motifs 4 + 5 correspond to the conserved functional important motifs of the TOPRIM domain [CDD 8119] defined by Aravind *et al.* [9]. Highly conserved residues are indicated by bold type.

α cannot suppress an *E. coli dnaG(ts)* strain [431]. The isolated N-terminus (253 res.) of α was shown to retain its primase function; note that in this truncated protein the TOPRIM motif 5 is located 5 residues downstream from the amber stop [505]. The primase(null) phenotype of the E214Q mutation suggests that the EGYATA motif downstream at short distance from motif 4 is functionally important (see Fig. C15). This motif is in most cases only poorly conserved in other phage-encoded DnaG-type primases, but several plasmid-encoded primases were detected that show conservation of this motif [431,335]. However, the functional role of this sequence motif is elusive.

BLAST searches revealed several phage-encoded proteins with significant similarity to DnaG (Table C22). In most of these proteins, the Zn-finger motif and the TOPRIM signature motifs 4 + 5 were readily detected, but for the sizes of the proteins there are clearly three different types: i. the proteins with a size of ~800 res. including φP4 α and two of its homologues, which also contain the φP4 α-type C-terminal helicase and DNA-binding domains (see previous chapter); ii. the proteins with a size of ~300 res. which probably represent the minimal version of a monofunctional primase; and iii. the large φN15 RepA-type proteins containing the TOPRIM in their N-termini, a variant of the φP4 α EGYATA motif, but lacking a detectable Zn-finger motif. The primase activity of the φN15 RepA protein has not yet been demonstrated, but the striking analogy of roles for the replication of their cognate replicons between the RepA and φP4 α suggests this function [507,364].

A unique organisation exists for the (putative) primase genes of four related Mycobacteriophages: a gene encoding a small polypeptide of ~130 res. containing a putative Zn-finger structural motif is located directly upstream of a gene encoding a ~150 res. polypeptide containing the primase signature motifs 4 + 5 (Table C23). These two open reading frames are separated by 1 nt in the genomes of φD29, φL5, and φBxz2. Also in the φBxb1 genome, gp51 and gp 50 are separated by the same 'reading-frame shift' but the stop codon present in the φD29 gp58 gene has been changed in the φBxb1 gp51 gene. This results in the 42 res. longer gp51 protein, and – as a consequence – the reading frames encoding gp51 and gp50 overlap. φBxb1 gp50 contains the signature motifs but is – as compared to its homologues in the other phages – considerably less well conserved, probably due to compensation for the overlap. Taken together, it is unlikely that the split of the genes is simply a sequencing artefact in all 4 sequences. Biochemical analysis would have to show that both polypeptides together exhibit primase activity but the genetic organisation is suggestive (see also Fig. C17). All four phages of this group encode highly similar DNA polymerases of the Pol I-type (Section C3.5.1.) and helicases (Section C3.3.). It is thus very likely that these phages inherited their unique primase genes together with the other replication genes from a common ancestor. A common origin of these mycobacterial phages was also proposed by the Pittsburgh phage group from a whole-genome comparison [125,342].

The φT7 gene 4A and φP4 α primase domains, the φT4 gp61 primase and *E. coli* DnaG all contain the 'TO-



**Fig. C17** Secondary structure prediction for DnaG-type primases/primase domains.

The (putative) primase/primase domain sequences, and the size/position of the signature motifs are drawn to scale. For signature motif numbering see Tables C23 + C24. Secondary structure predictions were obtained from the Jpred server [88]. Colour code: red = predicted  $\alpha$ -helical region; green = predicted  $\beta$ -stranded region. Dark blue = signature motifs as listed in the Tables; light blue = EGYATA motif (see Text for details) [431]. The split of the two (putative) primase genes of  $\phi$ D29 is indicated by a dotted line.

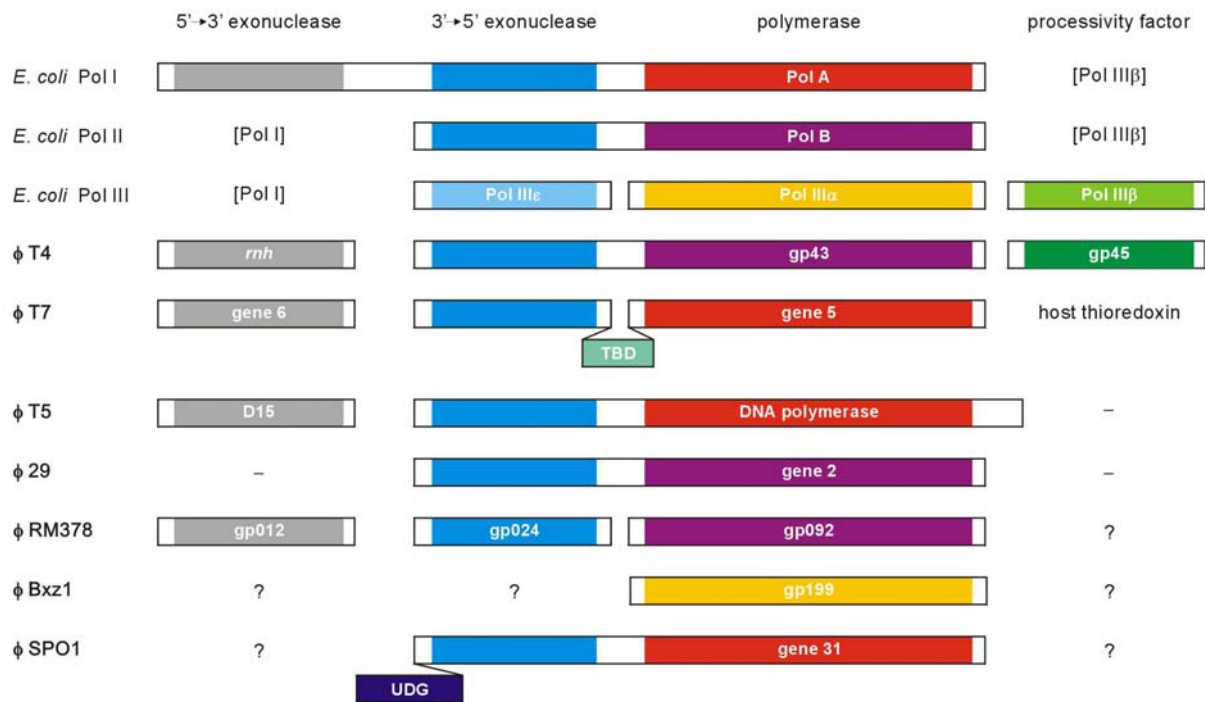
**Table C24** (Putative) primase domains of *Lactobacillus* sp. and *Burkholderia* sp. helicases.

motifs:					1	2	3	4
consensus					GhGFhh	hh.hGhDhD	xSGxGhHhhhRGx	hEhY. .xxGRFFxxTG
replicon	accession	gene	res.	ident.				
$\phi$ A2	NP_680517	orf35	770	[1-300]	GLGFFF	GY.VGIDVD	LSGEGIHIIIRGK	VEMY. .ESGRFFAMTG
$\phi$ 2389 / $\phi$ PSA	NP_511033	pri	757	50%	GLGFYF	PY.FGVDID	VSGTGIHIIAKGN	IEMY. .PDGRFFVMTG
$\phi$ adh	NP_050131	orf771	771	50%	GLAFYF	GY.VGLDVD	QSGKGIHAVFKGK	YEMY. .EAGRFFALTG
<i>L. innocua</i> p $\phi$	NP_471917	lin2587	757	49%	GLGFYF	PY.FGVDID	VSGTGIHIIAKGS	IEMY. .PDGRFFVMTG
$\phi$ 105	NP_690795	ORF11	806	37%	GIGFMF	PF.IGIDID	PSGEGVHIIAKGK	LEVY. .RHGRYFTFTG
$\phi$ BcepB1A	AAT37747	-	919	25%	GDGFTY	GFWLTADGG	SSGQGLHLFGRGV	LEFYTDKRGIAFGLSG

Motifs 1 - 4 were detected by BLAST alignment of the sequences as the regions with most significant conservation.

PRIM domain' (signature motifs 4 + 5; [9]) in addition to the N-terminal Zn-finger motif. However, none of the proteins shows any significant sequence similarity to the others, and their sizes differ considerably. The known crystal structures of the  $\phi$ T7 gene 4A primase domain and of a larger part of *E. coli* DnaG offered the possibility to compare these structures with each other, and to secondary predictions [205,204] (see Fig. C17). The

Jpred algorithm gives a reliability value of  $\sim 75\%$  for the predictions shown in Fig. C17, which seems reasonable. However, the algorithm is apparently biased towards the prediction of  $\alpha$ -helices and less suited to predict  $\beta$ -stranded regions. The embedding of the Zn-finger motif in several  $\beta$ -strands is reflected by the prediction, as is the positioning of motifs 4 + 5 in  $\alpha$ -helical regions. The region connecting motifs 4 + 5 differs in both crys-



**Fig. C18** Architecture of selected phage-encoded DNA polymerases.

Proteins are shown as open boxes (not to size; gene/protein name within box) with colours indicating the functional domain(s). Functional domains belonging to the same (sub)families are shown with identical colour (see text for details). 'UDG' uracil-d-glycosylase-like domain. 'TBD' thioredoxin-binding domain. '[Pol I]' indicates that the corresponding function is carried out by DNA Pol I (PolA). 'Pol IIIβ' indicates that the corresponding function can be carried out by the DNA Pol III β-subunit. '-' indicates that the corresponding function is not required for the respective replicon. A question mark indicates that a corresponding cognate gene of factor is not known yet for the respective replicon.

tal structures, and also the prediction cannot offer a plausible solution. In addition, the considerable length variation of the segment connecting the Zn-finger and motifs 4 + 5 makes it difficult to make any conclusive statement about the relationship of the proteins based on this type of analysis. Tentatively, we speculate that the ϕT7 gene 4A and ϕP4 α primase domain, ϕT4 gp61, and the split (putative) primase proteins gp58 and gp57 of ϕD29 are more related to each other than to any of the other two primases.

We have discussed the various types of primase-encoding genes found in phages – single, split, fused to helicase genes – that are all related to the prototype DnaG primase of *E. coli*, although rather distantly in most cases. We could not detect any phage-encoded primases with significant similarity to plasmid-encoded primases, e.g. to RepB of RSF1010 or to ColE2 Rep primase. Nevertheless, there may exist yet unknown primase genes: we were unable to detect – except for ϕT5 – primase-encoding genes in any of the (pro)phages that encode helicases (see previous chapter) and DNA poly-

merases of the ϕT5- and the ϕ12-subfamilies of Pol I-type DNA polymerases (BRM Section 3.4.).

Also, it is presently not possible to decide whether the N-terminal extensions of some very large F4-type and ϕP4 α-type helicases found in several phages can function as primases (see previous chapter; Fig. C15). E.g. several phages of *Lactobacilli* (Gram(+)) and *Burkholderia* sp. (Gram(-)) encode large proteins with significant similarity to other ϕP4 α-type helicases in their C-termini together with ~300 res. long N-termini of unknown function. These N-termini contain several well-conserved motifs in β-stranded regions (Jpred prediction; not shown) lacking any similarity to known primase motifs (Table C24). Whether these genes encode a novel type of primase-helicase enzyme needs to be clarified by biochemical experiments.

The observation that largely unrelated proteins could be classified by genetical and biochemical analysis as primases strongly suggests that such proteins evolved several times independently. This view is supported by the finding that for none of the known primases even a

distant relationship to DNA dependent RNA polymerases could be detected so far. For bacteriophage research, this situation is uncomfortable since it makes the prediction of gene functions for newly sequenced phage genomes difficult.

### C3.5. DNA polymerases and accessory proteins

#### C3.5.1. DNA polymerases

DNA polymerases are the key enzymes for the propagation of DNA genomes. All known DNA polymerases synthesise polynucleotide in the 5'→3' direction, and complementary to a single-stranded DNA template. Replicative DNA polymerases are unable to initiate *de novo* DNA synthesis on a DNA template but require the existence of a primer containing a free hydroxyl group to start DNA elongation [228]. Since Kornberg's purification of *E. coli* Pol I in the 1950s of the last century DNA polymerases are among the most intensely studied molecules, genetically, biochemically, and structurally [30]. A highly instructive review by Baker and Bell [18] might help the reader to find a way through the vast amount of literature, but other reviews are also suggested as starting points [201,422,428,381,27,186,150].

DNA polymerases are composed of structurally and functionally defined domains: i. a domain responsible for template binding, primer binding, and 5'→3' polynucleotide synthesis, ii. a proofreading or 3'→5' exonuclease domain responsible for the removal of misincorporated nucleotides [388, 366], iii. a 'structure specific 5'→3' exonuclease' or RNase H domain for the removal of RNA primers, and iv. a processivity factor in replicative DNA polymerases. *E. coli* DNA Pol I combines three functions in a single polypeptide, but lacks processivity. In contrast, polymerase functions are contributed by separate subunits in the *E. coli* Pol III holoenzyme which lacks a cognate 5'→3' exonuclease. Pol III $\alpha$  (*dnaE*) is the polymerase subunit, Pol III $\epsilon$  (*dnaQ*) the 3'→5' exonuclease subunit, and Pol III $\beta$  (*dnaN*) is the processivity factor, the sliding clamp (Section C3.5.2.). The Pol III  $\gamma/\tau$ ,  $\delta$ , and  $\delta'$  subunits are responsible for clamp-loading, and the  $\tau$ ,  $\chi$ ,  $\theta$ , and  $\psi$  small subunits of DNA Pol III mediate communication among the subunits and with other components of the replisome, e.g. SSB, the DnaG primase and the replicative helicase DnaB. Phage-encoded DNA polymerases provide further examples of DNA polymerase functions organised in individual subunits or combined in integral proteins.  $\phi$ T4 DNA polymerase gp43 contains 3'→5' exonuclease and polymerase domains, and requires the  $\phi$ T4 gp45 sliding clamp for processivity. The DNA polymerase of  $\phi$ 29 contains 3'→5' exonuclease and polymerase do-

main, and possess high intrinsic processivity.  $\phi$ T7 gene 5 DNA polymerase contains 3'→5' exonuclease and polymerase domains, and recruits host thioredoxin as processivity factor. In several phage replicases separate proteins carry the 3'→5' exonuclease and polymerase functions, e.g.  $\phi$ Bxz1,  $\phi$ RM378. In all phage replicons mentioned above, the 'structure specific 5'→3' exonuclease or RNase H is encoded by a separate gene (Fig. C18).

Most of the approximately 60 presently known phage-encoded (putative) DNA polymerases were detected by recent genome sequencing projects. All these putative DNA polymerases can be easily grouped into one of the known DNA polymerase families (see below). However, despite the wealth of structural data available for various DNA polymerases, four major topics exist where even a thorough comparison still cannot help to predict reliably the biochemical properties of a candidate DNA polymerase detected by genomics, emphasising the undiminished importance of biochemistry:

- i. it is not predictable whether a candidate polymerase is a high-fidelity or an error-prone polymerase. E.g. *E. coli* Pol II (*polB*) and  $\phi$ T4 gp43 DNA polymerase are members of the same polymerase family but differ with respect to error rate during synthesis.
- ii. the strand-displacement capacity is not predictable. E.g. the  $\phi$ 29 p2 and  $\phi$ T4 gp43 DNA polymerases are members of the same polymerase family, but  $\phi$ 29 p2 DNA polymerase is highly proficient in strand-displacement during processive DNA synthesis and can dispense with a helicase,  $\phi$ T4 gp43 polymerase in contrast requires its cognate processivity factor and helicase.
- iii. the processivity of a DNA polymerase is not predictable. E.g. the  $\phi$ T7 gene 5 and the  $\phi$ T5 DNA polymerases belong both to the Pol I-family of DNA polymerases,  $\phi$ T5 polymerase is inherently capable of high processive DNA synthesis, while  $\phi$ T7 gene 5 requires host thioredoxin as processivity factor. As a second example, *E. coli* PolIII $\alpha$  (*dnaE*) and  $\phi$ T4 gp43 polymerase belong to different polymerase families but both employ a sliding clamp as processivity factor [206].
- iv. it is not possible to predict which molecular mechanism governs the coupling of leading- and lagging-strand synthesis. This co-ordinated synthesis requires a functional asymmetry in the replisome. It is also not possible to predict from the sequence of a DNA polymerase how this asymmetry is put into practice for a given replisome [283].

We confine ourselves to a discussion of the different types of DNA polymerases encoded by bacteriophages, and their functional organisation in domains or separate

proteins. For practical reasons we chose to sort the phage DNA polymerases by their similarity to one of the five known DNA polymerases of *E. coli*: Pol I (*polA*), Pol II (*polB*), Pol III $\alpha$  (*dnaE*), Pol IV (*dinB*), and Pol V (*umuC*) [363].

***E. coli* Pol I (*polA*)-type DNA polymerases.** The crystal structure of the Klenow fragment of *E. coli* Pol I is known since 1991 [25]. For a thorough discussion of Pol I structure and functions we refer the reader to the review of Joyce and Steitz [201]. By the BLAST approach, altogether 40 phage-encoded (putative) DNA polymerases were detected (Table C25). These 40 DNA polymerases include 3 arbitrarily chosen prophage-encoded genes, which were however too numerous – particularly among *Xyella* strains – to include them all. None of these (mostly putative) phage DNA polymerases contains the N-terminal 5'→3' exonuclease domain characteristic for orthologues of *E. coli* Pol I. All protein sequences contain the 3'→5' exonuclease signature motifs defined by Salas and co-workers [28], and the polymerase signature motifs originally defined by Braithwaite and Ito [44]. This with almost perfect conservation of the functionally important residues – the few exceptions are discussed further below. Because these signature motifs were defined by combining the results of biochemical analyses of proteins mutated at catalytically important positions with known crystal structures the relevance of their detection in a candidate protein is considered superior to a calculated similarity (BLAST) value [201].

We found a reasonable correlation between the similarity of the proteins and the presence of several conserved individual residues within the 4 signature motifs of the polymerase domain. The degree of conservation of these residues within the signature motifs reflects the phylogenetic relationship of these DNA polymerases, an issue beyond the scope of this review. Here, these conserved residues allowed the loose definition of 4 subfamilies of the phage-encoded Pol I-type DNA polymerases: the  $\phi$ T7,  $\phi$ 12,  $\phi$ D29, and  $\phi$ T5 subfamilies. The considerable length variation among the 40 DNA polymerases (between 545 and 993 res.) is discussed in the subfamily paragraphs.

**Phage T7 gene 5-type subfamily.** The crystal structure of the  $\phi$ T7 gene 5 DNA polymerase is known [PDB 1T7P] [109]. It is the prototype phage-encoded DNA polymerase of the Pol I family, and contains all conserved signature motifs of the 3'→5' exonuclease and polymerase domains. Extensive biochemical analyses revealed that co-ordination of the highly processive leading- and lagging-strand synthesis by this polymerase requires an intricate interplay of the  $\phi$ T7 gene 4 helicase-primase, the gene 2.5 SSB, and 2 polymerases at a

specific DNA loop formed during synthesis of Okazaki fragments [248,247].

$\phi$ T7 gene 5 DNA polymerase requires host thioredoxin for processivity [441]. A domain (~76 res.) responsible for thioredoxin binding (TBD) could be identified within the region connecting the 3'→5' exonuclease and polymerase domains [175]. As inferred from sequence similarity, TBDs are also present in the DNA polymerases of phages  $\phi$ T3,  $\phi$ A1122,  $\phi$ YeO3-12, and  $\phi$ gh-1. The domain character of the closely related TBD of  $\phi$ T3 DNA polymerase (97% ident. res.) was demonstrated by its functional transfer to a corresponding position in an *E. coli* Pol I-type DNA polymerase lacking this region [94]. The DNA polymerases of phages  $\phi$ SP6,  $\phi$ K1-5, and  $\phi$ Felix01 also contain an extended region connecting the 3'→5' exonuclease and polymerase domains, but there is no primary sequence similarity to the  $\phi$ T7 gene 5 TBD, and its function therefore unclear. The DNA polymerases of Cyanophage  $\phi$ P60 may also use thioredoxin as processivity factor: the locus P60\_19 directly upstream of the polymerase gene encodes a putative thioredoxin. Roseophage  $\phi$ SIO1, and *Pseudomonas* sp. phage  $\phi$ PaP3 lack the extended interdomain region, and it is not known how these DNA polymerases maintain processivity.  $\phi$ PaP3 DNA polymerase seems to be encoded by two separate genes: while p39 is similar to the N-terminal ~190 res. of  $\phi$ T7 gene 5 DNA polymerase, the larger p32 protein is similar to the 'core' and contains the signature motifs.

The  $\phi$ T7 gene 5 DNA polymerase contains a small loop (EGDK motif) closely upstream of the signature motif 2a (res. 401-404). Deletion of this loop results in a mutant polymerase with a defect in primer utilisation and in its ability to contact the primase subunit of  $\phi$ T7 gene 4 helicase-primase [78]. The EGDK motif (EGDN in  $\phi$ gh-1 p18) is only conserved in those DNA polymerases of this subfamily that – like  $\phi$ T7 gene 5 – contain a TBD.

**Phage T5 DNA polymerase-type subfamily.** The DNA polymerase of  $\phi$ T5 has been characterised biochemically as a protein with high inherent processivity and a strong strand-displacement activity [245,73].  $\phi$ T5 DNA polymerase possesses a longer than usual C-terminal extension (~50 res.) also present in the  $\phi$ SPO1 gene 31 DNA polymerase but lacking in all other Pol I-type DNA polymerases. This region of the protein is suspected to be responsible for the high intrinsic processivity of the protein but direct evidence has not yet been presented [245]. This group of seven phage DNA polymerases is rather heterogeneous, and individual sequences have only a very low similarity with  $\phi$ T5 DNA polymerase. They were grouped mostly because the second conserved aspartic acid residue in polymerase motif C is replaced by either serine or alanine; conserved residues



**Table C25** Conserved motifs in the 3'→5' exonuclease and polymerase domains of *E. coli* Pol I-type polymerases.

domains				3'→5' exo domain		polymerase domain					
motifs:				ExoI	ExoII	ExoIII	2b	A 3	B 4	C 5	ident. res.
consensus				DxE	Nx <sub>2-3</sub> FD Y	YxxxD	TGRLSS	DhSxIELR	RxxAKxhNFghhYG	Rhhh-VHDEhkh K	
replicon	accession	gene	res.								
<i>E. coli</i>	P00582	<i>polA</i>	928	DTE	NL.KYD	YAAED	TGRLSS	DYSQIELR	RRSAKAINFGLIYG	RMIMQVHDELVVF	
φ T7	P00581	gene 5	704	DIE	NGHKYD	YNVQD	TGRATH	DASGLELR	RDNAKTTFIYGFILYG	AYMAVHDEIQV	[331-704]
φ A1122	NP_848283	p21	704	DIE	NGHKYD	YNVQD	TGRATH	DASGLELR	RDNAKTTFIYGFILYG	AYMAVHDEIQV	98%
φ T3	NP_523320	gene 5	704	DIE	NGHKYD	YNVQD	TGRATH	DASGLELR	RDNAKTTFIYGFILYG	AYMAVHDEIQV	97%
φ YeO3-12	NP_052093	gene 5	704	DIE	NGHKYD	YNVQD	TGRATH	DASGLELR	RDNAKTTFIYGFILYG	AQMGVHDEVQI	89%
φ gh-1	NP_813764	p18	709	DIE	NGIKYD	YCEQD	TGRATH	DASGLELR	RGIAKTTFIYAFILYG	CYMAVHDELQI	56%
φ P60	NP_570330	P60_20	587	DIE	NIVNYD	YCARD	TGRQAH	DASGLELR	RKSGKGVTYCLIYG	YPLAFVHDEQQL	37%
φ SIO1	NP_064753	p16	580	DIE	LKDVAE	GFMLD	TGAVTH	DAAAGIQLR	RPTAKTTFIYAFILYG	KQCAVHDEWQT	32%
φ PaP3	NP_775218	p39	175								30%
	NP_775225	p32	545	DLE	QWNIKD	TVQGD	TGRMRH	DGAGLELR	RDDAKTTFIYAFIYG	WKVLDIHDEGQW	25%
φ SP6	NP_853574	gp14	849	DAE	NFLGYD	SKLTD	TFRMRH	DGAGLELR	RDMAKTTFIYAFIYG	CGVANVHDEIQM	34%
φ K1-5	AAR90053	gp13	846	DAE	NFLGYD	SKLTD	TFRMRH	DGAGLELR	RDMAKTTFIYAFIYG	CGIANVHDEIQM	33%
φ Felix01	NP_944958	p179	906	DTE	NGLGYD	GWKAD	TGRMTQ	DQNSAQLV	RKKAKNGIYALLFG	RWVCSYHDEVSL	23%
φ 12	NP_803318	p12	650	DIE	NA.NFE	YCIRD	TGRWAG	DFSAIEAR	RQKGVSEELALCYQ	KIVGEVHDEVIV	[251-650]
<i>S. pyogenes</i> pφ	AAR83201	-	651	DIE	NA.SFE	YNRD	TGRWAG	DFSAIEAR	RQKGIKIELALCYG	KIVGEVHDEVII	48%
<i>S. pyogenes</i> SSI-1	NP_801684	sps0422	640	DIE	NA.QFE	YCIQD	TGRWAG	DFSAIEAR	RQKGIKIELALCYQ	GVVFEVHDEAII	55%
φ SPO2	P06225	gene L	648	DIE	NA.NFE	YCIQD	TGRWAG	DFSAIEAR	RQKGVKVELALCYQ	KTVMEVHDEAVL	56%
<i>X. fastidiosus</i> 9a5c	NP_299803	XF2525	726	DLE	NS.HFD	YAQRD	TGRWAG	DLSNIEGR	RQIGKVVQELALCYG	SIVLTVHDEIIT	32%
φ BPP-1	NP_958711	Bbp42	688	DLE	NS.HFD	YAGLD	TGRWAG	DLSNIEGR	RQIGKVMELMLCYE	QIVLSVHDEIIT	34%
φ APSE-1	NP_051006	P45	993	DLE	NGGMFD	YAGSD	TGRWAG	DLSNIEGR	RQIGKVMELGLCYG	EIVLTVHDEIIS	31%
φ BcepNazgul	NP_919002	gp49	693	DYE	NA.QFE	YNQD	TQRWAG	DLSSIEIV	RTMSKPAVLGAGYR	NIVLEVHDEIVT	31%
φ Vp16C	AAQ96570	orf39C	686	-	NI.TFE	YRYCD	TGRWISA	DFSAIEAV	KSLGKVAELASCYG	PVVMETHDELIA	33%
φ Vp16T	AAQ96505	orf38T	714	DVE	NI.TFE	YAYCD	TGRWISA	DFSAIEAV	KSLGKVAELASCYG	GKCYAVHDDVLY	31%
φ Bcep1	NP_944374	gp66	655	DVE	NA.PFD	YALID	TARGAG	DLSGIEAR	RQVGVKVDLALCFG	EIVFEVYDEALL	27%
φ Bcep781	NP_705685	gp57	655	DVE	NA.PFD	YALID	TARGAG	DVPKLEAF	RQVGVKAADLALCFG	EIVFEVYDEALL	26%
φ Bcep43	NP_958166	gp61	656	DVE	NA.PFD	YALID	TARGAG	DVPKLEAF	RQVGVKAADLALCFG	EIVFEVYDEALL	25%
φ D29	NP_046860	gp44	607	DTE	NA.SYD	YAGMD	TSRMSI	DYQAQELR	RKYAKTVNFGRVYG	VTLRLPIHDEIVA	[213-607]
φ L5	NP_039708	L5p41	595	DTE	NA.SFD	YAGMD	TSRMSI	DYQAQELR	RKYVGTANFQKVYG	VTLRLPIHDEIVA	71%
φ Bxz2	NP_817633	gp44	604	DTE	NA.SYD	YSGMD	TARMSI	DYQAQELR	RKYAKTVNFGRVYG	VTLRLPIHDEILA	71%
φ Bxb1	NP_075308	gp41	608	DTE	NA.AFD	YAGMD	TARMST	DYQAQELR	RKYVGMVNFAYVFG	MIRLVHDEVLA	52%
φ BT1	NP_813726	gp11	624	DTE	NA.MFD	YAGLD	TGRMSI	DFQAIEMR	RKYVFKGAGFGKVYG	VYMLPIHDEIVF	42%
φ C31	NP_047956	gp11	617	DTE	NA.PFD	YAGLD	TGRMSI	DFQAIEMR	RKYVFKGAGFGKVYG	VYMLPIHDEIVF	42%
φ Rosebush	NP_817817	gp56	609	DIE	NA.SFD	GAMAD	TGRMSY	DWAQIEPV	RPMKAVILLATMYG	HVQLAMHDELVV	32%
φ PG1	NP_943839	gp61	619	DIE	NA.PFD	YEGMD	TGRMSY	DWSQIEPV	RPTSKVLLSSMYG	ELYLAMHDEVVV	28%
pφ BC6A51	NP_831637	bc1864	802	DIE	NL.SFE	YSAED	TGRLNS	DFSGFELR	RTDAKAGNFGISYG	DMIAQIHDEIIF	27%
φ T5	P19822	-	829	DSE	NL.KFD	YAAKD	SGRLSS	DLTTAEVY	RQAAKAITFGILYG	KIVMLVHDSVVA	[301-829]
φ SPO1	AAA03732	gene 31	924	DLE	NG.KFD	YTLDD	TGRLSS	DYSQLELR	RTASKKIQFGIVYQ	RICITVHDSIVL	27%
φ VP5 / φ VP2	AAR92073	-	633	DTE	NA.SFD	YACPD	TGPALQ	DMASFEVR	QANAKQMNLMGIMFN	SMILNTHDSYSM	23%
φ VpV262	NP_640280	p19	661	DLE	NA.KFD	LKLRD	THRLSS	DGAQIEFR	RTDAKPDTFKPLYG	FLVNTVHDSVIS	-
φ Xp10	NP_858986	p39	794	DLE	NI.KFD	GTGED	TGRLSS	DFTALEIY	RSKAKGFSFQRAYG	LIVNTVHDAAYI	31%
φ KMV	NP_877458	gp19	807	DLE	NLDGVD	GPCGD	TTRLSS	DYTALEVV	RKNIKPKAFSAQYG	FLINNVDVAVYT	26%
φ 1 ( <i>V. harveyi</i> )	AAL85287	-	786	DLE	NF.KLD	YAGGD	TGRLAA	DFSQGELR	RNYGKPNFGLLYG	KFMAMIHDAVLS	23%

Initial BLAST searches with *E. coli* Pol I (*polA*) as query were refined in subsequent rounds: i. with parts of the Pol I sequence as query, ii. with *B. subtilis* Pol I sequence as query, iii. with known phage DNA polymerases as query (e.g. φT7, φT5), iv. with the least similar matches of preceding rounds as query, and v. by eye-scanning completely sequenced phage genomes for conserved motifs. The motifs ExoI, ExoII, and ExoIII were taken from Bernad *et al.* [28]. Polymerase domain motifs A - C for the Pol I family were taken from Joyce + Steitz [201] and [337], which correspond to motifs 3 - 5 of Delarue *et al.* [99]. The motifs 1 + 2 of Delarue *et al.* were omitted, and from the variation of motif 2 by Blanco *et al.* [38] only motif 2b included because motifs 1 + 2a could not be unambiguously identified in many phage DNA polymerase sequences. Bracketing indicates the proteins compared for similarity (% ident. res.) with φT7 gene 5, φ12 p12, φD29 gp44, and φT5 DNA polymerase, respectively. Conserved individual residues within the motifs are boxed. The micro-duplication in motif C of φVp16T orφ38T is indicated by underlining. Bold letters indicate the functionally important residues (see text).

within the signature motifs were not detected (see Table C25). The spacing of the signature motifs in the DNA polymerases of phages  $\phi$ Xp10 and  $\phi$ KMV suggests that these proteins contain an additional domain, in analogy to the TBD of  $\phi$ T7 gene 5 DNA polymerase. For these two proteins, however, the function of this interdomain-region is not known.

A rather unusual composition of functional domains is found in the  $\phi$ SPO1 gene 31 DNA polymerase: a stretch of ~200 residues N-terminal to the 3'→5' exonuclease and polymerase domains shows significant similarity with bacterial uracil-DNA glycosylases, e.g. 43% ident. res. over almost the entire length with *Thermus thermophilus* Hb8 UDG (205 res.).  $\phi$ SPO1 – together with  $\phi$ e,  $\phi$ HI,  $\phi$ 2C,  $\phi$ SP8, and  $\phi$ SP82 – belongs to a group of large *B. subtilis* phages (genome size ~150 kb), which contain 5-(hydroxymethyl)-2'-deoxyuridine instead of thymine in their linear dsDNA genomes [177]. It is not known, whether the UDG-like domain of  $\phi$ SPO1 has UDG activity, and whether such a function has any relevance for the replication of hmUra containing DNA. Alternatively, an UDG-like polypeptide could be simply employed as processivity factor for the polymerase by making use of its DNA-binding property. BLAST searches with  $\phi$ SPO1 gene 31 as query gave various matches in the NCBI database with proteins showing ~35% ident. res. confined to the N-terminal UDG-like domain (res. 1- ~200). Matching proteins were <250 res. in size and their entries assigned 'DNA polymerase, bacteriophage-type'; such falsely assigned proteins include *Archaeoglobus fulgidus* DSM 4304 NP\_071102, *Pyrococcus furiosus* DSM 3638 NP\_579114, *Halobacterium* sp. NRC-1 NP\_280750, *B. cereus* ATCC14579 NP\_831308, and *Clostridium tetani* E88 NP\_780781.

**Phage D29 gp44-type subfamily.** This group of nine (putative) DNA polymerases is characterised by a moderate to high similarity to Mycobacteriophage  $\phi$ D29 gp44 DNA polymerase. None of these proteins has yet been characterised biochemically. The proteins have a rather uniform size (~600 res.) and motif spacing, and several conserved residues within the signature motifs of the polymerase domain are specific for this subfamily. Only the DNA polymerase of  $\phi$ BC6A51 has a somewhat larger spacing between motifs ExoIII and 2a, but a BLAST search with this region only matched with the 3'→5' exonuclease-polymerase interdomain region of various bacterial Pol I homologues. It is therefore questionable whether this region has a special function.

**Phage 12 p12-type subfamily.** None of the 13 proteins of this group has been characterised biochemically. This group of (putative) DNA polymerases is very homogeneous with respect to size, sequence similarity (~

40% ident. res.) and spacing of the conserved signature motifs (including conserved residues within the signature motifs). This is somewhat surprising because they are encoded by a great variety of phages of Gram-negative as well as Gram-positive hosts. None of the proteins has an N- or C-terminal extension, and also the interdomain region separating the 3'→5' exonuclease from the polymerase domain has the same length found in *E. coli* Pol I. There are therefore no hints how these (putative) DNA polymerases gain processivity. The extremely wide spacing of polymerase signature motifs B and C in the  $\phi$ APSE-1 P45 DNA polymerase sequence (accession NP\_051006) is artefactual, an intron sequence has not been removed during 'translation' of the  $\phi$ APSE-1 genome sequence into accessible protein database entries.

The (putative) orf39C and orf38T DNA polymerases of *Vibrio* sp. phages  $\phi$ Vp16C and  $\phi$ Vp16T, respectively, deviate both from the canonical pattern.  $\phi$ Vp16C orf39C DNA polymerase lacks the ExoI motif due to a truncated N-terminus: homology (91% ident. res.) between orf39C and orf38T begins with res. 91 of the latter. A simple sequencing and reading-frame assignment error can be excluded because the DNA sequence upstream of the  $\phi$ Vp16C orf39C gene bears no similarity with the 5' region of the  $\phi$ Vp16T orf38T gene [402]. The essential acidic residues in the polymerase motif C of  $\phi$ Vp16T orf38T DNA polymerase cannot be unambiguously assigned because the region apparently contains a micro-duplication. Also the spacing of motifs B and C differs considerably from that observed in the orf39C and other proteins of this subfamily. The functionality of both DNA polymerases needs to be confirmed by biochemical analysis.

***E. coli* Pol II (*polB*)-type DNA polymerases.** DNA polymerase II of *E. coli* is involved in excision repair, translesion synthesis, and replication re-start, and encoded by the SOS-inducible *polB* gene (formerly *dinA*) [40,193,361]. The evolutionary origin of the *E. coli* *polB* gene is still subject to debate, but it does certainly not belong to the standard equipment of a bacterial genome, because orthologues are exclusively found in the genomes of several  $\beta$ - and  $\gamma$ -Proteobacteria. [217,122]. Structurally, Pol II is composed of an N-terminal 3'→5' exonuclease domain and a C-terminal polymerase domain. *E. coli* Pol II has a region of unknown function (~150 re.) preceding the 3'→5' exonuclease domain, and which is only found in the PolB orthologues. The C-terminal pentapeptide QLGLF was inferred to be involved in the interaction of Pol II with the  $\beta$ -clamp [90,262]. The three-dimensional structure of Pol II is not known, but all signature motifs characteristic for the Pol $\alpha$ -family are readily detected (Table C26) [488,201]. Two groups of phage-encoded DNA polymerases share the

**Table C26** Conserved motifs in the 3'→5' exonuclease and polymerase domains of *E. coli* Pol II-type polymerases.

domains				3'→5' exo domain			polymerase domain			
motifs:				ExoI	ExoII	ExoIII	A II	B III	C I	
consensus				DxE	Nx <sub>2-3</sub> FD	YxxxD	DhxSLYPS	KhxNShYG	hYGD <sup>T</sup> DShhh	ident. res.
replicon	accession	gene	res.							
<i>E. coli</i>	P21189	<i>polB</i>	783	DFE	NVVQFD	YNLKD	DYKSLYPS	KIIMNAFYG	IYGD <sup>T</sup> DSTFV	
φ 29 / φ PZA	NP_040719	p2	572	DFE	NL.KFD	YIKND	DVNSLYPA	KLMLNSLYG	IYCD <sup>T</sup> DSIHL	[1-572]
φ B103	Q37882	gp2	572	DFE	NL.KFD	YIKND	DVNSLYPS	KLMFDSLYG	IYCD <sup>T</sup> DSIHL	81%
φ M2	P19894	G	572	DFE	NL.KFD	YIKND	DVNSLYPS	KLMLNSLYG	IYCD <sup>T</sup> DSIHL	81%
φ GA-1	NP_073685	gene 2	578	DFE	NL.KFD	YLKHD	DVNSMYP A	KLMLNSLYG	IYAD <sup>T</sup> DSIHV	52%
φ Cp-1	Q37989	gene 5	568	DFE	NL.KFD	YIHVD	DINSMYP A	KIMLNSLYG	LYAD <sup>T</sup> DSLHL	33%
φ PRD1	P10479	gene I	553	DFE	NGGKFD	YLKGD	DVNSMYP H	KLILNSSYG	LYCD <sup>T</sup> DSIIC	24%
φ Bam35c	NP_943751	orf5	735	DTE	NL.DFD	YMEYD	DKNSLYPY	KLMQNALYG	AYCD <sup>T</sup> DSCAT	24%
φ 44AHJD / φ P68	NP_817305	orf10	755	DIE	NCNKYD	YIHND	DINSSYP Y	KVVLNGLYG	IYCD <sup>T</sup> DSL YM	21%
φ P1 <sup>1)</sup>	NP_064636	P82	694	DFE	NSKNFD	YFR.D	DINSAYPY	KIKLNSLYG	IYGD <sup>T</sup> DSIMF	<20%
φ C1	NP_852013	orf7	784	DIE	NGAKYD	YEMVD	DLNSSYPT	KHYSKLFYG	WYAD <sup>T</sup> DSL YM	<20%
φ T4	P04415	gp43	898	DIE	NIEGFD	YNIID	DLTSLYPS	KILINSLYG	AAGD <sup>T</sup> DSVYV	[1-898]
φ RB69	Q38087	gp43	903	DIE	NVESFD	YNIID	DLTSLYPS	KLILINSLYG	LYGD <sup>T</sup> DSIYV	61%
φ RB49	NP_891597	gp43	892	DIE	NTETFD	YNIGD	DLTSLYPS	KVLINSLYG	RYCD <sup>T</sup> DSVYV	55%
φ KVP40	NP_899330	gene 43	850	DIE	NIESFD	YQIQD	DLTSLYPS	KILINSLYG	IYGD <sup>T</sup> DSIYV	42%
φ Aeh1	NP_943895	gp43	919	DIE	NSEMF D	YLIRD	DLTSLYP H	KVLINSLYG	VYGD <sup>T</sup> DSIYL	38%
φ 25	AAP68686 AAP68688	gene 43A gene 43B	889	DIE	NSEGFD	YCVRD	DLTSLYPS	KILINSLYG	CYID <sup>T</sup> DSVYL	51%
φ 44RR2.8t	NP_932391 NP_932390	gp43A gp43B	889	DIE	NSEGFD	YCVRD	DLTSLYPS	KILINSLYG	CYID <sup>T</sup> DSVYL	50%
φ 65	AAR90916 AAR90917	gp43A gp43B	897	DIE	NSEGFD	YSIRD	DLTSLYP H	KVLINSLYG	IYCD <sup>T</sup> DSQYL	36%
φ RM378	NP_835611	p024	313	DIE	NV.ISD	YNAVD	-	-	-	-
φ RM378	NP_835679	p092	522	-	-	-	DFTSLYPS	KIMMNSMYG	IYSHTDSIFV	-

The motifs ExoI, ExoII, and ExoIII were taken from Bernad *et al.* [28]. Polymerase domain motifs A - C for the Pol $\alpha$  family were taken from Joyce + Steitz [201], which correspond to motifs I - III of Wong *et al.* [488]. Bracketing indicates the proteins compared for similarity (% ident. res.) with φT4 gp43, and φ29 p2, respectively. Bold letters indicate the functionally important residues (see text). 1) *Mycoplasma* sp. φ P1 must not be mistaken as *E. coli* φ P1.

**Table C27** Conserved motifs in the polymerase domain of *E. coli* Pol III $\alpha$ -type polymerases (DnaE); including the cognate 3'→5' exonuclease, DnaQ.

domain				3'→5' exo domain			polymerase		
motifs:				ExoI	ExoII	ExoIII $\epsilon$	C	A	
consensus				DxE	Nx <sub>2-3</sub> FD	HxxxxxD	hPDhDhDh	hhKhDhLG	ident. res.
replicon	accession	gene	res.						
<i>E. coli</i>	P10443	<i>dnaE</i>	1160	-	-	-	MPDFDVDF	LVKFDFLG	[1-1160]
<i>B. subtilis</i>	O34623	<i>dnaE</i>	1115	-	-	-	MPDIDIDF	LLKMDFLG	37%
<i>E. coli</i>	P03007	<i>dnaQ</i>	243	DTE	NA.AFD	HGALLD	-	-	-
<i>B. subtilis</i>	P13267	<i>polC</i>	1437	DVE	NA.SFD	HRAIYD	VPDIDLNF	LLKLDILG	23%
<i>B. subtilis</i> φ SPBc2	NP_046685	<i>yorL</i>	1305	-	-	-	LPDIDIDT	SVKFDLLT	22%
φ PIS136	AAL66178	DP1	1206	-	-	-	PPDIDIDF	LLKIDILG	37%
φ Bxz1	NP_818250	gp199	1131	-	-	-	FPDIDVDF	GVKDLLLA	26%
φ Barnyard	NP_818618	gp80	1111	-	-	-	YPDIDLDF	MLKVDFLG	26%

The motifs ExoI and ExoII were taken from Bernad *et al.* [28]; motif ExoIII $\epsilon$  from Moser *et al.* [304] and deRose *et al.* [100]. Polymerase domain motifs C and A were taken from Pritchard and McHenry [355]. Bracketing indicates the proteins compared for similarity (% ident. res.) with *E. coli* Pol III $\alpha$  (*dnaE*). Bold letters indicate the functionally important residues (see text).

conserved signature motifs with *E. coli* Pol II despite the lack of significant sequence similarity: the group of  $\phi 29$  p2-like DNA polymerases and the group of  $\phi T4$  gp43-like DNA polymerases (Table C26).

In marked contrast to *E. coli* Pol II and  $\phi T4$  gp43 DNA polymerase (see below), the p2 DNA polymerase of *B. subtilis* phage  $\phi 29$  is highly processive and exhibits a strong strand-displacement activity [372]. The latter property renders the replication of the 19.4 kb long dsDNA genome of  $\phi 29$  independent of a helicase.  $\phi 29$  DNA polymerase acts as monomer since the strand-displacement mechanism operating during  $\phi 29$  replication does not require a coupling of leading- and lagging-strand DNA synthesis. For the details of protein-primed replication and  $\phi 29$ -type DNA polymerases we refer the reader to the recent review by Meijer, Salas and co-workers, and primary literature cited therein [287].

The DNA polymerases of several Gram(+)-specific phages and the Gram(-)-specific phage  $\phi PRD1$  are similar in primary sequence, and of comparable size, as are the respective phage genomes [340]. It seems very likely, therefore, that also the biochemical properties of these rather compact DNA polymerases are very similar. This makes this group of enzymes excellent candidates for the elucidation of the molecular basis for inherent processivity and strand-displacement capability, not complicated by a coupling of leading- and lagging-strand DNA synthesis as in dimeric DNA polymerase holoenzymes.

The DNA polymerases of phages  $\phi Bam35C$ ,  $\phi 44AHJD/\phi P68$ , and  $\phi C1$  are considerably larger than  $\phi 29$  p2 DNA polymerase. In all three proteins, the 3'→5' exonuclease domain is separated from the polymerase domain by a ~200 res. long region of unknown function. In analogy to the structure of  $\phi T7$  gene 5 DNA polymerase this region might be required for binding of an (unknown) processivity factor. Alternatively it could be simply a flexible hinge region: BLAST searches with this region as query revealed limited similarity to putative hinge regions of several bi-modular transcription factors.

Pol II-type DNA polymerases involved in protein-primed replication are not confined to phage replicons: the linear mitochondrial plasmid pCIK1 (6.7 kb) of the phytopathogenic fungus *Claviceps purpurea* codes for protein with C-terminal similarity (31% ident. res.; ~200 res.) to orf7 DNA polymerase of  $\phi C1$ . Interestingly, pCIK1 replicates via protein priming like  $\phi 29$  [327].

Eight highly similar phage-encoded DNA polymerases make up the group of  $\phi T4$  gp43-type DNA polymerases (Table C26). The crystal structure of the  $\phi RB69$  gp43 protein has been solved [PDB 1CLQ], also in complex with the cognate clamp (see below) [407]. Thus this DNA polymerase is among the best understood structurally, and provides the prototype for the

*E. coli* Pol II-type DNA polymerases [203]. Within the gp43 structure, the conserved signature motifs were found at positions indicative of their importance for protein function. Benkovic and co-workers could demonstrate the interaction of the  $\phi T4$  gp43 C-terminus with the gp45 clamp [3]. The sequence of the extended C-terminal domain involved in the interaction with the cognate clamp is highly specific for this group of DNA polymerases. A coupling of leading- and lagging-strand synthesis by  $\phi T4$  gp43 DNA polymerase is achieved by a dimeric 'holoenzyme', and a direct interaction of the two polymerase subunits could already be shown to involve a region overlapping the polymerase domain [381]. The details of the  $\phi T4$  replisome were recently reviewed by Benkovic and co-workers [27].

To our knowledge, the (putative) p092 DNA polymerase of  $\phi RM378$  is the only *E. coli* Pol II-type DNA polymerase lacking an N-terminal 3'→5' exonuclease domain – a key feature of Pol III $\alpha$ -type DNA polymerases (see below). However, a 3'→5' exonuclease is encoded separately in the  $\phi RM378$  genome. The 3 polymerase signature motifs of p092 are located within the N-terminal half of the protein. It is not known whether the extended C-terminus (~200 res.) of p092 is necessary to achieve the topologically required orientation of the exonuclease and polymerase subunits within the DNA polymerase holoenzyme, or whether it is involved in contacting a yet unknown sliding clamp, as in  $\phi T4$  gp43-type DNA polymerases.

#### *E. coli* Pol III $\alpha$ (*dnaE*)-type DNA polymerases.

BLAST searches with *E. coli* Pol III $\alpha$  (*dnaE*) as query led unexpectedly to the detection of four phage-encoded homologues (Table C27). The three-dimensional structure of the  $\alpha$ -subunit of *E. coli* Pol III holoenzyme is not known but Pritchard and McHenry could identify functionally important residues by biochemical analysis of mutant proteins, and thus derive 2 signature motifs [355]. Despite their low overall similarity with Pol III $\alpha$ , the two signature motifs are well conserved and readily found at corresponding positions in all four phage-encoded proteins, which we consider *bona fide* DNA polymerases, therefore. Lacking an N-terminal 3'→5' exonuclease domain, all four DNA polymerases are structurally more related to DnaE than to PolC of *B. subtilis* (see Table C27). All four phages encoding Pol III $\alpha$ -like DNA polymerases infect Gram(+) hosts: *B. subtilis*, and three species of the Actinobacteria branch (high G+C Gram(+)-bacteria).  $\phi SPBc2$  is also found as curable prophage in the genome sequence of *B. subtilis* 168 [491]. The genome sequences of the mycobacterial phages  $\phi Barnyard$  and  $\phi Bxz1$  were determined just recently [342]. Therefore we expect that more phage-encoded DNA polymerases of the *E. coli* Pol III $\alpha$ -type will be found in the future, adding this type of DNA polymerase

**Table C28** Conserved motifs in the phage-encoded 5'→3' exonucleases.

motifs:				<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>
consensus				hhhh <b>D</b> G	hhhh <b>D</b> x	hRxExxxxxxx <b>YK</b> xxR	EADDhhG <b>x</b> h	hhT <b>x</b> D <b>K</b> D
replicon	accession	gene	res.					
<i>E. coli</i>	P00582	<i>polA</i>	928	LILVDG	AVVFDA	FRDE. .LFEHY <b>K</b> SHR	EADDVIGTL	I <b>S</b> TG <b>D</b> KD
<i>E. coli</i>	P38506	<i>exo</i>	281	LLIVDA	VAVFDD	WRHQ. .RLPDY <b>K</b> AGR	EADDLAATL	IV <b>S</b> T <b>D</b> KG
<i>A. aeolicus</i>	NP_214115	<i>pol</i>	289	LYILDG	VVVFDA	TKRE. KIYADY <b>K</b> QR	EADDVIAYL	IY <b>S</b> P <b>D</b> KD
φ T4	NP_049859	<i>mh</i>	305	ICLIDF	VLCIDN	WRRD. .FAYYY <b>K</b> KNR	EADDHIAVL	I <b>S</b> SS <b>D</b> GD
φ RB69	NP_861940	<i>mh</i>	290	I <b>A</b> LAD <b>F</b>	VLCMDN	WRRD. .FAYYY <b>K</b> KNR	EADDHIGVL	IV <b>A</b> SD <b>G</b> D
φ RB49	NP_891817	<i>mh</i>	315	VHLIDM	VLA <b>F</b> DS	RRDI. . .APYY <b>K</b> RNK	EADDIIAVL	IN <b>S</b> GD <b>G</b> D
φ KVP40	NP_899249	<i>mh</i>	310	TKSFDG	LLTVDL	LRSN. VLKN <b>K</b> IKYPR	EADDHIGVL	IT <b>S</b> SD <b>G</b> D
φ Aeh1	NP_944210	<i>mh</i>	306	FCIFDF	IIAVDN	WRRQ. .KYFY <b>K</b> KHR	EADDII <b>G</b> VL	IV <b>S</b> SD <b>G</b> D
φ 44RR2.8t	NP_932573	<i>mh</i>	307	VNIIDA	IIAFDD	WRRD. .LAWYY <b>K</b> KNR	EADDIIAIL	I <b>S</b> SS <b>D</b> SD
φ T7	P00638	gene 6	300	ILVMDG	LA <b>F</b> TDS	WRKE. LVDPNY <b>K</b> TNR	EGDDVMGVI	I <b>I</b> SC <b>D</b> KD
φ A1122	NP_848287	p25	300	ILVMDG	LA <b>F</b> TDS	WRKE. LVDPNY <b>K</b> ANR	EGDDVMGVI	I <b>I</b> SC <b>D</b> KD
φ T3	P20321	gene 6	302	ILVMDG	LA <b>F</b> TDS	WRKE. LVDPNY <b>K</b> ANR	EGDDVMGVI	I <b>I</b> SC <b>D</b> KD
φ YeO3-12	NP_052100	gene 6	303	VLVMDG	LA <b>F</b> TDD	WRKV. LVDETY <b>K</b> ENR	EGDDVMGII	LV <b>S</b> CD <b>K</b> D
φ gh-1	NP_813767	p21	314	ALDMDY	ILSGDD	WRKE. .VLETY <b>K</b> ANR	EGDDVCGIL	SV <b>S</b> CD <b>K</b> D
φ P60	NP_570332	P60_22	243	TL <b>L</b> IDA	LT <b>F</b> TDS	FRKD. .VEPTY <b>K</b> GNR	EADDVMGIL	L <b>I</b> SP <b>D</b> KD
φ PaP3	NP_775229	p28	294	IAGIDG	LTG <b>K</b> DN	FRLELATIRKY <b>K</b> GTR	EADDWL <b>G</b> VR	AC <b>S</b> R <b>D</b> KD
φ Felix01	NP_944975	p196	348	HVFIDS	GL <b>T</b> FDE	WERQ. .TWKE <b>A</b> KSEK	E <b>A</b> DSIVIAK	L <b>M</b> SID <b>K</b> D
φ Sp6	NP_853581	gp21	342	I <b>A</b> LIDG	GA <b>E</b> CDA	FRLRLA <b>F</b> TKPY <b>K</b> GTR	E <b>A</b> DDL <b>M</b> SIA	IV <b>S</b> AD <b>K</b> D
φ K1-5	AAR90062	gp 20	342	I <b>A</b> LIDA	AA <b>E</b> CDA	FRVRLA <b>F</b> TKPY <b>K</b> QR	E <b>A</b> DDL <b>M</b> SIA	IV <b>S</b> LD <b>K</b> D
<i>P. putida</i> φP	NP_744425	PP2276	199	[Q <b>N</b> LID <b>K</b> ]	LV <b>F</b> SD <b>K</b>	WRKT. .VLPT <b>Y</b> KHNR	EGDDVCGIL	V <b>A</b> SID <b>K</b> D
φ T5	P06229	D15	291	LMIVDG	IVLG <b>D</b> K	FRLE. .HLPEY <b>K</b> GNR	EADDMAAYI	L <b>I</b> ST <b>D</b> GD
φ Xp10	NP_858984	p37	309	PLHVDG	AAGADR	VHLS. .ASN <b>T</b> K <b>G</b> H <b>R</b>	EADDGIAYC	V <b>M</b> SAD <b>K</b> D
φ VpV262	NP_640291	p30	323	RSLVDA	VADVDI	WRNEDATI <b>Q</b> KY <b>K</b> GCR	EADDSIAAT	V <b>F</b> SG <b>D</b> KD
φ KMV	NP_877461	orf22	313	TLV <b>C</b> DA	TR <b>T</b> LD <b>T</b>	TRVHLTAAG <b>G</b> A <b>K</b> AYR	EADDG <b>M</b> MD	IR <b>S</b> DD <b>K</b> D
φ RM378	NP_835599	p012	318	VNLIDL	GMI <b>I</b> DD	LRKK. .LLP <b>Q</b> Y <b>K</b> EH <b>R</b>	EADDVIAHL	IV <b>S</b> T <b>D</b> KD

Motifs A - E correspond to the proposal of Gutman and Minton [156] except that motif F was omitted because it does not contain residues known to be involved in catalysis or DNA binding (see text). Bold letters indicate the functionally important residues (see text). Motif A of *P. putida* φP protein PP2276 is shown in parentheses because it is irregularly spaced to motif B.

ses to the well-known types of phage-encoded DNA polymerases of the Pol I- and Pol II-type.

None of the four phage-encoded Pol III $\alpha$ -like DNA polymerases has yet been analysed experimentally. It is therefore only possible to speculate about their points of conflict with or access to the replication machinery of their hosts. Some hints may be obtained from the present knowledge about the replicative DNA polymerase of *B. subtilis*, which also has not yet been thoroughly characterised. Since its early detection, the *B. subtilis* PolC protein was considered the functional equivalent of the  $\alpha$ -subunit (*dnaE*) of *E. coli* Pol III holoenzyme [24]. In contrast to *E. coli* Pol III $\alpha$ , the PolC<sub>Bsu</sub> polymerase contains an N-terminal 3'→5' exonuclease domain that belongs to the 'dnaQ-subfamily' of exonucleases

(see below) [21]. A homologue of the *E. coli* *dnaQ* gene is not found in the genome of *B. subtilis*. In addition to *polC*, a homologue of the *E. coli* *dnaE* gene is found in the *B. subtilis* genome, and Ehrlich and co-workers could show that both polymerases are essential for chromosome replication [101]. These authors also presented evidence suggesting that both polymerases form part of an asymmetric *B. subtilis* replisome, and that DnaE<sub>Bsu</sub> is responsible for lagging-strand synthesis. We hypothesise that the phage-encoded Pol III $\alpha$ -type DNA polymerases can replace the DnaE subunit in the asymmetric DNA polymerase holoenzyme of the host, thus linking the propagation of their cognate replicon to the host replication machinery.

**Table C29** Conserved motifs in the phage-encoded 3'→5' exonucleases.

motifs:				<b>ExoI</b>	<b>ExoII*</b>	<b>ExoIII</b>	<b>ExoIII<sub>ε</sub></b>
consensus				<b>DxE</b>	<b>Nx<sub>2-4</sub>F<sub>1</sub>D</b>	<b>YxxxD</b>	<b>HxxxxD</b>
replicon	accession	gene	res.		Y		
<i>E. coli</i>	P00582	<i>polA</i>	928	<b>DTE</b>	NL . . KYD	YAAED	
<i>E. coli</i>	P03007	<i>dnaQ</i>	243	<b>DTE</b>	NA . . AFD		HGALLD
<i>E. coli</i>	NP_416515	<i>sbcB</i>	475	<b>DYE</b>	NV . . RFD		HDAMAD
φ T4	NP_049629	<i>dexA</i>	227	<b>DFE</b>	NI . APSD		HDCAKD
φ RB69	NP_861704	<i>dexA</i>	225	<b>DFE</b>	NLAPSED		HDCAKD
φ RB49	NP_891579	<i>dexA</i>	222	<b>DFE</b>	[CEAVGVD]		HDAIHD
φ KVP40	NP_899271	<i>dexA</i>	230	<b>DYE</b>	LIPSAED		HDCAKD
φ Aeh1	NP_943887	<i>dexA</i>	226	<b>DME</b>	LKPSMND		HDCAKD
φ 44RR2.8t	NP_932376	<i>dexA</i>	221	<b>DWE</b>	NL . KPSPD		HDCAKA
φ 65	AAR90930	<i>dexA</i>	219	<b>DFE</b>	LLLKPSPD		HDCARD
φ 12	NP_803310	p04	310	<b>DFE</b>	NA . . LFD		HSAKFD
φ ST64B	NP_700403	sb30	189	<b>DLE</b>	NG . ANFD		HNALAD
φ PY54	NP_892096	p50	190	<b>DLE</b>	NG . ASFD		HNALAD
φ 186	NP_052286	orf81	194	<b>DTE</b>	NA . . DYD		HRALAD
φ PaP3	NP_775218	p39	175	<b>DLE</b>	NI . LDYD	YCLQD	
φ Xp10	NP_858982	p35	245	[DDI]	NG . KKFD	YNVVD	
φ KMV	NP_877463	orf24	348	<b>DIE</b>	NG . KRFD	YNIDD	
φ RM378	NP_835611	p024	313	<b>DIE</b>	NV . . ISD	YNAVD	
φ CJW1	NP_817562	gp115	261	<b>DIE</b>	NG . DRFD	YNIHD	
φ VWB	NP_958252	p10	263	<b>DLE</b>	NV . . PYD	FHELD	
φ Che8	NP_817375	gp37	273	<b>DVE</b>	NLLSQAD	YCAGD	

The ExoII motif Nx<sub>2-3</sub>(F/Y)D of Bernad *et al.* [28] was changed to ExoII\* motif Nx<sub>3-4</sub>(F/Y)D to allow the alignment of the φRB69 *dexA* sequence. The motif ExoIII<sub>ε</sub> was taken from Moser *et al.* [304] and deRose *et al.* [100]. Bold letters indicate the functionally important residues (see text). The ExoI motif is shown in parentheses for the φXp10 p35 protein because it deviates from the consensus, it is however regularly spaced to motif ExoII.

***E. coli* Pol IV (*dinB*)-type and *E. coli* Pol V (*umuC*)-type DNA polymerases.** These two polymerases belong to the 'Y class' of errorprone DNA polymerases involved in translesion synthesis in *E. coli* and possibly also in stationary-phase induced mutagenesis [331]. They have homologues in *B. subtilis* (*yqjH*, *yqjW*) and many other bacterial genomes from different phyla [147,439].

BLAST searches identified the YoIE protein of φSPBc2 as the only phage protein with significant similarity to *E. coli* DinB and UmuC. Because φSPB2c encodes the YorL DNA polymerase (see above), the involvement of YoIE in replication of this phage under 'normal' conditions seems unlikely. Interestingly, in contrast to their virtually complete absence from phage replicons, genes encoding proteins with significant similarity to *E. coli* DinB and UmuC polymerases are quite common among plasmid replicons from various incompatibility groups, e.g. IncL/M plasmid R471a, IncT

plasmid R394, IncN plasmid R46, IncJ plasmid R391 (see also [344]).

**5'→3' exonucleases.** A 5'→3' exonucleolytic or RNase H activity is necessary for the removal of RNA primers after replication of dsDNA, prior to strand sealing by DNA ligase. Because this type of enzymatic activity requires a 5'-end for recognition but cuts DNA near the junction of the single- and double-stranded DNA it has been termed 'structure specific' 5'→3' exonuclease [201]. The role of this specific exonuclease clearly differs from that of the 5'→3' exonuclease required to create single-stranded 3'-OH overhangs required for recombination although the proteins may be structurally related (Section C3.6.2.).

The 5'→3' exonucleolytic activity of *E. coli* Pol I is located within the N-terminal domain (res. 1-323), and fully retained in a truncated Pol I protein (*resA1*;

Q298→amber), showing that the isolated domain is structurally and functionally intact [200]. This N-terminal domain is found in all 124 known bacterial Pol I (*polA*) orthologues, but in none of the 40 phage-encoded DNA polymerases of the *E. coli* Pol I-type (see above). However, several phages that encode DNA polymerases also encode a protein with limited similarity to the *E. coli* Pol I N-terminus, and *E. coli* exonuclease IX. [46].

Six particularly well conserved sequence motifs A–F were found by Gutman and Minton by alignment of ten 5'→3' exonuclease sequences, including 4 phage-encoded proteins [156]. Since then, the crystal structures of the  $\phi$ T5 D15 protein (PDB 1XO1) and  $\phi$ T4 RNase H (PDB 1TFR) were solved [65,308]. The results of mutation analysis in combination with both crystal structures led to the identification of a conserved lysine residue involved in DNA binding (motif C), and 4 aspartate residues essential for Mg<sup>2+</sup>-binding, and thus catalysis (motifs A, B, D, E) [31,135]. The essential residues determined by Nossal and co-workers for  $\phi$ T4 RNase H are perfectly conserved in 21 phage-encoded genes, despite the lack of significant sequence similarity (>25% ident. res.) (Table C28) [31]. The region of similarity between the 21 proteins is in most cases confined to a stretch of ~50 residues encompassing motifs D and E. However, the spacing of the motifs is fairly uniform in all protein sequences. We infer from the latter and the perfect conservation of the important residues that the genes have been correctly assigned. Motif A is set in parentheses for the *P. putida* prophage locus PP2276 (Table C28) because the spacing to motif B is unusually short, and the protein might in fact be truncated at the N-terminus and non-functional. The phage-encoded exonucleases are found exclusively in three groups of phages: i. all phages encoding DNA polymerases of the  $\phi$ T4-subgroup of the *E. coli* Pol B-type encode a 5'→3' exonuclease/RNaseH (the complete DNA sequences of the *Aeromonas* sp. phages  $\phi$ 25 and  $\phi$  65 are not yet available); ii. within the  $\phi$ T7 subgroup of phages encoding *E. coli* Pol I-type DNA polymerases, a gene encoding a 5'→3' exonuclease/RNaseH with similarity to one of the proteins listed in Table C28 is only absent in the genome of phage  $\phi$ SIO1; iii. all phages of the heterogeneous  $\phi$ T5-subgroup encode a 5'→3' exonuclease (the  $\phi$ SPO1 genome has not yet been sequenced). Except for the functionally important residues, no other conserved group-specific residues could be detected.

13 bacterial genomes characterised by significant genome size reduction lack a *polA* gene (several *Mycoplasma* sp., *Ureaplasma* sp., *Wigglesworthia* sp., *Tropheryma* sp., *Buchnera* sp.). Nevertheless, all encode a protein similar to *E. coli* exoIX. Exo IX genes are, on the other hand, not present in all genomes encoding a *polA* orthologue. This observation emphasises the im-

portance of the 5'→3' exonuclease/ RNaseH enzyme for the replication of bacterial genomes. A unique situation is found in the genome of *A. aeolicus*: the *pol* gene (locus aq\_1628; falsely assigned 'DNA polymerase I 3'-5' exo domain') encodes a 5'→3' exonuclease/RNaseH enzyme (Table C28), the *polA* gene (locus aq\_1967) encodes a  $\phi$ T7-subfamily DNA polymerase, i.e. lacking the N-terminal 5'→3' exonuclease/RNaseH domain of *E. coli* Pol I. However, the *A. aeolicus* PolA protein shares ~30% ident. res. with proteobacterial Pol I proteins and has no detectable similarity with any DNA polymerase of the  $\phi$ T7-subfamily. The origin of the *A. aeolicus polA* gene remains enigmatic, therefore.

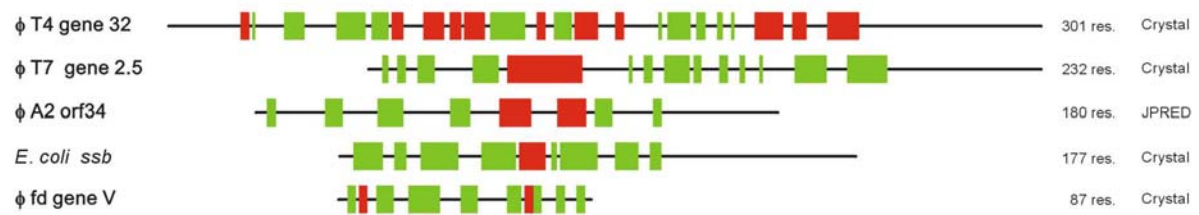
**3'→5' exonucleases.** The 3'→5' exonucleolytic activity required for proofreading, i.e. the removal of misincorporated bases in a growing DNA strand, is contributed by the  $\epsilon$ -subunit (*dnaQ*) in the *E. coli* DNA Pol III holoenzyme. Historically, the *dnaQ49* mutation was shown to result in mutator phenotype and led to the identification of the  $\epsilon$  subunit of *E. coli* DNA Pol III [393,348]. In contrast, the proofreading activity resides within the N-terminal domain of all known phage-encoded DNA polymerases of the *E. coli* Pol I- or Pol II-type, with the exceptions of the  $\phi$ RM378 DNA polymerase. It is therefore questionable whether the 3'→5' exonucleases encoded by phage replicons are important for their replication, or involved in recombination rather (Section C3.6.2.).

Mutational analysis of the *B. subtilis* PolC protein and the NMR analysis of the *E. coli* Pol III $\epsilon$  (*dnaQ*) protein suggest a significant structural similarity with the *E. coli* Pol I 3'→5' exonuclease domain [21,100]. Including the Pol III $\epsilon$  structure, the alignment of 3'→5' exonucleases led to the definition of a novel conserved motif ExoIII $\epsilon$ , specific for the 'DnaQ subfamily' of exonucleases, which deviates from the canonical ExoIII motif of *E. coli* Pol I- and Pol II-type DNA polymerases [28]. The recombination protein SbcB of *E. coli* also contains a N-terminal domain belonging to this 'DnaQ subfamily' (see Table C29).

The DexA protein of phage  $\phi$ T4 belongs to the 'DnaQ subfamily' of exonucleases, but it is not essential for replication under laboratory conditions [137]. DexA might instead be involved in the degradation of the host chromosome and  $\phi$ T4-specific recombination [185]. All  $\phi$ T4 group phages encode orthologues of DexA (Table C29), and the conservation of the entire set of replication proteins in this group makes it very likely that DexA is also non-essential in the other phages (the complete *Aeromonas* sp.  $\phi$ 25 genome sequence is not yet known).

The  $\phi$ 12 p04 protein may be the cognate 3'→5' exonuclease of the  $\phi$ 12 p12 DNA polymerase. No other phage encoding DNA polymerases of the  $\phi$ 12 subfamily





**Fig. C19** Secondary structures of SSB proteins.

The secondary structure elements of  $\phi$ T4 gene 32 [PDB 1GPC],  $\phi$ T7 gene 2.5 [PDB 1JE5],  $\phi$ fdgene V [PDB 1YHA], and *E. coli* *ssb* [PDB 1KAW] were taken from the known crystal structures, and obtained by Jpred prediction for  $\phi$ A2 orf34 [88]. Colour code: red = (predicted)  $\alpha$ -helical region; green = (predicted)  $\beta$ -stranded region; black line = unstructured/no prediction.

of Pol I-type polymerases encodes in addition a 3'→5' exonuclease, its occurrence in  $\phi$ 12 is unique, therefore. The  $\phi$ 12 p04 protein is special in a second aspect: it is rather large in size, the 3 signature motifs are found in the N-terminal half. BLAST searches revealed a significant similarity (~30% ident. res.) of its C-terminal half with C-termini of several NAD-dependent DNA ligases, while the N-terminal half gave matches exclusively with proteins of the 'DnaQ subfamily' of exonucleases. Any functional consequences of the apparent composite structure of this protein remain to be studied. Also the phages  $\phi$ ST64B,  $\phi$ PY54, and  $\phi$ 186 encode exonucleases of the 'DnaQ subfamily' (Table C29). These phage replicons are devoid of a cognate DNA polymerase and the participation of the sb30, orf181, and p50 proteins, respectively, in their replication unlikely.

The remaining 7 phage-encoded 3'→5' exonucleases belong to the subfamily of Pol I-type 3'→5' exonucleases (Table C29). We infer from the observation that the  $\phi$ RM378 p092 polymerase lacks the canonical exonuclease domain a functional role of  $\phi$ RM378 p024 exonuclease during replication of this phage. The role of the remaining proteins is unknown at present.

### C3.5.2. Sliding clamps, clamp loaders and DNA polymerase accessory proteins

Sliding clamps are processivity factors that tether the replicative DNA polymerase to the DNA template by encircling the DNA duplex with a ring-like structure. The sliding clamp of *E. coli* DNA polymerase III holoenzyme is formed by two DNA polymerase III $\beta$ -subunits. The  $\beta$ -clamp can attract other *E. coli* DNA polymerases in addition to DNA Pol III [263]. The archaeal and eukaryal PCNA clamps, and the  $\phi$ T4 clamp are homo-trimers. The proteins of the *E. coli*-type  $\beta$ -clamp, PCNA, and the  $\phi$ T4 gp45 clamp share no sequence similarity. Nevertheless, their three-dimensional ring structures derived from crystals can be readily superimposed, i.e. they are 'structural orthologues' (reviewed in [196]). Fueled by ATP hydrolysis, sliding clamps are

assembled around the DNA duplex by clamp loader protein complexes: the  $(\gamma/\tau)_3\delta\delta'$  complex in case of *E. coli* DNA polymerase III, the pentameric archaeal and eukaryal replication factor C complex (RFC) [324], and the pentameric  $\phi$ T4 gp44·gp62 (4:1) complex [453].

The crystal structures of the  $\phi$ T4 sliding clamp protein gp45 (PDB 1CZD) and of the highly related  $\phi$ RB69 gp45 protein (78% ident. res.; PDB 1B77, 1B8H) have been solved [407,298]. All phages encoding orthologues of the  $\phi$ T4 gp43 DNA polymerases also encode orthologues of the 228 res. long  $\phi$ T4 gp45 clamp:  $\phi$ RB69, *Aeromonas* sp.  $\phi$ 25 and  $\phi$ 65,  $\phi$ 44RR2.8t,  $\phi$ RB49,  $\phi$ Aeh1, and  $\phi$ KVP40. The gp45 proteins vary slightly in length (221–330 res.), and their similarities with the  $\phi$ T4 protein range from 27% ( $\phi$ KVP40 gp45) to 78% ( $\phi$ RB69 gp45). Five monomers of the  $\phi$ T4 gp44 protein and one monomer of gp62 form the clamp loader complex [349]. Orthologues of the  $\phi$ T4 gp62 and gp44 proteins are encoded by all phages of this group. Since also the helicases and SSBs of the  $\phi$ T4-like phages are highly similar it can be safely assumed that replisome formation in these phages closely resembles the  $\phi$ T4 pathway, including the coupling of leading- and lagging-strand synthesis [381,3,456]. The  $\phi$ T4 clamp – loaded preferentially to primer-template junctions by the pentameric clamp loader – also acts as transcriptional activator during late transcription of  $\phi$ T4 together with the gp33 co-activator and the late  $\sigma$ -factor gp55 [384,487]. It is not precisely known, however, how a direction of the  $\phi$ T4 clamp to its cognate replicon is achieved, avoiding its titration by the host chromosome. The RFC small subunits of various sequenced archaeal and eukaryal genomes show significant similarity with the  $\phi$ T4 gp44 clamp loader subunit, which might in part simply be due to the presence of a nucleotide-binding domain of the AAA-type [324,93]. Homologues of the  $\phi$ T4 DNA polymerase, or the clamp, or the clamp loader proteins are not found in the genomes of other phages.

Phage-encoded proteins similar to either the the eukaryal PCNA clamp protein(s) or the *E. coli*  $\beta$ -clamp

(*dnaN*) were not be detected by BLAST searches. Also searches for phage-encoded proteins similar to the other subunits of *E. coli* DNA polymerase III, i.e. the  $\gamma$ ( $\tau$ ),  $\delta$ , and  $\delta'$  clamp loader subunits, or the  $\chi$ ,  $\theta$ , and  $\psi$  subunits, were unsuccessful. It should be noted, however that genes encoding orthologues of the DNA Pol III  $\chi$ ,  $\theta$ , and  $\psi$  subunits are only found in *E. coli* and few closely related species. *Streptococcus pyogenes* MGAS8232 encodes a protein (locus SpyM3\_0690 as part of the p $\phi$  N315.1 genome, also present in other *S. pyogenes* genomes) showing similarity (22% ident. res.) in the N-terminal 180 res. with *E. coli* Pol III $\delta'$ . The assignment 'putative DNA polymerase III delta prime subunit - phage associated' is nevertheless misleading because the protein is a fulllength homologue of the DnaB<sub>Eco</sub>-type of helicases including all signature motifs (Section C3.3.; Table C28).

### C3.6. Single-strand DNA binding and recombination proteins

#### C3.6.1. Single-strand DNA binding proteins (SSBs)

Single-strand DNA binding proteins lack enzymatic activity but play important roles as accessory proteins in DNA replication, recombination, and repair (reviewed in [294, 260]). Prokaryotic SSBs bind cooperatively to single-stranded DNA, and non-specifically, albeit with a marked preference for poly(dT) *in vitro* [261,235]. SSBs prevent degradation of ssDNA by nucleases, keep unwound DNA in an open conformation for subsequent steps during the initiation of replication, and stimulate the DNA strand exchange activity of RecA or RecT (see following chapter). A unique feature is the requirement of *E. coli* SSB as activator of  $\phi$  N4 transcription [275]. SSBs of filamentous phages (e.g.  $\phi$ M13,  $\phi$ fd) or isometric phages (e.g.  $\phi$ X174,  $\phi$ G4) can inhibit replication, i.e. prevent the conversion of the ssDNA genomes into the double-stranded replicative form [45,130].

SSBs are essential for the replication of phages  $\phi$ T4 and  $\phi$ T7, the *E. coli* chromosomal replicon, and yeast, and it is assumed that SSBs are essential in all organisms (for references see [260,486]. Homologues of the *E. coli* *ssb* gene have been found in all presently known bacterial genomes (see below), in yeast mitochondria [461], and in the genomes of various phages and conjugative plasmids, e.g. F factor, R100, or Col1b. A specific role of the plasmid- and phage- encoded SSBs for their cognate replicons has been questioned because inactivation of the *ssb* gene did not prevent their propagation [146,184]. However, Lehnher and co-workers could recently demonstrate that the  $\phi$ P1 SSB provides a selective advantage for wildtype over  $\Delta$ *ssb-P1* phages in

infection and propagation during stationary phase growth of the *E. coli* host [26]. It seems thus possible that plasmid- and phage-encoded SSBs perform important functions under conditions that are not tested in the laboratory routine.

The crystal structures of the SSB proteins encoded by  $\phi$ T4 gene 32 [PDB 1GPC],  $\phi$ T7 gene 2.5 [PDB 1JE5],  $\phi$ fd gene V [PDB 1YHA], and *E. coli* *ssb* [PDB 1KAW] have been solved. The proteins show no sequence homology and are structurally only distantly related. Despite some similarity in their secondary structures, BLAST searches did not reveal primary sequence similarities between the  $\phi$ A2 orf34 SSB and *E. coli* SSB (Fig. C19).

The SSB proteins have distinct biochemical properties:  $\phi$ T4 SSB binds to DNA as monomer, SSBs of filamentous phages as dimers, and *E. coli* SSB as homotetramers. *E. coli* SSB monomers bind non-specifically to ~15 nt long oligonucleotides. For the SSB homotetramer, binding sites of  $35\pm 3$  nt and  $65\pm 5$  nt were determined *in vitro*, depending on salt and protein concentrations, as well as temperature and pH, suggesting that 2 or all 4 subunits of a single tetramer, respectively, bind to ssDNA (for references see [260]). Lohman and Ferrari assume that *E. coli* SSB binding to ssDNA *in vivo* occurs by the binding mode in which ~65 nt ssDNA are wrapped around one tetramer. This binding mode of 'limited cooperativity' results *in vitro* in the formation of beads consisting of double-tetramers separated by short protein- less DNA stretches that would allow other proteins to access the ssDNA [260]. The binding properties of SSB suggest that during the initiation of replication only very few tetramers actually bind to an unwound replication origin. A ~28 bp long stretch of the *B. subtilis* *oriC* was opened by the action of DnaA and HU *in vitro*, which was extended moderately by the addition of SSB [232].

The  $\phi$ T4,  $\phi$ T7,  $\phi$ A2, and *E. coli* SSBs possess highly acidic C-terminal tails of ~12 res. but not gene V SSB of  $\phi$ fd. It was shown that cooperative binding of *E. coli* SSB to ssDNA results in a conformational change of the protein exposing the C-terminus. SSB lacking up to 62 res. from the C-terminus is still able to form tetramers, and binds to ssDNA *in vitro* even tighter than wild-type SSB [482]. Comparable results were obtained also for  $\phi$ T4 gene 32 SSB [182,189], and Williams and co-workers proposed that these acidic C- termini are instrumental in signalling the SSB conformation to proteins involved in subsequent reactions [482]. For the acidic C-terminus of  $\phi$ T7 gene 2.5 SSB, an interaction with the gene 4A primase-helicase gene 5 polymerase could already be shown [213,214,212]. The *E. coli* *ssb-113* mutation results in a P→S exchange of the penultimate SSB residue, and – comparable to C-terminally truncated SSB – the SSB-113 mutant protein binds to ssDNA

**Table C30** Phage encoded SSB proteins.

<b>group 1: <math>\phi</math> T4 gene 32-type SSBs</b>	
genes	$\phi$ T4 gp32, $\phi$ T6 gp32, $\phi$ T2 gp32, $\phi$ RB69 gp32, $\phi$ 44RR2.8t gp32, $\phi$ RB49 p241, $\phi$ Aeh1 gp32, $\phi$ SV14 gp32, $\phi$ KVP40 gp32
lengths / similarity	295 - 322 res. / $\geq 48\%$ ident. res. with $\phi$ T4 gp32, full length
conserved C-tail	acidic tail; 4 - 6 D/E res. within the 12 C-terminal res.
<b>group 2: <math>\phi</math> T7 gene 2.5-type SSBs</b>	
genes	$\phi$ T7 gene 2.5, $\phi$ A1122 p12, $\phi$ T3 gene 2.5, $\phi$ YeO3-12 gene 2.5, $\phi$ gh-1 p10
lengths / similarity	232 - 233 / $\geq 50\%$ ident. res. with $\phi$ T7 gene 2.5, full length
conserved C-tail	12 res. acidic tail containing 4 - 8 E/D res.; C-terminal 3 res.: [G,E]DF
<b>group 3: <math>\phi</math> A2 orf34-type SSBs</b>	
genes	$\phi$ A2 orf34, $\phi$ adh orf175, $\phi$ PSA / $\phi$ 2389 orf49, $\phi$ DT1 orf34, $\phi$ O1205 orf11, $\phi$ sf19 orf151, $\phi$ 31 orf4, $\phi$ sf11 gp151, $\phi$ sf121 p37, $\phi$ mi7-9 gp18c, $\phi$ 105 orf10, $\phi$ ETA orf18
lengths / similarity	145 - 175 res. / $\geq 20\%$ ident. res. with $\phi$ A2 orf34, full length
conserved C-tail	12 res. acidic tail containing 3 - 5 D/E res.; C-terminal 9 res.: eI[S,e][D,e]x[D,L,v]PF;
<b>group 4: <math>\phi</math> M13 gene V-type SSBs</b>	
genes	$\phi$ M13 gene V, $\phi$ fd gene V, $\phi$ f1 gene V, $\phi$ lke gene V, $\phi$ Lf orf98, $\phi$ Cf1c p1
lengths	87-98 res. / $\geq 33\%$ ident. res. with $\phi$ M13 gene V, full length
conserved C-tail	no pronounced acidic tail, no consensus
<b>group 5: <math>\lambda</math> Ea10-type (putative) SSBs</b>	
genes	$\lambda$ <i>ssb</i> ( <i>ea10</i> ), $\phi$ VT2-Sa <i>ea10</i> , $\phi$ 933W <i>ea10</i> , $\phi$ 4795 <i>ea10</i>
lengths / similarity	122 res. / $\geq 99\%$ ident. res. with $\lambda$ <i>ea10</i> , full length
conserved C-tail	weakly acidic tail; C-terminal 12 res.: IFDSDDMTIKAA
<b>group 6: other SSBs</b>	
genes	$\phi$ lf1 geneV, $\phi$ Pf3 orf78; from filamentous phages of Proteobacteria $\phi$ N4 N4SSB, $\phi$ PRD1 gene XII; from Gram(-) hosts $\phi$ 29( $\phi$ PZA) gene 5A, $\phi$ B103 gene 5, $\phi$ Nf gene 5, $\phi$ GA-1 gene 5, $\phi$ P68 orf14, $\phi$ 44AHJD orf14; from Bacillales phages with genome sizes of ~20 kb $\phi$ BK5-T orf48; from <i>Lactococcus</i> sp., acidic C-tail: APMNISDEDLPP
<b>group 7: <i>E. coli</i>-type SSBs</b>	
Gram(-) specific	
genes	$\phi$ VP2 <i>ssb</i> , $\phi$ VP5 <i>ssb</i> , $\phi$ Bcep22 gp16, $\phi$ Sf6 gp27, $\phi$ P1 <i>ssb-P1</i> , $\phi$ T1 27
lengths / similarity	157 - 169 res. / $\geq 48\%$ ident. res. with <i>E. coli</i> <i>ssb</i> , full length
conserved C-tail	12 res. acidic tail containing 3 - 5 D/E res.; C-terminal 9 res.: m[D,e]fD[D,S][D,S]IPF;
Gram(+) specific	
genes	$\phi$ TP901-1 p12, $\phi$ LL-H orf299, $\phi$ SPP1 G36P, $\phi$ CJW1 gp102, $\phi$ SLT orf147, $\phi$ NIH1.1 gene 13, $\phi$ ul136 orf141, $\phi$ Tuc2009 orf15, $\phi$ MM1 p13, $\phi$ 7201 orf9, $\phi$ bIL285 orf15, $\phi$ PV83 orf18, $\phi$ bIL286 orf15, $\phi$ PVL orf45, $\phi$ VO1 orf1, $\phi$ 3626 p39, $\phi$ A118 ORF60
lengths / similarity	139 -181 res. / $\geq 28\%$ ident. res. with <i>E. coli</i> <i>ssb</i> , full length
conserved C-tail	12 res. acidic tail containing 3 - 5 D/E res.; C-terminal 8 res.: [d,e]ixd[D,q][L,v]P[F,v]

slightly tighter than wild-type SSB *in vitro* [72]. A possible interaction of *E. coli* SSB with the primosomal protein PriB was proposed based on biochemical experiments, but has not been substantiated [264]. More recently, a direct interaction of the acidic tail (15 res.) of *E. coli* SSB with the  $\psi$  subunit of DNA Pol III holoen-

zyme was shown by O'Donnell and coworkers to be necessary for efficient clamp loading [207]. They could also show that the SSB-113 mutant protein is defective in clamp loading. These findings support models for the initiation or the resumption of replication that consider SSB as an essential protein at the steps following helica-

se loading during replisome formation: SSB might be instrumental in co-ordinating the coupling of leading- and lagging-strand DNA synthesis [207,145]. Orthologues of the *E. coli* *holD* gene – encoding the  $\psi$  subunit of DNA Pol III holoenzyme – are only found in the genomes of the Enterobacteriales and a few other  $\gamma$ -subgroup Proteobacteria. However, the perfect conservation of the penultimate proline residue within a reasonably conserved acidic tail in all *E. coli*-type SSBs from Proteobacteria and Firmicutes, and in all  $\phi$ A2 orf34-type SSBs (see below) underscore the importance of the C-terminus of these SSBs. Database searches produced >50 matches for phage- encoded SSBs (prophages not included), which we sorted into 6 groups of similar proteins, and a seventh group containing sequences with no or few apparent homologues (Table C30). The largest group, group 7, containing the *E. coli*-type SSBs, is discussed last.

**Group 1: phage T4 gene 32 (*gp32*) SSB.** The 9 phages encoding homologues of the  $\phi$ T4 gene 32 single-strand DNA binding protein belong to the group of the T-even phages of Enterobacteriales, characterised by their intricate replication mode (BRM Section 2.5.). The high degree of similarity suggests that the biochemical properties of all proteins are very similar to those of  $\phi$ T4 *gp32* SSB. All proteins contain a highly acidic C-terminal tail, which however lacks a detectable consensus sequence. In addition to the full-length homologues, the data bases list numerous homologues of  $\phi$ T4 gene 32 for which only fragments are known from phages  $\phi$ FS- $\alpha$ ,  $\phi$ M1,  $\phi$ PST,  $\phi$ RB8,  $\phi$ RB10,  $\phi$ RB6,  $\phi$ RB15,  $\phi$ RB18,  $\phi$ SV76,  $\phi$ RB27,  $\phi$ RB32,  $\phi$ RB69,  $\phi$ RB3,  $\phi$ RB9, and  $\phi$ RB70.

**Group 2: phage T7 gene 2.5 SSB.** Like the SSBs of group 1, the 4 homologues of the  $\phi$ T7 gene 2.5 protein contain highly acidic C-terminal tails that lack a detectable consensus sequence. For the C-terminus this protein, specific interactions with gene 4 primase/helicase and gene 5 DNA polymerase was shown [223].

**Group 3: phage A2 orf34 SSB.** The 12 homologues of the  $\phi$ A2 orf34 SSB are exclusively found in the genomes of Bacilli phages (prophages not included). Although the similarity of the proteins among each other is rather low they all contain a well conserved acidic C-terminus. In addition to the overall structure (see above), the C-terminal tails of these SSBs are highly similar to those of *E. coli*-type SSBs of Gram(+)-specific phages (group 7). It is therefore reasonable to assume that the SSBs of both groups perform identical functions. The *ssb* genes of this group are in most cases closely located to (putative) replication genes of their cognate phage replicons, in the same transcription unit. The

cognate SSBs are probably essential for the replication of the phages of this group because Kim and Batt could demonstrate inhibition of phage replication in an antisense approach [211]. Since in these experiments a premature termination of transcription of the (essential)  $\phi$ P4 $\alpha$ -type helicase gene could have been the reason for the observed inhibition of phage replication, more experiments are warranted for a decisive statement. Notably, only 1 of the *ssb* genes of this group ( $\phi$ ETA orf18) is found associated with one of the set of 40 replication initiator genes of Gram(+)-specific phages discussed in Section C3.1.2. . Also *vice versa*: with one exception ( $\phi$ BK5-T, see group 6), only *ssb* genes of group 7 are found associated with the replication initiator genes of the set of 40. Apparently, mosaicism, i.e. HGT among phages, did not mess up phage phylogeny entirely.

**Group 4: phage M13 gene V SSB.** The 6 orthologues of the  $\phi$ M13 gene V SSB are found in the genomes of filamentous phages with a genome size of ~7 kb. Despite their high degree of homology, these proteins lack a consensus C-terminus.

**Group 5: phage  $\lambda$  Ea10 SSB.** The precise function of the  $\lambda$  Ea10 SSB is not known, it is however not required for initiation of  $\lambda$  replication. The almost complete identity – including the weakly acidic tail – of the 4 Ea10 orthologues suggests that these proteins have a common origin in one of the four phages and were spread only recently.

**Group 6: miscellaneous SSBs.** This group contains 11 SSBs, which did not show any similarity to SSBs of one of the other groups.  $\phi$ BK5-T orf48 encodes a SSB that has no sequence similarity with either group 3 nor group 7 SSBs, except for the acidic C-terminus, which is highly similar to that of the group 7 SSBs of Gram(+)-specific phages. The 6 homologues of  $\phi$ 29 gene 5 SSB seem to be a unique type of SSB involved in the protein-primed replication of the linear genomes of related Bacillus phages [155,287]. Also *E. coli* phage  $\phi$ PRD1 replicates via the protein-primed mode, and its SSB might be involved in this process. The specific function of  $\phi$ N4 SSB is not known. This phage replicon requires *E. coli* SSB for its transcription (see above). The  $\phi$ f1 geneV and the  $\phi$ Pf3 orf78 SSBs of filamentous phages have no apparent sequence similarity with group 4 SSBs, despite a comparable size.

The heterogeneity of the SSBs of this group suggests to us that in addition to the known SSB types more phage proteins with yet unknown function may be identified as SSBs. At present, it is clear that single-strand DNA binding proteins evolved several times, and that the *E. coli*-type SSB seems to be the evolutionary most successful, by prevalence.

*Group 7: E. coli SSB and phage-encoded homologues.* Homologues of *E. coli* SSB are encoded by 5 Gram(–)-specific phages and 17 Gram(+)-specific phages, and in addition by many prophages (see below). The 5 Gram(–)-specific phages possess different replication gene modules (BRM Section 3.1.). Most other phages of this host range having similar replication modules lack *ssb* genes. Therefore we presume that the *ssb* genes in the genomes of  $\phi$ VP2,  $\phi$ VP5,  $\phi$ Bcep22, and  $\phi$ Sf6 represent accessory genes required for specific growth conditions rather than for replication under standard laboratory conditions – similar to  $\phi$ P1 [26]. In contrast to the genomes of Proteobacteria-specific phages, *ssb* genes are more frequently found in the genomes of phages of Firmicutes hosts. As already observed for the SSBs of group 3, the *ssb* genes of the 17 Gram(+)-specific phages of group 7 are in many cases closely located to (putative) replication genes of their cognate phage replicons, probably in the same transcription unit. SSBs of this group might thus be essential for the replication of some of these phages even under normal conditions, but apparently not for SPP1 [475].

Two regions of the *E. coli*-type SSB proteins are conserved among all phage-encoded and also the chromosomally encoded SSBs: i. the N-terminal ~120 res. responsible for ssDNA-binding and tetramerisation (see above), and ii. the C-terminal ~15 res. . The connecting region is responsible for the observed variation in length of these SSBs, and poorly conserved. The short C-terminal consensus sequences determined for the SSBs encoded by the Gram(–)-specific and the Gram(+)-specific phages are slightly different but correspond well to the C-termini of the cognate SSBs of their hosts. They may be considered phylon-specific SSB- fingerprints. In addition to their conservation in the SSBs, C-terminal tails containing ~4 acidic residues and a terminal LPF motif are also found in the initiator proteins of the Gram(+)-specific phages  $\phi$ LL-H,  $\phi$ mv4,  $\phi$ r1t, and the *B. anthracis*  $\phi$ LambdaBa04. Although this feature is not found in other phage initiator proteins, and might be coincidental therefore. It could be rewarding to analyse the role of this putative signalling domain – by experiments.

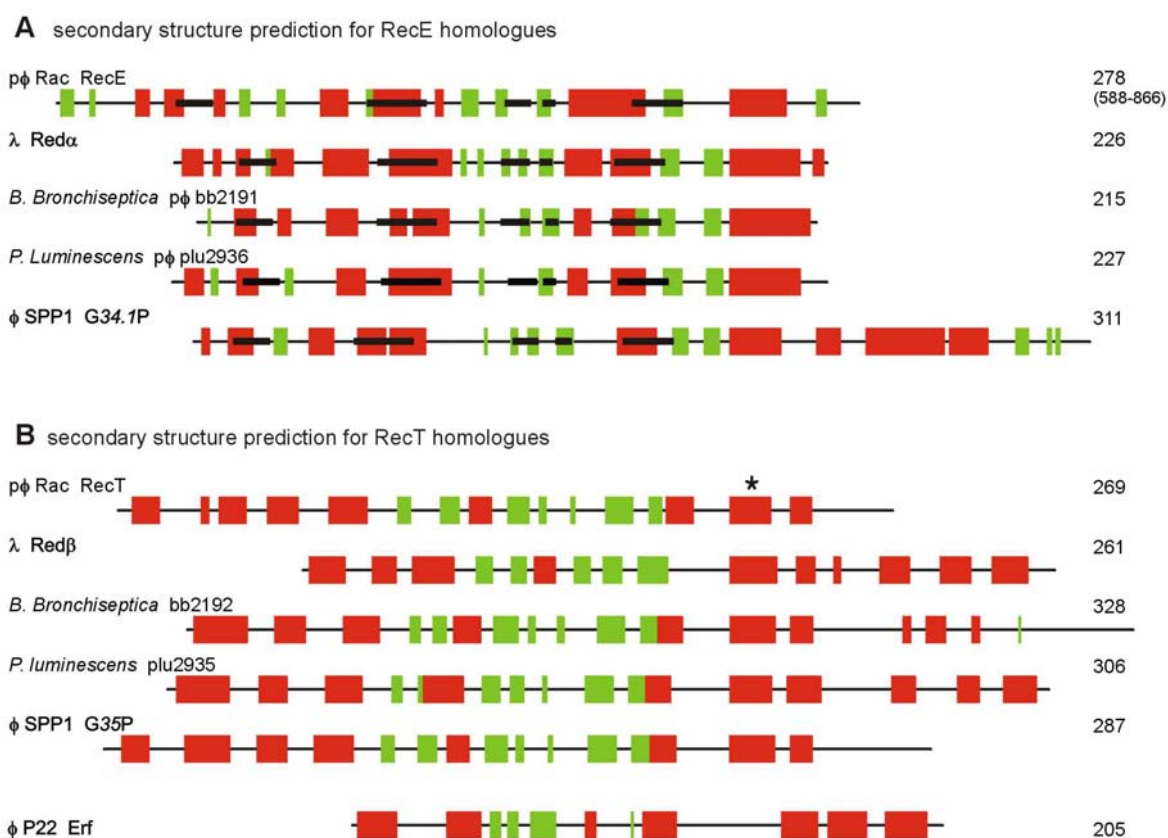
81 of the 136 presently known bacterial genomes contain a single *ssb* gene encoding a protein with significant similarity to *E. coli* SSB. 24 genomes code for either SSBs with an atypical (e.g. *Mycoplasma* sp.) or truncated C-terminus, or contain a C-terminally truncated in addition to a full-length *ssb* gene. However, 32 genomes code for two or more full-length SSBs with significant similarity to SSB<sub>Eco</sub>. Certainly, the 11 full-length and 4 truncated *ssb* genes in the genome of a *Phytoplasma* (sp. onion yellows) represents an extreme, but 3–4 *ssb* genes are frequently found in the genomes of Bacillales. On closer inspection, "surplus" *ssb* genes

are in most cases readily identified as parts of prophages. But the significant similarity of all *E. coli*-type SSB proteins has already led to confusion in gene assignment for some bacterial genomes containing multiple *ssb* gene(s). E.g. the locus sa0353 of the NCBI genome entry for *S. aureus* N315 is marked as "single-strand DNA-binding protein of phage phi PVL"; the gene is located between *rpsF* and *rpsR*, a feature also conserved in the genomes of *Bacillus* sp & *Listeria* sp., respectively. Locus sa1792, the second SSB-encoding gene in *S. aureus* N315 and marked "single-strand DNA-binding protein" is next to sa1791 encoding a prophage initiator gene, which suggests that sa1792 is a prophage *ssb* gene. Locus sa0353 represents the chromosomal *ssb* gene, therefore.

The presence of multiple *E. coli*-type *ssb* genes as part of prophages in ~25% of the presently known genomes is explained by the observation that temperate phage replicons preferably contain *E. coli*-type *ssb* genes. However, also homologues of the temperate phage  $\phi$ A2 orf34 gene are found in the genomes of e.g. *L. innocua* (locus lin2588) and *S. pyogenes* (locus SpyM3\_0960). There are, in addition, rare examples for the acquisition of *ssb* genes from virulent phages by chromosomal replicons. The *Xylella fastidiosa* locus XfasO1896 and the *X. campestris* locus XCC2059 encode proteins that are homologues of the  $\phi$ fd geneV SSB. Also the *rstB1* gene, part of the CTX $\phi$  prophage of *V. cholerae*, is a  $\phi$ fd geneV homologue [95]. *P. aeruginosa* strain PAO1 contains a chromosomal copy of the  $\phi$ Pf1 *ssb* gene (locus PA0720). The genome of *P. putida* KT2440 (locus PP2267) encodes a homologue of the  $\phi$ T7 gene2.5 SSB [397]. The phage-encoded *ssb* genes, which do not belong to the group 7, did not replace the chromosomal *E. coli*-type *ssb* gene as proposed by Forterre [126]. It seems more likely that phages acquired these *ssb* genes from their hosts by horizontal gene transfer (HGT) as suggested by Moreira [301] (see also BRM Section 4.). These observations emphasise the pivotal role of bacteriophages as vectors for HGT.

### C3.6.2. Recombination proteins

*Note in advance:* temperate phage replicons can integrate into the chromosomes of their hosts by site-specific recombination. Subsequently, the 'prophage' is propagated as part of the host chromosome. The shut-off of prophage replication by a  $\lambda$  *cI*-type repressor can not be compared to copy-number control mechanisms operating in plasmid replicons and chromosomes. Also, the recombination process resulting in prophage formation is mechanistically unrelated to replication. We will therefore not discuss integration here [153]. We note however that the highly conserved phage integrase



**Fig. C20** Secondary structure prediction for RecE/RecT analogues.

The secondary structure of *redα* is shown according to the crystal data [PDB 1AVQ]. For the other proteins, secondary structures are shown according to Jpred prediction [88]. The length of the individual protein sequences is indicated by the black line; all protein sequences are shown at the same scale. Red =  $\alpha$ -helices; green =  $\beta$ -strands. **A** Thicker black lines indicate the positions of the signature motifs defined by Vellani + Myers [463]. For the *E. coli* RecE sequence, only the functionally important res. 588-866 were included in the analysis [79]. **B** The helix containing the conserved MxxKTx<sub>2-3</sub>(R/K) motif (weakly conserved in *redβ*) is marked with an asterisk.

genes have been instrumental for the detection of prophages – or their remnants – in bacterial genomes [63].

The stringent coupling of replication and recombination was for many years considered a hallmark of the extravagant 'life-style' of bacteriophage T4 – without paradigmatic consequences for the studies of both nucleic acid processing pathways in bacteria and other phages [305,139]. The coupled replication and recombination of  $\phi$ T4 was addressed in BRM Section 2.5., and we will only discuss  $\phi$ T4 endonuclease VII here in the context of phage-encoded Holliday junction resolvases.

Also the recombination-dependent replication of bacteriophage Mu was taken as an exception, and comparable rather to the propagation of prokaryotic transposons. Because  $\phi$ Mu relies on the host replication/recombination machinery for propagation – except for its cognate transposase – we will not discuss it here, and

refer the reader to several instructive reviews and papers [67,161,314,420]. During the last decade, however, and due to an increasing awareness of the pioneering work of Kogoma and his collaborators [218], the importance of recombination for chromosome replication in *E. coli* came into the focus of replication research [238,85] (BRM Sections 2.5. + 2.6.).

In this section we shall discuss only those phage-encoded recombination pathways and proteins for which a direct functional connection to phage replication has been shown, or can be assumed. Proteins include 5'→3' exonucleases, singlestrand annealing proteins (SAPs), and Holliday junction resolvases. This section concludes with a brief survey of phage-encoded (putative) homologues of other known *E. coli* recombination protein. Phage recombination pathways and proteins have been reviewed extensively by Kuzminov [238].

**The  $\lambda$  Red $\alpha$ /Red $\beta$  (Exo/Bet) recombination pathway.** The observation that  $\lambda$  recombination was not affected by a *recA* mutation of the *E. coli* host led to the detection of phage-encoded recombination functions [113]. Accordingly,  $\lambda$  mutants were isolated that grew normally in a wild-type host, but failed to propagate in *rec<sup>-</sup>* hosts. The corresponding  $\lambda$  genes were identified as the genes coding for  $\lambda$  exonuclease (Red $\alpha$  or Exo) [254,414,64], and for beta protein (Red $\beta$  or Bet) [216,310]. Red $\alpha$  and Red $\beta$  occur *in vivo* as an equimolar complex [310]. During stepwise 5'→3' exonucleolytic degradation of one strand, Red $\alpha$  loads Red $\beta$  to the other strand. Red $\beta$  was shown *in vitro* to promote reannealing of complementary single strands. The structure of Red $\alpha$  [PDB AVQ] but not that of Red $\beta$  is known. More recently, Stahl and co-workers could show by genetic analysis that the Rap(NinG) protein (see below) is the 'cognate' Holliday junction resolvase for the  $\lambda$  Red $\alpha$ /Red $\beta$  recombination pathway [443].

Younger readers may wonder why decades passed between the detection of the  $\lambda$  encoded recombination genes by genetic experiments [127,49], early biochemistry [64], and the elucidation of the three-dimensional structure of red $\alpha$  [229]. The initially vague concepts of homologous recombination had to be transformed into robust experimental models describing the molecular mechanisms and the actors driving it. So, (almost) all of what is known today about homologous recombination in bacteria can be traced back to the ingenious genetic analyses of various combinations of *E. coli* *rec* and  $\lambda$  *red* mutants [238].

**The *E. coli*  $\phi$  Rac RecE/RecT recombination pathway.** Genetic analysis of *E. coli* *recBC* mutants (exonuclease V) revealed secondary suppressors termed *sbcA* that were able to (partly) restore recombination proficiency of the strain by activating the expression of exonuclease VIII [237,483]. The *sbcA* mutation(s) were mapped to the region of the cryptic Rac prophage of *E. coli* K12, and the gene encoding exonuclease VIII termed *recE*. Exonuclease VIII is a polypeptide of 866 res. in length but only the C-terminal 278 residues are required for exonuclease activity [79]. The function of the large N-terminal domain is not known. The *E. coli* *recBC sbcA* mutants not only allowed propagation of  $\lambda$  *red* phages – indicating that the activated function(s) are analogous to the  $\lambda$  *red* pathway – but also made the propagation of  $\lambda$  *red* phages independent of *recA*. The latter finding suggested that the *recE* region encoded a Red $\beta$ -like function in addition to exonuclease VIII [80]. Kolodner and co-workers determined the polypeptides expressed from the *recE* region: exonuclease VIII(RecE) and – from the same transcript – RecT with the characteristics of a SAP [159].

#### **The $\phi$ SPP1 G34.1P/G35P recombination pathway.**

As mentioned above, mutations in  $\phi$ SPP1 of either gene 34.1 or 35, respectively, result in a replication arrest phenotype, producing linear progeny molecules of sub-genome length [475]. The specific length of these molecules led Alonso and co-workers to hypothesise that unidirectional replication proceeds normally in such mutants from *oriL* to *oriR* – a duplicated origin structure though inactive for initiation. At *oriR*, bound G38P creates a road block for the replication fork, which cannot be overcome due to the lacking recombination activity of either G34.1P or G35P [12]. These authors could show that purified G35P protein is a SAP with properties comparable to RecT, and able to anneal single-stranded *oriL* DNA to the homologous part of a supercoiled template [12]. *Gene* 34.1 was recently shown to encode for the cognate 5'→3' exonuclease, termed G34.1P [280] or Chu [463]. Together, these studies demonstrate that the exonuclease/SAP recombination pathway is similar in Gram(–) and Gram(+)-specific phages with respect to the underlying molecular mechanism.

**The  $\phi$ P22 Erf recombination protein.** Early genetic analyses suggested that  $\phi$ P22 carries a RecA-like gene whose inactivation prevented the growth of the phage on an *E. coli* *recA* host [490,351,136]. The structure of  $\phi$ P22 Erf (essential recombination function) is not yet known. It was shown that the N-terminal ~140 res. of Erf are sufficient for ssDNA-binding and oligomerisation, the function of the C-terminus is not known [353,312]. The secondary structure prediction suggests that Erf is – at best – only distantly related with the RecT/Red $\beta$  SAPs although all known phage-encoded SAPs form fairly similar complexes with ssDNA (Fig. C20) [336]. In contrast to the phage recombination protein pairs discussed above, the  $\phi$ P22 Erf protein lacks a 'cognate' exonuclease. The *arf* gene downstream of *erf* in the P22 genome encodes a small protein (47 res.), which was shown by genetic analysis to be required for efficient recombination of  $\phi$ P22 [352]. Because of its small size it seems unlikely that *arf* possesses exonuclease activity but may instead recruit the yet unknown exonuclease. Since homologues of *arf* are not found in other phages this notion remains speculative.

**Evolutionary relationship of the RecE/RecT, G34.1P/G35P, and Red $\alpha$ /Red $\beta$  protein pairs.** The RecE/RecT recombination proteins encoded by the *E. coli* Rac prophage and the Red $\alpha$ /Red $\beta$  proteins of  $\lambda$  perform analogous functions. Despite subtle differences in protein activity and the failure to obtain a functional complementation in reciprocal assays [313], and despite the lack of significant primary sequence similarity, these protein pairs may be evolutionary related [219,61, 194].



**Table C31** Protein sequence similarity among RecE/RecT analogues.

		RecE	Red $\alpha$	bb2191	plu2936	G34.1P		RecT	Red $\beta$	bb2192	plu2935	G35P
$\phi$ Rac	RecE	x					RecT	x				
$\lambda$	Red $\alpha$	—	x				Red $\beta$	—	x			
<i>B. bronchiseptica</i> $\phi$	bb2191	—	29%	x			bb2192	31%	—	x		
<i>P. luminescens</i> $\phi$	plu2936	—	70%	28%	x		plu2935	28%	29%	41%	x	
$\phi$ SPP1	G34.1P	18% <sup>1)</sup>	—	25%	27%	x	G35P	40%	—	29%	28%	x

Only the functionally important res. 588-866 of *E. coli* RecE were included in the analysis [79]. BLAST analysis was performed using the NCBI 'bl2seq' software with default settings except that the 'low complexity' filter was disabled [445]. 1) value taken from Ayora *et al.* [12].

**Table C32** (Pro)phage-encoded  $\phi$ P22 Erf-like proteins.

A signature motifs:			A		B						
consensus			QxxGATSSYhRKYxhxxhFxx	x <sub>4-5</sub> DxD	KGxxSDAhKRAAVxWGIGRYLY						
			A	KR	L	G	CF	V	E		
B occurrence of motifs:			BLAST		BLAST		BLAST		motifs		
			$\phi$ P22		$\phi$ D3		$\phi$ ST64T		A B		
replicon	accession	gene	res.	1-150	151-205	1-200	201-275	1-180	181-235		
<i>C. acetobutylicum</i> $\phi$	NP_348558	cac1936	229	-	-	-	-	32	-		+
<i>T. thermophilus</i>	YP_005892	ttc1923	212	-	-	-	-	31	-		+
$\phi$ RM378	NP_835605	p018	164	-	-	-	-	39	-		+
$\phi$ bIL286	NP_076648	orf14	252	-	-	-	-	57	-		+
$\phi$ ul36	NP_663647	ORF252	252	-	-	-	-	57	-		+
$\phi$ ST64T	NP_720291	orf-235	235	-	76	-	-	x	x		+
$\phi$ Sf6	NP_958204	gene 28	235	-	63	-	-	94	81		+
$\phi$ ST104	YP_006369	ORF13	235	-	77	-	-	94	98		+
$\phi$ P22	NP_059596	erf	205	x	x	22	-	-	89		+
$\phi$ H-19B	AAD04639	erf	201	88	74	23	-	-	32		+
$\phi$ HK022	NP_037690	Erf	201	88	73	24	-	-	32		+
$\phi$ HK97	NP_037726	gp40	201	89	74	24	-	-	32		+
$\phi$ c2	NP_043534	e15	179	38	-	-	-	-	-		+
$\phi$ MM1	NP_150142	ORF9	235	36	-	-	-	-	-		+
$\phi$ 7201	NP_038308	ORF7	218	28	-	-	-	-	-		+
<i>S. agalactiae</i>	NP_687594	sag0565	224	23	-	30	-	-	-		+
$\phi$ NIH1.1	NP_438125	phiNIH1.1_12	224	23	-	30	-	-	-		+
<i>S. pyogenes</i> MGAS8232	NP_607818	spyM18_1796	224	23	-	30	-	-	-		+
$\phi$ PV83	NP_061607	orf17	212	24	-	26	-	-	-		+
<i>N. punctiforme</i>	ZP_00106022	npun0391	257	23	-	23	-	-	-		+
<i>L. johnsonii</i> $\phi$ Lj771	AAK27920	orf223	221	[42]	-	25	-	-	-		+
$\phi$ LL-H	AAL77543	orf178b	178	-	-	24	-	-	-		+
$\phi$ mv4	AAG31330	orf244	244	-	-	25	-	-	-		+
$\phi$ T1	YP_003917	28	226	-	-	23	-	-	-		+
$\phi$ SLT	NP_075480	orf216	216	-	-	-	-	-	-		+
pcp9 <i>B. burgdorferi</i>	AAG36971	orf4	302	-	-	27	-	-	-		+
<i>L. gasseri</i> $\phi$	ZP_00046420	lgas0587	270	-	-	-	-	-	-		+
<i>L. johnsonii</i> $\phi$ Lj928	NP_958522	ljo_1449	261	-	-	25	-	-	-		+
<i>D. vulgaris</i>	YP_011247	duv2032	237	-	-	27	-	-	-		+
<i>L. monocytogenes</i> $\phi$	NP_465842	lmo2318	232	-	-	26	-	-	-		+
$\phi$ D3	NP_061548	orf52	275	-	-	x	x	-	-		+
<i>B. cepacia</i> $\phi$	ZP_00219218	Bucepa02006186	215	-	-	30	-	-	-		+

Also The G34.*IP*/G35P protein pair encoded by *B. subtilis* phage  $\phi$ SPP1 belongs to this group of functional analogous exonuclease/SAP pairs, and Alonso and co-workers already pointed to the significant sequence similarity between RecT and G35P (40% ident. res.) [12]. The availability of many new genome sequences of (pro)phages and bacterial species made it attractive to re-evaluate the evolutionary relationship of these protein pairs. BLAST searches revealed protein sequences encoded by adjacent genes of (putative) prophages in the genomes of *B. bronchiseptica* ( $\beta$ - Proteobacteria) and *Photorhabdus luminescens* ( $\gamma$ -Proteobacteria) whose intermediate similarity values support the notion of a common origin of the phage- encoded exonuclease/SAP pairs (Table C31) [194]. This observation illustrates neatly the mutual benefits that phage research and genomic sequencing projects can achieve by combining their data.

The *P. luminescens* plu2935 protein shows significant similarity to RecT and Red $\beta$ . Likewise, the *P. luminescens* plu2936 protein is a homologue of Red $\alpha$ , but also similar to the *B. bronchiseptica* bb2191 protein, which in turn is similar to G34.*IP* of  $\phi$ SPP1. Because the sequences differ considerably in size we performed a secondary structure prediction analysis in search for related structures (Fig. C20). Both groups of proteins show a fairly similar 'core' of alternating  $\beta$ -stranded and  $\alpha$ -helical regions. We infer from this analysis that a common evolutionary origin of these phage-encoded exonuclease/SAP pairs is very likely. However, this analysis suggests that  $\phi$ P22 Erf is only very distantly related to the other 3 SAPs, if at all.

**Other phage-encoded (putative) exonuclease/SAP gene pairs.** As discussed above and also noted by others [463], phage-encoded exonucleases and SAPs form functional pairs and their genes are usually located close to each other in the genomes of their cognate replicons. A complex picture emerges when the BLAST similarities of the RecE/RecT, the G34.*IP*/G35P, the Red $\alpha$ /Red $\beta$  proteins, and the  $\phi$ P22 Erf protein to other phage-encoded (putative) exonuclease/SAP proteins are plotted in a matrix that includes the pairwise arrangement of the respective genes (Figs. C21 + C22). We analysed 91 phage and prophage sequences by this approach. A further known  $\sim$ 30 known sequences mostly from prophages were omitted because they did not form novel (sub)groups. Also we excluded the multiple exonuclease/SAP proteins encoded by *Phytoplasma* (sp. onion yellows) (see above for *Phytoplasma ssb* genes, Section C3.6.1). The matrix suggests the definition of six distinct groups:

*Group I* contains (putative) SAPs with similarity to  $\phi$ bIL286 orf 14 (Fig. C21; upper left part). The N-termi-

nal  $\sim$ 150 res. of these proteins show weak similarity to the eukaryotic Rad52-type SAPs. The bipartite human Rad52 protein [PDP 1H2I] contains within its N-terminal DNA-binding domain ( $\sim$ 200 res.) in a larger  $\alpha$ -helix the motif **KEAVTDGLKRALRSFGNALGNC**, which distantly resembles the motif B present in all proteins of this group at a corresponding position (Table C32) [416]. The Rad52 C-terminus might be responsible for the interaction with the Rad51 recombination protein [416]. Despite a reasonable conservation of the Ntermini, the SAPs of this group lack similarity in their C-termini. Moreover, The C-termini of the (putative) SAPs of  $\phi$ ST64T,  $\phi$ Sf6, and  $\phi$ ST104 are highly similar to the C-terminus of  $\phi$ P22 Erf (Table C32). Possible protein-protein interaction partner(s) of  $\phi$ P22 Erf are not known, therefore we can only mention here the intriguing C-terminal similarity of two groups of SAPs with completely unrelated N-terminal DNA-binding domains. A 'cognate' 5'→3' exonuclease could not be detected in the phage genomes of this group.

*Group II* contains SAPs with similarity to the  $\phi$ P22 Erf protein (Fig. C21; upper left to middle part). The crude BLAST approach did not reveal any similarity to one of the other SAPs. The initial sorting by eye already suggested that several subfamilies exist within this group. A more detailed BLAST search identified motif A (pos. 97-124 in  $\phi$ P22 Erf) as the most highly conserved region in all sequences of this group. Motif A is contained within the presumed DNA-binding region of  $\phi$ P22 Erf but its exact functional role is not known [312]. We tentatively use this motif here to discriminate sequences belonging to group I from those belonging to group II (see above). Motif A could also be identified unequivocally at corresponding positions in sequences that failed to give BLAST matches with  $\phi$ P22 Erf or  $\phi$ D3 orf52 as query, e.g.  $\phi$ SLT orf216 (Table C32). The C-terminal  $\sim$ 50 res. of the Erf-like proteins are highly divergent; only for the closest relatives of  $\phi$ P22 Erf similarity was found to three sequences from group I (see above).

Although Erf was originally identified in the  $\phi$ P22 phage of enteric hosts, it is now clear that *erf* genes are present in many (pro)phage genomes covering a wide host range. BLAST searches with red $\alpha$ , G34.*IP*, or RecE as query did not lead to the detection of a 'cognate' 5'→3' exonuclease for this group. There are three notable exceptions: i. the  $\phi$ T1 gene 29 – directly adjacent to gene 28 encoding an Erf-like protein – shows significant similarity to the RecE-like proteins gp60 of  $\phi$ Che9c and orfB of a *Legionella pneumophila* prophage [356], and ii. the dvu2033 protein of a *Desulfovibrio vulgaris* prophage is similar to the G34.*IP*-like protein encoded by the *C. tetani* E88 prophage locus ctc02146 (not shown). The third exception is  $\phi$ D3: while orf 51 is a

Red $\alpha$ -like exonuclease, orf 52 is a Erf-like SAP (Fig. C21). These three examples emphasise that the pairwise occurrence of exonuclease/SAP genes in phage genomes is a recurring feature although other than the prototype protein pairs are possible.

*Group III* contains gene pairs with similarity to the  $\lambda$  Red $\alpha$ /Red $\beta$  proteins (Fig. C21; lower right part). An unequivocal sorting of the protein pairs in this group is only possible for the very close relatives of  $\lambda$ . There exist three minor subgroups: i. the  $\phi$ D3-type with a Red $\alpha$ -like exonuclease and a Erf-like SAP (see above), ii. the *B. bronchiseptica* prophage-type with a Red $\alpha$ -like exonuclease (bb2191) and a RecT/G35P-like SAP (bb2192), and iii. the *P. rettgeri* prophage-type with a G34.IP-like exonuclease (orf59) and a Red $\beta$ -like SAP (orf68). A Red $\beta$ -like SAP gene but no 'cognate' exonuclease gene could be identified in the genomes of  $\phi$ LC3 and  $\phi$ 31.1.

*Group IV* contains proteins with similarity to the  $\phi$ SPP1 G34.IP/G35P protein pair (Fig. C22; upper left part). As can be deduced from the matrix, the G35P-like proteins of this group form a larger RecT-type family of phage encoded single-strand annealing proteins together with the SAPs of the  $\phi$ bIL309 and  $\phi$ Rac groups. The RecT family has no detectable structural similarity to proteins from the Erf family, despite their analogous function. There is however a distant sequence similarity and (predicted) structural relatedness between members of the RecT family and Red $\beta$  proteins: the (putative) *P. luminescens* SAP plu2935 shows similarity to Red $\beta$ -like as well as RecT-like proteins (see also above).

*Group V* contains SAPs with similarity to the  $\phi$ bIL309 RecT protein (Fig. C22; lower right part). All members of this group are encoded by Gram(+)-specific phages. The enzymatic activity of the  $\phi$ bIL309 RecT protein – or of any other protein from this group – has not yet been demonstrated, and the function assigned only by the convincing similarity to RecT and G35P. BLAST searches with Red $\alpha$ , G34.IP and RecE as query – or with any of the exonucleases of these groups – did not lead to the detection of a 'cognate' 5'→3' exonuclease for this group. For some of the homologues of the  $\phi$ bIL309 orf13 protein we could detect some weak similarities to RecB and RuvA proteins. It seems thus possible that the orf13-like proteins are involved in recombination, although their inclusion in the matrix is purely speculative at present. However, the 'canonical' gene arrangement – the exonuclease gene preceding the SAP gene within the same transcription unit – is reversed in the  $\phi$ bIL309 RecT/orf13-like gene pairs.

*Group VI* contains proteins with similarity to the *E. coli*  $\phi$ Rac RecE/RecT proteins (Fig. C22; lower right part). For several RecE-like proteins, a 'cognate' SAP could not be detected, e.g. for the *E. coli* prophage  $\phi$ CP-933P, or *Salmonella* sp. phage  $\phi$ ST64B. Interestingly, the putative exonucleases encoded by the  $\phi$ CP-933U, P, and M show significant similarity to the DexA exonuclease of the  $\phi$ T4 group of phages – which however is a 3'→5' exonuclease.

We can draw the following conclusions:

- i. within all groups, members with high, moderate and low similarity to each other were found. This suggests that the major possible variations within and among the groups defined here could be identified by analysing sequences from 91 (pro)phages.
- ii. in roughly one half of the cases, exonuclease genes and SAP genes are present as pairs in phage genomes. For the other half, it is presently not possible to decide whether the genes are present as single genes or whether the 'cognate' partner has not yet been identified.
- iii. the exonuclease and SAP genes occur in favoured combinations, but examples exist for any possible combination. This emphasises that the presence of both types of recombination proteins is more important for phage recombination than a particular exonuclease or SAP. We could however not detect any 'mixing' of the phage-encoded exonuclease/SAP gene pairs with the (functional analogous) chromosomal *recBCD/recA* genes.
- iv. with regard to sequence similarity, there is no apparent distinction between exonuclease/SAP pairs from Gram(+) and Gram(-)-specific phages, in contrast to the replication initiator proteins (Section C3.1.2.). This suggests that – consistent with experimental results (see above) – the phage-encoded recombination proteins perform their function independent from the host recombination machinery.
- v. in most cases, the gene encoding the exonuclease precedes the SAP-encoding gene within the same transcription unit. This gene arrangement is independent of the type of exonuclease and SAP, and may have enabled successful recombination events combining different types of exonuclease and SAP genes. In some cases, small intervening orfs (<100 res.) were found to separate the exonuclease and SAP genes.
- vi. the plasmid encoded Erf-like protein orf14 of the ~9 kb plasmid pcp9 of *B. burgdorferi* [430], and the (putative) 5'→3' exonuclease/SAP pairs of the 217 kb plasmid pCAR1 of *P. resinovorans* [269] and of the 199 kb plasmid pRts1 of *P. vulgaris* [311] are included in the matrix. No other plasmid-encoded proteins with similarity to either exonucleases or

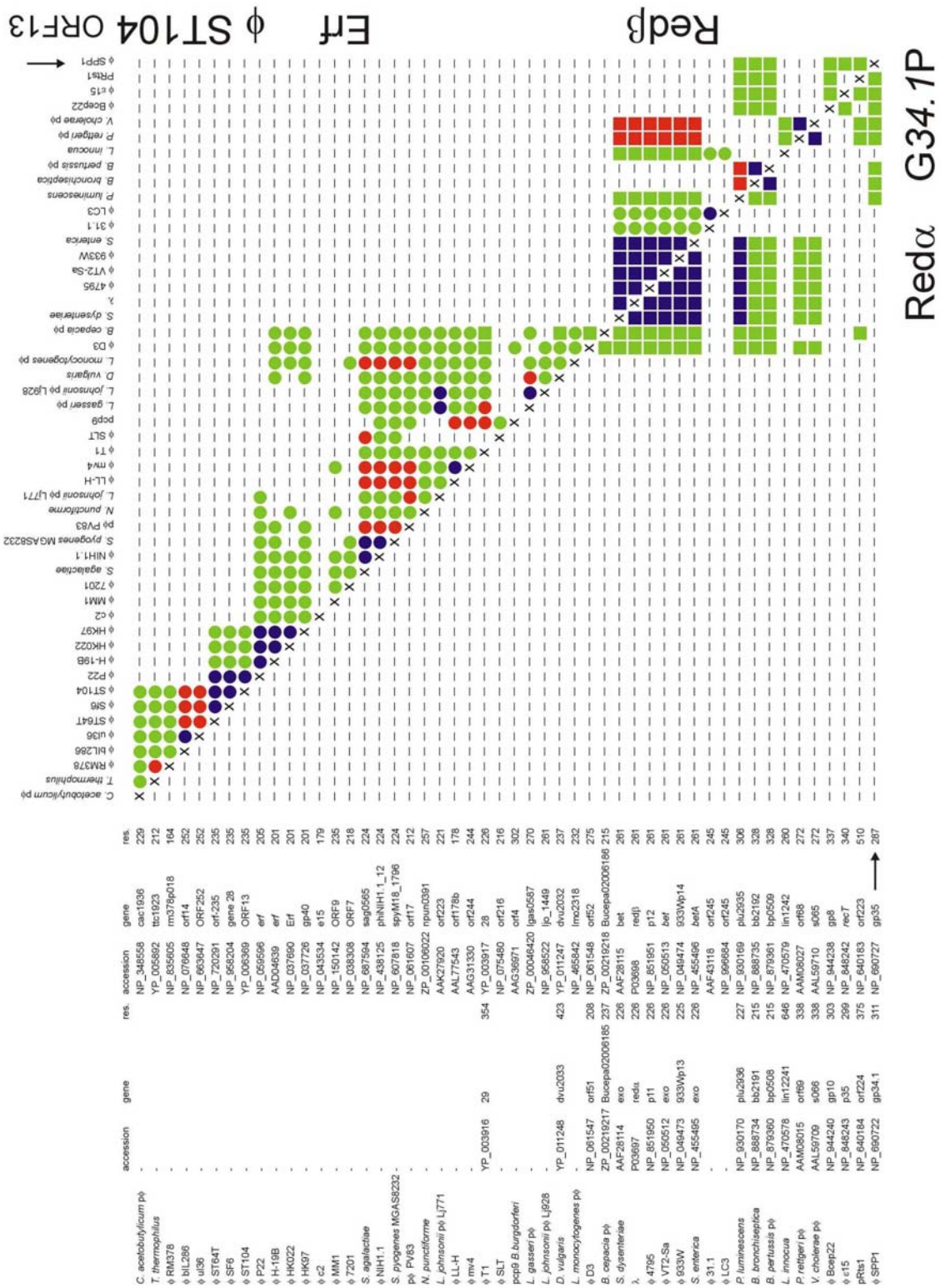


Fig. C21 Phage-encoded 5'→3' exonuclease / SAP protein pairs; part A.



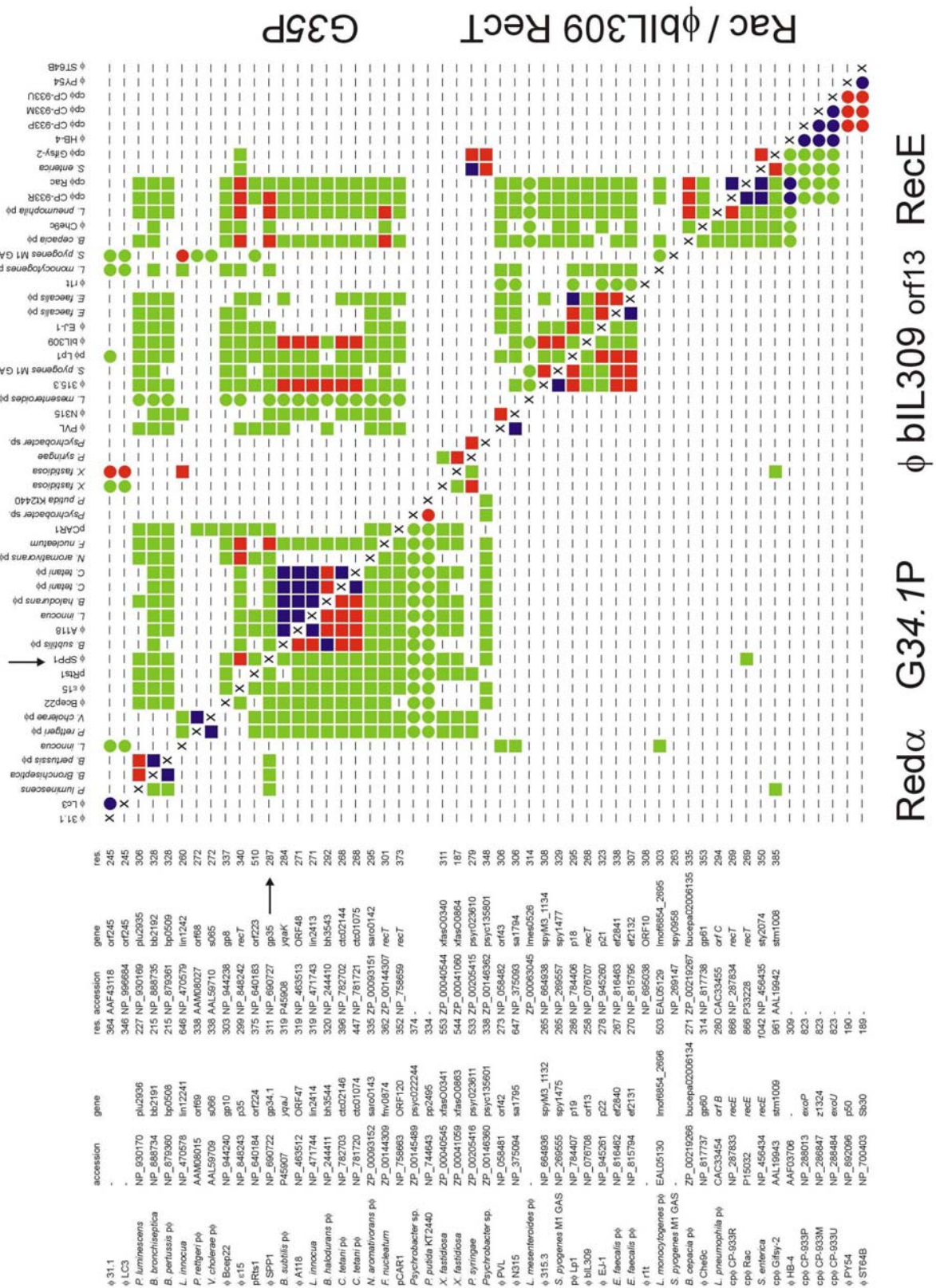


Fig. C22 Phage-encoded 5'→3' exonuclease / SAP protein pairs; part B.

Legend to Figs. C21 + C22. The BLAST similarity values for the pairwise comparison of the exonuclease proteins are shown in the lower left part, i.e. below the diagonal. The values for the pairwise comparison of the SAP proteins are shown in the upper right part, i.e. above the diagonal. Colour code: blue = similarity 61-100% ident. res.; red = similarity 41-60% ident. res.; green = ~20%-40% ident. res.. Squares indicate that in the pairwise comparison of exonuclease and SAP sequences the respective proteins have been detected in both compared sets. Circles indicate that in the pairwise comparison one protein was not detected in at least one of the compared sets. For reasons of space the original matrix was split into two separate figures. The region of overlap between both figures is indicated by vertical + horizontal arrows at the positions of the  $\phi$ SPP1 genes.

SAPs were detected by BLAST searches. Therefore, the presence of recombination genes in plasmids is exceptional, and it is not known whether the mentioned recombination proteins play any role during propagation of their plasmid replicons.

Despite the pairwise occurrence of exonuclease/SAP genes in many phage genomes, other genes encoding (putative) recombination proteins are rarely found close by. Among the few exceptions are the RdgC-like protein encoded by the gene upstream of the RecE-homologue of a *Burkholderia cepacia* prophage (locus Bucepa02006134, not included in Fig. C22), and a RusA-like protein encoded upstream of the Erf-like protein encoded by a *S. pyogenes* prophage (locus spyM18\_1796). SSB and replication genes, however, are found in many genomes of Gram(+)-specific phages in close vicinity of recombination genes, in some cases apparently within the same transcription unit, e.g. in  $\phi$ A118 (BRM Section 3.1.).

**Holliday junction resolvases.** DNA endonucleases that recognise and cleave the specific structure of four-way or crossover junctions have been termed Holliday junction resolvases. In the generally accepted model, the resolution of the Holliday junction(s) demarks the penultimate stage of homologous recombination, which is i. initiated by a 'loose end' and RecA-driven strand exchange/invasion, ii. leads to formation of a Holliday junction upon reciprocal strand exchange, iii. is continued by branch migration, and iv. is terminated by gap closure after resolution of Holliday junctions. In *E. coli*, branch migration is driven by the action of the RuvAB helicase, while Holliday junctions are resolved by the action of RuvC [61]. Recent experiments point to the possibility that Holliday junction resolution by RuvC may also trigger the initiation of RDR [403,281]. Replication fork arrest may lead to fork regression and the annealing of the newly synthesised strands, i.e. the formation of a Holliday junction as the 'ankle' of a chicken-foot structure, literally spoken. When the repair function of exonuclease V (RecBCD) is disabled, RuvC sets the strand scission(s) required to re-initiate replication via RDR. We speculate that a similar mechanism may also operate during the switch from  $\theta$ DR to  $\sigma$ DR required to obtain phage genome concatemers (see above).

Phage-encoded Holliday junction resolvases include  $\phi$ T4 endonuclease VII (gp49),  $\phi$ T7 endonuclease I (gene 3),  $\lambda$  Rap, orf3 protein with limited similarity to

*E. coli* RuvC from the lactococcal phage  $\phi$ bIL66, and RusA, encoded by the *E. coli* K12 cryptic prophage DLP12. Despite their highly divergent three-dimensional structures and unrelated protein sequences all these enzymes form dimers and require metal ions for enzymatic activity. The structure and function of the phage-encoded Holliday junction resolvases have been comprehensively reviewed by Sharples who also defined the five groups of phage-encoded Holliday junction resolvases [408]. Therefore we confine ourselves here to a discussion of the many new sequences detected by BLAST searches for the five junction resolvase groups.

*Gene 49 (gp49) endonuclease VII of phage T4.* The crystal structure of the dimeric  $\phi$ T4 endonuclease VII is known [PDB 1EN7], and functionally important residues have been determined: D40, H41, N62, and E65 are closely located to the catalytic core, residues F72 and W87 in the C-terminal region are involved in dimerisation [35,36]. Endonuclease VII contains four cysteine residues in a specific arrangement that is also found in several type II restriction endonucleases. All phages of the  $\phi$ T4-group encode a homologous endonuclease VII except for *Aeromonas* sp. phages  $\phi$ 65 and  $\phi$ 25 (sequences not yet completed) (Table C33). Motif E but not motif D is conserved in all endonuclease VII-type proteins detected by BLAST searches. The dimerisation of the proteins lacking a detectable motif C may follow a route different from that of  $\phi$ T4 gp49.

Notably, phages that encode a homologue of endonuclease VII but a DNA polymerase of the Pol I-type include  $\phi$ D29,  $\phi$ L5,  $\phi$ K1-5, and  $\phi$ SP6 (Section C3.5.1.). The *yajD* gene of the insect endosymbiont *Sodalis glossinidius* ( $\gamma$ -Proteobacteria) encodes a protein with similarity to  $\phi$ T4 gp49. By using *S. glossinidius* YajD as query, BLAST searches readily revealed homologues in the genomes of *E. coli*, *S. enterica*, *Yersinia pestis*, several other  $\gamma$ -Proteobacteria, and – notably – several *Streptomyces* species/strains. The function of these YajD proteins is not known but the conservation of the size and the structural motifs suggests that it may be closely related to the function of  $\phi$ T4 endonuclease VII. In the *E. coli* K12 genome, the *yajD* gene is located closely downstream of the *sbcC* and *sbcD* genes, which are related to the  $\phi$ T4 gp46/gp47 proteins. This suggests that these genes were originally acquired together from a  $\phi$ T4-like phage. The presence of a  $\phi$ T4 gp49 homologue in the genomes of  $\phi$ D29,  $\phi$ L5,  $\phi$ K1-5,  $\phi$ SP6, and in the genomes of the other phages listed in Table C33

**Table C34**  $\phi$ T7 endonuclease I (gene 3) homologues.

motifs:				A	B
consensus				xYTPDFhhP	hh <b>h</b> ETKGx
replicon	accession	gene	res. ident.		
$\phi$ T7	P00641	gene 3	149 x	TYTPDFLLP	IFV <b>ET</b> KGL
$\phi$ A1122	NP_848275	p13	151 87%	VYTPDFLLP	IF <b>I</b> ETKGL
$\phi$ T3	P20314	gene 3	152 81%	LYTPDFLLP	IF <b>I</b> ETKGL
$\phi$ YeO3-12	NP_052083	gene 3	153 86%	LYTPDFLLP	IF <b>I</b> ETKGL
$\phi$ gh-1	NP_813757	p11	147 60%	KYTPDFVLA	II <b>I</b> ETKGI
<i>P. putida</i> KT2440 p $\phi$	NP_744417	pp2268	141 50%	KYTPDFALA	II <b>V</b> ETKGR
$\phi$ P60	NP_570326	P60_orf16	66 47%	KYTPDFFLP	VILEVKGF
	NP_570327	P60_orf17	54 43%		
$\phi$ SIO1	NP_064756	p19	134 47%	TYTPDFVLP	VII <b>E</b> TKGR

The motifs A + B represent highly conserved stretches identified by BLAST searches with  $\phi$ T7 gene 3 protein as query; conserved functional residues are shown in bold type [157].

suggests, in addition, that a *yajD* gene was acquired more than once by phages via HGT from a chromosomal copy – possibly even three times (!) by  $\phi$ Omega. This applies also to  $\phi$ YeO3-12, which encodes for two different (putative) Holliday junction resolvases: p47 as a protein similar to  $\phi$ T4 gp49, and a homologue of  $\phi$ T7 gene 3 protein. However, the issue of the spread of the  $\phi$ T4 gp49- like proteins among bacterial chromosomes and phages requires a deeper analysis before a firm statement can be made.

*Gene 3 endonuclease I of phage T7.* The crystal structure of the dimeric  $\phi$ T7 endonuclease I is known [PDB 1MOD], and functionally important residues have been determined: D55, E65, and T66 are closely located to the catalytic core and involved in manganese ion binding (Table C34) [157]. With two exceptions, the phages encoding a resolvase similar to  $\phi$ T7 gene 3 protein also encode a DNA polymerase with significant similarity to  $\phi$ T7 DNA polymerase (Section C3.5.1 and Table C25). A contiguous endonuclease protein is probably synthesised from the overlapping  $\phi$ P60 genes orf16 (N-terminus) and orf 17 (C-terminus), and the split into two genes may have resulted from a sequencing error (F. Chen; pers. communication) [75]. Phages  $\phi$ SP6 and  $\phi$ K1-5 encode  $\phi$ T7-type DNA polymerases but their cognate resolvases are similar to  $\phi$ T4 endonuclease VII (see above). *Genes* encoding resolvases have not yet been identified in the genomes of  $\phi$ PaP3 and  $\phi$ Felix01, which also encode DNA polymerases of the  $\phi$ T7-type (see Table C25).

*Rap protein of phage Lambda.* The  $\lambda$  Rap (NinG) gene was identified by genetic analysis of recombina-

tion deficient  $\lambda$  mutants [178]. The function of Rap as a Holliday junction resolvase was confirmed by biochemical analysis [409]. The three-dimensional structure of  $\lambda$  Rap is not yet known. The particular arrangement of 8 cysteine residues in the  $\lambda$  Rap protein was proposed to mediate zinc ion binding, and is conserved throughout the entire group of Rap-like proteins (see Table C35) [408].

Two variants exist among the Rap-like proteins, which differ in the spacing of the two cysteine residues in motif A (Table C35). Because this difference in spacing is found also among otherwise highly conserved Rap homologues (> 90% ident. res.) it might not have functional consequences. In contrast to the RusA-like resolvases, the presence of a Rap encoding genes is entirely confined to Gram(-)-specific phages (see below). Interestingly – and also in contrast to the RusA proteins – only two distinct similarity groups exist among the Rap proteins: the proteins with almost perfect identity to  $\lambda$  Rap and the proteins with significant similarity, i.e. with identities ranging from >30% to <50% (listed separately in Table C35). A group of proteins with similarities lower than 30% is missing, which was found for almost all phage recombination and replication proteins analysed in this study. This observation indicates that either the expected low-similarity genes are still missing from the data sources because the candidate phages have not yet been identified, or that the Rap gene is a very recent evolutionary development for which the precursor has yet to be identified.

*Phage-encoded resolvases with similarity to E. coli RuvC.* Chopin and co-workers could show that the orf2 and orf3 proteins of  $\phi$ bIL66 are responsible for a struc-



**Table C33**  $\phi$ T4 endonuclease VII (gp49) homologues.

motifs:					A		B		C		D	E
consensus					xxChhCxx		xhDHDHxxx		hhCxxCNxxExxx		xxhxx	xhxahx
replicon	accession	gene	res.	ident.								
$\phi$ T4	P13340	gp49	157	x	GK <b>CL</b> ICQR	27	HL <b>DH</b> DHELN	LL <b>CN</b> LCNAAEGQM	HK <b>F</b> NR	YLEWLE		
$\phi$ RB69	NP_861776	gp49	157	89%	GK <b>CPI</b> CHR	27	HL <b>DH</b> DHELN	LL <b>CN</b> LCNAAEGQM	HK <b>F</b> NR	YLEWLE		
$\phi$ RB49	NP_891632	p061	157	44%	GIC <b>P</b> L <b>C</b> QL	27	HL <b>DH</b> DHDL	LL <b>H</b> SC <b>CN</b> RLEGLS	G <b>D</b> FKR	FVVWAT		
$\phi$ KVP40	NP_899379	49	151	38%	DR <b>C</b> V <b>L</b> CGT	22	HT <b>DH</b> DHNTG	VL <b>C</b> R <b>A</b> C <b>N</b> TYEGVV	HK <b>F</b> TR	YVQWLR		
$\phi$ Aeh1	NP_943932	gp49	161	49%	GIC <b>P</b> L <b>C</b> KR	27	HL <b>DH</b> DHELN	LL <b>C</b> CF <b>CN</b> KFEGMV	HE <b>F</b> MR	YITQLK		
$\phi$ 44RR2.8t	NP_932421	gp49	157	48%	G <b>C</b> CL <b>K</b> CKN	27	HL <b>DH</b> DHALE	LL <b>C</b> AR <b>CN</b> MLEGMI	HK <b>F</b> IR	YIEWMR		
$\phi$ D29	O64250	gp59	153	[34%]	GR <b>C</b> Y <b>I</b> CQR	22	S <b>V</b> D <b>H</b> DHKTG	LL <b>C</b> TM <b>CN</b> KYILGW	-	CIEMLQ		
$\phi$ L5	Q05272	L5p56	164	[41%]	GR <b>C</b> Y <b>I</b> CQR	21	S <b>V</b> D <b>H</b> DHKTG	LL <b>C</b> TM <b>CN</b> KYTLGW	-	CIEFFK		
$\phi$ K1-5	AAR90063	gp21	136	[36%]	WK <b>C</b> P <b>L</b> CGG	22	VL <b>DH</b> DHETG	V <b>V</b> CR <b>G</b> C <b>N</b> GAEGKI	G <b>V</b> ISG	QLQWLE		
$\phi$ SP6	NP_853582	gp22	136	[36%]	WK <b>C</b> P <b>L</b> CGG	22	VL <b>DH</b> DHETG	V <b>V</b> CR <b>G</b> C <b>N</b> GAEGKI	G <b>V</b> ISG	QLKWLE		
$\phi$ S-PM2	AAP92668	g49	142	[40%]	G <b>V</b> CA <b>I</b> CKG	22	C <b>V</b> D <b>H</b> DHETG	LL <b>C</b> R <b>N</b> C <b>N</b> MMLGQV	-	-		
$\phi$ Omega	NP_818437	gp138	160	[30%]	GR <b>C</b> AL <b>C</b> QR	21	AVE <b>H</b> DHKTG	I <b>C</b> CG <b>P</b> C <b>N</b> LGLVLGH	-	DIAFFE		
	NP_818388	gp89	166	[42%]	CAN <b>P</b> GCRA	29	HL <b>DH</b> DHDC	I <b>L</b> CP <b>A</b> C <b>N</b> LMLGSA	-	FRLYGK		
	NP_818499	gp199	229	[36%]	GK <b>C</b> AN <b>P</b> GC	33	P <b>V</b> D <b>H</b> DHACC	I <b>L</b> CP <b>P</b> C <b>N</b> TTLGQM	K <b>R</b> LRG	LRFWGV		
$\phi$ KMV	NP_877462	gp23	146	[30%]	K <b>L</b> C <b>P</b> L <b>C</b> GK	23	V <b>M</b> D <b>H</b> DHETG	VL <b>H</b> RS <b>C</b> NTAEGKI	-	IIPYLH		
<i>S. coelicolor</i> A3(2)	AAC64277	orf8	170	[38%]	GL <b>C</b> V <b>I</b> CLK	17	H <b>V</b> D <b>H</b> CHKTG	VL <b>C</b> F <b>N</b> C <b>S</b> AIKGL	-	GTSWKP		
$\phi$ Bx22	AAN01813	gp59	158	[43%]	GR <b>C</b> Y <b>G</b> CRR	21	S <b>V</b> D <b>H</b> DHETG	LL <b>C</b> T <b>A</b> C <b>N</b> RNVLGH	-	FIDYLD		
$\phi$ Bxb1	NP_075321	gp54	158	[35%]	GR <b>C</b> Y <b>I</b> CRK	21	AV <b>D</b> H <b>D</b> HRTG	LL <b>D</b> TP <b>C</b> NRNVLGH	HL <b>G</b> DD	GIDYLE		
$\phi$ Corndog	NP_817862	gp11	161	[38%]	GK <b>C</b> A <b>I</b> CQI	35	AV <b>D</b> H <b>D</b> HKMA	LL <b>C</b> GP <b>C</b> NQIGIRW	AL <b>V</b> RA	IQAYTR		
$\phi$ Xp10	NP_858983	35L	135	[26%]	GR <b>C</b> A <b>I</b> CQG	21	CL <b>D</b> H <b>N</b> HSTG	VL <b>H</b> RG <b>C</b> NSVLGLV	G <b>V</b> FAR	LAGYLR		
$\phi$ YeO3-12	NP_072072	p47	129	[31%]	G <b>V</b> CA <b>I</b> CGI	22	AV <b>D</b> H <b>C</b> HATG	LL <b>C</b> S <b>A</b> C <b>N</b> IALGKF	-	AITYLR		
<i>S. glossinidius</i>	AAS66876	yajD	113	[30%]	W <b>V</b> CG <b>R</b> CSR	28	H <b>I</b> D <b>H</b> DHSNN	ML <b>C</b> LY <b>C</b> H <b>D</b> HEHEK	-	ASLYGT		
<i>E. coli</i> K12	P19678	yajD	115	-	W <b>V</b> CG <b>R</b> CSR	28	H <b>I</b> D <b>H</b> DHTNN	LL <b>C</b> LY <b>C</b> H <b>D</b> HEHSK	-	ADQYGT		

The motifs A - C represent highly conserved stretches identified by BLAST searches with  $\phi$ T4 gp49 protein as query; conserved functional residues are shown in bold type [357]. the 4 conserved cysteine residues involved in metal ion binding in motifs A + C are indicated by gray shading, and spacing of the motifs shown in a separate column. Motifs C + D contain aromatic residues involved in protein dimerisation. In those cases where the region of similarity did not encompass the entire protein sequences, the identity values are shown in square brackets.

ture-specific endonuclease activity; while orf3 protein shows similarity to *E. coli* RuvC, the contribution of orf2 protein to this activity is still obscure [33]. BLAST searches revealed a group of orf3 homologues among lactococcal phages, and a second group of proteins with significant similarity (Table C36). All proteins share with *E. coli* RuvC the signature motifs A - D and F, but lack motif E\*.

As was pointed out by Sharples and co-workers, *E. coli* RuvC [PDB 1HJR] and  $\phi$ bIL66 orf3 protein share several conserved catalytically important residues in a reasonably related structural environment but differ considerably outside this 'core', therefore simple BLAST searches with  $\phi$ bIL66 orf3 as query fail to detect similarity to *E. coli* RuvC (Table C36) [410]. These authors suppose that the structural differences between the two proteins reflect their different substrate specifi-

cities: while RuvC is highly specific for 'X' structures, i.e. four-way junctions,  $\phi$ bIL66 orf3 protein can cleave branched 'Y' structure, in addition. It should be interesting to include in future studies aimed to clarify the molecular basis of this functional divergence two proteins detected by our BLAST searches: *Fusobacterium nucleatum* RuvC and the protein encoded by locus psyc163101 of *Psychrobacter* sp. (Table C36; third group). Both proteins possess motif E\* but show a reasonable similarity to  $\phi$ bIL66 orf3 protein (~25% ident. res.). However, the protein from *Psychrobacter* sp. is more closely related to *E. coli* RuvC (55% ident. res.) than to  $\phi$ bIL66 orf3 protein, but *F. nucleatum* RuvC is equally distantly related to both.

A fourth group of proteins with significant similarity to  $\phi$ bIL170 p63 is encoded exclusively by lactococcal phages. The signature motifs A - D are well conserved

**Table C35** Phage- and prophage-encoded  $\lambda$  Rap-like proteins.

motifs:				A				B				C				D			
				RCx <sub>2-4</sub> Cx				hCxxxCG				PChSCG				QCxxCN			
consensus	replicon	accession	gene	res.	ident.														
$\lambda$		P03770	Rap / NinG	204	x	RCKNDECR	11	WCSP	CG	66	PCIS	CG	30	QC	VVCN				
$\phi$ 21		Q9XJQ4	NinG	204	98%	RCKNDECR	11	WCSP	CG	66	PCIS	CG	30	QC	VVCN				
$\phi$ Sf6		NP_958228	gene 54	203	98%	RCKNDECR	11	WCSP	CG	66	PCIS	CG	30	QC	VVCN				
$\phi$ ST104		YP_006391	NinG	203	96%	RCKNEECR	11	WCSP	CG	66	PCVSCG	30	QC	VVCN					
$\phi$ P22		NP_059617	NinG	203	94%	RCKNEECR	11	WCSP	CG	66	PCIS	CG	30	QC	VVCN				
p $\phi$ CP-933K		NP_286501	z0953	207	89%	RCKNDECR	11	WCSP	CG	69	PCIS	CG	30	QC	VVCN				
$\phi$ Nil2		CAC95102	NinG	201	91%	KCKI . . CK	11	WCSP	EHG	66	PCIS	CG	30	QC	VVCN				
$\phi$ LC159		AAN59917	orf2	190	91%	KCKI . . CK	11	WCSP	EHG	66	PCIS	CG	30	QC	VVCN				
$\phi$ H-19B		O48427	NinG	201	91%	KCKI . . CK	11	WCSP	EHG	66	PCIS	CG	30	QC	VVCN				
$\phi$ 933W / $\phi$ VT2-Sa		NP_049497	p37	201	92%	KCKI . . CK	11	WC	CP	EHG	66	PCIS	CG	30	QC	VVCN			
$\phi$ 4795		NP_851974	ORF37	201	92%	KCKI . . CK	11	WC	CP	EHG	66	PCIS	CG	30	QC	VVCN			
p $\phi$ CP-933V		NP_288675	z3346	201	90%	KCKI . . CK	11	WC	CP	EHG	66	PCIS	CG	30	QC	VVCN			
p $\phi$ Gifsy-2		NP_459996	stm1021	203	47%	KCANKECR	11	VCSY	QCA	68	GCIS	CG	30	QC	DVYN				
<i>Y. pestis</i> p $\phi$		NP_405642	ypo2090	201	41%	NCKV . . CK	11	WC	DE	EHK	66	ACV	SCG	30	QC	APCN			
<i>H. somnus</i> p $\phi$		ZP_00122672	hsom0668	203	40%	KCKV . . CG	11	WC	SP	ECG	59	PCIS	CG	30	QC	SVCN			
<i>B. cepacia</i> p $\phi$ ?		ZP_00219525	bucepa02005864	190	37%	KCRE . . CG	11	VCS	PV	CA	55	PCIS	CG	28	QC	GPCN			
<i>P. syringae</i> p $\phi$		ZP_00127017	psyr023592	194	38%	TCDNPACG	11	VCG	W	ACG	55	PCV	SCG	30	QC	APCN			
$\phi$ D3		NP_061578	orf82	213	35%	KCQNTFCG	11	VCS	P	ACA	57	GCIS	CG	41	QC	KACN			
$\phi$ Aa $\phi$ 23		NP_852744	NinG	189	34%	KCKS . . CG	11	VCS	P	KCA	60	PCIA	CG	30	GC	IRCN			

The motifs A - D represent the conserved stretches containing the 8 cysteine residues identified by BLAST searches with  $\lambda$  Rap as query; the cysteine residues potentially involved in metal ion binding indicated by gray shading, and spacing of the motifs shown in separate columns.

in these proteins but they lack ~60–110 C-terminal residues as compared to  $\phi$ bIL170 p63 (Table C36). The function of these proteins remains to be studied but it seems doubtful whether they are functional Holliday junction resolvases. Interestingly,  $\phi$ vML3 contains genes for the longer as well as the shorter version of this protein; both proteins possess a truncated motif A at the extreme N-terminus.

*Phage-encoded resolvases with similarity to E. coli RusA.* In *E. coli* *ruvABC* mutants, homologous recombination can be rescued by activation of the *rusA* locus. The *rusA* gene is located in the region of the cryptic prophage DLP12 of *E. coli* K-12, and dimeric RusA has been shown by biochemical analysis to possess Holliday junction resolvase activity (reviewed in [408]). *E. coli* RusA has been crystallised [PDB 1Q8R], and the functionally important residues were determined: D70, D72 and D91 are likely to be involved in forming a magnesium ion-binding pocket involving both subunits of the RusA dimer [39,358]. In 2001, Sharples reported 11 homologues of RusA in the databases, but the number has increased in the meantime to >50 [408] (see Table C37).

We carried out the initial BLAST search for RusA-like proteins using *E. coli* RusA as query; in several subsequent steps, the most distantly related proteins were used as query until no further sequences with reasonable similarity were detected. With the exception of the RusA-like proteins of *A. aeolicus* (locus aq\_1953) and *Bacillus* sp. plasmid pSL32 (see below), all proteins similar to *E. coli* RusA detected by BLAST searches are either phage-encoded or contained within (putative) prophage regions in bacterial chromosomes from a wide range of species. All RusA-like proteins have a comparable size, and reasonably well conserved motifs A + B. In one group however, the spacing within motif A differs by 3 residues; the functional consequence of this different spacing is not known.

The low copy-number plasmid pLS32 (~70 kb) of *Bacillus natto* encodes a protein of unknown function, which is highly similar to *E. coli* RusA (28% ident. res.) despite its longer C-terminus [442]. The N-termini of the putative replication initiator RepN of pLS32 and the ORF4 replication initiator of  $\phi$ 7201 (Section C3.1.2.) share significant sequence similarity (38% ident. res.). The replication origin of pLS32 was mapped within the

**Table C36** Phage- and prophage-encoded RuvC-like proteins.

motifs:				A	B	C	D	E*	F
consensus				LxhDhS	xTGYAHx	xERxx	hxDYx <sub>1-4</sub> hxhExx	xQhKxxh	DDxxDAh
replicon	accession	gene	res. ident.						
φ bIL170	NP_047175	p63	160 x	L <b>A</b> IDFS	NTGYAFR	LERAK	LFDYF..IY <b>I</b> EEP	ISNSMWC	DD <b>M</b> ADAF
φ sk1	NP_044999	ORF53	160 91%	L <b>A</b> IDFS	NTGYAFR	LERAK	LFDYF..IY <b>I</b> EEP	VPNSKWC	DD <b>M</b> ADAF
φ bIL66 / φ bIL67	AAA99046	orf3	160 92%	L <b>A</b> IDFS	NTGYAFR	LERAK	LFDYF..IY <b>I</b> EEP	IPNSKWC	DD <b>Q</b> ADAF
φ 712	AAC63027	M3	161 95%	L <b>A</b> IDFS	NTGYAFR	LERAK	LFDYF..IY <b>I</b> EEP	ISNSMWC	DD <b>Q</b> ADAF
φ c2	NP_043550	gene 12	161 45%	L <b>A</b> IDFS	GTGYAFR	WERTF	LKDYH..MA <b>I</b> ETP	IDNSKWC	DN <b>M</b> ADAY
φ vML3	P13004	-	152 44%	EFPS	GTGYAFR	WERTF	LKDYH..MA <b>I</b> ETP	IDNSKWC	DN <b>M</b> ADAY
φ 315.3	NP_795506	SpyM3_1128	170 24%	L <b>S</b> LDVS	GTGWAIF	YERAK	PFEDI...V <b>I</b> EQN	VNVSTWR	D <b>D</b> EADAI
<i>S. pyogenes</i>	NP_607564	spym18_1491	170 23%	L <b>S</b> LDIS	GTGWALF	FERGR	LKKYY..CA <b>V</b> EKN	INVSTWR	D <b>D</b> EADAI
φ EJ-1	NP_945263	p24	171 22%	L <b>S</b> LDIS	ATGVAVF	FERGR	HFESI...V <b>V</b> EKN	WTWRKYW	D <b>D</b> EADAI
<i>F. nucleatum</i>	Q8RGS0	<i>ruvC</i>	190 24%	I <b>G</b> IDPG	IVGYGII	EERLE	YKPEF..MA <b>I</b> EDL	L <b>Q</b> VKIGI	D <b>D</b> AADAL
<i>Psychrobacter</i> sp.	ZP_00145849	psyc163101	201 23%	I <b>G</b> IDPG	MTGYGIL	PERLK	YA <b>D</b> EPIYTA <b>I</b> EQV	R <b>Q</b> IKQAV	Q <b>D</b> AADGL
<i>E. coli</i>	NP_047175	<i>ruvC</i>	160 -	L <b>G</b> IDPG	VTGYGVI	PSRLK	QP <b>D</b> YF...A <b>I</b> EQV	R <b>Q</b> VKQTV	A <b>D</b> AADAL
φ 923	AAP80756	orf L2	93 [41%]	L <b>A</b> IDFS	GTGYAFR	WERTF	LKDYH..MA <b>I</b> ETP	-	-
φ 5447 / φ 5469	AAN05716	L2	85 [41%]	L <b>A</b> IDFS	GTGYAFR	WERTF	LKDYH..MA <b>I</b> ETP	-	-
φ c6A	AAN05724	L2	84 [41%]	L <b>A</b> IDFS	GTGYAFR	WERTF	LKDYH..L <b>A</b> IETP	-	-
φ 5440	AAN05712	L2	84 [40%]	L <b>A</b> IDFS	GTGYAFR	WERTF	LKDYH..MA <b>I</b> ETP	-	-
φ 943	AAN05707	L2	84 [40%]	L <b>A</b> IDFS	GTGYAFR	WERTF	LKDYH..MA <b>I</b> ETP	-	-
φ vML3	S05334	-	97 [38%]	EFPS	GTGYAFR	WERTF	LKDYH..MA <b>I</b> ETP	-	-

Motif B + C represent highly conserved stretches identified by BLAST searches with φbIL170 p63 protein as query. Motifs A, D, E\* + F correspond to the RuvC signature motifs I - IV defined by Lilley + White [251]; the catalytically important residues are shown in bold type. The similarity values for protein sequences matching with the N-terminus of φbIL170 p63 only are shown in square brackets.

repN gene, and the plasmid may therefore represent a naturally occurring phage-derived plasmid replicon [442]. Alternatively – and in (weak) analogy to *E. coli* φP1 – pLS32 could be the prophage state of a yet unknown *Bacillus* sp. phage. It should be highly interesting to analyse the (putative) mechanism for copy number control in pLS32 because phage origins of the type apparently operating in this plasmid usually give rise to much higher copy numbers as have been determined for pSL32 (2-3 per chromosome)[442]. Other plasmid-encoded RusA-like proteins – or other Holliday junction resolvases – were not detected in extensive BLAST searches.

The discussion of phage-encoded Holliday junction resolvases shows that genes encoding such proteins are widespread among phages. Presently, five groups of resolvases are known, which share no structural features and are evolutionary only distantly related, if at all (but see [251]). This argues against the possibility that more and novel resolvases will be simply detectable by *in silico* analyses. Except for the Rap resolvase, which is only found in the genomes of phages and prophages

with significant similarity to λ, the resolvases possibly spread among the various phage families by HGT, but there remain gaps. The known resolvase genes do not cover the entire spectrum of known phages: resolvase genes have not yet identified φSPP1 (but in φA118), φPaP3 and φFelix01 (but in φT7).

**Other phage-encoded recombination proteins.** We conclude the discussion of recombination proteins with a brief survey of phage-encoded homologues of the 'canonical' *E. coli* recombination proteins [230,61]. Several of them are probably not directly involved in chromosome replication, but we were curious to learn more about the possible mutual exchange of recombination protein-encoding genes between chromosomes and phage replicons.

*RecA.* Among the most highly conserved bacterial proteins, RecA has to date only two detectable homologues in phage replicons [267]: the gp117 protein of Mycobacteriophage φCJW1 (≥60% ident. res. with RecA proteins of several *Mycobacterium* sp.), and the gp201

**Table C37** Phage- and prophage-encoded RusA-like proteins.

motifs:				A		B	
consensus				xxDxDX <sub>3/6</sub>		KxxxD xhhxDDxxhx	
replicon	accession	gene	res.	ident.			
pφ DLP12 <i>E. coli</i> K12	P40116	<i>rusA</i>	120	x	RRDLNLQ . . . KAAFD	GFWLDDAQVV	
pφ CP-933X <i>E. coli</i>	NP_287358	<i>z1873</i>	120	95%	RRDLNLQ . . . KAAFD	GFWLDDAQVV	
φ 82	Q37873	<i>rusA</i>	120	95%	RRDLNLQ . . . KASFD	GFWLDDAQVV	
φ ST64T	NP_720315	<i>Rus</i>	131	53%	KRDLNLP . . . KAVFD	GFWLDDGQID	
φ HK620	NP_112065	<i>rus</i>	120	43%	RRDLNLL . . . KGLLD	GFAEDDEQFD	
φ HK97	Q9MCN8	<i>gp67</i>	120	43%	RRDLNLL . . . KGLLD	GFAEDDEQFD	
<i>E. coli</i> pφ CP-933U	NP_288470	<i>z3115</i>	119	41%	RRDLNLL . . . KAPLD	GLLMDDQFD	
<i>E. coli</i> pφ CP-933P	NP_287995	<i>z6061</i>	119	40%	RRDLNLL . . . KAPLD	GLLIDDEQFD	
<i>E. coli</i> pφ CP-933O	NP_287519	<i>z2057</i>	119	39%	RRDLNLL . . . KAPLD	GVLMDDEQFD	
<i>E. coli</i> pφ CP-933M	NP_286866	<i>z1344</i>	124	41%	RRDLNLL . . . KAPLD	GLLMDDQFD	
<i>E. coli</i> pφ CP-933N	NP_282777	<i>z1785</i>	123	40%	RRDLNLL . . . KAPLD	EVLIDDEQFD	
<i>E. coli</i> pφ CP-933O	NP_287559	<i>z2101</i>	101	42%	RRDLNLL . . . KAPLD	GLLIDDEQFD	
φ SfV	NP_599075	<i>orf43</i>	129	38%	IRDLDNYN . . . KALFD	GVWEDDSQVK	
φ ST64B	NP_700418	<i>sb45</i>	129	41%	IRDLDNYN . . . KALFD	GVWEDDSQVK	
<i>L. pneumophila</i> pφ	CAC33459	<i>ORF G</i>	120	31%	RRDIDNLL . . . KCLLD	GVYDDDNQID	
φ Bcep1	NP_944331	<i>gp23</i>	131	41%	AWDVANRE . . . KLLSD	GFWRDDSQID	
φ Bcep781	NP_705648	<i>ORF20</i>	136	33%	AWDVANRE . . . KLICD	GFWRDDSQID	
φ Bcep43	NP_958127	<i>gp23</i>	128	33%	AWDVANRE . . . KLICD	GFWRDDSQID	
pLS32	BAA24869	-	185	28%	RKDPNNFL . . . KVPFD	GVYLDDVAL	
<i>A. aeolicus</i>	NP_214335	<i>aq_1953</i>	152	26%	KRDIDNML . . . KSLWD	GVIKNDNLIF	
<i>S. pyogenes</i> M1 GAS pφ	NP_268913	<i>spy0673</i>	131	[64%]	KPDTNLQ . . . KLLKD	GFWNDDAQVA	
φ TM4	NP_569805	<i>gp72</i>	213	[46%]	RPDLKLA . . . RAILD	VVWIDDSQVV	
φ NIH1.1	NP_438128	<i>15</i>	117	[34%]	KPDLNLE . . . KAVYD	VVWTDNIIIV	
φ 315.4	NP_665052	<i>SpyM3_1248</i>	146	[34%]	KPDLNLE . . . KAVYD	VVWTDNIIIV	
φ r1t	NP_695042	<i>orf14</i>	131	[42%]	RPDLNLM . . . KNLQD	RYYSDDSQIV	
φ bIL285 / φ LC3	NP_076590	<i>rusA</i>	129	[42%]	RPDLNLM . . . KNLQD	RYYSDDSQIV	
φ Tuc2009 / φ ul36	NP_108696	<i>orf19</i>	129	[42%]	RPDLNLM . . . KNLQD	RYYSDDSQIV	
φ BK5-T	NP_116543	<i>ORF51</i>	133	[42%]	RPDLNLM . . . KNLQD	RYYSDDSQIV	
<i>S. agalactiae</i> pφ	NP_687598	<i>SAG0569</i>	158	[33%]	KPDIDNLI . . . KAVFD	IVWSDDNIVC	
φ 7201	NP_038311	<i>orf10</i>	153	[33%]	KPDIDNLI . . . KALFD	IVWTDDNIVC	
<i>R. eutropha</i> pφ	ZP_00168225	<i>Raeut330401</i>	157	[40%]	KPDADNVL . . . KAVKD	IVWVDDAQAV	
φ PVL / <i>L. innocua</i> pφ	NP_058487	<i>orf 48</i>	145	23%	RPDIDNYV . . . KAILD	IMFSDDGKIV	
φ LL-H	AAL77552	<i>orf139</i>	139	25%	KPDTDNIE . . . KGIYD	LAYEDDKQIV	
<i>C. tetani</i> pφ	NP_782962	<i>ctc02433</i>	222	27%	HPDTINIT . . . KSIFD	GVIINDAQIR	
φ Bxz1	NP_818256	<i>gp205</i>	151	-	KPDLKLI . . . RAIGD	TVYTEDSRIV	
<i>N. aromatorans</i> pφ	ZP_00094972	<i>saro1995</i>	135	-	KPDGNIL . . . KALGD	IVWADDSQVA	
<i>B. cepacia</i> pφ	ZP_00218864	<i>Bucepa02006569</i>	158	-	KPDADNVV . . . KALKD	VVYGDDGQVV	
φ N315	NP_835540	<i>sa1789</i>	134	-	KPDIDNLI . . . KTVLD	HVWKDDNQIT	
<i>B. subtilis</i> pφ	P45911	<i>yqaN</i>	142	-	KPDVDNYV . . . KGVKD	LIYKDDSQVV	
<i>E. faecalis</i> pφ	NP_815791	<i>ef2128</i>	141	-	KPDLNLYF . . . KAVTD	LYKNDGQIAV	
φ 31.1	AAF43122	<i>orf139a</i>	139	29%	RSDPDNLQPTLKAIMD	GLWSDDNHEV	
φ ul36.1	AAF74065	<i>orf109</i>	109	[40%]	RSDPDNLQPTLKAIMD	GLWSDDNHEV	
φ TP901-1	NP_112678	<i>rus</i>	139	29%	RSDPDNLQPTLKAIMD	GLWSDDNHEV	
<i>X. fastidiosa</i> pφ	ZP_00038274	<i>xfasa0197</i>	153	29%	LPDDNMLARFKPYRD	ALGIDRRFV	
φ A118	NP_463528	<i>ORF63</i>	134	-	KKDPDNIAFAKFI	GFLENDNLNY	
φ MM1	NP_150151	<i>MM1p20</i>	133	[30%]	KLDPNLYPTVKAID	GIWTDNHNKV	
φ Che9c	NP_817741	<i>gp64</i>	127	-	RRDNHNLWPTVKALVD	GVVPDDTEH	
φ bIL309	NP_076712	<i>orf17</i>	136	-	KYDPPNYEPTSKALID	GIWNDDNLYNV	

For the prophages of the CP-933 series, only the genes of *E. coli* strain O157:H7 EDL933 but not for *S. flexneri* are listed.

protein of Mycobacteriophage  $\phi$ Bxz1. Gp117 can serve as example for recent horizontal gene transfer (HGT), i.e. the acquisition of a host gene by a phage replicon [301]. The origin of gp201 is more difficult to trace because it shows only ~30% similarity (ident. res.) to gp117 as well as to various bacterial RecA proteins of  $\phi$ Bxz1 hosts. The significant divergence of the gp201 sequence from other RecA proteins suggests that it was acquired by  $\phi$ Bxz1 much earlier than gp117 by  $\phi$ CJW1. Whether CJW1 gp117 and  $\phi$ Bxz1 gp201 still function as SAPs has not yet been analysed.

*Exonuclease V* (recB, recC, recD). No phage-encoded proteins with significant similarity to the *E. coli* RecB and RecC proteins were detected by BLAST searches, and also the RecD protein shares only some limited similarity with the Dda proteins of the  $\phi$ T4-group of phages in the helicase domain (SF1-type) [419] (Section C3.3.).

*RecJ*. The YorK protein of *B. subtilis*  $\phi$ SPBc2 is the only phage-encoded protein with similarity (23% ident. res.; full-length) to *E. coli* RecJ. The *york* gene is located adjacent to the *yorkL* gene encoding the DNA polymerase of  $\phi$ SPBc2; it has not been studied whether the putative 5'→3' exonuclease activity of the YorK protein is required for  $\phi$ SPBc2 replication, recombination, or both.

*RecQ*. Homologues of the *E. coli* RecQ helicase (superfamily II) are found in the genomes of bacteria from all phyla, and also in the genomes of archaea and eukarya. The only viral RecQ homologue (cac19131) is encoded by the fungal virus DpAV4, and was probably acquired from its host. The various phage-encoded SF2 helicases (Section C3.3.) may be functional analogues of RecQ but significant sequence similarity to *E. coli* RecQ could not be detected.

*RdgC*. Genes encoding homologues of the nucleoid-associated protein RdgC are only found in proteobacterial genomes. This protein modulates RecA binding to ssDNA in concert with SSB and the RecFOR recombination proteins, and has been characterised only recently [380,300]. Phage-encoded homologues of RdgC are the HrdC protein of  $\phi$ P1 (61% ident. res.), ORF10 protein of *Vibrio harveyi* phage  $\phi$ VHML (41% ident. res.), and probably p45 protein of *Yersinia* sp. phage  $\phi$ PY54 (28% ident. res. in the N-terminal half).

*SbcC, SbcD*. The *E. coli* SbcB and SbcD proteins form a heterodimer with a complex substrate-specific exo- and endonuclease activity, and with a particular role in the removal of hairpin DNA [83]. SbcC and SbcD homologues are well conserved in bacterial ge-

nomes from most phyla, and found also in archaeal and eukaryal genomes. The only detectable phage-encoded proteins with significant similarity to SbcC/ SbcD are the gp46/gp47 proteins of phages of the  $\phi$ T4-group (see above).

*RuvA, RuvB, and RuvC*. No phage-encoded homologues of the RuvA and RuvB proteins were found by BLAST searches. The Holliday junction resolvase encoded by *L. lactis* phage  $\phi$ IL66 shares the catalytic core residues with RuvC [33]; homologues of  $\phi$ IL66 orf3 protein are found in several other lactococcal phages and prophages of *Streptococcus* sp. (see above).

BLAST searches did not reveal proteins with significant similarity to the *E. coli* RecF, RecO, RecR, RecG, RecN, exonuclease I (*sbcB*), helicase II (*uvrD/recL*), and helicase IV (*helD*) recombination proteins. Although the functional equivalence or similarity of several recombination pathways in bacteria and phages is well documented, this survey suggests that the genes/proteins responsible for individual recombination pathways were kept largely separate during the long-lasting co-evolution of phages and their hosts. Clearly, the acquisition of bacterial genes encoding recombination proteins by phage replicons is the exception rather than the rule. Several prophage-encoded recombination proteins present in bacterial genomes were shown to constitute backup systems for mutated host recombination genes, but some require activation by mutation(s) in order to substitute for defective chromosomally encoded recombination proteins. There is thus far only one example where a phage recombination protein might have replaced a chromosomal recombination protein: the genome of *A. aeolicus* lacks detectable equivalents of the *E. coli* RuvABC proteins but encodes a homologue of the phage Holliday junction resolvase RusA [408].

#### C4. Concluding remarks

The still rapidly increasing number of completely sequenced genomes makes it impossible to predict a 'half-life' for the numerous figures and tables presented in this compendium. For quite a number of prophage genes present in bacterial chromosomes, however, we were able to deduce their function by (BLAST) sequence comparison combined with a search for signature motifs derived from homologues with known structure. In this respect, the compendium of phage replication origins and proteins may serve as a reliable starting-tool for function assignment for newly sequenced genomes at least for some time. On longer terms, it would be preferable to combine the knowledge-based approach presented here with automated approaches like the COG

(Clusters of Orthologous Groups of proteins) [444] or the CDD (Conserved Domains Database) [273] databases in an continuously updated online search tool.

In the introduction to the review (BRM Section 1.), we reason why a satisfactory definition of 'bacteriophages' is problematic. Also we mention that replication research has traditionally used phages as *experimental systems* rather than discussing their exact classification. Based on the results presented in this compendium, this somewhat sloppy approach seems largely justified: despite the enormous variability of phage-encoded replication functions, we rarely found plasmid- or chromosome-encoded homologues of phage replication proteins (except for prophages). There are only *very few* examples for gene exchange (HGT) between phage replicons and chromosomes. One example is the acquisition of the host *dnaB* gene encoding the replicative helicase by a number of lambdoid phages. Another example may be the (hypothetical) acquisition of *dnaC*-type helicase loader genes of lambdoid phages by several  $\gamma$ -proteobacteria and by the Bacillales. These issues are discussed in detail in BRM Sections 4.2. & 4.3., respectively. The activation of recombination functions encoded by resident prophages – e.g. *rusA*, *recET* – may serve as a third example, but in these cases one cannot speak of HGT in the strict sense.

There are also *very few* examples for an exchange of replication functions between phage and plasmid replicons. Clearly, plasmids and phages replicating by RCR share several features and their replication modules may have a common evolutionary origin. Phages however

lack the intricate copy-number control mechanisms operating in RCR plasmids. Also, a closer inspection of the initiator genes encoded by RCR plasmids and phages reveals conserved structural motifs that are – with few exceptions – confined to either the phage-encoded or the plasmid-encoded proteins (Tables C8 + C9).

Coliphages P1 and P4 are among the best-studied naturally occurring joint replicons, but homology searches suggest that both are solitaires in the phage world. Plasmids pXO1, pTAV320, and pTA1060 encode (putative) replication proteins sharing reasonable similarity in their N-termini with several  $\lambda$  O-type initiator proteins. Since this N-terminal DNA-binding domain contains a HTH-motif, which is also present in a number of transcription factors with N-terminal DNA-binding domains, it would require a detailed analysis to determine whether these plasmid genes and the initiator genes of lambdoid phages have a common origin. The replication of two naturally occurring linear plasmids of *Streptomyces rochei* may be driven by phage-derived replication modules (Section C2.2.). Both encode  $\lambda$  O-type initiators with the replication origin contained within the initiator gene [71]. Both plasmids encode, in addition, a helicase of the DnaB<sub>Eco</sub>-type in an arrangement typical for the 'IH-type' replication module (BRM Section 3.1.3.). Another example for a phage replication module driving plasmid replication may be the low copy-number plasmid pLS32 of *Bacillus natto* (Section C3.6.2.). Keeping these few exceptions in mind, we conclude that – with respect to replication – the term 'bacteriophage' can be used safely also in the future.

## References

- [1] Ajdic' D, McShan WM, McLaughlin RE, *et al.* (2002) Genome sequence of *Streptococcus mutans* UA159, a cariogenic dental pathogen. *Proc Natl Acad Sci USA* **99**: 14434-14439.
- [2] Alano P, Dehò G, Sironi G & Zangrossi S (1986) Regulation of the plasmid state of the genetic element P4. *Mol Gen Genet* **203**: 445-450.
- [3] Alley SC, Trakselis MA, Mayer MU, Ishmael FT, Jones AD & Benkovic SJ (2001) Building a replisome solution structure by elucidation of protein-protein interactions in the bacteriophage T4 DNA polymerase holoenzyme. *J Biol Chem* **276**: 39340-39349.
- [4] Altermann E, Klein JR & Henrich B (1999) Primary structure and features of the genome of the *Lactobacillus gasseri* temperate bacteriophage (phi)adh. *Gene* **236**: 333-346.
- [5] Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W & Lipman DJ (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* **25**: 3389-3402.
- [6] Arai K & Kornberg A (1979) A general priming system employing only dnaB protein and primase for DNA replication. *Proc Natl Acad Sci USA* **76**: 4308-4312.
- [7] Arai K & Kornberg A (1981) Mechanism of dnaB protein action. IV. General priming of DNA replication by dnaB protein and primase compared with RNA polymerase. *J Biol Chem* **256**: 5267-5272.
- [8] Arai K, Yasuda S & Kornberg A (1981) Mechanism of dnaB protein action. I. Crystallization and properties of dnaB protein, an essential replication protein in *Escherichia coli*. *J Biol Chem* **256**: 5247-5252.
- [9] Aravind L, Leipe DD & Koonin EV (1998) Toprim - a conserved catalytic domain in type IA and II topoisomerases, DnaG-type primases, OLD family nucleases and RecR proteins. *Nucleic Acids Res* **26**: 4205-4213.
- [10] Asano S, Higashitani A & Horiuchi K (1999) Filamentous phage replication initiator protein gpII forms a covalent complex with the 5' end of the nick it introduced. *Nucleic Acids Res* **27**: 1882-1889.
- [11] Ayora S, Langer U & Alonso JC (1998) *Bacillus subtilis* DnaG primase stabilises the bacteriophage SPP1 G40P helicase-ssDNA complex. *FEBS Lett* **439**: 59-62.
- [12] Ayora S, Missich R, Mesa P, Lurz R, Yang S, Egelman EH & Alonso JC (2002) Homologous-pairing activity of the *Bacillus subtilis* bacteriophage SPP1 replication protein G35P. *J Biol Chem* **277**: 35969-35979.

- [13] Ayora S, Stasiak A & Alonso JC (1999) The *Bacillus subtilis* bacteriophage SPP1 G39P delivers and activates the G40P DNA helicase upon interacting with the G38P-bound replication origin. *J Mol Biol* **288**: 71-85.
- [14] Baas PD & Jansz HS (1976) Bacteriophage phiX174 DNA synthesis in a replication-deficient host: determination of the origin of phiX DNA replication. *J Mol Biol* **102**: 633-656.
- [15] Backhaus H & Petri JB (1984) Sequence analysis of a region from the early right operon in phage P22 including the replication genes 18 and 12. *Gene* **32**: 289-303.
- [16] Bailey S, Sedelnikova SE, Mesa P, Ayora S, Alonso JC & Rafferty JB (2003) Crystallization of the *Bacillus subtilis* SPP1 bacteriophage helicase loader protein G39P. *Acta Crystallogr D Biol Crystallogr* **59**: 1090-1092.
- [17] Bailey S, Sedelnikova SE, Mesa P, *et al.* (2003) Structural analysis of *Bacillus subtilis* SPP1 phage helicase loader protein G39P. *J Biol Chem* **278**: 15304-15312.
- [18] Baker TA & Bell SP (1998) Polymerases and the replisome: machines within machines. *Cell* **92**: 295-305.
- [19] Bárcena M, Martin CS, Weise F, Ayora S, Alonso JC & Carazo JM (1998) Polymorphic quaternary organization of the *Bacillus subtilis* bacteriophage SPP1 replicative helicase (G40 P). *J Mol Biol* **283**: 809-819.
- [20] Bárcena M, Ruiz T, Donate LE, Brown SE, Dixon NE, Radermacher M & Carazo JM (2001) The DnaB.DnaC complex: a structure based on dimers assembled around an occluded channel. *EMBO J* **20**: 1462-1468.
- [21] Barnes MH, Spacciapoli P, Li DH & Brown NC (1995) The 3'-5' exonuclease site of DNA polymerase III from gram-positive bacteria: definition of a novel motif structure. *Gene* **165**: 45-50.
- [22] Barrett KJ, Gibbs W & Calendar R (1972) A transcribing activity induced by satellite phage P4. *Proc Natl Acad Sci USA* **69**: 2986-2990.
- [23] Bates DB, Asai T, Cao Y, Chambers MW, Cadwell GW, Boye E & Kogoma T (1995) The DnaA box R4 in the minimal *oriC* is dispensable for initiation of *Escherichia coli* chromosome replication. *Nucleic Acids Res* **23**: 3119-3125.
- [24] Bazill GW & Gross JD (1973) Mutagenic DNA polymerase in *B. subtilis*. *Nat New Biol* **243**: 241-243.
- [25] Beese LS & Steitz TA (1991) Structural basis for the 3'-5' exonuclease activity of *Escherichia coli* DNA polymerase I: a two metal ion mechanism. *EMBO J* **10**: 25-33.
- [26] Bendtsen JD, Nilsson AS & Lehnher H (2002) Phylogenetic and functional analysis of the bacteriophage P1 single-stranded DNA-binding protein. *J Virol* **76**: 9695-9701.
- [27] Benkovic SJ, Valentine AM & Salinas F (2001) Replisome-mediated DNA Replication. *Annu Rev Biochem* **70**: 181-208.
- [28] Bernad A, Blanco L, Lazaro JM, Martin G & Salas M (1989) A conserved 3'-5' exonuclease active site in prokaryotic and eukaryotic DNA polymerases. *Cell* **59**: 219-228.
- [29] Bernstein JA & Richardson CC (1988) A 7-kDa region of the bacteriophage T7 gene 4 protein is required for primase but not for helicase activity. *Proc Natl Acad Sci USA* **85**: 396-400.
- [30] Bessman M, Kornberg A, Lehman I & Simms E (1956) Enzymic synthesis of deoxyribonucleic acid. *Biochim Biophys Acta* **21**: 197-198.
- [31] Bhagwat M, Meara D & Nossal NG (1997) Identification of Residues of T4 RNase H Required for Catalysis and DNA Binding. *J Biol Chem* **272**: 28531-28538.
- [32] Bhattacharyya S & Griep MA (2000) DnaB helicase affects the initiation specificity of *Escherichia coli* primase on single-stranded DNA templates. *Biochemistry* **39**: 745-752.
- [33] Bidnenko E, Ehrlich SD & Chopin MC (1998) *Lactococcus lactis* phage operon coding for an endonuclease homologous to RuvC. *Mol Microbiol* **28**: 823-834.
- [34] Bird LE, Pan H, Soultanas P & Wigley, DB (2000) Mapping protein-protein interactions within a stable complex of DNA primase and DnaB helicase from *Bacillus stearothermophilus*. *Biochemistry* **39**: 171-182.
- [35] Birkenbihl RP & Kemper B (1998) Endonuclease VII has two DNA-binding sites each composed from one N- and one C-terminus provided by different subunits of the protein dimer. *EMBO J* **17**: 4527-4534.
- [36] Birkenbihl RP & Kemper B (1998) Localization and characterization of the dimerization domain of holliday structure resolving endonuclease VII of phage T4. *J Mol Biol* **280**: 73-83.
- [37] Biswas EE & Biswas SB (1999) Mechanism of DnaB helicase of *Escherichia coli*: Structural domains involved in ATP hydrolysis, DNA binding, and oligomerization. *Biochemistry* **38**: 10919-10928.
- [38] Blanco L, Bernad A, Blasco MA & Salas M (1991) A general structure for DNA-dependent DNA polymerases. *Gene* **100**: 27-38.
- [39] Bolt EL, Sharples GJ & Lloyd RG (1999) Identification of three aspartic acid residues essential for catalysis by the Rusa holliday junction resolvase. *J Mol Biol* **286**: 403-415.
- [40] Bonner CA, Randall SK, Rayssiguier C, Radman M, Eritja R, Kaplan BE, McEntee K & Goodman MF (1988) Purification and characterization of an inducible *Escherichia coli* DNA polymerase capable of insertion and bypass at abasic lesions in DNA. *J Biol Chem* **263**: 18946-18952.
- [41] Bouché JP, Rowen L & Kornberg A (1978) The RNA primer synthesized by primase to initiate phage G4 DNA replication. *J Biol Chem* **253**: 765-769.
- [42] Bouché JP, Zechel K & Kornberg A (1975) dnaG gene product, a rifampicin-resistant RNA polymerase, initiates the conversion of a single-stranded coliphage DNA to its duplex replicative form. *J Biol Chem* **250**: 5995-6001.
- [43] Bowden DW, Twersky RS & Calendar R (1975) *Escherichia coli* deoxyribonucleic acid synthesis mutants: their effect upon bacteriophage P2 and satellite bacteriophage P4 deoxyribonucleic acid synthesis. *J Bacteriol* **124**: 167-175.
- [44] Braithwaite DK & Ito J (1993) Compilation, alignment, and phylogenetic relationships of DNA polymerases. *Nucleic Acids Res* **21**: 787-802.
- [45] Brayer GD & McPherson A (1985) A model for intracellular complexation between gene-5 protein and bacteriophage fd DNA. *Eur J Biochem* **150**: 287-296.
- [46] Breyer WA & Matthews BW (2000) Structure of *Escherichia coli* exonuclease I suggests how processivity is achieved. *Nat Struct Biol* **7**: 1125-1128.
- [47] Briani F, Dehò G, Forti F & Ghisotti D (2001) The plasmid status of satellite bacteriophage P4. *Plasmid* **45**: 1-17.
- [48] Broker TR (1973) An electron microscopic analysis of pathways for bacteriophage T4 DNA recombination. *J Mol Biol* **81**: 1-16.
- [49] Brooks K & Clark AJ (1967) Behavior of lambda bacteriophage in a recombination deficient strain of *Escherichia coli*. *J Virol* **1**: 283-293.
- [50] Brown DR, Roth MJ, Reinberg D & Hurwitz J (1984) Analysis of bacteriophage phi X174 gene A protein-mediated termination and reinitiation of phi X DNA synthesis. I. Characterization of the termination and reinitiation reactions. *J Biol Chem* **259**: 10545-10555.
- [51] Bruand C & Ehrlich SD (2000) UvrD-dependent replication of rolling-circle plasmids in *Escherichia coli*. *Mol Microbiol* **35**: 204-210.
- [52] Bruand C, Velten M, McGovern S, Marsin S, Serena C, Ehrlich SD & Polard P (2005) Functional interplay between the *Bacillus subtilis* DnaD and DnaB proteins essential for initiation and re-initiation of DNA replication. *Mol Microbiol* **55**: 1138-1150.



- [53] Bruand C, Velten M, McGovern S, Marsin S, Serena C, Ehrlich SD & Polard P (2005) Functional interplay between the *Bacillus subtilis* DnaD and DnaB proteins essential for initiation and re-initiation of DNA replication. *Mol Microbiol* **55**: 1138-1150.
- [54] Bujalowski W, Klonowska MM & Jezewska MJ (1994) Oligomeric structure of *Escherichia coli* primary replicative helicase DnaB protein. *J Biol Chem* **269**: 31350-31358.
- [55] Burger KJ & Trautner TA. (1978) Specific labelling of replicating SPP1 DNA: analysis of viral DNA synthesis and identification of phage DNA-genes. *Mol Gen Genet* **166**: 277-285.
- [56] Burke RL, Alberts BM & Hosoda J (1980) Proteolytic removal of the COOH terminus of the T4 gene 32 helix-destabilizing protein alters the T4 in vitro replication complex. *J Biol Chem* **255**: 11484-11493.
- [57] Burke RL, Munn M, Barry J & Alberts BM (1985) Purification and properties of the bacteriophage T4 gene 61 RNA priming protein. *J Biol Chem* **260**: 1711-1722.
- [58] Burz DS, Beckett D, Benson N & Ackers GK (1994) Self-assembly of bacteriophage lambda cI repressor: effects of single-site mutations on the monomer-dimer equilibrium. *Biochemistry* **33**: 8399-8405.
- [59] Calendar R, Lindqvist B, Sironi G & Clark A.J (1970) Characterization of REP- mutants and their interaction with P2 phage. *Virology* **40**: 72-83.
- [60] Callanan MJ, O'Toole PW, Lubbers MW & Polzin KM (2001) Examination of lactococcal bacteriophage c2 DNA replication using two-dimensional agarose gel electrophoresis. *Gene* **278**: 101-106.
- [61] Camerini-Otero RD & Hsieh P (1995) Homologous recombination proteins in prokaryotes and eukaryotes. *Annu Rev Genet* **29**: 509-552.
- [62] Campbell A & Botstein D (1983) Evolution of the Lambdoid Phages. In: *Lambda II* (Hendrix RW, Roberts JW, Stahl FW & Weisberg RA., Eds.), pp. 365-380. Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, USA.
- [63] Canchaya C, Proux C, Fournous G, Bruttin A & Bruessow H (2003) Prophage genomics. *Microbiol Mol Biol Rev* **67**: 238-276.
- [64] Carter DM & Radding CM (1971) The role of exonuclease and beta protein of phage lambda in genetic recombination. II. Substrate specificity and the mode of action of lambda exonuclease. *J Biol Chem* **246**: 2502-2512.
- [65] Ceska TA., Sayers JR, Stier G & Suck D (1996) A helical arch allowing single-stranded DNA to thread through T5 5'-exonuclease. *Nature* **382**: 90-93.
- [66] Cha TA & Alberts BM (1990) Effects of the bacteriophage T4 gene 41 and gene 32 proteins on RNA primer synthesis: coupling of leading- and lag-ging-strand DNA synthesis at a replication fork. *Biochemistry* **29**: 1791-1798.
- [67] Chaconas G, Harshey RM, Sarvetnick N & Bukhari A.I (1981) Mechanism of bacteriophage Mu DNA transposition. *Cold Spring Harbor Symp Quant Biol* **45**(1): 311-322.
- [68] Chandry, PS, Moore SC, Boyce JD, Davidson BE & Hillier AJ (1997) Analysis of the DNA sequence, gene expression, origin of replication and modular structure of the *Lactococcus lactis* lytic bacteriophage sk1. *Mol Microbiol* **26**: 49-64.
- [69] Chang P & Marians KJ (2000) Identification of a region of *Escherichia coli* DnaB required for functional interaction with DnaG at the replication fork. *J Biol Chem* **275**: 26187-26195.
- [70] Chang PC & Cohen SN (1994) Bidirectional replication from an internal origin in a linear streptomyces plasmid. *Science* **265**: 952-954.
- [71] Chang PC, Kim ES & Cohen SN (1996) Streptomyces linear plasmids that contain a phage-like, centrally located, replication origin. *Mol Microbiol* **22**: 789-800.
- [72] Chase JW, L'Italien JJ, Murphy JB, Spicer EK & Williams KR (1984) Characterization of the *Escherichia coli* SSB-113 mutant single-stranded DNA-binding protein. Cloning of the gene, DNA and protein sequence analysis, high pressure liquid chromatography peptide mapping and DNA-binding studies. *J Biol Chem* **259**: 805-814.
- [73] Chatterjee DK, Fujimura RK, Campbell JH & Gerard GF (1991) Cloning and overexpression of the gene encoding bacteriophage T5 DNA polymerase. *Gene* **97**: 13-19.
- [74] Chatteraj DK (2000) Control of plasmid DNA replication by iterons: no longer paradoxical. *Mol Microbiol* **37**: 467-476.
- [75] Chen F & Lu J (2002) Genomic Sequence and Evolution of Marine Cyanophage P60: a New Insight on Lytic and Lysogenic Phages. *Appl Environ Microbiol* **68**: 2589-2594.
- [76] Chesney RH, Scott JR & Vapnek D (1979) Integration of the plasmid prophages P1 and P7 into the chromosome of *E. coli*. *J Mol Biol* **130**: 161-173.
- [77] Chopin MC, Rouault A, Ehrlich SD & Gautier M (2002) Filamentous phage active on the gram-positive bacterium *Propionibacterium freudenreichii*. *J Bacteriol* **184**: 2030-2033.
- [78] Chowdhury K, Tabor S & Richardson CC (2000) A unique loop in the DNA-binding crevice of bacteriophage T7 DNA polymerase influences primer utilization. *Proc Natl Acad Sci USA* **97**: 12469-12474.
- [79] Chu CC, Templin A & Clark AJ (1989) Suppression of a frameshift mutation in the *recE* gene of *Escherichia coli* K-12 occurs by gene fusion. *J Bacteriol* **171**: 2101-2109.
- [80] Clark AJ, Sharma V, Brenowitz S, et al. (1993) Genetic and molecular analyses of the C-terminal region of the *recE* gene from the Rac prophage of *Escherichia coli* K-12 reveal the *recT* gene. *J Bacteriol* **175**: 7673-7682.
- [81] Cohen G (1983) Electron microscopy study of early lytic replication forms of bacteriophage P1 DNA. *Virology* **131**: 159-170.
- [82] Cohen G & Sternberg N (1989) Genetic analysis of the lytic replicon of bacteriophage P1.I. Isolation and partial characterization. *J Mol Biol* **207**: 99-109.
- [83] Connelly JC, de Leau ES & Leach DR (1999) DNA cleavage and degradation by the SbcCD protein complex from *Escherichia coli*. *Nucleic Acids Res* **27**: 1039-1046.
- [84] Corpet F (1988) Multiple sequence alignment with hierarchical clustering. *Nucleic Acids Res* **16**: 10881-10890.
- [85] Cox MM, Goodman MF, Kreuzer KN, Sherratt DJ, Sandler SJ & Marians KJ (2000) The importance of repairing stalled replication forks. *Nature* **404**: 37-41.
- [86] Crooke E, Thresher R, Hwang DS, Griffith J & Kornberg A (1993) Replicatively active complexes of DnaA protein and the *Escherichia coli* chromosomal origin observed in the electron microscope. *J Mol Biol* **233**: 16-24.
- [87] Cuff JA & Barton GJ (2000) Application of Multiple Sequence Alignment Profiles to Improve Protein Secondary Structure Prediction. *Proteins* **40**: 502-511.
- [88] Cuff JA, Clamp ME, Siddiqui AS, Finlay M & Barton GJ (1998) JPred: a consensus secondary structure prediction server. *Bioinformatics* **14**: 892-893.
- [89] D'Ari R, Jaffé-Brachet A, Touati-Schwartz D & Yarmolinsky MB (1975) A dnaB analog specified by bacteriophage P1. *J Mol Biol* **94**: 341-366.
- [90] Dalrymple BP, Kongsuwan K, Wijffels G, Dixon NE & Jennings PA (2001) A universal protein-protein interaction motif in the eubacterial DNA replication and repair systems. *Proc Natl Acad Sci USA* **98**: 11627-11632.
- [91] Datta HJ, Khatri GS & Bastia D (1999) Mechanism of recruitment of DnaB helicase to the replication origin of the plasmid pSC101. *Proc Natl Acad Sci USA* **96**: 73-78.
- [92] Davey MJ, Fang L, McInerney P, Georgescu RE & O'Donnell M (2002) The DnaC helicase loader is a dual ATP/ADP switch protein. *EMBO J* **21**: 3148-3159.

- [93] Davey MJ, Jeruzalmi D, Kuriyan J & O'Donnell M (2002) Motors and switches: AAA+ machines within the replisome. *Nat Rev Mol Cell Biol* **3**: 826-835.
- [94] Davidson JF, Fox R, Harris DD, Lyons-Abbott S & Loeb LA (2003) Insertion of the T3 DNA polymerase thioredoxin binding domain enhances the processivity and fidelity of Taq DNA polymerase. *Nucleic Acids Res* **31**: 4702-4709.
- [95] Davis BM & Waldor MK (2003) Filamentous phages linked to virulence of *Vibrio cholerae*. *Curr Opin Microbiol* **6**: 35-42.
- [96] Dehò G, Zangrossi S, Sabbattini P, Sironi G & Ghisotti D (1992) Bacteriophage P4 immunity controlled by small RNAs via transcription termination. *Mol Microbiol* **6**: 3415-3425.
- [97] del Solar G & Espinosa M (2000) Plasmid copy number control: an ever-growing story. *Mol Microbiol* **37**: 492-500.
- [98] del Solar G, Giraldo R, Ruiz-Echevarria MJ, Espinosa M & Díaz Orejas R (1998) Replication and control of circular bacterial plasmids. *Microbiol Mol Biol Rev* **62**: 434-464.
- [99] Delarue M, Poch O, Tordo N, Moras D & Argos P (1990) An attempt to unify the structure of polymerases. *Protein Eng* **3**: 461-467.
- [100] DeRose EF, Li D, Darden T, Harvey S, Perrino FW, Schaaper RM & London RE (2002) Model for the catalytic domain of the proofreading epsilon subunit of *Escherichia coli* DNA polymerase III based on NMR structural data. *Biochemistry* **41**: 94-110.
- [101] Dervyn E, Suski C, Daniel R, Bruand C, Chapuis J, Errington J, Janniere L & Ehrlich SD (2001) Two Essential DNA Polymerases at the Bacterial Replication Fork. *Science* **294**: 1716-1719.
- [102] Desiere F, Lucchini S, Bruttin A, Zwahlen MC & Bruessow H (1997) A Highly Conserved DNA Replication Module from *Streptococcus thermophilus* Phages Is Similar in Sequence and Topology to a Module from *Lactococcus lactis* Phages. *Virology* **234**: 372-382.
- [103] Díaz Orejas R, Ziegelin G, Lurz R & Lanka E (1994) Phage P4 DNA replication in vitro. *Nucleic Acids Res* **22**: 2065-2070.
- [104] Díaz R, Barnsley P & Pritchard RH (1979) Location and characterisation of a new replication origin in the *E. coli* K12 chromosome. *Mol Gen Genet* **175**: 151-157.
- [105] Díaz R & Pritchard RH (1978) Cloning of replication origins from the *E. coli* K12 chromosome. *Nature* **275**: 561-564.
- [106] Dodson M, Echols H, Wickner S, *et al.* (1986) Specialized nucleoprotein structures at the origin of replication of bacteriophage lambda: localized unwinding of duplex DNA by a six-protein reaction. *Proc Natl Acad Sci USA* **83**: 7638-7642.
- [107] Donate LE, Llorca O, Bárcena M, Brown SE, Dixon NE & Carazo JM (2000) pH-controlled quaternary states of hexameric DnaB helicase. *J Mol Biol* **303**: 383-393.
- [108] Dong F, Gogol EP & von Hippel PH (1995) The phage T4-coded DNA replication helicase (gp41) forms a hexamer upon activation by nucleoside triphosphate. *J Biol Chem* **270**: 7462-7473.
- [109] Doublet S, Tabor S, Long AM, Richardson CC & Ellenberger T (1998) Crystal structure of a bacteriophage T7 DNA replication complex at 2.2 Å resolution. *Nature* **391**: 251-258.
- [110] Dubeau L, Hours C & Denhardt DT (1981) The mechanism of replication of phiX174. XVII. Purification and partial characterization of the gene A and A proteins. *Can J Biochem* **59**: 106-115.
- [111] Dunn JJ & Studier FW (1983) Complete nucleotide sequence of bacteriophage T7 DNA and the locations of T7 genetic elements. *J Mol Biol* **166**: 477-535.
- [112] Dunn K, Chrysogelos S & Griffith J (1982) Electron microscopic visualization of recA-DNA filaments: evidence for a cyclic extension of duplex DNA. *Cell* **28**: 757-765.
- [113] Echols H & Gingery R (1968) Mutants of bacteriophage lambda defective in vegetative genetic recombination. *J Mol Biol* **34**: 239-249.
- [114] Egelman EH, Yu X, Wild R, Hingorani MM & Patel SS (1995) Bacteriophage T7 helicase/primase proteins form rings around single-stranded DNA that suggest a general structure for hexameric helicases. *Proc Natl Acad Sci USA* **92**: 3869-3873.
- [115] Eisenberg S, Griffith J & Kornberg A (1977) phiX174 cistron A protein is a multifunctional enzyme in DNA replication. *Proc Natl Acad Sci USA* **74**: 3198-3202.
- [116] Eisenberg S & Kornberg A (1979) Purification and characterization of phiX174 gene A protein. A multifunctional enzyme of duplex DNA replication. *J Biol Chem* **254**: 5328-5332.
- [117] Eisenberg S, Scott JF & Kornberg A (1976) An enzyme system for replication of duplex circular DNA: the replicative form of phage phi X174. *Proc Natl Acad Sci USA* **73**: 1594-1597.
- [118] Eisenberg S, Scott JF & Kornberg A (1976) Enzymatic replication of viral and complementary strands of duplex DNA of phage phiX174 proceeds by separate mechanisms. *Proc Natl Acad Sci USA* **73**: 3151-3155.
- [119] Espinosa M, del Solar G, Rojo F & Alonso JC (1995) Plasmid rolling circle replication and its control. *FEMS Microbiol. Lett* **130**: 111-120.
- [120] Fang LH, Davey MJ & O'Donnell M (1999) Replisome assembly at oriC, the replication origin of *E. coli* reveals an explanation for initiation sites outside an origin. *Mol Cell* **4**: 541-553.
- [121] Felczak MM & Kaguni JM (2004) The Box VII motif of *Escherichia coli* DnaA protein is required for DnaA oligomerization at the *E. coli* replication origin. *J Biol Chem*
- [122] Filee J, Forterre P, Sen-Lin T & Laurent J (2002) Evolution of DNA polymerase families: evidences for multiple gene exchange between cellular and viral proteins. *J Mol Evol* **54**: 763-773.
- [123] Flensburg J & Calendar R (1987) Bacteriophage P4 DNA replication. Nucleotide sequence of the P4 replication gene and the cis replication region. *J Mol Biol* **195**: 439-445.
- [124] Foley S, Lucchini S, Zwahlen MC & Bruessow H (1998) A short noncoding viral DNA element showing characteristics of a replication origin confers bacteriophage resistance to *Streptococcus thermophilus*. *Virology* **250**: 377-387.
- [125] Ford ME, Sarkis GJ, Belanger AE, Hendrix RW & Hatfull GF (1998) Genome structure of mycobacteriophage D29: implications for phage evolution. *J Mol Biol* **279**: 143-164.
- [126] Forterre P (1999) Displacement of cellular proteins by functional analogues from plasmids or viruses could explain puzzling phylogenies of many DNA informational proteins. *Mol Microbiol* **33**: 457-465.
- [127] Franklin NC (1967) Deletions and functions of the center of the phi80 -lambda phage genome. Evidence for a phage function promoting genetic recombination. *Genetics* **57**: 301-318.
- [128] Frick DN, Baradaran K & Richardson CC (1998) An N-terminal fragment of the gene 4 helicase/primase of bacteriophage T7 retains primase activity in the absence of helicase activity. *Proc Natl Acad Sci USA* **95**: 7957-7962.
- [129] Frick DN & Richardson CC (2001) DNA Primases. *Annu Rev Biochem* **70**: 39-80.
- [130] Fulford W & Model P (1988) Bacteriophage f1 DNA replication genes. II. The roles of gene V protein and gene II protein in complementary strand synthesis. *J Mol Biol* **203**: 39-48.
- [131] Funnell BE, Baker TA & Kornberg A (1987) *In vitro* assembly of a prepriming complex at the origin of the *Escherichia coli* chromosome. *J Biol Chem* **262**: 10327-10334.

- [132] Funnell BE & Inman RB (1983) Bacteriophage P2 DNA replication. Characterization of the requirement of the gene B protein *in vivo*. *J Mol Biol* **167**: 311-334.
- [133] Furth ME, McLeester C & Dove WF (1978) Specificity determinants for bacteriophage lambda DNA replication. I. A chain of interactions that controls the initiation of replication. *J Mol Biol* **126**: 195-225.
- [134] Furth ME & Yates JL (1978) Specificity determinants for bacteriophage lambda DNA replication. II. Structure of O proteins of lambda-phi80 and lambda-82 hybrid phages and of a lambda mutant defective in the origin of replication. *J Mol Biol* **126**: 227-240.
- [135] Garforth SJ, Ceska TA, Suck D & Sayers JR (1999) Mutagenesis of conserved lysine residues in bacteriophage T5 5'-3' exonuclease suggests separate mechanisms of endo and exonucleolytic cleavage. *Proc Natl Acad Sci USA* **96**: 38-43.
- [136] Garzon A, Cano DA & Casades J (1998) The P22 Erf protein and host RecA provide alternative functions for transductional segregation of plasmid-borne duplications. *Mol Gen Genet* **259**: 39-45.
- [137] Gauss P, Gayle M, Winter RB & Gold L (1987) The bacteriophage T4 dexA gene: sequence and analysis of a gene conditionally required for DNA replication. *Mol Gen Genet* **206**: 24-34.
- [138] Geider K, Baumel I & Meyer TF (1982) Intermediate stages in enzymatic replication of bacteriophage fd duplex DNA. *J Biol Chem* **257**: 6488-6493.
- [139] George JW, Stohr BA, Tomso DJ & Kreuzer KN (2001) The tight linkage between DNA replication and double-strand break repair in bacteriophage T4. *Proc Natl Acad Sci USA* **98**: 8290-8297.
- [140] Gielow A, Diederich L & Messer W (1991) Characterization of a phage-plasmid hybrid (Phasyl) with two independent origins of replication isolated from *Escherichia coli*. *J Bacteriol* **173**: 73-79.
- [141] Gilbert W & Dressler D (1968) DNA replication: the rolling circle model. *Cold Spring Harbor Symp Quant Biol* **33**: 473-484.
- [142] Giraldo R (2003) Common domains in the initiators of DNA replication in Bacteria, Archaea and Eukarya: combined structural, functional and phylogenetic perspectives. *FEMS Microbiol Rev* **26**: 533-554.
- [143] Giraldo R & Fernandez-Tresguerres ME (2004) Twenty years of the pPS10 replicon: insights on the molecular mechanism for the activation of DNA replication in iteron-containing bacterial plasmids. *Plasmid* **52**: 69-83.
- [144] Girons IS, Bourhy P, Ottone C, Picardeau M, Yelton D, Hendrix RW, Glaser P & Charon N (2000) The LE1 bacteriophage replicates as a plasmid within *Leptospira biflexa*: construction of an *L. biflexa*-*Escherichia coli* shuttle vector. *J Bacteriol* **182**: 5700-5705.
- [145] Glover BP & McHenry CS (1998) The chi psi subunits of DNA polymerase III holoenzyme bind to single-stranded DNA-binding protein (SSB) and facilitate replication of an SSB-coated template. *J Biol Chem* **273**: 23476-23484.
- [146] Golub EI & Low KB (1986) Derepression of single-stranded DNA-binding protein genes on plasmids derepressed for conjugation, and complementation of an *E. coli* *ssb*- mutation by these genes. *Mol Gen Genet* **204**: 410-416.
- [147] Goodman MF (2002) Error-prone Repair DNA Polymerases in Prokaryotes and Eukaryotes. *Annu Rev Biochem* **71**: 17-50.
- [148] Gorbalenya AE, Koonin EV & Wolf YI (1990) A new superfamily of putative NTP-binding domains encoded by genomes of small DNA and RNA viruses. *Febs Lett* **262**: 145-148.
- [149] Gottesman S, Clark WP, Crecy-Lagard V & Maurizi MR (1993) ClpX, an alternative subunit for the ATP-dependent Clp protease of *Escherichia coli*. Sequence and *in vivo* activities. *J Biol Chem* **268**: 22618-22626.
- [150] Grabowski B & Kelman Z (2003) Archeal DNA replication: eukaryal proteins in a bacterial context. *Annu Rev Microbiol* **57**: 487-516.
- [151] Greenstein D & Horiuchi K (1987) Interaction between the replication origin and the initiator protein of the filamentous phage f1. Binding occurs in two steps. *J Mol Biol* **197**: 157-174.
- [152] Grigoriev A (1999) Strand-specific compositional asymmetries in double-stranded DNA viruses. *Virus Res* **60**: 1-19.
- [153] Groth AC & Calos MP (2004) Phage integrases: biology and applications. *J Mol Biol* **335**: 667-678.
- [154] Guenther B, Onrust R, Sali A, O'Donnell M & Kuriyan J (1997) Crystal structure of the delta' subunit of the clamp-loader complex of *E. coli* DNA polymerase III. *Cell* **91**: 335-345.
- [155] Gutiérrez C, Martin G, Sogo JM & Salas M (1991) Mechanism of stimulation of DNA replication by bacteriophage phi 29 single-stranded DNA-binding protein p5. *J Biol Chem* **266**: 2104-2111.
- [156] Gutman PD & Minton KW (1993) Conserved sites in the 5'-3' exonuclease domain of *Escherichia coli* DNA polymerase. *Nucleic Acids Res* **21**: 4406-4407.
- [157] Hadden JM, Declais AC, Phillips SEV & Lilley DMJ (2002) Metal ions bound at the active site of the junction-resolving enzyme T7 endonuclease I. *EMBO J* **21**: 3505-3515.
- [158] Hall MC & Matson SW (1999) Helicase motifs: the engine that powers DNA unwinding. *Mol Microbiol* **34**: 867-877.
- [159] Hall SD, Kane MF & Kolodner RD (1993) Identification and characterization of the *Escherichia coli* RecT protein, a protein encoded by the *recE* region that promotes renaturation of homologous single-stranded DNA. *J Bacteriol* **175**: 277-287.
- [160] Hamlett NV & Berger H (1975) Mutations altering genetic recombination and repair of DNA in bacteriophage T4. *Virology* **63**: 539-567.
- [161] Haniford DB & Chaconas G (1992) Mechanistic aspects of DNA transposition. *Curr Opin Genet Dev* **2**: 698-704.
- [162] Hansen EB (1989) Structure and regulation of the lytic replicon of phage P1. *J Mol Biol* **207**: 135-149.
- [163] Hansen EB & Yarmolinsky MB (1986) Host participation in plasmid maintenance: Dependence upon *dnaA* of replicons derived from P1 and F. *Proc Natl Acad Sci USA* **83**: 4423-4427.
- [164] Harth G, Baumel I, Meyer TF & Geider K (1981) Bacteriophage fd gene-2 protein. Processing of phage fd viral strands replicated by phage T7 enzymes. *Eur J Biochem* **119**: 663-668.
- [165] Hay N & Cohen G (1983) Requirement of *E. coli* DNA synthesis functions for the lytic replication of bacteriophage P1. *Virology* **131**: 193-206.
- [166] Heidekamp F, Langeveld SA, Baas PD & Jansz HS (1980) Studies of the recognition sequence of phi X174 gene A protein. Cleavage site of phi X gene A protein in St-I RFI DNA. *Nucleic Acids Res* **8**: 2009-2021.
- [167] Heinrich J, Riedel HD, Rückert B, Lurz R & Schuster H (1995) The lytic replicon of bacteriophage P1 is controlled by an antisense RNA. *Nucleic Acids Res* **23**: 1468-1474.
- [168] Heisig A, Severin I, Seefluth AK & Schuster H (1987) Regulation of the *ban* gene containing operon of prophage P1. *Mol Gen Genet* **206**: 368-376.
- [169] Hiasa H & Marians KJ (1999) Initiation of bidirectional replication at the chromosomal origin is directed by the interaction between helicase and primase. *J Biol Chem* **274**: 27244-27248.
- [170] Higashitani A, Greenstein D, Hirokawa H, Asano S & Horiuchi K (1994) Multiple DNA conformational changes induced by an initiator protein precede the nicking reaction in a rolling circle replication origin. *J Mol Biol* **237**: 388-400.

- [171] Higashitani A, Greenstein D & Horiuchi K (1992) A single amino acid substitution reduces the superhelicity requirement of a replication initiator protein. *Nucleic Acids Res* **20**: 2685-2691.
- [172] Higashitani A, Higashitani N & Horiuchi K (1997) Minus-strand origin of filamentous phage versus transcriptional promoters in recognition of RNA polymerase. *Proc Natl Acad Sci USA* **94**: 2909-2914.
- [173] Higashitani N, Higashitani A, Guan ZW & Horiuchi K (1996) Recognition mechanisms of the minus-strand origin of phage  $\phi$ 1 by *Escherichia coli* RNA polymerase. *Genes Cells* **1**: 829-841.
- [174] Higashitani N, Higashitani A & Horiuchi K (1993) Nucleotide sequence of the primer RNA for DNA replication of filamentous bacteriophages. *J Virol* **67**: 2175-2181.
- [175] Himawan JS & Richardson CC (1996) Amino acid residues critical for the interaction between bacteriophage T7 DNA polymerase and *Escherichia coli* thioredoxin. *J Biol Chem* **271**: 19999-20008.
- [176] Hinton DM & Nossal NG (1987) Bacteriophage T4 DNA primase-helicase. Characterization of oligomer synthesis by T4 61 protein alone and in conjunction with T4 41 protein. *J Biol Chem* **262**: 10873-10878.
- [177] Hoet PP, Coene MM & Cocito CG (1992) Replication Cycle of *Bacillus subtilis* Hydroxymethyluracil-Containing Phages. *Annu Rev Microbiol* **46**: 95-116.
- [178] Hollifield WC, Kaplan EN & Huang HV (1987) Efficient RecABC-dependent, homologous recombination between coliphage lambda and plasmids requires a phage ninR region gene. *Mol Gen Genet* **210**: 248-255.
- [179] Horiuchi K (1980) Origin of DNA replication of bacteriophage  $\phi$ 1 as the signal for termination. *Proc Natl Acad Sci USA* **77**: 5226-5229.
- [180] Horiuchi K (1986) Interaction between gene II protein and the DNA replication origin of bacteriophage  $\phi$ 1. *J Mol Biol* **188**: 215-223.
- [181] Horiuchi K (1997) Initiation mechanisms in replication of filamentous phage DNA. *Genes Cells* **2**: 425-432.
- [182] Hosoda J & Moise H (1978) Purification and physicochemical properties of limited proteolysis products of T4 helix destabilizing protein (gene 32 protein). *J Biol Chem* **253**: 7547-7558.
- [183] Hourcade D & Dressler D (1978) Site-specific initiation of a DNA fragment. *Proc Natl Acad Sci USA* **75**: 1652-1656.
- [184] Howland CJ, Rees CE, Barth PT & Wilkins BM (1989) The *ssb* gene of plasmid Collb-P9. *J Bacteriol* **171**: 2466-2473.
- [185] Huang YJ, Parker MM & Belfort M (1999) Role of exonucleolytic degradation in group I intron homing in phage T4. *Genetics* **153**: 1501-1512.
- [186] Huebscher U, Maga G & Spadari S (2002) Eukaryotic DNA Polymerases. *Annu Rev Biochem* **71**: 133-163.
- [187] Ikeda JE, Yudelevich A & Hurwitz J (1976) Isolation and characterization of the protein coded by gene A of bacteriophage  $\phi$ X174 DNA. *Proc Natl Acad Sci USA* **73**: 2669-2673.
- [188] Ilyina TV & Koonin EV (1992) Conserved sequence motifs in the initiator proteins for rolling circle DNA replication encoded by diverse replicons from eubacteria, eucaryotes and archaeobacteria. *Nucleic Acids Res* **20**: 3279-3285.
- [189] Imai Y, Ogasawara N, Ishigo-Oka D, Kadoya R, Daito T & Moriya S (2000) Subcellular localization of Dna-initiation proteins of *Bacillus subtilis*: evidence that chromosome replication begins at either edge of the nucleoids. *Mol Microbiol* **36**: 1037-1048.
- [190] Ishigo-Oka D, Ogasawara N & Moriya S (2001) DnaD protein of *Bacillus subtilis* interacts with DnaA, the initiator protein of replication. *J Bacteriol* **183**: 2148-2150.
- [191] Ishmael FT, Alley SC & Benkovic SJ (2001) Identification and Mapping of Protein-Protein Interactions between gp32 and gp59 by Cross-linking. *J Biol Chem* **276**: 25236-25242.
- [192] Ishmael FT, Alley SC & Benkovic SJ (2002) Assembly of the Bacteriophage T4 Helicase. Architecture and Stoichiometry of the gp41-gp59 Complex. *J Biol Chem* **277**: 20555-20562.
- [193] Iwasaki H, Nakata A, Walker GC & Shinagawa H (1990) The *Escherichia coli polB* gene, which encodes DNA polymerase II, is regulated by the SOS system. *J Bacteriol* **172**: 6268-6273.
- [194] Iyer LM, Koonin EV & Aravind L (2002) Classification and evolutionary history of the single-strand annealing proteins, RecT, Redbeta, ERF and RAD52. *BMC Genomics* **3**: 8-8.
- [195] Jacob F, Brenner S & Cuzin F (1963) On the regulation of DNA replication in bacteria. *Cold Spring Harbor Symp Quant Biol* **28**: 329-348.
- [196] Jeruzalmi D, O'Donnell M & Kuriyan J (2002) Clamp loaders and sliding clamps. *Curr Opin Struct Biol* **12**: 217-224.
- [197] Jezewska MJ & Bujalowski W (1996) Global conformational transitions in *Escherichia coli* primary replicative helicase DnaB protein induced by ATP, ADP, and single-stranded DNA binding. Multiple conformational states of the helicase hexamer. *J Biol Chem* **271**: 4261-4265.
- [198] Jones CE, Mueser TC, Dudas KC, Kreuzer KN & Nossal NG (2001) Bacteriophage T4 gene 41 helicase and gene 59 helicase-loading protein: a versatile couple with roles in replication and recombination. *Proc Natl Acad Sci USA* **98**: 8312-8318.
- [199] Jones CE, Mueser TC & Nossal NG (2000) Interaction of the Bacteriophage T4 Gene 59 Helicase Loading Protein and Gene 41 Helicase with Each Other and with Fork, Flap, and Cruciform DNA. *J Biol Chem* **275**: 27145-27154.
- [200] Joyce CM, Fujii DM, Laks HS, Hughes CM & Grindley ND (1985) Genetic mapping and DNA sequence analysis of mutations in the *polA* gene of *Escherichia coli*. *J Mol Biol* **186**: 283-293.
- [201] Joyce CM & Steitz TA (1994) Function and Structure Relationships in DNA Polymerases. *Annu Rev Biochem* **63**: 777-822.
- [202] Kadyrov FA & Drake JW (2004) UvsX recombinase and Dda helicase rescue stalled bacteriophage T4 DNA replication forks in vitro. *J Biol Chem* **279**(34): 35735-35740
- [203] Karam JD & Konigsberg WH (2000) DNA polymerase of the T4-related bacteriophages. *Prog Nucleic Acid Res Mol Biol* **64**: 65-96.
- [204] Kato M, Ito T, Wagner G, Richardson CC & Ellenberger T (2003) Modular architecture of the bacteriophage T7 primase couples RNA primer synthesis to DNA synthesis. *Mol Cell* **11**: 1349-1360.
- [205] Keck JL, Roche DD, Lynch AS & Berger JM (2000) Structure of the RNA polymerase domain of *E. coli* primase. *Science* **287**: 2482-2486.
- [206] Kelman Z, Hurwitz J & O'Donnell M (1998) Processivity of DNA polymerases: two mechanisms, one goal. *Structure* **6**: 121-125.
- [207] Kelman Z, Yuzhakov A, Andjelkovic J & O'Donnell M (1998) Devoted to the lagging strand-the subunit of DNA polymerase III holoenzyme contacts SSB to promote processive elongation and sliding clamp assembly. *EMBO J* **17**: 2436-2449.
- [208] Khan SA (1997) Rolling-circle replication of bacterial plasmids. *Microbiol Mol Biol Rev* **61**: 442-455.
- [209] Khan SA (2000) Plasmid rolling-circle replication: recent developments. *Mol Microbiol* **37**: 477-484.
- [210] Khan SA & Novick RP (1983) Complete nucleotide sequence of pT181, a tetracycline-resistance plasmid from *Staphylococcus aureus*. *Plasmid* **10**: 251-259.

- [211] Kim SG, Bor YC & Batt CA (1992) Bacteriophage resistance in *Lactococcus lactis* ssp. *lactis* using antisense ribonucleic acid. *J Dairy Sci* **75**: 1761-1767.
- [212] Kim YT & Richardson CC (1994) Acidic carboxyl-terminal domain of gene 2.5 protein of bacteriophage T7 is essential for protein-protein interactions. *J Biol Chem* **269**: 5270-5278.
- [213] Kim YT, Tabor S, Bortner C, Griffith JD & Richardson CC (1992) Purification and characterization of the bacteriophage T7 gene 2.5 protein. A single-stranded DNA-binding protein. *J Biol Chem* **267**: 15022-15031.
- [214] Kim YT, Tabor S, Churchich JE & Richardson CC (1992) Interactions of gene 2.5 protein and DNA polymerase of bacteriophage T7. *J Biol Chem* **267**: 15032-15040.
- [215] Kimura T, Asai T, Imai M & Takanami M (1989) Methylation strongly enhances DNA bending in the replication origin region of the *Escherichia coli* chromosome. *Mol Gen Genet* **219**: 69-74.
- [216] Kmiec E & Holloman WK (1981) Beta protein of bacteriophage lambda promotes renaturation of DNA. *J Biol Chem* **256**: 12636-12639.
- [217] Knopf CW (2000) Molecular mechanisms of replication of herpes simplex virus 1. *Acta Virol* **44**: 289-307.
- [218] Kogoma T (1997) Stable DNA replication: interplay between DNA replication, homologous recombination, and transcription. *Microbiol Mol Biol Rev* **61**: 212-238.
- [219] Kolodner R, Hall SD & Luisi-DeLuca C (1994) Homologous pairing proteins encoded by the *Escherichia coli* *recE* and *recT* genes. *Mol Microbiol* **11**: 23-30.
- [220] Kong D, Griffith JD & Richardson CC (1997) Gene 4 helicase of bacteriophage T7 mediates strand transfer through pyrimidine dimers, mismatches, and nonhomologous regions. *Proc Natl Acad Sci USA* **94**: 2987-2992.
- [221] Kong D, Nossal NG & Richardson CC (1997) Role of the bacteriophage T7 and T4 single-stranded DNA-binding proteins in the formation of joint molecules and DNA helicase-catalyzed polar branch migration. *J Biol Chem* **272**: 8380-8387.
- [222] Kong D & Richardson CC (1996) Single-stranded DNA binding protein and DNA helicase of bacteriophage T7 mediate homologous DNA strand exchange. *EMBO J* **15**: 2010-2019.
- [223] Kong D & Richardson CC (1998) Role of the acidic carboxyl-terminal domain of the single-stranded DNA-binding protein of bacteriophage T7 in specific protein-protein interactions. *J Biol Chem* **273**: 6556-6564.
- [224] Konieczny I & Marszalek J (1995) The requirement for molecular chaperones in lambda DNA replication is reduced by the mutation pi in lambda P gene, which weakens the interaction between lambda P protein and DnaB helicase. *J Biol Chem* **270**: 9792-9799.
- [225] Koonin EV (1992) DnaC protein contains a modified ATP-binding motif and belongs to a novel family of ATPases including also DnaA. *Nucleic Acids Res* **20**: 1997-1997.
- [226] Koonin EV & Ilyina TV (1992) Geminivirus replication proteins are related to prokaryotic plasmid rolling circle DNA replication initiator proteins. *J Gen Virol* **73**(10): 2763-2766.
- [227] Koonin EV & Ilyina TV (1993) Computer-assisted dissection of rolling circle DNA replication. *Biosystems* **30**: 241-268.
- [228] Kornberg A & Baker TA (1992) DNA Replication, WH Freeman and Company, New York, USA.
- [229] Koval R & Matthews BW (1997) Toroidal structure of lambda-dna-exonuclease. *Science* **277**: 1824-1827.
- [230] Kowalczykowski SC, Dixon DA, Eggleston AK, Lauder SD & Rehrauer WM (1994) Biochemistry of homologous recombination in *Escherichia coli*. *Microbiol Rev* **58**: 401-465.
- [231] Kowalski D & Eddy MJ (1989) The DNA unwinding element: A novel, cis-acting component that facilitates opening of the *Escherichia coli* replication origin. *EMBO J* **8**: 4335-4344.
- [232] Krause M & Messer W (1999) DnaA proteins of *Escherichia coli* and *Bacillus subtilis*: coordinate actions with single-stranded DNA-binding protein and interspecies inhibition during open complex formation at the replication origins. *Gene* **228**: 123-132.
- [233] Krevolin MD & Calendar R (1985) The replication of bacteriophage P4 DNA in vitro. Partial purification of the P4 alpha gene product. *J Mol Biol* **182**: 509-517.
- [234] Krevolin MD, Inman RB, Roof D, Kahn M & Calendar R (1985) Bacteriophage P4 DNA replication. Location of the P4 origin. *J Mol Biol* **182**: 519-527.
- [235] Kuil ME, Holmlund K, Vlaanderen CA & van Grondelle R (1990) Study of the binding of single-stranded DNA-binding protein to DNA and poly(rA) using electric field induced birefringence and circular dichroism spectroscopy. *Biochemistry* **29**: 8184-8189.
- [236] Kurosawa Y & Okazaki T (1979) Structure of the RNA portion of the RNA-linked DNA pieces in bacteriophage T4-infected *Escherichia coli* cells. *J Mol Biol* **135**: 841-861.
- [237] Kushner SR, Nagaishi H & Clark AJ (1974) Isolation of exonuclease VIII: the enzyme associated with sbcA indirect suppressor. *Proc Natl Acad Sci USA* **71**: 3593-3597.
- [238] Kuzminov A (1999) Recombinational repair of DNA damage in *Escherichia coli* and bacteriophage lambda. *Microbiol Mol Biol Rev* **63**: 751-813.
- [239] Langer U, Richter S, Roth A, Weigel C & Messer W (1996) A comprehensive set of DnaA box mutations in the replication origin, *oriC*, of *Escherichia coli*. *Mol Microbiol* **21**: 301-311.
- [240] Lanka E, Edelbluth C, Schlicht M & Schuster H (1978) *Escherichia coli* dnaB protein. Affinity chromatography on immobilized nucleotides. *J Biol Chem* **253**: 5847-5851.
- [241] Lanka E, Geschke B & Schuster H (1978) *Escherichia coli* dnaB mutant defective in DNA initiation: Isolation and properties of the dnaB protein. *Proc Natl Acad Sci USA* **75**: 799-803.
- [242] Lanka E, Mikolajczyk M, Schlicht M & Schuster H (1978) Association of the prophage P1ban protein with the dnaB protein of *Escherichia coli*. *J Biol Chem* **253**: 4746-4753.
- [243] Lanka E & Schuster H (1983) The DnaC protein of *Escherichia coli*. Purification, physical properties and interaction with DnaB protein. *Nucleic Acids Res* **11**: 987-997.
- [244] Learn BA, Um S-J, Huang L and McMacken R (1997) Cryptic single-stranded-DNA binding activities of the phage lambda P and *Escherichia coli* DnaC replication initiation proteins facilitate the transfer of *E. coli* DnaB helicase onto DNA. *Proc Natl Acad Sci USA* **94**: 1154-1159.
- [245] Leavitt MC & Ito J (1989) T5 DNA polymerase: structural-functional relationships to other DNA polymerases. *Proc Natl Acad Sci USA* **86**: 4465-4469.
- [246] LeBowitz JH & McMacken R (1986) The *Escherichia coli* dnaB replication protein is a DNA helicase. *J Biol Chem* **261**: 4738-4748.
- [247] Lee J, Chastain PD, Griffith JD & Richardson CC (2002) Lagging strand synthesis in coordinated DNA synthesis by bacteriophage t7 replication proteins. *J Mol Biol* **316**: 19-34.
- [248] Lee J, Chastain PD, Kusakabe T, Griffith JD & Richardson CC (1998) Coordinated leading and lagging strand DNA synthesis on a minicircular template. *Mol Cell* **1**: 1001-1010.
- [249] Lemonnier M, Ziegelin G, Reick T, Gomez AM, Diaz Orejas R & Lanka E (2003) Bacteriophage P1 Ban protein is a hexameric DNA helicase that interacts with and substitutes for *Escherichia coli* DnaB. *Nucleic Acids Res* **31**: 3918-3928.
- [250] Leonard AC & Grimwade JE (2005) Building a bacterial orisome: emergence of new regulatory features for replication origin unwinding. *Mol Microbiol* **55**: 978-985.
- [251] Lilley DM & White MF (2000) Resolving the relationships of resolving enzymes. *Proc Natl Acad Sci USA* **97**: 9351-9353.

- [252] Lin S & Kowalski D (1994) DNA helical instability facilitates initiation at the SV40 replication origin. *J Mol Biol* **235**: 496-507.
- [253] Lindqvist BH & Six EW (1971) Replication of bacteriophage P4 DNA in a nonlysogenic host. *Virology* **43**: 1-7.
- [254] Little JW (1967) An exonuclease induced by bacteriophage lambda. II. Nature of the enzymatic reaction. *J Biol Chem* **242**: 679-686.
- [255] Liu CC & Alberts BM (1980) Pentaribonucleotides of mixed sequence are synthesized and efficiently prime de novo DNA chain starts in the T4 bacteriophage DNA replication system. *Proc Natl Acad Sci USA* **77**: 5698-5702.
- [256] Liu Y & Haggård-Liungquist E (1994) Studies of bacteriophage P2 DNA replication: localization of the cleavage site of the A protein. *Nucleic Acids Res* **22**: 5204-5210.
- [257] Liu, Y & Haggård-Liungquist, E (1996) Functional characterization of the P2 A initiator protein and its DNA cleavage site. *Virology* **216**: 158-164.
- [258] Liu, Y, Saha, S & Haggård-Liungquist, E (1993) Studies of bacteriophage P2 DNA replication. The DNA sequence of the cis-acting gene A and ori region and construction of a P2 mini-chromosome. *J Mol Biol* **231**: 361-374.
- [259] Lobočka, MB, Rose, DJ, Plunkett, G, III, Rusin, M, Samojedny, A, Lehnerr, H, Yarmolinsky, MB & Blattner, FR (2004) Genome of bacteriophage P1. *J Bacteriol* **186**: 7032-7068.
- [260] Lohman, TM & Ferrari, ME (1994) *Escherichia coli* single-stranded DNA-binding protein: multiple DNA-binding modes and cooperativities. *Annu Rev Biochem* **63**: 527-570.
- [261] Lohman, TM & Overman, LB (1985) Two binding modes in *Escherichia coli* single strand binding protein-single stranded DNA complexes. Modulation by NaCl concentration. *J Biol Chem* **260**: 3594-3603.
- [262] Lopez de Saro, FJ, Georgescu, RE, Goodman, MF & O'Donnell, M (2003) Competitive processivity-clamp usage by DNA polymerases during DNA replication and repair. *EMBO J* **22**: 6408-6418.
- [263] Lopez de Saro FJ & O'Donnell M (2001) Interaction of the beta sliding clamp with MutS, ligase, and DNA polymerase I. *Proc Natl Acad Sci USA* **98**: 8376-8380.
- [264] Low RL, Shlomai J & Kornberg A (1982) Protein n, a primosomal DNA replication protein of *Escherichia coli*. Purification and characterization. *J Biol Chem* **257**: 6242-6250.
- [265] Lu YB, Ratnakar PVAL, Mohanty BK & Bastia D (1996) Direct physical interaction between DnaG primase and DnaB helicase of *Escherichia coli* is necessary for optimal synthesis of primer RNA. *Proc Natl Acad Sci USA* **93**: 12902-12907.
- [266] Ludlam AV, McNatt MW, Carr KM & Kaguni JM (2001) Essential amino acids of *Escherichia coli* DnaC protein in an N-terminal domain interact with DnaB helicase. *J Biol Chem* **276**: 27345-27353.
- [267] Lusetti SL & Cox MM (2002) The Bacterial RecA Protein and the Recombinational DNA Repair of Stalled Replication Forks. *Annu Rev Biochem* **71**: 71-100.
- [268] Madsen SM, Mills DJ, Djordjevic G, Israelsen H & Klaenhammer TR (2001) Analysis of the genetic switch and replication region of a P335-type bacteriophage with an obligate lytic lifestyle on *Lactococcus lactis*. *Appl Environ Microbiol* **67**: 1128-1139.
- [269] Maeda K, Nojiri H, Shintani M, Yoshida T, Habe H & Omori T (2003) Complete nucleotide sequence of carbazole/dioxin-degrading plasmid pCAR1 in *Pseudomonas resinovorans* strain CA10 indicates its mosaicity and the presence of large catabolic transposon Tn4676. *J Mol Biol* **326**: 21-33.
- [270] Maisnier-Patin S, Nordström K & Dasgupta S (2001) Replication arrests during a single round of replication of the *Escherichia coli* chromosome in the absence of DnaC activity. *Mol Microbiol* **42**: 1371-1382.
- [271] Mallory JB, Alfano C & McMacken R (1990) Host virus interactions in the initiation of bacteriophage lambda DNA replication. Recruitment of *Escherichia coli* DnaB helicase by lambda P replication protein. *J Biol Chem* **265**: 13297-13307.
- [272] Mannisto RH, Kivela HM, Paulin L, Bamford DH & Bamford JK (1999) The complete genome sequence of PM2, the first lipid-containing bacterial virus To Be isolated. *Virology* **262**: 355-363.
- [273] Marchler-Bauer A, Anderson JB, Cherukuri PF *et al.* (2005) CDD: a Conserved Domain Database for protein classification. *Nucleic Acids Res* **33**: D192-D196 (Database Issue).
- [274] Mardanov AV, Strakhova TS & Ravin NV (2004) RepA protein of the bacteriophage N15 exhibits activity of DNA helicase. *Dokl Biochem Biophys* **397**: 217-219.
- [275] Markiewicz P, Malone C, Chase JW & Rothman-Denes LB (1992) *Escherichia coli* single-stranded DNA-binding protein is a supercoiled template-dependent transcriptional activator of N4 virion RNA polymerase. *Genes Dev* **6**: 2010-2019.
- [276] Marsin S & Forterre P (1999) The active site of the rolling circle replication protein Rep75 is involved in site-specific nuclease, ligase and nucleotidyl transferase activities. *Mol Microbiol* **33**: 537-545.
- [277] Marsin S, McGovern S, Ehrlich SD, Bruand C & Polard P (2001) Early steps of *Bacillus subtilis* primosome assembly. *J Biol Chem* **276**: 45818-45825.
- [278] Marszalek J & Kaguni JM (1994) DnaA protein directs the binding of DnaB protein in initiation of DNA replication in *Escherichia coli*. *J Biol Chem* **269**: 4883-4890.
- [279] Marszalek J, Zhang WG, Hupp TR, Margulies C, Carr KM, Cherry S & Kaguni JM (1996) Domains of DnaA protein involved in interaction with DnaB protein, and in unwinding the *Escherichia coli* chromosomal origin. *J Biol Chem* **271**: 18535-18542.
- [280] Martinez-Jimenez MI, Alonso JC & Ayora S (2005) *Bacillus subtilis* Bacteriophage SPP1-encoded Gene 34.1 Product is a Recombination-dependent DNA Replication Protein. *J Mol Biol* **351**: 1007-1019.
- [281] McGlynn P & Lloyd RG (2000) Modulation of RNA polymerase by (p)ppGpp reveals a RecG-dependent mechanism for replication fork progression. *Cell* **101**: 35-45.
- [282] McGrath S, Seegers JF, Fitzgerald GF & van Sinderen D (1999) Molecular characterization of a phage-encoded resistance system in *Lactococcus lactis*. *Appl Environ Microbiol* **65**: 1891-1899.
- [283] McHenry CS (2003) Chromosomal replicases as asymmetric dimers: studies of subunit arrangement and functional consequences. *Mol Microbiol* **49**: 1157-1165.
- [284] McIntosh PK, Dunker R, Mulder C & Brown NC (1978) DNA of *Bacillus subtilis* bacteriophage SPP1: physical mapping and localization of the origin of replication. *J Virol* **28**: 865-876.
- [285] McMacken R, Ueda K & Kornberg A (1977) Migration of *Escherichia coli* DnaB protein on the template DNA strand as a mechanism in initiating DNA replication. *Proc Natl Acad Sci USA* **74**: 4190-4194.
- [286] Meijer M, Beck E, Hansen FG, Bergmans HE, Messer W, von Meyenburg K & Schaller H (1979) Nucleotide sequence of the origin of replication of the *E. coli* K-12 chromosome. *Proc Natl Acad Sci USA* **76**: 580-584.
- [287] Meijer WJ, Horcajadas JA & Salas M (2001) Phi29 family of phages. *Microbiol Mol Biol Rev* **65**: 261-287.
- [288] Mendelman LV, Notarnicola SM & Richardson CC (1992) Roles of bacteriophage T7 gene 4 proteins in providing primase and helicase functions in vivo. *Proc Natl Acad Sci USA* **89**: 10638-10642.
- [289] Mendelman LV & Richardson CC (1991) Requirements for primer synthesis by bacteriophage T7 63-kDa gene 4 protein.

- Roles of template sequence and T7 56-kDa gene 4 protein. *J Biol Chem* **266**: 23240-23250.
- [290] Messer W (2002) The bacterial replication initiator DnaA. DnaA and *oriC*, the bacterial mode to initiate DNA replication. *FEMS Microbiol Rev* **26**: 355-374.
- [291] Messer W, Hartmann-Kühlein H, Langer U, Mahlow E, Roth A, Schaper S, Urmoneit B & Woelker B (1992) The complex for replication initiation of *Escherichia coli*. *Chromosoma* **102**: S1-S6.
- [292] Messer W & Weigel C (1996) Initiation of chromosome replication. In: *Escherichia coli* and *Salmonella*, Cellular and Molecular Biology (Neidhardt FC, Curtiss III R, Ingraham J, Lin ECC, Low KB, Magasanik B, Reznikoff WS, Riley M, Schaechter M & Umbarger HE, Eds.), pp. 1579-1601. ASM Press, Washington, D.C., USA.
- [293] Messer W & Weigel C (1997) DnaA initiator - also a transcription factor. *Mol Microbiol* **24**: 1-6.
- [294] Meyer RR & Laine PS (1990) The single-stranded DNA-binding protein of *Escherichia coli*. *Microbiol Rev* **54**: 342-380.
- [295] Meyer TF & Geider K (1979) Bacteriophage fd gene II-protein. II. Specific cleavage and relaxation of supercoiled RF from filamentous phages. *J Biol Chem* **254**: 12642-12646.
- [296] Meyer TF, Geider K, Kurz C & Schaller H (1979) Cleavage site of bacteriophage fd gene II-protein in the origin of viral strand replication. *Nature* **278**: 365-367.
- [297] Missich R, Weise F, Chai S, Lurz R, Pedré X & Alonso JC (1997) The replisome organizer (G38P) of *Bacillus subtilis* bacteriophage SPP1 forms specialized nucleoprotein complexes with two discrete distant regions of the SPP1 genome. *J Mol Biol* **270**: 50-64.
- [298] Moarefi I, Jeruzalmi D, Turner J, O'Donnell M & Kuriyan J (2000) Crystal structure of the DNA polymerase processivity factor of T4 bacteriophage. *J Mol Biol* **296**: 1215-1223.
- [299] Moore DD, Denniston KJ & Blattner FR (1981) Sequence organization of the origins of DNA replication in lambdoid coliphages. *Gene* **14**: 91-101.
- [300] Moore T, McGlynn P, Ngo HP, Sharples GJ & Lloyd RG (2003) The RdgC protein of *Escherichia coli* binds DNA and counters a toxic effect of RecFOR in strains lacking the replication restart protein PriA. *EMBO J* **22**: 735-745.
- [301] Moreira D (2000) Multiple independent horizontal transfers of informational genes from bacteria to plasmids and phages: implications for the origin of bacterial replication machinery. *Mol Microbiol* **35**: 1-5.
- [302] Morris CF, Sinha NK & Alberts BM (1975) Reconstruction of bacteriophage T4 DNA replication apparatus from purified components: rolling circle replication following de novo chain initiation on a single-stranded circular DNA template. *Proc Natl Acad Sci USA* **72**: 4800-4804.
- [303] Moscoso M & Suarez JE (2000) Characterization of the DNA replication module of bacteriophage A2 and use of its origin of replication as a defense against infection during milk fermentation by *Lactobacillus casei*. *Virology* **273**: 101-111.
- [304] Moser MJ, Holley WR, Chatterjee A & Mian IS (1997) The proofreading domain of *Escherichia coli* DNA polymerase I and other DNA and/or RNA exonuclease domains. *Nucleic Acids Res* **25**: 5110-5118.
- [305] Mosig G (1998) Recombination and Recombination-dependent DNA Replication in Bacteriophage T4. *Annu Rev Genet* **32**: 379-413.
- [306] Moyer KE, Kimsey HH & Waldor MK (2001) Evidence for a rolling-circle mechanism of phage DNA synthesis from both replicative and integrated forms of CTXphi. *Mol Microbiol* **41**: 311-323.
- [307] Mueser TC, Jones CE, Nossal NG & Hyde CC (2000) Bacteriophage T4 gene 59 helicase assembly protein binds replication fork DNA. The 1.45 Å resolution crystal structure reveals a novel alpha-helical two-domain fold. *J Mol Biol* **296**: 597-612.
- [308] Mueser TC, Nossal NG & Hyde CC (1996) Structure of bacteriophage T4 RNase H, a 5' to 3' RNA-DNA and DNA-DNA exonuclease with sequence similarity to the RAD2 family of eukaryotic proteins. *Cell* **85**: 1101-1112.
- [309] Mufti S & Bernstein H (1974) The DNA-delay mutants of bacteriophage T4. *J Virol* **14**: 860-871.
- [310] Muniyappa K & Radding CM (1986) The homologous recombination system of phage lambda. Pairing activities of beta protein. *J Biol Chem* **261**: 7472-7478.
- [311] Murata T, Ohnishi M, Ara T, *et al.* (2002) Complete nucleotide sequence of plasmid Rts1: implications for evolution of large plasmid genomes. *J Bacteriol* **184**: 3194-3202.
- [312] Murphy KC, Casey L, Yannoutsos N, Potete AR & Hendrix RW (1987) Localization of a DNA-binding determinant in the bacteriophage P22 Erf protein. *J Mol Biol* **194**: 105-117.
- [313] Muyrers JP, Zhang Y, Buchholz F & Stewart AF (2000) RecE/RecT and Redalpha/Redbeta initiate double-stranded break repair by specifically interacting with their respective partners. *Genes Dev* **14**: 1971-1982.
- [314] Nakai H, Doseeva V & Jones JM (2001) Handoff from recombinase to replisome: insights from transposition. *Proc Natl Acad Sci USA* **98**: 8247-8254.
- [315] Nakai H & Richardson CC (1986) Interactions of the DNA polymerase and gene 4 protein of bacteriophage T7. Protein-protein and protein-DNA interactions involved in RNA-primed DNA synthesis. *J Biol Chem* **261**: 15208-15216.
- [316] Nakai H & Richardson CC (1988) The effect of the T7 and *Escherichia coli* DNA-binding proteins at the replication fork of bacteriophage T7. *J Biol Chem* **263**: 9831-9839.
- [317] Nanduri B, Byrd AK, Eoff RL, Tackett AJ & Raney KD (2002) Pre-steady-state DNA unwinding by bacteriophage T4 Dda helicase reveals a monomeric molecular motor. *Proc Natl Acad Sci USA* **99**: 14722-14727.
- [318] Neuwald AF, Aravind L, Spouge JL & Koonin EV (1999) AAA(+): A class of chaperone-like ATPases associated with the assembly, operation, and disassembly of protein complexes. *Genome Res* **9**: 27-43.
- [319] Ng JY & Mariani KJ (1996) The ordered assembly of the phiX174-type primosome. 2. Preservation of primosome composition from assembly through replication. *J Biol Chem* **271**: 15649-15655.
- [320] Niedenzu T, Roleke D, Bains G, Scherzinger E & Saenger W (2001) Crystal structure of the hexameric replicative helicase RepA of plasmid RSF1010. *J Mol Biol* **306**: 479-487.
- [321] Nossal NG (1980) RNA priming of DNA replication by bacteriophage T4 proteins. *J Biol Chem* **255**: 2176-2182.
- [322] Notarnicola SM, Mulcahy HL, Lee J & Richardson CC (1997) The acidic carboxyl terminus of the bacteriophage T7 gene 4 helicase/primase interacts with T7 DNA polymerase. *J Biol Chem* **272**: 18425-18433.
- [323] Notarnicola SM, Park K, Griffith JD & Richardson CC (1995) A domain of the gene 4 helicase/primase of bacteriophage T7 required for the formation of an active hexamer. *J Biol Chem* **270**: 20215-20224.
- [324] O'Donnell M, Jeruzalmi D & Kuriyan J (2001) Clamp loader structure predicts the architecture of DNA polymerase III holoenzyme and RFC. *Curr Biol* **11**: R935-R946.
- [325] Odegrip R & Haggård-Liungquist E (2001) The two active-site tyrosine residues of the a protein play non-equivalent roles during initiation of rolling circle replication of bacteriophage p2. *J Mol Biol* **308**: 147-163.
- [326] Odegrip R, Schoen S, Haggård-Liungquist E, Park K & Chatteraj DK (2000) The interaction of bacteriophage P2 B protein with *Escherichia coli* DnaB helicase. *J Virol* **74**: 4057-4063.



- [327] Oeser B, Gessner-Ulrich K, Deing P & Tudzynski P (1993) pCLK1 and pCIT5 - Two Linear Mitochondrial Plasmids from Unrelated *Claviceps purpurea* Strains: A Comparison. *Plasmid* **30**: 274-280.
- [328] Ogawa T (1975) Analysis of dnaB function of *Escherichia coli* K12 and the dnaB-like function of P1 prophage. *J Mol Biol* **94**: 327-340.
- [329] Ogawa T, Ogawa H & Tomizawa J (1988) Organization of the early region of bacteriophage phi 80. Genes and proteins. *J Mol Biol* **202**: 537-550.
- [330] Ogawa T & Okazaki T (1979) RNA-linked nascent DNA pieces in phage T7-infected *Escherichia coli*. III. Detection of intact primer RNA. *Nucleic Acids Res* **7**: 1621-1633.
- [331] Ohmori H, Friedberg EC, Fuchs RP, et al. (2001) The Y-family of DNA polymerases. *Mol Cell* **8**: 7-8.
- [332] Okazaki T, Kurosawa Y, Ogawa T, et al. (1979) Structure and metabolism of the RNA primer in the discontinuous replication of prokaryotic DNA. *Cold Spring Harbor Symp Quant Biol* **43(1)**: 203-219.
- [333] Ostergaard S, Brondsted L & Vogensen FK (2001) Identification of a replication protein and repeats essential for DNA replication of the temperate lactococcal bacteriophage TP901-1. *Appl Environ Microbiol* **67**: 774-781.
- [334] Pansegrau W & Lanka E (1991) Common sequence motifs in DNA relaxases and nick regions from a variety of DNA transfer systems. *Nucleic Acids Res* **19**: 3455
- [335] Pansegrau W & Lanka E (1992) A common sequence motif among prokaryotic DNA primases. *Nucleic Acids Res* **20**: 4931
- [336] Passy SI, Yu X, Li Z, Radding CM & Egelman EH (1999) Rings and filaments of beta protein from bacteriophage lambda suggest a superfamily of recombination proteins. *Proc Natl Acad Sci USA* **96**: 4279-4284.
- [337] Patel PH, Suzuki M, Adman E, Shinkai A & Loeb LA (2001) Prokaryotic DNA polymerase I: evolution, structure, and "base flipping" mechanism for nucleotide selection. *J Mol Biol* **308**: 823-837.
- [338] Patel SS & Hingorani MM (1993) Oligomeric structure of bacteriophage T7 DNA primase/helicase proteins. *J Biol Chem* **268**: 10668-10675.
- [339] Patel SS & Picha KM (2000) Structure and Function of Hexameric Helicases. *Annu Rev Biochem* **69**: 651-697.
- [340] Pecenkova T & Paces V (1999) Molecular phylogeny of phi29-like phages and their evolutionary relatedness to other protein-primed replicating phages and other phages hosted by gram-positive bacteria. *J Mol Evol* **48**: 197-208.
- [341] Pedré X, Weise F, Chai S, Lüder G & Alonso JC (1994) Analysis of cis and trans acting elements required for the initiation of DNA replication in the *Bacillus subtilis* bacteriophage SPP1. *J Mol Biol* **236**: 1324-1340.
- [342] Pedulla ML, Ford ME, Houtz JM, et al. (2003) Origins of highly mosaic mycobacteriophage genomes. *Cell* **113**: 171-182.
- [343] Peeters BP, Peters RM, Schoenmakers JG & Konings RN (1985) Nucleotide sequence and genetic organization of the genome of the N-specific filamentous bacteriophage IKE. Comparison with the genome of the F-specific filamentous phages M13, fd and fl. *J Mol Biol* **181**: 27-39.
- [344] Permina EA, Mironov AA & Gelfand MS (2002) Damage-repair error-prone polymerases of eubacteria: association with mobile genome elements. *Gene* **293**: 133-140.
- [345] Petit MA & Ehrlich SD (2002) Essential bacterial helicases that counteract the toxicity of recombination proteins. *EMBO J* **21**: 3137-3147.
- [346] Petri JB & Backhaus H (1984) Structural organization of the ori site of phage P22: comparison with other lambdoid ori sites. *Gene* **32**: 304-310.
- [347] Picha KM & Patel SS (1998) Bacteriophage T7 DNA helicase binds dTTP, forms hexamers, and binds DNA in the absence of Mg<sup>2+</sup>. The presence of dTTP is sufficient for hexamer formation and DNA binding. *J Biol Chem* **273**: 27315-27319.
- [348] Piechocki R, Kupper D, Quinones A & Langhammer R (1986) Mutational specificity of a proof-reading defective *Escherichia coli* dnaQ49 mutator. *Mol Gen Genet* **202**: 162-168.
- [349] Pietroni P, Young MC, Latham GJ & von Hippel PH (2001) Dissection of the ATP-driven reaction cycle of the bacteriophage T4 DNA replication processivity clamp loading system. *J Mol Biol* **309**: 869-891.
- [350] Polo S, Sturniolo T, Dehò G & Ghisotti D (1996) Identification of a phage-coded DNA-binding protein that regulates transcription from late promoters in bacteriophage P4. *J Mol Biol* **257**: 745-755.
- [351] Poteete AR & Fenton AC (1984) Lambda red-dependent growth and recombination of phage P22. *Virology* **134**: 161-167.
- [352] Poteete AR, Fenton AC & Semerjian AV (1991) Bacteriophage P22 accessory recombination function. *Virology* **182**: 316-323.
- [353] Poteete AR, Sauer RT & Hendrix RW (1983) Domain structure and quaternary organization of the bacteriophage P22 Erf protein. *J Mol Biol* **171**: 401-418.
- [354] Potrykus K, Wrobel B, Wegrzyn A & Wegrzyn G (2000) Replication of oriJ-based plasmid DNA during the stringent and relaxed responses of *Escherichia coli*. *Plasmid* **44**: 111-126.
- [355] Pritchard AE & McHenry CS (1999) Identification of the acidic residues in the active site of DNA polymerase III. *J Mol Biol* **285**: 1067-1080.
- [356] Pugh JC & Ritchie DA (1984) Formation of phage T1 concatemers by the RecE recombination pathway of *Escherichia coli*. *Virology* **135**: 200-206.
- [357] Raaijmakers H, Vix O, Toro I, Golz S, Kemper B & Suck D (1999) X-ray structure of T4 endonuclease VII: a DNA junction resolvase with a novel fold and unusual domain-swapped dimer architecture. *EMBO J* **18**: 1447-1458.
- [358] Rafferty JB, Bolt EL, Muranova TA, Sedelnikova SE, Leonard P, Pasquo A, Baker PJ, Rice DW, Sharples GJ & Lloyd RG (2003) The structure of *Escherichia coli* RusA endonuclease reveals a new Holliday junction DNA binding fold. *Structure (Camb)* **11**: 1557-1567.
- [359] Rakonjac J, Ward LJH, Schiemann AH, Gardner PP, Lubbers MW & O'Toole PW (2003) Sequence Diversity and Functional Conservation of the Origin of Replication in Lactococcal Prolate Phages. *Appl Environ Microbiol* **69**: 5104-5114.
- [360] Raney KD & Benkovic SJ (1995) Bacteriophage T4 Dda helicase translocates in a unidirectional fashion on single-stranded DNA. *J Biol Chem* **270**: 22236-22242.
- [361] Rangarajan S, Woodgate R & Goodman MF (2002) Replication restart in UV-irradiated *Escherichia coli* involving pols II, III, V, PriA, RecA and RecFOR proteins. *Mol Microbiol* **43**: 617-628.
- [362] Ratnakar PVAL, Mohanty BK, Lobert M & Bastia D (1996) The replication initiator protein pi of the plasmid R6K specifically interacts with the host-encoded helicase DnaB. *Proc Natl Acad Sci USA* **93**: 5522-5526.
- [363] Rattray AJ & Strathern JN (2003) Error-prone DNA polymerases: when making a mistake is the only way to get ahead. *Annu Rev Genet* **37**: 31-66.
- [364] Ravin NV, Kuprianov VV, Gilcrease EB & Casjens SR (2003) Bidirectional replication from an internal ori site of the linear N15 plasmid prophage. *Nucleic Acids Res* **31**: 6552-6560.
- [365] Ravin V, Ravin N, Casjens S, Ford ME, Hatfull GF & Hendrix RW (2000) Genomic sequence and analysis of the atypical temperate bacteriophage N15. *J Mol Biol* **299**: 53-73.

- [366] Reha-Krantz LJ (1998) Regulation of DNA polymerase exonucleolytic proofreading activity: studies of bacteriophage T4 "antimutator" DNA polymerases. *Genetics* **148**: 1551-1557.
- [367] Reha-Krantz LJ & Hurwitz J (1978) The *dnaB* gene product of *Escherichia coli*. II. Single stranded DNA-dependent ribonuclease triphosphatase activity. *J Biol Chem* **253**: 4051-4057.
- [368] Reiser W, Leibrecht I & Klein A (1983) Structure and function of mutants in the P gene of bacteriophage lambda leading to the pi phenotype. *Mol Gen Genet* **192**: 430-435.
- [369] Richardson CC (1983) Bacteriophage T7: minimal requirements for the replication of a duplex DNA molecule. *Cell* **33**: 315-317.
- [370] Richardson RW & Nossal NG (1989) Characterization of the bacteriophage T4 gene 41 DNA helicase. *J Biol Chem* **264**: 4725-4731.
- [371] Richardson RW & Nossal NG (1989) Trypsin cleavage in the COOH terminus of the bacteriophage T4 gene 41 DNA helicase alters the primase-helicase activities of the T4 replication complex in vitro. *J Biol Chem* **264**: 4732-4739.
- [372] Rodriguez I, Lazaro JM, Blanco L, *et al.* (2005) A specific subdomain in phi29 DNA polymerase confers both processivity and strand-displacement capacity. *Proc Natl Acad Sci USA* **102**: 6407-6412.
- [373] Rossmann MG, Moras D & Olsen KW (1974) Chemical and biological evolution of nucleotide-binding protein. *Nature* **250**: 194-199.
- [374] Rost B & Sander C (1994) Combining evolutionary information and neural networks to predict protein secondary structure. *Proteins* **19**: 55-72.
- [375] Roten CA, Gamba P, Barblan JL & Karamata D (2002) Comparative Genometrics (CG): a database dedicated to biometric comparisons of whole genomes. *Nucleic Acids Res* **30**: 142-144.
- [376] Roth A & Messer W (1995) The DNA binding domain of the initiator protein DnaA. *EMBO J* **14**: 2106-2111.
- [377] Roth MJ, Brown DR & Hurwitz J (1984) Analysis of bacteriophage phi X174 gene A protein-mediated termination and reinitiation of phi X DNA synthesis. II. Structural characterization of the covalent phi X A protein-DNA complex. *J Biol Chem* **259**: 10556-10568.
- [378] Rowen L & Kornberg A (1978) Primase, the DnaG protein of *Escherichia coli*. An enzyme which starts DNA chains. *J Biol Chem* **253**: 758-764.
- [379] Rubin EJ, Lin W, Mekalanos JJ & Waldor MK (1998) Replication and integration of a *Vibrio cholerae* cryptic plasmid linked to the CTX prophage. *Mol Microbiol* **28**: 1247-1254.
- [380] Ryder L, Sharples GJ & Lloyd RG (1996) Recombination-dependent growth in exonuclease-depleted *recBC sbcBC* strains of *Escherichia coli* K-12. *Genetics* **143**: 1101-1114.
- [381] Salinas F & Benkovic SJ (2000) Characterization of bacteriophage T4-coordinated leading- and lagging-strand synthesis on a minicircle substrate. *Proc Natl Acad Sci USA* **97**: 7196-7201.
- [382] San Martin C, Radermacher M, Wolpensinger B, Engel A, Miles CS, Dixon NE & Carazo JM (1998) Three-dimensional reconstructions from cryoelectron microscopy images reveal an intimate complex between helicase DnaB and its loading partner DnaC. *Structure* **6**: 501-509.
- [383] San Martin MC, Stamford NPJ, Dammerova N, Dixon NE & Carazo JM (1995) A structural model for the *Escherichia coli* DnaB helicase based on electron microscopy data. *J Struct Biol* **114**: 167-176.
- [384] Sanders GM, Kassavetis GA & Geiduschek EP (1995) Rules governing the efficiency and polarity of loading a tracking clamp protein onto DNA: determinants of enhancement in bacteriophage T4 late transcription. *EMBO J* **14**: 3966-3976.
- [385] Sandler SJ, Marians KJ, Zavitz KH, Couto J, Parent MA & Clark AJ (1999) *dnaC* mutations suppress defects in DNA replication- and recombination-associated functions in priB and priC double mutants in *Escherichia coli* K-12. *Mol Microbiol* **34**: 91-101.
- [386] Sanger F, Air GM, Barrell BG, *et al.* (1977) Nucleotide sequence of bacteriophage phi X174 DNA. *Nature* **265**: 687-695.
- [387] Sawaya MR, Guo S, Tabor S, Richardson CC & Ellenberger T (1999) Crystal structure of the helicase domain from the replicative helicase-primase of bacteriophage T7. *Cell* **99**: 167-177.
- [388] Schaaper RM (1993) Base selection, proofreading, and mismatch repair during DNA replication in *Escherichia coli*. *J Biol Chem* **268**: 23762-23765.
- [389] Schaper S & Messer W (1995) Interaction of the initiator protein DnaA of *Escherichia coli* with its DNA target. *J Biol Chem* **270**: 17622-17626.
- [390] Schekman R, Weiner JH, Weiner A & Kornberg A (1975) Ten proteins required for conversion of phiX174 single-stranded DNA to duplex form in vitro. Resolution and reconstitution. *J Biol Chem* **250**: 5859-5865.
- [391] Scherzinger E, Lanka E, Morelli G, Seiffert D & Yuki A (1977) Bacteriophage-T7-induced DNA-priming protein. A novel enzyme involved in DNA replication. *Eur J Biochem* **72**: 543-558.
- [392] Scherzinger E, Ziegelin G, Bárcena M, Carazo JM, Lurz R & Lanka E (1997) The RepA protein of plasmid RSF1010 is a replicative DNA helicase. *J Biol Chem* **272**: 30228-30236.
- [393] Scheuermann R, Tam S, Burgers PM, Lu C & Echols H (1983) Identification of the epsilon-subunit of *Escherichia coli* DNA polymerase III holoenzyme as the *dnaQ* gene product: a fidelity subunit for DNA replication. *Proc Natl Acad Sci USA* **80**: 7085-7089.
- [394] Schiemann AH, Rakonjac J, Callanan M, Gordon J, Polzin K, Lubbers MW & O'Toole PW (2004) Essentiality of the early transcript in the replication origin of the lactococcal prolate phage c2. *J Bacteriol* **186**: 8010-8017.
- [395] Schnos M, Zahn K, Blattner FR & Inman RB (1989) DNA looping induced by bacteriophage lambda O protein: implications for formation of higher order structures at the lambda origin of replication. *Virology* **168**: 370-377.
- [396] Schnos M, Zahn K, Inman RB & Blattner FR (1988) Initiation protein induced helix destabilization at the lambda origin: A prepriming step in DNA replication. *Cell* **52**: 385-395.
- [397] Scholl D, Kieczawa J, Kemp P, Rush J, Richardson CC, Merril C, Adhya S & Molineux IJ (2004) Genomic analysis of bacteriophages SP6 and K1-5, an estranged subgroup of the T7 supergroup. *J Mol Biol* **335**: 1151-1171.
- [398] Schouler C, Ehrlich SD & Chopin MC (1994) Sequence and organization of the lactococcal prolate-headed bIL67 phage genome. *Microbiology* **140**(11): 3061-3069.
- [399] Schrock RD & Alberts B (1996) Processivity of the gene 41 DNA helicase at the bacteriophage T4 DNA replication fork. *J Biol Chem* **271**: 16678-16682.
- [400] Schuster H, Lanka E, Edelbluth C, Geschke B, Mikolajczyk M, Schlicht M & Touati-Schwartz D (1979) A *dnaB*-analog DNA-replication protein of phage P1. *Cold Spring Harbor Symp Quant Biol* **43**(1): 551-557.
- [401] Scott JF, Eisenberg S, Bertsch LL & Kornberg A (1977) A mechanism of duplex DNA replication revealed by enzymatic studies of phage phi X174: catalytic strand separation in advance of replication. *Proc Natl Acad Sci USA* **74**: 193-197.
- [402] Seguritan V, Feng IW, Rohwer F, Swift M & Segall AM (2003) Genome sequences of two closely related *Vibrio parahaemolyticus* phages, VP16T and VP16C. *J Bacteriol* **185**: 6434-6447.
- [403] Seigneur M, Bidnenko V, Ehrlich SD & Michel B (1998) RuvAB acts at arrested replication forks. *Cell* **95**: 419-430.

- [404] Seitz H, Weigel C & Messer W (2000) The interaction domains of the DnaA and DnaB replication proteins of *Escherichia coli*. *Mol Microbiol* **37**: 1270-1279.
- [405] Seki T & Okazaki T (1979) RNA-linked nascent DNA pieces in phage T7-infected *Escherichia coli*. II. Primary structure of the RNA portion. *Nucleic Acids Res* **7**: 1603-1619.
- [406] Seufert W & Messer W (1987) Start sites for bidirectional *in vitro* replication inside the replication origin, *oriC*, of *Escherichia coli*. *EMBO J* **6**: 2469-2472.
- [407] Shamoo Y & Steitz TA (1999) Building a replisome from interacting pieces: sliding clamp complexed to a peptide from DNA polymerase and a polymerase editing complex. *Cell* **99**: 155-166.
- [408] Sharples GJ (2001) The X philes: structure-specific endonucleases that resolve Holliday junctions. *Mol Microbiol* **39**: 823-834.
- [409] Sharples GJ, Corbett LM & Graham IR (1998) lambda Rap protein is a structure-specific endonuclease involved in phage recombination. *Proc Natl Acad Sci USA* **95**: 13507-13512.
- [410] Sharples GJ, Ingleston SM & Lloyd RG (1999) Holliday junction processing in bacteria: insights from the evolutionary conservation of RuvABC, RecG, and RusA. *J Bacteriol* **181**: 5543-5550.
- [411] Shiffman D & Cohen SN (1992) Reconstruction of a *Streptomyces* linear replicon from separately cloned DNA fragments: existence of a cryptic origin of circular replication within the linear plasmid. *Proc Natl Acad Sci USA* **89**: 6129-6133.
- [412] Shinozaki K & Okazaki T (1977) RNA-linked nascent DNA pieces in T7 phage-infected *Escherichia coli* cells. I. Role of gene 6 exonuclease in removal of the linked RNA. *Mol Gen Genet* **154**: 263-267.
- [413] Shlomai J & Kornberg A (1980) An *E. coli* replication protein that recognizes a unique sequence within a hairpin region in  $\Phi$ X174 DNA. *Proc Natl Acad Sci USA* **77**: 799-803.
- [414] Shulman MJ, Hallick LM, Echols H & Signer ER (1970) Properties of recombination-deficient mutants of bacteriophage lambda. *J Mol Biol* **52**: 501-520.
- [415] Singleton MR, Sawaya MR, Ellenberger T & Wigley DB (2000) Crystal structure of T7 gene 4 ring helicase indicates a mechanism for sequential hydrolysis of nucleotides. *Cell* **101**: 589-600.
- [416] Singleton MR, Wentzell LM, Liu Y, West SC & Wigley DB (2002) Structure of the single-strand annealing domain of human RAD52 protein. *Proc Natl Acad Sci USA* **99**: 13492-13497.
- [417] Sivaprasad AV, Jarvinen R, Puspurs A & Egan JB (1990) DNA replication studies with coliphage 186. III. A single phage gene is required for phage 186 replication. *J Mol Biol* **213**: 449-463.
- [418] Skarstad K & Boye E (1994) The initiator protein DnaA: Evolution, properties and function. *Biochim Biophys Acta* **1217**: 111-130.
- [419] Smith GR (2001) Homologous recombination near and far from DNA breaks: alternative roles and contrasting views. *Annu Rev Genet* **35**: 243-274.
- [420] Sokolsky TD & Baker TA (2003) DNA gyrase requirements distinguish the alternate pathways of Mu transposition. *Mol Microbiol* **47**: 397-409.
- [421] Soultanas P (2002) A functional interaction between the putative primosomal protein DnaI and the main replicative DNA helicase DnaB in *Bacillus*. *Nucleic Acids Res* **30**: 966-974.
- [422] Sousa R (1996) Structural and mechanistic relationships between nucleic acid polymerases. *Trends Biochem Sci* **21**: 186-190.
- [423] Speck C & Messer W (2001) Mechanism of origin unwinding: sequential binding of DnaA initiator protein to double-stranded and single-stranded DNA in the AT-rich region of the replication origin. *EMBO J* **20**: 1469-1476.
- [424] Speck C, Weigel C & Messer W (1999) ATP- and ADP-DnaA protein, a molecular switch in gene regulation. *EMBO J* **18**: 6169-6176.
- [425] Staden R, Beal KF & Bonfield JK (1998) The Staden Package. In: *Bioinformatics Methods and Protocols*, Vol.132: Computer Methods in Molecular Biology (Misener S & Kraetz SA, Eds.), pp. 115-130. The Humana Press Inc., Totowa NJ, USA.
- [426] Stanley E, Walsh L, van der Zwet A, Fitzgerald GF & van Sinderen D (2000) Identification of four loci isolated from two *Streptococcus thermophilus* phage genomes responsible for mediating bacteriophage resistance. *FEMS Microbiol Lett* **182**: 271-277.
- [427] Stassen AP, Schoenmakers EF, Yu M, Schoenmakers JG & Konings RN (1992) Nucleotide sequence of the genome of the filamentous bacteriophage I2-2: module evolution of the filamentous phage genome. *J Mol Evol* **34**: 141-152.
- [428] Steitz TA (1999) DNA polymerases: structural diversity and common mechanisms. *J Biol Chem* **274**: 17395-17398.
- [429] Stenger DC, Revington GN, Stevenson MC & Bisaro DM (1991) Replicational release of geminivirus genomes from tandemly repeated copies: evidence for rolling-circle replication of a plant viral DNA. *Proc Natl Acad Sci USA* **88**: 8029-8033.
- [430] Stewart PE, Thalken R, Bono JL & Rosa P (2001) Isolation of a circular plasmid region sufficient for autonomous replication and transformation of infectious *Borrelia burgdorferi*. *Mol Microbiol* **39**: 714-721.
- [431] Strack B, Lessl M, Calendar R & Lanka E (1992) A common sequence motif, -E-G-Y-A-T-A-, identified within the primase domains of plasmid-encoded I- and P-type DNA primases and the alpha protein of the *Escherichia coli* satellite phage P4. *J Biol Chem* **267**: 13062-13072.
- [432] Stratling W & Knippers R (1973) Function and purification of gene 4 protein of phage T7. *Nature* **245**: 195-197.
- [433] Strick T, Allemand J, Croquette V & Bensimon D (2000) Twisting and stretching single DNA molecules. *Prog Biophys Mol Biol* **74**: 115-140.
- [434] Strick TR, Allemand JF, Bensimon D & Croquette V (2000) Stress-induced structural transitions in DNA and proteins. *Annu Rev Biophys Biomol Struct* **29**: 523-543.
- [435] Studier FW (1969) The genetics and physiology of bacteriophage T7. *Virology* **39**: 562-574.
- [436] Studier FW (1972) Bacteriophage T7. *Science* **176**: 367-376.
- [437] Sugimoto K, Oka A, Sugisaki H, Takanami M, Nishimura A, Yasuda Y & Hirota Y (1979) Nucleotide sequence of *Escherichia coli* replication origin. *Proc Natl Acad Sci USA* **76**: 575-579.
- [438] Sun W, Tormo J, Steitz TA & Godson GN (1994) Domains of *Escherichia coli* primase: functional activity of a 47-kDa N-terminal proteolytic fragment. *Proc Natl Acad Sci USA* **91**: 11462-11466.
- [439] Sung HM, Yeaman G, Ross CA & Yasbin RE (2003) Roles of YqjH and YqjW, homologs of the *Escherichia coli* UmuC/DinB or Y superfamily of DNA polymerases, in stationary-phase mutagenesis and UV-induced mutagenesis of *Bacillus subtilis*. *J Bacteriol* **185**: 2153-2160.
- [440] Szpirer J (1972) Control of development in temperate bacteriophages. IV. Specific action of N product at a transcription stop signal. *Mol Gen Genet* **114**: 297-304.
- [441] Tabor S, Huber HE & Richardson CC (1987) *Escherichia coli* thioredoxin confers processivity on the DNA polymerase activity of the gene 5 protein of bacteriophage T7. *J Biol Chem* **262**: 16212-16223.
- [442] Tanaka T & Ogura M (1998) A novel *Bacillus natto* plasmid pLS32 capable of replication in *Bacillus subtilis*. *FEBS Lett* **422**: 243-246.

- [443] Tarkowski TA, Mooney D, Thomason LC & Stahl FW (2002) Gene products encoded in the *ninR* region of phage lambda participate in Red-mediated recombination. *Genes Cells* **7**: 351-363.
- [444] Tatusov RL, Fedorova ND, Jackson JD, *et al.* (2003) The COG database: an updated version includes eukaryotes. *BMC Bioinformatics* **4**: 41-41.
- [445] Tatusova TA & Madden TL (1999) BLAST 2 Sequences, a new tool for comparing protein and nucleotide sequences. *FEMS Microbiol Lett* **174**: 247-250.
- [446] Taylor K & Wegrzyn G (1995) Replication of coliphage lambda DNA. *FEMS Microbiol Rev* **17**: 109-119.
- [447] Terzano S, Christian R, Espinoza FH, Calendar R, Dehò G & Ghisotti D (1994) A new gene of bacteriophage P4 that controls DNA replication. *J Bacteriol* **176**: 6059-6065.
- [448] Tessman ES (1966) Mutants of bacteriophage S13 blocked in infectious DNA synthesis. *J Mol Biol* **17**: 218-236.
- [449] Thirlway J, Turner IJ, Gibson CT, Gardiner L, Brady K, Allen S, Roberts CJ & Soultanas P (2004) DnaG interacts with a linker region that joins the N- and C-domains of DnaB and induces the formation of 3-fold symmetric rings. *Nucleic Acids Res* **32**: 2977-2986.
- [450] Thomson N, Sebahia M, Cerdano-Tarraga A & Parkhill J (2002) Sibling rivalry. *Trends Microbiol* **10**: 396-397.
- [451] Tocchetti A, Galimberti G, Dehò G & Ghisotti D (1999) Characterization of the *oriI* and *oriII* origins of replication in phage-plasmid P4. *J Virol* **73**: 7308-7316.
- [452] Tocchetti A, Serina S, Oliva I, Dehò G & Ghisotti D (2001) Cnr interferes with dimerization of the replication protein alpha in phage-plasmid P4. *Nucleic Acids Res* **29**: 536-544.
- [453] Torgov MY, Janzen DM & Reddy MK (1998) Efficiency and frequency of translational coupling between the bacteriophage T4 clamp loader genes. *J Bacteriol* **180**: 4339-4343.
- [454] Touati-Schwartz D (1979) A dnaB analog gene, specified by bacteriophage P1: genetic and physiological evidence for functional analogy and interactions between the two products. *Mol Gen Genet* **174**: 173-188.
- [455] Tougu K & Marians KJ (1996) The extreme C terminus of primase is required for interaction with DnaB at the replication fork. *J Biol Chem* **271**: 21391-21397.
- [456] Trakselis MA, Alley SC, Abel-Santos E & Benkovic SJ (2001) Creating a dynamic picture of the sliding clamp during T4 DNA polymerase holoenzyme assembly by using fluorescence resonance energy transfer *Proc Natl Acad Sci USA* **98**: 8368-8375.
- [457] Tseng TY, Frick DN & Richardson CC (2000) Characterization of a novel DNA primase from the *Salmonella typhimurium* bacteriophage SP6. *Biochemistry* **39**: 1643-1654.
- [458] Tsurimoto T & Matsubara K (1981) Purification of bacteriophage lambda-O protein that specifically binds to the origin of replication. *Mol Gen Genet* **181**: 325-331.
- [459] Tsurimoto T & Matsubara K (1981) Purified bacteriophage lambda O protein binds to four repeating sequences at the lambda replication origin. *Nucleic Acids Res* **9**: 1789-1799.
- [460] Umek RM & Kowalski D (1988) The ease of DNA unwinding as a determinant of initiation at yeast replication origins. *Cell* **52**: 559-567.
- [461] Van Dyck E, Foury F, Stillman B & Brill SJ (1992) A single-stranded DNA binding protein required for mitochondrial DNA replication in *S. cerevisiae* is homologous to *E. coli* SSB. *EMBO J* **11**: 3421-3430.
- [462] van Mansfeld AD, van Teeffelen HA, Baas PD & Jansz HS (1986) Two juxtaposed tyrosyl-OH groups participate in phi X174 gene A protein catalysed cleavage and ligation of DNA. *Nucleic Acids Res* **14**: 4229-4238.
- [463] Vellani TS & Myers RS (2003) Bacteriophage SPP1 Chu is an alkaline exonuclease in the SynExo family of viral two-component recombinases. *J Bacteriol* **185**: 2465-2474.
- [464] Velten M, McGovern S, Marsin S, Ehrlich SD, Noiro P & Polard P (2003) A two-protein strategy for the functional loading of a cellular replicative DNA helicase. *Mol Cell* **11**: 1009-1020.
- [465] Venkatesan M, Silver LL & Nossal NG (1982) Bacteriophage T4 gene 41 protein, required for the synthesis of RNA primers, is also a DNA helicase. *J Biol Chem* **257**: 12426-12434.
- [466] von Hippel PH & Delagoutte E (2003) Macromolecular complexes that unwind nucleic acids. *Bioessays* **25**: 1168-1177.
- [467] Wahle E, Lasken RS & Kornberg A (1989) The dnaB-dnaC replication protein complex of *Escherichia coli*. I. Formation and properties. *J Biol Chem* **264**: 2463-2468.
- [468] Wahle E, Lasken RS & Kornberg A (1989) The dnaB-dnaC replication protein complex of *Escherichia coli*. II. role of the complex in mobilizing dnaB functions. *J Biol Chem* **264**: 2469-2475.
- [469] Waterfield NR, Lubbers MW, Polzin KM, Le Page RW & Jarvis AW (1996) An origin of DNA replication from *Lactococcus lactis* bacteriophage c2. *Appl Environ Microbiol* **62**: 1452-1453.
- [470] Waters VL & Guiney DG (1993) Processes at the nick region link conjugation, T-DNA transfer and rolling circle replication. *Mol Microbiol* **9**: 1123-1130.
- [471] Wegrzyn A, Wegrzyn G & Taylor K (1995) Plasmid and host functions required for lambda plasmid replication carried out by the inherited replication complex. *Mol Gen Genet* **247**: 501-508.
- [472] Weigel C, Messer W, Preiss S, Welzeck M, Morigen and Boye E (2001) The sequence requirements for a functional *Escherichia coli* replication origin are different for the chromosome and a minichromosome. *Mol Microbiol* **40**: 498-507.
- [473] Weigel C, Schmidt A, Seitz H, Tuengler D, Welzeck M & Messer W (1999) The N-terminus promotes oligomerisation of the *Escherichia coli* initiator protein DnaA. *Mol Microbiol* **34**: 53-66.
- [474] Weigel C & Seitz H (2002) Strand-specific loading of DnaB helicase by DnaA to a substrate mimicking unwound *oriC*. *Mol Microbiol* **46**: 1149-1156.
- [475] Weise F, Chai S, Lüder G & Alonso JC (1994) Nucleotide sequence and complementation studies of the gene 35 region of the *Bacillus subtilis* bacteriophage SPP1. *Virology* **202**: 1046-1049.
- [476] Weng SF, Fan YF, Tseng YH & Lin JW (1997) Sequence analysis of the small cryptic *Xanthomonas campestris* pv. vesicatoria plasmid pXV64 encoding a Rep protein similar to gene II protein of phage 12-2. *Biochem Biophys Res Commun* **231**: 121-125.
- [477] Wickner RB, Wright M, Wickner S & Hurwitz J (1972) Conversion of phiX174 and fd single-stranded DNA to replicative forms in extracts of *Escherichia coli*. *Proc Natl Acad Sci USA* **69**: 3233-3237.
- [478] Wickner S & Hurwitz J (1974) Conversion of phiX174 viral DNA to double-stranded form by purified *Escherichia coli* proteins. *Proc Natl Acad Sci USA* **71**: 4120-4124.
- [479] Wickner S & Hurwitz J (1975) Interaction of *Escherichia coli* dnaB and dnaC(D) gene products *in vitro*. *Proc Natl Acad Sci USA* **72**: 921-925.
- [480] Wickner S & McKenney K (1987) Deletion analysis of the DNA sequence required for the *in vitro* initiation of replication of bacteriophage lambda. *J Biol Chem* **262**: 13163-13167.
- [481] Wickner SH & Zahn K (1986) Characterization of the DNA binding domain of bacteriophage lambda O protein. *J Biol Chem* **261**: 7537-7543.

- [482] Williams KR, Spicer EK, LoPresti MB, Guggenheimer RA & Chase JW (1983) Limited proteolysis studies on the *Escherichia coli* single-stranded DNA binding protein. Evidence for a functionally homologous domain in both the *Escherichia coli* and T4 DNA binding proteins. *J Biol Chem* **258**: 3346-3355.
- [483] Willis DK, Satin LH & Clark AJ (1985) Mutation-dependent suppression of *recB21 recC22* by a region cloned from the Rac prophage of *Escherichia coli* K-12. *J Bacteriol* **162**: 1166-1172.
- [484] Woelker B & Messer W (1993) The structure of the initiation complex at the replication origin, *oriC*, of *Escherichia coli*. *Nucleic Acids Res* **21**: 5025-5033.
- [485] Wojtkowiak D, Georgopoulos C & Zylicz M (1993) Isolation and characterization of ClpX, a new ATP-dependent specificity component of the Clp protease of *Escherichia coli*. *J Biol Chem* **268**: 22609-22617.
- [486] Wold MS (1997) Replication Protein A: A Heterotrimeric, Single-Stranded DNA-Binding Protein Required for Eukaryotic DNA Metabolism. *Annu Rev Biochem* **66**: 61-92.
- [487] Wong K & Geiduschek EP (1998) Activator-sigma interaction: A hydrophobic segment mediates the interaction of a sigma family promoter recognition protein with a sliding clamp transcription activator. *J Mol Biol* **284**: 195-203.
- [488] Wong SW, Wahl AF, Yuan PM, *et al.* (1988) Human DNA polymerase alpha gene expression is cell proliferation dependent and its primary structure is similar to both prokaryotic and eukaryotic replicative DNA polymerases. *EMBO J* **7**: 37-47.
- [489] Wrobel B & Wegrzyn G (2002) Evolution of lambdoid replication modules. *Virus Genes* **24**: 163-171.
- [490] Yamagami H & Yamamoto N (1970) Contribution of the bacterial recombination function to replication of bacteriophage P2. *J Mol Biol* **53**: 281-285.
- [491] Yasbin RE, Fields PI & Andersen BJ (1980) Properties of *Bacillus subtilis* 168 derivatives freed of their natural prophages. *Gene* **12**: 155-159.
- [492] Yeo HJ, Ziegelin G, Korolev S, Calendar R, Lanka E & Waksman G (2002) Phage P4 origin-binding domain structure reveals a mechanism for regulation of DNA-binding activity by homo- and heterodimerization of winged helix proteins. *Mol Microbiol* **43**: 855-867.
- [493] Yoda K & Okazaki T (1991) Specificity of recognition sequence for *Escherichia coli* primase. *Mol Gen Genet* **227**: 1-8.
- [494] Yonesaki T (1994) Involvement of a replicative DNA helicase of bacteriophage T4 in DNA recombination. *Genetics* **138**: 247-252.
- [495] Yonesaki T (1994) The purification and characterization of gene 59 protein from bacteriophage T4. *J Biol Chem* **269**: 1284-1289.
- [496] Yong Y & Romano LJ (1995) Nucleotide and DNA-induced conformational changes in the bacteriophage T7 gene 4 protein. *J Biol Chem* **270**: 24509-24517.
- [497] Yu X, Hingorani MM, Patel SS & Egelman EH (1996) DNA is bound within the central hole to one or two of the six subunits of the T7 DNA helicase. *Nat Struct Biol* **3**: 740-743.
- [498] Yu X, Jezewska MJ, Bujalowski W & Egelman EH (1996) The hexameric *E. coli* DnaB helicase can exist in different Quaternary states. *J Mol Biol* **259**: 7-14.
- [499] Zahn K & Blattner FR (1985) Binding and bending of the lambda replication origin by the phage O protein. *EMBO J* **4**: 3605-3616.
- [500] Zahn K & Blattner FR (1985) Sequence-induced DNA curvature at the bacteriophage lambda origin of replication. *Nature* **317**: 451-453.
- [501] Zahn K & Blattner FR (1987) Direct evidence for DNA bending at the lambda replication origin. *Science* **236**: 416-422.
- [502] Zechner EL, Wu CA & Marians KJ (1992) Coordinated leading- and lagging-strand synthesis at the *Escherichia coli* DNA replication fork. III. A polymerase-primase interaction governs primer size. *J Biol Chem* **267**: 4054-4063.
- [503] Ziegelin G, Calendar R, Ghisotti D, Terzano S & Lanka E (1997) Cnr protein, the negative regulator of bacteriophage P4 replication, stimulates specific DNA binding of its initiator protein alpha. *J Bacteriol* **179**: 2817-2822.
- [504] Ziegelin G, Calendar R, Lurz R & Lanka E (1997) The helicase domain of phage P4 alpha protein overlaps the specific DNA binding domain. *J Bacteriol* **179**: 4087-4095.
- [505] Ziegelin G, Linderoth NA, Calendar R & Lanka E (1995) Domain structure of phage P4 alpha protein deduced by mutational analysis. *J Bacteriol* **177**: 4333-4341.
- [506] Ziegelin G, Scherzinger E, Lurz R & Lanka E (1993) Phage P4 alpha protein is multifunctional with origin recognition, helicase and primase activities. *EMBO J* **12**: 3703-3708.
- [507] Ziegelin G, Tegtmeier N, Lurz R, Hertwig S, Hammerl J, Appel B & Lanka E (2005) The *repA* Gene of the Linear *Yersinia enterocolitica* Prophage PY54 Functions as a Circular Minimal Replicon in *Escherichia coli*. *J Bacteriol* **187**: 3445-3454.
- [508] Zuniga M, Franke-Fayard B, Venema G, Kok J & Nauta A (2002) Characterization of the putative replisome organizer of the lactococcal bacteriophage r1t. *J Virol* **76**: 10234-10244.
- [509] Zylicz M, Gorska I, Taylor K & Georgopoulos C (1984) Bacteriophage lambda replication proteins: formation of a mixed oligomer and binding to the origin of lambda DNA. *Mol Gen Genet* **196**: 401-406.

**Note added in proof:**

Due to the very recent change in the layout of FEMS Microbiology Reviews, the original numbering of Sections and Chapters is not preserved in the published main body of this review, but still referred to in this compendium. With the aim to minimise confusion resulting thereof, we present here a 'Table of Contents' of the main body with the original Section and Chapter numbering:

**bacteriophage replication modules****Contents: main body**

1. Introduction
2. Replication mechanisms
  - 2.1. Initiation by nicking: 'rolling circle'-type DNA replication (RCR) :  
*Phage fd - Phage P2*
  - 2.2. Initiation by melting: theta( $\theta$ )-type DNA replication ( $\theta$ DR) : *Phage  $\lambda$  - Phage SPP1 - Phage N15*
  - 2.3. Initiation at the ends of linear DNA: protein-primed DNA replication (ppDR) : *Phage  $\phi$ 29*
  - 2.4. Initiation of DNA replication by transcripts (tDR) : *Phage T7*
  - 2.5. Recombination-dependent DNA replication (RDR) : *Phage T4*
  - 2.6. Replication restart
3. Bacteriophage replication modules
  - 3.1. Phages encoding initiator proteins
    - 3.1.1. 'Initiator-solo' replication modules
    - 3.1.2. 'Initiator-helicase loader' replication modules
    - 3.1.3. 'Initiator-helicase' replication modules
    - 3.1.4. 'Initiator-helicase loader-helicase' replication modules
    - 3.1.5. Conclusions for Chapter 3.1.
  - 3.2. Replication module exchange among phages
  - 3.3.  $\phi$ P4 $\alpha$ -type helicase-primase encoding replication modules
  - 3.4. Phages encoding DNA polymerases  
*The phage T4-type replication module - The phage T7-type replication module - The phage D29-type replication module - The replication modules of the phages K, Bxz1, and T5 - The phage  $\phi$ 29-type replication module*
  - 3.5. Replication modules of phages replicating by RCR
  - 3.6. Phage replicons lacking replication protein genes
4. Evolutionary considerations
  - 4.1. The different types of phage-encoded helicases
  - 4.2. Phage-encoded homologues of the *E. coli* DnaB helicase
  - 4.3. Chromosomally encoded homologues of phage helicase loaders
5. Perspectives