



Beyond the Genome 2013

Informatics Challenge

Michael Schatz
Sven-Eric Schelhorn, October 2013



The prize

2

iPad mini



Reads

A metagenomic sample was generated by mixing portions of the reference sequences of several microbial species. Sequence reads were simulated from these portions. Within each portion of a reference sequence, a foreign insert (i.e, not originating from any of the microbial species) was placed. These inserts encode a message.



Message

One of the inserts corresponds to the 'wildtype', as deposited in public sequence databases such as NCBI nt. The other inserts are slight variations of this wildtype (>90% nuc. similarity).

The message we are seeking is encoded as nucleotide variants of the non-wildtype inserts with respect to the wildtype insert.

Consequently, there is one message for each non-wildtype insert in the read data. All messages together yield a quote that is the solution to the challenge.



How to encode a message into DNA

5

The screenshot displays the Science Magazine website interface. At the top, the Science logo is prominent, followed by navigation links for AAAS.ORG, FEEDBACK, HELP, and LIBRARIANS. A search bar is located on the right, with 'Science Magazine' entered. Below the search bar, a red navigation bar contains links for NEWS, SCIENCE JOURNALS, CAREERS, BLOGS & COMMUNITIES, MULTIMEDIA, and COLLECTIONS. A 'JOIN / SUBSCRIBE' button is also present. The main header area features the Science logo and the tagline 'The World's Leading Journal of Original Scientific Research, Global News, and Commentary.' Below this, a secondary navigation bar includes links for Science Home, Current Issue, Previous Issues, Science Express, Science Products, My Science, and About the Journal. The article page for 'Next-Generation Digital Information Storage in DNA' is displayed, published online on August 16, 2012. The article is by George M. Church, Yuan Gao, and Sriram Kosuri. The abstract describes a strategy to encode arbitrary digital information in DNA, write a 5.27-megabit book using DNA microchips, and read the book using next-generation DNA sequencing. The page also includes a sidebar with 'Article Views' (Abstract, Full Text, Full Text (PDF), Supplementary Materials) and 'Article Tools' (Save to My Folders, Download Citation, Alert Me When Article is Cited, Post to CiteULike, E-mail This Page, Get Permission, View PubMed Citation). A 'Related Content' section is at the bottom left. On the right, there are promotional banners for 'Get all of Science' and a 'WEBINAR' titled 'Clinical Validation of Cancer Biomarker Signatures using Array Technology'.

Science Magazine

AAAS.ORG | FEEDBACK | HELP | LIBRARIANS

Science Magazine

WASHINGTON UNIV SCH OF MED | ALERTS | ACCESS RIGHTS | MY ACCOUNT | SIGN IN

NEWS | SCIENCE JOURNALS | CAREERS | BLOGS & COMMUNITIES | MULTIMEDIA | COLLECTIONS | JOIN / SUBSCRIBE

Science The World's Leading Journal of Original Scientific Research, Global News, and Commentary.

Science Home | Current Issue | Previous Issues | Science Express | Science Products | My Science | About the Journal

Home > Science Magazine > Science Express > Church et al.

Article Views

- Abstract
- Full Text
- Full Text (PDF)
- Supplementary Materials

Article Tools

- Save to My Folders
- Download Citation
- Alert Me When Article is Cited
- Post to CiteULike
- E-mail This Page
- Get Permission
- View PubMed Citation

Related Content

Published Online August 16 2012

Science DOI: 10.1126/science.1226355

BREVIA

Next-Generation Digital Information Storage in DNA

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

¹ Author Affiliations

² To whom correspondence should be addressed. E-mail: sri.kosuri@wss.harvard.edu

ABSTRACT

Digital information is accumulating at an astounding rate, straining our ability to store and archive it. DNA is among the most dense and stable information media known. The development of new technologies in both DNA synthesis and sequencing make DNA an increasingly feasible digital storage medium. Here, we develop a strategy to encode arbitrary digital information in DNA, write a 5.27-megabit book using DNA microchips, and read the book using next-generation DNA sequencing.

Get all of Science

Join Now!

ADVERTISEMENT

WEBINAR

Clinical Validation of Cancer Biomarker Signatures using Array Technology



Numbers are bits

Bit: off or on, 0 or 1

Byte: eight bits, 01001000

Each position represents a power of 2 :

$$01001000 = 8 + 64 = 72$$

*All characters in the ASCII set can be represented by one byte
(0 – 127 (= 2^7-1))*



ASCII code table

Dec	Hex	Oct	Char	Dec	Hex	Oct	Char	Dec	Hex	Oct	Char	Dec	Hex	Oct	Char
0	0	000	NUL (null)	32	20	040	#32; space	64	40	100	#64; @	96	60	140	#96; `
1	1	001	SOH (start of heading)	33	21	041	#33; !	65	41	101	#65; A	97	61	141	#97; a
2	2	002	STX (start of text)	34	22	042	#34; "	66	42	102	#66; B	98	62	142	#98; b
3	3	003	ETX (end of text)	35	23	043	#35; #	67	43	103	#67; C	99	63	143	#99; c
4	4	004	EOF (end of transmission)	36	24	044	#36; \$	68	44	104	#68; D	100	64	144	#100; d
5	5	005	ENQ (enquiry)	37	25	045	#37; %	69	45	105	#69; E	101	65	145	#101; e
6	6	006	ACK (acknowledge)	38	26	046	#38; &	70	46	106	#70; F	102	66	146	#102; f
7	7	007	BEL (bell)	39	27	047	#39; '	71	47	107	#71; G	103	67	147	#103; g
8	8	010	BS (backspace)	40	28	050	#40; (72	48	110	#72; H	104	68	150	#104; h
9	9	011	TAB (horizontal tab)	41	29	051	#41;)	73	49	111	#73; I	105	69	151	#105; i
10	A	012	LF (NL line feed, new line)	42	2A	052	#42; *	74	4A	112	#74; J	106	6A	152	#106; j
11	B	013	VT (vertical tab)	43	2B	053	#43; +	75	4B	113	#75; K	107	6B	153	#107; k
12	C	014	FF (NP form feed, new page)	44	2C	054	#44; ,	76	4C	114	#76; L	108	6C	154	#108; l
13	D	015	CR (carriage return)	45	2D	055	#45; -	77	4D	115	#77; M	109	6D	155	#109; m
14	E	016	SO (shift out)	46	2E	056	#46; .	78	4E	116	#78; N	110	6E	156	#110; n
15	F	017	SI (shift in)	47	2F	057	#47; /	79	4F	117	#79; O	111	6F	157	#111; o
16	10	020	DLE (data link escape)	48	30	060	#48; 0	80	50	120	#80; P	112	70	160	#112; p
17	11	021	DC1 (device control 1)	49	31	061	#49; 1	81	51	121	#81; Q	113	71	161	#113; q
18	12	022	DC2 (device control 2)	50	32	062	#50; 2	82	52	122	#82; R	114	72	162	#114; r
19	13	023	DC3 (device control 3)	51	33	063	#51; 3	83	53	123	#83; S	115	73	163	#115; s
20	14	024	DC4 (device control 4)	52	34	064	#52; 4	84	54	124	#84; T	116	74	164	#116; t
21	15	025	NAK (negative acknowledge)	53	35	065	#53; 5	85	55	125	#85; U	117	75	165	#117; u
22	16	026	SYN (synchronous idle)	54	36	066	#54; 6	86	56	126	#86; V	118	76	166	#118; v
23	17	027	ETB (end of trans. block)	55	37	067	#55; 7	87	57	127	#87; W	119	77	167	#119; w
24	18	030	CAN (cancel)	56	38	070	#56; 8	88	58	130	#88; X	120	78	170	#120; x
25	19	031	EM (end of medium)	57	39	071	#57; 9	89	59	131	#89; Y	121	79	171	#121; y
26	1A	032	SUB (substitute)	58	3A	072	#58; :	90	5A	132	#90; Z	122	7A	172	#122; z
27	1B	033	ESC (escape)	59	3B	073	#59; ;	91	5B	133	#91; [123	7B	173	#123; {
28	1C	034	FS (file separator)	60	3C	074	#60; <	92	5C	134	#92; \	124	7C	174	#124;
29	1D	035	GS (group separator)	61	3D	075	#61; =	93	5D	135	#93;]	125	7D	175	#125; }
30	1E	036	RS (record separator)	62	3E	076	#62; >	94	5E	136	#94; ^	126	7E	176	#126; ~
31	1F	037	US (unit separator)	63	3F	077	#63; ?	95	5F	137	#95; _	127	7F	177	#127; DEL



Text Hello, world!



Text can be numbers

Text Hello, world!

ASCII 72 101 108 108 111 44 32 119 111 114
 108 100 33



Text Hello, world!

ASCII 72 101 108 108 111 44 32 119 111 114
 108 100 33

Binary 01001000 01100101 01101100 01101100
 01101111 00101100 00100000 01110111
 01101111 01110010 01101100 01100100



0 becomes A or C
1 becomes T or G



Full example

12

Text Hello, world!

ASCII 72 101 108 108 111 44 32 119 111 114
 108 100 33

Binary 01001000 01100101 01101100 01101100
 01101111 00101100 00100000 01110111
 01101111 01110010 01101100 01100100

DNA AGCAGCCC ATTCCGAT CTTATTAC CTTCTTCC
 CGGAGGGG AATATGCC ACTACCCA ATGTATTT
 ATTCTTGT ATTTAATC CTGCGGAA CGTAAGCC



How we prepared the data

13

1. *Downloaded a couple of reference sequences from NCBI*
2. *Excised out a chunk each*
3. *Took a DNA sequence (the 'wildtype insert') and inserted it into one of the chunks*
4. *Made copies of the wildtype insert and encoded messages as nucleotide variants with respect to the wildtype*
5. *Inserted one of the resulting variant inserts into each of the remaining reference sequences*
6. *Generated simulated reads from all new references at different coverages (metagenomics: non-uniform coverages)*



Variant encoding

14

(This example encodes only one letter, the real messages are longer parts of sentences)

OLD WAY: Full length encoding (absolute)



A	G	C	A	G	C	C	C
---	---	---	---	---	---	---	---

$= 01001000 = H...$

TODAY: variant encoding (relative)



A	G	C	A	G	C	C	C
---	---	---	---	---	---	---	---

$= 01001000 = H...$



What you get

15

dna-encode.pl	Perl script to encode/decode text to/ from DNA
sh_end_1.fastq.gz sh_end_2.fastq.gz	Paired end read data from the mixed references, fastq-format, 2x250bp from 1000+/-50bp fragments (inner distance 500+/-50bp)
lo_end_1.fastq.gz lo_end_2.fastq.gz	Paired end read data from the mixed references, fastq-format, 2x150bp from 5300+/-500bp fragments (inner distance 5000+/-500bp)



dna-encode.pl

NAME

dna-encode – encode and decode ASCII text into DNA

SYNOPSIS

dna-encode [OPTIONS]... [FILE]...

DESCRIPTION

This script encodes a string of characters first into big endian (network order) binary and then into DNA where zero become A or C and one becomes G or T.

This implementation is based on the algorithm described in George M. Church, Yuan Gao, and Sriram Kosuri. Next-Generation Digital Information Storage in DNA. Science 2012. DOI: [10.1126/science.1226355](https://doi.org/10.1126/science.1226355).



dna-encode.pl

OPTIONS

- d, --decode
Decode a DNA sequence into a message rather than the default of encoding a message into DNA.
- l, --little-endian
Encode/decode characters using little endian byte order rather than the default big endian byte order.
- r, --reverse-complement
Reverse complement the DNA after encoding (or before decoding).
- verbose
Output intermediate binary when encoding/decoding.



It's a *metagenomic* sample.
Choose your tools accordingly.



After you identified an insert, you need to identify the insert *wildtype*. There are several ways to distinguish it from the variants. BLAST, consensus, pairwise similarities...



*NCBI Blast may be unreliable due to
the Government Shutdown.
If yes, try to use the public BLAST
server at EBI/EMBL (WU-BLAST).*



Are there any questions?

*Otherwise, the link to the read data
and the email address to send the
answer to will follow next.*



Read data (~10mb)

Can be obtained now at:

<http://schatzlab.cshl.edu/btg2013.tgz>

Answer (quote and author of quote) should be sent to:

beyondthegenome2013@gmail.com

Winner is announced today at about 4pm

