







	Database	Records	Address
luct come of	BioGRID	>2 M protein interactions	https://thebiogrid.org
Just some of	EcoCyc/MetaCyc	>3,200 pathways from >3,400 organisms	http://www.ecocyc.org, http://www.metacyc.org
the resources	Ensembl (+ BioMart for easy sequence queries)	Major repository of DNA sequences, genomes, genes, proteins, and transcripts	http://useast.ensembl.org/index.html
	Entrez Genome	Millions of genome sequences	http://www.ncbi.nlm.nih.gov/genome?db=genome
bioinformatics	Expression Atlas	159K mRNA expression expts in 66 species	https://www.ebi.ac.uk/gxa/home/
	Genbank	>5 triillion bases sequenced; > 33 trillion bases as whole genome shotgun data	https://www.ncbi.nlm.nih.gov/genbank/
Think of these	Gene Expression Omnibus (GEO)	>7 M mRNA or protein expression expts	http://www.ncbi.nlm.nih.gov/geo/
data for new	Genomes Online Database (GOLD)	>500K genome sequences, many in progress	https://gold.jgi.doe.gov/index
discoveries	Human Protein Atlas	millions of high-res images of ~17K human proteins across tissues, cancers, & cell lines	http://www.proteinatlas.org/
	KEGG	Most known pathways, in >500 graphical diagrams and >10K organisms (via homology)	http://www.genome.ad.jp/kegg/
	Medline / PubMed	>37 million references	https://www.ncbi.nlm.nih.gov/PubMed/
	Mouse Genome Informatics	~20,000 mouse genes, diverse associated data & annotations	http://www.informatics.jax.org/
	Online Mendelian Inheritance in Man (OMIM)	Compendium of human genes and genetic phenotypes, data for >16,000 human genes	https://www.ncbi.nlm.nih.gov/omim/
	Pride	Hundreds of millions of peptide mass spectra from 10's of thousands of experiments	https://www.ebi.ac.uk/pride/archive/
	Reactome	>2K pathways involving >11K human proteins, also other organisms	https://www.reactome.org/
	SGD	~6,000 yeast genes, diverse associated data & annotations	https://www.yeastgenome.org/
	UniProtKB/SWISS-PROT	>570K hand-curated sequence entries from >14K organisms	https://www.uniprot.org/



Live demo Ensembl->BioMart->filter for [JOUBERT SYNDROME 5]-> CEP290, OMIM, Reactome, Human Protein Atlas & OpenCell, AlphaFold database (CEP290 is awful!)

It's nice to know that all of this exists, but ideally, you'd like to be able to so something constructive with the data.

That means getting the data inside your own programs.

All of these databases let you download data in big batches, but this isn't always the case, so....

Let's empower your Python scripts to grab data from	the web.
 For a number of specific biological databases, you can use BioPython BioPython lets you access sequence & structure databases, read fasta do simple sequence analyses, BLAST, etc, right from your Python co If you need to install it, just open an Anaconda prompt (on a PC) or la from Anaconda Navigator & type "pip install biopython" 	I/genome files, de unch a console window
e.g.	
from Bio import Entrez Entrez.email = "your_email@gmail.com" # Always tell NCBI who you handle = Entrez.efetch(db="nucleotide", id="EU490707", rettype="gb", print(handle.read()) LOCUS EU490707 1302 bp DNA linear PLN 26-JUL-2016	<mark>are</mark> , retmode="text")
DEFINITION Selenipedium aequinoctiale maturase K (matK) gene, partial cds; chloroplast. ACCESSION EU490707 VERSION EU490707.1	
KEYWORDS . SOURCE chloroplast Selenipedium aequinoctiale ORGANISM Selenipedium aequinoctiale ORIGIN	$\mathcal{N}\mathcal{N}$
1 attttttacg aacctgtgga aatttttggt tatgacaata aatctagttt agtacttgtg 61 aaacgtttaa ttactcgaat gtatcaacag aatttttga tttcttcggt taatgattct	biopython
 There's a complete pdf tutorial @ http://biopython.org/DIST/docs/tutori	al/Tutorial.pdf



The basic idea:

We first set up a "request" by opening a connection to the URL.

We then save the response in a variable and print it.

If it can't connect to the site, it'll print out a helpful error message instead of the page.

You can more or less use the commands in a cookbook fashion....

For example:	
import urllib.request	# include the urllib.request module
url = "https://www.utexas.edu/"	
x = urllib.request.urlopen(url) print(x.read())	# setup a request # read page and show the result to the user
	Python 3 version

We can be slightly fancier in order to handle different formats and the inevitable internet connection errors









If you run that program, you should get back				
>>>				
html				
lots of metadata				
OWN - NLM				
STAT- MEDLINE				
DCOM- 20010322				
LR - 20210108	the Medline entry for the human			
IS - 0028-0836 (Print)	genome sequence naner			
IS - 0028-0836 (Linking)	Senome sequence paper			
VI - 409 ID - 6822				
DP - 2001 Feb 15				
TI - Initial sequencing and analysis of the human ge	enome.			
PG - 860-921				
AB - The human genome holds an extraordinary trove of information about human				
development, physiology, medicine and evolution	on. Here we report the results of an			
international collaboration to produce and mak	e freely available a draft sequence of			
the human genome. We also present an initial analysis of the data, describing some				
of the insights that can be gleaned from the sec	juence.			
FAU - Lander, E S				
AD - Whitehead Institute for Biomedical Research	Center for Genome Research, Cambridge			
MA 02142. USA, lander@genome.wi.mit.edu				
,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,				
[and so on]				







Since this was PubMed, w	ve could have done i	t in BioPython as well:
from Bio import Entrez pmid = "11237011" Entrez.email = "your_email@gma handle = Entrez.efetch(db="pubn print(bandle read())	ail.com" # Always tell NCB ned", id=pmid, rettype="me	l who you are edline", retmode="text")
	<pre>print(handle.read().</pre>	count("AU - "))
 PMID- 11237011 OWN - NLM STAT- MEDLINE DCOM- 20010322 LR - 20240729 IS - 0028-0836 (Print) IS - 0028-0836 (Linking) VI - 409 IP - 6822 DP - 2001 Feb 15 TI - Initial sequencing and analysis of the human PG - 860-921 AB - The human genome holds an extraordinary development, physiology, medicine and evolu international collaboration to produce and ma of the human genome. We also present an ini some of the insights that can be gleaned from 	genome. trove of information about human tion. Here we report the results of an ike freely available a draft sequence tial analysis of the data, describing the sequence.	biopython

