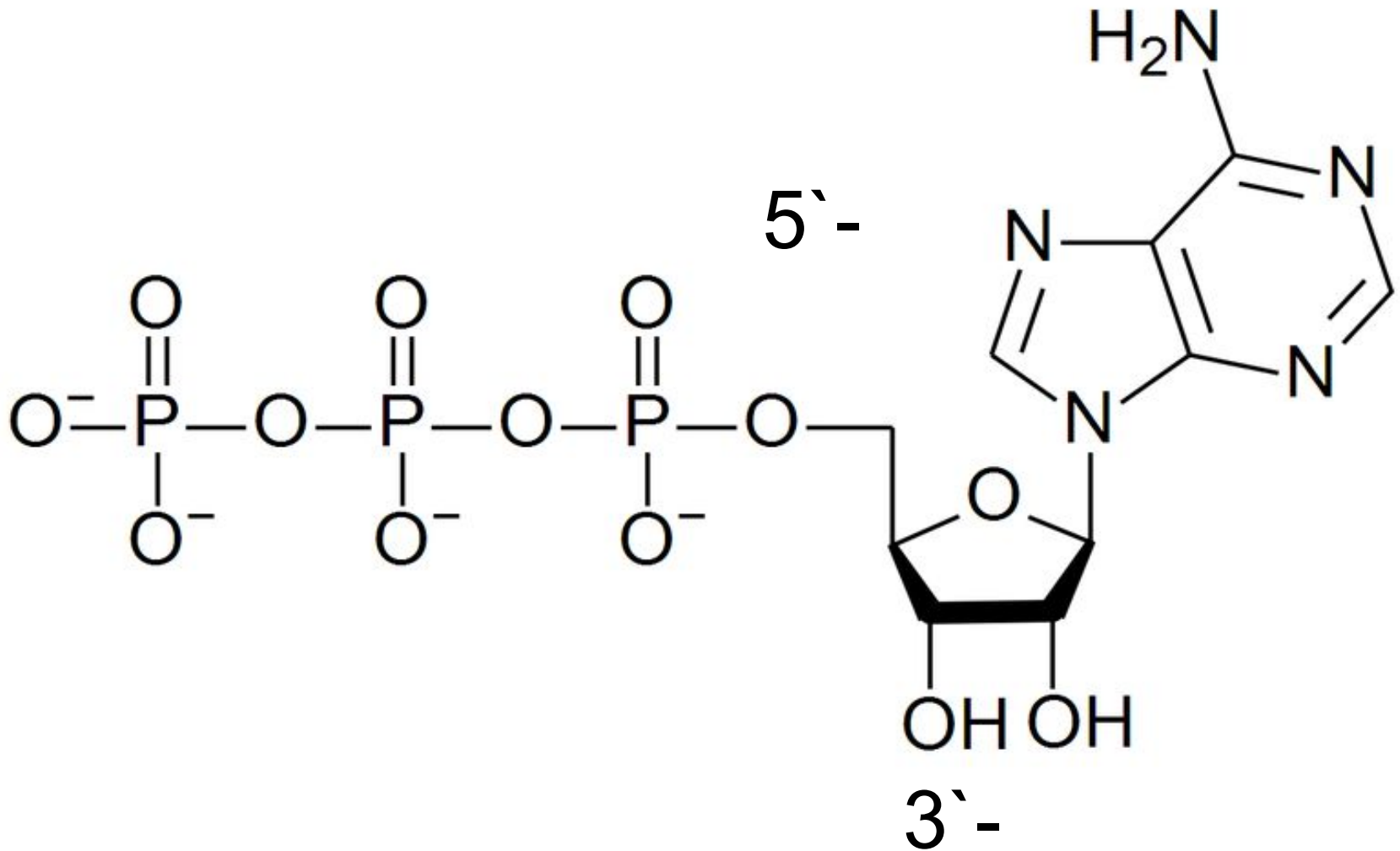


# DNA vs. computer

1. Про 5'-3' и всякую химию
2. Про банки данных (архивные vs. курируемые)
3. Святая троица EMBL – GenBank – DDBJ
4. Собственно EMBL, его разделы, классы данных и поля; CDS, кодирующие участки, ссылки из Swiss-Prot.

# АТФ

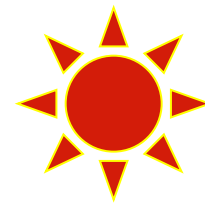


# Как записывают последовательности нуклеиновых кислот ?

1. Последовательность = последовательность однобуквенных символов.  
Никаких дефисов и обозначений фосфодиэфирных связей.
2. Одни и те же однобуквенные символы для последовательностей РНК и ДНК (при записи РНК обычно 'U'  $\Rightarrow$  'T').  
Любая последовательность по умолчанию считается ДНК (т.е. полимером 2'-дезоксирибонуклеотидов).
3. Одни и те же символы используются для обозначения азотистых оснований, нуклеозидов и нуклеотидов  
Допустимы заглавные и строчные буквы, хотя рекомендованы заглавные.
4. **Последовательность записывается в направлении 5'  $\rightarrow$  3'**

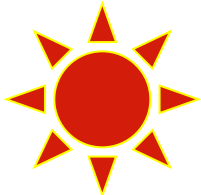
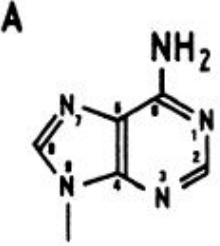
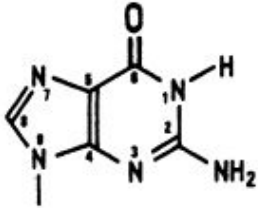

Пример:

5'-СТСГАС-3'



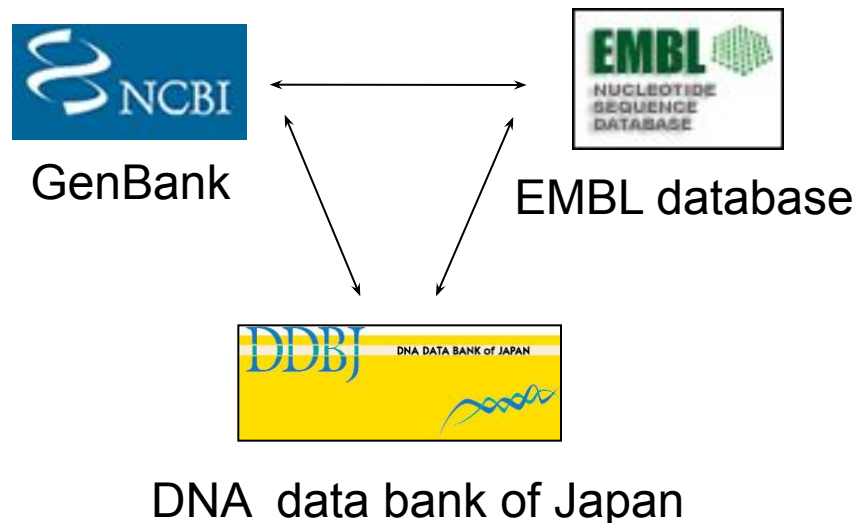
Nomenclature Committee of the International Union of Biochemistry (NC-IUB)  
Nomenclature for incompletely specified bases in nucleic acid sequences  
Recommendations 1984  
Biochem. J. (1985) 229, 281-286

# Общепринятые однобуквенные обозначения для стандартных азотистых оснований (остатков нуклеозидов и нуклеотидов) и вырожденных позиций в выравниваниях нуклеиновых кислот

Символ	Расшифровка	Происхождение обозначения	
<b>G</b>	G, guanine (гуанин)		
<b>A</b>	A, adenine, (аденин)		
<b>T</b>	T/U, thymine/uracyl (тимин в ДНК и урацил в РНК)		
<b>C</b>	C, cytosine (цитозин)		
<b>R</b>	A или G, purine(пурины)	pu <b>R</b> ine	
<b>Y</b>	C или T или U, pyrimidine (пиримидин)	p <b>Y</b> rimidine	
<b>M</b>	A или C	a <b>M</b> ino	
<b>K</b>	G или T	<b>K</b> eto	
<b>S</b>	G или C	<b>S</b> trong interaction (3 H bonds)	
<b>W</b>	A или T	<b>W</b> weak interaction (2 H bonds)	
<b>H</b>	A или C или T, но не G	в алфавите 'H' следует за 'G'	
<b>B</b>	G или T или C, но не A	'B' следует за 'A'	
<b>V</b>	G или C или A, но не T или U	'V' следует за 'U'	
<b>D</b>	G или A или T, но не C	'D' следует за 'C'	
<b>N</b>	G или A или T или C	a <b>N</b> y	

# NCBI и EBI

- National Center for Biotechnology Information и European Bioinformatics Institute (подразделение EMBL – European Molecular Biology Laboratory)
- Три базы данных – GenBank, EMBL и DDBJ (японская) – по сути, одно и то же.



# Что надо знать про банк EMBL

- что это архив (за содержание записи несёт ответственность её автор)
  - поэтому разноречивой терминологии
  - поэтому одно и то же по многу раз
  - поэтому много неисправленных ошибок
- что у последовательности из записи часто нет естественных границ
- что это часть триединства (EMBL, GenBank, DDBJ)
  - ежедневный обмен данными
- ... ну и смысл основных полей, конечно (особенно структуру поля FT!)

ttttacctcttttagtgatattgtgatagagcaaaaatcccgcacattgtgctcgggattgttttaaaccttgttgatttaattttcaatcgcttctttataaagaagtagtggtgccc  
acaacactcacattgcatatcaatacggcctttatgttcggctaataatttcgcaatttcttcatcagagatgagcagtagatgcagaactagaacgctcagcagagcagccaca  
gaaaaattgtacatcttgtgctggataaagattaacggtttctcgtgatataaacgataggagtaacttctcagggagaccaataattcttcatctttactgttgctgcgagc  
gtagttaaagtctcaaaatcttctggtgtaccagaaccatcaggcataattgtaataacatacctgctgcccactggctgcttcatattctccagtacgaataattaattgagttg



GenBank



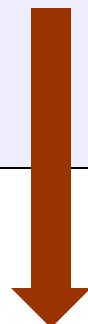
EMBL



DDBJ



компьютерный поиск гена, трансляция и компьютерная аннотация



Базы данных научной литературы



~2 500 000 последовательностей

Экспертиза



~200 000 последовательностей



UniParc  
(UniProt Archive)

UniRef  
(UniProt non-redundant Reference databases)



# Класс данных

## 3.1 Data Class

The data class of each entry, representing a methodological approach to the generation of the data or a type of data, is indicated on the first (ID) line of the entry. Each entry belongs to exactly one data class.

Class	Definition
CON	Entry constructed from segment entry sequences, drawing annotation from segment entries
ANN	Entry constructed from segment entry sequences with its own annotation
PAT	Patent
EST	Expressed Sequence Tag
GSS	Genome Survey Sequence
HTC	High Throughput CDNA sequencing
HTG	High Throughput Genome sequencing
MGA	Mass Genome Annotation
WGS	Whole Genome Shotgun
TPA	Third Party Annotation
STS	Sequence Tagged Site
STD	Standard (all entries not classified as above)



# Таксономический раздел

## 3.2 Taxonomic Division

The entries which constitute the database are grouped into taxonomic divisions, the object being to create subsets of the database which reflect areas of interest for many users.

In addition to the division, each entry contains a full taxonomic classification of the organism that was the source of the stored sequence, from kingdom down to genus and species (see below).

Each entry belongs to exactly one taxonomic division. The ID line of each entry indicates its taxonomic division, using the three letter codes shown below:

Division	Code
-----	----
Bacteriophage	PHG
Environmental Sample	ENV
Fungal	FUN
Human	HUM
Invertebrate	INV
Other Mammal	MAM
Other Vertebrate	VRT
Mus musculus	MUS
Plant	PLN
Prokaryote	PRO
Other Rodent	ROD
Synthetic	SYN
Transgenic	TGN
Unclassified	UNC
Viral	VRL

# Поле

ID - identification (begins each entry; 1 per entry)  
AC - accession number ( $\geq 1$  per entry)  
PR - project identifier (0 or 1 per entry)  
DT - date (2 per entry)  
DE - description ( $\geq 1$  per entry)  
KW - keyword ( $\geq 1$  per entry)  
OS - organism species ( $\geq 1$  per entry)  
OC - organism classification ( $\geq 1$  per entry)  
OG - organelle (0 or 1 per entry)  
RN - reference number ( $\geq 1$  per entry)  
RC - reference comment ( $\geq 0$  per entry)  
RP - reference positions ( $\geq 1$  per entry)  
RX - reference cross-reference ( $\geq 0$  per entry)  
RG - reference group ( $\geq 0$  per entry)  
RA - reference author(s) ( $\geq 0$  per entry)  
RT - reference title ( $\geq 1$  per entry)  
RL - reference location ( $\geq 1$  per entry)  
DR - database cross-reference ( $\geq 0$  per entry)  
CC - comments or notes ( $\geq 0$  per entry)  
AH - assembly header (0 or 1 per entry)  
AS - assembly information (0 or  $\geq 1$  per entry)  
FH - feature table header (2 per entry)  
FT - feature table data ( $\geq 2$  per entry)  
XX - spacer line (many per entry)  
SQ - sequence header (1 per entry)  
CO - contig/construct line (0 or  $\geq 1$  per entry)  
bb - (blanks) sequence data ( $\geq 1$  per entry)  
// - termination line (ends each entry; 1 per entry)

# FT

FT Key Location/Qualifiers=value

FT CDS 1..1000  
/codon=(seq:"cug",aa:Ser)  
/codon=(seq:"tga",aa:Trp)

<http://www.ebi.ac.uk/embl/WebFeat/index.html>



# CDS и exons

CDS – кодирующая последовательность, то есть ровно те нуклеотиды, что соответствуют белку (по крайней мере его основной форме).

Кодирующие участки – те фрагменты ДНК, из которых составлен CDS.

Exons – экзоны, то из чего будет составлена зрелая матричная РНК, они включают в себя 5` и 3` - нетранслируемые области – те части РНК, которые отвечают за регуляцию трансляции.

# Ссылки из записи Swiss-Prot на EMBL

(2001) <i>Proteins</i> <b>44</b> :270-281	
Position	X-RAY CRYSTALLOGRAPHY (1.89 ANGSTROMS).
Medline	<a href="#">21348730</a>
DOI	<a href="#">10.1002/prot.1092</a> ;
PubMed	<a href="#">11455600</a>   <a href="#">CiteXplore</a>

## Comments

<b>FUNCTION</b>	Decomposes hydrogen peroxide into water and oxygen; serves to protect cells from the toxic effects of hydrogen peroxide.
<b>CATALYTIC ACTIVITY</b>	2 H(2)O(2) = O(2) + 2 H(2)O.
<b>COFACTOR</b>	Heme group.
<b>SUBUNIT</b>	Homotetramer.
<b>SUBCELLULAR LOCATION</b>	Cytoplasm (Probable).
<b>INDUCTION</b>	By entry into stationary phase.
<b>SIMILARITY</b>	Belongs to the catalase family, HP11 subfamily.

## Copyright

Copyrighted by the UniProt Consortium, see <http://www.uniprot.org/terms> Distributed under the Creative Commons Attribution-NoDerivs License

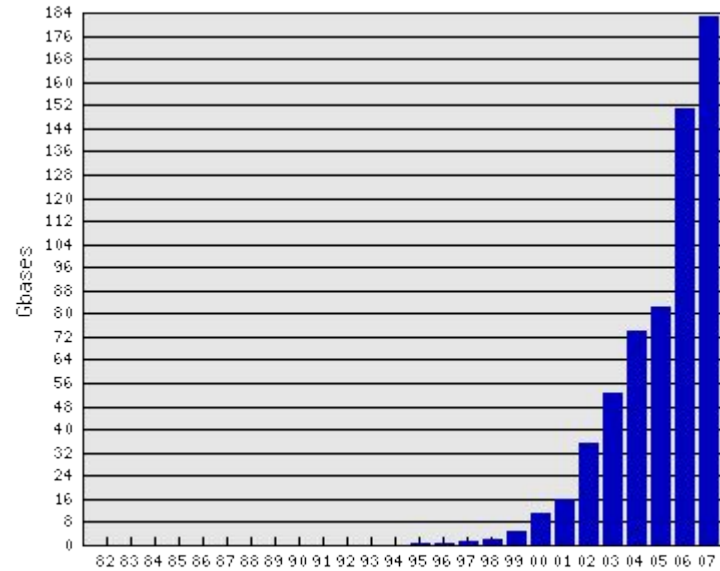
## Database cross-references

EMBL	<a href="#">M55161</a> ; <a href="#">AAA24039.1</a> ; -, Genomic_DNA. <a href="#">U00096</a> ; <a href="#">AAT48137.1</a> ; -, Genomic_DNA. <a href="#">AP009048</a> ; <a href="#">BAA15513.1</a> ; -, Genomic_DNA.
PIR	<a href="#">A39129</a> ; A39129.
RefSeq	<a href="#">AP_002351.1</a> ; -. <a href="#">YP_025308.1</a> ; -.
PDB	<a href="#">1CF9</a> ; X-ray; 1.80 A; A/B/C/D=27-753. <a href="#">1GG9</a> ; X-ray; 1.89 A; A/B/C/D=1-753. <a href="#">1GGE</a> ; X-ray; 1.89 A; A/B/C/D=1-753. <a href="#">1GGE</a> ; X-ray; 2.28 A; A/B/C/D=1-753. <a href="#">1GGH</a> ; X-ray; 2.15 A; A/B/C/D=1-753. <a href="#">1GGJ</a> ; X-ray; 1.92 A; A/B/C/D=1-753. <a href="#">1GGK</a> ; X-ray; 2.26 A; A/B/C/D=1-753. <a href="#">1IPH</a> ; X-ray; 2.80 A; A/B/C/D=1-753. <a href="#">1P7Y</a> ; X-ray; 2.40 A; A/B/C/D=1-753.

Каждая строка – отдельный сиквенс; первая ссылка в строке – запись в EMBL, вторая – CDS внутри этой записи (здесь идентификатор, например, AAA24039.1 – это идентификатор CDS в специальном дочернем банке данных EMBL-Coding sequences).

# Статистика EMBL

Total nucleotides



Number of entries

