

GigaScience

Carbon-based archiving: the current progress and future prospects of DNA-based data storage

--Manuscript Draft--

Manuscript Number:	GIGA-D-18-00466	
Full Title:	Carbon-based archiving: the current progress and future prospects of DNA-based data storage	
Article Type:	Review	
Funding Information:	Guangdong Provincial Academician Workstation of BGI Synthetic Genomics (2017B090904014)	Dr. YUE SHEN
Abstract:	The information explosion has led to rapid increases in the amount of data to be physically stored. Yet one of the essential challenges is still looking for a better solution: How to store these large amounts of data in a space-efficient and stable way? DNA-based storage is a promising approach for long-term digital information storage as DNA holds great potentials because of its unique bio-properties. This review summarizes the state-in-art methods including digit-to-DNA coding schemes and the media types used in DNA storage to provide a general overview of the most recent progress achieved in this field.	
Corresponding Author:	YUE SHEN CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Zhi Ping	
First Author Secondary Information:		
Order of Authors:	Zhi Ping Dongzhao Ma Xiaoluo Huang Shihong Chen Longying Liu YUE SHEN Sha Joe Zhu	
Order of Authors Secondary Information:		
Additional Information:		
Question	Response	
Are you submitting this manuscript to a special series or article collection?	No	
Experimental design and statistics	Yes	
Full details of the experimental design and statistical methods used should be given		

<p>in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

1 **Abstract**

2 The information explosion has led to rapid increases in the amount of data to be physically
3 stored. Yet one of the essential challenges is still looking for a better solution: How to store
4 these large amounts of data in a space-efficient and stable way? DNA-based storage is a
5 promising approach for long-term digital information storage as DNA holds great potentials
6 because of its unique bio-properties. This review summarizes the state-in-art methods
7 including digit-to-DNA coding schemes and the media types used in DNA storage to provide
8 a general overview of the most recent progress achieved in this field.

9

10 **Keywords:**

11 DNA digital storage, Binary-DNA encoding scheme, *in vivo/in vitro* DNA digital storage

1 Introduction to DNA-based storage

2 The concept of DNA-based storage was initially introduced by computer scientists and
3 engineers in 1960s [1]. One of the pioneering attempts was made in 1988 by Joe Davis. At his
4 seminal art work – “Microvenus” [2], Davis converted an icon into a string of binary digits,
5 encoded them into a 28 base-pair (bp) synthetic DNA and later successfully sequenced to
6 retrieve the “icon” [2]. Although Microvenus was originally designed for interstellar
7 communications, it demonstrated that non-biological information could be also stored in
8 DNA.

9 Three unique bio-features making DNA the focus of the next generation of digital-
10 information storage. First, DNA is remarkably stable comparing to other storage media. With
11 its double-helix-structure and base stacking interaction, DNA can last for a thousand times
12 longer than a silicon device [3] and survive at harsh conditions over millennia [4,5,6,7].
13 Second, DNA possesses high storage density. Intuitively, each gram(g) of single-stranded
14 DNA can maintain data up to 455 exabytes [8]. As the storage strategy is continuously
15 optimized, scientists have already achieved a density that is very close to this theoretic limit.
16 Last but not least, the biological property of DNA provides access to natural reading and
17 writing enzyme which enables information stored in it remains accessible for millennia [8]. A
18 recent announced project called “the Lunar Library™ project” aims to make a DNA archive
19 with the collection of 10,000 images and 20 books for long-term backup storage on the Moon.
20 This showcase suggests the potential and advantage of DNA as a medium in long-term digital
21 storage.

22 The accessibility of DNA-based storage is mainly driven by two enabling techniques - DNA
23 synthesis and DNA sequencing [9], of which the former serves for “encoding” and the later
24 for “decoding”. Typically, digital information is first transcoded into “ATCG” sequence using
25 developed coding scheme. These sequences are then synthesized into oligo-nucleotides(oligos)

1 or long DNA fragments to allow long-term storage. To retrieve data, DNA sequencing is
2 applied to obtain the original “ATCG” sequence from synthesized DNA and so the
3 information stored in DNA.

4 **Overview of current coding schemes for DNA storage**

5 According to previous studies, we can summarize that an optimal coding scheme usually
6 outperformances in achieving: 1. High fidelity. In data retrieval, there is an obvious trade-off
7 between accuracy and redundancy. Hence, to strike a balance, appropriate coding scheme and
8 error correcting strategy are applied to avoid and to correct errors induced by DNA synthesis
9 or sequencing. 2. High coding efficiency. With four elementary bases, DNA has the
10 theoretical coding potential to store information in quaternary scaffold at least twice as much
11 as that of binary codes. 3. Flexible accessibility. From a computer science standpoint, data be
12 stored is expected to have random access. Correspondingly, all proposed coding schemes are
13 designed to fulfill the above features.

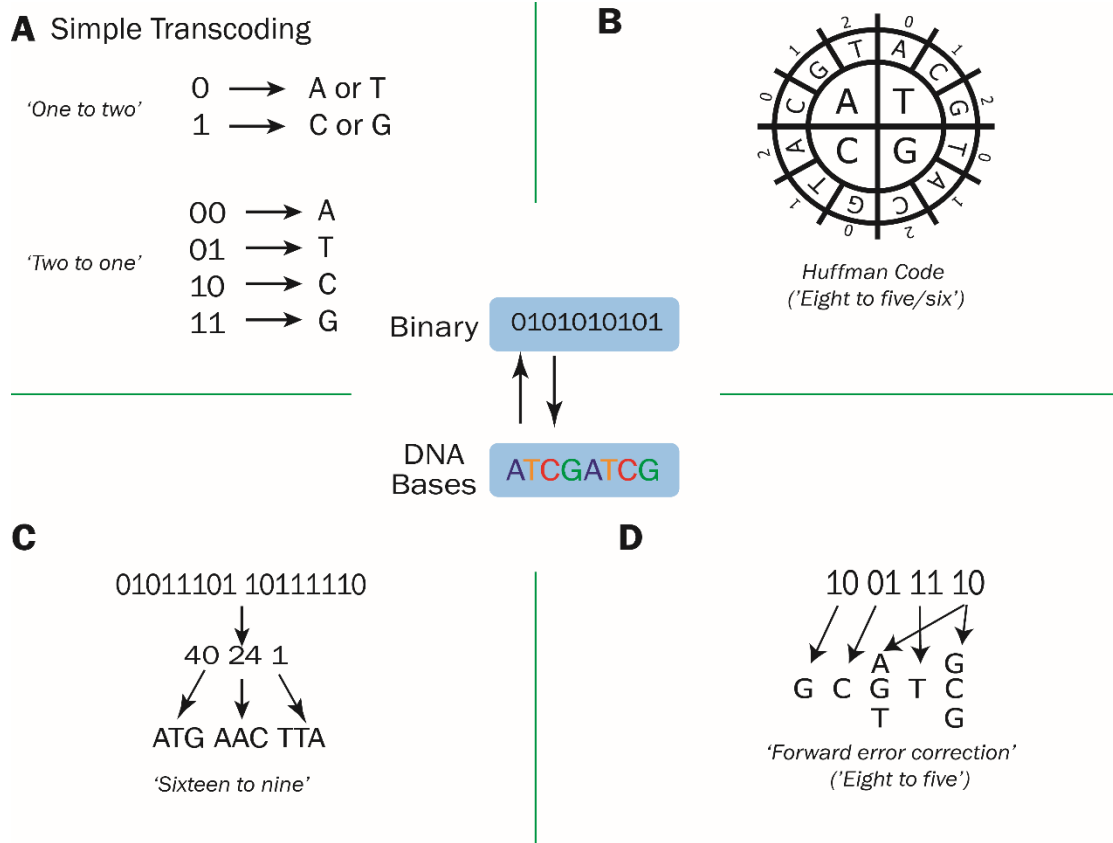


Figure 1 The different binary transcoding methods of reviewed DNA storage schemes. A) One binary bit is mapping to two optional bases [8]. B) Two binary bits are mapping to one fixed base [10]. C) Eight binary bits are transcoded through Huffman coding and then transcoded to five or six bases [11]. D) Two bytes (16 binary bits) are mapping to nine bases [12]. E) Eight binary bits are mapping to five bases [13].

• “Simple” code coding scheme

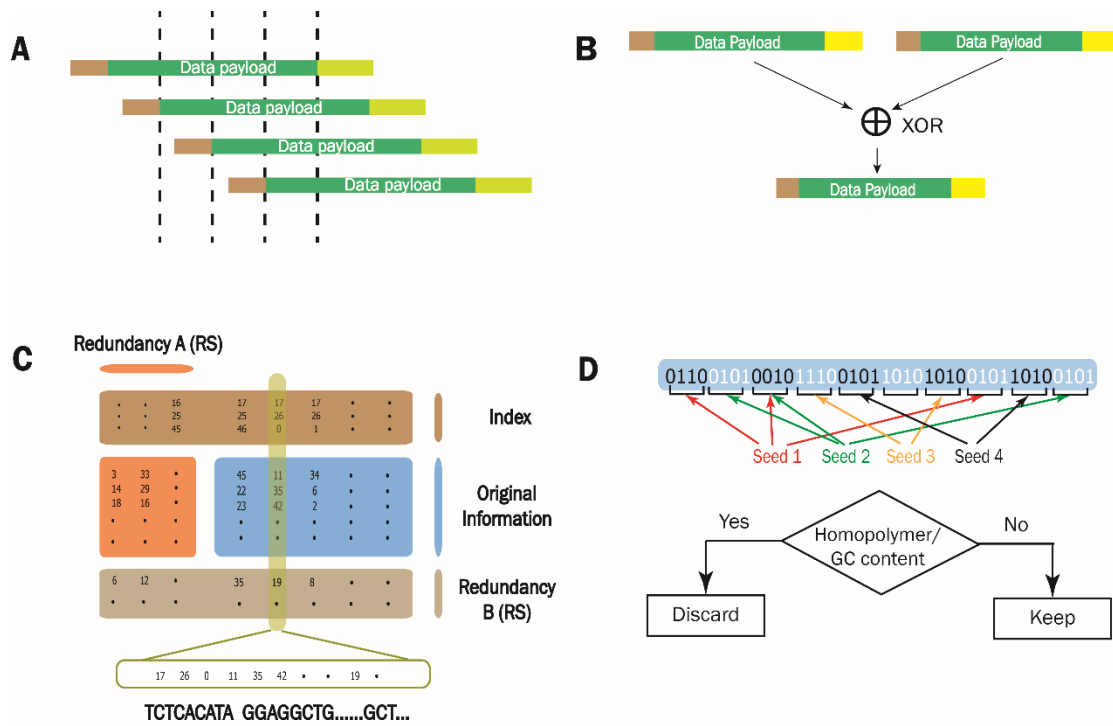
A “simple” code that aimed to tackle errors generated from DNA sequencing and synthesis (e.g. repeated sequences, secondary structure and abnormal GC content) was first proposed by Church et. al in 2012 [8]. By employing the free base swap strategy, Church and his colleagues encoded approximately 0.65 Mb data into ~8.8 Mb DNA oligos of 159 nt in length. It is considered as a milestone study in DNA storage given that large amount of digital data was successfully stored in DNA [14], which demonstrates the potential of DNA storage in coping with the challenge of information explosion. Yet, to allow its base swapping flexibility, this coding scheme sacrifices the information density - each binary code is transcoded into

1 one base (Fig. 1A). Researchers have later developed other coding strategies to overcome this
2 issue while maintain comparable performance.

3 • Huffman coding scheme

4 In 2013, Goldman and colleagues adopted Huffman code in their coding scheme, which
5 effectively improve the potential coding potential to 1.58 bits/nt [11]. Before transcribed into
6 DNA nucleotides, binary data was first converted into ternary Huffman code and then
7 transcribed to DNA sequence referred to a rotating encoding table (Fig. 1B). Every Byte of
8 the resulting data would be substituted by five or six ternary digits (comprises “0”, “1”, “2”
9 only), which can prevent generating mononucleotide repeats and compress the original data
10 by 25% to 37.5%. Besides, for ASCII text format files, compression further outperforms
11 by mapping the most common characters to five-digits ternary strings [11]. In addition, this
12 coding scheme employed simple parity-check coding for error detection and maintained a
13 four-fold coverage redundancy to prevent error and data loss (Fig. 2A). Nevertheless, it is
14 noted that the simple parity-check coding can only detect but not correct the errors and the
15 increased redundancy inevitably lower the coding efficiency.

1



2

3 Figure 2. The different redundancy types used in the reviewed DNA storage schemes: A) Increasing
 4 redundancy by repetition; B) Increasing redundancy by an exclusive-or (XOR) calculation; C)
 5 Increasing redundancy by using Reed-Solomon code for two rounds; D) Increasing redundancy by
 6 using fountain code.

7

8 • Improved Huffman coding scheme

9 In 2016, Bornholt et. al improved Goldman’s encoding scheme by an XOR encoding
 10 principle [12], which employed an exclusive-or (XOR, ‘ \oplus ’) operation to yield redundancy.

11 As shown in Fig. 2B, every two original sequences, A and B, will generate a redundant
 12 sequence C by $A \oplus B$. Therefore, with any two sequences (AB, AC or BC), one can easily
 13 recover the third sequence. Moreover, this coding scheme also provides the flexibility in
 14 providing redundancy according to the level of significance of particular data strands, namely
 15 “tunable redundancy”. This coding scheme successfully encoded 4 files with the total size of
 16 151 Kb and recovered 3 out of 4 files without manual intervention [12].

17 Moreover, the need of amplifying target files in large-scale database suggested the necessity
 18 of random-access in DNA storage. Therefore, in 2018, the same team put forward another

1 error-free coding scheme that allowed the users to randomly reach and recover individual files
2 in a large-scale system. In this coding scheme, unique polymerase chain reaction (PCR)
3 primers are assigned to individual files after rigorous screening, therefore, it allows users to
4 randomly access their target file(s). 200 Mb data was successfully stored and recovered in this
5 study, which set a new milestone for providing the feasibility of DNA storage of large scale
6 [13].

7 • A coding scheme based on Galois Field and Reed-Solomon Code
8 Focused on error detection and correction, a coding scheme based on Galois field and Reed-
9 Solomon (RS) code [14] was proposed by Grass and colleagues in 2015 [15]. Meanwhile, the
10 potential data density was improved to ~ 1.78 bits/nt. With the two-byte (8×2 bits) fundamental
11 information block, this coding scheme introduced a finite field (Galois field or GF) of DNA
12 nucleotide triplets as its elements (Fig. 1C). To prevent mononucleotide repeat > 3 nt during
13 encoding, the last two nucleotides of the triplet are varied, which can give 48 different triplets.
14 They indeed employed a GF(47) since 47 is the largest prime number smaller than 48 The
15 information block is then mapped to the three elements in GF(47), *i.e.* 256^2 to 47^3 . When
16 conducted error detection and correction, RS code was applied in this scheme. As shown in
17 Fig. 2C, two rounds of RS coding were applied horizontally and vertically to the matrix
18 generated by GF transcoding respectively.

19 In this pilot study, 83 kilobytes of text data were encoded *in silico* [15]. Although the data
20 size was not quite impressive it underlines the necessity of applying error-correction coding
21 and significantly enhances the coding efficiency.

22 • A “forward error correction” coding scheme
23 Blawat and colleagues proposed a coding scheme focusing on tackling errors generated from
24 DNA sequencing, amplification and synthesis (*e.g.* insertion, deletion and swapping). The
25 potential coding density was 1.6bits/nt. Two reference coding tables are specified in advance.
26 The one-byte (8 bits) fundamental information block is assigned to a 5 nt DNA sequence and

1 the 3rd and 4th nucleotide are swapped (Fig. 1D). Two other criteria are applied to prevent
2 mononucleotide repeat during this process: 1) the first three nucleotides should not be the
3 same; 2) the last two nucleotides should not be the same. Consequently, an 8-bits data block
4 (i.e. $2^8 = 256$ permutations for binary data) is transcoded into 704 different DNA blocks (4^5 -
5 4^3 - 4^4) [16]. They can be categorized into three clusters: clusters A & B of complete blocks,
6 256 per each and cluster C of 192 incomplete blocks. Data can then be mapped to DNA
7 blocks A and B as required, e.g. alternately mapped to A or B.

8 In this study, 22Mbytes of data were successfully encoded and stored in an oligo pool. The
9 data had been retrieved with no error, which proved the feasibility of this coding scheme. Yet,
10 this is not the case for detecting and correcting single-mutation. For example, “11100011”
11 could be mapped to a DNA block “TGTAG”. However, if an A-to-T transversion occurs, the
12 DNA block will change to “TGTTG”, which will give an error byte “11101111” after
13 decoding.

14 • Fountain code-based DNA storage coding scheme

15 In 2017, Erilich and Zielinski employed fountain code in their coding scheme [17]. Fountain
16 code is a widespread coding method of information communication system known for its
17 robustness and high efficiency [18]. Fountain code is also known as a rateless erasure code, in
18 which data to be stored is divided into k segments, namely resource packets. Potentially
19 limitless number of encoded packets can then derive from the resource packets. When it
20 returns n ($n > k$) encoded packets, the original resource data should be perfectly recovered. In
21 practice, n only need to be slightly larger than k to yield e great coding efficiency as well as
22 robustness for information communication [19].

23 Similarly, binary data-nucleotide sequence encryption is carried out. A fundamental two-bit to
24 one-nucleotide transcoding table is adopted, in which [00, 01, 10, 11] mapped to [A, C, G, T],
25 respectively (Fig. 1A). At first, original binary information is segmented to small blocks.
26 These blocks are chosen according to a pre-designed pseudorandom sequence of numbers. A

1 new data block is then created by bitwise addition of the selected blocks with random seeds
2 attached and transcoded to nucleotide blocks according to the transcoding table.

3 Mononucleotide repeats and abnormal GC content are prevented by a final verification (Fig.
4 2D) [17].

5 The oligos in this coding scheme are correlated and have grid-like topology to realize
6 extremely low but necessary redundancy. This study enables the theoretical limit of coding
7 potential unprecedentedly high, reaching 1.98bits/nt and remarkably reduces the requested
8 redundancy for an error-free recovery of source file. Moreover, the mechanism of random
9 selection and validity verification ensures that long single-nucleotide homopolymers would
10 not appear in the encoded sequence. However, in this coding scheme, the complexity level of
11 encoding and decoding is not linearly correlated to the data size. Thus, decoding could be
12 complicated and may require more resource and longer time for computation. However,
13 although it is claimed that a 4% loss of total packets would not affect the recovery of original
14 file in the report, in terms of the features of DNA Fountain code, loss of more packet may
15 cause the complete failure of recovery. If it is aiming to store for permanent preservation, the
16 amount of redundancy must be raised to ensure the information integrity when encounter
17 spoiled oligo pools.

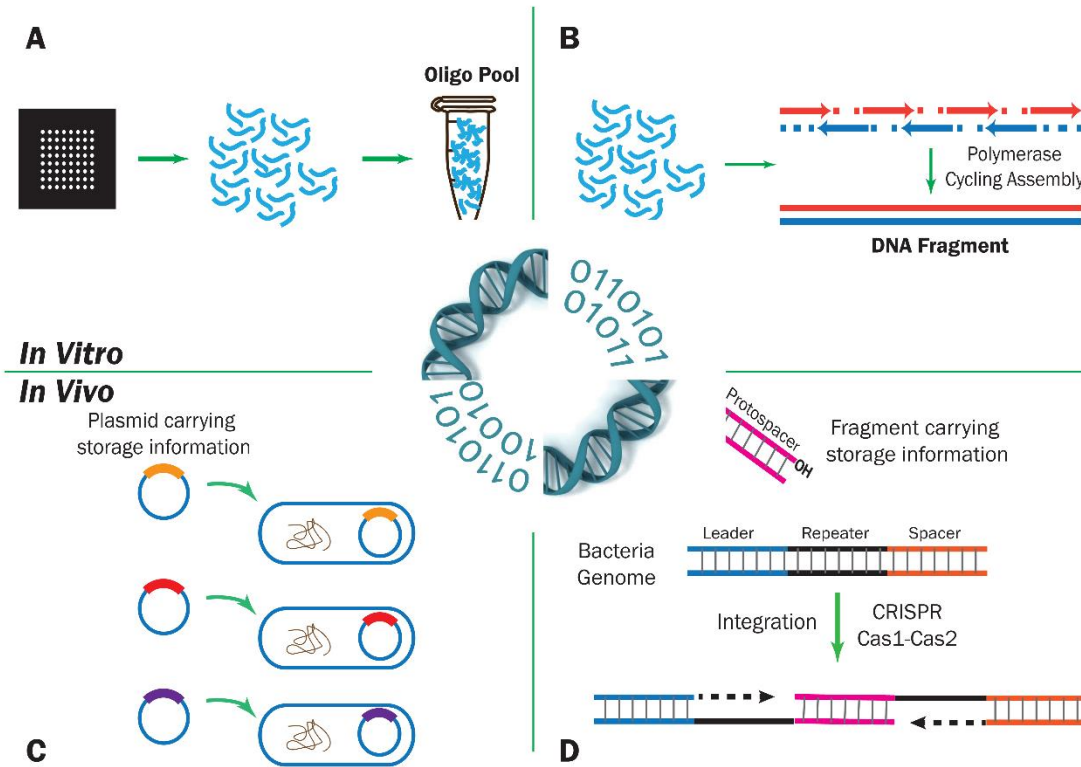
18 If DNA storage could be viewed as merely a storage process with high fidelity, DNA fountain
19 coding is the present only communication-based coding scheme. In DNA data storage and
20 retrieval, the most common error is caused by single nucleotide mutation. To address this
21 issue, most coding scheme will create high redundancy in order to tackle the mal-condition of
22 current communication channels, however, these error correction algorithms require complex
23 decoding procedures and much computing time. Here, fountain-coding scheme firstly show
24 that it is unnecessary to employ error detection/correction algorithms, which provide us an
25 alternative solution towards improving the performance of DNA coding.

26

1 **Overview of current media for DNA storage:**

2 Current DNA storage employed different media to store the encoded DNA sequences. In sum,
3 there are two types: *in vivo* and *in vitro*.

4



5

6 Figure 3. Two categories of DNA storage application. Panel A) and B) demonstrate the two ways of in
7 vitro DNA storage; panel C) and D) demonstrate two ways of in vivo DNA storage. A) Chip-based
8 high throughput DNA oligo analysis. DNA oligos carrying digital information are stored in the form of
9 oligo pool. B) DNA fragments synthesized by polymerase cycling assembly (PCA), the fragments will
10 carry the information to be stored. C) Digital information inserted into plasmid and then the plasmids
11 are transferred into bacteria cells. D) DNA fragments carrying digital information is inserted into
12 bacteria genome by employing CRISPR system using Cas1-Cas2 integrase.

13 ***In vivo* DNA storage**

14 *In vivo* DNA storage is commonly adopted in the pioneer works of DNA storage, such as the
15 *Microvenus* project, which used bacteria as the storage medium. Typically, encoded DNA
16 sequences are first cloned into plasmid and then transferred into the bacteria. Therefore, the

1 DNA sequences and so does the information it carries can be maintained in the tiny bacteria
2 and their billions of descendants.

3 Nevertheless, the capacity of bacteria for carrying plasmid is limited by the type of plasmids
4 and their corresponding size. In addition, the mutation of plasmid in bacteria is quite common.
5 During bacteria replication, the spontaneous mutation may ultimately alter the information
6 stored in them after a few years.

7 Recently, Church *et. al* demonstrated a novel method to encode an image and a short movie
8 clip into the bacteria genome using the CRISPR-Cas system with Cas1-Cas2 integrase [20].
9 Although it is reported that CRISPR-Cas system is not equally efficient to all the sequences,
10 this work greatly improved the capability of *in vivo* DNA storage.

11 ***In vitro* DNA storage**

12 Apart from *in vivo* DNA storage, *in vitro* DNA storage is more frequently seen in recent
13 studies. One of the most popular form is oligo library. This is largely due to the maturation of
14 chip-based high-throughput oligo synthesis technique [21], making the synthesis of large
15 amount of DNA oligos more cost-effective.

16 When synthesis, each oligo is given a short tag, or index, as all the oligos would be
17 completely mixed for high throughput synthesis and sequencing. Current oligo synthesis
18 technique is able to generate at most 200-mers in relatively high accuracy and purity [22].
19 Hence, the index should be as short as possible to save the information capacity in each oligo.
20 Apparently, much more indices will be needed if more DNA oligo sequences are generated
21 and mixed. However, similar to *in vivo* DNA storage, the larger the data size is, the more
22 DNA oligos is demanded for *in vitro* DNA storage, which will increase the size of indices in
23 oligo and thus lower the storing capacity and efficiency.

24 Alternatively, longer DNA fragments can be used instead of DNA oligos to avoid these
25 problems. In 2017, Yadzi *et. al* successfully encoded 3633 bytes of information (two images)

1 into 17 DNA fragments and recovered the image using homopolymer error correction [23].

2 Nevertheless, the current cost of DNA fragment synthesis is higher than that of oligo
3 synthesis, which increases the overall cost of DNA fragment-based storage.

4 Some other pioneer work also goes beyond our aforementioned DNA storage system. Song
5 and Zeng proposed a strategy which is claimed to be able to detect and correct error in each
6 byte [24]. They transformed short message into *E.coli* stellar competent cells and proved the
7 reliability of their strategy. Lee *et. al* incorporated enzymatic DNA synthesis and DNA
8 storage principles, reported an enzymatic-based DNA storage strategy [25]. All these
9 researches laid a sound foundation for world-wide application of this novel storage medium.

10

11 **Challenges of DNA-based storage**

12 **Limited size of synthetic DNA**

13 As mentioned above, information encoding in DNA depends on DNA synthesis. Based on the
14 final product size, DNA synthesis includes oligo synthesis (≤ 200 mer) and gene synthesis
15 (200-3,000 bp or above), while DNA oligos usually serve as basic building blocks for gene
16 synthesis. For cost saving purpose and to reduce complexity of DNA synthesis, primary
17 storage unit size is often limited below 200nt [21].

18 Due to this limit, information needs to be fragmented and indexed before encoded into DNA
19 to allow oligo synthesis (encoding) and pool sequencing (decoding) to reconstruct data in the
20 correct order. Thus, when the amount of information grows, not only the number of fragments
21 increase, but the indexing information also accumulates subsequently. Except for optimizing
22 the index length (see “DNA storage in beyond” below), techniques for synthesizing longer
23 oligo are considered to be the major challenge before we can push the envelope.

24 **DNA sequencing-induced errors**

1 Currently, there are two major types of DNA sequencing techniques: real-time, single-
2 molecule sequencing and massively parallel (or next generation) sequencing. The latter is a
3 high-throughput sequencing method and is dominant for short-read (<700bp, depending on
4 platform) sequencing while the former is on the opposite [9,26].

5 In DNA storage, massively parallel sequencing is widely used for data retrieval ever since it's
6 firstly employed by Church *et al.* in 2012. Two main reasons can explain this prevalence.
7 First, the length of the synthetic DNA generated from encoding is relatively short, which is
8 more cost-effective to sequence with massively parallel sequencing. Second, the throughput
9 and accuracy (~99.9%) of massively parallel sequencing still far surpass its counterparts [9].

10 However, this technique also comes with limitation. Most massively parallel sequencing
11 platforms require *in vitro* template amplification with primers to generate a complex template
12 library for sequencing. During this process, copying errors, sequence-dependent biases (for
13 example, in high- and low-GC regions and at long mononucleotide repeats) and information
14 loss (for example, methylation) are produced [9].

15 Nevertheless, sequencing with minimal biases and random errors in respect to accuracy and
16 contiguity is possible given that rapid progress is now achieved in real-time, single-molecule
17 sequencing. It is reported that this rising technique can tolerant high GC content and only
18 generates random errors [27], which is ideal in data retrieval. When it can also achieve high-
19 fidelity, the storing potential of DNA may be further unlocked.

20

21 **DNA storage in beyond**

22 In spite of all progresses been achieved, current DNA-based storage is still at the early stage
23 of its substantial applications. Moreover, its development will necessarily benefit from the
24 progress of the coding/decoding methods.

1 It could be foreseen that in the near future, DNA oligo synthesis could break the limitation of
2 200-mers, providing us longer primary storage unit. This will clearly improve the net coding
3 efficiency with same length of PCR primer and shorter index sequences. A simulation was
4 performed for DNA storage of 1GB file under theoretical limitation, i.e. one DNA base would
5 represent two binary bits. For each DNA oligo, the length of forward and reverse primers was
6 set as 20. Therefore, we could get:

$$i + d + 20 \times 2 = l$$

Equation (1)

9 where i is the length of index, d is the length of data payload and l is the length of DNA oligo.

10 As the file size is 1GB, we could get:

$$d \times 2^i = \frac{1024^3 \times 8bits}{2}$$

Equation (2)

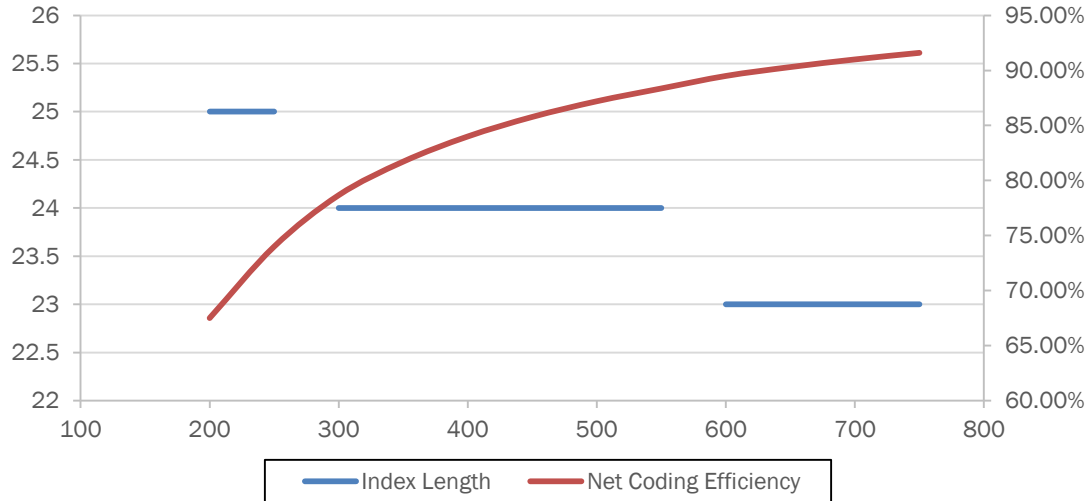
13 With combination of equation (1), therefore:

$$\log_2(l - 40 - i) + i = 32$$

Equation (3)

16 Hence, we could get an optimal index length with fixed DNA oligo length.

17 As Figure 4 shows, with the increasing of DNA oligo length, the index length decreases while
18 net coding efficiency increases. This calculation indicates that the efficiency of DNA storage
19 could be remarkably improved with the improvement of DNA oligo synthesis techniques.



1

2 Figure 4 A simulation of net coding efficiency for DNA storage of 1GB file. The x-axis represents the
 3 length of oligo to be synthesized, y-axis (left) represents the minimum length of index needed to record
 4 the information coordinates, y-axis (right) represents the net coding efficiency.

5 In addition, the scale of DNA synthesis also affects the information capacity of DNA storage
 6 per unit mass. High-throughput oligo synthesis is currently directed to microscale level with
 7 the development of chip-based DNA synthesis technology. In DNA storage, the information
 8 capacity of certain mass of DNA sequences also relates to the copy number of each DNA
 9 molecule. To date, the copy number of oligos is around 10^7 molecules in microchip high
 10 throughput synthesis without dilution according to Erlich *et. al* [17], which will give an
 11 information capacity level at $\sim 10^{13}$ bytes/g according to Equation (4). If the copy number
 12 decreased to 10^4 molecules per oligo, the information capacity will increase to $\sim 10^{16}$ bytes/g.
 13 Additionally, synthesis in microscale also reduces the cost by several orders of magnitude and
 14 saves the step of dilution.

15

$$C = n \times (N_m \mu \delta \gamma)^{-1}$$

16

Equation (4)

17 where C represents the information capacity; n represents the number of bytes each oligo
 18 carries, normally 10 – 20 bytes/molecule according to different coding schemes; N_m is the

1 number of molecules, μ is the number of nucleotides per molecule, δ is 320 Dalton/nucleotide;
 2 γ is 1.67×10^{-24} g/Dalton,

3 On the other hand, development of sequencing technique also significantly effects DNA
 4 storage. By summarizing the frequent-used sequencing platforms in DNA storage, we noticed
 5 that accuracy and cost of sequencing are no longer the only considerations for DNA storage.
 6 Portable, yet error-prone sequencing platforms like Oxford Nanopore MinION is gaining
 7 attention due to its potential for high-compactness and stand-alone DNA data storage systems
 8 [13]. This trend will be largely attributed to the emergence of a growing number of error-
 9 tolerant coding schemes, which enable recovery of data even using the error-prone
 10 sequencing platform. If the accuracy of these sequencing platforms can be further improved, a
 11 paradigm shift from next-generation sequencing to third generation sequencing may
 12 eventually take place in DNA storage.

13
 14

Platform	Error Rate	Runtime	Instrument Cost(US\$)	Cost per Gb (US\$)	Study
Illumina MiSeq	0.1%	21-56h*	\$99K	\$110-250	[12]Bornhol et al.,2016/[15]Grass et al.,2015/[17]Erlich Y and Zielinski D,2017/[20]Shipman et al.,2017
Illumina HiSeq 2000	2.0%	3-10d *	\$654K	\$41	[8]Church et al.,2012/[11]Goldman et al.,2013
Illumina HiSeq 2500	0.1%	7h-11d *	\$690	\$30-250	[16]Blawat et al.,2016
Illumina NextSeq	0.1%	<3d	\$1K	\$7	[13]Organick et al.,2018
Oxford Nanopore MinION	12.0%	up to 48h	\$1K	\$750	[13]Organick et al.,2018/[23]Yazdi et al.,2017

d, days; Gb, gigabase pairs; h, hours; K, thousand; * varied by read length and version of reagent kit

15
 16

Table 1. Summary of frequent-used sequencing platforms in DNA storage (data retrieved from [26]).

17 Taken together, DNA storage provides us the possibility to manipulate DNA as a carbon-
 18 based archive with an excellent storage density and stability. By combining the within-
 19 reached development of DNA synthesis and sequencing techniques, DNA might eventually

1 transform into the next generation digital information storage media with ideal volatility,
2 capacity, accessibility and reliability.

3 **Acknowledgement**

4 This work was supported by Guangdong Provincial Academician Workstation of BGI
5 Synthetic Genomics (No. 2017B090904014).

6

7

8

9

10

11

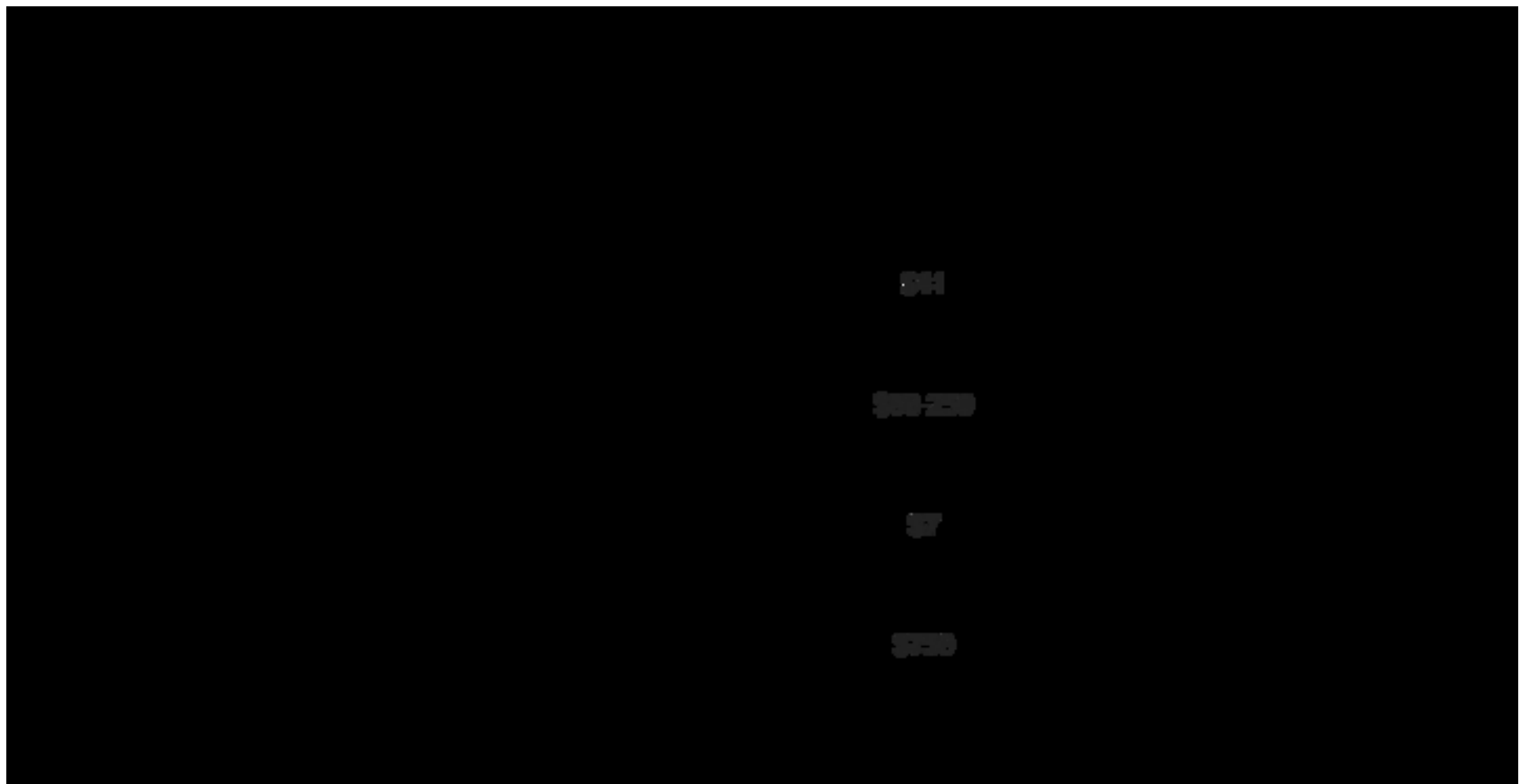
12 **References**

- 13 1. Neiman MS. Some fundamental issues of microminiaturization. *Radiotekhnika*.
14 1964;No. 1:3-12.
- 15 2. Joe Davis a. Microvenus. *Art Journal*. 1996; 1:70. doi:10.2307/777811.
- 16 3. Bonnet J, Colotte M, Coudy D, Couallier V, Portier J, Morin B, *et al*. Chain and
17 conformation stability of solid-state DNA: implications for room temperature storage.
18 *Nucleic Acids Research*. 2010;38 5:1531-46. doi:10.1093/nar/gkp1060.
- 19 4. Pääbo S, Poinar H, Serre D, Jaenicke-Després V, Hebler J, Rohland N, *et al*.
20 GENETIC ANALYSES FROM ANCIENT DNA. *Annual Review of Genetics*.
21 2004;38:645-79. doi:10.1146/annurev.genet.37.110801.143214.
- 22 5. Kool ET. Hydrogen bonding, base stacking, and steric effects in DNA replication.
23 *Annual Review of Biophysics & Biomolecular Structure*. 2001;30 1:1-22.

- 1 6. Nelson DL, Cox MM and Lehninger AL. *Lehninger principles of biochemistry*. New
2 York ; Basingstoke : W.H. Freeman, c2008. 5th ed.; 2008
- 3 7. Pierce BA. *Genetics : a conceptual approach*. New York, NY : W.H. Freeman, c2012.
4 4th ed., International ed.; 2012.
- 5 8. Church GM, Gao Y and Kosuri S. Next-generation digital information storage in
6 DNA. *Science*. 2012;337 6102:1628. doi:10.1126/science.1226355.
- 7 9. Shendure J, Balasubramanian S, Church GM, Gilbert W, Rogers J, Schloss JA, et al.
8 DNA sequencing at 40: past, present and future. *Nature*. 2017;550 7676:345-53.
9 doi:10.1038/nature24286.
- 10 10. De Silva PY and Ganegoda GU. New Trends of Digital Data Storage in DNA.
11 *Biomed Res Int*. 2016;2016:8072463. doi:10.1155/2016/8072463.
- 12 11. Goldman N, Bertone P, Chen S, Dessimoz C, LeProust EM, Sipos B, et al. Towards
13 practical, high-capacity, low-maintenance information storage in synthesized DNA.
14 *Nature*. 2013;494 7435:77-80. doi:10.1038/nature11875.
- 15 12. Bornholt J, Lopez R, Carmean DM, Ceze L, Seelig G and Strauss K. A DNA-Based
16 Archival Storage System. *SIGPLAN Not*. 2016;51 4:637-49.
17 doi:10.1145/2954679.2872397.
- 18 13. Organick L, Ang SD, Chen YJ, Lopez R, Yekhanin S, Makarychev K, et al. Random
19 access in large-scale DNA data storage. *Nat Biotechnol*. 2018;36 3:242-8.
20 doi:10.1038/nbt.4079.
- 21 14. Reed I and Solomon G. Polynomial Codes Over Certain Finite Fields. *Journal of the*
22 *Society for Industrial and Applied Mathematics*. 1960;8 2:300-4.
23 doi:10.1137/0108018.
- 24 15. Grass RN, Heckel R, Puddu M, Paunescu D and Stark WJ. Robust chemical
25 preservation of digital information on DNA in silica with error-correcting codes.
26 *Angew Chem Int Ed Engl*. 2015;54 8:2552-5. doi:10.1002/anie.201411378.

- 1 16. Blawat M, Gaedke K, Hütter I, Chen X-M, Turczyk B, Inverso S, et al. Forward Error
2 Correction for DNA Data Storage. *Procedia Computer Science*. 2016;80:1011-22.
3 doi:<https://doi.org/10.1016/j.procs.2016.05.398>.
4
5
6 17. Erlich Y and Zielinski D. DNA Fountain enables a robust and efficient storage
7 architecture. *Science*. 2017; 6328:950. doi:10.1126/science.aaj2038.
8
9
10 18. Byers JW, Luby M, Mitzenmacher M and Rege A. A digital fountain approach to
11 reliable distribution of bulk data. *Proceedings of the ACM SIGCOMM '98 conference*
12 *on Applications, technologies, architectures, and protocols for computer*
13 *communication*. Vancouver, British Columbia, Canada: ACM, 1998, p. 56-67.
14
15
16 19. MacKay DJ. Fountain codes. *IEE Proceedings-Communications*. 2005;152 6:1062-8.
17
18
19 20. Shipman SL, Nivala J, Macklis JD and Church GM. CRISPR-Cas encoding of a
20 digital movie into the genomes of a population of living bacteria. *Nature*. 2017;547
21 7663:345-9. doi:10.1038/nature23017.
22
23
24 21. Kosuri S and Church GM. Large-scale de novo DNA synthesis: technologies and
25 applications. *Nature methods*. 2014;11 5:499-507. doi:10.1038/nmeth.2918.
26
27
28 22. Ma S, Tang N and Tian J. DNA synthesis, assembly and applications in synthetic
29 biology. *Curr Opin Chem Biol*. 2012;16 3-4:260-7. doi:10.1016/j.cbpa.2012.05.001.
30
31
32 23. Yazdi S, Gabrys R and Milenkovic O. Portable and Error-Free DNA-Based Data
33 Storage. *Sci Rep*. 2017;7 1:5011. doi:10.1038/s41598-017-05188-1.
34
35
36 24. Song L and Zeng AP. Orthogonal Information Encoding in Living Cells with High
37 Error-Tolerance, Safety, and Fidelity. *ACS Synth Biol*. 2018;7 3:866-74.
38
39
40
41
42
43
44
45
46
47
48
49 25. Lee HH, Kalhor R, Goela N, Bolot J and Church GM. Enzymatic DNA synthesis for
50 digital information storage. *bioRxiv*. 2018; doi:10.1101/348987.
51
52
53
54 26. Goodwin S, McPherson JD and McCombie WR. Coming of age: ten years of next-
55 generation sequencing technologies. *Nat Rev Genet*. 2016;17 6:333-51.
56
57
58
59
60
61
62
63
64
65

- 1 27. Ross MG, Russ C, Costello M, Hollinger A, Lennon NJ, Hegarty R, *et al.*
2 Characterizing and measuring bias in sequence data. *Genome Biol.* 2013;14 5:R51.
3
4 3 doi:10.1186/gb-2013-14-5-r51.
5
6 4 28. van Dijk EL, Jaszczyszyn Y, Naquin D and Thermes C. The Third Revolution in
7 Sequencing Technology. *Trends Genet.* 2018;34 9:666-81.
8
9 5
10 6
11 doi:10.1016/j.tig.2018.05.008.
12
13
14 7
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65



A Simple Transcoding

'One to two'

0 → A or T
1 → C or G

'Two to one'

00 → A
01 → T
10 → C
11 → G

Binary 0101010101

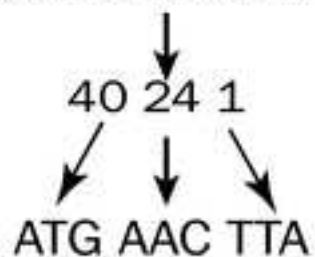
DNA Bases ATCGATCG

B

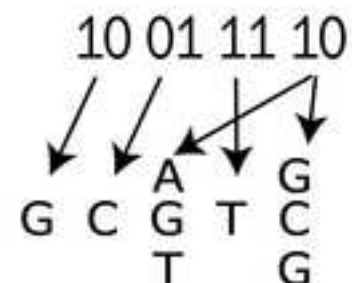
Huffman Code
'Eight to five/six'

C

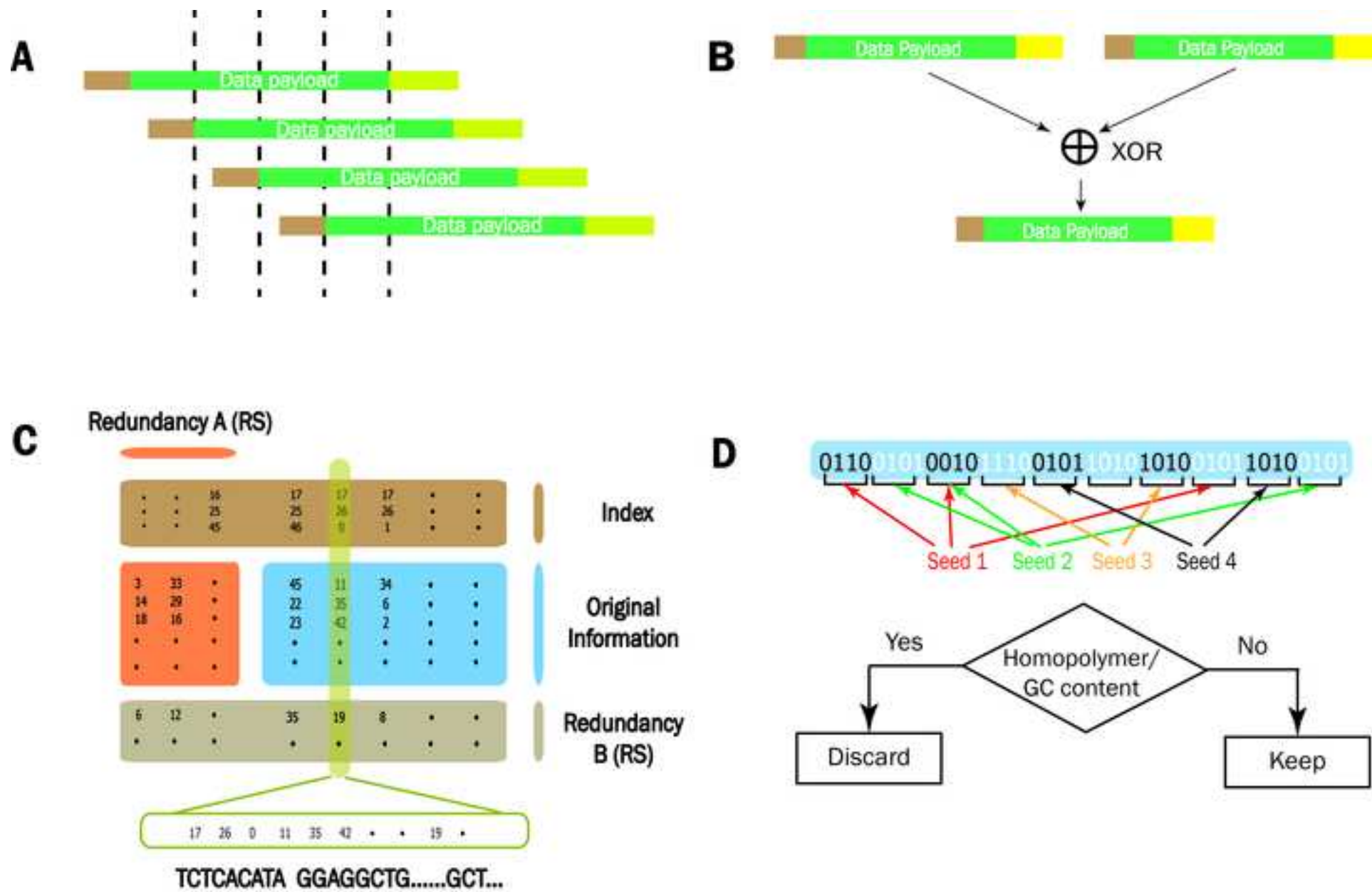
01011101 10111110



'Sixteen to nine'

D

'Forward error correction'
'Eight to five'



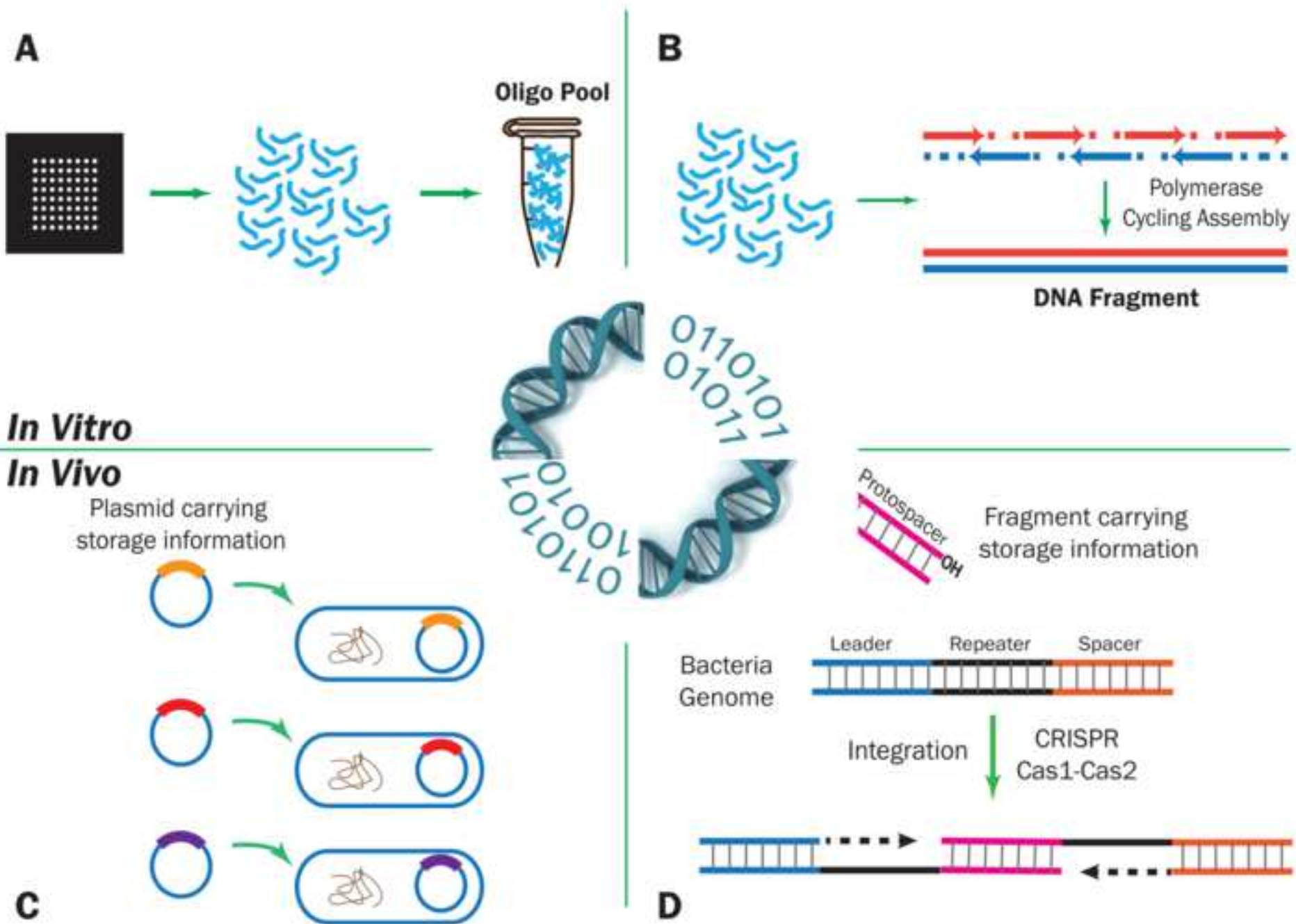
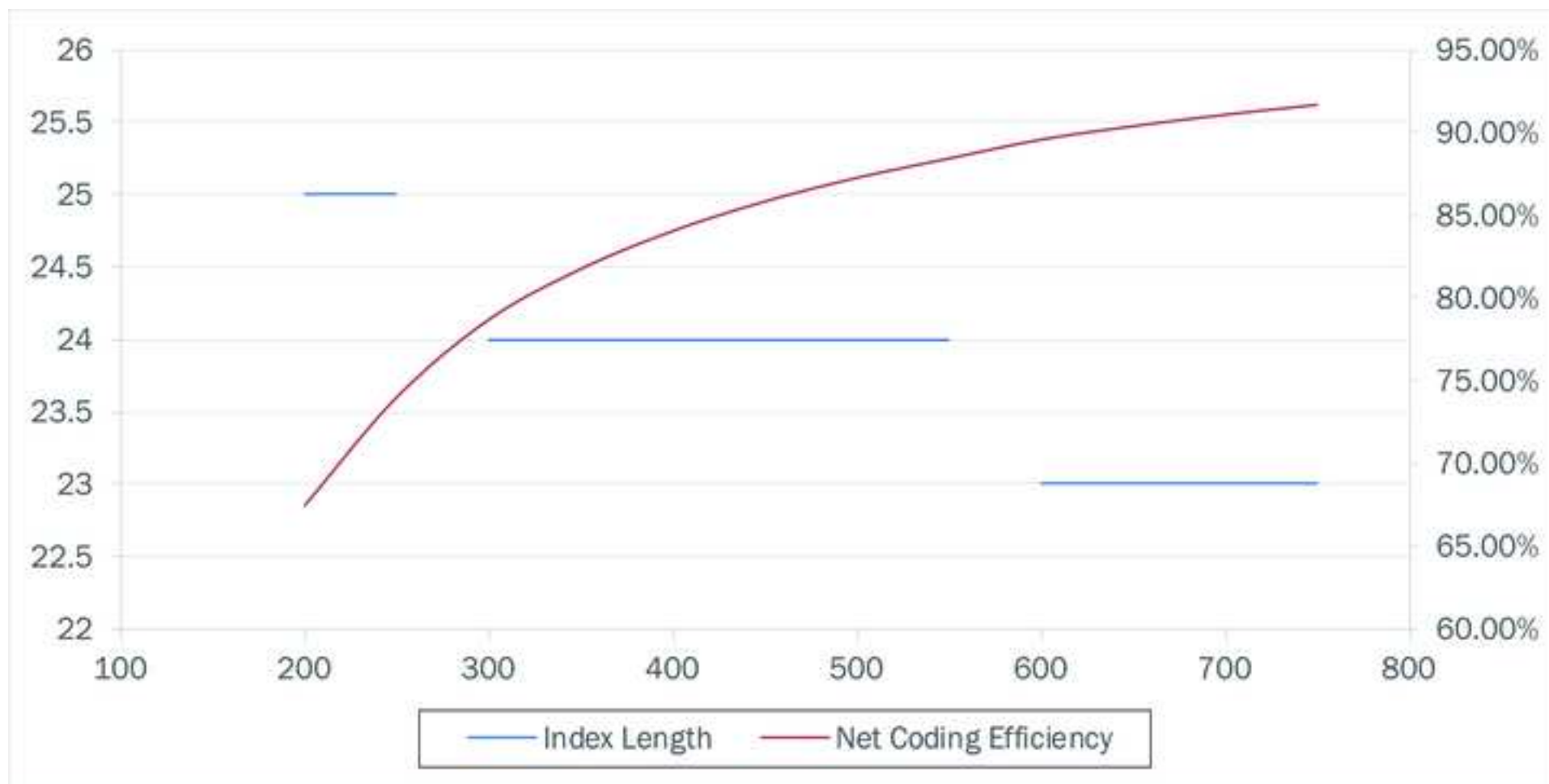


Figure 4



Dear Dr. Edmunds:

We submit our manuscript entitled "Carbon-based archiving: the current progress and future prospects of DNA-based data storage" to GigaScience for publication.

This manuscript is a timely review of DNA-based storage with focus on coding scheme and media type. We provide scalable measurements and technical opinions of this field, which we believe will be a great add on to people's current understanding and help promote its better development. As DNA-based storage is a promising bio-approach for large scale digital information storage, we consider it is well in scope of the GigaScience's publication criteria.

All authors have read and have abided by the publication ethics as set out by the Commission on Publication Ethics (COPE) for manuscripts submitted to GigaScience.

All authors declared that they have no conflicts of interest to this work.

The work described has not been submitted elsewhere for publication, in whole or in part, and all the authors listed have approved the manuscript that is enclosed.

Thank you very much for your attention and consideration.

Yours sincerely,

Yue (Chantal) Shen

Sha Joe Zhu