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

Presented by Divya and Natalie

March 13th, 2013

20.385

Nature (2013) 494: 77–80 doi:10.1038/nature11875

Who needs digital storage?





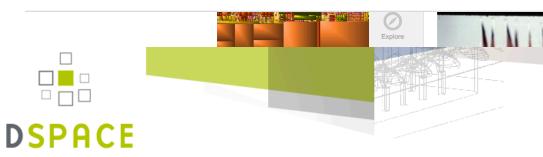
iCloud: Purchasing iCloud storage and billing

Summary

iCloud customers are provided with 5 GB of free cloud storage. Purchased music, movies, TV shows, apps, and books, as well as photos in your Photo Stream don't count against your 5 GB of free storage.

Products Affected

iCloud





ABOUT DSPACE

- About DSpace
- · Why Use DSpace?
- · Who's Using DSpace
- Use Case Examples
- Supporting Organization
- Service Providers

About DSpace

DSpace is the software of choice for academic, non-profit, and commercial organizations building open digital repositories. It is free and easy to install "out of the box" and completely customizable to fit the needs of any organization.

DSpace preserves and enables easy and open access to all types of digital content including text, images, moving images, mpegs and data sets. And with an ever-growing community of developers, committed to continuously expanding and improving the software, each DSpace installation benefits from the next.



Other examples of DNA storage

Eduardo Kac, 1999

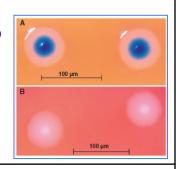
GENESIS

Bible quote to Morse code to DNA code to

expressed protein

Gibson et al, 2010 *Synthetic Genomics*

Watermarks with 46 authors and James Joyce and Richard Feynman quotes



George Church, 2012 *Regenesis*



Next-Generation Digital Information Storage in DNA

George M. Church 1,2, Yuan Gao3, Sriram Kosuri 1,2,*

- (-) living systems are unreliable, existing examples don't scale and \$\$\$\$

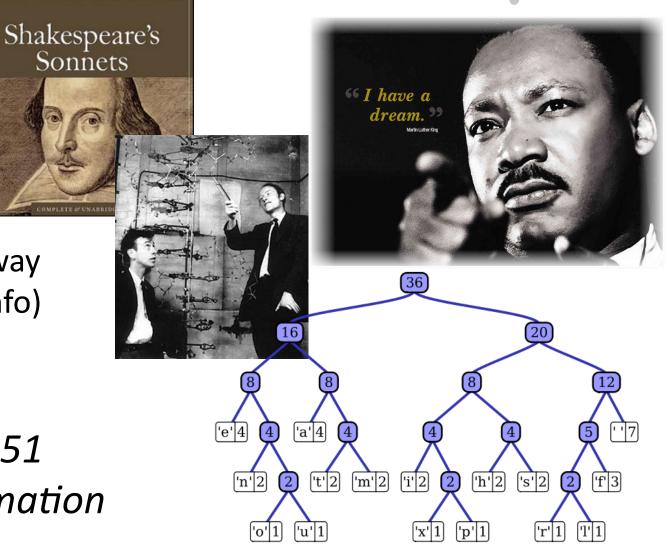
Proof of concept

Sonnets

EMBL-EBI

- ASCII txt
- •PDF
- •JPEG 2000
- •MP3
- Huffman code (way of compressing info)

TOTAL = 757,051bytes of information



An arbitrary computer file a string S_{\emptyset} of bytes, i.e. a value in the set $\{\emptyset \dots 255\}$

 $S_{\emptyset} = Birney \cup and \cup Goldman$

An arbitrary computer file a string S_{\emptyset} of bytes, i.e. a value in the set $\{\emptyset \dots 255\}$ given Huffman code, converting it to base-3. S_1 of characters in $\{\emptyset, 1, 2\}$. 'trit'. $S_1 = 20100 20210 10101 00021 20001 222111 02212 01112 00021$ 22100 02212 222212 02110 02101 22100 11021 01112 00021 d

then add "indexing" and "orientation" info

An arbitrary computer file

a string S_{\emptyset} of bytes, i.e. a value in the set $\{\emptyset \dots 255\}$

given Huffman code, converting it to base-3.

 S_1 of characters in $\{\emptyset, 1, 2\}$

'trit'.

converted	l to	a Dl	ŇΑ	string
-----------	------	------	----	--------

🕻 previous	next trit to encode		
nt written	Ø	1	2
Α	С	G	Т
С	G	Т	Α
G	Т	Α	C
Т	Α	С	G

An arbitrary computer file a string S_{\emptyset} of bytes, i.e. a value in the set $\{\emptyset \dots 255\}$

given Huffman code, converting it to base-3. S_1 of characters in $\{\emptyset, 1, 2\}$ 'trit'.

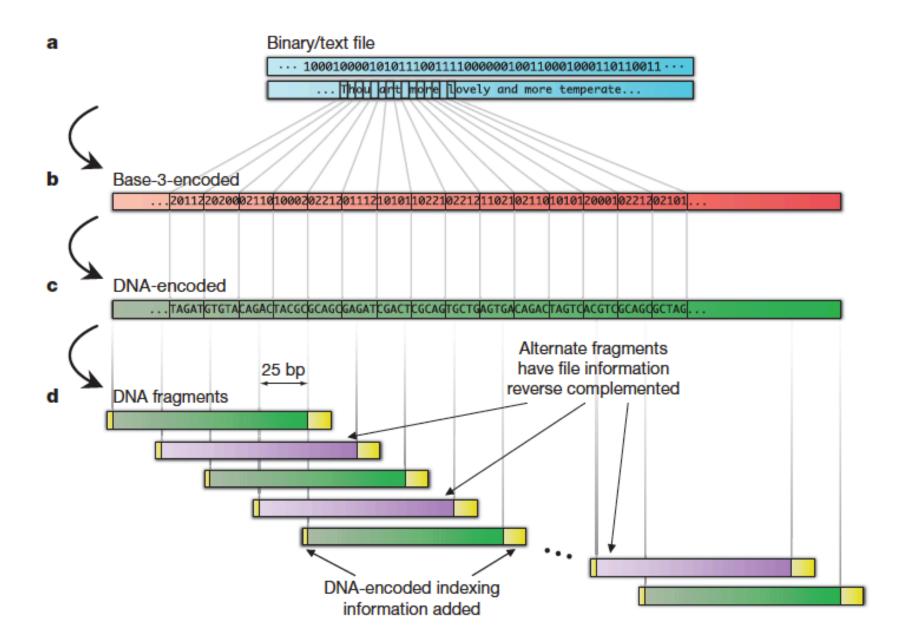
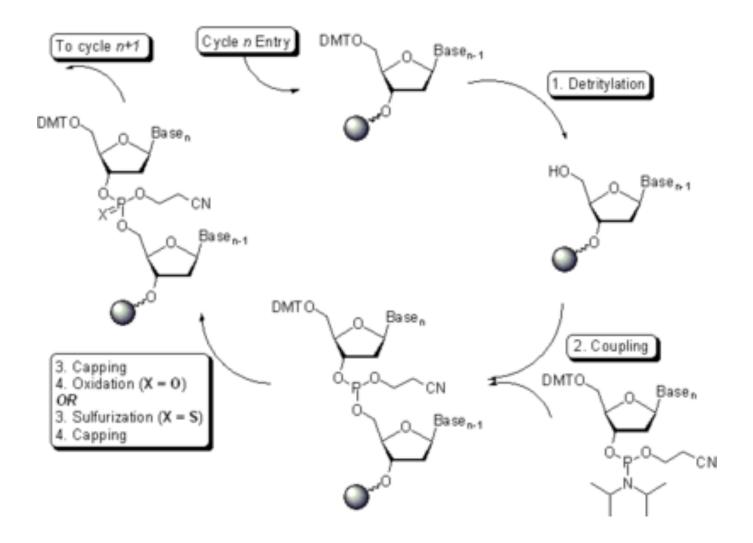


Fig. 1



Beaucage SL, Caruthers MH. (1981) Tetrahedron Lett. 22, 1859-62.

Precise details of the decoding procedure are left as an exercise for the reader.

Errors introduced during DNA synthesis, storage or sequencing could lead to various artefacts, particularly nt insertion, deletion or substitution. Recovery of information from fragments with such errors may be possible (details left as exercise)—we have not found this to be necessary due to the large numbers of perfectly-sequenced fragments available via high-throughput sequencing.

Key Questions

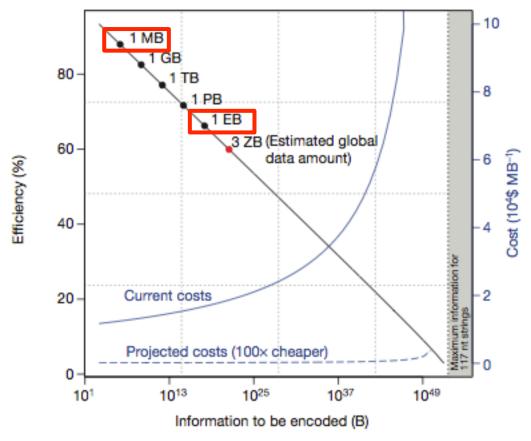
1. Can large amounts of data be stored using DNA?

2. Is decoding accurate?

3. Is DNA-storage cost-effective?

DNA-based Storage Feasible for Large Data

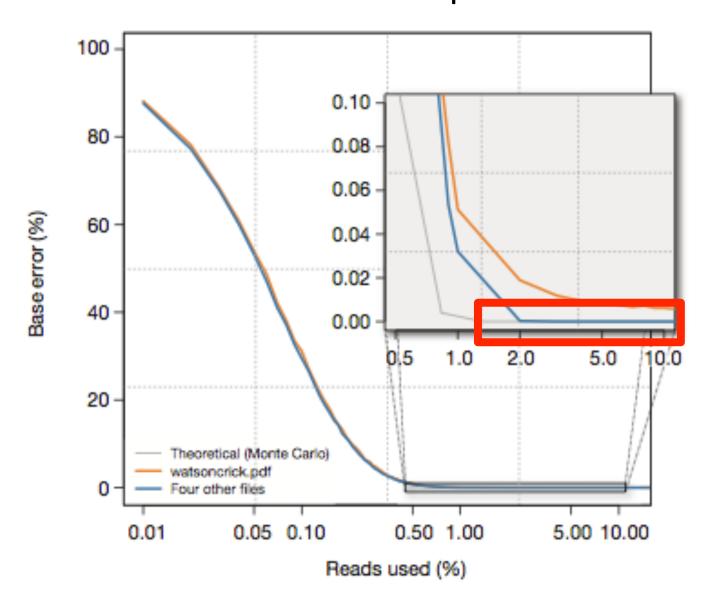




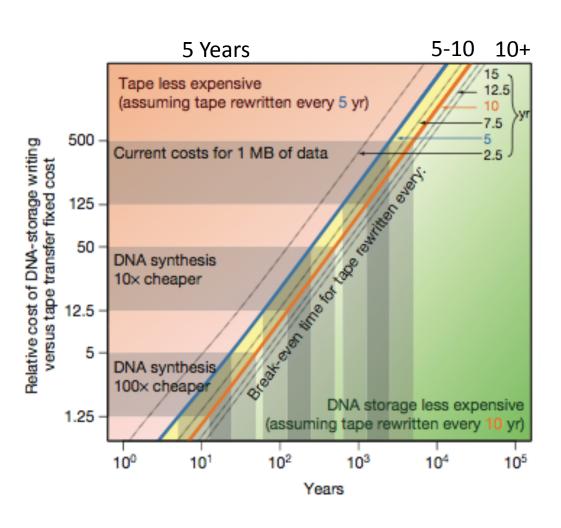
Bytes	Efficiency
10 ³	88%
10^{18}	>65%

Projected 2 magnitude reduction in cost

Zero Percent Error Rate after Reading More Than 2% of Sample

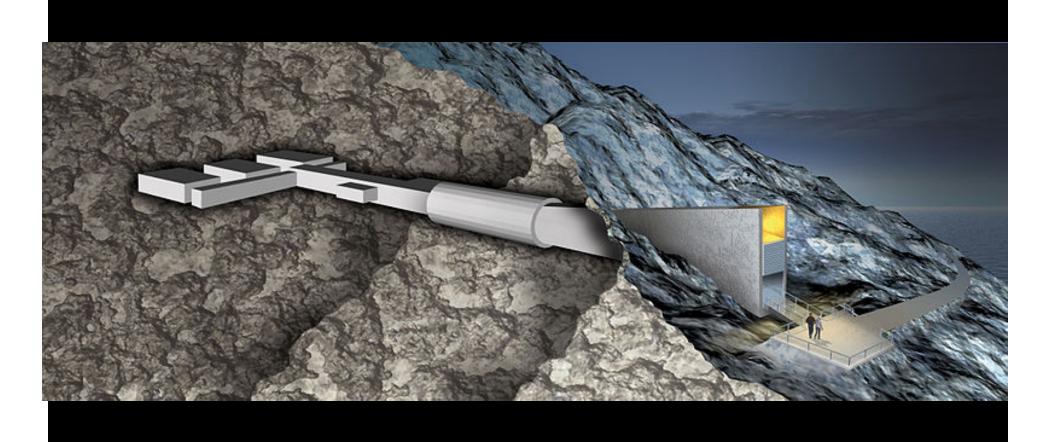


DNA-based Storage Ideal for Long-Horizon Archives with Low Access Rates



- LHC stores 90% of data on magnetic tape
- Magnetic tape needs to be rewritten periodically to counter data loss resulting from magnetic field breakdown
- Shortest breakeven time: ~100 years







Assumptions

- Cost of synthesis will decrease 100x
- No public wariness
- DNA will remain stable and fully readable over all periods of time
- Current decoding practices will still be in place and understandable in 500+ years
- Preferable to not maintain records

Limitations

- Accessibility
- Not widespread information storage plan
- Consume 10% of information every time you decode
- Not yet cost-effective
- Easily-susceptible to environmental damage

Significance

- Proof of concept of DNA as a dense and stable information storage system
- DNA as storage device for multiple file types
- DNA as a computer language
- Cool science!

Pros

- Low maintenance
- Low error-rate
- Addresses the problem of long-term data storage

Cons

- Cost-effectiveness
- Accessibility
- Public hesitation towards biology

Discussion Questions

- What are some concerns the general public would have about DNA-based storage?
- Would it be wise to put vital information such as government information and historical records in a relatively inaccessible place?
- How likely are people to adopt DNA-storage in the near future?