



BIST Founding Conference

Growing together, advancing science



Barcelona Institute of
Science and Technology

BIG DATA IN BIOINFORMATICS

MARCH 31, 2017

CCCB | Barcelona



BIST Founding Conference

Growing together, advancing science



Barcelona Institute of
Science and Technology

BIG DATA IS BIOINFORMATICS

MARCH 31, 2017

CCCB | Barcelona

BiG data is bioinformatics

- Heterogeneous data
 - numerical
 - non-numerical
 - Structures at different levels (from molecules to organisms)—images
 - Sequences
 - Longitudinal/dynamic—movies
 - Multi-dimensional
- Collected at multiple sites
 - Produced by individual small labs to large international consortiums
- Shared through the internet
 - Real time access
- Need of integrative analysis



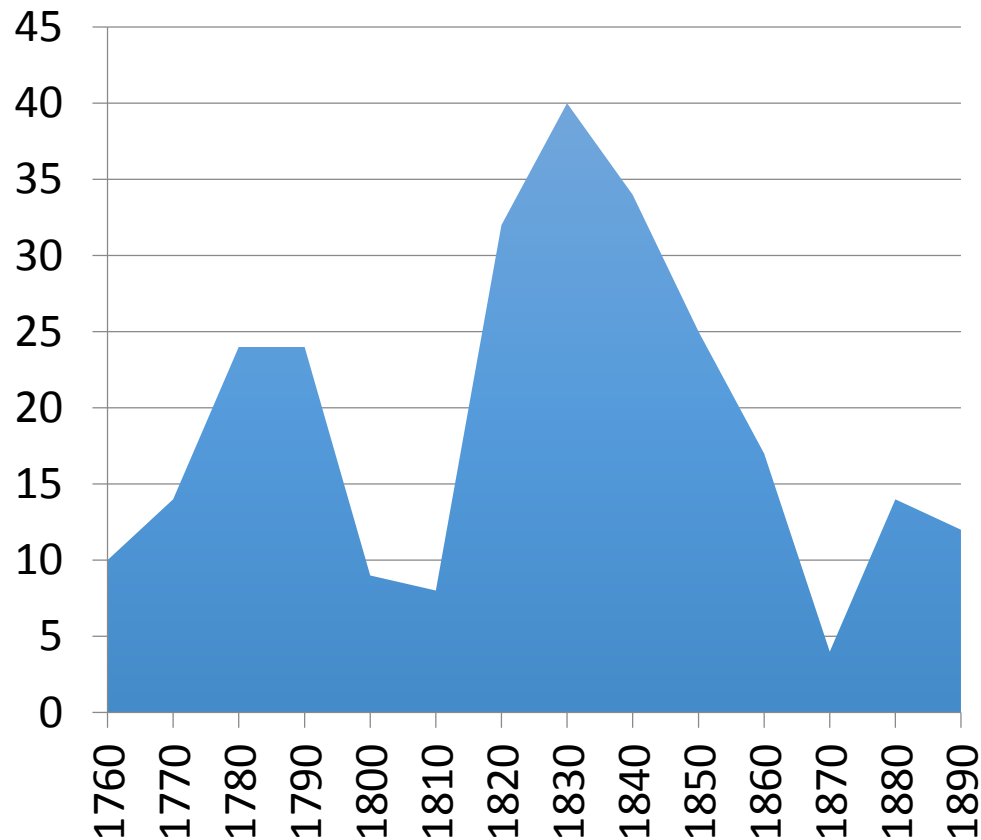
Founding Conference

Google search: X-informatics (june 4,2015)

bio informatics	24,500,000
chemo informatics	275,000
astro informatics	27,800
neuro informatics	331,000
socio informatics	14,100
geo informatics	548,000
meteo informatics	146
econo informatics	2,010
eco informatics	92,800
physico informatics	5,390

MARCH 31, 2017

CCCB | Barcelona



Number of scientific expeditions

■ # of commissioned years

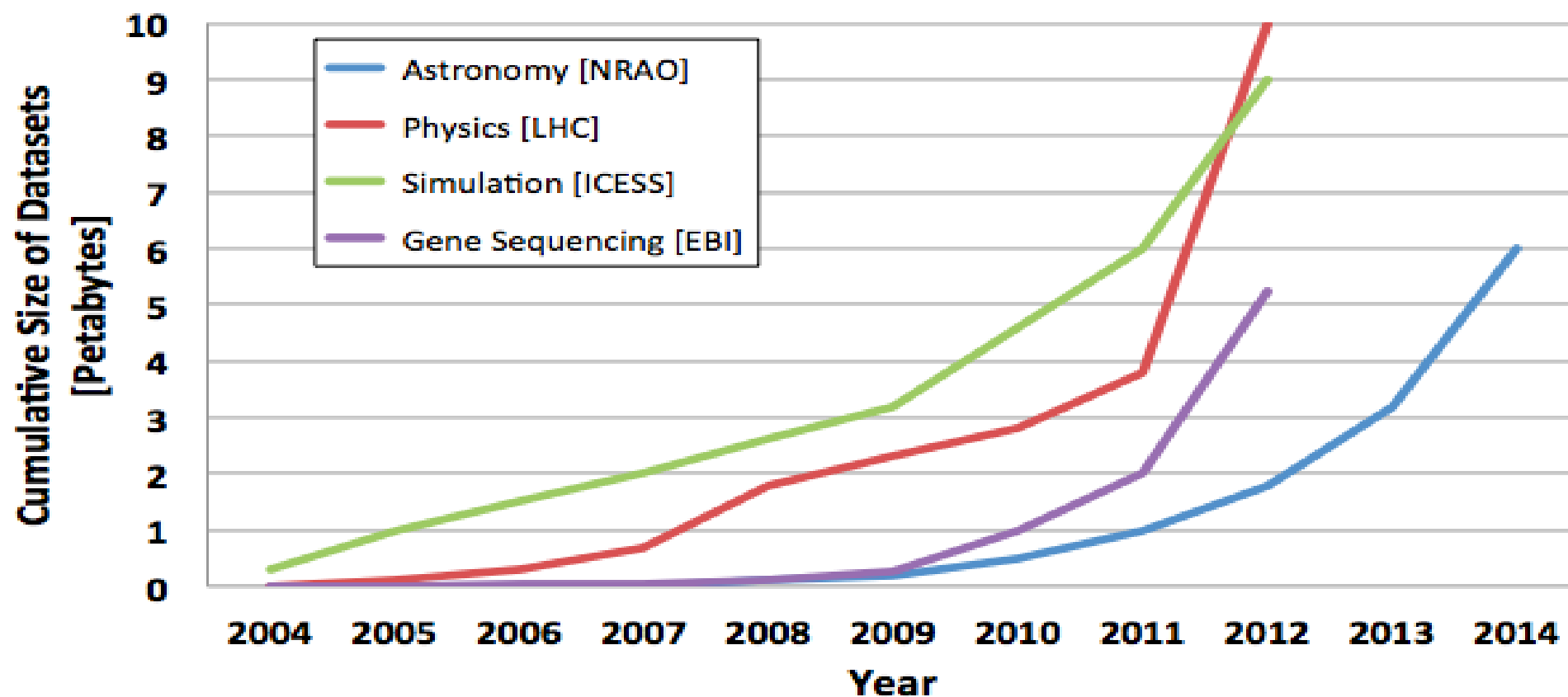


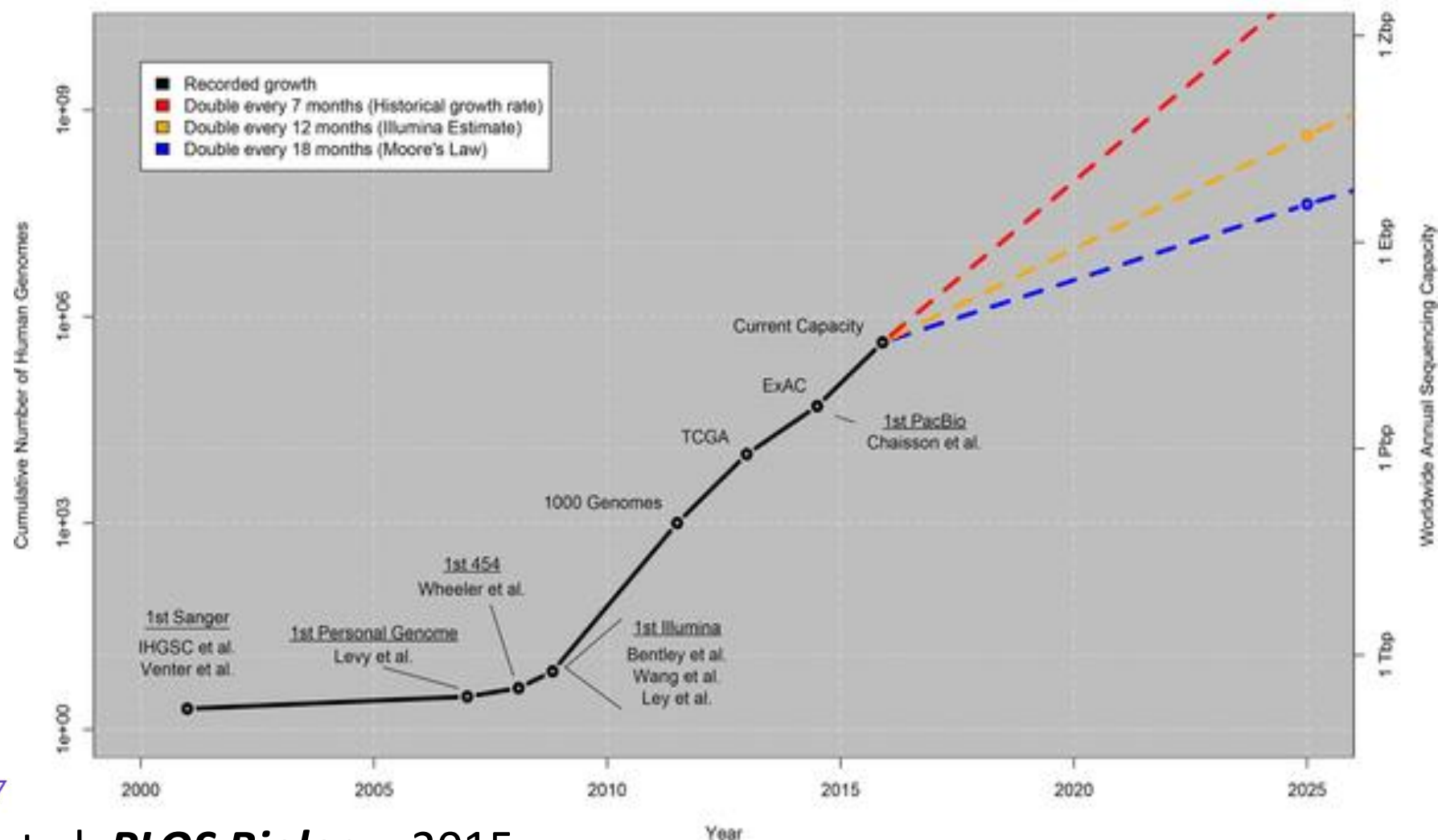
MARCH 31, 2017



CCCB | Barcelona

Cedric Notredame, CRG





MARCH 31, 2017

CCCB | Barcelona

Big Data: Astronomical or Genomical?

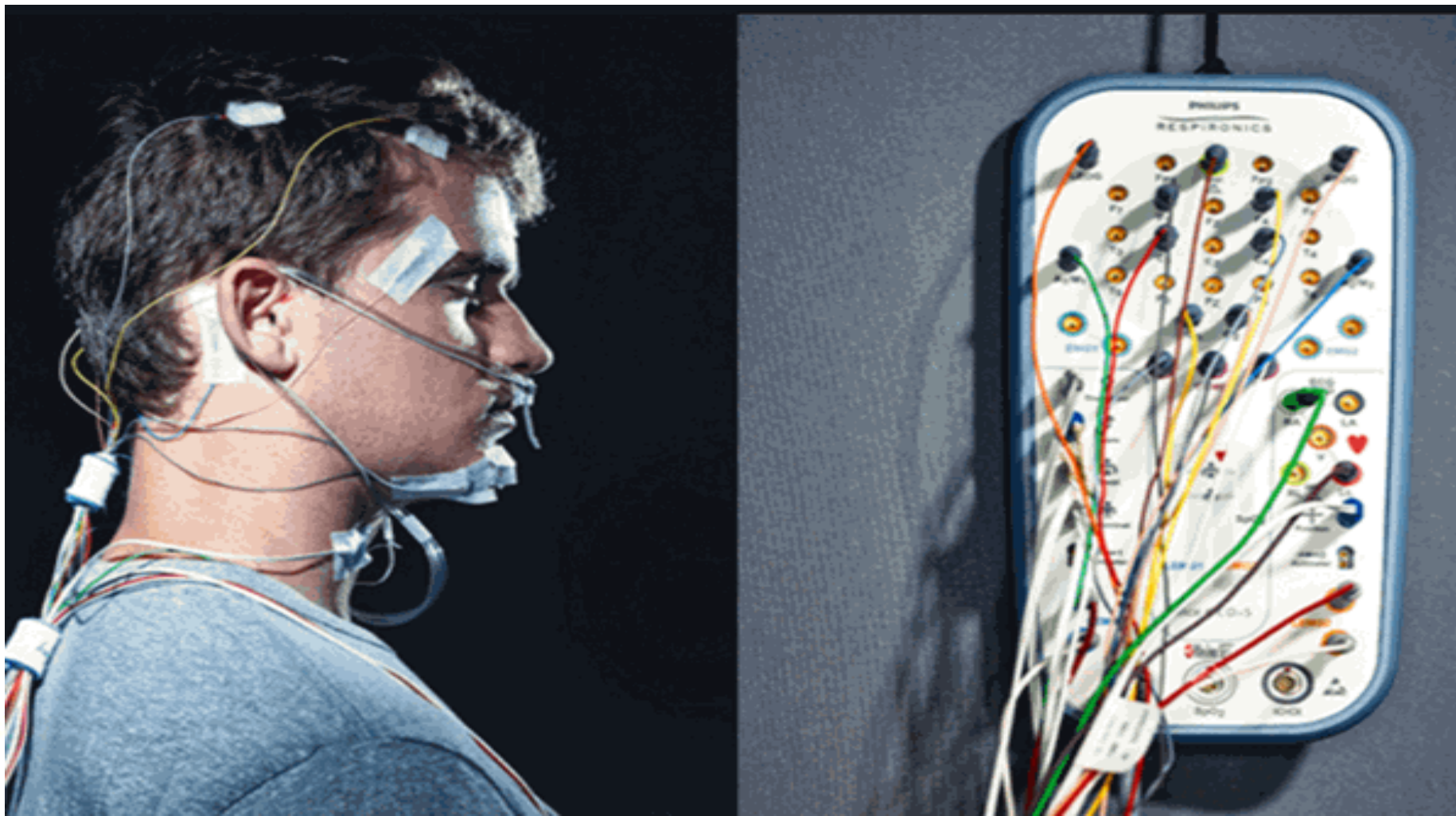
<u>Data Phase</u>	<u>Astronomy</u>	<u>Twitter</u>	<u>YouTube</u>	<u>Genomics</u>
Acquisition	25 zetta-bytes/year	0.5–15 billion tweets/year	500–900 million hours/year	1 zetta-bases/year
Storage	1 EB/year	1–17 PB/year	1–2 EB/year	2–40 EB/year
Analysis	In situ data reduction	Topic and sentiment mining	Limited requirements	Heterogeneous data and analysis
	Real-time processing	Metadata analysis		Variant calling, ~2 trillion central processing unit (CPU) hours
	Massive volumes			All-pairs genome alignments, ~10,000 trillion CPU hours
Distribution	Dedicated lines from antennae to server (600 TB/s)	Small units of distribution	Major component of modern user's bandwidth (10 MB/s)	Many small (10 MB/s) and fewer massive (10 TB/s) data movement

MARCH 31, 2017

CCCB | Barcelona

Table 1. Four domains of Big Data in 2025

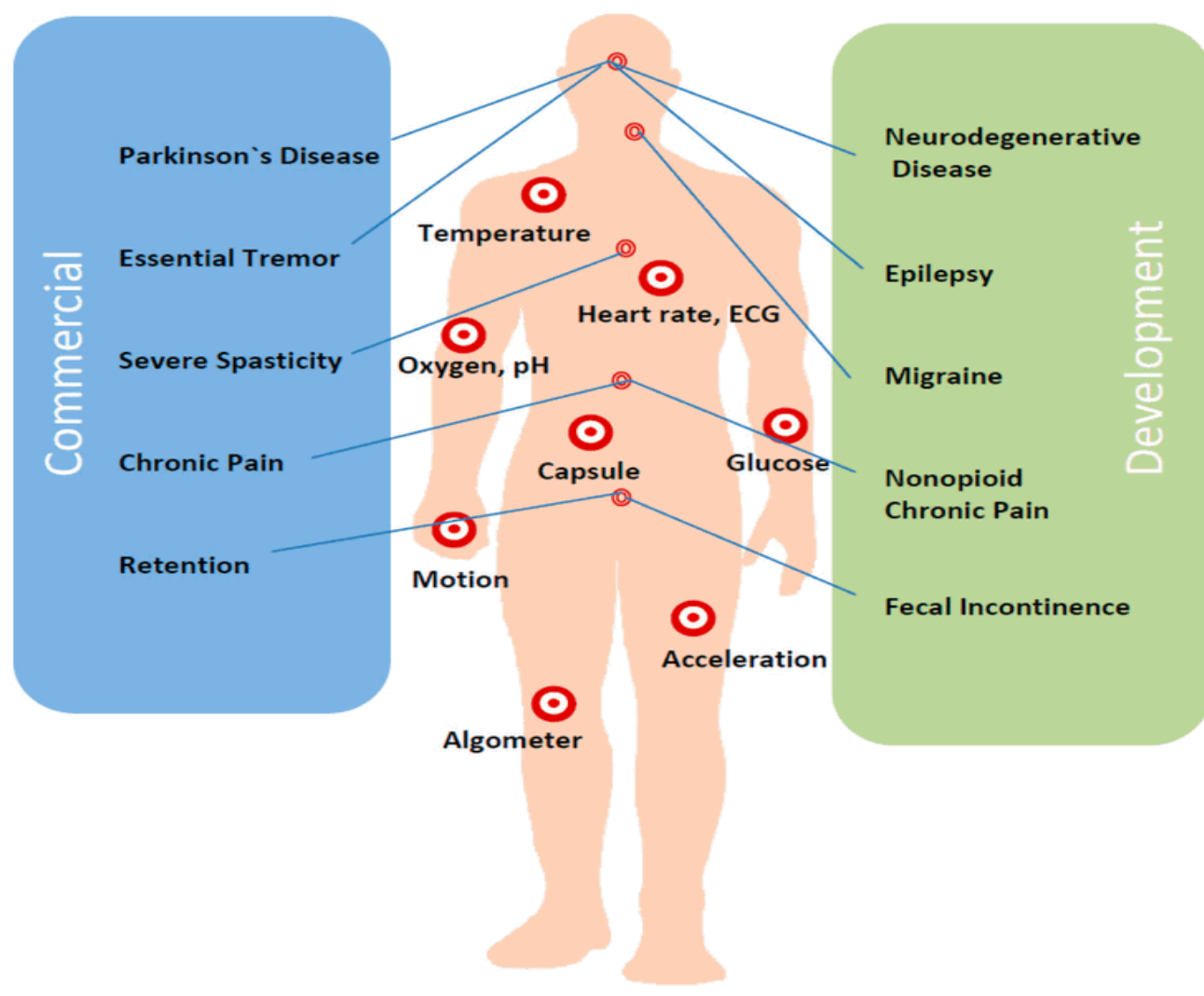
We are the Big Data



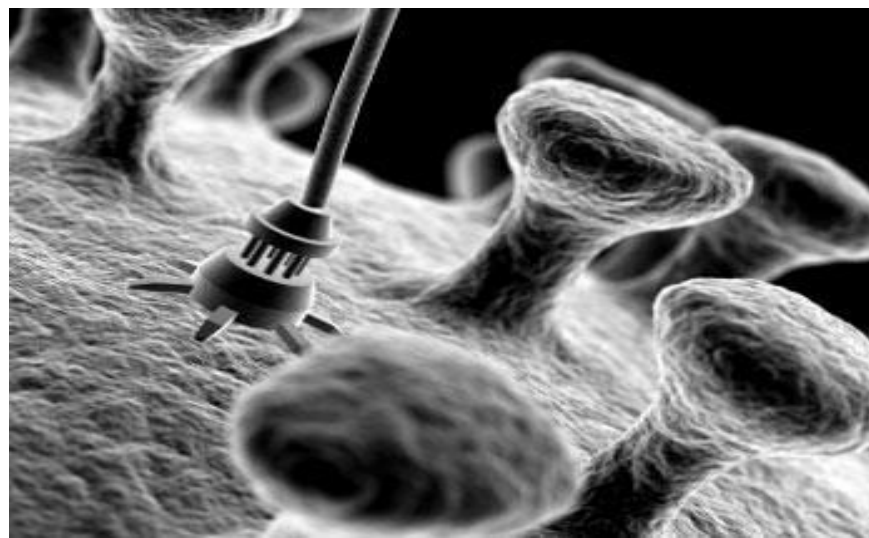
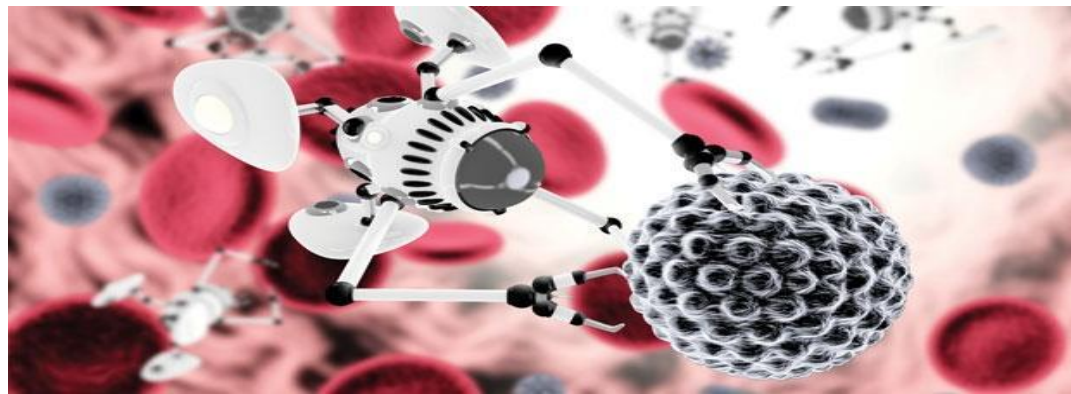
Wearable medical devices

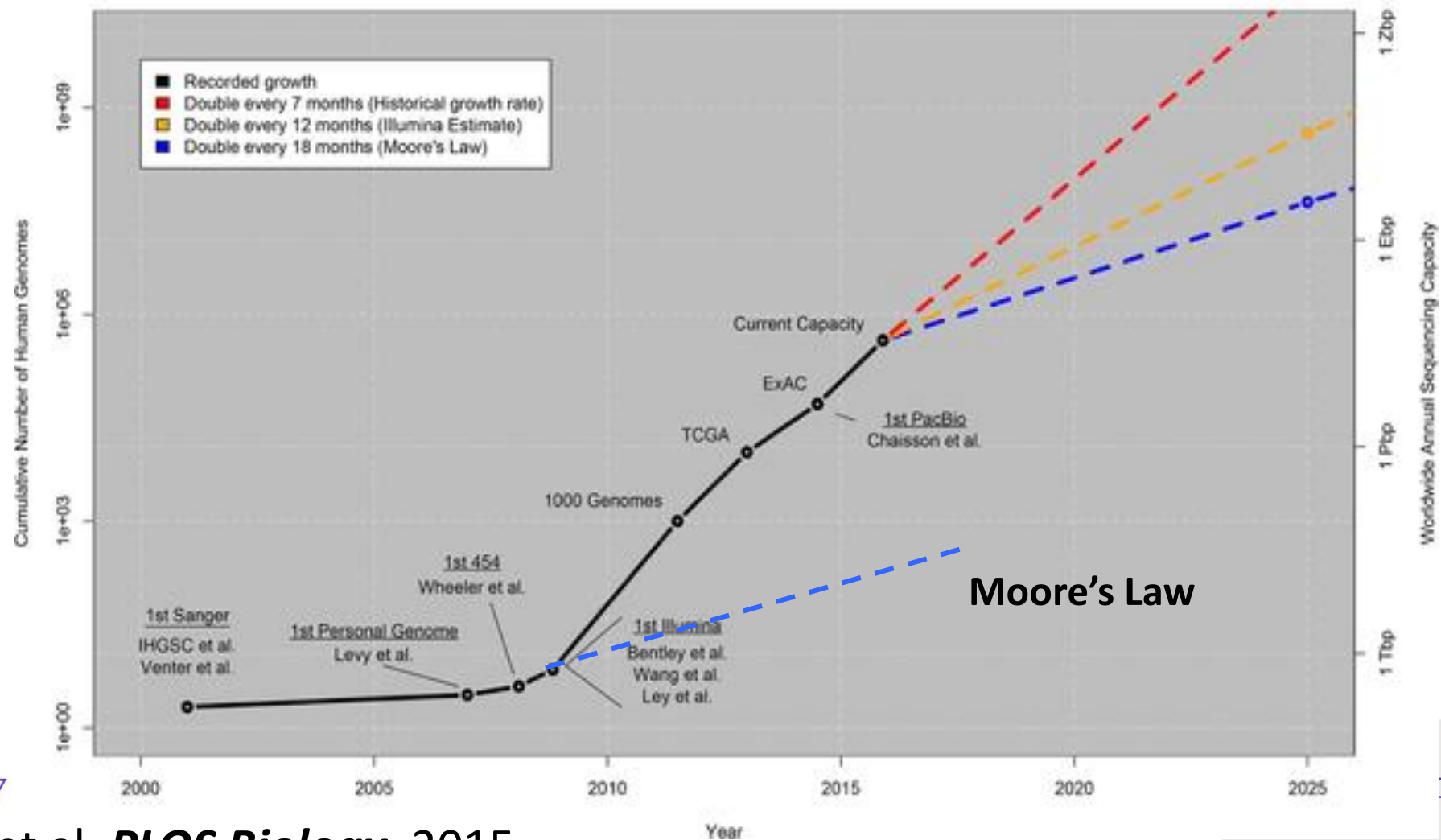


Implantable wearable devices



Nanowearables





MARCH 31, 2017

CCCB | Barcelona

Molecular Computation of Solutions to Combinatorial Problems

Leonard M. Adleman

The tools of molecular biology were used to solve an instance of the directed Hamiltonian path problem. A small graph was encoded in molecules of DNA, and the “operations” of the computation were performed with standard protocols and enzymes. This experiment demonstrates the feasibility of carrying out computations at the molecular level.

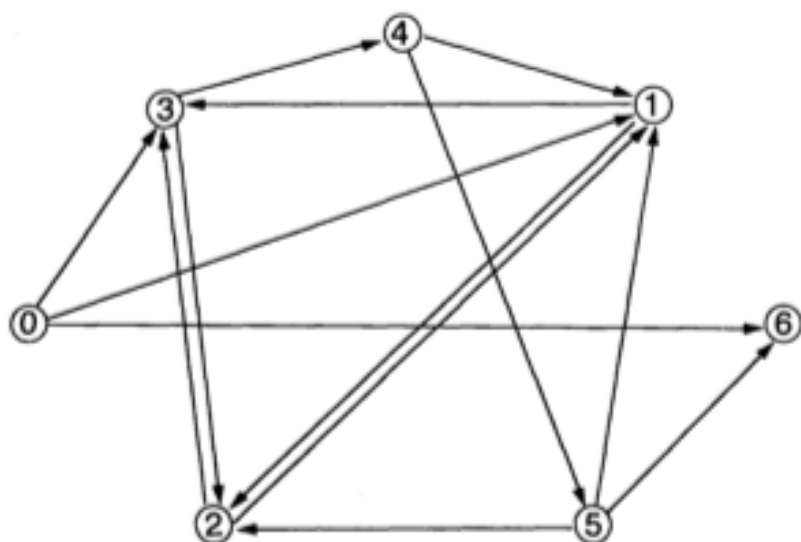
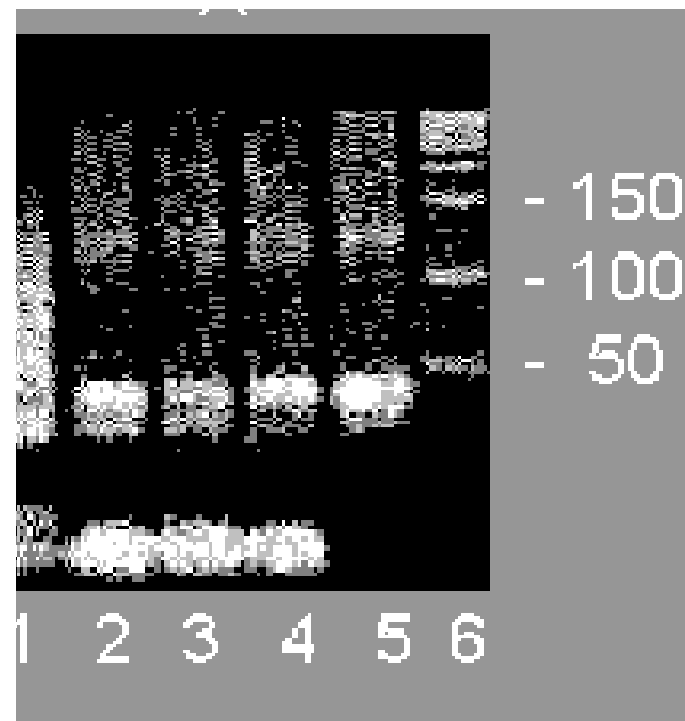


Fig. 1. Directed graph. When $v_{in} = 0$ and $v_{out} = 6$, a unique Hamiltonian path exists: $0 \rightarrow 1$, $1 \rightarrow 2$, $2 \rightarrow 3$, $3 \rightarrow 4$, $4 \rightarrow 5$, $5 \rightarrow 6$.

O_2 TATCGGATCGGTATATCCGA
 O_3 GCTATTCGAGCTTAAAGCTA
 O_4 GGCTAGGTACCAGCATGCTT
 $O_{2 \rightarrow 3}$ GTATATCCGAGCTATTCGAG
 $O_{3 \rightarrow 4}$ CTTAAAGCTAGGCTAGGTAC
 \bar{O}_3 CGATAAGCTCGAATTTTCGAT

$O_{2 \rightarrow 3}$ $O_{3 \rightarrow 4}$
 ↓
GTATATCCGAGCTATTCGAGCTTAAAGCTAGGCTAGGTAC
CGATAAGCTCGAATTTTCGAT
 \bar{O}_3

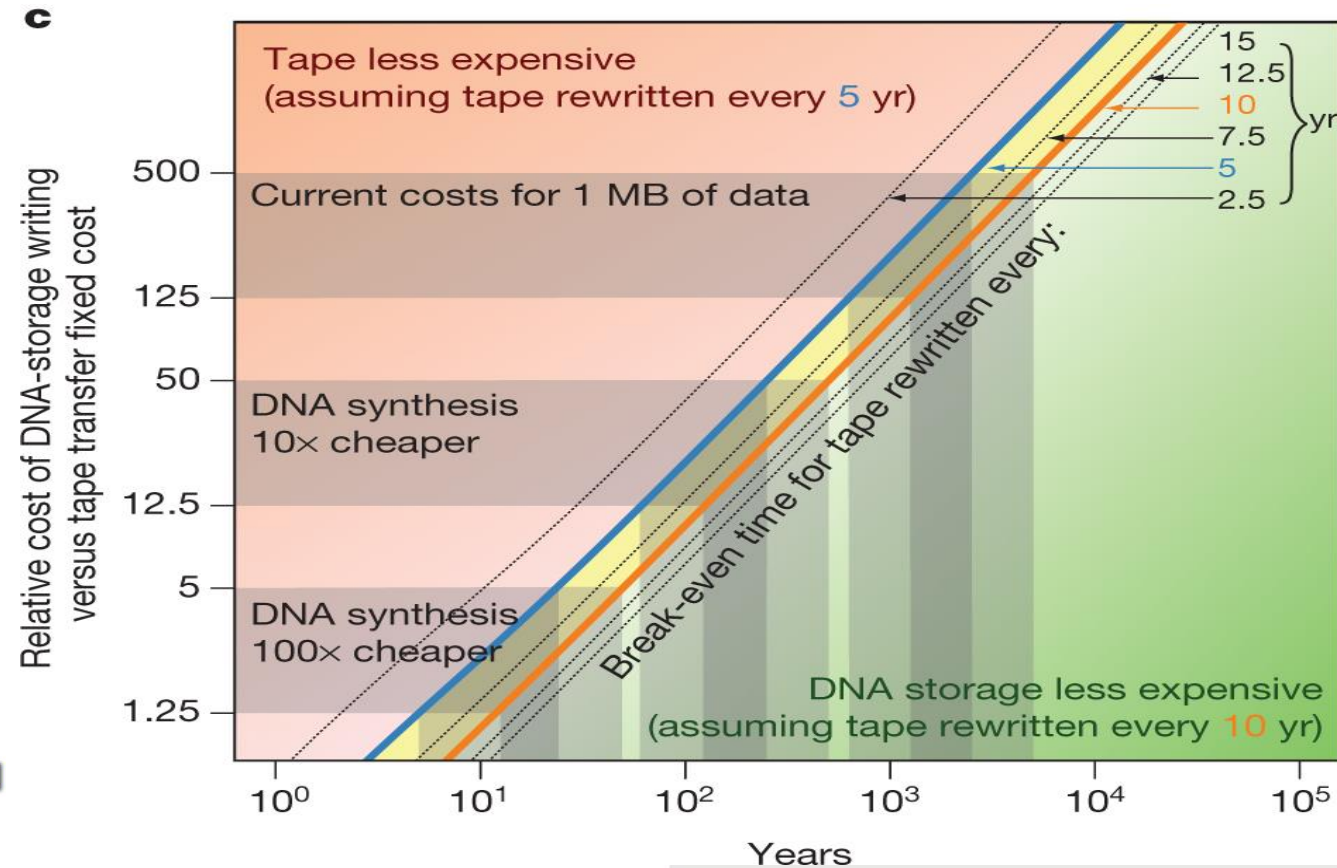
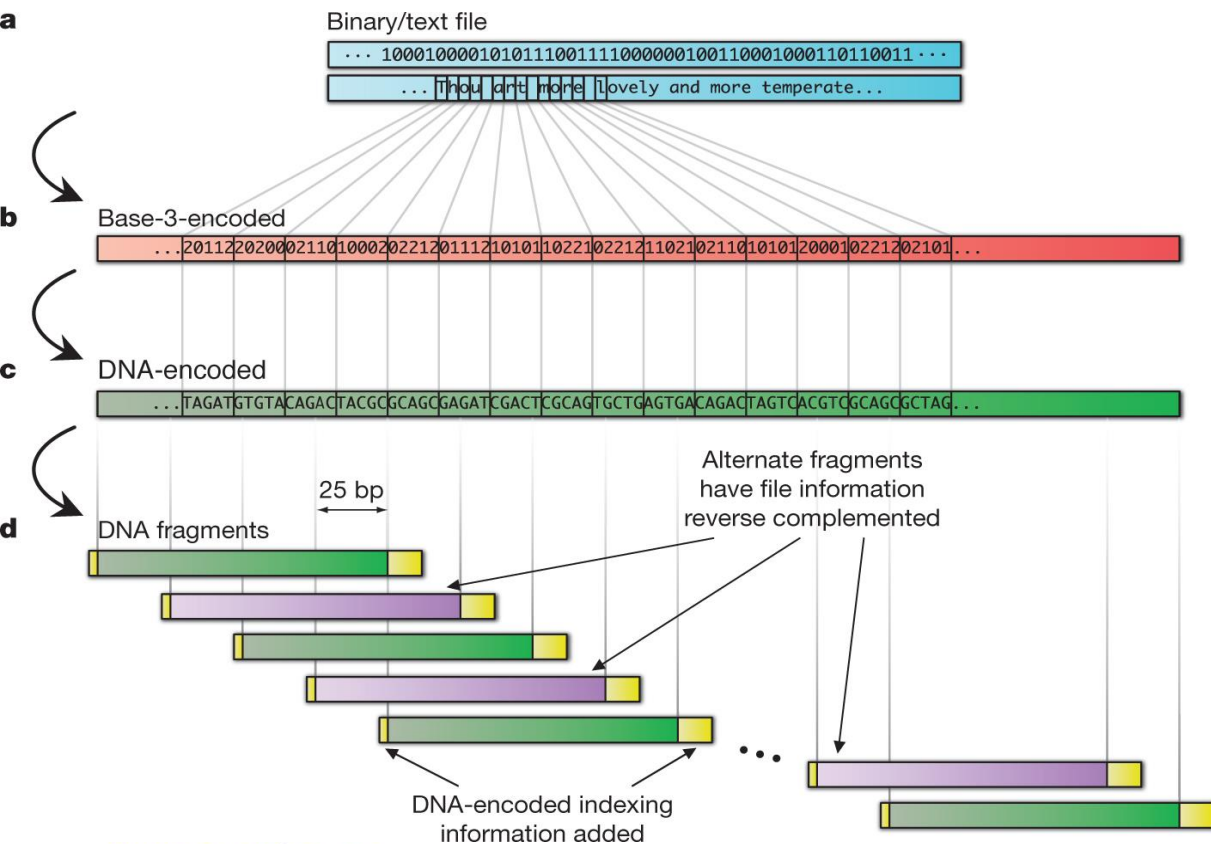
Fig. 2. Encoding a graph in DNA. For each vertex i in the graph, a random 20-mer oligonucleotide O_i is



Towards practical, high-capacity, low-maintenance information storage in synthesized DNA

Nick Goldman¹, Paul Bertone¹, Siyuan Chen², Christophe Dessimoz¹, Emily M. LeProust², Botond Sipos¹ & Ewan Birney¹

2 PB per 1 g DNA



“Goldman prediction”

- 2PB per 1g DNA (2×10^{15} bytes)
- Total world info (2013): 3ZBytes (3×10^{21} bytes)
 - Approximately 1.5×10^6 g (1,5 tonnes of DNA) to store all information
- Information doubling time: 2 years
- Mass of earth: 6×10^{27} g (google)
- $1.5 \times 10^6 \times 2^{x/2} \approx 6 \times 10^{27} \rightarrow x \approx 140$ years
- the mass of total info in the world stored in DNA exceeds the mass of the Earth in year 2157

REPORT

DNA STORAGE

DNA Fountain enables a robust and efficient storage architecture

215 PB per 1 g DNA

Yaniv Erlich^{1,2,3*} and Dina Zielinski¹

DNA is an attractive medium to store digital information. Here we report a storage strategy, called DNA Fountain, that is highly robust and approaches the information capacity per nucleotide. Using our approach, we stored a full computer operating system, movie, and other files with a total of 2.14×10^6 bytes in DNA oligonucleotides and perfectly retrieved the information from a sequencing coverage equivalent to a single tile of Illumina sequencing. We also tested a process that can allow 2.18×10^{15} retrievals using the original DNA sample and were able to perfectly decode the data. Finally, we explored the limit of our architecture in terms of bytes per molecule and obtained a perfect retrieval from a density of 215 petabytes per gram of DNA, orders of magnitude higher than previous reports.