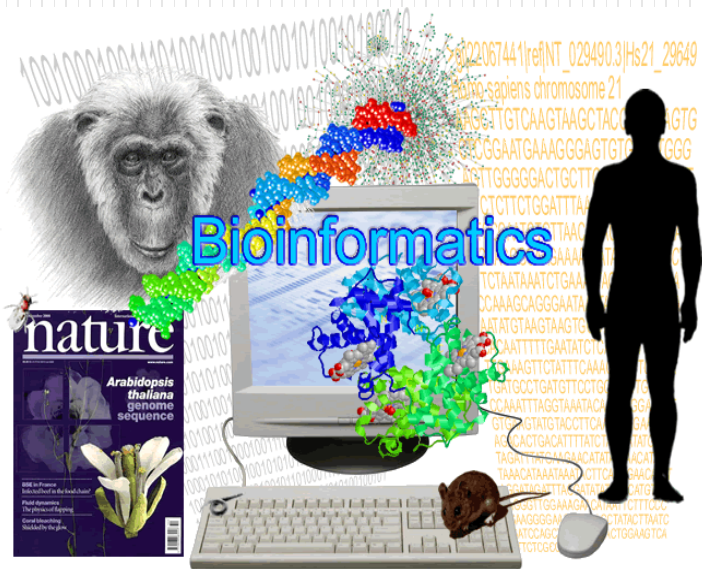


# Bioinformatics, Syntenic Biology & Genome Editing



薛佑玲 PhD

Institute of Biomedical Sciences

National Sun Yat-sen University

[ylshiu@mail.nsysu.edu.tw](mailto:ylshiu@mail.nsysu.edu.tw)

# Outline

---

**Introduction: a Short History About Bioinformatics**

---

**Bioinformatics Q & A**

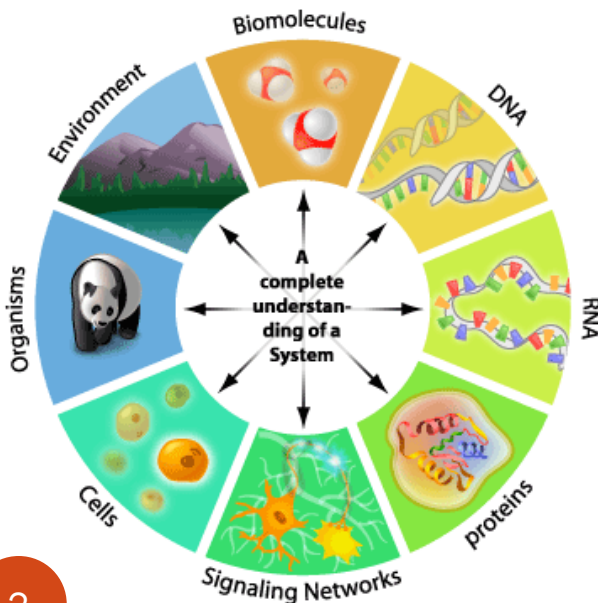
---

**Synthetic biology**

---

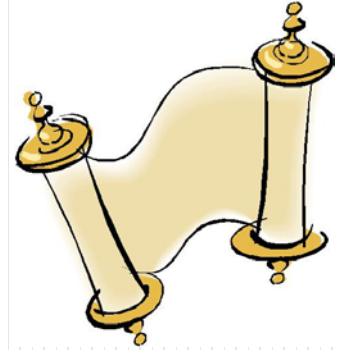
**Genome Editing**

---

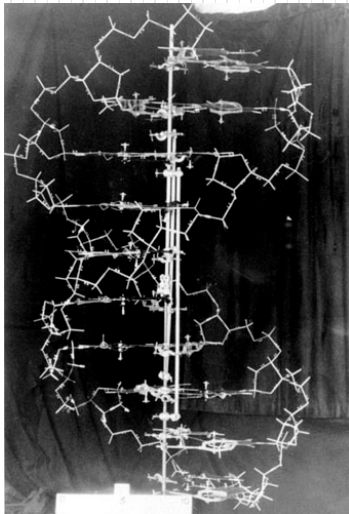


**b101nf0rmat1cs**

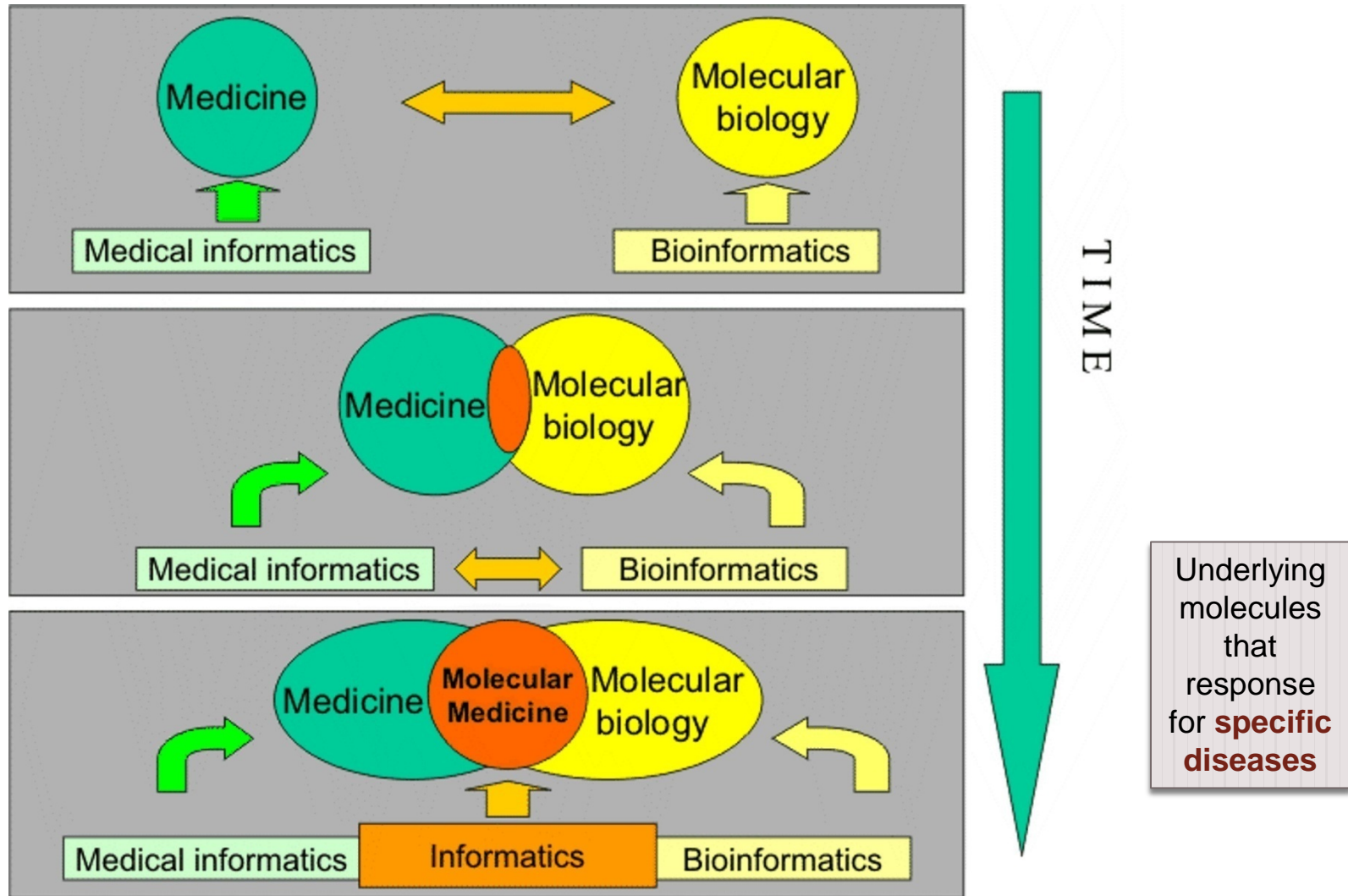
```
ACCATGGATTACATA00110110001101010  
GATTCCATTATAAGGA01100111000000100  
TGCCGGCAATAGGCA001110101000110101  
CAATAAGCATTCCAC001010101101011011
```



# A Short History about Bioinformatics



# The Convergence between MI & BI



# Top Ten Medical Breakthroughs – since 1840

Hygiene equipment

Antibiotics

Anesthetic

Vaccine

Discovery of DNA structure

Microbiology theory

'The Pill': the combined oral contraceptive pill

Evidence-based Medicine

Medical imaging ( e.g., X-ray, MRI...)

**Computer**

**Stem cell therapy**

根據British Medical Journal 線上意見調查，自1840年創刊以來，最重要的醫學里程碑

# Day 4: Computer Science and Medicine

Csedweek

11 部影片

訂閱



Video player controls including a play/pause button, a volume icon, a progress bar showing 0:04 / 2:03, a resolution dropdown set to 480p, and a full screen button.

# The Holy Grail of Bioinformatics

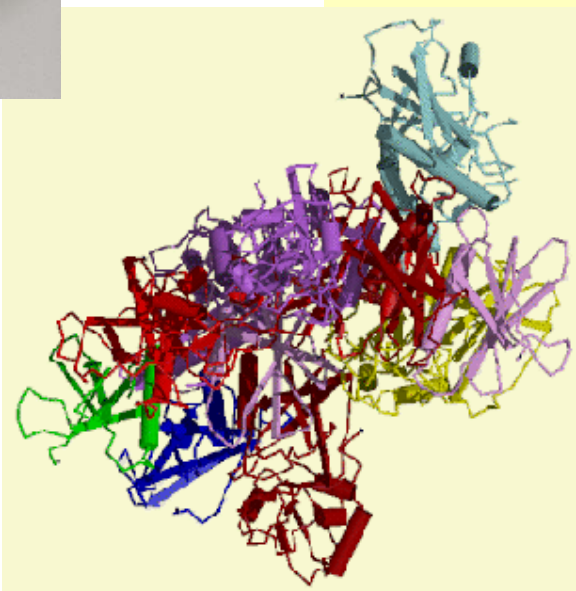
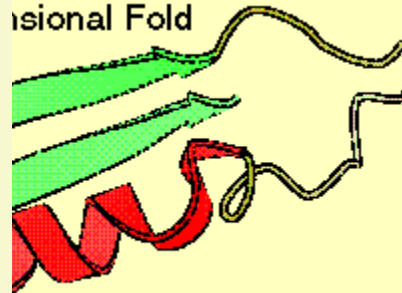


```
MNGTEGPNFYVPFSNKTGVVRS PF EAPQYYLAEPWQFSMLAAYMFL L I V L  
GFPINFLTLY V T V Q H K K L R T P L N Y I L L N L A V A D L F M V F G G F T T T L Y T S L H  
G Y F V F G P T G C N L E G F F A T L G G E I A L W S L V V L A I E R Y V V V C K P M S N F R F G E  
N H A I M G V A F T W V M A L A C A A P P L V G W S R Y I P Q G M Q C S C G A L Y F T L K P E I N N
```

Amino Acid Sequence



3D Structural Fold



...to be able to understand **the words in a sequence sentence** that form a particular protein **structure** (from Attwood & Parry-Smith 1999)

# A Short History Overview (I) - Wet

**1953:** Double helix of DNA (Waston & Crick)

**1954:** First protein sequence (**insulin** by **Sanger**)

**1958:** First X-ray 3D structure of a protein (**myoglobin** by Kendrew)

**1972:** First DNA sequencing

**1977:** Rapid **sequencing** techniques (**Gilbert & Sanger**)

**1986:** PCR (the photocopying machine of the biologist)

**1992:** Sequence of **yeast** chromosome III ( $3 \cdot 10^5$  bp)

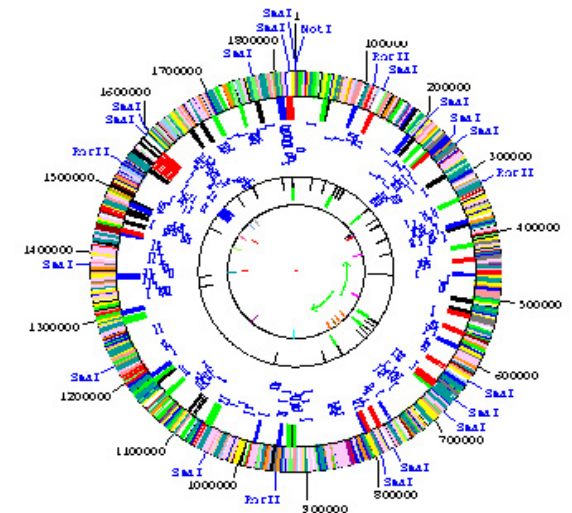
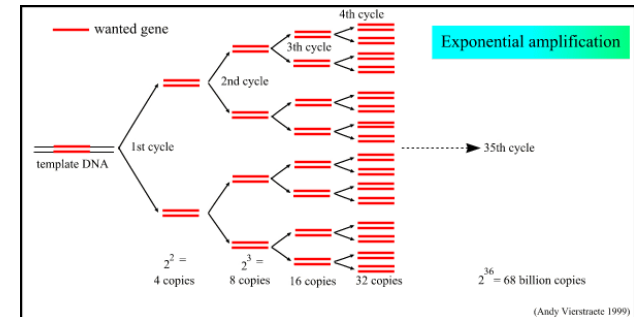
**1995:** Sequence of the genome of the bacteria: **Haemophilus influenzae** ( $2 \cdot 10^6$  bp)

**1999:** Sequence of the genome of a **multi-cellular organism**: *Caenorhabditis elegans* ( $10^8$  bp)

**2000:** Blue draft of the **human genome** ( $3 \cdot 10^9$  bp)

**2002:** Genome of *Ashbya gossypii* (**Saccharomycetes**)

**Recent:** [GOLD database](#)





# A Short History Overview (I) - Dry

**1965:** «Atlas of protein sequence and structure» (**Dayhoff**)

**1967:** Fitch WM (Phylogenetic trees)

**1970: Needleman/Wunsch** (1st similarity search algorithm)

**1971:** PDB (3D structure database)

**1977: Staden** (1st sequence analysis software suite)

**1980: EMBL Heidelberg**

**1980: Smith/Waterman algorithm**

**1982:** EMBL Nucleotide Sequence Database and GenBank

**1985: CABIOS** (1st scientific journal for bioinformatics)

**1985:** FASTP (ancestor of **FASTA**, Blast, etc.)

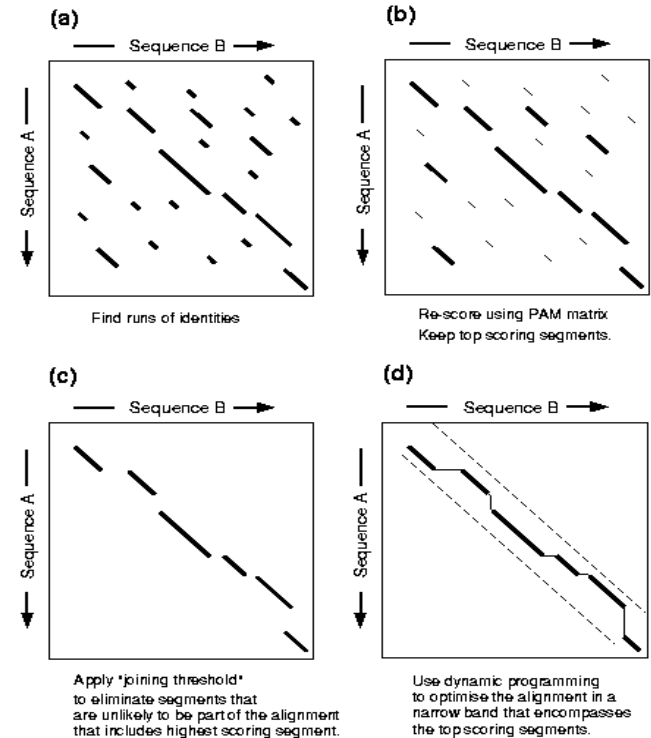
**1986:** Swiss-Prot (Protein Sequence Database)

**1988: Creation of the NCBI in the USA**

**1992:** EBI founded as EMBL outstation in **Hinxton** (Wellcome Trust Campus)

**1993: ExPASy** (1st WWW server for the life sciences)...

## FASTA Algorithm



# Early Bioinformatics: the birth of a discipline – Quzounis CA & Valencia A (2003)

**Table 2.** Twenty Publications that influenced our view of bioinformatics

Publication	Comments
Zuckerlandl and Pauling, 1965b	First use of molecular sequences for evolutionary studies
Fitch and Margoliash, 1967	Use of molecular sequences to build trees
Needleman and Wunsch, 1970	First implementation of dynamic programming for protein sequence comparison
Lee and Richards, 1971	Calculation of accessibility on protein structures
Chou and Fasman, 1974	First secondary structure prediction method
Tanaka and Scheraga, 1975	Simulation of protein folding
Dayhoff, 1978	First collection of protein sequences
Hagler and Honig, 1978	One of the first explicit attempts to simulate protein folding
Doolittle, 1981	Seminal paper examining divergence and convergence in protein evolution
Felsenstein, 1981	One of the first statistical treatments of evolutionary tree construction
Richardson, 1981a	The most comprehensive description of protein structure to that date
→ Kabsch and Sander, 1984	Discovery with profound implications for model building by homology and structure prediction
Novotny <i>et al.</i> , 1984	The inability of distinguishing correct from incorrect structures threw back structure prediction approaches for a long while
→ Chothia and Lesk, 1986	Examination of divergence between sequence and structure
Doolittle, 1986	Influential book on sequence analysis
Feng and Doolittle, 1987	The first approach for an efficient multiple sequence alignment procedure, later implemented in CLUSTAL
Lathrop <i>et al.</i> , 1987	One of the first applications of Artificial Intelligence in protein structure analysis and prediction
Ponder and Richards, 1987	The very first threading approach, using sequence enumeration
Altschul <i>et al.</i> , 1990	The implementation of a sequence matching algorithm based on Karlin's statistical work
Bowie <i>et al.</i> , 1991	The first implementation of protein structure prediction using threading

# Bioinformatics: A Snapshot 10 Years Ago

Pharmaceutical companies were **not interested**

**Life scientists** believed that it was an **outlet** for **failed biologists** that want to play around with computers

**Computer scientists** did not even consider it important, they confused it with **bio-inspired “computer sciences”**

*E.g., genetic algorithm, artificial life, ant algorithm, neural network*

DNA computers...

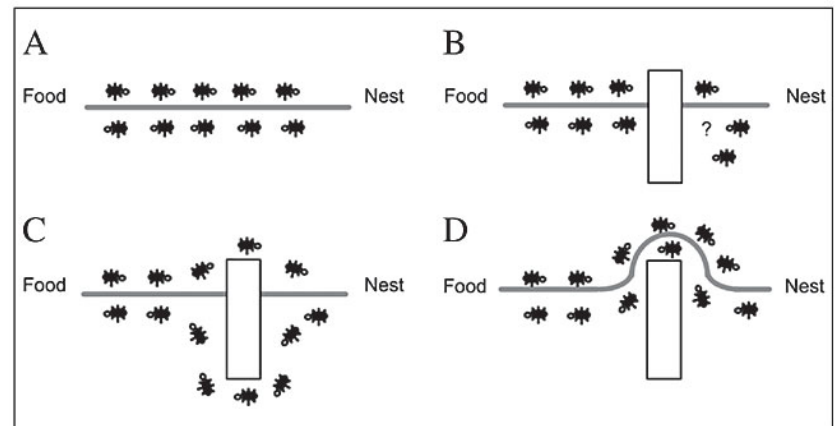
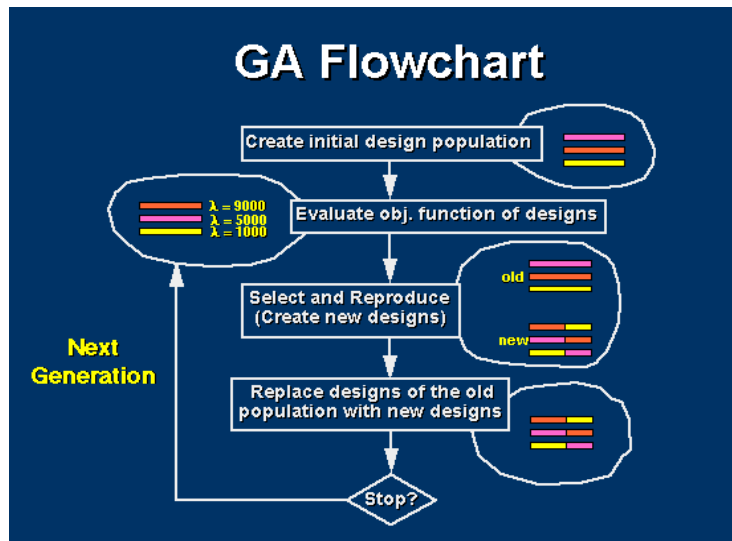


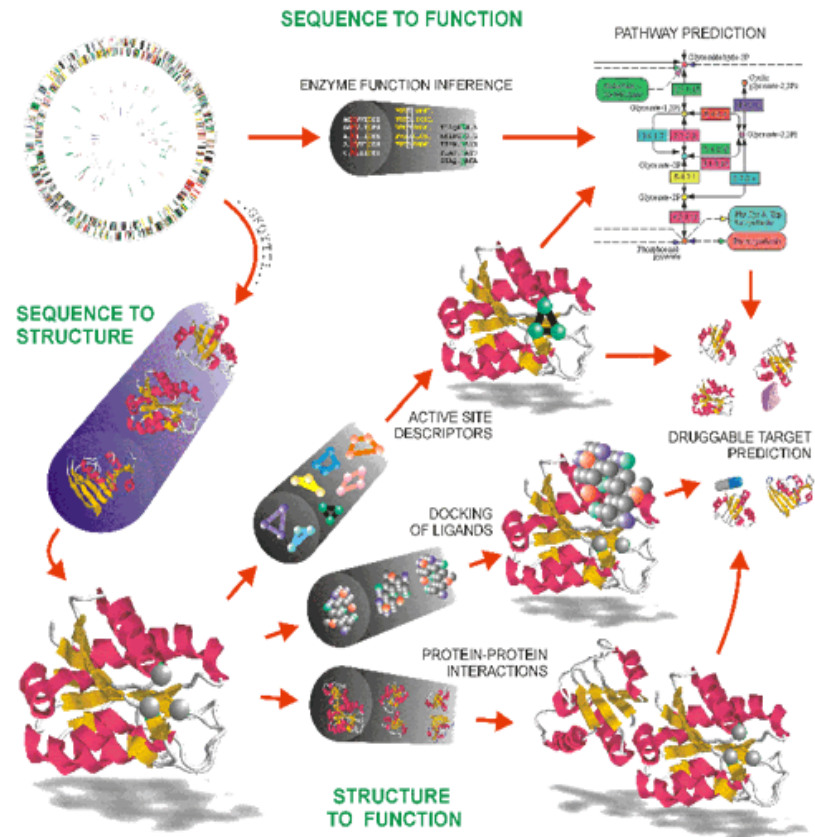
Figure 2. A, ants in a pheromone trail between nest and food; B, an obstacle interrupts the trail; C, ants find two paths to go around the obstacle; D, a new pheromone trail is formed along the shorter path.

# Bioinformatics in 2003

**Pharmaceutical companies** believe that it is **the most efficient way** to streamline the process of **drug discovery**

Some life scientists believe it is **the solution to all problems in life sciences** and that it will allow them **to avoid** doing **some experiments**

**Computer scientists** are very interested: **the scope and complexity** of the domain makes it the ideal field of application of **new software techniques** and specialized hardware developments



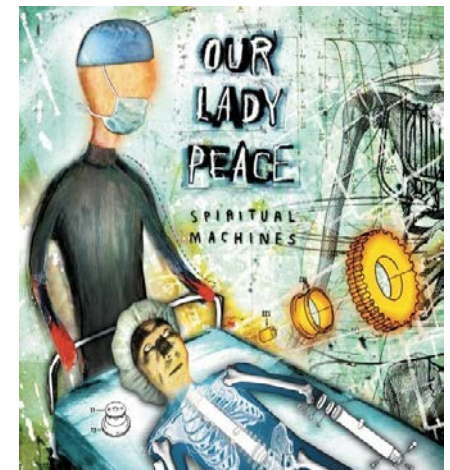
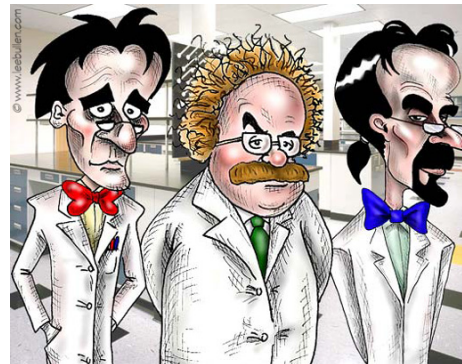


# Bioinformatics in 2010

Pharmaceutical companies use it **routinely**, but have realized that it **complements** rather than **replaces** experimental work

Life scientists use it **efficiently every day** and therefore **forget** that it **exists**

Computer scientists may have jumped on **another fancy subject**: Spiritual machines?



# Bioinformatics in 2020

## THEME: Innovation with AI and Cognitive Computing

### TOPICS OF INTEREST

Topics of interest include, but are not limited to:



Adaptive computation in bioinformatics  
Bio-data visualization  
Bio-inspired computing  
Biological network reconstruction and analysis  
Biomarker discovery  
Computational systems biology  
Coronavirus disease  
Disease classification  
DNA, RNA and protein sequence analysis



Drug discovery and validation  
Epigenetics/epigenomics  
Epidemiology  
Formal validation of biological systems  
Functional genomics  
Gene expression analysis  
Health informatics  
Human-centric applications  
Medical and biomedical informatics



Metagenomics data analysis  
Modeling and simulation of biological processes,  
pathways, etc.  
Molecular evolution and phylogeny  
Next-generation and Third-generation sequencing  
Parallel and distributed computing for life science  
Population genetics  
Proteomics & other omics  
Protein folding  
Translational bioinformatics

# Artificial Intelligence

- 一般稱的 AI 其實是 **Artificial Intelligence** 的縮寫，而這個名字也清楚地表達了它的涵義。
- 人工智慧的定義其實就是以「人工」編寫的**電腦程式**，去模擬出**人類的「智慧」行為**，其中包含模擬人類感官的「聽音辨讀、視覺辨識」、大腦的「推理決策、理解學習」、動作類的「移動、動作控制」等行為。

# ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



# MACHINE LEARNING

Machine learning begins to flourish.



# DEEP LEARNING

Deep learning breakthroughs drive AI boom.



1950's

1960's

1970's

1980's

1990's

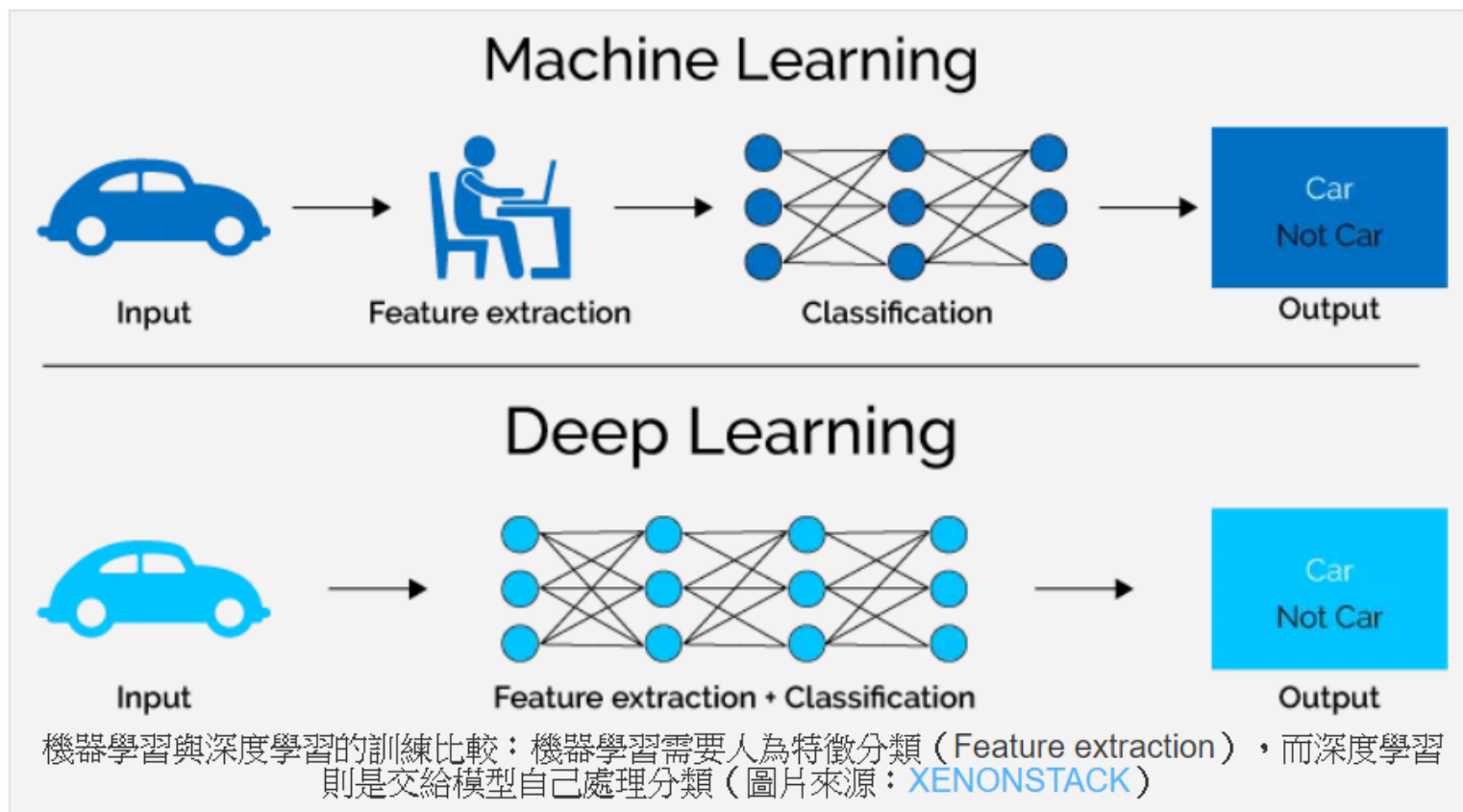
2000's

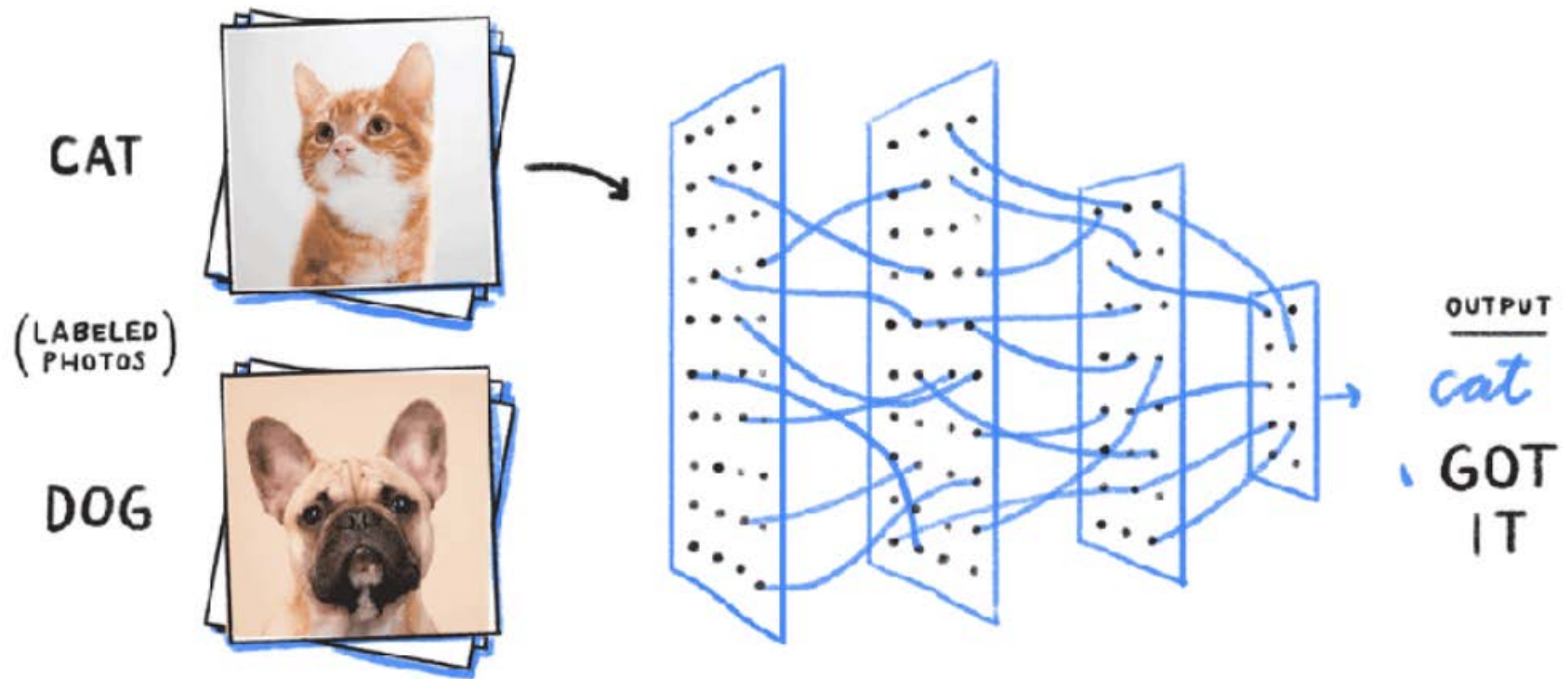
2010's

AI 演進 (圖片來源: [NVIDIA](#))



我們可以從上面這張圖清楚理解，AI、ML、DL 這三個名詞的關係就像洋蔥一樣層層遞進，機器學習（ML），是人工智慧（AI）底下的技術分支，而深度學習（DL）是近年才從機器學習衍伸出的領域，可以比喻為俄羅斯娃娃，一個子領域之中又有更深入的子領域。

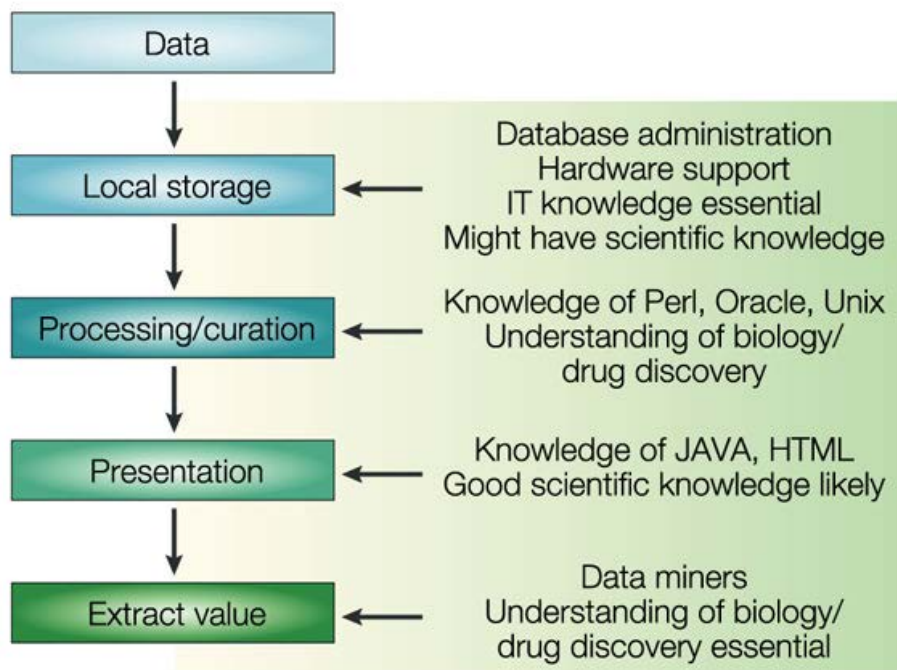




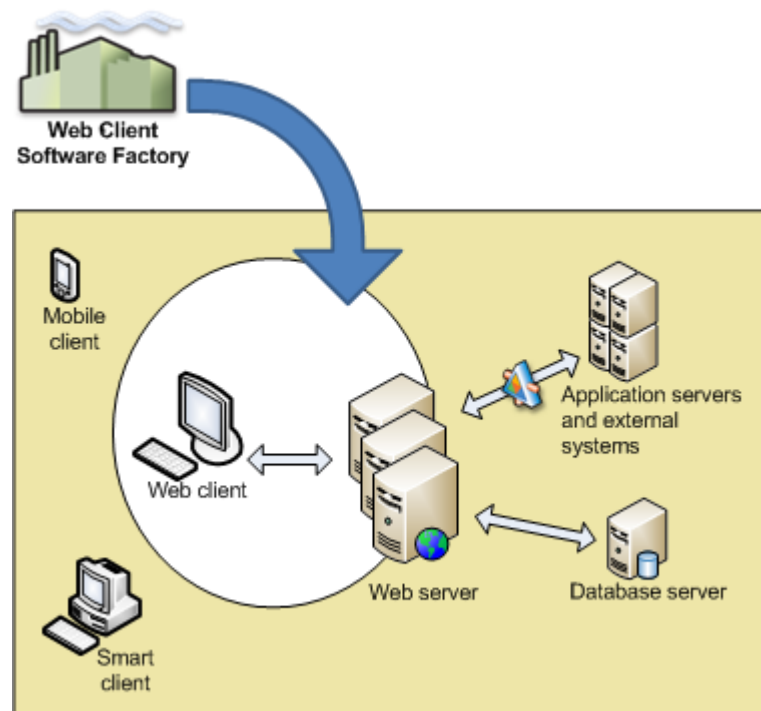
透過捲積神經網路訓練貓狗辨識 (圖片來源: [Medium](#))

## Convolutional Neural Network, CNN

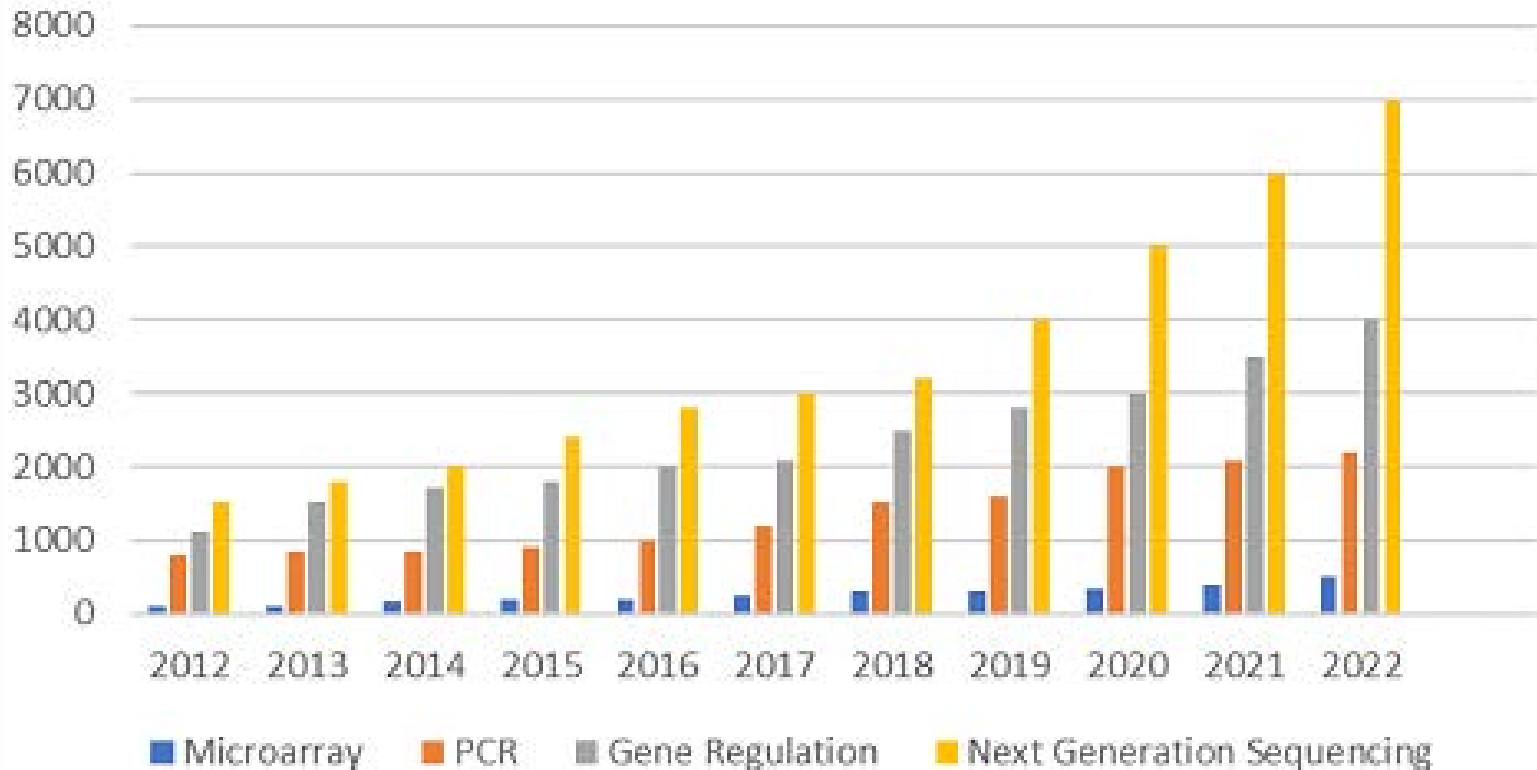
# Resources: databases & software



Nature Reviews | Drug Discovery




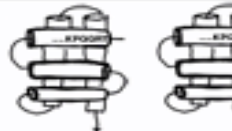





## Trascriptomics Technologies Market Analysis



# Breadth: Homologs, Large-scale Surveys, Informatics–

	pairwise comparison, sequence & structure alignment	multiple alignment, patterns, templates, trees	databases, scoring schemes, censuses
--	-----------------------------------------------------------	------------------------------------------------------	-----------------------------------------

<b>1</b>	<b>2</b>	<b>3-100</b>	<b>100+</b>
----------	----------	--------------	-------------

Depth: Rational Drug Design (physi		<b>Genome Sequence</b>	atc gatc gatattgggattgggga	atc gatc gatattgggattgggga atc gatc gatattgggattgggga	atc gatc gatattgggattgggga atc gatc gatattgggattgggga atc gatc gatattgggattgggga atc gatc gatattgggattgggga atc gatc gatattgggattgggga	atc gatc gatattgggattgggga atc gatc gatattgggattgggga atc gatc gatattgggattgggga atc gatc gatattgggattgggga atc gatc gatattgggattgggga
	gene finding	↓				
		<b>Protein Sequence</b>	ALMNAKKKPPQRT	ALMNAKKKPPQRT ALMNAKKKPPQRT	ALMNAKKKPPQRT ALMNAKKKPPQRT ALMNAKKKPPQRT ALMNAKKKPPQRT	ALMNAKKKPPQRT ALMNAKKKPPQRT ALMNAKKKPPQRT ALMNAKKKPPQRT ALMNAKKKPPQRT ALMNAKKKPPQRT ALMNAKKKPPQRT ALMNAKKKPPQRT ALMNAKKKPPQRT ALMNAKKKPPQRT ALMNAKKKPPQRT ALMNAKKKPPQRT
	structure prediction	↓				
		<b>Protein Structure</b>				
	geometry calculation	↓				
		<b>Protein Surface</b>				
	molecular simulation	↓				
		<b>Force Field</b>				
	structure docking	↓				
	<b>Ligand Complex</b>					



COFFEE



BREAK



*"Don't just sit there! If you've processed all the data there is, go out and find more data!"*

Reproduced in R.L. Weber, *"A random walk in science"*, IOP Publishing, 1973

# Case Study

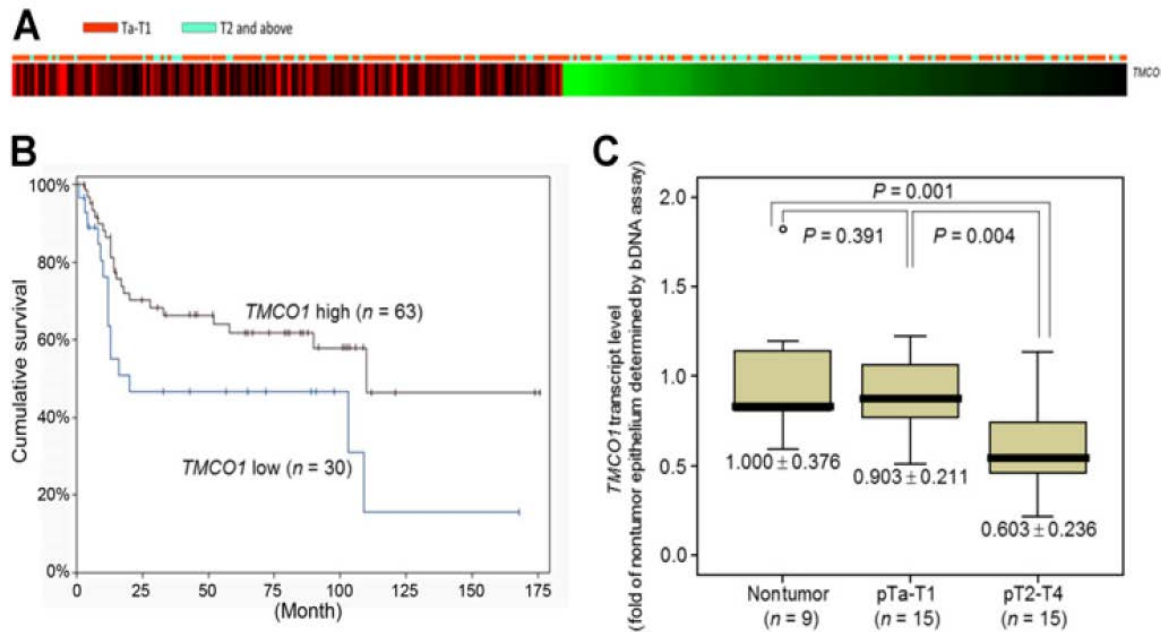


# Transmembrane and Coiled-Coil Domain 1 Impairs the AKT Signaling Pathway in Urinary Bladder Urothelial Carcinoma: A Characterization of a Tumor Suppressor



Chien-Feng Li<sup>1,2,3,4</sup>, Wen-Ren Wu<sup>5</sup>, Ti-Chun Chan<sup>1,5</sup>, Yu-Hui Wang<sup>1,6</sup>, Lih-Ren Chen<sup>4,7,8</sup>,  
Wen-Jeng Wu<sup>9,10,11,12,13,14,15</sup>, Bi-Wen Yeh<sup>9</sup>, Shih-Shin Liang<sup>5,16</sup>, and Yow-Ling Shiue<sup>5,17,18</sup>





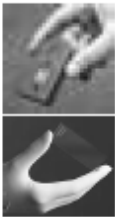
**Figure 1.**

Downregulation of the *TMCO1* protein predicts poor disease-specific and metastasis-free survivals. **A**, A heatmap shows the data analysis from GSE32894 (GEO dataset), which identified that the *TMCO1* transcript is significantly downregulated ( $P = 0.0009$ ) in muscle-invasive UBUC (blue bars). **B**, The downregulation of the *TMCO1* transcript was also predictive of poor overall survival in an independent dataset (GSE31684, GEO, NCBI;  $P = 0.0425$ ). **C**, Quantitative RT-PCR

# Gene Expression Omnibus (GEO)

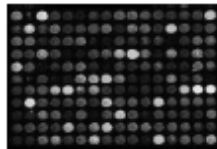
Submitted by  
Manufacturer\*

**GPL**  
Platform  
descriptions



Submitted by  
Experimentalists

**GSM**  
Raw/processed  
spot intensities  
from a single  
slide/chip



Entrez GEO

**GSE**  
Grouping of  
slide/chip data  
“a single experiment”

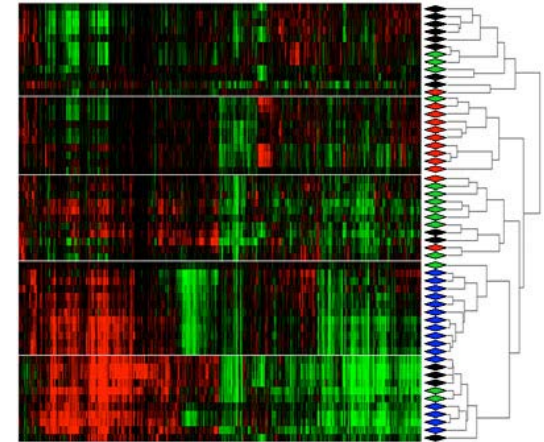


Curated by  
NCBI

**GDS**  
Grouping of  
experiments



Entrez  
GEO Datasets

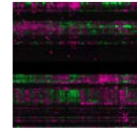




# CURATED DATASET BROWSER

Search for

Cluster Analysis



Download

- DataSet full SOFT file
- DataSet SOFT file
- Series family SOFT file
- Series family MINIML file
- Annotation SOFT file

## DataSet Record GDS4456: [Expression Profiles](#) [Data Analysis Tools](#) [Sample Subsets](#)

<b>Title:</b>	High-risk bladder cancer		
<b>Summary:</b>	Analysis of bladder cancer specimens from a high-risk population of patients who underwent radical cystectomy (Memorial Sloan-Kettering Cancer Center cohort). Results provide insight into the prediction of survival in high-risk urothelial carcinoma of the urinary bladder.		
<b>Organism:</b>	<i>Homo sapiens</i>		
<b>Platform:</b>	GPL570: [HG-U133_Plus_2] Affymetrix Human Genome U133 Plus 2.0 Array		
<b>Citations:</b>	Riester M, Taylor JM, Feifer A, Koppie T et al. Combination of a novel gene expression signature with a clinical nomogram improves the prediction of survival in high-risk bladder cancer. <i>Clin Cancer Res</i> 2012 Mar 1;18(5):1323-33. PMID: <a href="#">22228636</a> Riester M, Werner L, Bellmunt J, Selvarajah S et al. Integrative analysis of 1q23.3 copy-number gain in metastatic urothelial carcinoma. <i>Clin Cancer Res</i> 2014 Apr 1;20(7):1873-83. PMID: <a href="#">24486590</a>		
<b>Reference Series:</b>	<a href="#">GSE31684</a>	<b>Sample count:</b>	93
<b>Value type:</b>	transformed count	<b>Series published:</b>	2012/01/23

### Data Analysis Tools

Find genes [?](#)

Compare 2 sets of samples

Cluster heatmaps

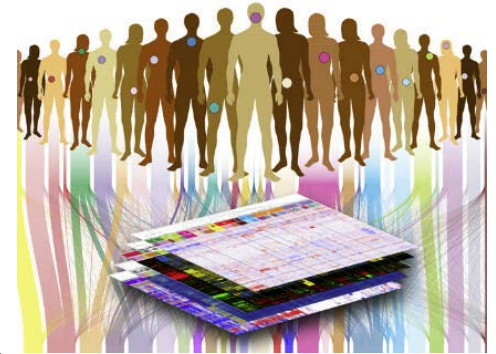
Find gene name or symbol:

Find genes that are up/down for this condition(s):  
 specimen  
 disease state

# TCGA: The Cancer Genome Atlas Program

- The Cancer Genome Atlas (TCGA), a landmark cancer genomics program, molecularly characterized over **20,000 primary cancer** and **matched normal samples** spanning **33 cancer types**. This joint effort between **NCI** and the **National Human Genome** Research Institute began in 2006, bringing together researchers from diverse disciplines and multiple institutions

- Over the next dozen years, TCGA generated over **2.5 petabytes** ( $2^{50}$ ) of **genomic**, **epigenomic**, **transcriptomic**, and **proteomic data**. The data, which has already led to improvements in our ability to **diagnose**, **treat**, and **prevent cancer**, will remain publicly available for anyone in the research community to use.





## Webinar 3: Expression Data Analysis & OQL

May 14, 2020



0:00 / 57:59



cBioPortal Webinar 3: Expression Data Analysis

SnapShot

# SnapShot: TCGA-Analyzed Tumors

Amy Blum, Peggy Wang, Jean C. Zenklusen

[Show more](#) ▾

Cancer type	Prevalence	TCGA cases assessed	Key findings
Breast lobular carcinoma	3,327,552	203	FOXA1 elevated in lobular carcinoma, GATA3 in ductal carcinoma; lobular enriched for PTEN loss and Akt activation
Breast ductal carcinoma		784	Four distinct subtypes: basal, Her2, luminal A, and luminal B; most common driver mutations: TP53, PIK3CA, GATA3; basal subtype similar to serous ovarian cancer
Prostate cancer	3,085,209	333	Highly heterogeneous with 26% driven by unknown alterations; ETS gene fusions or mutations in SPOP, FOXA1, or IDH1 define seven subtypes; actionable lesions in PI3K, MAPK, and DNA repair pathways
Colorectal adenocarcinoma	1,317,247	276	Colon and rectal cancers have similar genomic profiles; hypermutated subtype associated with favorable prognosis; new potential drivers: ARID1A, SOX9, FAM123B/WTX
Cutaneous melanoma	1,169,351	331	Established four subtypes: BRAF mutant, RAS mutant, NF1 mutant, and triple wild-type based on driver mutations; higher levels of immune lymphocyte infiltration correlated with better survival
Thyroid carcinoma	726,646	496	Majority driven by RAS or BRAFV600E mutations
Endometrial carcinoma	710,228	373	Classified endometrial cancers into four categories: POLE ultramutated, MSI (microsatellite instability) hypermutated, copy-number low, copy-number high
Uterine carcinosarcoma		57	Strong and varied degree of epithelial-mesenchymal transition; TP53 mutations in 91% of samples; PI3K alterations in half

## TCGA

Program History +

### TCGA Cancers Selected for Study

Publications by TCGA

Using TCGA +

Contact

## TCGA Cancers Selected for Study

The Cancer Genome Atlas (TCGA) selected the following cancers for study based on specific criteria that include:

- Poor prognosis
- Overall public health impact
- Availability of samples meeting standards for patient consent
- Availability of samples meeting standards for quality and quantity that include:
  - Primary, untreated tumor with a source of matched normal tissue or blood sample
  - Frozen, sufficiently sized, resection samples
  - Samples composed of at least 80% tumor nuclei (threshold later lowered to 60% with improved sequencing technology and computational methods)



TCGA cancers selected for study

Credit: National Cancer Institute

## TCGA's Cancers Selected for Study

Cancer Type Studied	# Cases Characterized (# Cases in Marker Paper)	Publication
Acute Myeloid Leukemia	200 (200)	NEJM 2013 <a href="#">↗</a>
Adrenocortical Carcinoma	92 (91)	Cancer Cell 2016 <a href="#">↗</a>
Bladder Urothelial Carcinoma	412 (131)	Nature 2014 <a href="#">↗</a> , Cell 2017 <a href="#">↗</a>
Breast Ductal Carcinoma	778 (430)	Nature 2012 <a href="#">↗</a>
Breast Lobular Carcinoma	201 (127)	Cell 2015 <a href="#">↗</a>
Cervical Carcinoma	307 (228)	Nature 2017 <a href="#">↗</a>
Cholangiocarcinoma	51 (38)	Cell Reports 2017 <a href="#">↗</a>
Colorectal Adenocarcinoma	633 (276)	Nature 2012 <a href="#">↗</a>
Esophageal Carcinoma	185 (164)	Nature 2017 <a href="#">↗</a>
Gastric Adenocarcinoma	443 (295)	Nature 2014 <a href="#">↗</a>
Glioblastoma Multiforme	617 (206)	Nature 2008 <a href="#">↗</a> , Cell 2013 <a href="#">↗</a>
Head and Neck Squamous Cell Carcinoma	528 (279)	Nature 2015 <a href="#">↗</a>





Analyze, Integrate, Discover

Home

Tutorial

TCGA


CPTAC

CBTTC

## Welcome to UALCAN

The **U**niversity of **AL**abama at Birmingham **CAN**cer data analysis Portal

UALCAN is a comprehensive, user-friendly, and interactive web resource for analyzing cancer OMICS data. It is built on PERL-CGI with high quality graphics using javascript and CSS. UALCAN is designed to, a) provide easy access to publicly available cancer OMICS data (TCGA, MET500, CPTAC and CBTTC), b) allow users to identify biomarkers or to perform in silico validation of potential genes of interest, c) provide graphs and plots depicting expression profile and patient survival information for protein-coding, miRNA-coding and lincRNA-coding genes, d) evaluate epigenetic regulation of gene expression by promoter methylation, e) perform pan-cancer gene expression analysis, f) Provide additional information about the selected genes/targets by linking to HPRD, GeneCards, Pubmed, TargetScan, The human protein atlas, DRUGBANK, Open Targets and the GTEx. These resources allow researchers to gather valuable information and data about the genes/targets of interest, g) provide clinical proteomic consortium data analysis including total/phospho-proteins and h) provide pediatric brain tumor gene expression and protein expression analysis.

Follow us on Twitter 

### Please cite:

- 1) Chandrashekar DS, Karthikeyan SK, Korla PK, Patel H, Shovon AR, Athar M, Netto GJ, Qin ZS, Kumar S, Manne U, Creighton CJ, Varambally S. UALCAN: An update to the integrated cancer data analysis platform. *Neoplasia*. 2022 Mar;25:18-27. doi: 10.1016/j.neo.2022.01.001 [PMID: 35078134]
- 2) Chandrashekar DS, Bachel B, Balasubramanya SAH, Creighton CJ, Rodriguez IP, Chakravarthi BVSK and Varambally S. UALCAN: A portal for facilitating tumor subgroup gene expression and survival analyses. *Neoplasia*. 2017 Aug;19(8):649-658. doi: 10.1016/j.neo.2017.05.002 [PMID:28732212]

Updated slides from a presentation about UALCAN cancer database at Dana-Farber Cancer Institute [Download](#)

## UALCAN: A Portal for Facilitating Tumor Subgroup Gene Expression and Survival Analyses<sup>1</sup>

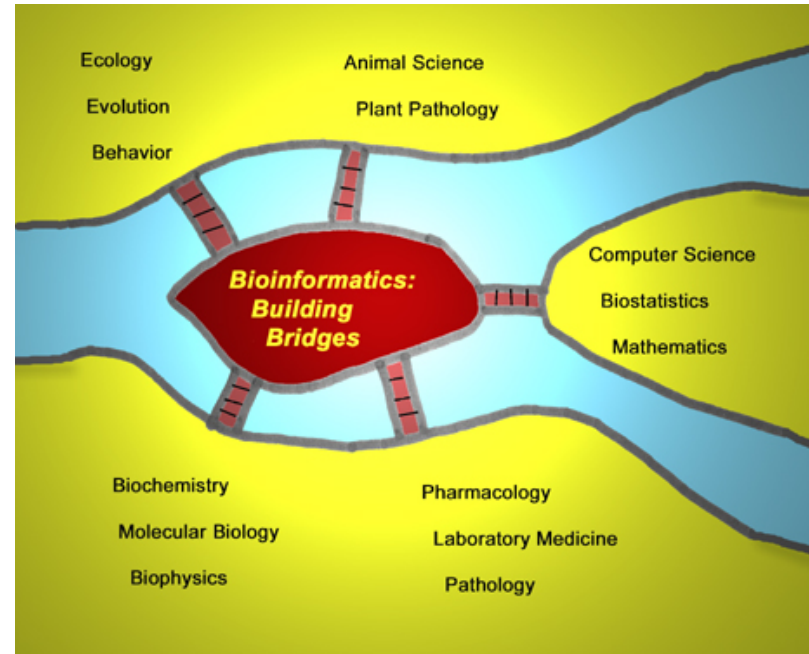


Darshan S. Chandrashekar<sup>\*,†</sup>, Bhuwan Bachel<sup>\*</sup>,  
Sai Akshaya Hodigere Balasubramanya<sup>\*,†</sup>,  
Chad J. Creighton<sup>†</sup>, Israel Ponce-Rodriguez<sup>\*</sup>,  
Balabhadrapatruni V.S.K. Chakravarthi<sup>\*,†</sup> and  
Sooryanarayana Varambally<sup>\*,†</sup>

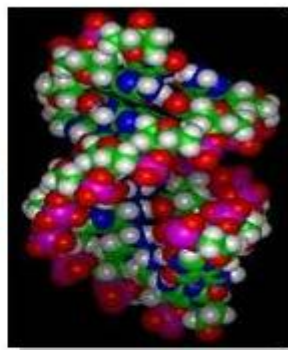
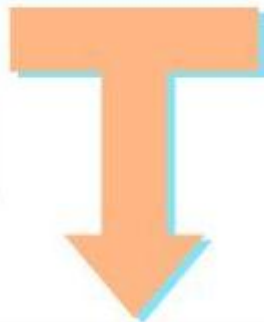
<sup>\*</sup>Molecular and Cellular Pathology, Department of Pathology, University of Alabama at Birmingham; <sup>†</sup>Comprehensive Cancer Center, University of Alabama at Birmingham, Birmingham, AL 35233, USA; <sup>‡</sup>Department of Medicine, Dan L. Duncan Comprehensive Cancer Center, and Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX 77030, USA

• Genomics data from The Cancer Genome Atlas (TCGA) project has led to the comprehensive molecular characterization of **multiple cancer types**. The large sample numbers in TCGA offer an excellent opportunity to address questions associated with **tumor heterogeneity**. Exploration of the data by cancer researchers and clinicians is imperative to unearth novel therapeutic/diagnostic biomarkers. Various computational tools have been developed to aid researchers in carrying out **specific TCGA data analyses**; however there is need for resources to facilitate the study of gene expression variations and survival associations across tumors. Here, we report UALCAN, an easy to use, interactive web-portal to perform to **in-depth analyses of TCGA gene expression data**. UALCAN uses **TCGA level 3 RNA-seq** and clinical data from **31 cancer types**. The portal's user-friendly features allow to perform: 1) analyze **relative expression of a query gene(s) across tumor and normal samples**, as well as in **various tumor sub-groups** based on individual cancer stages, tumor grade, race, body weight or other clinicopathologic features, 2) estimate the effect of **gene expression level and clinicopathologic features on patient survival**; and 3) identify the **top over- and under-expressed (up and down-regulated) genes** in individual cancer types. This resource serves as a platform for in silico validation of target genes and for identifying tumor subgroup specific candidate biomarkers. Thus, UALCAN web-portal could be extremely helpful in accelerating cancer research. UALCAN is publicly available at <http://ualcan.path.uab.edu>.

# Q & A



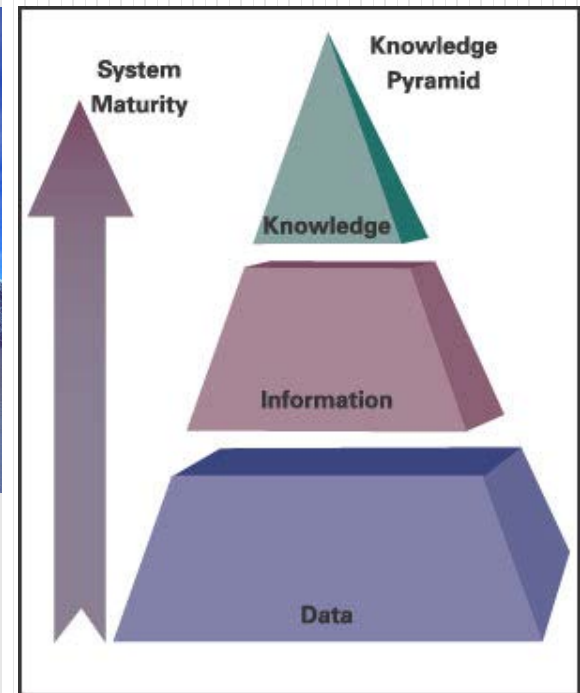
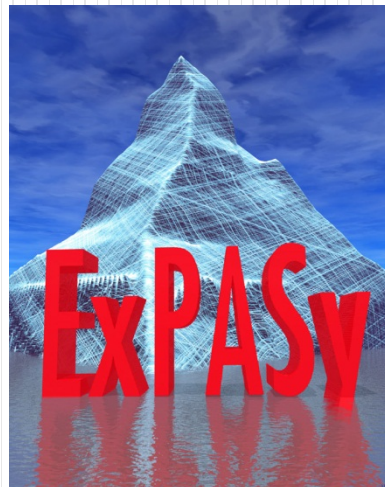
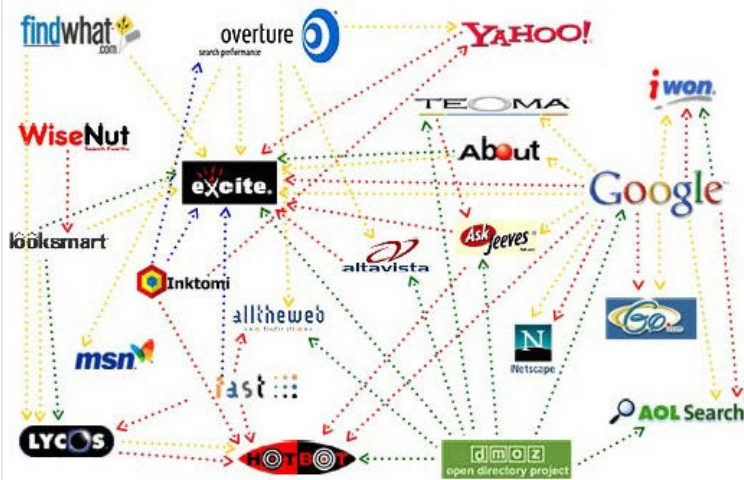
Computer systems



Biological systems

BIOINFORMATICS

# Q: How to Find the Right Stuffs?





Query all databases ▼

search

## Visual Guidance

## Categories

- proteomics
- genomics
- structural bioinformatics
- systems biology
- phylogeny/evolution
- population genetics
- transcriptomics
- biophysics
- imaging
- IT infrastructure
- drug design

## Resources A..Z

## Links/Documentation

ExPASy is the **SIB Bioinformatics Resource Portal** which provides access to scientific databases and software tools (i.e., *resources*) in different areas of life sciences including proteomics, genomics, phylogeny, systems biology, population genetics, transcriptomics etc. (see **Categories** in the left menu). On this portal you find resources from many different SIB groups as well as external institutions.

### Featuring today

#### STRING

Database of known and predicted protein-protein interactions  
[\[details\]](#)



# How to Find the Right Stuffs

Google Algorithm: **PageRank™**

PDF, 庫存頁面...

Askcom **ExpertRank** algorithm

Subject-specific popularity

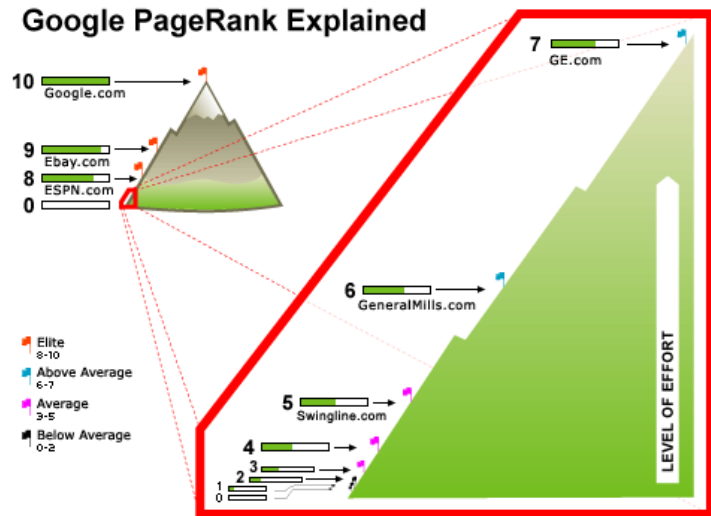
Use **the right key words**

PubMed: MeSH

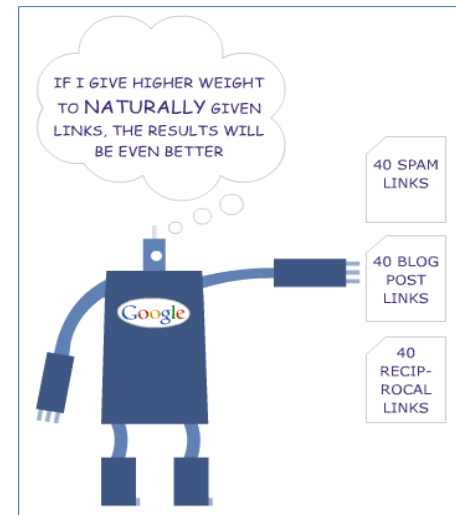
OMIM: index

Gene name: HUGO

Fidelity: edu > gov > org > com



©2007 Elliance, Inc.



# Search Efficiently

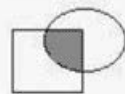
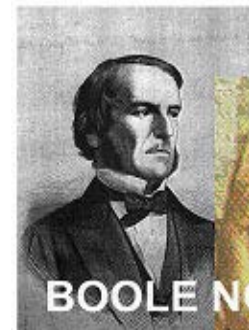
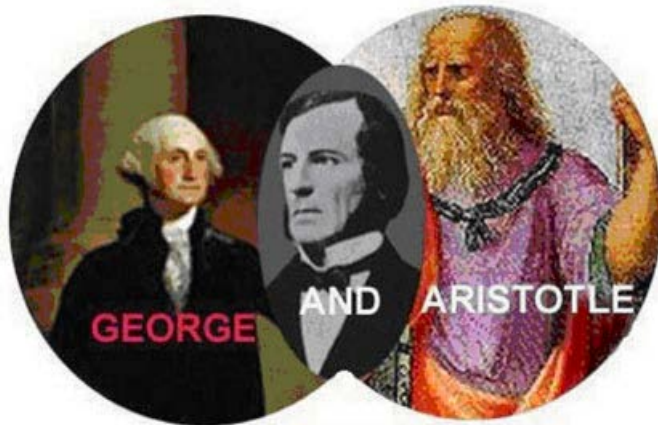
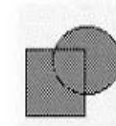
[Quick  
Tours](#)

[Search PubMed by Authors](#)

[My NCBI...](#)

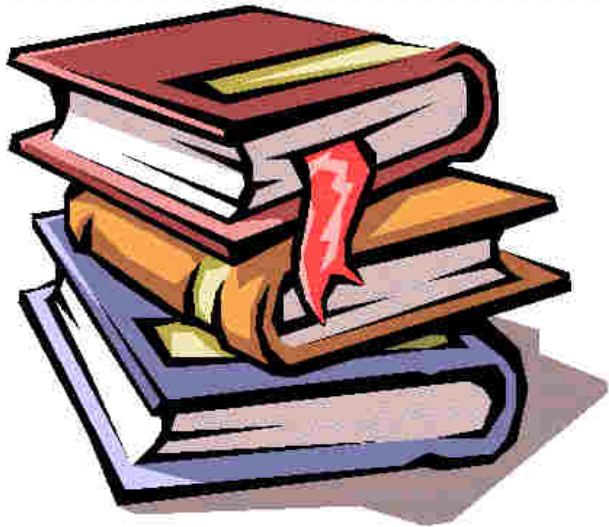
[Boolean  
operators](#)

**AND OR NOT**



This is a small search.  
Your results will  
include *both* words.

**Q: How to Find References Related  
to Your Favorite Gene (YFG)**





# Gene or Disease – Official Symbol



1----- (100000- )	Autosomal loci or phenotypes (entries created before May 15, 1994)
2----- (200000- )	
3----- (300000- )	X-linked loci or phenotypes
4----- (400000- )	Y-linked loci or phenotypes
5----- (500000- )	Mitochondrial loci or phenotypes
6----- (600000- )	Autosomal loci or phenotypes (entries created after May 15, 1994)



PubMed

- POU5F1



<b>POU5F1P8</b>	POU class 5 homeobox 1 pseudogene 8	Homo sapiens
<b>Pou5f1</b>	POU domain, class 5, transcription factor 1	Mus musculus
<b>Pou5f1-rs1</b>	POU domain, class 5, transcription factor 1, related sequence 1	Mus musculus
<b>Pou5f1-rs10</b>	POU domain, class 5, transcription factor 1, related sequence 10	Mus musculus
<b>Pou5f1-rs2</b>	POU domain, class 5, transcription factor 1, related sequence 2	Mus musculus
<b>Pou5f1-rs3</b>	POU domain, class 5, transcription factor 1, related sequence 3	Mus musculus
<b>Pou5f1-rs4</b>	POU domain, class 5, transcription factor 1, related sequence 4	Mus musculus
<b>Pou5f1-rs5</b>	POU domain, class 5, transcription factor 1, related sequence 5	Mus musculus
<b>Pou5f1-rs6</b>	POU domain, class 5, transcription factor 1, related sequence 6	Mus musculus
<b>Pou5f1-rs8</b>	POU domain, class 5, transcription factor 1, related sequence 8	Mus musculus
<b>Pou5f1-rs9</b>	POU domain, class 5, transcription factor 1, related sequence 9	Mus musculus
<b>Pou5f2</b>	POU domain class 5, transcription factor 2	Mus musculus
<b>POU5F1</b>	POU class 5 homeobox 1	Sus scrofa
<b>pou5f1</b>	POU domain, class 5, transcription factor 1	Danio rerio
<b>POU5F1</b>	POU class 5 homeobox 1	Pan troglodytes
<b>POU5F1</b>	POU class 5 homeobox 1	Pan troglodytes
<b>POU5F2</b>	POU domain class 5, transcription factor 2	Pan troglodytes

OMIM

- Preview and index

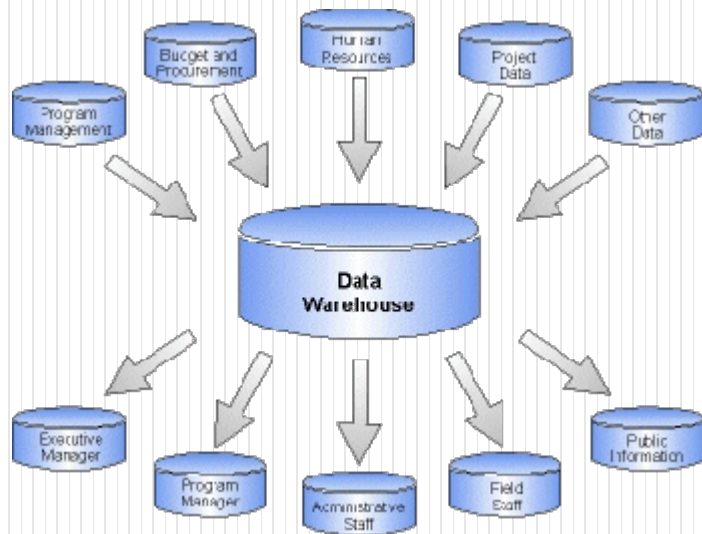
GeneCards/

human POU5F1 (symbol only)

Entrez Gene

- POU5F1

# Q: What is Derivative Databases?



# Leading Bioinformatic Centers

## NCBI, USA

- To develop **new methods** for integrative, **computer-based data analysis** to mine massive and complex **data sets**

## EBI, UK

- The EBI is a centre for **research** and **services** in **bioinformatics**
- The Institute manages **databases** of **biological data** including **nucleic acid, protein sequences & macromolecular structures**

## Tutorials

Training materials in HTML, PDF and Video formats

Filter this table

Type	Title and Description
Video	<b>A Guide to NCBI: Gene Expression, Part 1</b> Part 1 of the gene expression module from "A Librarian's Guide to NCBI," a workshop held at the National Library of Medicine in April 2013
Video	<b>A Guide to NCBI: Gene Expression, Part 2</b> Part 2 of the gene expression module from "A Librarian's Guide to NCBI," a workshop held at the National Library of Medicine in April 2013
Video	<b>A Guide to NCBI: Gene Expression, Part 3</b> Part 3 of the gene expression module from "A Librarian's Guide to NCBI," a workshop held at the National Library of Medicine in April 2013
PDF	<b>Align 2 Sequences</b> Aligning two groups of sequences and displaying the results in the NCBI sequence viewer
Video	<b>Assign Downloaders for dbGaP Data</b> Learn how an authorized user of controlled-access data can assign a downloader role to someone in his/her institution

## Online courses

Start now

[ArrayExpress: Discover functional genomics data quickly and easily](#)

Author: Anja Füllgrabe

ArrayExpress is a database of functional genomics data. This course will give you an overview of how these data are stored in ArrayExpress and will teach you how to effectively search and retrieve data from the [ArrayExpress website](#). [...]

Start now

[ArrayExpress: Quick tour](#)

Author: Melissa Burke

This quick tour provides an overview of EMBL-EBI's functional genomics database ArrayExpress. [...]

Start now

[Biocuration: An introduction](#)

Author:

Claire O'Donovan, leader of the Protein Function Content team at EMBL-EBI, gives an introduction into biocuration and talks about what it is like to work as a biocurator and the skill sets you need.[...]

# The National Center for Biotechnology Information (NCBI)

Founded **1988**

**NCBI** The **leading** American information provider; a division of the National Library of Medicine (NLM), NIH (Bethesda, USA)

**Roles** To develop **new information technologies** to aid our understanding of the **molecular** and **genetic processes** that underlie **health and disease**



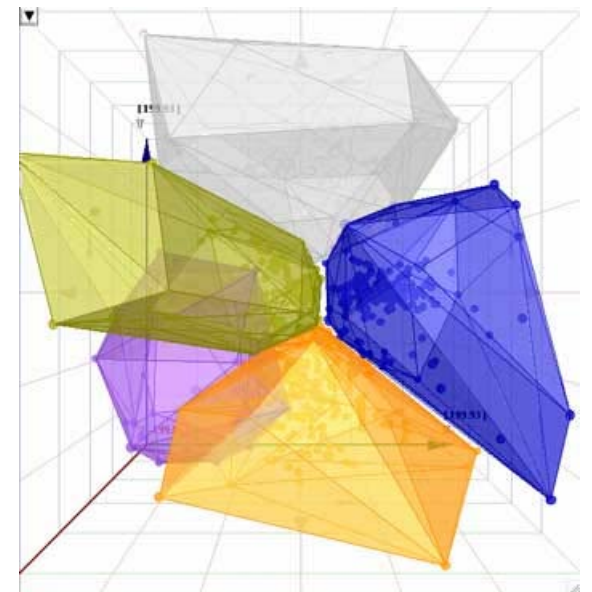
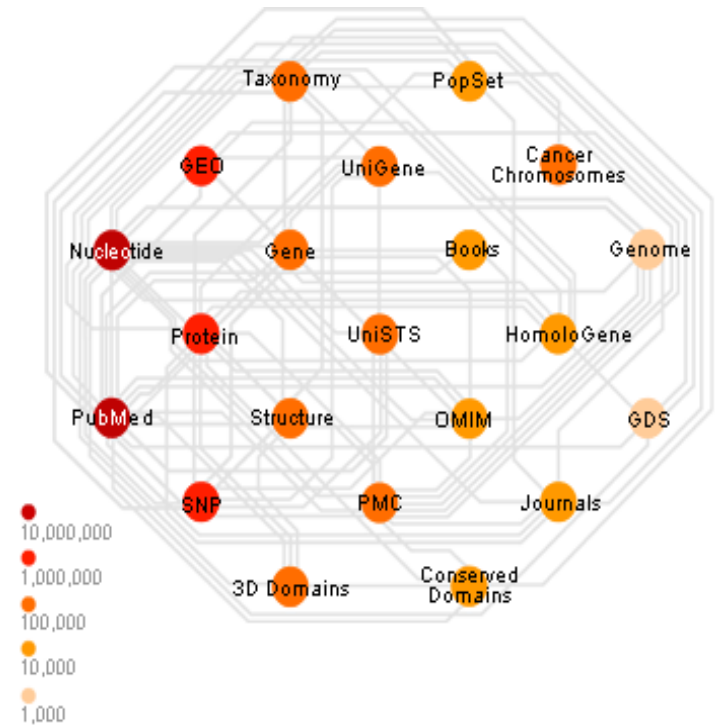
# Contents

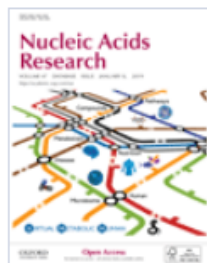
## Databases

- **Primary vs. derivative** databases
- **Value-added**

## Methodologies (tools)

- Tools: e.g., [BLAST](#), [NCBI](#)
- **Algorithms**
  - Neural network (NN)
    - Self-organizing map (SOM)
  - **Hidden Markov Model (HMM)**
  - K-means clustering





Volume 47, Issue D1  
08 January 2019

## Article Contents

Abstract

NEW AND UPDATED DATABASES

NAR ONLINE MOLECULAR  
BIOLOGY DATABASE  
COLLECTION


ACKNOWLEDGEMENTS

FUNDING

REFERENCES

—

## The 26th annual Nucleic Acids Research database issue and Molecular Biology Database Collection


Daniel J Rigden , Xosé M Fernández


*Nucleic Acids Research*, Volume 47, Issue D1, 08 January 2019, Pages D1–D7,  
<https://doi.org/10.1093/nar/gky1267>


**Published:** 29 December 2018




PDF

 Split View

 Cite

 Permissions

 Share ▼

### Abstract

The 2019 Nucleic Acids Research (NAR) Database Issue contains 168 papers spanning molecular biology. Among them, 64 are new and another 92 are updates describing resources that appeared in the Issue previously. The remaining 12 are updates on databases most recently published elsewhere. This Issue contains two Breakthrough articles, on the Virtual Metabolic Human (VMH) database which links human and gut microbiota metabolism with diet and disease, and Vibrism DB, a database of mouse brain anatomy and gene (co-)expression with sophisticated visualization and session sharing.

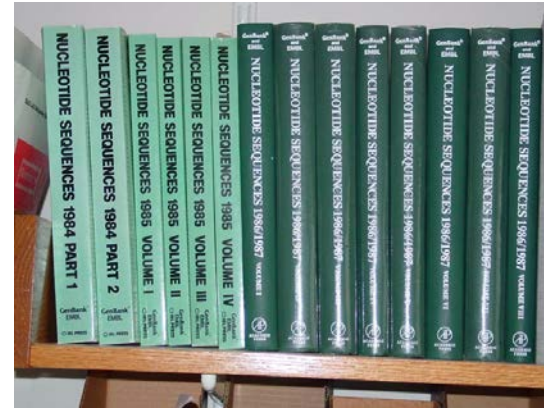
# Primary vs. Derivative Databases - NCBI

## Primary databases

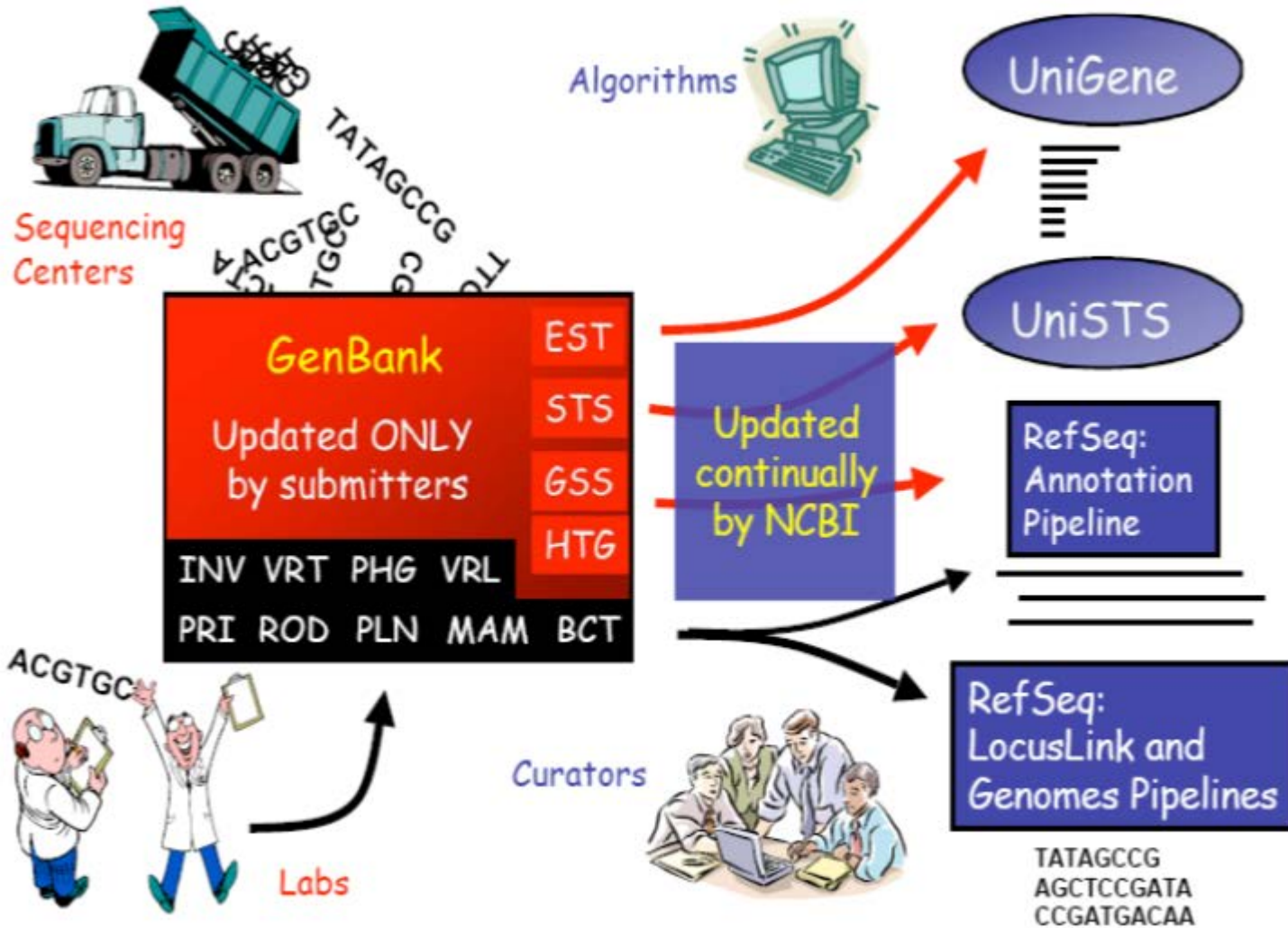
- **Original** submissions by **experimentalists**
- **Submitters** retain editorial control of records
- Archival in nature

## Derivative databases

- **Curated** by NCBI staffs
- **NCBI** retains **editorial control** of records
- Record content is **updated continually**

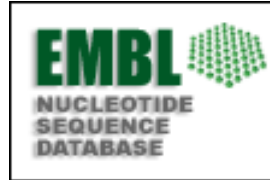


# Primary vs. Derivative Databases - NCBI





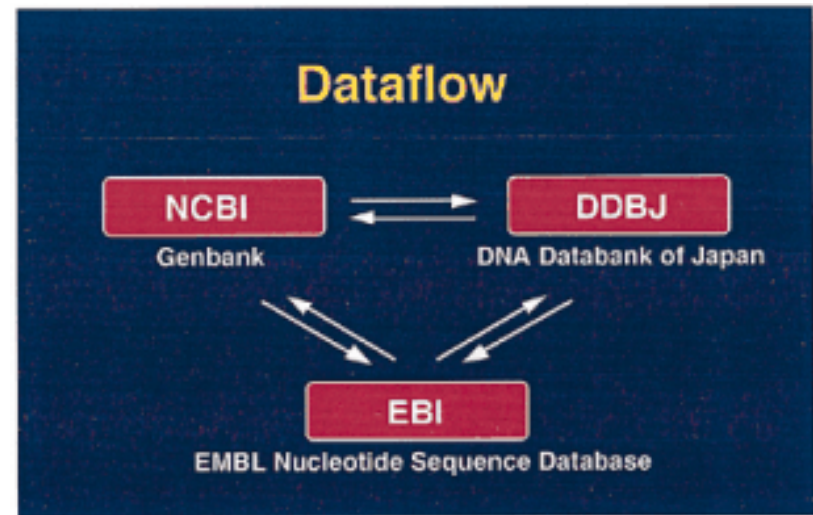
# Primary DNA Databases



GenBank (USA)

EMBL (Europe)

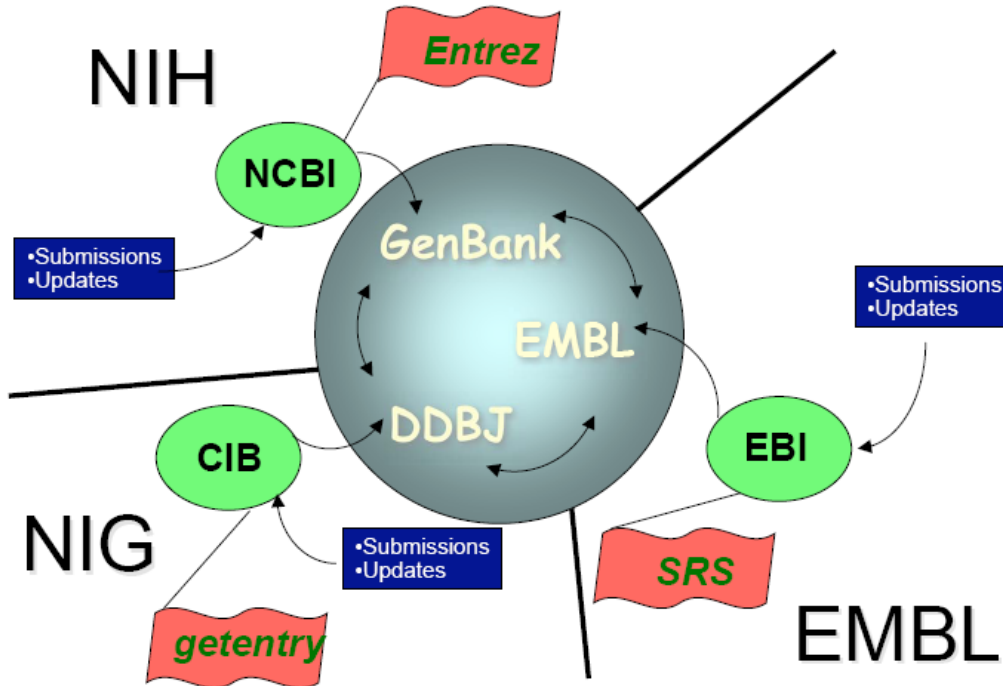
DDBJ (Japan)



National Institute of Health (**NIH**)

National Center for Biotechnology (**NCBI**)

Retrieval System Across all Databases in NCBI (**ENTREZ**)



National Institute of Genetics (NIG)

Center for Information Biology (CIB)



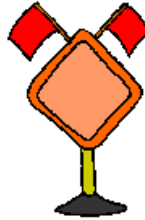
The European Bioinformatics Institute (**EBI**)

Sequence Retrieval System (**SRS**)

The *European Molecular Biology Laboratory* (**EMBL**)

# EMBL/GenBank/DDDBJ Annotations

Warning!!!



---

DNA  
data  
base  
annot  
ations  
are  
**full of  
errors**

In sequences, in annotations, in  
CDs attribution...

---

**No consistency** of annotations

---

Most annotations are done by the  
**submitters**

---

**Heterogeneity** of quality and  
updating

---



# Some Interesting Sequence Annotation

```
FT source      1..124
FT             /db_xref="taxon:4097"
FT             /organelle="plastid:chloroplast"
FT             /organism="Nicotiana tabacum"
FT             /isolate="Cuban cahibo cigar, gift from President Fidel
FT             Castro"
```

Or:

```
FT source      1..17084
FT             /chromosome="complete mitochondrial genome"
FT             /db_xref="taxon:9267"
FT             /organelle="mitochondrion"
FT             /organism="Didelphis virginiana"
FT             /dev_stage="adult"
FT             /isolate="fresh road killed individual"
FT             /tissue_type="liver"
```

???

# Organization of GenBank: Traditional Divisions

Records are divided into 18 Divisions.

- 12 Traditional
- 6 Bulk

## Traditional Divisions:

- Direct Submissions (Sequin and BankIt)
- Accurate
- Well characterized

PRI Primate  
PLN Plant and Fungal  
BCT Bacterial and Archeal  
INV Invertebrate  
ROD Rodent  
VRL Viral  
VRT Other Vertebrate  
MAM Mammalian  
PHG Phage  
SYN Synthetic (cloning vectors)  
ENV Environmental Samples  
UNA Unannotated

Entrez query: `gbdiv_xxx[Properties]`

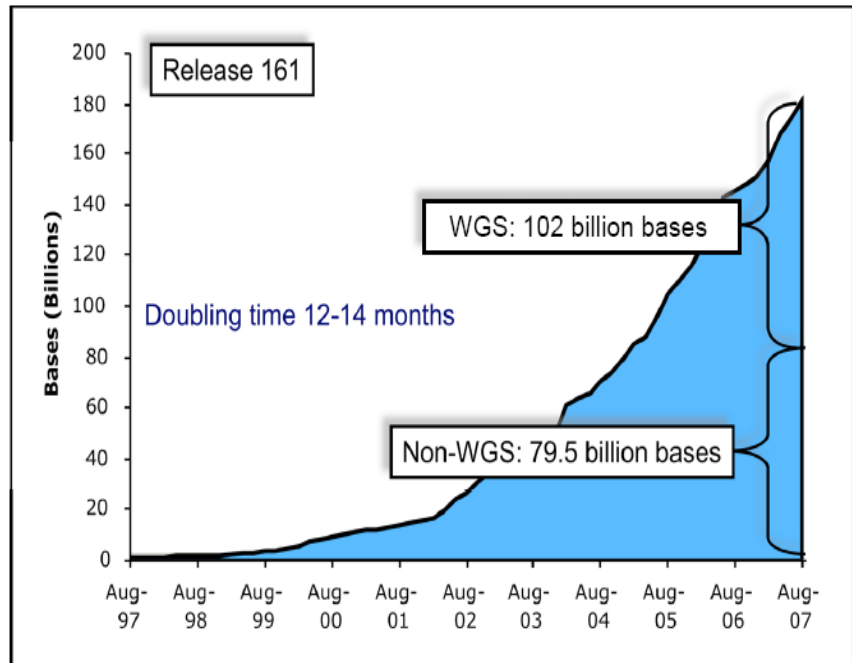
# Bulk GenBank Divisions

Batch submission & htg (email & ftp)

Inaccurate & poorly characterized

- **EST:** Expressed Sequence Tag
- **GSS:** Genome Survey Sequence
- **HTG:** High Throughput Genome
- **HTC:** High Throughput cDNA
- **STS:** Sequence Tagged Site

## The Growth of GenBank



# Organization of GenBank: Bulk Divisions

Records are divided into 18 Divisions.

- 12 Traditional
- 6 Bulk

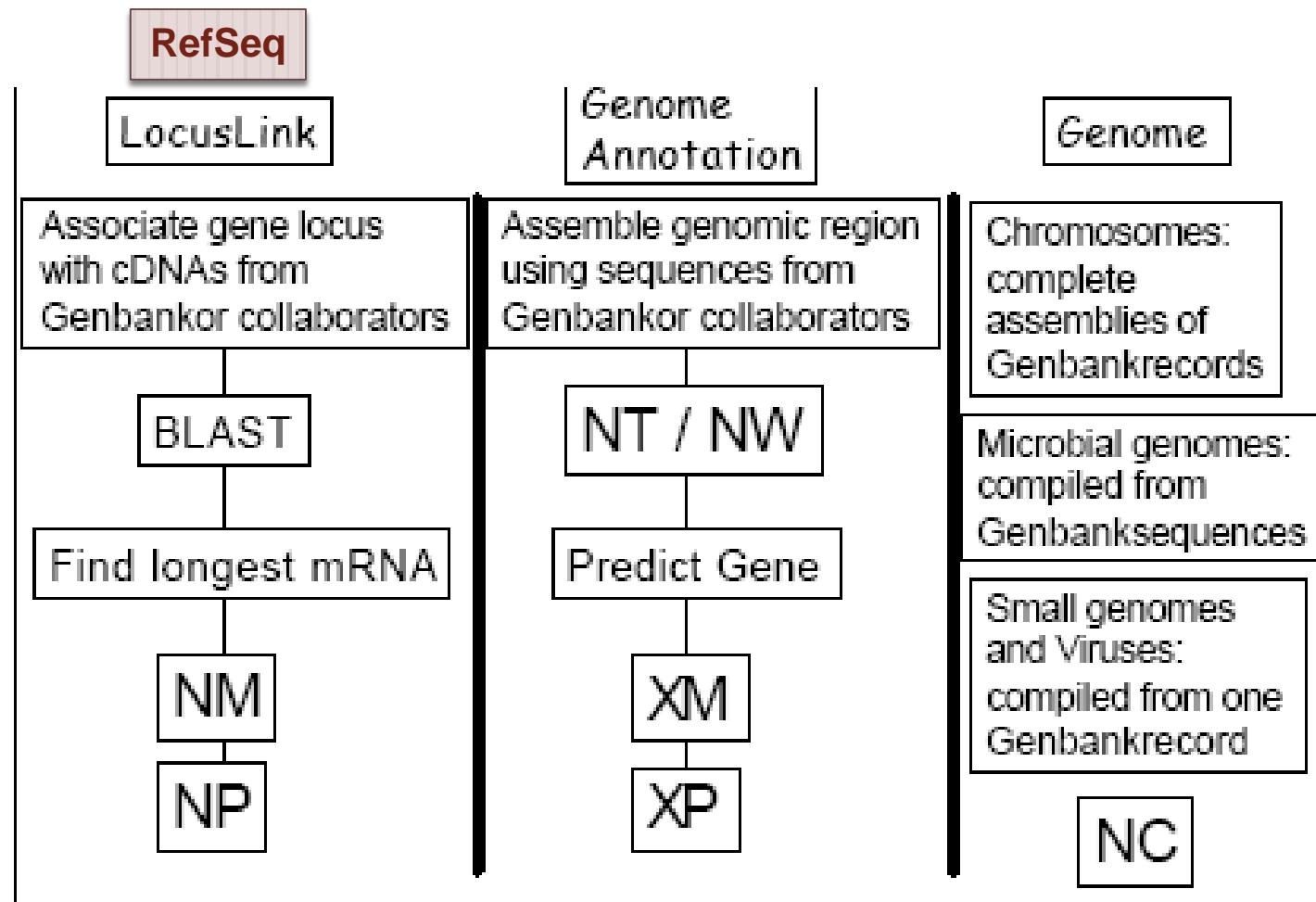
## **BULK Divisions:**

- Batch Submission  
(Email and FTP)
- Inaccurate
- Poorly characterized

EST Expressed Sequence Tag  
GSS Genome Survey Sequence  
HTG High Throughput Genomic  
STS Sequence Tagged Site  
HTC High Throughput cDNA  
PAT Patent

Entrez query: `gbdiv_xxx[Properties]`

# RefSeq Pipelines





# Selected RefSeq Accession Number

## mRNAs and Proteins

NM\_123456

Curated mRNA

NP\_123456

Curated Protein

NR\_123456

Curated non-coding RNA

XM\_123456

Predicted mRNA

XP\_123456

Predicted Protein

XR\_123456

Predicted non-coding RNA

## Gene Records

NG\_123456

Reference Genomic Sequence

## Chromosome

NC\_123455

Microbial replicons, organelle  
genomes, human chromosomes

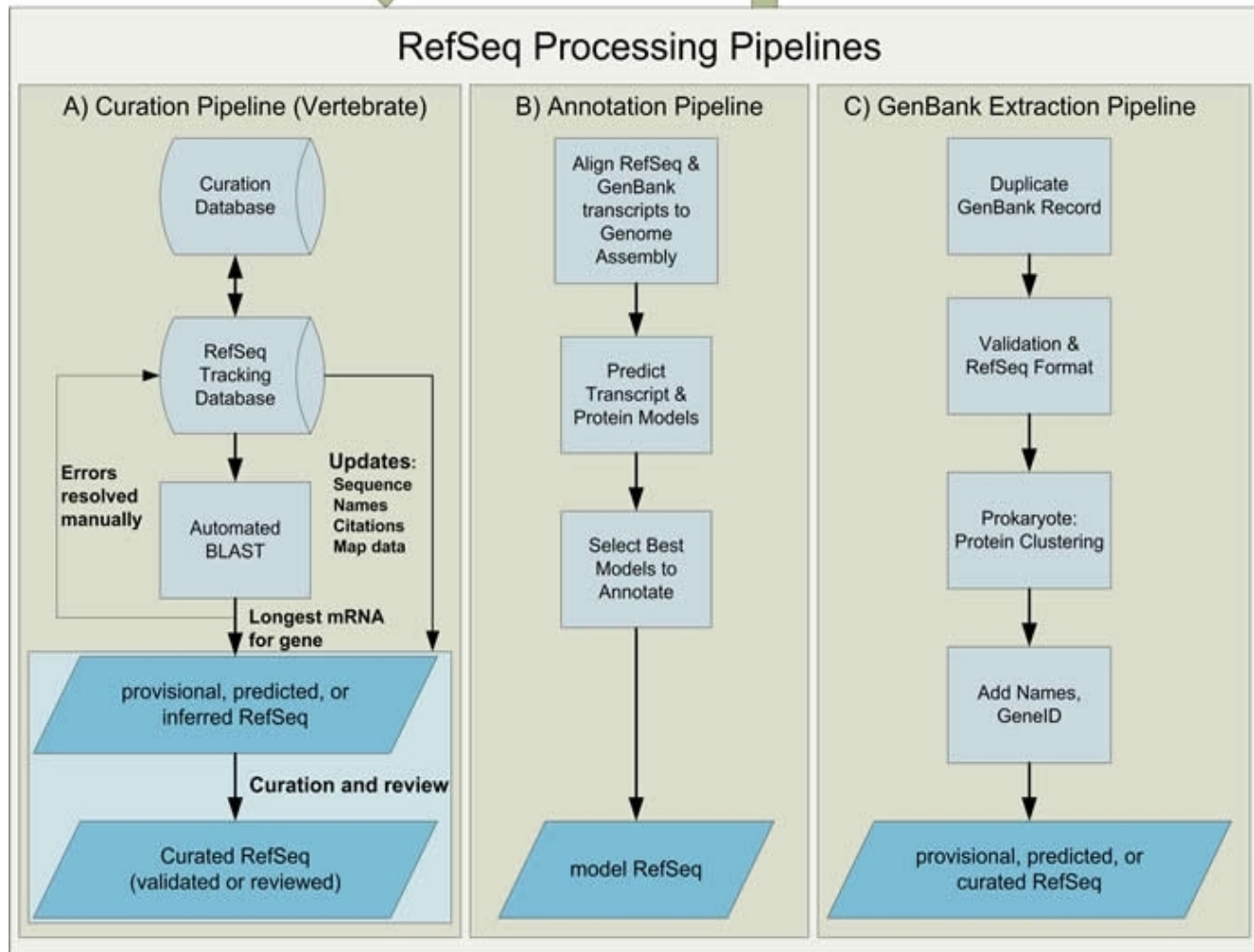
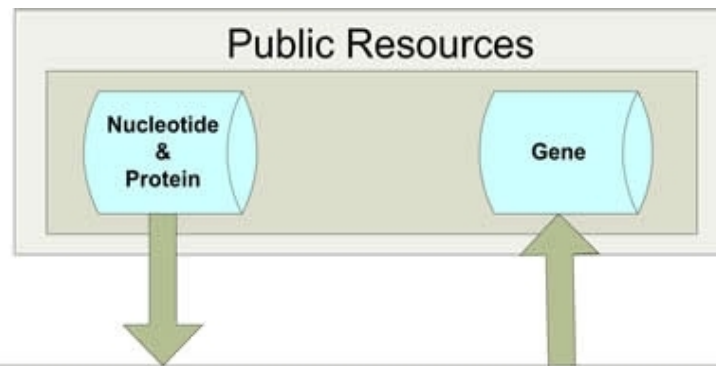
## Assemblies

NT\_123456

Contig

NW\_123456

WGS Supercontig



# RefSeq: NCBI's Derivative Sequence Database

- **Curated transcripts and proteins**
  - reviewed
  - human, mouse, rat, fruit fly, zebrafish, arabidopsis  
microbial genomes (proteins), and more
- **Model transcripts and proteins**
- **Assembled Genomic Regions (contigs)**
  - human genome      – chicken
  - mouse genome      – honeybee
  - rat genome          – sea urchin
- **Chromosome records**
  - Human genome
  - microbial
  - organelle

```
srcdb_refseq[Properties]
```

```
ftp://ftp.ncbi.nih.gov/refseq/release/
```

# RefSeq Benefits



---

**Non-redundancy**

---

**Explicitly** linked nucleotide & protein sequences

---

**Updates** to reflect current sequence data & biology

---

Data **validation**

---

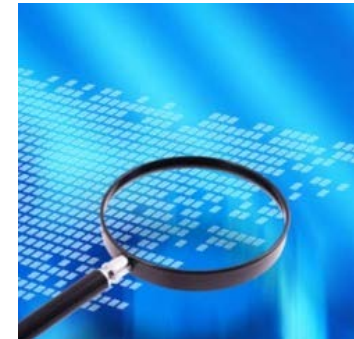
Format **consistency**

---

Distinct **accession** series

---

Stewardship by **NCBI staffs & collaborators**



interact with consistency

# Entrez Protein: Derivative Databases

## Example: CKS1B

[CDS](#)

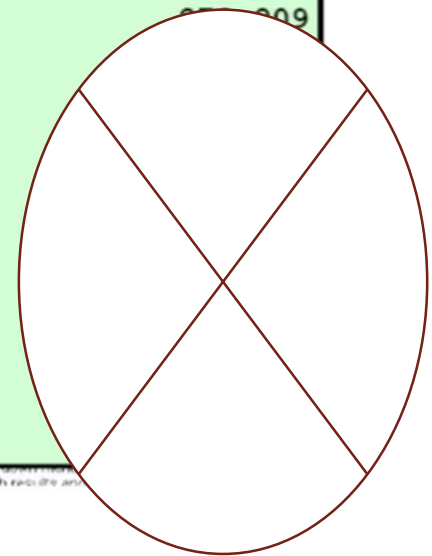
```

105..344
/gene="CKS1B"
/gene_synonym="CKS1; ckshs1; PNAS-16; PNAS-18"
/note="CDC28 protein kinase 1; CDC28 protein kinase 1B;
cell division control protein CKS1; NB4
apoptosis/differentiation related protein; PNAS-143;
CDC2-associated protein CKS1; CKS-1"
/codon_start=1
/product="cyclin-dependent kinases regulatory subunit 1"
/protein_id="NP_001817.1"
/db_xref="GI:4502857"
/db_xref="CCDS:CCDS1077.1"
/db_xref="GeneID:1163"
/db_xref="HGNC:19083"
/db_xref="HPRD:00299"
/db_xref="MIM:116900"
/translation="MSHKQIYYSDKYDDEEFYRHMVLPKDIAKLVPKTHLMSESEWR
NLGVQQSQGWVHYMIHEPEPHILLFRRLPKPKPKK"
164..291
/gene="CKS1B"
    
```



[exon](#)

Data Source	Sequences
GenPept	11,585,396
RefSeq	3,889,502
Third Party Annotation	5,263
Swiss Prot	2,009
PIR	
PRF	
PDB	
(PAT Division)	
<b>Total</b>	
<b>BLAST nr total</b>	
(no patents or env_nr -now 6 million)	



**PAT: patent**

# Search in NCBI Databases

---

**Searches** **Text:** e.g., *POU5F1* (Oct3/4);

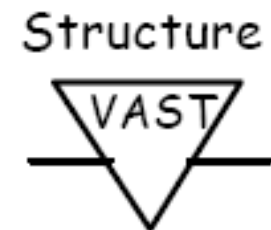
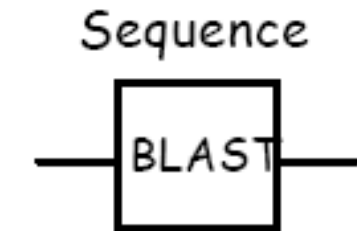
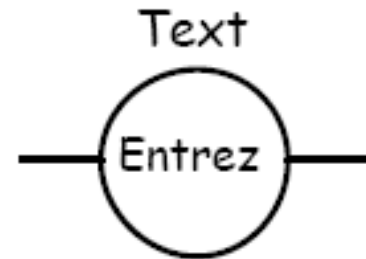
---

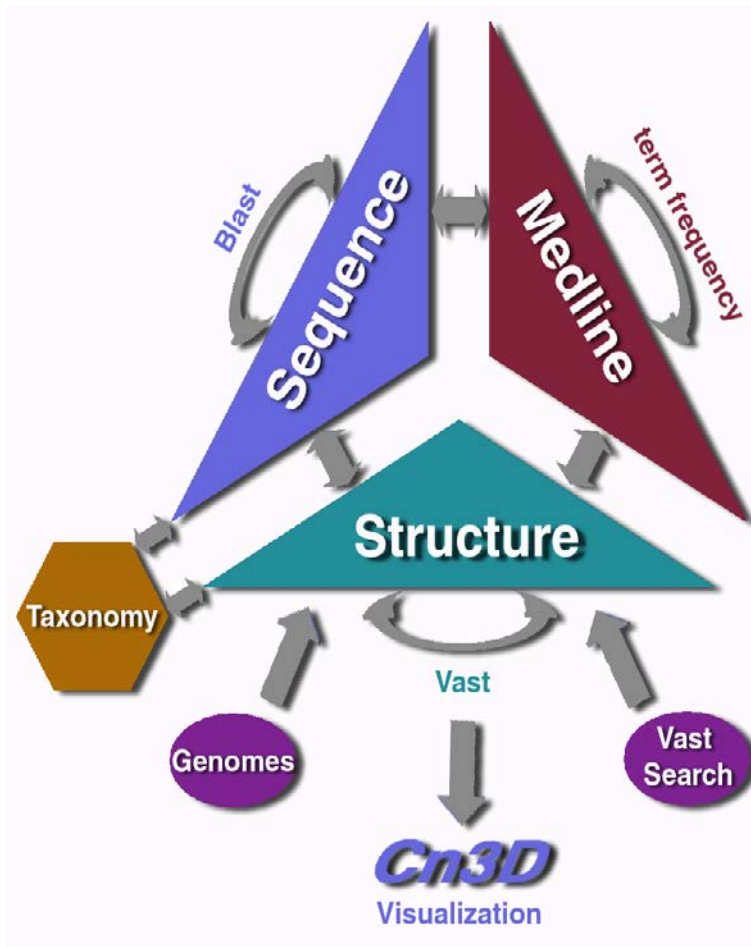
**Sequence:** e.g., [POU5F1](#)

---

**Structure:** e.g., [BRCA1](#)

---



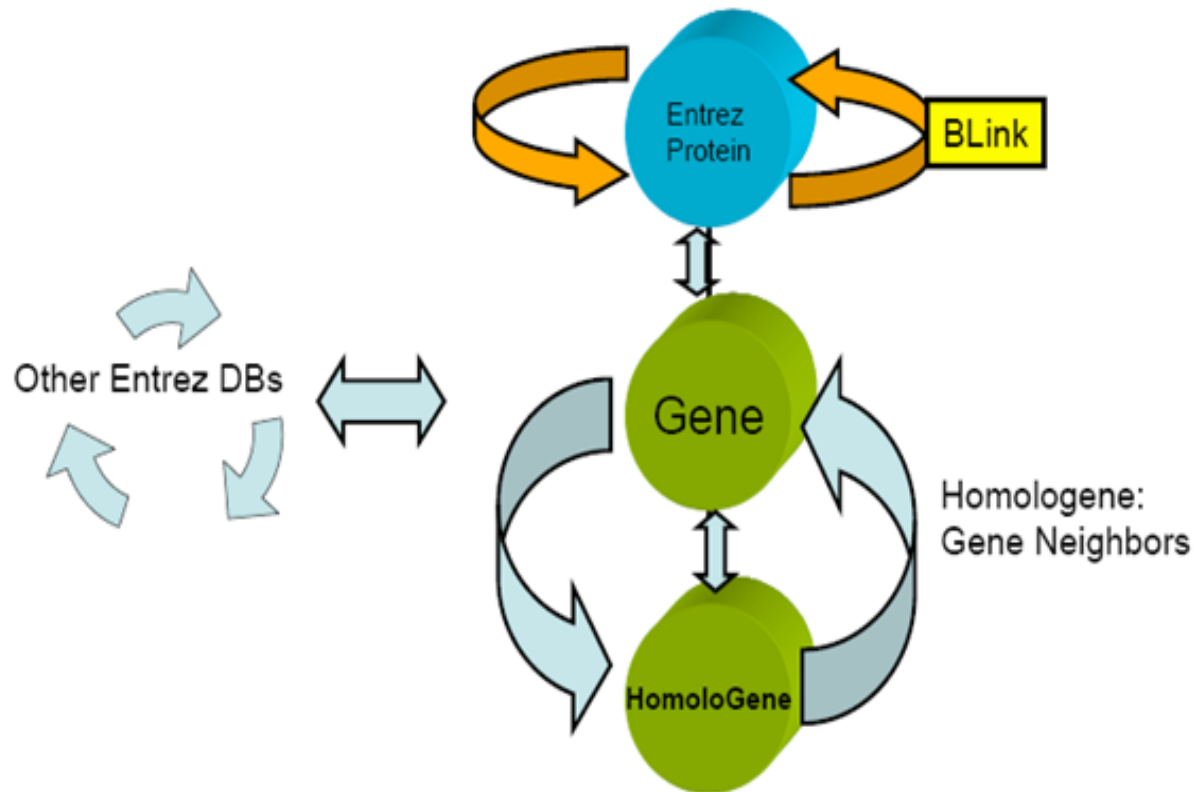


#### ▼ Links

[Explain](#)

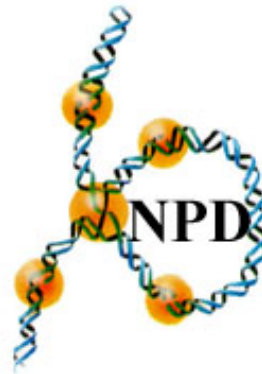
[Order cDNA clone](#)  
[Conserved Domains](#)  
[Genome](#)  
[GEO Profiles](#)  
[HomoloGene](#)  
[Map Viewer](#)  
[Nucleotide](#)  
[OMIM](#)  
[Full text in PMC](#)  
[Probe](#)  
[Protein](#)  
[PubMed](#)  
[PubMed \(OMIM\)](#)  
[PubMed \(GeneRIF\)](#)  
[SNP](#)  
[SNP: Genotype](#)  
[SNP: GeneView](#)  
[Taxonomy](#)  
[UniSTS](#)  
[AceView](#)  
[CCDS](#)  
[Ensembl](#)  
[Evidence Viewer](#)  
[HGNC](#)  
[HPRD](#)  
[KEGG](#)  
[MGC](#)  
[ModelMaker](#)  
[UniGene](#)  
[LinkOut](#)

# Entrez: Use Gene for everything





# Examples in Other Databases: Using the Official Symbol All the Time (except for protein structure)



The Nuclear Protein Database  
(e.g., TP53)



Home Genomes Blat Tables Gene Sorter PCR Session FAQ Help

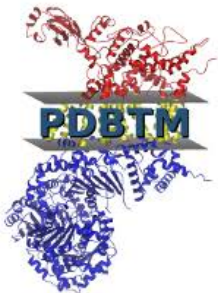
### Human (*Homo sapiens*) Genome Browser Gateway

The UCSC Genome Browser was created by the [Genome Bioinformatics Group of UC Santa Cruz](#).  
Software Copyright (c) The Regents of the University of California. All rights reserved.

clade	genome	assembly	position or search term	gene	image width	
Mammal	Human	Feb. 2009 (GRCh37/hg19)	chr6_mcf_hap5:2514038-2520393	POU5F1	800	<input type="button" value="submit"/>

[Click here to reset](#) the browser user interface settings to their defaults. [2011 ENCODE Usability Survey](#)

# Examples in Other Databases: Using the Official Symbol All the Time (except for protein structure)



## PDBTM: Protein Data Bank of Transmembrane Proteins

PDBTM version: 2019-02-22    Number of transmembrane proteins: 4084 (alpha: 3633 , beta: 427 )

all    << < 1a0s > >>

- Home
- Search
- Download
- Statistics
- Documents
- Help



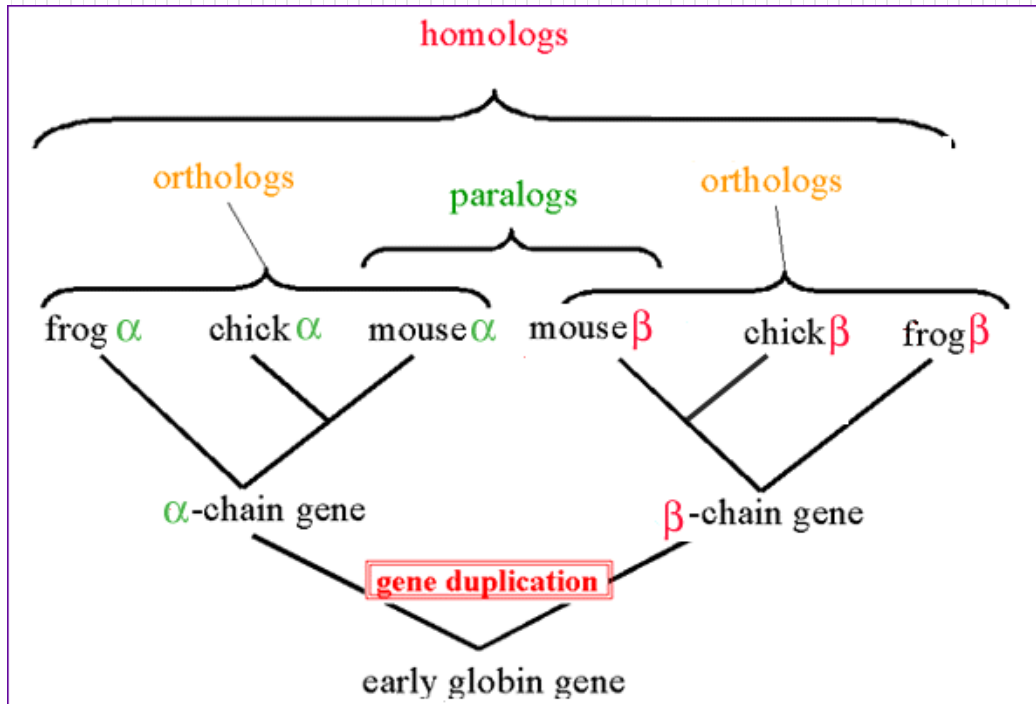
### Welcome to the PDBTM home page

PDBTM is the first comprehensive and up-to-date transmembrane protein selection of the Protein Data Bank (PDB). PDBTM database is maintained at the Institute of Enzymology by the Membrane Protein Bioinformatics Research Group. The PDBTM database was created by scanning all PDB entries with the TMDET algorithm. You can get more information about PDBTM in our articles and in the PDBTM manual. If you find PDBTM useful in your research, please cite our articles (Bioinformatics 20, 2964-2972; Nucleic Acids Research 33 Database Issue, D275-D278; Nucleic Acids Research 41 Database Issue, D524-D529 ).

6qex

PDBTM type: Tm\_Alpha  
Chain(s): A[12]

# Q: How Do You Find the Orthologs from Other Species



# Homologs (1)

## NCBI Homologene (links)

- A set of **maps** that shown **chromosomal regions** homologous between mouse, human & other species

## Example

- **POU5F1** (via ENTREZ\_GENE) **Links** to the “Homologene”
  - Protein: multiple alignment
  - Conserved domains
  - PubMed (references)
  - Protein → All links from this record → BLink

### 1: HomoloGene:8422. Gene conserved in Euteleostomi








#### Genes

Genes identified as putative homologs of one another during the construction of HomoloGene.

- POU5F1, *Homo sapiens*  
POU class 5 homeobox 1
- POU5F1L, *Pan troglodytes*  
POU domain, class 5, transcription factor 1-like
- POU5F1, *Canis lupus familiaris*  
POU class 5 homeobox 1
- POU5F1, *Bos taurus*  
POU class 5 homeobox 1
- Pou5f1, *Mus musculus*  
POU domain, class 5, transcription factor 1
- Pou5f1, *Rattus norvegicus*  
POU class 5 homeobox 1
- pou5f1, *Danio rerio*  
POU domain, class 5, transcription factor 1

#### Proteins

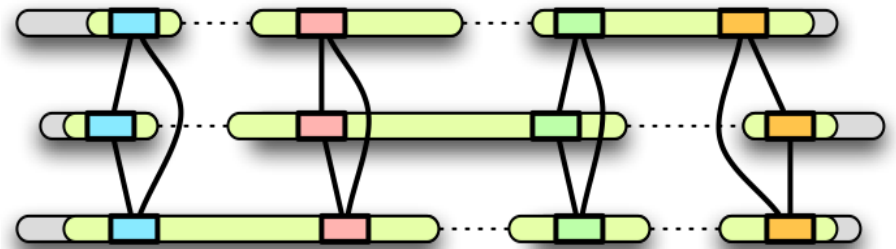
Proteins used in sequence comparisons and their conserved domain architectures.

- NP\_002692.2 360 aa 
- XP\_001135162.1 359 aa 
- XP\_538830.1 360 aa 
- NP\_777005.1 360 aa 
- NP\_038661.2 352 aa 
- NP\_001009178.1 352 aa 
- NP\_571187.1 472 aa 

# Homologs (2)

Hs and Mm links adjacent to each map name show the **mouse-human homology map** with the master chromosome as human or mouse

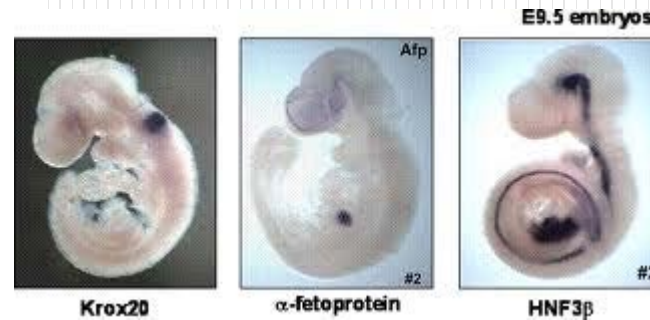
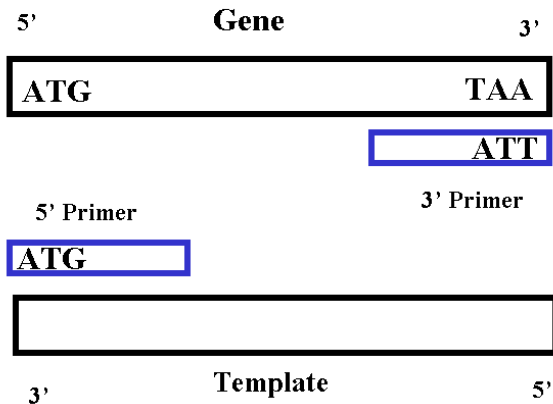
- [Mouse Genome Informatics](#)
- [Mm](#): *Pou5f1* (chr. 17; 19.23 cM)

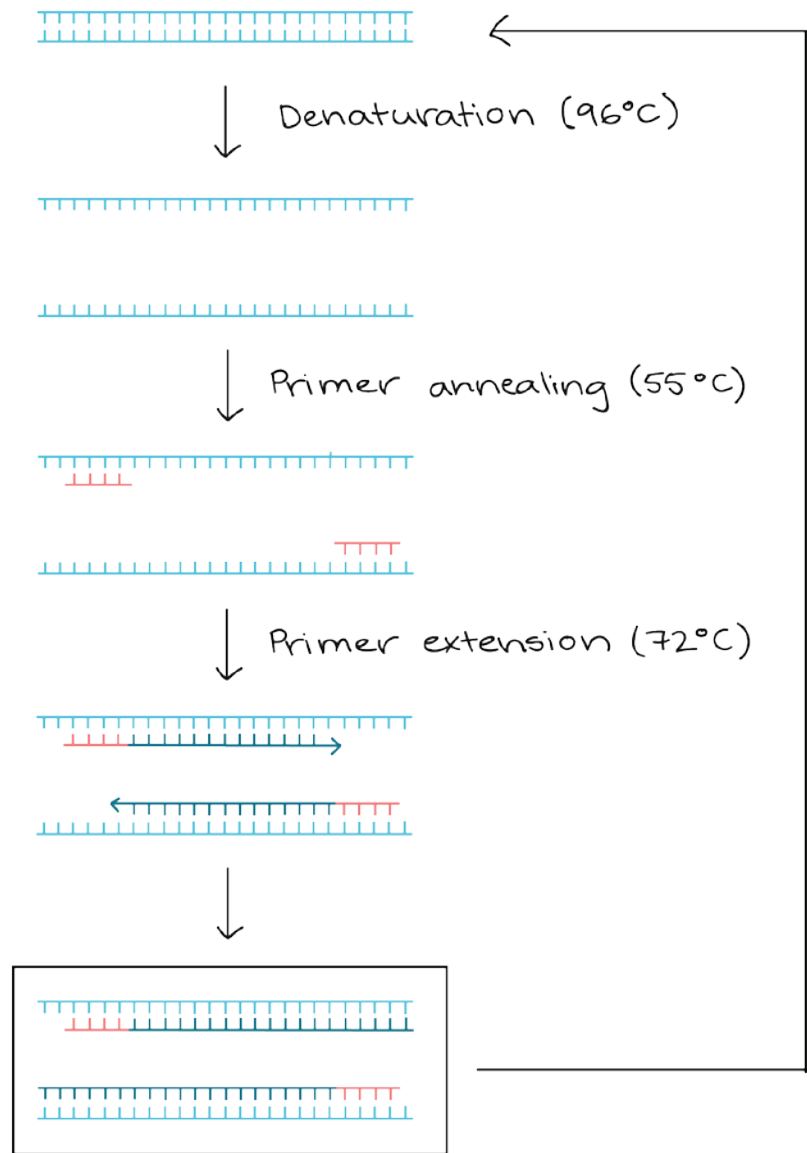


**Mercator**

Multiple Whole-Genome Orthology Map Construction

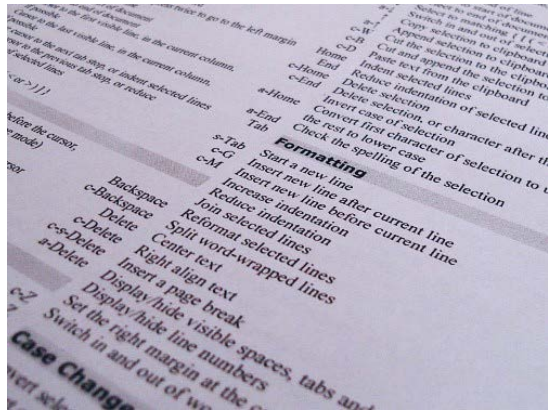
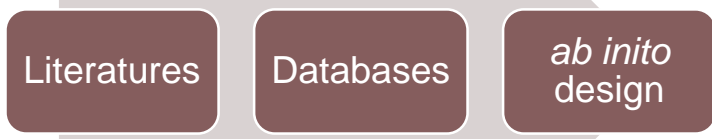
# Q: How to Design Primers/Probes for PCR/qPCR/Cloning/in situ hybridization





Repeat  
25-35X

Result after 1 cycle:  
# of DNA molecules  
doubled

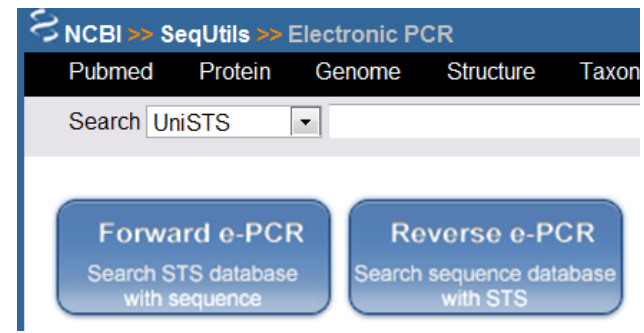


e.g., ACTB



BLAST

BLAST







► [NCBI/Primer-BLAST: Finding primers specific to your PCR template \(using Primer3 and BLAST\).](#) [more...](#) [Tips for finding specific primers](#)

[Reset page](#) [Save search parameters](#) [Retrieve recent results](#)

## PCR Template

Enter accession, gi, or FASTA sequence (A refseq record is preferred) [?](#) [Clear](#)

## Range

	From	To	
Forward primer	<input type="text"/>	<input type="text"/>	<a href="#">?</a> <a href="#">Clear</a>
Reverse primer	<input type="text"/>	<input type="text"/>	

Or, upload FASTA file

 [瀏覽...](#)

## Primer Parameters

Use my own forward primer  
(5'->3' on plus strand)

 [?](#) [Clear](#)

Use my own reverse primer  
(5'->3' on minus strand)

 [?](#) [Clear](#)

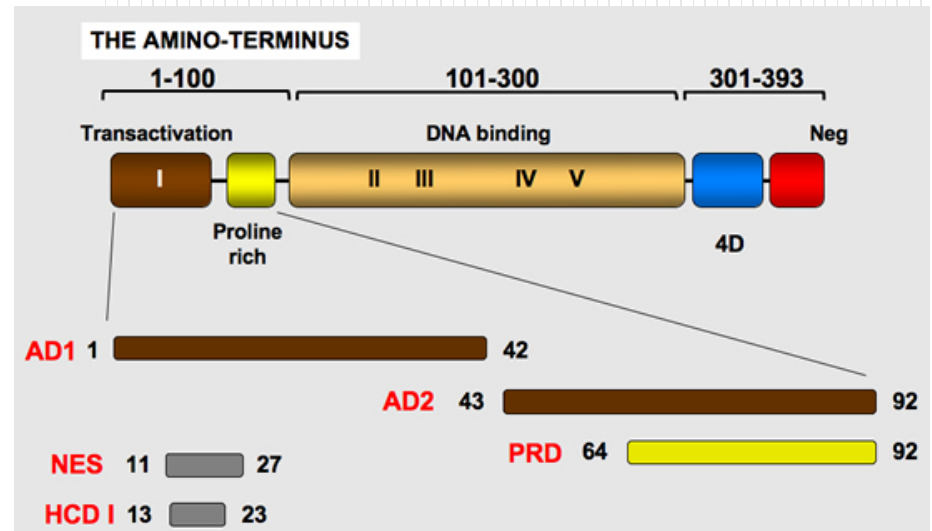
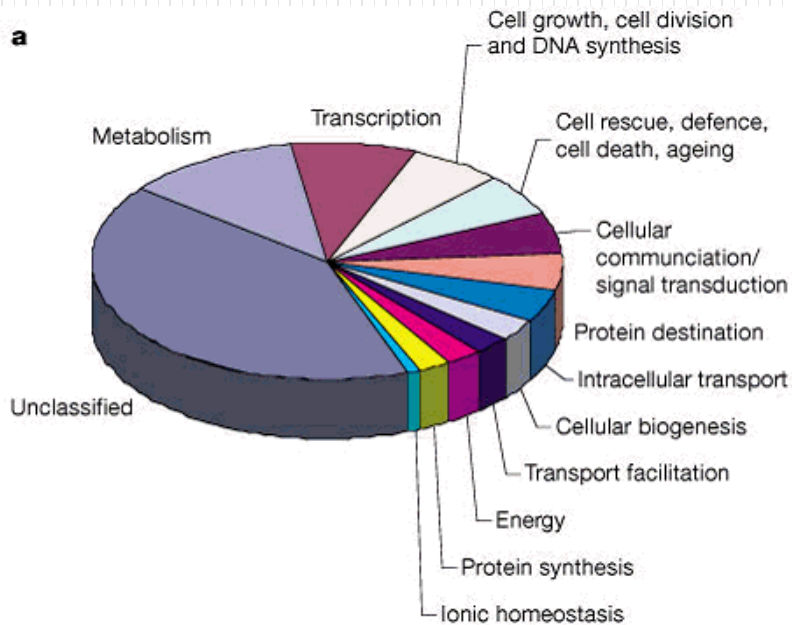
Min	Max
<input type="text" value="70"/>	<input type="text" value="1000"/>

PCR product size

# of primers to return

# Q: How to Find the Function and/or Structure of YFG

**a**



# 1. Gene Ontology

## Through integrated databases

- Entrez\_Gene
  - **GO terms**
- GeneCards
  - **GO terms**
- Uniprot/Swiss-Prot
  - POU5F1\_Human
  - General annotation (comments)
- Ontologies

Function	Evidence
<a href="#">DNA binding</a>	IDA <a href="#">PubMed</a>
<a href="#">miRNA binding</a>	IDA <a href="#">PubMed</a>
<a href="#">promoter binding</a>	IDA <a href="#">PubMed</a>
<a href="#">protein binding</a>	IPI <a href="#">PubMed</a>
<a href="#">sequence-specific DNA binding</a>	IEA
<a href="#">transcription factor activity</a>	IDA <a href="#">PubMed</a>
<a href="#">transcription factor binding</a>	IPI <a href="#">PubMed</a>

Process	Evidence
<a href="#">BMP signaling pathway involved in heart induction</a>	IMP <a href="#">PubMed</a>
<a href="#">anatomical structure morphogenesis</a>	TAS <a href="#">PubMed</a>
<a href="#">cardiac cell fate determination</a>	IDA <a href="#">PubMed</a>
<a href="#">cell fate commitment involved in the formation of primary germ layers</a>	IMP <a href="#">PubMed</a>
<a href="#">negative regulation of gene silencing by miRNA</a>	IMP <a href="#">PubMed</a>
<a href="#">positive regulation of SMAD protein nuclear translocation</a>	IDA <a href="#">PubMed</a>
<a href="#">positive regulation of catenin protein nuclear translocation</a>	IDA <a href="#">PubMed</a>
<a href="#">positive regulation of gene-specific transcription from RNA polymerase II promoter</a>	IDA <a href="#">PubMed</a>

# GO Evidence Code

## **Introduction**

### Experimental Evidence Codes

EXP: Inferred from Experiment

IDA: Inferred from Direct Assay

IPI: Inferred from Physical Interaction

IMP: Inferred from Mutant Phenotype

IGI: Inferred from Genetic Interaction

IEP: Inferred from Expression Pattern

### Computational Analysis Evidence Codes

ISS: Inferred from Sequence or Structural Similarity

ISO: Inferred from Sequence Orthology

ISA: Inferred from Sequence Alignment

ISM: Inferred from Sequence Model

IGC: Inferred from Genomic Context

RCA: inferred from Reviewed Computational Analysis

### Author Statement Evidence Codes

TAS: Traceable Author Statement

NAS: Non-traceable Author Statement

### Curator Statement Evidence Codes

IC: Inferred by Curator

ND: No biological Data available

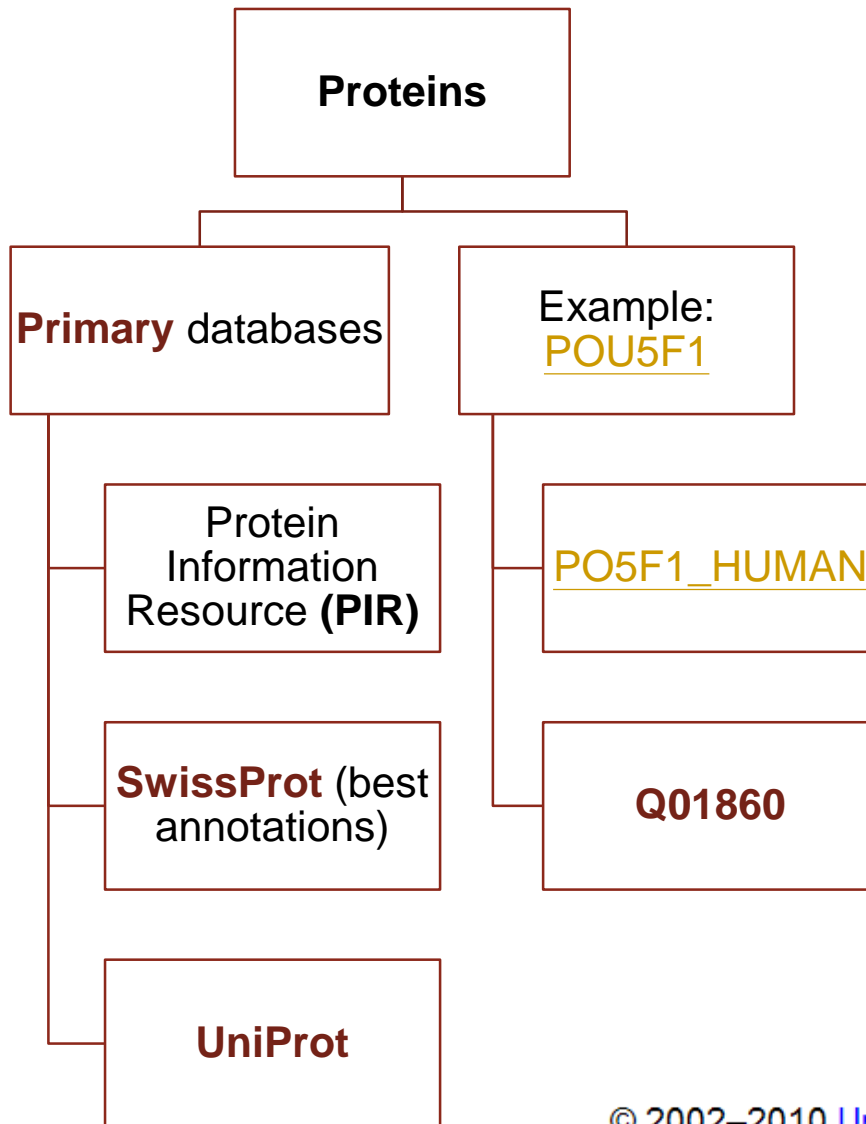
### **Automatically-assigned Evidence Codes**

IEA: Inferred from Electronic Annotation

### **Obsolete Evidence Codes**

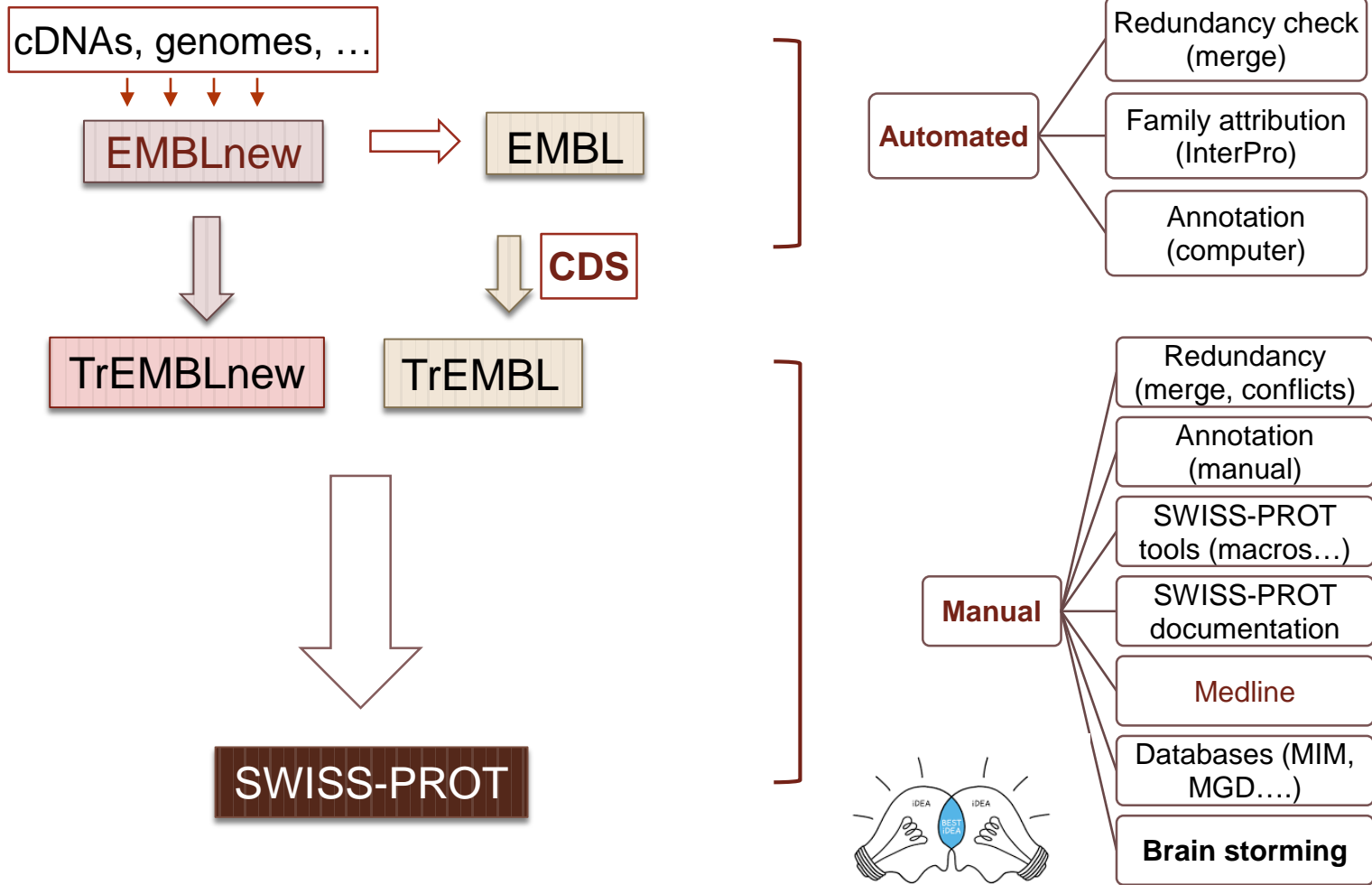
NR: Not Recorded

### **Note on Usage of the With/From Column**



© 2002–2010 UniProt Consortium | License & Disclaimer | Contact

# The Simplified Story of a SWISS-PROT Entry



Once in SWISS-PROT, the entry is no more in TrEMBL, **but still in EMBL (archive)**

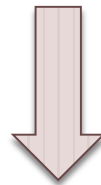
Domains, functional sites,  
protein families  
PROSITE  
InterPro  
Pfam  
PRINTS  
SMART  
Mendel-GFDb (plant gene  
families & EST annotations)

2D and 3D Structural dbs  
HSSP  
PDB

PTM  
CarbBank  
GlycoSuiteDB

2D-gel protein databases  
SWISS-2DPAGE  
ECO2DBASE  
HSC-2DPAGE  
Aarhus and Ghent  
MAIZE-2DPAGE

**SWISS-  
PROT**  
**UniProt  
KB**



Nucleotide sequence DB  
**EMBL, GeneBank, DDBJ**

Human diseases  
MIM

Protein-specific dbs  
GCRDb  
MEROPS (peptidase)  
REBASE  
TRANSFAC

Organism-spec. dbs  
DictyDb  
EcoGene  
FlyBase  
HIV  
MaizeDB  
MGD  
SGD  
StyGene (Salmonella)  
SubtiList  
TIGR  
TubercuList  
WormPep  
Zebrafish

# 2. UniProt/InterProt Annotations

UniProt  Advanced

BLAST Align Retrieve/ID mapping Help Contact

## UniProtKB results

[About UniProtKB](#) [Basket](#)

Filter by <sup>i</sup> BLAST Align Download Add to basket Columns Share 1 to 25 of 1,196 Show 25

Reviewed (814)

Unreviewed (382)

**Popular organisms**

- Human (230)
- Mouse (191)
- Rat (117)
- Bovine (63)
- Fruit fly (28)
- Other organisms

<input type="checkbox"/>	Entry	Entry name		Protein names	Gene names	Organism	Length	
<input type="checkbox"/>	P31750	AKT1_MOUSE		<b>RAC-alpha serine/threonine-protein ...</b>	<b>Akt1</b> Akt,Rac	Mus musculus (Mouse)	480	
<input type="checkbox"/>	P31749	AKT1_HUMAN		<b>RAC-alpha serine/threonine-protein ...</b>	<b>AKT1</b> PKB,RAC	Homo sapiens (Human)	480	
<input type="checkbox"/>	Q17941	AKT1_CAEEL		<b>Serine/threonine-protein kinase akt...</b>	<b>akt-1</b> C12D8.10	Caenorhabditis elegans	541	
<input type="checkbox"/>	P47196	AKT1_RAT		<b>RAC-alpha serine/threonine-protein ...</b>	<b>Akt1</b>	Rattus norvegicus (Rat)	480	
<input type="checkbox"/>	Q38998	AKT1_ARATH		<b>Potassium channel AKT1</b>	<b>AKT1</b> At2g26650,F18A8.2	Arabidopsis thaliana (Mouse-ear cress)	857	
<input type="checkbox"/>	Q8INB9	AKT1_DROME		<b>RAC serine/threonine-protein kinase</b>	<b>Akt1</b> CG4006	Drosophila melanogaster (Fruit fly)	611	



Entry

Feature viewer

Feature table

None

- Function
- Names & Taxonomy
- Subcellular location
- Pathology & Biotech
- PTM / Processing
- Expression
- Interaction
- Structure
- Family & Domains
- Sequences (2)
- Cross-references
- Publications
- Entry information
- Miscellaneous
- Similar proteins

▲ Top

## Domains and Repeats

Feature key	Position(s)	Length	Description	Graphical view	Feature identifier	Actions
Domain <sup>i</sup>	5 – 108	104	PH  PROSITE-ProRule annotation			Add BLAST
Domain <sup>i</sup>	150 – 408	259	Protein kinase PROSITE-ProRule annotation			Add BLAST
Domain <sup>i</sup>	409 – 480	72	AGC-kinase C-terminal			Add BLAST

## Region

Feature key	Position(s)	Length	Description	Graphical view	Feature identifier	Actions
Region <sup>i</sup>	14 – 19	6	Inositol-(1,3,4,5)-tetrakisphosphate binding			
Region <sup>i</sup>	23 – 25	3	Inositol-(1,3,4,5)-tetrakisphosphate binding			
Region <sup>i</sup>	228 – 230	3	Inhibitor binding			

Domain<sup>i</sup>

Binding of the PH domain to phosphatidylinositol 3,4,5-trisphosphate (PI(3,4,5)P<sub>3</sub>) following phosphatidylinositol 3-kinase alpha (PIK3CA) activity results in its targeting to the plasma membrane. The PH domain mediates interaction with TNK2 and Tyr-176 is also essential for this interaction. The AGC-kinase C-terminal mediates interaction with THEM4.

Sequence similarities<sup>i</sup>

Belongs to the [protein kinase superfamily](#). [AGC Ser/Thr protein kinase family](#). [RAC subfamily](#). Curated

Contains 1 [AGC-kinase C-terminal domain](#). Curated

Contains 1 [PH domain](#). PROSITE-ProRule annotation

Contains 1 [protein kinase domain](#). PROSITE-ProRule annotation

## Display



## PTM / Processing<sup>1</sup>

Entry

Feature viewer

Feature table

None

- Function
- Names & Taxonomy
- Subcellular location
- Pathology & Biotech
- PTM / Processing
- Expression
- Interaction
- Structure
- Family & Domains
- Sequences (2)
- Cross-references
- Publications
- Entry information
- Miscellaneous
- Similar proteins

▲ Top

### Molecule processing

Feature key	Position(s)	Length	Description	Graphical view	Feature identifier	Actions
Chain <sup>1</sup>	1 – 480	480	RAC-alpha serine/threonine-protein kinase		PRO_0000085605	<a href="#">Add</a> <a href="#">BLAST</a>

### Amino acid modifications

Feature key	Position(s)	Length	Description	Graphical view	Feature identifier	Actions
Modified residue <sup>1</sup>	14 – 14	1	N6-acetyllysine <a href="#">1 Publication</a>			
Modified residue <sup>1</sup>	20 – 20	1	N6-acetyllysine <a href="#">1 Publication</a>			
Disulfide bond <sup>1</sup>	60 ↔ 77		<a href="#">1 Publication</a>			
Modified residue <sup>1</sup>	124 – 124	1	Phosphoserine <a href="#">Combined sources</a>			
Modified residue <sup>1</sup>	126 – 126	1	Phosphoserine; alternate <a href="#">Combined sources</a>			
Glycosylation <sup>1</sup>	126 – 126	1	O-linked (GlcNAc); alternate <a href="#">1 Publication</a>			
Modified residue <sup>1</sup>	129 – 129	1	Phosphoserine; alternate <a href="#">Combined sources</a>			
Glycosylation <sup>1</sup>	129 – 129	1	O-linked (GlcNAc); alternate <a href="#">1 Publication</a>			
Modified residue <sup>1</sup>	176 – 176	1	Phosphotyrosine; by TNK2 <a href="#">1 Publication</a>			
Cross-link <sup>1</sup>	284 – 284		Glycyl lysine isopeptide (Lys-Gly) (interchain with G-Cter in ubiquitin) <a href="#">1 Publication</a>			
Disulfide bond <sup>1</sup>	296 ↔ 310		<a href="#">By similarity</a>			

# 3. If YFG Involves in Specific Function/Pathway? - through its interacted proteins

BioGRID 3.1

**CHD4**

*Mus musculus*

AA617397, mKIAA4075, D6ErtD380e, Mi-2beta, BC005710, KIAA4075, 9530019N15Rik, MGC11769

chromodomain helicase DNA binding protein 4

GO Process: 0 Terms

GO Function: 1 Terms

GO Component: 0 Terms

EXTERNAL DATABASE LINKOUTS

[MGI](#) | [Entrez Gene](#) | [RefSEQ](#) | [GenBank](#) | [UniprotKB](#)

Download 13 Associations For This Protein

Stats & Filters

Current Stati

High Throughput

10 (59%)

0 (0%)

Search Filter

No Filter: Show

Switch View:

Summary

Sortable Table

Displaying 13 total unique interactors




**POU5F1** | Otf-3, Oct3, Oct-3/4, Otf3, Oct3/4, Oct-3, Oct4, Otf-4, Oct-4, Otf3-rs7, Otf4, Otf3g

POU domain, class 5, transcription factor 1

**MTA2** | mmta2, Mta1l1, Mata1l1, AW550797

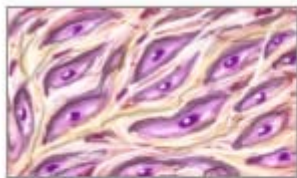
metastasis-associated gene family, member 2

# Databases for Protein – Protein Interaction

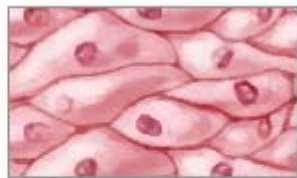
Resource	Comments
<a href="#">APID</a>	Agile Protein Interaction DataAnalyzer (Cancer Research Center, Salamanca, Spain)
<a href="#">BIND</a>	Biomolecular Interaction Network Database at the University of Toronto, Canada. No species restriction
<a href="#">CYGD</a>	PPI section of the Comprehensive Yeast Genome Database. Manually curated comprehensive <i>S. cerevisiae</i> PPI database at MIPS
<a href="#">DIP</a>	Database of Interacting Proteins at UCLA. No species restriction.
 <a href="#">GRID</a>	General Repository for Interaction Datasets. Mount Sinai Hospital, Toronto, Canada
<a href="#">HIV Interaction DB</a>	Interactions between HIV and host proteins.
 <a href="#">HPRD</a>	The Human Protein Reference Database. Institute of Bioinformatics, Bangalore, India and Johns Hopkins University, Baltimore, MD, USA.
<a href="#">HPID</a>	Human Protein Interaction Database. Department of computer Science and Information Engineering Inha University, Incheon, Korea
<a href="#">iHOP</a>	iHOP (Information Hyperlinked over Proteins). Protein association network built by literature mining
 <a href="#">IntAct</a>	Protein interaction database at EBI. No species restriction.
<a href="#">InterDom</a>	Database of putative interacting protein domains. Institute for InfoComm Research, Singapore.
<a href="#">JCB</a>	PPI site at the Jena Centre for Bioinformatics, Germany
<a href="#">MetaCore</a>	Commercial software suite and database. Manually curated human PPIs (among other things). GeneGo
<a href="#">MINT</a>	Molecular Interaction database at the Centro di Bioinformatica Molecolare, Universita di Roma, Italy.
<a href="#">MRC PPI links</a>	Commented list of links to PPI databases and resources maintained at the MRC Rosalind Franklin Centre for Genomics Research, Cambridge, UK
<a href="#">OPHID</a>	The Online Predicted Human Interaction Database. Ontario Cancer Institute and University of Toronto, Canada.
<a href="#">Pawson Lab</a>	Information on protein-interaction domains.
<a href="#">PPI</a>	...

# Q: What Kind of Cell Lines or Tissues I Should Use for PCR-based Cloning YFG?

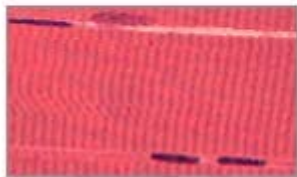
Four types of tissue



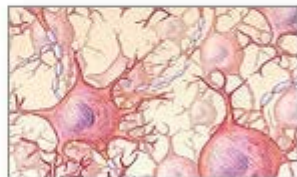
Connective tissue



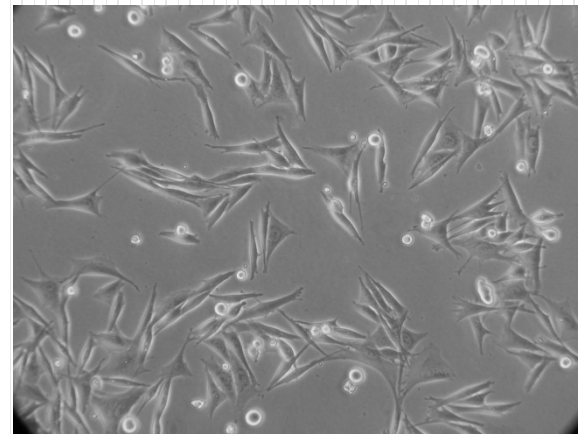
Epithelial tissue



Muscle tissue



Nervous tissue



Modify Query

### TCGA PanCancer Atlas Studies

User-defined Patient List (10953 patients / 10967 samples) - POU5F1

Queried gene is altered in • 161 (1%) of queried patients  
• 161 (1%) of queried samples



OncoPrint

Cancer Types Summary

Mutations

Survival

CN Segments

Expression

Download

POU5F1

Profile:

RNA Seq V2

Sort By:

Cancer Study



Log scale



Show mutations \*



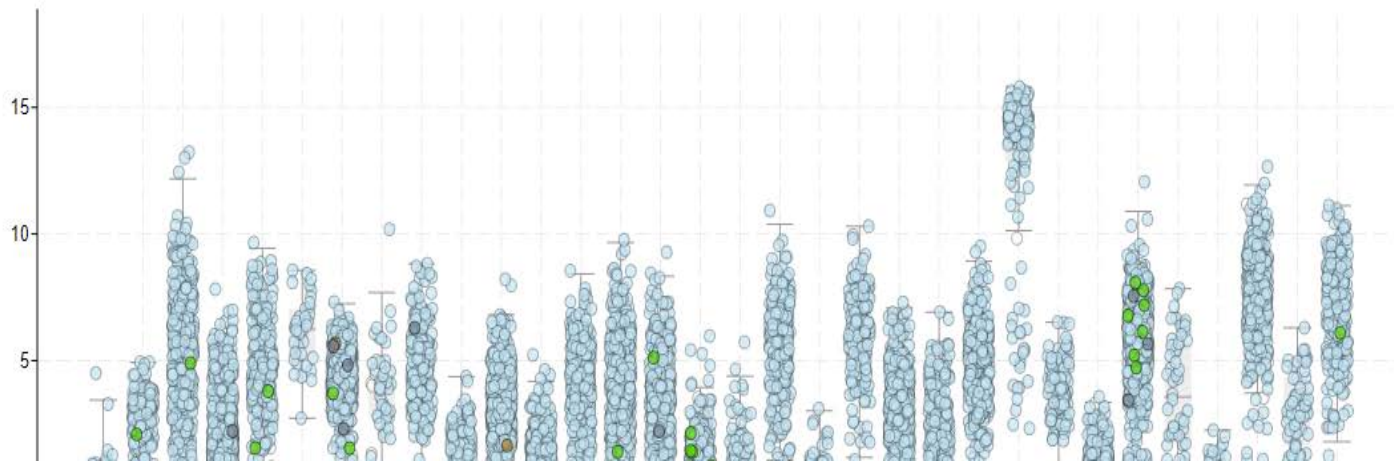
Show copy number alterations

Select studies:

TCGA Pan-Can Atlas (32)

Custom list

pression --- RNA Seq V2 (log2)



**Q: What Would I Do When I am Having  
Breakfast or a Coffee Break?**





## Coffee Break

### Tutorials for NCBI Tools

Edited by Laura Dean and Johanna McEntyre.

National Center for Biotechnology Information

Bethesda (MD): [National Center for Biotechnology Information \(US\)](#); 1999-.

[Copyright notice.](#)



*Coffee Break* is a resource at NCBI that combines reports on recent biomedical discoveries with use of NCBI tools. The result is an interactive tutorial that tells a biological story. Each report is based on a discovery reported in one or more articles from the recently published peer-reviewed literature. After a brief introduction that sets the work described into a broader context, the report focuses on how a molecular understanding can provide explanations of observed biology and lead to therapies for diseases.



# Bookshelf

U.S. National Library of Medicine  
National Institutes of Health

Search  ▾

[Limits](#) [Help](#)

Bookshelf ID: NBK1969



## NCBI News

Bethesda (MD): [National Center for Biotechnology Information \(US\)](#); 199

ISSN: 1060-8788

Publication No.: 94-3272

[Copyright notice.](#)

## Index of Issues

▣ [NCBI News, March 2011](#)

[Expand All](#)

[PubMed Interface for Mobile Devices Now Available](#)

[NCBI Bookshelf Updated to the New Entrez Design](#)

[New Organism Builds in UniGene](#)

[NCBI YouTube Video Update](#)



liver tumor mouse

## Seed tumor at liver of mouse-Surgery 種腫瘤在肝臟-開腹腔篇(一)

miss9ch

282 部影片

訂閱

This video contains animal  
experiment content,  
Viewer discretion is advise



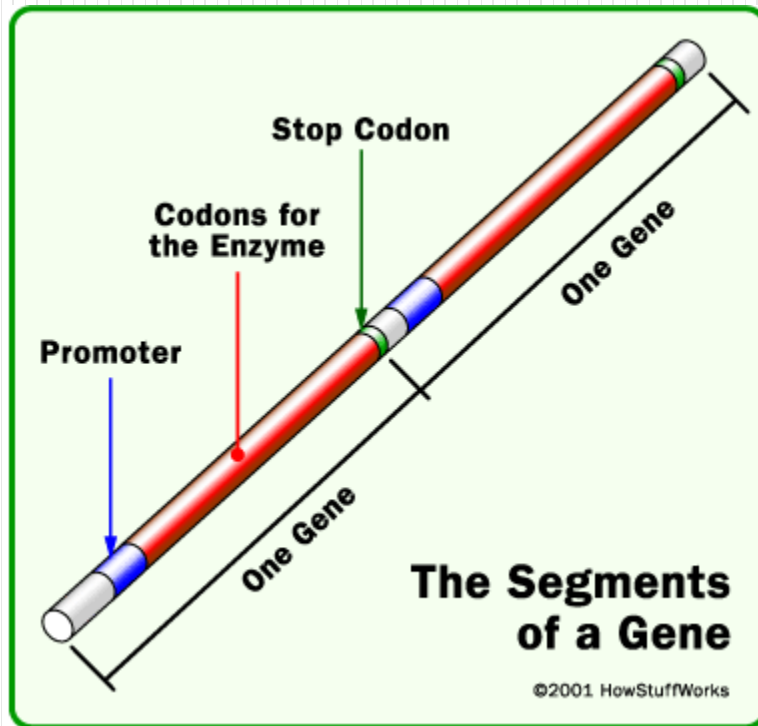
0:02 / 3:01



240p



# Q: How do You Know You've Cloned the Correct YGF? (Wild type vs. Mutant?)



NCBI/ BLAST/ blastn suite

blastn blastp blastx tblastn tblastx

BLASTN programs search nucleotide databases using a nucleotide query. [more...](#)

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s)

Clear

Query subrange

From

To

Genomic plus Transcript

Human genomic plus transcript (Human G+T)

Mouse genomic plus transcript (Mouse G+T)

Other Databases

Nucleotide collection (nr/nt)

Reference mRNA sequences (refseq\_rna)

Reference genomic sequences (refseq\_genomic)

NCBI Genomes (chromosome)

Expressed sequence tags (est)

Non-human, non-mouse ESTs (est\_others)

Genomic survey sequences (gss)

High throughput genomic sequences (HTGS)

Patent sequences(pat)

Protein Data Bank (pdb)

Human ALU repeat elements (alu\_repeats)

Sequence tagged sites (dbsts)

Whole-genome shotgun reads (wgs)

Environmental samples (env\_nt)

Human genomic plus transcript (Human G+T)

Or, upload file

Job Title

Align two or more sequences

Choose Search Set

Database

Exclude

Optional

Entrez Query

Models (XM/XP)  Uncultured/environmental sample sequences

[NCBI Homepage](#)

## Contamination

[Definition](#)  
[Sources](#)  
[Consequences](#)  
[Detection](#)

## VecScreen

[Overview](#)  
[Example](#)  
[Search Parameters](#)  
[Match Categories](#)  
[Interpretation](#)  
[Exceptions](#)

## UniVec Database

[Overview](#)  
[Redundancy](#)  
[Elimination](#)  
[Benefits](#)  
[Pseudo-](#)  
[Circularization](#)  
[Vector Representation](#)

## ▶ Screen a Sequence Using VecScreen

Enter your query sequence below as an Accession, GI, or **FASTA**.

## ▶ About VecScreen

[VecScreen](#) is a system for quickly identifying segments of a nucleic acid sequence that may be of vector origin. NCBI developed VecScreen to combat the problem of vector [contamination](#) in public sequence databases. This Web page is designed to help researchers identify and remove any segments of vector origin before sequence analysis or submission.

# ORF Finder (Open Reading Frame Finder)

PubMed

Entrez

BLAST

OMIM

Taxon

NCBI

**Tools**  
for data mining

**GenBank**  
sequence submission  
support and software

**FTP site**  
download data and  
software

The ORF Finder (Open Reading Frame Finder) is a graphical analysis tool which finds selectable minimum size in a user's sequence or in a sequence already in the data. This tool identifies all open reading frames using the standard or alternative genetic code. The ORF Finder should be helpful in preparing complete and accurate sequence data for the Sequin sequence submission software.

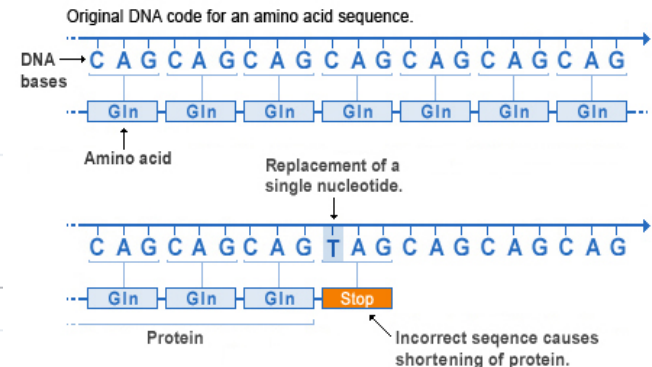
Enter GI or ACCESSION

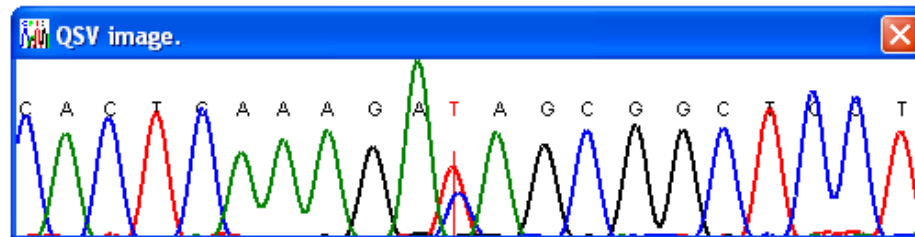
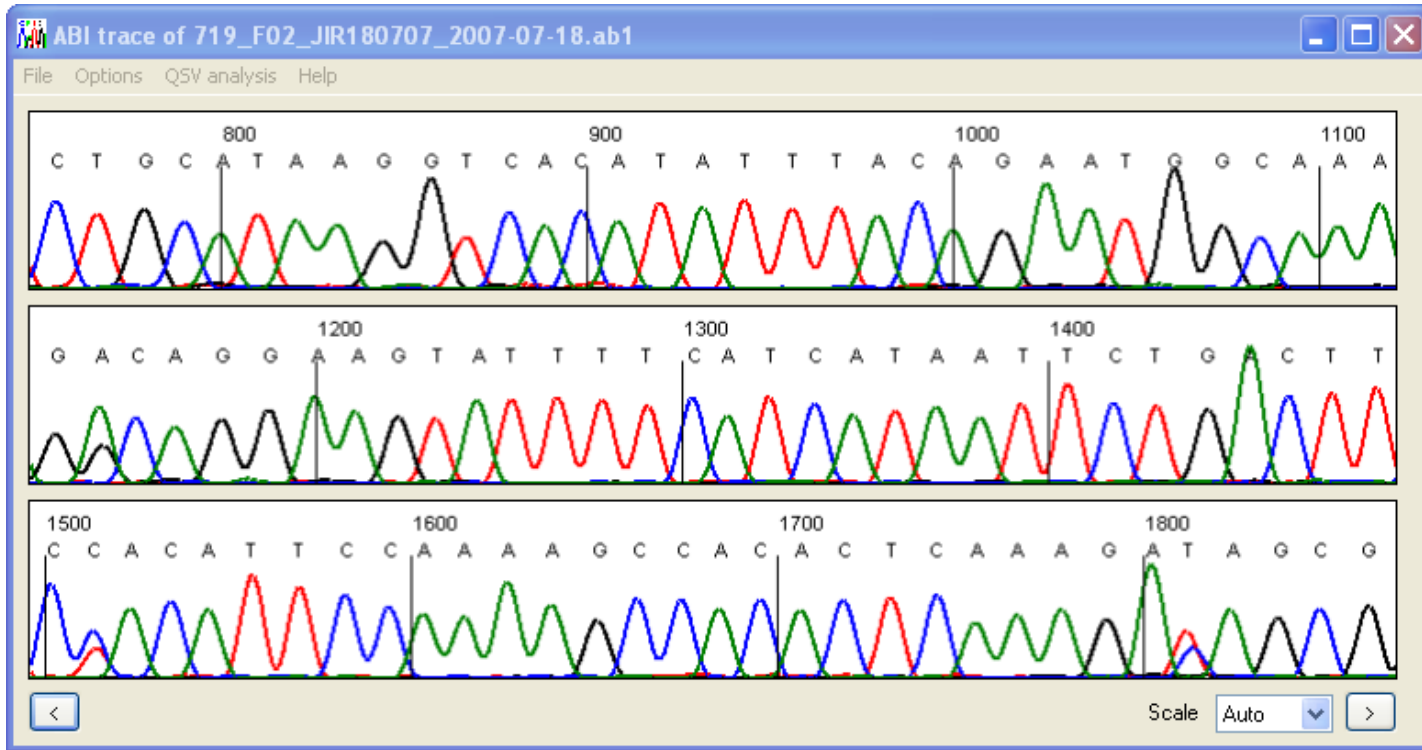
or sequence in FASTA format

FROM:  TO:

Genetic codes

Nonsense mutation





# When Cloned by Emails – get the map & confirmed

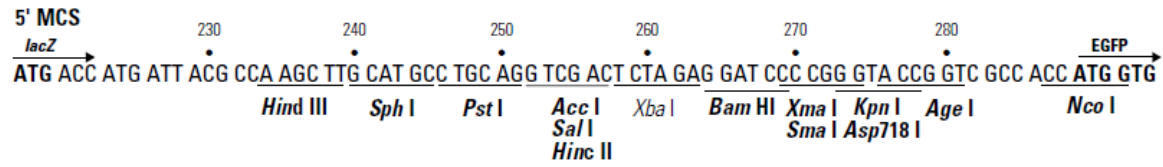
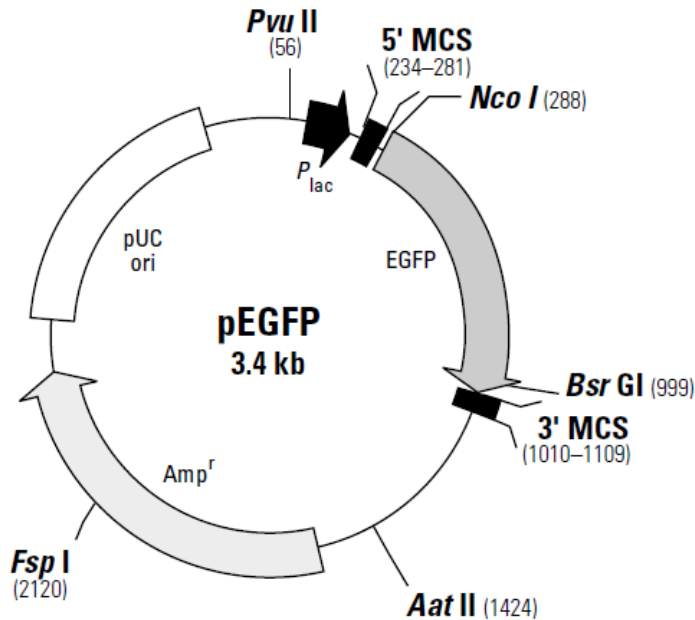
Specific EGFP Monoclonal Antibody for Westerns, IP and IC

Visit our website  
for more details!  
click here...

pEGFP Vector Information

PT3078-5

Catalog #6077-1





## Q: How to Get a Specific Sequence from Genome Databases



# Genome Biology

▼ Vertebrates	(17)
▼ Mammals	(14)
▼ Primates	(3)

[Map Viewer](#), NCBI

[Genome Browser](#), UCSC

[Ensembl Genome Browser](#), EBI



*e!*Ensembl

## Human (*Homo sapiens*) Genome Browser Gateway

The UCSC Genome Browser was created by the [Genome Bioinformatics Group of UC Santa Cruz](#).  
Software Copyright (c) The Regents of the University of California. All rights reserved.

clade	genome	assembly	position or search term	gene	image width	
Mammal	Human	Feb. 2009 (GRCh37/hg19)	chr6_mcf_hap5:2,514,038-2,520,39	POU5F1	800	submit

[Click here to reset](#) the browser user interface settings to their defaults. **2011 ENC**

track search    add custom tracks    configure tracks and display

- POU5F1
- POU5F1B
- POU5F1P1
- POU5F1P3
- POU5F1P4
- POU5F2
- POU5FLC12

**Survey**

## About the Human Feb. 2009 (GRCh37/hg19) assembly ([sequences](#))

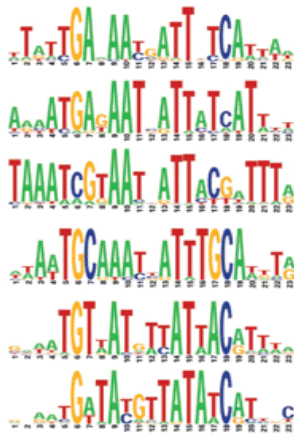
The February 2009 human reference sequence (GRCh37) was produced by the [Genome Reference](#)

## UCSC Genome Browser on Human Feb. 2009 (GRCh37/hg19) Assembly

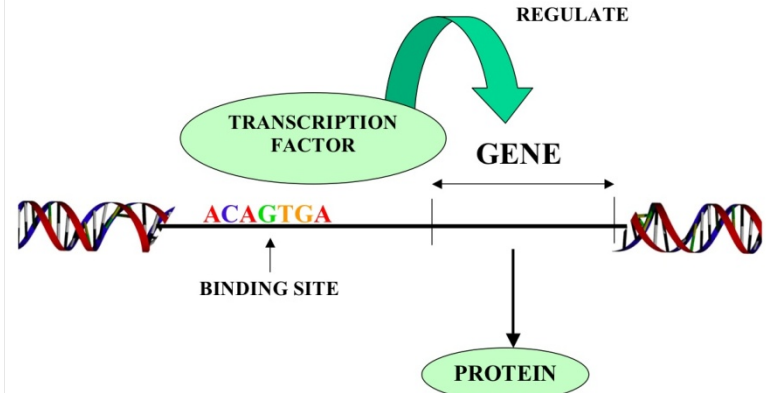
move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x

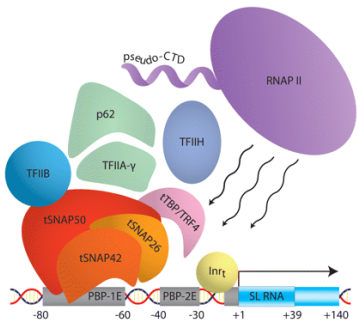
position/search chr6\_mcf\_hap5:2,514,038-2,520,39 [gene](#) jump clear size 6,356 bp. configure **2011 ENCODE Usability Survey**

# Q: How to Identify Potential Regulators?



Legend: A transcription factor molecule binds to the DNA at its binding site, and thereby regulates the production of a protein from a gene.





# Feature-Based Methods

Based on identifying **gene signals**

Promoter elements

Splice sites

Start/stop codons

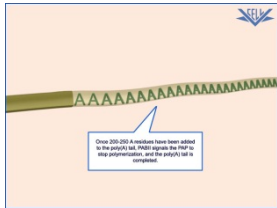
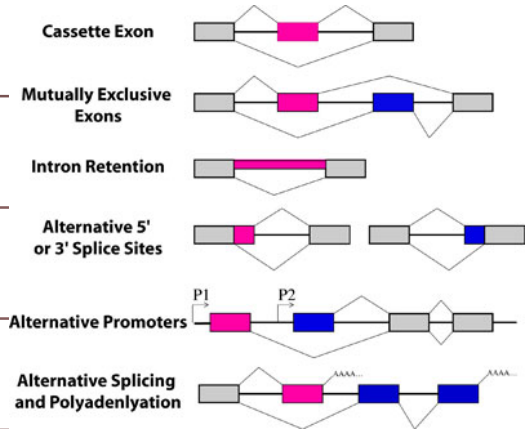
PolyA sites...

Consensus sequences

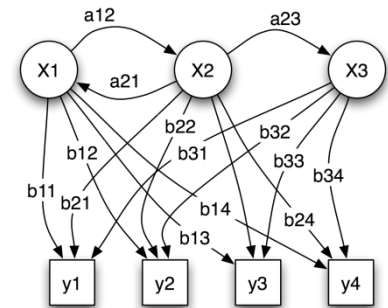
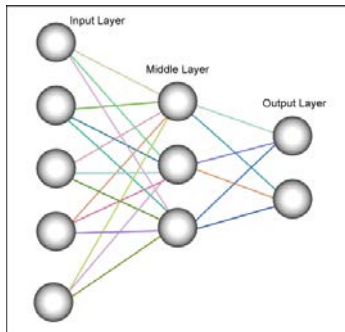
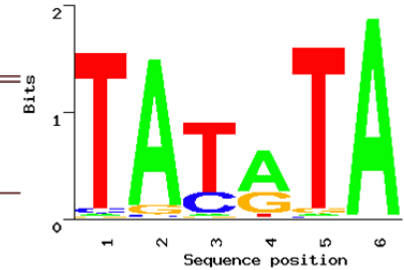
Weight matrices

Neural networks (NNs) Decision trees

Hidden Markov Models (HMMs)



Wide range of **methods**



# Promoter Databases and sites for analysis, prediction and search

[AlignACE](#)

motif-finding algorithm.

[Promoter Binding Element Database](#)  
[CpG promoter](#)

Arabidopsis thaliana promoter binding element database  
promoter mapping using CpG islands

[Core promoter](#)

to predict putative Transcriptional Start Site (TSS)

[dbtss](#)

Database of Transcriptional Start Sites

[Dragon Promoter Finder](#)

an advanced system for promoter recognition in vertebrates

[EPD](#)

an annotated non-redundant collection of eukaryotic POL II promoters

[FirstEF](#)

a 5' terminal exon and promoter prediction program

[Human Promoter Database](#)

Search for transcriptional start site

[Mcpromoter](#)

A statistical tool for the prediction of transcription start sites

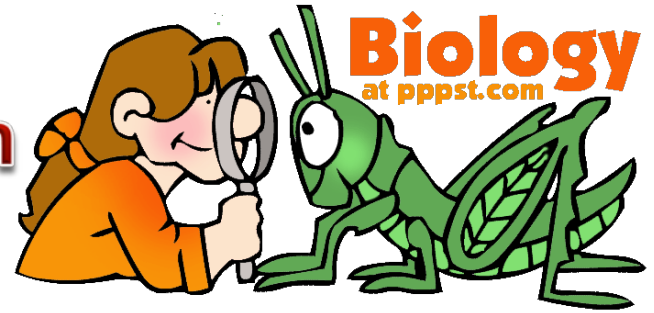
[Motif Explorer](#)

Motif & promoter visualization

[Neural Network Promoter Prediction](#)

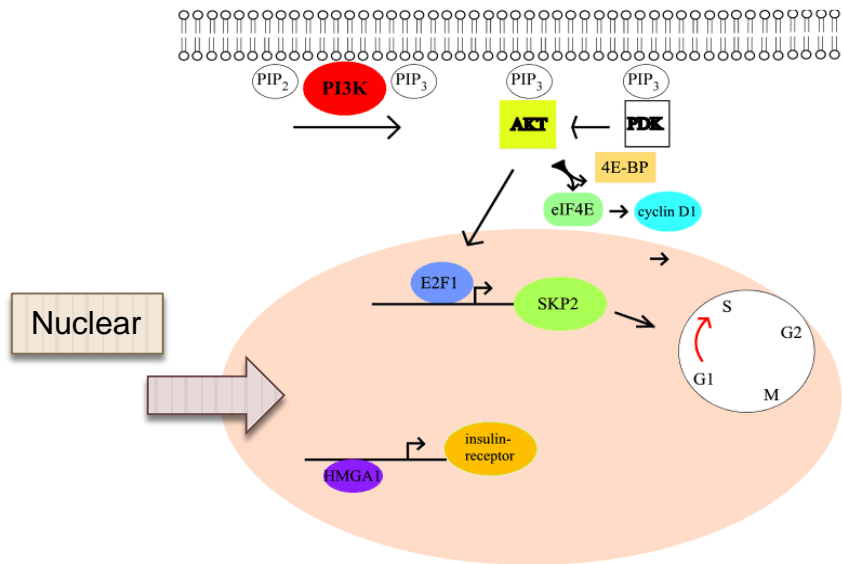
Neural Network Promoter Prediction

# Pattern-driven



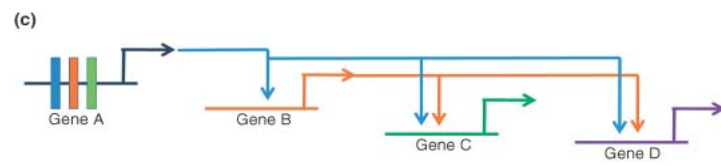
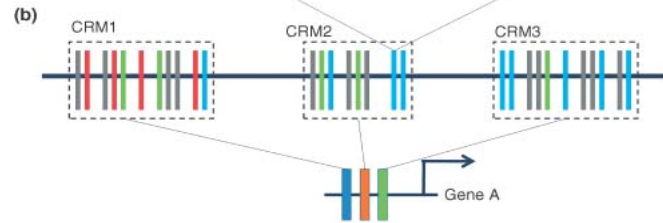
Success depends on **available of collections of annotated binding sites**

- Tend to produce huge numbers of **false-positive**
- **Reasons**
  - Binding sites (BS) for specific TFs often **variable**
  - Binding sites are short (typically **5-15 bp**)
  - **Interactions** between TFs (& other proteins) influence **affinity** & **specificity** of TF binding
  - One binding site often recognized by **multiple TFs**
  - **Biology is complex**: promoters often specific to **organism/cell/stage/environmental** condition



PI3K/AKT signaling in pancreatic cancer cells

	1	2	3	4	5	6	7	8	9	10	11	12	13	14
A	25	28	70	10	0	100	0	0	0	0	2	6	18	31
C	31	7	3	69	100	0	100	0	0	0	19	21	47	13
G	13	47	21	19	0	0	0	100	0	100	69	3	7	31
T	31	18	6	2	0	0	0	0	100	0	10	70	28	25



## Taking **sequence context/biology** into account (Do the **wet lab** experiments!!!)

**Eukaryotes:** clusters of TFBSs are common

**Probability** of “real” binding site increases if annotated **transcription start site (TSS) nearby**

- But **NOT** for enhancers
- Only a **small fraction of TSSs** have been experimentally mapped

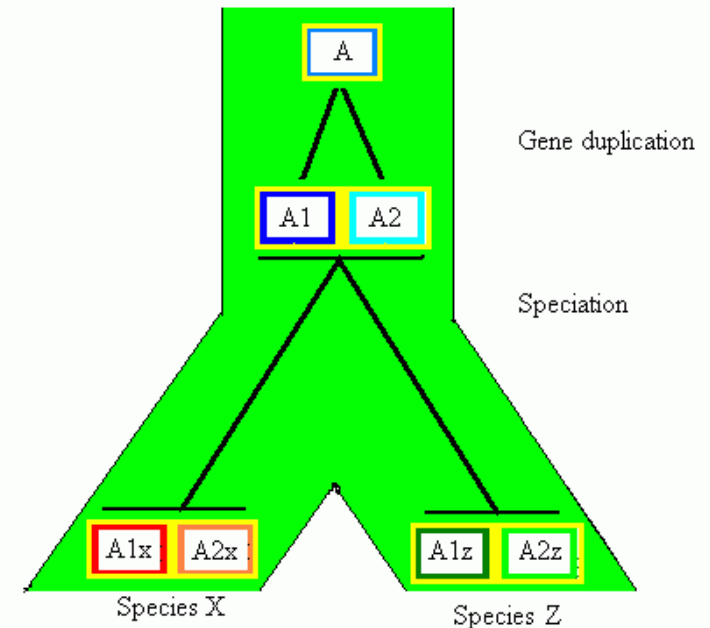
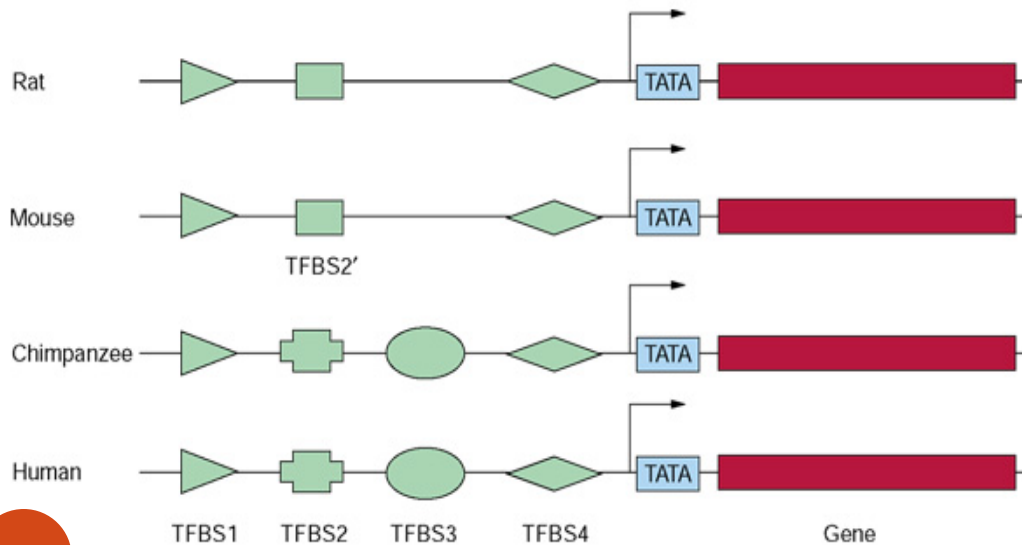
**Comparative** promoter mapping

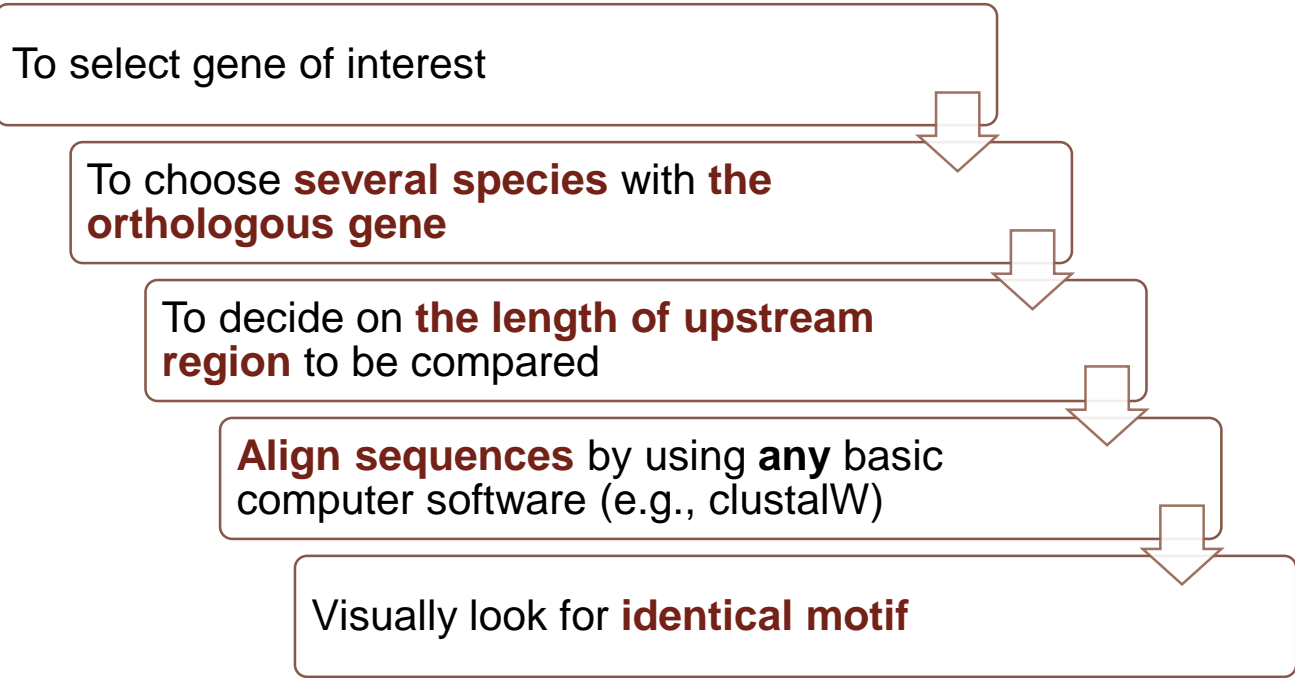


# Phylogenetic Footprinting

Patterns of gene regulation are often conserved across species

- Interspecies comparisons  $\Rightarrow$  to identify **common regulatory sequences** (Wasserman et al. 2000)
  - The selection of appropriate species, critical





```

Human  TAACAATGGTACATCCTAATGGAAGCTGOGAGGGAAATGCAATAATTTGCGGAAGCTGGGCATGAGCCCTGCCTCCAGCGGGTGGCGCTCGAGTCCGG 765
Dog     TAACAATGGTACATCCTAATGGAAGCTGOGAGGGAAATGCAATAATTTGCGGAAGCTGGGCATGAGCCCTGCCTCCAGCGGGTGGCGCTCGAGTCCGG 765
Mouse  TCACAATGGTACATCCTAATGGAAGCTGOGAGGGAAATGCAATAATTTGCGGAAGCGAAGCGATGCGCCAGTCTCCAGCGGGTGGCGCTCGAGTCCGA 941

Human  CTGAACGGCGGCAACTGGCGGCGGGCACGCGCCCGGGGCGCGCGCCACCCCTCTGCCTCCACCCAACTCCCTATTAGTGCAACGAGTTTACCTCTAG 865
Dog     CTGAACGGCGGCAACTGGCGGCGGGCACGCGCCCGGGGCGCGCGCCACCCCTCTGCCTCCACCCAACTCCCCATTAGTGCAACGAGTTTACCTCTAG 865
Mouse  CTGAACGGCGGCAACGGTGGCGGGGACGCGCCCGGGGCGCGCGCCACCCCTCTGCCTCCACCCAACTC----- 1014
  
```

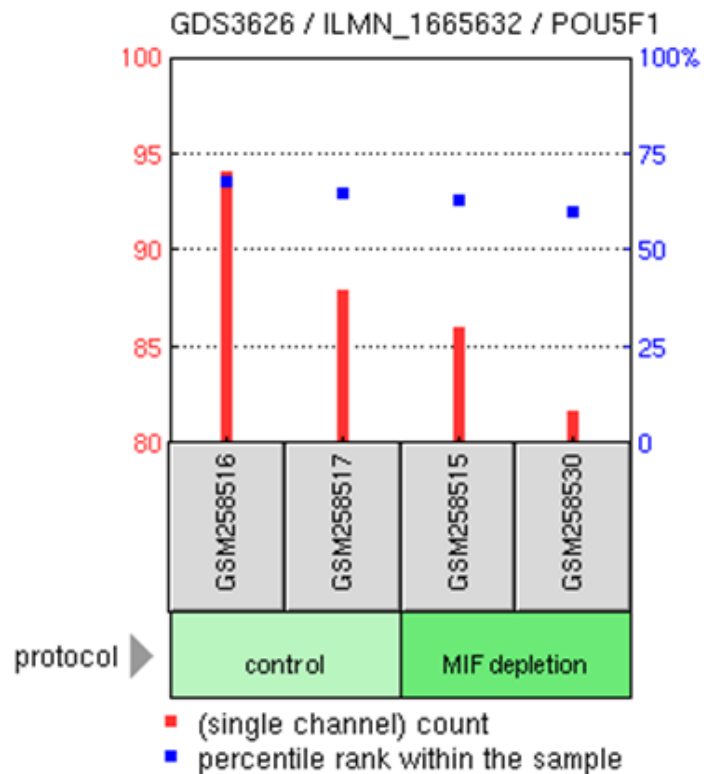
Potential TFBS: Ubx1 binding site  
NF- $\gamma$  binding site  
SP1 binding site  
GATA-1 binding site

\*All TF names are from human with orthologous TFs present in both dog and mouse.

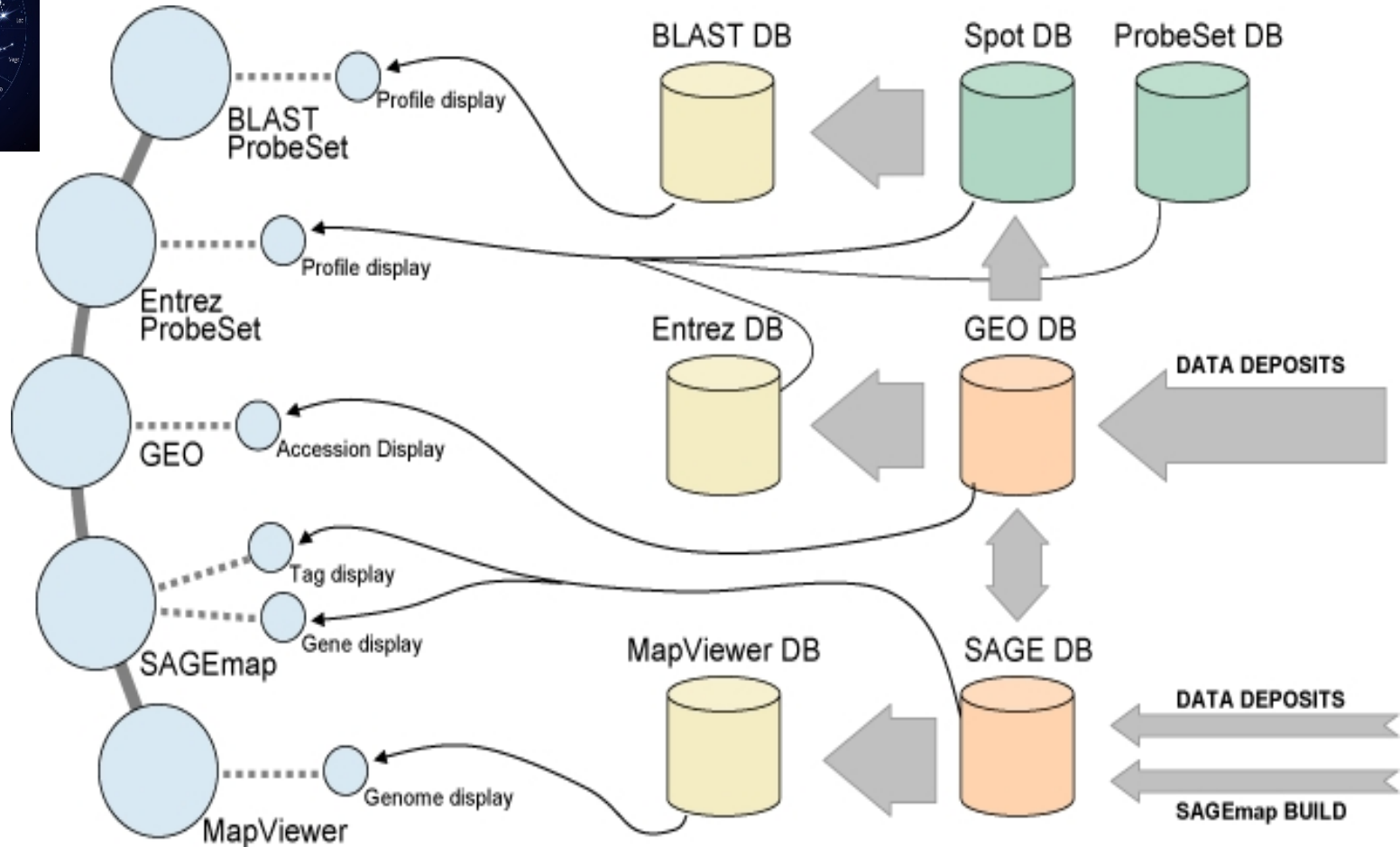
# One More Trick - Coregulation

**Title:** [GDS3626](#) / ILMN\_1665632 / POU5F1 / Homo sapiens

**Summary:** Analysis of HEK293 kidney cells depleted for the (0)/G(1) cell cycle arrest. Results provide insight into the mo



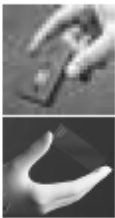
# Constellation of NCBI Gene Expression Resources



# Gene Expression Omnibus (GEO) (1)

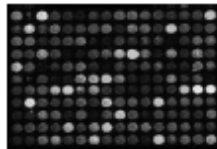
Submitted by  
Manufacturer\*

**GPL**  
Platform  
descriptions



Submitted by  
Experimentalists

**GSM**  
Raw/processed  
spot intensities  
from a single  
slide/chip



Entrez GEO

**GSE**  
Grouping of  
slide/chip data  
“a single experiment”

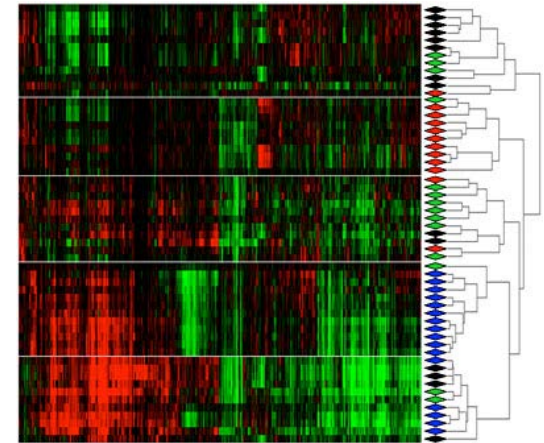


Curated by  
NCBI

**GDS**  
Grouping of  
experiments



Entrez  
GEO Datasets



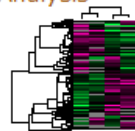
# Gene Expression Omnibus (GEO) (2)

- × Search GEO Profiles: POU5F1
  - × Or **Limit, Preview/Index**
- × GDS vs. GSE

Search for     [Advanced Search](#)

DataSet Record GDS46: <a href="#">Expression Profiles</a> <a href="#">Data Analysis Tools</a> <a href="#">Sample Subsets</a>			
<b>Title:</b>	E2F1-regulated genes		
<b>Summary:</b>	Identification of E2F1-regulated genes that modulate the transition from quiescence into DNA synthesis, or have roles in apoptosis, signal transduction, membrane biology, and transcription repression.		
<b>Organism:</b>	<i>Mus musculus</i>		
<b>Platform:</b>	GPL75: [Mu11KsubA] Affymetrix Murine 11K SubA Array		
<b>Citation:</b>	Ma Y, Croxton R, Moorer RL Jr, Cress WD. Identification of novel E2F1-regulated genes by microarray. <i>Arch Biochem Biophys</i> 2002 Mar 15;399(2):212-24. PMID: 11888208		
<b>Reference Series:</b>	<a href="#">GSE498</a>	<b>Sample count:</b>	4
<b>Value type:</b>	count	<b>Series published:</b>	2003/07/16

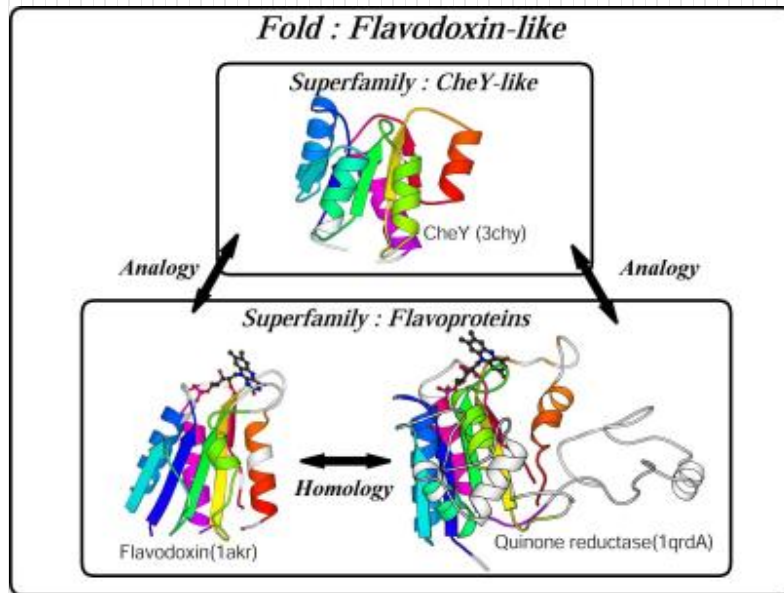
**Cluster Analysis**



**Download**

- DataSet SOFT file
- Series family SOFT file
- Series family MINiML file
- Annotation SOFT file

# Q: Can You Speculate the Function of YFG from Structure Similarity?

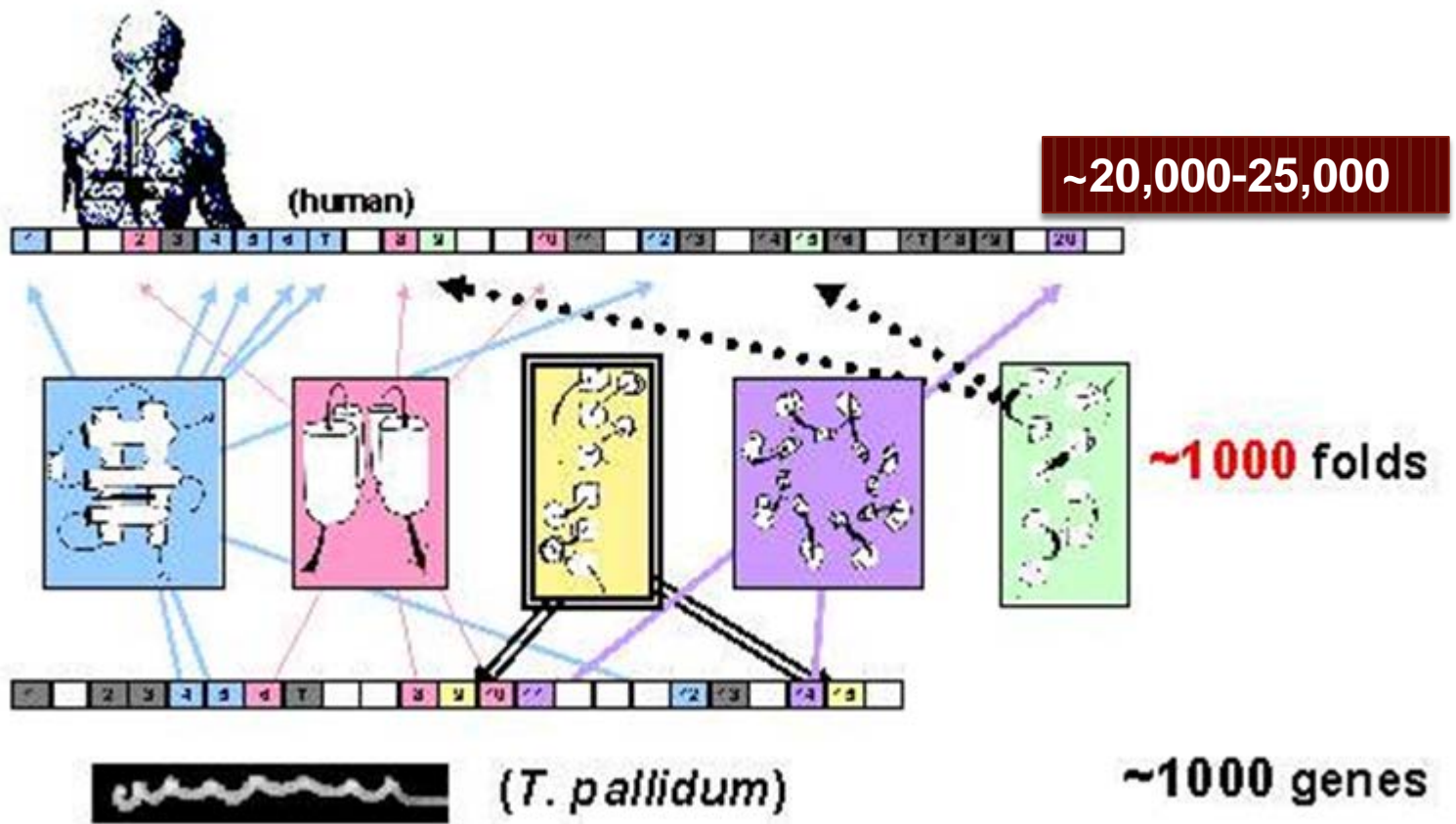


# Structures are More Conserved Than Sequences

Evolution	Homology	% Identity	Alignment Methods
Recent relationship - less divergence	Sequence alignments can be used to infer homology	100	Automatic Pairwise Alignment Methods
Increasing divergence		90	
	80		
	70		
	60		
Distant relationship	<b>Twilight Zone</b>	50	Consensus Methods
		40	
	30	Profile Methods	
	<b>Midnight Zone</b>	20	Structure Prediction
		10	
		0	



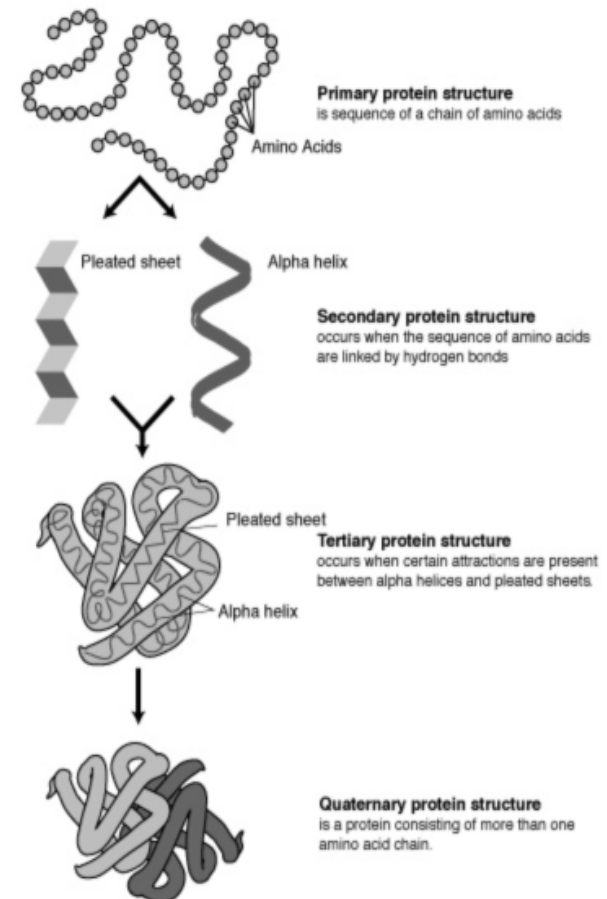
# Simplifying Genomes with Folds, Pathways



Significance: fold # << sequence ##

# Levels of Protein Sequence & Structure Organization

Level/ Database	Content	Example
Primary	Sequence	"AVILDRYFH"
Secondary	Motif	[AS]-[IL]2-X[DE]- R-[FYW]2-H
Tertiary	Domain/ module	a,b,c or @, *, #



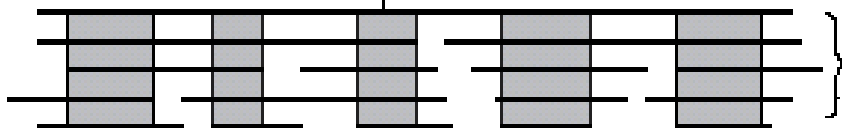
Single motif methods

permissive regular expression (IDENTIFY)

eMotif

exact regular expression (PROSITE)

XXXX  
XXXX  
XXXX



Full domain alignment methods

Profile (Profile library)

Hidden Markov Model (Pfam)

RWDAGCVN  
RWDSGCVN  
RWHHGCVQ  
RWKGACYN  
RWLVACEQ

XXXX  
XXXX  
XXXX

XXXX  
XXXX  
XXXX

XXXX  
XXXX  
XXXX

XXXX  
XXXX  
XXXX

frequency matrices (PRINTS)

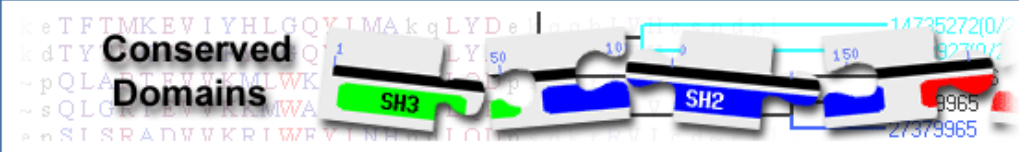
position-specific weight matrices (Blocks)

Multiple motif methods

Attwood 2000

# Major Secondary “Pattern” Database

2 <sup>nd</sup> Database	Primary Source	Stored Information
<u>PROSITE</u>	SWISS-PROT	Regular expression ( <b>pattern</b> )
<u>PROSITE</u>	BLOCKS+/Prints	<b>Fuzzy</b> expression ( <b>pattern</b> )
<u>PRINTS</u>	SWISS-PROT/ TrEMBL	Aligned motifs - fingerprints
Profiles ( <u>Prosite</u> )	SWISS-PROT	Weighted matrices ( <b>profiles</b> )
<u>Pfam/SMART</u>	SWISS-PROT	Hidden Markov Models ( <b>HMMs</b> )
Conserved Domain Database ( <u>CDD</u> )	NCBI	Position-specific scoring matrices ( <b>PSSMs</b> )



Search  for    Help

## Conserved Domains and Protein Classification

[RESOURCES](#) [SEARCH](#) [HOW](#)

### Resources

#### Conserved Domain Database (CDD)

CDD is a protein annotation resource that consists of a collection of well-annotated multiple sequence alignment models for ancient domains and full-length proteins. These are available as position-specific score matrices (PSSMs) for fast identification of conserved domains in protein sequences via RPS-BLAST. CDD content includes NCBI-curated domains, which use 3D-structure information to explicitly to define domain boundaries and provide insights into **sequence/structure/function relationships**, as well as domain models imported from a number of external source databases (Pfam, SMART, COG, PRK, TIGRFAM).

[Search](#) | [How To](#) | [Help](#) | [News](#) | [FTP](#) | [Publications](#)

#### CD-Search & Batch CD-Search

CD-Search is NCBI's interface to searching the Conserved Domain Database with protein query sequences. It uses RPS-BLAST, a variant of PSI-BLAST, to quickly scan a set of pre-calculated position-specific scoring matrices (PSSMs) with a protein query. The results of CD-Search are presented as an annotation of protein domains on the user query sequence (**illustrated example**), and can be visualized as domain multiple sequence alignments with embedded user queries. High confidence associations between a query sequence and conserved domains are shown as **specific hits**.

[CD-Search](#) | [Batch CD-Search](#) | [Help](#) | [FTP](#) | [Publications](#)

Search Database

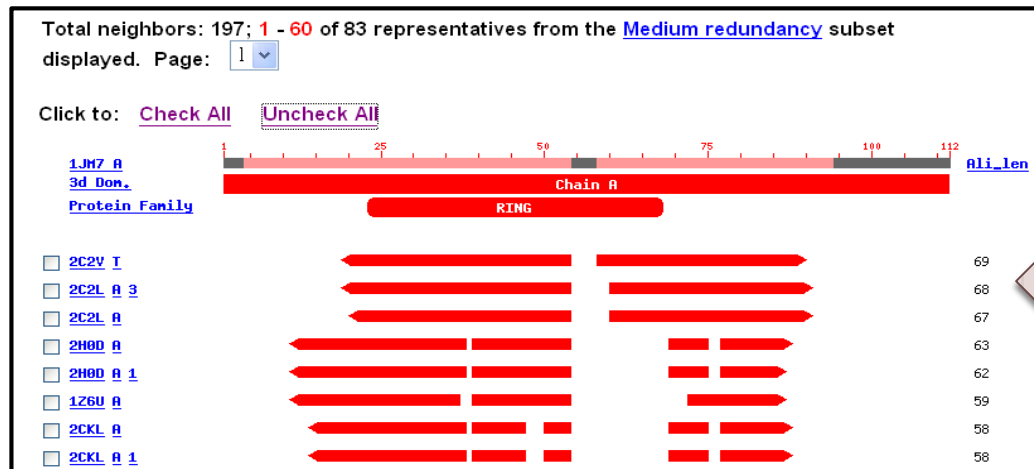
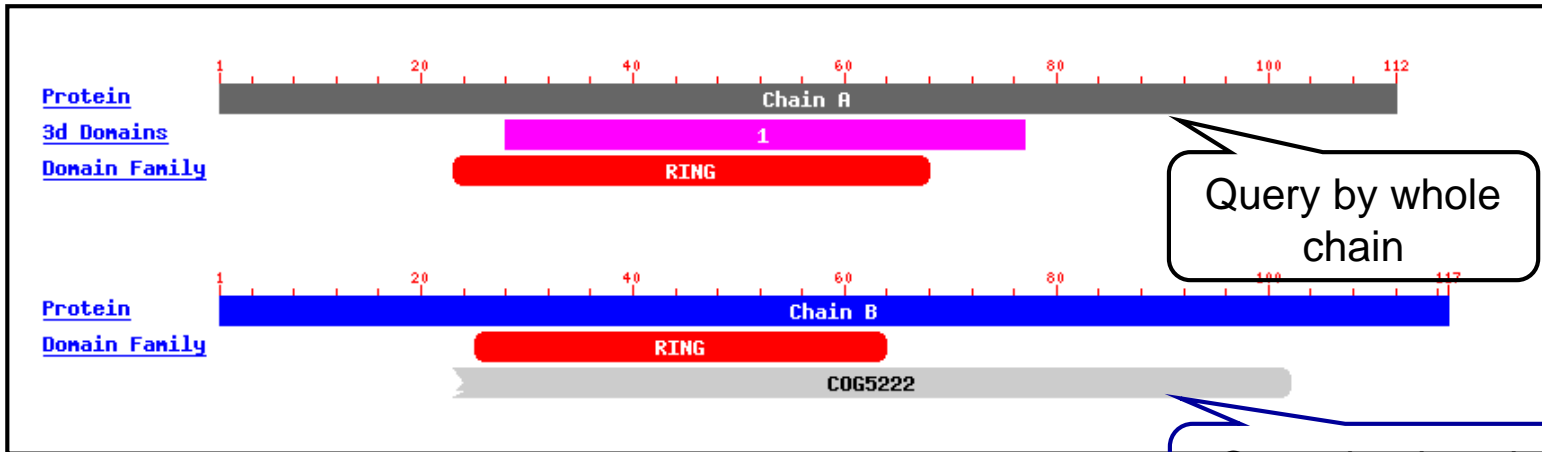
- CDD v2.28 - 39357 PSSMs
- SMART v5.1 - 791 PSSMs
- Pfam v24.0 - 11912 PSSMs
- COG v1.00 - 4873 PSSMs
- KOG v1.00 - 4825 PSSMs
- PRK v6.00 - 10885 PSSMs
- TIGR v10.00 - 4023 PSSMs

#### CDART: Conserved Domain Architectures

Conserved Domain Architecture database based on domain architecture queries. CDART finds protein similarities

forms similarity searches of the Entrez Protein database in the natural order of conserved domains in protein sequences. CDART finds evolutionary distances using sensitive domain

# VAST: Query by Chain or 3D Domain



Query by domain  
COG5222

Not found with  
chain query

# Synthetic Biology

## What is Synthetic Biology?

Biology,  
Engineering  
and  
Informatics

DNA



Proteins



Cells



High Value  
Applications

Human Therapeutics  
Industrial Products  
Agriculture  
Animal Sciences/  
Aquaculture  
Protein Production

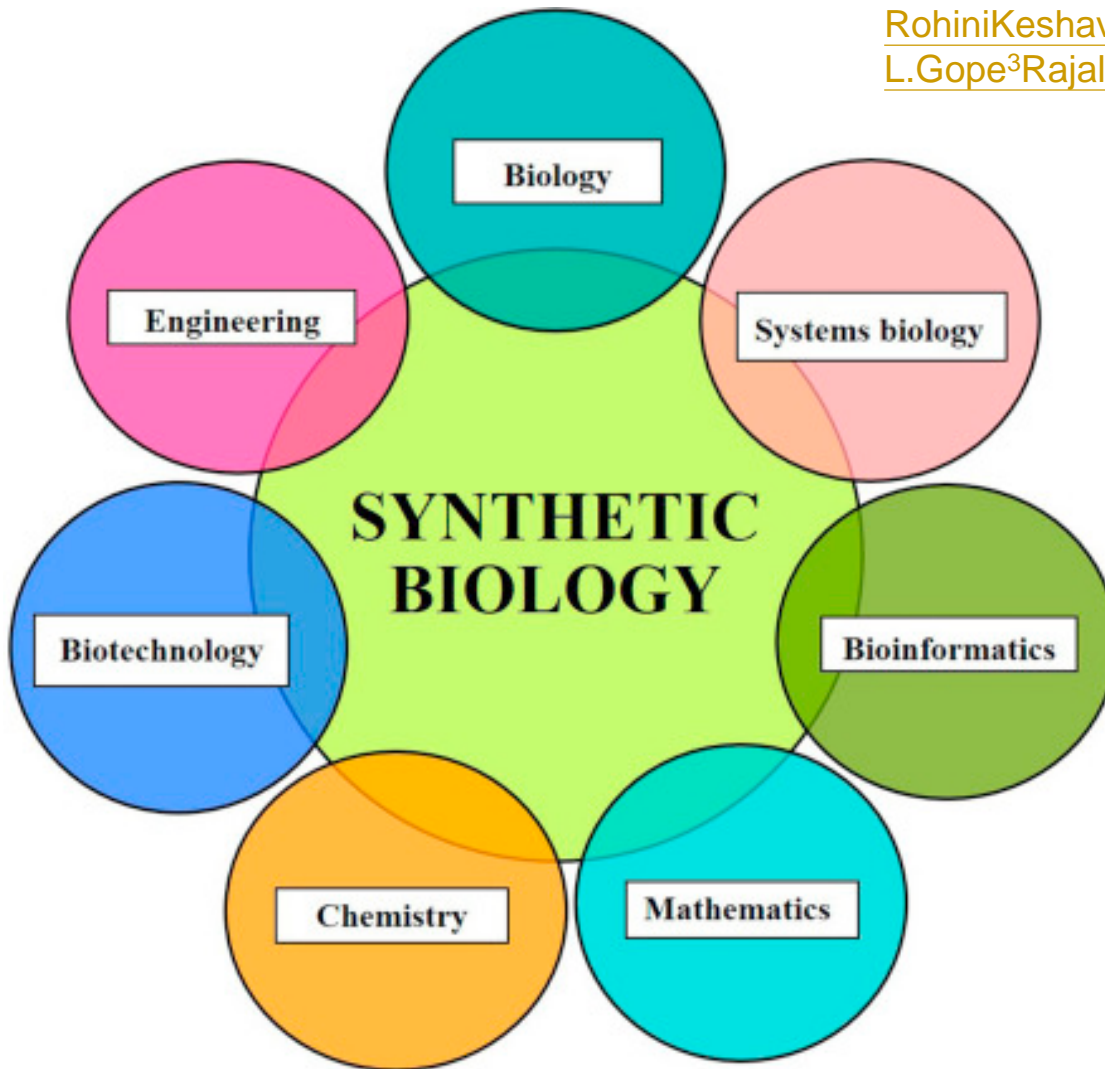




# Synthetic Biology

Synthetic biology is a field of science that involves redesigning organisms for useful purposes by engineering them to have new abilities. Synthetic biology researchers and companies around the world are harnessing the power of nature to solve problems in medicine, manufacturing and agriculture.



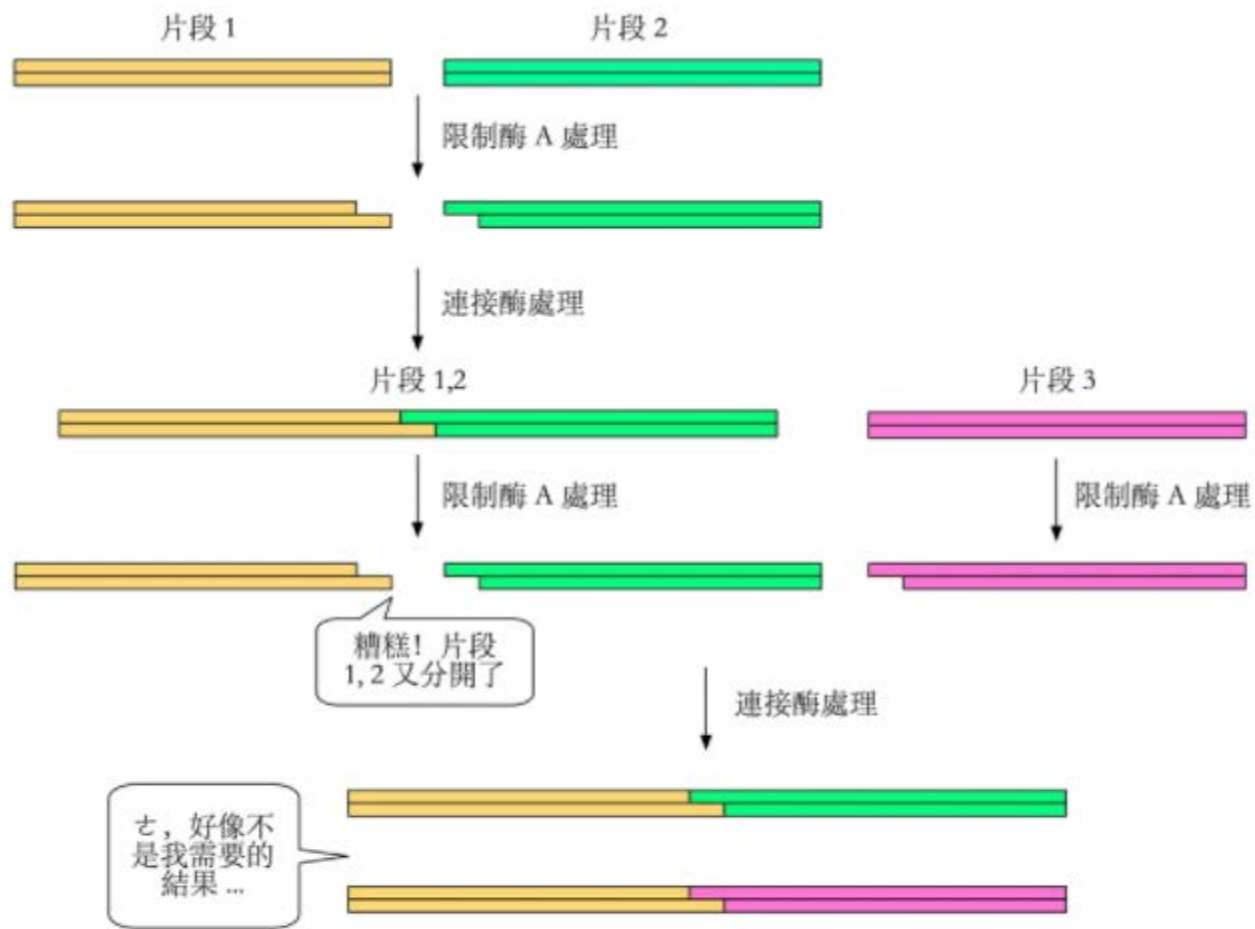


# Example 1

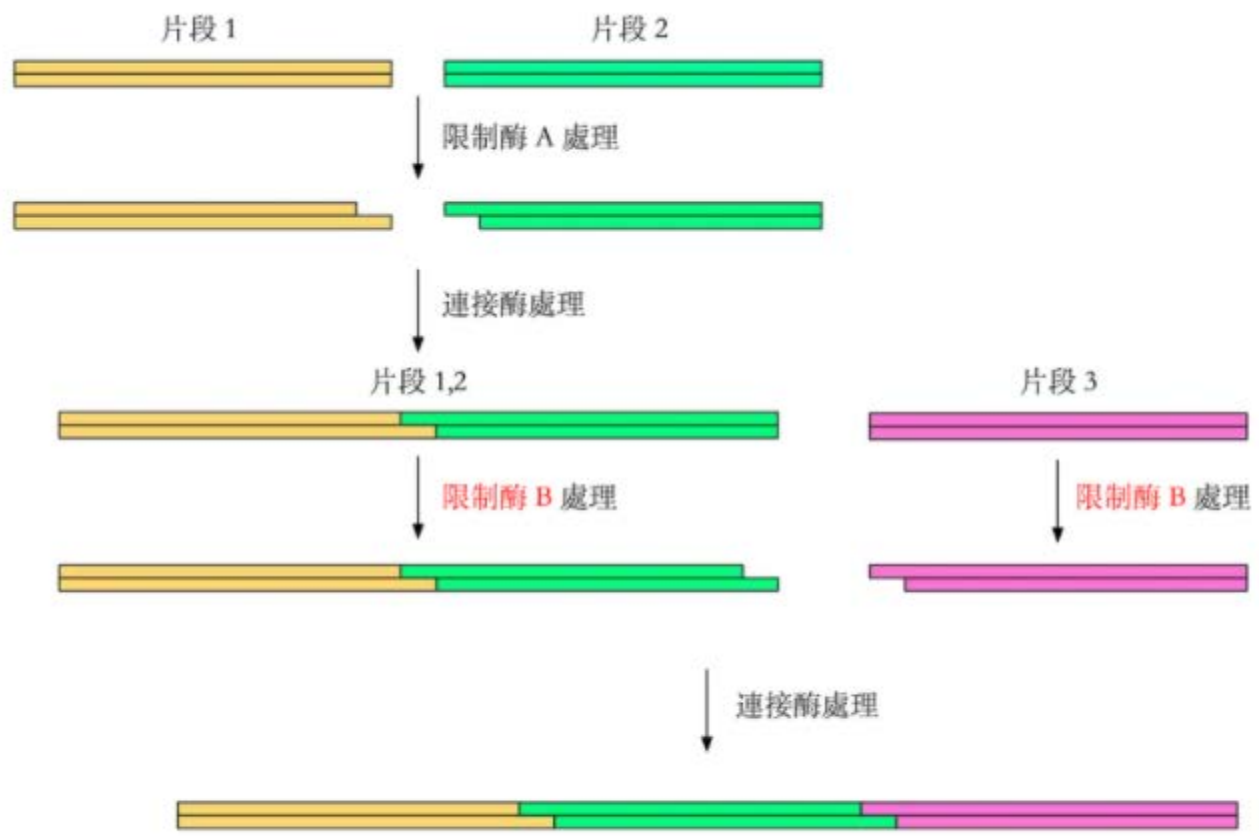
- 因為每合成一個鹼基對 (base pair, bp) DNA 的價格，三十年前要價數十至數百美元不等，而如今降低到只需要一美元或低於一美元，有人將這種現象比擬為生命科學研究上的摩爾定律。
- DNA 合成技術的成熟，大大降低了DNA 合成的經濟門檻，也預告著大尺度基因體工程與合成生物學研究時代的來臨。
- 2008 年，JCVI (J. Craig Venter Institute) 的研究人員用 5000 ~ 7000 bp 大小的化學合成DNA 片段 (chemically synthesized DNA fragments)，以人工方式兩兩相連接組裝成一個 582,970 bp 的 *Mycoplasma genitalium* 細菌基因體。

Features · iGEM合成生物學大賽 · 合成生物學 · 研究領域專題 · 編輯團隊的話

合成生物學專題



圖一 失敗的三段組裝



圖二 成功的三段組裝

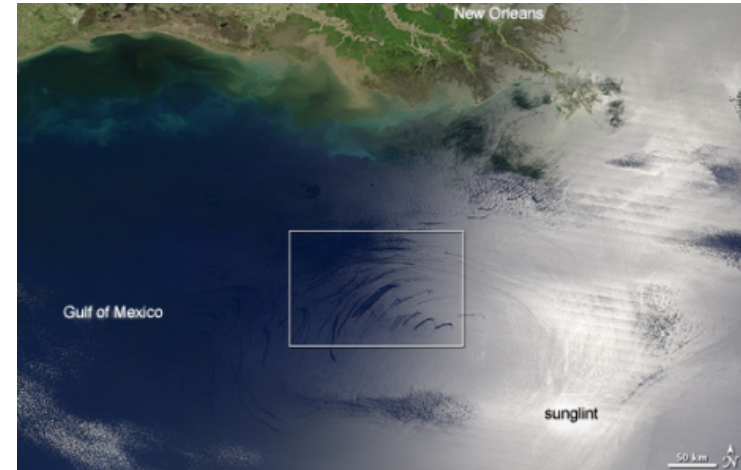
## Example 2



- Tom Knight 教授提出一種標準的 DNA 片段的組裝方式<sup>9,10</sup>，在每次的組裝可以使用相同的方式，不需要再費心選擇每次組裝使用的限制酶酵素。這樣的組裝方式，讓 DNA 片段可以像積木一樣，一個片段一個片段一直連續組裝下去，生物零件 (biological Part) 的概念就因此誕生。
- 將生物 DNA 片段零件化，是工程思維應用在分子生物學的一個重大發明。因此，透過生物零件的定義與標準化的組裝方式，我們可以進一步組裝生物設備 (biological device)，或更進一步可以組裝一個生物系統 (biological system)，形成一個由生物零件為基礎的工程框架<sup>11</sup>。

# Example 3: microorganisms harnessed for bioremediation

While invisible up close, microscopic oil **slicks** 浮油 from natural seeps 渗透 are visible from **space** because cohesion 凝聚 between oil molecules flattens wave action to **form smooth areas** on the water (2010, BP)



Petroleum-degrading microbes called *Oceanospirillales*

# Oil-eating microbes

Naturally occurring microbes in the ocean feed on the hydrocarbons in oil. Scientists hope to speed up the process for the large spill in the Gulf of Mexico, where warm temperatures also aid the reaction.



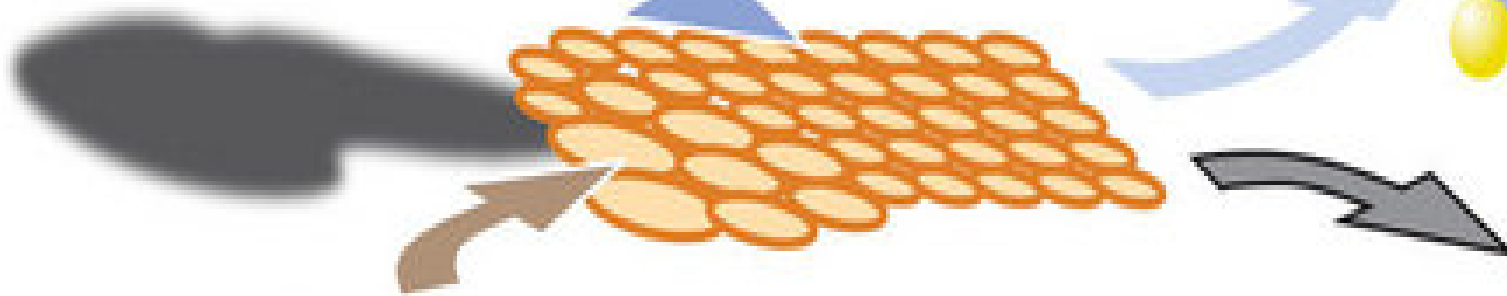
Oil contains hydrocarbons, which are made up of varying amounts of carbon and hydrogen



Oxygen is needed for the chemical reaction, but can be sparse at great ocean depths



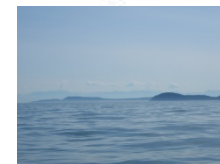
The microbes break apart the hydrocarbons and combine them with oxygen to create water and carbon dioxide



**Adding fertilizer** increases the size and number of the microbes so they can eat more oil; too much, however, can cause algae blooms, which starve the ecosystem of light and oxygen

**Not all of the oil** can be consumed, but what is left over is more easily dispersed by currents and wind

Source: Terry Hazen, Lawrence Berkeley National Lab  
Graphic: Miami Herald



© 2010 MCT

## Example 4: Rice modified to produce beta-carotene, a nutrient usually associated with carrots, that prevents vitamin A deficiency



Vitamin A deficiency causes blindness in 250,000 - 500,000 children every year and greatly increases a child's risk of death from infectious diseases.



# Example 5: Yeast engineered to produce rose oil as an eco-friendly and sustainable substitute for real roses that perfumers use to make luxury scents

## Engineered yeast could replace flowers in fragrances

By Michelle Yeomans [↗](#)

19-Mar-2015 - Last updated on 19-Mar-2015 at 13:43 GMT



RELATED TAGS: Synthetic biology, Dna

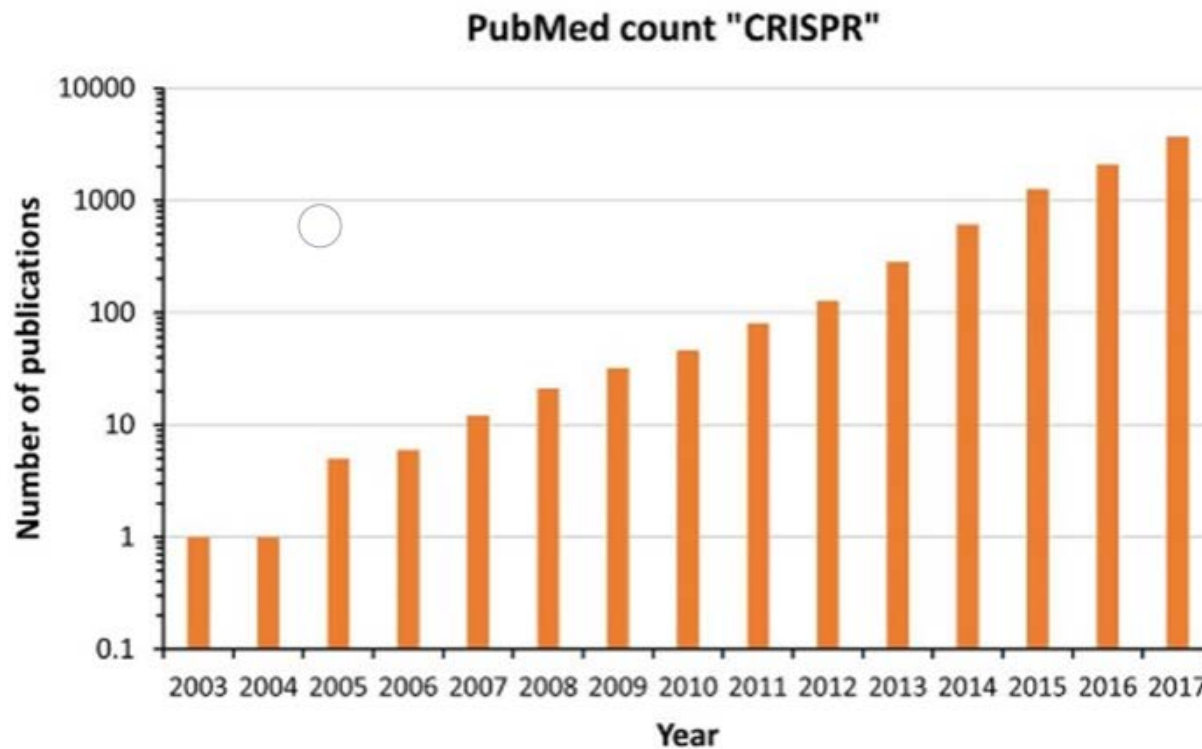
Boston-based specialists in synthetic biology Ginkgo Biowork is using yeast to produce fragrances that are cheaper than using naturally sourced ingredients.

# Genome Editing

# What is the difference between synthetic biology and genome editing?

- In some ways, synthetic biology is similar to another approach called "**genome editing**" because both involve **changing** an organism's genetic code; however, some people draw a distinction between these two approaches based on how that change is made
- In **synthetic biology**, scientists typically **stitch together** long stretches of DNA and insert them into an organism's genome.
  - These synthesized pieces of DNA could be genes that are found in other organisms or they could be entirely novel
- In **genome editing**, scientists typically use tools to **make smaller changes** to the organism's own DNA. Genome editing tools can also be used to delete or add small stretches of DNA in the genome.

# CRISPR/Cas9 Applications are Exploding and Revolutionize Molecular Biology



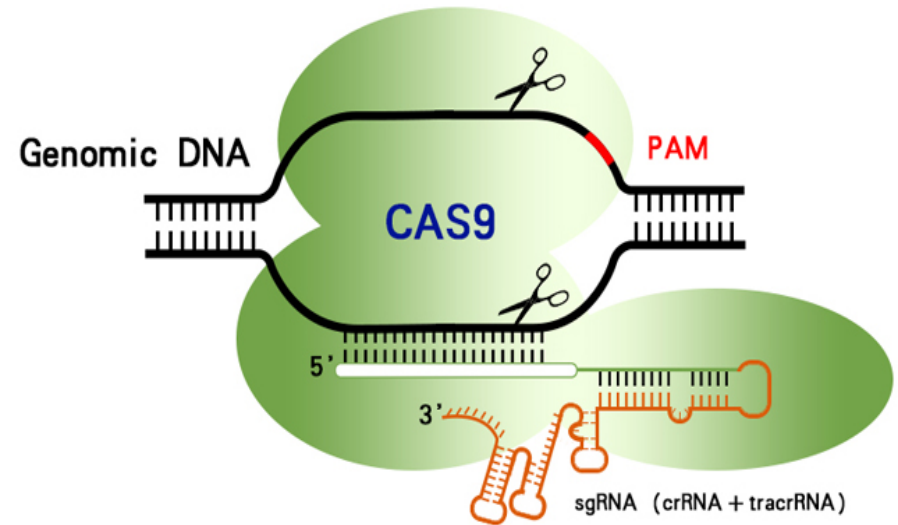
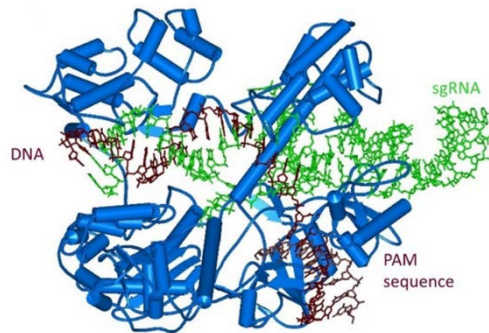
# Clustered Regularly Interspaced Short Palindromic Repeats (CRISPR)

- A genome editing technique that
  - Targets a specific section of DNA
  - Make a precise cut/break at the target site
  
- Applications
  1. To make a gene nonfunctional (knockout)
  2. Replace on version of a gene with another
    - E.g., gene therapy
    - David Vetter was born without a **functioning immune system** and spent his life in a bubble that protected him from germs. He died at age 12 in 1984. Scientists are using gene therapy to treat the disorder so that children can live normally.

**Adenosine Deaminase (ADA)**



# CRISPR/Cas9 Applications are Exploding and Revolutionize Molecular Biology



Structure of *staphylococcus aureus* Cas9 (blue) bound to single guide RNA (green) & targeted DNA (brown) (Nishimasu et al. 2015)

- Non-coding RNAs & Cas protein
- Protospacer adjacent motif (PAM) is a 2-6 base pair DNA sequence immediately following the DNA sequence targeted by the Cas9 nuclease in the CRISPR bacterial adaptive immune system
- sgRNA = single guide RNA = a targeting sequence (crRNA sequence) + (a Cas9 nuclease-recruiting sequence: tracrRNA)



CRISPR: Gene editing and beyond

135 [s://www.youtube.com/watch?v=4YKFw2KZA5o](https://www.youtube.com/watch?v=4YKFw2KZA5o)

A graphic featuring a bright sun with rays in the center, set against a blue sky with white clouds. The sun's rays create a lens flare effect. The word 'CRISPR' is written in large, bold, blue capital letters across the middle of the image. Below it, the acronym is expanded into its full name: 'Clustered Regularly Interspaced Short Palindromic Repeats', with each word starting with a red letter that corresponds to the letter in the acronym above.

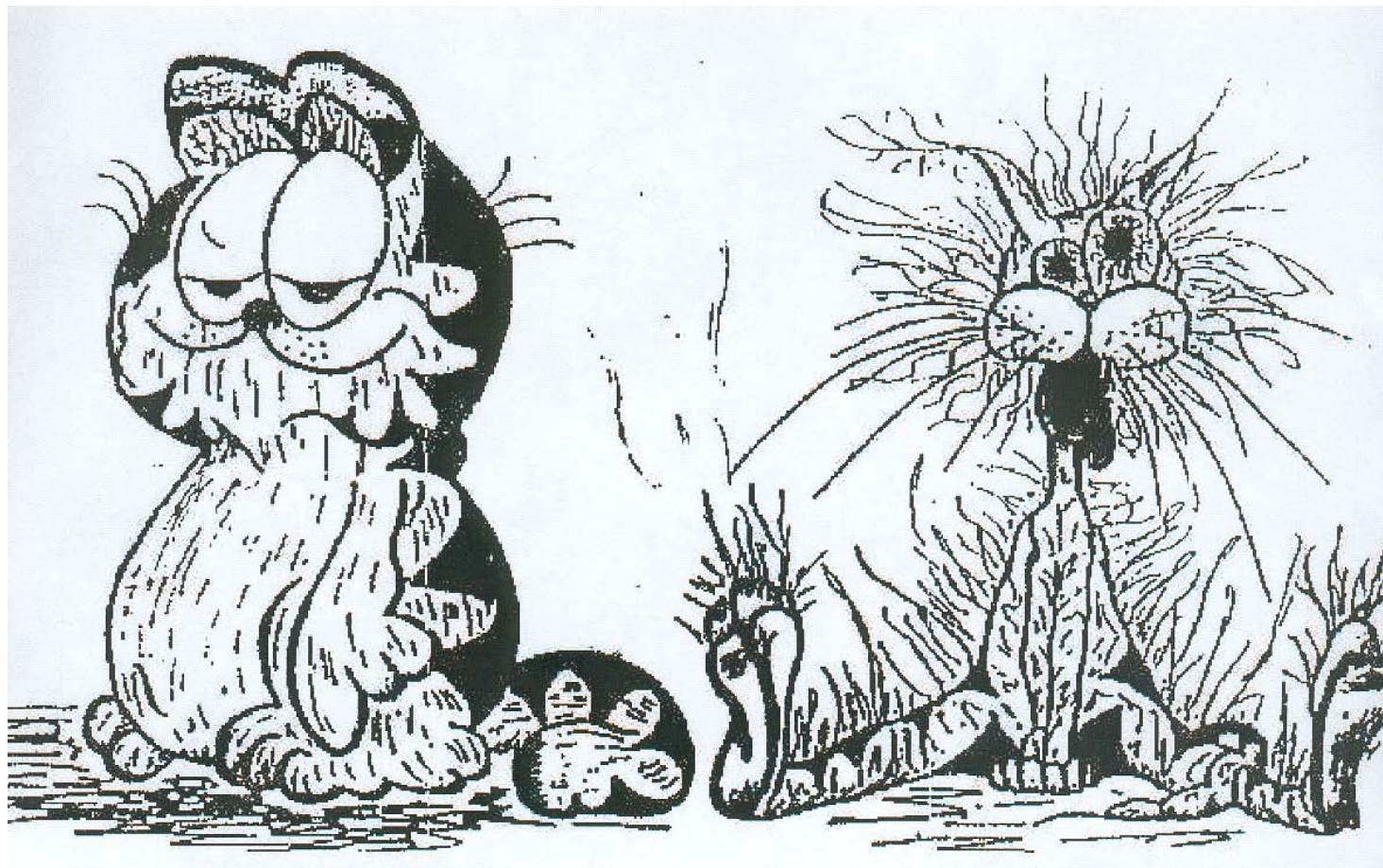
# CRISPR

*Clustered Regularly Interspaced Short Palindromic Repeats*

Genetic Engineering Will Change Everything Forever – CRISPR



**Before...**



**After...**