

## Random access DNA memory in a scalable, archival file storage system

James L. Banal<sup>1†</sup>, Tyson R. Shepherd<sup>1†</sup>, Joseph Berleant<sup>1†</sup>, Hellen Huang<sup>1</sup>, Miguel Reyes<sup>1,2</sup>,

Cheri M. Ackerman<sup>2</sup>, Paul C. Blainey<sup>1,2,3</sup>, and Mark Bathe<sup>1,2\*</sup>

<sup>1</sup>Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA  
02139 USA

<sup>2</sup>Broad Institute of MIT and Harvard, Cambridge, MA 02142 USA

<sup>3</sup>Koch Institute for Integrative Cancer Research at MIT, Cambridge, MA 02142 USA

\*Correspondence should be addressed to: [mark.bathe@mit.edu](mailto:mark.bathe@mit.edu)

†These authors contributed equally to this work

## 1 **ABSTRACT**

2 DNA is an ultra-high-density storage medium that could meet exponentially growing worldwide  
3 demand for archival data storage if DNA synthesis costs declined sufficiently and random access  
4 of files within exabyte-to-yottabyte-scale DNA data pools were feasible. To overcome the second  
5 barrier, here we encapsulate data-encoding DNA file sequences within impervious silica capsules  
6 that are surface-labeled with single-stranded DNA barcodes. Barcodes are chosen to represent file  
7 metadata, enabling efficient and direct selection of sets of files with Boolean logic. We  
8 demonstrate random access of image files from an image database using fluorescence sorting with  
9 selection sensitivity of 1 in  $10^6$  files, which thereby enables 1 in  $10^{6N}$  per  $N$  optical channels. Our  
10 strategy thereby offers retrieval of random file subsets from exabyte and larger-scale long-term  
11 DNA file storage databases, offering a scalable solution for random-access of archival files in  
12 massive molecular datasets.

13

## 14 **INTRODUCTION**

15 While DNA is conventionally the polymer used for storage and transmission of genetic  
16 information in biology, it can also be used for the storage of arbitrary digital information at  
17 densities far exceeding conventional data storage technologies such as flash and tape memory, at  
18 scales well beyond the capacity of the largest current data centers<sup>1,2</sup>. Recent progress in nucleic  
19 acid synthesis and sequencing technologies continue to reduce the cost of writing and reading  
20 DNA, thereby rendering DNA-based information storage potentially viable commercially in the  
21 future<sup>3-6</sup>. Demonstrations of its viability as a general information storage medium include  
22 numerous examples including the storage and retrieval of books, images, computer programs,  
23 audio clips, works of art, and Shakespeare's sonnets using a variety of encoding schemes<sup>7-13</sup>, with

24 data size limited primarily by the cost of DNA synthesis. In each case, digital information was  
25 converted to DNA sequences composed of ~100–200 nucleotide (nt) data blocks for ease of  
26 chemical synthesis and sequencing. Sequence fragments were then assembled to reconstruct the  
27 original, encoded information.

28         While significant effort in DNA data storage has focused on increasing the scale of DNA  
29 synthesis, as well as improving encoding schemes, an additional crucial aspect of a successful  
30 molecular data storage system is the ability to efficiently retrieve specific files, or random subsets  
31 of files, from a large-scale pool of DNA data on demand, without error, without data destruction,  
32 and ideally at low cost for a practical archival data storage and retrieval device. Toward this end,  
33 to date research has largely used conventional polymerase chain reaction (PCR)<sup>9,11,13</sup>, which uses  
34 up to 20–30 heating and cooling cycles with DNA polymerase to selectively amplify and extract  
35 specific DNA sequences from a DNA data pool using primers. Nested addressing barcodes<sup>14-16</sup>  
36 have also been used to uniquely identify a greater number of files, as well as biochemical affinity  
37 tags to selectively pull down oligos for targeted amplification<sup>17</sup>.

38         Major limitations of PCR-based approaches, however, include the length of DNA needed  
39 to uniquely label DNA data strands for file indexing, which dramatically reduces the DNA  
40 available for data storage. For example, for an exabyte-scale data pool, each file requires at least  
41 three barcodes<sup>17</sup>, or up to sixty nucleotides in total barcode sequence length, thereby reducing the  
42 number of nucleotides that can be used for data encoding. Further, selective amplification of a  
43 specific file using PCR requires access to the entire data pool for each query, which is destructive  
44 to the data pool, and intrinsically limited by the finite number orthogonal primers, e.g., 28,000 for  
45 previously demonstrated PCR-based random access system<sup>13</sup>, available to amplify target files  
46 without strand crosstalk due to non-specific hybridization. Finally, PCR-based approaches do not

47 allow for physical deletion of specific files from a data pool and require numerous heating and  
48 cooling cycles with DNA polymerase, which may be prohibitively costly, time-consuming, and  
49 impractical for random access memory in exabyte-to-yottabyte-scale data pools. While spatial  
50 segregation of data into distinct pools<sup>18</sup> and extraction of selected DNA using biochemical affinity  
51 pulldown have yielded significant improvements in PCR-based file selection strategies, these  
52 implementations vastly reduce data density<sup>17</sup>, and cannot access random subsets of files in this  
53 direct manner that is required for a truly scalable and deployable archival molecular file storage  
54 and retrieval system.

55 As an alternative to PCR-based approaches, here we focus on archival DNA data storage  
56 and retrieval by first encapsulating physically DNA-based files within discrete, impervious silica  
57 capsules, which we subsequently label with single-stranded DNA barcodes that enable direct,  
58 random access on the entire data pool via barcode hybridization, without need for amplification  
59 and without crosstalk with the physically isolated data-encoding DNA, followed by downstream  
60 selection that may be optical, physical, or biochemical. Each “unit of information” encoded in  
61 DNA we term a *file*, which includes both the DNA encoding the main data as well as any additional  
62 components used for addressing, storage, and retrieval. Each file contains a *file sequence*,  
63 consisting of the DNA encoding the main data, and *addressing barcodes*, or simply *barcodes*,  
64 which are additional short DNA sequences used to identify the file in solution using hybridization.  
65 We refer to a collection of files as a *data pool* or *database*, and the set of procedures for storing,  
66 retrieving, and reading out files is termed a *file system* (see **Supplementary Section S0** for a full  
67 list of terms).

68 As a proof-of-principle of our archival DNA file system, we encapsulated 20 image files,  
69 each composed of a ~0.1 kilobyte image file encoded in a 3,000-base-pair plasmid, within

70 monodisperse, 6- $\mu\text{m}$  silica particles that were chemically surface-labeled using up to three 25-mer  
71 single-stranded DNA (ssDNA) oligonucleotide barcodes chosen from a library of 240,000  
72 orthogonal primers, which allows for identification of up to  $\sim 10^{15}$  possible distinct files using only  
73 three unique barcodes per file<sup>19</sup> (**Fig. 1**). While we chose plasmids to encode DNA data in order  
74 to produce microgram quantities of DNA memory at low cost and to facilitate a renewable, closed-  
75 cycle write-store-access-read system using bacterial DNA data encoding and expression<sup>20-22</sup>, our  
76 file system is equally applicable to single-stranded DNA oligos produced using solid-phase  
77 chemical synthesis<sup>2,7,8,10-13,17</sup> or gene-length oligos produced enzymatically<sup>23-26</sup>, and larger file  
78 sizes on the megabyte to gigabyte scale. And while only twenty icon-resolution images were  
79 chosen as our image database, representing diverse subject matter including animals, plants,  
80 transportation, and buildings (**Supplementary Fig. 1**), our file system equally applies to  
81 thousands, billions, or larger sets of images, limited only by the cost of DNA synthesis, rather than  
82 any intrinsic property of our file system itself (**Supplementary Fig. 1**).

83 Fluorescence-activated sorting (FAS) was used to select target subsets of the complete data  
84 pool by first annealing fluorescent oligonucleotide probes that are complementary to the barcodes  
85 used to address the database<sup>27</sup>, enabling direct retrieval of specific, individual files from a pool of  
86  $(10^6)^N$  total files, where  $N$  is the number of fluorescence channels employed, without amplification  
87 required for PCR-based approaches, or loss of nucleotides available for data encoding. Further,  
88 our system enables direct, complex Boolean AND, OR, NOT logic to select random subsets of  
89 files with combinations of distinct barcodes to query the data pool, similar to conventional Boolean  
90 logic applied in text and file searches on solid-state silicon devices. And because physical  
91 encapsulation separates file sequences from external barcodes that are used to describe the  
92 encapsulated information, our file system offers long-term environmental protection of encoded

93 file sequences via silica encapsulation for permanent archival storage<sup>10,28,29</sup>, where external  
94 barcodes may be renewed periodically, further protected with secondary encapsulation, or replaced  
95 for more sophisticated file operations involving re-labeling of data pools. Taken together, our  
96 strategy presents a practical and scalable archival molecular file storage system with random  
97 access capability that applies to the exabyte-to-yottabyte scales, limited only by the current cost of  
98 DNA synthesis.

99

## 100 **File Synthesis**

101 Digital information in the form of 20 icon-resolution images was stored in a data pool, with each  
102 image encoded into DNA and synthesized on a plasmid. We selected images of broad diversity,  
103 representative of distinct and shared subject categories, which included several domestic and wild  
104 cats and dogs, US presidents, and several human-made objects such as an airplane, boats, and  
105 buildings (**Fig. 1** and **Supplementary Fig. 1**). To implement this image database, the images were  
106 substituted with black-and-white,  $26 \times 26$ -pixel images to minimize synthesis costs, compressed  
107 using run-length encoding, and converted to DNA (**Supplementary Fig. 1, 2**). Following  
108 synthesis, bacterial amplification, and sequencing validation (**Supplementary Fig. 3**), each  
109 plasmid DNA was separately encapsulated into silica particles containing a fluorescein dye core  
110 and a positively charged surface<sup>28,29</sup>. Because the negatively charged phosphate groups of the DNA  
111 interact with positively charged silica particles, plasmid DNA condensed on the silica surface,  
112 after which N-[3-(trimethoxysilyl)propyl]-N,N,N-trimethylammonium chloride (TMAPS) was  
113 co-condensed with tetraethoxysilane to form an encapsulation shell after four days of incubation  
114 at room-temperature<sup>10,29</sup> (**Fig. 2a**) to form discrete silica capsules containing the file sequence that  
115 encodes for the image file. Quantitative PCR (qPCR) of the reaction supernatant after

116 encapsulation (**Supplementary Fig. 4**) showed full encapsulation of plasmids without residual  
117 DNA in solution. To investigate the fraction of capsules that contained plasmid DNA, we  
118 compared the fluorescence intensity of the intercalating dye TO-PRO when added pre- versus post-  
119 encapsulation (**Supplementary Fig. 2**). All capsules synthesized in the presence of both DNA and  
120 TO-PRO showed a distinct fluorescence signal, consistent with the presence of plasmid DNA in  
121 the majority of capsules, compared with a silica particle negative control that contained no DNA.  
122 In order to test whether plasmid DNA was fully encapsulated versus partially exposed at the  
123 surface of capsules, capsules were also stained separately with TO-PRO post-encapsulation (**Fig.**  
124 **2b**). Using qPCR, we estimated  $10^6$  plasmids per capsule assuming quantitative recovery of DNA  
125 post-encapsulation (**Supplementary Fig. 5**). Because encapsulation of the DNA file sequence  
126 relies only on electrostatic interactions between positively-charged silica and the phosphate  
127 backbone of DNA, our approach can equally encapsulate any molecular weight of DNA molecule  
128 applicable to MB and larger file sizes, as demonstrated previously<sup>29</sup>, and is compatible with  
129 alternative DNA file compositions such as 100-200-mer oligonucleotides that are commonly  
130 used<sup>2,7,8,12,13,17</sup>.

131         Next, we chemically attached unique content addresses on the surfaces of silica capsules  
132 using orthogonal 25-mer ssDNA barcodes (**Supplementary Fig. 6**) describing selected features of  
133 the underlying image for file selection. For example, the image of an orange tabby house cat  
134 (**Supplementary Fig. 1**) was described with *cat*, *orange*, and *domestic*, whereas the image of a  
135 tiger was described with *cat*, *orange*, and *wild* (**Supplementary Fig. 1** and **Supplementary Table**  
136 **2**). To attach the barcodes, we activated the surface of the silica capsules through a series of  
137 chemical steps. Condensation of 3-aminopropyltriethoxysilane with the hydroxy-terminated  
138 surface of the encapsulated plasmid DNA provided a primary amine chemical handle that

139 supported further conjugation reactions (**Fig. 2c**). We modified the amino-modified surface of the  
140 silica capsules with 2-azidoacetic acid N-hydroxysuccinimide (NHS) ester followed by an  
141 oligo(ethylene glycol) that contained two chemically orthogonal functional groups: the  
142 dibenzocyclooctyne functional group reacted with the surface-attached azide through strain-  
143 promoted azide-alkyne cycloaddition while the NHS ester functional group was available for  
144 subsequent conjugation with a primary amine. Each of the associated barcodes contained a 5'-  
145 amino modification that could react with the NHS-ester groups on the surface of the silica capsules,  
146 thereby producing the complete form of our file. Notably, the sizes of bare, hydroxy-terminated  
147 silica particles representing capsules without barcodes were comparable with complete files  
148 consisting of capsules with barcodes attached, confirmed using scanning electron microscopy (**Fig.**  
149 **2d** and **2e**, left). These results were anticipated given that the encapsulation thickness was only on  
150 the order of 10 nm<sup>29</sup> and that additional steps to attach functional groups minimally increases the  
151 capsule diameter. We also observed systematic shifts in the surface charge of the silica particles  
152 as different functional groups were introduced onto their surfaces (**Fig. 2e**). Using hybridization  
153 assays with fluorescently-labelled probes<sup>30-32</sup>, we estimated the number of barcodes available for  
154 hybridization on each file to be on the order of 10<sup>8</sup> (**Supplementary Fig. 7**). Following synthesis,  
155 files were pooled and stored together for subsequent retrieval. Illumina MiSeq was used to read  
156 each file sequence and reconstruct the encoded image following selection and de-encapsulation,  
157 in order to validate the complete process of image file encoding, encapsulation, barcoding,  
158 selection, de-encapsulation, sequencing, and image file reconstruction (**Supplementary Figs. 9,**  
159 **10**).

160

161 **File Selection**



162 Following file synthesis and pooling, we used FAS to select specific targeted files from the  
163 complete data pool through the reversible binding of fluorescent probe molecules to the file  
164 barcodes (**Supplementary Fig. 6**). All files contained a fluorescent dye, fluorescein, in their core  
165 as a marker to distinguish files from other particulates such as spurious silica particles that  
166 nucleated in the absence of a core or insoluble salts that may have formed during the sorting  
167 process. Each detected fluorescein event was therefore interpreted to indicate the presence of a  
168 single file during FAS (**Supplementary Fig. 11**). To apply a query such as *flying* to the image  
169 database, the corresponding fluorescently labeled ssDNA probe was added, which hybridized to  
170 the complementary barcode displayed externally on the surface of a silica capsule for FAS  
171 selection (**Fig. 3a**).

172 We subjected the entire data pool to a series of experiments to test selection sensitivity of  
173 target subsets using distinct queries. First, we evaluated single-barcode selection of an individual  
174 file, specifically *Airplane*, out of a pool of varying concentrations of the nineteen other files as  
175 background (**Fig. 3b**). To select the *Airplane* file, we hybridized an AFDye 647-labelled ssDNA  
176 probe that is complementary to the barcode *flying*, which is unique to *Airplane*. We were able to  
177 detect and select the desired *Airplane* file through FAS even at a relative abundance of  $10^{-6}$   
178 compared with each other file (**Fig. 3c**). While comparable in sensitivity to a nested PCR barcoding  
179 data indexing approach<sup>17</sup>, unlike PCR that requires 20–30 of rounds of heating and cooling to  
180 selectively amplify the selected sequence, our approach selects files directly without need for  
181 thermal cycling and amplification. This strategy also applies to gating of  $N$  barcodes  
182 simultaneously in parallel optical channels, which offers file selection sensitivity of 1 in  $10^{6N}$  total  
183 files, where common commercial FAS systems offer up to  $N = 17$  channels<sup>33,34</sup>. For example,  
184 comparison of the retrieved sequences between the *flying* gate and the NOT *flying* gate after

185 chemical release of the file sequences from silica encapsulation revealed that 60–95% of the  
186 *Airplane* files were sorted into the `flying` gate (**Supplementary Figs. 18–21**), where we note  
187 that any sort probability above 50% indicates enrichment of *Airplane* within the correct population  
188 subset (`flying`) relative to the incorrect subset (`NOT flying`), while a sort probability of 100%  
189 would indicate ideal performance. Besides single file selection, our approach allows for repeated  
190 rounds of FAS selection, as well as Boolean logic, described below.

191

## 192 **Boolean Search**

193 Beyond direct selection of 1 in  $10^{6N}$  individual random files directly, without thermal cycling or  
194 loss of fidelity due to primer crosstalk, our system offers the ability to apply Boolean logic to select  
195 random file subsets from the data pool. AND, OR, and NOT logical operations were applied by  
196 first adding to the data pool fluorescently labeled ssDNA probes that were complementary to the  
197 barcodes (**Fig. 4**, left). This hybridization reaction was used to distinguish one or several files in  
198 the data pool, which were then sorted using FAS. We used two to four fluorescence channels  
199 simultaneously to create the FAS gates that executed the target Boolean logic queries (**Fig. 4**,  
200 middle). To demonstrate a NOT query, we added to the data pool an AFDye 647-labelled ssDNA  
201 probe that hybridized to files that contained the *cat* barcode. Files that did not show AFDye 647  
202 signal were sorted into the `NOT cat` subset (**Fig. 4a**). An example of an OR gate was applied to  
203 the data pool by simultaneously adding *dog* and *building* probes that both had the TAMRA label  
204 (**Fig. 4b**). All files that showed TAMRA signal were sorted into the `dog OR building` subset  
205 by the FAS. Finally, an example of an AND gate was achieved by adding *fruit* and *yellow* probes  
206 that were labelled with AFDye 647 and TAMRA, respectively. Files showing signal for both  
207 AFDye 647 and TAMRA were sorted into the `fruit AND yellow` subset in the FAS (**Fig. 4c**).

208 For each example query, we validated our sorting experiments by releasing the file sequence from  
209 silica encapsulation and sequencing the released DNA with Illumina MiniSeq (**Fig. 4**, right). Sort  
210 probabilities of each file for each search query are shown in **Supplementary Figs. S22–S24**.

211 The preceding demonstrations of Boolean logic gates enable file sorting with varying  
212 specificity of selection criteria for the retrieval of different subsets of the data pool. FAS can also  
213 be used to create multiple gating conditions simultaneously, thereby increasing the complexity of  
214 target file selection operations, as noted above. To demonstrate increasingly complex Boolean  
215 search queries, we selected the file containing the image of Abraham Lincoln from the data pool,  
216 which included images of two presidents, George Washington and Abraham Lincoln. The  
217 *president* ssDNA probe, fluorescently labeled with TAMRA, selected both *Lincoln* and  
218 *Washington* files from the data pool. The simultaneous addition of the *18<sup>th</sup> century* ssDNA probe,  
219 fluorescently labeled with AFDye 647 (**Fig. 5a**, left), discriminated *Washington*, which contained  
220 the *18<sup>th</sup> century* barcode, from the *Lincoln* file (**Fig. 5a**, middle). The combination of these two  
221 ssDNA probes permitted the complex search query `president AND (NOT 18th century)`.  
222 Sequencing analysis of the gated populations after reverse encapsulation validated that the sorted  
223 populations matched search queries for `president AND (NOT 18th century)`,  
224 `president AND 18th century`, and `NOT president` (**Fig. 5a**, right; **Supplementary**  
225 **Fig. 25**).

226 To demonstrate the feasibility of performing Boolean search using more than three  
227 fluorescence channels for sorting, we also selected the *Wolf* file from the data pool using the query  
228 `dog AND wild`, and used the *black & white* probe to validate the selected file (**Fig. 5b**, left).  
229 Because conventional FAS software is only capable of sorting using 1D and 2D gates, we first  
230 selected one out of the three possible 2D plots (**Fig. 5b**, left and bottom): *dog*-TAMRA against

231 *wild*-AFDye 647. We examined the *black & white*-TYE705 channel on members of the *dog* AND  
232 *wild* subset (**Fig. 5b**, left and bottom). Release of the encapsulated file sequence and subsequent  
233 sequencing of each gated population from the *dog* versus *wild* 2D plot validated sorting (**Fig. 5b**,  
234 right; **Supplementary Fig. 26**).

235 In contrast to single-stranded DNA oligos, our use of plasmids as a substrate for encoding  
236 information offered the ability to restore files into the data pool after retrieval. In cases where  
237 single images were sorted (**Figs. 4c, 5a, b**), we were able to transform competent bacteria from  
238 each search query that resulted in a single file (**Supplementary Fig. 27**). Amplified material was  
239 pure and ready for re-encapsulation into silica particles, which could be re-introduced directly back  
240 into the data pool. Importantly, our molecular file system and file selection process thereby  
241 represents a complete write-store-access-read cycle that in principle may be applied to exabyte and  
242 larger-scale datasets, with periodic renewal of single-stranded DNA barcodes and bacterial  
243 replication of DNA data following reading<sup>20-22</sup>. While sort probabilities were typically below the  
244 optimal 100% targeted for a specific file or file subset query, future work may characterize sources  
245 of error that could be due to sample contamination or random FAS errors. The latter type of error  
246 can be mitigated through repeated cycles of file selection in series. Our technical approach differs  
247 significantly from approaches that rely on selective PCR amplification for selection<sup>9,11,13,17,18</sup>, in  
248 which repeated amplifications may reduce fidelity of file selection.

249

## 250 **Discussion & Outlook**

251 We introduce a scalable, non-destructive, random access molecular file system for the direct access  
252 of arbitrary files and file-subsets from an archival DNA data store. The introduction of our file  
253 system overcomes former limitations of indirect, PCR-based file systems for the practical

254 implementation of archival DNA memory systems. This advance now leaves the high cost of DNA  
255 synthesis compared with alternative memory storage media as the primary remaining rate-limiting  
256 step for translation of this technology. While the overall data density of our file system is  
257 considerably lower than the theoretical limit of DNA data density due to the encapsulation of DNA  
258 files in silica particles, the physical size of exabyte-scale DNA data stored in our system is still  
259 orders of magnitude smaller than conventional archival file storage systems. For example,  
260 assuming 2 bits per base,  $10^{-21}$  grams per base, and a density of double-stranded DNA of 1.7 grams  
261 per cubic centimeter<sup>4</sup>, PCR-based random access approaches have a theoretical volumetric density  
262 limit of  $10^{27}$  bytes per  $m^3$ , compared with our approach of  $10^{24}$  bytes per  $m^3$  that is  $10^3$ -fold lower  
263 (**Supplementary Section S6**). However, PCR suffers from numerous issues such as enzyme cost,  
264 requirement of numerous heating and cooling cycles, and potential crosstalk between file  
265 sequences and barcodes<sup>17,18</sup>, which requires spatial segregation of file sequences in electrowetting  
266 devices<sup>18</sup> that reduced data density to  $\sim 10^{20}$  bytes per  $m^3$ , seven orders of magnitude below the  
267 theoretical limit for dry DNA (**Supplementary Section S6**).

268 In the current implementation of our file system, each file capsule contained  $10^6$  DNA  
269 plasmids, which could instead store multiple unique file-encoding plasmids or file fragments to  
270 increase data density to gigabyte-sized files per capsule, with an overall data density of  $10^{24}$  bytes  
271 per  $m^3$  (**Supplementary Section S6**), which is only three orders of magnitude lower than the  
272 theoretical data density limit of dry DNA, and four orders of magnitude higher than published  
273 approaches to storing and accessing DNA data with spatial segregation<sup>18</sup>. And equally important  
274 to data density per se is the physical size required to store an exabyte- or larger-scale DNA data  
275 pool. Using our approach,  $10^9$  gigabyte-sized files would still only require  $0.2 \text{ cm}^3$  of total dry  
276 volume, without any need for physically separated data pools. Notwithstanding, further increases

277 in data density could be achieved by using nanoparticles ~100–200 nm in diameter to encode  
278 files<sup>10,28,29</sup> sorted with higher sensitivity FAS systems<sup>35,36</sup>, or multiple layers of encapsulated  
279 DNA<sup>37</sup>.

280 In addition to data pool size and density, another crucial operating feature is the latency or  
281 time associated with DNA file retrieval. Because FAS scales linearly with the size of the data pool,  
282 retrieval time may still be limiting in an exabyte-scale data pool, even assuming gigabyte-sized  
283 files. To further reduce file selection time, future file system implementations may leverage  
284 parallel microfluidics-based optical sorting procedures, brighter fluorescent probes to increase  
285 selection throughput, alternative barcode implementations<sup>38-42</sup>, or physical sorting strategies such  
286 as direct biochemical pulldown<sup>17,43,44</sup>, such as recently implemented using direct magnetic  
287 extraction of files labelled with biochemical affinity tags<sup>17</sup>. Additional latency due to chemical  
288 deprotection of DNA from silica encapsulation renders our file system ideally suited to long-term,  
289 archival DNA storage at the exabyte-to-yottabyte scales.

290 Indeed, because we view our scalable file system as an alternative to tape-based, ‘cold’  
291 archival data storage systems rather than flash or other ‘hot’ memory, for which latency times may  
292 be tolerated on the time frame of several days to weeks, the foregoing latency limitations are of  
293 minimal importance compared with the transformative capability offered by our system to store  
294 exabyte-to-yottabyte-scale datasets with direct retrieval of arbitrary, random file subsets. Example  
295 applications include the retrieval of specific images from archival databases of astronomical image  
296 databases<sup>45</sup>, high-energy physics datasets<sup>46</sup>, or high-resolution deep ocean floor mapping<sup>47</sup>.

297 Finally, because our system is not limited to synthetic DNA, it applies equally to long-term  
298 archival storage of bacterial, human, and other genomes for archival sample preservation and  
299 retrieval<sup>23,48</sup>, forensic analysis, and retrospective analysis of pandemic outbreaks, as explored in

300 accompanying work<sup>49</sup>. Our demonstrated file system enables complex file search operations on  
301 underlying molecular data pools, moving us closer to realizing an economically viable, functional  
302 massive molecular file and operating system<sup>27,50,51</sup>.

303

## 304 **References**

- 305 1 Zhirnov, V., Zadegan, R. M., Sandhu, G. S., Church, G. M. & Hughes, W. L. Nucleic  
306 acid memory. *Nature Materials* **15**, 366-370, doi:10.1038/nmat4594 (2016).
- 307 2 Ceze, L., Nivala, J. & Strauss, K. Molecular digital data storage using DNA. *Nature*  
308 *Reviews Genetics* **20**, 456-466, doi:10.1038/s41576-019-0125-3 (2019).
- 309 3 Kosuri, S. & Church, G. M. Large-scale de novo DNA synthesis: technologies and  
310 applications. *Nature Methods* **11**, 499-507, doi:10.1038/nmeth.2918 (2014).
- 311 4 Zhirnov, V., Zadegan, R. M., Sandhu, G. S., Church, G. M. & Hughes, W. L. Nucleic  
312 acid memory. *Nature Materials* **15**, 366 (2016).
- 313 5 Palluk, S. *et al.* De novo DNA synthesis using polymerase-nucleotide conjugates. *Nature*  
314 *Biotechnology* **36**, 645-650, doi:10.1038/nbt.4173 (2018).
- 315 6 Lee, H. H., Kalhor, R., Goela, N., Bolot, J. & Church, G. M. Terminator-free template-  
316 independent enzymatic DNA synthesis for digital information storage. *Nature*  
317 *Communications* **10**, 2383, doi:10.1038/s41467-019-10258-1 (2019).
- 318 7 Church, G. M., Gao, Y. & Kosuri, S. Next-generation digital information storage in  
319 DNA. *Science* **337**, 1628-1628 (2012).
- 320 8 Goldman, N. *et al.* Towards practical, high-capacity, low-maintenance information  
321 storage in synthesized DNA. *Nature* **494**, 77 (2013).

- 322 9 Yazdi, S. M. H. T., Yuan, Y., Ma, J., Zhao, H. & Milenkovic, O. A rewritable, random-  
323 access DNA-based storage system. *Scientific Reports* **5**, 14138, doi:10.1038/srep14138  
324 (2015).
- 325 10 Grass, R. N., Heckel, R., Puddu, M., Paunescu, D. & Stark, W. J. Robust chemical  
326 preservation of digital information on DNA in silica with error-correcting codes.  
327 *Angewandte Chemie International Edition* **54**, 2552-2555 (2015).
- 328 11 Yazdi, S. M. H. T., Gabrys, R. & Milenkovic, O. Portable and error-free DNA-based data  
329 storage. *Scientific Reports* **7**, 5011, doi:10.1038/s41598-017-05188-1 (2017).
- 330 12 Erlich, Y. & Zielinski, D. DNA Fountain enables a robust and efficient storage  
331 architecture. *Science* **355**, 950-954 (2017).
- 332 13 Organick, L. *et al.* Random access in large-scale DNA data storage. *Nature*  
333 *Biotechnology* **36**, 242 (2018).
- 334 14 Kashiwamura, S., Yamamoto, M., Kameda, A., Shiba, T. & Ohuchi, A. in *8th*  
335 *International Workshop on DNA-Based Computers (DNA8)*. 112-123 (Springer).
- 336 15 Yamamoto, M., Kashiwamura, S., Ohuchi, A. & Furukawa, M. Large-scale DNA  
337 memory based on the nested PCR. *Natural Computing* **7**, 335-346 (2008).
- 338 16 Yamamoto, M., Kashiwamura, S. & Ohuchi, A. in *13th International Meeting on DNA*  
339 *Computing (DNA13)*. 99-108 (Springer).
- 340 17 Tomek, K. J. *et al.* Driving the scalability of DNA-based information storage systems.  
341 *ACS Synthetic Biology* **8**, 1241-1248, doi:10.1021/acssynbio.9b00100 (2019).
- 342 18 Newman, S. *et al.* High density DNA data storage library via dehydration with digital  
343 microfluidic retrieval. *Nature Communications* **10**, 1706 (2019).



- 344 19 Xu, Q., Schlabach, M. R., Hannon, G. J. & Elledge, S. J. Design of 240,000 orthogonal  
345 25mer DNA barcode probes. *Proceedings of the National Academy of Sciences* **106**,  
346 2289-2294, doi:10.1073/pnas.0812506106 (2009).
- 347 20 Farzadfard, F. *et al.* Single-nucleotide-resolution computing and memory in living cells.  
348 *Molecular Cell* **75**, 769-780. e764 (2019).
- 349 21 Farzadfard, F. & Lu, T. K. Genomically encoded analog memory with precise in vivo  
350 DNA writing in living cell populations. *Science* **346**, 1256272,  
351 doi:10.1126/science.1256272 (2014).
- 352 22 Farzadfard, F. & Lu, T. K. Emerging applications for DNA writers and molecular  
353 recorders. *Science* **361**, 870-875, doi:10.1126/science.aat9249 (2018).
- 354 23 Plesa, C., Sidore, A. M., Lubock, N. B., Zhang, D. & Kosuri, S. Multiplexed gene  
355 synthesis in emulsions for exploring protein functional landscapes. *Science* **359**, 343-347,  
356 doi:10.1126/science.aao5167 (2018).
- 357 24 Shepherd, T. R., Du, R. R., Huang, H., Wamhoff, E.-C. & Bathe, M. Bioproduction of  
358 pure, kilobase-scale single-stranded DNA. *Scientific Reports* **9**, 6121,  
359 doi:10.1038/s41598-019-42665-1 (2019).
- 360 25 Veneziano, R. *et al.* In vitro synthesis of gene-length single-stranded DNA. *Scientific*  
361 *Reports* **8**, 1-7 (2018).
- 362 26 Minev, D. *et al.* Rapid in vitro production of single-stranded DNA. *Nucleic Acids*  
363 *Research* **47**, 11956-11962, doi:10.1093/nar/gkz998 (2019).
- 364 27 Reif, J. H. *et al.* in *7th International Workshop on DNA-Based Computers (DNA7)*. 231-  
365 247 (Springer Berlin Heidelberg).

- 366 28 Paunescu, D., Fuhrer, R. & Grass, R. N. Protection and deprotection of DNA--high-  
367 temperature stability of nucleic acid barcodes for polymer labeling. *Angewandte Chemie*  
368 *International Edition* **52**, 4269-4272, doi:10.1002/anie.201208135 (2013).
- 369 29 Paunescu, D., Puddu, M., Soellner, J. O. B., Stoessel, P. R. & Grass, R. N. Reversible  
370 DNA encapsulation in silica to produce ROS-resistant and heat-resistant synthetic DNA  
371 "fossils". *Nature Protocols* **8**, 2440, doi:10.1038/nprot.2013.154 (2013).
- 372 30 Pillai, P. P., Reisewitz, S., Schroeder, H. & Niemeyer, C. M. Quantum-dot-encoded silica  
373 nanospheres for nucleic acid hybridization. *Small* **6**, 2130-2134,  
374 doi:10.1002/sml.201000949 (2010).
- 375 31 Leidner, A. *et al.* Biopebbles: DNA-functionalized core-shell silica nanospheres for  
376 cellular uptake and cell guidance studies. *Advanced Functional Materials* **28**, 1707572,  
377 doi:10.1002/adfm.201707572 (2018).
- 378 32 Sun, P. *et al.* Biopebble containers: DNA-directed surface assembly of mesoporous silica  
379 nanoparticles for cell studies. *Small* **15**, 1900083, doi:10.1002/sml.201900083 (2019).
- 380 33 Perfetto, S. P., Chattopadhyay, P. K. & Roederer, M. Seventeen-colour flow cytometry:  
381 unravelling the immune system. *Nature Reviews Immunology* **4**, 648-655,  
382 doi:10.1038/nri1416 (2004).
- 383 34 Chattopadhyay, P. K. *et al.* Quantum dot semiconductor nanocrystals for  
384 immunophenotyping by polychromatic flow cytometry. *Nature Medicine* **12**, 972-977,  
385 doi:10.1038/nm1371 (2006).
- 386 35 van Gaal, E. V. B., Spierenburg, G., Hennink, W. E., Crommelin, D. J. A. &  
387 Mastrobattista, E. Flow cytometry for rapid size determination and sorting of nucleic acid

- 388 containing nanoparticles in biological fluids. *Journal of Controlled Release* **141**, 328-  
389 338, doi:10.1016/j.jconrel.2009.09.009 (2010).
- 390 36 Lian, H., He, S., Chen, C. & Yan, X. Flow cytometric analysis of nanoscale biological  
391 particles and organelles. *Annual Review of Analytical Chemistry* **12**, 389-409,  
392 doi:10.1146/annurev-anchem-061318-115042 (2019).
- 393 37 Ablasser, A. & Chen, Z. J. cGAS in action: Expanding roles in immunity and  
394 inflammation. *Science* **363**, 1055+, doi:10.1126/science.aat8657 (2019).
- 395 38 Braeckmans, K. *et al.* Encoding microcarriers by spatial selective photobleaching. *Nature*  
396 *Materials* **2**, 169-173, doi:10.1038/nmat828 (2003).
- 397 39 Wilson, R., Cossins, A. R. & Spiller, D. G. Encoded microcarriers for high-throughput  
398 multiplexed detection. *Angewandte Chemie International Edition* **45**, 6104-6117,  
399 doi:10.1002/anie.200600288 (2006).
- 400 40 Pregibon, D. C., Toner, M. & Doyle, P. S. Multifunctional encoded particles for high-  
401 throughput biomolecule analysis. *Science* **315**, 1393-1396, doi:10.1126/science.1134929  
402 (2007).
- 403 41 Dagher, M., Kleinman, M., Ng, A. & Juncker, D. Ensemble multicolour FRET model  
404 enables barcoding at extreme FRET levels. *Nature Nanotechnology* **13**, 925-932,  
405 doi:10.1038/s41565-018-0205-0 (2018).
- 406 42 Martino, N. *et al.* Wavelength-encoded laser particles for massively multiplexed cell  
407 tagging. *Nature Photonics* **13**, 720-727, doi:10.1038/s41566-019-0489-0 (2019).
- 408 43 Lee, H., Kim, J., Kim, H., Kim, J. & Kwon, S. Colour-barcoded magnetic microparticles  
409 for multiplexed bioassays. *Nature Materials* **9**, 745-749, doi:10.1038/nmat2815 (2010).

- 410 44 Stewart, K. *et al.* in *24th International Conference on DNA Computing and Molecular*  
411 *Programming (DNA 24)*. 55-70 (Springer).
- 412 45 Broekema, P. C., Nieuwpoort, R. V. v. & Bal, H. E. in *Proceedings of the 2012 workshop*  
413 *on High-Performance Computing for Astronomy Data* 9–16 (Association for  
414 Computing Machinery, Delft, The Netherlands, 2012).
- 415 46 Gaillard, M. & Pandolfi, S. *CERN Data Centre passes the 200-petabyte milestone*,  
416 <https://cds.cern.ch/record/2276551> (2017).
- 417 47 Mayer, L. *et al.* The Nippon Foundation—GEBCO seabed 2030 project: The quest to see  
418 the world’s oceans completely mapped by 2030. *Geosciences* **8**, 63 (2018).
- 419 48 Breithoff, E. & Harrison, R. From ark to bank: extinction, proxies and biocapitals in ex-  
420 situ biodiversity conservation practices. *International Journal of Heritage Studies* **26**, 37-  
421 55 (2020).
- 422 49 Berleant, J., Banal, J. L., Schardl, T. B., Leiserson, C. E. & Bathe, M. Beyond Big Data:  
423 Transformative Capabilities of Archival DNA Storage and Retrieval. (2020).
- 424 50 Baum, E. B. Building an associative memory vastly larger than the brain. *Science* **268**,  
425 583-585 (1995).
- 426 51 Song, X. & Reif, J. Nucleic acid databases and molecular-scale computing. *ACS Nano*  
427 **13**, 6256-6268, doi:10.1021/acsnano.9b02562 (2019).

428

429 **Acknowledgments.** We gratefully acknowledge fruitful discussions with Charles Leiserson and  
430 Tao B. Schardl on the scalability and generalizability of our barcoding approach. We thank Glenn  
431 Paradis, Michael Jennings, and Michele Griffin of the Flow Cytometry Core at the Koch Institute  
432 in MIT and Patricia Rogers of the Flow Cytometry Facility at the Broad Institute of Harvard and

433 MIT for assistance and fruitful discussions in developing the flow cytometry workflow. We also  
434 thank David Mankus of the Nanotechnology Materials Core Facility at the Koch Institute in MIT  
435 for assistance in the imaging of the particles using the scanning electron microscope and Alla  
436 Leshinsky of the Biopolymer and Proteomics Core at the Koch Institute at MIT for assistance in  
437 mass spectrometry characterization.

438 **Funding.** M.B., J.L.B., T.R.S., and J.B. gratefully acknowledge funding from the Office of Naval  
439 Research N00014-17-1-2609, N00014-16-1-2506, N00014-12-1-0621, and N00014-18-1-2290  
440 and the National Science Foundation CCF-1564025 and CBET-1729397. Additional funding to  
441 J.B. was provided through an NSF Graduate Research Fellowship (Grant # 1122374). P.C.B. was  
442 supported by a Career Award at the Scientific Interface from the Burroughs Wellcome Fund.  
443 C.M.A. was supported by NIH grant F32CA236425.

444 **Author contributions.** J.L.B., T.R.S., and M.B. designed the file labeling and selection scheme.  
445 J.L.B, T.R.S., and C.M.A. implemented the file selection scheme using FAS. J.B. and T.R.S.  
446 developed the encoding scheme and metadata tagging of the images to DNA. T.R.S. designed the  
447 plasmid for encoding imaging. H.H. and T.R.S. performed the cloning, transformation, and  
448 purification of the plasmids. J.L.B. synthesized and purified all the TAMRA and AFDye 647-  
449 labelled DNA oligonucleotides. J.L.B. characterized the particles. J.L.B. developed the synthetic  
450 route to attach DNA barcodes on the surface of the particles. J.L.B. performed the encapsulation,  
451 barcoding, sorting, reverse encapsulation of the particles after sorting, and desalting. T.R.S., H.H.,  
452 and M.R. performed the sequencing. J.B. performed computational validation of the orthogonality  
453 of barcode sequences and J.L.B. performed the experimental validation of the orthogonality of  
454 barcode and probe sequences. J.B. developed the computational workflow to analyze the  
455 sequencing data, including statistical analyses. M.B. conceived of the file system and supervised

456 the entire project. P.C.B. supervised the FAS selection and supervised the sequencing workflow.

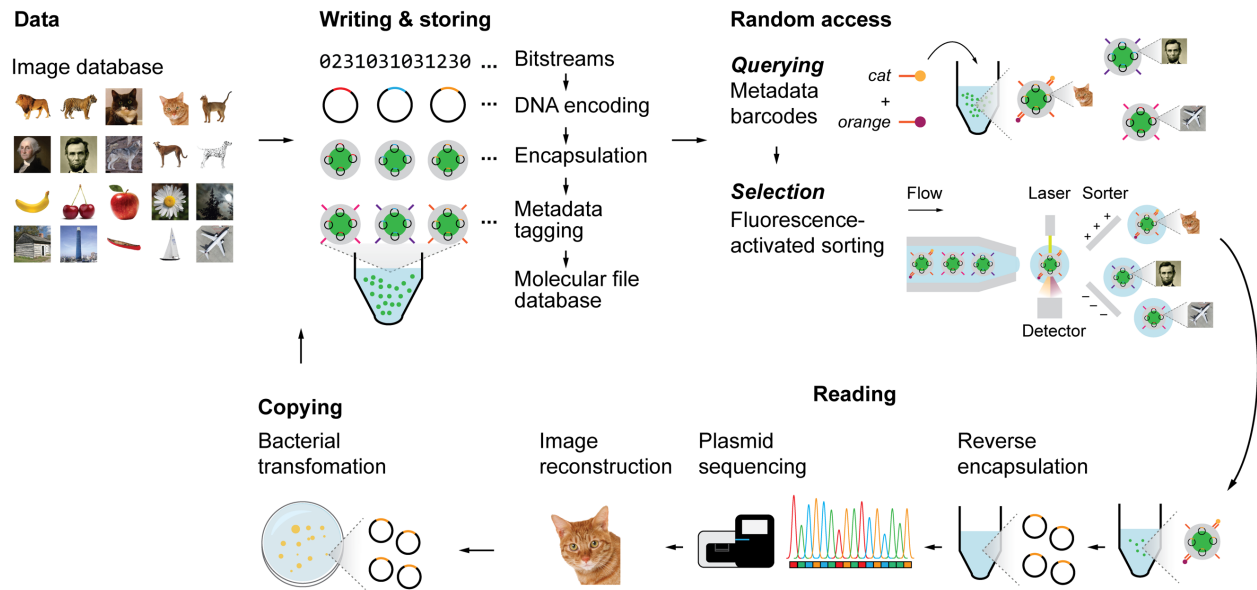
457 All authors analyzed the data and equally contributed to the writing of the manuscript.

458 **Competing interests.** T.R.S., J.L.B., J.B. & M.B. have filed provisional patents (17/029,948 and  
459 16/012,583) related to this work.

460 **Materials and correspondence.** Gene sequences and plasmid maps are available from AddGene  
461 (<https://www.addgene.org/depositing/77231/>). Software for sequence encoding and decoding is  
462 publicly available on GitHub (<https://github.com/lcbb/DNA-Memory-Blocks/>). All the data files  
463 used to generate the plots in this manuscript are available from M.B. upon request.

464 **Online content.** Any methods, additional references, and supplementary information are available  
465 at <https://doi.org/10.10XX/XXXXXX>.

466



467

468

469 **Figure 1 | Write-access-read cycle for a content-addressable molecular file system.** Colored

470 images were converted into  $26 \times 26$ -pixel, black-and-white icon bitmaps. The black-and-white

471 images were then converted into DNA sequences using ternary encoding scheme<sup>8</sup>. The DNA

472 sequences that encoded the images (file sequences) were inserted into a pUC19 plasmid vector

473 and encapsulated into silica particles using sol-gel chemistry. Silica capsules were then addressed

474 with content barcodes using orthogonal 25-mer single-stranded DNA strands, which were the final

475 forms of the files. Files were pooled to form the molecular file database. To query a file or several

476 files, fluorescently-labelled 15-mer ssDNA probes that are complementary to file barcodes were

477 added to the data pool. Particles were then sorted with fluorescence-activated sorting (FAS) using

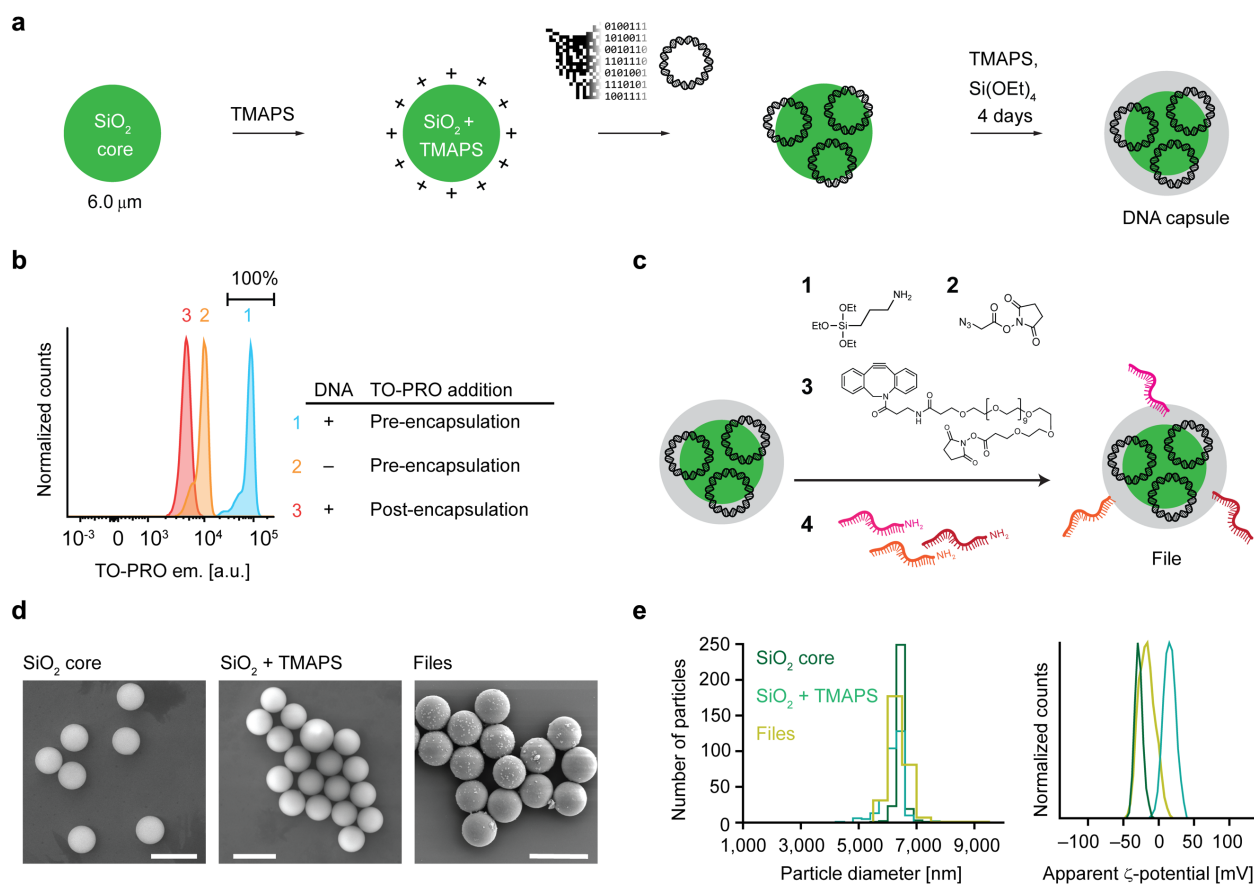
478 two to four fluorescence channels simultaneously. Addition of a chemical etching reagent into the

479 sorted populations released the encapsulated DNA plasmid. Sequences for the encoded images

480 were validated using Sanger sequencing or Illumina MiniSeq. Because plasmids were used to

481 encode information, re-transformation of the released plasmids into bacteria to replenish the

482 molecular file database thereby closed the write-access-read cycle.



483

484

485 **Figure 2 | Encapsulation of DNA plasmids into silica and surface barcoding. a**, Workflow of

486 silica encapsulation<sup>29</sup>. **b**, Raw fluorescence data from FAS experiments to detect DNA staining of

487 TO-PRO during or after encapsulation. **c**, Functionalization of encapsulated DNA particles. **d**,

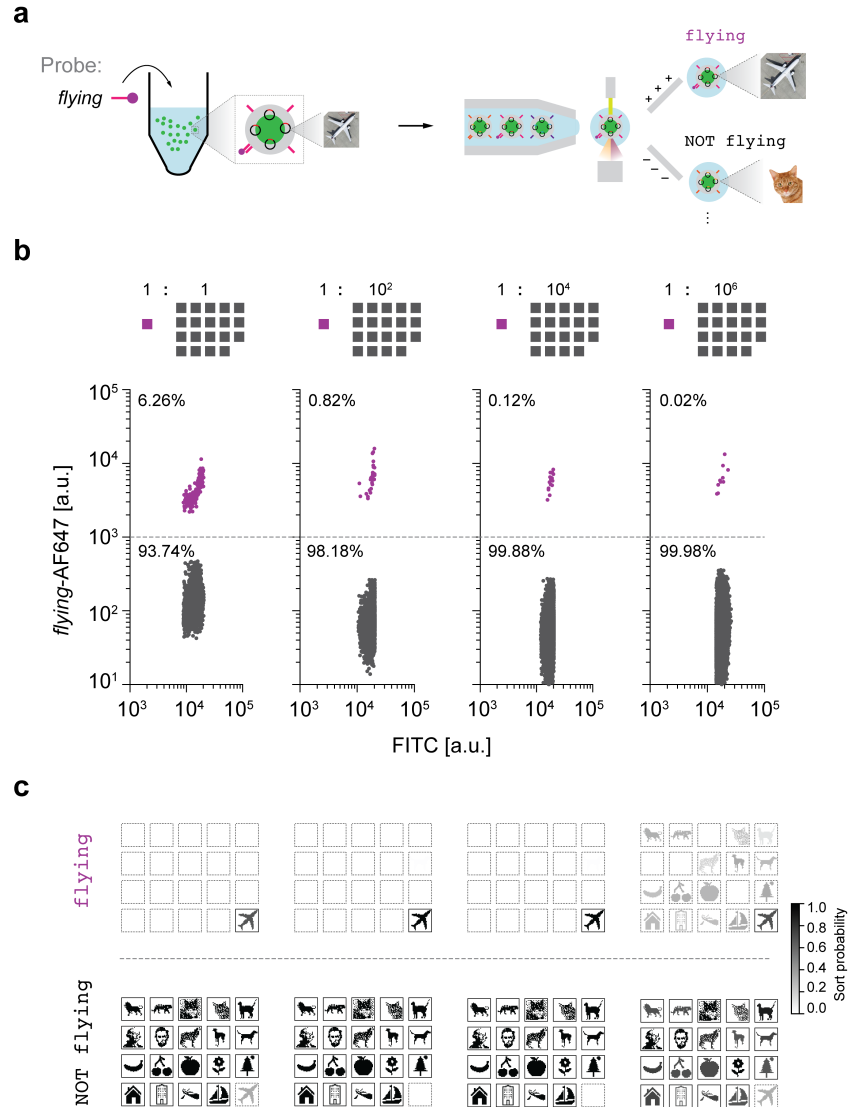
488 Scanning electron microscopy images of bare silica particles, silica particles functionalized with

489 TMAPS, and the file. **e**, Distribution of particle sizes determined from microscopy data (left) and

490 zeta potential analyses of silica particles and files.

491





492

493

494 **Figure 3 | Single-barcode sorting.** **a**, Schematic diagram of file sorting using FAS. **b**, Sorting of

495 *Airplane* from varying relative abundance of the other nineteen files as background. Percentages

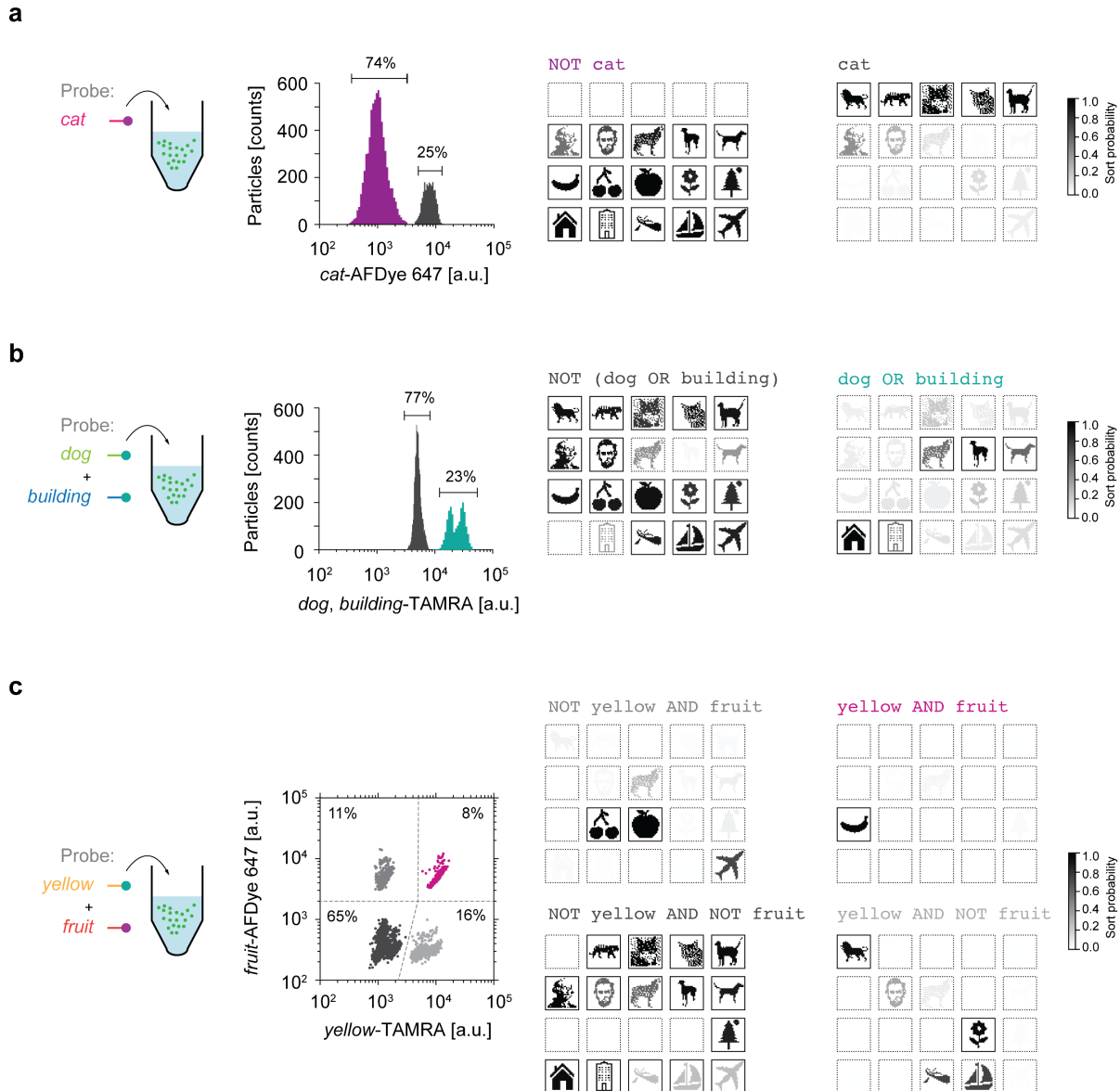
496 represent the numbers of particles that were sorted in the gate. Colored traces in each of the sorting

497 plots indicate the target population. **c**, Sequencing validation using Illumina MiniSeq. Sort

498 probability is the probability that a file is sorted into one gated population over the other gated

499 populations. Boxes with solid outlines indicate files that should be sorted into the specified gate.

500 Other files have dashed outlines.



501

502

503 **Figure 4 | Fundamental Boolean logic gates. a**, NOT cat selection. Raw fluorescence trace

504 from the FAS system (left) plotted on a 1D sorting plot showing the percent of particles that were

505 sorted in each gate. Sequencing using Illumina MiniSeq tested selection specificity (right). **b**, dog

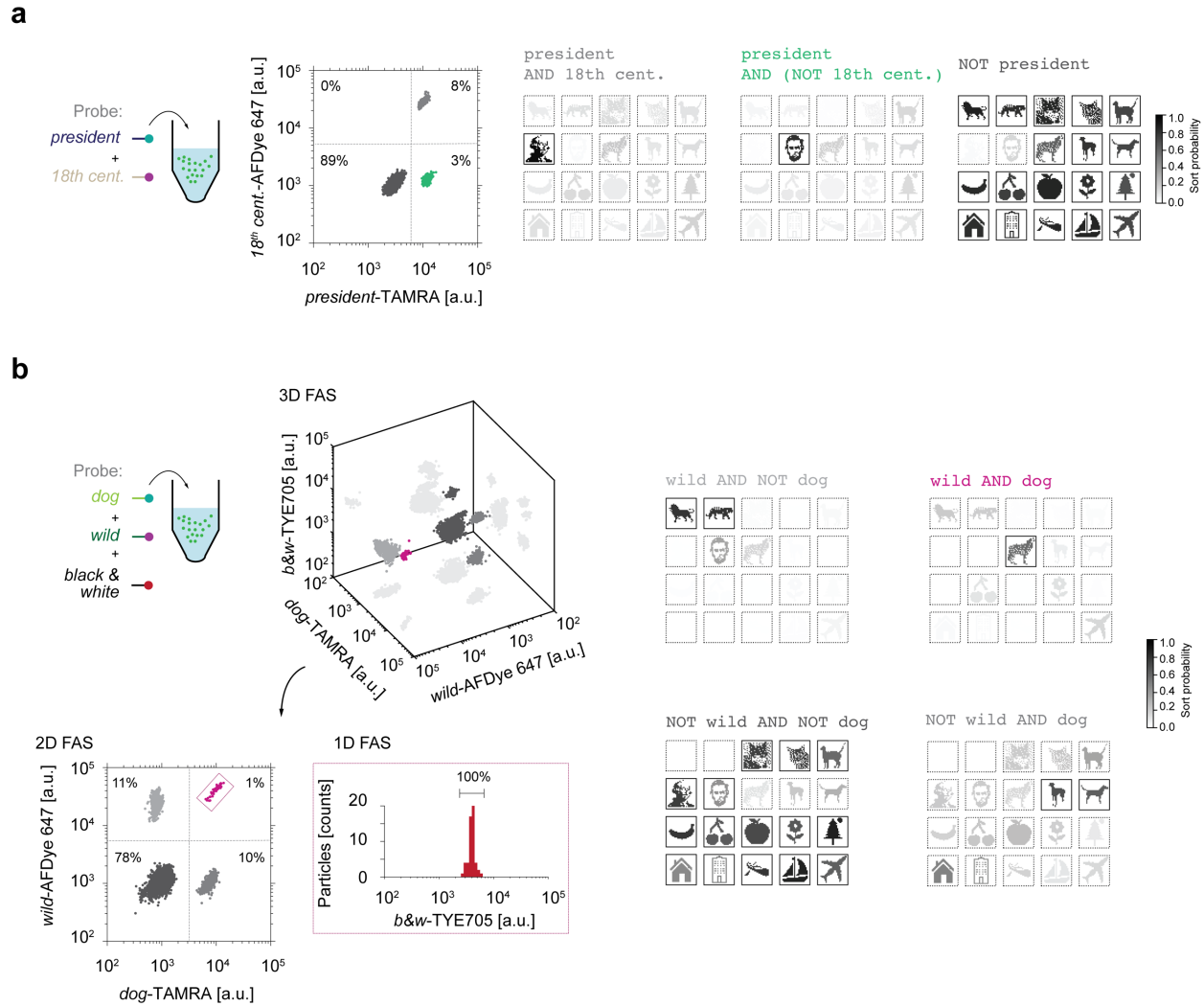
506 OR building selection. Raw fluorescence trace from the FAS system (left) plotted on a 1D

507 sorting plot showing the percent of particles that were sorted in each gate. Sequencing using

508 Illumina MiniSeq evaluated sorting using the OR gate (right). **c**, A 2D sorting plot to perform a

509 yellow AND fruit gate. Percentages in each quadrant show the percentages of particles that  
510 were sorted in each gate. Colored traces in all of the sorting plots indicate the target populations.  
511 Sort probability is the probability that a file is sorted into one gated population versus the other  
512 gated populations. Boxes with solid outlines indicate files that were intended to sort into the  
513 specified gate. Other files have dashed outlines.

514



515

516

517 **Figure 5 | Arbitrary logic searching. a, president AND (NOT 18<sup>th</sup> century) sorting.**

518 A 2D sorting plot (middle) was used to sort *Lincoln* by selecting a population that has high

519 TAMRA fluorescence but low AFDye 647 fluorescence. Sequencing using MiniSeq offered

520 quantitative evaluation of the sorted populations. **b, Multiple fluorescence channels were projected**

521 into a 3D FAS plot (left and top). There are three possible 2D plots that can be used for sorting.

522 To select the *Wolf* image using the query wild AND dog, a 2D plot of *wild* versus *dog* was first

523 selected and then populations selected using quadrant gates (left and bottom). One of the quadrants

524 were then selected where the *Wolf* image should belong based on the wild AND dog query in

525 order to test whether only a single population was present in the TYE705 fluorescence channel.  
526 Sequencing quantified the sorted populations (right) using Illumina MiniSeq. Sort probability is  
527 the probability that a file was sorted into one gated population over the other gated populations.  
528 Boxes with solid outlines indicate files that would ideally be sorted into the specified gate. Other  
529 files have dashed outlines.  
530