

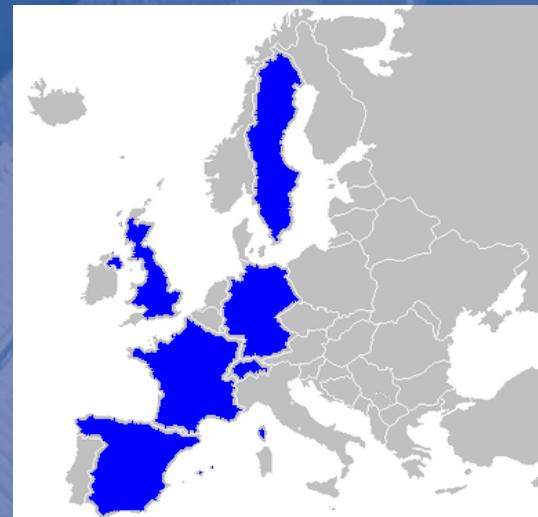


**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación

The HPC-Life Sciences scenario Byte vs Flop?

Modesto Orozco

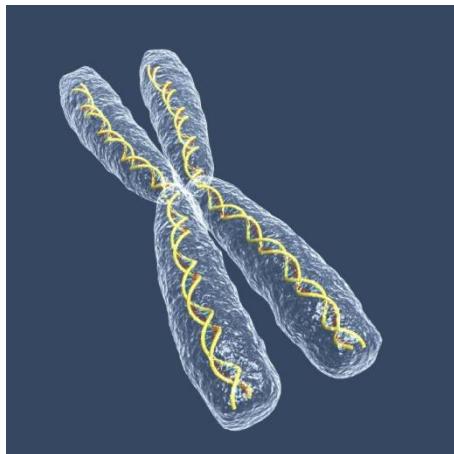
Henry Markram, Erik Lindahl, Paolo Carloni, Alfonso Valencia
Richard Lavery, Peter Coveney, Charles Laughton



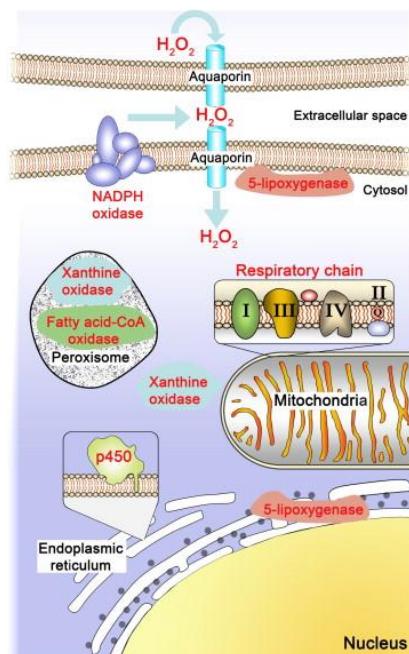
Exascale for Life Sciences applications



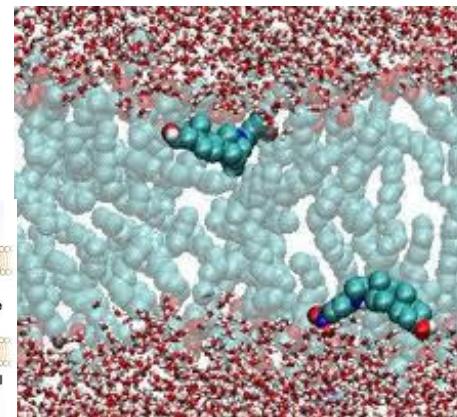
Omics



Systems Biology



Molecular Simulation



Organ Simulation



The dual nature of computers

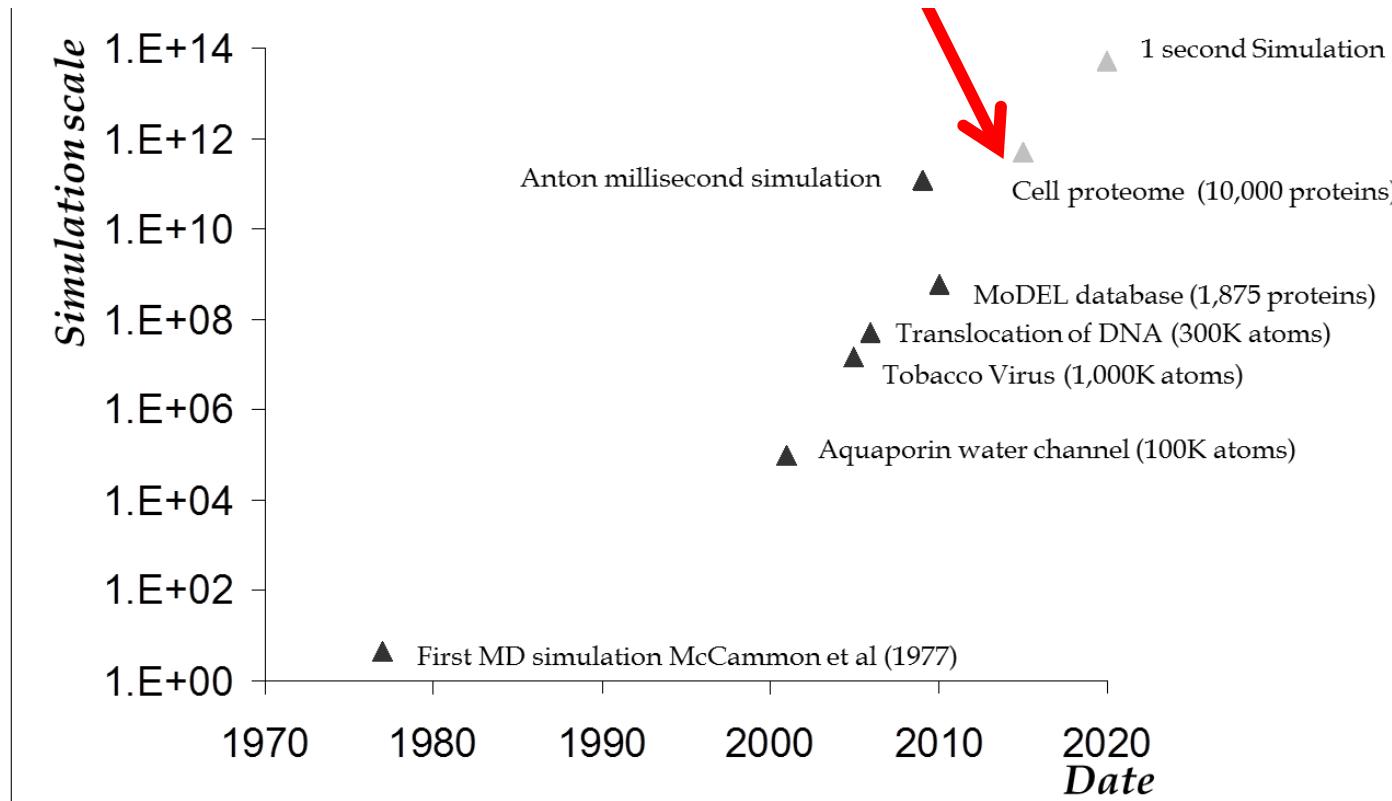
- ❑ Ordenador: a machine to manage data.
- ❑ Computador: a machine to do maths.



FLOPS and BYTES ARE EQUIALLY IMPORTANT
IN BIOLOGY?

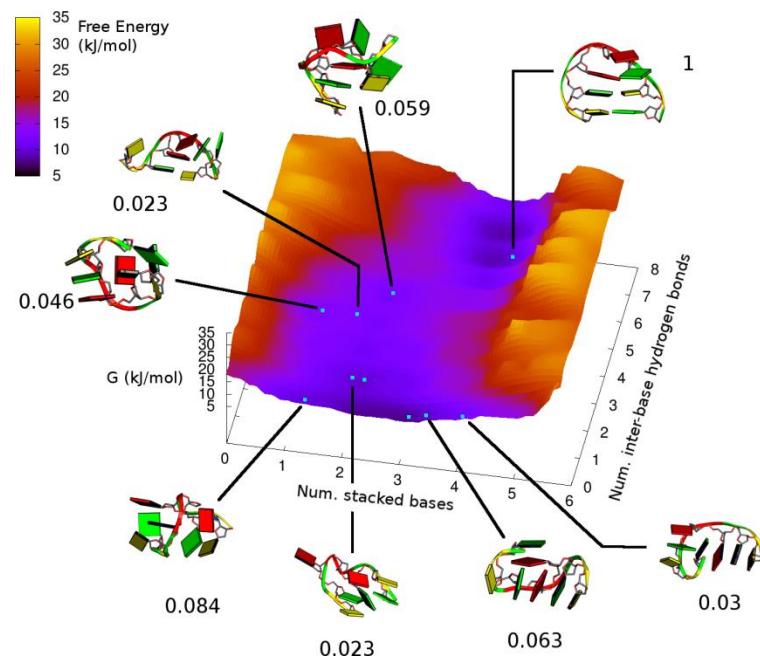
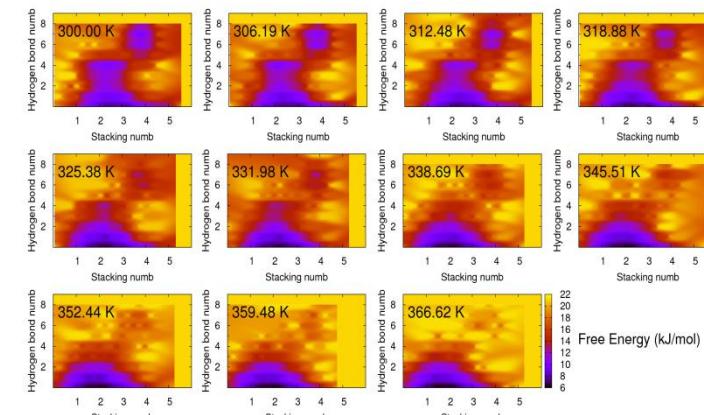
Evolution,...

Mature HIV-1 capsid (64M atoms) NAMD

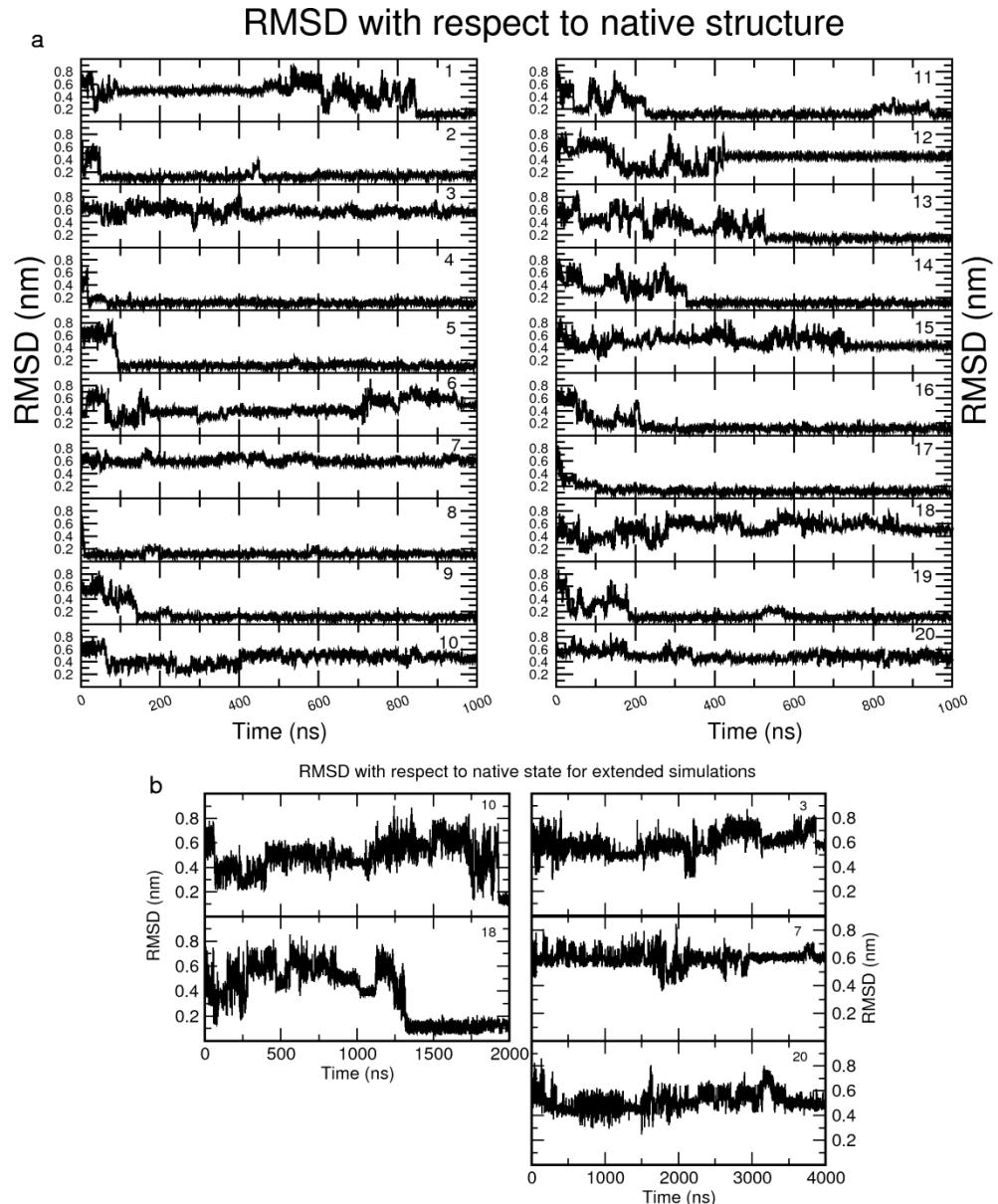


Simulation scale for long and large simulations is calculated multiplying simulation length (nanoseconds) with size (number of atoms).

Massive parallel simulations



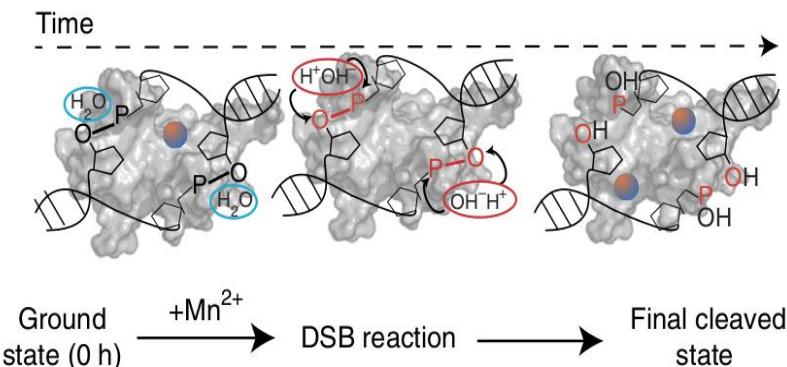
Portella & Orozco, Angew.Chem. 2010



Integrate experiments

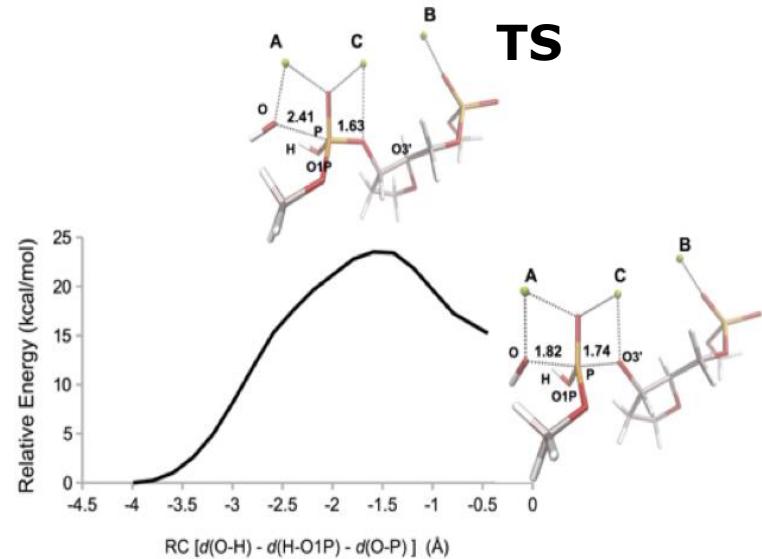
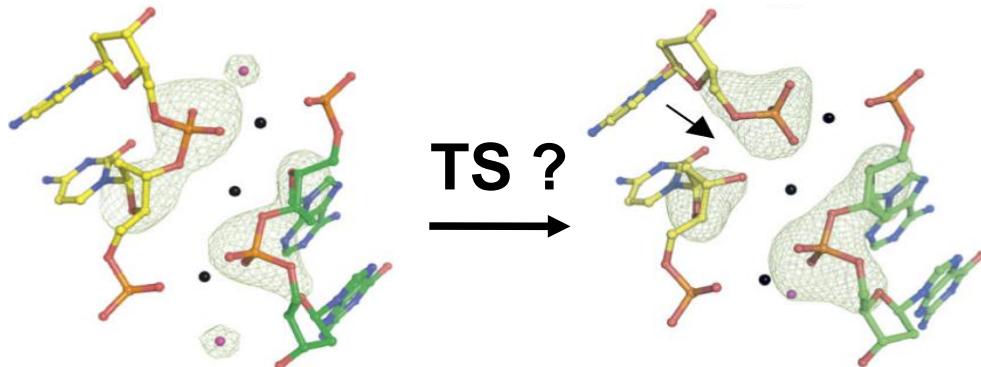


DNA double-strand break by the *I-Dmol* endonuclease



7 different catalytic intermediates were solved by X-ray, BUT reactive species cannot be trapped

We use a combination of **MD** and **QM** methods to fill the gaps and thus to complete the vision of the reaction, assign roles to the catalytic residues in the active site, estimate kinetic energy barriers ...

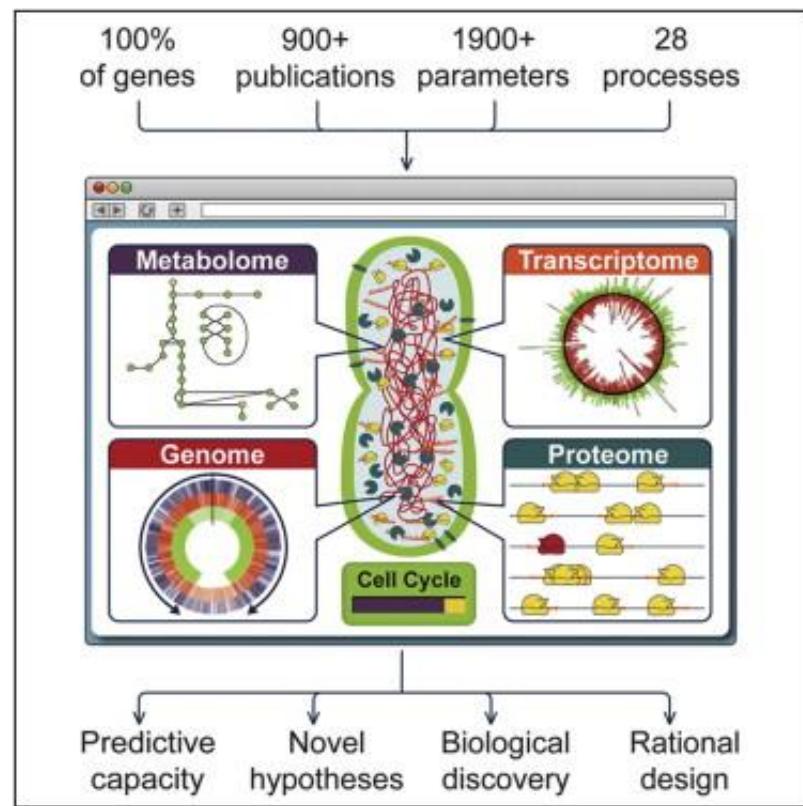


A reaction mechanism is proposed, involving the transient protonation of the PO_4 group by the catalytic water

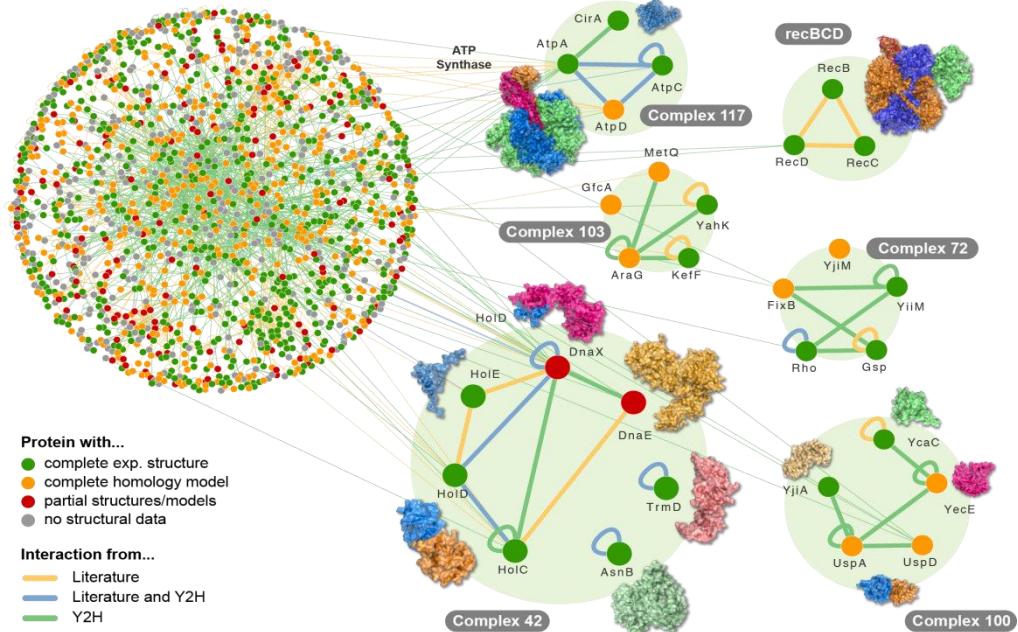
Whole-Cell simulation



- Whole-cell computational model of the life cycle of the human pathogen *Mycoplasma genitalium* from Genotype (Cell 2012)
- Entire organism (525 genes) modeled in terms of its molecular components
- Integration of multi-format data and very fragmented data
- experimental analysis directed by model predictions identified previously undetected kinetic parameters and biological functions
- 7 hours → 100 Gb of data

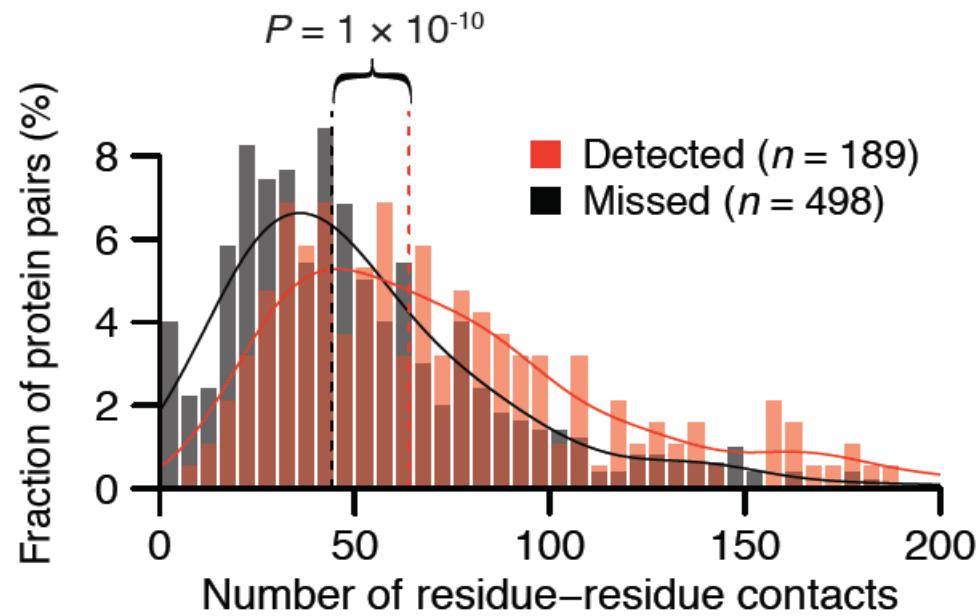
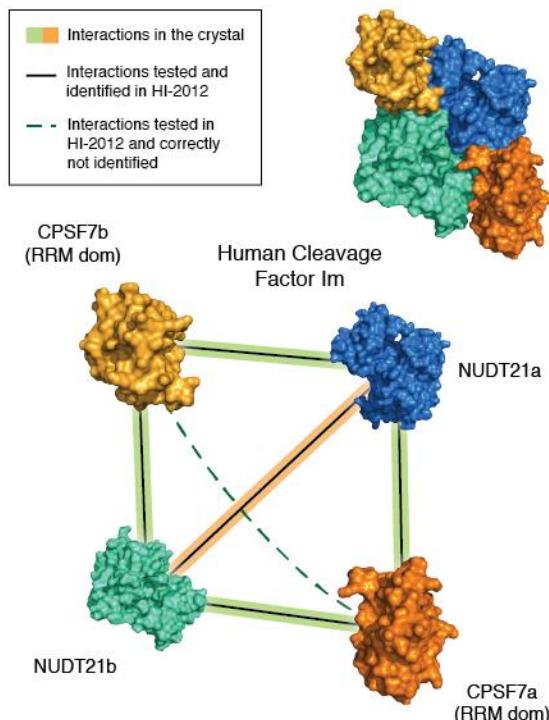


E. coli binary interactome network

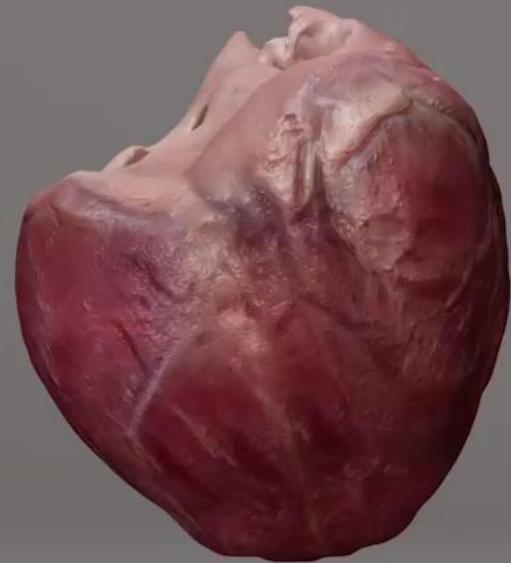


Human interactome network

Description	Tested		Identified		Not identified		P-value
	Direct	Indirect	TP	FP	FN	TN	
All interactions	687	56	189	2	498	54	1.01E-05



HEARTH SIMULATOR



BIOMECHANICS

J.M.Cela (BSC)

The dual nature of computers

- ❑ Ordenador: a machine to manage data.
- ❑ Computador: a machine to do maths.



FLOPS → BYTE but also BYTE→FLOPS

human genome project - Buscar con Google - Windows Internet Explorer

Favoritos IRB Barcelona Mail - Inbox ... Galería de Web Slice Sítios sugeridos

human genome project - Buscar con Google

Web Imágenes Videos Mapas Noticias Libros Correo Más ▾

Google

human genome project

Aproximadamente 2.080.000 resultados (0,26 segundos)

\$1,688 whole human exome
\$1,688/1st sample, \$1,999 for bulk HiSeq2000(2x100)&NimbleGen v2 cap
www.otogenetics.com

Artículos académicos para human genome project
express 1836 New goals for Implications

23andMe

elixir Data for Life

This article elsewhere

10 years ago, two fingers were enough to count the number of sequenced human genomes. Until last year, the fingers on two hands were enough. Today, the rate of such accounting is cascading so fast it is hard to keep track. Nature attempted nevertheless: we asked more than 90 genome centres and labs to estimate the number of human genome sequences they have in the works. Although far from comprehensive, the tally indicates that at least 1,700 human genomes will have been completed by the end of this month, and that the total will rise to more than 30,000 by the end of 2011.

KEY

500 genomes by the end of October 2010

500 genomes by the end of 2011

Number of high-throughput sequencers in the region

ScienceDaily®

Your source for the latest research news

News Articles Videos Images Books

Health & Medicine Mind & Brain Plants & Animals Earth & Climate Space & Time Matter &

Science News

Human Genome's Breaking Points: Genetic Sequence of Large-Scale Differences Between Human Genomes

ScienceDaily (Feb. 2, 2011) — A detailed analysis

Elixir is predicting 280,000 humans sequenced in 2014
(based on the analysis of 960000 SNPs)

> 20000 human genomes x year

Are you ready to learn how to make a difference?

Blog Cite

23andMe

Greenland ice-melt map gets the cold shoulder
20 September 2011

Texas prepares to fight for stem cells
20 September 2011

tes disgraced doctor
2011

arts spark desire for a

what's this?

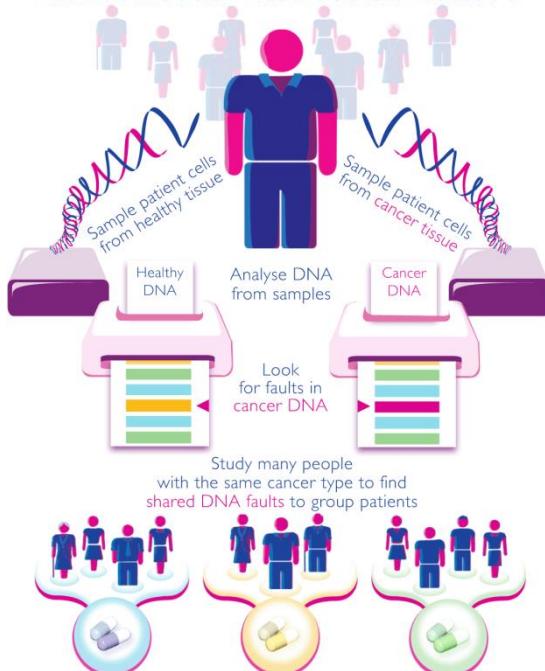
April 24th 2012

Genomic Grand Challenges



ICGL

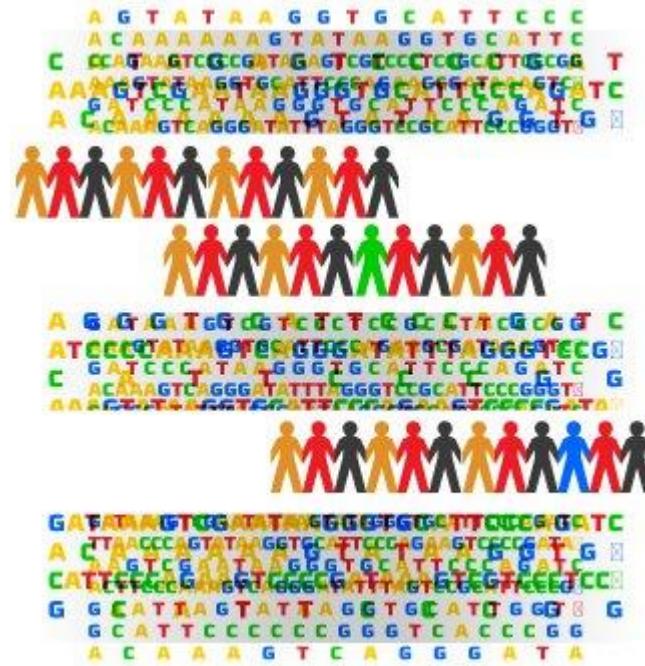
The International Cancer Genome Consortium



Within a decade, it will be possible to better tailor treatment

- Develop gene tests to routinely group patients
- Find new drugs that target specific groups better

1000G



Genomics-UK 100,000 genomes



Inbox (74) - modesto.or... Genomics England | 100,000 genomes www.genomicsengland.co.uk Modesto

Genomics england

Home About Genomics England The 100,000 Genomes Project GeCIP Library and resources News Contact us



Genomics England, with the consent of participants and the support of the public, is creating a lasting legacy for patients, the NHS and the UK economy through the sequencing of 100,000 genomes: [the 100,000 Genomes Project](#).

Genomics England was set up by the Department of Health to deliver the 100,000 Genomes Project. Initially the focus will be on rare disease, cancer and infectious disease. The project is currently in its pilot phase and will be completed by the end of 2017.

[Read more...](#)

News

NHS Genomic Medicine Centres announced for 100,000 Genomes Project

NHS England has announced eleven Genomic Medicine Centres that will lead the way in delivering the 100,000 Genomes Project. This marks the start of the main phase of the

Understanding genomics

What is genomics?

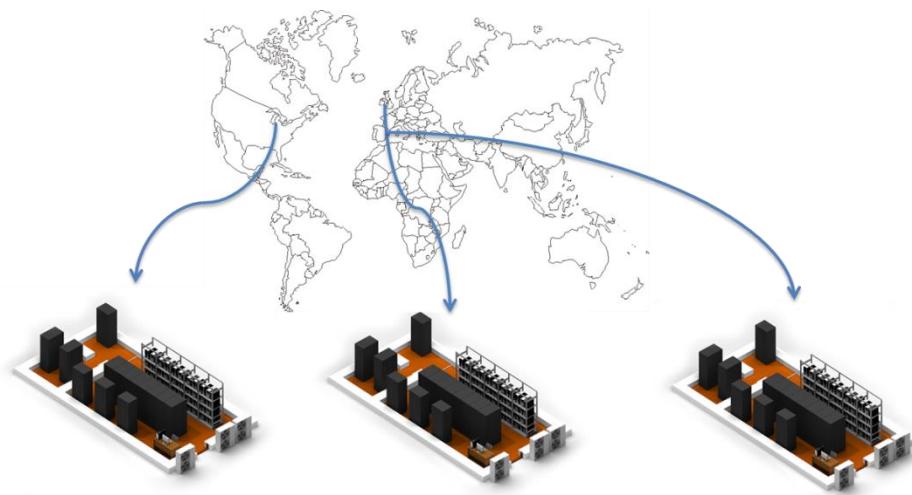
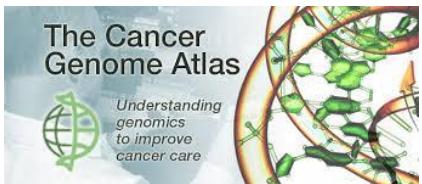
Our Head of Engagement, Vivienne Parry, explains more about genomics in this film courtesy of our partners at Health Education England.

Tweets

Genomics England @GenomicsEngland 7 Jan Check out this Storify about the day our new Genomic Medicine Centres were revealed: bit.ly/GMCStorify #Genomes100k #100KGPM Show Media

Department of Health 22 Dec

World-level genomic consortiums



University of Chicago
NIH Trusted Partner
Holds TCGA Data



EBI
Data Manager for 1KGP
Hosts EGA

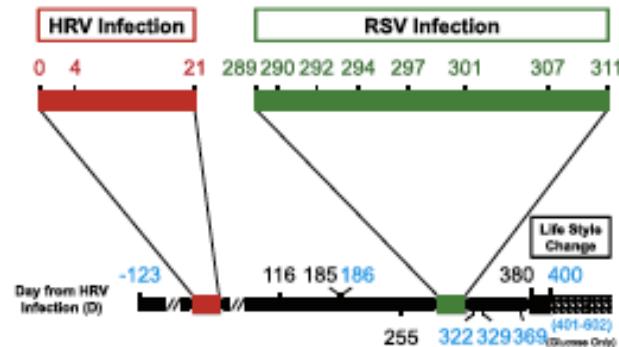


Barcelona Supercomputer Center
EGA Partner

- « European Genome-phenome Archive (EGA) repository will allow us to explore datasets from numerous genetic studies
- « Pan-cancer will produce rich data that will provide a major opportunity to develop an integrated picture of differences across tumor lineages

Personal Omics Profiling Reveals Dynamic Molecular and Medical Phenotypes

Rui Chen,^{1,11} George I. Mias,^{1,11} Jennifer Li-Pook-Than,^{1,11} Lihua Jiang,^{1,11} Hugo Y.K. Lam,^{1,12} Rong Chen,^{2,12} Elana Miriami,¹ Konrad J. Karczewski,¹ Manoj Hariharan,¹ Frederick E. Dewey,³ Yong Cheng,¹ Michael J. Clark,¹ Hogune Im,¹ Lukas Habegger,^{6,7} Suganthi Balasubramanian,^{6,7} Maeve O'Huallachain,¹ Joel T. Dudley,² Sara Hillenmeyer,¹ Rajini Haraksingh,¹ Donald Sharon,¹ Ghia Euskirchen,¹ Phil Lacroute,¹ Keith Bettinger,¹ Alan P. Boyle,¹ Maya Kasowski,¹ Fabian Grubert,¹ Scott Seki,² Marco Garcia,² Michelle Whirl-Carrillo,¹ Mercedes Gallardo,^{9,10} Maria A. Blasco,⁹ Peter L. Greenberg,⁴ Phyllis Snyder,¹ Teri E. Klein,¹ Russ B. Altman,^{1,5} Atul J. Butte,² Euan A. Ashley,³ Mark Gerstein,^{6,7,8} Kari C. Nadeau,² Hua Tang,¹ and Michael Snyder^{1,*}

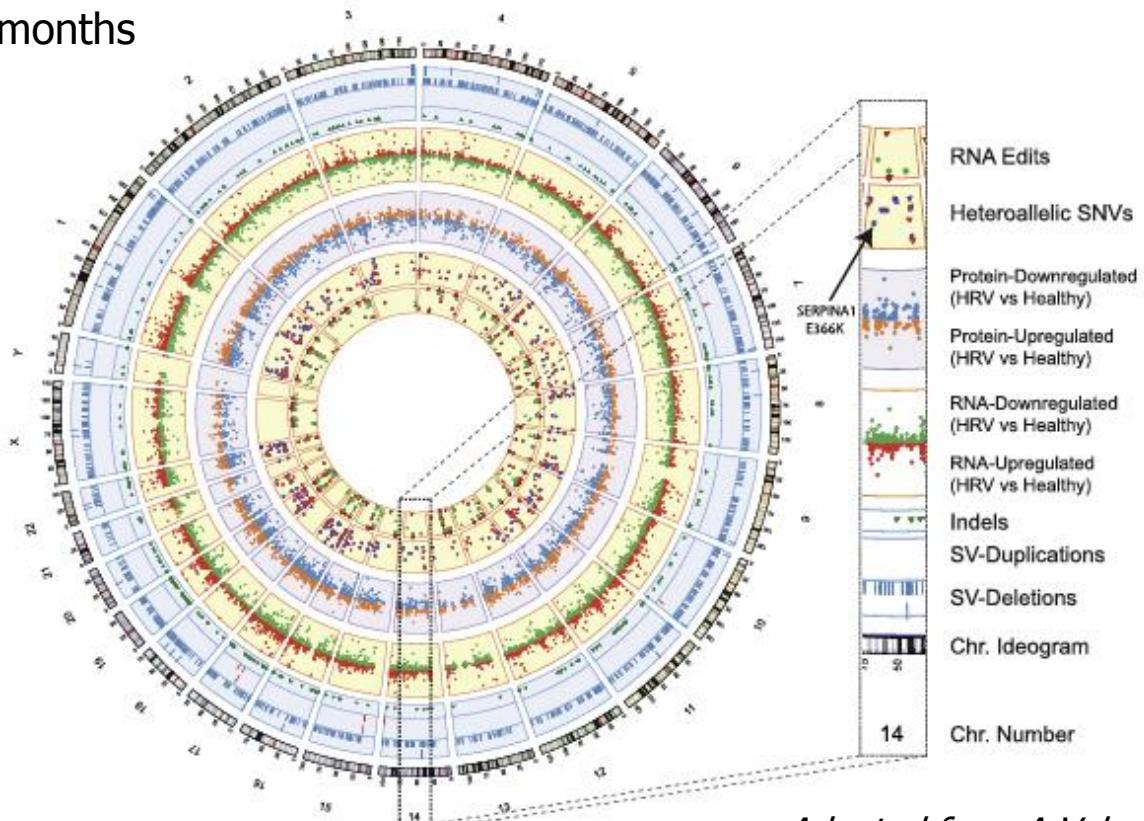


1 Individual followed for 14 months



Personalized Medicine

Source: New Yorker



Adapted from A. Valencia

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

Nick Goldman¹, Paul Bertone¹, Siyuan Chen², Christophe Dessimoz¹, Emily M. Le Proust², Botond Sipos¹ & Ewan Birney¹

Digital production, transmission and storage have revolutionized how we access and use information but have also made archiving an increasingly complex task that requires active, continuing maintenance of digital media. This challenge has focused some interest on DNA as an attractive target for information storage¹ because of its capacity for high-density information encoding, longevity under easily achieved conditions^{2–4} and proven track record as an information bearer. Previous DNA-based information storage approaches have encoded only trivial amounts of information^{2–7} or were not amenable to scaling-up⁸, and used no robust error-correction and lacked examination of their cost-efficiency for large-scale information archival⁹. Here we describe a scalable method that can reliably store more information than has been handled before. We encoded computer files totalling 739 kilobytes of hard-disk storage and with an estimated Shannon information¹⁰ of 5.2×10^6 bits into a DNA code, synthesized this DNA, sequenced it and reconstructed the original files with 100% accuracy. Theoretical analysis indicates that our DNA-based storage scheme could be scaled far beyond current global information volumes and offers a realistic technology for large-scale, long-term and infrequently accessed digital archiving. In fact, current trends in technological advances are reducing DNA synthesis costs at a pace that should make our scheme cost-effective for sub-50-year archiving within a decade.

Although techniques for manipulating, storing and copying large amounts of existing DNA have been established for many years^{11–13}, one of the main challenges for practical DNA-based information storage is the difficulty of synthesizing long sequences of DNA *de novo* to an exactly specified design. As in the approach of ref. 9, we represent the information being stored as a hypothetical long DNA molecule and encode this *in vitro* using shorter DNA fragments. This offers the benefit that isolated DNA fragments are easily manipulated *in vitro*^{11,13}, and that the routine recovery of intact fragments from samples that are tens of thousands of years old^{14,15} indicates that well-prepared synthetic DNA should have an exceptionally long lifespan in low-maintenance environments¹⁶. In contrast, approaches using living vectors^{16–18} are not as reliable, scalable or cost-efficient owing to disadvantages such as constraints on the genomic elements and locations that can be manipulated without affecting viability, the fact that mutation will cause the fidelity of stored and decoded information to reduce over time, and possibly the requirement for storage conditions to be carefully regulated. Existing schemes used for DNA computing in principle permit large-scale memory¹⁹, but data encoding in DNA computing is inextricably linked to the specific application or algorithm¹⁷ and no practical storage schemes have been realized.

As a proof of concept for practical DNA-based storage, we selected and encoded a range of common computer file formats to emphasize the ability to store arbitrary digital information. The five files comprised all 154 of Shakespeare's sonnets (ASCII text), a classic scientific paper¹⁸ (PDF format), a medium-resolution colour photograph of the European Bioinformatics Institute (JPEG 2000 format), a 2-s-excerpt from Martin Luther King's 1963 'I have a dream' speech (MP3 format) and a Huffman code²⁰ used in this study to convert bytes to base-3

digits (ASCII text), giving a total of 757,051 bytes or a Shannon information¹⁰ of 5.2×10^6 bits (see Supplementary Information and Supplementary Table 1 for full details).

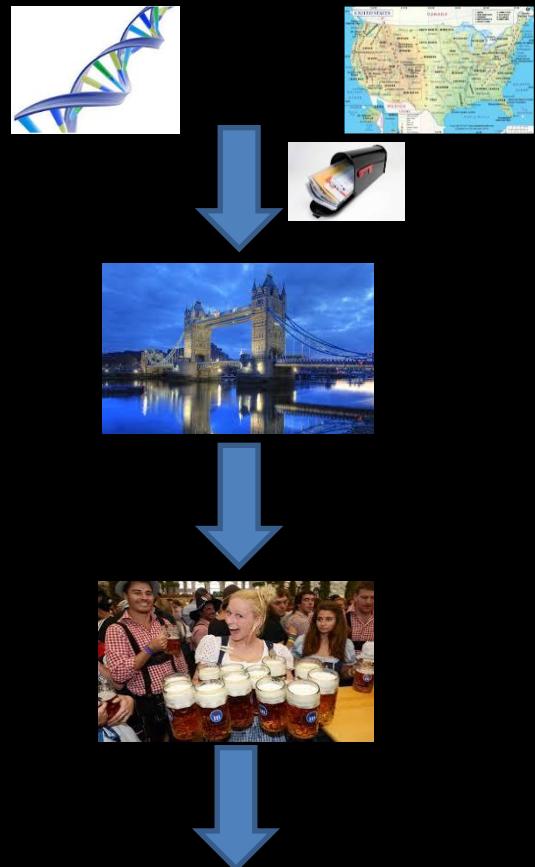
The bytes comprising each file were represented as single DNA sequences with no homopolymers (runs of ≥ 2 identical bases, which are associated with higher error rates in existing high-throughput sequencing technologies¹⁹ and led to errors in a recent DNA-storage experiment⁹). Each DNA sequence was split into overlapping segments, generating fourfold redundancy, and alternate segments were converted to their reverse complement (see Fig. 1 and Supplementary Information). These measures reduce the probability of systematic failure for any particular string, which could lead to uncorrectable errors and data loss. Each segment was then augmented with indexing information that permitted determination of the file from which it originated and its location within that file, and simple parity-check error-detection²⁰. In all, the five files were represented by a total of 153,335 strings of DNA, each comprising 117 nucleotides (nt). The perfectly uniform fragment lengths and absence of homopolymers make it obvious that the synthesized DNA does not have a natural (biological) origin, and so imply the presence of deliberate design and encoded information².

We synthesized oligonucleotides (oligos) corresponding to our designed DNA strings using an updated version of Agilent Technologies' OLS (oligo library synthesis) process²⁰, creating $\sim 1.2 \times 10^7$ copies of each DNA string. Errors occur only rarely (~ 1 error per 500 bases) and independently in the different copies of each string, again enhancing our method's error tolerance. We shipped the synthesized DNA in lyophilized form that is expected to have excellent long-term preservation characteristics²⁴, at ambient temperature and without specialized packaging, from the USA to Germany via the UK. After resuspension, amplification and purification, we sequenced a sample of the resulting library products at the EMBL Genomics Core Facility in paired-end mode on the Illumina HiSeq 2000. We transferred the remainder of the library to multiple aliquots and re-lyophilized these for long-term storage.

Our base calling using AYB²¹ yielded 79.6×10^6 read-pairs of 104 bases in length, from which we reconstructed full-length (117-nt) DNA strings *in silico*. Strings with uncertainties due to synthesis or sequencing errors were discarded and the remainder decoded using the reverse of the encoding procedure, with the error-detection bases and properties of the coding scheme allowing us to discard further strings containing errors. Although many discarded strings will have contained information that could have been recovered with more sophisticated decoding, the high level of redundancy and sequencing coverage rendered this unnecessary in our experiment. Full-length DNA sequences representing the original encoded files were then reconstructed *in silico*. The decoding process used no additional information derived from knowledge of the experimental design. Full details of the encoding, sequencing and decoding processes are given in Supplementary Information.

Four of the five resulting DNA sequences could be fully decoded without intervention. The fifth however contained two gaps, each a run

All 154 of Shakespeare's sonnets (ASCII text), WC paper, a EBI photo, Martin Luther King's 1963 speech and Huffman code.



All 154 of Shakespeare's sonnets (ASCII text), WC paper, a EBI photo, Martin Luther King's 1963 speech and Huffman code.

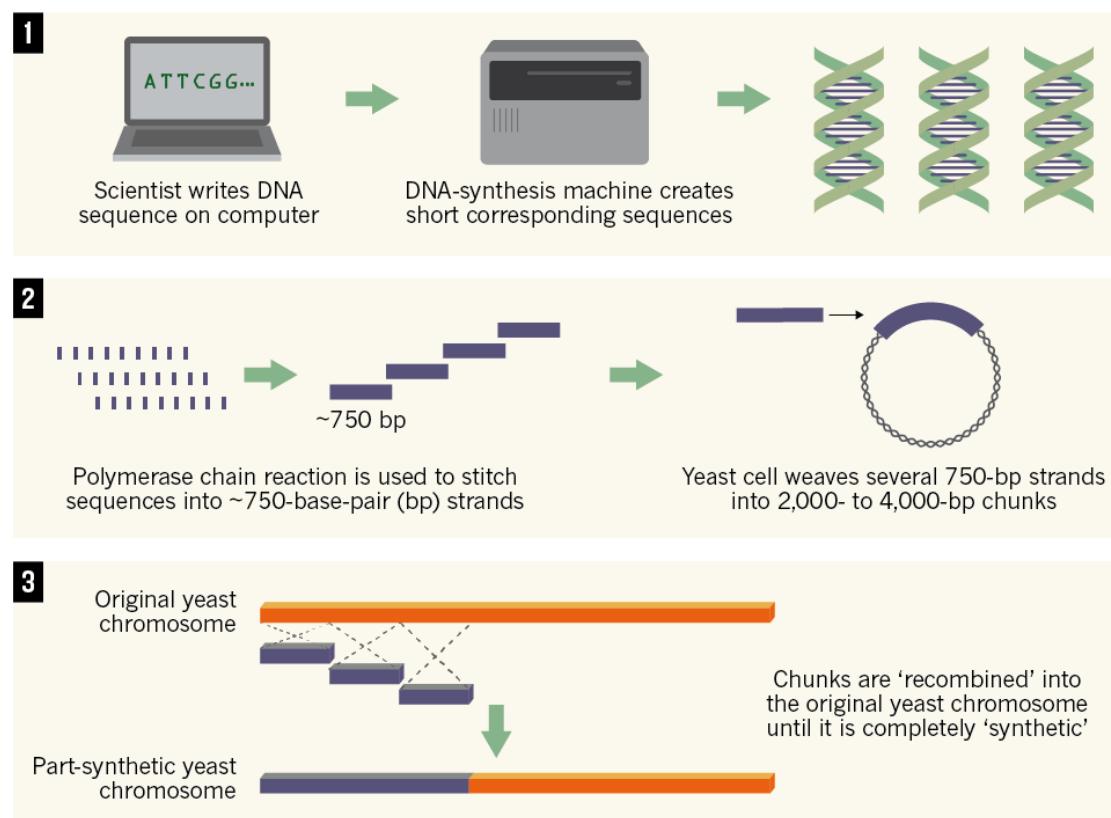
Synthetic biology



- 2010 -> parasite
- 2014: the synthesis of a functional 272,871–base pair designer eukaryotic chromosome
- 2.5% of the 12-million-base-pair *S. cerevisiae* genome

CONSTRUCTING LIFE

Researchers have synthesized a fully functional chromosome from the baker's yeast *Saccharomyces cerevisiae*. At 272,281 base pairs long, it represents about 2.5% of the organism's 12 million-base-pair genome.



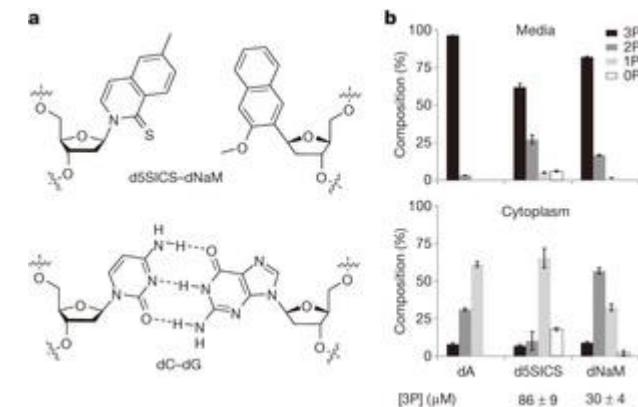
Science April 2014

A semi-synthetic organism with an expanded genetic alphabet



- ❑ Neither the presence of the unnatural triphosphates nor the replication of the UBP introduces a notable growth burden.
- ❑ Lastly, we find that the UBP is not efficiently excised by DNA repair pathways. Thus, the resulting bacterium is the first organism to propagate stably an expanded genetic alphabet.

Nature 509,385–388 (15 May 2014)



Data is extremely noisy
deriving information from data is very costly



TAAAT TAATAGT AT
 ATG AT
 ATAGT AT

30 x to 200 x coverage needed

Final results are strongly dependent on processing
No accepted standards!
Processing is very slow

example just CLL processing → 25% of Marenostrum



Smufin a reference-free method for somatic mutation detection

Normal Tumor

Normal Tumor

Time x 1 cancer sample

Sequencing: < 7 hours

Variant Calling: 1 week on a 512 core cluster

Variant Calling Smufin: 7 hours on a 16 core cluster

either point mutations or
structural variants

- Need to map reads to reference genome and discard variants.
 - Low sensitivity and/or specificity.
 - Detect either point mutation or structural variants
 - Best detection of structural variants of certain sizes.

Point mutations AND
structural variants of any size
(insertions, deletions, inversions)

- Detection of somatic variations through direct read comparison without mapping to a reference genome.
 - High sensitivity and/or specificity.
- Detects both point mutation or structural variants of any size