



Bioinformātika

Nils Rostoks

Latvijas Universitāte

Bioloģijas fakultāte

Lekcijas plāns

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101111001111001111010111100111101011111111001101100001

- Bioinformātika
- Bioloģiskā informācija - tās daudzveidība un apjoms
- Bioloģiskās informācijas datubāzes
- Genomu organizācija un evolūcija
- Salīdzinošā genomika

Mācību materiāli I

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011110111011110011110011110101111001111010111111111001101100001

- Lekcijas
- Lesk AM (2008) Introduction to Bioinformatics. 3rd ed. Oxford University Press, New York, USA
- Lesk AM (2005) Introduction to Bioinformatics. 2nd ed. Oxford University Press, New York, USA
- Pevzner P, Shamir R (2011) Bioinformatics for biologists. Cambridge University Press, Cambridge, UK
- Higgs PG, Attwood TK (2006) Bioinformatics and molecular evolution. Blackwell Publishing, Malden, USA, Oxford, UK
- Claverie, Jean-Michel Bioinformatics for dummies 2003 LUB:Biologijas-zin.-bibl., LUB:Centr.bibl.-krājums
- Higgs, Paul G. Bioinformatics and molecular evolution 2005 LUB:Biologijas-zin.-bibl., LUB:Juridisko-zin.-bibl.
- Mount, David W. Bioinformatics 2001 LUB:Centr.bibl.-krājums

Mācību materiāli II

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101110011110011110011110011110011110011110011100001

- Interneta resursi:

<http://www.ebi.ac.uk/training/online/>

<http://www.ncbi.nlm.nih.gov/education/>

<http://bioinformatics.oxfordjournals.org/>

GAACATTCCATGCTCATGGGTAGGAAGAATCAATATCGTGAAAATGGCCATACTGCCCAAGGTAA
TTTATAGATTCAATGCCATCCCCATCAAGCTACCAATGACTTTCTTCACAGAATTGGAAAAAACT
ACTTTAAAGTTCATATGGAACCAAAAAAGAGCCCGCATTGCCAAGTCAATCCTAAGCCAAAAGAA
CAAAGCTAGAGGCATCACACTACCTGACTTCAAACCTATACTACAAGGCTACAGTAACCAAAACAG
CATGGTACTGGTACCAAAGCAGAGATATAGACCAATGGAACAAAACAGTGCCCTCAGAAATAATA
CTGCATATCTACAACCATCTGATCTTTGACAAACCTGACAAAAACAAGCAATGGGGAAAGGATTC
CCTATTTAATAAATGGTGCTGGGAAAACCTGGCTAGCCATATGTAGAAAGCTGAAATTGGATCCCT
TCCTTACACCTTGTACAAAAATTAATTCAAGATGGATTACAGACTTAAATGTTAGACCTAAAACC
ATAAAAACCCTAGAAGAAAACCTAGGCAATACCATTACAGGACATAGGCATGGGCAAGAACTTCAT
GTCTAGAACACCAAAAGTAATGGCAACAAAAGCCAAAATTGACAAATGGGTCTAATTAAACTAAA
GAGCTTCTGCACAGCAAAGAAACTACCATCAGAGTGAAGAGGCAACCTACAGAATGGGAGAAAA
TTTTTGC AATCTGACAAAAGGGCTAATTTTTTGCATCTGACAAAGGGCTAATATCCAGAATCTACA
ATGAACTCAAACAAATTTACAAGAAAAAAACAAATTTACAAGAAAAAAACAAATTTACAAGAAA
AAAACAAATTTACAAGAAAAAAACAAACAACCCCATCAAAAAGTGGGCAAAGGATAAGAACAGTC
ACTTCTCAAAGAAGACATTTATGCAGCCAAAAGACACATGAAAAAATTCTCATCATCACTGGCC
ATCAGAGAAATGCAAATCAAACCCACAATGAGATACCATCTCACACCAGTTAGAATGGCGATCAT
TAAAAAGTCAAGAAACAACAGGTGCTGGAGAGGATGTGGAGAAATAGGAACACTTTTACACTGTT
AGTGGGACTGTAAACTAGTTCAACATTTGTGGAAGTCAGTATGGCGATTCTCAGGGATCTAGAAC
TAGAAATACCATTTGACCCAGCCATCCCATTACTGGGTATATACCCAAAGGATTATAAATCATGC
TGCTATAAAGACACTTGCACACATATGTTTATTGTGGTACTATTCACAATAGCAAAGACTTGGAA
CCAACCCAAATGTCCAACAATGATAGACTAGATTAAGAAAATGTGGCACATATACACTATGGAAT
ACTTTGCAGCCATAAAAAAGGATGAGTTCATGTCCTTTGTAGGGACATGGATGAAGCTGGAAACC
ATCATTCTCAGCAAACCTATCGCAAGGATAAAAAACCAAACACCGCATGTTCTCACTCATAGGTGG
TAACTGAACAATGAGAACACATGGTCACAGGAAGGGGAACATCACGCACTGGGGCCTGTTGTGGG
GTGGGGGGGAGTGGGGAGGGATAGCATTAGGAGACATACCTAATGTTAAATGACGATTTAATGGGT

Kas ir bioinformātika?

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011111011101101111001111010111111111001101100001

Bioinformatics derives knowledge from computer analysis of biological data. These can consist of the information stored in the genetic code, but also experimental results from various sources, patient statistics, and scientific literature. Research in bioinformatics includes method development for storage, retrieval, and analysis of the data. Bioinformatics is a rapidly developing branch of biology and is highly interdisciplinary, using techniques and concepts from informatics, statistics, mathematics, chemistry, biochemistry, physics, and linguistics. It has many practical applications in different areas of biology and medicine

M. Nilges, J.P. Linge, Institut Pasteur

Kas ir bioinformātika?

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011111011101111001111001111010111100111101011111111001101100001

Bioinformatics: Research, development, or application of computational tools and approaches for expanding the use of biological, medical, behavioral or health data, including those to acquire, store, organize, archive, analyze, or visualize such data.

Computational Biology: The development and application of data-analytical and theoretical methods, mathematical modeling and computational simulation techniques to the study of biological, behavioral, and social systems

*NIH Biomedical Information Science
and Technology Initiative Consortium*

Kas ir bioinformātika?

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011111011101111001111001111010111100111101011111111001101100001

Bioinformatics is the field of science in which biology, computer science, and information technology merge to form a single discipline. The ultimate goal of the field is to enable the discovery of new biological insights as well as to create a global perspective from which unifying principles in biology can be discerned

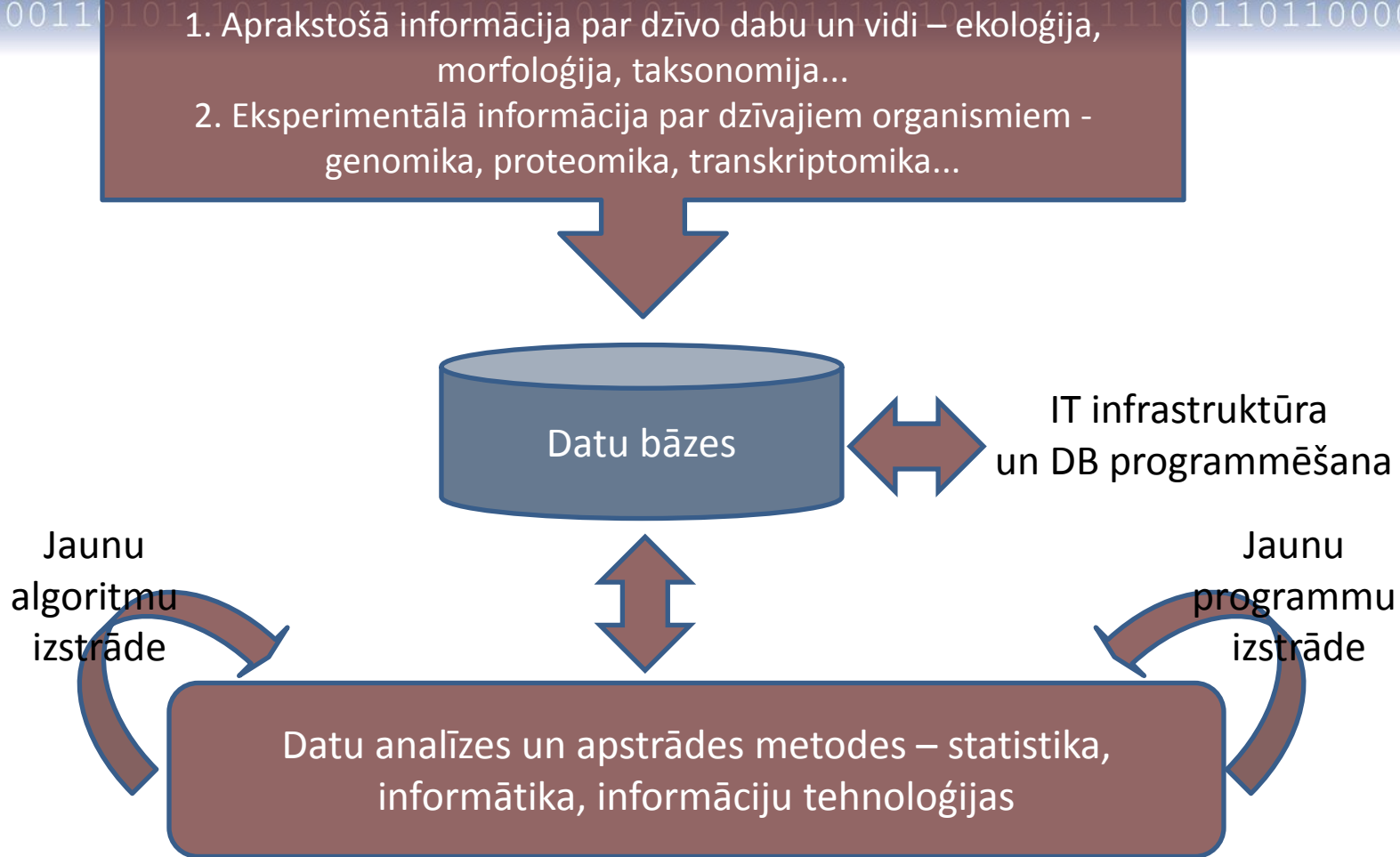
<http://www.ncbi.nlm.nih.gov/About/primer/bioinformatics.html>

Bioinformātikas pētījumu objekts

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111010111100111101011111111001101100001

- Genoma nukleotīdu secības (DNS/RNS) (genomika)
- Genoma ekspresija (dažādu RNS veidu analīze) (transkriptomika)
- Genoma kodētie proteīni (proteomika)
- Šūnas veidotie metabolīti (metabolomika)
- Mijiedarbības dažādu objektu starpā (sistēmbioloģija)

Bioloģija, statistika, programmēšana, informāciju tehnoloģijas





Bioloģiskā informācija - tās daudzveidība un apjoms

GGATCCATTAGAAACGACCAATTTATAGAGGAGAAAAAATGCTCCTCTAAAGGATACTCAACCGTTTACTGTCCCT
AAATGAGACCAACTAGAGTTGCTATGAGAACCTATGATCTGCATTGTGTGTCGCGAACACA...CAAG
GC...AACTATTATTTTATCAATGTATTACTTCTCCGTTCTAAAATATAAGAC...CA
CT...ATATACGGAACAAAATGAATGAATCTATACTTTAAAATATGTCTATGTACA...T
ATA...CTCTAAAAGTCTTATATTTAGGAACGGAGGGAGTATTAGAGAAGACAAC...T
ACTA...AAGAAACATTTTATACAGCACCTTCATCACCAGCACTGACACATACGAAAGCGTCT...GACCAG
CACACACCAGAACATCATTGACTATCTAGTGTTGGCTATGAGGATGAGGTGGGACACTTTCTTCTATTGTAACCG
CGGAAACATCGTGCCCTCTAGACGAGGTTCCGTAAAAAAGTAACAATAGTCCATGAAATTTTAAATTACAGAATA
TATAAAATCACTGGACGCAATCCCGTCACTATGTGTTGCTATCTCCAAATTCAAATTGGGTGCCATGAAAATTA
CGTTTTTCAGAGATGATTTACAAAATCTCATCATATAGAGAATTTTLAGATCGACTGAACACGCTAGAATCATGA
AATGTCATGAGAATGCAAATCATGTGCGAAAACCAGGTCCAAAATGCAATTTATGCACTCACCCAGTATCGTAC
TCCTGCAACCATCGTCAACAGATGTCATTCACCGTTGGGCCGGCCGGGCAGCATCTGCCAGAAGTAGGCAGTAGG
CTTGAATCATCTGGACATA...TAT...CTATCCCCCGTCCATCT
CGAATTTGTTGAATCAAATA...T...GCTCTTATCTACCCO

0000000110101100100101001100010100110001001011000101011100101010010101001
1100000101000100000000001000100000100110111001000000000101100000001000110
1001001010000001110010000010001000100001000000000101000100100111000100100
0001011000101000000101000000001100111011101000001010010100100000101001001
010001100001000100000000101101100100101101101010101000110001110101011011
101010110100100100010100010010100111000101100101101110101000100100001000111
1110001001101111010010110110011100000001000100001011001000000000000101000
00000000101011011100011111010100010100110001011000001000001110111001000000
110000101010010000010000010100100100010100000000100110101001011100100010010
000101001010001100001001011100001101101100000110000000110101011110100011001
011011001100110100101001010001011100111111111111101100101110100100111010011
100100010010110100000100110000010100100010010000111100000001000111111010010
110000010100010000011101001100100100111101110110001100001001110100001001111

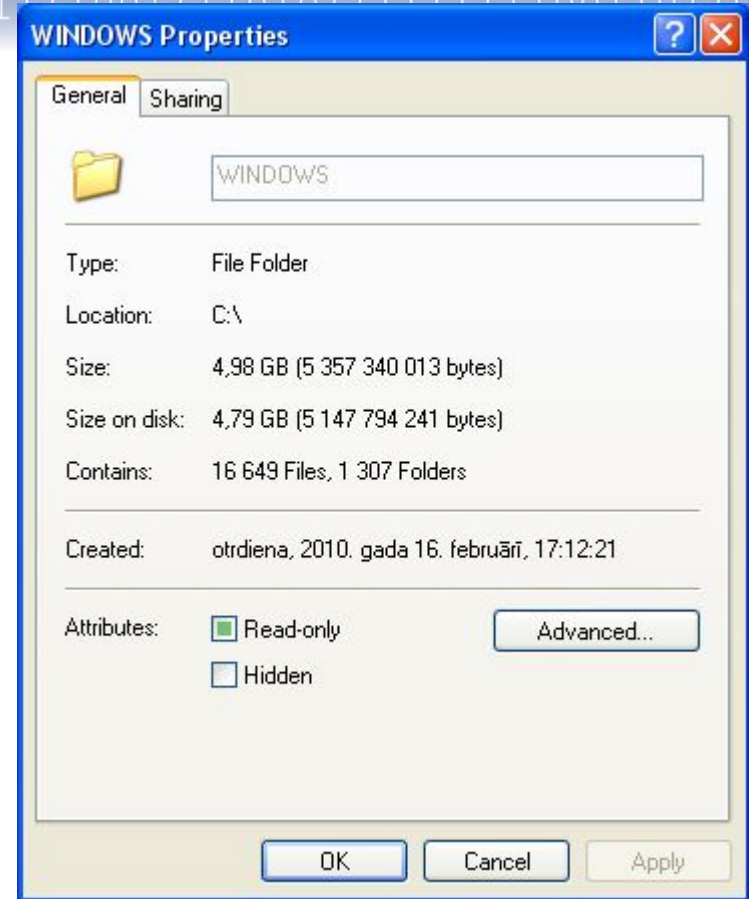
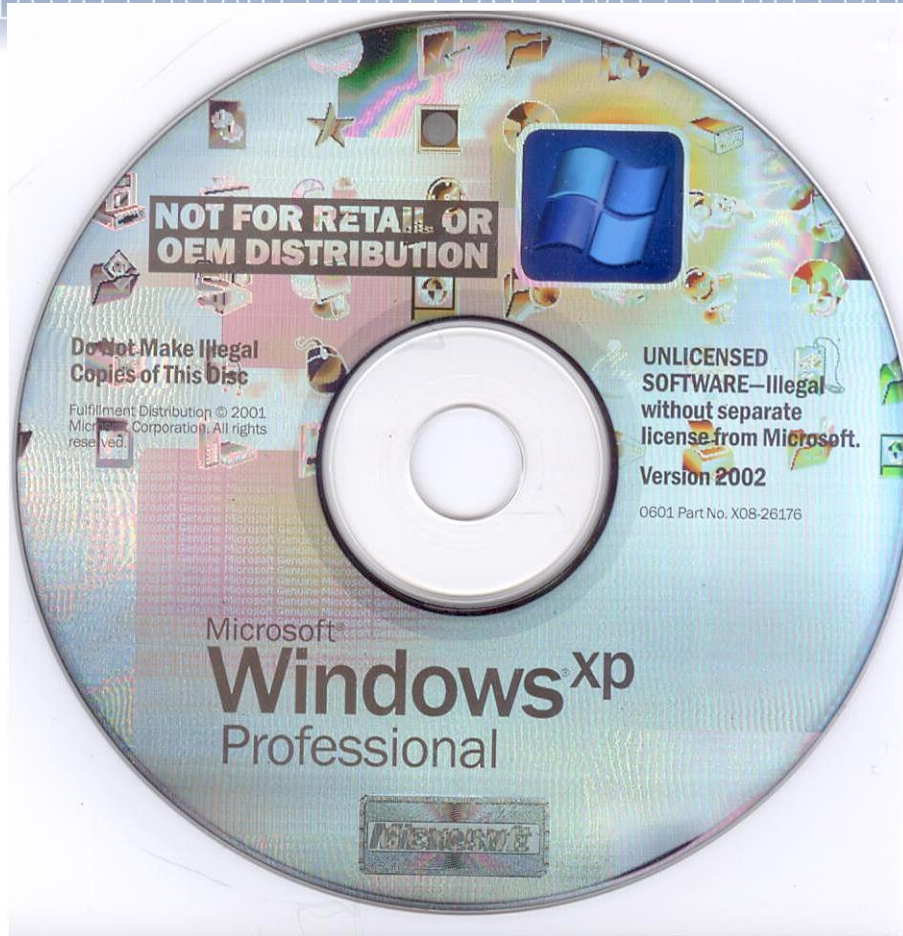
Digitālais cilvēka genoms

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111010111100111101011111111001101100001

- Diploīds genoms – apmēram 6×10^9 nukleotīdu
- Viena nukleotīda kodēšanai binārā formā nepieciešami vismaz 2 biti:
A = 00, C = 01, G = 10, T = 11
- 8 biti = 1 baits
- 1 baits var kodēt 4 nukleotīdus
- Diploīds cilvēka genoms binārā formā ir 1.5×10^9 baiti (datora darbībai vajag vairāk koda nekā dzīvības nodrošināšanai)

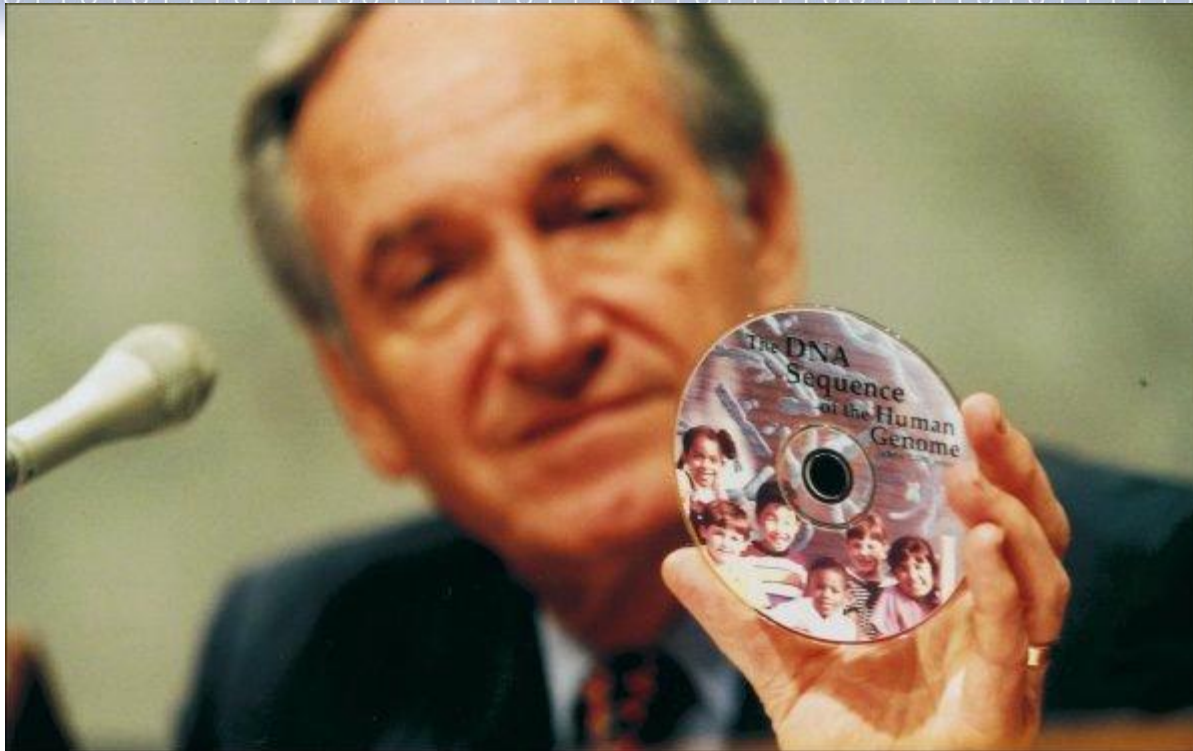
<http://www.tmssoft.com/article-genome.html>

Cilvēka genoms un MS Windows operētājsistēma?



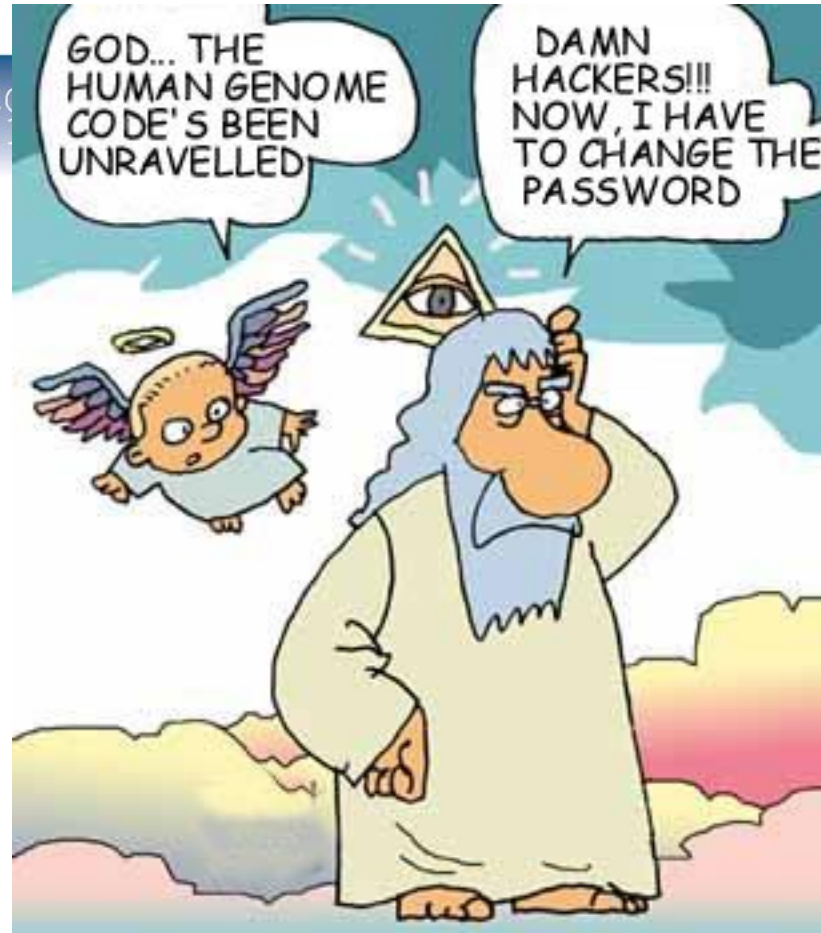
Cilvēka genoma sekvence

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011111011101101111001111010111111111001101100001



<http://www.genome.gov/dmd/img.cfm?node=Photos/People/People%20Outside%20NHGRI&id=79269>

ATGCAGAGCCAGATCGTGTGCC
001101011101001101011



AGGGGCGCCCATGCTGGAATTC
0111111111001101100001

Datorprogrammas DNS sekvences analīzei

Volume 4 Number 11 November 1977

Nucleic Acids Research

Sequence data handling by computer

R.Staden

MRC Laboratory of Molecular Biology, Hills Road, Cambridge CB2 2QH, UK

Received 10 October 1977

ABSTRACT

The speed of the new DNA sequencing techniques has created a need for computer programs to handle the data produced. This paper describes simple programs designed specifically for use by people with little or no computer experience. The programs are for use on small computers and provide facilities for storage, editing and analysis of both DNA and amino acid sequences. A magnetic tape containing these programs is available on request.

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111101110111100111101111001111010111111111001101100001

Genomā ietvertā informācija nosaka organisma identitāti

... vismaz baktērijās ...

Genoma transplantēšana baktērijās

Originally published in *Science Express* on 28 June 2007

Science 3 August 2007:

Vol. 317, no. 5838, pp. 632 - 638

DOI: 10.1126/science.1144622

[< Prev](#) | [Table of Contents](#) | [Next >](#)

CGCCCATGCTGGAATTC
11111001101100001

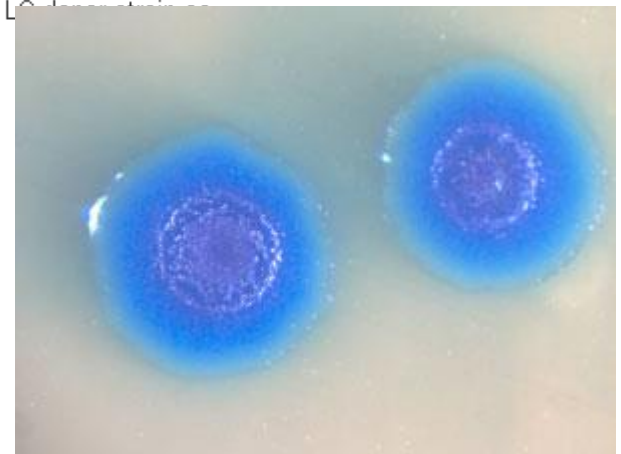
RESEARCH ARTICLES

Genome Transplantation in Bacteria: Changing One Species to Another

Carole Lartigue, John I. Glass,^{*} Nina Alperovich, Rembert Pieper, Prashanth P. Parmar,
Clyde A. Hutchison, III, Hamilton O. Smith, J. Craig Venter

As a step toward propagation of synthetic genomes, we completely replaced the genome of a bacterial cell with one from another species by transplanting a whole genome as naked DNA. Intact genomic DNA from *Mycoplasma mycoides* large colony (LC), virtually free of protein, was transplanted into *Mycoplasma capricolum* cells by polyethylene glycol-mediated transformation. Cells selected for tetracycline resistance, carried by the *M. mycoides* LC chromosome, contain the complete donor genome and are free of detectable recipient genomic sequences. These cells that result from genome transplantation are phenotypically identical to the *M. mycoides* LC donor strain as judged by several criteria.

The J. Craig Venter Institute, Rockville, MD 20850, USA.



Pilna baktērijas genoma ķīmiskā sintēze

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111001111010111100111101011111111001101100001

Originally published in *Science Express* on 24 January 2008
Science 29 February 2008:
Vol. 319, no. 5867, pp. 1215 - 1220
DOI: 10.1126/science.1151721

[< Prev](#) | [Table of Contents](#) | [Next >](#)

RESEARCH ARTICLES

Complete Chemical Synthesis, Assembly, and Cloning of a *Mycoplasma genitalium* Genome

Daniel G. Gibson, Gwynedd A. Benders, Cynthia Andrews-Pfannkoch, Evgeniya A. Denisova, Holly Baden-Tillson, Jayshree Zaveri, Timothy B. Stockwell, Anushka Brownley, David W. Thomas, Mikkel A. Algire, Chuck Merryman, Lei Young, Vladimir N. Noskov, John I. Glass, J. Craig Venter, Clyde A. Hutchison, III, Hamilton O. Smith*

We have synthesized a 582,970–base pair *Mycoplasma genitalium* genome. This synthetic genome, named *M. genitalium* JCVI-1.0, contains all the genes of wild-type *M. genitalium* G37 except MG408, which was disrupted by an antibiotic marker to block pathogenicity and to allow for selection. To identify the genome as synthetic, we inserted "watermarks" at intergenic sites known to tolerate transposon insertions. Overlapping "cassettes" of 5 to 7 kilobases (kb), assembled from chemically synthesized oligonucleotides, were joined by in vitro recombination to produce intermediate assemblies of approximately 24 kb, 72 kb ("1/8 genome"), and 144 kb ("1/4 genome"), which were all cloned as bacterial artificial chromosomes in *Escherichia coli*. Most of these intermediate clones were sequenced, and clones of all four 1/4 genomes with the correct sequence were identified. The complete synthetic genome was assembled by transformation-associated recombination cloning in the yeast *Saccharomyces cerevisiae*, then isolated and sequenced. A clone with the correct sequence was identified. The methods described here will be generally useful for constructing large DNA molecules from chemically synthesized pieces and also from combinations of natural and synthetic DNA segments.

Informācijas uzglabāšana DNS sekvencē

Scienceexpress

Brevia

Next-Generation Digital Information Storage in DNA

George M. Church,^{1,2} Yuan Gao,³ Sriram Kosuri^{1,2*}

¹Department of Genetics, Harvard Medical School, Boston, MA 02115, USA. ²Wyss Institute for Biologically Inspired Engineering, Boston, MA 02115, USA. ³Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD 21205, USA.

*To whom correspondence should be addressed. E-mail: sri.kosuri@wyss.harvard.edu

As digital information continues to accumulate, higher density and longer-term storage solutions are necessary (1). DNA has many potential advantages as a medium for immutable, high latency information storage needs (2). For example, DNA storage is very dense. At theoretical maximum, DNA can encode two bits per nucleotide (nt) or 455 exabytes per gram of ssDNA (3). Unlike most digital storage media, DNA storage is not restricted to a planar layer, and is often readable despite degradation in non-ideal conditions over millennia (4, 5). Finally, DNA's essential biological role provides access to natural reading and writing enzymes and ensures that DNA will remain a readable standard for the foreseeable future.

Storing messages in DNA was first demonstrated in 1988 (6) and the

for all but century-scale analysis and sequencing has become available, information (12). Our combined with library-compatible with future Reciprocally, large-scale could accelerate development (13). Future work could

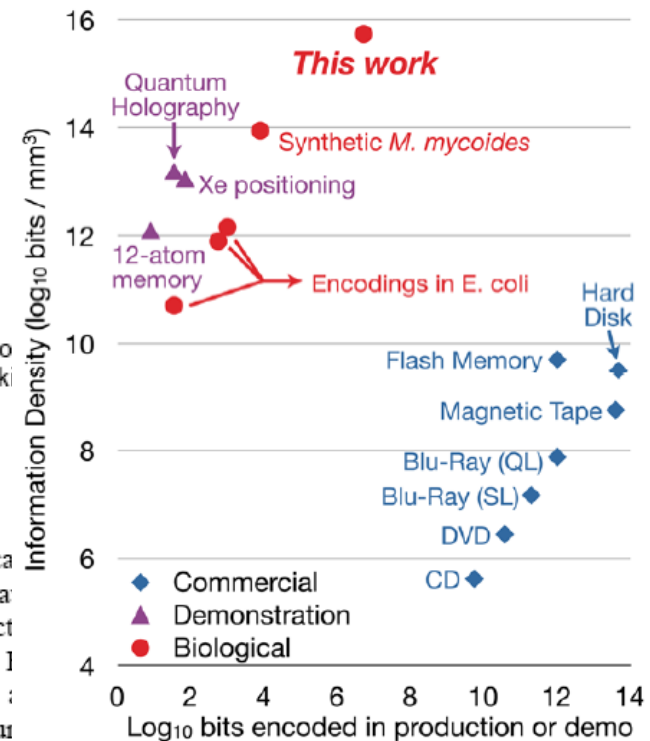


Fig. 1. Comparison to other measured by the log₁₀ of bits encoded in the report or commercial technologies. We plotted information density (log₁₀ of bits/mm²) versus current scalability as unit (3).

DNS sekvenēšana

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111110111011011110011110101111001111010111111111001101100001

- Divas metodes 1975. - 1977. gadā (*Allan Maxam un Walter Gilbert, Frederick Sanger*)

“These chemical procedures ... soon allowed the entire sequence of the plasmid cloning vector pBR322 (4362 bp) to be worked out by a single scientist in only one year.”

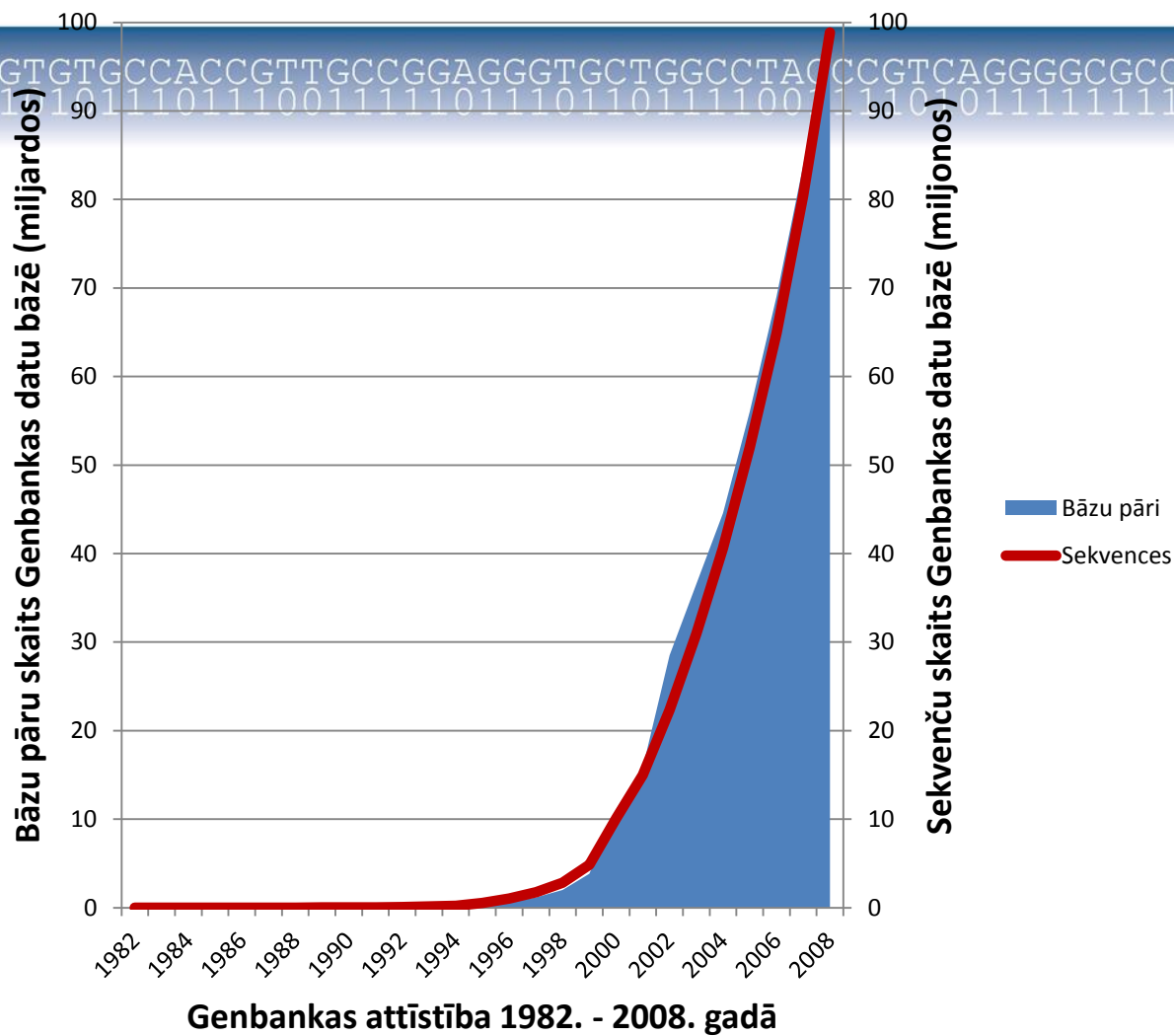
Molecular Biology of the Gene IVth ed. 1987

Genoma sekvenēšanas centri

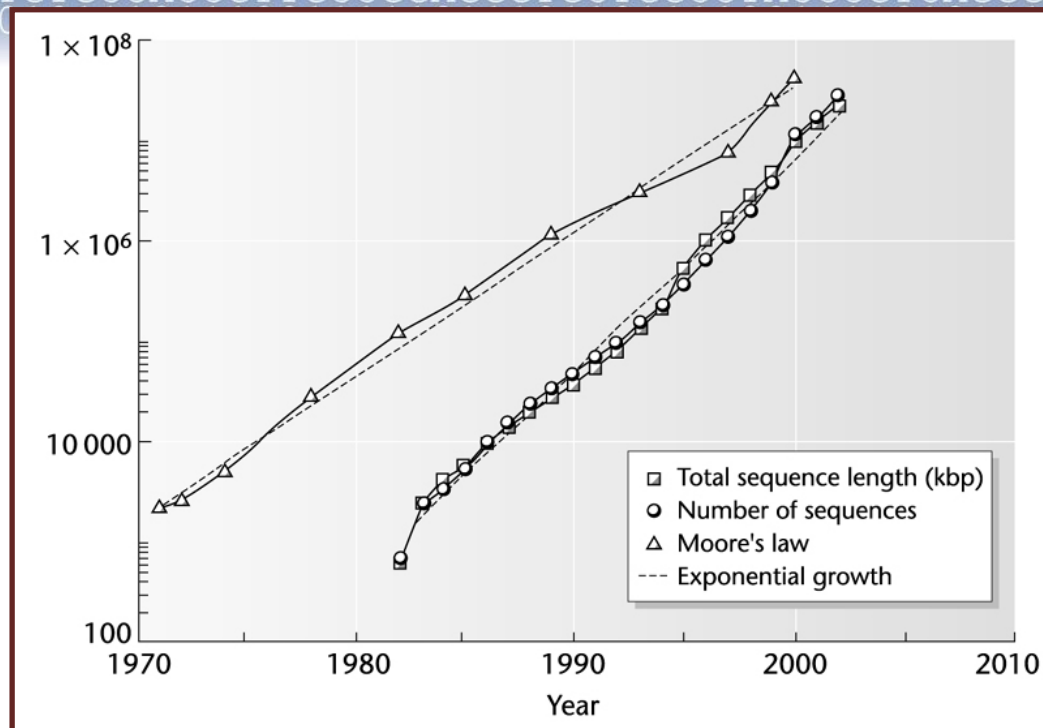
ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111110111011011110011110101111111100110110001



DNS sekvenču pieaugums GenBank

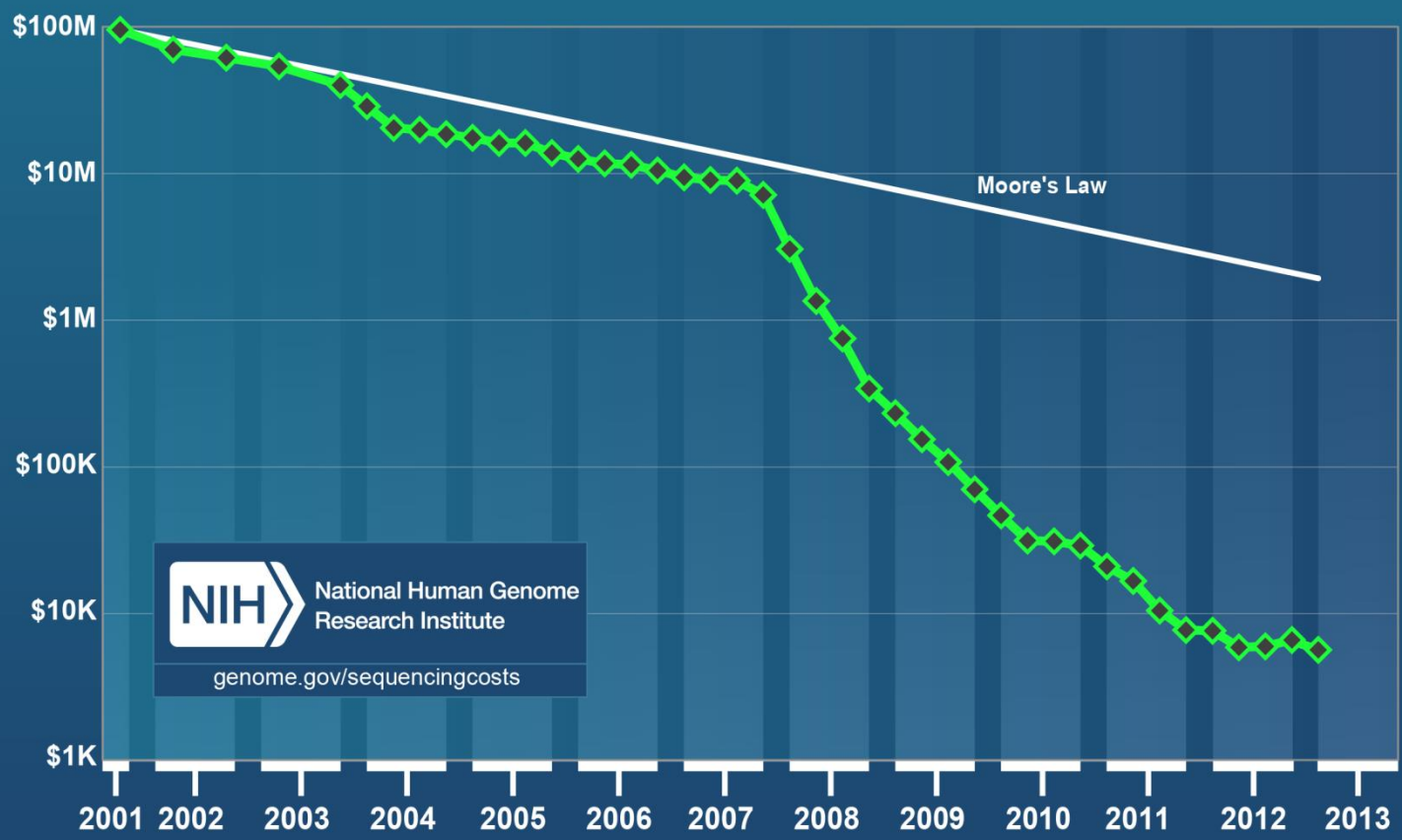


DNS sekvences un Mūra likums



Mūra likums – tranzistoru skaits datoru procesoros katru gadu dubultojas
Gan DNS sekvenču, gan tranzistoru skaita pieaugums uz datoru mikroshēmām ir eksponenciāls

Cost per Genome



Indivīda genoma sekvence

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011111011101101111001111010111111111111001101100001

- Cilvēka genoma projekts (1990. – 2003.) noteica pilnu cilvēka genoma sekвени, taču tā ir dažādu indivīdu genomu hibrīds
- Levy et al. (2007) The diploid genome sequence of an individual human. PLoS Biol 5: e254

Type	Number of Variants	bp Length	Min	Max	Mean
heterozygous SNP	1,762,541	1,762,541	1	1	1
homozygous SNP	1,450,860	1,450,860	1	1	1
heterozygous MNP	38,985	227,531	2	206	5.8
homozygous MNP	14,838	31,590	2	22	2.1
heterozygous indel	263,923	635,314	1	321	2.4
Complex	28,179	330,803	2	571	11.7
homozygous insertion	275,512	3,117,039	1	82,711	11.3
homozygous deletion	283,961	2,820,823	1	18,484	9.9
inversion	90	1,914,477	7	670,345	21,272
Total	4,118,889	12,290,978			

1000 genomu sekvences

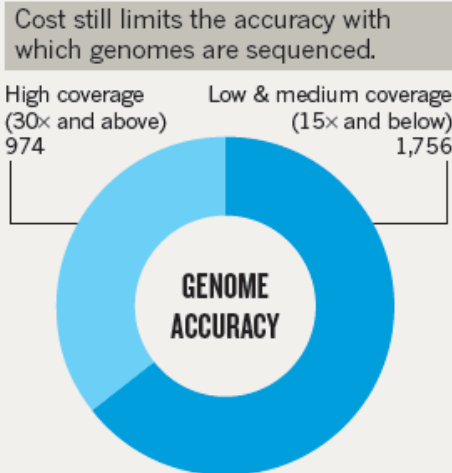
ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGTCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111101110110111100111101101111001111010111111111001101100001

- The Genomes Project (2010) A map of human genome variation from population-scale sequencing. Nature 467: 1061-1073

Why scientists want tens of thousands of genomes — and more

To understand populations

Comparing lots of genomes lets researchers identify points at which one genome differs from the next. Costs may be falling, but sequencing and data analysis are still pricey. So most researchers face a trade-off between the number of subjects and the accuracy in the sequences they can afford. For projects examining how populations commonly differ, sequencing a large number of individuals at relatively low accuracy or 'depth of coverage' is enough. About 900 genomes sequenced so far by the 1000 Genomes Project have been read three times on average.



To understand disease

Researchers trying to uncover rare disease-linked mutations — perhaps limited to just one family or an individual — need precision, typically sequencing each genome 30 times on average. Cancer genomes, many sequenced under the auspices of large collaborations, account for a sizeable chunk of high-coverage genome sequences completed to date. Projects scrutinizing people with diabetes, Crohn's disease and other disorders are starting to emerge. Analysing all the genome data is a huge challenge, as is turning genetic discoveries into clinical benefits.



Genomu struktūra un modernās metodes tās analīzei

000000001101011001001010011000101001100010010110001010111100101010010101001
1100000101000100000000001000100000100110111001000000000101100000001000110
1001001010000011100100000100010001000010000000001010001001001110001001000
0001011000101000000101000000001100111011101000001010010100100000101001001
01000110000100010000000101101100100101101101010101000110001110101011011011
101010110100100100010100010010100111000101100101101110101000100100001000111
111000100110111101001011011001110000001000100001011001000000000000101000
00000000101011011100011111010100010100110001011000001000001110111001000000
11000010101001000001000001010010010001010000000100110101001011100100010010
000101001010001100001001011100001101101100000110000000110101011110100011001
011011001100110100101001010001011100111111111111101100101110100100111010011
100100010010110100000100110000010100100010010000111100000001000111111010010
110000010100010000011101001100100100111101110110001100001001110100001001111

Genoms

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011111011101111001111001111010111100111101011111111001101100001

- Organisma genoms ir tā kopējā DNS molekulas nukleotīdu secība
- Prokariotu genomu veido hromosoma(s) un plazmīdas
- Eikariotu genoms sastāv no kodola genoma, kā arī mitohondriju un plastīdu (augiem) genoma

Genomika

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011110111011011110011110101111001111010111111111001101100001

- Bioloģijas zinātnes apakšnozare, kas nodarbojas ar sistemātisku un pilnīgu iedzimtības elementu pētīšanu molekulārā līmenī
- Genomika apraksta kā genomā kodētā informācija nosaka šūnu, organismu, populāciju un ekosistēmu attīstību

Bioinformātika un genomu struktūras anotācija

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111010111100111101011111111001101100001

```
1981 tatgtacgta cggaacgaaa cagtagtga caatgcgcga ggccacctag ctagctgcgc
2041 catgtatatg taagctatcc gtccacctca cgcgccaccg cgagctgatg cccgcggccg
2101 cgcgatcgac tcgaccacag gcctcgccgg cgcagccgta gccagctagc cgcgcgcgca
2161 ccgccgccat gtcgggagag ctctccgcgc gctcgtcgcc ctcgttcgcc gctcgcctt
2221 ccggcgccga cgacgcgcgg gcaggagagt cgacgccgct gcgggcccag gccgagagca
2281 gcaggcccgg cagcgccggc gtgaagctcc ggcggcggtg gcagcggcgg ctggggcggt
2341 ggcgcgtcgg ggacacgtgg gcgctggacc cgcgcgcgag gtgggtccgg gactggaacc
2401 gcgcctacct gctggcctgc gcggcggggc tgatggtgga cccgctcttc ctgtacgcgg
2461 tgtcggtgag cggcccgcgt atgtgcgtct tctcgcacgg ctggctcgcc gccgcggtca
2521 ccgcgctccg ctgcatggtg gacgccatgc acgcctggaa cctcctcag cagctccgcg
2581 tcgcgcgcgc ggcggcccgc gcggtggcgc gcggcccggc gggggtcgcg gacgaggagc
2641 aggccgaggc cgacgccgcg gccgcgcgaa gcctgcccgc gtacgccagg tccaggagag
2701 ggatggcgct cgacttcttc gtcatcctcc ccgtgatgca ggtgcacata tgccccctt
```


Eikariotu gēnu paredzēšana

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011110111011011110011110101111001111010111111111001101100001

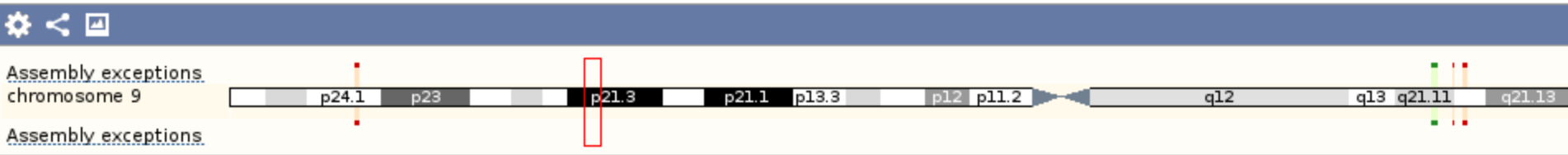
- Eikariotu genomi ir lieli un kompleksi
- Gēni ir lieli, sadalīti intronos un eksonos, pastāv alternatīvais splaisings
- Intronu sekvences evolucionē ātrāk par eksonu sekvencēm, tādēļ tās var stipri atšķirties pat starp evolucionāri tuviem organismiem
- Homoloģija ar citiem gēniem ļauj atrast tikai jau zināmus gēnus
- Eksperimentāla gēnu struktūras paredzēšana salīdzinot genomisko DNS ar cDNS
- Gēnu paredzēšana izmantojot bioinformātikas metodes

Eikariotu genomu komponenti

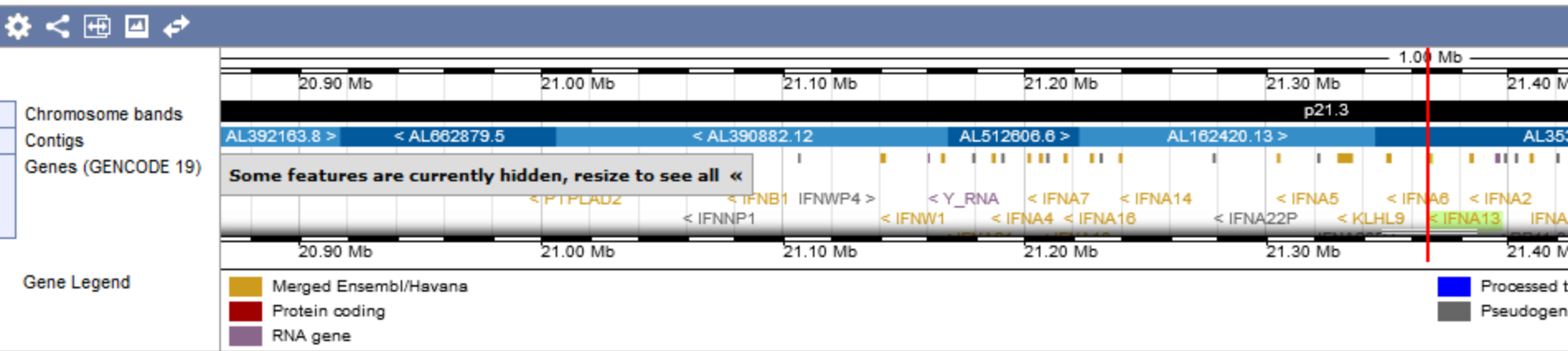
ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111110111011011110011110101111001111010111111111001101100001

- **Unikālas sekvences** - vienas kopijas gēni, unikālas regulatorās secības
- **Dažādo atkārtojumu anotācija ir ļoti svarīga, lai saprastu eikariotu genoma struktūru**
- **Bieži atkārtojumi** – minisatelīti, mikrosatelīti, telomēru atkārtojumi

Chromosome 9: 21,367,371-21,368,075



Region in detail i



<http://www.ensembl.org/index.html>

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101111001111001111010111100111101011111111001101100001

Bioloģiskās informācijas datubāzes. Informācijas meklēšanas un iegūšanas sistēmas

Bioloģiskā informācija

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011110111011110011110011110101111001111010111111111001101100001

- Molekulārās sekvences (nukleotīdu un aminoskābju)
- Proteīnu struktūras
- Gēnu ekspresijas dati
- Literatūra, kas saistīta ar bioloģisko informāciju

Bioinformātikas datu bāzes

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111010111100111101011111111001101100001

- Datubāze – organizēta informācijas kolekcija datoram nolasāmā formātā
- Dažādi datu bāzu veidi
Nukleotīdu, aminoskābju
Sekvenču, struktūras, literatūras, gēnu ekspresijas, molekulāro mijiedarbību, noteiktu organismu datu bāzes
- Informācijas meklēšana un iegūšana no dažādām datu bāzēm
- Internets un pieeja datu bāzēm

Arhīvs?

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111110111011011110011110101111001111010111111111001101100001

- Informācijas uzglabāšana ir tikai viena no datu bāzu funkcijām
- Datu bāze ir informācijas organizēšanas līdzeklis, informācijas meklēšanas efektivitāte datu bāzē ir atkarīga no tā, cik labi informācija ir organizēta
- Datu bāzes ir saistītas savā starpā un ļauj iegūt vispusīgu informāciju par meklēto jautājumu
- Datu bāzes piedāvā integrētus instrumentus, kas lietotājam ļauj patstāvīgi veikt datu analīzi

Kas kopīgs datu bāzēm?

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011110111011011110011110101111001111010111111111001101100001

- Dati ir sakārtoti tabulās, sadalīti pa grupām
- Katram datu bāzu ierakstam ir savs unikāls identifikators (ieraksta numurs – *Accession number* un GI)
- Dažādiem datu laukiem ir noteikts identifikators
- Datu bāzu ieraksti ir noteiktā standartizētā formātā, kas atvieglo to uzglabāšanu, meklēšanu un atrastās informācijas tālāku apstrādi

H.sapiens gene for interferon-alpha 13

GenBank: X75934.1

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

[Go to:](#)

LOCUS X75934 1539 bp DNA linear PRI 18-APR-2005

DEFINITION H.sapiens gene for interferon-alpha 13.

ACCESSION X75934

VERSION X75934.1 GI:439666

KEYWORDS interferon alpha.

SOURCE Homo sapiens (human)

ORGANISM [Homo sapiens](#)

Eukaryota; Metazoa; Chordata
Mammalia; Eutheria; Euarchonta
Catarrhini; Hominidae; Homo

REFERENCE 1

AUTHORS Henco,K., Brosius,J., Fujis
Hochstadt,J., Kovacic,T., P
Todokoro,K., Waelchli,M., N

TITLE Structural relationship of
pseudogenes

JOURNAL J. Mol. Biol. 185 (2), 227-

PUBMED [4057246](#)

REMARK (sites)

REFERENCE 2 (bases 1 to 1539)

AUTHORS Rostoks,N.

JOURNAL Unpublished

REFERENCE 3 (bases 1 to 1539)

AUTHORS Rostoks,N.

TITLE Direct Submission

JOURNAL Submitted (21-DEC-1993) Nil
of Genetics of Plants and M
Latvia LV-1842

FEATURES Location/Qualifier

source 1..1539

/organism="Homo sa

/mol_type="genomic

/db_xref="taxon:96

[TATA signal](#) 855..860

/citation=[1]

[mRNA](#) 887..>1539

/citation=[1]

[CDS](#) 956..1525

[CDS](#)

956..1525

/citation=[1]

/codon_start=1

/product="interferon-alpha 13"

/protein_id="CAA53538.1"

/db_xref="GI:439667"

/db_xref="GDB:136353"

/db_xref="GDB:136355"

/db_xref="GOA:P01562"

/db_xref="HGNC:5417"

/db_xref="HGNC:5419"

/db_xref="InterPro:IPR000471"

/db_xref="InterPro:IPR009079"

/db_xref="InterPro:IPR012351"

/db_xref="InterPro:IPR015589"

/db_xref="UniProtKB/Swiss-Prot:P01562"

/translation="MASPFALLMALVVLSCSSCSLGCGLPETHSLDNRRTLMLLAQM

SRISPSSCLMDRHDGFFPQEEFDGNQFQKAPAI SVLHELIIQQIFNLF TTKDSSAAWDE

DL LDKFCTELYQQLNDLEACVMQEERVGETPLMNADSI LAVKKYFRITL YL TEKKYS

PCAWEVVRAEIMRSLSLSTNLQERLRKE"

[sig peptide](#)

956..1024

/citation=[1]

[mat peptide](#)

1025..1522

/product="unnamed"

/citation=[1]

[variation](#)

984

ORIGIN

1 agatctcaaa gtcatatcat gagagggtgc ctctgcatac atatggtttg tcaactggcca

61 tcttatagat attgcttatg tttgatcctt agcatttctg tctgtgtttg gggctttgaa

121 atgaaatata aataatttat attttaacaa ttctactgaa gttgttcaac acatctatat

181 ttacgtcaag aattgaaga aeaattcttc acactctctc gctgactatc tagcagctac

CCGTCAGGGGCGCCCATGCTGGAATTC
111010111111111001101100001

Integrēta pieeja datu bāzēm

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011111011101111001111001111010111100111101011111111001101100001

- Zinot cik daudz un dažādu datu tipu pastāv (nukleotīdu un aminoskābju sekvences, literatūra, gēnu ekspresijas dati u.t.t.), būtu pamats bažām, kā atrast sakarības dažādu datu veidu starpā
- NCBI un EMBL piedāvā integrētas meklēšanas sistēmas – NCBI Entrez un EMBL SRS
- Abas meklēšanas sistēmas piedāvā gan teksta meklēšanu, gan molekulāro sekvenču meklēšanu izmantojot BLAST vai FASTA programmas

Būla operatori un ierobežotāji

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111010111100111101011111111001101100001

- ***Boolean operators*** – AND, OR un NOT

http://www.ncbi.nlm.nih.gov/books/NBK3837/#EntrezHelp.Entrez_Searching_Options

- ***Field*** – [Author], [Organism], [Journal], [Gene Name], bet noteikti *Field* attiecas uz noteiktām datu bāzēm

Meklēšanas sistēma paveic lielu darbu mūsu vietā, bet meklēšanas teksts vienmēr jāformulē pēc iespējas precīzi

Piemēram: «human BRCA1» (3846 rezultāti nukleotīdu datu bāzē, bet «Homo sapiens[Organism] AND BRCA1[Gene Name]» (147 rezultāti nukleotīdu datu bāzē)

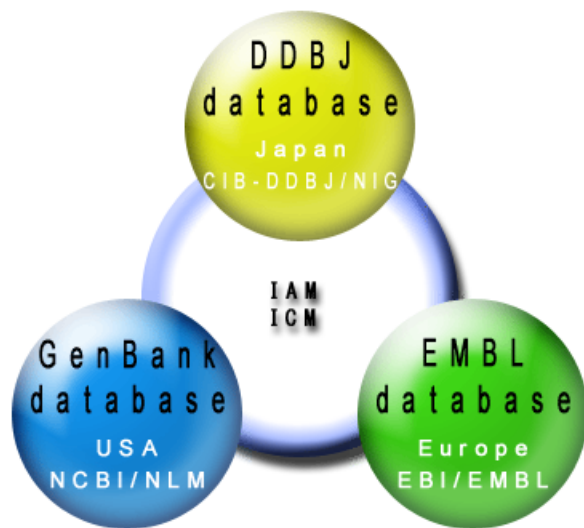
Primārās datu bāzes

- National Center for Biotechnology Information (NCBI)
GenBank
<http://www.ncbi.nlm.nih.gov>
- DNA Data Bank of Japan
<http://www.ddbj.nig.ac.jp/>
- The European Molecular Biology Laboratories (EMBL)
Nucleotide Archive <http://www.ebi.ac.uk/ena/>
- Primārās informācijas saturs visās 3 DB ir viens un tas pats, jo datu apmaiņa notiek katru dienu
- Sekvences identifikators (*accession number*) ir viens un tas pats visās DB

International Nucleotide Sequence Database Collaboration

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGTCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111010111100111101011111111001101100001

- <http://www.insdc.org/>
- INSDC ir sadarbības rezultāts starp GenBank, DDBJ, EMBL



**DDBJ/EMBL/GenBank
International
Nucleotide Sequence Database**

DDBJ: DNA Data Bank of Japan
CIB-DDBJ: Center for Information Biology and DNA Data Bank of Japan
NIG: National Institute of Genetics

EMBL: European Molecular Biology Laboratory
EBI: European Bioinformatics Institute

NCBI: National Center for Biotechnology Information
NLM: National Library of Medicine

IAM: International Advisory Meeting
ICM: International Collaborative Meeting

GenBank

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011110111011011110011110101111001111010111111111001101100001

- GenBank ir DNS sekvenču datu bāze
- GenBank datu bāze izveidota vēl pirms NCBI, kas ir struktūrvienība, kas šobrīd nodrošina šīs datu bāzes uzturēšanu
- Benson et al. (2013) GenBank. *Nucleic Acids Res* 41: D36-42

GenBank is a comprehensive public database of nucleotide sequences and supporting bibliographic and biological annotation

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111101110111100111101111001111010111111111001101100001

Genomu evolūcija un salīdzinošā genomika

Sugu un genomu evolūcija

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011111011101101111001111010111111111111001101100001

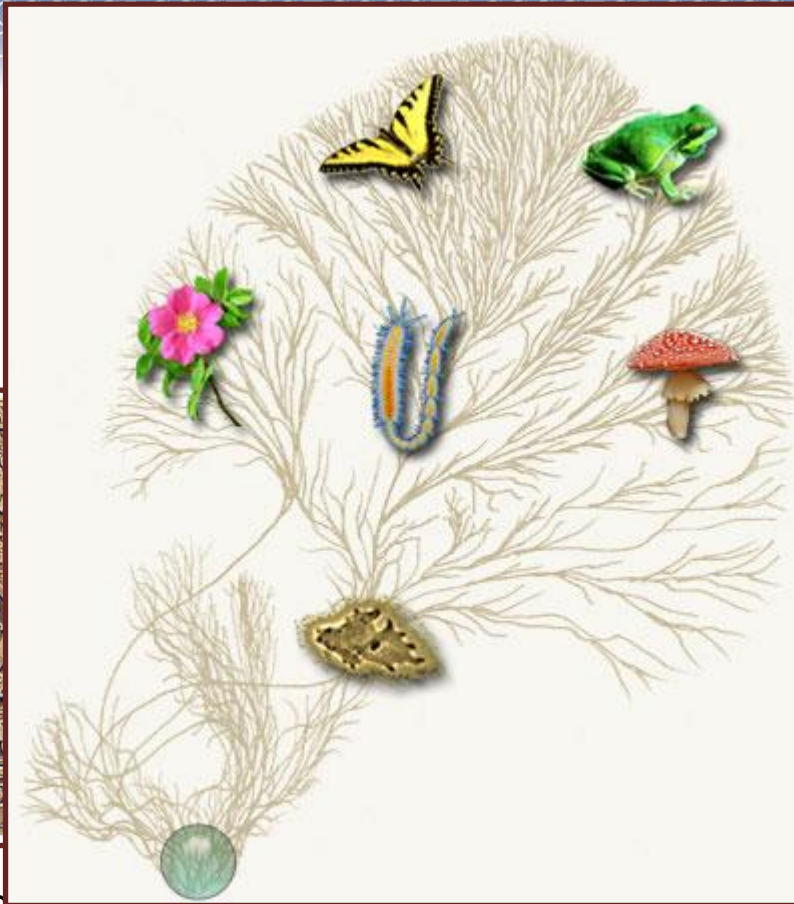
- Visas uz Zemes sastopamās organismu grupas ir cēlušās no viena kopēja senča, bet pēc tam evolucionējušas atsevišķi (izņēmums – horizontālā gēnu pārnese un organellu simbioze)
- DNS un proteīnu līdzības pamatā ir “*Identity by descent*”, vai arī konverģentā evolūcija
- Ja gēni dažādās sugās ir līdzīgi pēc DNS vai aminoskābju secības, tad iespējams tie veic līdzīgu funkciju
- Līdzīgas DNS un aminoskābju secības sauc par homologām (ar kopīgu izcelsmi)

Salīdzinošā genomika funkcionāli nozīmīgu elementu identifikācijai

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111101110111110011110101111001111010111111111001101100001

- Pēc būtības visspēcīgākā metode genomu anotācijai – ja DNS secības ir saglabājušas homologiju pēc miljoniem gadu ilgas neatkarīgas evolūcijas, tad tām droši vien ir funkcionāla nozīme

Dzīvības koks

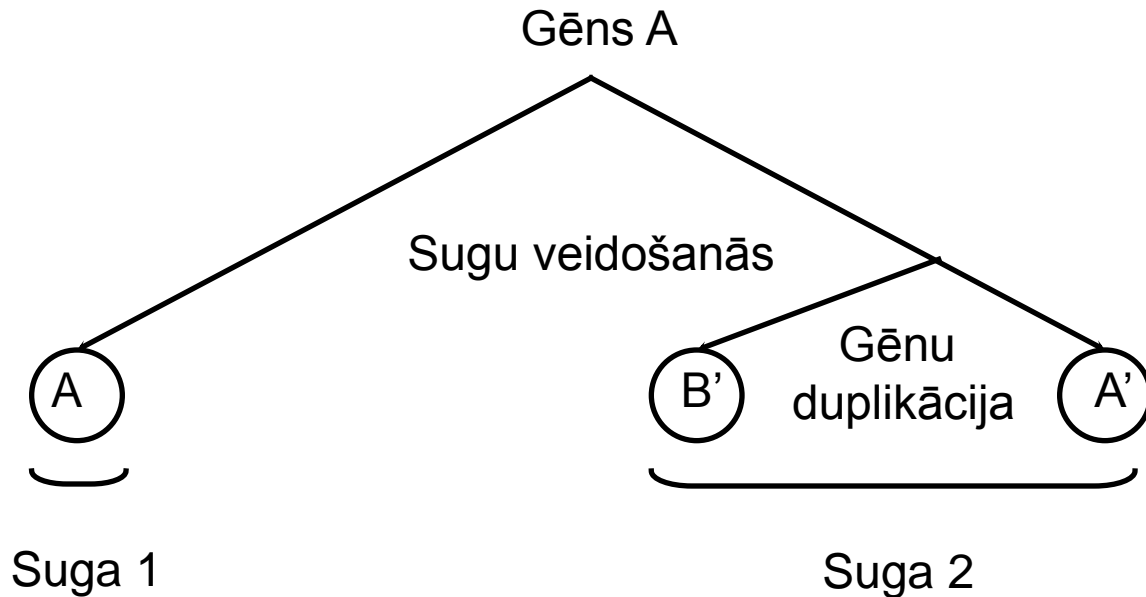


<http://www.tolweb.org/tree/>

Ortologi un paralogi

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111101110110111100111101011111111100110110001

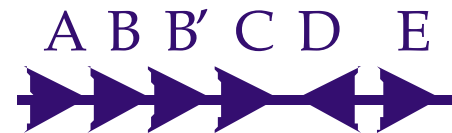
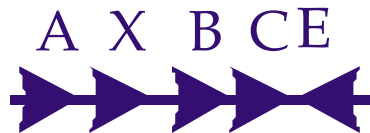
Pēdējais kopīgais sencis



A un A' ir ortologi. A' un B' ir paralogi. Pēc gēnu duplikācijas viena no gēna kopijām var turpināt pildīt iepriekšējo funkciju, kamēr otra gēna kopija mutāciju rezultātā var izmainīties un iegūt jaunu funkciju.

Gēnu kārtības saglabāšanās evolūcijas gaitā (sintēnija)

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111001111010111100111101011111111001101100001



Cilvēka un šimpanzes genomu salīdzinājums

ATGCAGAGCCAGA
0011010111010

Vol 437|1 September 2005|doi:10.1038/nature04072

nature

GAATTC
100001

ARTICLES

Initial sequence of the chimpanzee genome and comparison with the human genome

The Chimpanzee Sequencing and Analysis Consortium*

Here we present a draft genome sequence of the common chimpanzee (*Pan troglodytes*). Through comparison with the human genome, we have generated a largely complete catalogue of the genetic differences that have accumulated since the human and chimpanzee species diverged from our common ancestor, constituting approximately thirty-five million single-nucleotide changes, five million insertion/deletion events, and various chromosomal rearrangements. We use this catalogue to explore the magnitude and regional variation of mutational forces shaping these two genomes, and the strength of positive and negative selection acting on their genes. In particular, we find that the patterns of evolution in human and chimpanzee protein-coding genes are highly correlated and dominated by the fixation of neutral and slightly deleterious alleles. We also use the chimpanzee genome as an outgroup to investigate human population genetics and identify signatures of selective sweeps in recent human evolution.

cilvēks – šimpanze

cilvēks – pele

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGGTGCTGGCCTACCCGTCAGGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011110111011011110011110101111001111010111111111001101100001

- Genomu salīdzinājumi ļauj atbildēt uz fundamentāli atšķirīgiem jautājumiem
- **Cilvēks – pele.** Iespējams identificēt gēnus un kontroles elementus, kas saglabājušies nemainīgi evolūcijas gaitā. Var pētīt proteīnu funkcionālos domēnus, kas saglabājušies nemainīgi, vai arī gēnus, kas atšķir abus organismus.
- **Cilvēks – šimpanze.** Genomi, tai skaitā nekodējošā DNS, ļoti līdzīgi, kas ļauj identificēt mutācijas, kas radušās tieši cilvēka evolūcijas gaitā. Iespējams identificēt mutācijas gēnos un kodējošās daļās, kas potenciāli atšķir cilvēku no šimpanzes.

ENCODE project

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101111001111001111010111100111101011111111001101100001

- ENCyclopedia of DNA Elements (2003)
- Mērķis – identificēt visus funkcionālos elementus cilvēka genoma sekvencē
- Vairāk nekā 80% genoma sekvences noteikta sekvence, atklāti vairāk nekā 4 miljoni regulācijas rajonu
- Cilvēka genoma funkcionālo elementu analīze ar salīdzinošās genomikas palīdzību
- <http://www.genome.gov/10005107> un <http://www.nature.com/encode>

Genomu salīdzinājums no dažādām cilvēka populācijām

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGTCTGGCCTACCCGTCAGGGGCGCCCATGCTTGGAAATTC
00110101110100110101110111001111011101101111001111010111111111001101100001

Scienceexpress

Report

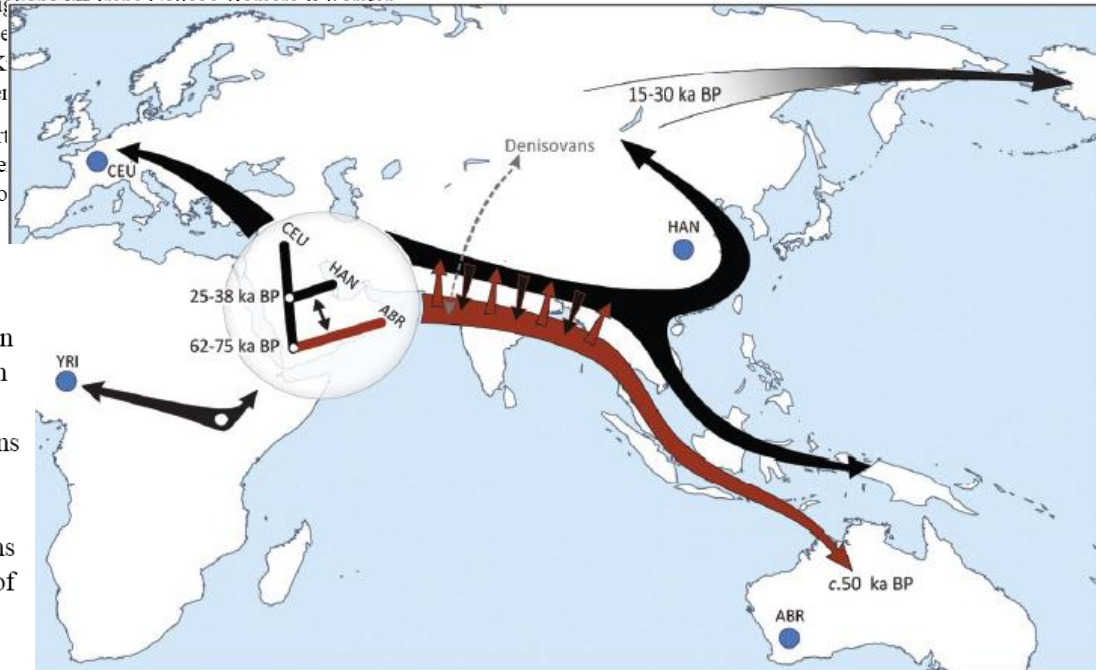
Rasmussen et al. (2011)
Science
10.1126/science.1211177

An Aboriginal Australian Genome Reveals Separate Human Dispersals into Asia

Morten Rasmussen,^{1,2*} Xiaosen Guo,^{2,3*} Yong Wang,^{4*} Kirk E. Lohmueller,^{4*} Simon Rasmussen,⁵ Anders Albrechtsen,⁶ Line Skotte,⁶ Stinus Lindgreen,^{1,6} Mait Metspalu,⁷ Thibaut Jombart,⁸ Toomas Kivisild,⁹ Weiwei Zhai,¹⁰ Anders Eriksson,¹¹ Andrea Manica,¹¹ Ludovic Orlando,¹ Francisco De La Vega,¹² Silvana Tridico,¹³ Ene Metspalu,⁷ Kasper Nielsen,⁵ María C. Ávila-Arcos,¹ J. Víctor Moreno-Mayar,^{1,14} Craig Muller,¹⁵ Joe Dortch,¹⁶ M. Thomas P. Gilbert,^{1,2} Ole Lund,⁵ Agata Wesolowska,⁵ Monika Karmin,⁷ Lucy A. Weinert,⁸ Bo Wang,³ Jun Li,³ Shuaishuai Tai,³ Fei Xiao,³ Tsunehiko Hanihara,¹⁷ George van Driem,¹⁸ Aashish R. Jha,¹⁹ François-Xavier Ricaut,²⁰ Peter de Knijff,²¹ Andrea B. Migliarese,^{2,3,6} David M. Lambert,²³ Søren Brunak,^{5,24} Peter Forster,^{25,26} Beata Keller,^{27,28} Ramneek Gupta,⁵ Carlos D. Bustamante,¹² Anders K. Balloux,⁸ Thomas Sicheritz-Pontén,^{5,29} Richard Villems,^{7,30} Rasmus Nielsen

¹Centre for GeoGenetics, Natural History Museum of Denmark, and Department of Geology, University of Copenhagen, Denmark. ²Sino-Danish Genomics Centre, University of Copenhagen, Denmark. ³Shenzhen Key Laboratory of Translational Genomics, Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China. ⁴Department of Biology, University of Copenhagen, Denmark. ⁵Department of Biology, University of Copenhagen, Denmark. ⁶Department of Biology, University of Copenhagen, Denmark. ⁷Department of Biology, University of Copenhagen, Denmark. ⁸Department of Biology, University of Copenhagen, Denmark. ⁹Department of Biology, University of Copenhagen, Denmark. ¹⁰Department of Biology, University of Copenhagen, Denmark. ¹¹Department of Biology, University of Copenhagen, Denmark. ¹²Department of Biology, University of Copenhagen, Denmark. ¹³Department of Biology, University of Copenhagen, Denmark. ¹⁴Department of Biology, University of Copenhagen, Denmark. ¹⁵Department of Biology, University of Copenhagen, Denmark. ¹⁶Department of Biology, University of Copenhagen, Denmark. ¹⁷Department of Biology, University of Copenhagen, Denmark. ¹⁸Department of Biology, University of Copenhagen, Denmark. ¹⁹Department of Biology, University of Copenhagen, Denmark. ²⁰Department of Biology, University of Copenhagen, Denmark. ²¹Department of Biology, University of Copenhagen, Denmark. ²²Department of Biology, University of Copenhagen, Denmark. ²³Department of Biology, University of Copenhagen, Denmark. ²⁴Department of Biology, University of Copenhagen, Denmark. ²⁵Department of Biology, University of Copenhagen, Denmark. ²⁶Department of Biology, University of Copenhagen, Denmark. ²⁷Department of Biology, University of Copenhagen, Denmark. ²⁸Department of Biology, University of Copenhagen, Denmark. ²⁹Department of Biology, University of Copenhagen, Denmark. ³⁰Department of Biology, University of Copenhagen, Denmark.

Fig. 2. Reconstruction of early spread of modern humans outside Africa. The tree shows the divergence of the Aboriginal Australian (ABR) relative to the CEPH European (CEU) and the Han Chinese (HAN) with gene flow between aboriginal Australasians and Asian ancestors. Purple arrow shows early spread of the ancestors of Aboriginal Australians into eastern Asia ~62,000 to 75,000 years B.P. (ka BP), exchanging genes with Denisovans, and reaching Australia ~50,000 years B.P. Black arrow shows spread of East Asians ~25,000 to 38,000 years B.P. and admixing with remnants of the early dispersal (red arrow) some time before the split between Asians and Native American ancestors ~15,000 to 30,000 years B.P. YRI, Yoruba.



ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111010111100111101011111111001101100001

Nukleīnskābju un proteīnu sekvenču līdzības pamatprincipi

**Dažādas salīdzināšanas metodes, to priekšrocības un
pielietošanas nosacījumi**

Sekvenču salīdzināšana

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011111011101101111001111010111111111111001101100001

- Angliski – *sequence alignment*
- Divi galvenie salīdzinājumu veidi – *pairwise sequence alignment* un *multiple sequence alignment*
- Latviskais tulkojums varētu būt “sekvenču pāru salīdzinājums” un “daudzkārtējs sekvenču salīdzinājums”

Sekvenču salīdzināšana

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111110111011110011110111100111111110011011011111111001101100001

- Ja mums ir divas vai vairākas DNS vai aminoskābju secības, tās analizējot vajadzētu:
 - izmērīt to vispārējo līdzību;
 - noteikt to līdzību katrā sekvences pozīcijā;
 - novērot līdzīgo un atšķirīgo rajonu izvietojumu sekvencēs;
 - novērtēt sekvenču evolucionārās attiecības

Sekvenču salīdzinājumu pielietojumi

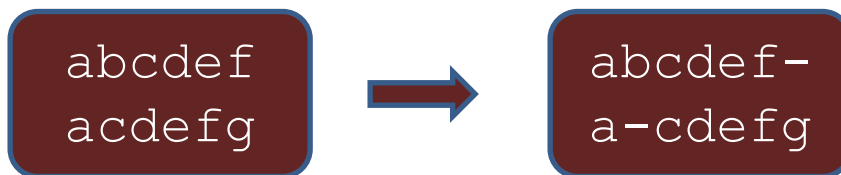
ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011111011101101111001111010111100111101011111111001101100001

- Genomu anotācijas pamatojas uz salīdzinājumu ar jau raksturotām sekvencēm (gēnu un atkārtojumu identifikācija)
- Funkcionālo domēnu noteikšana proteīnu sekvencēs
- DNS un proteīnu sekvenču radniecības pakāpes noteikšana, lai izvērtētu to filoģenētiskās attiecības un rekonstruētu šo secību evolūciju

Salīdzinājuma pamatprincipi

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGTCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111110111011110011110101111001111010111111111001101100001

- Sekvences salīdzina, nosakot to līdzību katrā sekvences pozīcijā
- Lai veidotu optimālu salīdzinājumu, iespējams sekvencēs ieviest pārtraukumus



- Nepieciešams izveidot kritērijus pēc kuriem spriest par salīdzinājuma kvalitāti

`acgtctga`
`agtttgat`

`acgtctga-`
`agttt-gat`

`acgtctga-`
`a-gtttgat`

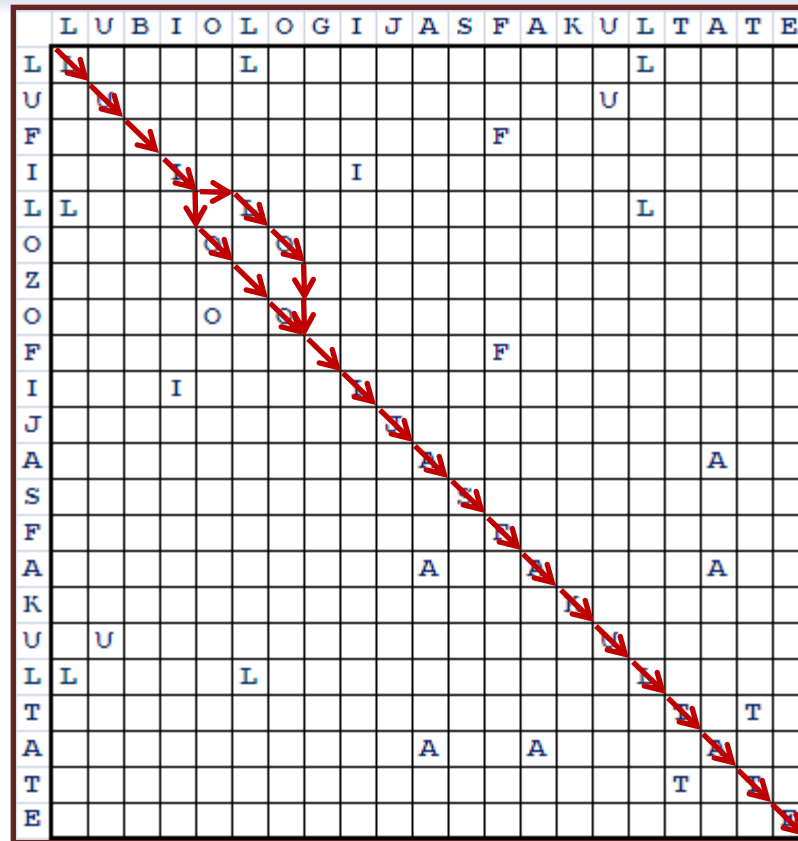
Dot plots

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111010111100111101011111111001101100001

- Visvienkāršākais divu sekvenču salīdzināšanas veids
- Dot plots ir kā tabula vai matrica, kurā rindas atbilst vienai sekvencai, bet kolonnas atbilst otrai sekvencai
- Vienkāršākajā veidā tiek atzīmēti tikai burti (nukleotīdi vai aminoskābes), kas katrā pozīcijā starp sekvencēm ir identiski

Dot plots un sekvenču salīdzinājums (*alignment*)

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111010111101011111111001101100001



Salīdzinājumu kvalitāte

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111110111011011110011110101111001111010111111111001101100001

L	U	B	I	-	O	L	O	G	I	J	A	S	F	A	K	U	L	T	A	T	E
L	U	F	I	L	O	Z	O	F	I	J	A	S	F	A	K	U	L	T	A	T	E
*	*		*		*		*	*	*	*	*	*	*	*	*	*	*	*	*	*	*

L	U	B	I	O	L	O	-	-	G	I	J	A	S	F	A	K	U	L	T	A	T	E
L	U	F	I	-	L	O	Z	O	F	I	J	A	S	F	A	K	U	L	T	A	T	E
*	*		*		*	*			*	*	*	*	*	*	*	*	*	*	*	*	*	*

Ceļš, kas tiek noiets dot plotā, atbilst sekvenču salīdzinājumam. Ja sekvences ir samērā līdzīgas, dot plots ļauj rekonstruēt sekvenču salīdzinājumu. Pastāv iespēja, ka vairāki ceļi dot plotā būs vienādi varbūtīgi. Nedrīkst pieņemt, ka ceļš dot plotā atbilst sekvenču molekulārajai evolūcijai

Divu sekvenču salīdzinājums

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111110111011110011110101111001111010111111111001101100001

- Kad izveidota salīdzinājumu novērtēšanas sistēma, iespējams meklēt **optimālu salīdzinājumu** – tādu, kas dod vislielāko punktu skaitu
- Ir situācijas, kad pastāv vairāki optimālie salīdzinājumi
- Globāls salīdzinājums – tiek salīdzinātas divas pilna garuma sekvences
- Lokāls salīdzinājums – salīdzina vienas sekvences segmentu ar otras sekvences segmentu
- Globāls salīdzinājums var būt ārkārtīgi laikietilpīgs (ja sekvenču garums ir n un m , tad matricas izmērs ir $n \times m$)

NCBI BLAST

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111110111011110011110101111001111010111111111001101100001

- ***Basic Local Alignment Search Tool***
- Līdzīgi kā globālajā sekvenču salīdzinājumā tiek izmantots dotplots, taču salīdzināšana notiek nevis ar pilna garuma sekvenci, bet gan tās noteikta garuma fragmentiem
- Altschul, S.F., Gish, W., Miller, W., Myers, E.W. & Lipman, D.J. (1990) Basic local alignment search tool. J. Mol. Biol. 215:403-410
- Altschul, S.F., Madden, T.L., Schäffer, A.A., Zhang, J., Zhang, Z., Miller, W. & Lipman, D.J. (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. Nucleic Acids Res. 25:3389-3402

BLAST mehānisms

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011110111011011110011110101111001111010111111111001101100001

- Visam pamatā – dotplots
- BLAST sadala meklēšanā izmantoto sekvenci noteikta garuma vārdos (piemēram, $k = 4$), tad atrod visas sekvences datu bāzē, kuras satur precīzi tādus vārdus
- BLAST cenšas pagarināt katru no sakrītošajiem vārdiem uz abām pusēm, nepieļaujot atšķirības un pārtraukumus (*mismatch, gap*)
- Pēc tam pagarinātie rajoni tiek apvienoti, pieļaujot atšķirības un pārtraukumus

BLAST

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011110111011011110011110101111001111010111111111001101100001

- BLASTN – nukleotīdu secība pret nukleotīdu datu bāzi
- BLASTP – proteīnu aminoskābju secība pret proteīnu datu bāzi
- BLASTX – translēta nukleotīdu secība pret proteīnu datu bāzi
- TBLASTN – aminoskābju secība pret translētu nukleotīdu datu bāzi
- TBLASTX – translēta nukleotīdu secība pret translētu nukleotīdu datu bāzi

Daudzkārtēja sekvenču salīdzināšana

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111110111011011110011110101111001111010111111111001101100001

- *Multiple sequence alignment (MSA)* – vairāku nukleotīdu vai aminoskābju secību salīdzinājums
- Viena no biežāk lietotajām bioinformātikas metodēm
- Plaši izmanto dažādu proteīnu domēnu meklēšanā, proteīnu struktūras modelēšanā un filoģenētiskajā analīzē

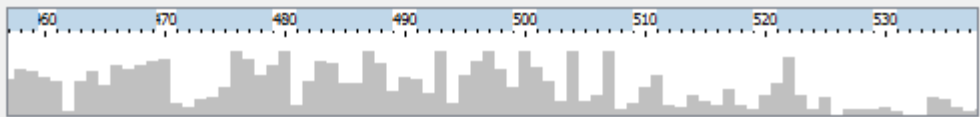
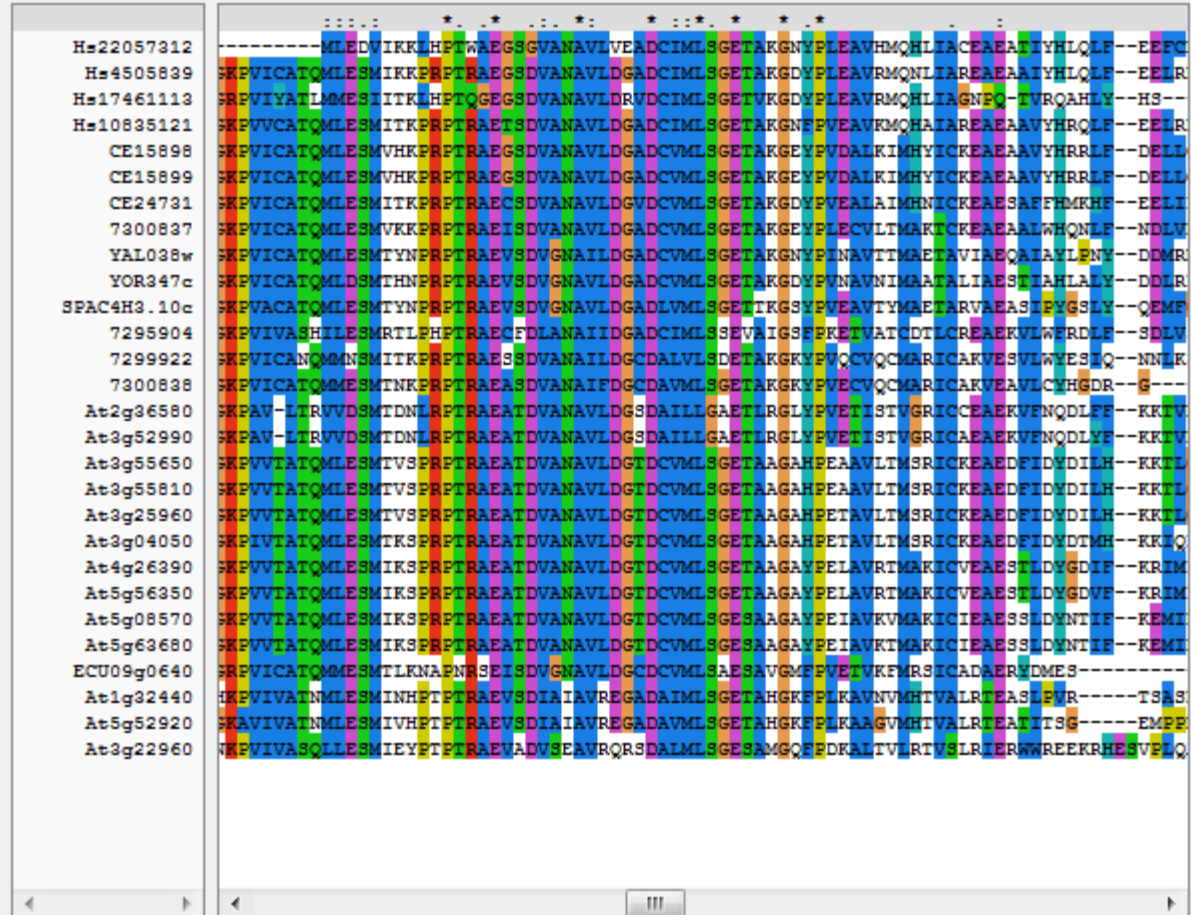
Mode: Multiple Alignment Mode Font: 10

>ref|XP_003428362.1| GM PREDICTED [Ornithorhynchus anatinus] Length=704

GENE ID: 100074825 LOC100074825 | pa [Ornithorhynchus anatinus]

Score = 644 bits (1662), Expect = Identities = 334/394 (85%), Positive

Query 48 VSNGCVSKILGRYYETGSIRPRAIG VSNGCVSKILGRYYETGSIRPRAIG
Sbjct 315 VSNGCVSKILGRYYETGSIRPRAIG
Query 108 SEGVCINDNIPSVSSINRVLRLAS SEGVCINDNIPSVSSINR R
Sbjct 375 SEGVCINDNIPSVSSINRGR--GE
Query 150 N-GQTGSWGTRPGWYPGTSVPGQPT G G G R + P P
Sbjct 433 ERGMEGCPGPRSKVWIDCPAPIIP-
Query 209 KLQRNRTSFTQEIEALEKEFERTH KLQRNRTSFTQEIEALEKEFERTH
Sbjct 491 KLQRNRTSFTQEIEALEKEFERTH
Query 269 EKLRNQRRQASNTPSHIPISSSFST EKLRNQRRQAS TPSHIPISSSFST
Sbjct 551 EKLRNQRRQASTTPSHIPISSSFST
Query 329 LPPMPSFTMANNLPMQPPVPSQTSS LPPMPSFTMANNLPMQPPVPSQTSS
Sbjct 611 LPPMPSFTMANNLPMQPPVPSQTSS
Query 389 SGTITSTGLISPGVSVPVQVPGSEPD SGTITSTGLISPGVSVPVQVPGSEPD
Sbjct 671 SGTITSTGLISPGVSVPVQVPGSEPD



Daudzkārtēju sekvenču salīdzinājumu veidošana

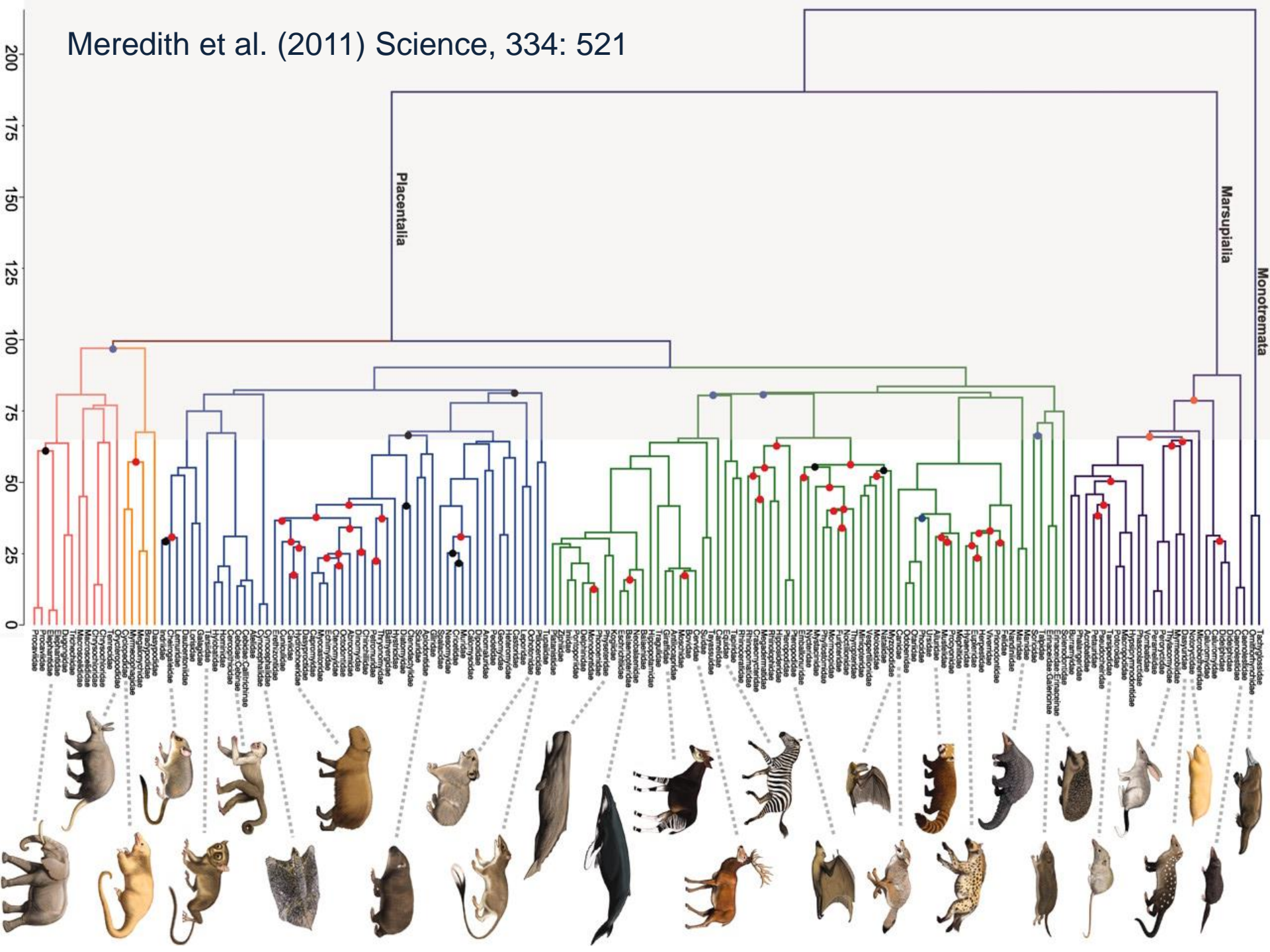
ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111010111100111101011111111001101100001

- Tāpat kā divu sekvenču salīdzināšanas gadījumā tiek meklēta līdzība starp nukleotīdiem vai aminoskābēm katrā salīdzinājuma pozīcijā
- Daudzkārtējā sekvenču salīdzinājumā var izdalīt sekojošus etapus:
 1. salīdzina savā starpā visus sekvenču pārus un izveido attālumu matricu;
 2. aprēķina dendrogrammu, kas balstīta uz attālumu matricas;
 3. Vienu pēc otras salīdzinājumam pievieno sekvences atbilstoši to secībai dendrogrammā sākot ar līdzīgākajām sekvencēm

Clustal X un Clustal W

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111010111100111101011111111001101100001

- Viena no populārākajām salīdzināšanas metodēm
- Oriģinālā publikācija - Higgins and Sharp (1988) **CLUSTAL**: a package for performing multiple sequence alignment on a microcomputer. Gene 73,237-244
- Thompson et al. (1994) **CLUSTAL W**: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, positions-specific gap penalties and weight matrix choice. Nucleic Acids Research, 22:4673-4680
- Thompson et al. (1997) The ClustalX windows interface: flexible strategies for multiple sequence alignment aided by quality analysis tools. Nucleic Acids Research, 24:4876-4882
- <http://www.clustal.org/clustal2/>



Evolūcija

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011111011101111001111010111100111101011111111001101100001

"Nothing in Biology Makes Sense Except in the
Light of Evolution"

Theodosius Dobzhansky (1973)

<http://www.2think.org/dobzhansky.shtml>

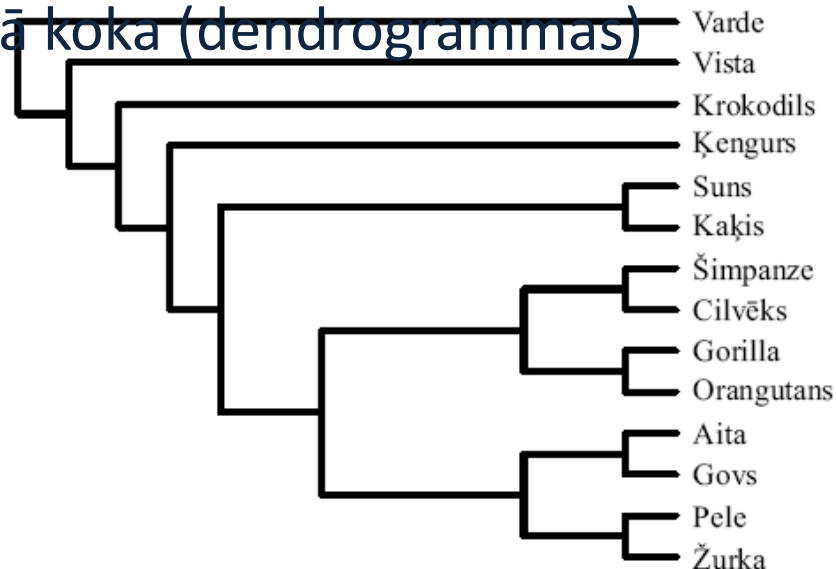
Filogenēze

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
00110101110100110101110111001111011101101111001111010111100111101011111111001101100001

- Filogenēze (*Phylogeny*) – sugas (vai citas taksonomiskas vienības) evolucionārā vēsture

Filogenētiskās analīzes pamatā ir ideja, ka visi dzīvie organismi ir cēlušies no kopīga senča

Evolucionārās attiecības starp taksonomiskajām vienībām parasti tiek attēlotas filogenētiskā koka (dendrogrammas) veidā



Molekulārā filoģenēze

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111110111011011110011110101111001111010111111111001101100001

- Evolūcijas kokus var veidot balstoties uz dažādām pazīmēm – morfoloģiskām, fizioloģiskām u.t.t.
- Molekulārā filoģenēze tāpat ir sugu (taksonomisko vienību) evolucionāro attiecību pētīšana izmantojot molekulārās sekvences

Terminoloģija

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111110111011011110011110101111001111010111111111001101100001

- Homoloģija (*Homology*) – līdzība, ko nosaka kopīga izcelšanās. Apgalvojums par kādu objektu homoloģiju nozīmē, ka tiek pieņemta to evolucionāra saistība
- Līdzība (*Similarity*) – objektu līdzības vai atšķirības mērs, neiedziļinoties to varbūtējā saistībā
- Klasteri – līdzīgu objektu grupas, kas tiek sakārtotas pēc līdzības pakāpes.
- Hierarhiskie klasteri – klasteri, kas veido klasterus, kas veido klasterus...

Sugu un gēnu dendrogrammas

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
0011010111010011010111011100111101110110111100111101111001111010111111111001101100001

- Parastākās dendrogrammas attēlo dažādu taksonomisko vienību evolūciju – sugu dendrogrammas
- Molekulārie dati (DNS un proteīnu sekvences) ļauj noteikt atsevišķu gēnu evolūciju
- Gēnu un sugu dendrogrammas ne vienmēr sakrīt

Dendrogrammas

ATGCAGAGCCAGATCGTGTGCCACCGTTGCCGGAGGGTGCTGGCCTACCCGTCAGGGGCGCCCATGCTGGAATTC
001101011101001101011101110011110111011011110011110101111001111010111111111001101100001

- Dendrogrammas topoloģija – norāda taksonomisko vienību radniecību
- Zaru garumi – norāda evolucionāro attālumu vai laiku kopš taksonomiskās vienības nodalījās.
- Dendrogrammas mezgli (*node*) – norāda kopējo priekšteci
- Molekulārais pulkstenis – atkarīgs no mutāciju ātruma. To kalibrē izmantojot fosīliju datus