



# Big Data kihívások a genomikában

Barta Endre

DEOEC, BMBI, Klinikai Genomikai Központ

MBK, Mezőgazdasági Genomikai és  
Bioinformatikai Csoport

[barta.endre@unideb.hu](mailto:barta.endre@unideb.hu)

[barta@abc.hu](mailto:barta@abc.hu)

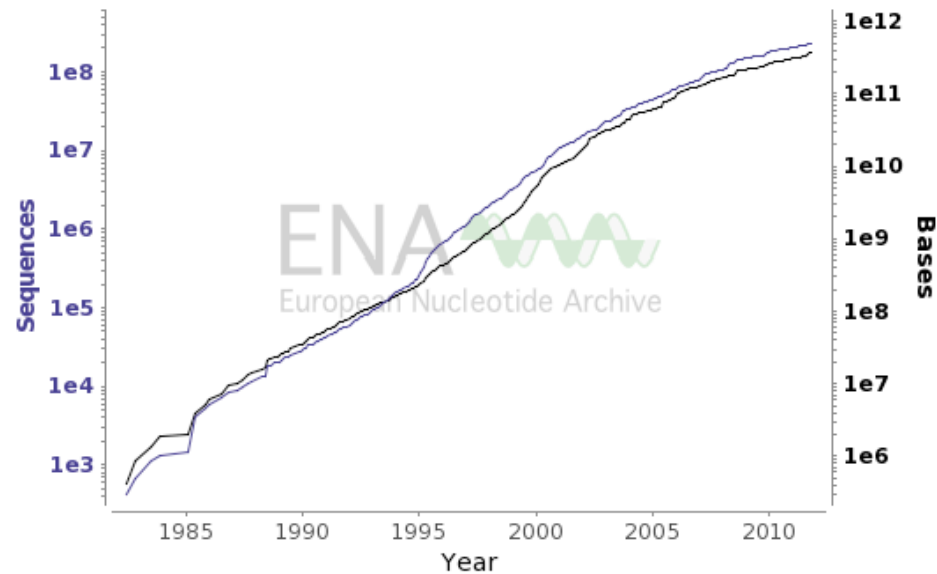


# Adatbázisok és a tárolókapacitás növekedése (MBK vs EMBL)

- 1990: MicroVax szerver
  - 2x 160 Mbyte HDD – 50 mbp
- 1993: SUN SparcServer 1000
  - 8x 512 Mbyte HDD – 150 mbp
- 1997: SUN Ultra Enterprise II
  - 4x 9 Gbyte HDD – 1 gbp
- 2002: SUN Vire480
  - 8x 180 Gbyte HDD – 23 gbp
  - Szekvencia + annotáció + index  
= ~ 100 Gbyte
- 2007: LINUX szerver
  - 6x 500 Gbyte HDD – 100 gbp
- 2011: LINUX klaszter
  - 8x 600 és 24x2000 Gbyte HDD – 370 gbp

## EMBL-Bank Growth

28-Nov-2011



— Sequences (229.7 millions) — Bases (373.0 billions)





# MBK külső hálózat

- 1991-ben 2400 bit/s modem Triesztbe
- 1992-ben 9600 bit/s X25 összeköttetés, HUNGARNET!
- 1993-ban 19.2 Kbit/s bérelt vonali összeköttetés
- 1994-ben 2 Mbit/s mikrohullámú összeköttetés
- 2000-ben 155 Mbit/s ATM összeköttetés (SZIE)
- 2002-ben 1Gbit/s optikai összeköttetés



# Adat transzfer szekvenátor és a szerver között



~1.4 Tbyte adat 11 nap futás alatt

Debreceni Egyetem Orvos-és Egészségtudományi Centrum  
Klinikai Genomikai és Személyre Szabott Orvoslási Központ

Vezető: Nagy László

Laborvezető: Bálint L. Bálint

Szekvenálás: Gyuris Tibor és Póliszka Szilárd

Szerver: Horváth Attila és Barta Endre

Bioinformatika: Nagy Gergely, Horváth Attila, Jónás Dávid  
és Barta Endre

## Prices:

Description	Qty	Unit Price	Discount	Total Price
Library preparation for HiSeq, Genomic Shotgun	1	800.00	10%	720.00
High-throughput sequencing on the illumina HiSeq, Paired-reads channel 2x100 bp	8	6,468.00	30%	36,220.80
Data delivery on external hard drive (OFFERED) (up to 500 GB, including handling and shipping)	1	250.00	100%	0.00

Sub-total: 36,940.80

Swiss VAT: EXPORT 0.00

**TOTAL (CHF) 36,940.80**

Prices are in Swiss Francs. In Switzerland, VAT of 8% is charged.  
Fasteris VAT number: 579 029 or CHE-110.264.966 TVA.

By sending us samples you agree with our Terms and Conditions and you are placing an order.



# Genomika, funkcionális genomika

- **Genomika:** A teljes örökítő anyag (magi, mitokondriális és ha van a kloroplasztisz DNS) szekvenálása, a szekvencia annotálása, elemzése
- **Funkcionális genomika:** A gének összességének, vagy egy-egy csoportjának a vizsgálata nagy áteresztőképességű kísérletekkel, amelyek vagy a genomszekvencia alapján lettek megtervezve (microarray, SNP), vagy szekvenálást használnak (ChIP-seq, RNA-seq stb.)

# Új generációs szekvenálás első publikáció 2005.

- A főbb modellszervezeteket már megszekvenálták hagyományos elsőgenerációs térképezős módszerekkel
- Különböző genomprogramok folyamatban

## Genome sequencing in microfabricated high-density picolitre reactors

Marcel Margulies<sup>1\*</sup>, Michael Egholm<sup>1\*</sup>, William E. Altman<sup>1</sup>, Said Attiya<sup>1</sup>, Joel S. Bader<sup>1</sup>, Lisa A. Bemben<sup>1</sup>, Jan Berka<sup>1</sup>, Michael S. Braverman<sup>1</sup>, Yi-Ju Chen<sup>1</sup>, Zhoutao Chen<sup>1</sup>, Scott B. Dewell<sup>1</sup>, Lei Du<sup>1</sup>, Joseph M. Fierro<sup>2</sup>, Xavier V. Gomes<sup>1</sup>, Brian C. Godwin<sup>1</sup>, Wen He<sup>1</sup>, Scott Helgesen<sup>1</sup>, Chun He Ho<sup>1</sup>, Gerard P. Irzyk<sup>1</sup>, Szilveszter C. Jando<sup>1</sup>, Maria L. I. Alenquer<sup>1</sup>, Thomas P. Jarvie<sup>1</sup>, Kshama B. Jirage<sup>1</sup>, Jong-Bum Kim<sup>1</sup>, James R. Knight<sup>1</sup>, Janna R. Lanza<sup>1</sup>, John H. Leamon<sup>1</sup>, Steven M. Lefkowitz<sup>1</sup>, Ming Lei<sup>1</sup>, Jing Li<sup>1</sup>, Kenton L. Lohman<sup>1</sup>, Hong Lu<sup>1</sup>, Vinod B. Makhijani<sup>1</sup>, Keith E. McDade<sup>3</sup>, Michael P. McKenna<sup>1</sup>, Eugene W. Myers<sup>2</sup>, Elizabeth Nickerson<sup>1</sup>, John R. Nobile<sup>1</sup>, Ramona Plant<sup>2</sup>, Bernard P. Puc<sup>1</sup>, Michael T. Ronan<sup>1</sup>, George T. Roth<sup>1</sup>, Gary J. Sarkis<sup>1</sup>, Jan Fredrik Simons<sup>2</sup>, John W. Simpson<sup>2</sup>, Maithreyan Srinivasan<sup>1</sup>, Karrie R. Tartaro<sup>1</sup>, Alexander Tomasz<sup>3</sup>, Kari A. Vogt<sup>1</sup>, Greg A. Volkmer<sup>1</sup>, Shally H. Wang<sup>1</sup>, Yong Wang<sup>1</sup>, Michael P. Weiner<sup>4</sup>, Pengguang Yu<sup>1</sup>, Richard F. Begley<sup>3</sup> & Jonathan M. Rothberg<sup>1</sup>

**454 Life Sciences Corporation**

Nature 437 (7057)  
376-380



# Moore törvény, Moore genom

- Újabb és újabb technológiák
  - Complete Genomics (hibridizálás)
  - Pacific Biosciences, Helicos egy molekula szekvenálása
  - Ion torrent: microchip alapú szekvenálás
  - Nanopore technológia

## Los Angeles Times

LOCAL U.S. WORLD BUSINESS SPORTS ENTERTAINMENT HEALTH LIVING TRAVEL OPINION Search GO

BREAKING CRIME L.A. APPS WEATHER TRAFFIC OBITS COMMUNITY CROSSWORDS COMICS

YOU ARE HERE: LAT Home → Collections → Chip

ADS BY GOOGLE

**Need Next Gen Sequencing?**  
We're an Experienced Sequencing Lab that can Help Advance Your Research  
[www.ambrigen.com/NextGenSequencing](http://www.ambrigen.com/NextGenSequencing)

**NextGen Sequencing**  
Ion Torrent PGM Sequencing Service  
Low Introductory Rate, Try Now!  
[www.EpochLifeScience.com](http://www.EpochLifeScience.com)

BOOSTER SHOTS: Oddities, musings and news from the health world

### Cheaper DNA sequencing possible with computer chip technology

July 22, 2011 | By Marissa Cevallos, HealthKey / For the Booster Shots blog

The race to decode a person's genome on the cheap got tighter this week. The sequencing company Ion Torrent announced this week in Nature that it used a \$49,500 machine, based on computer chip technology, to unravel a full human genome—apty, using the DNA of Gordon Moore, co-founder of Intel.

ADS BY GOOGLE



## Big Data in Healthcare – Saving Lives with Personalized Medicine

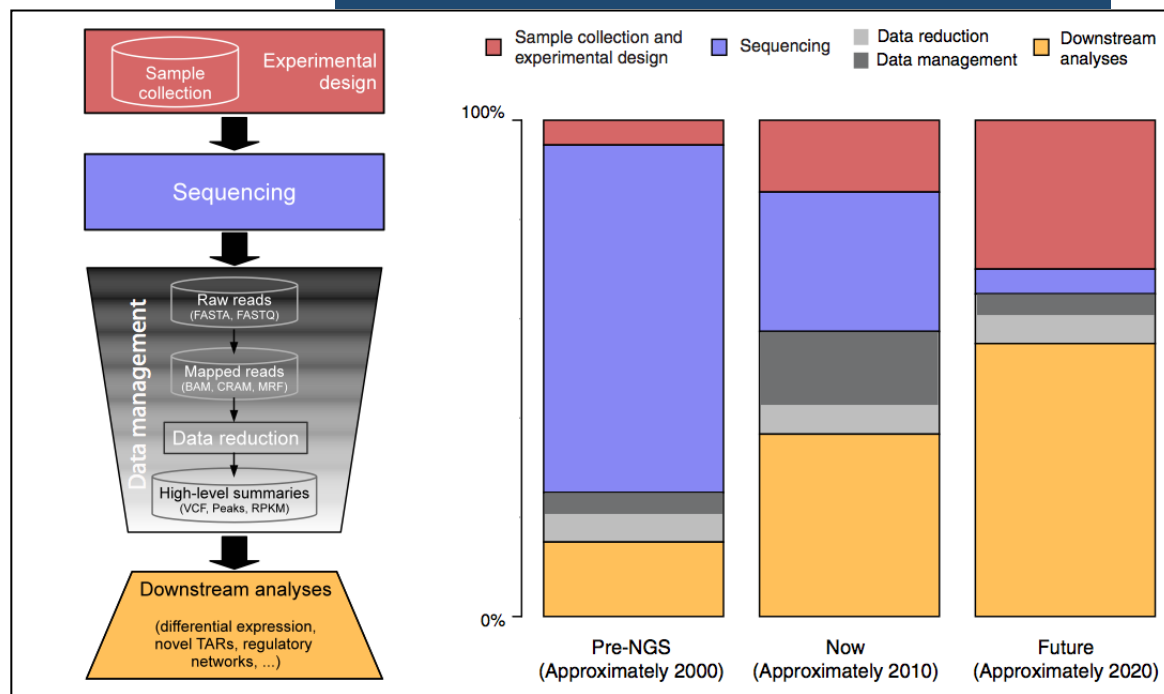
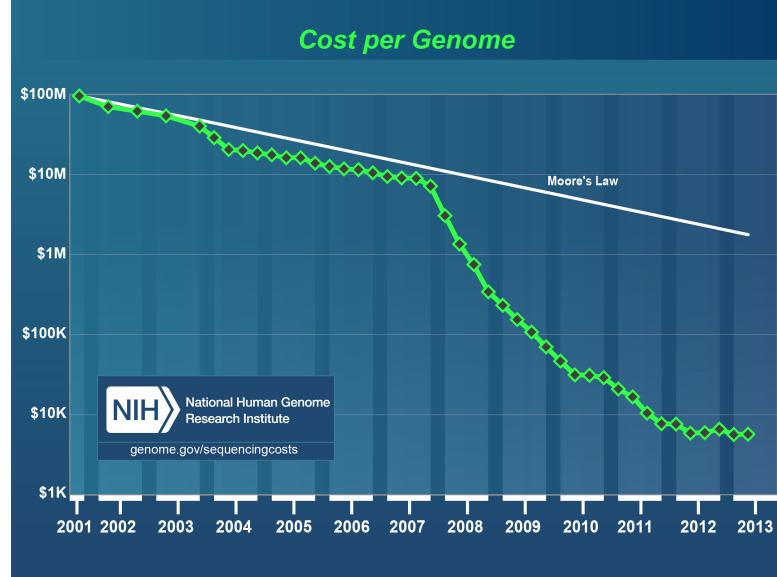
July 28, 2013 by [Rich Brueckner](#)

In this video, Intel Fellow Eric Dishman shares a personal story of how his 25 year struggle with kidney cancer was finally resolved through Big Data and personalized medicine. Dishman's doctors were able to treat him successfully after sequencing his complete genome.

For personalized medicine, only 50,000 people on earth have had their entire genome sequenced. As Dishman describes, this used to take months of supercomputing time and was simply out of reach for most people. This is starting to change thanks in part to advances in computing powered by Moore's Law.

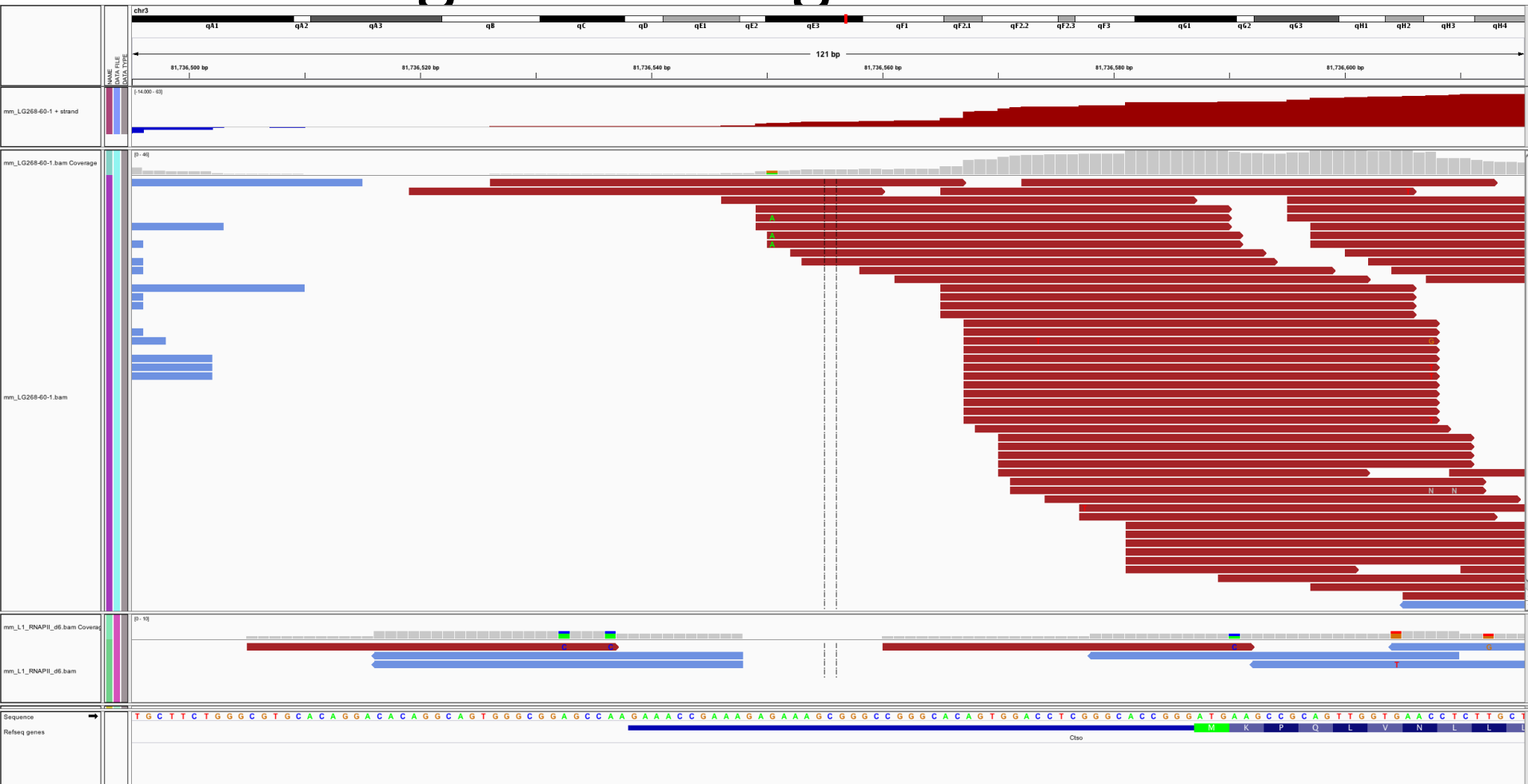
# Paradigmaváltás a kutatásban

- A „hagyományos” kutatási módszer hipotézis alapú
- A genomika lehetővé teszi az adatalapú kutatásokat
  - A kísérlet megtervezése a fenotípus alapján
  - Adatgenerálás
  - Adatfeldolgozás, szisztematikus analízis
  - Új hipotézisek felállítása
  - Kísérleti validáció

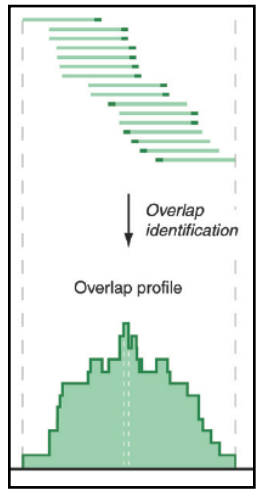
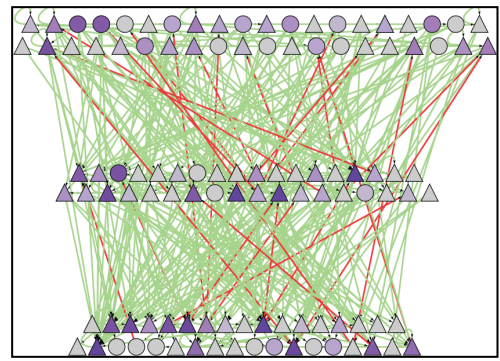
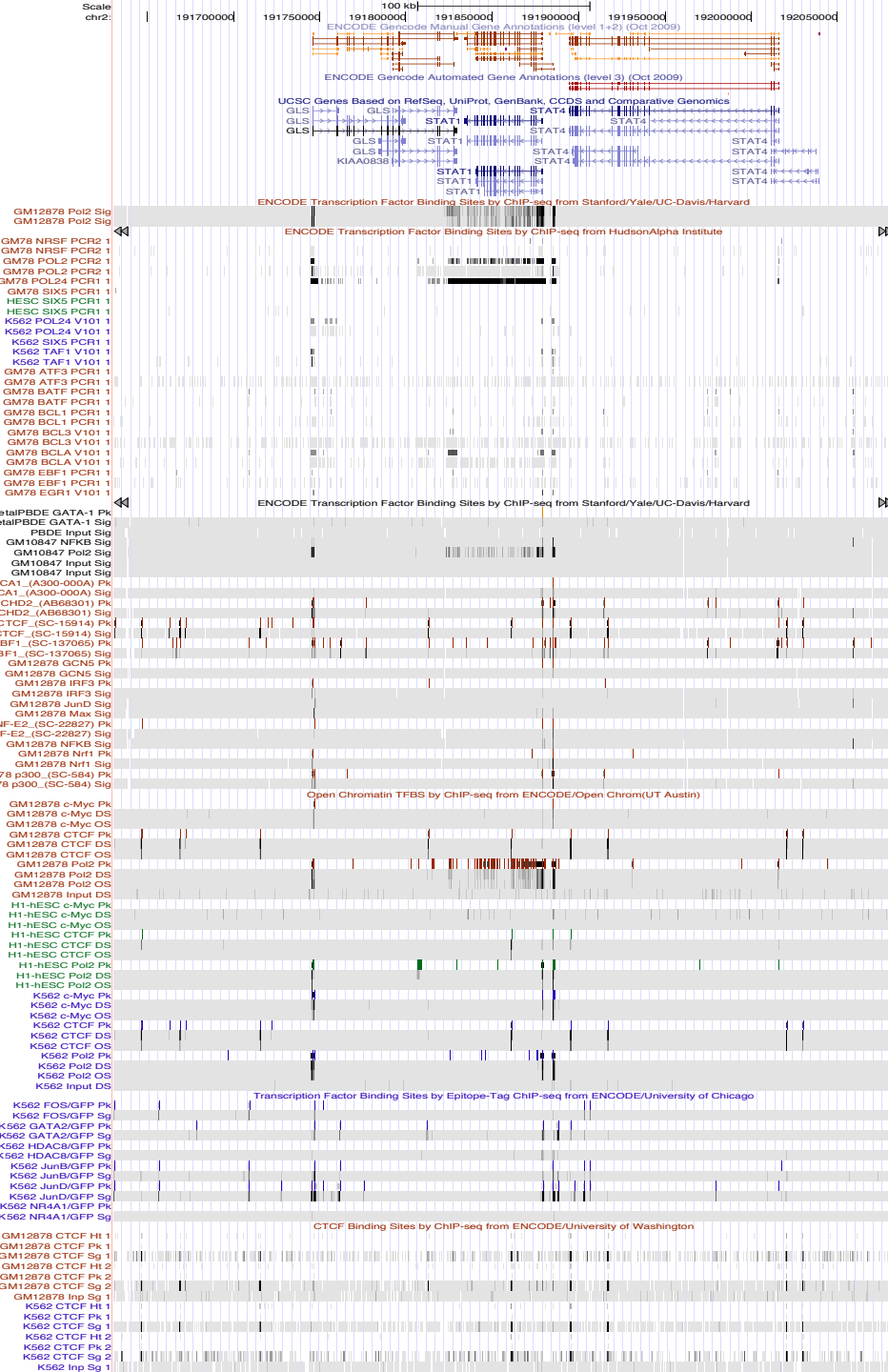




# Szekvenálás megjelenítése genomböngészőben



# Hierarchikus ENCODE adatmegjelenítés



nature ENCODE

Home Research Threads Additional Research News and Comment About Sponsor

nature ENCODE explorer

THREADS PAPERS

PRODUCTION WITH SUPPORT FROM illumina

THREAD OVERVIEW: Characterization of network topology

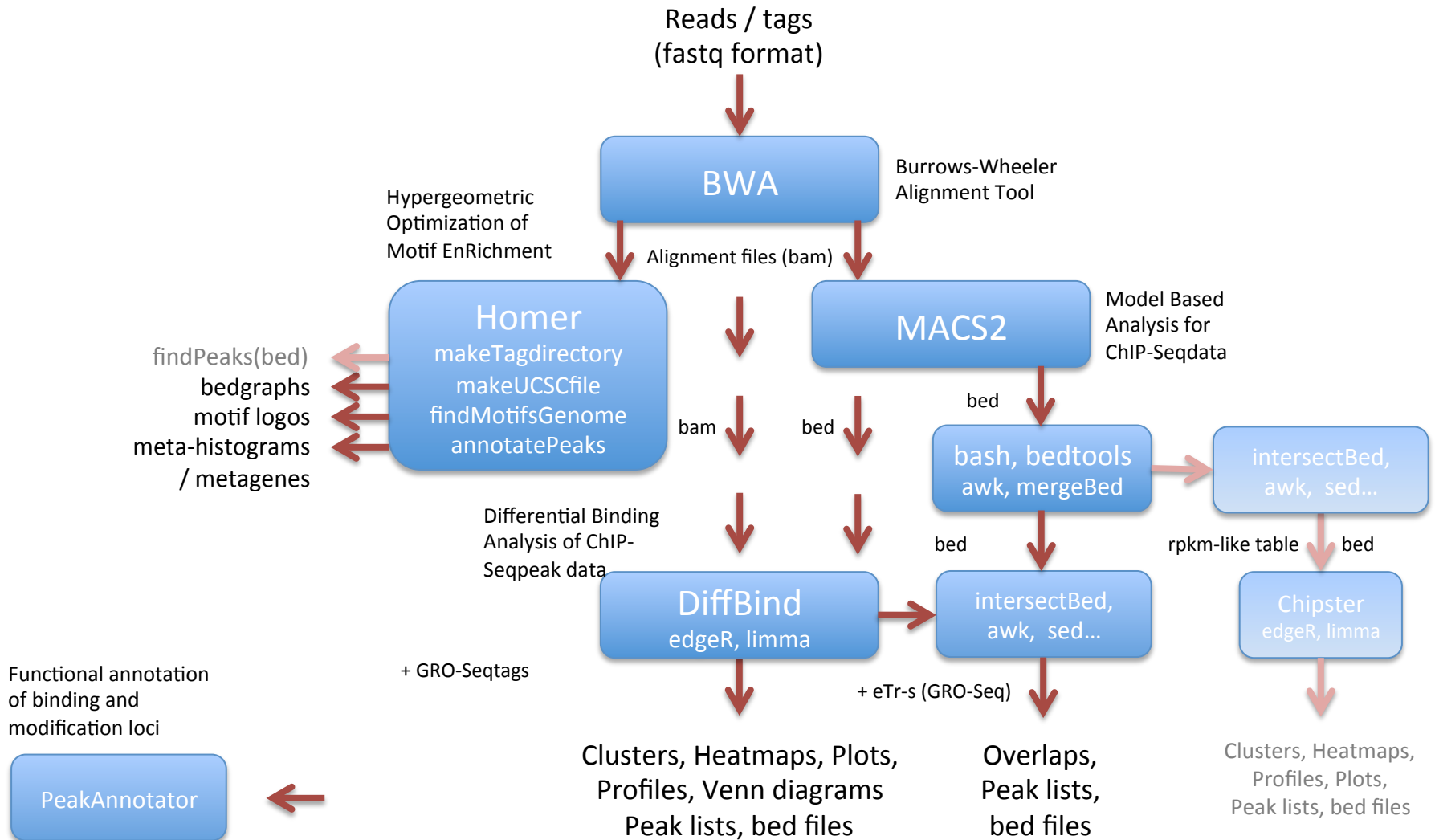
ENCODE data analysis helps to describe the various types of regulatory wiring implicit in the genome

Read Thread

The screenshot shows the Nature ENCODE explorer interface. It features a network visualization of regulatory wiring, with nodes representing genes and edges representing regulatory interactions. The interface includes a search bar, navigation links, and a detailed view of a specific thread titled "Characterization of network topology".



# ChIP-Seq análisis pipeline



# Big Data kihívások a genomikában

- Elsődleges szekvenálási adatok tárolása és feldolgozása
- Szekvencia archívumok (nincs adatmentés, párhuzamos adatbázisok vannak)
- Adatok feldolgozása 10-100x-os tárterületet igényel
- Szekvenálás berobbanása a humán diagnosztikába
- Adatok vizualizálása
- Speciális hardverigény (óriási I/O igény, esetenként nagy memóriaigény)



# Genomika és Big Data a kínai BGI-ban



**Table 1. BGI Computing Platform**

Location	Cores	Memory	Storage	T flops
Shenzhen	11,000	22TB	8.88PB	117T flops
Hong Kong	7,776	9.7TB	8.515PB	83T flops
Beijing	300	500GB	162TB	1.5T flops
Wuhan	1760	500GB	1PB	4T flops
US-CHOPS	300	800GB	500TB	2T flops
US-UC Davis	300	800GB	500TB	2T flops
EUR-Denmark	600	1.6TB	1PB	4T flops
<b>Total</b>	<b>22,036</b>	<b>35.9TB</b>	<b>20.6PB</b>	<b>214T flops</b>

## Project

- Million Plant & Animal Genomes Project
- Million Human Genomes Project
- Million Micro-ecosystem Genomes Project
- Article Published

## Industry

- BGI Tech
- Clinical Analysis Center
- Reproductive Health Center
- Agriculture Service

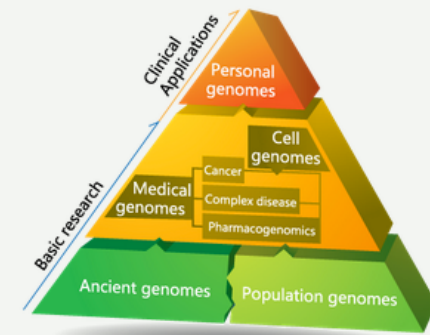


2013 Copyright BGI All Rights Reserved

Contact Us Investors Relationship Site map



## Project Overview



# Gímszarvas genom szekvenálása

- Nincsen referencia szekvencia (amihez illeszteni lehetne)
- 4 lane HiSeq 2000-es paired-end és 2 lane mate pair szekvenálás (~4 millió forint)
- 1.1 milliárd  $2 \times 100$  bp hosszú read  $\rightarrow$  220 milliárd bázispár  $\rightarrow$  ~80x lefedettség (elméletileg a genom minden pontjára 80 leolvasás jut)
- denovo genomösszerakás a pécsi szuperszámítógépen





# NIF szuperszámítógépek, Pécs

- **SGI UltraViolet1000** típusú, SMP (ccNUMA)
- **Intel Xeon X7542**(Nehalem EX) típusú 6 magos processzorok találhatóak, összesen **1152 mag (96 gép)**
- **6 Tbyte memória**- **Bármely processz egyben látja, meg tudja címezni!**
- **0,5 Pbyte diszkrendszer**
- ▶ **Suse enterprise operációs rendszer,**
- ▶ **SGE feladatütemező**





# denovo szekvencia összerakás Big Data vonzatai

- 772Gb  
memória
- 2.8Tb diszk
- Az összerakás  
3 hónapja  
folyik

```
[vserver] szarvas (0)$ tail -f ./logs/RunAllPathsLG_build1_8.log
sze szept 11 21:26:25 2013: building extenders_q_jump, memory usage = 61,544,034,304
cs szept 12 07:29:14 2013: now memory usage = 149,702,168,576
cs szept 12 07:29:52 2013: found extenders, memory usage = 45,580,173,312
cs szept 12 07:55:16 2013: reloaded extending data, memory usage = 298,075,394,048
cs szept 12 08:45:14 2013: join data created, memory usage = 17,702,809,600
cs szept 12 08:45:14 2013: considering 11157628 joins
cs szept 12 08:45:15 2013: forking 80 cottage workers to do 11157628 joins.
cs szept 12 08:45:58 2013: farming out work
Farming out 11157628 jobs to sub-processes.
```

```
[vserver] top (0)$ cat ebarta.top.2013-08-17_22:42:52.out
top - 22:42:55 up 5 days, 6:45, 0 users, load average: 607.30, 649.18, 642.81
Tasks: 20445 total, 502 running, 19943 sleeping, 0 stopped, 0 zombie
Cpu(s): 41.7%us, 4.4%sy, 0.0%ni, 53.9%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Mem: 6187177M total, 909152M used, 5278025M free, 388M buffers
Swap: 285945M total, 2048M used, 283897M free, 116378M cached
```

PID	USER	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TIME+	COMMAND
1173780	ebarta	20	0	772g	713g	2068	S	663	11.8	13745:16	exe
178350	ebarta	20	0	24032	15m	780	R	66	0.0	0:02.74	top
178346	ebarta	20	0	12628	1564	1232	S	0	0.0	0:00.10	sh
574656	ebarta	20	0	13016	1268	1264	S	0	0.0	0:00.02	57390
574662	ebarta	20	0	24952	2232	2228	S	0	0.0	0:00.16	exe
574663	ebarta	20	0	5632	636	632	S	0	0.0	0:00.00	tee
574671	ebarta	20	0	10388	1440	856	S	0	0.0	0:00.11	make
1173779	ebarta	20	0	12632	1228	1224	S	0	0.0	0:00.00	sh
1173783	ebarta	20	0	22220	1656	1532	S	0	0.0	15:46.61	MemMonitor
1173784	ebarta	20	0	5632	664	632	S	0	0.0	0:00.03	tee

```
[vserver] top (0)$ █
```

```
[vserver] ebarta (0)$ uname -a
Linux vserver 2.6.32.59-0.7-default #1 SMP 2012-07-13 15:50:56 +0200 x86_64 x86_64 x86_64 GNU/Linux
[vserver] ebarta (0)$ pwd
/scratch/ebarta
[vserver] ebarta (0)$ du -hs szarvas/
2,8T szarvas/
[vserver] ebarta (0)$ █
```

# Összefoglalás

- A biológia, molekuláris biológia mindig is élen járt a számítástechnika alkalmazásában
- A genomika a közeljövőben az egyik legnagyobb Big Data generáló tudomány lesz
- Megoldás: Infrastruktúra fejlesztés (ELIXIR)



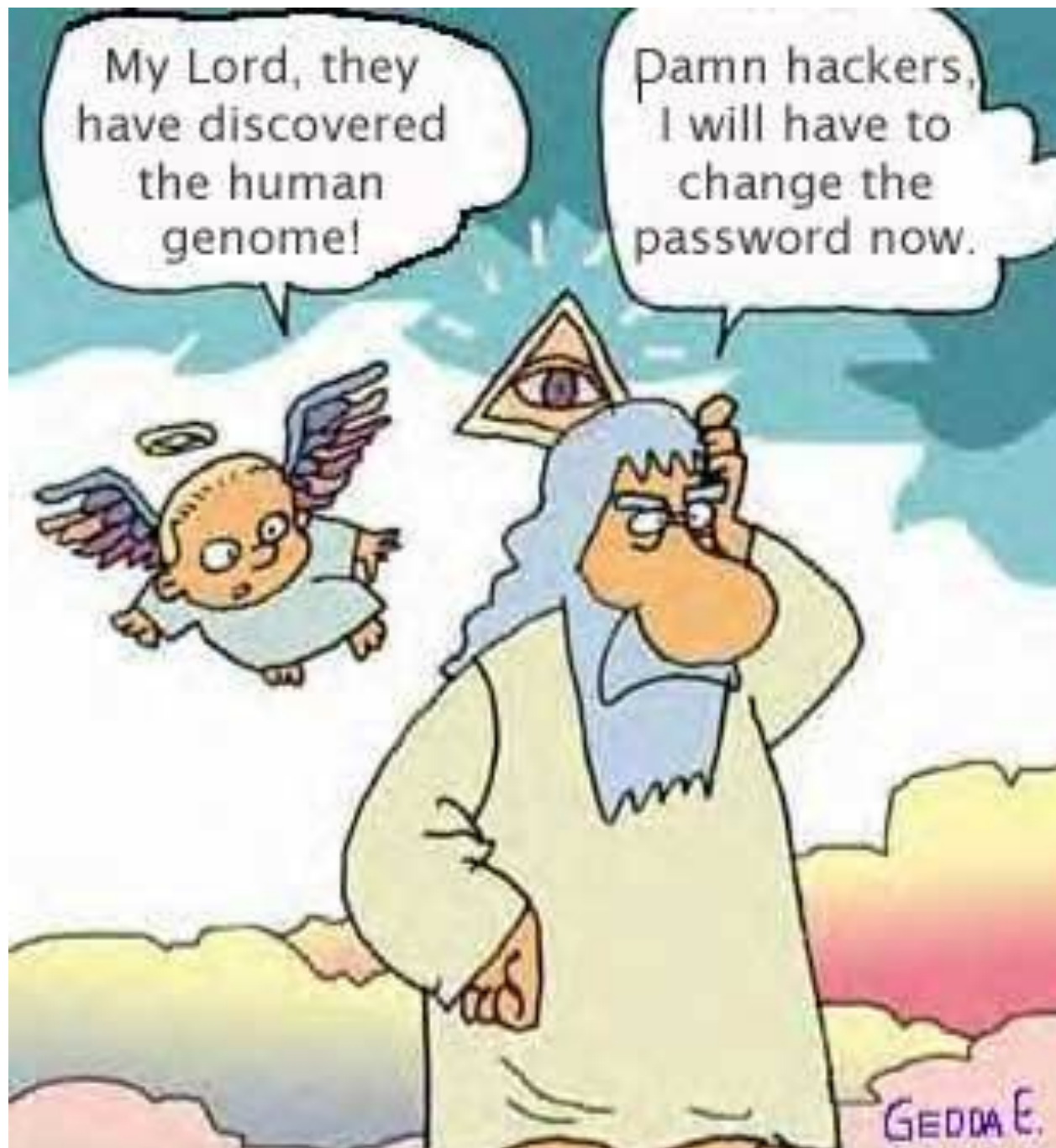
[about](#) [for funders](#) [for researchers](#) [for industry](#) [events](#) [news & media](#) [members area](#)

## Welcome to ELIXIR

*"ELIXIR unites Europe's leading life science organisations in managing and safeguarding the massive amounts of data being generated every day by publicly funded research. It is a pan-European research infrastructure for biological information. ELIXIR will provide the facilities necessary for life science researchers - from bench biologists to cheminformaticians - to make the most of our rapidly growing store of information about living systems, which is the foundation on which our understanding of life is built."*



Dr Niklas Blomberg  
ELIXIR Director



**Köszönöm a  
figyelmet!**