

**INFORMATION
PROCESSING RETRIEVAL
(DMLS03)
(MLISC)**



ACHARYA NAGARJUNA UNIVERSITY

CENTRE FOR DISTANCE EDUCATION

NAGARJUNA NAGAR,

GUNTUR

ANDHRA PRADESH

CONTENTS

Lesson 2 :	Universal Decimal Classification	1-19
Lesson 3 :	BSO – the Broad System of Ordering	1-43
Lesson 4 :	Indexing Languages	1-7
Lesson 5 :	ISBD	1-3
Lesson 6 :	AACR2	1-5
Lesson 7 :	Common Communication Format	1-11
Lesson 8 :	Machine Readable Catalogue	1-6
Lesson 9 :	ISO 2709	1-3
Lesson 10 :	Metadata	1-7
Lesson 11:	Dublin Core	1-4
Lesson 12 :	Evaluation of Information Retrieval	1-9
Lesson 13 :	Information Retrieval Models	1-16

**INFORMATION
PROCESSING RETRIEVAL
- PAPER – III
MLISc.**

Lesson 2: UNIVERSAL DECIMAL CLASSIFICATION

The **Universal Decimal Classification** (UDC) is a bibliographic and library classification developed by the Belgian bibliographers Paul Otlet and Henri La Fontaine at the end of the 19th century. UDC provides a systematic arrangement of all branches of human knowledge organized as a coherent system in which knowledge fields are related and inter-linked.

Originally based on the Dewey Decimal Classification, the UDC was developed as a new analytico-synthetic classification system with a significantly larger vocabulary and syntax that enables very detailed content indexing and information retrieval in large collections.^[1] In its first edition in 1905, the UDC already included many features that were revolutionary in the context of knowledge classifications: tables of generally applicable (aspect-free) concepts - called common auxiliary tables; a series of special auxiliary tables with specific but re-usable attributes in a particular field of knowledge; an expressive notational system with connecting symbols and syntax rules to enable coordination of subjects and the creation of a documentation language proper. Although originally designed as an indexing and retrieval system, due to its logical structure and scalability, UDC has become one of the most widely used knowledge organization systems in libraries, where it is used for either shelf arrangement, content indexing or both.^[7] UDC codes can describe any type of document or object to any desired level of detail. These can include textual documents and other media such as films, video and soundrecordings, illustrations, maps as well as realia such as museum objects.

Since the first edition in French "Manuel du Répertoire bibliographique universel" (1905), UDC has been translated and published in various editions in 40 languages. UDC Summary, an abridged Web version of the scheme is available in over 50 languages.^[10] The classification has been modified and extended over the years to cope with increasing output in all areas of human knowledge, and is still under continuous review to take account of new developments.

Application of UDC

UDC is used in around 150,000 libraries in 130 countries and in many bibliographical services which require detailed content indexing. In a number of countries it is the main classification system for information exchange and is used in all type of libraries: public, school, academic and special libraries. UDC is also used in national bibliographies of around 30 countries. Examples of large databases indexed by UDC include:

NEBIS (The Network of Libraries and Information Centers in Switzerland) - 2.6 million records

COBIB.SI (Slovenian National Union Catalogue) - 3.5 million records

Hungarian National Union Catalogue (MOKKA) - 2.9 million records

VINITI RAS database (All-Russian Scientific and Technical Information Institute of Russian Academy of Science) with 28 million records

Meteorological & Geostrophysical Abstracts (MGA) with 600 journal titles

PORBASE (Portuguese National Bibliography) with 1.5 million records

UDC has traditionally been used for the indexing of scientific articles which was an important source of information of scientific output in the period predating electronic publishing. Collections of research articles in many countries covering decades of scientific output contain UDC codes. Examples of journal articles indexed by UDC

UDC Structure

Notation

A notation is a code commonly used in classification schemes to represent a class, i.e. a subject and its position in the hierarchy, to enable mechanical sorting and filing of subjects. UDC uses Arabic numerals arranged decimally. Every number is thought of as a decimal fraction with the initial decimal point omitted, which determines the filing order. An advantage of decimal notational systems is that they are infinitely extensible, and when new subdivisions are introduced, they need not disturb the existing allocation of numbers. For ease of reading, a UDC notation is usually punctuated after every third digit:

Notation	Caption (Class description)
539.120	Theoretical problems of elementary particles physics. Theories and models of fundamental interactions
539.120.2	Symmetries of quantum physics
539.120.22	Conservation laws
539.120.222	Translations. Rotations
539.120.224	Reflection in time and space
539.120.226	Space-time symmetries
539.120.23	Internal symmetries
539.120.3	Currents
539.120.4	Unified field theories
539.120.5	Strings

In UDC the notation has two features that make the scheme easier to browse and work with:

- **hierarchically expressive** - the longer the notation, the more specific the class: removing the final digit automatically produces a broader class code.

- **syntactically expressive** - when UDC codes are combined, the sequence of digits is interrupted by a precise type of punctuation sign which indicates that the expression is a combination of classes rather than a simple class e.g. the colon in 34:32 indicates that there are two distinct notational elements: 34 Law. Jurisprudence and 32 Politics; the closing and opening parentheses and double quotes in the following code 913(574.22)"19"(084.3) indicate four separate notational elements: 913 Regional geography, (574.22) North Kazakhstan (Soltüstik Qazaqstan); "19" 20th century and (084.3) Maps (document form)

Basic features and syntax

UDC is an analytico-synthetic and faceted classification. It allows an unlimited combination of attributes of a subject and relationships between subjects to be expressed. UDC codes from different tables can be combined to present various aspects of document content and form, e.g. 94(410)"19"(075) History (*main subject*) of United Kingdom (*place*) in 20th century (*time*), a textbook (*document form*). Or: 37:2 Relationship between Education and Religion. Complex UDC expressions can be accurately parsed into constituent elements.

UDC is also a disciplinary classification covering the entire universe of knowledge.^[23] This type of classification can also be described as *aspect* or *perspective*, which means that concepts are subsumed and placed under the field in which they are studied. Thus, the same concept can appear in different fields of knowledge. This particular feature is usually implemented in UDC by re-using the same concept in various combinations with the main subject, e.g. a code for language in common auxiliaries of language is used to derive numbers for ethnic grouping, individual languages in linguistics and individual literatures. Or, a code from the auxiliaries of place, e.g.(410) *United Kingdom*, uniquely representing the concept of United Kingdom can be used to express 911(410) *Regional geography of United Kingdom* and 94(410) *History of United Kingdom*.

Organization of classes

Concepts are organized in two kinds of tables in UDC:

- **Common auxiliary tables** (including certain auxiliary signs). These tables contain facets of concepts representing, general recurrent characteristics, applicable over a range of subjects throughout the main tables, including notions such as place, language of the text and physical form of the document, which may occur in almost any subject. UDC numbers from these tables, called common auxiliaries are simply added at the end of the number for the subject taken from the main tables. There are over 15,000 of common auxiliaries in UDC.
- **The main tables or main schedules** containing the various disciplines and branches of knowledge, arranged in 9 main classes, numbered from 0 to 9 (with class 4 being vacant). At the beginning of each class there are also series of special auxiliaries, which express aspects that are recurrent within this specific class. Main tables in UDC contain more than 60,000 subdivisions.

Main classes

- 0 Science and Knowledge. Organization. Computer Science. Information Science. Documentation. Librarianship. Institutions. Publications
- 1 Philosophy. Psychology
- 2 Religion. Theology
- 3 Social Sciences
- 4 *vacant*
- 5 Mathematics. Natural Sciences
- 6 Applied Sciences. Medicine. Technology
- 7 The Arts. Entertainment. Sport
- 8 Linguistics. Literature
- 9 Geography. History

The vacant class 4 is the result of a planned schedule expansion. This class was freed by moving linguistics into class 8 in the 1960s to make space for future developments in the rapidly expanding fields of knowledge; primarily natural sciences and technology.

Common auxiliary tables

Common auxiliaries are aspect-free concepts that can be used in combination with any other UDC code from the main classes or with other common auxiliaries. They have unique notational representations that makes them stand out in complex expressions. Common auxiliary numbers always begin with a certain symbol known as a facet indicator, e.g. = (equal sign) always introduces concepts representing the language of a document; (0...) numbers enclosed in parentheses starting with zero always represent a concept designating document form. Thus (075) Textbook and =111 English can be combined to express, e.g.(075)=111 Textbooks in English, and when combined with numbers from the main UDC tables they can be used as follows: 2(075)=111 Religion textbooks in English, 51(075)=111 Mathematics textbooks in English etc.

- =... Common auxiliaries of language. Table 1c
- (0...) Common auxiliaries of form. Table 1d
- (1/9) Common auxiliaries of place. Table 1e
- (=...) Common auxiliaries of human ancestry, ethnic grouping and nationality. Table 1f
- "... " Common auxiliaries of time. Table 1g helps to make minute division of time e.g.: "1993-1996
- -0... Common auxiliaries of general characteristics: Properties, Materials, Relations/Processes and Persons. Table 1k.
- -02 Common auxiliaries of properties. Table 1k
- -03 Common auxiliaries of materials. Table 1k

- -04 Common auxiliaries of relations, processes and operations. Table 1k
- -05 Common auxiliaries of persons and personal characteristics. Table 1k this table is repeated

Connecting signs

In order to preserve the precise meaning and enable accurate parsing of complex UDC expressions, a number of connecting symbols are made available to relate and extend UDC numbers. These are:

Symbol	Symbol name	Meaning	Example
+	<u>plus</u>	coordination, addition	e.g. 59+636 <u>zoology</u> and <u>animal breeding</u>
/	<u>stroke</u>	consecutive extension	e.g. 592/599 Systematic zoology (everything from 592 to 599 inclusive)
:	<u>colon</u>	relation	e.g. 17:7 Relation of <u>ethics</u> to <u>art</u>
[]	<u>squarebrackets</u>	subgrouping	e.g. 311:[622+669](485) <u>statistics</u> of <u>mining</u> and <u>metallurgy</u> in <u>Sweden</u> (the auxiliary qualifiers 622+669 considered as a unit)
*	asterisk	Introduces non-UDC notation	e.g. 523.4*433 Planetology, minor planet Eros (IAU authorized number after the asterisk)
A/Z	alphabetical extension	Direct alphabetical specification	e.g. 821.133.1MOL French literature, works of Molière

UDC classes in this outline are taken from the Multilingual Universal Decimal Classification Summary (UDCC Publication No. 088) released on the by the UDC Consortium under the Creative Commons Attribution Share Alike 3.0 license (first release 2009, subsequent update 2012).

Main Tables[

0 Science and knowledge. Organization. Computer science. Information. Documentation. Librarianship. Institution. Publications[\[edit source\]](#)

- 00 Prolegomena. Fundamentals of knowledge and culture. Propaedeutics
- 001 Science and knowledge in general. Organization of intellectual work
- 002 Documentation. Books. Writings. Authorship
- 003 Writing systems and scripts
- 004 Computer science and technology. Computing
- 004.2 Computer architecture
- 004.3 Computer hardware
- 004.4 Software
- 004.5 Human-computer interaction
- 004.6 Data
- 004.7 Computer communication
- 004.8 Artificial intelligence
- 004.9 Application-oriented computer-based techniques
- 005 Management
- 005.1 Management Theory
- 005.2 Management agents. Mechanisms. Measures
- 005.3 Management activities
- 005.5 Management operations. Direction
- 005.6 Quality management. Total quality management (TQM)
- 005.7 Organizational management (OM)
- 005.9 Fields of management
- 005.92 Records management
- 005.93 Plant management. Physical resources management
- 005.94 Knowledge management
- 005.95/.96 Personnel management. Human Resources management
- 006 Standardization of products, operations, weights, measures and time
- 007 Activity and organizing. Information. Communication and control theory generally (cybernetics)
- 008 Civilization. Culture. Progress
- 01 Bibliography and bibliographies. Catalogues
- 02 Librarianship
- 030 General reference works (as subject)

- 050 Serial publications, periodicals (as subject)
- 06 Organizations of a general nature
- 069 Museums
- 070 Newspapers (as subject). The Press. Outline of journalism
- 08 Polygraphies. Collective works (as subject)
- 09 Manuscripts. Rare and remarkable works (as subject)

1 Philosophy. Psychology[\[edit source\]](#)

- 101 Nature and role of philosophy
- 11 Metaphysics
- 111 General metaphysics. Ontology
- 122/129 Special Metaphysics
- 13 Philosophy of mind and spirit. Metaphysics of spiritual life
- 14 Philosophical systems and points of view
- 159.9 Psychology
- 159.91 Psychophysiology (physiological psychology). Mental physiology
- 159.92 Mental development and capacity. Comparative psychology
- 159.93 Sensation. Sensory perception
- 159.94 Executive functions
- 159.95 Higher mental processes
- 159.96 Special mental states and processes
- 159.97 Abnormal psychology
- 159.98 Applied psychology (psychotechnology) in general
- 16 Logic. Epistemology. Theory of knowledge. Methodology of logic
- 17 Moral philosophy. Ethics. Practical philosophy

2 Religion. Theology[\[edit source\]](#)

The UDC tables for religion are fully faceted. Indicated in italics below, are special auxiliary numbers that can be used to express attributes (facets) of any specific faith. Any special number can be combined with any religion e.g. *-5 Worship* can be used to express e.g. *26-5 Worship in Judaism, 27-5 Worship in Christianity, 24-5 Worship in Buddhism*. The complete special auxiliary tables contain around 2000 subdivisions of various attributes that can be attached to express various aspects of individual faiths to a great level of specificity allowing equal level of detail for every religion.

2-1/-9 Special auxiliary subdivision for religion

2-1 *Theory and philosophy of religion. Nature of religion. Phenomenon of religion*

2-2 *Evidences of religion*

2-3 *Persons in religion*

2-4 *Religious activities. Religious practice*

2-5 *Worship broadly. Cult. Rites and ceremonies*

2-6 *Processes in religion*

2-7 *Religious organization and administration*

2-8 *Religions characterised by various properties*

2-9 *History of the faith, religion, denomination or church*

21/29 *Religious systems. Religions and faiths*

21 *Prehistoric and primitive religions*

22 *Religions originating in the Far East*

23 *Religions originating in Indian sub-continent. Hindu religion in the broad sense*

24 Buddhism

25 *Religions of antiquity. Minor cults and religions*

26 Judaism

27 Christianity

28 Islam

29 *Modern spiritual movements*

3 Social Sciences[\[edit source\]](#)

303 *Methods of the social sciences*

304 *Social questions. Social practice. Cultural practice. Way of life (Lebensweise)*

305 *Gender studies*

308 *Sociography. Descriptive studies of society (both qualitative and quantitative)*

311 Statistics as a science. Statistical theory

314/316 Society

314 *Demography. Population studies*

316 Sociology

32 Politics

33 Economics. Economic science

34 Law. Jurisprudence

35 *Public administration. Government. Military affairs*

36 *Safeguarding the mental and material necessities of life*

37 Education

39 *Cultural anthropology. Ethnography. Customs. Manners. Traditions. Way of life*

4 Vacant[edit source]

This section is currently vacant.

5 Mathematics. Natural sciences[edit source]

502/504 Environmental science. Conservation of natural resources. Threats to the environment and protection against them

502 The environment and its protection

504 Threats to the environment

51 Mathematics

510 Fundamental and general considerations of mathematics

511 Number theory

512 Algebra

514 Geometry

517 Analysis

519.1 Combinatorial analysis. Graph theory

519.2 Probability. Mathematical statistics

519.6 Computational mathematics. Numerical analysis

519.7 Mathematical cybernetics

519.8 Operational research (OR): mathematical theories and methods

52 Astronomy. Astrophysics. Space research. Geodesy

53 Physics

531/534 Mechanics

535 Optics

536 Heat. Thermodynamics. Statistical physics

537 Electricity. Magnetism. Electromagnetism

538.9 Condensed matter physics. Solid state physics

539 Physical nature of matter

54 Chemistry. Crystallography. Mineralogy

542 Practical laboratory chemistry. Preparative and experimental chemistry

543 Analytical chemistry

544 Physical chemistry

546 Inorganic chemistry

547 Organic chemistry

548/549 Mineralogical sciences. Crystallography. Mineralogy

55 Earth Sciences. Geological sciences

56 Palaeontology

57 Biological sciences in general

58 Botany

59 Zoology

6 Applied sciences. Medicine. Technology[\[edit source\]](#)

Class 6 occupies the largest proportion of UDC schedules. It contains over 44,000 subdivisions. Each specific field of technology or industry usually contains more than one special auxiliary table with concepts needed to express operations, processes, materials and products. As a result, UDC codes are often created through the combination of various attributes. Equally, some parts of this class enumerate concepts to a great level of detail e.g. *621.882.212 Hexagon screws with additional shapes. Including: Flank screws. Collar screws. Cap screws*

- 60 Biotechnology
- 61 Medical sciences
- 611/612 Human biology
- 613 Hygiene generally. Personal health and hygiene
- 614 Public health and hygiene. Accident prevention
- 615 Pharmacology. Therapeutics. Toxicology
- 616 Pathology. Clinical medicine
- 617 Surgery. Orthopaedics. Ophthalmology
- 618 Gynaecology. Obstetrics
- 62 Engineering. Technology in general
- 620 Materials testing. Commercial materials. Power stations. Economics of energy
- 621 Mechanical engineering in general. Nuclear technology. Electrical engineering. Machinery
- 622 Mining
- 623 Military engineering
- 624 Civil and structural engineering in general
- 625 Civil engineering of land transport. Railway engineering. Highway engineering
- 626/627 Hydraulic engineering and construction. Water (aquatic) structures
- 629 Transport vehicle engineering
- 63 Agriculture and related sciences and techniques. Forestry. Farming. Wildlife exploitation
- 630 Forestry
- 631/635 Farm management. Agronomy. Horticulture
- 633/635 Horticulture in general. Specific crops
- 636 Animal husbandry and breeding in general. Livestock rearing. Breeding of domestic animals
- 64 Home economics. Domestic science. Housekeeping
- 65 Communication and transport industries. Accountancy. Business management. Public relations
- 654 Telecommunication and telecontrol (organization, services)
- 655 Graphic industries. Printing. Publishing. Book trade
- 656 Transport and postal services. Traffic organization and control
- 657 Accountancy
- 658 Business management, administration. Commercial organization

- 659 Publicity. Information work. Public relations
- 66 Chemical technology. Chemical and related industries
- 67 Various industries, trades and crafts
- 68 Industries, crafts and trades for finished or assembled articles
- 69 Building (construction) trade. Building materials. Building practice and procedure

7 The arts. Recreation. Entertainment. Sport[\[edit source\]](#)

- 7.01/09 *Special auxiliary subdivision for the arts*
- 7.01 *Theory and philosophy of art. Principles of design, proportion, optical effect*
- 7.02 *Art technique. Craftsmanship*
- 7.03 *Artistic periods and phases. Schools, styles, influences*
- 7.04 *Subjects for artistic representation. Iconography. Iconology*
- 7.05 *Applications of art (in industry, trade, the home, everyday life)*
- 7.06 *Various questions concerning art*
- 7.07 *Occupations and activities associated with the arts and entertainment*
- 7.08 *Characteristic features, forms, combinations etc. (in art, entertainment and sport)*
- 7.091 *Performance, presentation (in original medium)*
- 71 Physical planning. Regional, town and country planning. Landscapes, parks, gardens
- 72 Architecture
- 73 Plastic arts
- 74 Drawing. Design. Applied arts and crafts
- 745/749 Industrial and domestic arts and crafts. Applied arts
- 75 Painting
- 76 Graphic art, printmaking. Graphics
- 77 Photography and similar processes
- 78 Music
- 79 Recreation. Entertainment. Games. Sport
- 791 Cinema. Films (motion pictures)
- 792 Theatre. Stagecraft. Dramatic performances
- 793 Social entertainments and recreations. Art of movement. Dance
- 794 Board and table games (of thought, skill and chance)
- 796 Sport. Games. Physical exercises
- 797 Water sports. Aerial sports
- 798 Riding and driving. Horse and other animal sports
- 799 Sport fishing. Sport hunting. Shooting and target sports

8 Language. Linguistics. Literature

Tables for class 8 are fully faceted and details are expressed through combination with common auxiliaries of language (Table 1c) and a series of special auxiliary tables to indicate other facets or attributes in Linguistics or Literature. As a result, this class allows for great specificity in indexing although the schedules themselves occupy very little space in UDC. The subdivisions of e.g. *811 Languages* or *821 Literature* are derived from common auxiliaries of language =1/=9 (Table 1c) by substituting a point for the equals sign, e.g. *811.111 English language* (as a subject of a linguistic study) and *821.111 English literature* derives from =111 *English language*. Common auxiliaries of place and time are also frequently used in this class to express place and time facets of Linguistics or Literature, e.g. *821.111(71)"18" English literature of Canada in 19th century*

80 General questions relating to both linguistics and literature. Philology

801 Prosody. Auxiliary sciences and sources of philology

808 Rhetoric. The effective use of language

81 Linguistics and languages

81`1/4 Special auxiliary subdivision for subject fields and facets of linguistics and languages

81`1 General linguistics

81`2 Theory of signs. Theory of translation. Standardization. Usage. Geographical linguistics

81`3 Mathematical and applied linguistics. Phonetics. Graphemics. Grammar. Semantics. Stylistics

81`4 Text linguistics, Discourse analysis. Typological linguistics

81`42 Text linguistics. Discourse analysis

81`44 Typological linguistics

811 Languages

Derived from the common auxiliaries of language =1/=9 (Table 1c) by replacing the equal sign = with prefix *811*. e.g. =111 *English* becomes *811.111 Linguistics of English language*

811.1/.9 All languages natural or artificial

811.1/.8 Individual natural languages

811.1/.2 Indo-European languages

811.21/.22 Indo-Iranian languages

811.3 Dead languages of unknown affiliation. Caucasian languages

811.4 Afro-Asiatic, Nilo-Saharan, Congo-Kordofanian, Khoisan languages

811.5 Ural-Altaic, Palaeo-Siberian, Eskimo-Aleut, Dravidian and Sino-Tibetan languages. Japanese. Korean. Ainu

811.6 Austro-Asiatic languages. Austronesian languages

811.7 Indo-Pacific (non-Austronesian) languages. Australian languages

811.8 American indigenous languages

811.9 Artificial languages

82 Literature

82-1/-9 Special auxiliary subdivision for literary forms, genres

82-1 Poetry. Poems. Verse

- 82-2 *Drama. Plays*
- 82-3 *Fiction. Prose narrative*
- 82-31 *Novels. Full-length stories*
- 82-32 *Short stories. Novellas*
- 82-4 *Essays*
- 82-5 *Oratory. Speeches*
- 82-6 *Letters. Art of letter-writing. Correspondence. Genuine letters*
- 82-7 *Prose satire. Humour, epigram, parody*
- 82-8 *Miscellanea. Polygraphies. Selections*
- 82-9 *Various other literary forms*
- 82-92 *Periodical literature. Writings in serials, journals, reviews*
- 82-94 *History as literary genre. Historical writing. Historiography. Chronicles. Annals.*

Memoirs

- 82.02/.09 *Special auxiliary subdivision for theory, study and technique of literature*
- 82.02 *Literary schools, trends and movements*
- 82.09 *Literary criticism. Literary studies*
- 82.091 *Comparative literary studies. Comparative literature*
- 821 *Literatures of individual languages and language families*

Derived from the common auxiliaries of language =1/=9 (Table 1c) by replacing the equal sign = with prefix 821. e.g. =111 *English* becomes 821.111 *English literature*

9 Geography. Biography. History[\[edit source\]](#)

Tables for Geography and History in UDC are fully faceted and place, time and ethnic grouping facets are expressed through combination with common auxiliaries of place (Table 1d), ethnic grouping (Table 1f) and time (Table 1g)

- 902/908 *Archaeology. Prehistory. Cultural remains. Area studies*
- 902 *Archaeology*
- 903 *Prehistory. Prehistoric remains, artefacts, antiquities*
- 904 *Cultural remains of historical times*
- 908 *Area studies. Study of a locality*
- 91 *Geography. Exploration of the Earth and of individual countries. Travel. Regional geography*
- 910 *General questions. Geography as a science. Exploration. Travel*
- 911 *General geography. Science of geographical factors (systematic geography). Theoretical geography*
- 911.2 *Physical geography*
- 911.3 *Human geography (cultural geography). Geography of cultural factors*
- 911.5/.9 *Theoretical geography*

- 912 Nonliterary, nontextual representations of a region
- 913 Regional geography
- 92 Biographical studies. Genealogy. Heraldry. Flags
- 929 Biographical studies
- 929.5 Genealogy
- 929.6 Heraldry
- 929.7 Nobility. Titles. Peerage
- 929.9 Flags. Standards. Banners
- 93/94 History
- 930 Science of history. Historiography
- 930.1 History as a science
- 930.2 Methodology of history. Ancillary historical sciences
- 930.25 Archivistcs. Archives (including public and other records)
- 930.85 History of civilization. Cultural history
- 94 General

Common Auxiliary Tables

Common auxiliaries of language. Table 1c

- =1/=9 Languages (natural and artificial)
- =1/=8 Natural languages
- =1/=2 Indo-European languages
- =1 Indo-European languages of Europe
- =11 Germanic languages
- =12 Italic languages
- =13 Romance languages
- =14 Greek (Hellenic)
- =15 Celtic languages
- =16 Slavic languages
- =17 Baltic languages
- =18 Albanian
- =19 Armenian
- =2 Indo-Iranian, Nuristani (Kafiri) and dead Indo-European languages
- =21/=22 Indo-Iranian languages
- =21 Indic languages
- =22 Iranian languages
- =29 Dead Indo-European languages (not listed elsewhere)
- =3 Dead languages of unknown affiliation. Caucasian languages

- =34 Dead languages of unknown affiliation, spoken in the Mediterranean and Near East (except Semitic)
- =35 Caucasian languages
- =4 Afro-Asiatic, Nilo-Saharan, Congo-Kordofanian, Khoisan languages
- =41 Afro-Asiatic (Hamito-Semitic) languages
- =42 Nilo-Saharan languages
- =43 Congo-Kordofanian (Niger-Kordofanian) languages
- =45 Khoisan languages
- =5 Ural-Altaiic, Palaeo-Siberian, Eskimo-Aleut, Dravidian and Sino-Tibetan languages. Japanese. Korean. Ainu
- =51 Ural-Altaiic languages
- =521 Japanese
- =531 Korean
- =541 Ainu
- =55 Palaeo-Siberian languages
- =56 Eskimo-Aleut languages
- =58 Sino-Tibetan languages
- =6 Austro-Asiatic languages. Austronesian languages
- =61 Austro-Asiatic languages
- =62 Austronesian languages
- =7 Indo-Pacific (non-Austronesian) languages. Australian languages
- =71 Indo-Pacific (non-Austronesian) languages
- =72 Australian languages
- =8 American indigenous languages
- =81 Indigenous languages of Canada, USA and Northern-Central Mexico
- =82 Indigenous languages of western North American Coast, Mexico and Yucatán
- =84/=88 Central and South American indigenous languages
- =84 Ge-Pano-Carib languages. Macro-Chibchan languages
- =85 Andean languages. Equatorial languages
- =86 Chaco languages. Patagonian and Fuegian languages
- =88 Isolated, unclassified Central and South American indigenous languages
- =9 Artificial languages
- =92 Artificial languages for use among human beings. International auxiliary languages (interlanguages)
- =93 Artificial languages used to instruct machines. Programming languages. Computer languages

(0...) Common auxiliaries of form. Table 1d

- (0.02/.08) Special auxiliary subdivision for document form*
- (0.02) Documents according to physical, external form*
- (0.03) Documents according to method of production*
- (0.032) Handwritten documents (autograph, holograph copies). Manuscripts. Pictorial documents (drawings, paintings)*
- (0.034) Machine-readable documents*
- (0.04) Documents according to stage of production*
- (0.05) Documents for particular kinds of user*
- (0.06) Documents according to level of presentation and availability*
- (0.07) Supplementary matter issued with a document*
- (0.08) Separately issued supplements or parts of documents*
- (01) Bibliographies
- (02) Books in general
- (03) Reference works
- (04) Non-serial separates. Separata
- (041) Pamphlets. Brochures
- (042) Addresses. Lectures. Speeches
- (043) Theses. Dissertations
- (044) Personal documents. Correspondence. Letters. Circulars
- (045) Articles in serials, collections etc. Contributions
- (046) Newspaper articles
- (047) Reports. Notices. Bulletins
- (048) Bibliographic descriptions. Abstracts. Summaries. Surveys
- (049) Other non-serial separates
- (05) Serial publications. Periodicals
- (06) Documents relating to societies, associations, organizations
- (07) Documents for instruction, teaching, study, training
- (08) Collected and polygraphic works. Forms. Lists. Illustrations. Business publications
- (09) Presentation in historical form. Legal and historical sources
- (091) Presentation in chronological, historical form. Historical presentation in the strict sense
- (092) Biographical presentation
- (093) Historical sources
- (094) Legal sources. Legal documents

(1/9) Common auxiliaries of place. Table 1e

- (1) Place and space in general. Localization. Orientation
- (1-0/-9) Special auxiliary subdivision for boundaries and spatial forms of various kinds*
- (1-0) Zones*

- (1-1) *Orientation. Points of the compass. Relative position*
- (1-11) *East. Eastern*
- (1-13) *South. Southern*
- (1-14) *South-west. South-western*
- (1-15) *West. Western*
- (1-17) *North. Northern*
- (1-19) *Relative location, direction and orientation*
- (1-2) *Lowest administrative units. Localities*
- (1-5) *Dependent or semi-dependent territories*
- (1-6) *States or groupings of states from various points of view*
- (1-7) *Places and areas according to privacy, publicness and other special features*
- (1-8) *Location. Source. Transit. Destination*
- (1-9) *Regionalization according to specialized points of view*
- (100) Universal as to place. International. All countries in general
- (2) Physiographic designation
- (20) Ecosphere
- (21) Surface of the Earth in general. Land areas in particular. Natural zones and regions
- (23) Above sea level. Surface relief. Above ground generally. Mountains
- (24) Below sea level. Underground. Subterranean
- (25) Natural flat ground (at, above or below sea level). The ground in its natural condition, cultivated or inhabited
- (26) Oceans, seas and interconnections
- (28) Inland waters
- (29) The world according to physiographic features
- (3) Places of the ancient and mediaeval world
- (31) Ancient China and Japan
- (32) Ancient Egypt
- (33) Ancient Roman Province of Judaea. The Holy Land. Region of the Israelites
- (34) Ancient India
- (35) Medo-Persia
- (36) Regions of the so-called barbarians
- (37) Italia. Ancient Rome and Italy
- (38) Ancient Greece
- (399) Other regions. Ancient geographical divisions other than those of classical antiquity
- (4/9) Countries and places of the modern world
- (4) Europe
- (5) Asia
- (6) Africa
- (7) North and Central America
- (8) South America

(9) States and regions of the South Pacific and Australia. Arctic. Antarctic

(=...) Common auxiliaries of human ancestry, ethnic grouping and nationality. Table 1f

They are derived mainly from the common auxiliaries of language =... (Table 1c) and so may also usefully distinguish linguistic-cultural groups, e.g. =111 English is used to represent (=111) English speaking peoples

(=01) Human ancestry groups

(=011) European Continental Ancestry Group

(=012) Asian Continental Ancestry Group

(=013) African Continental Ancestry Group

(=014) Oceanic Ancestry Group

(=017) American Native Continental Ancestry Group

(=1/=8) Linguistic-cultural groups, ethnic groups, peoples **[derived from Table 1c]**

(=1:1/9) Peoples associated with particular places

e.g. (=111:71) Anglophone population of Canada

"..." Common auxiliaries of time. Table 1g

"0/2" Dates and ranges of time (CE or AD) in conventional Christian (Gregorian) reckoning

"0" First millennium CE

"1" Second millennium CE

"2" Third millennium CE

"3/7" Time divisions other than dates in Christian (Gregorian) reckoning

"3" Conventional time divisions and subdivisions: numbered, named, etc.

"4" Duration. Time-span. Period. Term. Ages and age-groups

"5" Periodicity. Frequency. Recurrence at specified intervals.

"6" Geological, archaeological and cultural time divisions

"61/62" Geological time division

"63" Archaeological, prehistoric, protohistoric periods and ages

"67/69" Time reckonings: universal, secular, non-Christian religious

"67" Universal time reckoning. Before Present

"68" Secular time reckonings other than universal and the Christian (Gregorian) calendar

"69" Dates and time units in non-Christian (non-Gregorian) religious time reckonings

"7" Phenomena in time. Phenomenology of time

-0 Common auxiliaries of general characteristics. Table 1k

-02 Common auxiliaries of properties

- 021 Properties of existence
- 022 Properties of magnitude, degree, quantity, number, temporal values, dimension, size
- 023 Properties of shape
- 024 Properties of structure. Properties of position
- 025 Properties of arrangement
- 026 Properties of action and movement
- 027 Operational properties
- 028 Properties of style and presentation
- 029 Properties derived from other main classes

-03 Common auxiliaries of materials

- 032 Naturally occurring mineral materials
- 033 Manufactured mineral-based materials
- 034 Metals
- 035 Materials of mainly organic origin
- 036 Macromolecular materials. Rubbers and plastics
- 037 Textiles. Fibres. Yarns. Fabrics. Cloth
- 039 Other materials

-04 Common auxiliaries of relations, processes and operations

- 042 Phase relations
- 043 General processes
- 043.8/.9 Processes of existence
- 045 Processes related to position, arrangement, movement, physical properties, states of matter
- 047/-049 General operations and activities

-05 Common auxiliaries of persons and personal characteristics

- 051 Persons as agents, doers, practitioners (studying, making, serving etc.)
- 052 Persons as targets, clients, users (studied, served etc.)
- 053 Persons according to age or age-groups
- 054 Persons according to ethnic characteristics, nationality, citizenship etc.
- 055 Persons according to gender and kinship
- 056 Persons according to constitution, health, disposition, hereditary or other traits
- 057 Persons according to occupation, work, livelihood, education
- 058 Persons according to social class, civil status

Lesson 3: BSO - the Broad System of Ordering

BSO - the Broad System of Ordering - is a modern machine-held classification system embracing all fields of knowledge.

- Usage of BSO
- **Needling the haystack:** subject searching in wide angled information systems, such as the Internet
- **A systematic overview of knowledge:** what does it mean for subject searching?
- **BSO as a systematic overview of knowledge.**
- **BSO cycle of knowledge-the conceptual skeleton.**
- BSO cycle of knowledge.**Sector descriptions.**
- **Starting to clad the skeleton**
- **Adding more cladding**
- **Full detail excerpt**
- **Index excerpt**

Usage of BSO

1. As a fall-back aid to subject searching on the Net or any miscellaneous compilation or collection covering many subject fields
2. As a subject tagging code applied to individual items or records in wide angle collections or compilations. In this application the result is an orderly and easily grasped subject arrangement of the items
3. As a mediating tool in changing over from one subject indication system to another

Needling the haystack: subject searching in wide angled information systems, such as the Internet

When your search question dives straight into the haystack and instantly retrieves the sought needle, then you don't need BSO. At other times searching in a large information store or network may all too readily bring the needle-in-haystack problem to mind. You may draw a blank or be presented with an offering which is not quite what you were looking for. At such times BSO could often help in suggesting alternative search approaches. Similarly, if you are not quite sure of the appropriate subject word to begin your subject search with, a glance at the index and systematic schedules of BSO could help to set you on your way.

A systematic overview of knowledge: what does it mean for subject searching?

When you think of the total universe of knowledge as covered by Internet do you see it as an intricate maze of subjects in a tangled net of relationships, or as an orderly map enabling one to navigate with confidence? Most subject searchers fall somewhere between these two extremes. Real polymaths are rare. On the other hand, most of us with interest or close involvement with particular subject areas tend to think about the contents of these limited areas in a more or less structured way. Problems begin when our search needs take us to the very fringes of our familiar subject areas and beyond.

BSO as a systematic overview of knowledge

There are all sorts of possible overviews of knowledge, systematic at various levels and in varying degrees. BSO is a compilation of about 6800 terms, arranged in an order which is systematic or structured at the level of meaning. It is a knowledge classification which attempts to reflect a modern consensual world outlook, and because it is essentially concept-oriented rather than word-oriented, its captions are supplemented by definition notes whenever the scope and meaning of the caption is not clear from its context in the systematic schedules. The system gives solutions to many dilemmas met by searchers looking for "composite" subjects which can be represented in language only by phrases. Familiarisation with the system is made easy by the fact that its arrangement is governed by a repeated and widely applied underlying pattern which enhances its overall simplicity

BSO Cycle of knowledge - the conceptual skeleton

To get a quick general idea of the way BSO arranges the broadest categories of knowledge read the brief skeletal **sector descriptions** below in clockwise order starting at the bottom left text block (marked with **) and proceeding upwards block by block to the top, then across to the top right block, and then downwards.

View the same cyclic arrangement with the **skeleton clad** with subject names.

View **more detail** in the linear form of an indented hierarchical outline subject list, derived by cutting and unrolling the cycle.

BSO Cycle of knowledge. Sector descriptions.

Read blocks of text clockwise. Start at **

: 410 to 480		-----
: Human beings	>->	: 500 to 588
as		: Humanities and
subject of		: society
knowledge		

: 359 to 397
: Human interaction
with
: non-human
organism

: 600 to 899
: Material products of
: human societies,
: Artefacts

_____ to 345
 :Natural sciences in
 order of
 :levels of
 increasing complexity

**

_____ to 188
 : Formal and
 general sciences
 : techniques and
 activities
 : which are
 utilised in
 : connection with
 a variety of
 : more
 specialised areas
 of
 : knowledge
 pursuit and
 : handling

_____ to 397
 : Non-utilitarian
 : products of human
 societies,
 :Mentefacts
 :(Language, Arts,
 :Religion, Esoterica)

^

BSO cuts cycle here

Starting to clad the skeleton

A linear sequence of classes is obtained by cutting into this cyclic structure at a chosen point. The classes separated at the point of division are the first and last classes, of the linear sequence. The BSO point of division is at the base of the diagram below

480 Sports/games	>Humanities/social studies	500
470 Human needs	History/related sciences	510
460 Education	Area studies	520
450 Psychology	Society	527
445 Behavioural sciences	Social sciences	530
420 Medicine	Sociology	535
410 Biomedical sciences	Demography	537
390 Environment	Politics	540
380 Wildlife exploitation	Public administration	550
370 Forestry	Law	560
366 Animal husbandry	Social welfare	570
360 Agriculture	Economics	580
359 Applications of life sciences	Enterprise management	588
340 Zoology	Technology	600
330 Botany	Production technology	620
320 Microbiology	Materials handling	625
310 Biological sciences	Packaging/storage	627
300 Life sciences	Energy technology	631
290 Geography	Materials technology	635
270 Geology	Nuclear technology	640
260 Earth sciences	Electrotechnology	650
250 Space and earth sciences	Thermal engineering	670
230 Chemistry	Mechanical engineering	680
228 Crystallography	Construction technology	710
210 Physics	Environmental technology	730
205 Physical sciences	Transport technology	740
203 Natural sciences	Military sci./technology	760
200 Science and technology	Mining	780
188 Metrology	Process industries	800
186 Testing and trials	Metal technology	860
182 Research	Wood/pulp/paper technology	871,95
166 Standardisation	Textiles technology	877
165 Management	Particular products manf.	890
160 Systemology/cybernetics	Language/literature	910
150 Communication sciences	Art	940
140 Information sciences	Religion/atheism	970
120 Mathematics	Esoteric practices	992
118 Logic		
** 112 Philosophy		

Adding more cladding

- 088 PHENOMENA & ENTITIES FROM A MULTI- or NON-DISCIPLINARY POINT OF VIEW
- 100 KNOWLEDGE GENERALLY
- 112 PHILOSOPHY
- 116 SCIENCE OF SCIENCE
- 118 LOGIC
- 120 MATHEMATICS
- 125 STATISTICS & PROBABILITY
- 127 DATA PROCESSING
- 128 COMPUTER SCIENCE
- 140 INFORMATION SCIENCES
- 143 Libraries
- 144 Archives
- 146 Museums
- 146,80 Exhibitions
- 148 MEETINGS
- 148,80 CONSULTANCY
- 148,85 INTERVIEWING
- 150 COMMUNICATION SCIENCES
- 152 Reprography & printing
- 152,60 Book trade
- 152,80 Intellectual property
- 155 Destination-directed communication
- 156 Mass communication
- 158,20 Publicity
- 160 SYSTEMOLOGY & CYBERNETICS
- 163 Operations research
- 165 MANAGEMENT
- 166 STANDARDISATION & STANDARDS
- 168 ORGANISATIONS
- 182 RESEARCH
- 184 DISCOVERIES, INVENTIONS & PATENTS
- 186 TESTING & TRIALS
- 188 METROLOGY

- 200 SCIENCE & TECHNOLOGY (TOGETHER)
- 203 NATURAL SCIENCES
- 205 PHYSICAL SCIENCES
- 210 Physics
- 212 Energy interactions & forms
- 214 Particle & high-energy physics
- 215 Nuclear physics
- 217 Atomic, molecular & ion physics

219	Vacuum physics
222	Bulk matter physics
224	Plasma & fluid physics
225	Condensed matter physics
226	Physics of solids
228	Crystallography
230	Chemistry
232	Physical chemistry
234	Chemistry of particular substances
235	Inorganic chemistry
237	Organic chemistry
238	Polymer chemistry
250	Space & earth sciences
252	Astronomy & astrophysics
258	Space research
260	Earth sciences
262	Geodesy & surveying
263	Geophysics
265	Atmospheric sciences
267	Hydrospheric sciences
270	Geology
290	Geography
300	LIFE SCIENCES
310	BIOLOGICAL SCIENCES
312	Biophysics & biochemistry
313,20	Molecular biology
313,30	Cell & tissue biology
313,70	Genetics
313,75	Evolution
315,19,34	Biological materials
315,20	Developmental biology
315,54	Morphology
315,55	Physiology
315,56	Pathology
315,58	Immunology
315,59,30	Biological material structures
315,60	Biological parts, organs & functional systems
318	Ecology
318,50	Ethology
320	Microbiology
330	Botany
340	Zoology
359	APPLICATIONS OF LIFE SCIENCES
360	AGRICULTURE
363	Plant crop production
363,70	Horticulture

- 363,80 Production of specific plant crops
- 366 Animal husbandry
- 368 VETERINARY SCIENCE
- 370 FORESTRY
- 380 WILDLIFE EXPLOITATION
- 390 ENVIRONMENT
- 397 Natural resources
- 410 BIOMEDICAL SCIENCES
- 420 MEDICINE
- 422 Preventive medicine
- 423 Social medicine
- 425 Clinical & internal medicine
- 426 Surgery
- 428 Diseases
- 430 PARTS, ORGANS & SYSTEMS OF THE HUMAN BODY
- 432 Body parts & regional specialties
- 433 Body organs & systems
- 433,50 Defence system
- 434 Integumentary & musculoskeletal system
- 435 Visceral systems & organs
- 437 Nervous system & sense organs
- 438 Mental health & disorders
- 439 BIOMEDICAL SPECIALTIES, BY HUMAN SUBJECT OR PATIENT,
& BY ENVIRONMENT
- 445 BEHAVIOURAL SCIENCES
- 450 PSYCHOLOGY
- 460 EDUCATION
- 470 HUMAN NEEDS
- 475 Household science
- 477 Work & leisure occupations
- 477,50 Leisure & recreation
- 478 Tourism & travel
- 480 Sports & games
- 482 Individual prowess & athletic sports
- 483 Ball games
- 485 Special environment sports
- 486,40 Wheel vehicle sports
- 486,50 Animal sports
- 486,70 Target & quarry sports
- 488,30 Board & piece games
- 489 Social diversions & pastimes

- 500 HUMANITIES & SOCIAL STUDIES
- 510 HISTORY & RELATED SCIENCES
- 512 Archaeology & prehistory
- 513 History of particular epochs & periods

515	History of main world areas & continents
516	History of particular countries (including regions & localities within a particular country)
520	AREA STUDIES
526	Area studies of particular countries (including regions & localities within a particular country)
527	SOCIETY
528	SOCIAL GROUPS & COMMUNITIES
529	Ethnic & linguistic & religious groups
530	SOCIAL SCIENCES
532	Culture
533	CULTURAL ANTHROPOLOGY
535	SOCIOLOGY
537	DEMOGRAPHY
540	POLITICAL SCIENCE & POLITICS
542	Political institutions & organisations
543	Political organisational patterns & systems
544	Political history
545	Politics of particular groupings of states
546	Politics of particular states & countries
550	PUBLIC ADMINISTRATION
560	LAW
562	Civil law
563	Public law
565	International law
567	Systems of law (by origin)
568	Law of particular countries
570	SOCIAL WELFARE
580	ECONOMICS
581,80	Microeconomics
582	Macroeconomics
584	Economic organisation
586	Sectorial economics
588	MANAGEMENT OF ENTERPRISES
600	TECHNOLOGY
610	SCIENTIFIC BASIS OF TECHNOLOGY
611	EQUIPMENT & PLANT
612	SYSTEMS ENGINEERING
612,55	COMPUTER TECHNOLOGY
615	TECHNICAL TESTING
617	MAINTENANCE & SERVICING ENGINEERING
618	TECHNICAL & INDUSTRIAL DESIGN
620	PRODUCTION TECHNOLOGY
625	MATERIALS HANDLING
627	PACKAGING & STORAGE & DISPATCH

631	ENERGY TECHNOLOGY
635	MATERIALS TECHNOLOGY
640	NUCLEAR TECHNOLOGY
642	Nuclear reactor technology
645	Isotope technology
650	ELECTRONIC & ELECTRICAL TECHNOLOGIES
653	ELECTRONIC ENGINEERING
655	TELECOMMUNICATION ENGINEERING
657	ELECTRICAL ENGINEERING
670	THERMAL ENGINEERING & APPLIED THERMODYNAMICS
673	Heat engines
678	Refrigeration technology
680	MECHANICAL ENGINEERING
684	FLUID ENGINEERING
686	VACUUM TECHNOLOGY
688	VIBRATION & ACOUSTIC ENGINEERING
710	CONSTRUCTION TECHNOLOGY
712	CIVIL ENGINEERING
714	ILLUMINATING ENGINEERING
716	BUILDING CONSTRUCTION & SERVICES
716,70	Types of buildings
720	ARCHITECTURE
724	LANDSCAPE DESIGN
726	PHYSICAL PLANNING
730	ENVIRONMENTAL TECHNOLOGY
730,30	Pollution control
734	Public & industrial health engineering
736	Safety engineering
738	Rescue & salvage operations
740	TRANSPORT TECHNOLOGY & SERVICES
742	Road transport technology & services
743	Railway transport technology & services
745	Water transport technology & services
747	Air transport technology & services
748	Space transport technology
760	MILITARY SCIENCE & TECHNOLOGY
764	Warfare
767	Weapons in warfare
780	MINING
782	Solid fuel extraction
783	Metal mining
784	Non-metallic mineral mining
786	Oil & gas extraction technology
788	Ore & mineral dressing
800	PROCESS INDUSTRIES
810	Chemical technology & engineering

811	Chemical engineering
812	Chemical technology
813	Chemical agents & basic industrial chemicals
815	Industrial gases technology
816,90	Biotechnology
823	Technology of particular groups of chemicals
825	Petroleum technology
826	Natural oils, fats & waxes technology
828	Polymer technology
831,30	Chemical agents for materials processing
831,40	Surface-active agents
831,44	Adhesives & sealants
831,47	Lubricants
831,60	Surface finishing agents
832	Fuels & explosives technology
834	Colour industry technology
836	Pharmaceuticals & related technologies
838	Agricultural chemicals technology
840	FOOD & DRINK TECHNOLOGY
841	Food technology
845	Drinks technology
847,50	TOBACCO PROCESSING
849,13	MINERAL PROCESSING TECHNOLOGY
850	Non-metallic mineral technologies
852,32	Lime & lime products
852,34	Gypsum & anhydrite products
852,50	Cement & mortar technology
852,60	Concrete technology
854	Ceramics & clayware technology
856	Glass & glass ceramics technology
858	Fibrous & layered silicate mineral technologies
860	METAL TECHNOLOGY
864	Metal products
865,50	Metallurgy of particular & kinds of metals
865,60	Alloys
866,60	Ferrous metallurgy
867	Non-ferrous metallurgy
871,95	WOOD, PULP & PAPER TECHNOLOGY
872	Wood technology
873	Pulp & paper technology
875	LEATHER & OTHER ANIMAL PRODUCTS TECHNOLOGY
877	TEXTILES TECHNOLOGY
878	CORDAGE & WIRE ROPE MAKING
890	MANUFACTURE & TECHNOLOGY OF PARTICULAR PRODUCTS NOT SCHEDULED IN BSO AREA 600 TO 878

- 910 LANGUAGE & LITERATURE
 - 911 LINGUISTICS
 - 912 USE OF LANGUAGE
 - 912,20 Authorship
 - 912,30 Reading
 - 915 LITERATURE
 - 920 SPECIAL PHILOLOGICAL STUDIES
 - 921 Indo-European languages & literatures
 - 923,20 Afro-Asian languages & literatures
 - 923,60 Caucasian languages & literatures
 - 923,70 Basque language & literature
 - 924 Eurasian & North Asian languages & literatures
 - 925,20 Dravidian languages
 - 925,51 Sino-Tibetan languages & literatures
 - 925,70 Austronesian & Oceanic languages & literatures
 - 927 African languages & literatures
 - 928 Amerindian languages & literatures

 - 940 ARTS
 - 943 PLASTIC ARTS
 - 945 GRAPHIC FINE ARTS
 - 947 PHOTOGRAPHY AS ART
 - 949 DECORATIVE ARTS & HANDICRAFTS
 - 950 MUSIC & PERFORMING ARTS
 - 951 Music
 - 952 Vocal music
 - 953 Instrumental music
 - 955 PERFORMING ARTS
 - 957 Cinema

 - 970 RELIGION & ATHEISM
 - 971 Atheism & rationalism
 - 972 Religion
 - 973,90 Particular religions
 - 974 Religions of Indian origin
 - 975 Religions of Far Eastern origin
 - 975,70 Religions of Iranian origin
 - 976 Judaism
 - 977 Christianity
 - 978 Islam
 - 979 Other religions & semi-religious cults

 - 992 ESOTERIC PRACTICES & MOVEMENTS
-

Full detail excerpt

The complete full detail BSO classification expands the preceding list by an overall factor of about 25. The following excerpt is extracted from the Biomedical Sciences area of BSO:

439,31,40	Foetal biomedicine
439,32	Paediatrics = Pediatrics = Infancy & childhood
439,32,12,33	Metabolism
439,32,13,70	Genetics
439,32,15,20	Development
439,32,20,38	Radiology
439,32,20,58	Immunology
439,32,22	Preventive paediatrics
439,32,22,60	Food health & hygiene
439,32,23,70	Hospitals = Children's hospitals
439,32,25,30	Diagnosis
439,32,25,40	Therapeutics
439,32,25,70	Anaesthesiology
439,32,26	Surgery
439,32,28	Diseases
439,32,28,41	Toxicology
439,32,28,43	Allergy & anaphylaxis
439,32,28,80	Communicable diseases
439,32,34,31	Dermatology = Skin biomedicine
439,32,34,50	Musculoskeletal system
439,32,34,75	Connective tissue
439,32,35,30	Alimentary system
439,32,35,31	Nutrition process
439,32,35,32	Dentistry
439,32,35,40	Respiratory system
439,32,35,51	Cardiovascular system
439,32,35,58	Lymphology
439,32,35,60	Endocrinology

Index excerpt

This excerpt is an arbitrarily chosen slice from the index to the full BSO classification. It is not an index to the terms in the preceding excerpt from the full BSO classification which do not fall within the alphabetic range Foetuses-Food, Health & hygiene.

Foetuses, Biomedical sciences 439,31,40
 Fog, Meteorology 265,67,49
 Folacin, Biochemistry 312,37,87,36,80
 Folacin, Biochemistry, Biomedical sciences 412,37,87,36,80

Folds, Geology 273,34,50
Folk & traditional art 940,52
Folk & traditional music 951,52
Folk literature 915,52
Folklore 533,60
Folklore, Nature 533,62
Food, Additives & flavouring materials, Technology 841,71
Food, Allergies, Medicine 428,43,60
Food, Animal derived, Human needs 471,51
Food, Animal husbandry products 366,81
Food, Contamination, Health & hygiene 422,60,40
Food, Crops 363,64
Food, Disease causation, Medicine 420,52,60
Food, Fishery products 387,87
Food, Flavouring materials, Technology 841,72
Food, Health & hygiene 422,60
Food, Health & hygiene, Genital system 435,75,22,60

BSO - ORIGIN AND PRE-DEVELOPMENT PHASE

Subject indication for an information network

'Subject indication' is the phrase used in this manual to refer to those facilities of an information system which enable it to be interrogated by queries which have a subject as their point of departure. The user supplies to the system the name of a subject with the aim of extracting information on that subject from the system's store. Frequently the system contains in its store, not the information ultimately required, but records of the names and addresses for documents in which the information is likely to be found. In such a case subject indication has as its aim the identification of *documents* carrying the required information. The tools or languages of subject indication include indexing languages, classification systems, controlled term or keyword lists and thesauri.

The Broad System of Ordering is such a subject indication language, or, more specifically, a classification system, developed for a proposed world-wide information network covering the whole field of knowledge. At first sight there appears to be a little reason for supposing that a subject indication language for a network should be fundamentally different from a subject information language for information system generally, and it is arguable that the schedules of *BSO* bear this out. However, there are areas of such uncertainty surrounding subject indication languages, that it would have been rash indeed not to have put the matter to the test by information requirements in view. It will be seen that the result of the exercise bears a family resemblance to some of the document classifications which have preceded it, despite the fact that during the exercise no reference or recourse was made to literary or documentary warrant in the direct sense. Whether the differences between *BSO* and the document classifications are considered significant or trivial in themselves, they could possibly prove to be essential to the network application.

Origin of BSO

BSO originated in the context of the idea which emerged in the 1960s that consideration should be given to the possibility of a global network of scientific information centres, taking into account particularly the needs of developing countries. The network idea was itself triggered by a technological development, not at that time generally available, but certainly upon the near horizon. This was the possibility of cheap and fast data transmission links. It is notable that thinking about the information network, involving the first steps towards system definition began about a decade before the hardware became generally available. This was an honourable exception to the more usual situation in the mechanisation of information services in which the computer hardware was available well advance of system planning.

The subject indication sub-system of the network was seen as an important part of the whole system. It was vital in such a network that information on the subjects of documentary resources held by any one participating centre should be accessible to all the centres in the network. There were two closely interlocked, but still separable, problems here. The first was that, despite - and perhaps because of - the growth of mechanisation of information services, which in the late 1960s was just getting under way on a substantial scale, a greater amount of subject indication activity depended upon human intuitive skill and know-how than in the pre-mechanisation period ending about 1955. There was, for instance, an unprecedented proliferation of controlled keyword lists and thesauri, for use with mechanised systems, but little sign of common logical rationale in their construction which might otherwise itself be amenable

to mechanisation. These human skills produced indexing tools of great diversity for particular subject areas. Did these tools, the detailed construction principles of which were usually not fully communicable, offer a suitable model for the subject indication language of the proposed network?

The second problem arose for the realisation that for often good and sufficient reasons, centres representing the various subject fields would continue to use a variety of subject indication languages, corresponding to a variety of needs. Accordingly, communication through the use of a possible standard indexing language, which all participating centres would use for subject description of documents, was ruled out. On the other hand a solution seemed to lie in a procedure whereby subject information coded in one local indexing language could be converted by clerical means into the codes of another language conveying the same subject information.

Switching indexing languages

It so happened that this solution involving interconnection of individual local indexing languages, by a mediating or switching language had been under study by the Groupe d'Etude sur l'Information scientifique based at Marseilles, since 1963. Unlike the proposed global scientific information network, the system envisaged by GEIS was composed of centres dealing with the same subject discipline, and the particular discipline used as a study sample was the Science of Scientific Information itself. This difference is of some importance when considering the transfer *en bloc*, of the conclusions of GEIS in their 'Intermediate Lexicon' project to the context of the global scientific information network. A key feature in the GEIS scheme was the 'equivalence' or 'conversion' table in which the code for a given concept as rendered in one indexing language was coupled with the code of the same concept in another indexing language. It was assumed that such coupling was practicable to the extent that both indexing languages were, in fact as well as in pretension, lists of terms each of which corresponded with a definite and unambiguous concept. In fact the term 'indexation' was reserved for the process of concept analysing a document and assigning a code accordingly. (The code could be a notation symbol of a classification, or a descriptor or authorised term drawn from a thesaurus or subject heading list).

For the simplest case of a network in which the participating centres taken in aggregate used only two local indexing languages, all that was required was a pair of 'equivalence tables', one leading from language A to B, and the other from language B to A. If there were more than three indexing languages represented in the network, then it became more economical, in terms of the number of pairs of 'equivalence tables' required, to employ what has been variously termed as switching language, a mediating language, or a communication indexing language. A message (i.e. a subject request, or the answer to a subject request) would thus proceed from centre A (using local indexing language A) via an 'equivalence table' to the switching language, and then outward to a further 'equivalence table' reaching its destination, centre B, coded in the form in which the same subject is rendered in local indexing language B. This system, which is exactly analogous to a telephone network, would require one pair of 'equivalence tables' (one subscriber's line in the telephone analogy) between each centre's local indexing language and the switching language, ; whereas, if there were no switching language employed, each centre would need to construct, and of course maintain, as many pairs of 'equivalence tables' as there were languages in use in the network, minus one. It has been mentioned that such a system was expected to work, subject to the condition that the various local languages were in fact concept controlled. By the same token, the switching language itself would need to be one in

which each representation (notation symbol, or term) corresponded to one, and one only, concept; and in which for every concept there was one, and one only preferred representation. The form and arrangement of the switching language could be considered a function of the kind of use for which it was intended. If it were required only for simple matching, its arrangement would be immaterial to those engaged in day-to-day operations of sending and receiving messages. If additionally it were intended to employ such a switching language for hierarchical search, it would be necessary to incorporate the necessary hierarchical linkages into the language. However from the point of view of constructing and updating a switching language under controlled vocabulary conditions, some form of schematic arrangement, on the lines of a classification, is probably mandatory. This schematic arrangement might not, however, be the form in which the language was most conveniently held in a computer store.

It is to be supposed that the idea emanating from GEIS (which was later given quantitative elaboration in a research study carried out at the Polytechnic of North London School of Librarianship) because it was the only one of its kind available, must have affected the thinking of members of the first bodies charged with the task of considering the global science information programme, when they turned to the question of subject indication.

This text is followed by the last section of the chapter *Steps towards clarification* (pp. 3-8)

CHAPTER 2: APPLICATION OF BSO

The limits of broadness

The first task of the FID/SRC Working Group was that of sharpening the somewhat indefinite terms of the remit entrusted to it. The central question here was to try to decide in the most concrete possible manner what was to be understood by broadness as a feature of the proposed *Broad System of Ordering*. What should the determining principle be, which would cause some terms to be included in the scheme as sufficiently broad, and others to be rejected on the ground that they were too specialised?

Arithmetical Approach

Several possible approaches to an answer to this question could be foreseen. The answer could be purely arithmetical - a stated total number of terms in the system could be settled in advance. The answer could be based on a particular property of terms in relation to the classification structure yet to be devised - namely the hierarchical level of the term within the structure. Or, it could be based upon some inherent semantic property which a term might or might not possess. Or, yet again, it could be based upon some formal linguistic property of terms. Finally, it might be possibly based upon some sufficiently objective social property or phenomenon associated with the term or with the concept denoted by it. An obvious thought here was that a social property useful as marking cutoff of detail might well be one closely related with the purpose which it was hoped BSO would serve.

The first approach to the problem of defining broadness, or cut-off point, for BSO - the laying down in advance of the total number of terms to be used - had some special attractions in relation to cost predictability, especially in the context of mechanised exchange of information within a network. Clearly the cost of the computer processing of such exchange would fairly closely depend upon the size of the interconnection language to be traversed in the passing of each message. This dependence is probably less significant now (1978) than at the outset of

BSO development in 1973. With the expected future use of microprocessor elements as customary computing hardware, it is likely to be even less significant in future. At the beginning of the development of BSO, it was provisionally assumed that the full scheme might contain 2000 terms, and in the first draft submitted for comment in 1975 there were in fact 2100 terms. This draft was faulted both on account of its omissions and on account of alleged over-development of detail. Such criticism on mutually opposed grounds might have been crudely interpreted as a justification for the middle position taken by the BSO draft. However, the tenor of the comments themselves pointed to a great weakness of any solution to the cutoff problem based upon a prescribed maximum number of terms. As the approach to the maximum is reached, the question of what is, or is not, to be included in the system, comes to depend upon refined judgments of the relative importance of candidate subjects. Reliability in such judgments or judgments reflecting a real consensus are hardly to be expected from practitioners in the borderline specialties themselves - these specialties are, after all, often in competition among themselves for social recognition and funding. For this reason any purely arithmetical characterisation of cutoff in terms of the total number of subject-terms which the system is to contain is likely to be unsatisfactory and tendentious.

A brief side-glance at the arithmetical size of the scheme as a result of a cutoff criterion to be described later may be in order here. The 3rd revised draft of BSO (1978) contains about 4000 terms. The 18th edition of DDC, an established general classification for books has about 80,000 terms, so in approximate terms an average 'broad block' of information which can be designated by BSO is 20 times 'broader' than typical information at book level, and 3 times 'broader' than the information units which can be designated by the abridged UDC.

Hierarchical Approach

Many comments received on the earlier drafts of the scheme assumed without question that cutoff could appropriately be defined by reference to some-hierarchical level in the scheme. The arguments against such a basis for setting the limits of detail of the scheme are formidable. It can be contended that the policy on limit of detail, far from being derivative from the exigencies of the structure of the ordering system itself, should be independent of that structure. The structure is for the purpose of ordering, not for delimitation of acceptable detail. Furthermore hierarchical level of a given term is one of the most unstable features of all classifications in face of necessary changes required by the arrival of new knowledge. Much new knowledge arises by the fusion, following the discovery of common properties, of two or more hitherto separate subjects on the same hierarchical level. Whenever this occurs the separate subjects and all other subjects subsumed by them change their hierarchical level. Another consideration is that a statement of a hierarchical level is often made for explanatory or presentational purposes (for example in BSO 212 ENERGY INTERACTIONS & FORMS (ANY STATE OF MATTER). Obviously alternative presentational strategies are possible, and they will to some extent depend on available type variations for display. Also a chosen strategy may at some time have to be modified because of the appearance of a new subject remote from the hierarchical statement in question. Hierarchical levels are thus determined both by logical imperatives and presentational nuances. Both factors are subject to necessary change, and their states at a particular moment should not be the determinants of system detail cutoff. Finally the practical importance of a subject is by no means necessarily correlated with its hierarchical level. For instance 923,70 BASQUE LANGUAGE, being a unique member of a set is on the same hierarchical level as 921 INDO-EUROPEAN LANGUAGES.

The lack of agreement between natural languages as to the incidence of 'logies' and 'graphies' probably reflects the fact that mental organisation - the central characteristic of the kind of knowledge which constitutes a discipline - is not an all-or-nothing property. While one can perceive intuitively that, for instance, Chemistry is a more highly organised system of thought than Reprography, this is not to say that Reprography is not a discipline. Indeed, it would be quite hard to identify any subject matter which has generated literature, which can confidently be said to possess zero mental organisation. On the practical plane of handling subjects found in documents, no hard and fast line can be drawn between disciplines and non-disciplines from the standpoint of mental organisation of the material. It is only a question of more or less, and no simple method of scaling this spectrum of more or less mental organisation was available which would have enabled the FID/SRC Working Party to apply a quantitative criterion for cutoff of detail.

This text is followed by the sections:
Linguistic Approach: Subjects v. subject fields (pp 10-11),
Subjects as targets of organised information sources (pp. 11-13)

CHAPTER 3: DEVELOPING BSO - COLLECTING, STRUCTURING AND FEEDBACK (pp. 14-22)

CHAPTER 4: THE FIELD TEST OF BSO (pp. 23-33)

CHAPTER 5: APPLICATION OF BSO

The primary purpose for which BSO has been compiled is to serve as an exchange or switching language for use in an information network covering all subjects and in principle extending to users anywhere in the world.

Concept representation as the basis of switching

Behind the surface idea of subject indication switching between different indexing languages lies the assumption that despite the fact that individual centres participating in a network may differ from one another in the formalisms of their local indexing languages, there is between them an underlying agreement as to the nature and relations of the concepts represented in the local indexing languages. In other words, diversity belongs to the plane of language and terminology, but agreement to the plane of thought and idea. Switching is accordingly feasible on the plane of thought and idea on which agreements exists.

Different sets of indexing terms, descriptors, or notation symbols, used in different indexing languages to represent the same idea can be made to switch their idea-content between centres, provided that

- a) each local indexing language consists of terms and symbols, each of which is the sole representation, in the language, of a particular idea, and also represents that idea alone
- b) some neutral representation of the idea, agreed by all concerned, becomes the medium for clerical linkage for switching purposes. The neutral or switching language of concept representation must, like the local languages involved in the switching process be a controlled language.

Does this mean that a centre using free-text indexing cannot participate in switching? The answer is that in formal terms such a centre could participate, but practically it is unlikely to do so, because, in preparing the necessary concordance tables between its own input and the switching language, it would need to embark upon a vocabulary control exercise no less onerous than the control of the local indexing language itself: this is, however, the burden from which free-text indexing seeks to escape.

Form of switching language

The next question which arises is: what form of controlled indexing language is appropriate for the switching duty? Should it be arbitrary identifying code, a thesaurus, or a classification? An arbitrary code which carries no implicit or explicit information upon relations between vocabulary control itself - namely, the selection of codes to represent concepts uniquely - depends upon prior process of clustering concepts in order to establish near relationships and actual identity. An arbitrary code is no aid to such clustering. ON the other hand thesauri and classifications do display semantic relations - relations between ideas on the plane of meaning.

The choice between universal classification and universal thesaurus for the switching language role follows from the manner in which each displays relationships. A classification attempts to display relationships as a totality by means of tabulation. A thesaurus depicts relationships in a fragmentary manner, in the form of binary linkages, each of which is probably separated from semantically 'next neighbour' binary linkages by the accident of the alphabet. There is very little question as to which manner of relational display is the more useful for the purpose of controlling the vocabulary in face of an incoming flow of candidate new terms. Indeed, it is becoming increasingly common for thesauri themselves to supplement the fragmentary manner of showing semantic relationships, by adding to the alphabetical sequence of keywords, ancillary sections of grouped, categorised, or fully classified terms. Indeed it is becoming increasingly common for thesauri themselves to supplement the fragmentary manner of showing semantic relationships, by adding to the alphabetical sequence of keywords, ancillary section of grouped, categorised, of fully classified terms. Fro a small thesaurus the clustering process essential to vocabulary control in admitting new terms may be undertaken informally as a purely mental activity. If the thesaurus is large and of wide subject scope, then reliable and economic control of the vocabulary requires that the clustering should be externally formalised as a classification structure. It has been argued earlier that the practicability of a universal switching language depends critically upon its ability to be controlled, revised, and updated with minimum effort. A classification , more than any other form of indexing language, is amenable to easy, predictable, yet at the same time fully controlled updating. This is the essential ground upon which it is the preferred form of indexing language for the universal switching application. That existing universal classifications have failed, or are visibly failing, precisely in this respect does not vitiate the argument. The theoretical developments in classification of the last half-century have been preferentially applied to special subject classifications. BSO is in one sense an attempt to bring many of these developments into the sphere of general classification. It seems likely that these developments, all in the direction of bringing pervasive structural patterns into general classification, may hold the key to resolving the updating/keeping-up-with-knowledge problem which besets the established systems of universal classification and their users.

New knowledge, new technology and universal classification

On the broadest perspective, the UNISIST requirement of a classification, covering all fields, for exchange or switching purposes, may be seen as a particular concrete manifestation of a more general new need for a universal classification which has emerged only in the present decade. This need has arisen from the conjunction of three separate factors. The first of these concerns the process by which growing points of new knowledge often appear astride of discipline boundaries, and in aggregate have the effect of diminishing the practical significance of these boundaries. This process has been well recognised for many years but its impact has only recently been fully felt. The fringe or marginal subjects of specialised information services are spreading ever more widely over the total field of knowledge. It is not only that some of the socially more significant of the new technologies are of mixed scientific parentage. There is at present a considerable emphasis on what may be termed holistic approaches to all departments of human affairs. The ground 'between' technology, economics and apparently more distantly related social sciences is at present receiving unprecedented attention, as may be seen from the appearance of such interdisciplinary information services as SPLINES. Equally the boundaries between technology and social sciences have become blurred by the integrated concept 'Environment' which ultimately stems from the realm of biology and psychology. The rise of this holistic standpoint has on the one side strained the capacity of the established general classifications for accommodation to near breaking point, and on the other stimulated a new need for a universal classification.

The second factor contributing to a new need for universal classifications is directly technical in character. The limitations of clerical manual methods of manipulation and transfer of information records tended to confine such activities to single disciplines, within which quantities of material to be processed were sometimes manageable. Electronic data processing has vastly relaxed these limitations, and accordingly the significance of the discipline boundaries themselves has relaxed.

The third factor contributing to a new need for universal classifications is directly technical in character. The limitations of clerical manual methods of manipulation and transfer of information records tended to confine such activities to single disciplines, within which quantities of material to be processed were sometimes manageable. Electronic data processing has vastly relaxed these, limitations, and accordingly the significance of the discipline boundaries themselves has relaxed.

The third factor leading to a renewed need for a universal classification has been the internationalisation of information processing activities. Access to information is far less than hitherto the prerogative of advanced countries alone. In the developing countries there is at present great activity in the setting up of information centres covering all fields of knowledge, and collecting or arranging access to information from all sources. This flow of information on a global scale has re-animated the whole issue of a universal classification, particularly in its role for indicating the nature of the subject-content of information requests and documents.

BSO for switching and mediating

All of the three above factors are clearly related to the universal switching language application of BSO. An information network should be capable of connecting centres individually oriented to different focal disciplines. Its practicability on a large scale depends substantially upon exploiting

data processing and transmission technology, and the associated switching language has to be capable of surmounting linguistic and cultural barriers.

There are other possible applications, essentially of the same operational type as the switching language, which may be envisaged for BSO. In all cases they are products of the first and third of the three general factors mentioned above which seem to require a new universal classification, but the second factor concerning the liberation of earlier restraints owing to data processing technology is generally less significant than in the switching language application and may be absent altogether. In most cases, though not exclusively, these additional applications involve users who potentially may be in any part of the world.

Networking is not the only context within which neutral mediating languages come into play. Considerable financial resources are at the moment being applied to the translation, harmonisation, and interconversion of thesauri. For the minority of thesauri which themselves are no more detailed than BSO it is possible to conceive of BSO as a clerical switching language. For larger thesauri involved in interconversion projects designed to give users of one indexing language access to documents indexed in a different language, the use of BSO is a mediating, or common reference, language would achieve the necessary preliminary clustering of related terms from both thesauri, and would provide a framework for the higher organisation of the formed clusters, which might not be carried through to the finished product, but in any case would be useful as a provisional concept-holding device while the conversion work was in progress. The advantage of this approach would be both to eliminate decision process in the preliminary clustering, and to enable the broadest view of the overall subject-structure of the thesauri to be available from a very early stage in the project. Thus costly looping back whereby an early decision has to be modified to conform with the implications of a decision taken later - very characteristic of piecemeal operations on a structure of which the integrity is for the time being invisible - would be eliminated. The preliminary clusters thus formed would of course require to be broken down further by human intelligence - this being the inevitable limitation of a 'coarse' ordering system.

Other applications

Another example of a possible application of BSO is as an aid in the routing operations of referral centres and clearinghouses in dealing with inquiries. Compared with the switching application, this use of BSO would exploit communications technology equally, but its involvement with data processing would be less sophisticated.

Ultimately serving the same purposes as the referral centre, but serving individual demand by the mass medium of an older form of communication - the printed world - is the comprehensive directory of specialist organisations and specialist information sources. From the point of view of subject indication, present standards in publications of this kind could be improved to the substantial benefit of users. Such improvement could be realised either by arranging the material by BSO codes or by providing and index from BSO codes to page or item serial numbers.

It is also possible to envisage the use of BSO in purely disseminative modes of communication. As a subject tag supplied on copies of distributed reports and separates of all kinds, BSO codes would serve recipients of this material both as a 'coarse' interest filter, and secondly as a temporary filing system both for purposes of retrieval and subsequent control of disposal of little used material.

In these latter applications, and in some others such as the possible use of BSO codes as subject indicators in machine readable records. BSO would to some extent be competitive with existing established general classifications. The seriousness of this competition would perhaps depend upon

a) the inherent advantages and disadvantages, input cost-wise and user-wise or relatively 'coarse' subject specification versus the more detailed specification aiming at book level or documentation level in the established schemes

b) the relative merits of BSO and the established schemes in providing unequivocal placing for subjects, and thus in ease of decision effort in indexing

c) achievement by BSO of a new style of updating arrangement which would permit prompt assimilation of new knowledge into the scheme at a cost to the user which would be found acceptable

One final question which arises here is whether BSO might conceivably in future infiltrate or invade the territory proper of the established document classifications. In other words, will it ever be used for shelving books or filing documents in libraries? The answer to this question depends upon established systems rather than upon BSO. All that can be said is that if the established systems are found wanting on either of the two issues labeled b) and c) in the foregoing paragraph, then this same question will doubtless be raised repeatedly. It is not entirely unusual for tools of this kind to be used for purposes other than those for which they were originally intended. Furthermore there is nothing in the design of BSO which would inhibit elaboration to a greater depth of detail.

CHAPTER 6: BSO - DESCRIPTION OF THE SCHEME

One System design and user effort

It has been suggested in the preceding chapter that a switching indexing language needs to be economical in usage. The benefits of networking are not obtainable entirely without cost. The indexing of material by a switching language at a centre would, after all, be an addition to indexing effort normally put forth for local purposes. It is therefore essential that the additional cost of communication with other centres in the network should not contain any unnecessary element. It is against this background that the question of the cost of BSO to the user, both in day-to-day operation and in making changes consequent upon changes in the content and structure of knowledge, has been a matter of primary concern at every step in designing the scheme.

A classification user's unnecessary costs arise mainly in two ways. First, day-to-day application of the scheme may demand more decision effort than is necessary. Second, the local implementation of update amendments to the scheme may involve unnecessary effort.

Unnecessary decision effort is the result either of gross mismatch between the subjects found in the material to which the classification is to be applied and the concepts represented in the classification itself, or to lack of structural homogeneity in the scheme itself. It should be noted that mismatch is the result not only of initial shortcomings of the scheme but also of delays in updating. Lack of structural homogeneity may be paraphrased as unnecessary complexity in the scheme due to absence of overall pattern. An example of an inhomogenous general

classification would be one which was prepared simply by bringing together the special classifications corresponding to each included subject area, and listing them sequentially (possibly in some logical or otherwise helpful order). Any discipline, almost by definition, represents a particular viewpoint. A series of classifications, each optimal for the needs of a particular viewpoint, form, when added together, a general classification of great complexity, and consequently demand excessive decision effort in being applied.

Unnecessary effort in implementing updating, both on the part of the updater and of the user, is demanded when the insertion of a new subject requires not only an addition to the schedule but also a re-notation of adjacent terms representing old knowledge. This may arise either because the area involved was in the first place inadequately structured or because of a constraint offered by the notation.

These considerations are reflected in the general features of BSO, which include a marked incidence of structured pattern, both within and transcending discipline boundaries. The system is also highly prescriptive. There are no alternative placings offered. Completely definitive and embracing procedures are laid down by which indexers deal with the necessary factor of cross-classification in the schedule, which is therefore expected to be non-ambiguous in use and predictable in updating.

Neutrality and value judgments

After the question of the economics of the system comes the matter of its neutrality. All special classifications reflect the special viewpoint partly inherent in the discipline concerned and partly conventional among specialists within the discipline. Likewise, all general classifications are vulnerable to the charge that they reflect some particular world outlook or philosophy. This is obviously a question with potentially serious implications for a scheme intended for global use. Like the material which will be subject to switching in the foreseeable future, BSO reflects in many ways the standpoint of European tradition and culture. Within this limitation, the compilers have tried to stand outside sectional philosophies and to avoid decisions which have sectional philosophical implications. It is perhaps necessary to insist that neither the hierarchical nor ordinal position of any term carries any implication as to the importance of the associated concept.

Outline of BSO

The outline of the system is as follows:

FIRST OUTLINE OF BSO

088	Phenomena & entities from a multi or non-disciplinary point of view	460 EDUCATION
		470 HUMAN NEEDS
		475 Household science
		477 Work & leisure
		480 Sports & games
100	KNOWLEDGE GENERALLY	
112	Philosophy	500 HUMANITIES, CULTURAL & SOCIAL

		SCIENCES	
116	Science of science	510	History
118	Logic	526	Area studies
120	Mathematics	530	Social sciences
128	Computer science	533	Cultural anthropology
140	Information sciences	535	Sociology
150	Communication sciences	537	Demography
160	Systemology	540	Political science & politics
165	Management	550	Public administration
182	Research	560	Law
188	Metrology	570	Social welfare
200	SCIENCE AND TECHNOLOGY	580	Economics
203	Natural sciences	588	Management of enterprises
205	Physical sciences	600	TECHNOLOGY
210	Physics	910	LANGUAGE, LINGUISTICS & LITERATURE
230	Chemistry		
250	Space& earth sciences		
300	Life sciences		
300/439	Application of life science	940	ARTS
360	Agriculture	943	Plastic arts
368	Veterinary science	945	Graphic fine arts
368	Forestry	949	Decorative arts & handicrafts
380	Wild life exploitation	950	Music & performance arts
390	Environment & natural resources	970	RELIGION & ATHEISM
410	Biomedical sciences		
445	Behavioural sciences		
450	Psychology		

It is perhaps instructive to compare this outline with those of two other systems to which it has some resemblance. The first of these is the draft arrangement of subject fields according to Object Areas, devised by Dr. Dahlberg and detailed on pages 16-17. As has been made clear in Chapter 3, this was undoubtedly a major germinal influence upon BSO. The second scheme played no special part in the development of BSO; rather, likenesses to it gradually emerged as BSO was progressively elaborated. This is the outline of the Bliss Bibliographic Classification, which has become of topical interest since the first volumes of the 2nd edition of this classification (BC2) appeared in 1977. While a draft outline of BC2 had been available during most of the BSO development period, no special significance was attributed to it during the Working Group's review of sources of terms.

All three schemes (BSO, Object Area Scheme, and BC2) are fairly similar down to point 300 in the BSO outline. Both BSO and BC2 begin with Generalia and Phenomena from a multi-disciplinary point of view, though they disagree as to which of these two should appear first. The preliminary sciences (112 to 188 in BSO) are virtually identical in content with a similar group in BC2, though the three subgroups into which these fields may be divided are arranged somewhat differently in the two schemes. Only one of these subgroups occupies a similar preliminary position (Area 1) in the Object Area Scheme. The other two subgroups are located in Area f1 to which there is nothing comparable either in BSO or in BC2.

Both BSO and BC2 have Science and technology (taken together) immediately following the above-mentioned 'preliminary' sciences; all three schemes have General technology remotely located from Natural Sciences and closely following Economics and Management of Enterprises. BC2 and the Object Area Scheme are alike in intercalating Physics-based technology and Chemistry-based technology with the respective broad sciences. The Object Area Scheme goes further in placing Mining and Metallurgy within the Area of the Cosmic and Earth Sciences. BSO, alone of the three, after assigning a place near Natural Sciences for Science and technology taken together and generally, places all Technology based on Physics and Chemistry in a position remote from Physical Sciences. This position corresponds to that of general Technology in BC2 and the Object Area Scheme.

All three schemes agree in forming a ladder of the sciences, exhibiting Henry Bliss's 'gradation of specialties' and corresponding fairly closely to an arrangement of entities illustrating integrative levels of increasing complexity. This sequence begins at Physics, passes through Chemistry, Space and Earth sciences to Biology and Medicine, Psychology and Education. BSO and the Object Area Scheme go a little further together by adding Sports, Games and Leisure after Education. In BSO alone Sports, Games and Leisure are comprehended under the more general idea of Human Needs.

In the Life Sciences area BSO and the Object Area Scheme are in approximate agreement in placing Agriculture and Animal Husbandry in the same general region as the parent sciences Botany and Zoology. BC2, on the other hand, separates Agriculture and places it under remotely located General Technology.

In the area of the Humanities and Social Sciences the divergence between the three schemes are greater than in the Natural Science area. BC2 and the Object Area Scheme both place some, but not all, of the Social Sciences before History. BSO reverses this, thereby keeping all Social Sciences together, with the exception of Education. The end of the BSO outline is generally similar to the corresponding part of the Object Area Scheme. BC2 is here sharply different from the other two schemes in placing Religion, the Occult, and Ethics before the more 'practical' Social Sciences, and also in locating Recreative Arts immediately before Fine Arts.

Taking now the broadest view of the order of the BSO outline, the 'preliminary' sciences (112 to 188) are essentially methodological sciences and techniques, applicable to many fields, and necessary tools for activity in the subject fields 200 to 890 (with the probable exception of 510 History and 520 Area Studies). It has already been noted that the sequence from 210 Physics onwards is one of increasing complexity. As this is of practical consequence in connection with matters of citation order in classifying composite subjects, it may bear restatement in different terms. Each science in this sequence has methodological and phenomenal aspects which when taken in isolation belong to preceding sciences in the sequence and not to following ones. Conversely each science in the series may contribute 'aspects' to sciences following it, but not to those preceding it.

The fact has already been mentioned that BSO follows a more traditional line than either BC2 or the Object Area Scheme in completely separating the applications of physics, chemistry, and of the space and earth sciences, from the parent sciences. It should be noted that no general classification since the Subject Classification of JD Brown has attempted a completely intimate collocation of sciences with their applications. No one now advocates this intimate collocation which fragments the whole of science and the whole of technology. The BSO Panel saw virtually no advantage and very many disadvantages in separating the steadily converging

physical sciences by inserting their associated technologies between them, nor in the concomitant scatter of technology, with General Technology occurring later in the schedule than many individual technologies. However, in the area of the Biological Sciences both BSO and the Object Area Scheme collocate sciences and their applications at a very broad level only. For instance BSO has Botany and Zoology followed by Agriculture comprising Plant Crops and Animal Husbandry. The more intimate collocation giving the sequence Botany - Plant Crops - Zoology - Animal Husbandry is rejected by both schemes. When, on the other hand, we reach the Biomedical Sciences both BSO and BC2 have the science and technology closely intermixed, so that, for instance, the physiology and the clinical medicine of a particular body organ are placed together

Syntactic aspects and combinatory facilities

The discussion of subject sequence in the BSO schedules so far undertaken in this chapter has been against the background of semantic relationships - the closeness or distance between terms on the plane of meaning, when they are considered simply as isolated terms. However, so-called cross-classification has been mentioned in passing. It has also been stated that some subjects may be used as 'tools' in other subjects, and that some contribute 'aspects' to others. 'Tools' and 'aspects' represent certain kinds of syntactic relations. These are relations, at the concept level, between terms which stand together to denote compound or composite subjects. 'Cross-classification' is frequently used to refer to the dilemmas experienced by classifiers attempting to assign places for composite subjects in classification schemes which are inadequately prescriptive on the handling of syntactic relations.

Among organised information sources there are some which are devoted to subjects which are composite in nature. Accordingly BSO has comprehensive facilities for combining notational elements to represent composite subjects. It is, in fact, a fully synthesising or faceted system, though it has not been thought necessary or even desirable to label facets as such.

Combinatory facilities in classification systems inevitably raise the issue of order in which the elements are combined, also called citation order or facet order. In some working situations this issue may be bypassed or left to intuitive judgment, but for a neutral mediating indexing language covering all subject fields a completely fixed or prescriptive citation order appears to be necessary to ensure reasonably 'noise-free' transmission of information.

In BSO the order in which notational elements are combined to form codes for composite subjects is in the majority of cases the reverse of the order in which the elements are set down in the classification schedule. Without the qualification in the majority of cases citation order problems would be reduced to purely clerical procedures, and if we can specify those situations to which the reverse-schedule-sequence rule applies without exception, we still have a highly time and effort saving feature of BSO.

It is first of all useful to categorise combinations into internal combinations which comprise notational elements drawn from the same subject field (e.g. 575,32,0,73,50 Child welfare in disaster relief, constructed from the elements 575,32 Child welfare and 573,50 Disaster relief and aid) on the one hand, and external combinations constructed from notational elements taken from different subject fields (e.g. 550-163 Operations research in public administration, using elements 550 Public administration and 163 Operations research) on the other.

In order to make this categorisation completely explicit it is necessary to state unambiguously what is meant in this context by a subject field. Subject fields for defining internal vs. external combinations are enumerated as 'Combination areas' on page xi of the published BSO. A combination with both elements drawn from one of these 'Combination areas' is an internal combination.

Internal combinations without exception obey the reverse-schedule sequence rule for combination order. (In the above example the leading element in the combination, 575,32 is later in the schedule than the second element 573,50).

The structural background to this combination rule, is that each subject field is elaborated according to a facet pattern, which, with very slight variations, is repeated over many fields. The following is the commonest facet pattern, given in schedule sequence which would be reversed for combination order:

- 1) Tools or equipment for carrying out operations
- 2) Operations (i.e. purposive activities by people)
- 3) Processes, interactions
- 4) Parts, subsystems of objects of action or study, or of products
- 5) Objects of action or study, or products, or total systems

(In the example above the first element in the combination order, namely the concept Child belongs to facet (5), the second element, the process which requires a welfare operation to be undertaken, namely the concept Disaster, belongs to facet (3). Facet (4) is inapplicable to this subject field. Facet (2) is applicable but has no role in this combination because the operation, Welfare already defines the whole 'combination area'. Facet (1) would be applicable if a particular kind of welfare agency were to be specified). Such regularity of underlying pattern covering the whole scheme is conducive to economy both in the day-to-day use of the scheme by indexers or searchers and in predictable updating.

The reverse-schedule-sequence rule cannot be used in the same clerical or mechanical manner in deciding combination order for external combinations, though more often than otherwise it would give correct and consistent results. The reason why it cannot be employed reliably for external combinations can be shown from a single example. Let us assume that reverse-schedule-sequence is being used as the basis for combination order in the case of Educational psychology. The rule will then give 460-450 (460 is Education, the hyphen or dash is the connecting symbol for external combinations, 450 is Psychology). Educational psychology, may be approximately factored as the Psychological aspects of the Education process. How then do we code Psychological education, the teaching and training in the subject Psychology? If the reverse-schedule-sequence rule were used we should arrive again at 460-450 as for Educational psychology. For any indexing system covering the whole of knowledge this would produce unacceptable noise at output. The example leads to two further considerations. The first is that external combinations should not (in the manner of the UDC colon connecting symbol) be used to indicate any relationship. This would lead not only to output noise, but also to anomalies in file sequence of classified material. The second consideration is that both in Educational psychology and Psychological education, one of the subjects (or rather the phenomena of one of the subjects) is the 'recipient' or 'target' to which the other subject contributes a set of aspects or properties. Thus psychological viewpoints are contributed to the education process in Educational Psychology, and an education process is contributed or

applied to the realm of psychology in Psychological Education. An interesting difference to be noted in passing is that the 'recipient' in Educational psychology is the primary phenomena of education, i.e. the education process, while the 'recipient' in the case Psychological education is not the primary phenomena of psychology - there is no reference here to the educating of psychological processes - but the second-order phenomena of people involved in psychology as a field of interest or profession.

The upshot of these considerations is that combination order for external combinations in BSO needs to be determined by reference to the relation between the elements which require connection. The following rule which emphasises the directionality of the 'recipient' element in the relation and the 'aspect contributing' element is believed to be unambiguously applicable to situations which can be represented by external combinations, and is recommended:

Cite first the notation for the element denoting application area, mission, purpose, end-product or whole system: more generally the subject which 'receives' an action or effect, or is seen according to a particular viewpoint, or has a property attributed to it

Cite second the notation for the element denoting aspect, approach, action applied, agent, or part of a stated whole: more generally the subject which 'contributes' an aspect, approach or action.

Use of the above relational formula where the 'aspect contribution' element belongs to the area 210 to 450 will normally produce combination orders which reverse the schedule order, as in the case of internal combinations throughout the schedule. This is because in this area the entities and phenomena studied by a particular science include aspects and properties which essentially belong to other sciences located earlier in the schedule sequence. For instance biological entities may have physical or chemical properties: medical, psychological and social phenomena may have biological aspects. In these cases the roles of 'aspect contributor' and 'recipient' elements cannot be reversed, as long as the 'recipient' element is the primary phenomena of the subject field concerned. If, as has been shown in the case of Psychological Education the 'recipient' element is the second order phenomena associated with the subject field an apparent reversal of roles results in a combination order which is identical with the schedule order of the elements.

Cases in which the 'aspect contributed' element belongs to the area 460 to 992, and to which the relational formula is applied, more frequently produce combination orders which break the reverse-schedule sequence rule. However, a glance at the outline suggests that the main exceptions to the reverse-schedule-sequence rule fall into a few categories. These are

1. Any social or historical aspect of any subject field 112 to 480 ('social' here is to be understood as including any aspect corresponding to subject fields 530 to 588)
2. Terminological aspect of any subject field 112 to 890.

Finally it should be noted that some departures from the reverse-schedule-sequence rule occur when the relational formula is applied to composite subjects of which both elements are drawn from the area 112 to 188.

Discipline and phenomena classes

One further feature of the BSO outline is worthy of mention. General classifications are based primarily upon subject disciplines which are methodologies and special points of view usually, but not necessarily, focussing upon a definite set of entities or phenomena. A consequence is that in conventional general classifications there is no way of classing entities or phenomena as such, merely described, or treated from many points of view. An institution dealing with, for example, Fish in all their aspects, zoological, economic, aquacultural, technical, mythological, and as quarry in a pastime, is not appropriately placed in Zoology (345,62). Furthermore there are organised information sources dealing in a multi-disciplinary manner with such topics as Food and Housing. In virtually every general classification before BSO such vital topics of everyday life have been misleadingly assigned to the discipline Sociology. Sociology is also an invariable dumping ground for multi-aspect studies, also reflected in institutional warrant, of social groups, such as Women, Racial Minorities, the Aged and the Disabled. These studies are by no means primarily sociological in viewpoint or treatment. BSO has attempted to deal with this problem by including a few phenomenon- or entity-based classes in its main outline, all containing a marked human reference, and by supplying a special location 088 at the beginning of the classification, for other phenomena or entities not included in the above-mentioned phenomena/entity-based classes. The enumerated phenomena/entity-based classes are 470 Human Needs, covering Food, Clothing and Shelter in their most extended aspects, together with Leisure, 520 Area Studies which are multidisciplinary in character, and 528 Social Groups. In connection with the residual phenomenon/entity class at 088, the problem arises as to how the multifarious phenomena and entities which might need to be assigned here should be individualised and ordered. The solution to this problem, as used in BSO, may be illustrated by the example of Fish already given. It was pointed out that Zoology offers one point of view upon Fish and that it would therefore be wrong to assign multidisciplinary (or non-disciplinary) material on Fish within the discipline of Zoology. Yet despite this mismatch in aspect or point of view Fish have a relation with the discipline Zoology which is not of the same Kind as their relation to Economics, Aquaculture, Food technology, Sport, or Mythology. The special relation with Zoology consists of the fact that the concept Fish is uniquely defined by the zoological characteristics of fish, namely their anatomical and physiological features. The concept Fish is not similarly defined - though it may well be described - in terms of the characteristics peculiar to Economics, Aquaculture, Food technology, Sport or Mythology. When treated in multidisciplinary manner any entity such as Fish may be linked - though not subordinated - to the discipline within which it is uniquely defined, and this circumstance make available a mechanism whereby an entity may be individualised and ordered at 088. The notation for the entity within the discipline which defines it is simply added to 088,. Thus the notion of Fish is uniquely defined in zoological terms. The notation for zoological aspects of Fish is 345,62. The same notation added to 088 as 088,345,62 then signifies Fish in all their aspects, zoological and other.

A somewhat different basic view, but a similar mechanism, is applied to the problem of individualising technical products. There is no question here of multiplicity of points of view. The point of view is assumed to be technical, embracing manufacture and the technique for using and maintaining the product. The problem is simply one of individualising the great number of kinds of products which emerge from technical processes. In BSO products defined by purpose or designed for a particular purpose are classed at the end of the Technology schedule at 890, and individualised by reference to the BSO code for the particular purpose, elsewhere in the scheme. It is necessary to emphasise 'elsewhere in the scheme' as the purpose of some products is simply to contribute to more complex technology. Such products (e.g. Switchgear)

with a role internal to technology are normally enumerated in the BSO Technology schedules. The scheduled heading covers both their manufacture and use (Manufacture can be distinguished from use by employment of the suffix ,06,20 taken from 620 Production technology). One consequence of the policy for individualising by purpose those products with purposes external to technology is that 877,60 Cloth and fabric technology does not schedule manufacture of clothing as a product. The technology of the purpose-defined product Clothing is classed at 890,472. The 472 is taken from the root Human Needs code for Clothing.

Common Facets

BSO has Time and Place facets, introduced by notation -01 and -02 respectively, which are functionally similar to Time and Place divisions provided in other general classification schemes. They are applicable to every subject field except those such as 510 History, 520 Area studies, 544 to 546 Political history and Politics of individual states and groupings of states, where Place and Time are specially scheduled facets. The Place facet makes use of ISO two-character alphabetical codes, and can also specify transnational political areas (e.g. EEC countries) areas defined by language, race or religion, and areas defined by the usual physical geographical factors (e.g. Tropical areas).

An Optional facet enabling the type of information source to be specified has also been included as a result of the field test. It was found that data as to type of information source is often given prominence in descriptive material upon which indexers rely in order to establish the subject field for classifying purposes. However this data does not form part of the subject description and failure to realise this results in codes being applied which give misleading information. For example, lack of attention to this factor could cause such an information source title as *British Technology Index* to be wrongly coded as

600-026,GB An information source on the technology of Great Britain

whereas the correct coding is

600 33-026,GB An index, originating in Great Britain, on technology

Were it decided not to use the Optional facet, the correct coding in this case would be 600. Though the above example is taken from the area of technology, ambiguity in the use of place designation in the titles of services and institutions is even more commonly encountered in the social sciences. The use of the Optional facet compels the classifier to penetrate such ambiguity in searching for the correct subject description of an item.

Notation

Notation is the last feature of the BSO to be dealt with in this descriptive account of the scheme, and this perhaps reflects the view of the compilers that notation is at all times to be regarded as an ancillary to the structure of the classification. The scheme was in the first place constructed independently of any notation. The present notation could be uncoupled and another used in its place without changing the character of the system, always assuming that any new notation would be no less able than the present one to handle combinations and produce the required order.

The notation given in the published BSO is intended to be read and carried in mind by human users. There could be good reasons why a notation intended to be read and stored by a machine might be rather different. The human user reacts negatively to overlength and over-complexity arising from the appearance of symbols from different sets (e.g. alphabetical and numerical). Within the necessarily prescribed size limits of even variable length records, the computer is not seriously troubled by length of notation, and as all species of digits are in any case converted to numerical values for processing, a superficially mixed notation has no terrors for it. The BSO notation is believed to be tolerably brief: over 90% of the uncompounded terms cited in the schedules have codes of 5 numerical symbols length. By relying on the use of numbers as the main symbol set, and using other symbols only sparingly it manages not to be over-obtrusive. Also, by eschewing the secondary functions often accorded to notations it is capable of admitting new subjects, without limit, at their logically correct positions. It fulfills its primary function of mechanising the sequence of subjects in the schedule or in a user's file but it gives no structural information apart from that necessarily implied by the order of subjects alone. Notations of this kind are often termed 'non-expressive' notations, though it should not be overlooked that such notations do express syntactic relations. A 'non-expressive' notation was devised for 850, because all experience of the earlier established general classifications goes to show that notations which express structure, particularly hierarchy, create very difficult and sometimes insoluble problems in the insertion of new subjects in their correct places. Too often new subjects are inserted in the wrong place because of the presence of a notational gap, or the process of inserting it in the correct place involves the re-notating of the neighbouring part of the schedule. It is the dilemma embodied in these two alternatives which causes most of the decision effort and cost entailed in revising the established general classifications. This effort, and the associated delay and cost to users, should not, it was felt, be accepted as a necessity in connection with an ongoing universal switching indexing language. It was here, more perhaps than anywhere else, that the requirements of the UNISIST switching language demanded a complete break with tradition.

It was stated above that the published notation was intended for the human user. This is not quite the same as saying that it is intended only for manual switching systems. Computer processed switching systems also have human users of switching languages at both input and output ends of the switching system. Also, it is not to say that the notation could not be fed to a computer for switching between symbols having the same meaning in different local indexing languages. However, if facilities for combining switching with interactive computer-aided search were required, it would be preferable to employ for this purpose another notation containing built-in cues enabling the machine to traverse requested search paths. For the human user, the schedule of terms itself, the conceptual pattern implicit in the manner of their ordering, their hierarchical status, and the cross-references directing to related locations in the schedule, together constitute the search aid. For computer search all these matters must be explicit in the notation. Such a fully-expressive computer-oriented notation would be far too long and complex for direct use by human beings. However, given removal of the constraints upon length and complexity necessary for the human user, such a computer-oriented notation could be as hospitable to new knowledge as is the present BSO human-oriented notation.

Arabic numerals were chosen as the base symbol-set of the BSO notation because they are the best known set with elements carrying well-understood sequence values, and because they are invariable throughout the world. The numerical characters are supplemented by two punctuation signs, the hyphen and comma, and by the Roman alphabet A to Z for occasional situations where individualisation rather than grouping is required, as for instance in specifying the names of individual artists. Some notation elements are drawn from outside coding systems, such as

the ISO code for names of countries, and the Groups of the Periodic table also employ Roman numerals. The use of characters supplementary to numerals demands a fixed system of ordinal values as between the supplementary characters and numerals. The following sequence gives the recommended ordinal value system for files organised by BSO:

Spaces after last symbol of notation

Two spaces, followed by further numerical characters This occurs when the Optional facet for type of source is used.

- followed by further numerical characters This is the connecting symbol for external combinations of notation.

, followed by further characters This is a semantically empty character which introduces intercalated numbers filing between consecutive members of a notational array.

00 to 99 or 000 to 999 These two set never occur together in a file in such a way as to require ordinal preference between them.

A to Z

Turning now from the ordinal value of individual symbols to the make-up of a notation or code for a given subject, all codes begin with a member of the millesimal array 000 to 999. Between any two subjects represented by consecutive members of this millesimal array, further subjects may be interposed by adding to the first of the two consecutive numbers concerned a comma followed by a member of the two-digit centesimal array 00 to 99. In similar fashion further subjects may be interposed between consecutive numbers of the 00 to 99 array, by adding a comma followed by members of a further 00 to 99 array. Accordingly a typical code structure comprises a single group of 3 numbers followed by an indefinite number of groups of 2 numbers, all groups being separated by commas (e.g. 915,15,50.....)

In a few well-defined situations this typical 3,2,2.... pattern may be varied. In notations which contain the hyphen (external combinations, and Time and Place Facets) the 3,2,2.... pattern may appear on both sides of the hyphen. However, in many cases the hyphen links two groups of 3 numbers (e.g. 642-580 Nuclear reactor economics). Internal notation combinations contain the single separated number a as a connecting symbol (e.g. 978,0,72,37 The Koran). The numbers 088 and 890 are the leading number groups of untypical 3,3,2,2,... patterns.

Notational combination

A tabulation of the procedures for combining notation is given on page xiii of BSO, and is reproduced here on page 51. It should be added that in some fields where composite subjects are expected to arise frequently by comparison with unitary enumerated subjects, the schedules themselves provide for intersecting concepts by Expand like instructions. Notations produced by this mechanism are always shorter than the combined notation in which elements are linked by the connecting symbols mentioned in the tabulation. In deciding when to use Expand like instructions, the BSO Panel were obliged to balance the advantage of brevity against the two disadvantages that the Expand like mechanism is more likely to lead to indexing errors than the connection of two notational elements by a set of connecting symbols, and will also use up more of the available brief notation, thus in the long run causing an increase in the length of notation of future enumerated subjects. At the level of notation manipulation the difference between an internal combination (outside the area 600 to 890) and an Expand like instruction is that while the internal combination adds a connecting symbol and deletes the first numeral of the second notational element, an Expand like instruction omits any connecting symbol and at the same time deletes two or more of the first numerals of the second notational element. At the concept level Expand like instructions can be more versatile than internal combinations. This versatility is manifested in the BSO schedules where an Expand like instruction adds the legend 'with incorporated additions marked +'. These 'incorporated additions' are concepts which arise only in subordination to a facet combination. For instance the BSO Physics schedule consists essentially of a list of energy interactions and forms (from which the notion of a specific medium is absent), followed by a list of media or forms of matter. The Expand like instructions cause notions of forms of matter to be combined with notions of energy. Thus for instance we have 224,25 Plasmas and fluids, Mechanics. An important branch of Plasma and Fluid Mechanics is Magnetohydrodynamics. This branch of mechanics is logically dependent upon the combination of Plasmas and fluids with Mechanics. At the level of generalised energy interactions there is probably no term to comprehend the abstract idea of mechanical motion, magnetic fields and electric fields in triangular interaction. Accordingly 224,34 Magnetohydrodynamics is entered as an 'incorporated addition' subordinate to 224,25 Mechanics of plasmas and fluids.

REFERENCES

- 1 BLISS BIBLIOGRAPHIC CLASSIFICATION: 2nd ed. Edited by V. Broughton, and J. Mills. London, Butterworth, 1977- in progress
- 2 BLISS, H.E. Organization of knowledge in libraries and the subject approach to books. New York, H.W. Wilson, 1933
- 3 BROWN, J.D. Subject classification. 3rd ed. revised and enlarged by J.D. Stewart. London, Grafton, 1939.

Procedures for combining notation

For composite subjects not given in Schedule and not derivable from 'Expand like' notes in schedule

Internal combinations (Both elements scheduled in the same Combination Area)		External combinations (Elements scheduled in different Combination Areas)	
Combination other than 600 to 890	Areas than	Combination Area 600 to 890	
1. Write notations for separate elements in reverse schedule order, side by side, with space for 2 characters between them		1. Write notations for separate elements in reverse schedule order, side by side, with space for 2 characters between them	1. Decide combination order according to citation rule (see p. 45) 2. Write down separate notations in chosen combination order 3. Insert Dash (or hyphen) between the 2 notation elements
2. Delete 1st digit of the 2nd notation element		2. Insert a comma after the 1st digit of the 2nd notation element	
3. Insert ,0, in the 3 character space created	3	3. Insert ,0 in the 2 character space	

CHAPTER 7: PRACTICAL SUBJECT INDICATION WITH BSO

This chapter will deal with some general issues which arise in the practical use of BSO as an indexing language. Initially its application as a subject tagging or ordering code for information-bearing items recorded in a static medium, such as a directory, will be considered. In applications of this kind the linking of BSO codes to local indexing language codes is not a factor which has to be taken into account. When such linking is necessary, as for instance when BSO codes need to be linked to indexing languages of specialist information centres in an exchange network, fresh considerations come into play. These are discussed in the latter part of the chapter. [...]

Using BSO as a direct classification

BSO is controlled indexing and retrieval language. [...]

Background knowledge for concept analysis

How much background knowledge of the subjects concerned is required for concept analysis? [...]

Issues of policy in classifying

Decisions are needed upon certain areas of indexing or classification policy, which the classification system itself should not attempt to take, nor even to advise upon, because the answers are strictly related to the particular uses to which the scheme may be put. Answers are, however, needed if consistency and uniformity of practice are to be achieved. These qualities are at an especially high premium in any scheme of co-operative subject indication, of which a switching network is but one example. [...]

OUTLINE OF BSO

SUBJECT FIELDS

088	<u>Phenomena & entitites from a multi or non-disciplinary point of view</u>		
100	<u>KNOWLEDGE GENERALLY</u>	500	<u>HUMANITIES, CULTURAL & SOCIAL SCIENCES</u>
200	<u>SCIENCE AND TECHNOLOGY</u>	510	<u>History and related sciences</u>
203	<u>Natural sciences</u>	526	<u>Area studies</u>
300	<u>Life sciences</u>	530	<u>Social sciences</u>
410	<u>Biomedical sciences</u>	600	<u>TECHNOLOGY</u>
445	<u>Behavioural sciences</u>	910	<u>LANGUAGE, LINGUISTICS & LITERATURE</u>
460	<u>EDUCATION</u>	940	<u>ARTS</u>
470	<u>HUMAN NEEDS</u>	970	<u>RELIGION & ATHEISM</u>

088 PHENOMENA & ENTITIES FROM A MULTI- or NON-DISCIPLINARY POINT OF VIEW

100	KNOWLEDGE GENERALLY
112	PHILOSOPHY
116	SCIENCE OF SCIENCE
118	LOGIC
120	MATHEMATICS
125	STATISTICS & PROBABILITY
127	DATA PROCESSING
128	COMPUTER SCIENCE
140	INFORMATION SCIENCES
143	Libraries
144	Archives
146	Museums
146,80	Exhibitions
148	MEETINGS
148,80	CONSULTANCY
148,85	INTERVIEWING
150	COMMUNICATION SCIENCES
152	Reprography & printing
152,60	Book trade
152,80	Intellectual property
155	Destination-directed communication
156	Mass communication
158,20	Publicity
160	SYSTEMOLOGY & CYBERNETICS

163	Operations research
165	MANAGEMENT
166	STANDARDISATION & STANDARDS
168	ORGANISATIONS
182	RESEARCH
184	DISCOVERIES, INVENTIONS & PATENTS
186	TESTING & TRIALS
188	METROLOGY
200	SCIENCE & TECHNOLOGY (TOGETHER)
203	NATURAL SCIENCES
205	PHYSICAL SCIENCES
210	Physics
212	Energy interactions & forms
214	Particle & high-energy physics
215	Nuclear physics
217	Atomic, molecular & ion physics
219	Vacuum physics
222	Bulk matter physics
224	Plasma & fluid physics
225	Condensed matter physics
226	Physics of solids
228	Crystallography
230	Chemistry
232	Physical chemistry
234	Chemistry of particular substances
235	Inorganic chemistry
237	Organic chemistry
238	Polymer chemistry
250	Space & earth sciences
252	Astronomy & astrophysics
258	Space research
260	Earth sciences
262	Geodesy & surveying
263	Geophysics
265	Atmospheric sciences
267	Hydrospheric sciences
270	Geology
290	Geography
300	LIFE SCIENCES
310	BIOLOGICAL SCIENCES
312	Biophysics & biochemistry
313,20	Molecular biology
313,30	Cell & tissue biology
313,70	Genetics
313,75	Evolution
315,19,34	Biological materials
315,20	Developmental biology
315,54	Morphology
315,55	Physiology

315,56	Pathology
315,58	Immunology
315,59,30	Biological material structures
315,60	Biological parts, organs & functional systems
318	Ecology
318,50	Ethology
320	Microbiology
330	Botany
340	Zoology
359	APPLICATIONS OF LIFE SCIENCES
360	AGRICULTURE
363	Plant crop production
363,70	Horticulture
363,80	Production of specific plant crops
366	Animal husbandry
368	VETERINARY SCIENCE
370	FORESTRY
380	WILDLIFE EXPLOITATION
390	ENVIRONMENT
397	Natural resources
410	BIOMEDICAL SCIENCES
420	MEDICINE
422	Preventive medicine
423	Social medicine
425	Clinical & internal medicine
426	Surgery
428	Diseases
430	PARTS, ORGANS & SYSTEMS OF THE HUMAN BODY
432	Body parts & regional specialties
433	Body organs & systems
433,50	Defence system
434	Integumentary & musculoskeletal system
435	Visceral systems & organs
437	Nervous system & sense organs
438	Mental health & disorders
439	BIOMEDICAL SPECIALTIES, BY HUMAN SUBJECT OR PATIENT, & BY ENVIRONMENT
445	BEHAVIOURAL SCIENCES
450	PSYCHOLOGY
460	EDUCATION
470	HUMAN NEEDS
475	Household science
477	Work & leisure occupations
477,50	Leisure & recreation
478	Tourism & travel
480	Sports & games
482	Individual prowess & athletic sports
483	Ball games
485	Special environment sports
486,40	Wheel vehicle sports
486,50	Animal sports

486,70	Target & quarry sports
488,30	Board & piece games
489	Social diversions & pastimes
500	HUMANITIES & SOCIAL STUDIES
510	HISTORY & RELATED SCIENCES
512	Archaeology & prehistory
513	History of particular epochs & periods
515	History of main world areas & continents
516	History of particular countries (including regions & localities within a particular country)
520	AREA STUDIES
526	Area studies of particular countries (including regions & localities within a particular country)
527	SOCIETY
528	SOCIAL GROUPS & COMMUNITIES
529	Ethnic & linguistic & religious groups
530	SOCIAL SCIENCES
532	Culture
533	CULTURAL ANTHROPOLOGY
535	SOCIOLOGY
537	DEMOGRAPHY
540	POLITICAL SCIENCE & POLITICS
542	Political institutions & organisations
543	Political organisational patterns & systems
544	Political history
545	Politics of particular groupings of states
546	Politics of particular states & countries
550	PUBLIC ADMINISTRATION
560	LAW
562	Civil law
563	Public law
565	International law
567	Systems of law (by origin)
568	Law of particular countries
570	SOCIAL WELFARE
580	ECONOMICS
581,80	Microeconomics
582	Macroeconomics
584	Economic organisation
586	Sectorial economics
588	MANAGEMENT OF ENTERPRISES
600	TECHNOLOGY
610	SCIENTIFIC BASIS OF TECHNOLOGY
611	EQUIPMENT & PLANT
612	SYSTEMS ENGINEERING
612,55	COMPUTER TECHNOLOGY
615	TECHNICAL TESTING
617	MAINTENANCE & SERVICING ENGINEERING
618	TECHNICAL & INDUSTRIAL DESIGN

620	PRODUCTION TECHNOLOGY
625	MATERIALS HANDLING
627	PACKAGING & STORAGE & DISPATCH
631	ENERGY TECHNOLOGY
635	MATERIALS TECHNOLOGY
640	NUCLEAR TECHNOLOGY
642	Nuclear reactor technology
645	Isotope technology
650	ELECTRONIC & ELECTRICAL TECHNOLOGIES
653	ELECTRONIC ENGINEERING
655	TELECOMMUNICATION ENGINEERING
657	ELECTRICAL ENGINEERING
670	THERMAL ENGINEERING & APPLIED THERMODYNAMICS
673	Heat engines
678	Refrigeration technology
680	MECHANICAL ENGINEERING
684	FLUID ENGINEERING
686	VACUUM TECHNOLOGY
688	VIBRATION & ACOUSTIC ENGINEERING
710	CONSTRUCTION TECHNOLOGY
712	CIVIL ENGINEERING
714	ILLUMINATING ENGINEERING
716	BUILDING CONSTRUCTION & SERVICES
716,70	Types of buildings
720	ARCHITECTURE
724	LANDSCAPE DESIGN
726	PHYSICAL PLANNING
730	ENVIRONMENTAL TECHNOLOGY
730,30	Pollution control
734	Public & industrial health engineering
736	Safety engineering
738	Rescue & salvage operations
740	TRANSPORT TECHNOLOGY & SERVICES
742	Road transport technology & services
743	Railway transport technology & services
745	Water transport technology & services
747	Air transport technology & services
748	Space transport technology
760	MILITARY SCIENCE & TECHNOLOGY
764	Warfare
767	Weapons in warfare
780	MINING
782	Solid fuel extraction
783	Metal mining
784	Non-metallic mineral mining
786	Oil & gas extraction technology
788	Ore & mineral dressing
800	PROCESS INDUSTRIES
810	Chemical technology & engineering
811	Chemical engineering
812	Chemical technology

813	Chemical agents & basic industrial chemicals
815	Industrial gases technology
816,90	Biotechnology
823	Technology of particular groups of chemicals
825	Petroleum technology
826	Natural oils, fats & waxes technology
828	Polymer technology
831,30	Chemical agents for materials processing
831,40	Surface-active agents
831,44	Adhesives & sealants
831,47	Lubricants
831,60	Surface finishing agents
832	Fuels & explosives technology
834	Colour industry technology
836	Pharmaceuticals & related technologies
838	Agricultural chemicals technology
840	FOOD & DRINK TECHNOLOGY
841	Food technology
845	Drinks technology
847,50	TOBACCO PROCESSING
849,13	MINERAL PROCESSING TECHNOLOGY
850	Non-metallic mineral technologies
852,32	Lime & lime products
852,34	Gypsum & anhydrite products
852,50	Cement & mortar technology
852,60	Concrete technology
854	Ceramics & clayware technology
856	Glass & glass ceramics technology
858	Fibrous & layered silicate mineral technologies
860	METAL TECHNOLOGY
864	Metal products
865,50	Metallurgy of particular & kinds of metals
865,60	Alloys
866,60	Ferrous metallurgy
867	Non-ferrous metallurgy
871,95	WOOD, PULP & PAPER TECHNOLOGY
872	Wood technology
873	Pulp & paper technology
875	LEATHER & OTHER ANIMAL PRODUCTS TECHNOLOGY
877	TEXTILES TECHNOLOGY
878	CORDAGE & WIRE ROPE MAKING
890	MANUFACTURE & TECHNOLOGY OF PARTICULAR PRODUCTS NOT SCHEDULED IN BSO AREA 600 TO 878
910	LANGUAGE & LITERATURE
911	LINGUISTICS
912	USE OF LANGUAGE
912,20	Authorship
912,30	Reading
915	LITERATURE

- 920 SPECIAL PHILOLOGICAL STUDIES**
- 921 Indo-European languages & literatures**
- 923,20 Afro-Asian languages & literatures**
- 923,60 Caucasian languages & literatures**
- 923,70 Basque language & literature**
- 924 Eurasian & North Asian languages & literatures**
- 925,20 Dravidian languages**
- 925,51 Sino-Tibetan languages & literatures**
- 925,70 Austronesian & Oceanic languages & literatures**
- 927 African languages & literatures**
- 928 Amerindian languages & literatures**

- 940 ARTS**
- 943 PLASTIC ARTS**
- 945 GRAPHIC FINE ARTS**
- 947 PHOTOGRAPHY AS ART**
- 949 DECORATIVE ARTS & HANDICRAFTS**
- 950 MUSIC & PERFORMING ARTS**
- 951 Music**
- 952 Vocal music**
- 953 Instrumental music**
- 955 PERFORMING ARTS**
- 957 Cinema**

- 970 RELIGION & ATHEISM**
- 971 Atheism & rationalism**
- 972 Religion**
- 973,90 Particular religions**
- 974 Religions of Indian origin**
- 975 Religions of Far Eastern origin**
- 975,70 Religions of Iranian origin**
- 976 Judaism**
- 977 Christianity**
- 978 Islam**
- 979 Other religions & semi-religious cults**

- 992 ESOTERIC PRACTICES & MOVEMENTS**

BSO, also termed SRC (Subject-field Reference Code), is a classification system developed within the UNISIST-program for the purpose of interconnection of information systems. It is a disciplinary organized system founded in 1972 as a UNESCO project in the UNISIST-program "World Science Information System" in cooperation with FID. The idea behind BSO is related to the idea of networks and probably represents the latest attempt to create a new universal classification. It is developed by the Englishman Eric Coates in cooperation with others, including the Bliss Classification Association.

BSO was meant to be an international switching language, an overall information retrieval language to transfer blocks of information in coarse subject groups between information systems applying different indexing languages.

Coates (1979) states: "theoretically the switching operation requires nothing more than a neutral code system in which concepts are represented". Dahlberg (1978) regards it as a positive step towards standardization: "A standard classification assists library rationalization and national and international cooperation on statistics, research and cataloguing".

As referred in the entry on switching language, has critical voices been raised towards the theoretical assumptions behind such systems. BSO did not develop to fulfill the wishes behind its construction. This may have a connection to such inherent problematic assumptions in the concept of switching languages.

BSO (Broad system of ordering)

BSO, also termed SRC (Subject-field Reference Code), is a classification system developed within the [UNISIST](#)-program for the purpose of interconnection of information systems. It is a disciplinary organized system founded in 1972 as a UNESCO project in the UNISIST-program "World Science Information System" in cooperation with [FID](#). The idea behind BSO is related to the idea of networks and probably represents the latest attempt to create a new universal classification. It is developed by the Englishman Eric Coates in cooperation with others, including the Bliss Classification Association.

BSO was meant to be an international [switching language](#), an overall [information retrieval language](#) to transfer blocks of information in coarse subject groups between information systems applying different indexing languages.

Coates (1979) states: "theoretically the switching operation requires nothing more than a neutral code system in which concepts are represented". Dahlberg (1978) regards it as a positive step towards standardization: "A standard classification assists library rationalization and national and international cooperation on statistics, research and cataloguing".

As referred in the entry on [switching language](#), has critical voices been raised towards the theoretical assumptions behind such systems. BSO did not develop to fulfill the wishes behind its construction. This may have a connection to such inherent problematic assumptions in the concept of switching languages.

Literature:

Coates, E. J. (1979). The Broad System of Ordering. *International Forum for Information and Documentation*, 4(3), 3-6.

Coates, E. J. (1980). The Broad System of Ordering (BSO). IN: *New trends in documentation and information: proceedings of the 39th FID Congress, University of Edinburgh, 25-28 September 1978*. Edited by Peter J. Taylor. London, Aslib, 259-273.

Coates, E. J. (1981). The Broad System of Ordering: the compilers reply to their critics. *International Forum on Information and Documentation*, 6(1), 24-30.

Coates, E.; Lloyd, G. & Simandl, D. (1978). BSO: Broad System of Ordering: schedule and index. Prepared by the FID/BSO Panel. 3rd revision. The Hague : Federation Internationale de Documentation (FID) : United Nations Educational, Scientific and Cultural Organization (UNESCO). (FID Publication 564).

Coates, E.; Lloyd, G. & Simandl, D. (1979). *The BSO Manual: the development, rationale and use of the Broad System of Ordering*. The Hague : Federation Internationale de Documentation (FID). (FID Publication 580).

Dahlberg, I. (1975). The terminology of subject-fields. *International Classification*, 2(1), 31-37.

Dahlberg, I. (1977). Major developments in classification. IN: *Advances in librarianship, volume 7*. Edited by M. J. Voigt & M. H. Harris. New York: Academic Press, 41-103.

Dahlberg, I. (1978). Normung und Klassifikation. *DK-Mitteilungen*, 22(5-6), 13-17.

Dahlberg, I. (1980). The Broad System of Ordering (BSO) as a basis for an Integrated Social Science Thesauri? *International Classification*, 7(2), 66-72.

Dahlberg, I. (1982). The Broad System of Ordering (BSO) as a Basis for an Integrated Social Science Thesaurus? *International Classification*, 7, 66-72.

DeHart, F. E. (1982). Topic relevance and BSO switching effectiveness. *International Classification*, 9(2), 71-76.

Foskett, D. J. (1975). Classification. IN: *Handbook of special librarianship and information work, 4th ed.* Edited by W. E. Batten. London, Aslib, 153-197.

Foskett, A. C. (1979). The Broad System of Ordering: old wine into new bottles? *International Forum for Information and Documentation*, 4(3), 7-12.

Interdepartmental Commission on Classification at the USSR State Committee for Science and Technology: Comments on the Broad System of Ordering. *International Forum for Information and Documentation*, 4(3), 1979, 25-27.

Lloyd, G. (1979). *BSO-Broad System of Ordering*. Arlington, VA.: Educational Resources Information Center. ERIC report. ED-186 005.

Madeley, H. (1983). The Broad System of Ordering. *Australian Academic and Research Libraries*, 14(4), 235-246.

Perreault, J. M. (1979). Some problems in the BSO. *International Forum for Information and Documentation*. 4(3), 16-20.

Rybatchenkov, V. (1974). Development of a Broad System of Ordering for UNISIST purposes. *International Classification*, 1(1), 20-21.

Soergel, Dagobert: The Broad System of Ordering - a critique. *International Forum for Information and Documentation*, 4(3), 1979, 21-24.

Svenonius, Elaine: Compatibility of retrieval languages: introduction to a forum. *International Classification*, 10(1), 1983, 2-4.

Toman, J. & Lloyd, G. A. (1975). Introduction to the Subject-field Reference Code (SRC) or Broad System of Ordering (BSO) for UNISIST. IN: *Ordering systems for global information*

newworks. Proceedings of the 3rd international study conference on classification research. Bombay, India., 6-11, January, 1975. p. 321-326.

Vickery, B. C. & McIlwaine, I. (1979). Structuring and switching; a discussion of the Broad System of Ordering. *International Forum for Information and Documentation*, 4 (3), 13-15.

Wellisch, H. H. (1976). Dewey Decimal Classification, Universal Decimal Classification, and the Broad System of Ordering: the evolution of universal ordering systems. IN: *Major classification systems-the Dewey centennial: papers presented at the Allerton Park Institute Number 21, held November 9-12, 1975, Allerton Park, Monticello, Illinois*. Edited by K. L. Henderson. Urbana-Champaign, University of Illinois Graduate School of Library Science, 113-124.

B S O. Broad System of Ordering. A general, faceted classification scheme for information exchange and switching.

<http://www.ucl.ac.uk/fatks/bsol/>

Lesson 4: Indexing Languages

Introduction:

One of the important functions of an information retrieval system is to match the content of documents with the user's requirements or queries. When the librarians or document lists make subject approach to information, they are confronted with the difficult task of subject indexing. Subject indexing is a method of information retrieval. The subject index helps the searcher from an unclear or a rough statement to an extensive standard one. The basic objective of subject indexing is to match the content of documents with the user's queries. A number of systems namely chain indexing, POPSI, PRECIS etc have been developed for preparing subject index entries of documents.

In order to darksome various complex indexing problems many forms of controlled vocabulary have been developed such as thesaurus, thesauri fact et.

Indexing Language:

The phrase "Indexing Language" is generally defined as all the words permitted either to describe a specific document or to construct a query to search a document file, along with the rules describing how the terms are to be used and in what relation to each other. In other words allowable in that indexing language. It is the complete set of terms in the natural language that are employed in the collection of documents. The list includes all required synonyms that are used in the process of indexing a set of documents. But it does not mean that all the terms in the list can be used to actually index the documents.

Types of indexing Languages:

Indexing language have been categorized into a number of fundamental types. A well known category is designated as assigned and derived indexing systems. In the former an indexer assigns terms or descriptors on the basis of subjective interpretation of the concepts implied in the documents. It is an intellectual method involving the finding out of specific subject of the document and assigning an appropriate subject heading.

In derived indexing system all terms or descriptors are taken from the document it self. In other words derived term systems are almost clerical and can be easily mechanized. Author indexes, title indexes, citation indexes and natural language indexes are derived term systems, where as all indexing languages with vocabulary control devices such as subject heading lists, thesauri and classification schemes are assigned term systems.

Natural Language Indexing:

Any information retrieval system without vocabulary control, if referred to as a "Natural Language" or sometimes as a "free – text" system because the system allows the indexer to select the term to be used directly from the text being indexed, or in automatic systems, the terms are selected by the computer. The terms are chosen from the text in self this approach may also be called indexing by extraction. The uniterm systems in the early days are the example of natural language system.

Natural language has some advantages. Its vocabulary is uptodate and its keeps on growing assimilating new concepts as soon as they come into being. If we want to conduct a highly specific search, the natural language system is more useful.

Controlled vocabulary or Artificial Language.

The lists of subject heading classification schemes and thesauri are representatives of indexing languages. The vocabulary (Indexing Language) of a retrieval system exerts a considerable influence on its performance. It greatly influence of the construction of each search strategy. Because of the controlled vocabulary of indexing languages. It is exact and precise. And there is no problem of synonyms and homonyms references taken case of these.

Advantages of controlled vocabulary:

The controlled vocabulary has indeed a number of obvious advantages.

It controls synonyms and near synonyms and brings semantically related terms together.

If properly constructed control vocabulary avoids many problems of false coordination or in correct term relationship.

Thesaurus:

The word "thesaurus" is derived from Greek and Latin words which means a "Treasury" according to the Oxford English dictionary the earliest usage of work thesaurus was known in 1565 from the title. Thesaurus Language Romance et Britanie. Modern usage may be said to date from 1852 when the first edition of the thesaurus of English words and phrases was published by Peter Mark Roget. The addition of "and phrases" in the title had great significance. It is considered as the "Father of all thesauri". An interesting and entertaining historical account has been given by Karen Spark Jones. Who traces the origin of "Synonymy" in dictionaries and identifies the main differences from natural language: the source involves "vocabulary normalization".

The concept of thesaurus in Library and information Science. Used by Helen Brouson for the first time in connection with information retrieval at the dorking conference on classification on May 14th 1957.

Thesaurus – definitions:

Here some of the important definitions are quoted on thesaurus and understand the meaning of them. The Oxford English Dictionary defines the thesaurus as a "Treasury" or "Store house of Knowledge."

English Language thesaurus means "a book contain a store of words of information about a particular field or set of concepts, specifically, a dictionary of synonym."

The second revised edition of the guidelines for the establishment and development of monolingual thesauri defines the thesaurus as the vocabulary of a controlled indexing language, formerly organized so that the a priori relationship between the concepts are made explicit."

All the above definitions show that the thesauri mean a treasury or store house in its general usage. As the word entered the literary field, through the well known regents thesaurus the usage has been changed to mean a lexicon or a dictionary or encyclopedia.

The usage of the term thesaurus, in Library and Information Science reference tools denote a special function in information retrieval through controlling the terminology.

Purpose of thesaurus:

The major purposes of a thesaurus include the following:

1. To provide a map of a given field of knowledge, indicating how concepts or ideas about concepts are related to one other, which helps on indexer or a searcher to understand the structure of the field.
2. To provide a standard vocabulary for a given subject field.
3. To provide a guide for users of the system so that they choose the correct term for a subject search; this stress the importance of cross references. It an interest uses more than one synonym in the same index for example; “abroad” “Foreign” and “overseas” – then documents are liable to be indexed haphazardly under all these; a searcher who chooses one and finds document indexed there, will assume that he has found the correct term and will stop his search without knowing that there are other useful documents indexed under the other synonyms.
4. To provide classified hierarchies so that a search can be broadened or narrowed systematically.

Indexing terms and their relationships:

The basic elements in a thesaurus are the individual words, terms, or phrases and these are offer called “descriptors” or Keywords” the descriptors are arranged in thesaurus an alphabetical order. An indexer assign the descriptor terms to describe the content of documents.

The terms which are not proffered to be used in indexing are “ Non-descriptors” they are proper names of corporate bodies, government agencies institutions and firms, geographical names etc.

Relation ships:

One of the key points of a thesaurus is that it indicates the relationship among terms.

They are :

Hierarchical or structural Relationship.

Eqatence or preferential Relation.

Associate or Affinitive Relation.

The symbols to express the relationships in thesauri have become more or less standardized as follows;

SN Scope Note.

USE Equivalent to “see” reference.

UF Use for, the reciprocal of use.

BT Broader Term.

NT Narrow Term.

RT Related Term.

For example:

INTELLIGENCE

BT: Ability.

NT: Comprehension

RT: Aptitude

By having these relationships displayed, both the indexers and the searchers are in a better position to cover the full range of options that may be possible in either indexing or searching.

Hierarchical relationships:

The hierarchical relation expresses super/subordinate of concepts. There are two types of relationships in this category; Genes – Species and Part – whole relationships. This relationship is indicated in order that users may make the transition from a first access point treated terms or access point. It helps the indexes to select the most specific terms to describe a concept that is available in the thesaurus. The relationships are displayed in the thesaurus by using the symbol BT and NT.

Preferential Relationships:

These are generally indicate preferred terms or descriptors and thus perform the necessary function of controlling index vocabulary. By specifying which terms are to be used in the index and which are not. Preferential relationship in thesaurus is indicated by UF (use for) and USE.

Associative or Affinitive Relationships:

Affinitive relationships exist between terms that are not necessarily connected to one another in any fixed hierarchical manner. This relationship in a thesaurus is displayed by the symbol RT Related Term.

For example, instruction is a related concept to education Teaching, and Courses:

Instruction

RT Education

Teaching

Courses

The UNISIST Guidelines State that “the associative relation is usually employed to cover the other relations between concepts that are related but are neither consistently hierarchical nor equivalent “.

An alternative means of expressing an affinitive relationship is found in MESH, i.e. Medical Subject headings the term “See related” is the equivalent of “RT”.

Most thesauri make a clear distinction. But when BT and NT are reasonably simple to identify, related terms present a much more complex issue and no system has yet succeeded in defining precisely which terms should be enumerated as RT to any other terms.

Thus a thesaurus provide, the control of terminology by showing a structured display of concepts supplying for each concept all terms that might express that concept and presenting the associative and hierarchical relationships of the vocabulary. A list of terms which does not include structural and relational information is not a thesaurus. It is merely an alphabetical list of descriptors or subject heading.

Parts of Thesaurus:

A thesaurus usually has at least two major parts.

They are main part and an auxiliary part. The main part in a thesaurus is normal alphabetical list of all descriptors giving complete information each descriptor including the concept relationship. This part includes both descriptors and non-descriptors along with scope notes and definitions.

Auxiliary part:

In order to improve the access to the main part a thesaurus may contain several auxiliary parts. i.e., permutation subject index, systematic listings like Hierarchical Index, and subject category Index as in TEST Thesaurus of Engineering and Scientific Terms. The thesaurus may also contain a faceted classification along with alphabetical thesaurus as in Thesaurofacet.

Construction of thesaurus:

Thesauri have been compiled in a variety of different ways.

Aitcheson and Gilchrist have given same practical steps for construction of a thesaurus.

They are :

- Identify the needs of the users.
- Define the subject field.
- Decide the type and design of thesaurus layout.
- Collect terms from subject literature, users, specialist.
- Screen and edit the terms as per the rules of the thesaurus.
- Record the terms in the term cards thesaurus form.
- Sorting and grouping of thesaurus cards.
- Prepare a hierarchical structure and other associated parts.
- Test the thesaurus against a selected collection of document.
- Get the thesaurus evaluated by subject specialists and users.

Let us further discuss these ten steps in detail:

1. Identification of the subject field: The nature of the users, their, type, the needs of a library and information system influence the construction of terminology.

2. Definition of subject field: the boundaries of a topic coverage and the depth to which various aspects of the subject are to be treated may be decide. Study of the general or technical thesauri already available may give some idea about the main sections of the subject. Form and geographical divisions could be taken from any major classification system or thesauri.

3. Decide the type and design of the thesaurus layout: what amount of natal is expected from thesaurus? At this stage, the compiler should clarify his ideas on the type of thesaurus to be constructed. Decision should be taken on the characteristics of thesaurus hierarchical and other interrelation ships of the terms.

4. Collect the terms from subject literature, users and subject specialists. Once the design and layout are decided, the collection of terms for these subject areas may be taken up. Sources like descriptor lists or thesauri, subject heading lists, classification schemes, nomenclatures of indexes of journals and obstructing periodical etc. should be scanned for an initial list of the terms.

5. Screening and editing the terms as per rules of thesaurus soergel suggests that the indexing would be most effective if done by a number of different subject experts in the same field using terms of their choice. The subject experts may be shown the list of terms in their own subject fields and asked to comment, making amendments and adding term. This will help to screen and edit the terminology. The rules of thesaurus construction (Ex. UNISIST Guidelines for the establishment and Development of monolingual Thesauri) as decided in the beginning should be followed consistently.

6. Recording of terms:

In a machine-selected thesaurus, the listing of terms will be printed or displayed by the computer. If the thesaurus includes humanly selected terms, it is convenient to record term on a card. Each term written on card should show RT BT NT VF etc and SN when necessary to determine the status of term for inclusion.

7. Sorting and Grouping of thesaurus cards:

All the thesaurus cards are to be sorted out and grouped according to their subject groups and sub-groups. Duplicate entries are to be eliminated in the preliminary scanning.

8. Prepare the Hierarchical structure and other Associated posts: From the group of terms, inter term relationships are to be identified and hierarchical structures are to be developed among the descriptors. Facet analysis helps to identity and display of underlying structures.

9. Test the thesaurus against a selected collection of Documents: any thesaurus, thus compiled should be tested against a selected collection of documents of the concerned subject field to examine its efficiency and use in information retrieval.

10. Evaluation of thesaurus: user satisfaction may be helpful in determining the quality of a thesaurus; based on the feedback of the user the thesaurus could be refined.

Salton's five principles of thesaurus construction:

1. No very rare concepts should be included in the thesaurus since they could be expected to produce many matches between documents and search requests.

2. Very common high frequency terms should also be excluded from the thesaurus.

3. Non-significant words should be studied carefully before they are included in the list of words to be eliminated.

4. Ambiguous terms should be included only for the senses that are likely to be present in the document collection to be treated.

5. Each concept class should include only terms of roughly equivalent frequency so that the matching characteristics are approximately the same for each term within a category.

Relationships among terms:

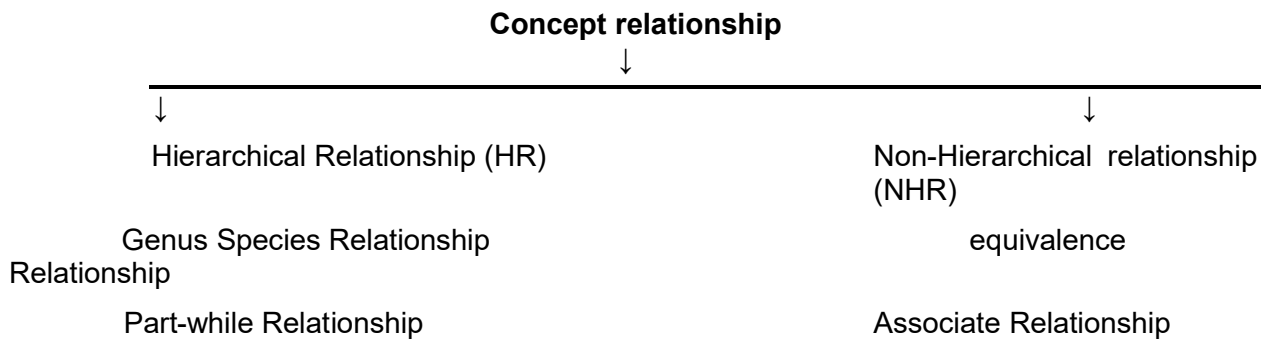
Thesaurus is a controlled list terminology for a given subject area and shows relationships among the terms.

These may be categorized as

- i. hierarchical Relationship; and
- ii. Non- hierarchical relationship.

In hierarchical relationship , there are two types, namely, genus species and part-whole. The non- hierarchical relations can be further divided into equivalence and associative relationship.

These relationships can be shown in the following chart;



Hierarchical Relationships:

The HR for a given concept arises from its super ordinate and subordinate links, that is either genus special and whole- part. In the thesaurus the super ordinate link is represented as Broader Term (BT) and the subordinate link as Narrower Term(NT).

Example:

Gens – species hierarchical relationship:

RADIATION

NT Electromagnetic radiation

ELECTRONIC RADIATION

BT Radiation

Whole – part hierarchical relationship:

INDIA

NT Tamilnadu

TAMILNADU

BT India

Non – Hierarchical Relationship (NHR)

Equivalence Relationship:

The equivalence relationship implies the control of the synonymy through preferred and non- preferred terms. The preferred terms are the terms used consistently to represent the concepts when indexing. The non-preferred terms are the synonyms or vasi – synonyms of terms, which are not is used for preferred terms and UF (Used for) for non- preferred terms

Ex: DEVELOPING COUNTRIES

UF Underdeveloped Countries.

Underdeveloped Countries

USE DEVELOPING COUNTRIES.

Associate Relationship:

The association relationship is the non- hierarchical relationships among the preferred terms.

Two kinds of terms can be linked by the associative relationship:

- a) Those that belong to same category, and
- b) Those belonging to different categories.

a) Terms belonging to the same category:

Example: SHIPS

BOATS

BT: VEHICLES

BT:VEHICLES

RT: BOATS

RT SHIPS

B) Terms belonging to different categories:

1. Process and resulting product.
 - i. Ex Cooking
 - ii. RT Food.
2. Situation or condition and what may occur in that situation
 - i. Ex: Inflation
 - ii. RT Price rise.
3. An action and the product of the action:
 - i. Ex: WEAVING
 - ii. RT: Cloth
4. Effect and cause
 - i. Ex. Poverty
 - ii. RT Unemployment

5. Process and person usually associated with the process

Ex. Administration of Justice

RT Judge

Role of thesaurus in information retrieval:

Thesauri have been used for over decades to improve precision and recall in information retrieval. The performance of an information retrieval system can be improved by the use of controlled vocabularies from a thesaurus in the concerned field of knowledge.

Thesauri will have a keystone role to play in systems of integrated database. Databases cannot be regarded as integrated unless vocabulary can be used to achieve a degree of compatibility.

Example

Following is a short list of word famous thesaurus:

1. Thesaurus of engineering and scientific terms (TEST).
2. Thesauro facet. A thesaurus and façade classification for engineering and related subjected.
3. Information retrieval thesaurus of education terms.
4. Roots thesaurus.
5. OECD macro thesaurus.
6. UNESCO thesaurus.
7. INIS thesaurus.
8. MESH Medical subject heading.
9. SPINES thesaurus published by UNESCO.
10. INSPEC Thesaurus.

References:

1. Aitchison, J and Gilchrist. Thesaurus construction: a practical manual. London: ASLIB 1972.
2. Fosketac The subject approach to information, 4 th ed. London: clive Bingley.1982
3. Guha, B . Documentation and information; services, techniques and systems, 2nd ed. Calcutta; the work press, 1983.
4. Tonley, Helen M. and Ralph.D Gee. Thasaurus making, London; Andrew devtsch 1980.
5. Indiragandhi National open University, Information procession and retrieval (MLIS-3 Block -1; Unit 4; thesaurus) New Delhi IGNOU, 1985.
6. BR Ambedkar Open University, Information Process and Retrival (MLIS Block-II, Unit 8 ; Thesourus – its structure, functions and construct) BRA, Hyderabad 1998
7. Riaz, Muhammad, Advanced Indexing and abstracting practices New Delhi: Atlantic publisher 1989.
8. Lancaster, F.W. Vocabulary contron for information Retrivel. Washing for D.C; Information resource pross, 1972.
9. Vickery, B-C Technique of information retrieval, London: Butter Worts, 1970.
10. UNESCO: Guidelines for the Establishment and development of monolingual thesauri, 2nd Rev. ed. Paris: UNESCO, 1973.

Lesson 5: ISBD

Structure

0. Objective
1. Introduction
2. History of ISBD
 - 2.1 Revision
3. Scope, Purpose and Use
 - 3.1 Scope
 - 3.2 Purpose
 - 3.3 Use
4. Structure of an ISBD Record
 - 4.1 Outline of the ISBD
 - 4.2 Punctuation
 - 4.3 Sources of Information
 - 4.4 Example Records of ISBD (M)
5. Conclusion
6. References

0. Objective

The objective of this lesson is to explain the basic features of International Standard Bibliographic Description (ISBD). It explains the evolution of ISBD for different forms of library materials. The basic structure of the ISBD is clearly dealt with in this lesson. The organization of different bibliographic elements and their order with prescribed punctuation is explained in detail.

After studying this lesson you will be able to

- What is ISBD
- History of ISBD
- Outline of ISBD

1. Introduction

The **International Standard Bibliographic Description (ISBD)** is intended to serve as a principal standard to promote universal bibliographic control. ISBD is a set of rules produced by the International Federation of Library Associations and Institutions (IFLA) to create a bibliographic description of all published literature, in a standard, Internationally acceptable human-readable form, especially for use in a bibliography or a library catalog. The ISBD main goal is to offer consistency when sharing bibliographic information. The ISBD is the standard that determines the data elements to be recorded or transcribed in a specific sequence as the basis of the description of the resource being catalogued. In addition, it employs prescribed punctuation as a means of recognizing and displaying these data elements and making them understandable independently of the language of the description. A consolidated edition of the ISBD was published in 2007 and revised in 2011, superseding earlier separate ISBDs for different types of documents such as monographs, serials, cartographic materials, electronic resources, non-book materials, and printed music. IFLA's ISBD Review Group is responsible for maintaining the ISBD.

2. History of ISBD

The International Standard Bibliographic Description (ISBD) dates back to 1969, when IFLA Committee on Cataloguing sponsored an International Meeting of Cataloguing Experts. The first ISBD(M) for Monographs appeared in 1971. The ISBD(S) for Serials was published in 1974. In the same year 'First standard edition' of ISBD(M) was also published incorporating the suggestions received from experts. The ISBD(G) a general International standard bibliographic description suitable for all types of materials was published in the year 1977. The ISBD(M) was also revised and published as 'First standard edition revised' in the year 1978. The other ISBDs as detailed below are published subsequently:

ISBD(CM) for cartographic materials	1977
ISBD(NBM) for non-book materials	1977
ISBD(S) revised for Serials	1977
ISBD(A) for Antiquarian (older monographs)	1980
ISBD(PM) for printed music	1980

2.1 Revision

The IFLA committee met in 1977 and decided to review all ISBDs for every five years and an ISBD Review committee was formed. This committee met in 1981 to plan for reviewing and revising the ISBDs. Consequently the ISBDs were republished as follows:

ISBD(M)	1987
ISBD(CM)	1987
ISBD(NBM)	1987
ISBD(S)	1988
ISBD(CF) for computer file	1990
ISBD(A)	1991
ISBD(PM)	1991
ISBD(G)	1992
ISBD(ER) for Electronic resources	1997 as replacement of ISBD(CF)

Following publications are the outcomes of "Second General Review project"

ISBD(CR) for serials and other continuing resources	2002 as a replacement of ISBD(S)
ISBD(M)	2002
ISBD(G)	2004

A consolidated edition of the ISBD was published in 2007 and revised in 2011, superseding earlier separate ISBDs for different types of documents such as monographs, serials, cartographic materials, electronic resources, non-book materials, and printed music. IFLA's ISBD Review Group is responsible for maintaining the ISBD.

3. Scope, Purpose and Use

3.1 Scope

ISBD specifies the requirements for the description and identification of the most common types of published resources that are likely to appear in library collections. It also assigns an order to the elements of description and specifies a system of punctuation for description of resources.

3.2 Purpose

The primary purpose of the ISBD is to provide stipulations for compatible descriptive cataloguing worldwide in order to facilitate International exchange of bibliographic records.

ISBD aims to:

- make records from different sources interchangeable
- assist in the interpretation of records across language barriers
- assist in the conversion of bibliographic records to electronic form
- enhance interoperability with other content standards

3.3 Use

The ISBD provides maximum number of elements to describe the resources for various bibliographic activities. However, ISBD designated three categories of elements as detailed below:

Mandatory elements: These are required in all bibliographic activities and presence of these elements is compulsory.

Optional elements: These elements may be included or omitted depending up on the discretion of the bibliographic agency.

Conditional: These elements are required under certain conditions. If the condition does not warrant the presence of these elements, their use is optional.

4. Structure of an ISBD record

The ISBD prescribes eight areas of description. Each area, except area 7, is composed of multiple elements with structured classifications. Elements and areas that do not apply to a particular resource are omitted from the description. Standardized punctuation (colons, semicolons, slashes, dashes, commas, and periods) is used to identify and separate the elements and areas. The order of elements and standardized punctuation make it easier to interpret bibliographic records when one does not understand the language of the description.

4.1 Outline of the ISBD

Area	Prescribed punctuation	Element	Usage	Repeatability
1. Title and statement of responsibility area	[] = : / ;	1.1 Title proper	M	R
		1.2 General material designation. GMDs	O	
		1.3 Parallel title	C	
		1.4 Other title information	C	
		1.5 Statements of responsibility	C	
2: Edition area	=	First Statement	M	R
		Subsequent statement	C	
2: Edition area	=	2.1 Edition statement	M	R
		2.2 Parallel edition	O	

	/ ;	2.3 Statement of responsibility relating to edition First statement Subsequent statement	M O	R
3. Material or type of resource specific area			M	R
4: Publication, production, distribution, etc., area	: ,	4.1 Place of publication 4.2 Name of publisher 4.4 Date of publication	M M M	
5. Physical description area	: , +	5.1 Specific material designation and extent 5.2 Other physical details 5.3 Dimensions 5.4 Accompanying material	M M M O	R
6: Series area	= : / , ;	6.1 Series 6.2 Parallel title of the series 6.3 Other title of the series 6.4 Statement of responsibility 6.5 ISSN 6.6 Numbering within series	M C C C M C	R R R
7: Notes area			C	R
8: Resource identifier (e.g. ISBN, ISSN) and terms of availability area	= : ()	8.1 Resource identifier 8.2 Key title 8.3 Terms of availability and price 8.4 Qualifications	M C O O	R R

4.2 Punctuation

Each element of the description, except the first element of area 1, is either preceded or enclosed by prescribed punctuation. Prescribed punctuation is preceded and followed by space except for coma (,) and point (.) which are only followed by space.

Each area of the ISBD other than area 1 is preceded by a point, space, dash, space (. -). Each area can also be written as separate paragraph without ant preceding punctuation.

When the first element of an area is not present in a description, the prescribed punctuation of the first element that is present is replaced by a point, space, dash, space (. -), preceding the area.

When an element or area ends with a point and the prescribed punctuation for the element or area that follows begins with a point, both points are to be given

Eg. 3rd ed.. –
Not 3rd ed. –

The mark omission of information is indicated by three points (...). The mark of omission is preceded and followed by a space.

4.3 Sources of Information

For all types of material the resource itself constitutes the basis of the description. The ISBDs provided the preferred sources of information for describing all types of resources. For mono graphs the title page is the preferred source of information.

4.4 Example Records of ISBD (M)

A typical ISBD records:

Record 1

A manual for writers of research papers, theses, and dissertations : Chicago style for students and researchers / Kate L. Turabian ; revised by Wayne C. Booth, Gregory G. Colomb, Joseph M. Williams, and University of Chicago Press editorial staff. — 7th ed. — Chicago : University of Chicago Press, 2007. — xviii, 466 p. : ill. ; 23 cm. — (Chicago guides to writing, editing, and publishing). — Includes bibliographical references (p. 409-435) and index. — ISBN 978-0-226-82336-2(cloth : alk. paper) : USD35.00. — ISBN 978-0-226-82337-9 (pbk. : alk. paper) : USD17.00

Record 2

Theory of classification / Krishan Kumar. —
Delhi : Vikas, 1979. — xii, 510p ; 22 cm.
ISBN 0 7069 0797 3 Cloth: Rs. 60

Record 3

Education and the social order / Bertrand Russel. — Revised and enlarged ed. —
London : George Allen & Unwin, 1975. — viii, 301p ; 22 cm.
Previous edition published as: 'Education and the modern world.'
London : Van Nostrand, 1932
ISBN 0 486 22876 9 Paperback: \$5.00

5. Conclusion

The application of computer has brought home the great importance of uniformity; precision; and compatibility of bibliographic tools like library catalogues. This emphasizes the need for standardization. International Standard Bibliographic Descriptions (ISBDs) developed by IFLA serve a good example of an attempt towards uniform cataloguing practices and achieving successful and convenient international exchange of bibliographic information in written as well as machine-readable form. The acceptance of ISBD by all libraries at national and international level is one step towards making Universal Bibliographic Control (UBC) a reality.

6. References

- Chan, Lois Mai. *Cataloging and Classification: an introduction*. New York: McGraw-Hill Humanities, 1994.
- *International Standard Bibliographic Description (ISBD)*. Preliminary consolidated ed. München: K.G. Saur, 2007. (IFLA series on bibliographic control, vol. 31)
- Svenonius, Elaine. *The Intellectual Foundation of Information Organization*. Boston: The MIT Press, 2000.
- Willer, Mirna; Dunsire, Gordon; Bosancic, Boris (2010). "ISBD and the Semantic Web". *JLIS.it* (University of Florence) 1 (2).doi:10.4403/jlis.it-4536. Retrieved 29 June 2013.

Lesson 6: AACR2

Structure

0. Objective
1. Introduction
2. AACR Governance Structure
3. Organization of description
4. Levels of Detail in Description
5. Structure of AACR2
 - 5.1 Part I of AACR2
 - 5.2 Part II of AACR2
6. Appendices:
7. Conclusion

0. Objective

The objective of this lesson is to explain the basic features of Anglo-American Cataloguing Rules 2. It explains the basic structure of the code, organization of description of the library materials and levels of description etc.

After studying this lesson you will be able to

- What is AACR2
- Governance of AACR
- Organization of description of documents
- Levels of description recommended by AACR2
- structure of AACR2

1. Introduction:

The **Anglo-American Cataloguing Rules** (AACR) are a national cataloging code first published in 1967. The rules in the code are based on "Statement of Principles" adopted by the International Conference of Cataloguing Principles in 1961. The rules in AACR1 have been formulated primarily to meet the requirements of general research libraries. Though it is not developed as an International code, it is being widely used in various countries. AACR2 stands for the *Anglo-American Cataloguing Rules, Second Edition*. Despite the claim to be 'Anglo-American', the first edition of AACR was published in 1967 in somewhat distinct North American and British texts. The second edition of 1978 unified the two sets of rules (adopting the British spelling 'cataloguing') and brought them in line with the International Standard Bibliographic Description.

AACR2 exists in several print versions, as well as an online version. Gorman has edited several revisions of AACR2 including a concise edition. Print versions are available from the publishers. The online version is available only via Cataloger's Desktop from the Library of Congress. Various translations are also available from other sources.

Principles of AACR include cataloguing from the item 'in hand' rather than inferring information from external sources and the concept of the 'chief source of information' which is preferred where conflicts exist.

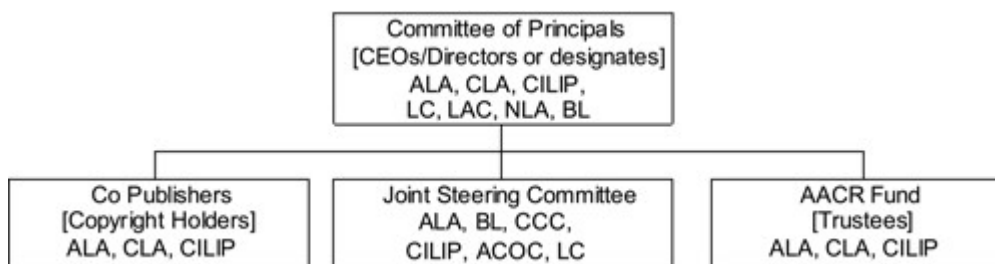
Over the years AACR2 has been updated by occasional amendments, and was significantly revised in 1988 (2nd edition, 1988 revision) and 2002 (2nd edition, 2002 revision). The 2002 revision included substantial changes to sections for non-book materials. A schedule of annual updates began in 2003 and ceased with 2005. AACR is published in English and has been translated into other languages

AACR2 has been succeeded by Resource Description and Access (commonly referred to as RDA), which was released in June 2010. This new code is informed by the Functional Requirements for Bibliographic Records and was conceived to be a framework more flexible and suitable for use in a digital environment. In the fall of 2010, the Library of Congress, National Library of Medicine, National Agricultural Library, and several other institutions and national libraries of other English-speaking countries performed a formal test of RDA, the results of which were released in June 2011.

2. AACR Governance Structure

AACR2 is published under the auspices of the **AACR Fund**. The publication of the Anglo-American Cataloguing Rules (AACR) is governed by the Committee of Principals, which coordinates three subordinate groups: The Co-Publishers of AACR, The Joint Steering Committee for Revision of AACR, and The AACR Fund Committee (Trustees). Following is the AACR governance structure:

AACR Governance Structure



The **Anglo-American Cataloguing Rules2** (AACR2) are designed for use in the construction of catalogues and other lists in general libraries of all sizes. The rules cover the description of, and the provision of access points, for all library materials commonly collected at the present time. The second edition of the Anglo-American cataloguing Rules appeared in 1978 is based on a reconciliation of the British and North American texts of the 1967 edition. It is published jointly by the American Library Association, the Canadian Library Association, and the Chartered Institute of Library and Information Professionals in the UK. The editor is Michael Gorman. AACR2 is designed for use in the construction of catalogues and other lists in

general libraries of all sizes. The rules cover the description of, and the provision of access points for, all library materials commonly collected at the present time.

3. Organization of description:

The description is divided into the following areas:

Title and statement of responsibility

Edition

Material Specific details

Publication, distribution, etc

Physical description

Series

Note

Standard number and terms of availability

Each are is further divided into a number of elements.

4. Levels of Detail in Description

The AACR2 code has prescribed three levels of description. First level provides the minimum information which is necessary to identify a given document. Second level can be called standard description. It provides all the data which may be considered necessary for the description of documents. The third level provides information covering every descriptive element described in the code. The choice of level of description would depend uponj the purpose to be satisfied by a given catalogue in the library.

5. Structure of AACR2

The structure of AACR2 is based on ISBD, which make distinction between bibliographic description and the access points. The AACR2 consists of two parts. Part I deals with the provision of information describing the item being catalogued, and Part II deals with the determination and establishment of headings (access points) under which the descriptive information is to be presented to catalogue users, and with the making of references to those headings. This part contains rules for choice of main and added entries in chapter 21 and form of headings and uniform titles in chapters 22 to 25 and references in chapter 26 of the code. In both parts the rules proceed from the general to the specific.

5.1 Part I of AACR2

AACR2 has used general framework for all bibliographic descriptions called ISBD(G). This has enabled it to achieve an effective and consistent approach to bibliographic description. All the rules for description confirm to a single set of punctuation conventions. The same principles have been adopted for description of different kinds of library materials. The rules in part I deal with description of all types of library information materials.

Part I of AACR2 consists of 13 chapters. Chapter 1 provides general rules for description of all materials collected now-a-days by libraries. Chapters 2 to 10 provide rules for describing specific types of materials as given below:

Chapter 2	: Books, pamphlets and printed materials
Chapter 3	: Cartographic materials
Chapter 4	: Manuscripts
Chapter 5	: Music
Chapter 6	: Sound Recordings
Chapter 7	: Motion pictures and video recordings
Chapter 8	: Graphic materials
Chapter 9	: Machine-readable data files
Chapter 10	: Three dimensional artifacts and Realia
Chapter 11	: Microforms
Chapter 12	: Serials
Chapter 13	: Analysis

In chapters 2-10 the areas and elements prescribed by ISBD(G) and their prescribed punctuations have been described in terms of the particular type of library material. Content of an element which is special to a particular type of material has been dealt in detail in the concerned chapter. In case, general rule is applicable reference is made to that of chapter 1.

Chapters 11-13 deal with microforms, serial and analysis which are of partial general in nature. In some cases the rules provided in these chapters have modified the provisions given in preceding chapters. In other cases the rules given in these chapters are supposed to be used along with the rules described in preceding chapters. Analysis is conceived as a process by which a part of publication is described is related to the whole document. The rules for preparing analytical entries and multi-level descriptions are provided in chapter 13 which is entitled 'Analysis'.

The rules in each chapter consist of:

0. Preliminary rules
1. Title and statement of responsibility
2. Edition
3. Material Specific details
4. Publication, distribution, etc
5. Physical description
6. Series
7. Note
8. Standard number and terms of availability
9. Supplementary items
10. Items made up of several types of material
11. Facsimiles, photocopies, and other reproductions

5.2 Part II of AACR2

Part II deals with Headings, Uniform titles, and References. It consists of 6 chapters numbered 21 to 26. Chapter 21 deals with access points. Chapters 22-24 deal with headings for persons, Geographic names and corporate bodies respectively. Chapter 25 deal with Uniform titles and rules for giving different types of references are dealt in chapter 26. In each chapter general rules precede special rules. In case, no specific rules exist for a given specific problem, then general rules are to be applied.

Rules in chapter 21 provide an integrated approach to all library materials in deciding access points. The concept of mixed responsibility is also dealt in detail in this chapter. Rules for added entries also provided in this chapter.

6. Appendices:

AACR2 code provide four appendices dealing with capitalization, abbreviations, the treatment of numerals and glossary of terms used in cataloguing.

7. Conclusion:

Anglo-American Cataloguing Rules2 (AACR2) is widely used catalogue code throughout the world for describing the documents for the purpose of catalogue. It aims to respond to the changes that have been taking place in different aspects of librarianship. The code made an attempt to make the rules as contemporary as could be possible in the circumstances. The entire code is based on ISBD and in accordance with "Paris Principles" on cataloguing rules. An attempt has been made to make the code more accessible to cataloguers and bibliographers in language and articulation. AACR2 has not been able to make all the provisions required for library automation. However, built on foundations established by the Anglo-American Cataloguing Rules, RDA is being developed as a new standard designed for use in a digital environment. RDA will be co-published by the American Library Association, the Canadian Library Association and the Chartered Institute of Library and Information Professionals (CILIP).

LESSON 7: COMMON COMMUNICATION FORMAT (CCF)

STRUCTURE

0. Objective
1. Introduction
2. Origin of CCF
3. CCF Record Structure
 - 3.1 Record Label
 - 3.2 Directory
 - 3.3 Data Fields
 - 3.4 Record Separator
 - 3.5 Sample CCF Record
4. Criticism on CCF
5. CCF Tags
6. Conclusion
7. References

0. OBJECTIVE

The objective of this lesson is to explain the basic features of Common Communication Format (CCF). It explains the basic structure of the format, organization of description of the library materials in the format.

After studying this lesson you will be able to

- What is CCF
- Origin of CCF
- Structure of CCF
- Criticism on CCF

1. INTRODUCTION

CCF (Common Communication Format) is a format for the exchange of bibliographic records. CCF is devised by taking into consideration the major existing International exchange formats. There were different practices in record creation using different formats such as UNIMARC, ISDS, MEKOF-2, ASIDIC, USSR-US CCF etc., resulting in records which, when merged into a database, will show their different origins. To achieve homogeneity and avoid ambiguity a small number of data elements were identified which were used by virtually all information-handling communities, including both libraries and abstracting and indexing organizations. These commonly used data elements formed the core of the CCF. The Common Communication Format (CCF) is to provide a detailed and structured method for recording a number of mandatory and optional data elements in a computer-readable bibliographic record for exchange purposes between two or more computer-based systems. However, it can also be useful within non-computerized bibliographic systems.

A technique was developed to show relationships between bibliographic records, and between elements within bibliographic records. The concept of the record segment was developed and refined, and a method for designating relationships between records, segments,

and fields was accepted. The first edition of CCF: The Common Communication Format was published by UNESCO in 1984.

CCF provides rules to achieve consistency, uniformity and compatibility between more than one computer systems. Within an information system, the record which forms the database will usually exist in a number of separate but highly compatible formats. At the very least there will be:-

- a format in which records will be input to the system.
- a format best suited to long term storage.
- a format to facilitate retrieval, and
- a format in which records will be displayed.

In addition, if two or more organizations wish to exchange records with one another, it will be necessary for each of these organizations to agree upon a common standard format for exchange purposes. Each must be able to convert to an exchange format record from an internal-format record, and vice-versa.

2. ORIGIN OF CCF

In April 1978 the Unesco/PGI (i.e. Unesco General Information Programme) sponsored an International symposium on Bibliographic exchange formats, which was held in Taormina, Sicily. The symposium was convened to study the desirability and feasibility of establishing the maximum compatibility between existing bibliographic exchange formats.

Immediately after this symposium UNESCO/PGI formed an Adhoc group to establish a Common Communication Format. At the start itself, the Adhoc group has decided certain principles which the CCF still follows like:-

1. The structure of the format to conform the international standard ISO – 2709.
2. The core record to consist of a small number of mandatory data elements essential for bibliographic description are identified in a standard manner.
3. The mandatory elements are augmented by the well identified optional data elements.

3. CCF RECORD STRUCTURE

Each CCF record consists of four major parts, i.e.

1. Record Label (having 24 Characters);
2. Directory (having Variable Length);
3. Data fields (having Variable Length);
4. Record Separator (having 1 Character).

3.1 Record Label

Each CCF record begins with a fixed label of 24 characters. The contents of which are as follows:-

Character Positions	Assigned No. of Characters	Contents
0-4 th	5 Chrs	Record Length- The length of the record includes the label, directory, data fields and

		record separator. (Use of 5 characters for the record length permits records as long as 99,999 characters)
5 th	1 Chr.	Record Status- Using a code taken from the list of Record Status Codes from CCF manual (Pg. No.143)
6 th	1 Chr.	Blank
7 th	1 Chr.	Bibliographic Level- Codes are given in CCF manual on Page No. 144.(Character position for bibliographical level is not used in factual record & is filled with space).
8-9 th	2 Chrs	Blank
10 th	1 Chr.	Indicator Length- Used to fix the length of the indicator.
11 th	1 Chr.	Subfield Identifier Length- eg., ^a, ^b, ^c, ^d, ... etc.
12-16 th	5 Chrs	Base Address of Data- the location within the record at which the first data field begins.
17-19 th	3 Chrs	Blank
20 th	1 Chr.	Length of the 'Length of data field' in each directory entry- use of 4 characters here permits data field as long as 9,999 characters.
21 st	1 Chr.	Length of 'Starting Character position in each directory entry- normally we use 5 because for getting 10,000 th (9,999+1th) position we need 5 characters only.
22 nd	1 Chr.	Length of 'Implementation defined' section of each directory entry- Normally unused (Blank i.e., Zero)
23 rd place	1 Chr.	Blank.

3.2 Directory

A single Directory Entry Contains:

24-26 th	3 Chrs	Tag- By using 3 Characters we can use 999 possible tags.
27-30 th	4 Chrs	Length of data field- A four digit number showing how many characters are occupied

by the data field, including indicators & data field separator but excluding the record separator code if the data field is the last field in the record.

31-35 th	5 Chrs	Starting Character Position- It gives the position of the first character of the next data field relative to the base address of the data.
36 th	1 Chr.	Segment Identifier- Chosen from 0-9 &/or A-Z to designate the data field as being a member of particular segment(Normally not used).
37 th	1 Chr.	Occurrence Identifier- Chosen from 0-9 &/or A-Z, which differentiate multiple occurrences of data fields that carry the same tag within same record (Normally not used).

3.3 Data Fields

A single Data Field Contains:

38-39 th	2 Chrs	For Indicators
40-41 st	2 Chrs	For Subfield Identifiers
Variable	Variable	Subfield
	1 Chr	Field Separator

3.4 Record Separator

Each record is separated by one character record separator

At End	1 Chr.	Record Separator.
--------	--------	-------------------

3.5 Sample CCF Record

For Example (CCF Record Structure):-

001610000022000720004500 □ Record Label200004000000 □ Tag(200), Length of Data field (0040), Starting Character Position of Data(00000).300002000040 Tag(300), Length of Data field (0020), Starting Character Position of Data(00040).400001600060 □ Tag(400), Length of Data field (0016), Starting Character Position of data(00060).440001100076 □ Tag(440), Length of Data field (0011), Starting Character Position of Data(00076).# □ Field Separator
^aProlegomena to Library Classification#^aRanganathan^bS.R.#^aDelhi
^bJaypee#^a19990000# □ Field Separator# □ Record Separator.

4. CRITICISM ON CCF:

Many of the disadvantages of CCF are based on the disadvantages of the Cataloguing Codes. The CCF is basically a tag code to facilitate data exchange between two or more systems. It should be independent of the cataloguing codes.

The field Physical description (460) which describes the physical attributes of the item cannot be used to produce catalogue cards following different codes. For example: AACR insists that if there are papers numbered in roman numerals they should be taken into consideration. However, CCC does not insist on this. The problem is that the catalogue codes not only prescribe what should be descriptive element, they have formulated rules on how they should be represented.

While exchanging the data the field Place and Name of the Distributor (420) is of no use to the other system, because the distributor can vary from place to place. It can be argued that this field is provided for internal use. But there is no field to provide accession number which is frequently used for internal purposes.

The subfield "B" which describes statement of responsibility in the field Title and Statement of responsibility (200) and the field Name of Person (300) is an overlapping concept to each other.

Field Segment linking fields (080,081,082,083,084,085) does not give any clue about the linking of documents or records. Only Field to Field linking (086) makes some sense.

The Field Source of Record (020) is a non-repeatable one. But if the database is merged with a master database, then it may be repeated disputing the concept of non-repeatability.

The major disadvantage in CCF is the different codes used for data elements. In this case CCF seems to be consistently inconsistent.

For example: for the subfield 'A' which describes language of the record in Field Language and Script of the record (031) they have followed a code list from ISDS manual. The codes are in alphabets. The same problem surfaces in the fields with tag numbers 040, 041, 200,201, 210, 220, 221, etc., where ever the question of language arises. The same type of problem is there in Record Status Codes, bibliographic data level codes, completeness of record codes. But these problems does not appear in character set codes like, physical medium codes, role codes, type of material codes, by using numerals instead of alphabets. Here the consistency fails.

Again no Code is given for subfield 'B' of tag 110 which describes national bibliographic agency code in the field National Bibliography Number (110) and sub field 'B' which describes legal deposit agency code in field Legal Deposit Number (111).

5. CCF TAGS

Following are the Common Communication Format (CCF) Tag Numbers

Tag Name

001 Record Identifier
 010 RECORD IDENTIFIER USED IN A SECONDARY SEGMENTS
 010A Identifier

 011 ALTERNATIVE RECORD IDENTIFIER (R)
 011A Alternative identifier
 011B Identification of Agency in coded form
 011C Name of agency

 015 BIBLIOGRAPHIC LEVEL OF SECONDARY SEGMENT
 015A Bibliographic level

 020 SOURCE OF RECORD
 020A Identification of agency in coded form
 020B Name of agency
 020C Name of code set
 020D Rules for bibliographic description
 020L Language of name of agency

 021 COMPLETENESS OF RECORD
 021A Level of completeness code

 022 DATE ENTERED ON FILE
 022A Date

 023 DATE AND NUMBER OF RECORD VERSION
 023A Version date
 023B Version number

 030 CHARACTER SETS USED IN RECORD
 030A Alternative Control set (C1)
 030B Default Graphic set (G0)
 030C Second Graphic Set (G1)
 030D Third Graphic Set (G2)
 030E Fourth Graphic Set (G3)
 030F Additional Control Set (R)
 030G Additional Graphic Set (R)

 031 LANGUAGE AND SCRIPT OF RECORD (R)
 031A Language of the record(R)
 031B Script of the record(R)

 040 **LANGUAGE AND SCRIPT OF ITEM(R)**
 040A Language of item (R)
 040B Script of item
 041 Language and Script of Summary (R)
 041A Language of the summary (R)
 041B Script of the summary

- 050 **PHYSICAL MEDIUM**
050A Physical medium code (R)
- 060 **TYPE OF MATERIAL**
060A Type of material code (R)
- 080 **SEGMENT LINKING FIELD: GENERAL VERTICAL RELATIONSHIP (R)**
080A Segment relationship code
080B Segment indicator code
- 081 **SEGMENT LINKING FIELD: VERTICAL RELATIONSHIP FROM MONOGRAPH**
081A Segment relationship code
081B Segment indicator code
- 082 **SEGMENT LINKING FIELD: VERTICAL RELATIONSHIP FROM MULTI-VOLUME MONOGRAPHIC**
082A Segment relationship code
082B Segment indicator code
- 083 **SEGMENT LINKING FIELD: VERTICAL RELATIONSHIP FROM SERIAL**
083A Segment relationship code
083B Segment indicator code
- 085 **SEGMENT LINKING FIELD: HORIZONTAL(R)**
085A Segment relationship code
085B Segment indicator code
- 086 **FIELD TO FIELD LINKING(R)**
086A field linked from
086B Field relationship code
086C Field linked to
- 100 **INTERNATIONAL STANDARD BOOK NUMBER(R)**
100A ISBN
100B Invalid ISBN (R)
100C Qualification (R)
- 101 **INTERNATIONAL STANDARD SERIAL NUMBER (ISSN)**
101A ISSN
101B Invalid ISSN
101C Cancelled ISSN (R)
- 102 **CODEN**
102A Coden
- 110 **NATIONAL BIBLIOGRAPHIC NUMBER (R)**
110A National Bibliographic Number
110B National Bibliographic Agency Code
- 111 **LEGAL DEPOSIT NUMBER(R)**

- 111A Legal deposit number
- 111B Legal deposit agency

- 120 DOCUMENT IDENTIFICATION NUMBER(R)
- 120A Document identification number
- 120B Type of number

- 200 TITLE AND ASSOCIATED STATEMENT(S) OF RESPONSIBILITY (R)
- 200A Title (R)
- 200B Statement of responsibility associated with title (R)
- 200L Language of title
- 200S Script of title

- 201 KEY TITLE
- 201A Key title
- 201B Abbreviated key title
- 201L Language of key title
- 201S Script of key title

- 210 PARALLEL TITLE AND ASSOCIATED STATEMENT(S) OF
RESPONSIBILITY (R)
- 210A Parallel title
- 210B Statement of responsibility associated with parallel title (R)
- 210L Language of parallel title
- 210S Script of parallel title

- 220 SPINE TITLE (R)
- 220A Spine title
- 220L Language of spine title

- 221 COVER TITLE (R)
- 221A Cover title
- 221L Language of cover title

- 222 ADDED TITLE PAGE TITLE (R)
- 222A Added title page title
- 222L Language of added title page title

- 223 RUNNING TITLE (R)
- 223A Running Title
- 223L Language of running title

- 230 OTHER TITLE (R)
- 230A Other variant title
- 230L Language of title
- 240 UNIFORM TITLE (R)
- 240A Uniform title
- 240B Number of part(s) (R)
- 240C Name of part(s) (R)
- 240D Form subheading (R)
- 240E Language of item (as part of uniform title) (R)

- 240F version
- 240G Date of version
- 240L Language of uniform title
- 240Z Authority number

- 260 EDITION STATEMENT AND ASSOCIATED STATEMENT(S) OF RESPONSIBILITY (R)
 - 260A Edition Statement
 - 260B Statement of responsibility associated with edition (R)
 - 260L Language of edition statement

- 300 NAME OF PERSON (R)
 - 300A Entry element
 - 300B Other name elements
 - 300C Additional elements to name
 - 300D Date(s)
 - 300E Role (coded) (R)
 - 300F Role (non-coded) (R)
 - 300Z Authority number

- 310 NAME OF CORPORATE BODY (R)
 - 310A Entry element
 - 310B Other part(s) of name (R)
 - 310C Qualifier (R)
 - 310D Address of corporate body
 - 310E Country of corporate body
 - 310F Role (coded) (R)
 - 310G Role (non-coded) (R)
 - 310L Language of entry element
 - 310S Script of entry element
 - 310Z Authority number

- 320 NAME OF MEETING (R)
 - 320A Entry element
 - 320B Other part(s) of name (R)
 - 320C Qualifier (R)
 - 320E Country
 - 320G Location of meeting
 - 320H Date of meeting (in ISO format)
 - 320I Date of meeting (in free format)
 - 320J Number of meeting
 - 320L Language of entry element
 - 320S Script of entry element
 - 320Z Authority number

- 330 AFFILIATION (R)
 - 330A Entry element
 - 330B Other parts of the name (R)
 - 330C Qualifier (R)
 - 330D Address (R)
 - 330E Country of affiliation

- 330L Language of entry element

- 400 PLACE OF PUBLICATION AND PUBLISHER (R)
 - 400A Place of publication (R)
 - 400B Name of publisher
 - 400C Full address of publisher (R)
 - 400D Country of publisher (R)

- 410 PLACE OF MANUFACTURE AND NAME OF MANUFACTURER (R)
 - 410A Place of manufacture (R)
 - 410B Name of manufacturer
 - 410C Full address of manufacturer (R)
 - 410D Country of manufacturer (R)

- 420 PLACE AND NAME OF DISTRIBUTOR (R)
 - 420A Place of distributor (R)
 - 420B Name of distributor
 - 420C Full address of distributor (R)
 - 420D Country of distributor (R)

- 440 DATE OF PUBLICATION (R)
 - 440A date in formalized form
 - 440B date in non-formalized form

- 441 DATE OF LEGAL DEPOSIT
 - 441A Date legal deposit

- 450 SERIAL NUMBERING
 - 450A Serial numbering and date

- 460 PHYSICAL DESCRIPTION
 - 460A Number of pieces and designation
 - 460B Other descriptive details
 - 460C Dimensions
 - 460D Accompanying material (R)

- 480 SERIES STATEMENT AND ASSOCIATED STATEMENT(S) OF RESPONSIBILITY (R)
 - 480A Series Statement
 - 480B Statement of responsibility associated with series statement
 - 480C Part statement
 - 480D ISSN
 - 480L Language of title
 - 480S Script of title

- 490 PART STATEMENT (R)
 - 490A Volume / part numeration and designation (R)
 - 490B Pagination defining a part
 - 490C Other identifying data defining a part

- 500 NOTE (R)

- 500A Note
- 510 NOTE ON BIBLIOGRAPHICAL RELATIONSHIP (R)
- 510A Note
- 520 SERIAL FREQUENCY NOTE (R)
- 520A Frequency
- 520B Dates of frequency
- 530 CONTENTS NOTE (R)
- 530A Note
- 600 ABSTRACT (R)
- 600A Abstract
- 600L Language of abstract
- 610 CLASSIFICATION SCHEME NOTATION
- 610A Notation (R)
- 610B Identification of classification scheme
- 620 SUBJECT DESCRIPTOR (R)
- 620A Subject descriptor
- 620B Identification of subject system

6. CONCLUSION:

CCF is devised by taking into consideration the major existing International exchange formats and was intended to be used for the transfer of records between systems. This is purely an exchange format. It does not give any information about circulation system of any particular library. In spite of its disadvantages it is one of the most widely used formats especially in developing countries. But still some serious, fruitful steps should be taken to overcome its problems.

7. REFERENCES:

1. CCF: the Common Communication Format, 2nd Ed., Paris, UNESCO, 1988 (PGI-88/WS/2).
2. Ellen, Gradley and Hopkinson, Alan, Exchanging Bibliographic data: MARC and other international format. Library Association Publishing Ltd., London, 1990, pp. 209-222.
3. International Standard ISO 2709(E): Documentation- Format for Bibliographic Information Interchange on Magnetic Tape.(In Handbook on International Standards Governing Information Transfer by International Organization for Standards, 1977. pp. 291-294.

Lesson 8: Machine Readable Catalogue 21(MARC21)

Structure

0. Objective
1. Introduction
2. Record Structure
 - 2.1 Outline of Leader
 - 2.2 Outline of Record Directory
 - 2.3 Outline of Control Fields
 - 2.4 Outline of Variable fields
3. Field Designations
4. Content in a MARC Record
5. MARC formats
6. MARC 21
7. MARCXML
 - 7.1 MARCXML primary design goals included
8. Future
9. Conclusion
10. References

0. Objective

The objective of this lesson is to explain the basic features of Machine Readable Cataloguing (MARC). It explains the evolution of MARC. The basic structure of the MARC Record is clearly dealt with in this lesson. The Field designators and subfield codes and tags for the variable fields are explained. MARCXML version which is developed for web applications is also explained.

After studying this lesson you will be able to know

- What is MARC
- History of MARC
- Structure of MARC
- field designators and tags

1. Introduction

MARC (MACHINE-Readable Cataloging) standards are a set of digital formats for the description of items, such as books, patents, serials etc. catalogued by libraries. It was developed by Henriette Avram at the US Library of Congress during the 1960s to create records that can be used by computers, and to share those records among libraries. By 1971, MARC formats had become the national standard for dissemination of bibliographic data in the United States. MARC standard became the international standard by 1973. There are several versions of MARC in use around the world, the most predominant being MARC 21, created in 1999. The MARC 21 family of standards now includes formats for authority records, holdings

records, classification schedules, and community information, in addition to the format for bibliographic records.

2. Record Structure

MARC records are typically stored and transmitted as binary files. MARC uses the ISO 2709 standard to define the structure of each record. This includes a marker to indicate where each record begins and ends, as well as a set of characters at the beginning of each record that provide a directory for locating the fields and subfields within the record.

The basic machine readable catalogue record on a MARC tape consists of the Leader, the Record Directory, the Control Fields and the Variable Fields.

Leader	Record Directory	Control Fields	Variable Fields
--------	------------------	----------------	-----------------

The control field consists of both variable control number and Variable Fixed Fields. The Leader is fixed in length for all records contain 24 characters. It is a set of fields describing the general structure of the individual entry. The Record Directory is an index to the location of the Control and Variable Fields in the record. It consists of a series of fixed length entries, one for each variable field in the record.

An entry in the Record Directory contains the identification tag, the length and starting character position of each variable field in the record. The record Directory will end with a field-terminator code. Since the number of variable fields in a record can vary, the total length of the Record Directory is also variable. All fields end with field-terminator code except the last variable field which ends with record-terminator code. All Variable Fields are made up of variable length alphanumeric data. Each variable field is identified by three character numeric tag in the record directory. Tags may be repeated as required in a logical record. However, tags associated with control fields will not be repeated in a logical record.

2.1 Outline of Leader

The total number of characters in the Leader is 24, and there are nine data elements in the Leader as described below:

Name of Data Element	No. of Characters
Record Length	5
Record Status	1
Type of Record	1
Bibliographic Level	1
Blanks	2
Indicator Count	1
Subfield Count	1
Base Address of Data	5
Blank Character	7

2.2 Outline of the Record Directory

Each record directory consist of three elements viz, Tag, Length of the field and Starting character position of the field in the record. Each entry in the directory is 12 character entries as detailed below. The number of entries in the record directory corresponds to the number of variable fields present in the record. The record directory is terminated by field-terminator code.

The sample record directory is given below:

Entry	Element in the entry	No. Characters
Field 1	Tag of the field	3
	Length of the Field	4
	Starting Character Position of the field	5
Field 2	Tag of the field	3
	Length of the Field	4
	Starting Character Position of the field	5

Field n	Tag of the field	3
	Length of the Field	4
	Starting Character Position of the field	5
	Field Terminator (F/T)	1

2.3 Outline of Control Fields

Data Element - 1	Data Element - 2	Data Element - 3	-----	Data Element - n	F/T
------------------	------------------	------------------	-------	------------------	-----

Example:

Tag	Name of the Control Field	Data Elements	No. of Characters	Character position in the field
001	Library of Congress Card Number	Year	2	3-4
		Number	6	5-10
		supplement etc.	1	11
008	Fixed length Data Elements	1 data entered on file	6	0-5
		2 Type of publication	1	6
		3 Date 1	4	7-10
		.		
		.		
		10 Govt. pub. indicat	1	28
		16 Biography code	1	34
17 Language code	3	35-37		

2.4 Outline of Variable Fields

Variable field consists of indicators, subfield codes, data elements and the field terminator. Further each variable field is assigned a tag and the tag is stored in the directory. The directory, control fields and variable fields are always terminated by a field terminator. Finally the last character in the record is a Record Terminator.

Indicator
Subfield code
Data Element
Subfield code
Data Element
.
.
.
Subfield code
Data Element
F/T
Indicator
Subfield code
Data Element
.
.
.
Subfield code
Data Element
F/T
.
.
.
F/T
R/T

3. Field Designations

Indicators: Each variable field will begin with 2 character code which provides descriptive information about the field. The contents of the indicators are specified for the fields in which they are used, If the indicators are not used with a particular field, they will contain blanks.

Subfield codes: Variable fields are made up of a single data element or a group of data elements. The subfield code precedes each data element in a field and identifies the data element. The subfield code consists of 2 characters. For the purpose of these specifications, the delimiter will be represented by "\$".

Data Elements: All the data elements in the variable fields may have variable lengths. Each variable field or data element has a tag. Some fields are repeatable.

Subfield Codes in Variable Fields: The subfield code identifies the constituent data elements of a variable field. For example the imprint field, tag 260, may have the following 3 data elements in its respective subfield codes:

Place	\$a
Publisher	\$b
Date	\$c

Variable Fields: Each data element in the variable field has tag of three characters. Following are few fields with tags for the purpose of illustration.

010	LC Card number
100	Main entry Personal name
245	Title
300	Collation
650	Topical subject heading
700	Personal name added entry

4. Content in a MARC Record

MARC transmits information about a bibliographic item, not the content of that item; this means it is a metadata transmission standard, not a content standard. The actual content a cataloger will place in each MARC field is usually governed and defined by standards outside of MARC, except for a handful of fixed fields defined by the MARC standards themselves. The Anglo-American Cataloguing Rules, for example, define how the physical characteristics of books and other item should be expressed. The Library of Congress Subject Headings (LCSH) provides a list of authorized subject terms to describe the main content of the item. Other cataloging rules, subject thesauri, and classification schedules can also be used.

5. MARC formats

MARC formats

Name	Description
Authority records	Provide information about individual names, subjects, and uniform titles. An authority record establishes an authorized form of each heading, with references as appropriate from other forms of the heading.
Bibliographic records	Describe the intellectual and physical characteristics of bibliographic resources (books, sound recordings, video recordings, and so forth).
Classification records	MARC records containing classification data. For example, the Library of Congress Classification has been encoded using the MARC 21 Classification format.
Community Information records	MARC records describing a service providing agency. For example, the local homeless shelter or tax assistance provider.
Holdings records	Provide copy-specific information on a library resource (call number, shelf location, volumes held, and so forth).

6. MARC 21

MARC 21 was designed to redefine the original MARC record format for the 21st century and to make it more accessible to the international community. MARC 21 has formats for the following five types of data: Bibliographic Format, Authority Format, Holdings Format, Community Format, and Classification Data Format. Currently MARC 21 has been implemented successfully by The British Library, the European Institutions and the major library institutions in the United States, and Canada.

MARC 21 is a result of the combination of the United States and Canadian MARC formats (USMARC and CAN/MARC). MARC21 is based on the ANSI standard Z39.2, which allows users of different software products to communicate with each other and to exchange data. MARC 21 in UTF-8 format allows all the languages supported by Unicode.

7. MARCXML

In 2002, the Library of Congress developed the MARC XML schema as an alternative record structure, allowing MARC records to be represented in XML. Libraries typically expose their records as MARC XML via a web service, often following the SRU or OAI-PMH standards.

MARC XML is an XML schema base on the common MARC 21 standards. MARC XML was developed by the Library of Congress and adopted by it and others as a means of facilitating the sharing of, and networked access to, bibliographic information.

7.1 MARCXML primary design goals included:

- Simplicity of the schema
- Flexibility and extensibility
- Lossless and reversible conversion from MARC
- Data presentation through XML style sheets
- MARC records updates and data conversions through XML transformations
- Existence of validation tools

8. Future

The future of the MARC formats is a matter of some debate among libraries. On the one hand, the storage formats are quite complex and are based on outdated technology. On the other, there is no alternative bibliographic format with an equivalent degree of granularity. The billions of MARC records in tens of thousands of individual libraries (including over 50,000,000 belonging to the OCLC consortium alone) create inertia. The Library of Congress has launched the Bibliographic Framework Initiative (BIBFRAME), that aims at providing a replacement for MARC that provides greater granularity and easier re-use of the data expressed in multiple catalogs.

Lesson 9: ISO 2709

1. Introduction

ISO 2709 is a standard for Record Structure of machine readable bibliographic record. It is an ISO standard for bibliographic descriptions, titled “Information and documentation—Format for information exchange”. It is maintained by the Technical Committee for Information and Documentation

2. History

ISO 2709 is a format for the exchange of bibliographic information. It was developed in the 1960s under the direction of Henriette Avram of the Library of Congress to encode the information printed on library cards. It was first created as ANSI Standard Z39.2, and called Information Interchange Format. The 1981 version of the standard was titled Documentation Format for bibliographic information interchange on magnetic tape. The latest edition of that standard is Z39.2 published in 1994 (ISSN: 1041-5653). The ISO standard supersedes Z39.2. The current standard is ISO 2709 released in December 2008.

3. Basic structure

ISO 2709 record has three sections:

- **Record label:** The first 24 characters of the record. This is the only portion of the record that is fixed in length. The record label includes the record length and the base address of the data contained in the record. It also has data elements that indicate how many characters are used for indicators and subfield identifiers.

Name of Data Element	No. of Characters
Record Length	5
Record Status	1
Type of Record	1
Bibliographic Level	1
Blanks	2
Indicator Count	1
Subfield Count	1
Base Address of Data	5
Blank Character	7

- **Directory:** The directory provides the entry positions to the fields in the record, along with the field tags. A directory entry has four parts and cannot exceed nine characters in length:

- Field tag (3 characters)
- Length of the field
- Starting character position of the field
- (Optional) Implementation-defined part

Entry	Element in the entry	No. Characters
Field 1	Tag of the field	3
	Length of the Field	4
	Starting Character Position of the field	5
Field 2	Tag of the field	3
	Length of the Field	4
	Starting Character Position of the field	5

Field n	Tag of the field	3
	Length of the Field	4
	Starting Character Position of the field	5
	Field Terminator (F/T)	1

- **Data fields (Variable fields):** A string containing all field and subfield data in the record

Indicator
Subfield code
Data Element
Subfield code
Data Element
.
.
.
Subfield code
Data Element
F/T
Indicator
Subfield code
Data Element
.
.
.
Subfield code
Data Element
F/T
.
.
.
F/T
R/T

- **Record separator:** A single character.

The tags are often displayed as labels on bibliographic fields. Each bibliographic field has an associated tag. The tags are stored in the directory not in the bibliographic field.

Fields

There are three kinds of fields in the ISO 2709 record:

- **Record identifier field:** Identifying the record and assigned by the organization that creates the record. The record identifier field has tag 001.
- **Reserved fields:** Reserved fields supply data which may be required for the processing of the record. Reserved fields always have a tag in the range 002–009 and 00A–ZZZ.
- **Bibliographic Fields:** These are in the range 010–999 and 0AA–ZZZ. The bibliographic fields contain data and a field separator. They can also have these optional sub-parts:
 - **Indicator** (0–9 characters, as coded in the Leader): Indicators generally provide further information about the contents of the field, the relationship between the field and other fields in the record, or about action required in certain data manipulation processes (including display labels).
 - **Identifier** (0–9 characters): This identifies data within the bibliographic field. Where used, identifiers are composed of a delimiter (1 char, IS₁ of ISO 646) and an identifying code (1–9 chars, as defined in the leader), plus a variable length string containing the data.

Example

MARC21 library cataloging record coded in ISO 2709 format. MARC21 is an instance of ISO 2709 that has the following characteristics:

- tags are in the range 002–999 only
- there is a two-character indicator on each field, and each character is a separately defined data element
- the identifier within data fields (called "subfield code" in MARC21) is a single ASCII character preceded by IS₁ of ISO 646.

References

1. "ISO 2709:2008 - Information and documentation -- Format for information exchange". Retrieved 21 January 2011.
2. "ISO 2709:1981 - Documentation -- Format for bibliographic information interchange on magnetic tape". Retrieved 21 January 2011.

Lesson 10: METADATA

Structure

0. Objective
1. Introduction
2. Definition of Metadata
3. Organization of description
4. Levels of Detail in Description
5. Structure of AACR2
 - 5.1 Part I of AACR2
 - 5.2 Part II of AACR2
6. Appendices:
7. Conclusion

0. Objective

The objective of this lesson is to explain the basic features of Anglo-American Cataloguing Rules 2. It explains the basic structure of the code, organization of description of the library materials and levels of description etc.

After studying this lesson you will be able to

- What is AACR2
- Governance of AACR
- Organization of description of documents
- Levels of description recommended by AACR2
- structure of AACR2

1. Introduction

The term **metadata** refers to "data about data". The term is ambiguous, as it is used for two fundamentally different concepts (types). **Structural metadata** is about the design and specification of data structures and is more properly called "data about the containers of data; **descriptive metadata**, on the other hand, is about individual instances of application data, the data content. Metadata (Meta content) are traditionally found in the card catalogs of libraries. As information has become increasingly digital, metadata are also used to describe digital data using metadata standards specific to a particular discipline. By describing the contents and context of data files, the quality of the original data/files is greatly increased. For example, a webpage may include metadata specifying what language it is written in, what tools were used to create it, and where to go for more on the subject, allowing browsers to automatically improve the experience of users.

2. Definition of Metadata

Metadata (Meta content) are defined as the data providing information about one or more aspects of the data, such as:

- Means of creation of the data
- Purpose of the data

- Time and date of creation
- Creator or author of the data
- Location on a computer network where the data were created
- Standards used

For example, a digital image may include metadata that describe how large the picture is, the color depth, the image resolution, when the image was created, and other data. A text document's metadata may contain information about how long the document is, who the author is, when the document was written, and a short summary of the document.

Metadata are data. As such, metadata can be stored and managed in a database, often called a Metadata registry or Metadata repository. However, without context and a point of reference, it might be impossible to identify metadata just by looking at them. For example: by itself, a database containing several numbers, all 13 digits long could be the results of calculations or a list of numbers to plug into an equation - without any other context, the numbers themselves can be perceived as the data. But if given the context that this database is a log of a book collection, those 13-digit numbers may now be identified as ISBNs - information that refers to the book, but is not itself the information within the book.

The term "metadata" was coined in 1968 by Philip Bagley, in his book "Extension of programming language concepts", where it is clear that he uses the term in the ISO 11179 "traditional" sense, which is "structural metadata" i.e. "data about the containers of data"; rather than the alternate sense "content about individual instances of data content" or meta content, the type of data usually found in library catalogues. Since then the fields of information management, information science, information technology, and librarianship have widely adopted the term. In these fields the word *metadata* is defined as "data about data". While this is the generally accepted definition, various disciplines have adopted their own more specific explanation and uses of the term.

Libraries

Metadata have been used in various forms as a means of cataloging archived information. The Dewey Decimal System employed by libraries for the classification of library materials is an early example of metadata usage. Library catalogues used 3x5 inch cards to display a book's title, author, subject matter, and a brief plot synopsis along with an abbreviated alpha-numeric identification system which indicated the physical location of the book within the library's shelves. Such data help classify, aggregate, identify, and locate a particular book.

Photographs

Metadata may be written into a digital photo file that will identify who owns it, copyright & contact information, what camera created the file, along with exposure information and descriptive information such as keywords about the photo, making the file searchable on the computer and/or the Internet. Some metadata are written by the camera and some is input by the photographer and/or software after downloading to a computer. However, not all digital cameras enable you to edit metadata. This functionality has been available on most Nikon DSLRs since the Nikon D3 and on most new Canon cameras since the Canon EOS 7D. Photographic Metadata Standards are governed by organizations that develop the following standards. They include, but are not limited to:

- IPTC Information Interchange Model IIM (International Press Telecommunications Council),
- IPTC Core Schema for XMP

- XMP – Extensible Metadata Platform (an ISO standard)
- Exif – Exchangeable image file format, Maintained by CIPA (Camera & Imaging Products Association) and published by JEITA (Japan Electronics and Information Technology Industries Association)
- Dublin Core (Dublin Core Metadata Initiative – DCMI)
- PLUS (Picture Licensing Universal System).

Video

Metadata are particularly useful in video, where information about its contents (such as transcripts of conversations and text descriptions of its scenes) are not directly understandable by a computer, but where efficient search is desirable.

Web pages

Web pages often include metadata in the form of Meta tags. Description and keywords Meta tags are commonly used to describe the Web page's content. Most search engines use these data when adding pages to their search index.

3. Creation of metadata

Metadata can be created either by automated information processing or by manual work. Elementary metadata captured by computers can include information about when an object was created, who created it, when it was last updated, file size and file extension.

For the purposes of this article, an "object" refers to any of the following:

- A physical item such as a book, CD, DVD, map, chair, table, flower pot, etc.
- An electronic file such as a digital image, digital photo, document, program file, database table, etc.

4. Metadata types

The metadata application is many fold covering a large variety of fields of application there are nothing but specialized and well accepted models to specify types of metadata. Bretheron & Singley (1994) distinguish between two distinct classes: structural/control metadata and guide metadata. **Structural metadata** are used to describe the structure of computer systems such as tables, columns and indexes. **Guide metadata** are used to help humans find specific items and are usually expressed as a set of keywords in a natural language. According to Ralph Kimball metadata can be divided into 2 similar categories: technical metadata and business metadata. **Technical metadata** correspond to internal metadata, *business metadata* - to external metadata. Kimball adds a third category named **Process metadata**. On the other hand, NISO distinguishes among three types of metadata: descriptive, structural and administrative. **Descriptive metadata** are the information used to search and locate an object such as title, author, subjects, keywords, publisher; **Structural metadata** give a description of how the components of the object are organized; and **Administrative metadata** refer to the technical information including file type. Two sub-types of administrative metadata are rights management metadata and preservation metadata.

5. Metadata structures

Metadata (Meta content), or more correctly, the vocabularies used to assemble metadata (Meta content) statements, are typically structured according to a standardized

concept using a well-defined metadata scheme, including: metadata standards and metadata models. Tools such as controlled vocabularies, taxonomies, thesauri, data dictionaries and metadata registries can be used to apply further standardization to the metadata. Structural metadata commonality is also of paramount importance in data model development and in database design.

5.1 Metadata syntax

Metadata (Meta content) syntax refers to the rules created to structure the fields or elements of metadata (Meta content). A single metadata scheme may be expressed in a number of different markups or programming languages, each of which requires a different syntax. For example, Dublin Core may be expressed in plain text, HTML, XML and RDF.

A common example of (guide) Meta content is the bibliographic classification, the subject, the Dewey Decimal class number. There is always an implied statement in any "classification" of some object. To classify an object as, for example, Dewey class number 514 (Topology) (i.e. books having the number 514 on their spine) the implied statement is: "<book><subject heading><514>". This is a subject-predicate-object triple, or more importantly, a class-attribute-value triple. The first two elements of the triple (class, attribute) are pieces of some structural metadata having a defined semantic. The third element is a value, preferably from some controlled vocabulary, some reference (master) data. The combination of the metadata and master data elements results in a statement which is a Meta content statement i.e. "Meta content = metadata + master data". All these elements can be thought of as "vocabulary". Both metadata and master data are vocabularies which can be assembled into Meta content statements. There are many sources of these vocabularies, both Meta and master data: UML, EDIFACT, XSD, Dewey/UDC/LoC, SKOS, ISO-25964, Pantone, Linnaean Binomial Nomenclature etc. Using controlled vocabularies for the components of Meta content statements, whether for indexing or finding, is endorsed by ISO-25964: "If both the indexer and the searcher are guided to choose the same term for the same concept, then relevant documents will be retrieved."

5.2 Hierarchical, linear and planar schemata

Metadata schema can be hierarchical in nature where relationships exist between metadata elements and elements are nested so that parent-child relationships exist between the elements. An example of a hierarchical metadata schema is the IEEE LOM schema where metadata elements may belong to a parent metadata element. Metadata schema can also be one-dimensional, or linear, where each element is completely discrete from other elements and classified according to one dimension only. An example of a linear metadata schema is Dublin Core schema which is one dimensional. Metadata schema is often two dimensional, or planar, where each element is completely discrete from other elements but classified according to two orthogonal dimensions.

5.3 Metadata hyper mapping

In all cases where the metadata schemata exceed the planar depiction, some type of hyper mapping is required to enable display and view of metadata according to chosen aspect and to serve special views. Hyper mapping frequently applies to layering of geographical and geological information overlays.

5.4 Granularity

The degree to which the data or metadata are structured is referred to as their granularity. Metadata with a high granularity allow for deeper structured information and enable greater levels of technical manipulation however, a lower level of granularity means that

metadata can be created for considerably lower costs but will not provide as detailed information. The major impact of granularity is not only on creation and capture, but moreover on maintenance. As soon as the metadata structures get outdated, the access to the referred data will get outdated. Hence granularity shall take into account the effort to create as well as the effort to maintain.

6. Metadata standards

International standards apply to metadata. Much work is being accomplished in the national and international standards communities, especially ANSI (American National Standards Institute) and ISO (International Organization for Standardization) to reach consensus on standardizing metadata and registries.

The core standard is ISO/IEC 11179-1:2004 and subsequent standards (see ISO/IEC 11179). All yet published registrations according to this standard cover just the definition of metadata and do not serve the structuring of metadata storage or retrieval neither any administrative standardization. It is important to note that this standard refers to metadata as the data about containers of the data and not to metadata (Meta content) as the data about the data contents. It should also be noted that this standard describes itself originally as a "data element" registry, describing disembodied data elements, and explicitly disavows the capability of containing complex structures. Thus the original term "data element" is more applicable than the later applied buzzword "metadata".

The Dublin Core metadata terms are a set of vocabulary terms which can be used to describe resources for the purposes of discovery. The original set of 15 classic metadata terms, known as the Dublin Core Metadata Element Set are endorsed in the following standards documents:

- IETF RFC 5013
- ISO Standard 15836-2009
- NISO Standard Z39.85

6.1 Library and information science

Libraries employ metadata in library catalogues, most commonly as part of an Integrated Library Management System. Metadata are obtained by cataloguing resources such as books, periodicals, DVDs, web pages or digital images. These data are stored in the integrated library management system, ILMS, using the MARC metadata standard. The purpose is to direct patrons to the physical or electronic location of items or areas they seek as well as to provide a description of the item/s in question.

More recent and specialized instances of library metadata include the establishment of digital libraries including e-print repositories and digital image libraries. While often based on library principles, the focus on non-librarian use, especially in providing metadata, means they do not follow traditional or common cataloging approaches. Given the custom nature of included materials, metadata fields are often specially created e.g. taxonomic classification fields, location fields, keywords or copyright statement. Standard file information such as file size and format are usually automatically included.

Standardization for library operation has been a key topic in international standardization (ISO) for decades. Standards for metadata in digital libraries include Dublin Core, METS, MODS, DDI, ISO standard Digital Object Identifier (DOI), ISO standard Uniform

Resource Name (URN), PREMIS schema, Ecological Metadata Language, and OAI-PMH. Leading libraries in the world give hints on their metadata standards strategies.

Kimball et al. refers to three main categories of metadata: Technical metadata, business metadata and process metadata. Technical metadata are primarily definitional, while business metadata and process metadata are primarily descriptive. Keep in mind that the categories sometimes overlap.

- **Technical metadata** define the objects and processes in a DW/BI system, as seen from a technical point of view. The technical metadata include the system metadata which define the data structures such as: tables, fields, data types, indexes and partitions in the relational engine, and databases, dimensions, measures, and data mining models. Technical metadata define the data model and the way it is displayed for the users, with the reports, schedules, distribution lists and user security rights.
- **Business metadata** are a content from the data warehouse described in more user-friendly terms. The business metadata tell you what data you have, where they come from, what they mean and what is their relationship is to other data in the data warehouse. Business metadata may also serve as a documentation for the DW/BI system. Users who browse the data warehouse are primarily viewing the business metadata.
- **Process metadata** are used to describe the results of various operations in the data warehouse. Within the ETL process, all key data from tasks are logged on execution. This includes start time, end time, CPU seconds used, disk reads, disk writes and rows processed. When troubleshooting the ETL or query process, this sort of data becomes valuable. Process metadata are the fact measurement when building and using a DW/BI system. Some organizations make a living out of collecting and selling this sort of data to companies - in that case the process metadata become the business metadata for the fact and dimension tables. Collecting process metadata is in the interest of business people who can use the data to identify the users of their products, which products they are using and what level of service they are receiving.

6.2 Metadata on the Internet

The HTML format used to define web pages allows for the inclusion of a variety of types of metadata, from basic descriptive text, dates and keywords to further advance metadata schemes such as the Dublin Core, e-GMS, and AGLS standards. Pages can also be geo-tagged with coordinates. Metadata may be included in the page's header or in a separate file. Micro formats allow metadata to be added to on-page data in a way that users do not see, but computers can readily access.

7. Metadata management

Metadata management is the end-to-end process and governance framework for creating, controlling, enhancing, attributing, defining and managing a metadata schema, model or other structured aggregation, either independently or within a repository and the associated supporting processes (often to enable the management of content). The World Wide Web Consortium (W3C) has identified Governance as a key challenge in the advancement of third generation Web Technologies (Web 3.0, Semantic Web), and a number of research prototypes, such as S3DB, explore the use of semantic modeling to identify practical solutions.

7.1 Database management

Each relational database system has its own mechanisms for storing metadata. Examples of relational-database metadata include:

- Tables of all tables in a database, their names, sizes and number of rows in each table.
- Tables of columns in each database, what tables they are used in, and the type of data stored in each column.

In database terminology, this set of metadata is referred to as the catalog. The SQL standard specifies a uniform means to access the catalog, called the information schema, but not all databases implement it, even if they implement other aspects of the SQL standard. For an example of database-specific metadata access methods, see Oracle metadata. Programmatic access to metadata is possible using APIs such as JDBC, or Schema Crawler.

Lesson 11: Dublin Core

Structure

0. Objective
1. Introduction
2. Background of Dublin Core
3. Levels of the Dublin Core Standard
 - 3.1 Simple Dublin Core
 - 3.2 Qualified Dublin Core
4. Syntaxes of Dublin Core Meta data
5. Conclusion
6. References

0. Objective

The objective of this lesson is to explain the basic features of Dublin core Meta data standards for cataloguing web resources. It explains the basic structure of the Dublin core code, organization of description of the web resources of information and levels of description etc.

After studying this lesson you will be able to

- What is Dublin Core
- HISTORY OF Dublin Core
- Levels of the Dublin Core Standard
- Syntax of Dublin Core

1. Introduction

The Dublin Core metadata terms are a set of vocabulary terms which can be used to describe resources for the purposes of identification/discovery. The terms can be used to describe a full range of web resources (video, images, web pages, etc.), physical resources such as books and objects like artworks. The original set of 15 classic metadata terms, known as the Dublin Core Metadata Element Set are endorsed in the following standards documents:

- IETF RFC 5013
- ISO Standard 15836-2009
- NISO Standard Z39.85

Dublin Core Metadata can be used for multiple purposes, from simple resource description, to combining metadata vocabularies of different metadata standards, to providing interoperability for metadata vocabularies in the Linked data cloud and Semantic web implementations.

2. Background

"Dublin" refers to Dublin, Ohio, USA where the work originated during the 1995 invitational OCLC/NCSA Metadata Workshop, hosted in by Online Computer Library Center (OCLC), a library consortium based there, and the National Center for Supercomputing Applications (NCSA).

"Core" refers to the metadata terms as "broad and generic being usable for describing a wide range of resources".

The semantics of Dublin Core were established and are maintained by an international, cross-disciplinary group of professionals from librarianship, computer science, text encoding, museums, and other related fields of scholarship and practice.

The Dublin Core Metadata Initiative (DCMI) provides an open forum for the development of interoperable online metadata standards for a broad range of purposes and of business models. DCMI's activities include consensus-driven working groups, global conferences and workshops, standardization, and educational efforts to promote widespread acceptance of metadata standards and practices. In 2008, DCMI separated from OCLC and incorporated as an independent entity.

3. Levels of the standard

The Dublin Core standard includes two levels viz. Simple and Qualified.

3.1 Simple Dublin Core

The Simple **Dublin Core Metadata Element Set (DCMES)** consists of 15 metadata elements:

1. Title
2. Creator
3. Subject
4. Description
5. Publisher
6. Contributor
7. Date
8. Type
9. Format
10. Identifier
11. Source
12. Language
13. Relation
14. Coverage
15. Rights

Each Dublin Core element is optional and may be repeated. The DCMI has established standard ways to refine elements and encourage the use of encoding and vocabulary schemes. There is no prescribed order in Dublin Core for presenting or using the elements. The Dublin Core became ISO 15836 standard in 2006 and is used as a base-level data element set for the description of learning resources in the ISO/IEC 19788-2 Metadata for learning resources (MLR) -- Part 2: Dublin Core elements, prepared by the ISO/IEC JTC1 SC36.

EXAMPLE OF CODE

```
<meta name="DC.Publisher" content="publisher-name" >
```

AN EXAMPLE OF USE [AND MENTION] OF D.C. (BY WEBCITE)

At the web page which serves as the "archive" form for WebCite, it says, in part: "Metadata (optional). These are Dublin Core elements. [...]".

3.2 Qualified Dublin Core

Qualified Dublin Core includes three additional elements (Audience, Provenance and Rights Holder), as well as a group of element refinements (also called qualifiers) that refine the semantics of the elements in ways that may be useful in resource discovery. Subsequent to the specification of the original 15 elements, an ongoing process to develop exemplary terms extending or refining the Dublin Core Metadata Element Set (DCMES) was begun. The additional terms were identified, generally in working groups of the Dublin Core Metadata Initiative, and judged by the DCMI Usage Board to be in conformance with principles of good practice for the qualification of Dublin Core metadata elements.

Elements refinements make the meaning of an element narrower or more specific. A refined element shares the meaning of the unqualified element, but with a more restricted scope. The guiding principle for the qualification of Dublin Core elements, colloquially known as the *Dumb-Down Principle*, states that an application that does not understand a specific element refinement term should be able to ignore the qualifier and treat the metadata value as if it were an unqualified (broader) element. While this may result in some loss of specificity, the remaining element value (without the qualifier) should continue to be generally correct and useful for discovery.

In addition to element refinements, Qualified Dublin Core includes a set of recommended encoding schemes, designed to aid in the interpretation of an element value. These schemes include controlled vocabularies and formal notations or parsing rules. A value expressed using an encoding scheme may thus be a token selected from a controlled vocabulary (for example, a term from a classification system or set of subject headings) or a string formatted in accordance with a formal notation, for example, "2000-12-31" as the ISO standard expression of a date. If an encoding scheme is not understood by an application, the value may still be useful to human reader.

Audience, Provenance and RightsHolder are elements, but not part of the Simple Dublin Core 15 elements. Use Audience, Provenance and RightsHolder only when using Qualified Dublin Core. DCMI also maintains a small, general vocabulary recommended for use within the element Type. This vocabulary currently consists of 12 terms.

4. Syntaxes of Dublin Core Metadata

Syntax choices for Dublin Core metadata depend on a number of variables, and "one size fits all" prescriptions rarely apply. When considering an appropriate syntax, it is important to note that Dublin Core concepts and semantics are designed to be syntax independent, are equally

applicable in a variety of contexts, as long as the metadata is in a form suitable for interpretation both by machines and by human beings.

The Dublin Core Abstract Model provides a reference model against which particular Dublin Core encoding guidelines can be compared, independent of any particular encoding syntax. Such a reference model allows implementers to gain a better understanding of the kinds of descriptions they are trying to encode and facilitates the development of better mappings and translations between different syntaxes.

5. Conclusion

The Dublin Core metadata terms are a set of vocabulary terms used to describe a full range of web resources (video, images, web pages, etc.), physical resources such as books etc. It is an open source standard which can be improved by any one. The Dublin Core Metadata initiative (DCMI) provides an open forum for the development of interoperable online metadata standards for a broad range of purposes and of business models.

Lesson 12: Evaluation of Information retrieval

Information retrieval is the activity of obtaining information resources relevant to an information need from a collection of information resources. Searches can be based on metadata or on full-text (or other content-based) indexing.

Automated information retrieval systems are used to reduce what has been called "information overload". Many universities and public libraries use IR systems to provide access to books, journals and other documents. Web search engines are the most visible IR applications.

History

The idea of using computers to search for relevant pieces of information was popularized in the article *As We May Think* by Vannevar Bush in 1945.^[1] The first automated information retrieval systems were introduced in the 1950s and 1960s. By 1970 several different techniques had been shown to perform well on small text corpora such as the Cranfield collection (several thousand documents).^[1] Large-scale retrieval systems, such as the Lockheed Dialog system, came into use early in the 1970s.

In 1992, the US Department of Defense along with the National Institute of Standards and Technology (NIST), cosponsored the Text Retrieval Conference (TREC) as part of the TIPSTER text program. The aim of this was to look into the information retrieval community by supplying the infrastructure that was needed for evaluation of text retrieval methodologies on a very large text collection. This catalyzed research on methods that scale to huge corpora. The introduction of web search engines has boosted the need for very large scale retrieval systems even further.

Timeline

Before the 1900s

1801: Joseph Marie Jacquard invents the Jacquard loom, the first machine to use punched cards to control a sequence of operations.

1880s: Herman Hollerith invents an electro-mechanical data tabulator using punch cards as a machine readable medium.

1890 Hollerith cards, keypunches and tabulators used to process the 1890 US Census data.

- **1920s-1930s**

Emanuel Goldberg submits patents for his "Statistical Machine" a document search engine that used photoelectric cells and pattern recognition to search the metadata on rolls of microfilmed documents.

- **1940s–1950s**

late 1940s: The US military confronted problems of indexing and retrieval of wartime scientific research documents captured from Germans.

1945: Vannevar Bush's *As We May Think* appeared in *Atlantic Monthly*.

1947: Hans Peter Luhn (research engineer at IBM since 1941) began work on a mechanized punch card-based system for searching chemical compounds.

1950s: Growing concern in the US for a "science gap" with the USSR motivated, encouraged funding and provided a backdrop for mechanized literature searching systems (Allen Kent *et al.*) and the invention of citation indexing (Eugene Garfield).

1950: The term "information retrieval" appears to have been coined by Calvin Mooers.^[2]

1951: Philip Bagley conducted the earliest experiment in computerized document retrieval in a master thesis at MIT.^[3]

1955: Allen Kent joined Case Western Reserve University, and eventually became associate director of the Center for Documentation and Communications Research. That same year, Kent and colleagues published a paper in *American Documentation* describing the precision and recall measures as well as detailing a proposed "framework" for evaluating an IR system which included statistical sampling methods for determining the number of relevant documents not retrieved.

1958: International Conference on Scientific Information Washington DC included consideration of IR systems as a solution to problems identified. See: *Proceedings of the International Conference on Scientific Information, 1958* (National Academy of Sciences, Washington, DC, 1959)

1959: Hans Peter Luhn published "Auto-encoding of documents for information retrieval."

- **1960s:**

early 1960s: Gerard Salton began work on IR at Harvard, later moved to Cornell.

1960: Melvin Earl Maron and John Lary Kuhns^[4] published "On relevance, probabilistic indexing, and information retrieval" in the *Journal of the ACM* 7(3):216–244, July 1960.

1962:

- Cyril W. Cleverdon published early findings of the Cranfield studies, developing a model for IR system evaluation. See: Cyril W. Cleverdon, "Report on the Testing and Analysis of an Investigation into the Comparative Efficiency of Indexing Systems". Cranfield Collection of Aeronautics, Cranfield, England, 1962.
- Kent published *Information Analysis and Retrieval*.

1963:

- Weinberg report "Science, Government and Information" gave a full articulation of the idea of a "crisis of scientific information." The report was named after Dr. Alvin Weinberg.
- Joseph Becker and Robert M. Hayes published text on information retrieval. Becker, Joseph; Hayes, Robert Mayo. *Information storage and retrieval: tools, elements, theories*. New York, Wiley (1963).

1964:

- Karen Spärck Jones finished her thesis at Cambridge, *Synonymy and Semantic Classification*, and continued work on computational linguistics as it applies to IR.
- The National Bureau of Standards sponsored a symposium titled "Statistical Association Methods for Mechanized Documentation." Several highly significant papers, including G. Salton's first published reference (we believe) to the SMART system.

mid-1960s:

- National Library of Medicine developed MEDLARS Medical Literature Analysis and Retrieval System, the first major machine-readable database and batch-retrieval system.
- Project Intrex at MIT.

1965: J. C. R. Licklider published *Libraries of the Future*.

1966: Don Swanson was involved in studies at University of Chicago on Requirements for Future Catalogs.

late 1960s: F. Wilfrid Lancaster completed evaluation studies of the MEDLARS system and published the first edition of his text on information retrieval.

1968:

- Gerard Salton published *Automatic Information Organization and Retrieval*.
- John W. Sammon, Jr.'s RADC Tech report "Some Mathematics of Information Storage and Retrieval..." outlined the vector model.

1969: Sammon's "A nonlinear mapping for data structure analysis" (IEEE Transactions on Computers) was the first proposal for visualization interface to an IR system.

• 1970s

early 1970s:

- First online systems—NLM's AIM-TWX, MEDLINE; Lockheed's Dialog; SDC's ORBIT.
- Theodor Nelson promoting concept of hypertext, published *Computer Lib/Dream Machines*.

1971: Nicholas Jardine and Cornelis J. van Rijsbergen published "The use of hierarchic clustering in information retrieval", which articulated the "cluster hypothesis."^[5]

1975: Three highly influential publications by Salton fully articulated his vector processing framework and term discrimination model:

- *A Theory of Indexing* (Society for Industrial and Applied Mathematics)
- *A Theory of Term Importance in Automatic Text Analysis* (JASIS v. 26)
- *A Vector Space Model for Automatic Indexing* (CACM 18:11)

1978: The First ACM SIGIR conference.

1979: C. J. van Rijsbergen published *Information Retrieval* (Butterworths). Heavy emphasis on probabilistic models.

• 1980s

1980: First international ACM SIGIR conference, joint with British Computer Society IR group in Cambridge.

1982: Nicholas J. Belkin, Robert N. Oddy, and Helen M. Brooks proposed the ASK (Anomalous State of Knowledge) viewpoint for information retrieval. This was an important concept, though their automated analysis tool proved ultimately disappointing.

1983: Salton (and Michael J. McGill) published *Introduction to Modern Information Retrieval* (McGraw-Hill), with heavy emphasis on vector models.

1985: David Blair and Bill Maron publish: *An Evaluation of Retrieval Effectiveness for a Full-Text Document-Retrieval System*

mid-1980s: Efforts to develop end-user versions of commercial IR systems.

1985–1993: Key papers on and experimental systems for visualization interfaces.

Work by Donald B. Crouch, Robert R. Korfhage, Matthew Chalmers, Anselm Spoerri and others.

1989: First World Wide Web proposals by Tim Berners-Lee at CERN.

• 1990s

1992: First TREC conference.

1997: Publication of Korfhage's *Information Storage and Retrieval*^[6] with emphasis on visualization and multi-reference point systems.

late 1990s: Web search engines implementation of many features formerly found only in experimental IR systems. Search engines become the most common and maybe best instantiation of IR models, research, and implementation.

Overview

An information retrieval process begins when a user enters a query into the system. Queries are formal statements of information needs, for example search strings in web search engines. In information retrieval a query does not uniquely identify a single object in the collection. Instead, several objects may match the query, perhaps with different degrees of relevancy.

An object is an entity that is represented by information in a database. User queries are matched against the database information. Depending on the application the data objects may be, for example, text documents, images,^[7] audio,^[8] mind maps^[9] or videos. Often the documents themselves are not kept or stored directly in the IR system, but are instead represented in the system by document surrogates or metadata.

Most IR systems compute a numeric score on how well each object in the database matches the query, and rank the objects according to this value. The top ranking objects are then shown to the user. The process may then be iterated if the user wishes to refine the query.^[10]

Performance and correctness measures

Many different measures for evaluating the performance of information retrieval systems have been proposed. The measures require a collection of documents and a query. All common measures described here assume a ground truth notion of relevancy: every document is known to be either relevant or non-relevant to a particular query. In practice queries may be ill-posed and there may be different shades of relevancy.

Precision

Precision is the fraction of the documents retrieved that are relevant to the user's information need.

$$\text{precision} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{retrieved documents}\}|}$$

In binary classification, precision is analogous to positive predictive value. Precision takes all retrieved documents into account. It can also be evaluated at a given cut-off rank, considering only the topmost results returned by the system. This measure is called *precision at n* or *P@n*.

Note that the meaning and usage of "precision" in the field of Information Retrieval differs from the definition of accuracy and precision within other branches of science and technology.

Recall[edit]

Recall is the fraction of the documents that are relevant to the query that are successfully retrieved.

$$\text{recall} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{relevant documents}\}|}$$

In binary classification, recall is often called sensitivity. So it can be looked at as *the probability that a relevant document is retrieved by the query*.

It is trivial to achieve recall of 100% by returning all documents in response to any query. Therefore recall alone is not enough but one needs to measure the number of non-relevant documents also, for example by computing the precision.

Fall-out

The proportion of non-relevant documents that are retrieved, out of all non-relevant documents available:

$$\text{fall-out} = \frac{|\{\text{non-relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{non-relevant documents}\}|}$$

In binary classification, fall-out is closely related to specificity and is equal to $(1 - \text{specificity})$. It can be looked at as *the probability that a non-relevant document is retrieved by the query*.

It is trivial to achieve fall-out of 0% by returning zero documents in response to any query.

F-measure[edit]

The weighted harmonic mean of precision and recall, the traditional F-measure or balanced F-score is:

$$F = \frac{2 \cdot \text{precision} \cdot \text{recall}}{(\text{precision} + \text{recall})}$$

This is also known as the F_1 measure, because recall and precision are evenly weighted.

The general formula for non-negative real β is:

$$F_\beta = \frac{(1 + \beta^2) \cdot (\text{precision} \cdot \text{recall})}{(\beta^2 \cdot \text{precision} + \text{recall})}$$

Two other commonly used F measures are the F_2 measure, which weights recall twice as much as precision, and the $F_{0.5}$ measure, which weights precision twice as much as recall.

The F-measure was derived by van Rijsbergen (1979) so that F_β "measures the effectiveness of retrieval with respect to a user who attaches β times as much importance to recall as

precision". It is based on van Rijsbergen's effectiveness measure $E = 1 - \frac{1}{\frac{\alpha}{P} + \frac{1-\alpha}{R}}$. Their relationship is $P^\beta = 1 - E'$ where $\alpha = \frac{1}{1 + \beta^2}$.

Average precision[edit]

Precision and recall are single-value metrics based on the whole list of documents returned by the system. For systems that return a ranked sequence of documents, it is desirable to also consider the order in which the returned documents are presented. By computing a precision and recall at every position in the ranked sequence of documents, one can plot a precision-recall curve, plotting precision $P(r)$ as a function of recall r . Average precision computes the average value of $P(r)$ over the interval from $r = 0$ to $r = 1$.^[11]

$$\text{AveP} = \int_0^1 p(r) dr$$

That is the area under the precision-recall curve. This integral is in practice replaced with a finite sum over every position in the ranked sequence of documents:

$$\text{AveP} = \sum_{k=1}^n P(k) \Delta r(k)$$

where k is the rank in the sequence of retrieved documents, n is the number of retrieved documents, $P(k)$ is the precision at cut-off k in the list, and $\Delta r(k)$ is the change in recall from items $k - 1$ to k .^[11]

This finite sum is equivalent to:

$$\text{AveP} = \frac{\sum_{k=1}^n (P(k) \times \text{rel}(k))}{\text{number of relevant documents}}$$

where $\text{rel}(k)$ is an indicator function equaling 1 if the item at rank k is a relevant document, zero otherwise.^[12] Note that the average is over all relevant documents and the relevant documents not retrieved get a precision score of zero.

Some authors choose to interpolate the $P(r)$ function to reduce the impact of "wiggles" in the curve.^{[13][14]} For example, the PASCAL Visual Object Classes challenge (a benchmark for computer vision object detection) computes average precision by averaging the precision over a set of evenly spaced recall levels $\{0, 0.1, 0.2, \dots, 1.0\}$.^{[13][14]}

$$\text{AveP} = \frac{1}{11} \sum_{r \in \{0, 0.1, \dots, 1.0\}} p_{\text{interp}}(r)$$

where $p_{\text{interp}}(r)$ is an interpolated precision that takes the maximum precision over all recalls greater than r :

$$p_{\text{interp}}(r) = \max_{\tilde{r}: \tilde{r} \geq r} p(\tilde{r}).$$

An alternative is to derive an analytical $p(r)$ function by assuming a particular parametric distribution for the underlying decision values. For example, a *binormal precision-recall curve* can be obtained by assuming decision values in both classes to follow a Gaussian distribution.^[15]

R-Precision[edit]

Precision at **R**-th position in the ranking of results for a query that has **R** relevant documents. This measure is highly correlated to Average Precision. Also, Precision is equal to Recall at the **R**-th position.

Mean average precision[edit]

Mean average precision for a set of queries is the mean of the average precision scores for each query.

$$\text{MAP} = \frac{\sum_{q=1}^Q \text{AveP}(q)}{Q}$$

where Q is the number of queries.

Discounted cumulative gain

DCG uses a graded relevance scale of documents from the result set to evaluate the usefulness, or gain, of a document based on its position in the result list. The premise of DCG is that highly relevant documents appearing lower in a search result list should be penalized as the graded relevance value is reduced logarithmically proportional to the position of the result.

The DCG accumulated at a particular rank position P is defined as:

$$\text{DCG}_P = rel_1 + \sum_{i=2}^P \frac{rel_i}{\log_2 i}.$$

Since result set may vary in size among different queries or systems, to compare performances the normalised version of DCG uses an ideal DCG. To this end, it sorts

documents of a result list by relevance, producing an ideal DCG at position p ($IDCG_p$), which normalizes the score:

$$nDCG_p = \frac{DCG_p}{IDCG_p}.$$

The nDCG values for all queries can be averaged to obtain a measure of the average performance of a ranking algorithm. Note that in a perfect ranking algorithm, the DCG_p will be the same as the $IDCG_p$ producing an nDCG of 1.0. All nDCG calculations are then relative values on the interval 0.0 to 1.0 and so are cross-query comparable.

Other Measures

Mean reciprocal rank

- Spearman's rank correlation coefficient

Model types



For effectively retrieving relevant documents by IR strategies, the documents are typically transformed into a suitable representation. Each retrieval strategy incorporates a specific model for its document representation purposes. The picture on the right illustrates the relationship of some common models. In the picture, the models are categorized according to two dimensions: the mathematical basis and the properties of the model.

First dimension: mathematical basis

Set-theoretic models represent documents as sets of words or phrases. Similarities are usually derived from set-theoretic operations on those sets. Common models are:

- Standard Boolean model
- Extended Boolean model
- Fuzzy retrieval
- *Algebraic models* represent documents and queries usually as vectors, matrices, or tuples. The similarity of the query vector and document vector is represented as a scalar value.
- Vector space model
- Generalized vector space model
- (Enhanced) Topic-based Vector Space Model
- Extended Boolean model
- Latent semantic indexing aka latent semantic analysis
- *Probabilistic models* treat the process of document retrieval as a probabilistic inference. Similarities are computed as probabilities that a document is relevant for a given query. Probabilistic theorems like the Bayes' theorem are often used in these models.
- Binary Independence Model
- Probabilistic relevance model on which is based the okapi (BM25) relevance function
- Uncertain inference
- Language models

- Divergence-from-randomness model
- Latent Dirichlet allocation
- *Feature-based retrieval models* view documents as vectors of values of *feature functions* (or just *features*) and seek the best way to combine these features into a single relevance score, typically by learning to rank methods. Feature functions are arbitrary functions of document and query, and as such can easily incorporate almost any other retrieval model as just a yet another feature.

Second dimension: properties of the model

Models without term-interdependencies treat different terms/words as independent. This fact is usually represented in vector space models by the orthogonality assumption of term vectors or in probabilistic models by an independency assumption for term variables.

- *Models with immanent term interdependencies* allow a representation of interdependencies between terms. However the degree of the interdependency between two terms is defined by the model itself. It is usually directly or indirectly derived (e.g. by dimensional reduction) from the co-occurrence of those terms in the whole set of documents.
- *Models with transcendent term interdependencies* allow a representation of interdependencies between terms, but they do not allege how the interdependency between two terms is defined. They relay an external source for the degree of interdependency between two terms. (For example a human or sophisticated algorithms.)

Lesson 13: Information Retrieval Models

1 Introduction

A quick overview of the major textual retrieval methods were discussed in this lesson. First a general model of the information retrieval process was discussed and then briefly describe the major retrieval methods and characterize them in terms of their strengths and shortcomings.

2 General Model of Information Retrieval

The goal of **information retrieval** (IR) is to provide users with those documents that will satisfy their information need. We use the word "document" as a general term that could also include non-textual information, such as multimedia objects. Figure 4.1 provides a general overview of the information retrieval process, which has been adapted from Lancaster and Warner (1993). Users have to formulate their information need in a form that can be understood by the retrieval mechanism. There are several steps involved in this translation process that we will briefly discuss below. Likewise, the contents of large document collections need to be described in a form that allows the retrieval mechanism to identify the potentially relevant documents quickly. In both cases, information may be lost in the transformation process leading to a computer-usable representation. Hence, the matching process is inherently imperfect.

Information seeking is a form of problem solving [Marcus 1994, Marchionini 1992]. It proceeds according to the interaction among eight subprocesses: problem recognition and acceptance, problem definition, search system selection, query formulation, query execution, examination of results (including relevance feedback), information extraction, and reflection/iteration/termination. To be able to perform effective searches, users have to develop the following expertise: knowledge about various sources of information, skills in defining search problems and applying search strategies, and competence in using electronic search tools.

Marchionini (1992) contends that some sort of spreadsheet is needed that supports users in the problem definition as well as other information seeking tasks. The InfoCrystal is such a spreadsheet because it assists users in the formulation of their information needs and the exploration of the retrieved documents, using the a visual interface that supports a "what-if" functionality. He further predicts that advances in computing power and speed, together with improved information retrieval procedures, will continue to blur the distinctions between problem articulation and examination of results. The InfoCrystal is both a visual query language and a tool for visualizing retrieval results.

The information need can be understood as forming a pyramid, where only its peak is made visible by users in the form of a conceptual query (see Figure 2.1). The conceptual query captures the key concepts and the relationships among them. It is the result of a conceptual analysis that operates on the information need, which may be well or vaguely defined in the user's mind. This analysis can be challenging, because users are faced with the general "vocabulary problem" as they are trying to translate their information need into a conceptual query. This problem refers to the fact that a single word can have more than one meaning, and,

conversely, the same concept can be described by surprisingly many different words. Furnas, Landauer, Gomez and Dumais (1983) have shown that two people use the same main word to describe an object only 10 to 20% of the time. Further, the concepts used to represent the documents can be different from the concepts used by the user. The conceptual query can take the form of a natural language statement, a list of concepts that can have degrees of importance assigned to them, or it can be statement that coordinates the concepts using Boolean operators. Finally, the conceptual query has to be translated into a query surrogate that can be understood by the retrieval system.

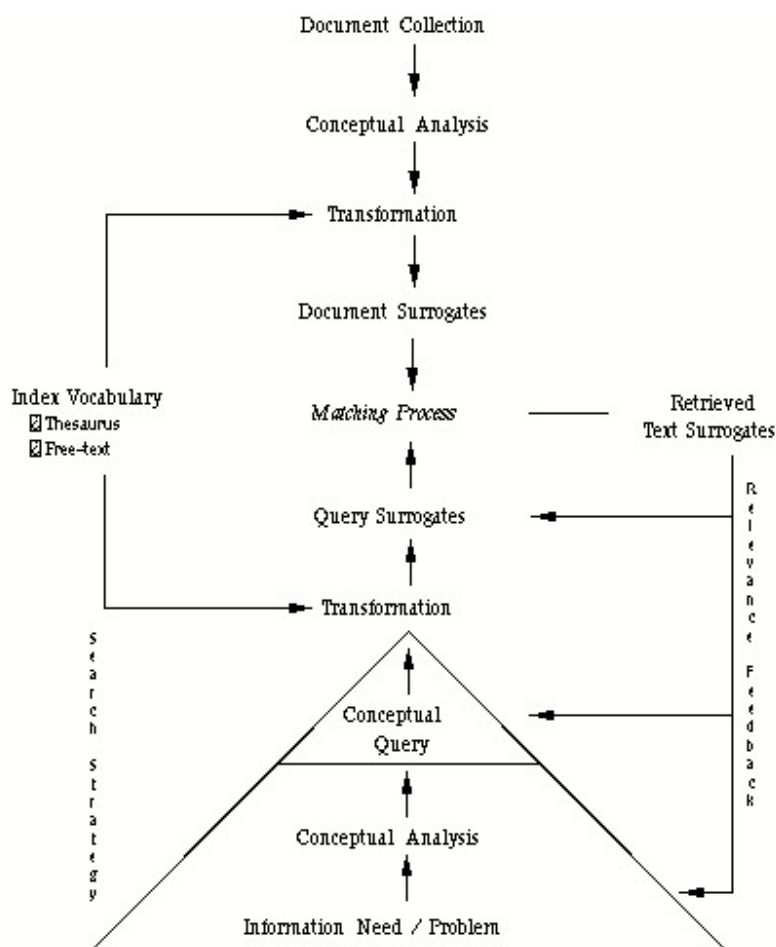


Figure 1: represents a general model of the information retrieval process, where both the user's information need and the document collection have to be translated into the form of surrogates to enable the matching process to be performed. This figure has been adapted from Lancaster and Warner (1993).

Similarly, the meanings of documents need to be represented in the form of text surrogates that can be processed by computer. A typical surrogate can consist of a set of index terms or descriptors. The text surrogate can consist of multiple fields, such as the title, abstract, descriptor fields to capture the meaning of a document at different levels of resolution or focusing on different characteristic aspects of a document. Once the specified query has been

executed by IR system, a user is presented with the retrieved document surrogates. Either the user is satisfied by the retrieved information or he will evaluate the retrieved documents and modify the query to initiate a further search. The process of query modification based on user evaluation of the retrieved documents is known as relevance feedback [Lancaster and Warner 1993]. Information retrieval is an inherently interactive process, and the users can change direction by modifying the query surrogate, the conceptual query or their understanding of their information need.

It is worth noting here the results, which have been obtained in studies investigating the information-seeking process, that describe information retrieval in terms of the cognitive and affective symptoms commonly experienced by a library user. The findings by Kuhlthau et al. (1990) indicate that thoughts about the information need become clearer and more focused as users move through the search process. Similarly, uncertainty, confusion, and frustration are nearly universal experiences in the early stages of the search process, and they decrease as the search process progresses and feelings of being confident, satisfied, sure and relieved increase. The studies also indicate that cognitive attributes may affect the search process. User's expectations of the information system and the search process may influence the way they approach searching and therefore affect the intellectual access to information.

Analytical search strategies require the formulation of specific, well-structured queries and a systematic, iterative search for information, whereas browsing involves the generation of broad query terms and a scanning of much larger sets of information in a relatively unstructured fashion. Campagnoni et al. (1989) have found in information retrieval studies in hypertext systems that the predominant search strategy is "browsing" rather than "analytical search". Many users, especially novices, are unwilling or unable to precisely formulate their search objectives, and browsing places less cognitive load on them. Furthermore, their research showed that search strategy is only one dimension of effective information retrieval; individual differences in visual skill appear to play an equally important role.

These two studies argue for information displays that provide a spatial overview of the data elements and that simultaneously provide rich visual cues about the content of the individual data elements. Such a representation is less likely to increase the anxiety that is a natural part of the early stages of the search process and it caters for a browsing interaction style, which is appropriate especially in the beginning, when many users are unable to precisely formulate their search objectives.

3 Major Information Retrieval Models

The following major models have been developed to retrieve information: the **Boolean** model, the **Statistical** model, which includes the vector space and the probabilistic retrieval model, and the **Linguistic and Knowledge-based** models. The first model is often referred to as the "exact match" model; the latter ones as the "best match" models [Belkin and Croft 1992]. The material presented here is based on the textbooks by Lancaster and Warner (1992) as well as Frakes and Baeza-Yates (1992), the review article by Belkin and Croft (1992), and discussions with Richard Marcus, my thesis advisor and mentor in the field of information retrieval.

Queries generally are less than perfect in two respects: First, they retrieve some irrelevant documents. Second, they do not retrieve all the relevant documents. The following two measures are usually used to evaluate the effectiveness of a retrieval method. The first one,

called the *precision rate*, is equal to the proportion of the retrieved documents that are actually relevant. The second one, called the *recall rate*, is equal to the proportion of all relevant documents that are actually retrieved. If searchers want to raise precision, then they have to narrow their queries. If searchers want to raise recall, then they broaden their query. In general, there is an inverse relationship between precision and recall. Users need help to become knowledgeable in how to manage the precision and recall trade-off for their particular information need [Marcus 1991].

3.1.1 Standard Boolean

In Table 2.1 we summarize the defining characteristics of the standard Boolean approach and list its key advantages and disadvantages. It has the following strengths: 1) It is easy to implement and it is computationally efficient [Frakes and Baeza-Yates 1992]. Hence, it is the standard model for the current large-scale, operational retrieval systems and many of the major on-line information services use it. 2) It enables users to express structural and conceptual constraints to describe important linguistic features [Marcus 1991]. Users find that synonym specifications (reflected by OR-clauses) and phrases (represented by proximity relations) are useful in the formulation of queries [Cooper 1988, Marcus 1991]. 3) The Boolean approach possesses a great expressive power and clarity. Boolean retrieval is very effective if a query requires an exhaustive and unambiguous selection. 4) The Boolean method offers a multitude of techniques to broaden or narrow a query. 5) The Boolean approach can be especially effective in the later stages of the search process, because of the clarity and exactness with which relationships between concepts can be represented.

The standard Boolean approach has the following shortcomings: 1) Users find it difficult to construct effective Boolean queries for several reasons [Cooper 1988, Fox and Koll 1988, Belkin and Croft 1992]. Users are using the natural language terms AND, OR or NOT that have a different meaning when used in a query. Thus, users will make errors when they form a Boolean query, because they resort to their knowledge of English.

	Standard Boolean
Goal	<ul style="list-style-type: none"> • Capture conceptual structure and contextual information
Methods	<ul style="list-style-type: none"> • Coordination: AND, OR, NOT • Proximity • Fields • Stemming / Truncation
(+)	<ul style="list-style-type: none"> • Easy to implement • Computationally efficient => all the major on-line databases use it • Expressiveness and Clarity Synonym specifications (OR-clauses) and phrases (AND-clauses).
(-)	<ul style="list-style-type: none"> • Difficult to construct Boolean queries. • All or nothing AND too severe, and OR does not differentiate enough. • Difficult to control output: Null output <-> Overload. • No ranking • No weighting of index or query terms • No uncertainty measure

Table 1: summarizes the defining characteristics of the standard Boolean approach and list the its key advantages and disadvantages.

For example, in ordinary conversation a noun phrase of the form "A and B" usually refers to more entities than would "A" alone, whereas when used in the context of information retrieval it refers to fewer documents than would be retrieved by "A" alone. Hence, one of the common mistakes made by users is to substitute the AND logical operator for the OR logical operator when translating an English sentence to a Boolean query. Furthermore, to form complex queries, users must be familiar with the rules of precedence and the use of parentheses. Novice users have difficulty using parentheses, especially nested parentheses. Finally, users are overwhelmed by the multitude of ways a query can be structured or modified, because of the combinatorial explosion of feasible queries as the number of concepts increases. In particular, users have difficulty identifying and applying the different strategies that are available for narrowing or broadening a Boolean query [Marcus 1991, Lancaster and Warner 1993]. 2) Only documents that satisfy a query exactly are retrieved. On the one hand, the AND operator is too severe because it does not distinguish between the case when none of the concepts are satisfied and the case where all except one are satisfied. Hence, no or very few documents are retrieved when more than three and four criteria are combined with the Boolean operator AND (referred to as the Null Output problem). On the other hand, the OR operator does not reflect how many concepts have been satisfied. Hence, often too many documents are retrieved (the Output Overload problem). 3) It is difficult to control the number of retrieved documents. Users

are often faced with the null-output or the information overload problem and they are at loss of how to modify the query to retrieve the reasonable number documents. 4) The traditional Boolean approach does not provide a relevance ranking of the retrieved documents, although modern Boolean approaches can make use of the degree of coordination, field level and degree of stemming present to rank them [Marcus 1991]. 5) It does not represent the degree of uncertainty or error due the vocabulary problem [Belkin and Croft 1992].

3.1.2 Narrowing and Broadening Techniques

As mentioned earlier, a Boolean query can be described in terms of the following four operations: degree and type of coordination, proximity constraints, field specifications and degree of stemming as expressed in terms of word/string specifications. If users want to (re)formulate a Boolean query then they need to make informed choices along these four dimensions to create a query that is sufficiently broad or narrow depending on their information needs. Most narrowing techniques lower recall as well as raise precision, and most broadening techniques lower precision as well as raise recall. Any query can be reformulated to achieve the desired precision or recall characteristics, but generally it is difficult to achieve both. Each of the four kinds of operations in the query formulation has particular operators, some of which tend to have a narrowing or broadening effect. For each operator with a narrowing effect, there is one or more inverse operators with a broadening effect [Marcus 1991]. Hence, users require help to gain an understanding of how changes along these four dimensions will affect the broadness or narrowness of a query.

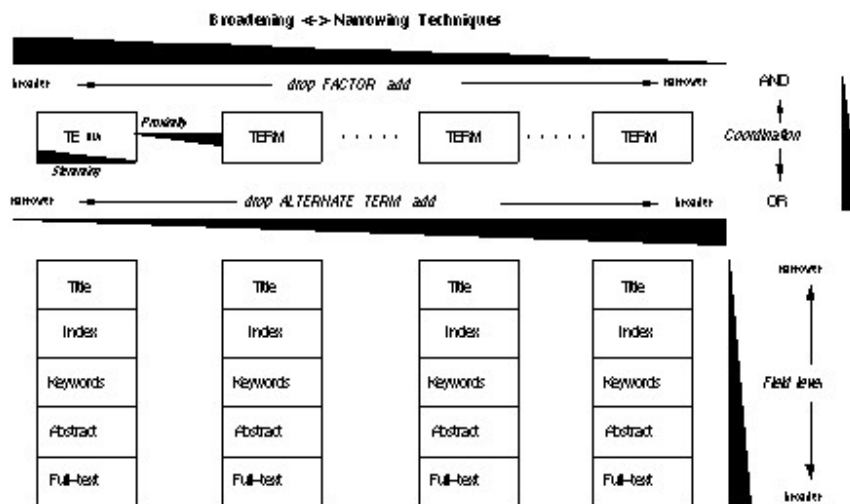


Figure 2: captures how coordination, proximity, field level and stemming affect the broadness or narrowness of a Boolean query. By moving in the direction in which the wedges are expanding the query is broadened.

Figure 2.2 shows how the four dimensions affect the broadness or narrowness of a query: 1) *Coordination*: the different Boolean operators AND, OR and NOT have the following effects when used to add a further concept to a query: a) the AND operator narrows a query; b)

the OR broadens it; c) the effect of the NOT depends on whether it is combined with an AND or OR operator. Typically, in searching textual databases, the NOT is connected to the AND, in which case it has a narrowing effect like the AND operator. 2) *Proximity*: The closer together two terms have to appear in a document, the more narrow and precise the query. The most stringent proximity constraint requires the two terms to be adjacent. 3) *Field level*: current document records have fields associated with them, such as the "Title", "Index", "Abstract" or "Full-text" field: a) the more fields that are searched, the broader the query; b) the individual fields have varying degrees of precision associated with them, where the "title" field is the most specific and the "full-text" field is the most general. 4) *Stemming*: The shorter the prefix that is used in truncation-based searching, the broader the query. By reducing a term to its morphological stem and using it as a prefix, users can retrieve many terms that are conceptually related to the original term [Marcus 1991].

Using Figure 2.2, we can easily read off how to broaden query. We just need to move in the direction in which the wedges are expanding: we use the OR operator (rather than the AND), impose no proximity constraints, search over all fields and apply a great deal of stemming. Similarly, we can formulate a very narrow query by moving in the direction in which the wedges are contracting: we use the AND operator (rather than the OR), impose proximity constraints, restrict the search to the title field and perform exact rather than truncated word matches. In Chapter 4 we will show how Figure 2.2 indicates how the broadness or narrowness of a Boolean query could be visualized.

3.1.3 Smart Boolean

There have been attempts to help users overcome some of the disadvantages of the traditional Boolean discussed above. We will now describe such a method, called *Smart Boolean*, developed by Marcus [1991, 1994] that tries to help users construct and modify a Boolean query as well as make better choices along the four dimensions that characterize a Boolean query. We are not attempting to provide an in-depth description of the Smart Boolean method, but to use it as a good example that illustrates some of the possible ways to make Boolean retrieval more user-friendly and effective. Table 2.2 provides a summary of the key features of the Smart Boolean approach.

Users start by specifying a natural language statement that is automatically translated into a Boolean Topic representation that consists of a list of factors or concepts, which are automatically coordinated using the AND operator. If the user at the initial stage can or wants to include synonyms, then they are coordinated using the OR operator. Hence, the Boolean Topic representation connects the different factors using the AND operator, where the factors can consist of single terms or several synonyms connected by the OR operator. One of the goals of the Smart Boolean approach is to make use of the structural knowledge contained in the text surrogates, where the different fields represent contexts of useful information. Further, the Smart Boolean approach wants to use the fact that related concepts can share a common stem. For example, the concepts "computers" and "computing" have the common stem comput*.

	Smart Boolean
Goal	<ul style="list-style-type: none"> • Structure search (re-)formulation process. • Use structural and contextual knowledge-bases and clarity of Boolean expressions.
Methods	<ul style="list-style-type: none"> • Natural language statement is automatically translated into Boolean Topic Representation • Boolean Topic Representation: <ul style="list-style-type: none"> ANDs of ORs of concepts Keyword /stem, all fields • Conceptual info. -> Coordination and Add/Drop Factor • Contextual info. -> Proximity • Structural info. -> Field levels • Synonym or word relationships -> Stemming/Truncation overlap => all this information can be used to rank documents • Techniques to Broaden and Narrow query
(+)	<ul style="list-style-type: none"> • No need for Boolean operators <ul style="list-style-type: none"> => Convert operator-free statement into ANDs of ORs • Assist user in query (re)formulation: <ul style="list-style-type: none"> by asking users targeted questions to automatically modify the query. • "Why irrelevant?" -> activates narrowing methods. • "Broaden by Dropping Factors" to estimate recall.
(-)	<ul style="list-style-type: none"> • How to visualize ? <ul style="list-style-type: none"> • Conceptual query representation (BTR) • Query modification techniques and their effects • Structured relevance feedback

Table 2: summarizes the defining characteristics of the Smart Boolean approach and list the its key advantages and disadvantages.

The initial strategy of the Smart Boolean approach is to start out with the broadest possible query within the constraints of how the factors and their synonyms have been coordinated. Hence, it modifies the Boolean Topic representation into the query surrogate by using only the stems of the concepts and searches for them over all the fields. Once the query surrogate has been performed, users are guided in the process of evaluating the retrieved document surrogates. They choose from a list of reasons to indicate why they consider certain documents as relevant. Similarly, they can indicate why other documents are not relevant by interacting with a list of possible reasons. This user feedback is used by the Smart Boolean system to automatically modify the Boolean Topic representation or the query surrogate, whatever is more appropriate. The Smart Boolean approach offers a rich set of strategies for modifying a query based on the received relevance feedback or the expressed need to narrow or broaden the query. The Smart Boolean retrieval paradigm has been implemented in the form of a system called CONIT, which is one of the earliest expert retrieval systems that was able to demonstrate that ordinary users, assisted by such a system, could perform equally well as experienced search intermediaries [Marcus 1983]. However, users have to navigate through a

series of menus listing different choices, where it might be hard for them to appreciate the implications of some of these choices. A key limitation of the previous versions of the CONIT system has been that lacked a visual interface. The most recent version has a graphical interface and it uses the tiling metaphor suggested by Anick et al. (1991), and discussed in section 10.4, to visualize Boolean coordination [Marcus 1994]. This visualization approach suffers from the limitation that it enables users to visualize specific queries, whereas we will propose a visual interface that represents all whole range of related Boolean queries in a single display, making changes in Boolean coordination more user-friendly. Further, the different strategies of modifying a query in CONIT require a better visualization metaphor to enable users to make use these search heuristics. In Chapter 4 we show how some of these modification techniques can be visualized.

3.1.4 Extended Boolean Models

Several methods have been developed to extend the Boolean model to address the following issues: 1) The Boolean operators are too strict and ways need to be found to soften them. 2) The standard Boolean approach has no provision for ranking. The Smart Boolean approach and the methods described in this section provide users with relevance ranking [Fox and Koll 1988, Marcus 1991]. 3) The Boolean model does not support the assignment of weights to the query or document terms. We will briefly discuss the *P-norm* and the *Fuzzy Logic* approaches that extend the Boolean model to address the above issues.

Extended Boolean Models	
Goal	<ul style="list-style-type: none"> • Less strict Boolean operators • Ranked output
Methods	<div style="border: 1px solid black; width: 20px; height: 20px; margin-bottom: 10px;"></div> <ul style="list-style-type: none"> • Fuzzy logic <p>[OR -> max], [AND -> min] and [NOT -> 1 - max]</p> <p>(-) Lack of sensitivity of min and max: $\min(0.2, 0.8) = \min(0.2, 0.3)$.</p>

Table 3: summarizes the defining characteristics of the Extended Boolean approach and list the its key advantages and disadvantages.

The **P-norm** method developed by Fox (1983) allows query and document terms to have weights, which have been computed by using term frequency statistics with the proper normalization procedures. These normalized weights can be used to rank the documents in the order of decreasing distance from the point (0, 0, ..., 0) for an OR query, and in order of increasing distance from the point (1, 1, ..., 1) for an AND query. Further, the Boolean operators have a coefficient P associated with them to indicate the degree of strictness of the operator

(from 1 for least strict to infinity for most strict, i.e., the Boolean case). The P-norm uses a distance-based measure and the coefficient P determines the degree of exponentiation to be used. The exponentiation is an expensive computation, especially for P-values greater than one.

In **Fuzzy Set theory**, an element has a varying degree of membership to a set instead of the traditional binary membership choice. The weight of an index term for a given document reflects the degree to which this term describes the content of a document. Hence, this weight reflects the degree of membership of the document in the fuzzy set associated with the term in question. The degree of membership for union and intersection of two fuzzy sets is equal to the maximum and minimum, respectively, of the degrees of membership of the elements of the two sets. In the "Mixed Min and Max" model developed by Fox and Sharat (1986) the Boolean operators are softened by considering the query-document similarity to be a linear combination of the min and max weights of the documents.

3.2 Statistical Model

The *vector space* and *probabilistic* models are the two major examples of the statistical retrieval approach. Both models use statistical information in the form of term frequencies to determine the relevance of documents with respect to a query. Although they differ in the way they use the term frequencies, both produce as their output a list of documents ranked by their estimated relevance. The statistical retrieval models address some of the problems of Boolean retrieval methods, but they have disadvantages of their own. Table 2.4 provides summary of the key features of the vector space and probabilistic approaches. We will also describe *Latent Semantic Indexing* and *clustering* approaches that are based on statistical retrieval approaches, but their objective is to respond to what the user's query did not say, could not say, but somehow made manifest [Furnas et al. 1983, Cutting et al. 1991].

3.2.1 Vector Space Model

The **vector space model** represents the documents and queries as vectors in a multidimensional space, whose dimensions are the terms used to build an index to represent the documents [Salton 1983]. The creation of an index involves lexical scanning to identify the significant terms, where morphological analysis reduces different word forms to common "stems", and the occurrence of those stems is computed. Query and document surrogates are compared by comparing their vectors, using, for example, the cosine similarity measure. In this model, the terms of a query surrogate can be weighted to take into account their importance, and they are computed by using the statistical distributions of the terms in the collection and in the documents [Salton 1983]. The vector space model can assign a high ranking score to a document that contains only a few of the query terms if these terms occur infrequently in the collection but frequently in the document. The vector space model makes the following assumptions: 1) The more similar a document vector is to a query vector, the more likely it is that the document is relevant to that query. 2) The words used to define the dimensions of the space are orthogonal or independent. While it is a reasonable first approximation, the assumption that words are pairwise independent is not realistic.

3.2.2 Probabilistic Model

The **probabilistic retrieval** model is based on the Probability Ranking Principle, which states that an information retrieval system is supposed to rank the documents based on their

probability of relevance to the query, given all the evidence available [Belkin and Croft 1992]. The principle takes into account that there is uncertainty in the representation of the information need and the documents. There can be a variety of sources of evidence that are used by the probabilistic retrieval methods, and the most common one is the statistical distribution of the terms in both the relevant and non-relevant documents.

We will now describe the state-of-art system developed by Turtle and Croft (1991) that uses Bayesian inference networks to rank documents by using multiple sources of evidence to compute the conditional probability $P(\text{Info need}|\text{document})$ that an information need is satisfied by a given document. An inference network consists of a directed acyclic dependency graph, where edges represent conditional dependency or causal relations between propositions represented by the nodes. The inference network consists of a document network, a concept representation network that represents indexing vocabulary, and a query network representing the information need. The concept representation network is the interface between documents and queries. To compute the rank of a document, the inference network is instantiated and the resulting probabilities are propagated through the network to derive a probability associated with the node representing the information need. These probabilities are used to rank documents.

The statistical approaches have the following strengths: 1) They provide users with a relevance ranking of the retrieved documents. Hence, they enable users to control the output by setting a relevance threshold or by specifying a certain number of documents to display. 2) Queries can be easier to formulate because users do not have to learn a query language and can use natural language. 3) The uncertainty inherent in the choice of query concepts can be represented. However, the statistical approaches have the following shortcomings: 1) They have a limited expressive power. For example, the NOT operation can not be represented because only positive weights are used. It can be proven that only 2^N of the 2^{2N} possible Boolean queries can be generated by the statistical approaches that use weighted linear sums to rank the documents. This result follows from the analysis of Linear Threshold Networks or Boolean Perceptrons [Anthony and Biggs 1992]. For example, the very common and important Boolean query $((A \text{ and } B) \text{ or } (C \text{ and } D))$ can not be represented by a vector space query (see section 5.4 for a proof). Hence, the statistical approaches do not have the expressive power of the Boolean approach. 3) The statistical approach lacks the structure to express important linguistic features such as phrases. Proximity constraints are also difficult to express, a feature that is of great use for experienced searchers. 4) The computation of the relevance scores can be computationally expensive. 5) A ranked linear list provides users with a limited view of the information space and it does not directly suggest how to modify a query if the need arises [Spoerri 1993, Hearst 1994]. 6) The queries have to contain a large number of words to improve the retrieval performance. As is the case for the Boolean approach, users are faced with the problem of having to choose the appropriate words that are also used in the relevant documents.

Table 4 summarizes the advantages and disadvantages that are specific to the vector space and probabilistic model, respectively. This table also shows the formulas that are commonly used to compute the term weights. The two central quantities used are the inverse term frequency in a collection (*idf*), and the frequencies of a term *i* in a document *j* (*freq(i,j)*). In the probabilistic model, the weight computation also considers how often a term appears in the relevant and irrelevant documents, but this presupposes that the relevant documents are known or that these frequencies can be reliably estimated.

<i>Statistical</i>	Vector Space	Probabilistic
Motivation	Simplify query formulation Ability to control output	Address uncertainty in query representations
Goal	Rank the output based on <div style="display: flex; justify-content: space-around;"> Similarity Probability of Relevance </div>	
Methods	Cosine measure	Use of different models
Source	Query Term Statistics <u>Vector-Space:</u> <ul style="list-style-type: none"> • $\text{similarity}(Q,D) = \sum (w_{iq} \times w_{ij}) / \text{"normalizer"}$ where $w_{iq} = (0.5 + 0.5 \text{freq}_{iq} / \text{maxfreq}_{iq}) \times \text{idf}(i)$ $w_{ij} = \text{freq}_{ij} \times \text{idf}(i)$ • inverse term freq. in collection $\text{idf}(i) = \log_2 (N-n(i)) / n(i)$. <u>Probabilistic:</u> <ul style="list-style-type: none"> • term weight $= \log [(r_t / R-r_t) / ((n_t - r_t) / ((N-n_t) - (R-r_t)))]$ ="hits / misses) / (false alarms/correct misses)" • $\text{similarity}_{jk} = \sum (C + \text{idf}(i)) \times \text{tf}(i,j)$ where $\text{tf}(i,j) = K + (1-K) (\text{freq}(i,j) / \text{maxfreq}(j))$. 	
Issues	<ul style="list-style-type: none"> • How to express NOT ? • Proximity searches ? • Limited expressive power • Computationally intensive • Assumes that terms are independent. • Lack of structure to represent important linguistic features • How to better visualize the retrieved set ? 	<ul style="list-style-type: none"> • Estimation of needed probabilities • Prior knowledge needed. • Independence assumption • Boolean relations lost. • Which model is best ?

Table 4: summarizes the defining characteristics of the statistical retrieval approach, which includes the vector space and the probabilistic model and we list the their key advantages and disadvantages.

If users provide the retrieval system with relevance feedback, then this information is used by the statistical approaches to recompute the weights as follows: the weights of the query terms in the relevant documents are increased, whereas the weights of the query terms that do not appear in the relevant documents are decreased [Salton and Buckley 1990]. There are multiple ways of computing and updating the weights, where each has its advantages and disadvantages. We do not discuss these formulas in more detail, because research on relevance feedback has shown that significant effectiveness improvements can be gained by using quite simple feedback techniques [Salton and Buckley 1990]. Furthermore, what is

important to this thesis is that the statistical retrieval approach generates a ranked list, however how this ranking has been computed in detail is immaterial for the purpose of this thesis.

3.2.3 Latent Semantic Indexing

Several statistical and AI techniques have been used in association with domain semantics to extend the vector space model to help overcome some of the retrieval problems described above, such as the "dependence problem" or the "vocabulary problem". One such method is **Latent Semantic Indexing (LSI)**. In LSI the associations among terms and documents are calculated and exploited in the retrieval process. The assumption is that there is some "latent" structure in the pattern of word usage across documents and that statistical techniques can be used to estimate this latent structure. An advantage of this approach is that queries can retrieve documents even if they have no words in common. The LSI technique captures deeper associative structure than simple term-to-term correlations and is completely automatic. The only difference between LSI and vector space methods is that LSI represents terms and documents in a reduced dimensional space of the derived indexing dimensions. As with the vector space method, differential term weighting and relevance feedback can improve LSI performance substantially.

Foltz and Dumais (1992) compared four retrieval methods that are based on the vector-space model. The four methods were the result of crossing two factors, the first factor being whether the retrieval method used Latent Semantic Indexing or keyword matching, and the second factor being whether the profile was based on words or phrases provided by the user (Word profile), or documents that the user had previously rated as relevant (Document profile). The LSI match-document profile method proved to be the most successful of the four methods. This method combines the advantages of both LSI and the document profile. The document profile provides a simple, but effective, representation of the user's interests. Indicating just a few documents that are of interest is as effective as generating a long list of words and phrases that describe one's interest. Document profiles have an added advantage over word profiles: users can just indicate documents they find relevant without having to generate a description of their interests.

3.3 Linguistic and Knowledge-based Approaches

In the simplest form of automatic text retrieval, users enter a string of keywords that are used to search the inverted indexes of the document keywords. This approach retrieves documents based solely on the presence or absence of exact single word strings as specified by the logical representation of the query. Clearly this approach will miss many relevant documents because it does not capture the complete or deep meaning of the user's query. The Smart Boolean approach and the statistical retrieval approaches, each in their specific way, try to address this problem (see Table 2.5). Linguistic and knowledge-based approaches have also been developed to address this problem by performing a morphological, syntactic and semantic analysis to retrieve documents more effectively [Lancaster and Warner 1993]. In a morphological analysis, roots and affixes are analyzed to determine the part of speech (noun, verb, adjective etc.) of the words. Next complete phrases have to be parsed using some form of syntactic analysis. Finally, the linguistic methods have to resolve word ambiguities and/or generate relevant synonyms or quasi-synonyms based on the semantic relationships between words. The development of a sophisticated linguistic retrieval system is difficult and it requires complex knowledge bases of semantic information and retrieval heuristics. Hence these

systems often require techniques that are commonly referred to as artificial intelligence or expert systems techniques.

3.3.1 DR-LINK Retrieval System

We will now describe in some detail the DR-LINK system developed by Liddy et al., because it represents an exemplary linguistic retrieval system. DR-LINK is based on the principle that retrieval should take place at the conceptual level and not at the word level. Liddy et al. attempt to retrieve documents on the basis of what people mean in their query and not just what they say in their query. DR-LINK system employs sophisticated, linguistic text processing techniques to capture the conceptual information in documents. Liddy et al. have developed a modular system that represents and matches text at the lexical, syntactic, semantic, and the discourse levels of language. Some of the modules that have been incorporated are: The Text Structurer is based on discourse linguistic theory that suggests that texts of a particular type have a predictable structure which serves as an indication where certain information can be found. The Subject Field Coder uses an established semantic coding scheme from a machine-readable dictionary to tag each word with its disambiguated subject code (e.g., computer science, economics) and to then produce a fixed-length, subject-based vector representation of the document and the query. The Proper Noun Interpreter uses a variety of processing heuristics and knowledge bases to produce: a canonical representation of each proper noun; a classification of each proper noun into thirty-seven categories; and an expansion of group nouns into their constituent proper noun members. The Complex Nominal Phraser provides means for precise matching of complex semantic constructs when expressed as either adjacent nouns or a non-predicating adjective and noun pair. Finally, The Natural Language Query Constructor takes as input a natural language query and produces a formal query that reflects the appropriate logical combination of text structure, proper noun, and complex nominal requirements of the user's information need. This module interprets a query into pattern-action rules that translate each sentence into a first-order logic assertion, reflecting the Boolean-like requirements of queries.

Linguistic Level	Boolean Retrieval	Statistical	Linguistic and Knowledge-based
Lexical	Stop word list	Stop word list	Lexicon
Morphological	Truncation symbol	Stemming	Morphological analysis
Syntactic	Proximity operators	Statistical phrases	Grammatical phrases
Semantic	Thesaurus	Clusters of co-occurring words	Network of words/phrases in semantic relationships

Table 5: characterizes the major retrieval methods in terms of how deal with lexical, morphological, syntactic and semantic issues.

To summarize, the DR-LINK retrieval system represents content at the conceptual level rather than at the word level to reflect the multiple levels of human language comprehension. The text representation combines the lexical, syntactic, semantic, and discourse levels of understanding to predict the relevance of a document. DR-LINK accepts natural language statements, which it translates into a precise Boolean representation of the user's relevance requirements. It also produces a summary-level, semantic vector representations of queries and documents to provide a ranking of the documents.

4 Conclusion

There is a growing discrepancy between the retrieval approach used by existing commercial retrieval systems and the approaches investigated and promoted by a large segment of the information retrieval research community. The former is based on the Boolean or Exact Matching retrieval model, whereas the latter ones subscribe to statistical and linguistic approaches, also referred to as the Partial Matching approaches. First, the major criticism leveled against the Boolean approach is that its queries are difficult to formulate. Second, the Boolean approach makes it possible to represent structural and contextual information that would be very difficult to represent using the statistical approaches. Third, the Partial Matching approaches provide users with a ranked output, but these ranked lists obscure

Key Problems	Possible Solutions
Selection of Search Vocabulary	<ul style="list-style-type: none"> • Thesaurus • Latent Semantic Indexing
Search strategy (re)formulation	<ul style="list-style-type: none"> • Smart Boolean • Statistical & Linguistic Approaches • Thesaurus • Graphical Interfaces
Information Overload	<ul style="list-style-type: none"> • Ranking • Clustering • Visualization

Table 6: lists some of the key problems in the field of information retrieval and possible solutions.

valuable information. Fourth, recent retrieval experiments have shown that the Exact and Partial matching approaches are complementary and should therefore be combined [Belkin et al. 1993].

In Table 2.6 we summarize some of the key problems in the field of information retrieval and possible solutions to them. We will attempt to show in this thesis: 1) how visualization can offer ways to address these problems; 2) how to formulate and modify a query; 3) how to deal with large sets of retrieved documents, commonly referred to as the information overload

problem. In particular, this thesis overcomes one of the major "bottlenecks" of the Boolean approach by showing how Boolean coordination and its diverse narrowing and broadening techniques can be visualized, thereby making it more user-friendly without limiting its expressive power. Further, this thesis shows how both the Exact and Partial Matching approaches can be visualized in the same visual framework to enable users to make effective use of their respective strengths.