

1 Low bias DNA for sustainable and efficient data storage

2 Yanmin Gao^{1,2}, Xin Chen³, Jianye Hao⁴, Chengwei Zhang⁵, Hongyan Qiao^{1,2} and Hao Qi^{1,2*}

3 ¹ School of Chemical Engineering and Technology, Tianjin University, Tianjin, China.

4 ² Key Laboratory of Systems Bioengineering (Ministry of Education), Tianjin University,
5 Tianjin, China.

6 ³ Center for Applied Mathematics, Tianjin University, Tianjin, China.

7 ⁴ College of intelligence and computing, Tianjin University, Tianjin, China.

8 ⁵ College of Information Science and Technology, Dalian Maritime University, Dalian, China.

9 * Correspondence should be addressed to H.Q. (haoq@tju.edu.cn)

10

11 **Abstract:**

12 In DNA data storage, the huge sequence complexity is challenging for repeatable and efficient
13 information reading from massive DNA oligo pool. Here, we demonstrated that synthetic oligo
14 pool comprising over ten thousand strands was largely skewed by PCR process due to its inherent
15 mechanism of inefficient priming, product-as-template, error-spreading prone, which caused
16 serious oligos dropout, over 50% oligos lost in repeated PCR amplification, much more fatal than
17 base mutant for correct information retrieve. Therefore, we developed a new biochemical
18 framework isothermal DNA reading (iDR) from large-scale oligo pool normalization (OPN) and
19 isothermal amplification. Due to its priming-free and error-spreading proof, it achieved low-biased
20 and stable amplification even for successive 10 times deep-reading without performance lost, the

21 first repeatable DNA storage. Furthermore, the skewed oligo pool with uneven molecule copy
22 number of each oligos was rectified by a total 1280 chemically synthesized OPN-probes, the
23 largest scale oligo normalization so far, and the necessary amount of sequencing reads for perfect
24 oligos retrieve was further largely decreased. These advanced features enable the iDR chemical
25 framework being ideal for manipulating the huge sequence complexity and building the
26 sustainable and efficient DNA storage “hardware”.

27

28 **Introduction:**

29 In DNA data storage, technology originally developed for bioengineering approach including array
30 oligo synthesis, PCR and DNA sequencing were integrated to construct the “hardware” for DNA
31 storage¹⁻⁵. Array synthesized DNA oligo pool comprising of from thousands up to millions of oligo
32 strands with several hundred bases in length has been utilized in many advanced bioengineering
33 applications, e.g., probe blend, DNA origami assembly, and genome synthesis¹. Due to both the
34 location on microchip and DNA sequence difference⁵, the copy number for each oligo strand from
35 array synthesis is distinct. Furthermore, the sheer number of synthesized oligos from array on
36 microchip is very small, roughly from 10^5 to 10^{12} at the concentration of femtomolar depending
37 on the synthesis platform^{6, 7}. Generally, amount of a few hundred nanograms DNA at the
38 concentration of micromolar is necessary for high quality DNA sequencing covering all oligos in
39 the pool on commercial high throughput DNA sequencing platform, e.g., Illumina⁸. Therefore,
40 amplification is crucial to boost the signal for the subsequent DNA sequencing for reading data.

41 The oligo pool size for DNA data storage is much larger at least several orders of magnitude than
42 that in other bioengineering applications⁹⁻¹¹. Furthermore, the unevenness of copy number
43 generated a huge complexity that caused serious problem for DNA molecule retrieval and data
44 decoding^{4, 5, 12, 13}. In current reported systems, the information reading was achieved by PCR
45 amplification and NGS, but the copy number unevenness originally stemmed from microchip
46 synthesis was further skewed by highly biased amplification process^{5, 12}. Therefore, more
47 amplified DNA materials was necessary to fetch the minor oligo strands from the skewed oligo
48 pool. The length and sequence context, GC content and secondary structure of DNA molecule are
49 well known to introduce large amplification bias in PCR for amplification of multiple templates
50 in parallel¹⁴. Minor oligos could be excluded easily in a few amplification cycles because of its

51 replication disadvantage. However, DNA storage requires amplification of from tens of thousands
52 up to million distinct DNA strands at the same time. Thus far, this problem was dealt with by high
53 encoding, physical and sequencing redundancy^{13, 15}, but paid price of losing storage density and
54 increasing cost in all the respects of synthesis, sequencing and decoding calculation. In previous
55 studies, it has been demonstrated that deep PCR amplification (over 60 thermal amplification
56 cycles) increased the unevenness of copy number largely and required several orders of magnitude
57 more sequencing reads for decoding¹². For extreme case, only a few oligos (less than 10 % of the
58 total oligos) wiped out almost all others after PCR amplification¹⁶. Moreover, successive repeated
59 PCR amplification significantly skewed the copy number distribution, caused large positive shift¹².
60 Even the sequences were carefully designed to minimize the sequence context difference, PCR
61 still caused significant chaos from the huge sequence complexity of oligo pool for DNA storage
62 and then imperiled the information decoding.

63 For practical data storage, crucial issues have to be addressed. First, information reading method
64 should be redesigned to handle amplification of oligo pool with high complexity of DNA sequence
65 at low bias and also support the repeated reading for long-term storage. Second, it is required to
66 flatten the unevenness of oligo copy number following microchip synthesis for decoding with less
67 sequencing resource. Last, the thermal cycling process in PCR not only consumes energy but also
68 will trouble the operation of storage device. It has been reported that long thermal treatment at
69 around 65°C resulted to oligo decay and higher temperature caused more damage to longer DNA
70 molecule^{16, 17}. All these issues come down to developing biosystem which could stably and
71 repeatedly handle oligo pool with high sequence complexity and uneven copy number.

72 Here, we adapted a BASIC code system^{18, 19} recently developed for digital distributed file system
73 for DNA storage, in which a high information density was achieved with a low encoding

74 redundancy of 1.56%, the lowest one ever reported so far. Above all, especially for amplification
75 of oligo pool with high sequence complexity, advanced isothermal DNA reading, termed as iDR,
76 was designed from a novel strand displacement isothermal reaction that drives amplification with
77 distinct mechanism to PCR. Deep statistics analysis demonstrated that oligo pool was more prone
78 to dropout than base error during biased amplification. Oligo pool comprising over ten thousand
79 strands was skewed even in a very light 10 cycles PCR amplification resulting to almost 4 times
80 of oligo dropout than iDR. Moreover, dropout enlarged with the amplification depth and repeated
81 times increased, 13.90% and 53.19% of oligo were wiped out in a 60-cycle deep and 10-times
82 repeated PCR amplification respectively. In contrast, the letter error, either base substitution or
83 indel mutation, remained relatively stable, but the depth of same error increased, namely error-
84 spreading prone. These observations demonstrated that the biased amplification skewed the
85 evenness of oligo pool largely, but base error majorly depended on the total amount of molecular
86 replication. However, due to its molecular nature, iDR addressed these crucial issues, inefficient-
87 priming, product-as-template that caused the bias and error-spreading in PCR amplification. For
88 iDR, oligo dropout remained very stable at low level of around 1.33%, even after successive 10
89 times repeated amplification. Over 70% of uniformity of oligo pool has been lost, decreased to
90 0.15 from 0.50 in multiple PCR, but iDR remained almost same at 0.58. To the best of our
91 knowledge, this is the first DNA storage system successfully achieved the deep repeated
92 information reading.

93 Furthermore, we developed large-scale oligos pool normalization (OPN) to stoichiometrically
94 rectify the oligo copy number unevenness for further improvement of information decoding. Total
95 1280 OPN probes were designed from 5 distinct anchors and one set of 256 barcode. The large
96 OPN probe mixture was demonstrated being able to improve the Gini index of oligo pool

97 comprising over thousands oligo stands, the most large-scale oligo normalization reported so far.
98 Oligos was 100% retrieved from sequencing reads much less than original oligo pool read by PCR.
99 The new biochemical framework was demonstrated being able to handle high sequence complexity,
100 by which stable multiple reading was achieved based on error-spreading proof and low bias
101 isothermal amplification, large scale oligo normalization and information could be decoded from
102 at least two orders of magnitude less necessary DNA material and sequencing resource and also
103 highly compatible with any other encoding software system for sustainable and efficient DNA data
104 storage.

105

106 **Results:**

107 **BASIC code for DNA mediated distributed storage.**

108 In current DNA storage, the digital file was divided and written into a large group of small piece
109 DNA oligos. Every single oligo function as individual information carrying unit. And then, the
110 entire file can be read out from sequencing all the oligos (Fig. 1a and Supplementary Note 1).
111 Therefore, DNA storage can be considered as a biomolecular distributed storage system, the
112 “Software-Defined Storage” in semiconductor hard drive. Due to this nature, we adapted a BASIC
113 coding system that is one well optimized regenerating code with reducing computational
114 complexity for distributed storage. It is flexible for various DNA oligo length from different
115 commercial synthesis platform. Generally, the encoding process started by dividing the target file
116 into non-overlapping groups. And then, the split information was encoded into DNA sequence
117 following an optimized encoding process (Supplementary Figs. 1-3) with two adjustable parameter
118 K (corresponding to oligo number in one non-overlapping group) and L (corresponding to the
119 length of oligo). Besides the normal file types, for the first time we tested the storage of genome
120 sequence including human mitochondrion and one artificial bacteria cell with a different encoding
121 strategy (Supplementary Fig. 4 and Supplementary Note 2). In genome sequence, there are many
122 complicate structures hard for correct sequencing. Through the encoding process, the genome
123 sequence was rewritten into nucleotide sequence again but with complicated sequence avoided and
124 accurate sequencing guaranteed by error correcting code. Furthermore, oligo pools with different
125 payload length have been designed to store 2.85 MB files in totally 109,568 oligo strands (Fig. 1a
126 and Supplementary Figs. 6 and 8). Reed-Solomon code was used for error correction and coding
127 redundancy for tolerance of missing entire oligo. In comparing with previous reported systems,
128 relative high information density 1.65 bits/nt was achieved with a 1.56% of coding redundancy,

129 which allowing success decoding as long as oligo dropout less than 1.56%, randomly losing 4
130 strands from 256 (Supplementary Fig. 5 and Supplementary Note 6). Technically, higher encoding
131 redundancy tolerates losing more oligos and success decoding can be achieved from lower
132 sequencing coverage. However, higher redundancy requires more DNA synthesis. Considering the
133 synthesis cost is higher than sequencing, represented over 90% of the total cost of DNA storage,
134 it is more practical to trade-off encoding redundancy for synthesis cost.

135 Chip-based synthesis only produces very small amount of oligo, roughly from 10^5 to 10^{12} at the
136 concentration of femtomolar depending on the synthesis platform⁷. In the current workflow of
137 DNA data storage, the quality of oligo pool majorly impacted the storage performance⁴.
138 Particularly, both the heterogenous oligo sequences and the unevenness of oligo copy number
139 generated a huge sequence complexity and then caused the PCR amplification bias. Minor oligo
140 molecules are prone to drop-out and more sequencing coverage was required for decoding oligo
141 pool with largely skewed copy number distribution (Fig. 1b). The material of skewed oligo pool
142 made stable and repeatable information decoding a huge challenge.

143 **Low-bias and error-spreading proof isothermal amplification**

144 It well known that PCR generated biased amplification from its inherent mechanism, i.e., product-
145 as-template, priming and thermal cycling dependent amplification²⁰. In order to address these
146 problems that easily skewed the oligo copy number distribution, we designed a method from novel
147 isothermal DNA replication. In comparing with PCR depending on thermal cycling to drive DNA
148 replication, sequence specific nickase and DNA polymerase with processive strand displacement
149 activity were recruited for DNA amplification under consistent low temperature²¹. After
150 systematical optimization (Supplementary Figs. 12-24 and Supplementary Note 4), Nt.BbvCI and
151 exonuclease deficient DNAP I Klenow fragment ($3' \rightarrow 5'$ exo⁻) were used to amplify oligos

152 immobilized on magnetic microbeads. The immobilized isothermal DNA replication system was
153 designated as iDR, isothermal DNA reading (Fig. 2a).

154 The intrinsic features of iDR specifically facilitate DNA storage application. In theory, for one 50
155 μ l reaction, 10 PCR thermal cycles required about 177.8J energy only for thermal regulation at
156 least two orders of magnitude more than 30 mins iDR of 2.52J (Fig. 2b and Supplementary Note
157 7) and it may be a huge issue at large scale of operation. Real time monitoring indicated that
158 amplification rate of iDR is very close to PCR, even generally it was considered as a linear
159 replication (Supplementary Figs. 25 and 26). Single-stranded or double-stranded DNA can be
160 produced in a controlled manner (Fig. 2c). Specially, ssDNA was amplified in a primer-free
161 manner and allows iDR to be a universal method for reading information with no sequence
162 information needed in advance. The nickase mediated site-specific phosphodiester bond cleavage
163 initiated the iDR amplification and generated a 5' terminal phosphate group, its function was
164 verified by direct ligation to a FAM labeled probe (Fig. 2d and Supplementary Figs. 27 and 28).
165 This inborn phosphate group is very convenient for subsequent functional adapter linking.

166 Amplified oligos were sequenced on commercial Illumina Hiseq 4000 platform with 150 paired-
167 end cycles and then deeply analyzed by a set of statistics methods developed from bioinformatic
168 BLAST program (Supplementary Fig. 30). Sequenced reads with various number of letter error
169 including substitution and indel were counted with significant different amount. Among them,
170 single letter error accounted for the vast majority in mutant sequenced reads for both PCR (80.1%)
171 and iDR (81.7%) (Fig. 2e, Supplementary Figs. 31-35 and Supplementary Note 11). The total indel
172 (0.03%) and substitution (0.2%) base rate were consistent with previous studies¹³. Although, RS
173 code is able to correct multiple errors in same DNA strand but will increase the computation
174 complexity largely^{17,22}. Therefore, both sequenced reads with no or single substitution/indel error

175 were collected as valid reads for further analysis. The distribution of the number of reads per each
176 given sequence were different between iDR and light PCR of just 10 thermal cycles. The
177 distribution normality was quantified by a modified function (Supplementary Note 9). In theory,
178 oligo pool of higher value distribution normality require less sequencing reads to recover all
179 synthesized oligos¹². The coverage distribution of both 10 cycles PCR and iDR were positively
180 skewed with a long tail, which comprising of high copy number oligos, but the normality of 10
181 cycles PCR decreased about 16% than iDR (Fig. 2f). The proportion of oligos with high copy
182 number in the tail, the top 30% of high coverage, enlarged with the PCR cycles increased
183 (Supplementary Fig. 38 and Supplementary Note 9). No obvious difference was observed between
184 the coverage distribution of iDR amplification from free oligo pool and oligo pool immobilized on
185 magnetic beads (Supplementary Figs. 39 and 40). These results demonstrated that even just 10
186 thermal cycles largely skewed the oligo pool due to the PCR amplification bias and the
187 amplification of iDR was much low-biased.

188 Deep errors, reads with mutant in high copy number, impeded information decoding. Sequenced
189 reads were further sorted out as group of M0G0 (with no letter error) and M1G1 (with single letter
190 substitution or indel error) by developed BLAST programs. The distribution of the number of reads
191 in M0G0 and M1G1 overlapped for both PCR and iDR (Fig. 2g). The extent of overlap between
192 coverage distribution of M0G0 and M1G1 correlated with the potential of mutant reads affecting
193 retrieval of correct reads for decoding. The max coverage for M1G1 of 10 cycles PCR was counted
194 as 63 and 6956 oligos in PCR M0G0 was identified with copy number lower than it, accounting
195 for 60.38% of total oligos. However, the corresponding number for iDR were counted respectively
196 as 28 and 2011, accounting for 17.46% of total oligos, 71% less deep error than PCR. It indicted
197 that in PCR 60.38% of oligos will be infected if retrieved by principle of law of large number-

198 makers, but only 17.46% for iDR. Because same master pool used as template, it demonstrated
199 that error accumulated much more in PCR than iDR. And even only 10 thermal cycles of PCR
200 caused deep error-spreading that would cause huge calculation in identifying the majority (over
201 60%) of correct oligos and the proportion enlarged with PCR cycle number increased
202 (Supplementary Fig. 45). Moreover, both the total rate for substitution or indel error and proportion
203 of valid reads (combined M0G0 and M1G1 group) in all noisy sequenced reads remained stable
204 for both iDR and PCR of thermal cycles from 10 to 60 (Fig. 2h and Supplementary Figs. 47 and
205 48). In comparison, the dropout rate increased significantly from 4.18% to 13.80% with PCR
206 amplification becoming deeper, much higher than iDR of 1.33%. Low replication fidelity of DNA
207 polymerases and around 1% miss reading coming along with NGS sequencing process²³ largely
208 contributed to the massive error reads with low copy number, e.g. 1-2 copy number, but it is
209 relatively easy to identify them from correct reads in high copy number. Considering the
210 mechanism, product-as-template and inefficient priming process caused the high amplification
211 bias and made PCR amplification prone to error-spreading and oligo dropout. In contrast, iDR was
212 designed to synthesize new oligo only from the original templates without priming process for
213 replication initiation and therefore iDR achieved a low biased and error-spreading proof
214 amplification and will require much less calculation resource in decoding process. Additionally,
215 high temperature treatment resulted to DNA oligo decay especially for long strands^{17,24}. The decay
216 lost rate in a 10 cycles PCR was calculated as 21.8% and 0.035% for iDR (Fig. 2b and
217 Supplementary Note 8) from a plotted DNA half-life graph (Supplementary Fig. 29).

218 **Multiple-repeated information reading.**

219 The amplification capability for successive deep reading oligo pool containing from 11,520 to
220 89,088 DNA strands was examined. PCR amplification was successively performed 10 times from

221 aliquot of previous reaction and iDR amplification was repeated 10 times from immobilized oligos
222 pool (Fig. 3a, Supplementary Figs. 49-51). In #1, #5 and #10 of successive PCR, the proportion of
223 amplified oligos with up to total 10 substitution or indel letter error decreased from 89.09% (± 10
224 of PCR #1) to 49.97% (± 10 of PCR #10) and from 86.77% (± 1 of PCR #1) to 48.58% (± 1 of PCR
225 #10). In comparison, the repeated iDR remained very consistent over 90% (Fig. 3b). These
226 statistics results indicated that large noise was introduced during PCR procedure and amplified
227 oligo with imperfect length increased from inefficient replication or miss priming. However, no
228 obvious difference was observed in proportion of M1G1 in total sequenced reads between PCR
229 and iDR, but the mean copy number in M1G1 reads increased from 1.10 of #1 PCR to 1.95 of #10
230 PCR, but only slightly changed for iDR, 1.08 of #1, 1.08 of #5 and 1.05 of #10 (Fig. 3c). This
231 result was in agreement with previous experiment in figure 2g and indicated that error accumulated
232 to high copy number during PCR and iDR achieved error-spreading proof amplification.

233 Interestingly, it was observed that the distribution of the number of reads per each given sequence
234 changed differently. Normality of #1, #5 and #10 PCR decreased significantly from 0.50 to 0.15.
235 In contrast, #1, #5 and #10 iDR gave a consistent normality of 0.58 (Fig. 3d, Supplementary Figs.
236 52-55). For successive PCR, the coverage distribution was largely positively skewed and the
237 proportion of both low copy number and high copy number oligos significantly increased. It
238 indicated that large enrichment driven by amplification bias for part of oligos efficiently occurred
239 with the successive PCR, but not in repeated iDR. Furthermore, it was observed that only top 1%
240 oligo of high coverage increased its proportion significantly and both 1% oligos of middle and low
241 coverage decreased while the oligo pool was successively read by PCR. In contrast, all the 1%
242 oligos remained steady in repeated iDR (Supplementary Fig. 56 and Supplementary Note 9). The
243 dropout rate in random valid reads set with 10x coverage was quantified to further assess the

244 amplification bias. For PCR, the dropout rate increased sharply from 4.18% (#1 PCR) to 53.19%
245 (#10 PCR) (Fig.3e), but iDR remained steady at about 2%. Due to the tolerance for 1.56% dropout
246 in BASIC encoding algorithms, we also calculated the coverage depth for random valid reads set
247 with 1.56% dropout (Supplementary Note 6), the crucial parameter for the theoretical minimum
248 decoding coverage. For PCR, the minimum coverage depth was quantified as 17.2 (#1 PCR) and
249 167 (#5 PCR). Because #10 PCR lost too many oligos, the minimum coverage depth was
250 calculated as 426 from calculation (Supplementary Fig. 61). For iDR, the minimum coverage
251 depth was quantified as 11 for #1, 12 for #5 and 12.5 for #10. Thus far, for PCR based sustainable
252 DNA media reading, one strategy is deep amplification from trace DNA material with large
253 number of thermal cycles and another is successive amplification from aliquot of previous reaction.
254 However, we demonstrated that both strategies are not practical. The nature of PCR including
255 especially inefficient priming and product-as-template generated large amplification bias that
256 significantly skewed the copy number distribution and resulted huge dropout in both 60 cycles
257 deep amplification and 10 times successive amplification, especially for oligo pool with huge
258 sequence complexity. Based on these deep statistical analyses, we pointed out that iDR
259 amplification was more stable, robust and suitable for repetitive read than PCR in DNA data
260 storage.

261 **Large-scale oligo pool normalization improved information reading.**

262 Besides the amplification procedure, the unevenness of oligo copy number originally stem from
263 microchip-based synthesis. However, there is still short of handy technology to normalize oligo
264 pool with high sequence complexity. To address this problem, we developed a simple oligo pool
265 method based on a previous reported study²⁵ but with more simple and flexible procedure and
266 without expensive DNA chemical modification. Two-parts probe was designed, in which constant

267 anchor part provide basic efficient binding and short variable 9nts barcode part generating 256
268 specific targets recognition with no expensive modification. In particular, oligo pool normalizing
269 (OPN) sequence was synthesized at the 3' terminal end, which comprising of barcode (9 nts) and
270 universal fragment (17 nts) (Fig. 4b and Supplementary Fig. 8). The barcode guided the OPN probe
271 to identify its target while universal fragment binding as an anchor (Fig. 4b). The unique barcode
272 sequence "GWSWSWS", alternating strong (G or C) and weak (A or T) bases (e.g. CACTGT or
273 GTCTGA), has been proven being able to generate 256 high specific binding (Supplementary Note
274 5). After perfect binding, OPN probe was extended by DNA polymerase, turning oligo into double-
275 stranded form and then the remaining ssDNA oligos were removed by exonuclease I degradation
276 (Fig. 4a). Due to the precisely selected degradation of normalization and the low productivity of
277 microchip oligo synthesis, it is necessary to amplify the oligo pool to a large quantity but keeping
278 its single-stranded form. Therefore, we developed a simple single-stranded oligo pool
279 amplification protocol (SOA), by which ssDNA oligo pool was amplified by PCR with normal
280 forward primer and 5'-phosphated reverse primer and then only strand extended from forward
281 primer was removed by lambda exonuclease and the intact complementary strand was left
282 (Supplementary Figs. 62 and 63).

283 First, we tested the normalization of three oligos of different length with 3' terminal FAM label,
284 which was mixed at input molar ratio of 1:5:25. After OPN normalization, the output oligos were
285 quantified as molar ratio of 0.95:1:1.25 (Fig. 4d and Supplementary Fig. 64). Next, we proceeded
286 to test normalization improved information decoding. As proof-of-concept, OPN 1.0 with 256
287 probes, which were synthesized separately and mixed at equal molar ratio, were used to normalize
288 256 oligos synthesized from microchip with other 10752 oligos together (Fig. 4c). Following the
289 SOA protocol, 256 single-stranded oligos pool was prepared and normalized, each OPN probe was

290 equal to the average molar concentration of the oligo pool. The coverage depth of each oligo in
291 OPN 1.0-iDR was significantly improved in comparing with iDR amplified oligo pool without
292 OPN normalization (Fig. 4e). The PCR-amplified oligo pool showed skewed normality than OPN-
293 iDR (Supplementary Figs. 65 and 67) and high copy number oligo in PCR was obviously
294 normalized in OPN-iDR (Supplementary Fig. 66) with decreased standard deviation in
295 quantification of molecule number of each oligo strands. And then, dropout rate for random valid
296 reads set with various coverage depth was plotted (Supplementary Fig. 68). All of the dropout rate
297 decreased as coverage increased. At the coverage of about 25, the dropout rate of PCR became
298 lower than 1.56%, the limit for decoding. However, all of dropout rate of OPN-iDR was lower
299 than 1.56% and became 0 at around coverage of 20. The dropout rate of iDR was always higher
300 than OPN-iDR but much lower than PCR. At coverage of 10, the dropout rate of OPN-iDR was
301 almost one order of magnitude lower than PCR. The recognition capability of OPN1.0 was
302 extended simply by recruiting more anchor design. In OPN2.0 (Fig. 4f), 1024 specific target oligo
303 recognition was generated from 4 distinct anchor sequences and the one set 256 barcode. All of
304 these 1024 oligos was perfect retrieved from 87x sequencing reads less than that 223x of PCR
305 amplified oligos (Fig. 4g and Supplementary Figs.69 and 71). Technically, this method can be
306 expanded by using a combination of a group of specific universal probes¹³ and the 256 barcodes
307 to manipulate oligo pool comprising of up to 3 Million oligos (Supplementary note 5). Although
308 minimum decoding coverage depth previously reported (Supplementary Fig. 72) is lower than the
309 result presented here, we have to point out that the minimum decoding coverage is highly
310 dependent on both the encoding system with different encoding redundancy and the quality of
311 oligo pool from different synthesis platform (Supplementary Figs. 73-76) and it is hard to compare
312 the results cross different systems.

313 **Discussion:**

314 DNA data storage is one very artificial application and actually the concept has been proposed for
315 a long time²⁶. But till lately, the significantly increasing capability of DNA synthesis and
316 sequencing start making it possible. In many respects, practical DNA data storage requires more
317 powerful capability for all related biotechnologies, including synthesis, sequencing and
318 manipulation of oligo pool with high sequence complexity. The size of oligo pool used in DNA
319 storage is already several orders of magnitude larger than that in other applications. However, there
320 is still short of hardware technology which is designed for DNA storage, especially manipulation
321 of DNA oligo pool with huge sequence complexity and unevenness of molecule copy number. In
322 studies to date, PCR is still the only method for amplification of oligo pool for reading information.
323 However, with deep bioinformatic and statistics analysis, we demonstrated that PCR amplification
324 skewed the oligo pool even after a very few thermal cycles (10 cycles) and the skewness largely
325 increased, deep error spreading and massive dropout, as the function of amplification cycle
326 numbers. Deep error, mutant sequenced reads in high coverage, interfered decoding and caused
327 significant increased calculation, but decoding fatally crashed due to massive dropout. We contend
328 that the features of low temperature, priming-free, and enzyme-mediated double helix unzipping
329 in iDR amplification presented here, overcame the major amplification bias related issues of PCR
330 including the inefficient priming, product-as-template, sequence context dependent, and high
331 temperature heating. iDR achieved very stable amplification performance, which efficiently
332 prevented the mutant error from spreading, decreased over 70% deep error, and achieved
333 successive repeated deep reading with consistent outcome quality. Actually, the physical storage
334 density, the most significant advantage of DNA storage with a calculation of a few kilograms DNA
335 material for storage of data from the whole human society²⁷, highly depend on the quality of DNA

336 reading. Therefore, higher physical storage density could be achieved by iDR system, which could
337 store information with at least two orders of magnitude less DNA material and sequencing resource
338 than the current PCR method (Fig. 5 and Supplementary Note 12).

339 Although, there is many variants of PCR, such as emulsion PCR, digital PCR, and multiplex PCR,
340 but none of them can avoid the crucial issues of product-as-template, inefficient priming, and
341 complex thermal regulation. Furthermore, Replication fidelity of DNAP I Klenow fragment used
342 in iDR system is much lower than Q5 DNA polymerase in PCR. There is still much room for
343 improvement. Additionally, practical DNA storage system must move out of biochemical test tube
344 to build device by highly integrating all related biochemical processes. Prototypes of DNA storage
345 hardware have already been proposed with high density DNA material in microfluidic device^{28,29}.
346 We contend that besides some crucial features of iDR make it more fit for hardware construction.
347 In iDR, only low temperature, slightly higher than room temperature, was required. Considering
348 the operation of DNA storage up to large scale, there will be a huge difference in energy
349 consumption. Furthermore, iDR is able to work in a primer-free fashion, only defined protein
350 enzyme mix is required no matter what information was encoded, which makes it possible for
351 universal information reading without any sequence information required in advance. To the best
352 of our knowledge, iDR is the first system for stable repeated DNA information reading which is
353 crucial and necessary feature for practical and sustainable storage hardware.

354 The large-scale oligo normalization is another advanced feature for DNA storage. Comparing with
355 the previous reported system, OPN was developed as a simple and very economical process
356 without any expensive modification directly on oligo pool or OPN probe, such as dexoyuracl (dU)
357 and biotin, none of them has been reported to be synthesized directly on chip-array. Then OPN
358 could be applied to any chip-array oligo synthesis platform and the OPN probe synthesis cost is

359 also acceptable (Supplementary Fig.77). Additionally, SOA protocol was developed enabling
360 OPN to manipulate oligo pool with very small amount DNA molecule, generally chip-array
361 synthesized oligo pool is very small amount. Moreover, the capability of OPN could be easily
362 extended. OPN 2.0 achieved specifically targeting 1024 oligos, 4 times of OPN 1.0, simply by
363 combining 4 distinct anchor sequences, which is the most large-scale oligo pool normalization
364 reported so far. 1024 array synthesized oligos was revised with improved Gini index
365 (Supplementary Fig.70), and was perfect identified from sequencing reads (87x) less than that
366 (223x) for PCR amplified oligos. Therefore, the necessary DNA material and sequencing resource
367 for information decoding could be decreased at least two orders of magnitude (Fig. 5), which is a
368 huge advantage in considering application of DNA storage in large scale comparing to current big
369 data center. Following the parallel extension strategy, target recognition could be leveraged to
370 million in an economical way (Supplementary note 5 and Supplementary Fig.63). Therefore, we
371 believed that this new biochemical framework with stable repeated reading and large-scale
372 normalization lays a foundation for development of practical and sustainable DNA storage.

373

374 **ACKNOWLEDGMENTS**

375 We would like to thank Professor Hanxu Hou from Dongguan University of Technology for advice
376 and assistance with designing algorithm for BASIC code. We also thank Yixi Wang from College
377 of intelligence and computing at Tianjin University for her help in test of encoding system. This
378 work was supported by the National Science Foundation of China (Grant No.21476167,
379 No.21778039 and No.21621004).

380

381 **AUTHOR CONTRIBUTIONS**

382 Y.G. and H.Q. designed and performed all experiments. X.C., J.H. and C.Z. designed and
383 developed the encoding and decoding program. X.C. developed program for the bioinformatics
384 statistics analysis. Y.G., H.Qiao and H.Q. collected and analyzed all experiment data. G.Y. T.H.
385 and H.Q. wrote the manuscript. H.Q. designed experiments, analyzed data and supervised this
386 work.

387

388 **COMPETING FINANCIAL INTERESTS**

389 H.Q. is the inventor of two patents application for the biochemical method described in this article.
390 The initial filings were assigned Chinese patent application (201911086860.0 and
391 201911087247.0) and international patent application (PCT/CN2019/123916). The remaining
392 authors declare no competing financial interests.

393 **References:**

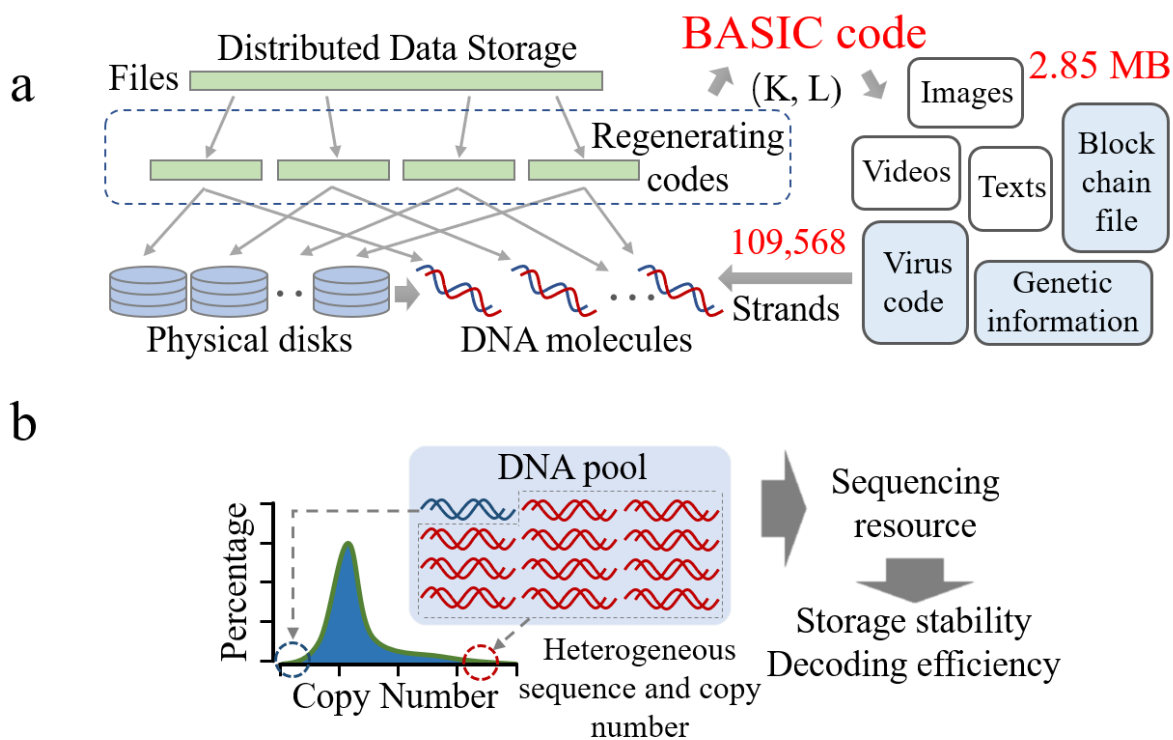
- 394 1. Kosuri, S. & Church, G.M. Large-scale de novo DNA synthesis: technologies and
395 applications. *Nat Methods* **11**, 499-507 (2014).
- 396 2. Goodwin, S., McPherson, J.D. & McCombie, W.R. Coming of age: ten years of next-
397 generation sequencing technologies. *Nat Rev Genet* **17**, 333-351 (2016).
- 398 3. Kozarewa, I. et al. Amplification-free Illumina sequencing-library preparation facilitates
399 improved mapping and assembly of (G+C)-biased genomes. *Nat Methods* **6**, 291-295
400 (2009).
- 401 4. Organick, L. et al. Experimental Assessment of PCR Specificity and Copy Number for
402 Reliable Data Retrieval in DNA Storage. (2019).
- 403 5. Chen, Y.-J. et al. Quantifying Molecular Bias in DNA Data Storage. (2019).
- 404 6. Klein, J.C. et al. Multiplex pairwise assembly of array-derived DNA oligonucleotides.
405 *Nucleic Acids Res* **44**, e43 (2016).
- 406 7. Jingdong Tian, H.G., Nijing Sheng, Xiaochuan Zhou, Erdogan Gulari, Xiaolian Gao &
407 George Church Accurate multiplex gene synthesis from programmable DNA microchips.
408 *Nature* **432**, 1050-1054 (2004).
- 409 8. Linnarsson, S. Recent advances in DNA sequencing methods - general principles of sample
410 preparation. *Exp Cell Res* **316**, 1339-1343 (2010).
- 411 9. Kosuri, S. et al. Scalable gene synthesis by selective amplification of DNA pools from
412 high-fidelity microchips. *Nat Biotechnol* **28**, 1295-1299 (2010).
- 413 10. Lipshutz, R.J., Fodor, S.P., Gingeras, T.R. & Lockhart, D.J. High density synthetic
414 oligonucleotide arrays. *Nat Genet* **21**, 20-24 (1999).
- 415 11. Bonde, M.T. et al. Direct mutagenesis of thousands of genomic targets using microarray-
416 derived oligonucleotides. *ACS Synth Biol* **4**, 17-22 (2015).
- 417 12. Erlich, Y. & Zielinski, D. DNA Fountain enables a robust and efficient storage architecture.
418 *Science* **355**, 950-954 (2017).
- 419 13. Organick, L. et al. Random access in large-scale DNA data storage. *Nat Biotechnol* **36**,
420 242-248 (2018).
- 421 14. Aird, D. et al. Analyzing and minimizing PCR amplification bias in Illumina sequencing
422 libraries. *Genome Biol* **12**, R18 (2011).
- 423 15. Goldman, N. et al. Towards practical, high-capacity, low-maintenance information storage
424 in synthesized DNA. *Nature* **494**, 77-80 (2013).
- 425 16. Heckel, R., Mikutis, G. & Grass, R.N. A Characterization of the DNA Data Storage
426 Channel. *Sci Rep* **9**, 9663 (2019).
- 427 17. Grass, R.N., Heckel, R., Puddu, M., Paunescu, D. & Stark, W.J. Robust chemical
428 preservation of digital information on DNA in silica with error-correcting codes.
429 *Angewandte Chemie* **54**, 2552-2555 (2015).
- 430 18. Hou, H., Shum, K.W., Chen, M. & Li, H. BASIC Codes: Low-Complexity Regenerating
431 Codes for Distributed Storage Systems. *IEEE Transactions on Information Theory* **62**,
432 3053-3069 (2016).
- 433 19. Hou, H., Han, Y.S., Shum, K.W. & Li, H. A Unified Form of EVENODD and RDP Codes
434 and Their Efficient Decoding. *IEEE Transactions on Communications* **66**, 5053-5066
435 (2018).
- 436 20. Mullis, K.B. & Faloona, F.A. Specific synthesis of DNA in vitro via a polymerase-
437 catalyzed chain reaction. *Methods Enzymol* **155**, 335-350 (1987).

- 438 21. Joneja, A. & Huang, X. Linear nicking endonuclease-mediated strand-displacement DNA
439 amplification. *Anal Biochem* **414**, 58-69 (2011).
- 440 22. Anavy, L., Vaknin, I., Atar, O., Amit, R. & Yakhini, Z. Data storage in DNA with fewer
441 synthesis cycles using composite DNA letters. *Nat Biotechnol* (2019).
- 442 23. Yazdi, S.M., Yuan, Y., Ma, J., Zhao, H. & Milenkovic, O. A Rewritable, Random-Access
443 DNA-Based Storage System. *Sci Rep* **5**, 14138 (2015).
- 444 24. A Characterization of the DNA Data Storage Channel. (2018).
- 445 25. Pinto, A., Chen, S.X. & Zhang, D.Y. Simultaneous and stoichiometric purification of
446 hundreds of oligonucleotides. *Nat Commun* **9**, 2467 (2018).
- 447 26. Wallace, M.R. Molecular Cybernetics: The Next Step? *Kybernetes* **7**, 265-268 (1978).
- 448 27. Extance, A. How DNA could store all the world's data. *Nature* **537**, 22-24 (2016).
- 449 28. Newman, S. et al. High density DNA data storage library via dehydration with digital
450 microfluidic retrieval. *Nat Commun* **10**, 1706 (2019).
- 451 29. Takahashi, C.N., Nguyen, B.H., Strauss, K. & Ceze, L. Demonstration of End-to-End
452 Automation of DNA Data Storage. *Sci Rep* **9**, 4998 (2019).

453

454

455



456

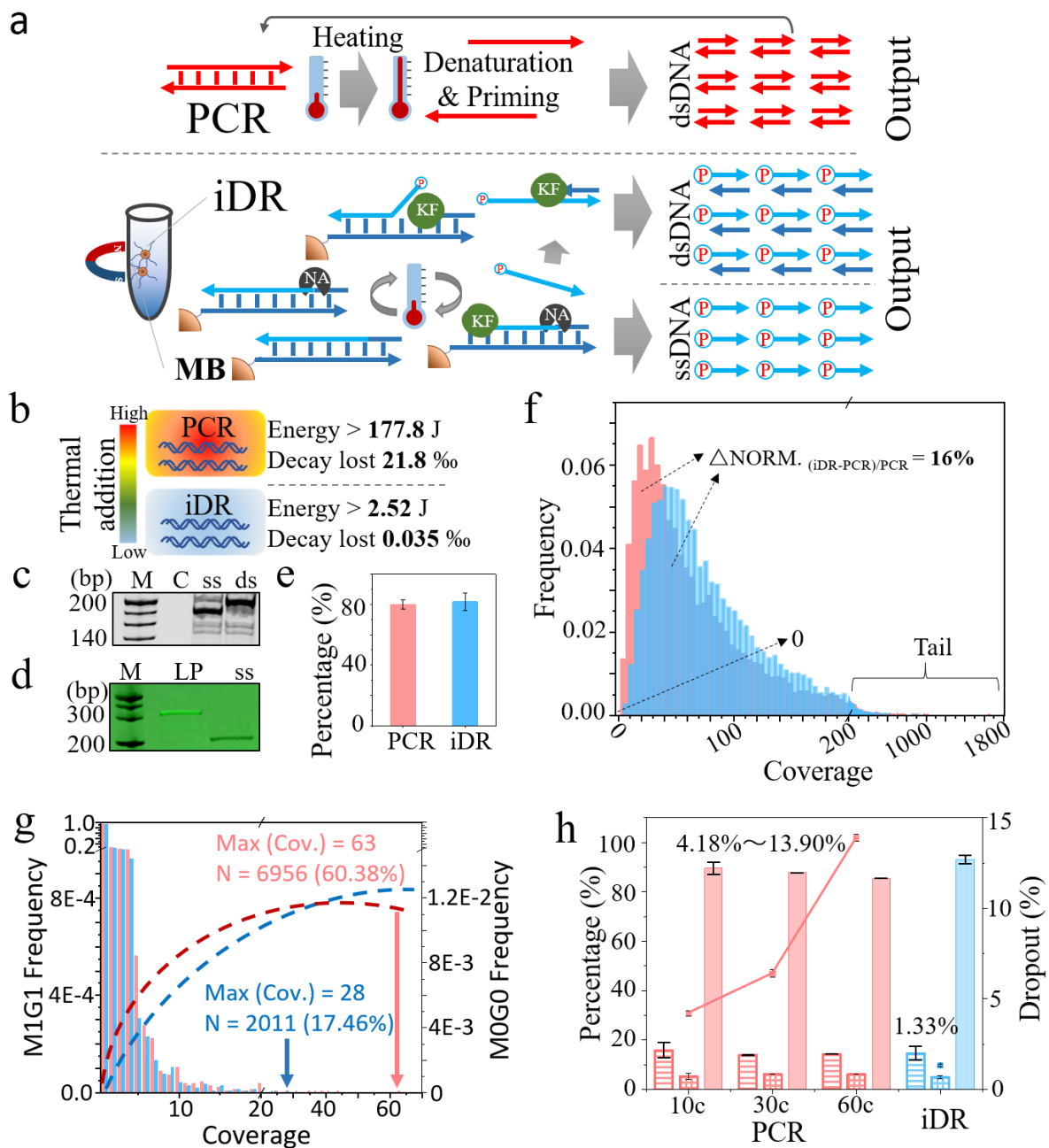
457

458 **Figure 1:** Overview of the DNA data storage.

459 **(a)** Schematic of the DNA mediated distributed data storage and totally 2.85 MB data, including
460 text,s, images, videos, and block chain file, computer virus code, genetic information, were
461 encoded into total 109,568 synthesis oligos using a adapted BASIC code system.

462 **(b)** Illustration of data storage in synthesis oligo pool. The huge sequence complexity and
463 unevenness of oligo copy number is challenging the stable storage and precise decoding.

464



465

466

467 **Figure 2:** Designed isothermal DNA reading (iDR) for DNA data storage.

468 (a) Illustration of the DNA amplification of PCR (upper) and iDR (lower). PCR depend on thermal

469 cycling to drive replication. By contrast, isothermal DNA reading was developed from

470 systematically optimized strand displacement mediated replication, in which nickase and specific
471 DNA polymerase collaborate to drive amplification under consistent low temperature.

472 **(b)** Calculation of energy consumption and the oligo degradation decay for one 50ul liquid reaction
473 amplified by 10 thermal cycles PCR or iDR. PCR and iDR required 177.8 J and 2.52 J respectively
474 only for thermal regulation theoretically. Thermal treatment caused oligo decay rate was calculated
475 as 21.8‰ for 10 cycles PCR and 0.035‰ for iDR respectively.

476 **(c)** Single-stranded and double-stranded oligo amplified by iDR from a 218nt dsDNA template
477 were analyzed on 10% native PAGE gel. M: 20 bp DNA Ladder; C: negative control of
478 amplification without input template; ss: single-stranded DNA product; ds: double-stranded DNA
479 product.

480 **(d)** A 5' terminal FAM labeled 30nt single-stranded probe ligated to 198nt single-stranded DNA
481 with an inborn 5' terminal phosphate group directly from iDR amplification was analyzed and
482 imaged on a 12% UREA denature PAGE. M: 20 bp DNA Ladder; LP: ligation product; ss: ssDNA
483 product of iDR.

484 **(e)** Sequenced reads with single letter error, substitution or indel, accounted for the vast majority
485 in total mutant sequenced reads for both PCR (80.1%) and iDR (81.7%). Error bars represent the
486 mean \pm s.d., where $n = 3$.

487 **(f)** Coverage depth distribution of sequenced reads for 10 cycles PCR and iDR amplification were
488 positively skewed but with 0 oligo dropout, but the distribution normality of PCR decreased about
489 16% than iDR.

490 **(g)** Sequenced reads were grouped as M0G0 (with no letter error) and M1G1 (with single letter
491 substitution or indel error) by developed BLAST programs. In M1G1 of 10 cycles PCR the max

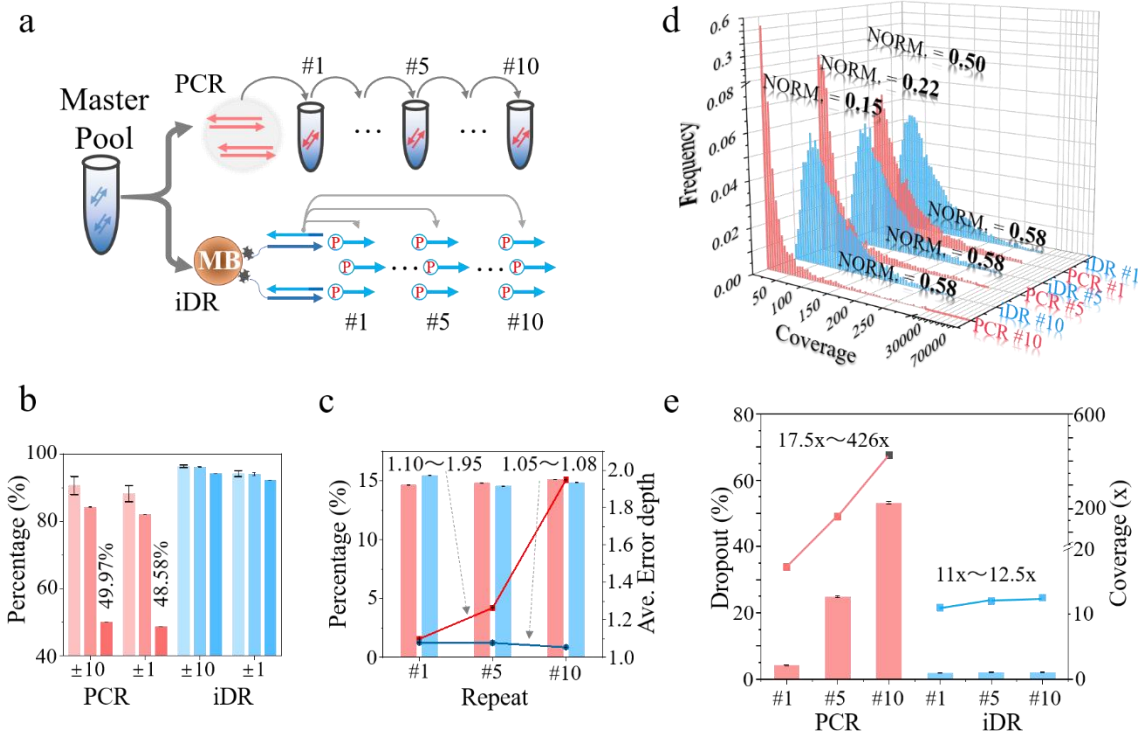
492 coverage was counted as 63 and 6956 oligos overlapped with M0G0 group, accounting for 60.38%
493 of total oligos. For iDR it was 28 and 2011 respectively, accounting for 17.46% of total oligos
494 with 71% of deep error decreased.

495 **(h)** The proportion of error reads and valid reads (total reads with no or single substitution/indel
496 error) per million noise sequenced reads were plotted, striped column for base substitution; grid
497 column for indel, and solid bar for valid reads. For 10, 30 and 60 cycles PCR, the substitution were
498 15.92%, 13.90% and 14.27%, and indel were 5.28%, 6.16% and 6.15%, and valid reads were
499 89.58%, 87.73% and 85.53% respectively. For iDR, 14.69% for substitution error, 4.99 for indel
500 and 93.15% for valid reads. Error bars represent the mean \pm s.d., where $n = 3$. Oligo dropout rate
501 calculated from random reads set with a mean 10x coverage depth was plotted. For 10, 30 and 60
502 thermal cycles PCR (red line), it was counted as 4.18%, 6.41% and 13.90% respectively, and 1.33%
503 for iDR (blue dot). Error bars represent the mean \pm s.d., where $n = 10$.

504

505

506



507

508

509 **Figure 3: Multiple information reading from DNA storage.**

510 **(a)** Illustration of successive DNA reading. Repeated PCR amplification was successively
 511 performed 10 times with aliquot of previous reaction as template and the oligos pool immobilized
 512 on magnetic beads was successive amplified 10 times by iDR. The same master oligo pool was
 513 used as the initial template and #1, #5 and #10 amplified oligos were analyzed.

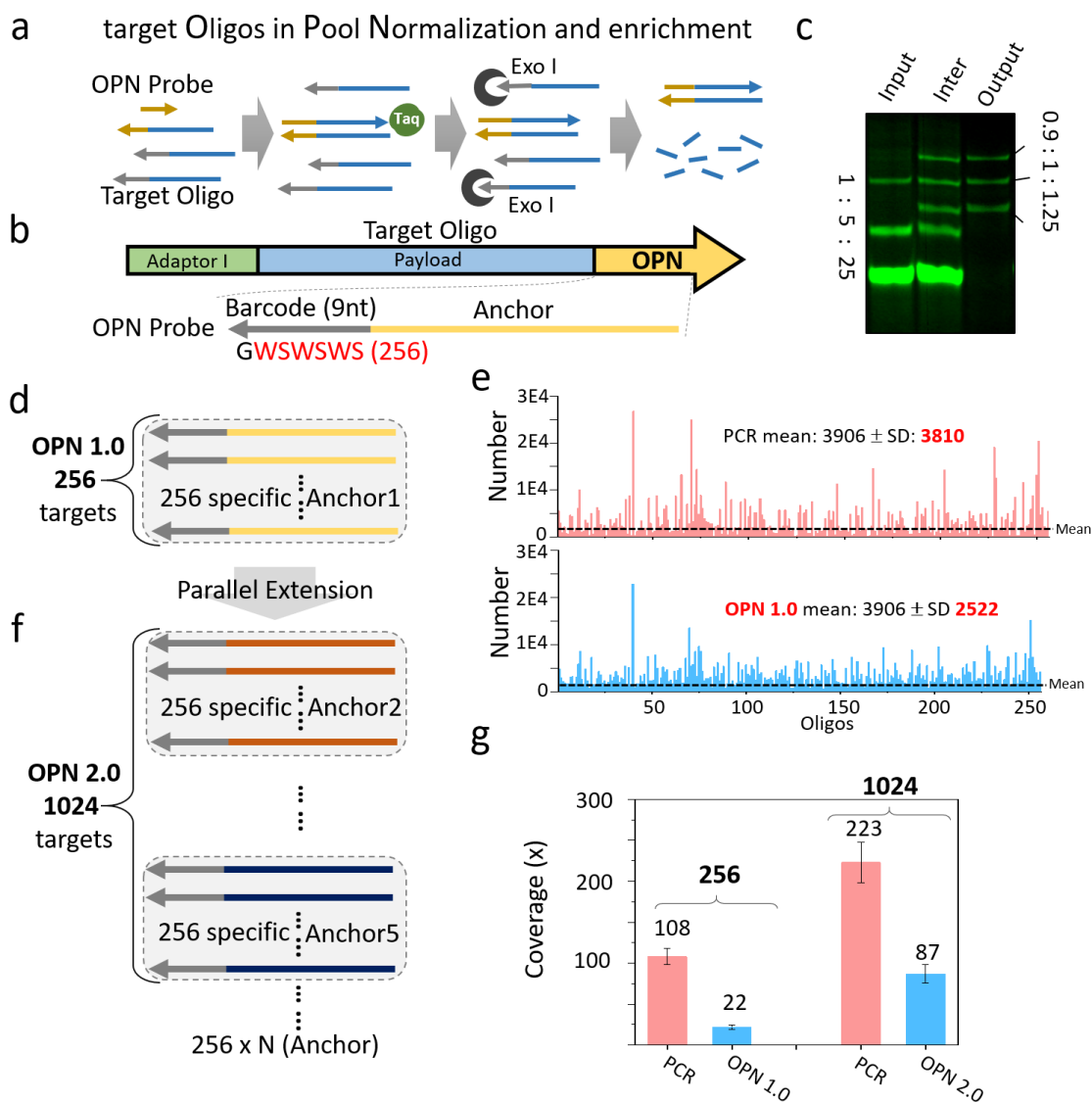
514 **(b)** The proportion of reads with up to 10 substitution or indel letter error per million sequenced
 515 reads were counted as 90.65% (#1 PCR, light red), 84.26% (#5 PCR, red), 49.97% (#10 PCR, dark
 516 red), 96.35% (#1 iDR, light blue), 96.14% (#5 iDR, blue), and 94.12% (#10 iDR, dark blue), and

517 error reads with up to 1 error were counted as 88.18% (#1 PCR), 82.05% (#5 PCR), 48.58% (#1
518 PCR), 94.18% (#1 iDR), 94.04% (#5 iDR), and 92.21% (#10 iDR). Error bars represent the mean
519 \pm s.d., where $n = 3$.

520 (c) The proportion of M1G1 reads per million valid reads were, PCR (red column), 14.65% (#1
521 PCR), 14.82% (#5 PCR), 15.10% (#10 PCR), and iDR (blue column) 15.46% (#1 iDR), 14.55%
522 (#5 iDR), and 14.86% (#10 iDR). The average error reads depth of PCR (red line, 0.5 for #1, 0.22
523 for #5, 0.15 for #10) increased and iDR (blue), but iDR (blue line, 1.08 for #1), 1.08 for #5), 1.05
524 for #10) remained stable. Error bars represent the mean \pm s.d., where $n = 3$.

525 (d) The distribution of reads number per each given sequence per million sequenced reads of #1,
526 #5, and #10 of PCR and iDR. The distribution normality was 0.5 (#1 PCR), 0.22 (#5 PCR), 0.15
527 (#10 PCR), 0.58 (#1, #5, #10 iDR) respectively.

528 (e) The dropout rate for random sequenced reads set with 10x coverage depth was plotted, PCR
529 (red column, 4.18% for #1, 22.89% for #5 and 53.19% for #10) and iDR (blue column, 1.86% for
530 #1, 2.08% for #5 and 2.09% for #10). Error bars represent the mean \pm s.d., where $n = 3$. The
531 coverage depth for random sequenced reads with 1.56% dropout was calculated as 17.2 (#1 PCR),
532 167 (#5 PCR), 426 (#10 PCR) and 11 (#1 iDR), 12 (#5 iDR), and 12.5 (#10 iDR) respectively.
533 Error bars represent the mean \pm s.d., where $n = 10$.



534

535

536 **Figure 4** Large scale normalized oligo pool improved information read-out.

537 **(a)** Brief workflow for OPN probe mediated oligo pool normalization. Oligos are captured by the
 538 corresponding OPN probe and extended to dsDNA by Taq DNA polymerase. Then the remaining
 539 ssDNA oligos and probes are digested by exonuclease I.

540 **(b)** Structure of oligo for oligo pool normalizing (OPN). The adaptor I region, R (nickase
541 recognition sequence) region, payload, and OPN region. In OPN region, barcode sequence with 9
542 nts in length and a unique “GWSWSWS” pattern with alternating strong (G or C) and weak (A or
543 T) bases (e.g. CACTGT or GTCTGA), by which 256 high specific binding could be generated and
544 universal fragment (17 nts).

545 **(c)** Normalization of 3’ terminal FAM labeled three oligos, 60 nts, 73 nts and 90 nts in length, with
546 input molar ratio of 1:5:25 and the output oligos were quantified as molar ratio of 0.95:1:1.25 on
547 15% native PAGE after OPN.

548 **(d)** OPN 1.0 constructed on one anchor sequence with 256 specific target oligo recognition.

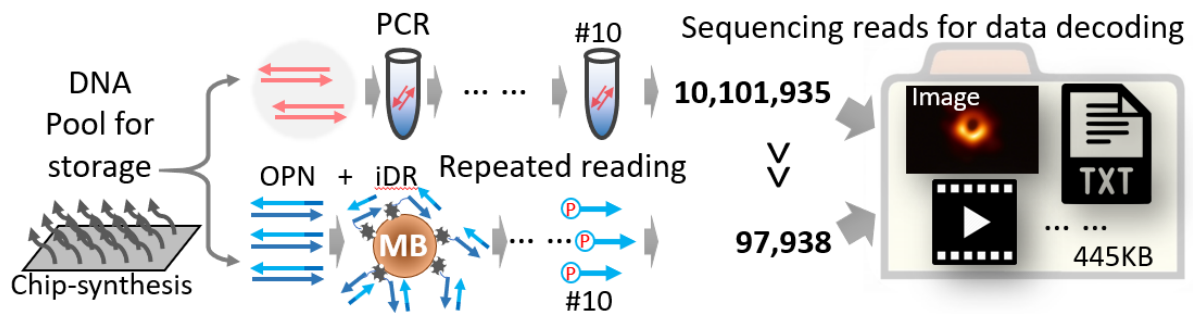
549 **(e)** The copy number of each oligo of 256 oligo pool per million valid sequenced reads. 256 oligo
550 pool amplified by PCR (red) was counted with mean copy number 3906 ± 3810 and (blue) 3906
551 ± 2522 for oligo pool revised by OPN 1.0 and then read by iDR.

552 **(f)** OPN 2.0 constructed from 4 distinct anchor sequences and one set of 256 barcode with 1024
553 specific target oligo recognition.

554 **(g)** The minimal necessary sequencing reads for complete oligos retrieve without dropout for 254
555 oligos pool read by PCR and OPN 1.0-iDR respectively, and 1024 oligo pool read by PCR and
556 OPN 2.0-iDR respectively.

557

558



559

560 **Figure 5.** Sustainable and repeated DNA storage.

561 Chip-synthesized DNA oligo pool, in which digital information including image, video and text,
562 were stored, was repeated read by PCR and OPN-iDR respectively. For current used PCR,
563 10,101,935 NGS noise reads was necessary for perfect decoding 445KB digital files, but two
564 orders of magnitude less only 97,938 NGS noise reads were necessary for OPN-iDR.

565

566 **ONLINE METHODS:**

567 **DNA Master Oligo Pool.** The pool (Pool 1/2--Twist Bioscience; Pool 3/4--CustomArray) was
568 resuspended in 1x TE buffer for a final concentration of 2 ng/uL. PCR was performed using Q5
569 High-Fidelity DNA Polymerase (NEB). Mix 10 ng of ssDNA pool (5 uL) with 2 µL of 100 µM of
570 the forward primer and 2 µL of 100 µM of the Adaptor 2 (Adaptor 2-1), 10 µL 5x Q5 Reaction
571 Buffer, 4 µL of 2.5 mM dNTPs, 0.5 µL Q5 High-Fidelity DNA Polymerase. Thermocycling
572 conditions were as follows: 5 min at 98°C; 10 cycles of: 30 s at 98°C, 30 s at 56°C, 15 s at 72°C;
573 followed by a 5 min extension at 72°C. The reaction was then purified according to the instructions
574 in the Eastep Gel and PCR Cleanup Kit and eluted in 50 µL DNase/RNase-free water. This library
575 was considered the master pool. All primers we used are in Supplementary Table 2.

576

577 **iDR reaction.** 10 ng of DNA oligo out of the master pool was attached to 1 µL of Streptavidin
578 Magnetic Beads (NEB). The iDR reaction mixtures contained 1 µL of the DNA template attached
579 to the beads (10 ng/µL), 0.25 mM dNTPs, 2.5 µL 10x NEBuffer 2, 0.08 U/µL Nt.BbvCI (NEB),
580 0.16 U/µL KF polymerase (exo⁻) (Vazyme), 4 µM T4 Gene 32 Protein, 0.2 mg/mL BSA, 0.5 µM
581 Adaptor 2 (For production of ssDNA, adaptor 2 was not added.). The mixtures were incubated at
582 37°C for 30 min. After amplification, the specific amplicon and the template was isolated through
583 magnetic pull-down. In the process of repeated iDR, the template attached to magnetic beads was
584 washed by Wash/Binding buffer (0.5 M NaCl, 20 mM Tris-HCl (pH 7.5), 1 mM EDTA) twice and
585 mixed with the above mentioned iDR reaction mixtures except for the DNA template. The process
586 was proceeded for 10 times. To retrieve the information, the amplified products were purified and
587 then sequenced on one Illumina Hiseq 4000 platform with 150 paired-end cycles in Novogene
588 (Supplementary Note 3). The process of optimization in detail are given in Supplementary Note 4.

589

590 **PCR reaction.** PCR was performed using Q5 High-Fidelity DNA Polymerase and forward
591 primer/adaptor 2 (10ng DNA master pool, 2 μ L of forward primer (100 μ M); 2 μ L of adaptor 2
592 (100 μ M)), 10 μ L 5x Q5 reaction buffer in a 50 μ L reaction. Thermocycling conditions were as
593 follows: 5 min at 98°C; 10 cycles of: 30 s at 98, 30 s at 58°C, 10 s at 72°C, followed by extension
594 at 72°C for 5 min. In the repeated PCR, each subsequent PCR reaction consumed 1 μ L of the prior
595 PCR reaction and employed 10 cycles in each 50- μ L reaction. The PCR product was purified by
596 Eastep Gel and PCR Cleanup Kit and eluted in 50 μ L DNase/RNase-free water. Then we
597 sequenced the PCR product on Illumina Hiseq 4000 platform.

598

599 **OPN probe.** Each probe contains two parts, from 5' to 3': a universal sequence (adaptor 2) and a
600 oligo-specific barcode sequence. Each barcode sequence is comprised of a number of commutative
601 strong (C or G) and weak (T or A) nucleotides according to earlier report²⁵. Here, the length of
602 barcode is 8 nucleotides, corresponding to a total of $2^8 = 256$ barcode instances (Supplementary
603 Note 5).

604 **Single-stranded oligo pool amplification (SOA).** The schematic of SOA is illustrated
605 (Supplementary Fig. 61). We used PCR with reverse primer (adaptor 2) modified with 5'
606 phosphate to amplify an oligonucleotide library with specific barcodes sequences. Then PCR
607 product was degraded from 5' phosphate groups to 3' direction by lambda exonuclease, thus
608 conversion of linear double-stranded DNA to single-stranded DNA (ssDNA). The mixture was
609 purified by Eastep Gel and PCR Cleanup Kit. Then 10% denaturing (7 mol/L urea) PAGE was
610 used to analyze the degraded products. Gel band quantitation was used to assess the yield of

611 ssDNA. Azurespot software was subsequently used to perform band detection, background
612 subtraction and band quantitation. More detailed steps are in Supplementary Note 10.

613 **Oligo pool normalizing (OPN).** The schematic of OPN is illustrated (Supplementary Fig. 61 and
614 62). The OPN probes were synthesized respectively (Supplementary Table 4 and Table 5). An
615 equimolar mixture of OPN probes (256 or 1024), which the number of each OPN probe was
616 equivalent to the average molar concentration of the oligo pool, was applied to capture the
617 corresponding oligo separately. The 256/1024 oligos were mixed with 256/1024 OPN probes and
618 hybridization buffer (10 mM Tris-EDTA, 0.5 M NaCl, and 0.05% Tween-20 (volume / volume))
619 in 20 μ L reaction. The mixture was denatured at 95°C for 3 min and slowly cooled to 60°C at a
620 ramp of 0.1°C/s, following kept for 2 h at 60°C using an Eppendorf Mastercycler instrument. Then,
621 extension reaction was carried out when the target was captured by corresponding OPN probe. To
622 ensure temperature uniformity, the pre-reaction mixture containing Taq DNA polymerase and
623 dNTPs was also pre-heated to 60°C before adding to the ssDNA/OPN probe mixture. The resulting
624 mixture was incubated for another 15 min at 60°C to obtain dsDNA product. Further, Exo I was
625 added to the resulting mixture to digest the remaining ssDNA and OPN probes. After Exo I was
626 inactivated, Streptavidin Magnetic Beads were added and incubated for another 30 min at 37°C in
627 shaker to isolate dsDNA product. The dsDNA product attached to magnetic beads was the template
628 for OPN-iDR. More detailed steps are in Supplementary Note 10.

629 **Data availability.** The BASIC code for encoding and decoding for both Linux and Windows and
630 bioinformatic analysis programs may be obtained via ([https://github.com/xiaomingao/DNA-](https://github.com/xiaomingao/DNA-information-storage)
631 [information-storage](https://github.com/xiaomingao/DNA-information-storage)). Furthermore, the original sequencing FASTQ file and the designed sequence
632 file may be obtained via
633 (<http://pan.tju.edu.cn:80/#/link/627DB5C9EB819984F1183D8D4A0B72E3>, Code: oZ3D).