How feasible is DNA as a long-term digital storage system?

Nick Goldman

EMBL-European Bioinformatics Institute

EMBL-EBI







155 petabytes

of storage capacity in our data centres

EMBL-EBI delivered 140 million jobs to its users in 2017 ~38 million

requests to EMBL-EBI websites every day

Requests from

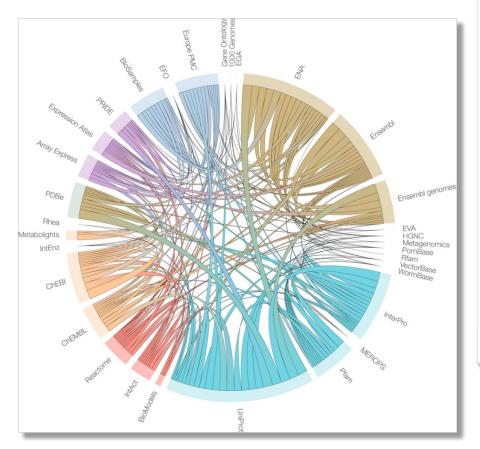
3.3 million

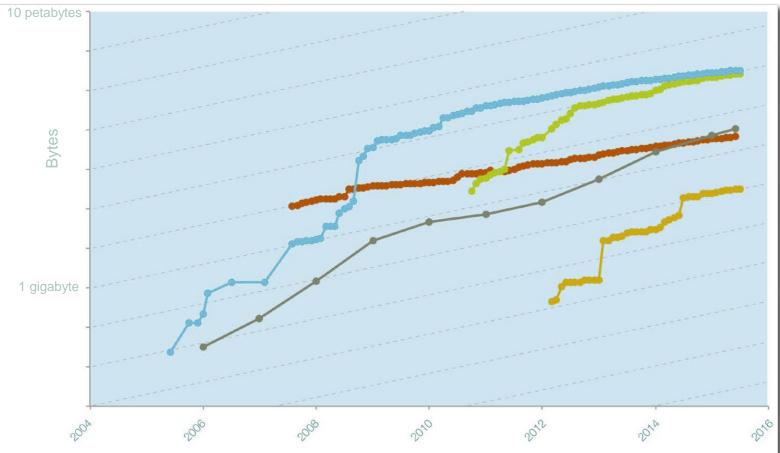
EMBL-EBI websites, each month

EBI is a major data provider to the life sciences research community

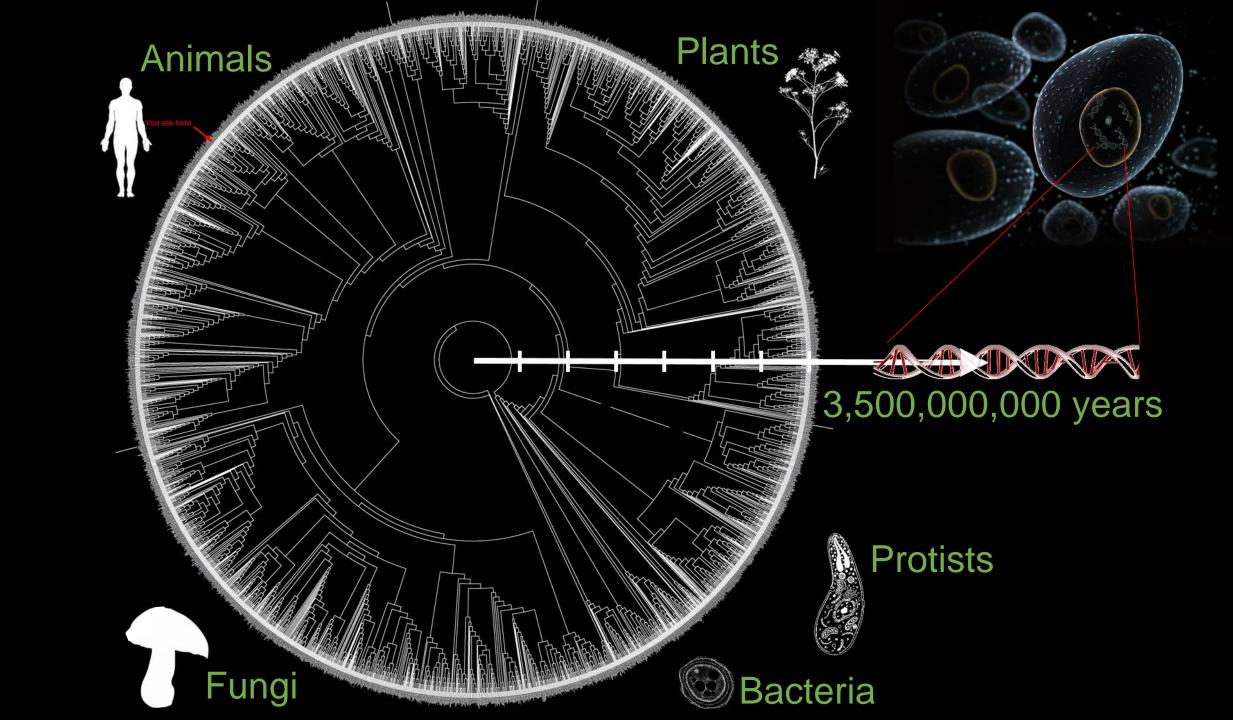


Databases at EBI

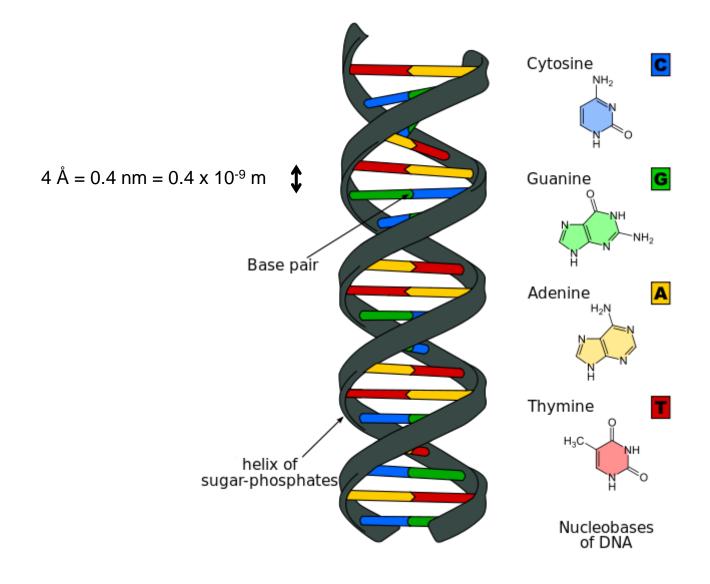




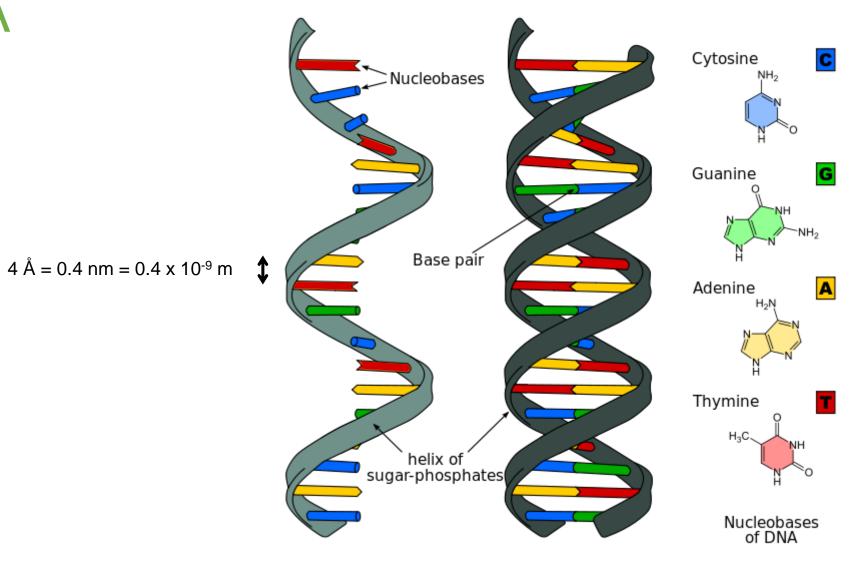




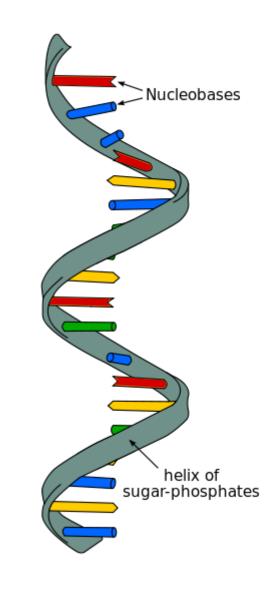
DNA



DNA



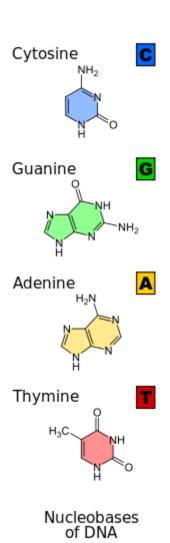
DNA



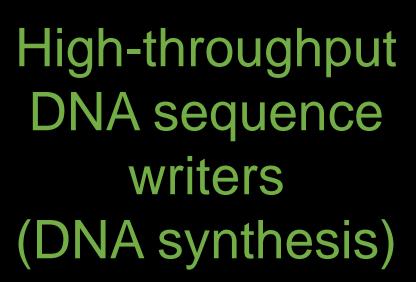
G G

G

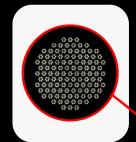
G



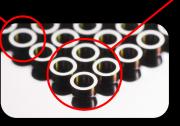














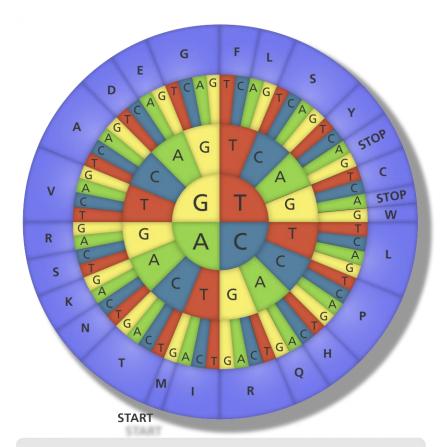
High-throughput
DNA sequence
readers
(DNA sequencing)



	T			С			A			G		
T	TTT	Phe	F	TCT	Ser	S	TAT	Tyr	Y	TGT	Cys	С
	TTC			TCC			TAC			TGC		
	TTA	Leu	L	TCA			TAA	STOP		TGA	STOP	
	TTG			TCG			TAG			TGG	Trp	W
С	CTT	Leu	L	CCT	Pro	P	CAT	His	Н	CGT	Arg	R
	CTC			CCC			CAC			CGC		
	CTA			CCA			CAA	Gln	Q	CGA		
	CTG			CCG			CAG			CGG		
A	ATT	Ile	I	ACT	Thr	T	AAT	Asn	N	AGT	Ser	S
	ATC			ACC			AAC			AGC		
	ATA			ACA			AAA	Lys	K	AGA	Arg	R
	ATG	Met	M	ACG			AAG			AGG		
G	GTT	Val	V	GCT	Ala	Α	GAT	Asp	D	GGT	Gly	G
	GTC			GCC			GAC			GGC		
	GTA			GCA			GAA	Glu	E	GGA		
	GTG			GCG			GAG			GGG		

the Genetic code of life...

...codes for proteins, but not mp3s or pdfs

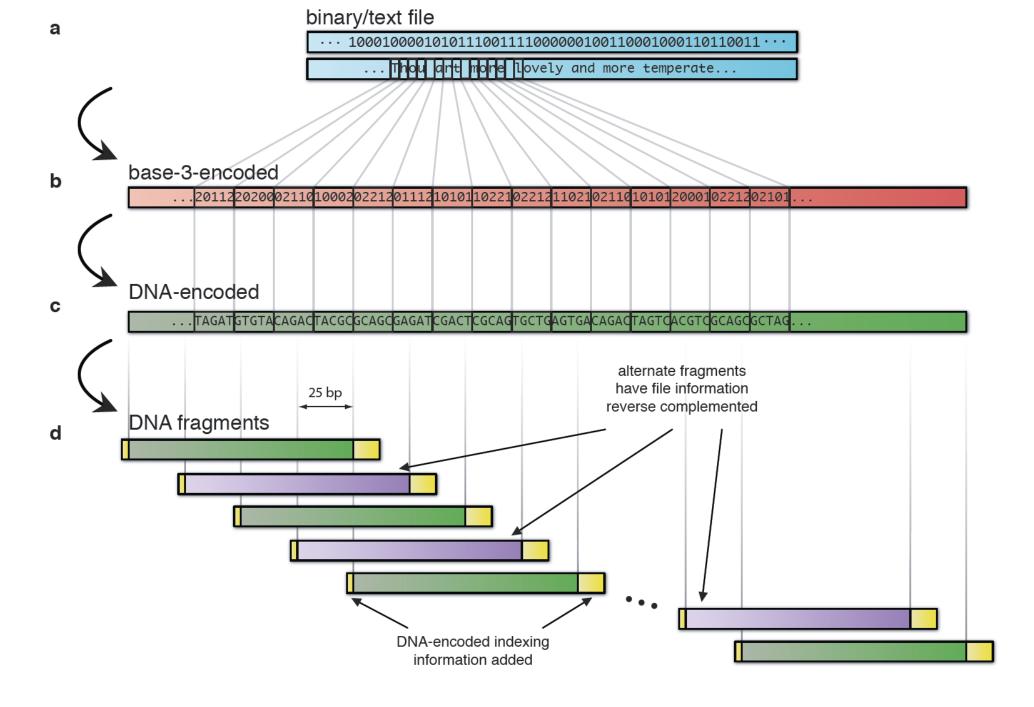


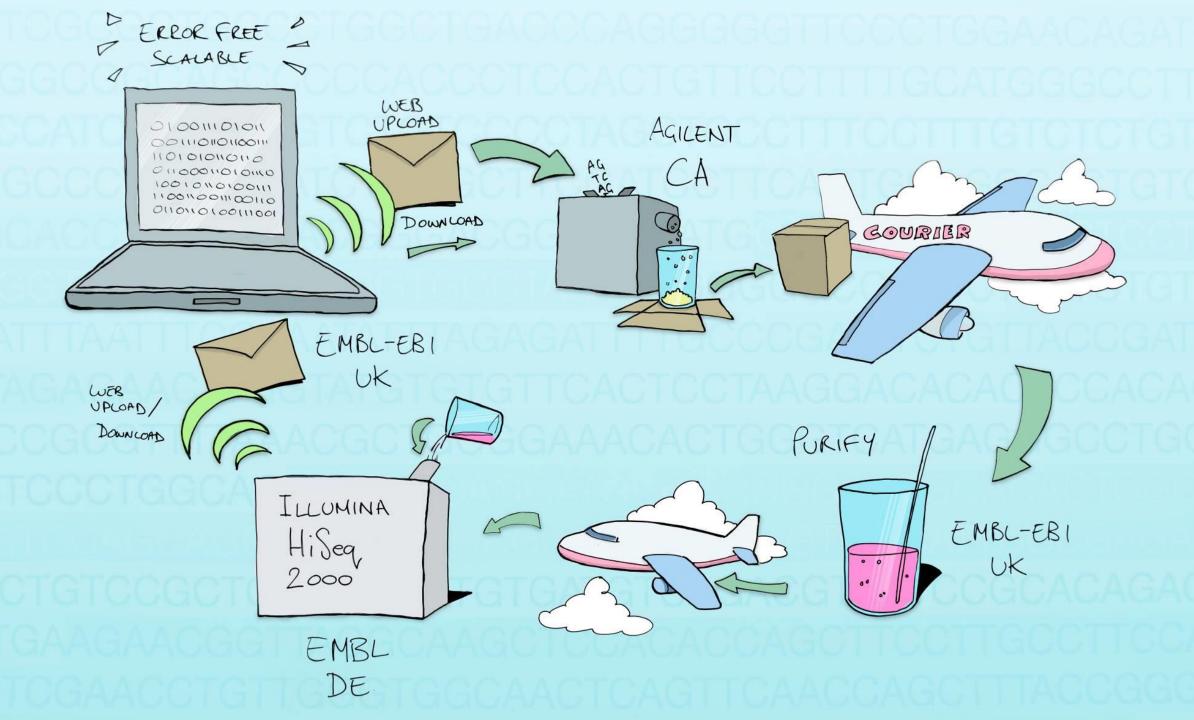
Amino acid code

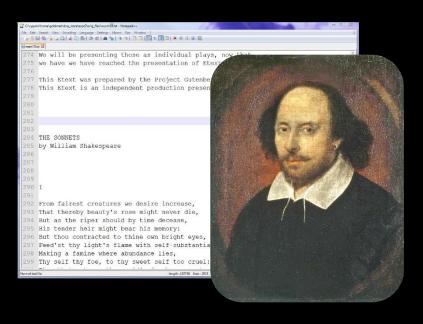
A - Alanine G - Glycine M - Methionine S - Serine **C** - Cysteine H - Histidine N - Asparagine T - Threonine D - Aspartic acid I - Isoleucine P - Proline V - Valine E - Glutamic acid K - Lysine Q - Glutamine W - Tryptophan Y - Tyrosine F - Phenylalanine L - Leucine R - Arginine













No. 4356 April 25, 1953

NATURE

equipment, and to Dr. G. E. R. Deacon and the is a residue on each chain every 3.4 A. in the z-direccaptain and officers of R.R.S. Discovery II for their part in making the observations.

- Young, F. B., Gerrard, H., and Jevons, W., Phil. Mag., 40, 149
- * Longuet-Higgins, M. S., Mon. Not. Roy. Astro. Soc., Geophys. Supp., 8, 285 (1949). Von Arx, W. S., Woods Hole Papers in Phys. Oceanog. Meteor., 11
- *Ekman, V. W., Arkiv. Mat. Astron. Fysik. (Steckholm), 2 (11) (1905).

MOLECULAR STRUCTURE OF NUCLEIC ACIDS

A Structure for Deoxyribose Nucleic Acid

WE wish to suggest a structure for the salt of deoxyribose nucleic acid (D.N.A.). This structure has novel features which are of considerable biological interest.

A structure for nucleic acid has already been proposed by Pauling and Corey1. They kindly made manuscript available to us in advance of publication. Their model consists of three intertwined chains, with the phosphates near the fibre axis, and the bases on the outside. In our opinion, this structure is unsatisfactory for two reasons (1) We believe that the material which gives the X ray diagrams is the salt, not the free acid. Without the acidic hydrogen atoms it is not clear what forces would hold the structure together, especially as the negatively charged phosphates near the axis will repel each other. (2) Some of the van der Waals distances appear to be too small.

Another three-chain structure has also been suggested by Fraser (in the press). In his model the phosphates are on the outside and the bases on the inside, linked together by hydrogen bonds. This structure as described is rather ill-defined, and for

this reason we shall not comment

on it. We wish to put forward a radically different structure for the salt of deoxyribose nucleic acid. This structure has two helical chains each coiled round the same axis (see diagram). We have made the usual chemical assumptions, namely, that each chain consists of phosphate diester groups joining B-D-deoxyribofuranose residues with 3',5' linkages. The two chains (but Both chains follow righthanded helices, but owing to chemical arguments. in opposite directions. Each loosely resembles Furberg's model No. I; that is, the bases are on the inside of the helix and the phosphates on elsewhere. the outside. The configuration of the sugar and the atoms sugar being roughly perpendi-

tion. We have assumed an angle of 36° between adjacent residues in the same chain, so that the structure repeats after 10 residues on each chain, that is, after 34 A. The distance of a phosphorus atom from the fibre axis is 10 A. As the phosphates are on the outside, cations have easy access to them.

The structure is an open one, and its water content is rather high. At lower water contents we would expect the bases to tilt so that the structure could become more compact.

The novel feature of the structure is the manner in which the two chains are held together by the purine and pyrimidine bases. The planes of the bases are perpendicular to the fibre axis. They are joined together in pairs, a single base from one chain being hydrogen-bonded to a single base from the other chain, so that the two lie side by side with identical z-co-ordinates. One of the pair must be a purine and the other a pyrimidine for bonding to occur. The hydrogen bonds are made as follows : purine position to pyrimidine position 1; purine position 6 to pyrimidine position 6.

If it is assumed that the bases only occur in the structure in the most plausible tautomeric forms (that is, with the keto rather than the enol configurations) it is found that only specific pairs of bases can bond together. These pairs are : adenine (purine) with thymine (pyrimidine), and guanine (purine) with cytosine (pyrimidine).

In other words, if an adenine forms one member of a pair, on either chain, then on these assumptions the other member must be thymine; similarly for guanine and cytosine. The sequence of bases on a single chain does not appear to be restricted in any way. However, if only specific pairs of bases can be formed, it follows that if the sequence of bases on one chain is given, then the sequence on the other chain is automatically determined.

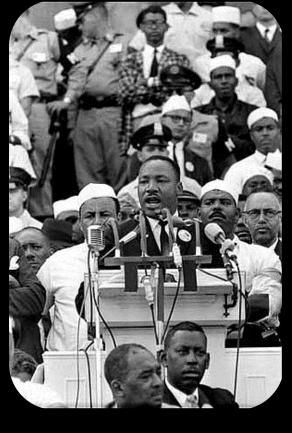
It has been found experimentally3,4 that the ratio of the amounts of adenine to thymine, and the ratio of guanine to cytosine, are always very close to unity for deoxyribose nucleic acid.

It is probably impossible to build this structure with a ribose sugar in place of the deoxyribose, as the extra oxygen atom would make too close a van der Waals contact.

The previously published X-ray datas, on deoxyribose nucleic acid are insufficient for a rigorous test of our structure. So far as we can tell, it is roughly compatible with the experimental data, but it must be regarded as unproved until it has been checked against more exact results. Some of these are given in the following communications. We were not aware not their bases) are related by a of the details of the results presented there when we dyad perpendicular to the fibre devised our structure, which rests mainly though not entirely on published experimental data and stereo-

the dyad the sequences of the atoms in the two chains run pairing we have postulated immediately suggests a possible copying mechanism for the genetic material. Full details of the structure, including the conditions assumed in building it, together with a set of co-ordinates for the atoms, will be published

We are much indebted to Dr. Jerry Donohue for constant advice and criticism, especially on internear it is close to Furberg's atomic distances. We have also been stimulated by 'standard configuration', the a knowledge of the general nature of the unpublished experimental results and ideas of Dr. M. H. F. cular to the attached base. There Wilkins, Dr. R. E. Franklin and their co-workers a







1 MB of computer files = almost invisible dust of DNA







MOLECULAR STRUCTURE OF **NUCLEIC ACIDS**

A Structure for Deoxyribose Nucleic Acid

WE wish to suggest a structure for the salt of deoxyribose nucleic acid (D.N.A.). This structure has novel features which are of considerable biological interest.

A structure for nucleic acid has already been proposed by Pauling and Corey1. They kindly made their manuscript available to us in advance of publication. Their model consists of three intertwined chains, with the phosphates near the fibre axis, and the bases on the outside. In our opinion, this structure is unsatisfactory for two reasons: (1) We believe that the material which gives the X-ray diagrams is the salt, not the free acid. Without the acidic hydrogen atoms it is not clear what forces would hold the structure together, especially as the negatively charged phosphates near the axis will repel each other. (2) Some of the van der Waals distances appear to be too small.

Another three-chain structure has also been suggested by Fraser (in the press). In his model the phosphates are on the outside and the bases on the inside, linked together by hydrogen bonds. This structure as described is rather ill-defined, and for

This figure is purely

diagrammatic. The two ribbons symbolize the

two phosphate—sugar chains, and the hori-

bases holding the chains.

ine marks the fibre axis.

this reason we shall not comment

on it. We wish to put forward a radically different structure for the salt of deoxyribose nucleic acid. This structure has two helical chains each coiled round the same axis (see diagram). We have made the usual chemical assumptions, namely, that each chain consists of phosphate diester groups joining \$\beta-D-deoxyribofuranose residues with 3',5' linkages. The two chains (but not their bases) are related by a dyad perpendicular to the fibre axis. Both chains follow righthanded helices, but owing to the dyad the sequences of the atoms in the two chains run in opposite directions. Each chain loosely resembles Furberg's2 model No. I; that is, the bases are on the inside of the helix and the phosphates on the outside. The configuration of the sugar and the atoms near it is close to Furberg's 'standard configuration', the sugar being roughly perpendicular to the attached base. There

MOLECULAR STRUCTURE OF NUCLEIC ACIDS

A Structure for Deoxyribose Nucleic Acid

WE wish to suggest a structure for the salt of deoxyribose nucleic acid (D.N.A.). This structure has novel features which are of considerable biological interest.

A structure for nucleic acid has already been proposed by Pauling and Corey¹. They kindly made their manuscript available to us in advance of publication. Their model consists of three intertwined chains, with the phosphates near the fibre axis, and the bases on the outside. In our opinion. this structure is unsatisfactory for two reasons: (1) We believe that the material which gives the X-ray diagrams is the salt, not the free acid. Without the acidic hydrogen atoms it is not clear what forces would hold the structure together, especially as the negatively charged phosphates near the axis will repel each other. (2) Some of the van der Waals distances appear to be too small.

Another three-chain structure has also been suggested by Fraser (in the press). In his model the phosphates are on the outside and the bases on the inside, linked together by hydrogen bonds. This structure as described is rather ill-defined, and for

this reason we shall not comment



This figure is purely diagrammatic. The two ribbons symbolize the two phosphate—sugar zontal rods the pairs of bases holding the chains cular to the attached base. There line marks the fibre axis

We wish to put forward a radically different structure for the salt of deoxyribose nucleic acid. This structure has two helical chains each coiled round the same axis (see diagram). We have made the usual chemical assumptions, namely, that each chain consists of phosphate diester groups joining \$-D-deoxyribofuranose residues with 3',5' linkages. The two chains (but not their bases) are related by a dyad perpendicular to the fibre axis. Both chains follow righthanded helices, but owing to the dyad the sequences of the atoms in the two chains run in opposite directions. Each chain loosely resembles Furberg's2 model No. I; that is, the bases are on the inside of the helix and the phosphates on the outside. The of the sugar and near it is close 'standard configu sugar being roug

What has happened since?

"NAM" (Nucleic Acid Memory) research established

















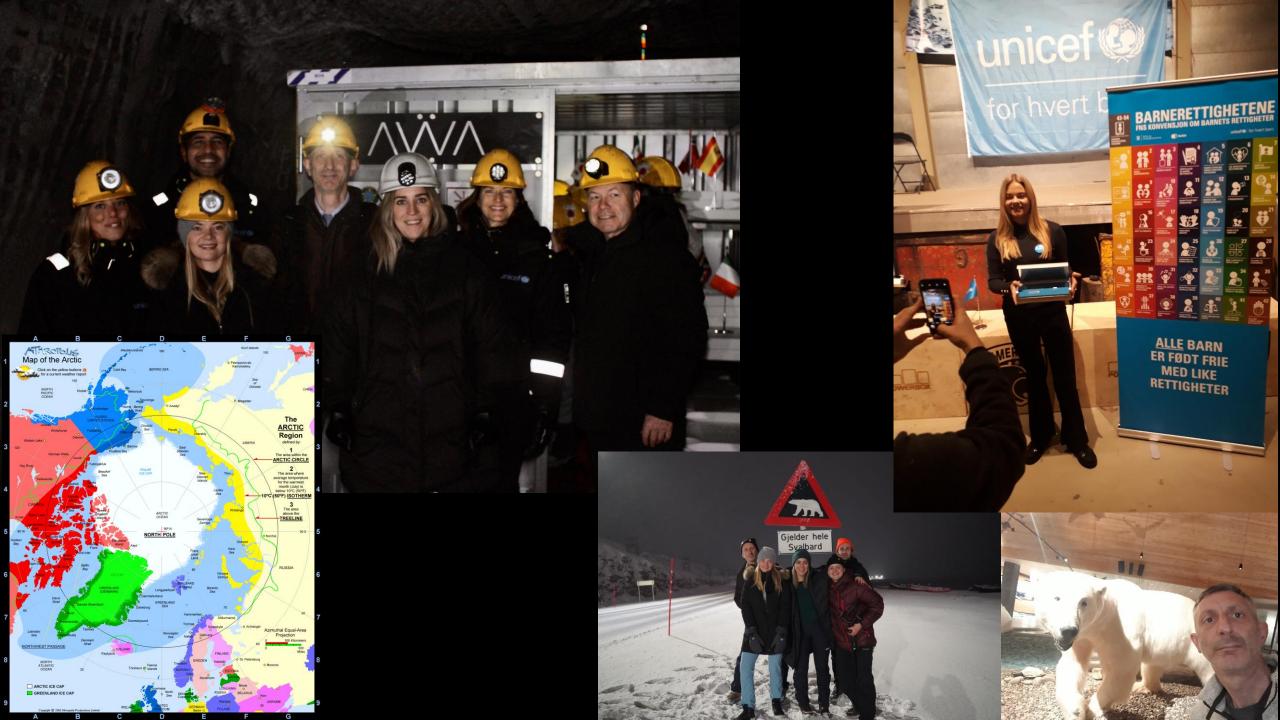






Bitcoin Challenge Davos 2015 Sequence the DNA, decode the message, claim the prize of 1 Bitcoin







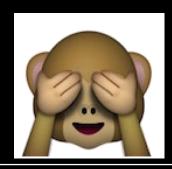
Technicolor and Harvard University

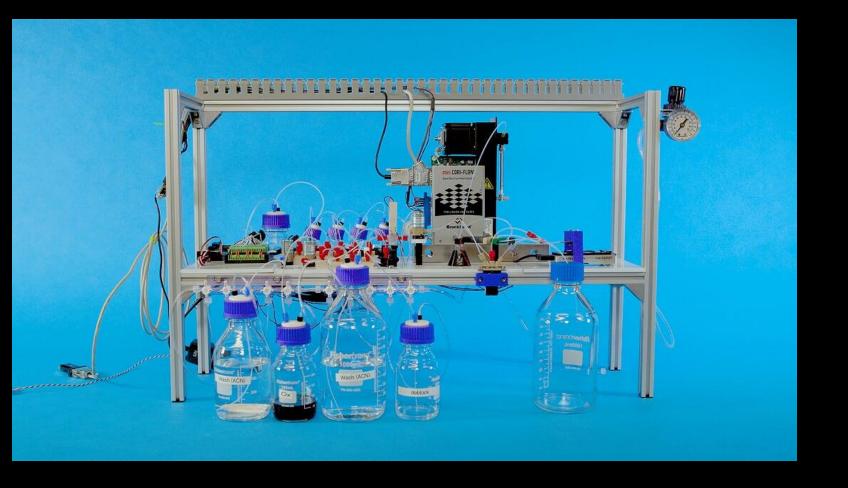
University of Washington and Microsoft Research

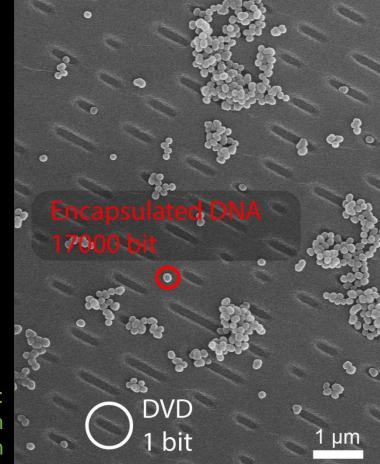












Picture credit:
Robert Grass, Julian Koch
ETH Zürich

6 areas where progress is being made:

- phosphoramidite synthesis continuing to improve (accuracy, speed, cost); enzymatic synthesis companies appearing regularly (various companies)
- new data encoding techniques reduce the amount of DNA needed (multiple teams)
- DNA random access memory (RAM University of Illinois)
- first 'table-top' end-to-end device prototype exists (University of Washington/Microsoft)
- can data encoding permit 'computing in DNA'? (image search University of Washington/Microsoft)
- integration of information in DNA into materials (ETH Zürich)



Why use DNA for storing information?

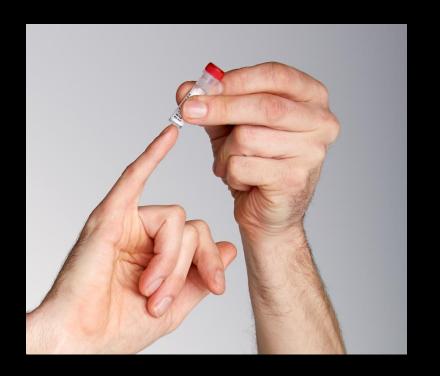
Why use DNA for storing information?

First: why not to use DNA for storing



Why use DNA for storing information?





1 test-tube full of DNA = 1000000 CD ROMs





x 10 000

MB GB 80-3 ZB (Estimated global 60-Cost (10⁴\$ MB⁻¹) data amount) Efficiency (%) 40-Current costs 20-Projected costs (100× cheaper) 10^{13} 1025 1037 1049 10^{1} Information to be encoded (B)

Encoding technique scales up beyond global data volumes

Archives made of DNA would take up very little space



7 MAY 2010 VOL 328 SCIENCE www.sciencemag.org Vol 456 20 November 2008 doi:10.1038/nature07446

RESEARCH ARTICLE

40 000 years old A Draft Sequence of the **Neandertal Genome**

Richard E. Green, *† Johannes Krause, † Adrian Udo Stenzel, † Martin Kircher, † Nick Patterson, 2

]effrey Hernár Barbar Eric S. Christi Vladin lavier Daniel Janet

DNA lasts for a very long time

Sequencing the nuclear genome of the extinct woolly mammoth

Webb Miller¹, Daniela I. Drautz¹, Aakrosh Ratan¹, Barbara Pusey¹, Ji Qi¹, Arthur M. Lesk¹, Lynn P. Tomsho¹, Michael D. Packard¹, Fangqing Zhao¹, Andrei Sher²†, Alexei Tikhonov³, Brian Ranev⁴, Nick Patterson⁵ in Lindblad-Toh⁵, Eric S. Lander⁵, Jai

20 000 years old

Tom Pringle⁸ & Steph

www.sciencemag.org SCIENCE VOL 306 26 NOVEMBER 2004

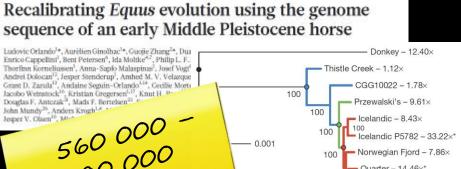
Rise and Fall of the

Beth Shapiro, 1,2 Alexei J. Drummond, Andrew Rambaut, 2 Michael C. Wilson.³ Paul E. Matheus.⁴ Andrei V. Sher.⁵ Oliver G. Pybus, M. Thomas P. Gilbert, 1,2 Ian Barnes, 6 Jonas Binladen, Eske Willerslev, Anders J. Hansen,

Jonathan C. Driver, 11 Duane G. Fro Grant Keddie, ¹⁴ Pavel Kosint Larry D. Martin, ¹⁷ Robert O. St Richard Tedford, 20 Sergei Z

Beringian Steppe Bison

Gennady F. Baryshnikov, ⁸ James A. Burns, ⁹ Sergei Davydov, ¹⁰



63 — Thoroughbred – 21.08×

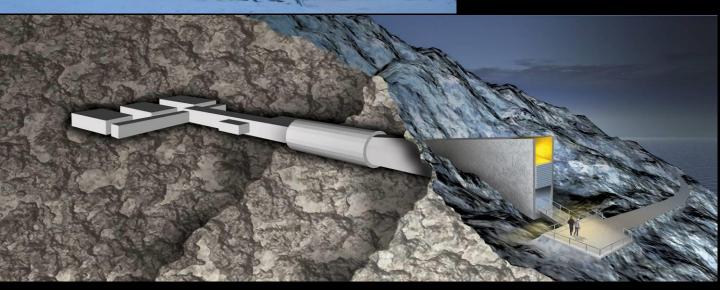
Ludovic Orlando¹*, Aurèlien Ginolhac¹*, Guojie Zhang²*, Dua Enrico Cappellini¹, Bent Petersen⁶, Ida Moltke^{4,7}, Philip L. F. Thorfinn Korneliussen¹, Anna-Sapío Malaspinas², Josef Vogt⁴ Thistle Creek - 1.12× Andrei Dolocan¹², Jesper Stenderup¹, Amhed M. V. Velazque Grant D. Zazula¹¹, Andaine Seguin-Orlando^{1,14}, Gecilie Morte Jacobo Weinstock 16, Kristian Gregersen 1, 17, Knut H. P. Douglas F. Antezak²⁸, Mads F. Bertelsen²² John Mundy26, Anders Krogh14 years old Quarter - 14.46× Standardbred - 12.16× Myr (0.341-0.375) 0.287 Myr (0.274-0.307) Arabian - 11.03×

74 | NATURE | VOL 499 | 4 JULY 20





DNA lasts for a very long time in easy-to-create conditions









With DNA it is easy to make many copies of your data — exponentially!











Card Reader Service for 80-Column IBM Punch Cards http://PunchCardReader.com

7 | 2227777277 | 2277277 | 22772777 | 2277277 | 2277277777 | 227727777 | 2277277777 | 22777 | 22777 | 22777 | 2

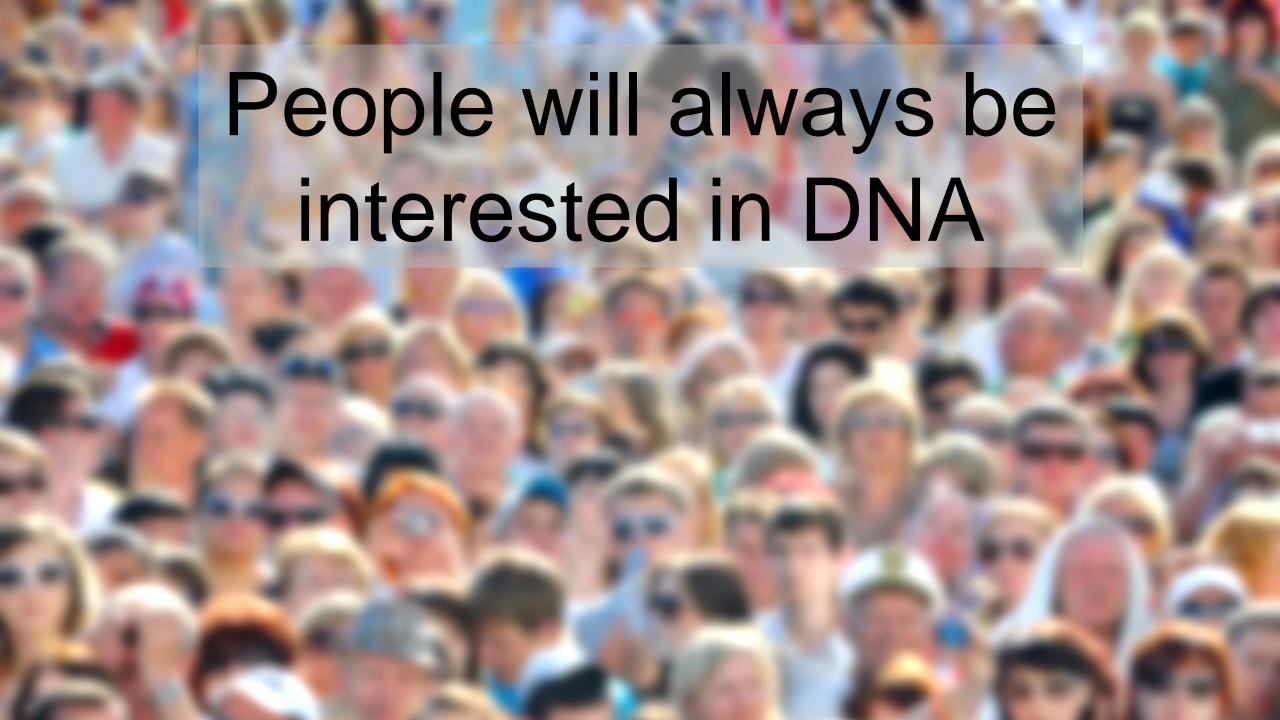








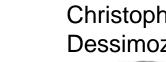


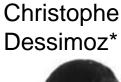


Ewan Birney



Paul Bertone*





Agilent

Technologies



Botond

Jossy Sayir*











Emily Hesketh*

and additional thanks to:

David MacKay[†], Graeme Mitchison[†] **Cambridge:**

EBI: Kevin Gori*, Daniel Henk*, Remco Loos*, Ari Löytynoja*, Hazel Marsden*, Tim Massingham*,

Sarah Parks*, Roland Schwarz*

EMBL-Heidelberg: Vladimir Benes, Dinko Pavlinić, Jonathon Blake*

EMBLEM: Birgit Kerber, Boris Bryk

Art for Eating: Charlotte Jarvis Echidna: Oliver Bradley * = now moved to new positions