# EKSPERIMENTALNE METODE V BIOKEMIJI IN MOLEKULARNI BIOLOGIJI

**Določanje nukleotidnih zaporedij**

**Sangerjeva metoda, metode druge in tretje generacije**

**Genomika**

**Podatkovne zbirke nukleotidnih zaporedij: GenBank**

**EST kloni**
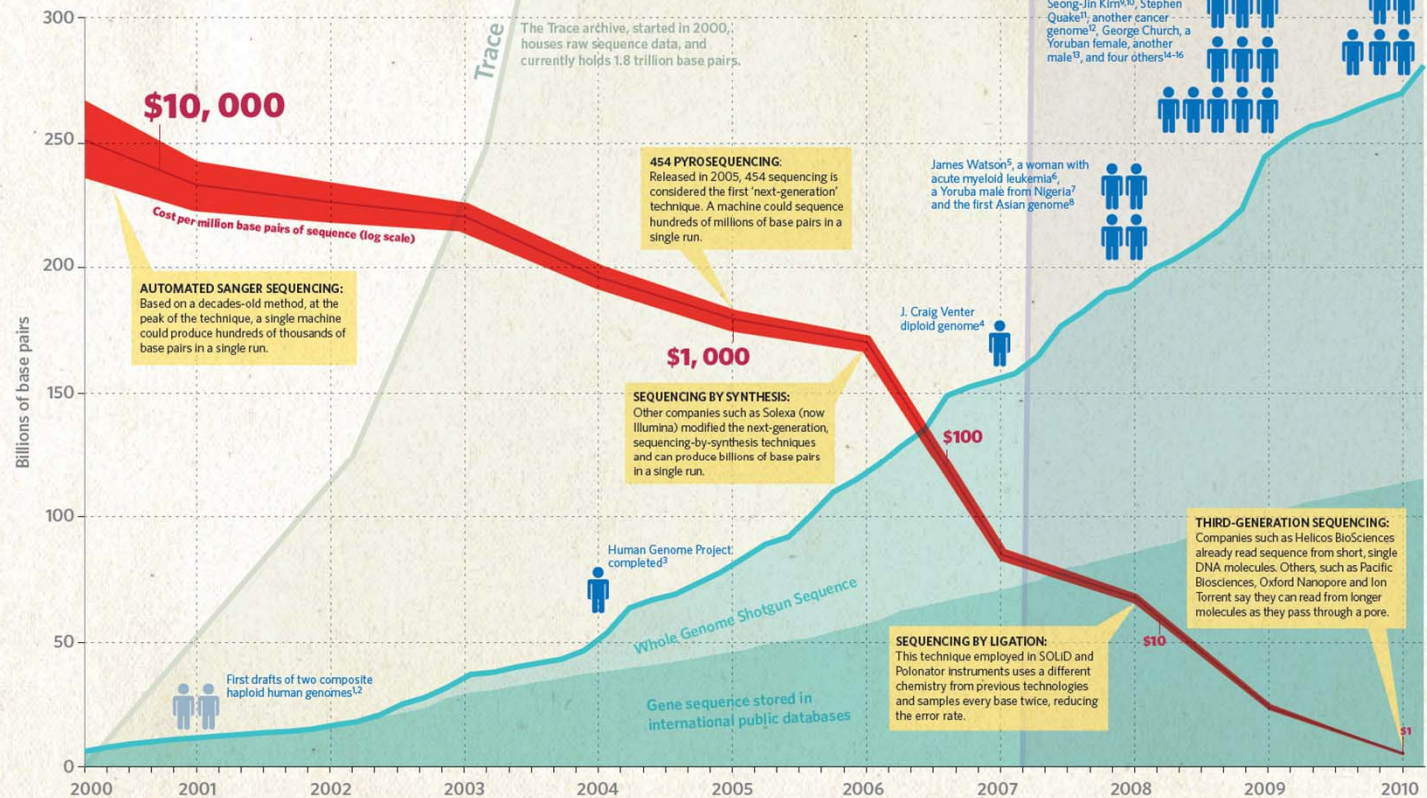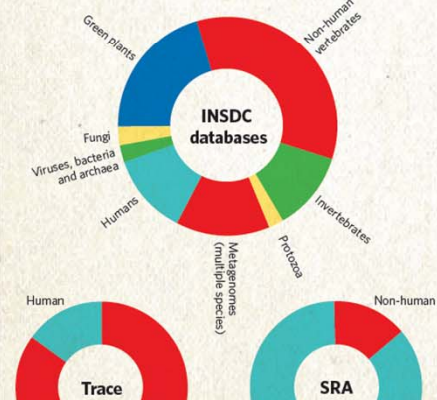
**Projekt človeški genom**

THE SEQUENCE EXPLOSION

A t the time of the announcement of the first drafts of the human genome in 2000, there were 8 billion base pairs of sequence in the three main databases for 'finished' sequence: GenBank, run by the US National Center for Biotechnology Information; the DNA Databank of Japan; and the European Molecular Biology Laboratory (EMBL) Nucleotide Sequence Database. The databases share their data regularly as part of the International Nucleotide Sequence Database Collaboration (INSDC). In the subsequent first post-genome decade, they have added another 270 billion bases to the collection of finished sequence, doubling the size of the database roughly every 18 months. But this number is dwarfed by the amount of raw sequence that has been created and stored by researchers around the world in the Trace archive and Sequence Read Archive (SRA). See Editorial, page 649, and human genome special at www.nature.com/humangenome

DNA SEQUENCES BY TAXONOMY

International Nucleotide Sequence Database Collaboration: The main repositories of 'finished' sequence span a wide range of organisms, representing the many priorities of scientists worldwide.

INSDC databases

Green plants
Non-human vertebrates
Fungi
Viruses, bacteria and archaea
Humans
Invertebrates
Protozoa
Metagenomes (multiple species)

Human
Non-human

Trace
SRA

The Sequence Read Archive (SRA) houses raw data from next-generation sequencing and has grown to 25 trillion base pairs. If this chart were to accommodate it, it would stretch to more than 12 metres — twice the height of an average giraffe.

A glioma cell line[17], Inuk[18], !Gubi and Archbishop Desmond Tutu[19], James Lupski[20], and a family of four[21]

Two Korean males including Seong-Jin Kim[9,10], Stephen Quake[11], another cancer genome[12], George Church, a Yoruban female, another male[13], and four others[14–16]

Trace

The Trace archive, started in 2000, houses raw sequence data, and currently holds 1.8 trillion base pairs.

$10,000

Cost per million base pairs of sequence (log scale)

454 PYROSEQUENCING: Released in 2005, 454 sequencing is considered the first 'next-generation' technique. A machine could sequence hundreds of millions of base pairs in a single run.

James Watson[5], a woman with acute myeloid leukemia[6], a Yoruba male from Nigeria[7] and the first Asian genome[8]

AUTOMATED SANGER SEQUENCING: Based on a decades-old method, at the peak of the technique, a single machine could produce hundreds of thousands of base pairs in a single run.

$1,000

J. Craig Venter diploid genome[4]

SEQUENCING BY SYNTHESIS: Other companies such as Solexa (now Illumina) modified the next-generation, sequencing-by-synthesis techniques and can produce billions of base pairs in a single run.

$100

THIRD-GENERATION SEQUENCING: Companies such as Helicos BioSciences already read sequence from short, single DNA molecules. Others, such as Pacific Biosciences, Oxford Nanopore and Ion Torrent say they can read from longer molecules as they pass through a pore.

Human Genome Project completed[3]

Whole Genome Shotgun Sequence

SEQUENCING BY LIGATION: This technique employed in SOLiD and Polonator instruments uses a different chemistry from previous technologies and samples every base twice, reducing the error rate.

$10

First drafts of two composite haploid human genomes[1,2]

Gene sequence stored in international public databases

$1

Billions of base pairs

300
250
200
150
100
50

2000  2001  2002  2003  2004  2005  2006  2007  2008  2009  2010

Nature, 1. 4. 2010

# Primerjava metod za določanje nukleotidnih zaporedij
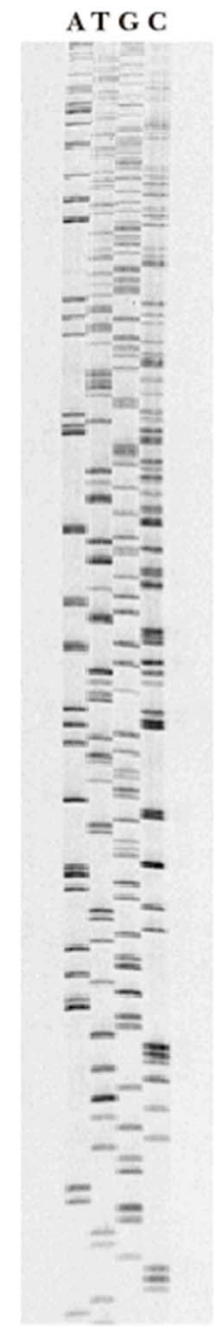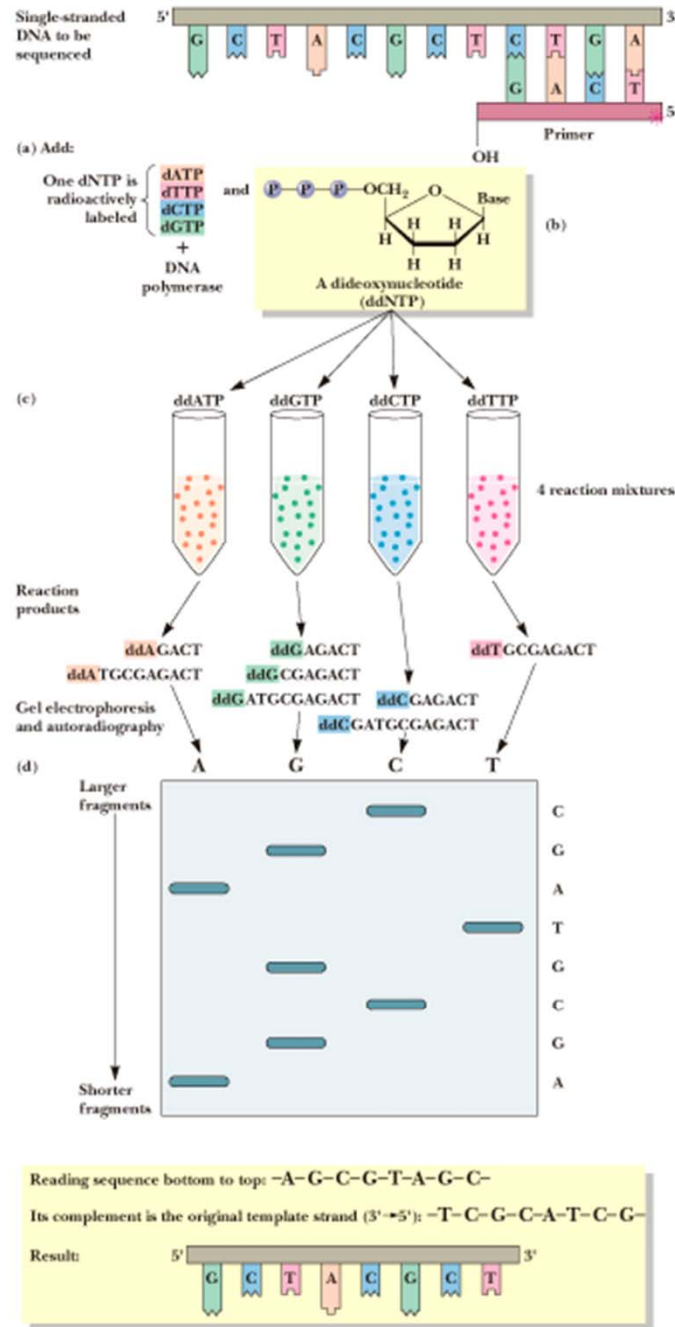
**Metode 2. generacije**

**Metoda 1. generacije**

Comparison of next-generation sequencing methods [37][38]

| Method | Single-molecule real-time sequencing (Pacific Bio) | Ion semiconductor (Ion Torrent sequencing) | Pyrosequencing (454) | Sequencing by synthesis (Illumina) | Sequencing by ligation (SOLiD sequencing) | Chain termination (Sanger sequencing) |
|---|---|---|---|---|---|---|
| Read length | 5,500 bp to 8,500 bp avg (10,000 bp N50); maximum read length >30,000 bases[39][40][41] | up to 400 bp | 700 bp | 50 to 300 bp | 50+35 or 50+50 bp | 400 to 900 bp |
| Accuracy | 99.999% consensus accuracy; 87% single-read accuracy[42] | 98% | 99.9% | 98% | 99.9% | 99.9% |
| Reads per run | 50,000 per SMRT cell, or ~400 megabases[43][44] | up to 80 million | 1 million | up to 3 billion | 1.2 to 1.4 billion | N/A |
| Time per run | 30 minutes to 2 hours[45] | 2 hours | 24 hours | 1 to 10 days, depending upon sequencer and specified read length[46] | 1 to 2 weeks | 20 minutes to 3 hours |
| Cost per 1 million bases (in US$) | $0.33-$1.00 | $1 | $10 | $0.05 to $0.15 | $0.13 | $2400 |
| Advantages | Longest read length. Fast. Detects 4mC, 5mC, 6mA.[47] | Less expensive equipment. Fast. | Long read size. Fast. | Potential for high sequence yield, depending upon sequencer model and desired application. | Low cost per base. | Long individual reads. Useful for many applications. |
| Disadvantages | Moderate throughput. Equipment can be very expensive. | Homopolymer errors. | Runs are expensive. Homopolymer errors. | Equipment can be very expensive. Requires high concentrations of DNA. | Slower than other methods. Have issue sequencing palindromic sequence.[48] | More expensive and impractical for larger sequencing projects. |

**Metode 3. generacije: uporaba nanopor**

Wikipedia: DNA sequencing

# „Sangerjeva reakcija" 1. generacija določanja zaporedij („chain terminating")

**The Nobel Prize in Chemistry 1980**

# Walter Gilbert

# Frederick Sanger

"for their contributions concerning the determination of base sequences in nucleic acids"

Sanger is responsible for the first complete determination of the sequence of a DNA molecule. He has established the sequence of the 5375 building blocks in DNA from a bacterial virus called phi-X174.

Sequence investigations with the methods of Gilbert and Sanger together with the recombinant-DNA technique make excellent tools for continued investigations of the structure and function of the genetic material.

# 1986

Avtomatizacija določanja nukleotidnih zaporedij in uporaba fluorescentnih barvil

Leroy Hood in Mike Hunkapiller

CELERA
an Applera Corporation Business



300 ABI PRISM® 3700 DNA Analyzers, Applied Biosystems



800 povezanih Compaq Alpha-based 64-bit postaj, vsaka sposobna več kot 250 bilijonov primerjav zaporedij na uro.



GAP          SINGLE STRANDED

# „Pyrosequencing" 2. generacija določanja zaporedij („by synthesis")



Gharizedah B et al. (2001) Lab. Investig. 81:673–679

Margulies M et al. (2005) Nature 437, 376-380

**Tehnologija vzporednega določanja zaporedij**

- Masovno določanje zaporedij
- Hitro
- Kloniranje ni potrebno
- „Emulsion PCR"
- Cenejše
- Bolj zanesljivo določanje zaporedij

Npr. sistem 454 Roche

Steen JA and Cooper MA (2011) Nature Methods 8,548–549

# „Nanopore sequencing" 3. generacija določanja zaporedij



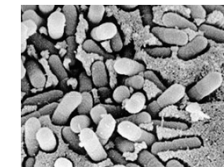Branton D *et al.* (2008) *Nat. Biotechnol.* 26: 1146 – 1153.

## Tehnologija določanja zaporedij s pomočjo proteinskih nanopor

- Masovno določanje zaporedij
- Hitro
- Kloniranje in PCR ni potrebno
- Določitev veliko daljših zaporedij, npr. kbp

Npr. Oxford Nanopore Technologies

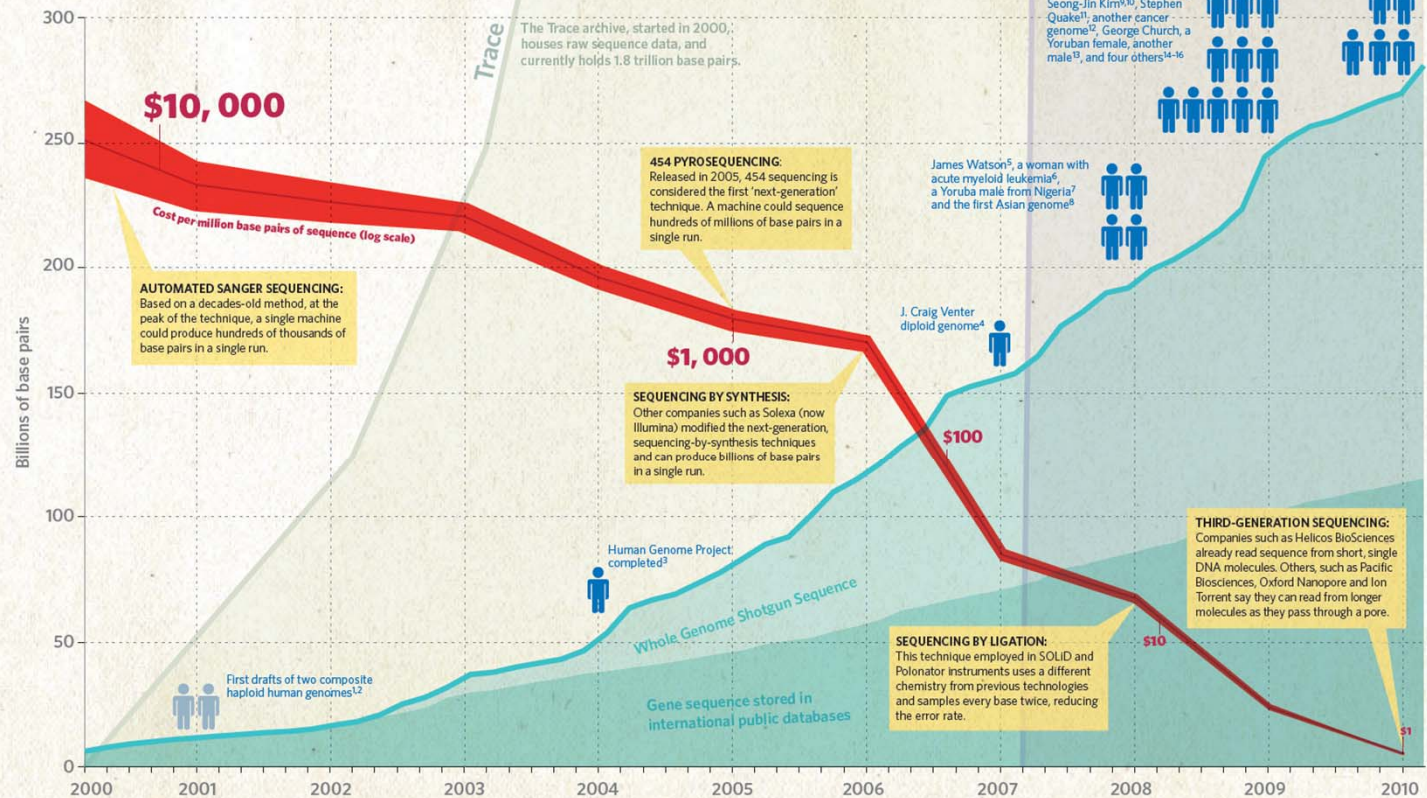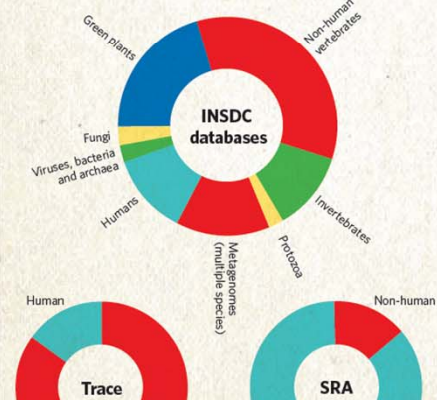| | Število baznih parov $\times 10^6$ | Število genov | Določen |
|---|---|---|---|
| Bakteriofag $\phi$X174 | 0.005 | 10 | 1977 |
| *Mycoplasma genitalium* | 0.58 | 483 | 1995 |
| *Hemophilus influenzae* | 1.83 | 1738 | 1995 |
| *M. tuberculosis* | 4.41 | 3959 | 1998 |
| *Escherichia coli* | 4.6 | 4377 | 1997 |
| *Saccharomyces cerevisiae* | 12.00 | 5885 | 1996 |
| *Caenorhabditis elegans* | 95.50 | 19.820 | 1998 |
| *Drosophila melanogaster* | 180.00 | 13.601 | 2000 |
| *Arabidopsis thaliana* | 117.00 | 25.498 | 2000 |
| Človek | 3300.00 | ≈ 34.000 | 2001 |

HUMAN GENOME

# THE SEQUENCE EXPLOSION

A t the time of the announcement of the first drafts of the human genome in 2000, there were 8 billion base pairs of sequence in the three main databases for 'finished' sequence: GenBank, run by the US National Center for Biotechnology Information; the DNA Databank of Japan; and the European Molecular Biology Laboratory (EMBL) Nucleotide Sequence Database. The databases share their data regularly as part of the International Nucleotide Sequence Database Collaboration (INSDC). In the subsequent first post-genome decade, they have added another 270 billion bases to the collection of finished sequence, doubling the size of the database roughly every 18 months. But this number is dwarfed by the amount of raw sequence that has been created and stored by researchers around the world in the Trace archive and Sequence Read Archive (SRA).
See Editorial, page 649, and human genome special at www.nature.com/humangenome

## DNA SEQUENCES BY TAXONOMY

**International Nucleotide Sequence Database Collaboration:**
The main repositories of 'finished' sequence span a wide range of organisms, representing the many priorities of scientists worldwide.

Green plants
Non-human vertebrates
INSDC databases
Fungi
Viruses, bacteria and archaea
Humans
Invertebrates
Protozoa
Metagenomes (multiple species)

Human
Non-human

**Trace**

**SRA**

The Sequence Read Archive (SRA) houses raw data from next-generation sequencing and has grown to 25 trillion base pairs. If this chart were to accommodate it, it would stretch to more than 12 metres — twice the height of an average giraffe.

Trace

The Trace archive, started in 2000, houses raw sequence data, and currently holds 1.8 trillion base pairs.

SRA

A glioma cell line[17], Inuk[18], !Gubi and Archbishop Desmond Tutu[19], James Lupski[20], and a family of four[21]

Two Korean males including Seong-Jin Kim[9,10], Stephen Quake[11]; another cancer genome[12], George Church, a Yoruban female, another male[13], and four others[14–16]

$10,000

Cost per million base pairs of sequence (log scale)

**454 PYROSEQUENCING:**
Released in 2005, 454 sequencing is considered the first 'next-generation' technique. A machine could sequence hundreds of millions of base pairs in a single run.

James Watson[5], a woman with acute myeloid leukemia[6], a Yoruba male from Nigeria[7] and the first Asian genome[8]

**AUTOMATED SANGER SEQUENCING:**
Based on a decades-old method, at the peak of the technique, a single machine could produce hundreds of thousands of base pairs in a single run.

J. Craig Venter diploid genome[4]

$1,000

**SEQUENCING BY SYNTHESIS:**
Other companies such as Solexa (now Illumina) modified the next-generation, sequencing-by-synthesis techniques and can produce billions of base pairs in a single run.

$100

**THIRD-GENERATION SEQUENCING:**
Companies such as Helicos BioSciences already read sequence from short, single DNA molecules. Others, such as Pacific Biosciences, Oxford Nanopore and Ion Torrent say they can read from longer molecules as they pass through a pore.

Human Genome Project completed[3]

Whole Genome Shotgun Sequence

**SEQUENCING BY LIGATION:**
This technique employed in SOLiD and Polonator instruments uses a different chemistry from previous technologies and samples every base twice, reducing the error rate.

$10

First drafts of two composite haploid human genomes[1,2]

Gene sequence stored in international public databases

$1

Billions of base pairs

300
250
200
150
100
50

2000  2001  2002  2003  2004  2005  2006  2007  2008  2009  2010

Nature, 1. 4. 2010

# NUKLEOTIDNE PODATKOVNE ZBIRKE

Dnevna izmenjava zaporedij!

**DDBJ**    *DNA Data Bank of Japan*
Zaporedja javno dostopnih virov in večjih genomskih projektov

**EMBL**

Podatkovna zbirka EBI. Vsebuje direktno vpisana zaporedja, rezultate določevanj zaporedij genomov, zaporedja iz literature in patentov. Iskanje in primerjave zaporedij preko vmesnikov

**GenBank**

Zaporedja javno dostopnih virov in večjih genomskih projektov.
Dnevna izmenjava podatkov z DDBJ in EMBL.
Razdeljena na posamezne odseke (nivo organizmov, EST, PAT,  STS, GSS, HTG).
**157.943.793.171** baz v **171.123.749** zaporedjih (14. 2. 2014)
iz več kot **165 000** organizmov
Dostop preko Entreza (Medline ali BLAST)
**dbEST**    podatkovna zbirka oznak izraženih zaporedij

**Ostale specializirane podatkovne zbirke**
GSDB (The Genome Sequence DataBase), SGD (*Saccharomyces* Genome Database), UniGene, TDB (TIGR podatkovna zbirka), ACeDB (A *C. elegans* DataBase)

**Table 1.**
**Growth of GenBank divisions (nucleotide base pairs)**

| Division | Description | Release 191 (8/2012) | Annual increase (%) [a] |
|---|---|---|---|
| **Taxnomic divisons** | | | |
| SYN | Synthetic | 928 200 038 | 494.2% |
| PHG | Phages | 84 079 451 | 34.4% |
| ENV | Environmental samples | 3 374 433 548 | 32.1% |
| VRL | Viruses | 1 429 464 786 | 21.1% |
| BCT | Bacteria | 8 439 854 434 | 21.0% |
| PLN | Plants | 5 481 470 133 | 15.6% |
| MAM | Other mammals | 863 036 872 | 6.9% |
| VRT | Other vertebrates | 2 886 594 595 | 6.7% |
| PRI | Primates | 6 317 656 773 | 3.3% |
| UNA | Unannotated | 127 803 | 1.5% |
| ROD | Rodents | 4 435 106 948 | 0.9% |
| INV | Invertebrates | 2 493 058 927 | −1.7% |
| **Functional divisions** | | | |
| TSA | Transcriptome shotgun data | 5 759 588 580 | 207.3% |
| WGS | Whole-genome shotgun data | 308 196 411 905 | 47.9% |
| PAT | Patented sequences | 12 118 622 726 | 8.6% |
| GSS | Genome survey sequences | 21 947 780 105 | 5.7% |
| EST | Expressed sequence tags | 40 888 051 100 | 4.8% |
| HTG | High-throughput genomic | 24 359 210 558 | 0.1% |
| STS | Sequence tagged sites | 636 262 446 | 0.1% |
| HTC | High-throughput cDNA | 639 165 410 | −3.5% |
| TOTAL | All GenBank sequences | 451 278 177 138 | 33.1% |



Benson DA et al. (2013) Nucl. Acids Res. 41: D36-D42

**Iskanje preko Entrez Nucleotide**
**CoreNucleotide** (the main collection)- osnovna zbirka
**dbEST** (Expressed Sequence Tags)- zbirka izraženih oznak zaporedij
**dbGSS** (Genome Survey Sequences)- neanotirana zaporedja („single-read")

**EST kloni**

cDNA knjižnica kot osnova

•Iskanje novih genov
•Pomoč pri določanju genov v genomih
•Kvantifikacija izražanja genov
•Primerjava med celicami/tkivi

Z modernimi metodami določanja zaporedij ni več aktualno-
**RNASeq**



mRNA — AAAAAA TTTTTT *Not* I — *Not* I primer-adapter

First strand synthesis — AAAAAA TTTTTT *Not* I

Second strand synthesis — AAAAAA TTTTTT *Not* I

*Sal* I adapter addition — *Sal* I — AAAAAA TTTTTT *Not* I *Sal* I

*Not* I digestion / Size fractionation — *Sal* I — AAAAAA TTTTTT *Not* I

Ligation to *Not* I-*Sal* I-cut vector — AAAAAA TTTTTT

cDNA ready for transformation

# EST kloni kot pomoč pri določanju genov



**TIGR**

**T**he **I**nstitute for **G**enomic **R**esearch

http://www.tigr.org/

Tudi ostali: Merck/IMAGE, Incyte....

**dbEST**

**dbEST (Nature Genetics 4:332-3;1993) is a division of GenBank that contains sequence data and other information on "single-pass" cDNA sequences, or "Expressed Sequence Tags", from a number of organisms.** A brief account of the history of human ESTs in GenBank is available (Trends Biochem. Sci. 20:295-6;1995). Also, consult the special "Genome Directory" issue of Nature (vol. 377, issue 6547S, 28 September 1995).

Adams MD *et al.* Complementary DNA sequencing: expressed sequence tags and human genome project. Science. 1991 21;252(5013):1651-6.

GenBank    [Nucleotide ▼]    [                    ]    Searc

GenBank ▼ | Submit ▼ | Genomes ▼ | WGS ▼ | HTGs ▼ | EST/GSS ▼ | Metagenomes ▼ | TPA ▼ | TSA ▼ | INSDC ▼

## dbEST release 130101

Summary by Organism - 01 January 2013

**Number of public entries: 74,186,692**

| Organism | Entries |
|---|---|
| Homo sapiens (human) | 8,704,790 |
| Mus musculus + domesticus (mouse) | 4,853,570 |
| Zea mays (maize) | 2,019,137 |
| Sus scrofa (pig) | 1,669,337 |
| Bos taurus (cattle) | 1,559,495 |
| Arabidopsis thaliana (thale cress) | 1,529,700 |
| Danio rerio (zebrafish) | 1,488,275 |
| Glycine max (soybean) | 1,461,722 |
| Triticum aestivum (wheat) | 1,286,372 |
| Xenopus (Silurana) tropicalis (western clawed frog) | 1,271,480 |
| Oryza sativa (rice) | 1,253,557 |
| Ciona intestinalis | 1,205,674 |
| Rattus norvegicus + sp. (rat) | 1,162,136 |
| Drosophila melanogaster (fruit fly) | 821,005 |
| Panicum virgatum (switchgrass) | 720,590 |
| Xenopus laevis (African clawed frog) | 677,911 |
| Oryzias latipes (Japanese medaka) | 666,891 |
| Brassica napus (oilseed rape) | 643,881 |
| Gallus gallus (chicken) | 600,434 |

EST        | EST ▼ |  human pancreas                                          ⊗   | Search |

Save search  Limits  Advanced                                                          Help

Display Settings: ⊡ Summary, 20 per page, Sorted by Default order          Send to: ⊡

ⓘ Found 2957530 nucleotide sequences.  Nucleotide (11127)  EST (157261)  GSS (2789142)

**Filter your results:**

All (157261)

Bacteria (0)

mRNA (157261)

Manage Filters

**Results: 1 to 20 of 157261**                    << First  < Prev   Page 1 of 7864   Next >  Last >>

▼ Top Organisms [Tree]
Homo sapiens (154723)
Mus musculus (2537)
Necator americanus (1)

☐ ig26h09.y5 **Human** Fetal **Pancreas** 1B **Homo sapiens** cDNA clone IMAGE: 5-, mRNA sequence
1. 594 bp linear mRNA
Accession: CK824495.1  GI: 44841420
EST    GenBank    FASTA

☐ ig26h09.x5 **Human** Fetal **Pancreas** 1B **Homo sapiens** cDNA clone IMAGE: 3-, mRNA sequence
2. 579 bp linear mRNA
Accession: CK824494.1  GI: 44841419
EST    GenBank    FASTA

**Find related data**
Database: Select ▼
Find items

☐ ig26h08.y5 **Human** Fetal **Pancreas** 1B **Homo sapiens** cDNA clone IMAGE: 5-, mRNA sequence
3. 582 bp linear mRNA
Accession: CK824493.1  GI: 44841418
EST    GenBank    FASTA

☐ ig26h08.x5 **Human** Fetal **Pancreas** 1B **Homo sapiens** cDNA clone IMAGE: 3-, mRNA sequence
4. 570 bp linear mRNA
Accession: CK824492.1  GI: 44841417
EST    GenBank    FASTA

**Search details**

```
("Homo sapiens"[Organism] OR
human[All Fields]) AND
pancreas[All Fields]
```

☐ ig26h07.y5 **Human** Fetal **Pancreas** 1B **Homo sapiens** cDNA clone IMAGE: 5- similar to TR:Q9ULJ9 Q9ULJ9 KIAA1221
5. PROTEIN ;, mRNA sequence
589 bp linear mRNA
Accession: CK824491.1  GI: 44841416
EST    GenBank    FASTA

Search                                See more...

☐ ig26h07.x5 **Human** Fetal **Pancreas** 1B **Homo sapiens** cDNA clone IMAGE: 3- similar to TR:Q9ULJ9 Q9ULJ9 KIAA1221
6. PROTEIN ;, mRNA sequence
599 bp linear mRNA
Accession: CK824490.1  GI: 44841415
EST    GenBank    FASTA

**Recent activity**
                                        Turn Off   Clear

🔍 human pancreas (157261)                          EST

☐ ig26h02.y5 **Human** Fetal **Pancreas** 1B **Homo sapiens** cDNA clone IMAGE: 5- similar to TR:O94941 O94941 KIAA0860 PROTEIN.
7. ;, mRNA sequence
593 bp linear mRNA
Accession: CK824489.1  GI: 44841414
EST    GenBank    FASTA

📄 hx26h11.y1 Human primary human ocular
pericytes. Equalized (hx) Homo sapiens  EST

🔍 human (10463910)                                EST

📄 Landscape of transcription in human cells.
                                          PubMed

📄 An integrated encyclopedia of DNA elements
in the human genome.                      PubMed

☐ ig26h02.x5 **Human** Fetal **Pancreas** 1B **Homo sapiens** cDNA clone IMAGE: 3- similar to TR:O94941 O94941 KIAA0860 PROTEIN.

See more...

Display Settings: ☑ EST                                                    Send to: ☑

# ig26h09.y5 Human Fetal Pancreas 1B Homo sapiens cDNA clone IMAGE: 5-, mRNA sequence

GenBank: CK824495.1

GenBank    FASTA

**IDENTIFIERS**

```
dbEST Id:       21765372
EST name:       ig26h09.y5
GenBank Acc:    CK824495
GenBank gi:     44841420
```

**CLONE INFO**
```
Clone Id:       IMAGE: (5')
Source:         Harvard University (HHMI) & Washington University (GSC)
DNA type:       cDNA
```

**PRIMERS**
```
Sequencing:     -40UP from Gibco
PolyA Tail:     Unknown
```

**SEQUENCE**
```
TTTTTTTTTTTCTTAATCCAGTCTATCAATAGACTTTATTTTAAAGCAGTTGTTGGTTCA
CAGAAAAACTGAGCAGAAGGTGCAGAGATATCCCAGAGGACCCCAACCCCCTGGGCCATG
GACCAGTACCAGCTGGTGGCCTGTTAGAAATTGGGCCACACAGTGGGCGAGTTAGCATTA
CCAACTGAGCCCCGCCTTCTGTCAGATCAGTGGTGGCATTAGATTCTCATAGGAGTGCGA
ATCCTATTGTTAACTGTGCATGCAAGAGATCTAGGTTGCATGTTCCTTATGAGAACCTAG
CTAATGCTTTTTGATCTGAGGTGGAACAGTTTCACAGTTTCATCCCCAAATCATATCCAC
ATTCCAACCACTATCTGTGGAAAAATTGTCTTCCATAAAACTGGTCCCTGGTGCCAAAAA
GGTTGGGGACTGCTGCCAGATACACTTGGCTCCCACACATGTGTAGCCTCCCCCATTAGC
AACATCTCCCACCAGAGTGGTACATTTGCTATAATTTATGTACTTACATTGGCACATCAT
TATCACCCAAAGTCCATAGTTTACATTAAGTCTCACTCTTGGTGTTGTACATTC
```

```
Entry Created:  Mar 1 2004
Last Updated:   Mar 11 2004
```

**COMMENTS**
```
                This read is a 5' RESEQUENCE of a previously sequenced
                pancreas clone
                Good hit to opposite strand read...wrong orientation BUT
                PASSED FOR MOUSE-PANCREAS VERIFICATION
```

**LIBRARY**
```
Lib Name:       LIBEST_009966 Human Fetal Pancreas 1B
Organism:       Homo sapiens
Tissue type:    Fetal Pancreas (4 Pooled Donors, 18 - 20 weeks, Stratagene
                #738023)
Develop. stage: Fetal Pancreas
Vector:         pBluescript SK(-)
R. Site 1:      NotI
```

**Analyze this sequence**
Run BLAST
Pick Primers

**Related information**
BioSample
Taxonomy
Map Viewer

>gi|44841420|gb|CK824495.1|CK824495 ig26h09.y5 Human Fetal Pancreas 1B Homo
sapiens cDNA clone IMAGE: 5', mRNA sequence
TTTTTTTTTTTCTTAATCCAGTCTATCAATAGACTTTATTTTAAAGCAGTTGTTGGTTCACAGAAAAACT
GAGCAGAAGGTGCAGAGATATCCCAGAGGACCCCAACCCCCTGGGCCATGGACCAGTACCAGCTGGTGGC
CTGTTAGAAATTGGGCCACACAGTGGGCGAGTTAGCATTACCAACTGAGCCCCGCCTTCTGTCAGATCAG
TGGTGGCATTAGATTCTCATAGGAGTGCGAATCCTATTGTTAACTGTGCATGCAAGAGATCTAGGTTGCA
TGTTCCTTATGAGAACCTAGCTAATGCTTTTTGATCTGAGGTGGAACAGTTTCACAGTTTCATCCCCAAA
TCATATCCACATTCCAACCACTATCTGTGGAAAAATTGTCTTCCATAAAACTGGTCCCTGGTGCCAAAAA
GGTTGGGGACTGCTGCCAGATACACTTGGCTCCCACACATGTGTAGCCTCCCCCATTAGCAACATCTCCC
ACCAGAGTGGTACATTTGCTATAATTTATGTACTTACATTGGCACATCATTATCACCCAAAGTCCATAGT
TTACATTAAGTCTCACTCTTGGTGTTGTACATTC

UCSC Genome browser
http://genome.ucsc.edu/goldenPath/aug2001Tracks.html

International Human Genome Sequencing Consortium
(2001) *Initial sequencing and analysis of the human genome*.
Nature 409, 860-921

Craig Venter



Francis Collins



- International Human Genome Sequencing Consortium (2001) *Initial sequencing and analysis of the human genome*. Nature 409, 860-921
- J. Craig Venter et al. (2001) *The Sequence of the Human Genome*. Science, 1304-1135

Prve verziji človeškega genoma.

- The Human Genome 10th Anniversary: http://www.sciencemag.org/site/extra/genomeanniversary/

Tematska številka ob 10. obletnici objave človeškega genoma.

- The 1000 Genomes Project Consortium (2012) *An integrated map of genetic variation from 1,092 human genomes*. Nature 491, 56-65

By characterizing the geographic and functional spectrum of human genetic variation, the 1000 Genomes Project aims to build a resource to help to understand the genetic contribution to disease. **Here we describe the genomes of 1,092 individuals from 14 populations**, constructed using a combination of low-coverage whole-genome and exome sequencing. By developing methods to integrate information across several algorithms and diverse data sources, we provide a validated haplotype map of **38 million single nucleotide polymorphisms, 1.4 million short insertions and deletions, and more than 14,000 larger deletions**. We show that individuals from different populations carry different profiles of rare and common variants, and that low-frequency variants show substantial geographic differentiation, which is further increased by the action of purifying selection. We show that evolutionary conservation and coding consequence are key determinants of the strength of purifying selection, that rare-variant load varies substantially across biological pathways, and that each individual contains hundreds of rare non-coding variants at conserved sites, such as motif-disrupting changes in transcription-factor-binding sites. This resource, which captures up to 98 % of accessible single nucleotide polymorphisms at a frequency of 1 % in related populations, enables analysis of common and low-frequency variants in individuals from diverse, including admixed, populations.

- The ENCODE Project Consortium (2012) *An integrated encyclopedia of DNA elements in the human genome*. Nature 489, 57–74

**The human genome encodes the blueprint of life, but the function of the vast majority of its nearly three billion bases is unknown. The Encyclopedia of DNA Elements (ENCODE) project has systematically mapped regions of transcription, transcription factor association, chromatin structure and histone modification.** These data enabled us to assign biochemical functions for 80% of the genome, in particular outside of the well-studied protein-coding regions. Many discovered candidate regulatory elements are physically associated with one another and with expressed genes, providing new insights into the mechanisms of gene regulation. The newly identified elements also show a statistical correspondence to sequence variants linked to human disease, and can thereby guide interpretation of this variation. Overall, the project provides new insights into the organization and regulation of our genes and genome, and is an expansive resource of functional annotations for biomedical research.
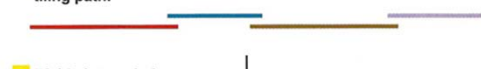
# Map-based genome sequencing



# Shotgun genome sequencing



**Prekrito zaporedje** ("contig")- prebrana zaporedja, ki se prekrivajo brez vrzeli in imajo visoko zanesljivost določitve.
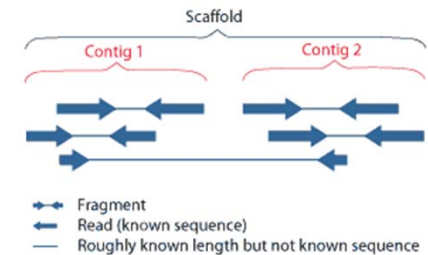
**Ogrodje** ("scaffold")- del genoma rekonstruiran iz odčitkov. Vsebuje prekrita zaporedja in vrzeli. Navadno se uporablja parne odčitke ("pair-end reads")
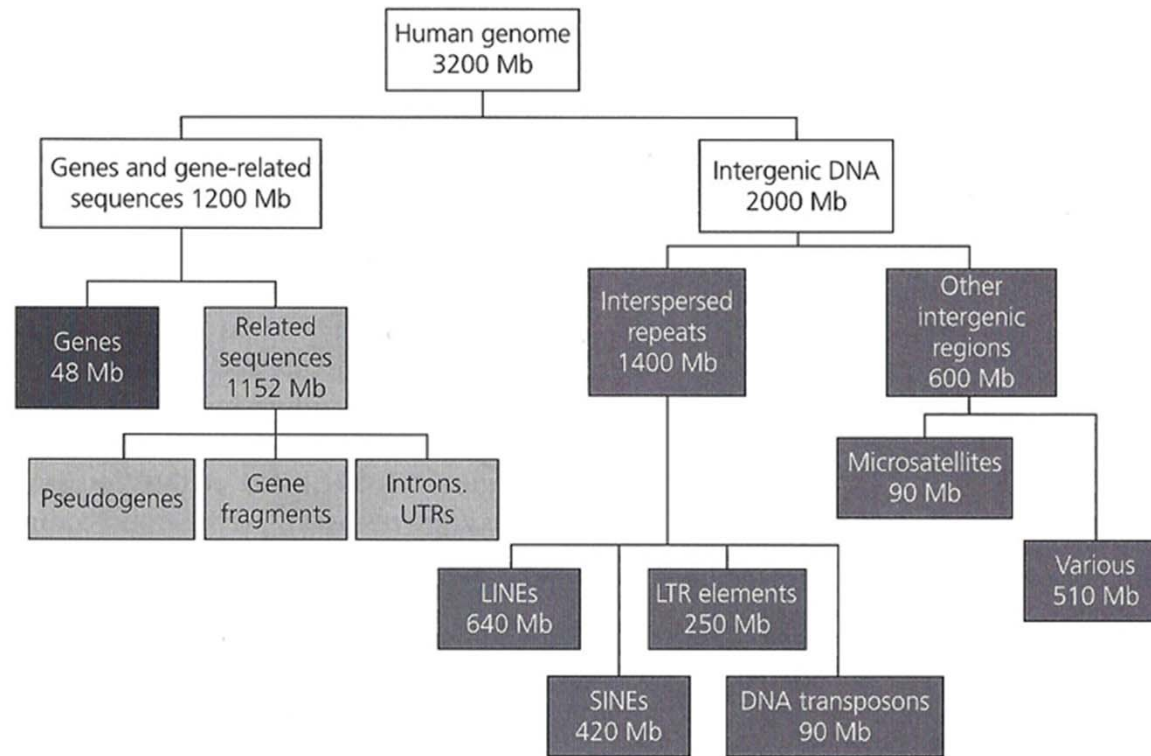
**Odčitek** ("read")- eno prebrano zaporedje

Vir: http://www.discoveryandinnovation.com/BIOL202

| Chromosome | Length (mm) | Base pairs | Variations | Confirmed proteins | Putative proteins | Pseudogenes | miRNA | rRNA | snRNA | snoRNA | Misc ncRNA | Links | Centromere position (Mbp) | Cumulative (%) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 85 | 249,250,621 | 4,401,091 | 2,012 | 31 | 1,130 | 134 | 66 | 221 | 145 | 106 | EBI | 125.0 | 7.9 |
| 2 | 83 | 243,199,373 | 4,607,702 | 1,203 | 50 | 948 | 115 | 40 | 161 | 117 | 93 | EBI | 93.3 | 16.2 |
| 3 | 67 | 198,022,430 | 3,894,345 | 1,040 | 25 | 719 | 99 | 29 | 138 | 87 | 77 | EBI | 91.0 | 23.0 |
| 4 | 65 | 191,154,276 | 3,673,892 | 718 | 39 | 698 | 92 | 24 | 120 | 56 | 71 | EBI | 50.4 | 29.6 |
| 5 | 62 | 180,915,260 | 3,436,667 | 849 | 24 | 676 | 83 | 25 | 106 | 61 | 68 | EBI | 48.4 | 35.8 |
| 6 | 58 | 171,115,067 | 3,360,890 | 1,002 | 39 | 731 | 81 | 26 | 111 | 73 | 67 | EBI | 61.0 | 41.6 |
| 7 | 54 | 159,138,663 | 3,045,992 | 866 | 34 | 803 | 90 | 24 | 90 | 76 | 70 | EBI | 59.9 | 47.1 |
| 8 | 50 | 146,364,022 | 2,890,692 | 659 | 39 | 568 | 80 | 28 | 86 | 52 | 42 | EBI | 45.6 | 52.0 |
| 9 | 48 | 141,213,431 | 2,581,827 | 785 | 15 | 714 | 69 | 19 | 66 | 51 | 55 | EBI | 49.0 | 56.3 |
| 10 | 46 | 135,534,747 | 2,609,802 | 745 | 18 | 500 | 64 | 32 | 87 | 56 | 56 | EBI | 40.2 | 60.9 |
| 11 | 46 | 135,006,516 | 2,607,254 | 1,258 | 48 | 775 | 63 | 24 | 74 | 76 | 53 | EBI | 53.7 | 65.4 |
| 12 | 45 | 133,851,895 | 2,482,194 | 1,003 | 47 | 582 | 72 | 27 | 106 | 62 | 69 | EBI | 35.8 | 70.0 |
| 13 | 39 | 115,169,878 | 1,814,242 | 318 | 8 | 323 | 42 | 16 | 45 | 34 | 36 | EBI | 17.9 | 73.4 |
| 14 | 36 | 107,349,540 | 1,712,799 | 601 | 50 | 472 | 92 | 10 | 65 | 97 | 46 | EBI | 17.6 | 76.4 |
| 15 | 35 | 102,531,392 | 1,577,346 | 562 | 43 | 473 | 78 | 13 | 63 | 136 | 39 | EBI | 19.0 | 79.3 |
| 16 | 31 | 90,354,753 | 1,747,136 | 805 | 65 | 429 | 52 | 32 | 53 | 58 | 34 | EBI | 36.6 | 82.0 |
| 17 | 28 | 81,195,210 | 1,491,841 | 1,158 | 44 | 300 | 61 | 15 | 80 | 71 | 46 | EBI | 24.0 | 84.8 |
| 18 | 27 | 78,077,248 | 1,448,602 | 268 | 20 | 59 | 32 | 13 | 51 | 36 | 25 | EBI | 17.2 | 87.4 |
| 19 | 20 | 59,128,983 | 1,171,356 | 1,399 | 26 | 181 | 110 | 13 | 29 | 31 | 15 | EBI | 26.5 | 89.3 |
| 20 | 21 | 63,025,520 | 1,206,753 | 533 | 13 | 213 | 57 | 15 | 46 | 37 | 34 | EBI | 27.5 | 91.4 |
| 21 | 16 | 48,129,895 | 787,784 | 225 | 8 | 150 | 16 | 5 | 21 | 19 | 8 | EBI | 13.2 | 92.6 |
| 22 | 17 | 51,304,566 | 745,778 | 431 | 21 | 308 | 31 | 5 | 23 | 23 | 23 | EBI | 14.7 | 93.8 |
| X | 53 | 155,270,560 | 2,174,952 | 815 | 23 | 780 | 128 | 22 | 85 | 64 | 52 | EBI | 60.6 | 99.1 |
| Y | 20 | 59,373,566 | 286,812 | 45 | 8 | 327 | 15 | 7 | 17 | 3 | 2 | EBI | 12.5 | 100.0 |
| mtDNA | 0.0054 | 16,569 | 929 | 13 | 0 | 0 | 0 | 2 | 0 | 0 | 22 | EBI | N/A | 100.0 |

Vir: Wikipedia

# Overview of human genome



**Ponovitve zaporedij v genomu:**

- transpozoni
- enostavne ponovitve, npr. 3-100, zaporedja nukleotidov (mikrosateliti). 2-6 nukleotidov, npr. $(CA)_n$
- večje ponovitve (10-300 kb)

# 20.687 genov, ki kodirajo proteine

Funkcijska delitev človeških genov. Največje družine: vezavni proteini za nukleinske kisline transkripcijski faktorji, encimi, transporterji, receptorji...



isomerases; 94; 0,5%
receptors; 1076; 6,3%
storage proteins; 15; 0,1%
structural proteins; 280; 1,6%
surfactants; 15; 0,1%
cell junction proteins; 67; 0,4%
chaperones; 130; 0,8%
transcription factors; 2067; 12,0%
phosphatases; 230; 1,3%
membrane traffic proteins; 321; 1,9%
transfer/carrier proteins; 248; 1,4%
hydrolases; 454; 2,6%
defense/immunity proteins; 107; 0,6%
calcium-binding proteins; 63; 0,4%
viral proteins; 7; 0,0%

unclassified; 4061; 23,6%

extracellular matrix proteins; 72; 0,4%
proteases; 476; 2,8%
cytoskeletal proteins; 441; 2,6%
transporters; 1098; 6,4%
transmembrane receptor regulatory/adaptor proteins; 84; 0,5%
transferases; 1512; 8,8%
oxidoreductases; 550; 3,2%
lyases; 104; 0,6%
cell adhesion molecules; 93; 0,5%
ligases; 260; 1,5%
nucleic acid binding; 1466; 8,5%
signaling molecules; 961; 5,6%
enzyme modulators; 857; 5,0%

Vir: Wikipedia

## dbSNP

Database of single nucleotide polymorphisms (SNPs) and multiple small-scale variations that include insertions/deletions, microsatellites, and non-polymorphic variants.

### Getting Started

Overview of dbSNP

FAQ

Factsheet

### Submit Data

Clinically Associated Human Variations

All Other Variations

Hold Until Published (HUP) Policies

### Access Data

Important RefSNP (RS) Attributes

Web Search

Batch Query

The NCBI Short Genetic Variations (SNV) database, also known as dbSNP, catalogs short variations in
nucleotide sequences from a wide range of organisms. **These variations include single nucleotide variations, short nucleotide insertions and deletions, short tandem repeats and microsatellites.**
Some of these rare human entries
have additional information associated with them, including disease associations, genotype information
and allele origin, as some variations are somatic rather than germline events.

The 1000 Genomes Project Consortium (2012) *An integrated map of genetic variation from 1,092 human genomes*. Nature 491, 56-65

...By developing methods to integrate information across several algorithms and diverse data sources, we provide a validated haplotype map of **38 million single nucleotide polymorphisms, 1.4 million short insertions and deletions, and more than 14,000 larger deletions**. ...

# Genomski brskalnik- Ensembl
http://www.ensembl.org/Homo_sapiens