

## Summary

MIT startup from the laboratory of Prof. Mark Bathe in the Department of Biological Engineering at MIT seeking a part or full-time business lead for the next phase of customer discovery and business launch. Our technology has the potential to transform the way DNA and RNA are stored and managed, unlocking waves of innovation in the high-growth biotechnology industry.

## What is Cache DNA?

Our mission is to provide a low-cost platform to store nucleic acids that are mission-critical to a number of areas such as preserving genetic material from critically-endangered species, understanding biological threats, and upgrading DNA-based molecular file systems. By putting samples in barcoded, micron-scale silica capsules, The technology of Cache DNA allows biobanks to replace large-footprint freezer systems with portable, room temperature vials that fit in the palm of your hand. Cache's approach to random access allows for the same experience of querying and retrieval that we expect from digital databases. Our vision is to transform the way biological samples are stored, providing as critical infrastructure for the biotechnology revolution.

## Who are we looking for?

We are seeking a passionate, creative, and empathetic entrepreneur to help our early-stage venture realize the potential of its promising science and technology. Your immediate goal will be to lead our ongoing conversations with leading contract research organizations, biotech companies, and government agencies as we forge groundbreaking partnerships. You will also help draft our business plan using market discovery and research on product-market fit. We are committed to diversity in our founding team. Please contact us if you are interested in learning more about this opportunity.

## What is the IP portfolio of Cache DNA?

Cache DNA has two pending patents that are central to its business development. These patents are intended to cover the end-to-end solution for DNA data storage and retrieval, and broadly apply to storage and retrieval of any nucleic acid samples.

- [In one patent](#), we have developed a zettabyte-scale archival DNA data storage platform that provides random-access memory (RAM) retrieval of arbitrary megabyte-to-terabyte datasets from this zettabyte archival data pool. Our approach importantly offers the ability of selection of subsets of these data in a manner that mimics searching for arbitrary text and image queries on Google from a desktop computer or handheld device. Beyond archival data storage and retrieval, our approach offers the opportunity to perform large-scale computational operations by leveraging massively parallel interactions of molecules in solution, approaching  $10^{23}$  (or Avogadro's number) processes at any given time, with minimal energy use.
- [A second patent covers](#) the ability to write DNA using nano-to-femtoliter droplets to reduce waste materials and actuated using microfluidics. This technology allows the parallel writing of 1,000–100,000 unique DNA sequences, chemically or enzymatically, thus offering a competitive approach to existing microarray-based chemical synthesis. These patents offer a competitive technological advantage for Cache to build a successful business model for DNA storage and retrieval.

## Who are we?

### Technical Lead / Co-Founder

[James L. Banal, Ph.D.](#), Postdoctoral Scholar, Dept. of Biological Engineering, MIT ([jbanal@mit.edu](mailto:jbanal@mit.edu))

### Scientific Co-Founder

[Mark Bathe, Ph.D.](#), Professor, Dept. of Biological Engineering, MIT ([mark.bathe@mit.edu](mailto:mark.bathe@mit.edu))<sup>‡</sup>

### Business Advisors

[Howard Bornstein, M.B.A.](#), Principal, InnoSpark Ventures

[Fernando Rodriguez-Villa](#), Director, International Strategy, Indigo Ag; GM/Founding Team, TellusLabs

### Scientific Advisory Board

[Paul C. Blainey, Ph.D.](#), Associate Professor, Dept. of Biological Engineering, MIT

[George M. Church, Ph.D.](#), Professor, Dept. of Genetics, Harvard Medical School

[Jeremiah A. Johnson, Ph.D.](#), Professor, Dept. of Chemistry, MIT

<sup>‡</sup>Point of contact

# Arbitrary Boolean logical search operations on massive molecular file systems

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)

<sup>†</sup>These authors contributed equally to this work.

**DNA is an ultra-high-density storage medium that could meet exponentially growing worldwide data storage demand. However, accessing arbitrary data subsets within exabyte-scale DNA data pools is limited by the finite addressing space for individual DNA-based blocks of data. Here, we form files by encapsulating data-encoding DNA within silica capsules that are surface-labeled with multiple unique barcodes. Barcoding is performed with single-stranded DNA representing file metadata that enables Boolean logic selection on the entire pool of data. We demonstrate encapsulation and Boolean selection of sub-pools of image files using fluorescence-activated sorting, with selection sensitivity of 1 in  $10^6$  files per channel. Our strategy in principle enables retrieval of targeted data subsets from exabyte- and larger-scale data pools, thereby offering a random access file system for massive molecular data sets.**

DNA is the polymer used for storage and transmission of genetic information in biology. In principle, DNA can also be used as a medium for the storage of arbitrary digital information at densities far exceeding existing commercial data storage technologies and at scales well beyond the capacity of current data centers <sup>1</sup>. Ongoing advances in nucleic acid synthesis and sequencing technologies also continue to reduce dramatically the cost of writing and reading DNA, thereby rendering DNA-based digital information storage potentially viable economically in the near future <sup>2-5</sup>. As demonstrations of its viability as a general information storage medium, to date books, images, computer programs, audio clips, works of art, and Shakespeare's sonnets have all been stored in DNA using a variety of encoding schemes <sup>6-12</sup>. In each case, digital information was converted to DNA sequences and typically fragmented into 100–200 nucleotide (nt) blocks of data

for ease of chemical synthesis and sequencing. Sequence fragments were then assembled to reconstruct the original, encoded information.

While significant research effort has focused on improving DNA synthesis and encoding schemes, an additional, crucial aspect of digital data storage and retrieval is the ability to access specific subsets of a data pool on demand, which is conventionally achieved using polymerase chain reaction (PCR)<sup>8,10,12</sup>. PCR-based strategies take advantage of the ease of replication of DNA to extract specific DNA sequences from a DNA data pool using custom-designed forward and reverse primers that are complementary to the flanking sequences of interest. Nested addressing barcodes<sup>13-15</sup> can also be used to uniquely identify files using multiple barcodes. For an exabyte-scale data pool, each file requires at least four barcodes, or up to one hundred nucleotides in total barcode sequence length, thereby nearly eliminating the number of nucleotides that can be used for data encoding. Further, orthogonality of barcodes to other barcodes and file sequences present in the data pool is essential for reliable data access. To overcome these limitations, previous approaches have used spatial segregation of data into distinct pools<sup>16</sup>. While PCR is typically known for its ease of amplifying specific DNA sequences, errors in priming via strand crosstalk can lead to information loss. In addition, selective amplification of a specific file using PCR requires access to the entire data pool for each query, which is also destructive to the sample queried. Finally, PCR-based approaches do not allow for physical deletion of specific files from a data pool, other than implementing an address overwrite<sup>10</sup>.

As an alternative to PCR-based data access, inspired by genomic segmentation within biological cells, here we physically encapsulate and thereby isolate DNA-based molecular data within discrete silica capsules, which we subsequently label to enable random access of the data pool via hybridization and subsequent optical selection. Each unit of information encoded in DNA

we term a *file*, which includes both the DNA encoding the main data as well as any additional components used for addressing, storage, and retrieval. Each file contains a *file sequence*, consisting of the DNA encoding the main data, and *addressing barcodes*, or simply *barcodes*, which are additional short DNA sequences used to identify the file in solution using hybridization. We refer to a collection of files as a *data pool* or *database*, and the procedures for storing, retrieving, and reading out files is termed a *file system* (see **Supplementary Section S0** for a full list of terms).

As a proof-of-principle of our file system, we encapsulated 3,000-nt plasmids encoding 85-byte images, the files, within monodisperse, 6- $\mu$ m spherical silica particles that were chemically surface-labeled using up to three 25-mer single-stranded DNA (ssDNA) oligonucleotides, the barcodes, chosen from a library of 240,000 orthogonal primers, allowing identification of up to  $\sim 10^{15}$  possible distinct files using only three unique barcodes per file <sup>17</sup> (**Fig. 1**). Twenty icon-resolution images were chosen in the data pool to represent diverse subject matter including animals, plants, transportation, and buildings, and labeled with DNA barcodes that represent the categories to which each image belongs (**Supplementary Fig. 1**). Fluorescence-activated sorting (FAS) was used to select target subsets of the complete data pool by first annealing fluorescent oligonucleotide probes that are complementary to the barcodes, in order to address the DNA database <sup>18</sup>. Retrieval of specific, individual files and collections of files described by Boolean AND, OR, and NOT logic was achieved using combinations of distinct barcodes to query the data pool. Because physical encapsulation separates file sequences from barcodes used to describe the encapsulated information, our file system offers highly specific, robust data retrieval operations; the ability to delete specific subsets of data; in addition to long-term environmental protection of encoded file sequences via silica encapsulation <sup>9,19,20</sup>. While we apply our proposed file system to

a prototypical kilobyte-scale image database here, our approach is fully scalable to massive molecular data pools at the exabyte- and larger-scales, as well as alternative encapsulation strategies<sup>21,22</sup>, barcode implementations<sup>23-27</sup>, and physical or other sorting strategies using biochemical affinity, optical, or other labeling approaches<sup>28-30</sup>.

## File Synthesis

Digital information in the form of 20 icon-resolution images was stored in a data pool, with each image encoded into DNA and synthesized on a plasmid. We selected images of broad diversity, representative of distinct and shared subject categories, which included several domestic and wild cats and dogs, US presidents, and several human-made objects such as an airplane, boats, and buildings (**Fig. 1** and **Supplementary Fig. 1**). To implement this image database, the images were substituted with black-and-white,  $26 \times 26$ -pixel images to minimize synthesis costs, compressed using run-length encoding, and converted to DNA (**Supplementary Fig. 1, 2**). Following synthesis, bacterial amplification, and sequencing validation (**Supplementary Fig. 3**), each plasmid DNA was separately encapsulated into silica particles containing a fluorescein dye core and a positively charged surface<sup>19,20</sup>. Because the negatively charged phosphate groups of the DNA interact with positively charged silica particles, plasmid DNA condensed on the silica surface, after which N-[3-(trimethoxysilyl)propyl]-N,N,N-trimethylammonium chloride (TMAPS) was co-condensed with tetraethoxysilane to form an encapsulation shell after four days of incubation at room-temperature<sup>9,20</sup> (**Fig. 2a**) to form discrete silica capsules containing the file sequence that encodes for the image file. Quantitative PCR (qPCR) of the reaction supernatant after encapsulation (**Supplementary Fig. 4**) showed full encapsulation of plasmids without residual DNA in solution. To investigate the fraction of capsules that contained plasmid DNA, we

compared the fluorescence intensity of the intercalating dye TO-PRO when added pre- versus post-encapsulation (**Supplementary Fig. 2**). All capsules synthesized in the presence of both DNA and TO-PRO showed a distinct fluorescence signal, consistent with the presence of plasmid DNA in the majority of capsules, compared with a silica particle negative control that contained no DNA. In order to test whether plasmid DNA was fully encapsulated versus partially exposed at the surface of capsules, capsules were also stained separately with TO-PRO post-encapsulation (**Fig. 2b**). Using qPCR, we estimated  $10^6$  plasmids per capsule assuming quantitative recovery of DNA post-encapsulation (**Supplementary Fig. 5**).

Next, we chemically attached unique content addresses on the surfaces of silica capsules using orthogonal 25-mer ssDNA barcodes (**Supplementary Fig. 6**) describing selected features of the underlying image. For example, the image of an orange tabby house cat (**Supplementary Fig. 1**) was described with *cat*, *orange*, and *domestic*, whereas the image of a tiger was described with *cat*, *orange*, and *wild* (**Supplementary Fig. 1** and **Supplementary Table 2**). To attach the barcodes, we activated the surface of the silica capsules through a series of chemical steps. Condensation of  $\gamma$ -aminopropyltriethoxysilane with the hydroxy-terminated surface of the encapsulated plasmid DNA provided a primary amine chemical handle that supported further conjugation reactions (**Fig. 2c**). We modified the amino-modified surface of the silica capsules with  $\beta$ -azidoacetic acid N-hydroxysuccinimide (NHS) ester followed by an oligo(ethylene glycol) that contained two chemically orthogonal functional groups: the dibenzocyclooctyne functional group reacted with the surface-attached azide through strain-promoted azide-alkyne cycloaddition while the NHS ester functional group was available for subsequent conjugation with a primary amine. Each of the associated barcodes contained a 5'-amino modification that could react with the NHS-ester groups on the surface of the silica capsules, thereby producing the complete form

of our file. Notably, the sizes of bare, hydroxy-terminated silica particles representing capsules without barcodes were comparable with complete files consisting of capsules with barcodes attached, confirmed using scanning electron microscopy (**Fig. 2d** and **2e**, left). These results were anticipated given that the encapsulation thickness was only on the order of 10 nm<sup>20</sup> and that additional steps to attach functional groups minimally increases the capsule diameter. We also observed systematic shifts in the surface charge of the silica particles as different functional groups were introduced onto their surfaces (**Fig. 2e**). Using hybridization assays with fluorescently-labelled probes<sup>31-33</sup>, we estimated the number of barcodes available for hybridization on our files to be on the order of 10<sup>8</sup> (**Supplementary Fig. 7**). Following synthesis, files were pooled and stored together for subsequent retrieval. Illumina MiSeq was used to read each file sequence and reconstruct the encoded image following selection and de-encapsulation, in order to validate the complete process of image file encoding, encapsulation, barcoding, selection, de-encapsulation, sequencing, and image file reconstruction (**Supplementary Figs. 9, 10**).

## File Selection

Following file synthesis and pooling, we used FAS to select specific targeted file subsets from the entire data pool. All files contained a fluorescent dye, fluorescein, in their core as a marker to distinguish files from other particulates such as spurious silica particles that nucleated in the absence of a core or insoluble salts that may have formed during the sorting process. Each detected fluorescein event was therefore interpreted to indicate the presence of an individual file at sufficiently low concentrations queried using FAS (**Supplementary Fig. 11**). For any query applied to the entire image database, a fluorescently-labelled ssDNA probe hybridized to its complementary barcode displayed externally on the surface of the silica capsule (**Fig. 3a**).



We subjected the entire data pool to a series of experiments to test selection sensitivity of target subsets using distinct queries. First, we evaluated single-barcode selection of an individual file, specifically *Airplane*, out of a pool of varying concentrations of the nineteen other files as background (**Fig. 3b**). To select the *Airplane* file, we hybridized an AFDye 647-labelled ssDNA probe that is complementary to the barcode *flying*, which is unique to *Airplane*. We were able to detect and select the desired *Airplane* file through FAS even at a relative abundance of  $10^{-6}$  compared with each other file (**Fig. 3c**). Comparison of the retrieved sequences between the flying gate and the NOT flying gate after chemical release of the file sequences from silica encapsulation revealed that 60–95% of the *Airplane* files were sorted into the flying gate (**Supplementary Figs. 18–21**). Note that any sort probability above 50% indicates enrichment of *Airplane* within the correct population subset (flying) relative to the incorrect subset (NOT flying), while a sort probability of 100% would indicate ideal performance.

## Boolean Search

Aside from selecting single files, Boolean logic can be used to select a specific subset of the data pool. We demonstrated AND, OR, and NOT logical operations by first adding to the data pool fluorescently-labelled ssDNA probes that were complementary to the barcodes (**Fig. 4**, left). This hybridization reaction was used to distinguish one or several files in the data pool, which were then sorted using FAS. We used two to four fluorescence channels simultaneously to create the FAS gates that executed the target Boolean logic queries (**Fig. 4**, middle). To demonstrate a NOT query, we added to the data pool an AFDye 647-labelled ssDNA probe that hybridized to files that contained the *cat* barcode. Files that did not show AFDye 647 signal were sorted into the NOT *cat* subset (**Fig. 4a**). An example of an OR gate was applied to the data pool by simultaneously adding

*dog* and *building* probes that both had the TAMRA label (**Fig. 4b**). All files that showed TAMRA signal were sorted into the dog OR building subset by the FAS. Finally, an example of an AND gate was achieved by adding *fruit* and *yellow* probes that were labelled with AFDye 647 and TAMRA, respectively. Files showing signal for both AFDye 647 and TAMRA were sorted into the fruit AND yellow subset in the FAS (**Fig. 4c**). For each example query, we validated our sorting experiments by releasing the file sequence from silica encapsulation and sequencing the released DNA with Illumina MiniSeq (**Fig. 4**, right). Sort probabilities of each file for each search query are shown in **Supplementary Figs. S22–S24**.

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

To demonstrate the possibility of performing Boolean search using more than three fluorescence channels for sorting, we selected the *Wolf* file from the data pool using the query dog AND wild, and used the *black & white* probe to validate the selected file (**Fig. 5b**, left). Because conventional FAS software is only capable of sorting using 1D and 2D gates, we first selected one out of the three possible 2D plots (**Fig. 5b**, left and bottom): *dog*-TAMRA against *wild*-AFDye 647. We examined the *black & white*-TYE705 channel on members of the dog AND wild subset (**Fig. 5b**, left and bottom). Release of the encapsulated file sequence and subsequent sequencing of each gated population from the *dog* versus *wild* 2D plot validated sorting (**Fig. 5b**, right; **Supplementary Fig. 26**).

The use of plasmids as a substrate for encoding information offered a convenient workflow for restoring files into the data pool after retrieval. In cases where single images were sorted (**Figs. 4c, 5a, b**), we were able to transform competent bacteria from each search query that resulted in a single file (**Supplementary Fig. 27**). Amplified material was pure and ready for re-encapsulation into silica particles, which could be re-introduced directly back into the data pool. Importantly, our molecular file system and file selection process thereby represents a complete write-access-read cycle that can in principle be applied to exabyte and larger-scale datasets. While sort probabilities were typically below the perfect 100% targeted for a specific file or file subset query, future work would be required to better characterize sources of error that may be due to sample contamination, FAS error, or imperfect orthogonality of barcode sequences employed (**Supplementary Fig. 6**)<sup>17</sup>.

## Outlook

We present a non-destructive molecular file system that is capable of both specific file selection and Boolean logic search operations for random access of single files or file subsets in a data pool.

Our implementation easily scales by increasing the numbers of barcodes per file and query fluorophores used for file selection, which can thereby address files in a larger-scale database for random access and computation. For example, labeling each file using four distinct barcodes instead of only the three used here renders it possible to label  $\binom{2.4 \times 10^5}{4} \approx 10^{20}$  files uniquely using the existing pool of  $\sim 10^5$  orthogonal barcodes<sup>17</sup>. Assuming an FAS system is capable of sorting a single file from  $10^6$  others using each fluorescent channel alone, as demonstrated in this work using a commercial FAS, one may theoretically sort a single file from  $10^{24}$  others using a conventional four-channel FAS system. This file system would then in principle offer sufficient sensitivity and specificity to select a single file from an exabyte or even yottabyte data pool. However, the time needed to perform FAS scales linearly with the size of the data pool, which may be prohibitively long even for exabyte-scale data pools using only 10–100 bytes per file. For example, 12 minutes was required to select at least one hundred copies of the *Airplane* file in a data pool in which this file has a relative abundance of  $10^{-6}$  compared with other files (**Fig. 3**). This is in contrast to selecting one hundred copies of the *Airplane* file in a data pool that contained equivalent numbers of nineteen other files, which required only  $\sim 30$  seconds. Thus, in order to search through an entire exabyte-scale data pool within 24 hours, each file should consist of approximately 100 gigabytes, assuming a typical commercial FAS device that searches at 10,000 files per second. In order to reduce file selection time, future implementations of our molecular file system should therefore leverage parallel microfluidics-based optical sorting procedures and brighter fluorescence probes to increase selection throughput and sensitivity, and thereby reduce the pool search time. Alternatively, direct magnetic pulldown of files labelled with biochemical or affinity tags may be employed<sup>30</sup>.

Aside from speed and specificity of data access, data density is also of importance to DNA data storage. Notably, both file size and data density can be tuned independently in our file system by changing the information content of loaded DNA and the size of the silica particles employed for encapsulation. While we used 6- $\mu\text{m}$  silica core particles here in order to maximize fluorescence signal-to-noise ratios for a commercial FAS instrument, this also limited volumetric density of our DNA file system<sup>3</sup>. Specifically, using this approach an exabyte-scale data pool consisting of a 100-byte file per particle would require approximately  $10^{16}$  files and  $1\text{ m}^3$  total dry volume, or  $10^{18}$  bytes per  $\text{m}^3$ . In comparison, PCR-based random access has a theoretical volumetric density limit of  $10^{24}$  bytes per  $\text{m}^3$ <sup>3</sup>, although additional methods are required to prevent crosstalk between file sequences and barcodes. To further increase the data density of our file system, future implementations may benefit from using nanoparticles  $\sim 100\text{--}200\text{ nm}$  in diameter to encode files<sup>9,19,20</sup> and higher sensitivity FAS systems<sup>34,35</sup> or direct biochemical, magnetic, or other pulldown for file and data subset selection from the data pool.

Beyond increasing file selection speed and data density, utilization of spectrally distinct fluorescent probes and discrete labeling intensities<sup>26</sup> would allow for far more complex and efficient logical operations than demonstrated here<sup>36</sup>. Physical particle parameters including forward and side-scatter could additionally be used to perform multi-dimensional sorting of particles with different scattering cross-sections, with or without additional fluorescence channels<sup>37</sup>. Repeated cycles of file selection in series could also further increase selection fidelity. While our technical approach differs significantly from approaches that rely on selective amplification for block selection<sup>8,12,16</sup>, in which amplifications may reduce fidelity of file selection, PCR-based random access approaches will typically have faster read-write times because they forgo encapsulation and de-encapsulation steps required by our approach<sup>9,19,20</sup>, which is therefore ideally

suited to long-term, archival data storage and retrieval with periodic file and barcode renewal. Aside from DNA data storage, population enrichment on our prototypical database of 20 unique files encoded in DNA plasmids with silica encapsulation and retrieval demonstrated using barcodes labels may alternatively be applied directly to biological DNA and other nanoscale sample management, such as genomic samples in biobanking or protein-encoding databases<sup>38</sup>. In either case, subsets of data or genomic sample pools may be enriched using Boolean AND, OR, and NOT logic, which complements existing PCR-based approaches. These operations enrich the capabilities of performing computation and sorting on underlying molecular data pools, moving us closer to realizing an economically viable, functional, massive molecular file and operating system<sup>18,39,40</sup>.

## References

- 1 Ceze, L., Nivala, J. & Strauss, K. Molecular digital data storage using DNA. *Nature Reviews Genetics* **20**, 456-466, doi:10.1038/s41576-019-0125-3 (2019).
- 2 Kosuri, S. & Church, G. M. Large-scale de novo DNA synthesis: technologies and applications. *Nature Methods* **11**, 499-507, doi:10.1038/nmeth.2918 (2014).
- 3 Zhirnov, V., Zadegan, R. M., Sandhu, G. S., Church, G. M. & Hughes, W. L. Nucleic acid memory. *Nature Materials* **15**, 366, doi:10.1038/nmat4594 (2016).
- 4 Palluk, S. *et al.* De novo DNA synthesis using polymerase-nucleotide conjugates. *Nature Biotechnology* **36**, 645-650, doi:10.1038/nbt.4173 (2018).
- 5 Lee, H. H., Kalhor, R., Goela, N., Bolot, J. & Church, G. M. Terminator-free template-independent enzymatic DNA synthesis for digital information storage. *Nature Communications* **10**, 2383, doi:10.1038/s41467-019-10258-1 (2019).

274 6 Church, G. M., Gao, Y. & Kosuri, S. Next-generation digital information storage in  
275 DNA. *Science* **337**, 1628, doi:10.1126/science.1226355 (2012).

276 7 Goldman, N. *et al.* Towards practical, high-capacity, low-maintenance information  
277 storage in synthesized DNA. *Nature* **494**, 77-80, doi:10.1038/nature11875 (2013).

278 8 Yazdi, S. M. H. T., Yuan, Y., Ma, J., Zhao, H. & Milenkovic, O. A rewritable, random-  
279 access DNA-based storage system. *Scientific Reports* **5**, 14138, doi:10.1038/srep14138  
280 (2015).

281 9 Grass, R. N., Heckel, R., Puddu, M., Paunescu, D. & Stark, W. J. Robust chemical  
282 preservation of digital information on DNA in silica with error-correcting codes.  
283 *Angewandte Chemie International Edition* **54**, 2552-2555, doi:10.1002/anie.201411378  
284 (2015).

285 10 Yazdi, S. M. H. T., Gabrys, R. & Milenkovic, O. Portable and error-free DNA-based data  
286 storage. *Scientific Reports* **7**, 5011, doi:10.1038/s41598-017-05188-1 (2017).

287 11 Erlich, Y. & Zielinski, D. DNA Fountain enables a robust and efficient storage  
288 architecture. *Science* **355**, 950-954, doi:10.1126/science.aaj2038 (2017).

289 12 Organick, L. *et al.* Random access in large-scale DNA data storage. *Nature*  
290 *Biotechnology* **36**, 242–248, doi:10.1038/nbt.4079 (2018).

291 13 Kashiwamura, S., Yamamoto, M., Kameda, A., Shiba, T. & Ohuchi, A. in *DNA 2002:*  
292 *DNA Computing*. (eds M. Hagiya & A. Ohuchi) 112–123 (Lecture Notes in Computer  
293 Science, Vol. 2568, Springer, 2003).

294 14 Yamamoto, M., Kashiwamura, S., Ohuchi, A. & Furukawa, M. Large-scale DNA  
295 memory based on the nested PCR. *Natural Computing* **7**, 335-346 (2008).

296 15 Yamamoto, M., Kashiwamura, S. & Ohuchi, A. in *DNA 2007: DNA Computing*. (eds  
297 M.H. Garzon & H. Yan) 99–108 (Lecture Notes in Computer Science, Vol. 4848,  
298 Springer).

299 16 Newman, S. *et al.* High density DNA data storage library via dehydration with digital  
300 microfluidic retrieval. *Nature Communications* **10**, 1706 (2019).

301 17 Xu, Q., Schlabach, M. R., Hannon, G. J. & Elledge, S. J. Design of 240,000 orthogonal  
302 25mer DNA barcode probes. *Proceedings of the National Academy of Sciences* **106**,  
303 2289–2294, doi:10.1073/pnas.0812506106 (2009).

304 18 Reif, J. H. *et al.* in *DNA 2001: DNA Computing*. (eds N. Jonoska & N.C. Seeman) 231–  
305 247 (Lecture Notes in Computer Science, Vol. 2340, Springer, 2002).

306 19 Paunescu, D., Fuhrer, R. & Grass, R. N. Protection and deprotection of DNA--high-  
307 temperature stability of nucleic acid barcodes for polymer labeling. *Angewandte Chemie*  
308 *International Edition* **52**, 4269–4272, doi:10.1002/anie.201208135 (2013).

309 20 Paunescu, D., Puddu, M., Soellner, J. O. B., Stoessel, P. R. & Grass, R. N. Reversible  
310 DNA encapsulation in silica to produce ROS-resistant and heat-resistant synthetic DNA  
311 "fossils". *Nature Protocols* **8**, 2440, doi:10.1038/nprot.2013.154 (2013).

312 21 Alexakis, T. *et al.* Microencapsulation of DNA within alginate microspheres and  
313 crosslinked chitosan membranes for in vivo application. *Applied Biochemistry and*  
314 *Biotechnology* **50**, 93–106 (1995).

315 22 Borodina, T. *et al.* Controlled release of DNA from self-degrading microcapsules.  
316 *Macromolecular Rapid Communications* **28**, 1894–1899, doi:10.1002/marc.200700409  
317 (2007).



- 23 Braeckmans, K. *et al.* Encoding microcarriers by spatial selective photobleaching. *Nature Materials* **2**, 169-173, doi:10.1038/nmat828 (2003).
- 24 Wilson, R., Cossins, A. R. & Spiller, D. G. Encoded microcarriers for high-throughput multiplexed detection. *Angewandte Chemie International Edition* **45**, 6104-6117, doi:10.1002/anie.200600288 (2006).
- 25 Pregibon, D. C., Toner, M. & Doyle, P. S. Multifunctional encoded particles for high-throughput biomolecule analysis. *Science* **315**, 1393-1396, doi:10.1126/science.1134929 (2007).
- 26 Dagher, M., Kleinman, M., Ng, A. & Juncker, D. Ensemble multicolour FRET model enables barcoding at extreme FRET levels. *Nature Nanotechnology* **13**, 925-932, doi:10.1038/s41565-018-0205-0 (2018).
- 27 Martino, N. *et al.* Wavelength-encoded laser particles for massively multiplexed cell tagging. *Nature Photonics* **13**, 720-727, doi:10.1038/s41566-019-0489-0 (2019).
- 28 Lee, H., Kim, J., Kim, H., Kim, J. & Kwon, S. Colour-barcode magnetic microparticles for multiplexed bioassays. *Nature Materials* **9**, 745-749, doi:10.1038/nmat2815 (2010).
- 29 Stewart, K. *et al.* in *International Conference on DNA Computing and Molecular Programming*. 55-70 (Vol. Springer).
- 30 Tomek, K. J. *et al.* Driving the scalability of DNA-based information storage systems. *ACS Synthetic Biology* **8**, 1241-1248, doi:10.1021/acssynbio.9b00100 (2019).
- 31 Pillai, P. P., Reisewitz, S., Schroeder, H. & Niemeyer, C. M. Quantum-dot-encoded silica nanospheres for nucleic acid hybridization. *Small* **6**, 2130-2134, doi:10.1002/sml.201000949 (2010).

- 32 Leidner, A. *et al.* Biopebbles: DNA-functionalized core–shell silica nanospheres for cellular uptake and cell guidance studies. *Advanced Functional Materials* **28**, 1707572, doi:10.1002/adfm.201707572 (2018).
- 33 Sun, P. *et al.* Biopebble containers: DNA-directed surface assembly of mesoporous silica nanoparticles for cell studies. *Small* **15**, 1900083, doi:10.1002/smll.201900083 (2019).
- 34 van Gaal, E. V. B., Spierenburg, G., Hennink, W. E., Crommelin, D. J. A. & Mastrobattista, E. Flow cytometry for rapid size determination and sorting of nucleic acid containing nanoparticles in biological fluids. *Journal of Controlled Release* **141**, 328-338, doi:10.1016/j.jconrel.2009.09.009 (2010).
- 35 Lian, H., He, S., Chen, C. & Yan, X. Flow cytometric analysis of nanoscale biological particles and organelles. *Annual Review of Analytical Chemistry* **12**, 389-409, doi:10.1146/annurev-anchem-061318-115042 (2019).
- 36 Perfetto, S. P., Chattopadhyay, P. K. & Roederer, M. Seventeen-colour flow cytometry: unravelling the immune system. *Nature Reviews Immunology* **4**, 648-655, doi:10.1038/nri1416 (2004).
- 37 Mage, P. L. *et al.* Shape-based separation of synthetic microparticles. *Nature Materials* **18**, 82-89, doi:10.1038/s41563-018-0244-9 (2019).
- 38 Plesa, C., Sidore, A. M., Lubock, N. B., Zhang, D. & Kosuri, S. Multiplexed gene synthesis in emulsions for exploring protein functional landscapes. *Science* **359**, 343-347, doi:10.1126/science.aao5167 (2018).
- 39 Baum, E. B. Building an associative memory vastly larger than the brain. *Science* **268**, 583-585 (1995).

40 Song, X. & Reif, J. Nucleic acid databases and molecular-scale computing. *ACS Nano*  
13, 6256-6268, doi:10.1021/acsnano.9b02562 (2019).

**Acknowledgments.** We gratefully acknowledge fruitful discussions with Charles Leiserson and Tao B. Schardl on the scalability and generalizability of our barcoding approach. We thank Glenn Paradis, Michael Jennings, and Michele Griffin of the Flow Cytometry Core at the Koch Institute in MIT and Patricia Rogers of the Flow Cytometry Facility at the Broad Institute of Harvard and MIT for assistance and fruitful discussions in developing the flow cytometry workflow. We also thank David Mankus of the Nanotechnology Materials Core Facility at the Koch Institute in MIT for assistance in the imaging of the particles using the scanning electron microscope and Alla Leshinsky of the Biopolymer and Proteomics Core at the Koch Institute at MIT for assistance in mass spectrometry.

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

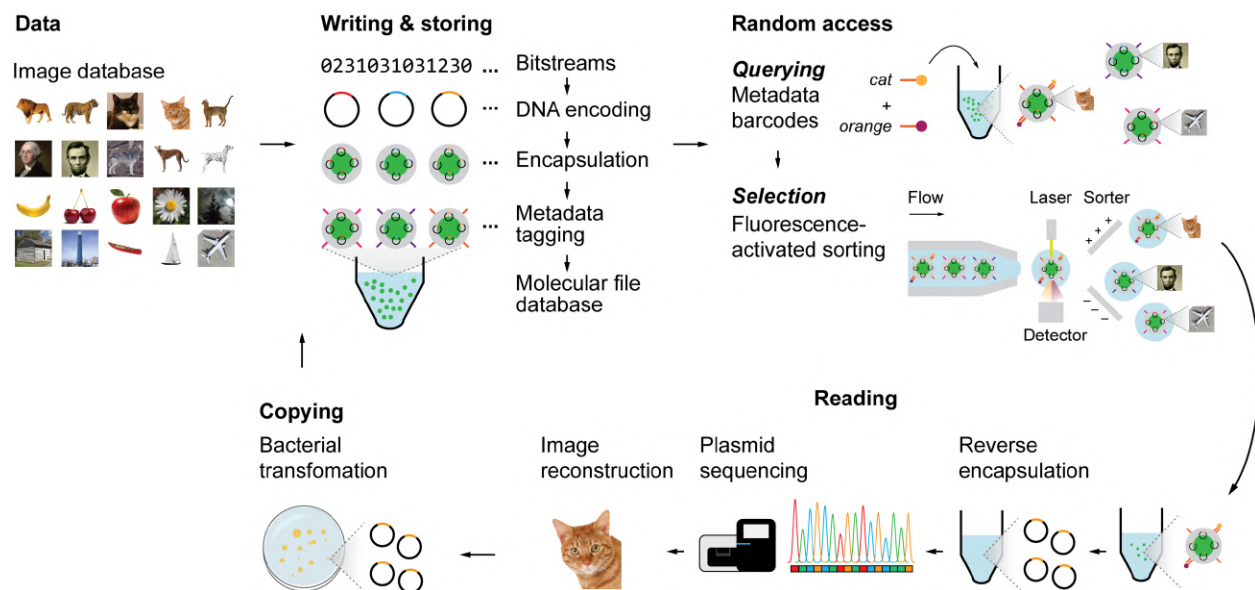
**Author contributions.** J.L.B., T.R.S., and M.B. designed the file labeling and selection scheme. J.L.B., T.R.S., and C.M.A. implemented the file selection scheme using FAS. J.B. and T.R.S. developed the encoding scheme and metadata tagging of the images to DNA. T.R.S. designed the plasmid for encoding imaging. H.H. and T.R.S. performed the cloning, transformation, and purification of the plasmids. J.L.B. synthesized and purified all the TAMRA and AFDye 647-

labelled DNA oligonucleotides. J.L.B. characterized the particles. J.L.B. developed the synthetic route to attach DNA barcodes on the surface of the particles. J.L.B. performed the encapsulation, barcoding, sorting, reverse encapsulation of the particles after sorting, and desalting. T.R.S., H.H., and M.R. performed the sequencing. J.B. performed computational validation of the orthogonality of barcode sequences. J.B. developed the computational workflow to analyze the sequencing data, including statistical analyses. M.B. conceived of the file system and supervised the entire project. P.C.B. supervised the FAS selection and supervised the sequencing workflow. All authors analyzed the data and equally contributed to the writing of the manuscript.

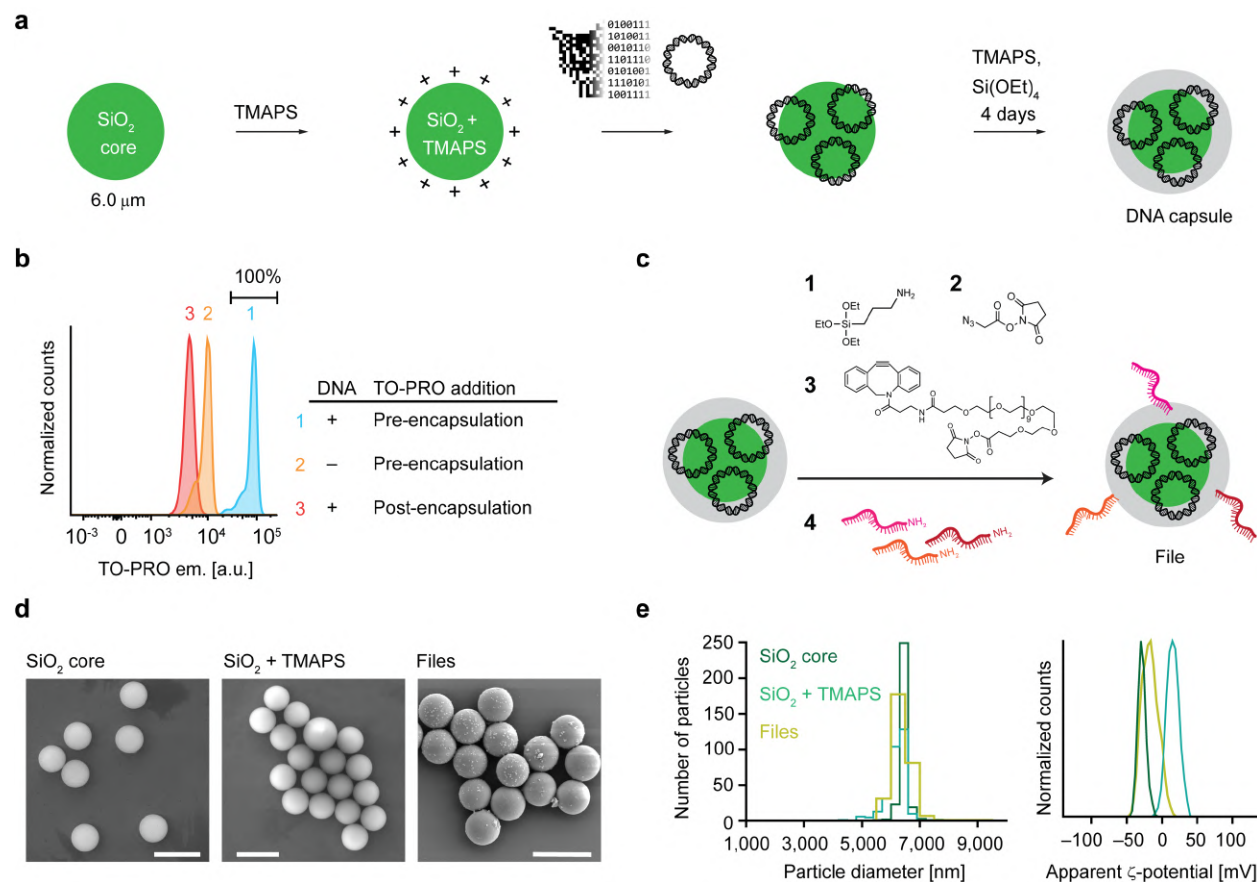
**Competing interests.** T.R.S., J.L.B., J.B. & M.B. have filed provisional patents (17/029,948 and 16/012,583) related to this work.

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

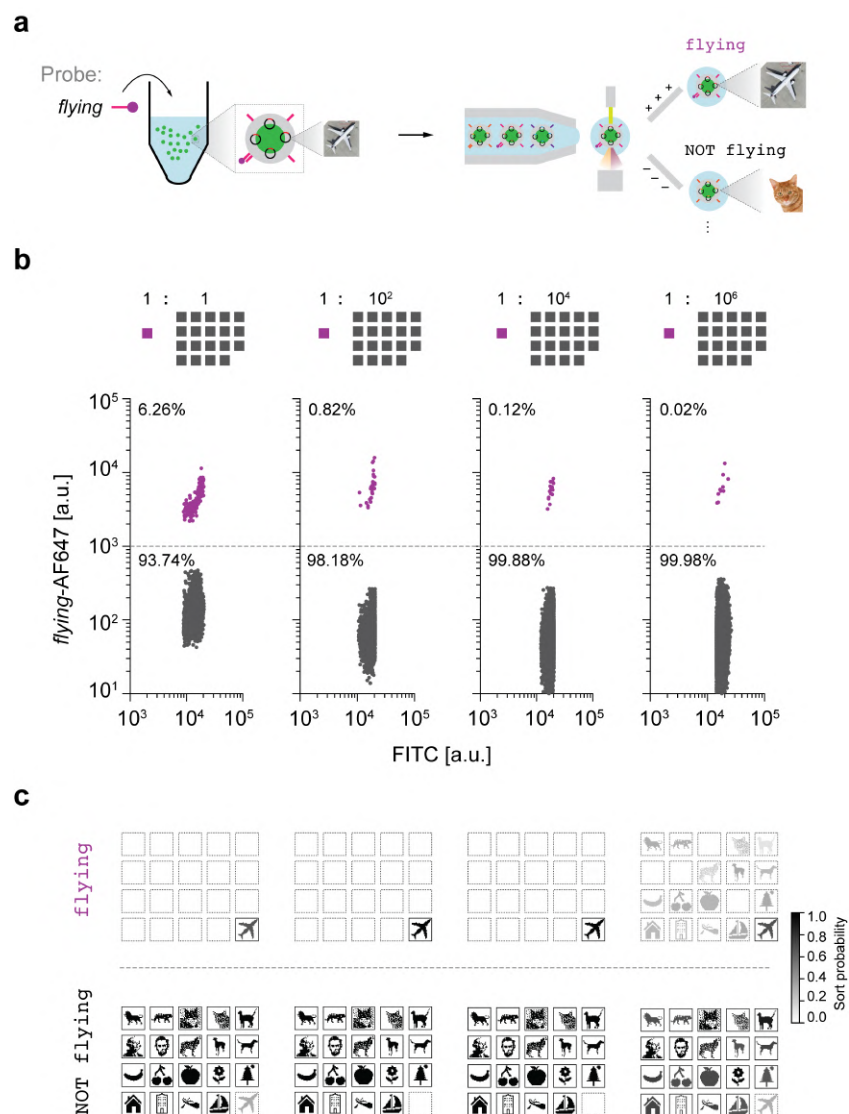
**Online content.** Any methods, additional references, and supplementary information are available at <https://doi.org/10.10XX/XXXXX>.



**Figure 1 | Write-access-read cycle for a content-addressable molecular filesystem.** Colored images were converted into  $26 \times 26$ -pixel, black-and-white icon bitmaps. The black-and-white images were then converted into DNA sequences using ternary encoding scheme <sup>7</sup>. The DNA sequences that encoded the images (file sequences) were inserted into a pUC19 plasmid vector and encapsulated into silica particles using sol-gel chemistry. Silica capsules were then addressed with content barcodes using orthogonal 25-mer single-stranded DNA strands, which were the final forms of the files. Files were pooled to form the molecular file database. To query a file or several files, fluorescently-labelled 15-mer ssDNA probes that are complementary to file barcodes were added to the data pool. Particles were then sorted with fluorescence-activated sorting (FAS) using two to four fluorescence channels simultaneously. Addition of a chemical etching reagent into the sorted populations released the encapsulated DNA plasmid. Sequences for the encoded images were validated using Sanger sequencing or Illumina MiniSeq. Because plasmids were used to encode information, re-transformation of the released plasmids into bacteria to replenish the molecular file database thereby closed the write-access-read cycle.

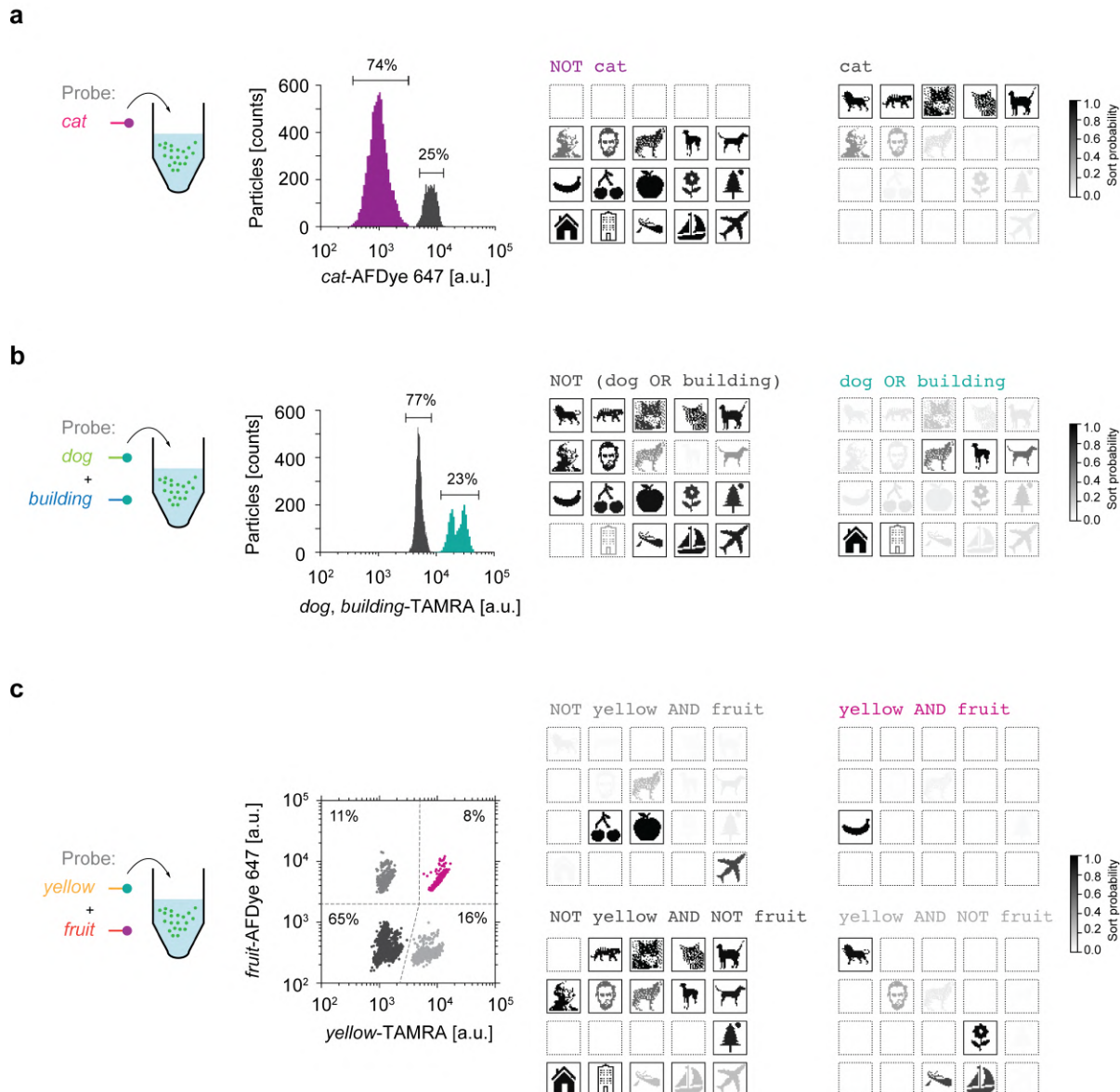


**Figure 2 | Encapsulation of DNA plasmids into silica and surface barcoding. a**, Workflow of silica encapsulation<sup>20</sup>. **b**, Raw fluorescence data from FAS experiments to detect DNA staining of TO-PRO during or after encapsulation. **c**, Functionalization of encapsulated DNA particles. **d**, Scanning electron microscopy images of bare silica particles, silica particles functionalized with TMAPS, and the file. **e**, Distribution of particle sizes determined from microscopy data (left) and zeta potential analyses of silica particles and files.



**Figure 3 | Single-barcode sorting.** **a**, Schematic diagram of file sorting using FAS. **b**, Sorting of *Airplane* from varying relative abundance of the other nineteen files as background. Percentages represent the numbers of particles that were sorted in the gate. Colored traces in each of the sorting plots indicate the target population. **c**, Sequencing validation using Illumina MiniSeq. Sort probability is the probability that a file is sorted into one gated population over the other gated populations. Boxes with solid outlines indicate files that should be sorted into the specified gate. Other files have dashed outlines.

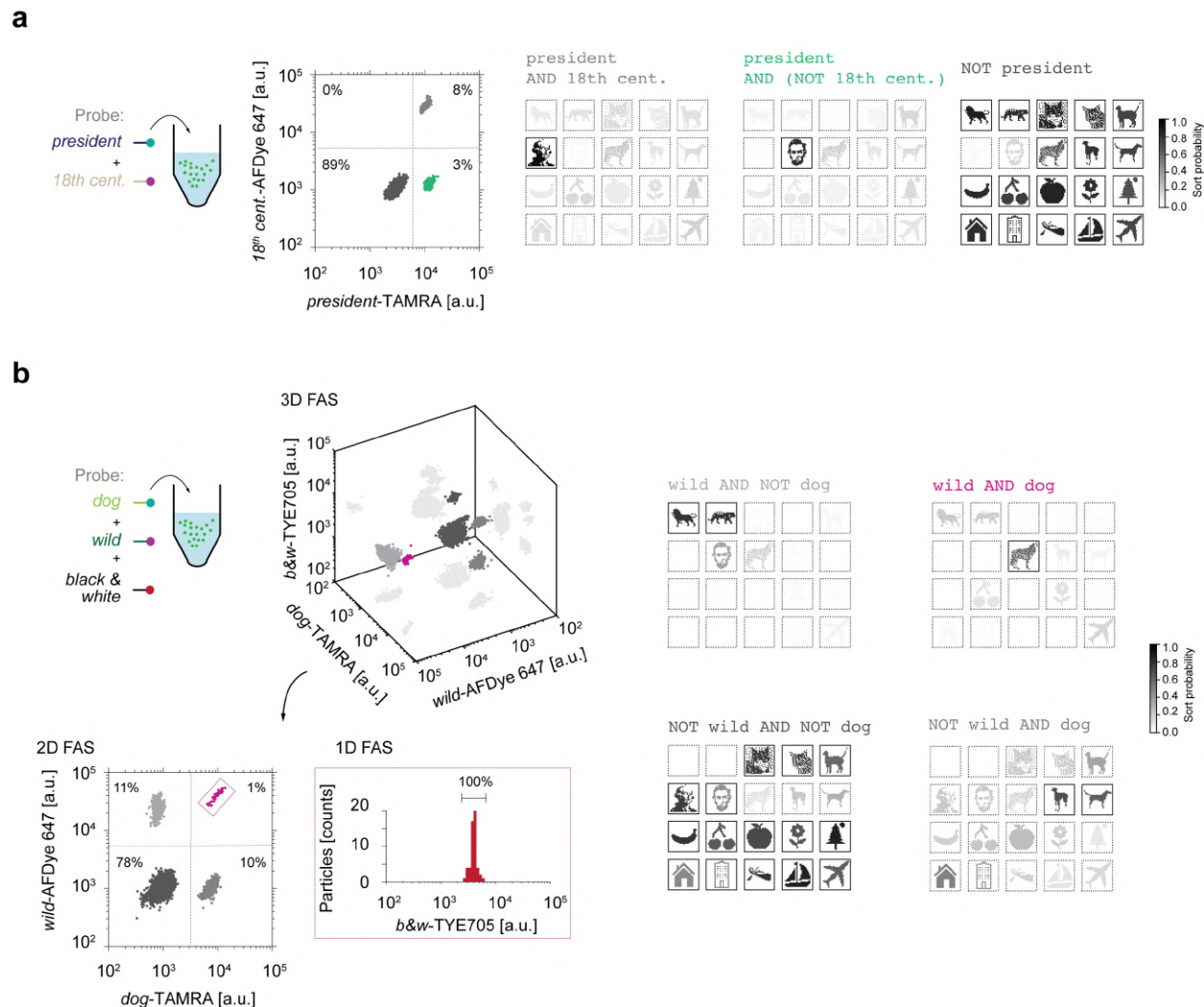




**Figure 4 | Fundamental Boolean logic gates. a**, NOT *cat* selection. Raw fluorescence trace from the FAS system (left) plotted on a 1D sorting plot showing the percent of particles that were sorted in each gate. Sequencing using Illumina MiniSeq tested selection specificity (right). **b**, *dog* OR *building* selection. Raw fluorescence trace from the FAS system (left) plotted on a 1D sorting plot showing the percent of particles that were sorted in each gate. Sequencing using Illumina MiniSeq evaluated sorting using the OR gate (right). **c**, A 2D sorting plot to perform a



444 yellow AND fruit gate. Percentages in each quadrant show the percentages of particles that  
445 were sorted in each gate. Colored traces in all of the sorting plots indicate the target populations.  
446 Sort probability is the probability that a file is sorted into one gated population versus the other  
447 gated populations. Boxes with solid outlines indicate files that were intended to sort into the  
448 specified gate. Other files have dashed outlines.  
449



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

A 2D sorting plot (middle) was used to sort *Lincoln* by selecting a population that has high TAMRA fluorescence but low AFDye 647 fluorescence. Sequencing using MiniSeq offered quantitative evaluation of the sorted populations. **b**, Multiple fluorescence channels were projected into a 3D FAS plot (left and top). There are three possible 2D plots that can be used for sorting. To select the *Wolf* image using the query wild AND dog, a 2D plot of *wild* versus *dog* was first selected and then populations selected using quadrant gates (left and bottom). One of the quadrants were then selected where the *Wolf* image should belong based on the wild AND dog query in

460 order to test whether only a single population was present in the TYE705 fluorescence channel.  
461 Sequencing quantified the sorted populations (right) using Illumina MiniSeq. Sort probability is  
462 the probability that a file was sorted into one gated population over the other gated populations.  
463 Boxes with solid outlines indicate files that would ideally be sorted into the specified gate. Other  
464 files have dashed outlines.  
465

# Supplementary Materials for

## Arbitrary Boolean logical search operations on massive molecular file systems

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 to: [mark.bathe@mit.edu](mailto:mark.bathe@mit.edu)

<sup>†</sup>These authors contributed equally to this work.

## Table of Contents

S0.	Glossary of terms .....	3
S1.	Materials and methods .....	4
S2.	Core memory plasmid sequences.....	5
S3.	DNA encapsulation .....	10
S4.	Barcoding DNA capsules.....	12
S5.	Estimating plasmid copy numbers in DNA files.....	14
S6.	Query probes .....	15
S7.	Estimating surface-accessible DNA barcodes using DNA hybridization assay .....	17
S8.	Sequencing analysis .....	19
S9.	Fluorescence sorting of files .....	23
S10.	Verification of DNA retrieval from sorted sequences .....	28
S11.	Bacterial transformation of sorted sequences .....	39
S12.	References .....	40

## S0. Glossary of terms

Terms	Definition
file	The most basic unit of a file system, consisting of the DNA encoding the main data (the file sequence), addressing barcodes, and any other components necessary for storage and/or retrieval. In our particular file system, each file is a silica particle that contains the file sequence and displays on its surface DNA barcodes describing features of the data. File names in this paper are italicized and the first letter is capitalized. Example: <i>Cat2</i>
file sequence	A DNA sequence that encodes the data in the file.
barcode/addressing barcode	A 25-mer single-stranded DNA sequence that is used to describe a single feature of the data. Barcode names in this paper are italicized and lower case. Example: <i>cat</i>
capsule	A silica particle that contains encapsulated DNA; a capsule with barcodes added to the surface constitutes a file.
data pool/database	A collection of files.
probe	A fluorescently-labelled 15-mer single-stranded DNA sequence that is complementary to a barcode. Probe names in this paper are italicized and lower case. Example: <i>cat</i>
query	A request for a particular subset of a database. Queries in this paper are written in monospaced Courier font. Example: <code>cat AND (NOT wild)</code>
file system	A system for storing and organizing data.

## S1. Materials and methods

**General materials.** All DNA oligonucleotides (oligos), including 5'-amino-modified DNA oligos and TYE705-modified oligos, were purchased from Integrated DNA Technologies (IDT; Coralville, IA) with standard desalting as purification method and were received as dry pellets. Upon receipt of the DNA oligonucleotides, the pellets were dissolved in 1× PBS (catalog number: 79378) from Millipore Sigma (Milwaukee, WI) and kept at 4 °C until further use.

Fluorescein-core silica particles with 5-μm diameter and hydroxyl-terminated surfaces (catalog number: DNG-L034) were obtained from Creative Diagnostics (Shirley, NY). N-hydroxysuccinimide (NHS) ester of TAMRA (catalog number: 1255-25) and AFDye 647 (catalog number: 1121-5) was purchased from Fluoroprobes (Scottsdale, AZ). DBCO-PEG13-NHS ester (catalog number: 1015), and azidoacetic acid NHS ester (catalog number: 1070) were purchased from Click Chemistry Tools (Scottsdale, AZ). TEOS (catalog number: 131903), and APTS (catalog number: 440140) were purchased from Millipore Sigma (Milwaukee, WI). Ammonia in water (28% NH<sub>3</sub>; catalog number: 338818) was purchased from Millipore Sigma and stored at 4 °C. TMAPS (50% in methanol) was obtained through Alfa Aesar (Haverhill, MA) and stored at 4 °C. Anhydrous organic solvents, dimethyl sulfoxide (DMSO; catalog number: 276855), *N*-methyl-2-pyrrolidone (catalog number: 270458), isopropanol (catalog number: 278475), and ethanol (catalog number: 459836), were purchased from Millipore Sigma (Milwaukee, WI). Activated molecular sieves (3 Å; Millipore Sigma; catalog number: 208574) were added to anhydrous DMSO and DMF upon opening.

Gene sequences and all oligonucleotides were purchased as specified from IDT. Gene sequences and plasmid maps are available from AddGene (<https://www.addgene.org/depositing/77231/>). Plasmids were verified by IDT using Illumina MiSeq and by Sanger Sequencing by GeneWiz (South Plainfield, NJ) and Illumina MiniSeq. SeaKem agarose was purchased from Lonza (Basel, Switzerland). SybrSafe was purchased from ThermoFisher (Waltham, MA). Luna universal qPCR Master Mix was purchased from New England Biolabs (NEB, Ipswich, MA). Qiagen (Venlo, Netherlands) HiSpeed Plasmid Midi and Maxi Kits were used for plasmid purification after amplification in 100 mL of LB (Sigma; St. Louis, MO) of transformed DH5α *Escherichia coli* cells (NEB). Illustra S-200 HR spin columns (GE Healthcare; Boston, MA) were used for buffer-exchanged salt removal. Qiagen Spin Miniprep kits were used for small-scale cleanup and concentration of DNAs.

**Characterization of materials.** Dynamic light scattering and surface zeta-potentials were measured using a Malvern Zetasizer Nano ZSP. All samples for surface zeta-potentials were prepared and measured in a standard fluorescence quartz cuvette (catalog number: 3-Q-10) from Starna Cells, Inc. (Atascadero, CA) at a concentration of 0.1 mg mL<sup>-1</sup> with a volume of 700 μL. A universal 'dip' probe (catalog number: ZEN1002) from Malvern was used to measure zeta potential of particles. Scanning electron microscopy of the particles were performed using a Zeiss Gemini 2 Field Emission Scanning Electron Microscope. Samples were mounted on silicon substrates or glass. For glass-mounted samples, 5 nm of gold was sputter-coated to make the samples conductive.

## S2. Core memory plasmid sequences

Twenty  $26 \times 26$ -pixel, black-and-white icon bitmaps were generated as representations of 20 high-resolution color images (**Supplementary Fig. 1**) encompassing a broad range of subject matters. Each black-and-white icon was converted to a length-676 bitstring in a column-first order, with each black or white pixel encoded as a 0 or 1, respectively. This bitstring was compressed via run length encoding, replacing each stretch of consecutive 0s or 1s with a 2-tuple (value, length) to generate a list of 2-tuples describing the entire bitstring. The maximum length is 15; runs longer than 15 bits are encoded as multiple consecutive runs of the same value. This run length encoding was converted to a sequence of base-4 digits as follows:

- 1) Begin with an empty string, and set the current run to the first run of the list.
- 2) Append the value of the current run (0 or 1).
- 3) Using 2 base-4 digits, append the length of the current run.
- 4) If the length of this run was 15, encode the next run starting with Step (2). Otherwise, encode starting at Step (3). If no runs remain, the process is complete.

This process produces a quaternary string describing the entire run length encoding of the image. To avoid homopolymer runs and repeated subsequences in the final nucleotide sequence, each digit is offset by a random number generated from a linear congruential random number generator (LCG) beginning with a random seed (i.e. the  $i^{\text{th}}$  value generated by the LCG is added, modulo 4, to the  $i^{\text{th}}$  base-4 digit of the quaternary string). We used an LCG of multiplier 5, modulus  $2^{31}-1$ , and increment 0, although any LCG parameters with period longer than the length of the sequence would be suitable.

The final “randomized” quaternary string is converted to nucleotides by a direct mapping (0 = G, 1 = A, 2 = T, 3 = C). The number used to seed the LCG is prepended to this sequence by converting it into a ternary string of length 20, whose digits are encoded in nucleotides via a base transition table, as done previously by Goldman et al.<sup>1</sup>: (0 = GA, AT, TC, CG; 1 = GT, AC, TG, CA; 2 = GC, AG, TA, CT). The first digit is encoded directly (0 = A, 1 = T, 2 = C).

This sequence was modified with additional flanking sequences added to the beginning and end. A 64-bit wavelet hash of each bitmap was calculated using the whash function provided by the ImageHash Python package, available through the Python Package Index (<https://pypi.org/project/ImageHash/>). The 64-bit hash was split into two 32-bit halves, each of which was represented in a length-24 ternary string that was subsequently converted to nucleotides through the same process as applied to the seed. The two 24-nt regions were appended to the beginning and end of the sequence (**Supplementary Table 1**, orange text). The sequence containing image hash, seed, and image encoding, was additionally flanked on the 5' and 3' ends by sequences (5'-CGTCGTCGTCCTCAAAC-3' and 5'-GCTGAAAAGGTGGCATCAAT-3', respectively) that allow amplification from a “master primer” pair that would amplify every sequence in the molecular plasmid database (**Supplementary Table 1**, purple text; **Supplementary Fig. 2**).

The final sequence was checked for problematic subsequences, specifically, GGGG, CCCC, AAAAA, TTTTT, and the restriction enzyme recognition sites GAATTC and CTGCAG. If any of these subsequences were found outside of expected occurrences in the constant flanking regions, the entire sequence was regenerated with a new random seed until no such subsequences appeared.

The generated sequences were cloned on a pUC19-based vector. The software for sequence encoding and decoding is publicly available on GitHub at <https://github.com/lcbb/DNA-Memory-Blocks/> and the plasmids are publicly available from AddGene (<https://www.addgene.org/depositing/77231/>). Each master primer and hash barcode pairs were verified by PCR and agarose gel analysis and PCR bias was checked by qPCR (**Supplementary Fig. 3**).



## A Direct conversion to bitmap icon

Cat2



- cat
- domestic
- orange



black & white



26×26 bitmap  
(diffusion dither)

Cat1



- cat
- domestic
- black & white

Cat3



- cat
- domestic
- brown

Wolf



- dog
- wild
- black & white

## B Bitmap icon representations of associated images

Airplane



- man-made
- air
- flying

Apple



- fruit
- red
- seeds

Banana



- fruit
- yellow
- seeds

Canoe



- man-made
- water
- oars

Cherries



- fruit
- red
- pit

Dog1



- dog
- domestic
- brown

Dog2



- dog
- domestic
- black & white

Flower



- plant
- white
- yellow

House



- man-made
- building
- wood

Lion



- cat
- wild
- yellow

Lincoln



- human
- 19th century
- president

Sailboat



- man-made
- water
- sails

Skyscraper



- man-made
- building
- steel

Tiger



- cat
- wild
- orange

Tree



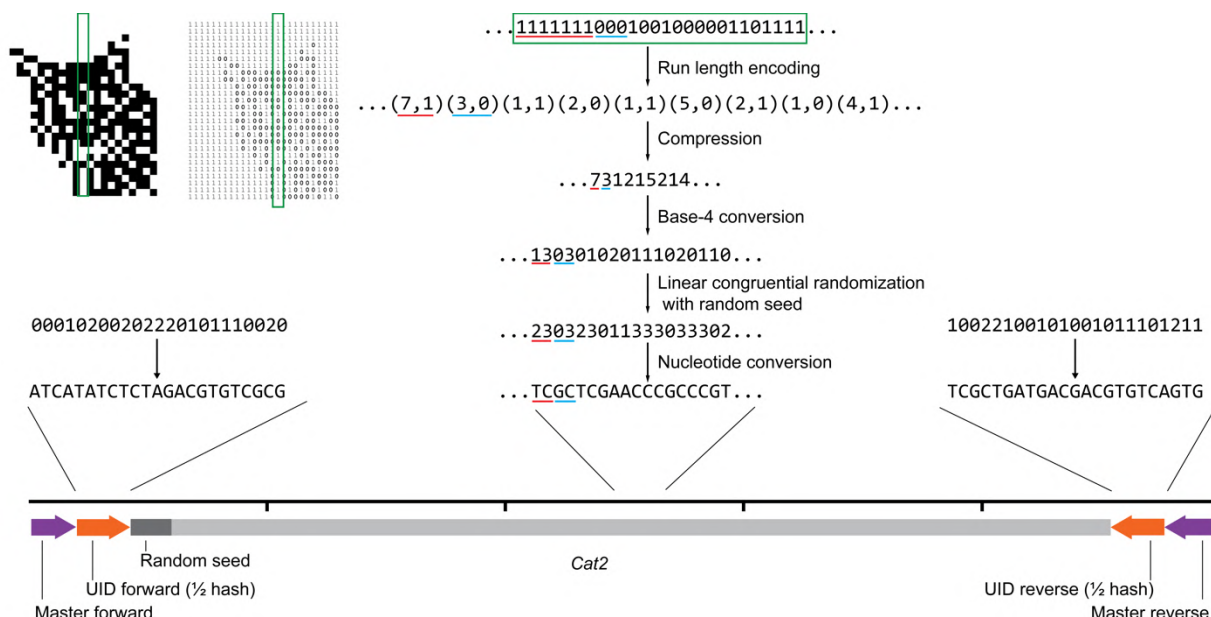
- plant
- tree
- moon

Washington



- human
- 18th century
- president

**Supplementary Figure 1. Image database with icon representations and content descriptors of the original images.** (A) Three pictures of cats and one picture of a wolf were directly converted to  $26 \times 26$ -pixel icons using Adobe Photoshop, by first changing the image to black and white, and then reducing the resolution to 26 pixels per inch with a 1-inch $\times$ 1-inch image using diffusion dithering. (B) Sixteen other images were additionally selected of broad subject matter and icon images were chosen for image representation and reduced to  $26 \times 26$ -pixel sizes. Icon images were used to reduce sequence size and therefore cost of synthesis.



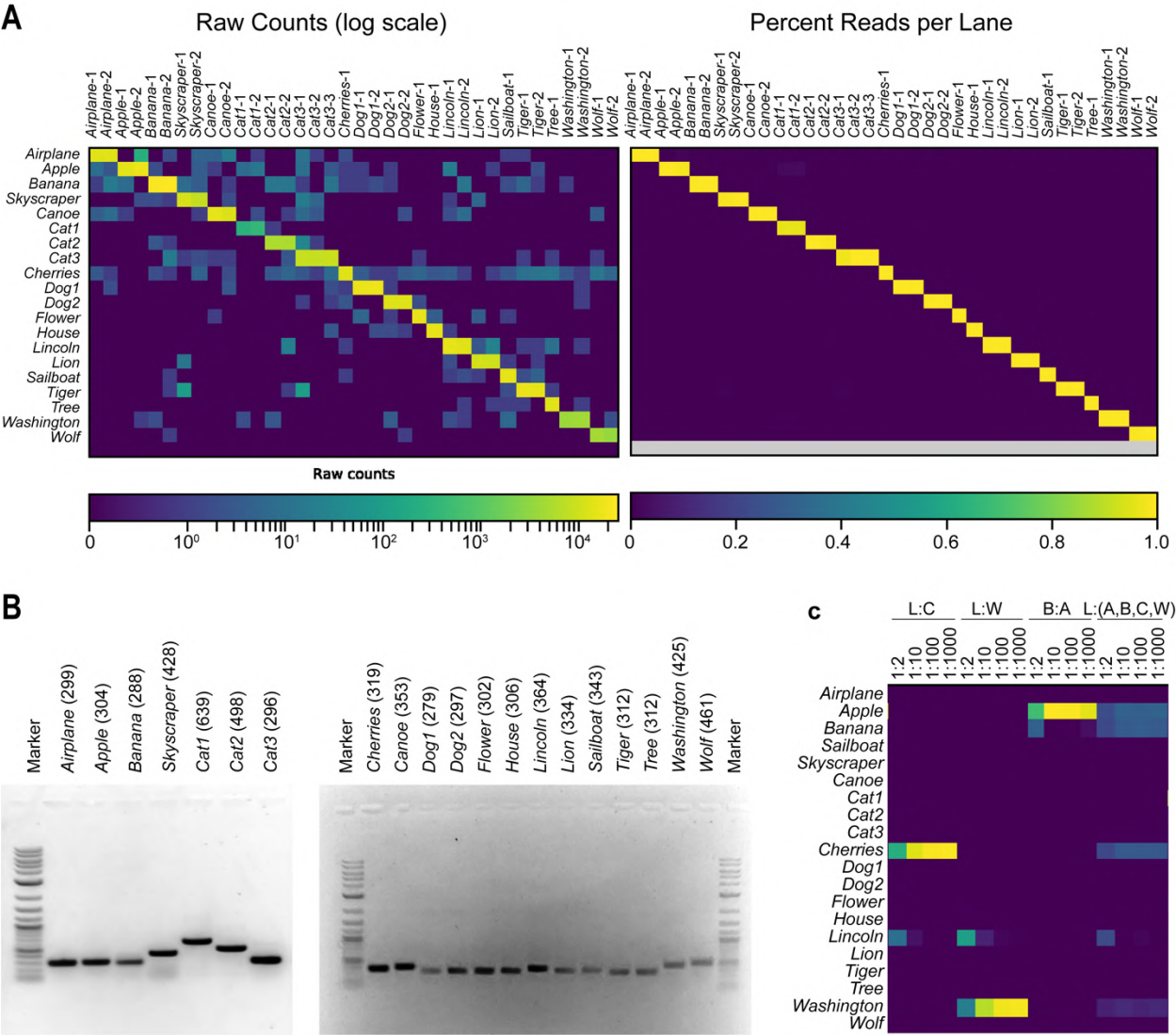
**Supplementary Figure 2. DNA encoding of representative icons.** The black and white 26 × 26-pixel icon is converted to 1 (white) or 0 (black) and the bitmap is generated, followed by run length encoding compression, a base-4 conversion, encryption by simple randomization with a seed (encoded in the DNA), and converted to nucleotides. A binary 64-bit wavelet hash is calculated from the icon image and split in half and converted to DNA to act as a UID primer pair (orange). A master primer pair is added to flank the construct. *Cat2* is shown as an example, but all icon images go through the same workflow.

**Supplementary Table 1. DNA sequences for encoded plasmids.** Purple text indicates master primer regions and orange text indicates hash primer regions.

Image	Insert sequence (5' to 3' direction)
<i>Airplane</i>	CGTCGTCGTCCTCAAACATCATGCTACACTGTACAGTGCATACACTGAGACTGTATCGCACGAAATTTGAGATATTTTCGTATC CTTGCCCTTTGCTAGGCTGAATTTGTTGGGCAGGCTCCTCGCACCAGATGATATTCGGGATAGCCTCCAGAATCTTAAGATTAAAAATCC AAAGGCTCCAGCCCTGACATCTTCCAGAGATGCTTCCCTCCAGAGACTCGCTGGTGGTTATAGGCTACCTTTGTTATAGAACGAT CGCTGCTCTGACTGTAGATATCGCTGAAAAGGTGGCATCAAT
<i>Apple</i>	CGTCGTCGTCCTCAAACATCATATCTCGTCGATCGCTCAGAATGAGCTGTGACGCGTGATCTAACATGATTTTGAACGGGA TCGGGAACGGTGGTCAACCTGCTTTTCTCTTCGGCTTAGCATGCAGTCTCATTCCGACGATCATCGCGTCCCATAGGATATACCA AATTATCATGAACAGTTTCTGCTATAGGTCCATAATTCAGTTCTTACATGCTGAGAATTGAATGAGGCTTTCACCTTTCTGCGGA TATCGCTGACACTATAGACGTGCTCGCTGAAAAGGTGGCATCAAT
<i>Banana</i>	CGTCGTCGTCCTCAAACATCATCTGCGCTGACATGTAGCTTATAGACTACACATGCACACTCCGTGTGTCATTGGAGGTGCT CAGTCCGGTGCAACTTTTCCGAAGTACTTTTGGTGGGCCAGATTTGACGGAAGCTTGGTGTACACTTAACCTGTGTTTATTGTCT ACCAGGAATTAACAGACTTACGTTGGGACACTAGTCAAGTAACTAGCCTCAACAAGTATCAGGCCAAACACATCATGTGAGATCGT ATCTATCTCGCTGAAAAGGTGGCATCAAT
<i>Sailboat</i>	CGTCGTCGTCCTCAAACATCATATCTCTATCTACTACTGTGTATATACTAGAGCTGCGACTCGCCCTACCGATTAGGAGTGCC ACTGGGAGCAACCGCTAACTGTAAGGTGGGTACCCGACGGTGCCGACCAAGAGCGAGGGCGGATGGAATATGAGGCGATGCCGACTT TGGCTACTGCACAACCGTAGTTTAGCCGCGCCCGAGGACTCTAGAGTATAGAACCAGTACGACAGAAGCTAGTACTAAGTTTGCC GATATGGCGTAGTGAGCGGTGCAAGGAGACGTCCTGGCTAATCATGTGAGACGATGTGAGCTCGCTGAAAAGGTGGCATCAAT
<i>Skyscraper</i>	CGTCGTCGTCCTCAAACATCATATCTCAGCTCATCTACTAGACATATGCATCGTGTGACAGTGTGTCAGAGTGTGTCAGGACTATGAGT AGGTGTTATCACTACGCACACCAACTTGTGGTAGCGGTAAGACATAGACGCGGATCCTCCTCGGCCCTACGTAAAAATTCCTGCTG CGTGACTTAAGACTCCTCGCAAGACGGAGATCTCGCGTTTGGCGGAAGAAGCGCCGACCGTGGTGGAAAGGTATCGTTACTAAC GAGGCAGGCGCAAGATGCGTGTCCGTAGGCCGAGTGGATAATGACCGTGGCAGTGATCTACAGTGTCTAGCGGTACTCTGTACCG GCGGTTCCCGATTGGTAAGTTACCACCATTTGTACACTATCGTCAGCGTAGCATGCGCTATAGCTGAAAAGGTGGCATCAAT
<i>Canoe</i>	CGTCGTCGTCCTCAAACATCATCACACTAGCGCTGTGCGACACTAGCAGTGTGACGTGATTTCCACTCATATCAACTCATGTT CTGGTGTGCGTGTGAGCTGAGCTACACATGAATACCATTAGACTTTCCGGGCACGAATAAACCGCCTGAGCACGAGAACCTAAGT TCGTCAAACCACTTCCAAGTCTGGGAACCAAAATATCCAAAACCGCTTCTCTTTGACTGCTAGGATCCAAGCAATGGAAC GCAAAATGCCAATAAGATGTGACATGTATGCTTGCATACGCGGTAGCATCGCTGACAGATCGTCTGCTAGCTGAAAAGGTG GCATCAAT

Image	Insert sequence (5' to 3' direction)
<i>Cat1</i>	CGTCGTCGTCCCTCAAACATCATGTCTATGAGTACACTCGCGATGTGACAGTCACAGATCGTCGCATAAGTGCTCCACCCTGTGTCATTGGCGCCAGTAATGTGCCATGTCACTCACTAAGAATAAGTAAAAATGAAGCTGGGAGCGCCCGGAGATTAGTCTGCCACGACTCAACCTTAGATCGAACGTCCTCTCTCCGATAGTTAGGCCATACTGGTACGTATCAAGGCTTCGGGTTGAATGAGCACCAAGTCCTCGCCTTAACTCGTCCGGCCTCTTGAGGATTATCTTTAGTTACTGGAAGTAGGTACTGAAAAGCATCTGCGCCCTAGCAAGGCTTATTTTAGGATCTCAGTGGGATGGAAGCGTATGCCACATGGTAGCGAAAAAGTCGTCTTGCTGTGCCGGAGTCGGCTACGGCCTGATGAGCTGAGGCGGAGGGCTTGCCCTCGAAAGCTCTACTAACAATTCAATAAATGTGGGAATGCTACATAAATGTCGTAGTACCGTCAGAGATAGGAGACGGGTCGATTAACACTTCTACGCAGGGTTATAATAGTTTGACACTCAGCAATCGCGACACTACACGCTGAAAAGGTGGCATCAAT
<i>Cat2</i>	CGTCGTCGTCCCTCAAACATCATATCTCTAGACGTGTGCGAAGTCGACACAGCGTACACGCTGTGCGTTACAACGAGTGAGCTAACATATGCAGATCACGTTTGACGGGAGTACATTATAAATCACC CGCAGCTTAGTCTCGAGCTCGCCGTGGAGATAGATAGTGTGCCAGCATCGACGCTCGCTCTCGTCAGCAGTAACGCGCTCGAAACCAGTCTGAGACACGAAATCGTATATCTGTCAATTCGTCGATGGGCTATTGGACCGGACAGCTCTCGCTCGAACCCGCCGTGCTAAAAATTTGCGGACTATTTAAACAAGAGTCATCCTGTTCTACTACAGTGAAGCCCTGTTGATGGGCGCTGCTGGAGAGTAGTACGGATGACAACATTAAACGACGAGGCCGCTACGAGATATAC TAGGAATCGAGGTCGTCCCGAGCATCGCTGATGACGACGTGTGAGTGGCTGAAAAGGTGGCATCAAT
<i>Cat3</i>	CGTCGTCGTCCCTCAAACATCGCGCTCACACGCGAGACAGAAATGCGAGCTATACGACATCTAACACCCGTACCGTCACTCCAGGCGGATGTAACGTTGTCTGGAGAGTTGCGTGTCTCCTCGCATCGAAGCAGCAGCATAGCTTGTCTGGTTATCGCGCGTAAACAGAGGCCACGCGGCTGAACCCGCTCGCCGCTTGCCACGAAATCTGTGATTACAAATTCGTCTCTCTGTAAGGTGCTATCGCTGAGCAGACGTGCACACTAGCTGAAAAGGTGGCATCAAT
<i>Cherries</i>	CGTCGTCGTCCCTCAAACATCATATGCGCTAGCTCACAGAGATCTATATGTGCGTGAGCGAGTCGAGTTGCGCGGGTCACGTGCTAACCTAGAGTAATATATAGCGAAATTCACAAGTTCCTGATCATTAGTTAAATGCATAATCCGATCCGCTACGCCGACTCGCGGAGGAGCGAGAGCTGCTGATGTTTTCAGGTGATGCTCAAGGCCAGTACAAGATTGCTTGAGAGTTCCACCGAGCGGGCCGGGTTCTGTTCCGTCACTTGGCGTAGATCGCACGAGACTATGTAGAGTAGGCTGAAAAGGTGGCATCAAT
<i>Dog1</i>	CGTCGTCGTCCCTCAAACATCATATATCTCACTGCGACTATCGATGAGTGTGCTCGCAACGGGACGTGCAAGGTCCAACGTATAAGGTTCCCTATACCACAGAACACGGACGGGTATGCTACACCGGGATTCTAACGCCCGTTCGGCCGCACATCTGGGACTTAATTGTCTAGCAACTTAGCTGTTTATGTGCTTAGTGACGATATGTCACGTTGTCCTGTTGGGATCGCTGACACTATAGACTGCTGCGCTGAAAAGGTGGCATCAAT
<i>Dog2</i>	CGTCGTCGTCCCTCAAACATCGCGCTCTGACGCGTCATCGACAGAGCACGTCGTGTGCGCTCTTCCATAATTAACACGAGGACTCGCGGGAACACCCATCGGATCGCATCCGAAGTATGGATAGGACTAAAGGAAACCGCGTGTGCTGCAAGTGAACCTCTCCACCTTCGCCAATGTTAAATGACCTAACATTGACAGAAATAGCGCAGCTCTTGACATGCCCCGTGCAACGCATCAAGCGTGGGCTATCGTAGATATACAGCAGAGCTGCGCTGAAAAGGTGGCATCAAT
<i>Flower</i>	CGTCGTCGTCCCTCAAACATCATATCTCATCTCTATCAGACGACACGCTGAGACGACATATGTAGCGGCTGTGGTTTATCACAGCCTATCATTCAACCTATTATGGGTGAGTCGTATATGAGCGAACTGTTGATGCGCTCCGAGGCGCTCAAATAACAGCCGACAGAGTGGCCGTTCTCAGATCAGCTCCAGTTTAACTGTGAGCGGGAGATAATAGCAGGCCAGGTACAGAGGAGCTGTTGACTAATGACTAATCGATCGACTGAGAGCATGTGCGAGCTGAAAAGGTGGCATCAAT
<i>House</i>	CGTCGTCGTCCCTCAAACATCATATCTCACTGCGACTATCGATGAGACATACTCAGCTCGTAACCGTCCGGTGTGCAAGGAAGCTCTTTGTATTGGCGGTGTACGCGGTACAAGGCCGTATTATCGGTATGACACTCGAAAGTGAAAATATATATGGTGATACCTTAATGGCGTAACGCGACCCCTACGTAGCTATGAGCGCAGCAGCTCTTAGTGTTCTTGAGCGCCCTCATTACATCTGCCGTACCTACTATCGTAGTAGACAGCGCACGTGCGCTGAAAAGGTGGCATCAAT
<i>Lincoln</i>	CGTCGTCGTCCCTCAAACATCATATCAGAGCACACACGCTCGATGATCTGTGCTGCTCACTCTGATACTCACACTGCGCACGAGGTGGGACCTCGAGATGGTGAGCTCACTGGCTCCGGAACGTGCCGCTTACCAGTCCACTTTCTAACAAAGGTGCGCTCATGCAAGCTATCCTTGTGCTCGCGGAGCGCGCAGCGATTGTGAGTTAGTCTCGTTAGTCAACACTGGCTATGATTAGTATTGCGACC TTGTGCCGTATAGTAGTTACCCACTGTTAACACCGTGGGCATTGTGCGACAATAGTAAACATCGACTCTGATATATGATGCTCAGCTGAAAAGGTGGCATCAAT
<i>Lion</i>	CGTCGTCGTCCCTCAAACATCGCTAGCAGCGCTCTGCTCACAATACTCTATGAGCGCTCTACGAAAGTCTCCGTCTCTGAGAGCGACATAAACTGCTCGCTACACGCTCTCGGATCGATTTTACATCACACCGCAGCGATGGTTGAGCTCGCAGCTATCAATCTTAGCGAGGAAAGCGTCGACCCCTATGCATCACAAGGGCTATTACCGGTGACTCACCTTCAACTGGGCATCGGGCGTTGCATCATAGGTTACAATCAGTAGATATGTTATGCTTCGTAATGTCCCATCATCTGACAGATGCAAGTCGCTAGCTGAAAAGGTGGCATCAAT
<i>Tiger</i>	CGTCGTCGTCCCTCAAACATCATCTGCGCTGACGATCGTATATGATGACTGTACATAGACATTCTTGTCAGCGGGCGGACTGCTGGGTTACGAGACAGGTATGGGTTGCACGTGATATGACCATGCCTAGATAGTGCCTGCTCGGGTGTGCGCAATGTTGATTGCTTGCTCATTACTTCCGGAGGAATCCCGACTGTATTGTTAAAGGTGGGCGGCATAGCTGTTTAAATGCTGCCTATAACCCGAAGGGTTCCGTCCATTGGATCGCTGTATCAGAGCTCTACGCTGAAAAGGTGGCATCAAT
<i>Tree</i>	CGTCGTCGTCCCTCAAACATCATATCTCTGCTACTCGCGTCACTAGTCGTGCACTGCTATTCCCAACAACCATGCCATAGTTTTTATGATGCATACGGAGAGAGTCATGTGAGTTGGGCCGATTATCCGAGACAATCGTTTCATACCTAAGGAACTGGTCACCTTACATTCCTCGGTCTAAGAGCCAACACTTGATACTACACTAGCAACAGCGGGTAGACCAATATGATTTTGATCAAAAGTACTTGCGTAACCTGTTATCGATCGATCGACATCTCTGACAGCTGAAAAGGTGGCATCAAT
<i>Washington</i>	CGTCGTCGTCCCTCAAACATCATATCAGCTCGCTACGTCAGATTATGAGCAGTCATCACGACAGATGTCCATTCCAATGGTACACATAGTGTGTCCTAGATCGACAGTATTCGCGGCGAGTATTATTAGCGGCTCTGTTTCAGCACCCGGAACAATAAGGACGAATAGTATACATATGCCTAATACGTTTTCACGTGCGGCACAACCTATAGAGAATATCCGGGATGACGCGTGACAATCAGTGATCAATGTAAGCTTTAACCGATTAGGTAGCGTCGTGGTTAACCGAGACTGTGAGTTACTCTATGTTCTGCAATAGCCGAGGTACCTTTGTTCTCAACCTGGGCGGCAATCGTAAATCCTGCTTATAGCTAGGATCATCGACGCGAGAGCTATCTAGGCTGAAAAGGTGGCATCAAT

Image	Insert sequence (5' to 3' direction)
<i>Wolf</i>	<p>CGTCGTCGTCCTCAAATATCATCTCTAGCGAGCTACAGCACTCACACACGATGACGATCATGTACGAGCTCGGCACTGATATG</p> <p>GAGACTCTCTATCTCAGGCCTCTACTTCTACACTGAGAGATTATCTGGTCTTTCCGCCGAGATTAGCTAAGCTTAAAAATCGTTCA</p> <p>ACGGTAACTAACCTGCTTTCTTGAGGAATTTCCGGCCGCTCCCGGTACATCGTACCCGTTTACTACCAACTTACAACACATGGC</p> <p>GCAGATCTAAGCATCGAGTTTCAAATGCAATTTCTATGATACCCGCTGTACTCCACCAAAGATACGAAATGCCACCTTCTAAA</p> <p>AGTTATCTCCACATATAGAGCGACTACGCAGAAGGCAGGAATAACTCCTATATTTAGGTTTATGCTCGGTTTATCGCTGACGCGCT</p> <p>CAGTACAGCAGCTGAAAAGGTGGCATCAAT</p>



**Supplementary Figure 3. DNA plasmid database characterization.** (A) Illumina MiniSeq was used to validate the plasmids after DNA generation, showing baseline purity for each. (B) Each plasmid was amplified using the master primers and agarose gel was used to validate sizes, shown in parentheses in nt. (C) Dilutions of plasmid databases were generated, *Lincoln* (L) to *Cherries* (C), *Lincoln* to *Washington* (W), *Banana* (B) to *Apple* (A), and *Lincoln* to *Apple*, *Banana*, *Cherries*, and *Washington*. Illumina MiniSeq was used to read out the population after amplification with the master primers and addition of adaptors and sequencing barcodes.

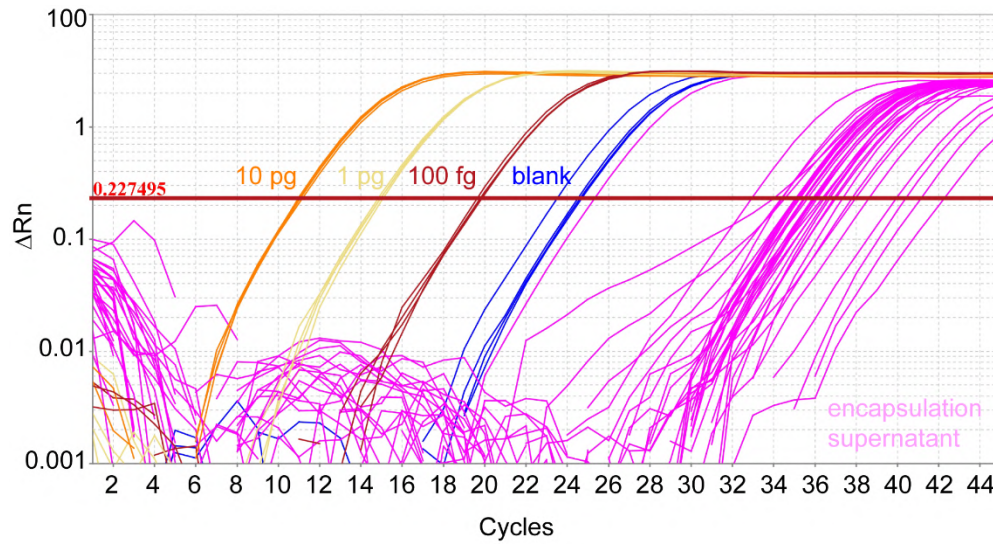
### S3. DNA encapsulation

**Functionalizing 5- $\mu$ m silica particle with TMAPS.** We adapted a procedure published previously<sup>2</sup>. A volume of 1.0 mL of 50 mg mL<sup>-1</sup> fluorescein-core 5- $\mu$ m silica particles was added into a 2.0 mL DNA/RNA LoBind Eppendorf tube. The particles were centrifuged at 1,000 rpm for 10 seconds using a benchtop centrifuge. The particles were redispersed in 1.0 mL anhydrous ethanol with vigorous vortexing. The particles were centrifuged and redispersed in ethanol five times. We then added 50  $\mu$ L of 50% TMAPS in methanol to the dispersed 5  $\mu$ m silica particles (50 mg mL<sup>-1</sup> in ethanol). The mixture was stirred overnight at room temperature using a thermal mixer from Thermo-Fisher (Waltham, MA) at 1,200 rpm. The mixture was centrifuged at 1,000 rpm and washed with ethanol five times to remove any unreacted TMAPS. The functionalized particles were finally redispersed in 1.0 mL DNase/RNase-free water. The particles were stored at room-temperature until further use.

**Encapsulating plasmid DNA.** We adapted a procedure published previously<sup>2</sup>. For each data encoding plasmid, a mass of 1.0 mg of TMAPS-functionalized, fluorescent 5  $\mu$ m particle was added into a 2 mL LoBind Eppendorf tube containing 15  $\mu$ g of plasmid DNA dissolved in 1 mL of water. The mixture was mixed gently using a tube revolver (Thermo Fisher) at 30 rpm and at room-temperature for 5 minutes. A volume of 10  $\mu$ L of 50% TMAPS in methanol was then added to mixture and stirred for 10 minutes, at 1000-rpm, and at 25 °C using a thermal mixer. After 10 minutes, a volume of 2  $\mu$ L of TEOS was added and the mixture is stirred for 24 hours, at 1000 rpm, and at 25 °C using a thermal mixer (Thermo-Fisher). An additional 5  $\mu$ L of TEOS was then added and the mixture is stirred for 4 days, at 1000 rpm, and at 25 °C which forms our encapsulated DNA particles or DNA capsules. The mixture is centrifuged at 2,000  $\times$  g for 3 minutes to sediment the DNA capsules then the supernatant was removed without disturbing the particle sediment pellet. The particles were washed repeatedly for 5 times by re-dispersing the particles with 1 mL of water, sedimenting the particles with a centrifuge at 2,000  $\times$ g for 3 minutes, and removing the supernatant. After the final wash, the DNA capsules were re-dispersed in 1 mL of ethanol with 30 seconds of vortex mixing. A volume of 20  $\mu$ L of  $\gamma$ -aminopropyltriethoxysilane was then added and the mixture was stirred for 18 hours, at 1000 rpm, and at 25 °C using a thermal mixer. The mixture is centrifuged at 2,000  $\times$  g for 3 minutes to sediment the amino-modified, DNA capsules then the supernatant was removed without disturbing the particle sediment pellet. The particles were washed repeatedly for 5 times by re-dispersing the particles with 1 mL of *N*-methyl-2-pyrrolidone, sedimenting the particles with a centrifuge at 2,000  $\times$  g for 3 minutes, and removing the supernatant. After the final wash, the DNA capsules were re-dispersed in 1 mL of *N*-methyl-2-pyrrolidone with 30 seconds of vortex mixing and the resulting colloidal suspension was then transferred into a clean 2 mL Eppendorf LoBind tube.

**Encapsulation efficiency.** After four days of encapsulation, all aqueous washes were collected. The amount of DNA that remained in each was estimated using qPCR (**Supplementary Fig. 4**). The washes were amplified using the master primer pair. All of the washes show plasmid concentrations that are below blank, suggesting that our encapsulation efficiency is quantitative.





**Supplementary Figure 4.** PCR amplification curves ( $\Delta Rn$  vs Cycle, log scale) shown for the encapsulation supernatant for all plasmids after four days of encapsulation, with 10 pg (orange), 1 pg (yellow), and 100 fg (red) of the *Cat2* plasmid used for calibration. Overlay of 20 plasmids individually amplified with master primers (magenta) and blank (blue).

#### S4. Barcoding DNA capsules

**Metadata descriptors DNA hash Supplementary Table chema.** The subject matter of each of the original high-resolution images was associated with metadata. The subject matter included sets of cats and dogs, both wild and domestic, and of a variety of colors: a domestic black-and-white cat (*Cat1*), a domestic orange cat (*Cat2*), domestic brown cat (*Cat3*), a domestic brown dog (*Dog1*), a domestic black-and-white dog (*Dog2*), a wild yellow cat (*Lion*), a wild orange cat (*Tiger*), and a wild black-and-white dog (*Wolf*). Also included in the database were historical US presidents (*Washington*, an 18<sup>th</sup> century president, and *Lincoln*, a 19<sup>th</sup> century president); man-made objects (*Canoe*, *Skyscraper*, *Airplane*, *Sailboat*, and *House*); fruits (*Apple*, *Banana*, and *Cherries*); and plants (*Tree* and *Flower*) (**Supplementary Fig. 1**). From these descriptions, each encoded image was annotated with three semantic metadata descriptors (**Supplementary Table 2**) associated with the original image. A table was then generated to associate each descriptor with a unique sequence chosen from a list of 240,000 orthogonal barcode sequences <sup>3</sup>.

**Supplementary Table 2. Metadata single-stranded key-value pairs.** Single-stranded DNA sequences were purchased from IDT. /5AmMC6/ denotes an amino hexyl modifier on the 5' end of each ssDNA sequence.

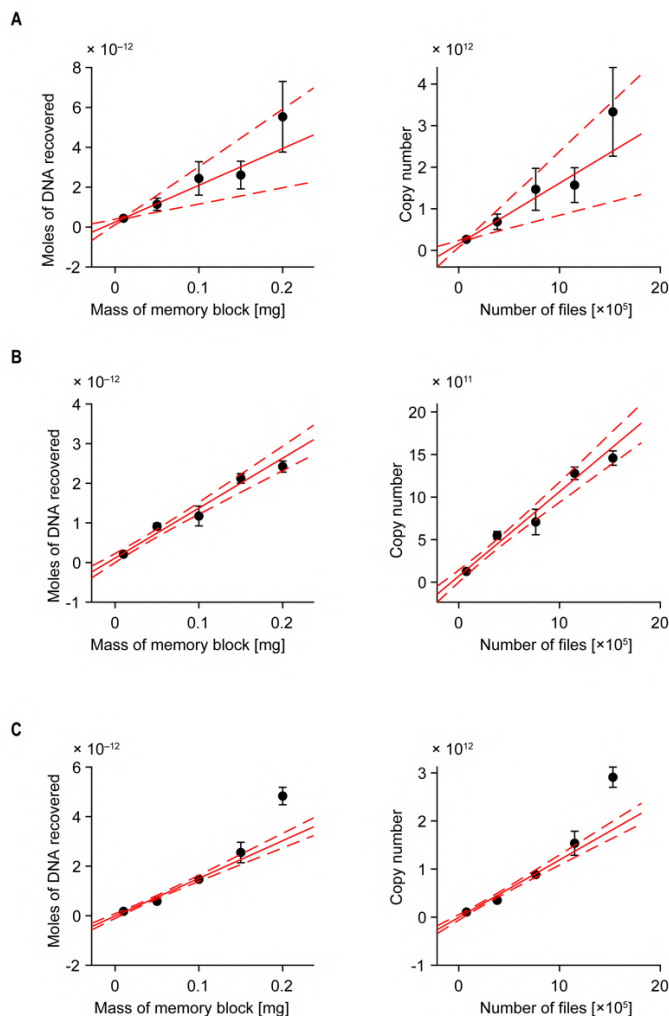
Images	Barcode 1	Sequence	Barcode 2	Sequence	Barcode 3	Sequence
<i>Airplane</i>	<i>man-made</i>	/5AmMC6/ATGGATGCACGT CCACAAGAAGCAG	<i>air</i>	/5AmMC6/TAGAAGCGTTCC GACGAAGTTACCT	<i>flying</i>	/5AmMC6/GAGATTATTTC TCGTTCCGCCAG
<i>Apple</i>	<i>fruit</i>	/5AmMC6/GAAGTTATTCGG TATCTGTGCCGCT	<i>red</i>	/5AmMC6/AGCGCTTGGGA CACGTGAAGTAAC	<i>seeds</i>	/5AmMC6/ACGTCACGTCGC CTATGGCGTTATT
<i>Banana</i>	<i>fruit</i>	/5AmMC6/GAAGTTATTCGG TATCTGTGCCGCT	<i>seeds</i>	/5AmMC6/ACGTCACGTCGC CTATGGCGTTATT	<i>yellow</i>	/5AmMC6/TAATGTGGCTTG GCTCACCCTAGG
<i>Canoe</i>	<i>man-made</i>	/5AmMC6/ATGGATGCACGT CCACAAGAAGCAG	<i>water</i>	/5AmMC6/CTGGTTTGATCC GACACATTGATTC	<i>oars</i>	/5AmMC6/GTTTCCGCATAA ACTCAGGGGAGTC
<i>Cat1</i>	<i>cat</i>	/5AmMC6/AACGATTGTTAT GCCCCCTAACTCAG	<i>domestic</i>	/5AmMC6/TCTTAACAAAGG ATGGGCAGGTCGC	<i>black &amp; white</i>	/5AmMC6/TTCAGGGTGGAA GTACCTCCAGAT
<i>Cat2</i>	<i>cat</i>	/5AmMC6/AACGATTGTTAT GCCCCCTAACTCAG	<i>domestic</i>	/5AmMC6/TCTTAACAAAGG ATGGGCAGGTCGC	<i>orange</i>	/5AmMC6/CTGAATACTACA CGCCGTGGTGAAG
<i>Cat3</i>	<i>cat</i>	/5AmMC6/AACGATTGTTAT GCCCCCTAACTCAG	<i>domestic</i>	/5AmMC6/TCTTAACAAAGG ATGGGCAGGTCGC	<i>brown</i>	/5AmMC6/ATCTATCTGTTG GAGTTAACGTACC
<i>Cherry</i>	<i>fruit</i>	/5AmMC6/GAAGTTATTCGG TATCTGTGCCGCT	<i>red</i>	/5AmMC6/AGCGCTTGGGA CACGTGAAGTAAC	<i>pit</i>	/5AmMC6/GAGCCGATTAG TAGCAGTGTCCA
<i>Dog1</i>	<i>dog</i>	/5AmMC6/AAAAGCAAGGTC GTTACATGGAGTT	<i>domestic</i>	/5AmMC6/TCTTAACAAAGG ATGGGCAGGTCGC	<i>brown</i>	/5AmMC6/ATCTATCTGTTG GAGTTAACGTACC
<i>Dog2</i>	<i>dog</i>	/5AmMC6/AAAAGCAAGGTC GTTACATGGAGTT	<i>domestic</i>	/5AmMC6/TCTTAACAAAGG ATGGGCAGGTCGC	<i>black &amp; white</i>	/5AmMC6/TTCAGGGTGGAA GTACCTCCAGAT
<i>Flower</i>	<i>plant</i>	/5AmMC6/TAAGCAATGGGT TCCACACTACGTA	<i>white</i>	/5AmMC6/TTTTATGCCGTG TTGTTGCGCGTAC	<i>yellow</i>	/5AmMC6/TAATGTGGCTTG GCTCACCCTAGG
<i>House</i>	<i>man-made</i>	/5AmMC6/ATGGATGCACGT CCACAAGAAGCAG	<i>building</i>	/5AmMC6/GTAGTTCGGGT GCATACTACCTGA	<i>wood</i>	/5AmMC6/GGGCGCAGAAGT CTCTATTCTAGAA
<i>Lincoln</i>	<i>human</i>	/5AmMC6/CATCGTAGGAAT GCGGCCGAGAATC	<i>19th century</i>	/5AmMC6/CGATGTAGTCAT CCCGATGTGCTGG	<i>president</i>	/5AmMC6/ATGGACGACTTG GGACGGGTATCAA
<i>Lion</i>	<i>cat</i>	/5AmMC6/AACGATTGTTAT GCCCCCTAACTCAG	<i>wild</i>	/5AmMC6/ACTCCGAGGAAC TTCGTGCTTAGTG	<i>yellow</i>	/5AmMC6/TAATGTGGCTTG GCTCACCCTAGG
<i>Sailboat</i>	<i>man-made</i>	/5AmMC6/ATGGATGCACGT CCACAAGAAGCAG	<i>water</i>	/5AmMC6/CTGGTTTGATCC GACACATTGATTC	<i>sails</i>	/5AmMC6/CTTACTTTCTTA CTCACTTCTCCAG
<i>Skyscraper</i>	<i>man-made</i>	/5AmMC6/ATGGATGCACGT CCACAAGAAGCAG	<i>building</i>	/5AmMC6/GTAGTTCGGGT GCATACTACCTGA	<i>steel</i>	/5AmMC6/TAGTGTGTGCCC ACTGTAGCCGTGA
<i>Tiger</i>	<i>cat</i>	/5AmMC6/AACGATTGTTAT GCCCCCTAACTCAG	<i>wild</i>	/5AmMC6/ACTCCGAGGAAC TTCGTGCTTAGTG	<i>orange</i>	/5AmMC6/CTGAATACTACA CGCCGTGGTGAAG
<i>Tree</i>	<i>plant</i>	/5AmMC6/TAAGCAATGGGT TCCACACTACGTA	<i>tree</i>	/5AmMC6/GATCAGAATCTA CTCGCATAGCCTC	<i>moon</i>	/5AmMC6/AGTTAAATGTCC CAGGCTTGTACC
<i>Washington</i>	<i>human</i>	/5AmMC6/CATCGTAGGAAT GCGGCCGAGAATC	<i>18th century</i>	/5AmMC6/GCAGTAAAGCTC GGTCCGATCTTCA	<i>president</i>	/5AmMC6/ATGGACGACTTG GGACGGGTATCAA
<i>Wolf</i>	<i>dog</i>	/5AmMC6/GAGTATCCGTTT GATTGTGCTCGC	<i>wild</i>	/5AmMC6/AGTCGTCCGAAA TATTGCATTCTTG	<i>black &amp; white</i>	/5AmMC6/TTCAGGGTGGAA GTACCTCCAGAT

**Chemical attachment of DNA barcodes on DNA capsules.** Using all the DNA capsules from the previous step, a mass of 5 mg of  $\beta$ -azido acetic acid *N*-hydroxysuccinimide ester and 5  $\mu$ L *N,N*-diisopropylethylamine as catalyst were added and the mixture was stirred for 2 hours, at 1000 rpm, and at 25 °C using a thermal mixer. The azide-modified DNA capsules were washed repeatedly for 5 times by re-dispersing the particles with 1 mL of *N*-methyl-2-pyrrolidone, sedimenting the azide-modified DNA capsules with a centrifuge at  $2,000 \times g$  for 3 minutes, and removing the supernatant. After the final wash, the azide-modified DNA capsules were re-dispersed in 1 mL of *N*-methyl-2-pyrrolidone. A mass of 2-mg of DBCO-PEG13-NHS ester was added and the mixture was stirred for 30 minutes, at 1,000 rpm, and at 25 °C using a thermal mixer. The particles were washed repeatedly for 5 times by re-dispersing the PEG-modified DNA capsules with 1 mL of *N*-methyl-2-pyrrolidone, sedimenting the PEG-modified DNA capsules with a centrifuge at  $2,000 \times g$  for 3 minutes, and removing the supernatant. After the final wash, the PEG-modified DNA capsules were re-dispersed in 200  $\mu$ L of *N*-methyl-2-pyrrolidone with 30-seconds of vortex mixing and 1 minute sonication (Cole Parmer; Vernon Hills, IL). A volume 10  $\mu$ L of each ssDNA barcode (500  $\mu$ M in nuclease-free water) and 200  $\mu$ L of PEG-modified DNA capsules were added to 770  $\mu$ L of 0.1 M bicarbonate buffer (pH 9.2) in a 1.5 mL Eppendorf LoBind tube. The mixture was stirred for 2 hours, at 1000 rpm, and at 25 °C using a thermal mixer to produce the final form of our data blocks or “files”. The files were washed repeatedly for 5 times by re-dispersing the particles with 1 mL of saline Tris-acetate-EDTA buffer with surfactants (40 mM Tris, 20 mM acetate, 2 mM EDTA, 1.0% Tween-20, 1.0% sodium dodecyl sulfate, and 500 mM NaCl), sedimenting the particles with a centrifuge at  $2,000 \times g$  for 3 minutes, and removing the supernatant. After the final wash, the particles were re-dispersed in 500  $\mu$ L of saline Tris-acetate-EDTA (40 mM Tris, 20 mM acetate, 2 mM EDTA, 1.0% Tween-20, 1.0% sodium dodecyl sulfate, and 500 mM NaCl) with 30 seconds of vortex mixing and 1 minute sonication. All the files were then pooled together which forms the file pool or molecular file database with an estimated final concentration of 2.0 mg mL<sup>-1</sup> in 10.0 mL of 40 mM Tris, 20 mM acetate, 2 mM EDTA, 1.0% Tween-20, 1.0% sodium dodecyl sulfate, and 500 mM NaCl.



## S5. Estimating plasmid copy numbers in DNA files

We took different volumes, 10  $\mu\text{L}$ , 50  $\mu\text{L}$ , 100  $\mu\text{L}$ , 150  $\mu\text{L}$ , and 200  $\mu\text{L}$ , from three different files that have particle concentrations at approximately  $1 \text{ mg mL}^{-1}$ . The particles were centrifuged for  $10,000 \times g$  for 1 minute and the supernatant was carefully removed using a pipette. The residue particles were dissolved using 45  $\mu\text{L}$  of 5:1 buffered oxide etch and incubating the mixture for 5 minutes. The mixture was vortexed for 5 seconds to re-suspend the pellet and the mixture was statically incubated at room temperature for 5 minutes. A volume of 5  $\mu\text{L}$  of 1 M phosphate buffer (0.75 M  $\text{Na}_2\text{HPO}_4$ ; 0.25 M  $\text{NaH}_2\text{PO}_4$ ; pH 7.5 at 0.1 M) was then added, vortexed for 1 second, and desalted twice through an Illustra MicroSpin S-200 HR column (GE Healthcare) before analyzing the samples through qPCR.



**Supplementary Fig. 5. Concentration and copy number of plasmid DNA for different memory files.** (A) *Cat2*, (B) *Sailboat*, (C) *Washington*. Solid red line denotes the best linear fit determined using Deming regression. Broken red lines are fit uncertainties.

**Supplementary Table 3. Fit results for Supplementary Fig. 5.**

File	Slope [moles of DNA recovered per mass of files]	Slope [copy number per file]
<i>Cat2</i>	$1.84 \times 10^{-11} \pm 1.04 \times 10^{-11}$	$1.46 \times 10^6 \pm 0.83 \times 10^6$
<i>Sailboat</i>	$1.26 \times 10^{-11} \pm 0.18 \times 10^{-11}$	$0.99 \times 10^6 \pm 0.14 \times 10^6$
<i>Washington</i>	$1.52 \times 10^{-11} \pm 0.17 \times 10^{-11}$	$1.20 \times 10^6 \pm 0.13 \times 10^6$

## S6. Query probes

An external table of probe sequences was generated associating each of these content descriptors to a DNA sequence database of reverse complements from the sequences displayed on the files and truncated to 15 nucleotides to maintain approximately 50 °C annealing temperatures using the IDT OligoAnalyzer tool (<https://www.idtdna.com/pages/tools/oligoanalyzer>). A 50 °C annealing temperature was chosen as the target temperature to easily de-hybridize the probe strands. A full-length 25-mer sequence would require annealing down from 95 °C and maintaining a temperature of 72 °C as designed for orthogonality and would introduce non-specific interactions at room temperature annealing, which would complicate sorting. Orthogonality of the truncated probe strands was tested computationally using the Nucleic Acids Package (NUPACK)<sup>4-6</sup>.

**Synthesis and purification of fluorophore-labelled single-stranded DNA query probes.** A volume of 200 µL of a 1 mM (200 nmol) solution of hexylamine-modified single-stranded DNA reconstituted in 100 mM sodium bicarbonate buffer (pH 9.2) and 20 µL of 50 mM stock solution of either TAMRA or AFDye 647 NHS ester in DMSO (10,000 nmol, 50 equivalents) were added sequentially in a 1.5-mL LoBind Eppendorf tube. The reaction mixture was mixed at 25 °C using a thermal mixer for 2 hours after which it was passed through an Illustra NAP-5 gel filtration column (GE Healthcare; Marlborough, MA) for desalting and removal of residual small molecules, such as NHS, unreacted dye NHS esters, or dye acids that were formed due to hydrolysis.

The desalted reaction mixture was further purified with ion-pairing, reverse-phase high performance liquid chromatography (IP-RP-HPLC) using a Waters (Milford, MA) Alliance HPLC e2695 system equipped with a Waters 2998 photodiode array detector, Waters Fraction Manager–Analytical, and an XBridge Oligonucleotide BEH C18 2.1 mm × 50 mm column with a particle size of 2.5 µm. The aqueous mobile phase for IP-RP-HPLC is composed of 0.1 M triethylammonium acetate in HPLC-grade water (Millipore Sigma) with pH 7.0 while the organic mobile phase is composed of 0.1 M triethylammonium acetate in 90:10 (w/w) acetonitrile and HPLC-grade water. A focused gradient was optimized and used for all purification methods (**Supplementary Table 4**). All purification runs were run at a flow rate of 1.0 mL min<sup>-1</sup>.

**Supplementary Table 4. Focused gradient table for IP-RP-HPLC purification of DNA-dye conjugates.**

Dye conjugate	Hold	Linear gradient	Hold	Clean-up linear gradient	Clean-up hold	Re-equilibration linear gradient	Re-equilibration hold
AFDye 647	Aqueous: 87%, Organic: 13% Time: 0–1 minute	Aqueous: 87% → 85% Organic: 13% → 15% Time: 1–16 minutes	Aqueous: 85% Organic: 15% Time: 16–17 minutes	Aqueous: 85% → 50% Organic: 50% → 50% Time: 17–18 minutes	Aqueous: 50% Organic: 50% Time: 18–19 minutes	Aqueous: 50% → 87% Organic: 50% → 13% Time: 19–20 minutes	Aqueous: 87% Organic: 13% Time: 20–23 minutes
TAMRA	Aqueous: 83%, Organic: 17% Time: 0–1 minute	Aqueous: 83% → 81% Organic: 17% → 19% Time: 1–16 minutes	None	Aqueous: 81% → 50% Organic: 19% → 50% Time: 16–17 minutes	Aqueous: 50% Organic: 50% Time: 17–18 minutes	Aqueous: 50% → 83% Organic: 50% → 17% Time: 18–19 minutes	Aqueous: 83% Organic: 17% Time: 19–23 minutes

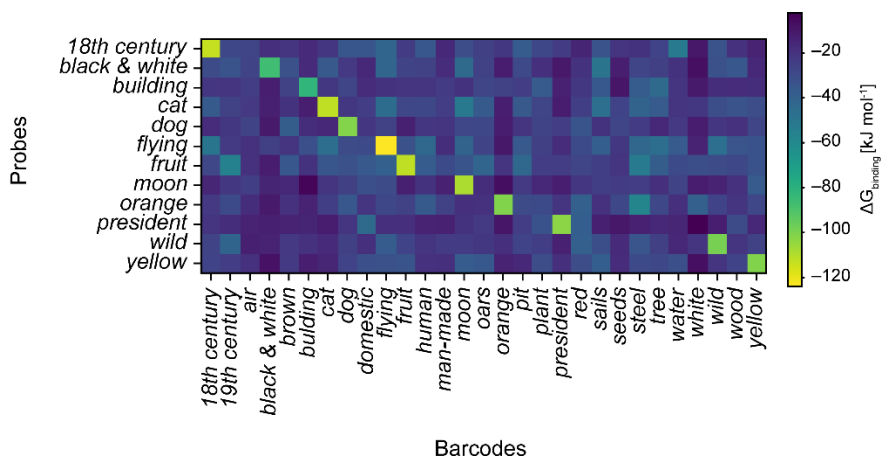
The column was heated and maintained at 50 °C using a column temperature controller and thermostat. Absorption intensities at 260 nm and 550 nm for TAMRA-DNA conjugates and 260 nm and 647 nm for AFDye 647-DNA conjugates were measured and used to automatically determine when to collect fractions. Collected purified fractions were dried to pellet form using a SpeedVac SPD300 (Thermo Fisher). The pelleted dye-modified DNA strands were reconstituted in 50 µL of HPLC-grade water and characterized with matrix-assisted laser desorption/ionization time-of-flight mass spectrometry using a Bruker (Billerica,

MA) microflex to confirm conjugation of the dye. Oligothymidine single-stranded oligonucleotide standards (Waters MassPREP OST) were used to calibrate the mass spectrometry measurements. The final concentrations of the dye-labelled single-stranded DNA probes were determined by measuring the absorbance spectra from 190–840 nm using a Nanodrop 2000 (Thermo Fisher) and using the molar absorption coefficients of the dyes at the maximum absorbance peak (AFDye 647: 270,000 M<sup>-1</sup> cm<sup>-1</sup>; TAMRA: 92,000 M<sup>-1</sup> cm<sup>-1</sup>) to calculate the concentrations.

**Supplementary Table 5. Sequences and mass characterization of dye-labelled query DNA probes.** MW = molecular weight.

Query barcode	Sequence	Expected MW (Da)	Measured MW (Da)
18 <sup>th</sup> century-AFDye 647	AFDye647-C6-TGAAGATCGGACCGA	5660.0	5660.2
cat-AFDye 647	AFDye647-C6-CTGAGTTAGGGGCAT	5682.2	5688.1
flyng-AFDye 647	AFDye647-C6-CTGGGCCGAACGAGG	5677.2	5675.3
fruit-AFDye 647	AFDye647-C6-AGCGGCACAGATACC	5605.2	5607.9
wild-AFDye 647	AFDye647-C6-CACTAAGCACGAAGT	5604.2	5601.0
building-TAMRA	TAMRA-C6-TCAGGTAGTATGCAC	5183.6	5176.5
dog-TAMRA	TAMRA-C6-AACTCCATGTAACGA	5136.6	5134.7
president-TAMRA	TAMRA-C6-TTGATACCCGTCCCA	5079.5	5077.1
yellow-TAMRA	TAMRA-C6-CCTAGCGGTGAGCCA	5169.6	5167.1

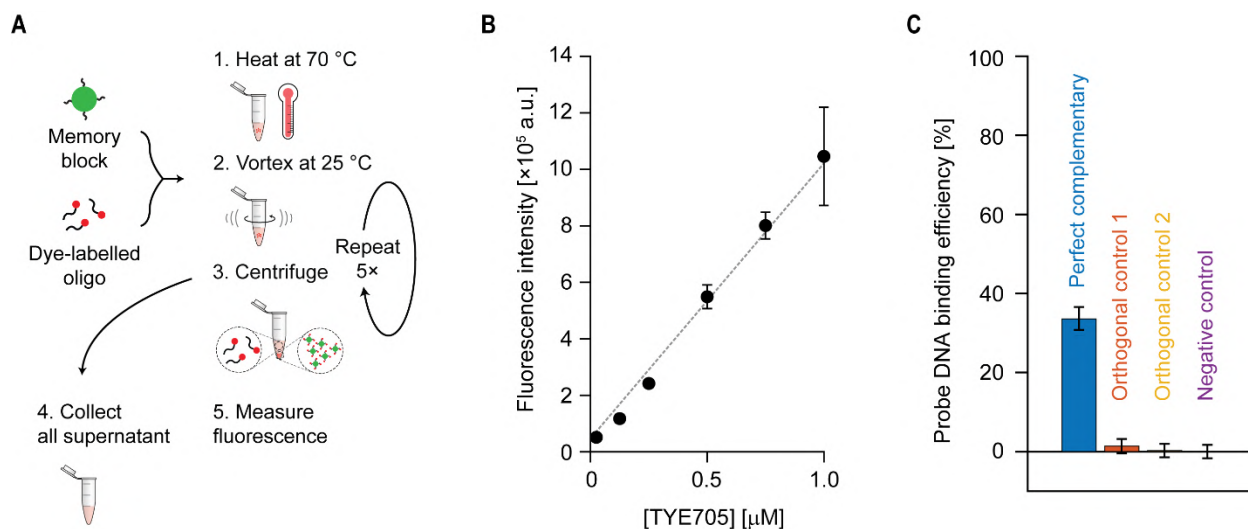
**Computational barcode validation.** The orthogonality of barcode and probe sequences was confirmed using NUPACK<sup>4-6</sup> to estimate the energetic favorability of binding between a single-stranded barcode and a single-stranded probe molecule in solution. These estimates do not account for the surface effects of having multiple barcodes densely packed on the surface of the silica bead. Using NUPACK's complexes function, the difference in energy between the double-stranded probe-barcode complex and the two single-stranded complexes at room temperature (25 °C) was calculated as the free energy of binding between the two DNA sequences. The binding energy for all probe-barcode pairs used during our experiments is shown in **Supplementary Fig. 6**. Overall, these results show that the binding affinity is much stronger between correct probe-barcode pairs as compared to incorrect probe-barcode pairs, which should be orthogonal to one another.



**Supplementary Figure 6. NUPACK estimates of binding affinity at 25 °C between barcode-probe pairs.** Affinities were estimated at 1M NaCl and without MgCl<sub>2</sub>. Only probes actually used for sorting were tested. Binding affinity is strongest for correct pairs, although some interactions between non-orthogonal pairs exist.

## S7. Estimating surface-accessible DNA barcodes using DNA hybridization assay

DNA hybridization assay was used to estimate the number of surface-accessible DNA barcodes<sup>7-9</sup>. In a typical experiment, we add 1  $\mu\text{L}$  of 500  $\mu\text{M}$  of TYE705-modified single-stranded DNA probe to a 1.5 mL Eppendorf LoBind tube that contains 50  $\mu\text{L}$  of 2 mg  $\text{mL}^{-1}$  of files that has surface-attached barcodes that are complementary to the TYE705-modified DNA probe sequence. Negative controls, with sample volumes of 50  $\mu\text{L}$  of 2 mg  $\text{mL}^{-1}$  in 1.5-mL Eppendorf LoBind tubes, were measured simultaneously with the test files. These negative controls either have surface-attached barcodes that are orthogonal to the TYE705-modified DNA probe sequence or hydroxy-terminated silica surface. Upon addition of the TYE705-modified DNA probe sequence into the file solution, the mixtures were mixed at 70  $^{\circ}\text{C}$  at 1,200 rpm using a thermal mixer for 5 minutes. The mixtures were cooled to 20  $^{\circ}\text{C}$  at 1,200 rpm using thermal mixer over 20 minutes and then centrifuged at  $10,000 \times g$  for 1 minute. The supernatant was collected, and the pelleted particles were washed with 80  $\mu\text{L}$  of 40 mM Tris, 20 mM acetate, 2 mM EDTA, 1.0% Tween-20, 1.0% sodium dodecyl sulfate, 500 mM NaCl. The sedimentation and washing process was repeated for five additional times while collecting the supernatant each cycle and pooling all the collected supernatant. A calibration curve using the TYE705-modified single-stranded DNA probe was used to determine the concentration of unhybridized TYE705-modified single-stranded DNA probe that remained in the supernatant solution (**Supplementary Fig. 7**).



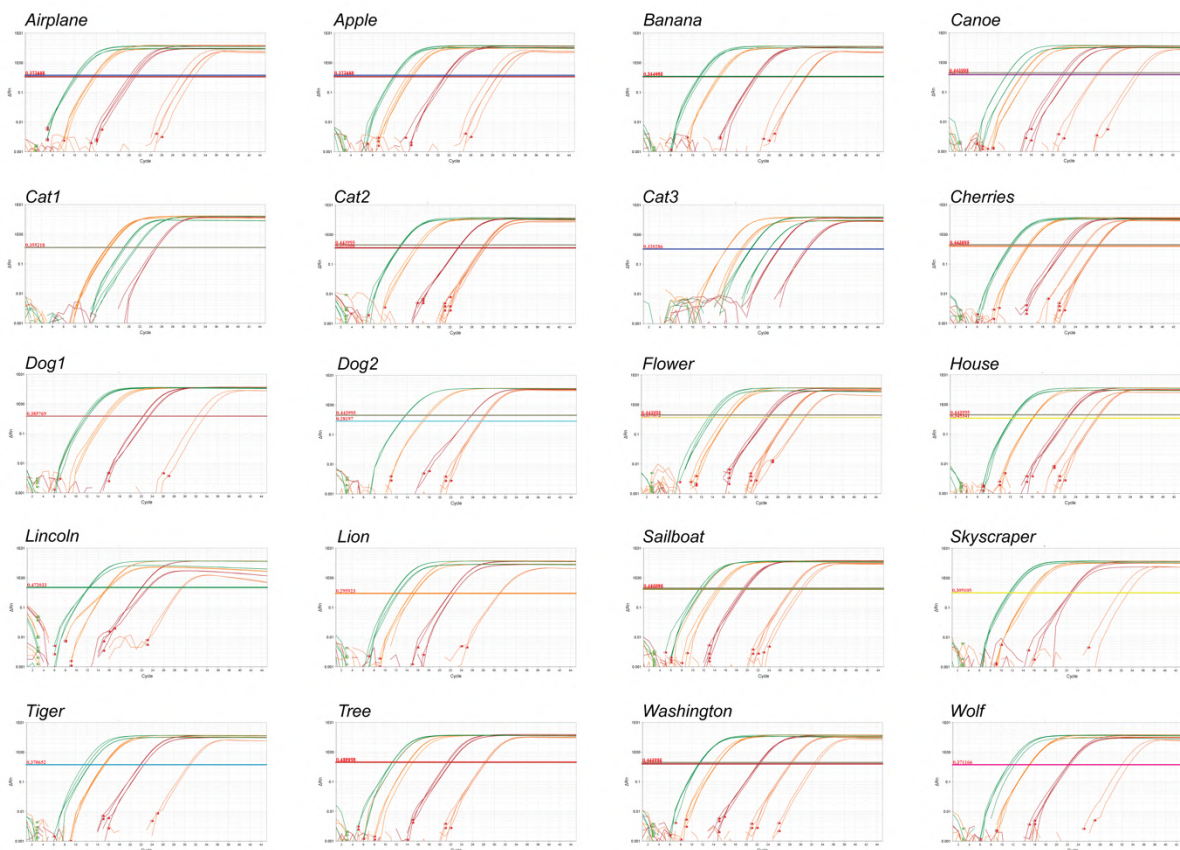
**Supplementary Figure 7. Hybridization assay to estimate the number of surface-accessible DNA barcodes.** (A) Sampling and analysis workflow. (B) Calibration curve used to determine the concentration of remaining unbound TYE705-modified DNA probes which encodes for the *black & white* barcode. Buffer for all dilutions: 40 mM Tris, 20 mM acetate, 2 mM EDTA, 1.0% Tween-20, 1.0% sodium dodecyl sulfate. (C) Probe DNA binding efficiency determined from the concentration of unbound TYE705-modified DNA probes that are left in the supernatant. Error bars are standard deviations from three independent replicates. Perfect complementary: File (*Dog2*) that has a *black & white* barcode. Orthogonal control 1: *Cat1*. Orthogonal control 2: *Airplane*. Negative control: silica core with hydroxy-terminated surface.

We used the density of silica and mass of silica that was sampled to estimate the number of silica particles in solution, and subtract the difference of initial TYE705-DNA concentration and the concentration of TYE705-DNA that remained in the supernatant to determine the concentration of hybridized TYE705-DNA. The number of TYE705-DNA that were hybridized onto the file surface can be calculated by taking the product of the volume of the supernatant, the concentration of the hybridized TYE705-DNA, and Avogadro's number ( $6.022 \times 10^{23}$  objects  $\text{mole}^{-1}$ ). Assuming that the hybridization efficiency is unity, the ratio of the number of hybridized TYE705-DNA and number of particles in solution provides the estimate of the number of surface-accessible barcodes per particle. The calculation is outlined as follows:

<b>Number of silica particles</b>	<p>Concentration of silica particles (<math>C_{\text{silica}}</math>) = 2.0 mg mL<sup>-1</sup>  Volume sampled (<math>V_{\text{silica}}</math>) = 0.050 mL  Mass of silica particles in solution (<math>M_{\text{silica}}</math>) = <math>C_{\text{silica}} \times V_{\text{silica}} = 1 \times 10^{-4}</math> g</p> <p>Density of a silica particle (<math>D_{\text{particle}}</math>) = 2.0 g mL<sup>-1</sup>  Volume of a silica particle (<math>V_{\text{particle}}</math>) = <math>4 \times \pi \times (2.5 \times 10^4 \text{ cm})^3 / 3 = 6.5 \times 10^{-11} \text{ cm}^3</math>  Mass of a silica particle (<math>M_{\text{particle}}</math>) = <math>D_{\text{particle}} \times V_{\text{particle}} = 1.3 \times 10^{-10}</math> g</p> <p>Number of silica particles (<math>N_{\text{particle}}</math>) = <math>M_{\text{silica}} / M_{\text{particle}} = 7.6 \times 10^5</math> particles</p>
<b>Number of hybridized TYE705-DNA</b>	<p>Initial concentration of TYE705-DNA (<math>C_{\text{initial}}</math>) = 500 <math>\mu</math>M  Volume of TYE705-DNA used (<math>V_{\text{used}}</math>) = 1 <math>\mu</math>L  Moles of TYE705 used (<math>m_{\text{initial}}</math>) = <math>C_{\text{initial}} \times V_{\text{used}} = 5 \times 10^{-10}</math> moles</p> <p>Measured concentration of TYE705-DNA from calibration curve (<math>C_{\text{supernatant}}</math>) = <math>0.740 \pm 0.05 \text{ } \mu\text{M}</math> (mean <math>\pm</math> s.d., n = 8)  Total volume of supernatant (<math>V_{\text{supernatant}}</math>) = <math>450 \pm 11 \text{ } \mu\text{L}</math> (mean <math>\pm</math> s.d., n = 8)  Moles of TYE705-DNA in supernatant (<math>m_{\text{supernatant}}</math>) = <math>C_{\text{supernatant}} \times V_{\text{supernatant}} = 3.3 \pm 0.2 \times 10^{-10}</math> moles (propagated uncertainty: mean <math>m_{\text{supernatant}} \times \sqrt{(\text{s.d. } V_{\text{supernatant}} / \text{mean } V_{\text{supernatant}})^2 + (\text{s.d. } m_{\text{supernatant}} / \text{mean } m_{\text{supernatant}})^2}</math>)</p> <p>Moles of hybridized DNA (<math>m_{\text{hybridized}}</math>) = <math>m_{\text{initial}} - m_{\text{supernatant}} = 1.7 \pm 0.2 \times 10^{-10}</math> moles (propagated uncertainty: <math>\sqrt{(\text{s.d. } m_{\text{supernatant}})^2}</math>)</p> <p>Number of hybridized DNA (<math>n_{\text{hybridized}}</math>) = <math>m_{\text{hybridized}} \times 6.022 \times 10^{23}</math> objects mole<sup>-1</sup> = <math>1.0 \pm 0.1 \times 10^{14}</math> hybridized DNA (propagated uncertainty: mean <math>n_{\text{hybridized}} \times \sqrt{(\text{s.d. } m_{\text{hybridized}} / \text{mean } m_{\text{hybridized}})^2}</math>)</p>
<b>Number of surface-accessible DNA barcodes</b>	<p>Number of surface-accessible DNA barcodes per particle (<math>n_{\text{surface}}</math>) = <math>n_{\text{hybridized}} / N_{\text{particle}} = 1.3 \pm 0.1 \times 10^8</math> surface-accessible DNA barcodes per particle (propagated uncertainty: mean <math>n_{\text{surface}} \times \sqrt{(\text{s.d. } m_{\text{hybridized}} / \text{mean } m_{\text{hybridized}})^2}</math>)</p>

## S8. Sequencing analysis

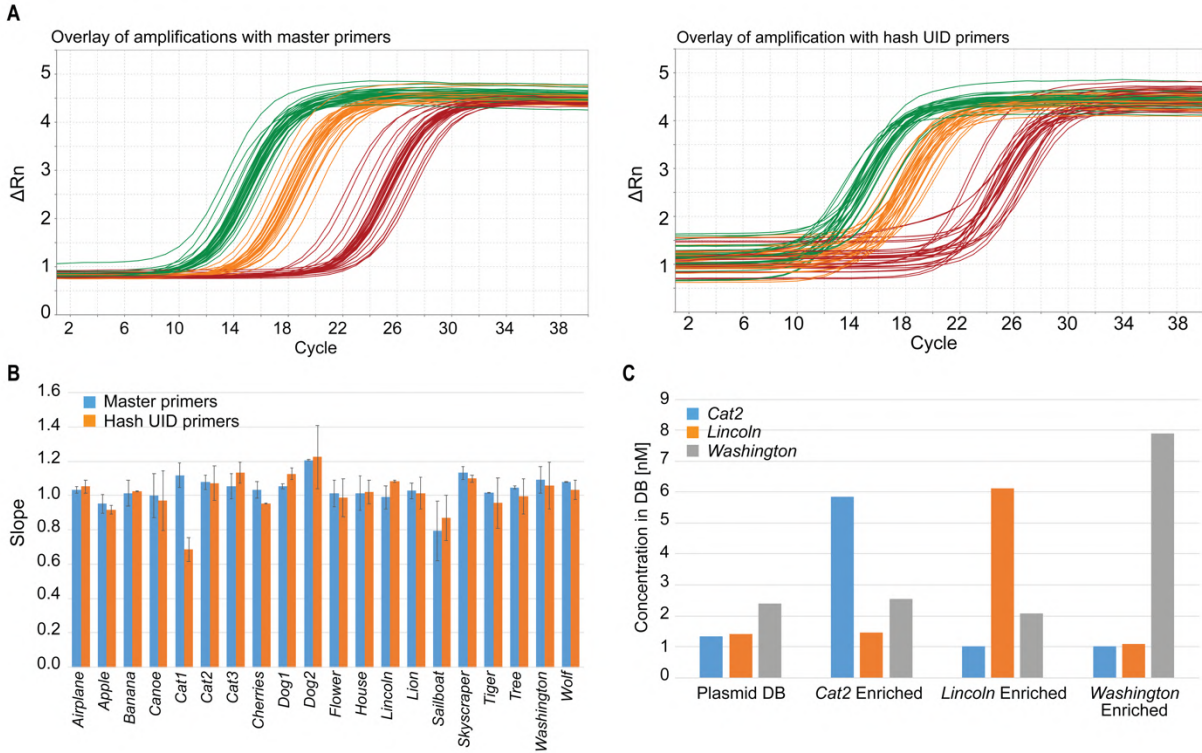
Three methods were used for verification of retrieval of DNA, which included quantitative PCR (qPCR), next-generation sequencing, and bacterial transformation and Sanger sequencing. For qPCR, standard curves were generated for each of the 20 plasmids using the 20 hash barcode primer pairs for 100 fg, 1 pg, and 10 pg of each, as judged by dilution from absorbance at 260 nm using the Nanodrop.



**Supplementary Figure 8. Plasmid standard curves.** PCR amplification curves ( $\Delta R_n$  vs. Cycle, log scale) shown for each plasmid DNA, for both master and hash UID primers, with 100 fg (red), 10 pg (orange), and 100 pg (green) of the indicated plasmid.

Gram weights were converted by molecular weight to moles to copy number, and the stand curve was generated as the log of the amounts vs. the threshold cycle. Fit curves were generated and the slopes are reported in **Supplementary Fig. 9**. Each file copy number approximately doubling per cycle. A simple plasmid mixture composed of equimolar amounts of *Cat2*, *Lincoln*, and *Washington*, and 5 $\times$  enriched of each of the 3 plasmids, with qPCR using the hash barcodes for amplification, showing capabilities of qPCR in isolating enriched populations.





**Supplementary Figure 9. Quantitative PCR analysis of plasmid.** (A) Overlay of 20 plasmids individually amplified with master primers (left) and hash UID primers (right), with 100 fg (red), 10 pg (orange), and 100 pg (green). (B) Calculated slopes from each standard curve from each master (blue) and hash UID (orange) primer amplification series, error bars are standard deviation of triplicate measurements. (C) Plasmid database with approximately equimolar concentrations, or with *Cat2*, *Lincoln*, or *Washington* 5× enriched, as shown, with qPCR with hash UID primers were used to quantify each of the three memory plasmids for each of the four databases.

For Illumina MiniSeq and MiSeq sequencing, the master primer pair with 5' extensions matching Illumina Nextera sequencing adapters were used to amplify all plasmids simultaneously (**Supplementary Fig. 9**). Template amounts were adjusted based on concentrations determined with Qubit fluorescence assay (Thermo Fisher) or qPCR. If required, the amplification was simultaneously followed by qPCR and enough cycles were used to rise above the Ct, or alternatively obtain a final concentration of 2 ng  $\mu\text{L}^{-1}$ . Dual sequence indices were then added to the adaptor-modified inserts at the 5' and 3' ends, associating the sequencing lane with a particular logic-gated pull, which was followed by SPRIbead cleanup. A 25  $\mu\text{L}$  PCR reaction amplified the material over 8–10 cycles using Kapa HiFi polymerase with 1 ng of template and 1  $\mu\text{M}$  forward and reverse primers. After amplification, this was combined with 20  $\mu\text{L}$  of SPRIselect beads, mixed, and let stand for 5 min. The mix was then separated by magnetic plates, and washed twice with 150  $\mu\text{L}$  80% ethanol, dried for 2 min, and eluted in 20  $\mu\text{L}$  Qiagen TE buffer. Samples were quantified using the Qubit fluorescence assay with the provided high-sensitivity buffer and standards. A sequencing pool was generated to approximately equimolar amount per index pair. Illumina MiniSeq with 150 × 150 read lengths was used to read out the start and end of each sequence. Sequences were demultiplexed, and sequence clustering was used to count the number of occurrences of each image.

**Supplementary Table 6. Primer sequences for sequence adaptor addition.** Master primer sequence is underlined.

Mem_R1	TCGTCGGCAGCGTCAGATGTGTATAAGAGACAG <u>CGTCGTCGCCCTCAA</u> ACT
Mem_R2	GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGGCAGATTGATGCCACCTTTTCAGC

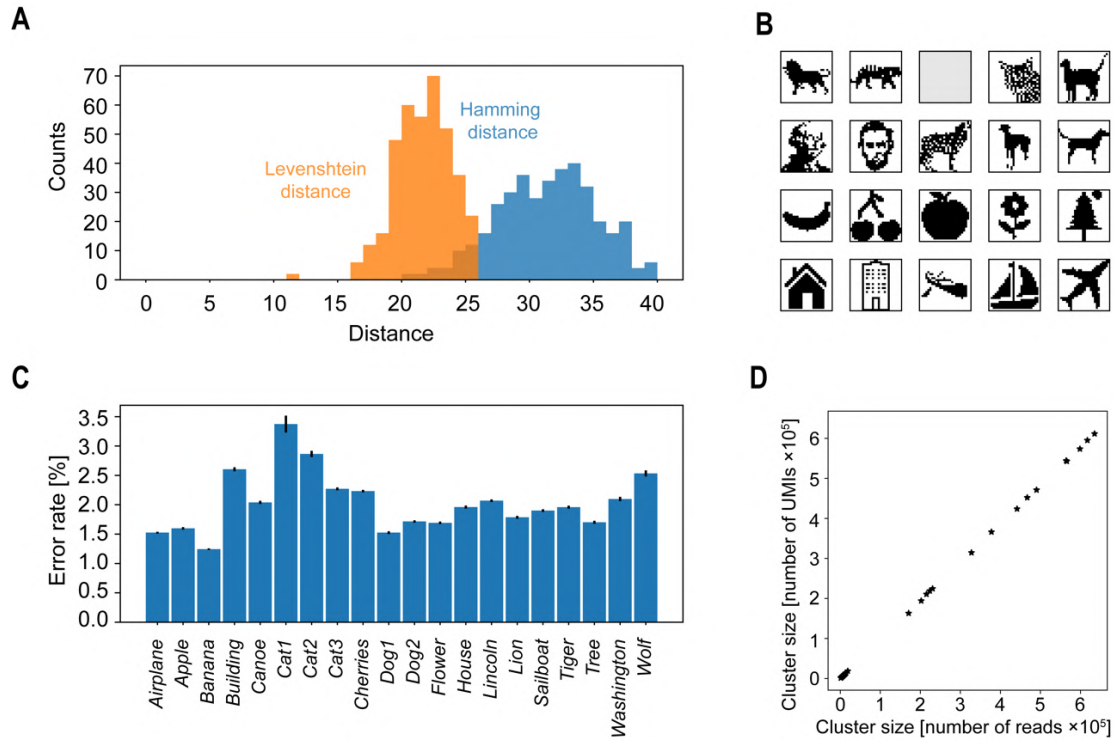
MiniSeq and MiSeq reads were clustered through a two-pass procedure that focused solely on the “hash” sequences flanking the image encoding internal regions of each sequence, because each hash sequence provided a unique identifier distinct between all sequences in our test set. UID’s were extracted from each read and clustering was subsequently performed with the following algorithm:

- 1) Create an empty list to store all observed clusters.
- 2) For each read do the following:
  - a. Check if the Hamming distance between the read’s UID and any cluster UID is at or below a predetermined threshold of 5. This distance threshold was determined as significantly below the Hamming distance between any pair of correct UIDs corresponding to one of the expected sequences (see **Supplementary Fig. 10**).
  - b. For any clusters satisfying this criterion, update the nucleotide counts observed at each position in the UID. Update the consensus UID of this cluster by a simple majority vote of the base identity at each position.
  - c. Periodically perform the following cluster clean-up check: If a cluster UID has a Hamming distance to one or more other clusters at or below the distance threshold, merge these clusters and update their nucleotide counts at each position and the consensus sequence. For all analyses, this was performed every 20,000 reads.
- 3) After processing all reads, remove any clusters that make up less than 0.02% of the total reads observed.
- 4) For each read, assign it to all clusters for which its UID satisfies the Hamming distance threshold. Some reads may not be assigned to any clusters or to multiple clusters, although in practice the latter occurrence was quite rare. Do not update any clusters during this step.
- 5) Each cluster’s UID was compared against the correct hash for each of the expected sequences and assigned to the one with the lowest Hamming distance, if this distance was less than the threshold. These assignments were used as the counts for each file in that sample.

As a control, clustering on some samples was also performed using an internal region of each sequence rather than the hash sequence, with minimal change to the results (data not shown). Code used to perform clustering is available on Github (<https://github.com/lcbb/DNA-Memory-Blocks/>). The sort probability for a file into a particular fraction was calculated as the count associated with that file divided by the sum of the counts for that file over all fractions generated from an initial sample. This metric is also referred to as *enrichment* throughout the text. In **Figs. 3–5** in the main text, the enrichment of each file is indicated by the percent opacity of the images displayed on the grid.

**Image reconstruction with Illumina MiSeq sequencing.** Reconstruction of the original images was carried out using Illumina MiSeq ran with  $300 \times 300$  read lengths, which span the plasmids sequences that encode images completely except for *Cat1*, which was 649 nucleotides in length. For this image, the clustering statistics were used although the image shown is not reconstructed from the sequencing results. To perform reconstruction, each pair of forward and reverse reads were aligned relative to each other to obtain a complete sequence for that read. From all sequences within each cluster, a consensus sequence was generated by determining the nucleotide identity that was most commonly observed at each nucleotide position. From this consensus sequence, the image was reconstructed by reversing the encoding process described in **Section S2**. With the exception of *Cat1*, all images were reconstructed successfully. Error rates were low enough that image reconstruction was almost always possible with as few as three reads. Results of the image reconstruction are shown in **Supplementary Fig. 10**.



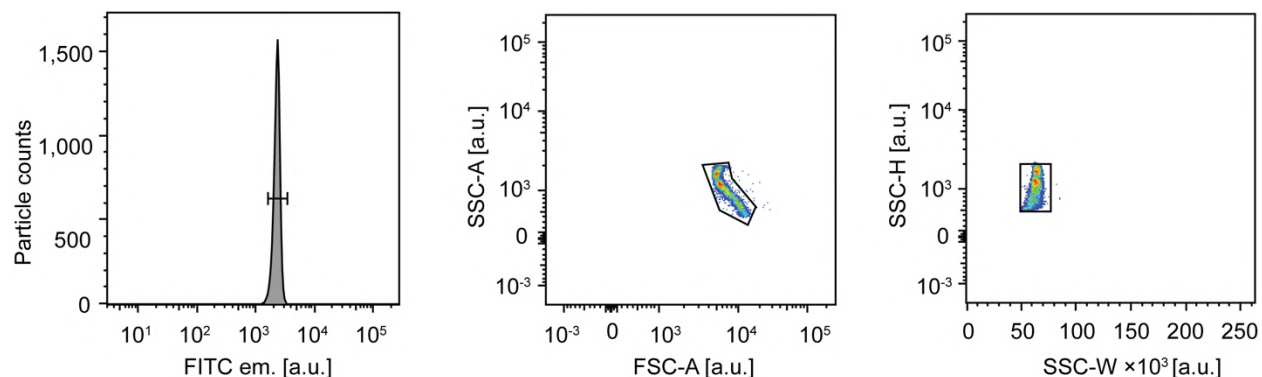


**Supplementary Figure 10.** (A) Hamming and Levenshtein distances between all pairs of hashes taken from file DNA sequences. Hamming distance was used rather than Levenshtein distance during clustering to reduce processing time. The minimum Hamming distance between any pair was 20, indicating that the distance threshold of 5 used during clustering is sufficient to avoid clustering correct sequences. (B) Image reconstructions from Illumina MiSeq  $300 \times 300$  sequencing, performed on a pool containing all files. Images were successfully reconstructed for all templates with the exception of *Cat1*, whose length (649 nucleotides) prevented full sequencing. (C) Sequencing error rates per base for each template. Error rates ranged from 1% to 3.5%, which is consistent with previous literature on sequencing of DNA de-encapsulated from silica particles<sup>10</sup>. Error bars on the error rates show standard errors of the mean. (D) To determine if PCR amplification bias could affect the relative counts of each file sequence during sequencing, universal molecular identifiers (UMIs) were added to some samples prior to PCR amplification. The size of each cluster was recalculated as the number of unique UMIs in that cluster. UMIs were 12-nt long random sequences added to the 3' and 5' ends of the file sequences, and two UMIs with a Levenshtein distance less than or equal to 1 were considered equivalent. The data show that the number of UMIs in each cluster was scaled linearly with the number of reads in that cluster, indicating that UMIs were not necessary for accurately measuring relative cluster size.

## S9. Fluorescence sorting of files

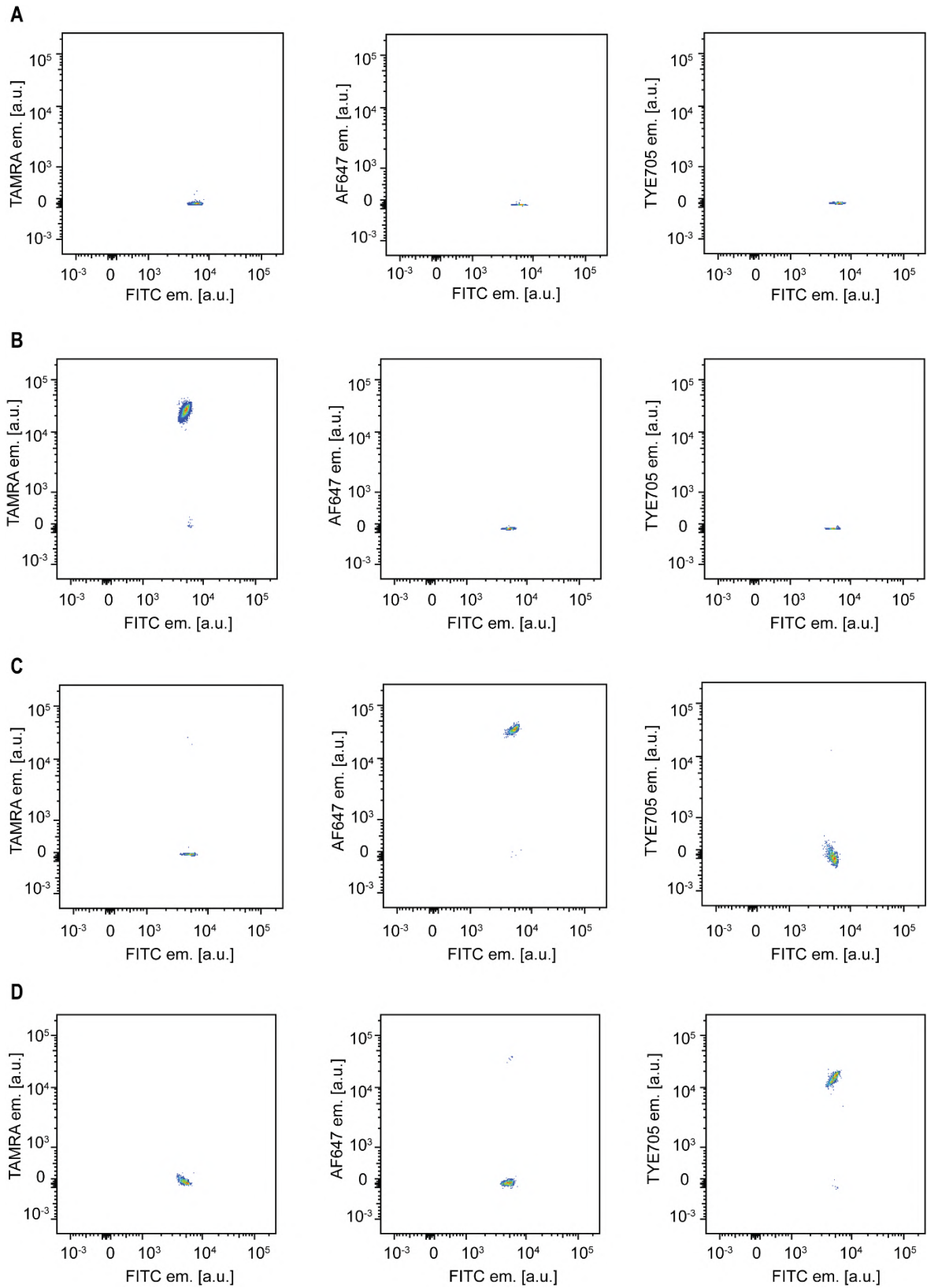
**Querying molecular file database using fluorescently-labelled probes.** The molecular file database was vortexed for 10 seconds, sonicated for two minutes, and re-vortexed for another 10 seconds to re-disperse the settled particles. A volume of 100- $\mu$ L of the molecular database ( $2 \text{ mg mL}^{-1}$ ) is added into a 1.5 mL Eppendorf LoBind tube. Dye-labelled probes for querying the molecular file database were added such that the final concentration of the DNA-dye single-stranded DNA in solution is  $5 \mu\text{M}$ . The resulting mixtures were mixed at  $70^\circ\text{C}$  at 1,200-rpm using a thermal mixer for 5 minutes. The mixtures were then cooled to  $20^\circ\text{C}$  at 1,200 rpm using a thermal mixer over 20 minutes and then centrifuged at  $10,000 \times g$  for 1 minute. The supernatant was discarded, and the pelleted particles were washed with 500- $\mu$ L of 40 mM Tris, 20 mM acetate, 2 mM EDTA, 1.0% Tween-20, 1.0% sodium dodecyl sulfate, 500 mM NaCl. The sedimentation and washing process was repeated for five additional times to remove non-specifically bound dye-DNA. The particles are finally re-suspended in 500  $\mu$ L of 40 mM Tris, 20 mM acetate, 2 mM EDTA, 1.0% Tween-20, 1.0% sodium dodecyl sulfate, 500 mM NaCl.

**Fluorescence-activated sorting.** All fluorescence-activated sorting (FAS) experiments were performed on a BD FACS Aria III flow cytometer. Samples were filtered through a Corning® 70- $\mu\text{m}$  cell strainer (Fisher Scientific) prior to particle sorts. Samples are flowed into the instrument with  $1\times$  PBS as sheath fluid at a flow rate that maintains an events detection rate of 1,200 events per second and below. We found that performing sorting at a flow rate that exceeds this events rate clogged the FAS instrument intermittently. All sorts were accomplished with a standard 70  $\mu\text{m}$  nozzle. The sample was held at room-temperature and agitated periodically every 5 minutes by stopping the sort and vortexing the sample vigorously with a vortex mixer. We note that the internal agitator in the flow cytometer with a 300 rpm agitation speed was not sufficient to prevent the silica particles from sedimenting over time and we found that periodically agitating the sample tube every 5 minutes with a vortex mixer was more effective. Since all files must contain a fluorescein core, all particles were gated by default using the 'FITC' laser and detector settings, which is defined by gating the majority population in the 'FITC-A' channel histogram, in addition to standard FSC and SSC gates to minimize sorting of doublets (**Supplementary Fig. 11**). All FAS experiments were performed at room-temperature.

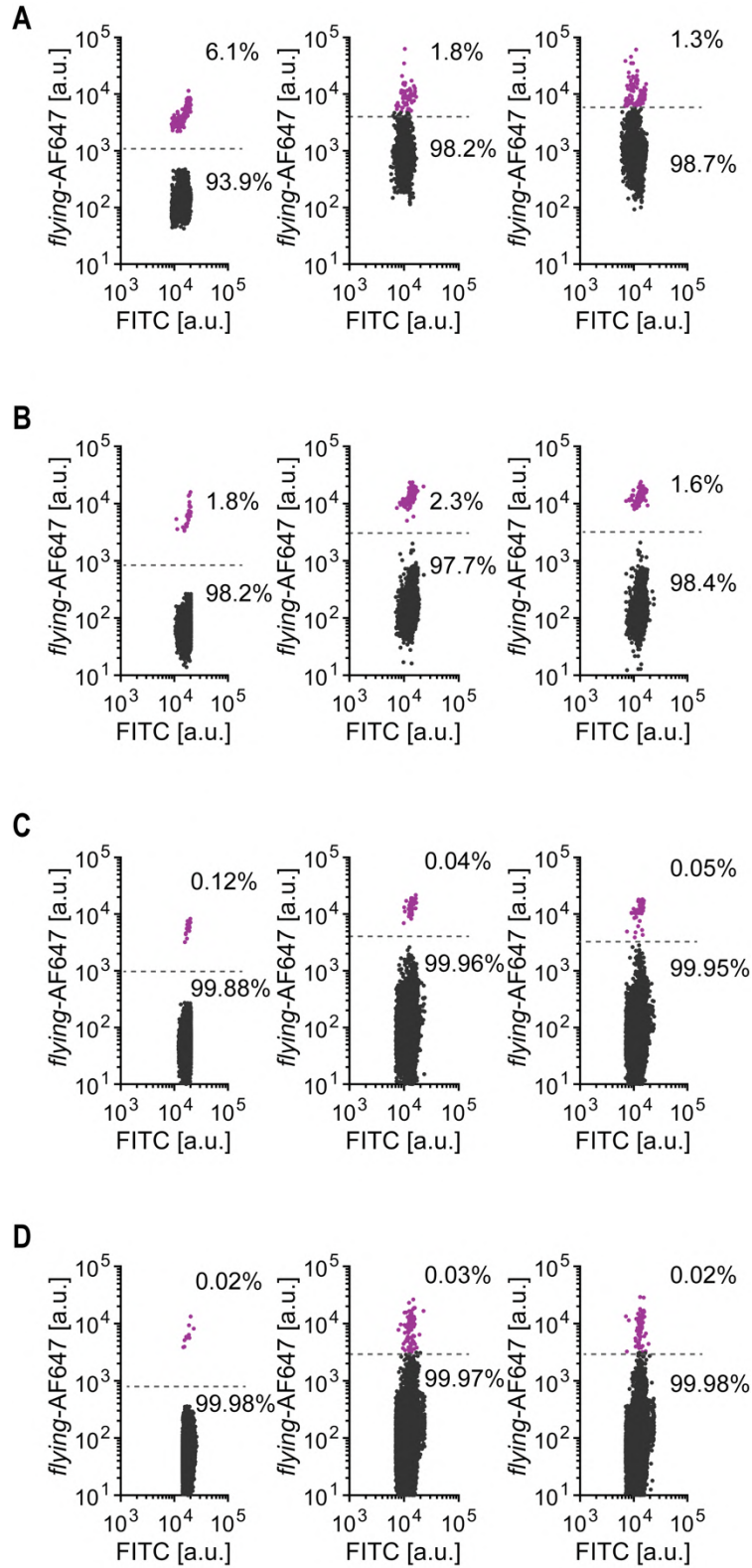


**Supplementary Figure 11.** Example of a standard set of gates used in all particle sorts. Colors indicate number density.

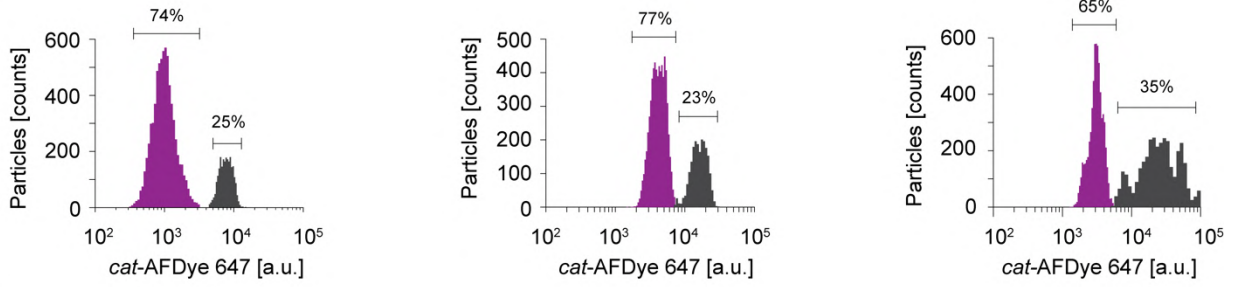
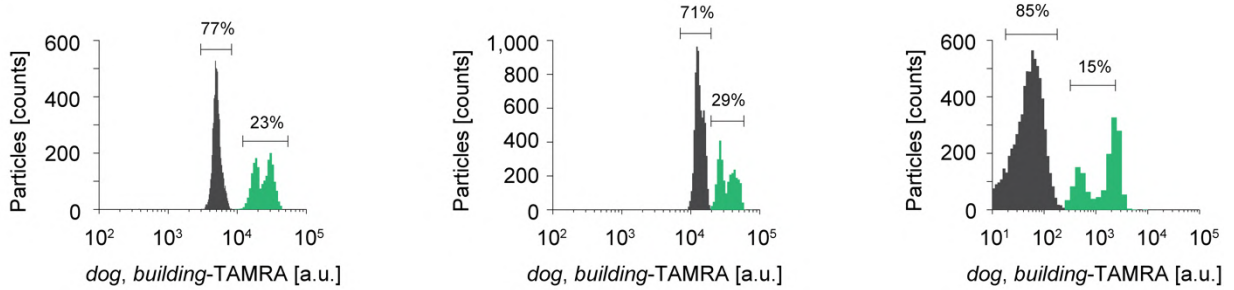
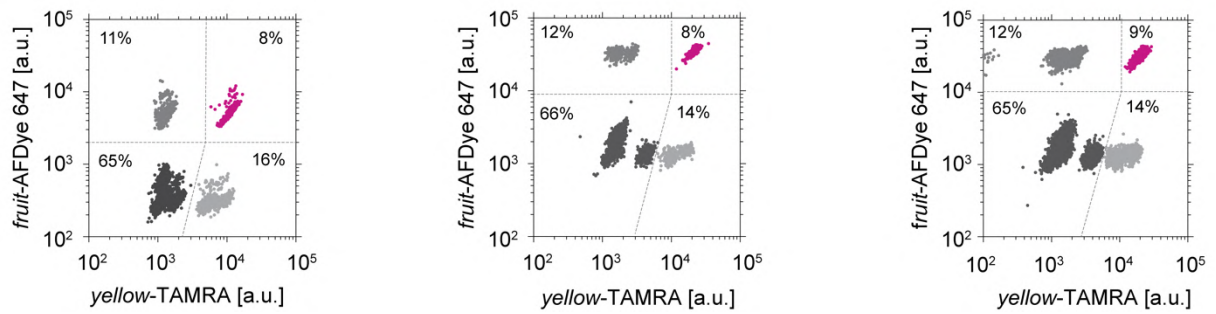
We also ran positive controls for every dye before every FAS experiment to validate that there is no significant spectral crosstalk during the sorting process and to validate that we have a distinguishable fluorescence signal in the presence of other fluorescent dyes. For example, because all our files have FITC, we validated that there is no fluorescence spillover of FITC in the TAMRA (PE-Texas Red channel), AF647 (APC channel), or TYE705 (Alexa Fluor 700 channel) that would otherwise make it difficult to distinguish the different particle populations. **Supplementary Fig. 12** summarizes the fluorescence traces of single and two-color controls.



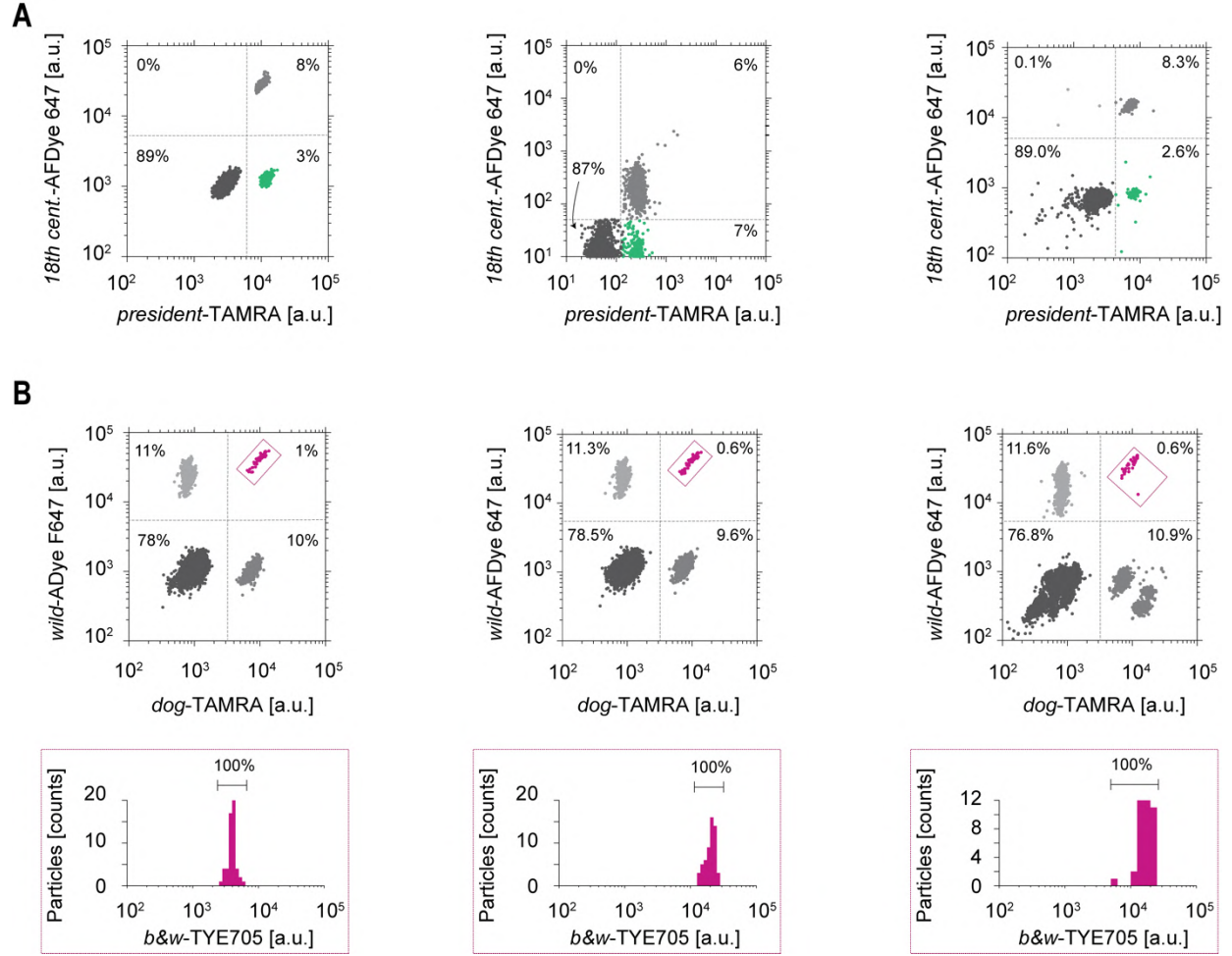
**Supplementary Figure 12. Single and two-color barcode positive controls.** (A) FITC only, (B) FITC + TAMRA, (C) FITC + AF647, and (D) FITC + TYE705. Colors indicate number density.



**Supplementary Figure 13. Raw data of FAS replicates from single barcode *flying* selections at different relative abundance of Airplane files compared to the other nineteen files. (A) 1:1 ratio of *Airplane* to each other file, (B) 1:10<sup>2</sup> ratio, (C) 1:10<sup>4</sup> ratio, (D) 1:10<sup>6</sup> ratio.**

**A****B****C**

**Supplementary Figure 14. Raw data of FAS replicates from Boolean logic selections. (A) NOT cat, (B) dog OR building, (C) fruit AND yellow.**

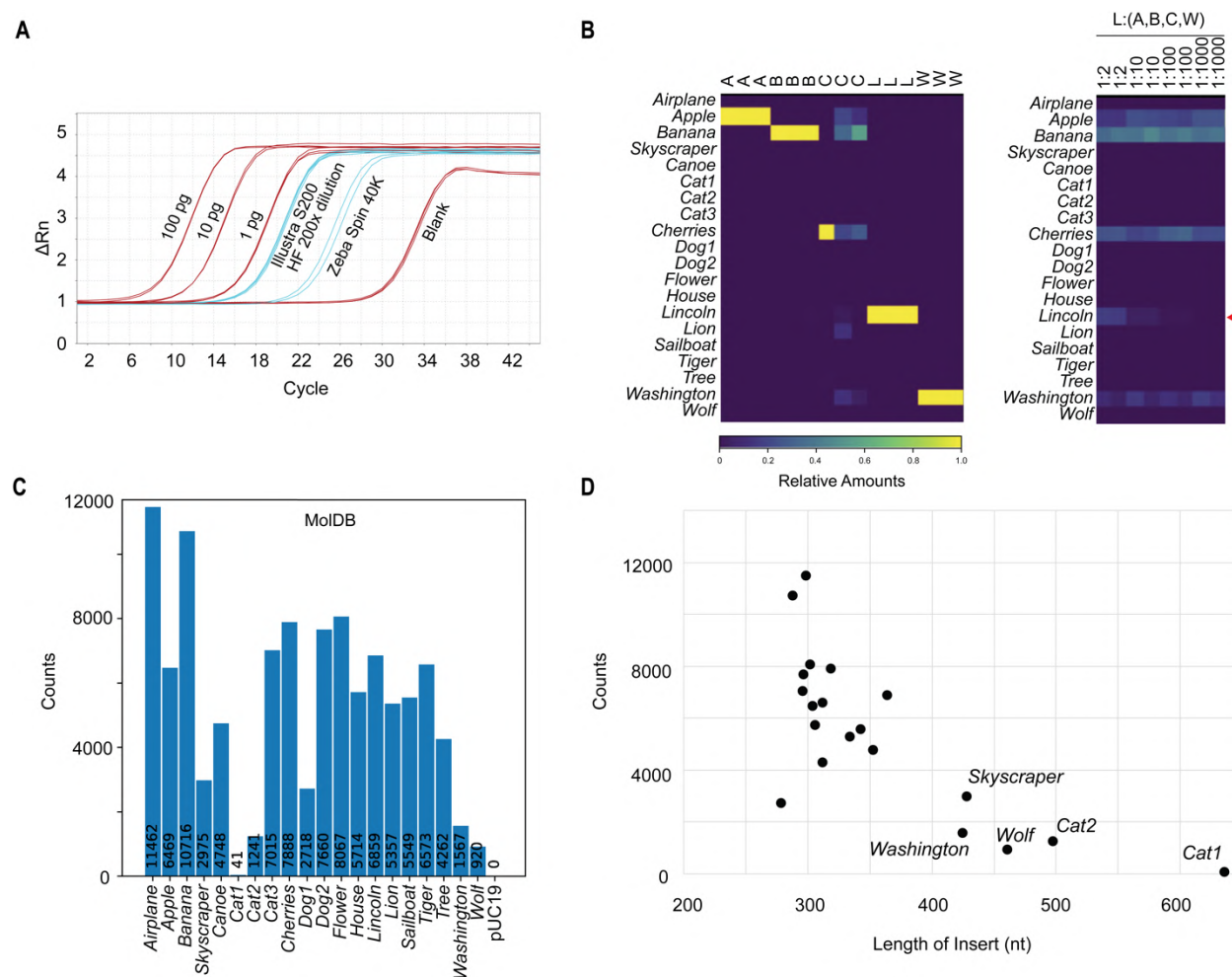


**Supplementary Figure 15. Raw data of FAS replicates from combined Boolean logic operations. (A) president AND (NOT 18<sup>th</sup> century). (B) dog AND wild.**

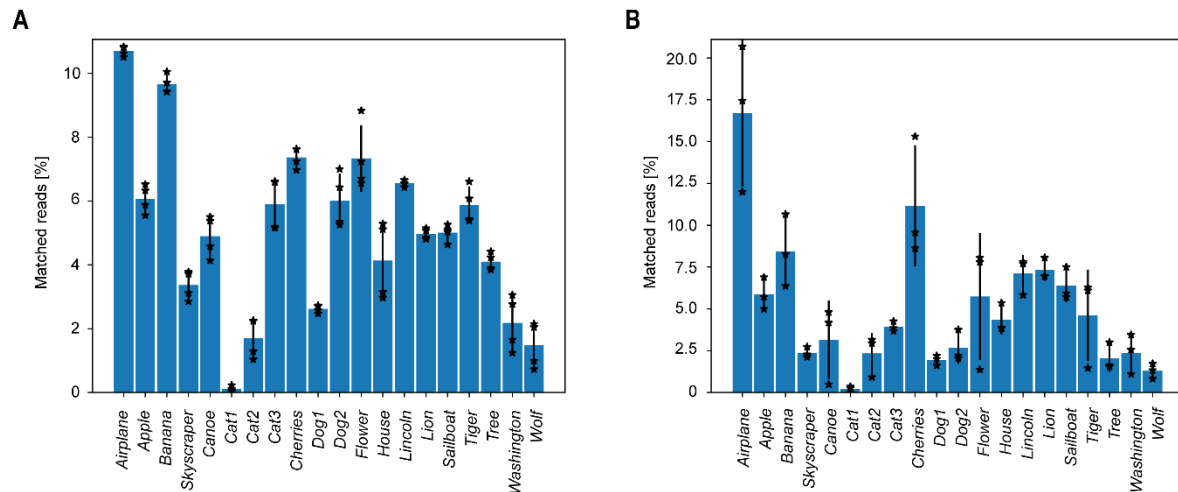


## S10. Verification of DNA retrieval from sorted sequences

**Release of DNA from sorted files.** Sorted populations were centrifuged at  $10,000 \times g$  for 1 minute. The supernatant was carefully removed with a pipette to avoid disturbing the silica pellets. A volume of 45  $\mu\text{L}$  of CMOS-grade 5:1 buffered oxide etch (Avantor; Visalia, CA) was then added. The mixture was vortexed for 5 seconds to re-suspend the pellet and the mixture was statically incubated at room temperature for 5 minutes. A volume of 5  $\mu\text{L}$  of 1 M phosphate buffer (0.75 M  $\text{Na}_2\text{HPO}_4$ ; 0.25 M  $\text{NaH}_2\text{PO}_4$ ; pH 7.5 at 0.1 M) was then added, vortexed for 1 second, and desalted twice through an Illustra MicroSpin S-200 HR column, which we found to be most effective compared to other clean-up methods.



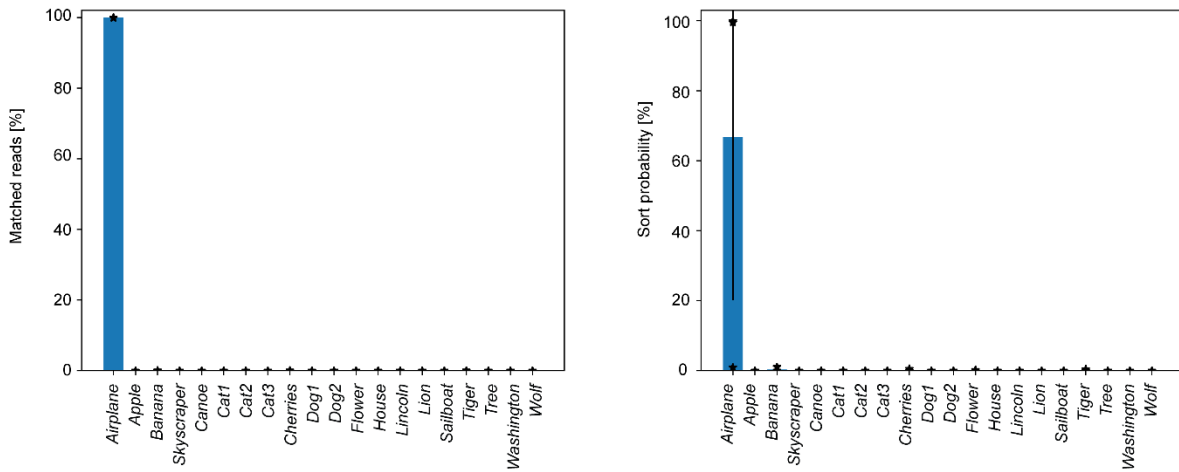
**Supplementary Figure 16. Cleanup and quality control on the release of DNA from the silica.** (A) Cleanup of release solution and additional salts from the release were tested by several methods of buffer exchange. Dilution of the released DNA mixture by 200-fold allowed for qPCR detection, as did cleanup by Illustra MicroSpin S-200 HR (GE Healthcare) and PCR Kleen Purification spin columns (Bio-Rad) with minimal sample loss. Zeba Spin 40K MWCO did not yield as much DNA. Illustra MicroSpin S-200 HR was used for all subsequent cleanup. (B) Release and purification of a subset of the molecular file database containing *Apple* (A), *Banana* (B), *Cherries* (C), *Lincoln* (L), and *Washington* (W) characterized from individual releases (left), and when part of a pool with *Lincoln* diluted (right), quantified by amplification, barcoding, and Illumina MiniSeq sequencing. (C) Release and purification of the entire molecular file database followed by amplification and barcoding, and sequencing from MiniSeq shows broadly similar profile of counts, with low detection of *Cat1*, and moderately low amplification of *Cat2*, *Dog1*, *Skyscraper*, *Washington*, and *Wolf*. (D) Low sequencing counts from the plasmid database can be explained by the variable length of the inserts being assayed. *Cat1*, *Cat2*, *Skyscraper*, *Washington*, *Wolf* are longer than 400 nucleotides.



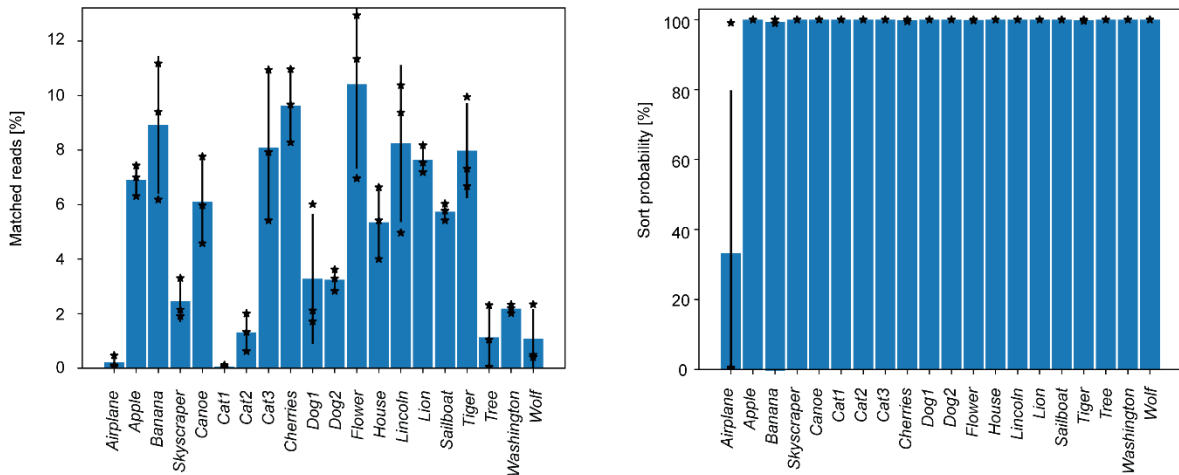
**Supplementary Figure 17. Count and sort probability statistics of sequencing reads from molecular file database. (A)** Molecular file database that did not pass through the FAS and were released directly. Mean and standard deviations were calculated from four independent replicates. **(B)** Molecular file database that was sorted into the FAS instrument using 'FITC' as the only sorting gate (all files contain FITC by default). Mean and standard deviations were calculated from three independent replicates. Matched reads are the number of reads matching each template divided by the number of reads matched to any template.



**A**

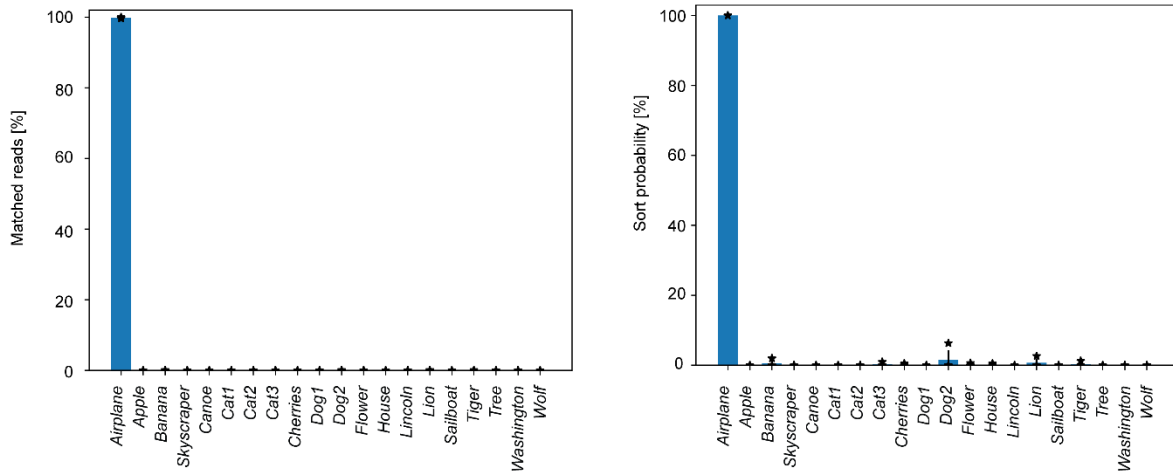


**B**

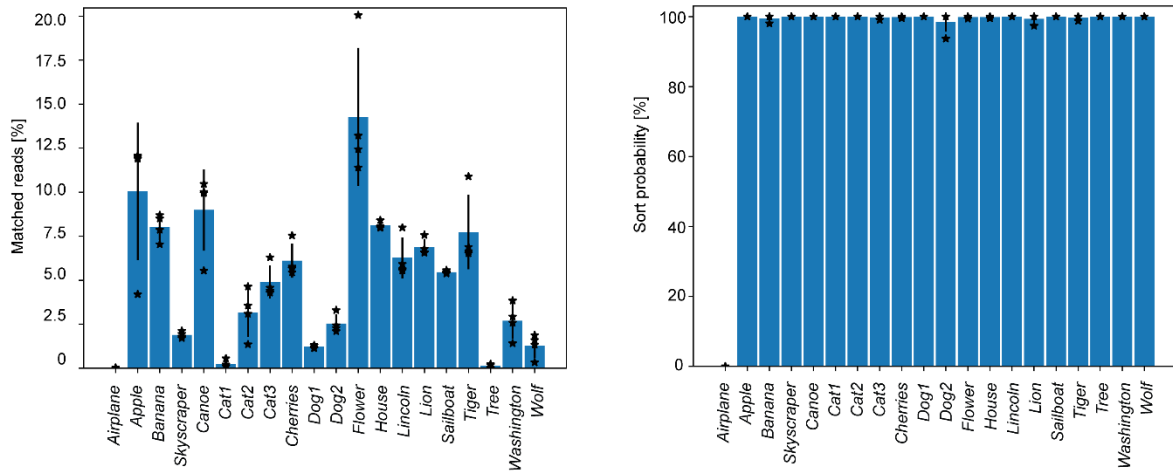


**Supplementary Figure 18. Count and sort probability statistics of sequencing reads from single-barcode sorts of *Airplane* in molecular database containing a 1:1 ratio of *Airplane* to each other file (i.e. equimolar concentration). (A) Sorted populations from the flying gate. Left: matched reads, Right: sort probabilities. (B) Sorted populations from the NOT flying gate. Left: matched reads, Right: sort probabilities. Matched reads are the number of reads matching each template divided by the number of reads matched to any template. Mean and standard deviations were calculated from four independent replicates.**

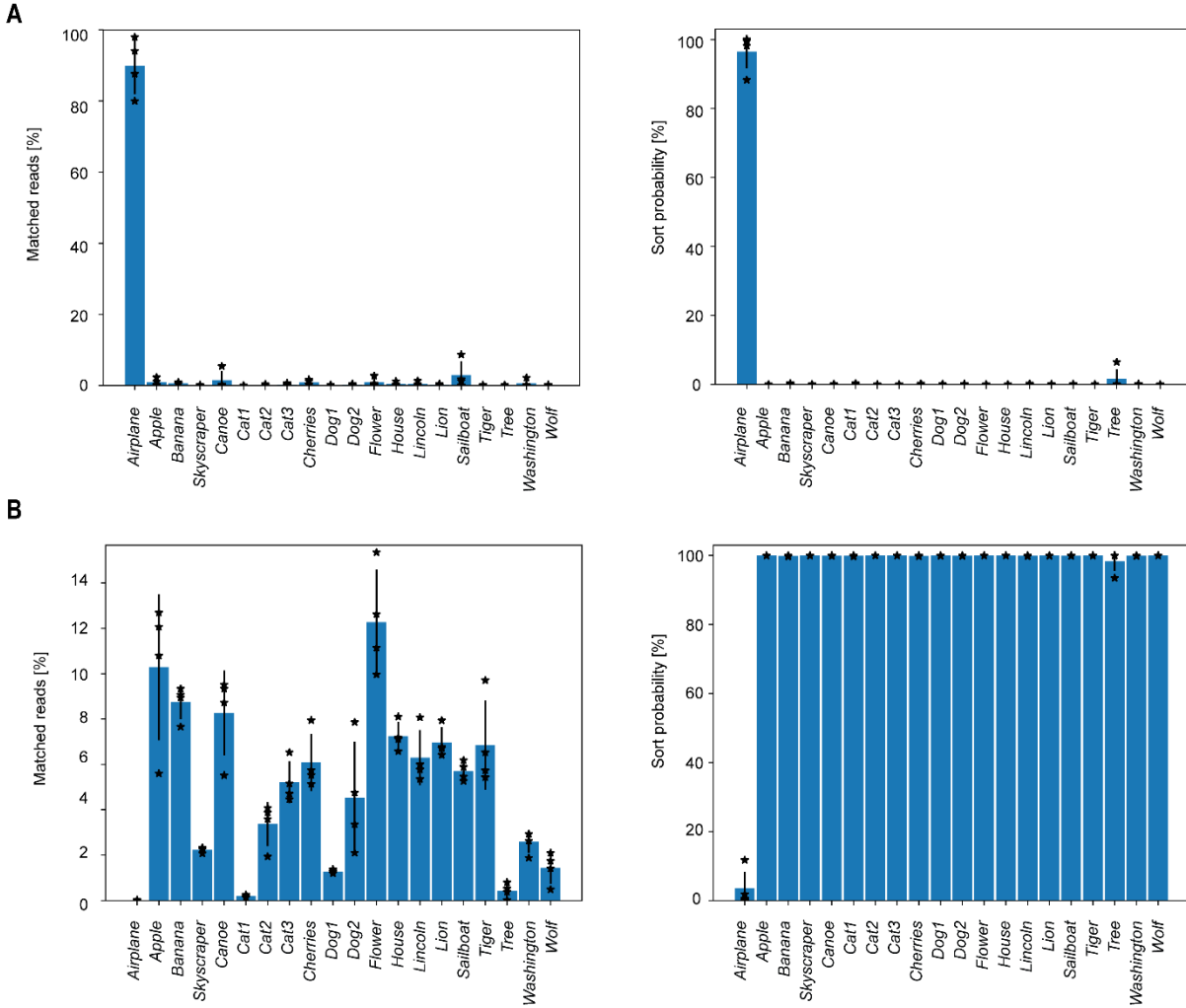
**A**



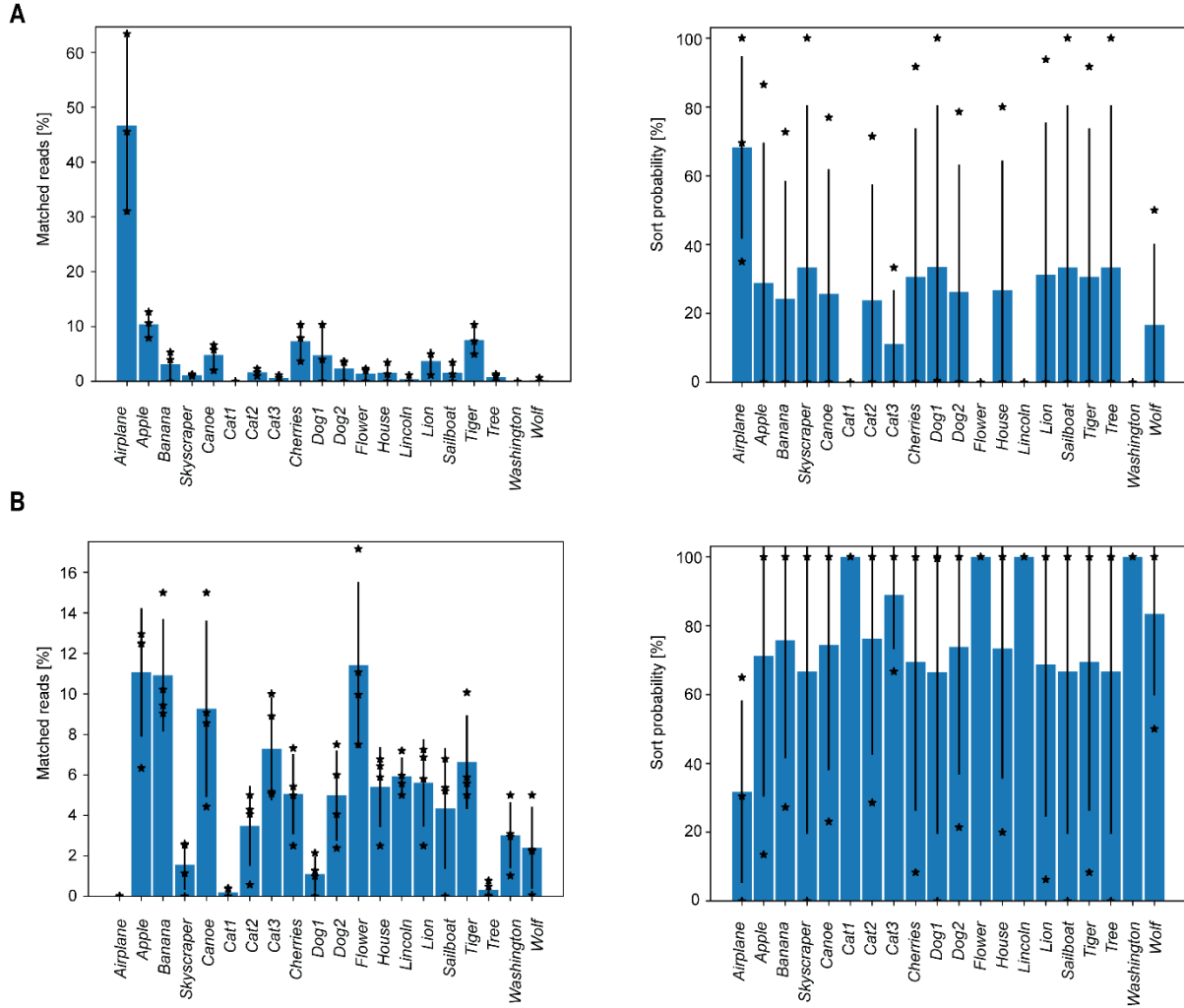
**B**



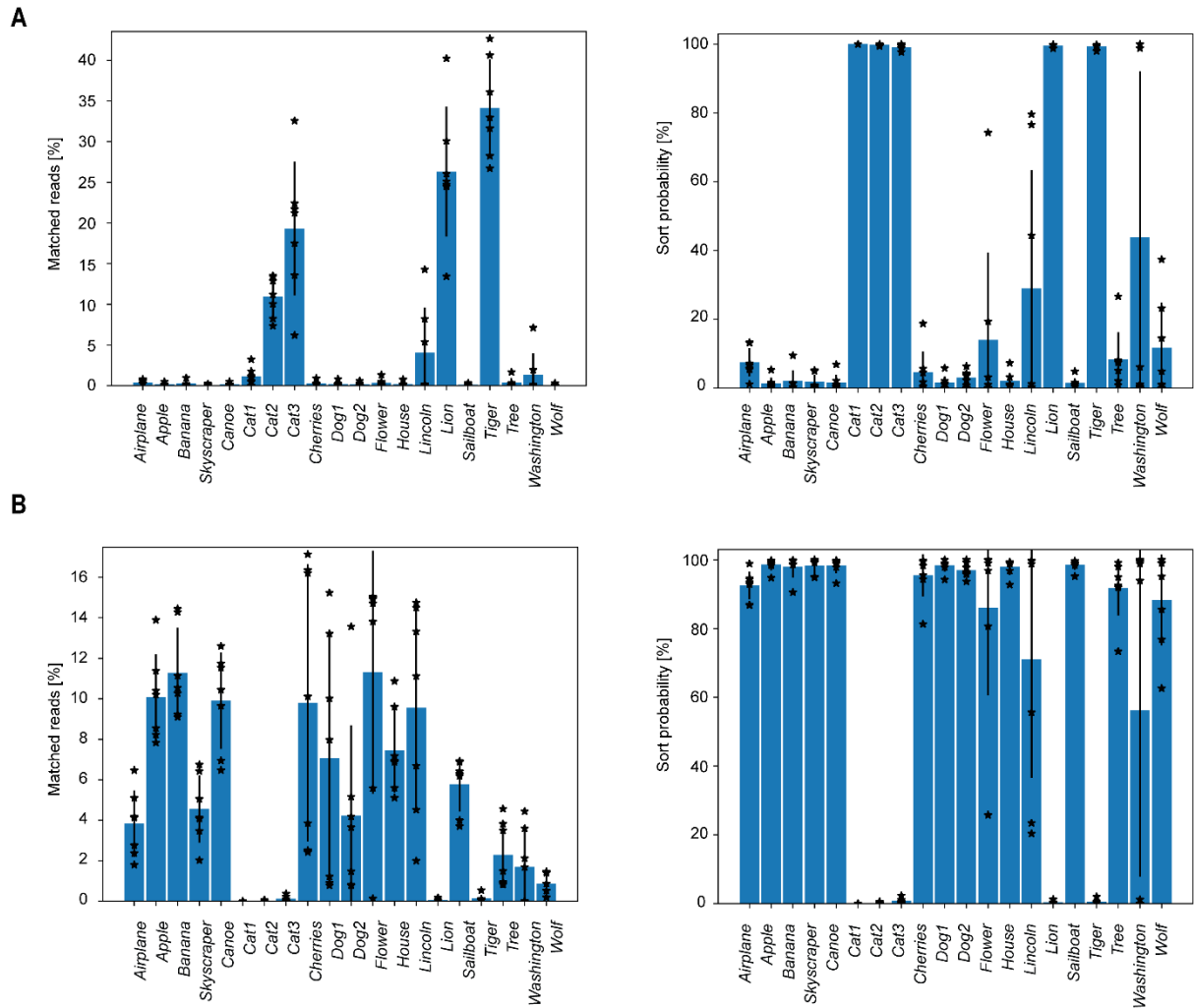
**Supplementary Figure 19. Count and sort probability statistics of sequencing reads from single-barcode sorts of *Airplane* in molecular database containing a 1:100 ratio of *Airplane* to each other file. (A) Sorted populations from the flying gate. Left: matched reads, Right: sort probabilities. (B) Sorted populations from the NOT flying gate. Left: matched reads, Right: sort probabilities. Matched reads are the number of reads matching each template divided by the number of reads matched to any template. Mean and standard deviations were calculated from four independent replicates.**



**Supplementary Figure 20. Count and sort probability statistics of sequencing reads from single-barcode sorts of *Airplane* in molecular database containing a 1:10,000 of *Airplane* to each other file. (A) Sorted populations from the flying gate. Left: matched reads, Right: sort probabilities. (B) Sorted populations from the NOT flying gate. Left: matched reads, Right: sort probabilities. Matched reads are the number of reads matching each template divided by the number of reads matched to any template. Mean and standard deviations were calculated from four independent replicates.**

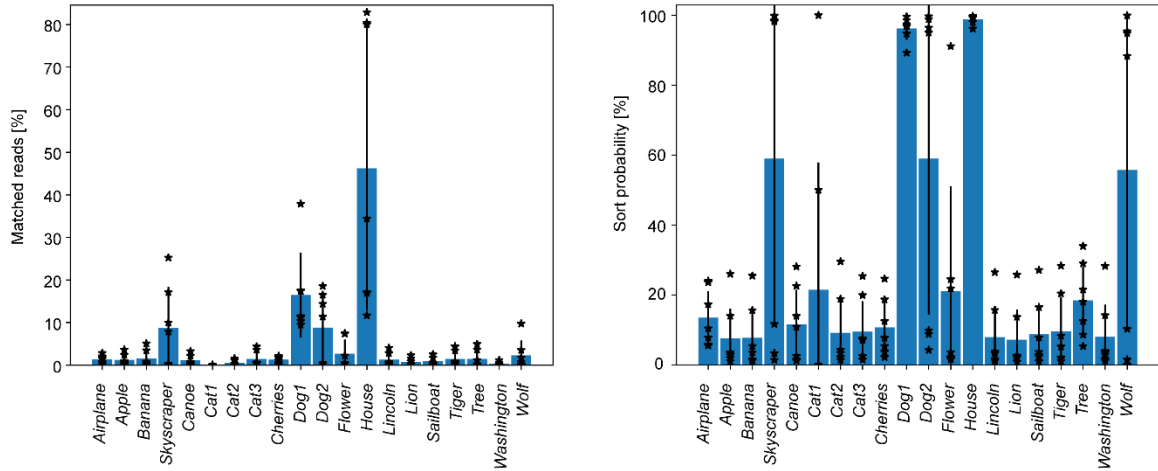


**Supplementary Figure 21. Count and sort probability statistics of sequencing reads from single-barcode sorts of *Airplane* in molecular database containing a 1:1,000,000 ratio of *Airplane* to each other file. (A) Sorted populations from the flying gate. Left: matched reads, Right: sort probabilities. (B) Sorted populations from the NOT flying gate. Left: matched reads, Right: sort probabilities. Matched reads are the number of reads matching each template divided by the number of reads matched to any template. Mean and standard deviations were calculated from three independent replicates.**

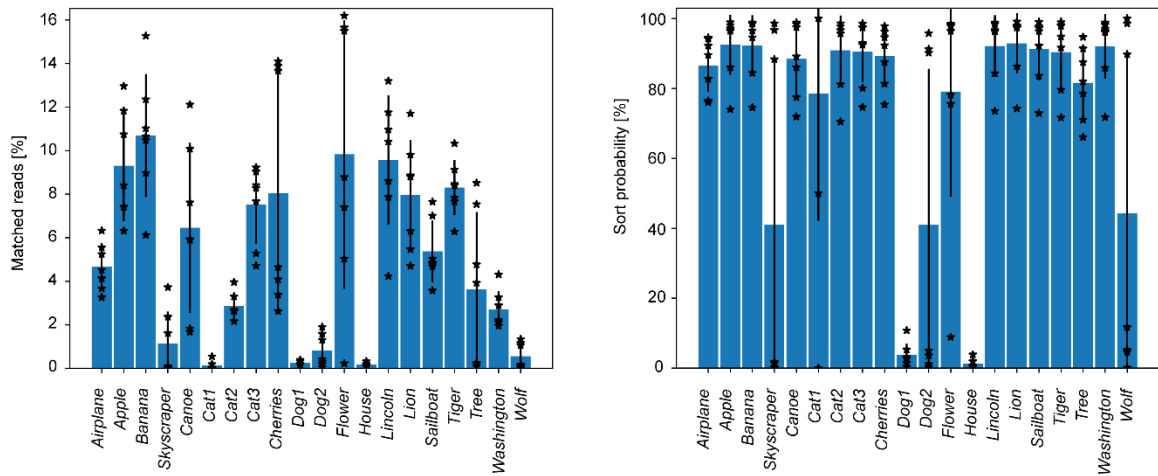


**Supplementary Figure 22. Count and sort probability statistics of sequencing reads from NOT cat sorts. (A)** Sorted populations from the cat gate. Left: raw counts, Right: sort probabilities. **(B)** Sorted populations from the NOT cat gate. Left: matched reads, Right: sort probabilities. Matched reads are the number of reads matching each template divided by the number of reads matched to any template. Mean and standard deviations were calculated from seven independent replicates.

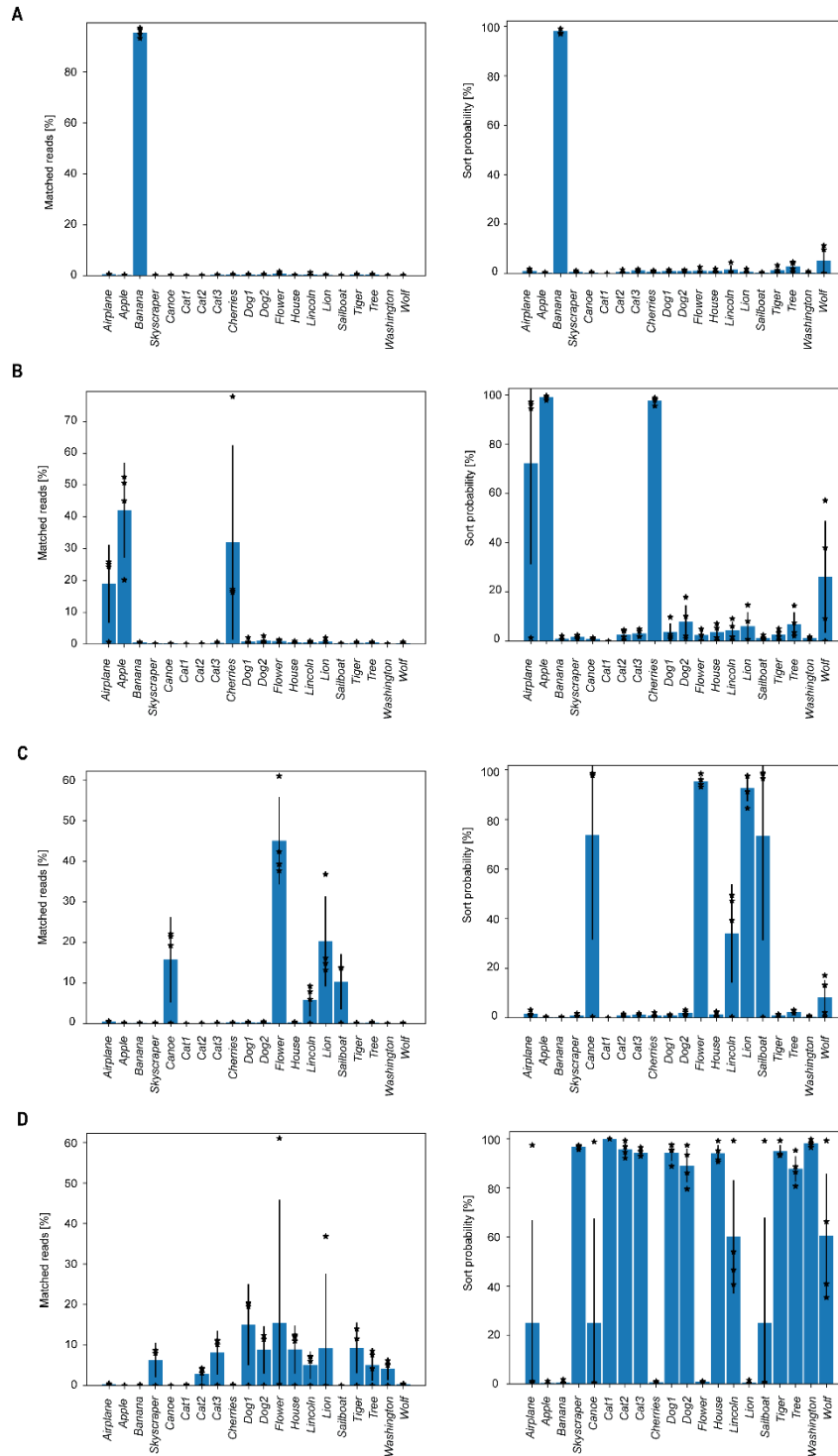
**A**



**B**

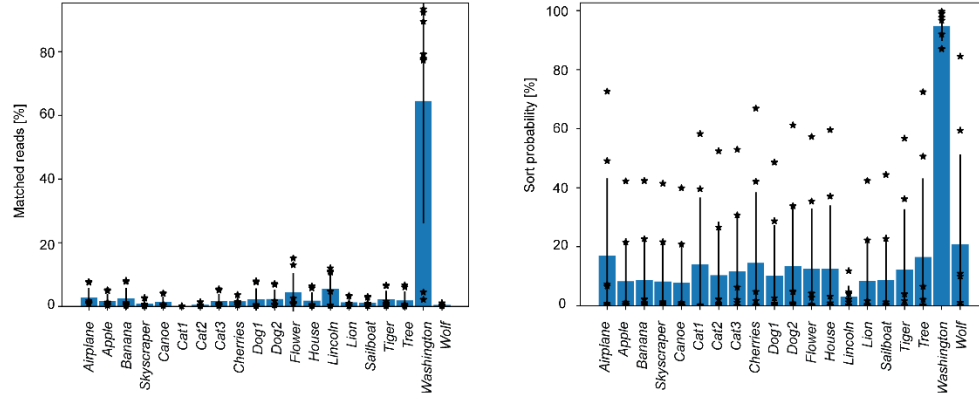


**Supplementary Figure 23. Count and sort probability statistics of sequencing reads from dog OR building sorts. (A)** Sorted populations from the dog OR building gate. Left: raw counts, Right: sort probabilities. **(B)** Sorted populations from the NOT (dog OR building) gate. Left: matched reads, Right: sort probabilities. Matched reads are the number of reads matching each template divided by the number of reads matched to any template. Mean and standard deviations were calculated from seven independent replicates.

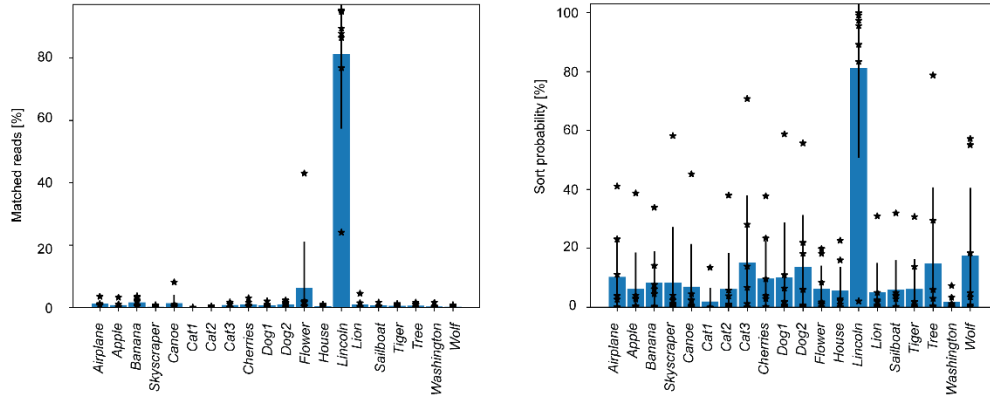


**Supplementary Figure 24. Count and sort probability statistics of sequencing reads from yellow AND fruit sorts.** (A) Sorted populations from the yellow AND fruit gate. Left: matched reads, Right: sort probabilities. (B) Sorted populations from the NOT yellow AND fruit gate. Left: matched reads, Right: sort probabilities. (C) Sorted populations from the yellow AND NOT fruit gate. Left: matched reads, Right: sort probabilities. (D) Sorted populations from the NOT fruit AND NOT yellow gate. Left: matched reads, Right: sort probabilities. Mean and standard deviations were calculated from four independent replicates. Matched reads are the number of reads matching each template divided by the number of reads matched to any template.

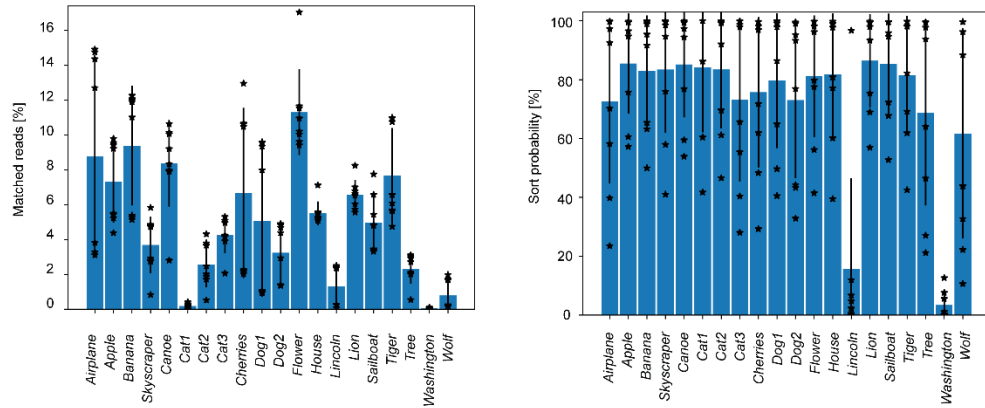
**A**



**B**

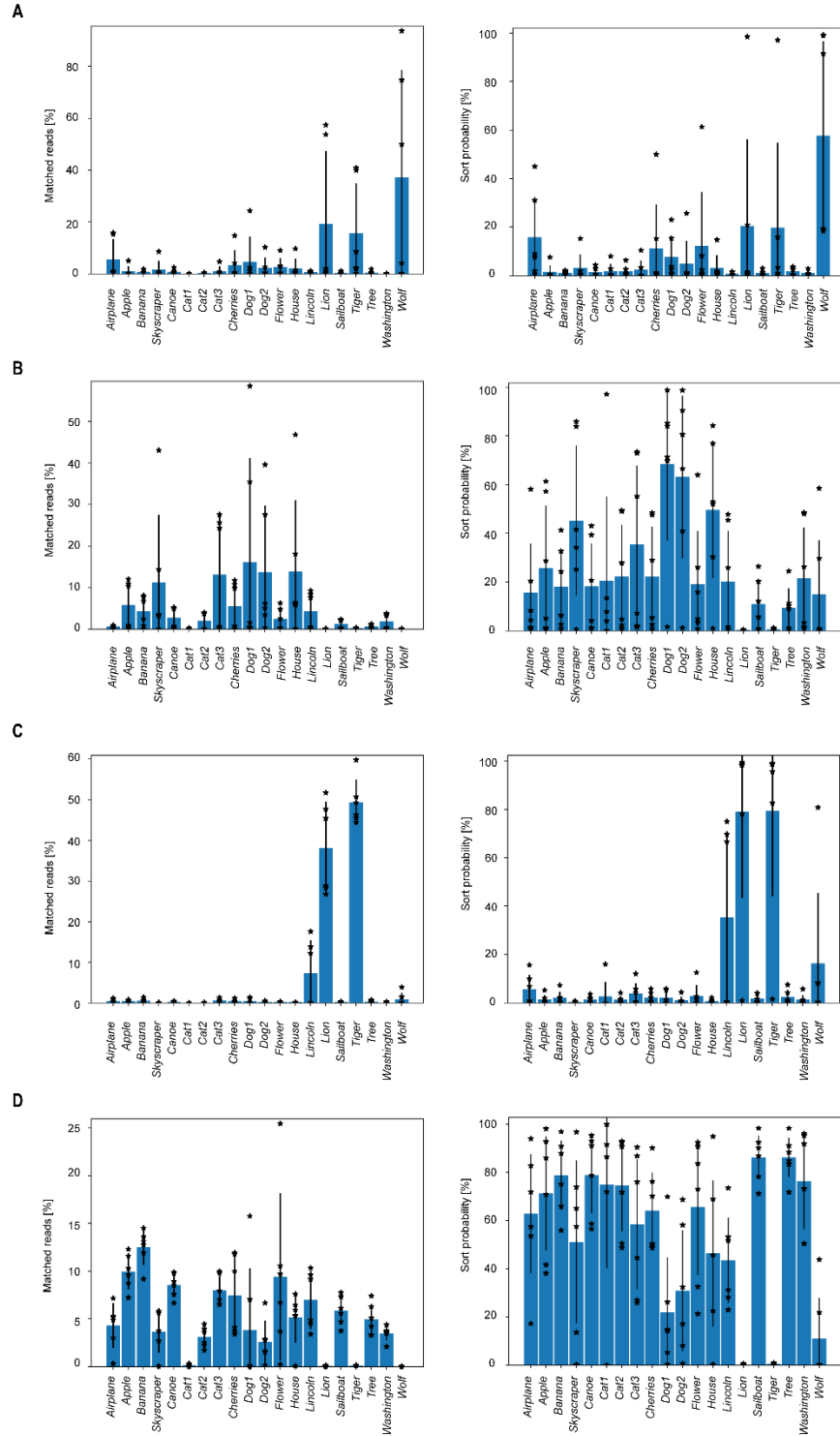


**C**



**Supplementary Figure 25. Count and sort probability statistics of sequencing reads from president AND 18th century sorts.** (A) Sorted populations from the president AND 18th century gate. Left: matched reads, Right: sort probabilities. (B) Sorted populations from the president AND NOT 18th century gate. Left: matched reads, Right: sort probabilities. (C) Sorted populations from the NOT president gate. Left: matched reads, Right: sort probabilities. Matched reads are the number of reads matching each template divided by the number of reads matched to any template. Mean and standard deviations were calculated from eight independent replicates.

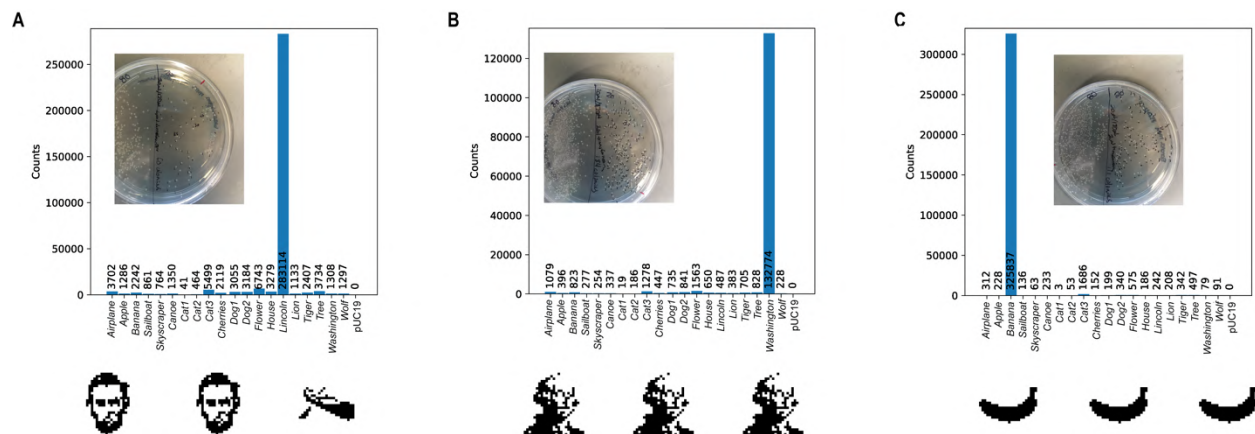




**Supplementary Figure 26. Count and sort probability statistics of sequencing reads from dog AND wild sorts.** (A) Sorted populations from the dog AND wild gate. Left: matched reads, Right: sort probabilities. (B) Sorted populations from the dog AND NOT wild gate. Left: matched reads, Right: sort probabilities. (C) Sorted populations from the NOT dog AND wild gate. Left: matched reads, Right: sort probabilities. (D) Sorted populations from the NOT dog AND NOT wild gate. Left: matched reads, Right: sort probabilities. Matched reads are the number of reads matching each template divided by the number of reads matched to any template. Mean and standard deviations were calculated from six independent replicates.

## S11. Bacterial transformation of sorted sequences

Samples that were sorted to single populations (yellow AND fruit: *Banana*; president AND 18th century: *Washington*; president AND NOT 18th century: *Lincoln*) were transformed to chemically competent *E. coli* 10 $\beta$  cells (NEB). Of the transformed cells, three colonies from each were grown in 4-mL LB overnight at 37 °C and Qiagen miniprep spin purification was used to retrieve the plasmid. Each of the three plasmids, as well as the PCR amplified release from each of the three populations, were sent for Sanger sequencing. All three amplicons showed primarily expected DNA sequences and each of the three images were retrieved from the sequencing of the sorts. Of the transformed colonies, 2 out of the 3 of the *Lincoln* colonies were verified to be *Lincoln* plasmids (one *Canoe* colony was also retrieved from this sample), while 3 out of the 3 *Banana* and *Washington* colonies were recovered. Sanger sequencing of the nine colonies showed no errors in results and all images were retrieved by inverse DNA to image processing. The bacterially amplified DNA was pure and readily available for re-encapsulation in a closed read-write cycle.



**Supplementary Figure 27. Bacterial transformation with sorted and cleaned DNA.** Cleanup of DNA release solution and additional salts away from the DNA allowed for transformation of NEB DH10 $\beta$  cells with the purified DNA. DNA sorted to single populations from (A) president AND NOT 18<sup>th</sup> century, (B) president AND 18<sup>th</sup> century, and (C) yellow AND fruit were transformed and grown to single colonies, and 3 colonies were selected and grown for DNA preparation and Sanger sequencing was applied to each. The expected Lincoln sort yielded two positive colonies and one colony encoding *Canoe* (A, bottom), with *Washington* (B, bottom) and *Banana* (C, bottom) sorts showed all three colonies returning the expected encoded image.

## S12. References

- 1 Goldman, N. *et al.* Towards practical, high-capacity, low-maintenance information storage in synthesized DNA. *Nature* **494**, 77-80, doi:10.1038/nature11875 (2013).
- 2 Paunescu, D., Puddu, M., Soellner, J. O. B., Stoessel, P. R. & Grass, R. N. Reversible DNA encapsulation in silica to produce ROS-resistant and heat-resistant synthetic DNA "fossils". *Nature Protocols* **8**, 2440, doi:10.1038/nprot.2013.154 (2013).
- 3 Xu, Q., Schlabach, M. R., Hannon, G. J. & Elledge, S. J. Design of 240,000 orthogonal 25mer DNA barcode probes. *Proceedings of the National Academy of Sciences* **106**, 2289-2294, doi:10.1073/pnas.0812506106 (2009).
- 4 Dirks, R. M. & Pierce, N. A. A partition function algorithm for nucleic acid secondary structure including pseudoknots. *Journal of Computational Chemistry* **24**, 1664-1677, doi:10.1002/jcc.10296 (2003).
- 5 Dirks, R. M. & Pierce, N. A. An algorithm for computing nucleic acid base-pairing probabilities including pseudoknots. *Journal of Computational Chemistry* **25**, 1295-1304, doi:10.1002/jcc.20057 (2004).
- 6 Dirks, R. M., Bois, J. S., Schaeffer, J. M., Winfree, E. & Pierce, N. A. Thermodynamic analysis of interacting nucleic acid strands. *SIAM Review* **49**, 65-88, doi:10.1137/060651100 (2007).
- 7 Pillai, P. P., Reisewitz, S., Schroeder, H. & Niemeyer, C. M. Quantum-dot-encoded silica nanospheres for nucleic acid hybridization. *Small* **6**, 2130-2134, doi:10.1002/smll.201000949 (2010).
- 8 Leidner, A. *et al.* Biopebbles: DNA-functionalized core-shell silica nanospheres for cellular uptake and cell guidance studies. *Advanced Functional Materials* **28**, 1707572, doi:10.1002/adfm.201707572 (2018).
- 9 Sun, P. *et al.* Biopebble containers: DNA-directed surface assembly of mesoporous silica nanoparticles for cell studies. *Small* **15**, 1900083, doi:10.1002/smll.201900083 (2019).
- 10 Heckel, R., Mikutis, G. & Grass, R. N. A characterization of the DNA data storage channel. *Scientific Reports* **9**, 9663, doi:10.1038/s41598-019-45832-6 (2019).