



# Beyond the Genome 2013

## *Informatics Challenge*

*Michael Schatz  
Sven-Eric Schelhorn, October 2013*



The prize

# 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.*



# The problem, 2

## 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

AAAS.ORG | FEEDBACK | HELP | LIBRARIANS

Science Magazine Enter Search Term SEARCH ADVANCED

WASHINGTON UNIV SCHL OF MED ALERTS ACCESS RIGHTS MY ACCOUNT SIGN IN

AAAS 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 Published Online August 16 2012 < Science Express Index

Science DOI: 10.1126/science.1226355 Read Full Text to Comment (0)

Full Text BREVIA

Full Text (PDF) Supplementary Materials

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

George M. Church<sup>1,2</sup>, Yuan Gao<sup>3</sup>, Sriram Kosuri<sup>1,2,\*</sup>

+ Author Affiliations

\* To whom correspondence should be addressed. E-mail: [sri.kosuri@wyss.harvard.edu](mailto:sri.kosuri@wyss.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

CANCER CRUSADE #40

Join Now!

ADVERTISEMENT

WEBINAR

Clinical Validation of Cancer Biomarker Signatures using Array Technology



# Numbers are bits

6

*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<sup>7</sup>-1))*





# Example

8

Text

Hello, world!



# Text can be numbers

9

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



*O becomes A or C  
I 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 ATTTAAC 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)

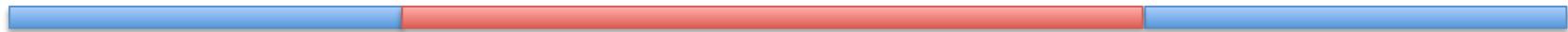


A G C A G C C C

= 01001000 = H...

TODAY: variant encoding (relative)

wild



var |



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.



# 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...



## Tip 3

20

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



# Questions?

21

*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](mailto:beyondthegenome2013@gmail.com)

*Winner is announced today at about 4pm*

