

**BIOC8142**

Unix II – Scripting, web clients, databases  
and formats

BIOC 8142

February 13, 2013

Bill Pearson [wrp@virginia.edu](mailto:wrp@virginia.edu) 4-2818 Jordan 6-057

Goals of today's lecture:

- Creating simple bash scripts
- Survey of Bioinformatics databases (Ouellette)
  - Primary vs reference
  - Annotations and cross-references
  - Survey of file formats
- Scripts as web browsers

1

What should you do to reinforce the  
lecture material?

- Scripting the bash shell (Google "bash introduction", focus on variables, flow control)
  - [tldp.org/LDP/abs/html/](http://tldp.org/LDP/abs/html/) (concise intro)
  - Learning the Bash Shell, 3<sup>rd</sup> edition (Ch 4 and 5)  
[proquest.safaribooksonline.com/book/operating-systems-and-server-administration/unix/0596009658](http://proquest.safaribooksonline.com/book/operating-systems-and-server-administration/unix/0596009658)
- Bioinformatics databases:  
Pevsner (2004) "Bioinformatics and Functional Genomics 2<sup>nd</sup> ed" Wiley-Blackwell, Ch. 1 (on reserve, HSL)
- Web clients – `curl`, `wget` (`man curl`, `man wget`)

How will you be graded on this material:

- Homework on collab

2



## You are free to:

- Copy, share, adapt, or re-mix;
- Photograph, film, or broadcast;
- Blog, live-blog, or post video of;

## This presentation. Provided that:

- You attribute the work to its author and respect the rights and licenses associated with its components.

Slide Concept by Cameron Neylon, who has waived all copyright and related or neighbouring rights. This slide only [ccZero](https://creativecommons.org/licenses/by-sa/4.0/).  
 Social Media Icons adapted with permission from originals by Christopher Ross. Original images are available under GPL at:  
<http://www.thisismyurl.com/free-downloads/15-free-speech-bubble-icons-for-popular-websites>

## Unix II – scripting, web clients, databases

- Scripting – putting commands in a file
  - bash commands:  

```
for $file in ( *.fasta ); do ... Done
```
  - Essential for reproducibility – your electronic lab notebook
  - Automation of repetitive tasks (run blast search using 20 files)
- Web clients – `curl/wget` – allow scripting of web access
  - Download a list of protein sequences using accessions
  - Homework – (a) do a blast search with tabular output; (b) extract accessions of hits; (c) download those sequences; (d) search with them

## (bash) shell scripts

- files ending with `.sh` suffix
- shebang: `#!/bin/bash` or `#!/bin/sh`
- useful to capture (potentially long) history of UNIX commands into a reproducible analysis
  - you will always need to repeat your analysis
  - you will never remember all the necessary steps
- with some modification, your script can be made generic, and reusable for other data

## shell scripts contain commands

```
franklin: 1 $ echo $PATH # a simple command
/home/wrp/bin:/usr/local/bin:/bin:/usr/bin:./seqprg/bin

franklin: 2 $ echo_path.sh
# echo_path.sh contains "echo $PATH"
bash: ./echo_path.sh: Permission denied
# cannot execute because -rw-r--r--

franklin: 3 $ sh echo_path.sh # can execute with 'sh'
/home/wrp/bin:/usr/local/bin:/bin:/usr/bin:./seqprg/bin

franklin: 4 $ chmod +x echo_path.sh # make executable

franklin: 5 $ echo_path.sh # now it works
/home/wrp/bin:/usr/local/bin:/bin:/usr/bin:./seqprg/bin
```

## (bash) shell variables

- Your unix session has two kinds of variables, env (environment) variables, and SHELL variables, refer to them with \$NAME
  - Individual variables can be seen with 'echo'
 

```
echo $PATH
```
  - All environment variables are listed with 'env'
- You can make your own variables for a command as well:
 

```
FILES=$(ls *.aa)
echo $FILES
```

  - shell variables never have a '\$' on the left of the '=', and ALWAYS have a '\$' on the right side.
 

```
NEW_FILES=$FILES
```
- \$SHELL variables are transient; to make them permanent, use:
 

```
export PATH=$PATH:/seqprg/bin
```

## (bash) shell flow control

- `for name in [...] ; do [...] ; done`  
– do something for each item in a list
- `if [...] ; then [...] ;`  
`elif [...] ; then [...] ;`  
`else [...] fi`  
– specify behavior depending on conditions
- `';'` are only necessary when putting multiple commands on one line.  
`for ... ; do ... ; done`

## Producing new filenames

```
franklin$ for f in *.aa; # file glob (*)
> do
> n=$(basename $f .aa) # $(command) makes output
                        into a string
> new=$n.new # ${n} if no '.' or '/'
> new2="this${n}that"
> echo $f $new $new2
> done
gstml_human.aa gstml_human.new thisgstml_humanthat
sequence.aa sequence.new thissequencethat
```



## COMPUTATIONAL & COMPARATIVE GENOMICS: Understanding and Using Biological Databases

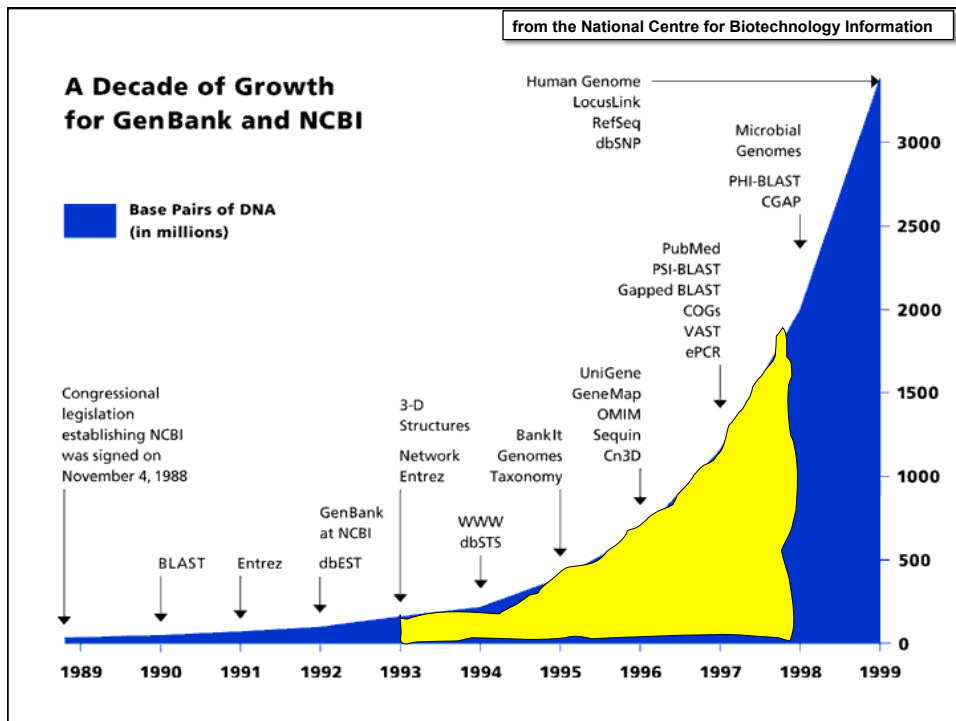
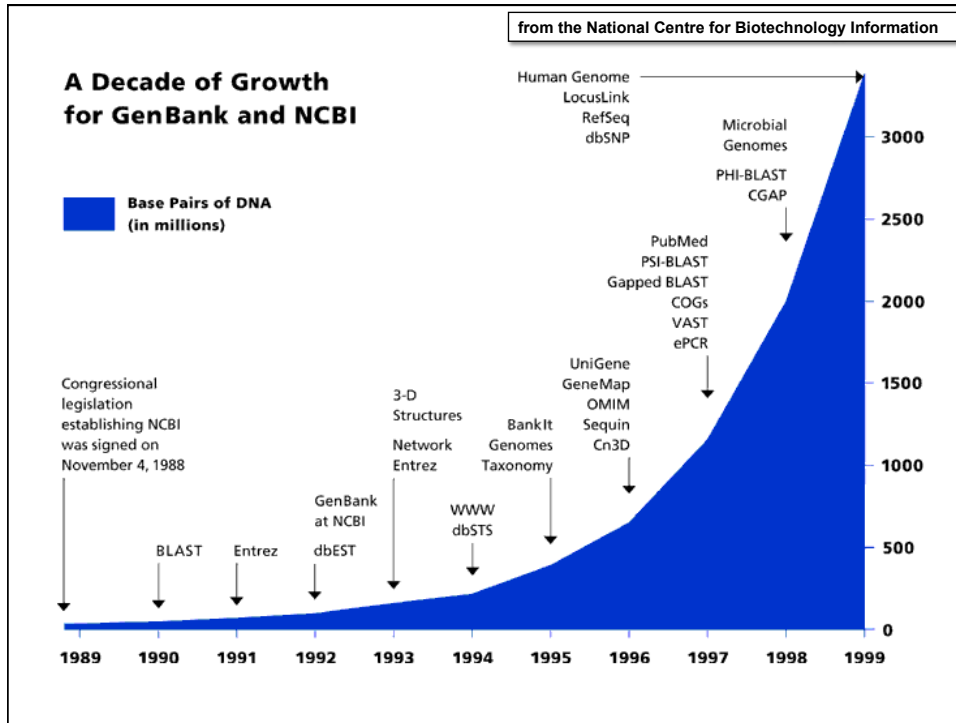
November 30<sup>th</sup>, 2012

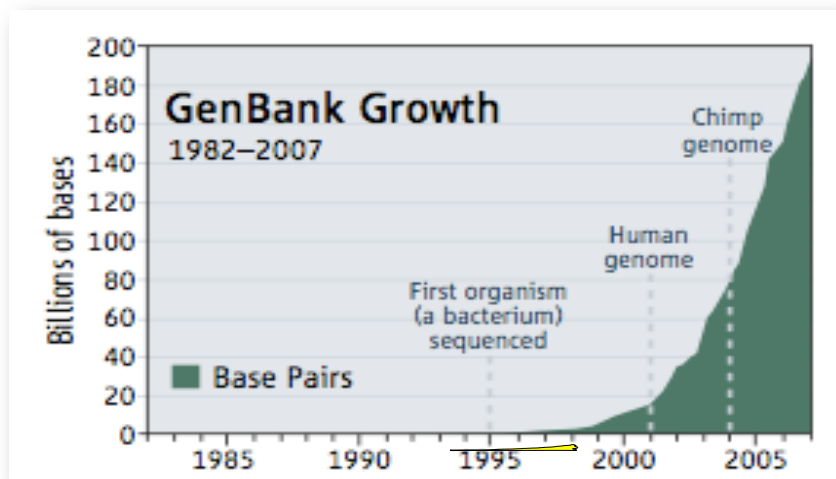
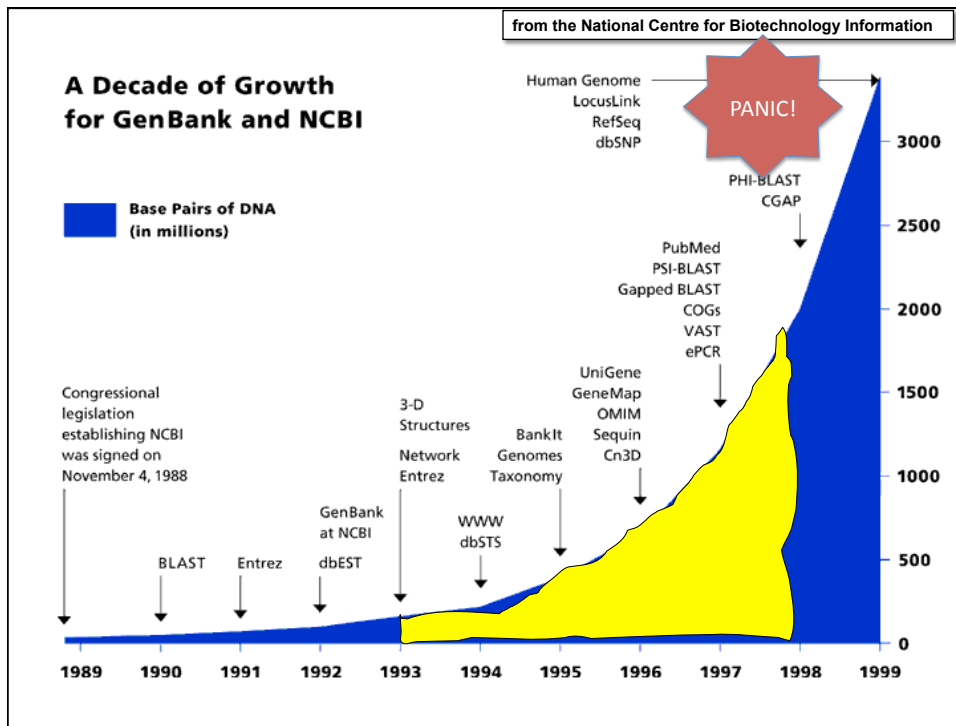
B.F. Francis Ouellette      francis@oicr.on.ca

- Associate Director, Informatics and Biocomputing, Ontario Institute for Cancer Research, Toronto, ON
- Associate Professor, Department of Cell and Systems Biology, University of Toronto, Toronto, ON.



<http://ncbi.nlm.nih.gov/>



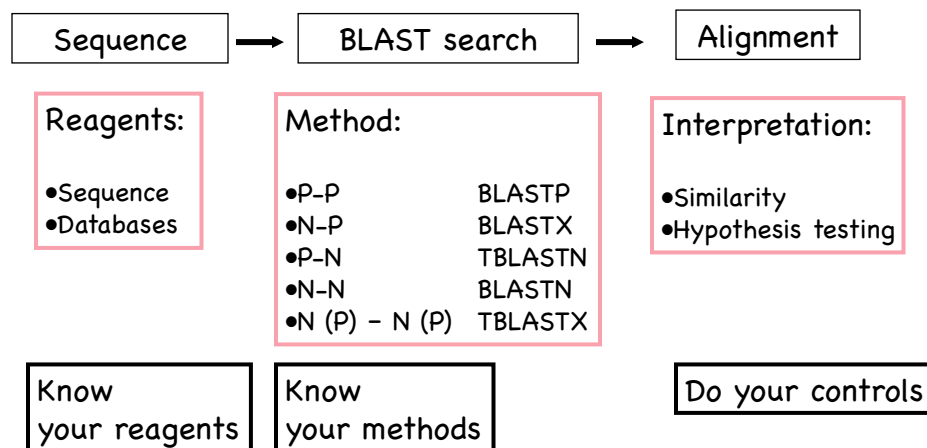




## Bioinformatics reagent: **Databases**

- Organized array of information
- Place where you put things in, and (if all is well) you should be able to get them out again.
- Resource for other databases and tools.
- Simplify the information space by specialization.
- Bonus: Allows you to make discoveries.
- Important question to ask:  
**what is the data model?**

## Bioinformatics experiments:



## Bioinformatics Citizenship: What it means, and what does it cost?

### Users must help to keep public databases correct

*Sir*— With the continued growth of the public DNA sequence databases, and the recent addition of the 11,000,000,000th nucleotide to GenBank (including DDBJ, EMBL and GenBank), it is timely to assess how we use these databases.

GenBank is the archive of all publicly available DNA, RNA and protein sequences. Upon publication a new sequence and its annotations appear in it. Investigators use GenBank in many ways, most commonly for similarity searches such as BLAST; to retrieve records; and for sequence analysis, multiple sequence alignment or pattern finding search. Errors sometimes occur in GenBank, ranging from the trivial (incorrect postal codes), to the misleading (30 nucleotides of vector left on the ends of a record), to the mission-critical (a full length mRNA without a coding sequence (CDS) annotated on it). Also very common are incomplete references that prevent researchers from linking the GenBank record to the publication that refers to it first.

Over the years some people have chosen to report these errors, but in most cases

© 2001 Macmillan Magazines Ltd

they are left unmodified. An uncorrected 'discovered' error is one of the worst possible failings in GenBank, so if you discover an error, report it to the database ([update@ncbi.nlm.nih.gov](mailto:update@ncbi.nlm.nih.gov)) and it should be rectified — although a follow-up is advised to make sure this gets done.

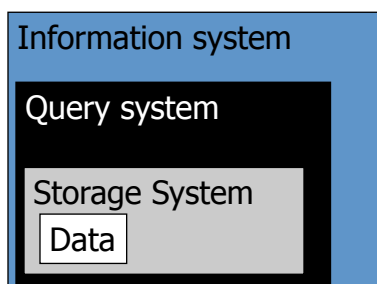
If you are a submitter, look at the record you submitted a few years ago: is it still correct? Was the citation ever updated? Take pride in the sequences that carry your name! Our ability to interpret genomes depends on all of these records being as accurate as possible. This is a task for all users of the databases.

#### Francis Ouellette

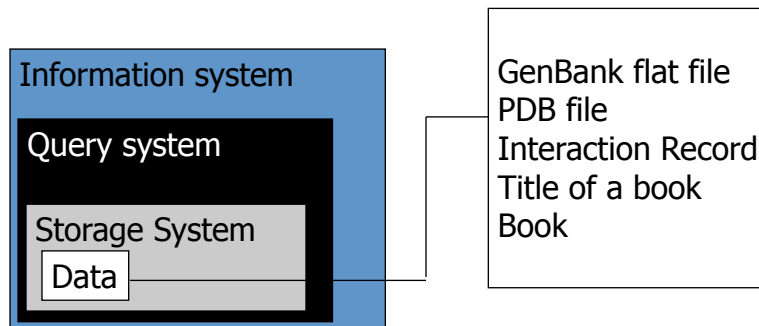
*Bioinformatics Core Facility, Centre for Molecular Medicine and Therapeutics, University of British Columbia, 950 West 28th Avenue, Vancouver, British Columbia V5Z 4H4, Canada*

Nature 409:452

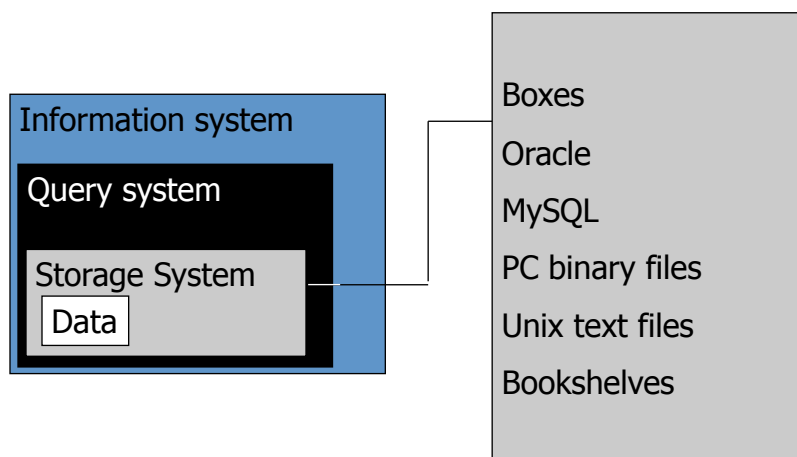
## Databases



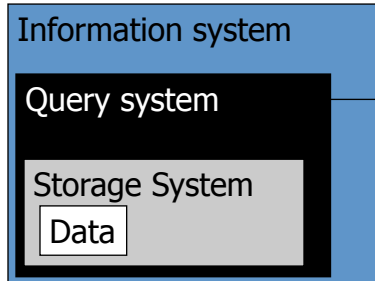
## Databases



## Databases

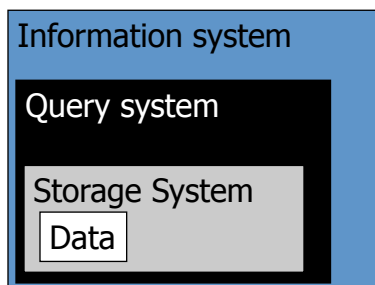


## Databases



A List you look at  
A catalogue  
indexed files  
SQL  
grep

## Databases



The library of Congress  
Google  
Entrez  
EnSEMBL  
UCSC genome browser

	Oct 15th, 2011
Nucleotide records	144,458,648
WGS records	68,330,215
Protein sequences	44,764,043
3D structures in PDB	76,973
Interactions and pathways	299,756
Human Unigene Cluster	122,727
Completed Genome projects	14,054
Different taxonomy Nodes	835,872
dbSNP records	149,211,539
Genes	9,096,508
RefSeq Genomic records	2,419,419
RefSeq RNA Records	2,679,762
RefSeq Protein Records	13,137,813
GenSAT images	106,414
GEO profiles	63,811,486
Homologene gene	128,030
PubChem compounds	28,485,889
PubMed records	21,298,405
Free PubMed records	3,466,553
PubMed Central records	2,303,372
OMIM records	21,992

	Oct 15th, 2011	Jul 17, 1999
Nucleotide records	144,458,648	4,456,822
WGS records	68,330,215	
Protein sequences	44,764,043	706,862
3D structures in PDB	76,973	9,780
Interactions and pathways	299,756	
Human Unigene Cluster	122,727	75,832
Completed Genome projects	14,054	10,870
Different taxonomy Nodes	835,872	52,889
dbSNP records	149,211,539	6,377
Genes	9,096,508	515
RefSeq Genomic records	2,419,419	
RefSeq RNA Records	2,679,762	
RefSeq Protein Records	13,137,813	
GenSAT images	106,414	
GEO profiles	63,811,486	
Homologene gene	128,030	
PubChem compounds	28,485,889	
PubMed records	21,298,405	10,372,886
Free PubMed records	3,466,553	
PubMed Central records	2,303,372	
OMIM records	21,992	10,695

	Oct 15th, 2011	Jul 17, 1999	fold differences
Nucleotide records	144,458,648	4,456,822	32
WGS records	68,330,215		
Protein sequences	44,764,043	706,862	63
3D structures in PDB	76,973	9,780	8
Interactions and pathways	299,756		
Human Unigene Cluster	122,727	75,832	2
Completed Genome projects	14,054	10,870	1
Different taxonomy Nodes	835,872	52,889	16
dbSNP records	149,211,539	6,377	23,398
Genes	9,096,508	515	17,663
RefSeq Genomic records	2,419,419		
RefSeq RNA Records	2,679,762		
RefSeq Protein Records	13,137,813		
GenSAT images	106,414		
GEO profiles	63,811,486		
Homologene gene	128,030		
PubChem compounds	28,485,889		
PubMed records	21,298,405	10,372,886	2
Free PubMed records	3,466,553		
PubMed Central records	2,303,372		
OMIM records	21,992	10,695	2

<http://www.ncbi.nlm.nih.gov/gquery/>

All [filter]

Search across databases

Result counts displayed in gray indicate one or more terms not found

22494752	PubMed: biomedical literature citations and abstracts	204913	Books: online books
2652821	PubMed Central: free, full text journal articles	22755	OMIM: online Mendelian Inheritance in Man
30184	Site Search: NCBI web and FTP sites		
75810005	Nucleotide: Core subset of nucleotide sequence records	147550	dbGaP: genotype and phenotype
74204953	EST: Expressed Sequence Tag records	6057967	UniGene: gene-oriented clusters of transcript sequences
35574863	GSS: Genome Survey Sequence records	46629	CDD: conserved protein domain database
74949497	Protein: sequence database	30181863	Clone: integrated data for clone resources
9649	Genome: whole genome sequences	545525	UniSTS: markers and mapping data
87647	Structure: three-dimensional macromolecular structures	170801	PopSet: population study data sets
1069505	Taxonomy: organisms in GenBank	78654154	GEO Profiles: expression and molecular abundance profiles
237508610	SNP: short genetic variations	930248	GEO DataSets: experimental sets of GEO data
3556617	dbVar: Genomic structural variation	6634	Epigenomics: Epigenetic maps and data sets
12286070	Gene: gene-centered information	649042	PubChem BioAssays: bioactivity screens of chemical substances
303184	SRA: Sequence Read Archive	46722996	PubChem Compound: unique small molecule chemical structures
435917	BioSystems: Pathways and systems of interacting molecules	116650376	PubChem Substances: deposited chemical substance records
133548	HomoloGene: eukaryotic homology groups	794663	Protein Clusters: a collection of related protein sequences
16627460	Probe: sequence-specific reagents	2837	OMIA: online Mendelian Inheritance in Animals
70315	BioProject: aggregated biological research project data	1768156	BioSample: biological material descriptions

Feb, 2013

# Formats

- DNA sequence (GenBank Flat Files)
- Protein Sequences
- Other formats to know about
  - FASTA
  - GFF3
  - XML

# GenBank Flat File (GBFF)

```

LOCUS       JM675711                1704 bp    mRNA    linear    PLM 01-000-2011
DEFINITION  Prunus salicina mitogen-activated protein kinase 1 mRNA, complete
            cds.
ACCESSION   JM675711
VERSION    JM675711.1  GI:136528442
KEYWORDS   .
SOURCE     Prunus salicina
            ORGANISM  Prunus salicina
                        Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta;
                        Spermatophyta; Magnoliophyta; eudicotyledons; core eudicotyledons;
                        rosids; Malvales; Rosales; Malvaceae; Myricales; Prunus.
REFERENCE   1  (bases 1 to 1704)
AUTHORS    Jiang,C., Pan,D. and Chen,G.
TITLE      Construction and Analysis of a Normalized Full-length cDNA Library
            of Sweet Prunus salicina
JOURNAL    Unpublished
REFERENCE   2  (bases 1 to 1704)
AUTHORS    Jiang,C., Pan,D. and Chen,G.
TITLE      Direct Subcloning
JOURNAL    Horticulture (18-400-2011) Horticulture, Fujian Agriculture and
            Forestry University, Fuzhou, Fu Jian 350002, China
            Location:Qualifiers
FEATURES             source          1..1704
                     /organism="Prunus salicina"
                     /mol_type="mRNA"
                     /architecture="1"
                     /db_xref="taxon:8122"
                     138..1249
                     /note="WAK1"
                     /function="kinase"
                     /product="Mitogen-activated protein kinase 1"
                     /protein_id="AJ282831.1"
                     /db_xref="GI:136528442"
                     /transcript_id="JM675711.1"
ORIGIN
1 ccaattacggp ctggtatcagp gggatcattt ttgctcctca tatcctccct tccaaatg
61 ttgactcttt agaacacaa tccatctggt tcttgcaac atttggatt ggcactcgt
121 ctactctggp gattctgatt tigtctcagp aggtcacaa atctatccg aggtcctp
181 ccacattgpc cgtatcgtc agatcaattt gtagcaaac tcttttggp ttctatgaa
241 gactctcttc tccaaagp cgttctcagp agttctcct gttacttct gttcctctp
301 gactctcttc tccaaagp cgttctcagp agttctcct gttacttct gttcctctp
361 actctctctc tccaaagp cgttctcagp agttctcct gttacttct gttcctctp
421 tttatctctc tccaaagp cgttctcagp agttctcct gttacttct gttcctctp
481 actctctctc tccaaagp cgttctcagp agttctcct gttacttct gttcctctp
541 gactctcttc tccaaagp cgttctcagp agttctcct gttacttct gttcctctp
601 ttctctcttc tccaaagp cgttctcagp agttctcct gttacttct gttcctctp
661 gactctcttc tccaaagp cgttctcagp agttctcct gttacttct gttcctctp
721 ttgactcttt gttactctt ttgactcagp agtgactc tctctcaat tttgactc
781 cttctctcag agttactct ttgactcagp agtgactc tctctcaat tttgactc
841 gttactcttc cttctctcag agttactct ttgactcagp agtgactc tctctcaat
901 tttgactcag agttactct ttgactcagp agtgactc tctctcaat tttgactc
961 aaactactct cttctctcag agttactct ttgactcagp agtgactc tctctcaat
1021 gttactcttc cttctctcag agttactct ttgactcagp agtgactc tctctcaat
1081 ttgactcttc cttctctcag agttactct ttgactcagp agtgactc tctctcaat
1141 gttactcttc cttctctcag agttactct ttgactcagp agtgactc tctctcaat
1201 ttgactcttc cttctctcag agttactct ttgactcagp agtgactc tctctcaat
1261 ttgactcttc cttctctcag agttactct ttgactcagp agtgactc tctctcaat
1321 ttgactcttc cttctctcag agttactct ttgactcagp agtgactc tctctcaat
1381 ttgactcttc cttctctcag agttactct ttgactcagp agtgactc tctctcaat
1441 ttgactcttc cttctctcag agttactct ttgactcagp agtgactc tctctcaat
1501 gttactcttc cttctctcag agttactct ttgactcagp agtgactc tctctcaat
1561 ttgactcttc cttctctcag agttactct ttgactcagp agtgactc tctctcaat
1621 gttactcttc cttctctcag agttactct ttgactcagp agtgactc tctctcaat
1681 aaaaaaaaaa aaaaaaaaaa
    
```

**Header**

- Title
- Taxonomy
- Citation

**Features (AA seq)**

**DNA Sequence**

## FASTA

```
>
MSEYQPSL FALNPMGF SPLDGSKSTNENVSASTSTAKPMVGLIFDKFIKTEEDI
IKQDTPSNLDFDFALPQTATAPDAKTVLPPELDDAVVESFFSSSTDSTPMFEYEN
LEDNSKEWTS LFDNDIPVTTDDVSLADKAIESTEEVSLVPSNLEVSTTSFLPTPV
EDAKLTQTRKVKKPN SVVKKSHHVGKDDERLDH LGVVAYNRKQRSIPLSPIVPES
SDPAALKRARNT EAARRSRARKLQRMKQLEDKVEELLSKNYHLENEVARLKKLVGE
R
```

## Databases

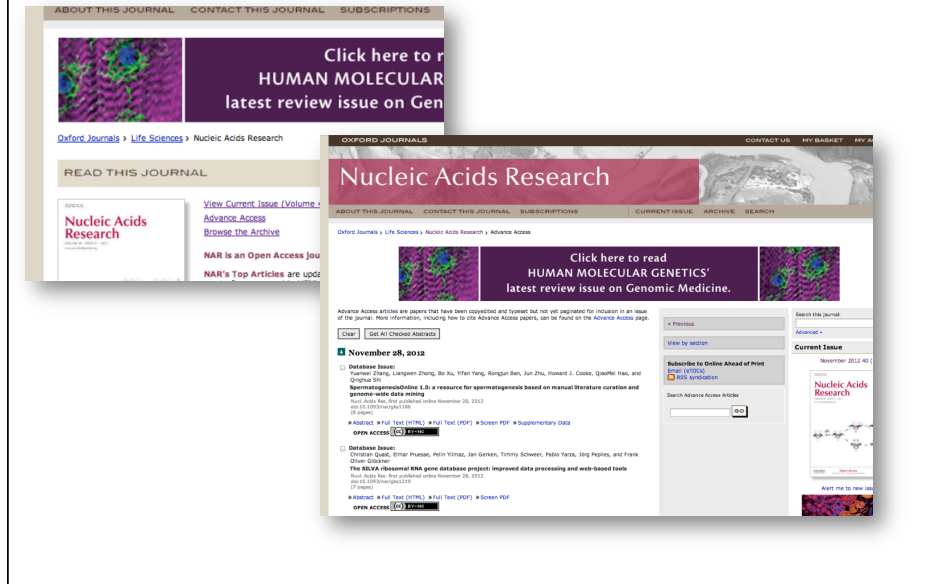
- Primary (archival)
  - GenBank/EMBL/DDBJ
  - UniProt
  - PDB
  - Medline (PubMed)
  - Intact
- Secondary (curated)
  - RefSeq
  - Taxon
  - UniProt
  - OMIM
  - SGD
  - Biosamples/  
Bioprojects



[http://nar.oxfordjournals.org/content/39/suppl\\_1](http://nar.oxfordjournals.org/content/39/suppl_1) January 2011

[http://nar.oxfordjournals.org/content/38/suppl\\_1](http://nar.oxfordjournals.org/content/38/suppl_1) January 2010

<http://nar.oxfordjournals.org/>



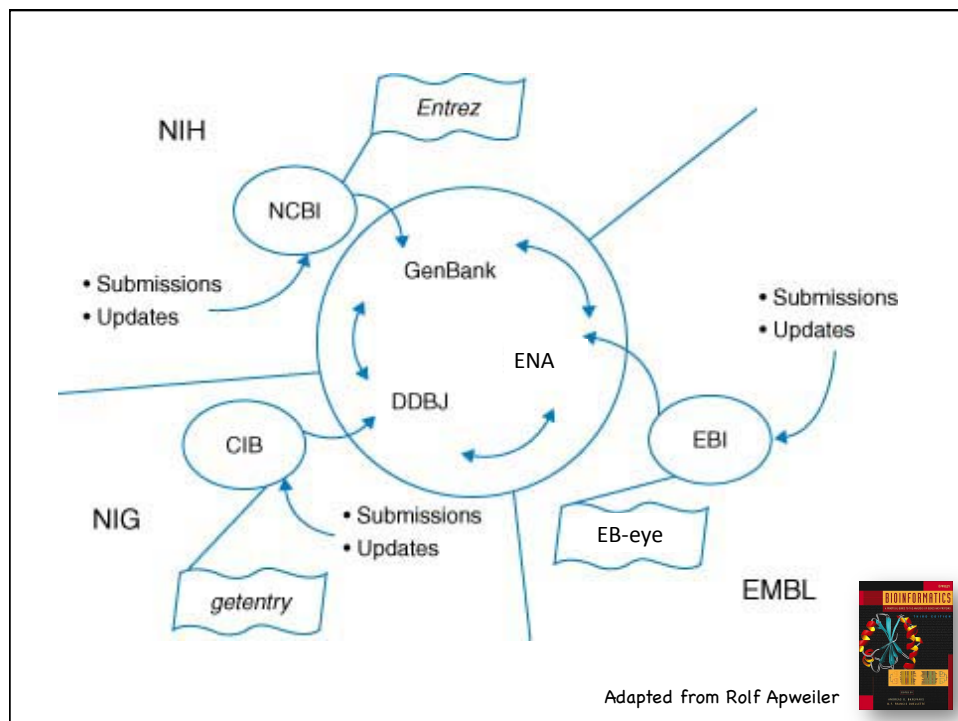
## Sequence Databases

- Primary DNA (archive)
  - DDBJ/ENA/GenBank
- Primary protein (curated/automation)
  - UniProtKB
- Curated Databases (lots of human labour)
  - RefSeq (Genomic, mRNA and protein)
  - UniProtKB/SwissProt and neXtprot

## What is GenBank?

GenBank is the NIH genetic sequence database of all publicly available DNA and derived protein sequences, with annotations describing the biological information these records contain.

<http://www.ncbi.nlm.nih.gov/genbank/>  
 Benson et al., Nucleic Acids Res. 2012 (out in 2011)  
<http://www.ncbi.nlm.nih.gov/pubmed/22144687>



## Types of files in GenBank

- From one-gene investigators
  - Often a very well annotated cDNA
  - A genomic segment from an new invertebrate
  - A mitochondria or virus
- From population/phylogenetic analysis
  - rRNA amplicon from environmental sampling
- From Genome Centers:
  - Gene expression:
    - Expressed Sequence Tags
    - Full Length Insert cDNA
    - TSA
  - Genome sequencing projects
    - HTG
    - CON

## Organismal Divisions

		Used in which database?
BCT	Bacterial	DDBJ - GenBank
FUN	Fungal	ENA
HUM	Homo sapiens	DDBJ - ENA
INV	Invertebrate	all
MAM	Other mammalian	all
ORG	Organelle	ENA
PHG	Phage	all
PLN	Plant	all (not same data in all)
PRI	Primate (also see HUM)	all (not same data in all)
PRO	Prokaryotic	ENA
ROD	Rodent	all
SYN	Synthetic and chimeric	all
VRL	Viral	all
VRT	Other vertebrate	all

## Functional Divisions

**PAT** Patent  
**EST** Expressed Sequence Tags  
**TSA** Transcriptome Shotgun Assembly  
**STS** Sequence Tagged Site  
**GSS** Genome Survey Sequence  
**HTG** High Throughput Genome (unfinished)  
**HTC** High throughput cDNA (unfinished)  
**CON** Contig assembly instructions  
**ENV** Environmental sampling methods

Organismal divisions:

<b>BCT</b>	<b>FUN</b>	<b>INV</b>	<b>MAM</b>	<b>PHG</b>	<b>PLN</b>
<b>PRI</b>	<b>ROD</b>	<b>SYN</b>	<b>VRL</b>	<b>VRT</b>	

## Guiding Principals

In GenBank, records are grouped for various reasons: understand this is key to using and fully taking advantage of this database.

## Identifiers

- You need identifiers which are stable through time
- Need identifiers which will always refer to specific sequences
- Need these identifiers to track history of **sequence** updates
- Also need feature and annotation identifiers (need to track important things)
  - Genes
  - Transcripts
  - Proteins
  - ((( Phenotype )))

## LOCUS, Accession, NID and protein\_id

**LOCUS:** Unique string of 10 letters and numbers in the database. Not maintained amongst databases, and is therefore a poor sequence identifier.

**ACCESSION:** A unique identifier to that record, citable entity; does not change when record is updated. A good record identifier, ideal for citation in publication.

**VERSION:** ID system where the accession and version play the same function as the accession and gi number.

**Nucleotide gi:** Geninfo identifier (gi), a unique integer which will change every time the sequence changes.

**Protein gi:** Geninfo identifier (gi), a unique integer which will change every time the sequence changes.

**protein\_id:** Identifier which has the same structure and function as the nucleotide Accession.version numbers, but slightly different format.

## LOCUS, Accession, gi and PID

```

LOCUS      HSU40282      1789 bp      mRNA      PRI      21-MAY-1998
DEFINITION Homo sapiens integrin-linked kinase (ILK) mRNA, complete cds.
ACCESSION  U40282
VERSION    U40282.1      GI:3150001
    
```

```

LOCUS: HSU40282
ACCESSION: U40282
VERSION: U40282.1
GI: 3150001
Protein gi: 3150002
protein_id: AAC16892.1
    
```



```

CDS          157..1515
              /gene="ILK"
              /note="protein serine/threonine kinase"
              /codon_start=1
              /product="integrin-linked kinase"
              /protein_id="AAC16892.1"
              /db_xref="GI:3150002"
    
```

Lecture 1.3

NCBI Resources How To

Nucleotide

Format

- Summary
- GenBank
- GenBank (full)
- FASTA
- FASTA (text)
- Graphics
- ASN.1
- Revision History
- Accession List
- GI List

Apply

chromosome 17, clone RP11-726012, complete sequence

106210 bp DNA linear PRI 03-MAR-2001  
 chromosome 17, clone RP11-726012, complete sequence.

ORGANISM [Homo sapiens](#) (human)

Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominidae; Homo.

REFERENCE 1 (bases 1 to 106210)  
 AUTHORS Birren,B., Linton,L., Nusbaum,C. and Lander,E.  
 TITLE Homo sapiens chromosome 17, clone RP11-726012  
 JOURNAL Unpublished

REFERENCE 2 (bases 1 to 106210)  
 AUTHORS Birren,B., Fasman,K., Linton,L., Nusbaum,C., Lander,E., Allen,N., Anderson,M., Baker,J., Baldwin,J., Barna,N., Beckerly,R., Benn,J., Boutwell,C., Brown,A., Castle,A., Cerny,J., Colangelo,M., Collins,S., Collymore,A., Cooke,P., Corliss,D., Depayre,E., Devon,K., Dewar,K., Donelan,L., Ferreira,P., FitzHugh,W., Forrest,C., Funke,R., Gage,D., Gardyna,S., Geraigery,K., Grant,G., Hayes,B., Heaford,A., Herena,L., Horton,L., Howland,L.C.,

Display Settings:  Revision History

Show difference between I and II as GenBank/GenPept

Gi	Version	Update Date
7670171	6	Apr 29, 2000 05:07 AM
13194375	7	Mar 3, 2001 05:06 AM

LOCUS AC005517 106330 bp DNA linear PRI 29-APR-2000  
 LOCUS AC005517 106210 bp DNA linear PRI 03-MAR-2001  
 DEFINITION Homo sapiens chromosome 17, clone RP11-726012, complete sequence.  
 ACCESSION AC005517  
 VERSION AC005517.6 GI:7670171  
 VERSION AC005517.7 GI:13194375  
 KEYWORDS HTG.  
 SOURCE Homo sapiens (human)  
 ORGANISM Homo sapiens  
 Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;  
 Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo.  
 REFERENCE 1 (bases 1 to 106330)  
 Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;  
 Catarrhini; Hominidae; Homo.

[Open Comparison in separate window](#)

[Homo sapiens chromosome 17, clone RP11-726012, complete sequence](#)  
 106,210 bp linear DNA  
 Accession: AC005517.7 GI: 13194375  
 Current status: live

I	II	Version	Gi	Update Date
<input type="radio"/>	<input type="radio"/>	7	13194375	Nov 30, 2009 04:42 PM
<input type="radio"/>	<input checked="" type="radio"/>	7	13194375	Mar 3, 2001 05:06 AM
<input checked="" type="radio"/>	<input type="radio"/>	6	7670171	Apr 29, 2000 05:07 AM
<input type="radio"/>	<input type="radio"/>	5	7574846	Apr 15, 2000 05:15 AM

## Accession number “space”

- GenBank:
  - 1+5 (L12345, U00001)
  - 2+6 (AF000001, AC000003)
- WGS (Not distributed with GenBank)
  - 4+2+6 (AAAA01000001, AAAD01000001)
- Protein:
  - 1+5 or 3+5
- All have “accession.version”



## Secondary Accession Numbers

- When you 'retire' accession numbers, these often are put in the secondary accession number space. (e.g GenBank Accession number L05146)
- With the removal of sequence length limits, GenBank will now allow continuous ranges of secondary accessions.
- As of GenBank Release 146.0 (February 2005), it is legal to represent continuous ranges of secondary accessions by a start accession, a dash character, and an end accession. (e.g. for the *E. coli* genome)

**ACCESSION U00096 AE000111-AE000510**

```

LOCUS      U00096                4639675 bp    DNA    circular BCT 01-SEP-2011
DEFINITION Escherichia coli str. K-12 substr. MG1655, complete genome.
ACCESSION  U00096 AE000111-AE000510
VERSION    U00096.2  GI:48994873
DBLINK     Project: 225
KEYWORDS   .
SOURCE     Escherichia coli str. K-12 substr. MG1655
  ORGANISM Escherichia coli str. K-12 substr. MG1655
            Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacteriales;
            Enterobacteriaceae; Escherichia.
REFERENCE  1 (bases 1 to 4639675)
AUTHORS    Blattner,F.R., Plunkett,G. III, Bloch,C.A., Perna,N.T., Burland,V.,
            Riley,M., Collado-Vides,J., Glasner,J.D., Rode,C.K., Mayhew,G.F.,
            Gregor,J., Davis,N.W., Kirkpatrick,H.A., Goeden,M.A., Rose,D.J.,
            Mau,B. and Shao,Y.
TITLE      The complete genome sequence of Escherichia coli K-12
JOURNAL    Science 277 (5331), 1453-1474 (1997)
PUBMED    9278503
  
```

## WGS: Whole Genome Shotgun (Not in GenBank release, not shared with ENA/DDBJ)

- Contigs from ongoing Whole Genome Shotgun sequencing projects
- The nucleotides from WGS projects go into the BLAST 'wgs' database, whereas the proteins go into the BLAST nr database.
- More info, and how to submit to this division:  
<http://www.ncbi.nlm.nih.gov/Genbank/wgs.html>
- Accession format is 4+2+6

<http://www.ncbi.nlm.nih.gov/projects/WGS/WGSprojectlist.cgi>

NCBI WGS Sequencing Projects

PubMed Entrez BLAST OMM Books Taxonomy Structure

WGS Home  
MapViewer

GenBank WGS Projects:

The accession number prefix connects to the summary page of the master record of the WGS project. The organism name connects to the Taxonomy Browser, where links to other resources are present. The each project is listed. If the WGS contigs are assembled into scaffolds or chromosomes, then the number of GenBank CON records representing these scaffolds or chromosomes is listed. The WGS contig records are available as links from the WGS master record. Go to [GenomeProject](#) or [MapViewer](#) for additional resources.

Prefix	GPID	Organism	# contigs	# CONs	Annotation	Complete
AAAA	261	<i>Oryza sativa</i> (indica cultivar-group)	50,231	3,095	Y-s	---
AAAB	1438	<i>Anopheles gambiae</i> str. PEST	69,724	5	Y-c	---
AAAC	299	<i>Bacillus anthracis</i> A2012	1	-	-	---
AAAD	11785	<i>Mus musculus</i>	20	---	-	CM000224
AAAE	56	<i>Rhodobacter sphaeroides</i>	159	---	Y	CP000143-CP000147_DQ232586-DQ232587
AAAF	57	<i>Rhodospirillum rubrum</i>	15	---	-	BX571963
AAAG	58	<i>Rhodospirillum rubrum</i>	10	---	Y	CP000230
AAAH	59	<i>Chloroflexus aurantiacus</i> J-10-fl	77	---	Y	CP000999
AAAI	259	<i>Ralstonia metallidurans</i> CH34	29	---	Y	CP000352-CP000355
AAAJ	254	<i>Burkholderia xenovorans</i> LB400	58	---	Y	CP000270-CP000272
AAAK	71	<i>Enterococcus faecium</i> DO	163	---	Y-c	---
AAAL	269	<i>Xylella fastidiosa</i> Dixon	32	---	Y-c	---
AAAM	261	<i>Xylella fastidiosa</i> A-m-1	219	---	Y-c	---
AAAN	262	<i>Magnetococcus</i> sp. MC-1	108	---	Y	CP000471
AAAO	84	<i>Lactobacillus gasserii</i> ATCC 33323	46	---	Y	CP000413
AAAP	6	<i>Magnetospirillum magnetotacticum</i> MS-1	3,880	---	-	---
AAAQ	94	<i>Thermobifida fusca</i>	19	---	Y	CP000088
AAAR	103	<i>Methanosarcina barkeri</i> str. fusaro	69	---	Y	CP000098-CP000099
AAAS	177	<i>Geobacter metallireducens</i> GS-15	11	---	Y	CP000148
AAAT	12	<i>Pseudomonas fluorescens</i> Pf0-1	24	---	Y	CP000094

## WGS record (not in GenBank)

```

LOCUS      AAC002000001      55046 bp      DNA      linear      PLN 03-OCT-2006
DEFINITION Cryptococcus neoformans var. grubii H99 contg1.1, whole genome
            shotgun sequence.
ACCESSION  AAC002000001 AAC002000000
VERSION    AAC002000001.1 GI:115505304
KEYWORDS   WGS.
SOURCE     Cryptococcus neoformans var. grubii H99
ORGANISM   Cryptococcus neoformans var. grubii H99
            Eukaryota; Fungi; Basidiomycota; Hymenomyces;
            Heterobasidiomycetes; Tremellomycetidae; Tremellales; Tremellaceae;
            Filobasidiella.
  
```

```

REFERENCE  3 (bases 1 to 55046)
AUTHORS    Birren,B., Lander,E., Galagan,J., Nusbaum,C., Devon,K., Ma,L.-J.,
            Jaffe,D., Butler,J., Alvarez,P., Gnerre,S., Grabherr,M., Kleber,M.,
            Mauceli,E., Brockman,W., Rounsley,S., Young,S., LaButti,K.,
            Pushparaj,V., DeCaprio,D., Crawford,M., Koehrsen,M., Engels,R.,
            Montgomery,P., Pearson,M., Howarth,C., Larson,L., Luoma,S.,
            White,J., Kodira,C., O'Leary,S., Yandava,C., Zeng,Q., Alvarado,L.,
            Dietrich,F., Heitman,J. and Kronstad,J.
CONSRM     The Broad Institute Genome Sequencing Platform
TITLE      Direct Submission
JOURNAL    Submitted (26-SEP-2006) Broad Institute of MIT and Harvard, 320
            Charles Street, Cambridge, MA 02141, USA
FEATURES   Location/Qualifiers
            source                1..55046
                                     /organism="Cryptococcus neoformans var. grubii H99"
                                     /mol_type="genomic DNA"
                                     /strain="H99"
                                     /variety="grubii"
                                     /serotype="A"
                                     /db_xref="taxon:235443"
ORIGIN
1  gtcaagagat aaacgccat gctcaccata cacagactat gatccgtcag cattgtcttt
61  tgtcacgagg ataacaatgg acgttacgc aaactcctct ccatccacac cctatattct
121 ccatcgccag cttttccatc ttgccacct ctctgccat cctgctttcg acaccgtctg
181 ctcttctggc gtccttccgc actctggtat ctttccctt cggcgtggtt gactatcgac
241 aacagtcagc gagagtaaca ctcagagcaa cggcgacaaa cgaatcagca actcacagtc
301 tgtacagcgt actatatctc cagcatccac cgccctgccc tccgcgtgtg accataacctc
361 ctccctgctg tctcatccac acgcaaccaa acgttgcttg cttgcttgcg tcaagagtag
  
```

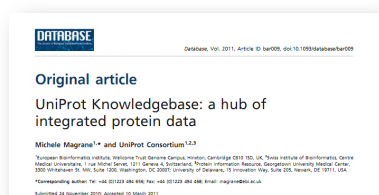
## Sequences NOT in GenBank

- WGS: whole genome shotgun
- TPA: third party annotations
- SNPs
- SAGE tags (serial analysis of gene expression)
- RefSeq (Genomic, mRNA, or protein)
- Consensus sequences

## What is UniProtKB?

UniProt is a protein sequence database that is the result of a merge from SWISS-PROT and PIR and is funded by the NIH, EMBL and the Swiss Gov't. It is the main distributed, annotated, and curated protein sequence database. Data in UniProt is derived from coding sequence annotations in ENA (GenBank/ENA/DDBJ) nucleic acid sequence data, and from sequences in PIR and SP. UniProt is a Flat-File database just like GenBank or ENA/EMBL

- <http://www.uniprot.org/>
- <http://database.oxfordjournals.org/content/2011/bar009.long>



# Uniprot Curation

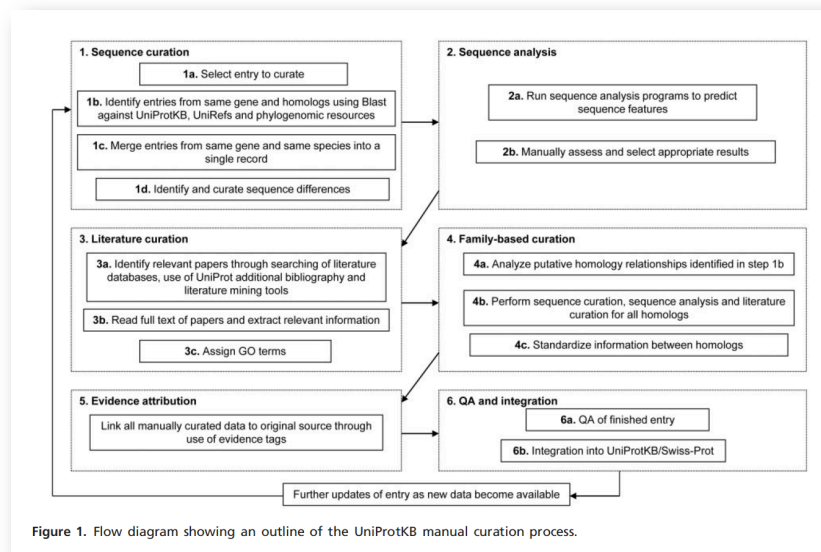


Figure 1. Flow diagram showing an outline of the UniProtKB manual curation process.

```

ID CYS3_YEAST Reviewed; 394 AA.
AC F31373; D6VFK6;
DT 01-JUL-1993, integrated into UniProtKB/Swiss-Prot.
DT 23-JAN-2007, sequence version 2.
DT 19-OCT-2011, entry version 115.
DE RecName: Full=Cystathionine gamma-lyase;
DE EC=4.4.1.1;
DE AltName: Full=Gamma-cystathionase;
DE AltName: Full=Sulfur transfer protein 1;
GN Name=CYS3; Synonyms=CY11, STR1; OrderedLocusNames=YAL012W;
GN ORFNames=FUN35;
OS Saccharomyces cerevisiae (strain ATCC 204508 / S288c) (Baker's yeast).
OC Eukaryota; Fungi; Dikarya; Ascomycota; Saccharomycotina;
OC Saccharomycetes; Saccharomycetales; Saccharomycetaceae; Saccharomyces.
OX NCBI_TaxID=559292;
RX [1]
RP NUCLEOTIDE SEQUENCE [GENOMIC DNA], AND PARTIAL PROTEIN SEQUENCE.
RX MEDLINE=92250430; PubMed=1577698;
RA Ono B., Tanaka K., Naito K., Heike C., Shinoda S., Yamamoto S.,
RA Ohmori S., Oshima T., Toh-e A.;
RT "Cloning and characterization of the CYS3 (CY11) gene of Saccharomyces
RT cerevisiae.";
RL J. Bacteriol. 174:3339-3347(1992).
RX [2]
RP NUCLEOTIDE SEQUENCE [GENOMIC DNA], AND CHARACTERIZATION.
RC STRAIN=DBY939;
RX MEDLINE=93328685; PubMed=8335636;
RA Yamagata S., D'Andrea R.J., Fujisaki S., Isaji M., Nakamura K.;
RT "Cloning and bacterial expression of the CYS3 gene encoding
RT cystathionine gamma-lyase of Saccharomyces cerevisiae and the
RT physicochemical and enzymatic properties of the protein.";
RL J. Bacteriol. 175:4800-4808(1993).
RX [3]
RP NUCLEOTIDE SEQUENCE [GENOMIC DNA].
RC STRAIN=ATCC 204511 / S288c / AB972;
RX MEDLINE=93289814; PubMed=8511966; DOI=10.1002/yea.320090406;
RA Barton A.B., Kaback D.B., Clark M.W., Keng T., Ouellette B.F.F.,
RA Storms R.K., Zeng B., Zhong W.W., Fortin N., Delaney S., Bussey H.;
RT "Physical localization of yeast CYS3, a gene whose product resembles
RT the rat gamma-cystathionase and Escherichia coli cystathionine gamma-
RT synthase enzymes.";
RL Yeast 9:363-369(1993).
  
```

```

FT INIT_MET 1 1 Removed.
FT CHAIN 2 394 Cystathionine gamma-lyase.
FT BINDING 52 52 /PTD-PRO_000114754.
FT BINDING 104 104 Substrate (By similarity).
FT BINDING 109 109 Substrate (By similarity).
FT BINDING 334 334 Substrate (By similarity).
FT MOD_RES 2 2 Phosphothreonine.
FT MOD_RES 39 39 Phosphoserine.
FT MOD_RES 40 40 Phosphoserine.
FT MOD_RES 204 204 N6-(pyridoxal phosphate)lysine (By
similarity).
FT MOD_RES 362 362 Phosphoserine.
FT HELIX 9 15
FT STRAND 38 47
FT TURN 51 53
FT HELIX 56 68
FT STRAND 72 78
FT HELIX 80 89
FT STRAND 96 102
FT HELIX 105 113
FT TURN 116 118
FT STRAND 123 127
FT HELIX 128 135
FT STRAND 138 145
FT TURN 150 152
FT HELIX 158 168
FT TURN 169 173
FT STRAND 175 179
FT TURN 181 183
FT HELIX 184 187
FT HELIX 190 193
FT STRAND 196 201
FT TURN 202 207
FT STRAND 215 220
FT HELIX 222 235
FT HELIX 241 251
FT HELIX 254 272
FT TURN 276 278
FT STRAND 279 283
FT HELIX 293 299
FT HELIX 301 303
FT STRAND 307 315
FT HELIX 317 326
FT STRAND 328 332
FT STRAND 342 344
FT TURN 346 350
FT TURN 356 360
FT STRAND 361 363
FT STRAND 368 372
FT HELIX 378 393
SQ SEQUENCE 394 AA; 42542 MW; 861A5AA7557697FC CRC64;
MTLQSDKFA TKAHAGEHV DVHGSVIEPI LSTTFKQSS PANPIGTIYEY SRSQNPNREN
LERAVALENA QYGLAFSSGS ATTTATLQSL PQGSHAVSIG DVGVGTHRYF TKVANAHVGE
TSFTNDLLND LPQLIKENTK LWIETPTNP TLKVTDIQKV ADLKKHAAQ QDVLVVDNT
FLSPYISNPL NFGADIVVHS ATKYINGHSD VLVGLVATNN KPLYERLQFL QNAIGAIKPS
PDAWLTHRGL KTLHLRVRQA ALSANKIAEF VLAADKENVVA VNYVGLKTHP NYDVVLKQHR
DALGGGIMSF RIKGGAEAAAS KPASSTRLPF LAESLGGIEE LLEVAVMTH GGIPEKARE
AGVFDDLVRI SVGIEDTDDL LEDIKALKQK ATN
//

```

**UniprotKB/SwissProt**

```

ID CYS3_YEAST STANDARD; PRT; 393 AA.
AC P31373;
DT 01-JUL-1993 (REL. 26, CREATED)
DE CYSTATHIONINE GAMMA-LYASE (EC 4.4.1.1) (GAMMA-CYSTATHIONASE).
GN CYS3 OR CY11 OR STR1 OR YAL012W OR FUN35.
OS TAXONOMY
OC SACCHAROMYCETACEAE; SACCHAROMYCES.

RX CITATION
CC -!- CATALYTIC ACTIVITY: L-CYSTATHIONINE + H(2)O = L-CYSTEINE +
CC NH(3) + 2-OXOBUTANOATE.
CC -!- COFACTOR: PYRIDOXAL PHOSPHATE.
CC -!- PATHWAY: FINAL STEP IN THE TRANS-SULFURATION PATHWAY SYNTHESIZING
CC L-CYSTEINE FROM L-METHIONINE.
CC -!- SUBUNIT: HOMOTETRAMER.
CC -!- SUBCELLULAR LOCATION: CYTOPLASMIC.
CC -!- SIMILARITY: BELONGS TO THE TRANS-SULFURATION ENZYMES FAMILY.
CC -----
CC DISCLAMOR
CC -----
DR DATABASE cross-reference
KW CYSTEINE BIOSYNTHESIS; LYASE; PYRIDOXAL PHOSPHATE.
FT INIT_MET 0 0
FT BINDING 203 203 PYRIDOXAL PHOSPHATE (BY SIMILARITY).
SQ SEQUENCE 393 AA; 42411 MW; 55BA2771 CRC32;
TLQSDKFA TKAHAGEHVD VHGVSVIEPI LSTTFKQSS ANPIGTIYEYS RSQNPENREN
ERAVALENA QYGLAFSSGS ATTTATLQSL PQGSHAVSIG DVGVGTHRYF TKVANAHVGE
TSFTNDLLND LPQLIKENTK LWIETPTNP TLKVTDIQKV ADLKKHAAQ QDVLVVDNT
FLSPYISNPL NFGADIVVHS ATKYINGHSD VLVGLVATNN KPLYERLQFL QNAIGAIKPS
PDAWLTHRGL KTLHLRVRQA ALSANKIAEF VLAADKENVVA VNYVGLKTHP NYDVVLKQHR
DALGGGIMSF RIKGGAEAAAS KPASSTRLPF LAESLGGIEE LLEVAVMTH GGIPEKARE
AGVFDDLVRI SVGIEDTDDL LEDIKALKQK ATN
//

```

## UniProtKB

UniProt incorporates:

- Function of the protein
- Post-translational modification
- Domains and sites.
- Secondary structure.
- Quaternary structure.
- Similarities to other proteins;
- Diseases associated with deficiencies in the protein
- Sequence conflicts, variants, etc.

## In closing ...

- Often only use FASTA files (eg for BLAST)
- Using any sequence where the coordinates are important, need an accession.version OR a gi number
- GBFF and neXtprot are simply human readable versions of these records
- GBFF have become a vehicle for a lot more information than they were meant to do when they were created!
- Keep in mind that GenBank is DNA centric and is a poor vehicle for protein and mRNA expression/interaction information: NCBI (and others) have other databases for these entities.
- All databases I mentioned today are fully “open” ...

## Users must help to keep public databases correct

<http://goo.gl/bdA1q>

*Sir*—With the continued growth of the public DNA sequence databases, and the recent addition of the 11,000,000,000th nucleotide to GenBank (including DDBJ, EMBL and GenBank), it is timely to assess how we use these databases.

GenBank is the archive of all publicly available DNA, RNA and protein sequences. Upon publication a new sequence and its annotations appear in it. Investigators use GenBank in many ways, most commonly for similarity searches such as BLAST; to retrieve records; and for sequence analysis, multiple sequence alignment or pattern finding search. Errors sometimes occur in GenBank, ranging from the trivial (incorrect postal codes), to the misleading (30 nucleotides of vector left on the ends of a record), to the mission-critical (a full length mRNA without a coding sequence (CDS) annotated on it). Also very common are incomplete references that prevent researchers from linking the GenBank record to the publication that refers to it first.

Over the years some people have chosen to report these errors, but in most cases

they are left unmodified. An uncorrected 'discovered' error is one of the worst possible failings in GenBank, so if you discover an error, report it to the database ([update@ncbi.nlm.nih.gov](mailto:update@ncbi.nlm.nih.gov)) and it should be rectified — although a follow-up is advised to make sure this gets done.

If you are a submitter, look at the record you submitted a few years ago: is it still correct? Was the citation ever updated? Take pride in the sequences that carry your name! Our ability to interpret genomes depends on all of these records being as accurate as possible. This is a task for all users of the databases.

**Francis Ouellette**

*Bioinformatics Core Facility, Centre for Molecular Medicine and Therapeutics, University of British Columbia, 950 West 28th Avenue, Vancouver, British Columbia V5Z 4H4, Canada*

© 2001 Macmillan Magazines Ltd

## Scripting from the WWW: `wget/curl`

- Most bioinformatic analyses require resources from the web, e.g. sequences, domain information, datasets, etc.
  - The NCBI and EBI resources are usually scriptable; e.g. write a script that takes a set of accessions from a file and get the sequences
  - Often all that is required is to recognize the URL of the information desired
    - <http://www.ncbi.nlm.nih.gov/protein/P09488>
  - Sometimes, you will need more information to get the desired format (e.g. FASTA)
    - <http://www.ncbi.nlm.nih.gov/protein/121735?report=fasta>
- `curl` and `wget` allow you to pull a web page into a file from the command line:
 

```
curl http://uniprot.org/uniprot/P09488.fasta > p09488.fasta
```
- Sometimes this is what you need; other times more work is required



## Finding a URL (www.uniprot.org)

The first screenshot shows the UniProt search interface. The search bar contains the query 'P09488', which is circled in red. The search results page shows the entry for 'P09488 (GSTM1\_HUMAN)' with the URL 'www.uniprot.org/uniprot/P09488' circled in red.

## Finding a URL to download (uniprot)

The first screenshot shows the UniProt search results for 'P09488 (GSTM1\_HUMAN)'. The 'fasta' download link is circled in red. The second screenshot shows the FASTA file download page with the URL 'www.uniprot.org/uniprot/P09488.fasta' circled in red.

```
curl http://www.uniprot.org/uniprot/P09488.fasta
```

## Finding a URL to download (NCBI)

RecName: Full=Glutathione S-transferase Mu 1; AltName: Full=GST HB sub - Protein - NCBI

Protein  Search

RecName: Full=Glutathione S-transferase Mu 1; AltName: Full=GST HB sub - Protein - NCBI

Protein  Limits Advanced

Display Settings:  GenPept Send to:

**RecName: Full=Glutathione S-transferase Mu 1; AltName: Full=GST HB subunit 4; AltName: Full=GST class-mu 1; AltName: Full=GSTM1-1; AltName: Full=GSTM1a-1a; AltName: Full=GSTM1b-1b; AltName: Full=GTH4**

UniProtKB/Swiss-Prot: P09488.3  
[FASTA](#) [Graphics](#)

[Change region](#)  
[Customize view](#)  
[Analyze this se](#)  
[Run BLAST](#)  
[Identify Conserved](#)

## Finding a URL to download (NCBI)

```
curl http://www.ncbi.nlm.nih.gov/protein/121735?report=fasta
```

RecName: Full=Glutathione S-transferase Mu 1; AltName: Full=GST HB sub - Protein - NCBI

Protein  Limits Advanced

Display Settings:  FASTA Send to:

**We are sorry, but NCBI web applications do not support your browser and may not function properly. [More information](#)**

**RecName: Full=Glutathione S-transferase Mu 1; AltName: Full=GST HB subunit 4; AltName: Full=GST class-mu 1; AltName: Full=GSTM1-1; AltName: Full=GSTM1a-1a; AltName: Full=GSTM1b-1b; AltName: Full=GTH4**

[Change region](#)  
[Analyze this se](#)  
[Run BLAST](#)  
[Identify Conserved](#)  
[Highlight Sequences](#)

## NCBI e-utilities

- The NCBI does not allow their web server to be used for large-scale, automated downloads (unlike Uniprot)  
[www.ncbi.nlm.nih.gov/guide/howto/dwn-records/](http://www.ncbi.nlm.nih.gov/guide/howto/dwn-records/)
- NCBI provides e-utilities (esearch.cgi, efetch.cgi) for programmatic access to ALL NCBI databases (proteins, DNA, also PubMed)  
[www.ncbi.nlm.nih.gov/books/NBK25500/](http://www.ncbi.nlm.nih.gov/books/NBK25500/)
- e-utilities need an Entrez GI number (e.g. 121735 for P09488.1)
- With a GI number, downloading a fasta file is easy (possible):  

```
curl 'http://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=protein&id=121735&rettype=fasta&retmode=text'
```

  
Quotes are required to protect '&' and '?' from shell
- (to keep things simple, use UniProt)

## Homework

From the tabular blastp search you did earlier (-outfmt 6) against swissprot:

1. Edit the file to save 10 lines with E(-)values > 0.1
2. isolate the library (subject) accession information for each of the 10 lines in the edited file, and save the accession in a new file
3. For each accession, split it into its component parts (hit 'man cut' to see how to change the delimiter).
  - Save the gi numbers for use with NCBI in one file
  - Isolate the accessions (P12345.3) in another file.
  - For this second file, isolate only the accessions without the version information.
4. For of the 10 gi numbers or accessions
  - Use the gi numbers to get the sequences from the NCBI
  - Use the protein accessions to get the sequences from UniProt
5. For each of the 10 UniProt accessions, run another blastp search against swissprot, saving the output to accession.bp  
(e.g. `blastp -db swissprot -query P12345.aa > P12345.bp`)