



## Cool New Stuff!

NCBI Resources for Phylogenically-Defined Next Generation Analysis in and out of the Cloud  
 Ben Busby (ben.busby@nih.gov)  
 NCBI Genome Resources Workshop  
 PAG XXVII January 14, 2019



Identical Protein Groups

**Argonaute family protein**

Gene/Fold FASTA GenInfo BLAST

Name: Argonaute family protein  
 RefSeq Identifier/Protein: NP\_858118.1, NP\_858118.1, NP\_858118.1, NP\_858118.1, NP\_858118.1  
 Assembly Accessions: 2  
 Protein Accessions: 6  
 CDS Regions: 6  
 Total Rows: 7

#	Accession	CDS Region In	Accession	Protein	Name	Organism	Strain	Assembly
1	NP_858118.1	1377446..1377732 (+)	NP_858118.1	Argonaute family protein	Arabidopsis thaliana			NCBI_200907128.1
2	NP_858118.1	1377446..1377732 (+)	NP_858118.1	Argonaute family protein	Arabidopsis thaliana			NCBI_200907128.1
3	NP_858118.1	1377446..1377732 (+)	NP_858118.1	Argonaute family protein	Arabidopsis thaliana			NCBI_200907128.1
4	NP_858118.1	1377446..1377732 (+)	NP_858118.1	Argonaute family protein	Arabidopsis thaliana			NCBI_200907128.1
5	NP_858118.1	1377446..1377732 (+)	NP_858118.1	Argonaute family protein	Arabidopsis thaliana			NCBI_200907128.1
6	NP_858118.1	1377446..1377732 (+)	NP_858118.1	Argonaute family protein	Arabidopsis thaliana			NCBI_200907128.1
7	NP_858118.1	1377446..1377732 (+)	NP_858118.1	Argonaute family protein	Arabidopsis thaliana			NCBI_200907128.1



Tom Madden's BLAST\_v5 Talk (slightly modified)

If you'd like to see the original slides from Tom, or this slide deck, check out <https://www.slideshare.net/benbusby>

## Search the Argonaute protein family with BLAST!

- Against the nr protein database
- Limit the search to Schizosaccharomyces pombe



## BLASTDBv5: a better BLAST database!

- Limit search by **taxonomy** with a command-line parameter `-taxids`
- Improved **performance** when limiting BLAST search with accessions.
  - **Faster!**
- **Retrieve sequences** by taxonomy from the BLAST database.
  - Extract data!



Database: All non-redundant GenBank CDS translations+PDB+SwissProt+PIR+PRF excluding environmental samples from WGS projects  
 148,264,992 sequences; 54,322,627,498 total letters

Query= NP\_858118.1 Argonaute family protein [Arabidopsis thaliana]  
 Length=997

Sequences producing significant alignments:

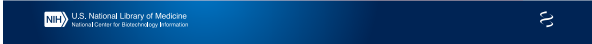
Accession	Organism	Score	E (Bits)	Value
NP_587782.1	argonaute [Schizosaccharomyces pombe]	487	2e-128	

NP\_587782.1 argonaute [Schizosaccharomyces pombe]  
 CAEL9275.1 argonaute [Schizosaccharomyces pombe]  
 Length=834

Score = 487 bits (1047), Expect = 2e-128, Method: Compositional matrix adjust.  
 Identities = 284/892 (32%), Positives = 449/892 (50%), Gaps = 91/892 (10%)

Query 136 SLPPASSKAVTFPPVRGRGTLGKQVRRANHF-LVQVADRLYHYDVSINPEVJSKTVNR 194  
 S: P+S Ae BRG G LGK++ +HW F ++ + + Y V +  
 Sbjct 2 SYKPSSEIAL----RPGYGLGKQITLKANFFQIISLPMETINQVHYVIGD----GSRVPR 54

Query 195 NVMKLLWNKYDSHLGGKS---PAYDGRKSLYTAGLPFDSKEFVNLAKRADGSSGDK 251  
 +L + + G S YDGR ++ G + + + VH+ GS  
 Sbjct 55 KQSLILNSKVKQYFGSSMWSYDGRKQKQKGGIADGTLK--VLI-----GESHSP 106



## BLAST command-line parameters

```
blastp
  -db nr
  -query NP_850110.fsa
  -outfmt "7 qaccver saccver pident length evaluate ssciname" # tabular output
  -out NP_850110.4896.tab
  -task blastp-fast # runs faster!
  -num_threads 8 # use 8 CPUs
  -evalue 0.05 # limit to more significant matches
  -taxids 4896 # limit search to Schizosaccharomyces pombe
```

## Need taxids for a whole Phylum?

```
bash-4.2$ get_species_taxids.sh -n Ascomycota
Taxid: 4890
rank: phylum
division: ascomycetes
scientific name: Ascomycota
common name: ascomycetes
1 match(es) found.
bash-4.2$ get_species_taxids.sh -t 4890 > Ascomycota.taxidlist
4754
4896
4897
4899
4983
4989
4911
4914
```

## A script to find taxid

Distributed with the BLAST package!

```
bash-4.2$ get_species_taxids.sh -n "Schizosaccharomyces pombe"
Taxid: 4896
rank: species
division: ascomycetes
scientific name: Schizosaccharomyces pombe
common name: fission yeast
1 match(es) found.
```

## BLAST against Ascomycota sequences!

```
blastp
  -db nr
  -query NP_850110.fsa
  -outfmt "7 qaccver saccver pident length evaluate ssciname" # tabular output
  -out NP_850110.Ascomycota.tab
  -task blastp-fast # runs faster!
  -num_threads 8 # use 8 CPUs
  -evalue 0.05 # limit to more significant matches
  -taxidlist Ascomycota.taxids # limit search to Ascomycota
```

## Tabular output

```
# BLASTP 2.8.0+
# Query: NP_850110.1 Argonaute family protein [Arabidopsis thaliana]
# Database: nr
# Fields: query acc.ver, subject acc.ver, % identity, alignment length, evaluate, subject sci name
# 1 hits found
NP_850110.1 NP_587782.1 31.839 892 1.57e-128 Schizosaccharomyces pombe
# BLAST processed 1 queries
```

## Ascomycota results

```
# BLASTP 2.8.0+
# Query: NP_850110.1 Argonaute family protein [Arabidopsis thaliana]
# Database: nr
# Fields: query acc.ver, subject acc.ver, % identity, alignment length, evaluate, subject sci name
# 581 hits found
NP_850110.1 GAD47878.1 36.139 891 1.93e-164 Saitoella complicata NRRL Y-17804
NP_850110.1 XP_019807722.1 36.558 829 1.39e-156 Saitoella complicata NRRL Y-17804
NP_850110.1 O0075158.1 34.228 855 1.83e-146 Ligomyces starknyi NRRL Y-21557
NP_850110.1 CU989782.1 31.767 894 6.68e-134 Tuber aestivum
NP_850110.1 XP_006233847.1 38.624 898 3.85e-130 Tuber melanosporum Mel2b
NP_850110.1 XP_013825347.1 32.548 882 1.01e-127 Schizosaccharomyces cryophilus OY26
NP_850110.1 OR019337.1 32.993 882 2.93e-127 Protomyces lactucaedebilis
NP_850110.1 NP_587782.1 31.839 892 1.45e-125 Schizosaccharomyces pombe
NP_850110.1 XP_013819123.1 33.333 852 5.49e-125 Schizosaccharomyces octosporus yf5286
NP_850110.1 EP544588.1 38.718 982 7.32e-125 Dactyliellina haptotyla CBS 298-58
NP_850110.1 CU54138.1 33.684 855 1.87e-124 Tuber aestivum
NP_850110.1 XP_011108241.1 31.188 981 2.47e-124 Arthrobotrys oligospora ATCC 24927
NP_850110.1 EHC45585.1 31.823 985 1.42e-122 Drechslerella stenobrocha 248
NP_850110.1 EHC45585.1 31.842 931 3.45e-128 Pseudomyces pannonum VM F-4518 (FW-2643)
NP_850110.1 KF208212.1 38.444 959 1.22e-116 Drechslerella stenobrocha 248
NP_850110.1 XP_002173826.1 31.818 858 4.95e-116 Schizosaccharomyces japonicus yf5275
NP_850110.1 XP_011128792.1 31.973 882 8.11e-116 Arthrobotrys oligospora ATCC 24927
NP_850110.1 KF213812.1 38.341 959 9.18e-116 Pseudomyces pannonum VM F-4520 (FW-2644)
NP_850110.1 XP_011119832.1 31.585 877 2.56e-114 Arthrobotrys oligospora ATCC 24927
NP_850110.1 KF111959.1 38.714 957 1.34e-112 Pseudomyces pannonum VM F-4281 (FW-2241)
NP_850110.1 EP543673.1 29.814 966 1.75e-111 Dactyliellina haptotyla CBS 288-58
```

## BLAST excluding Viridiplantae

```
blastp
-db nr
-query NP_850110.fsa
-outfmt "7 qcovser pident length evalui ssciname"
-out NP_850110.noGreenPlants.tab
-task blastp-fast
-num_threads 8
-evalue 0.05
-negative_taxidlist Viridiplantae.taxids # exclude green plants
```

## Everything

```
# BLASTP 2.6.0+
# Query: NP_850110.1 Argonaute family protein [Arabidopsis thaliana]
# Database: nr
# Fields: query acc.ver, subject acc.ver, % identity, alignment length, evalue, subject sci name
# 521 hits found
NP_850110.1 NP_850110.1 100.000 997 0.0 Arabidopsis thaliana
NP_850110.1 GAF68978.1 99.988 997 0.0 Arabidopsis thaliana
NP_850110.1 AA021514.1 99.988 997 0.0 Arabidopsis thaliana
NP_850110.1 XP_002588982.1 89.634 984 0.0 Arabidopsis lyrata subsp. lyrata
NP_850110.1 XP_018473185.1 82.513 995 0.0 Camelina sativa
NP_850110.1 XP_015017939.1 81.508 994 0.0 Camelina sativa
NP_850110.1 XP_015018467.1 81.094 961 0.0 Camelina sativa
NP_850110.1 XP_018518727.1 81.094 961 0.0 Camelina sativa
NP_850110.1 XP_086296291.1 77.756 1036 0.0 Capsella rubella
NP_850110.1 XP_018439599.1 74.926 1069 0.0 Camelina sativa
NP_850110.1 XP_013692156.1 73.548 973 0.0 Brassica napus
NP_850110.1 XP_012636559.1 75.548 973 0.0 Brassica oleracea var. oleracea
NP_850110.1 XP_009183826.1 74.705 976 0.0 Brassica rapa
NP_850110.1 KFK44936.1 73.053 979 0.0 Arabis alpina
NP_850110.1 XP_015647627.1 74.603 976 0.0 Brassica napus
NP_850110.1 XP_018448495.1 79.482 869 0.0 Raphanus sativus
NP_850110.1 XP_004489883.1 72.788 989 0.0 Eutrema sativagineum
NP_850110.1 XP_018498899.1 72.019 986 0.0 Camelina sativa
NP_850110.1 XP_018521782.1 67.286 991 0.0 Cleome hassleriana
NP_850110.1 XP_021818669.1 65.119 926 0.0 Manihot esculenta
NP_850110.1 PNT55858.1 67.086 875 0.0 Populus trichocarpa
```



## No Viridiplantae

```
# BLASTP 2.6.0+
# Query: NP_850110.1 Argonaute family protein [Arabidopsis thaliana]
# Database: nr
# Fields: query acc.ver, subject acc.ver, % identity, alignment length, evalue, subject sci name
# 588 hits found
NP_850110.1 XP_021356153.1 42.881 884 0.0 Mizuhopecten yessoensis
NP_850110.1 QW49186.1 41.968 884 0.0 Mizuhopecten yessoensis
NP_850110.1 XP_021356154.1 41.968 884 0.0 Mizuhopecten yessoensis
NP_850110.1 XP_021356147.1 41.742 884 0.0 Mizuhopecten yessoensis
NP_850110.1 XP_021965366.1 41.619 877 0.0 Folsomia candida
NP_850110.1 G0N88002.1 41.686 878 0.0 Orchesella cincta
NP_850110.1 G0A8021.1 41.723 877 0.0 Folsomia candida
NP_850110.1 XP_014666293.1 42.648 856 0.0 Priapulus caudatus
NP_850110.1 XP_082413188.1 42.272 854 0.0 Ixodes scapularis
NP_850110.1 ANH83651.1 48.785 915 0.0 Ptilius walprovincialis
NP_850110.1 XP_014666292.1 42.748 854 0.0 Priapulus caudatus
NP_850110.1 XP_021965365.1 41.563 883 0.0 Folsomia candida
NP_850110.1 XP_021965367.1 41.563 883 0.0 Folsomia candida
NP_850110.1 XP_082413231.1 48.781 893 0.0 Drosophila sechellia
NP_850110.1 EL199798.1 41.344 878 0.0 Capitella teleta
NP_850110.1 XP_012178889.1 48.252 954 0.0 Bombyx terrestris
NP_850110.1 XP_021356148.1 41.298 884 0.0 Mizuhopecten yessoensis
NP_850110.1 XP_021356151.1 41.516 884 0.0 Mizuhopecten yessoensis
NP_850110.1 G0F97116.1 48.488 981 0.0 Drosophila willistoni
NP_850110.1 XP_015833121.1 48.488 981 0.0 Drosophila willistoni
NP_850110.1 XP_022667621.1 48.488 981 0.0 Drosophila willistoni
```

## Timings

Search set	Run time (seconds)
Schizosaccharomyces pombe	8.7
Ascomycota	11.7
No Viridiplantae	101
Everything	91

8 CPUs, best of three runs using time command.



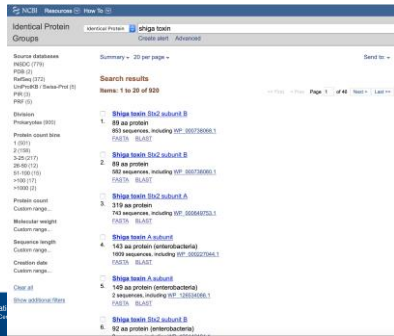
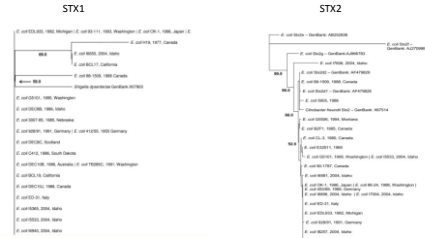
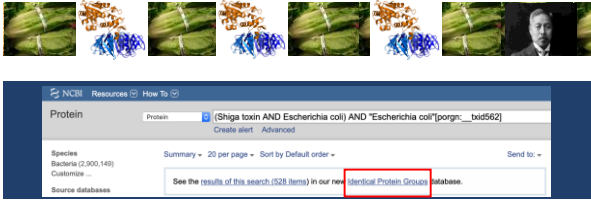
Tax BLAST report

Query= NP\_850110.1 Argonaute family protein [Arabidopsis thaliana]

Length=997

Accession	Description	Score	E-value
XP_021356153	protein argonaute-2-like isoform X6 [Mizuhopecten yessoensis]	668	0.0
QW49186	protein argonaute-2 [Mizuhopecten yessoensis]	666	0.0
XP_021356154	protein argonaute-2-like isoform X7 [Mizuhopecten yessoensis]	665	0.0
XP_021965367	protein argonaute-2-like isoform X2 [Mizuhopecten yessoensis]	665	0.0
XP_021356148	protein argonaute-2-like isoform X1 [Mizuhopecten yessoensis]	664	0.0
XP_021356151	protein argonaute-2-like isoform X3 [Mizuhopecten yessoensis]	664	0.0
XP_021356150	protein argonaute-2-like isoform X4 [Mizuhopecten yessoensis]	662	0.0
XP_021356149	protein argonaute-2-like isoform X5 [Mizuhopecten yessoensis]	662	0.0
XP_021356152	protein argonaute-2-like isoform X8 [Mizuhopecten yessoensis]	662	0.0
XP_021356155	protein argonaute-2-like isoform X9 [Mizuhopecten yessoensis]	662	0.0
XP_021356146	protein argonaute-2-like isoform X1 [Mizuhopecten yessoensis]	662	0.0
XP_021356149	protein argonaute-2-like isoform X1 [Mizuhopecten yessoensis]	662	0.0
XP_021965366	protein argonaute-2-like isoform X2 [Folsomia candida]	661	0.0
G0A8021	protein argonaute-2 [Folsomia candida]	660	0.0
XP_021965365	protein argonaute-2-like isoform X1 [Folsomia candida]	657	0.0
XP_021965367	protein argonaute-2-like isoform X2 [Folsomia candida]	657	0.0
G0N88002	protein argonaute-2 [Orchesella cincta]	660	0.0
XP_014666293	PRSDICTED: protein argonaute-2-like isoform X5 [Priapulus caudatus]	659	0.0
XP_014666292	PRSDICTED: protein argonaute-2-like isoform X6 [Priapulus caudatus]	658	0.0

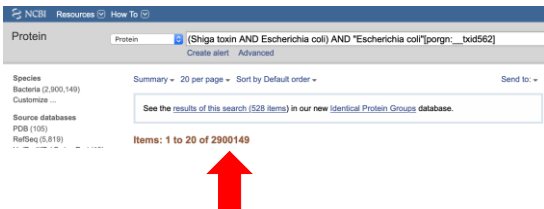




### Search the Shiga toxin family protein with BLAST

- Against the nr protein database
- Limit the search to E coli

### Why Use Identical Protein Groups?



### Search the Shiga toxin family protein with BLAST

- Against the nr protein database
- Limit the search to E coli

But which Shiga toxin?

### Stx1 and Stx2 Example Sequences

```

Shiga toxin Stx1 subunit A [Escherichia coli]
Shiga toxin Stx2 subunit A [Escherichia coli]

```

### But what did I need to do to make my demo?

Grab sequences!

```

ben.busby@known-blast-workbench-sbx-camacho:~$ search -db protein -query WP_123135120 | efetch -format fasta
>WP_123135120.1 Shiga toxin Stx1 subunit A [Escherichia coli]
MELIIFNLVFFRFVIFSNVVAKEFTLDFSTACTVYSELVYKSAIETPLDTISSGDTLLMIDSGTDGN
LFAVDVPRGIDPPEGRFNRLIVERNLYVTFGNVNRINNVYRFAADFSHTFPGTAVTLSDGSSYTLQ
RYVADISQNDQINRHSLLTLYLDLMSHSDTSLTQSVFARMRLRVYVTAALRFGRIGQRFTLLDLSGR
SYMTAFQVQLIANKRLSGLPVDYHSDVRFVRFISFDNMLLGEVALILNCHHHSRVAELVPEEFP
SMCPVDRVRVRIITHNKLWDSSTLGAILLRRAISS

```

Do I need to get all of nr?!

update\_blast.pl -- decompress works

But theres a better way!



### Stx1 and Stx2 Example Sequences

```

Shiga toxin Stx1 subunit A [Escherichia coli]
Shiga toxin Stx2 subunit A [Escherichia coli]

```

Representative samples. How do I know I have the right ones? Why if they are misannotated? Answer: BLAST ALL of them!

### CloudBLAST and Docker Images to the Rescue!

- <http://ncbi.github.io/blast-cloud/>
- <https://github.com/ncbi/makeblastdb4cloud>
- Also, check out magicBLAST!



### BLAST command-line parameters

```

blastp
-db refseq-protein_v5
-query WP_123135120_123130813.fasta
-outfmt "7 qaccver saccver pident length eval evalue ssciname" # tabular output
-out Shiga_Ecoli.tab
-task blastp-fast # runs faster!
-num_threads 32 # use 8 CPUs
-evalue 0.00005 # limit to more significant matches
-taxids 562 # limit search to Escherichia coli

```

### But how do I use the BLAST Docker(s)?!

<https://github.com/ncbi/docker/blob/master/blast/README.md>

```

Install NCBI-provided BLAST databases
The $BLASTDB_S3R environment variable refers to an existing, writable directory on the local host. The following command will download the swissprot_v5 BLAST database from OCP into the $BLASTDB_S3R directory (notice the -w flag argument, which sets the working directory for that command):

docker run --rm \
  -v $BLASTDB_S3R:/blast/blastdb:rw \
  -w /blast/blastdb \
  ncbi/blast \
  update_blastdb.pl --source gpc swissprot_v5

Make and install my own BLAST databases
If you have your own sequence data in a file called sequences.fasta and want to make a BLAST database, please run the command below:

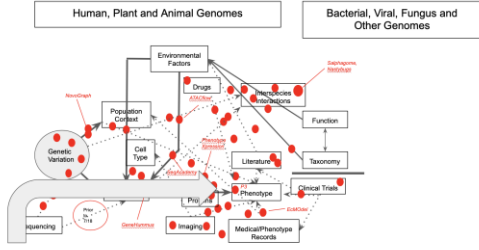
docker run --rm \
  -v $HOME/blastdb_custom:/blast/blastdb_custom:rw \
  -w $HOME/blastdb_custom \
  -v /blast/blastdb_custom \
  ncbi/blast \
  makeblastdb -in /blast/fasta/sequences.fasta -dbtype prot -out proteins -title 'My BLASTDB title'

```

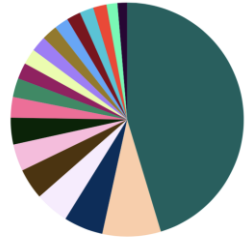




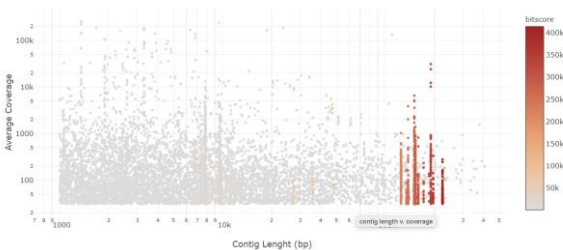
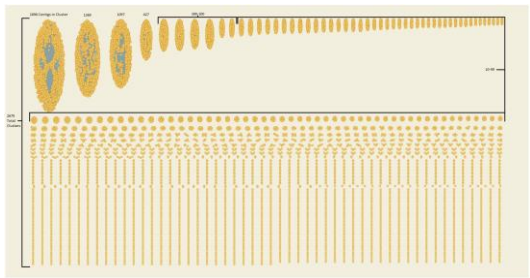
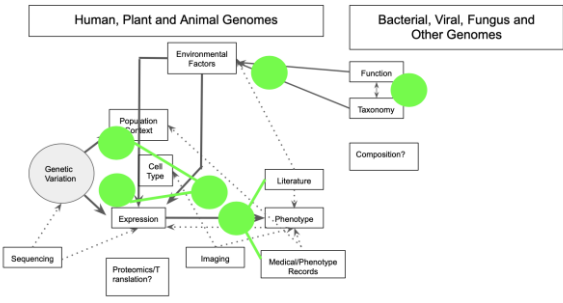
Previous Hackathons (commits 7/1/2018 - 1/1/19 only)



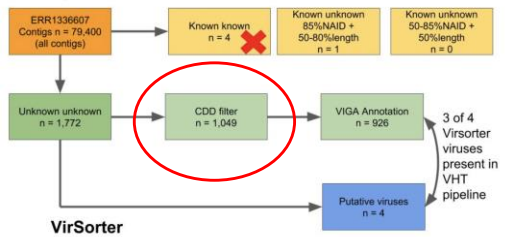
- uncultured crAssphage - 3398 contigs
- Enterobacteria phage HK630 - 608 contigs
- Enterobacteria phage P88 - 409 contigs
- Stx2-converting phage 1717 - 351 contigs
- Enterobacteria phage mEp460 - 320 contigs
- Enterobacteria phage cdtI - 279 contigs
- Enterobacteria phage phiP27 - 277 contigs
- Enterobacteria phage SFV - 214 contigs
- Escherichia virus P1 - 196 contigs
- Shigella phage SfV - 167 contigs
- Escherichia phage APCEc01 - 163 contigs
- Escherichia phage 1210 - 162 contigs
- Escherichia phage PBECO 4 - 161 contigs
- Salmonella phage 118970\_sas3 - 147 contigs
- Enterobacteria phage fIAA91-ss - 146 contigs
- Escherichia phage TL-2011b - 145 contigs
- Enterobacteria phage YYZ-2008 - 139 contigs
- Enterobacteria phage BP-4795 - 114 contigs
- Escherichia phage D108 - 101 contigs



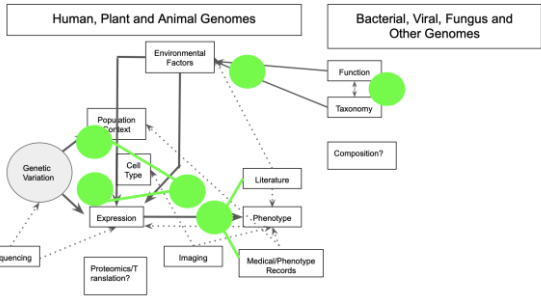
"STRIDES" Hackathons!



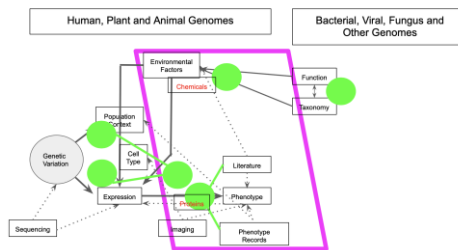
Example SRA dataset



**"STRIDES" Hackathons!**



**Future Hackathons: Pipelining and Community Tools!**



44

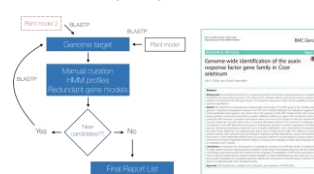
**Data Analysis Hackathons!**



- Planned Events**
- RNAseq – UNG, March 11-13 2019; Palo Alto, June 2019  
Smaller events at USF and Pittsburgh
  - Human Pan Genomics – UCSC, 25-27 March 2019; Haplotype Association, Baltimore, June 2019
  - Variant Synthetic Dataset Generation – Boston, April 15-16, 2019
  - AMR and Prokaryotic Genome Annotate-athon NIH Campus, August 2019

**GeneHumus**

Automated Gene Family Discovery!



**Creating a Community**

Upcoming Hackathons

- ADISTRIE: 10/11 HACKATHON ON 1/16/18, ABAP
- Anti-Hunting Hackathon – San Diego, ECHO – ATTY/ICM, January 9-11
- Antibiotic Discovery Research (ADRI), TBD
- Health of Food/Food/NCBI – Applications Open, APPI/HGM, February 4-6, 2019
- NIH Bethesda Campus NCBI – Team Processes New Code, February 20-22, 2019
- Genomics – History of US/NCBI – Team Processes New Code, February 26-27, 2019
- Research Triangle at UNC/NCBI – Plus IPHases workshop phase 2, March 11-13, 2019
- Perigonne and HackSpace Hackathon at UCSD/NCBI, March 26-27, 2019
- NIH/NIH/NIH/NIH Hackathon, April 10-16, 2019
- Innovation @ BCM/NIH/NCBI, October 13-14, 2019

<https://biohackathons.github.io>

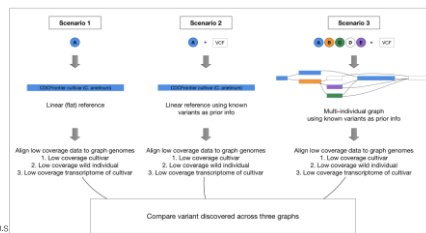
Previous Hackathons

- Saliva at UTSW and SDC/NCBI Community Panel, November 9-10
- 2018 Spring/Summer Bioinformatics Data Science, November 15-17

**MACHINE LEARNING PROJECT MANAGEMENT PRODUCT MANAGEMENT**

**HummusGraph**

Because graph genomes aren't just for humans!





Thank you.

This research was supported by the Intramural Research Program of the NIH, National Library of Medicine.

Jan Buchmann  
Chaochih Liu  
**Tom Madden**  
Moamen Elmassy  
Kim LeBlanc  
Allissa Dillman  
Olatun Awé  
Anjana Raina  
Peter Meric  
Francoise Thibaud-Nissen  
Douglas Sletta  
Michael Cruseo

**Hackathon  
Participants!**  
And visiting bioinformaticians!

Eric Cox  
Vamsi Kodali  
Brian Smith-White  
Bart Tawick  
Kim Pruitt

 U.S. National Library of Medicine  
National Center for Biotechnology Information

Watch NCBI News for updates!

<http://www.ncbi.nlm.nih.gov/news/>  
<https://www.youtube.com/user/NCBINLM>



## NCBI Genome Resources Workshop

Time	Topic
12:50 – 1:10	Submission of Genomes to GenBank <i>Karen Clark</i>
1:10 – 1:30	GEO Submissions and Usage <i>Steve Wilhite</i>
1:30 – 1:55	From Annotation to Visualization: Exploring Genes and Genomes with NCBI Tools <i>Eric Cox</i>
1:55 – 2:15	Programmatic Access to Genomic Data: E-Utilities and FTP <i>Vamsi K. Kodali</i>
2:15 – 2:35	NCBI Resources for Phyletically-Defined Next Generation Analysis in and out of the Cloud (a.k.a. Cool New Stuff!) <i>Ben Busby</i>
2:35 – 3:00	Q & A session

 U.S. National Library of Medicine  
National Center for Biotechnology Information

Visit NCBI Booth **223**

Contact us: [info@ncbi.nlm.nih.gov](mailto:info@ncbi.nlm.nih.gov)