

Mon, Oct 5

# Lecture 4A:

# Molecular sequences



*Asimina triloba*

© Matilda Adams/  
Missouri Botanical Garden

Practical Bioinformatics (Biol 4220)  
Instructor: Michael Landis  
Email: [michael.landis@wustl.edu](mailto:michael.landis@wustl.edu)



# Lecture 4A outline

1. Molecular sequence features
2. Molecular sequence download
3. Some commands
4. Lab 4A overview

A ***molecular sequence*** is a linear sequence from a molecular alphabet

*DNA sequences* are most widely used molecular sequences; *RNA sequences* and *amino acid (AA) sequences* are not uncommon

Sequence data informs all aspects of biology,  
*e.g.*

- mechanism for disease
- gene expression and development
- heredity and ancestry
- protein and cell function
- biodiversity surveys

# Example DNA sequence

sequence length is  
15 base pairs (bp)

Four states in DNA alphabet:  
ACGT

ACGGCTTCTAGCGAT

fifth position  
is in state C



# Example amino acid (AA) sequence

sequence length  
is 15 residues

Twenty states in AA alphabet:  
ACDEFGHIKLMNPQRSTVWY

Y G H I K M A P Q T A E F G H

ninth position  
is in state Q



# Genome

A genome is the complete set of sequences for a single sample (species, individual, tissue)

Genomic data has a *hierarchical structure*

- ***samples*** from a population have different genomes
- each genome is broken into multiple molecules, called ***chromosomes***
- each chromosome may contain some number of ***genes***
- each gene contains ***molecular characters***, and vary in total length and content

# Genes

Genes are discrete heritable subunits of the genome

Some parts of the genome form natural subunits;  
others are named by researchers

Common genes (or loci)

- protein-coding sequences (exons, introns)
- transcription factor binding sites
- repeat regions (microsatellites)
- transposons, retrotransposons
- many loci have no known function

# DNA to mRNA to AA

...ATGCGACGATGGATACCATAG...

# DNA to mRNA to AA

...ATGCGACGATGGATACCATAG...



1. DNA *transcribed*  
into mRNA

AUGCGACGAUGGAAUACCAUAG

# DNA to mRNA to AA

...ATGCGACGATGGATACCATAG...

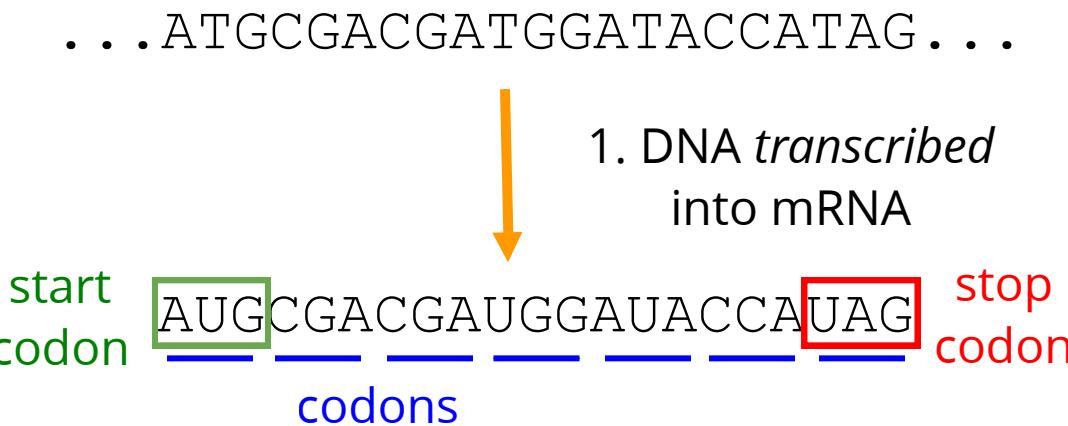


1. DNA *transcribed*  
into mRNA

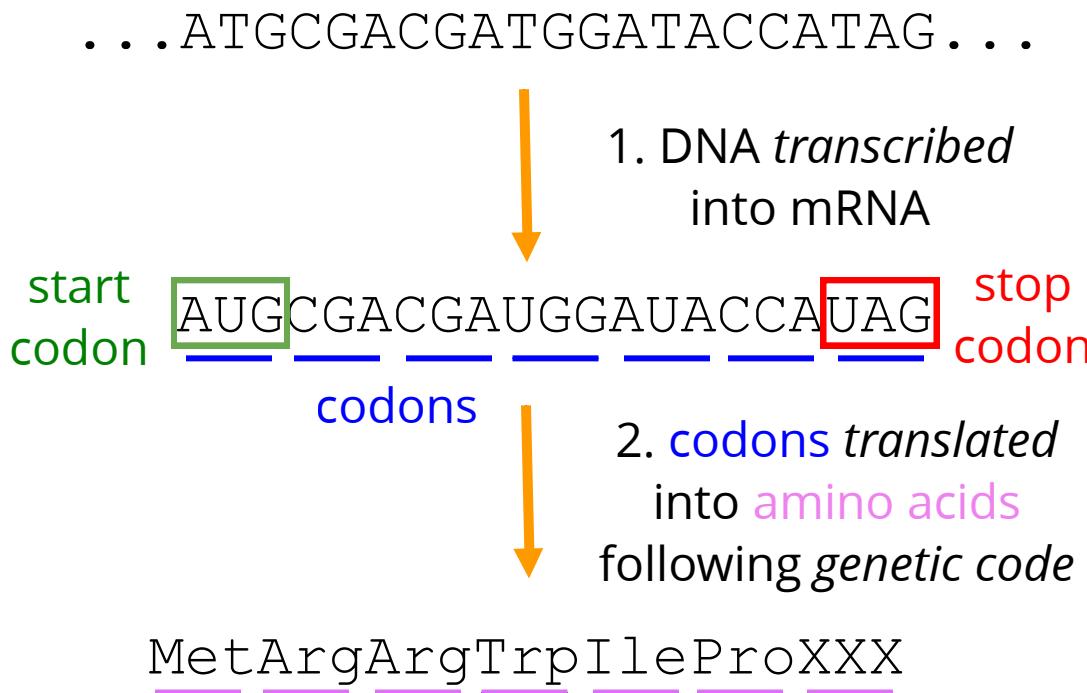
AUGCGACGAUGGAAUACCAUAG

—  
codons  
—

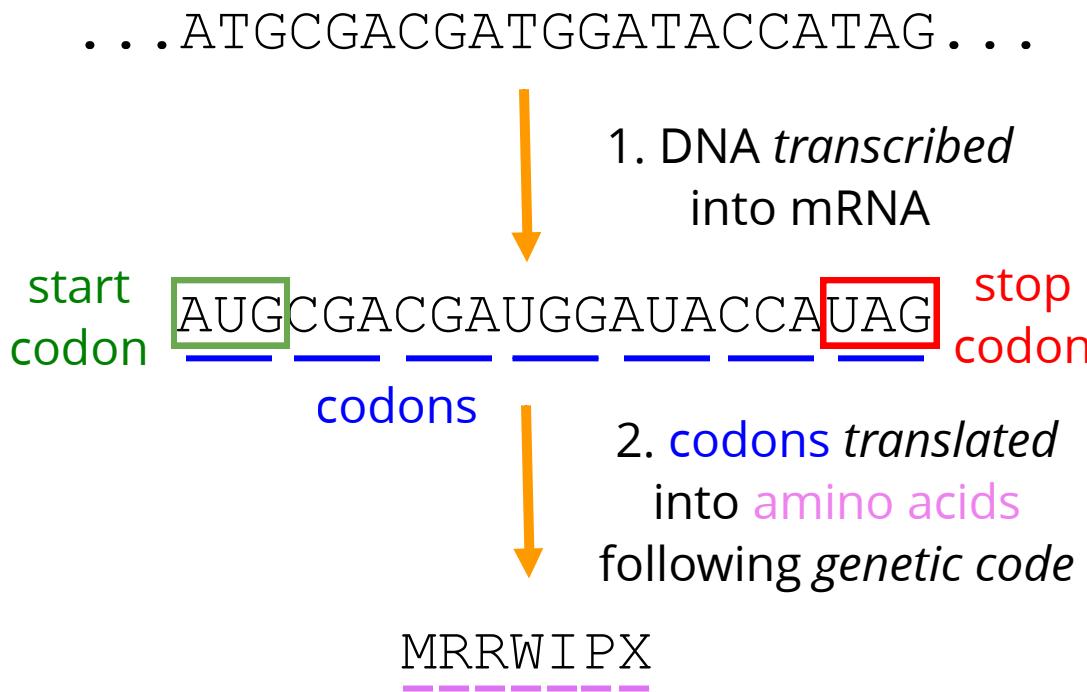
# DNA to mRNA to AA



# DNA to mRNA to AA



# DNA to mRNA to AA



# Genetic code

		Second letter								
		U	C	A	G					
First letter	U	UUU UUC UUA UUG	Phenylalanine (Phe)	UCU UCC UCA UCG	Serine (Ser)	UAU UAC UAA UAG	Tyrosine (Tyr)	UGU UGC UGA UGG	Cysteine (Cys)	U C A G
	C	CUU CUC CUA CUG	Leucine (Leu)	CCU CCC CCA CCG	Proline (Pro)	CAU CAC CAA CAG	Histidine (His)	CGU CGC CGA CGG	Arginine (Arg)	U C A G
	A	AUU AUC AUA AUG	Isoleucine (Ile)	ACU ACC ACA ACG	Threonine (Thr)	AAU AAC AAA AAG	Asparagine (Asn)	AGU AGC AGA AGG	Serine (Ser)	U C A G
	G	GUU GUC GUA GUG	Methionine (Met)	GCU GCC GCA GCG	Alanine (Ala)	GAU GAC GAA GAG	Aspartic acid (Asp)	GGU GGC GGA GGG	Arginine (Arg)	U C A G
			Valine (Val)						Glycine (Gly)	

© Copyright, 2014, University of Waikato. All rights reserved.  
[www.biotechlearn.org.nz](http://www.biotechlearn.org.nz)

Codons determine AA translation

# Genetic code

		Second letter								
		U	C	A	G					
First letter	U	UUU UUC UUA UUG	Phenylalanine (Phe)	UCU UCC UCA UCG	Serine (Ser)	UAU UAC	Tyrosine (Tyr)	UGU UGC	Cysteine (Cys)	U C A G
	C	CUU CUC CUA CUG	Leucine (Leu)	CCU CCC CCA CCG	Proline (Pro)	CAU CAC CAA CAG	Histidine (His)	CGU CGC CGA CGG	Arginine (Arg)	U C A G
	A	AUU AUC AUU AUG	Isoleucine (Ile)	ACU ACC ACA ACG	Threonine (Thr)	AAU AAC AAA AAG	Asparagine (Asn)	AGU AGC AGA AGG	Serine (Ser)	U C A G
	G	GUU GUC GUA GUG	Methionine (Met)	GCU GCC GCA GCG	Valine (Val)	GAA GAC GAA GAG	Lysine (Lys)	GGU GGC GGA GGG	Arginine (Arg)	U C A G

start codon

stop codons

# Genetic code

		Second letter								
		U	C	A	G					
First letter	U	UUU UUC UUA UUG	Phenylalanine (Phe)	UCU UCC UCA UCG	Serine (Ser)	UAU UAC	Tyrosine (Tyr)	UGU UGC	Cysteine (Cys)	U C A G
	C	CUU CUC CUA CUG	Leucine (Leu)	CCU CCC CCA CCG	Proline (Pro)	CAU CAC CAA CAG	Histidine (His)	CGU CGC CGA CGG	Arginine (Arg)	U C A G
	A	AUU AUC AUU AUG	Isoleucine (Ile)	ACU ACC ACA ACG	Threonine (Thr)	AAU AAC	Asparagine (Asn)	AGU AGC AGA AGG	Serine (Ser)	U C A G
	G	GUU GUC GUA GUG	Methionine (Met)	GCU GCC GCA GCG	Alanine (Ala)	GAU GAC GAA GAG	Aspartic acid (Asp)	GGU GGC GGA GGG	Arginine (Arg)	U C A G
						Glutamic acid (Glu)			Glycine (Gly)	

start codon

stop codons

© Copyright, 2014. University of Waikato. All rights reserved.  
[www.biotechlearn.org.nz](http://www.biotechlearn.org.nz)

GUU, GUC, GUA, and GUG  
all encode Valine

AAA and AAG  
encode Lysine

# Protein structure

Amino acid  
sequences

**Primary Structure** = sequence  
of amino acids

3-letter code

**Lys**-Thr-Tyr-Phe-Pro-His-  
Phe-Asp-Leu-Ser-His-**Gly** ...

1-letter code

**K**TYFPHFDLSH**G**

# Protein structure

Amino acid sequences

**Primary Structure** = sequence of amino acids

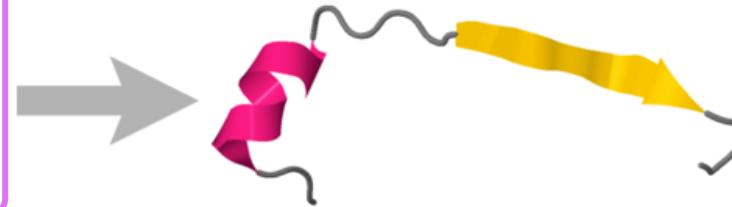
3-letter code

**Lys**-Thr-Tyr-Phe-Pro-His-Phe-Asp-Leu-Ser-His-**Gly** ...

1-letter code

**K**TYFPHFDLSH**G**

**Secondary Structure** = **alpha helices, beta strands**



# Protein structure

Amino acid sequences

**Primary Structure** = sequence of amino acids

3-letter code

Lys-Thr-Tyr-Phe-Pro-His-Phe-Asp-Leu-Ser-His-Gly ...

1-letter code

KTYFPHFDLSH**G**

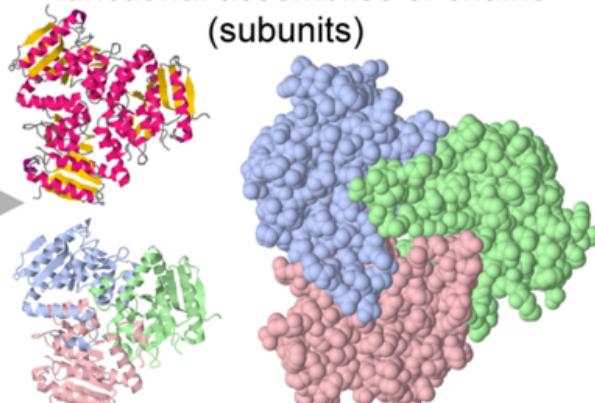
**Secondary Structure** =  
**alpha helices, beta strands**



**Tertiary Structure** = fold helices and strands into domains

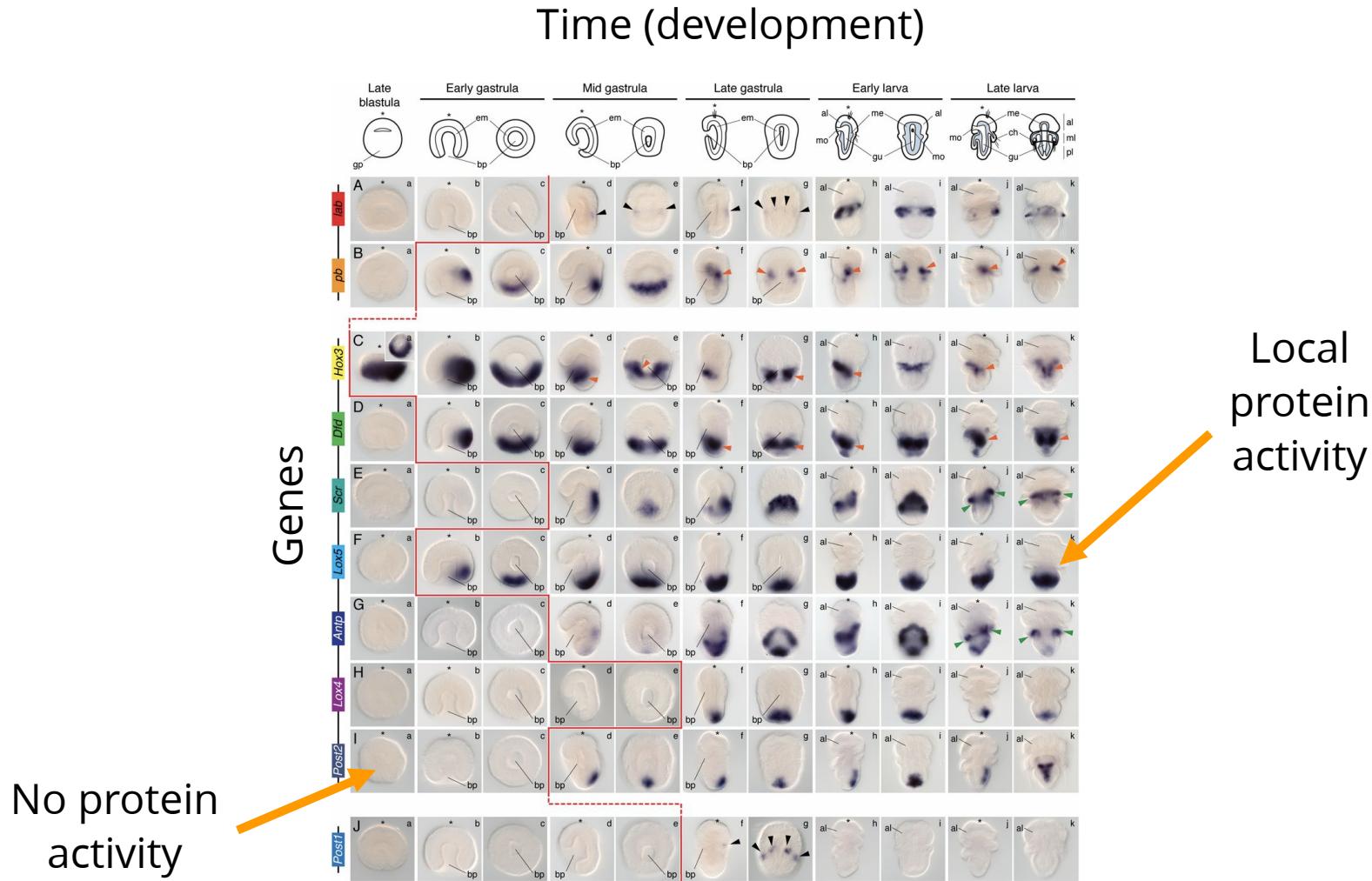


**Quaternary Structure (Biological Units)** = functional assemblies of chains (subunits)

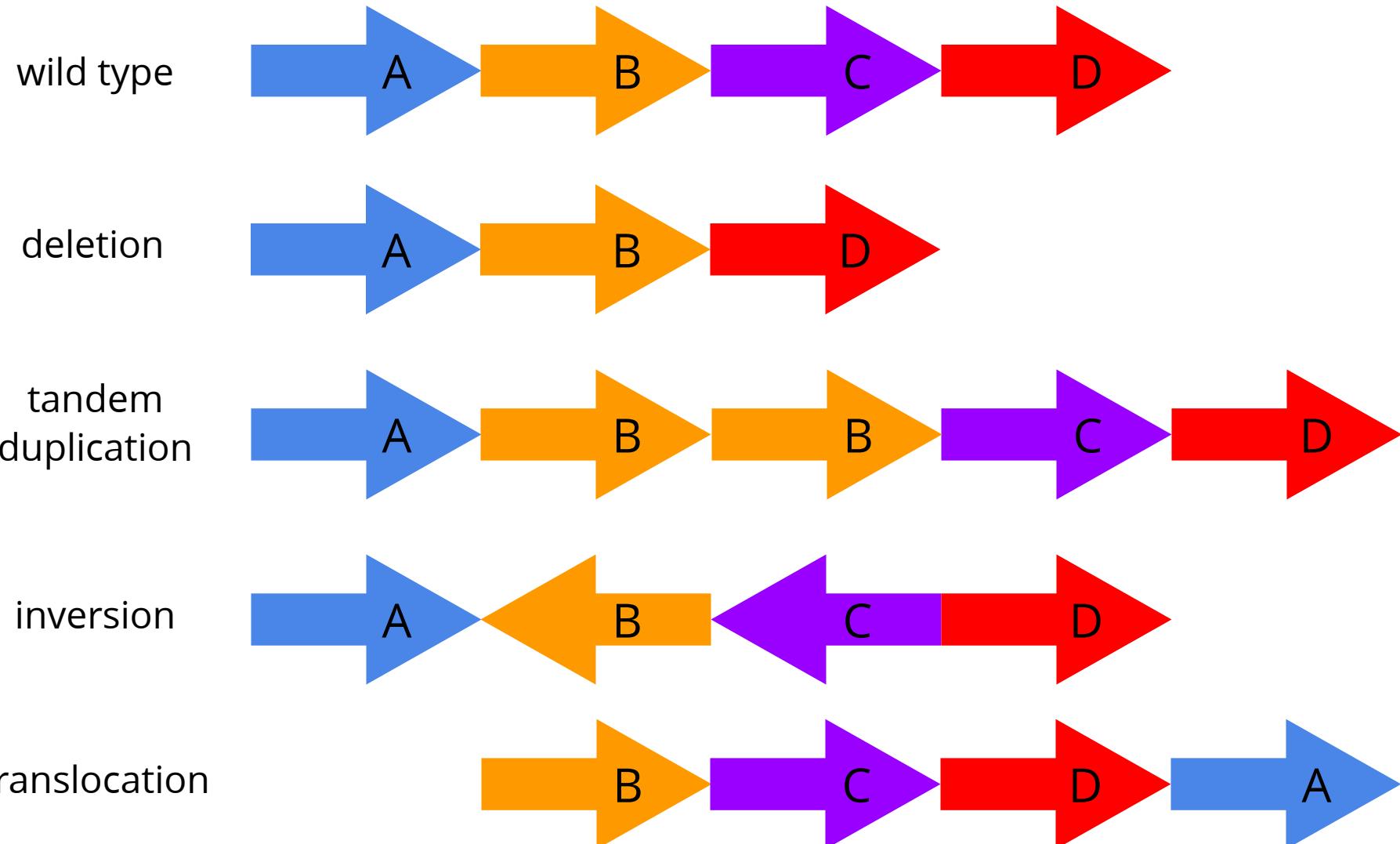


higher order **protein structure**  
increasingly determines  
protein function

# Measuring *gene expression* (using mRNA) relies on sequence data analysis



# ***Structural variation*** of genes among chromosomes studied by analyzing sequence data



# Sequence variation

	*	*   ***
Felis_cattus_cytB	tccgttattcat <b>t</b> tcaatc	
Mus_musculus_cytB	tccgttat <b>c</b> cac <b>a</b> caatc	
Homo_sapiens_cytB	tccgttatt <b>c</b> tctcaatc	
Bos_taurus_cytB	t <b>c</b> tgttattcactcaatc	

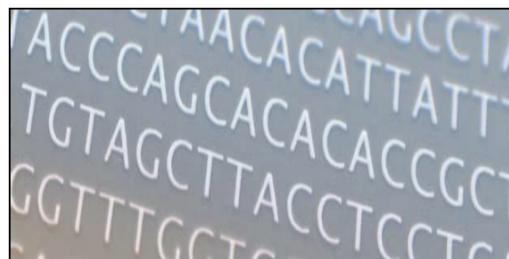
*samples across species*

# NCBI GenBank

NCBI Resources How To Sign in to NCBI

Nucleotide Nucleotide Advanced Search Help

COVID-19 is an emerging, rapidly evolving situation.  
Get the latest public health information from CDC: <https://www.coronavirus.gov>.  
Get the latest research from NIH: <https://www.nih.gov/coronavirus>.  
Find NCBI SARS-CoV-2 literature, sequence, and clinical content: <https://www.ncbi.nlm.nih.gov/sars-cov-2/>.



## Nucleotide

The Nucleotide database is a collection of sequences from several sources, including GenBank, RefSeq, TPA and PDB. Genome, gene and transcript sequence data provide the foundation for biomedical research and discovery.

### Using Nucleotide

[Quick Start Guide](#)

[FAQ](#)

[Help](#)

[GenBank FTP](#)

[RefSeq FTP](#)

### Nucleotide Tools

[Submit to GenBank](#)

[LinkOut](#)

[E-Utilities](#)

[BLAST](#)

[Batch Entrez](#)

### Other Resources

[GenBank Home](#)

[RefSeq Home](#)

[Gene Home](#)

[SRA Home](#)

[INSDC](#)

<https://www.ncbi.nlm.nih.gov>

# GenBank sequences

GenBank publicly hosts annotated DNA sequences

- 218,642,238 sequences
- 654,057,069,549 bases

GenBank sequences are used extensively in research

- identifying anonymous sequences
- inferring gene function
- searching for drug targets
- expanding datasets
- meta-analyses

# GenBank accession

## Homo sapiens isolate C15542T cytochrome b (cytb) gene, partial cds; mitochondrial

GenBank: JN034136.1

[FASTA](#) [Graphics](#) [PopSet](#)

---

Go to:

LOCUS JN034136 1141 bp DNA linear PRI 07-JUL-2012  
DEFINITION Homo sapiens isolate C15542T cytochrome b (cytb) gene, partial cds; mitochondrial.  
ACCESSION JN034136  
VERSION JN034136.1  
KEYWORDS .  
SOURCE mitochondrion Homo sapiens (human)  
ORGANISM [Homo sapiens](#)  
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;  
Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;  
Catarrhini; Hominidae; Homo.  
REFERENCE 1 (bases 1 to 1141)  
AUTHORS Lingala,H., Kumaraswamy,T. and Penagaluru,P.  
TITLE Novel variants in CytB mitochondrial gene in children with non syndromic congenital deafness in south India  
JOURNAL Unpublished  
REFERENCE 2 (bases 1 to 1141)  
AUTHORS Lingala,H., Kumaraswamy,T. and Penagaluru,P.  
TITLE Direct Submission  
JOURNAL Submitted (29-MAY-2011) Molecular Biology, Research Scholar,  
Hubsiguda, Hyderabad, Andhra Pradesh 500007, India

# GenBank accession (cont'd)

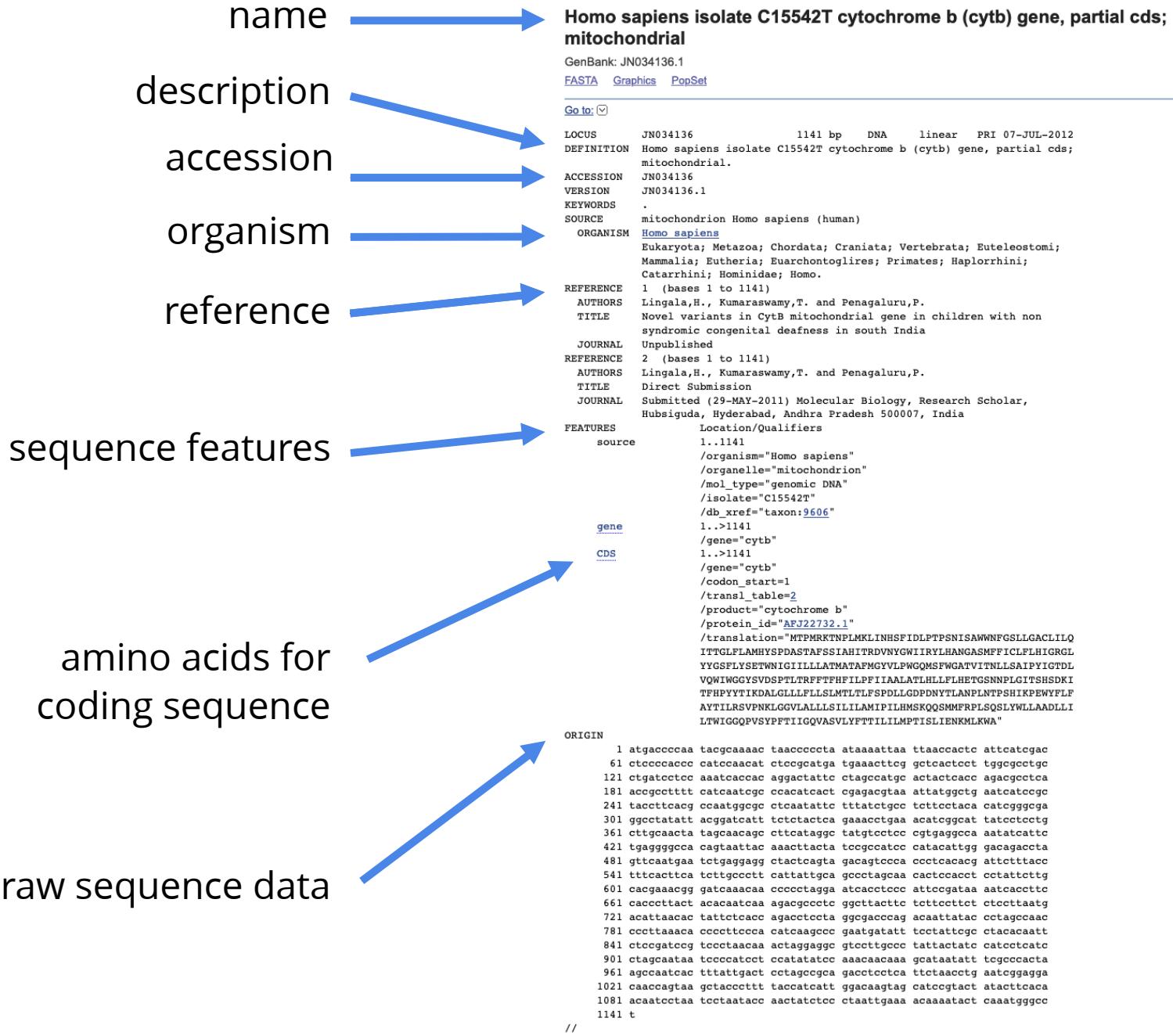
FEATURES	Location/Qualifiers
source	1..1141 /organism="Homo sapiens" /organelle="mitochondrion" /mol_type="genomic DNA" /isolate="C15542T" /db_xref="taxon: <a href="#">9606</a> "
<u>gene</u>	1..>1141 /gene="cytb"
<u>CDS</u>	1..>1141 /gene="cytb" /codon_start=1 /transl_table= <a href="#">2</a> /product="cytochrome b" /protein_id=" <a href="#">AFJ22732.1</a> " /translation="MTPMRKTNP <span style="font-family: monospace;">L</span> MKLINHSFIDLPTPSNISAWWNFGSLLGACLILQ ITTGLFLAMHYSPDASTAFSSIAHITRDVNYGWIIRYLHANGASMFFICLFLHIGRGL YYGSFLYSETWNIGIILLATMATAFMGYVLPWGQMSFWGATVITNLLSAIPYIGTDL VQWIWGGYSVDSPTLTRFFTfhFILPFIIAALATLHLLFLHETGSNNPLGITSHSDKI TFHPYYTIKDALG <span style="font-family: monospace;">L</span> LLFLLSLMTLTLFSPDLLGD <span style="font-family: monospace;">P</span> DNYTLANPLNTPSHIKPEWYFLF AYTILRSVPNKLG <span style="font-family: monospace;">V</span> LALLSILILAMIPILHMSKQQSMMFRPLSQSLYWL <span style="font-family: monospace;">A</span> ADLLI LTWIGGQPVSYPFTIIGQVASVLYFTTILILMPTISLIENKMLKWA"

# GenBank accession (cont'd)

## ORIGIN

```
1 atgaccccaa tacgcaaaac taacccccta ataaaattaa ttaaccactc attcatcgac
61 ctccccaccc catccaacat ctccgcata tgaaacttcg gctcactcct tggcgccctgc
121 ctgatcctcc aaatcaccac aggactattc ctagccatgc actactcacc agacgcctca
181 accgcctttt catcaatcgc ccacatcaact cgagacgtaa attatggctg aatcatccgc
241 tacccctacg ccaatggcgc ctcaatattc tttatctgcc tcttcctaca catcggcga
301 ggccttatatt acggatcatt tctctactca gaaacctgaa acatcgcat tatcctcctg
361 cttgcaacta tagcaacagc cttcataggc tatgtcctcc cgtgaggcca aatatcattc
421 tgaggggcca cagtaattac aaacttacta tccgccatcc catacattgg gacagaccta
481 gttcaatgaa tctgaggagg ctactcagta gacagtccca ccctcacacg attctttacc
541 tttcacttca tcttgccctt cattattgca gccctagcaa cactccacact cctattctg
601 cacgaaacgg gatcaaacaa ccccttagga atcacctccc attccgataa aatcaccttc
661 cacccttact acacaatcaa agacgcctc ggcttacttc tcttccttct ctccttaatg
721 acattaacac tatttcacc agacccctta ggcgacccag acaattatac cctagccaac
781 cccttaaaca ccccttccca catcaagccc gaatgatatt tcctattcgc ctacacaatt
841 ctccgatccg tccctaaacaa actaggaggc gtccttgccc tattactatc catcctcatc
901 ctagcaataa tccccatcct ccatatatcc aaacaacaaa gcataatatt tcgcccacta
961 agccaatcac ttatttact cctagccgca gacccctca ttctaacctg aatcgagga
1021 caaccagtaa gctacccttt taccatcatt ggacaagtag catccgtact atacttcaca
1081 acaatcctaa tccttaatacc aactatctcc ctaattgaaa acaaaaatact caaatgggcc
1141 t
```

//



# NCBI collections

[Viruses](#) > [Riboviria](#) > [Orthornavirae](#) > [Kitrinoviricota](#) > [Alsuviricetes](#) > [Martellivirales](#) > [Togaviridae](#) >

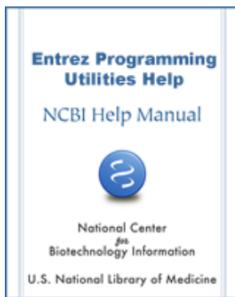
## Alphavirus - 33 complete genomes

Retrieve sequences:

\* The list view for each taxonomy node shows only the next level of sublineages.  
\* Unclassified/unassigned names are written in copper

Species [32] unclassified Alphavirus [1]

Genome	Accession	RefSeq type	Source information	Segm	Length	Protein	Neighbors	Host	Created	Updated
<input type="checkbox"/> Show / hide all segment lists										
Aura virus	<a href="#">NC_003900</a>	complete		-	11824 nt	3	1		02/09/1999	08/13/2018
Barmah Forest virus	<a href="#">NC_001786</a>	complete	strain:BH2193	-	11488 nt	4	32	human, invertebrates, vertebrates	01/14/1997	08/13/2018
Bebaru virus	<a href="#">NC_016962</a>	complete		-	11877 nt	3	-		03/09/2012	08/13/2018
Cabassou virus	<a href="#">NC_038670</a>	complete	strain:CaAr 508	-	11385 nt	2	-		08/24/2018	08/24/2018
Chikungunya virus	<a href="#">NC_004162</a>	complete	strain:S27-African prototype	-	11826 nt	2	876	human, invertebrates, vertebrates	09/06/2002	10/29/2018
Eastern equine encephalitis virus	<a href="#">NC_003899</a>	complete	strain:ssp. North American variant	-	11675 nt	4	453	human, invertebrates, vertebrates	12/16/1991	08/13/2018
Eilat virus	<a href="#">NC_018615</a>	complete	isolate:EO329	-	11634 nt	3	1	invertebrates	09/20/2012	08/13/2018
Everglades virus	<a href="#">NC_038671</a>	complete	strain:Everglades Fe3-7c	-	11395 nt	2	-		08/24/2018	08/24/2018
Fort Morgan virus	<a href="#">NC_013528</a>	complete	isolate:CM4-146	-	11381 nt	3	1	invertebrates, vertebrates	11/25/2009	08/13/2018
Getah virus	<a href="#">NC_006558</a>	complete	isolate:swine	-	11597 nt	3	36	invertebrates, vertebrates	12/17/2004	08/13/2018
Highlands J virus	<a href="#">NC_012561</a>	complete	isolate:585-01	-	11526 nt	3	9	invertebrates, vertebrates	04/15/2009	08/13/2018
Madariaga virus	<a href="#">NC_023812</a>	complete	strain:MADV/Cebus apella/BRA/BEAN5122/1956	-	11624 nt	2	29	human, invertebrates, vertebrates	03/20/2014	08/13/2018
Mayaro virus	<a href="#">NC_003417</a>	complete		-	11411 nt	4	40	human, invertebrates, vertebrates	02/22/2002	08/13/2018
Middelburg virus	<a href="#">NC_024887</a>	complete	isolate:ArB-8422	-	11550 nt	2	3	invertebrates, vertebrates	09/17/2014	08/13/2018
Mosso das Pedras virus	<a href="#">NC_038857</a>	complete	strain:78V-3531	-	11465 nt	2	-		08/24/2018	08/24/2018
Mucambo virus	<a href="#">NC_038672</a>	complete	strain:Mucambo BeAn 8	-	11391 nt	2	1	invertebrates	08/24/2018	08/24/2018
Ndumu virus	<a href="#">NC_016959</a>	complete		-	11688 nt	4	-	invertebrates, vertebrates	03/09/2012	08/13/2018
Onyong-nyong virus	<a href="#">NC_001512</a>	complete		-	11835 nt	3	3	human, invertebrates	08/02/1993	08/13/2018
Pixuna virus	<a href="#">NC_038673</a>	complete	strain:Pixuna BeAr 35645	-	11344 nt	2	-	vertebrates	08/24/2018	08/24/2018
Rio Negro virus	<a href="#">NC_038674</a>	complete	strain:AG80-663	-	11494 nt	2	-	invertebrates	08/24/2018	08/24/2018



# Entrez Direct: E-utilities on the Unix Command Line

Jonathan Kans, PhD<sup>✉1</sup>

Created: April 23, 2013; Updated: September 14, 2020.

## Getting Started

### Introduction

Entrez Direct (EDirect) provides access to the NCBI's suite of interconnected databases (publication, sequence, structure, gene, variation, expression, etc.) from a Unix terminal window. Search terms are entered as command-line arguments. Individual operations are connected with Unix pipes to allow construction of multi-step queries. Selected records can then be retrieved in a variety of formats.

EDirect also includes an argument-driven function that simplifies the extraction of data from document summaries or other results that are returned in structured XML format. This can eliminate the need for writing custom software to answer ad hoc questions. Queries can move seamlessly between EDirect commands and Unix utilities or scripts to perform actions that cannot be accomplished entirely within Entrez.

[www.ncbi.nlm.nih.gov/books/NBK179288/](http://www.ncbi.nlm.nih.gov/books/NBK179288/)

# *Wget, web-get*

Downloads files over HTTP/S and FTP;  
simple interface, allows for recursive copy

```
# download a file located at the target URL
$ wget https://website.com/download_me.txt
# save the file to a new location (`-O`)
$ wget -O my_downloaded_file.txt https://website.com/download_me.txt
# save all files listed in target file
$ wget -i MultipleDownloads.txt
# recursively download all files and subdirectories located
# in the target directory (`-r -np`)
$ wget -r -np http://example.com/configs/.vim/
# ...same as above, but ignore some files (`-R`)
$ wget -r -np -R "index.html*" http://example.com/configs/.vim/
```

# *curl*, transfer a URL

Downloads files over many ports;  
many features, including resume download

download  
interrupted

```
# download file from web address
$ curl -O https://releases.ubuntu.com/20.04.1/ubuntu-20.04.1-desktop-amd64.iso
  % Total    % Received % Xferd  Average Speed   Time     Time     Time  Current
               Dload  Upload   Total   Spent    Left  Speed
  3 2656M    3 89.8M    0      0  11.2M       0  0:03:56  0:00:08  0:03:48 16.9M^C
# resume download using `^C` -
$ curl -C - -O https://releases.ubuntu.com/20.04.1/ubuntu-20.04.1-desktop-amd64.iso
** Resuming transfer from byte position 104992768
  % Total    % Received % Xferd  Average Speed   Time     Time     Time  Current
               Dload  Upload   Total   Spent    Left  Speed
 10 2555M   10  255M    0      0  24.0M       0  0:01:46  0:00:10  0:01:36 27.4M
# file transfer FTP
$ curl -u demo:password -O ftp://test.rebex.net/readme.txt
```

# Lab 4A

[github.com/WUSTL-Biol4220/home/labs/lab\\_03A.md](https://github.com/WUSTL-Biol4220/home/labs/lab_03A.md)