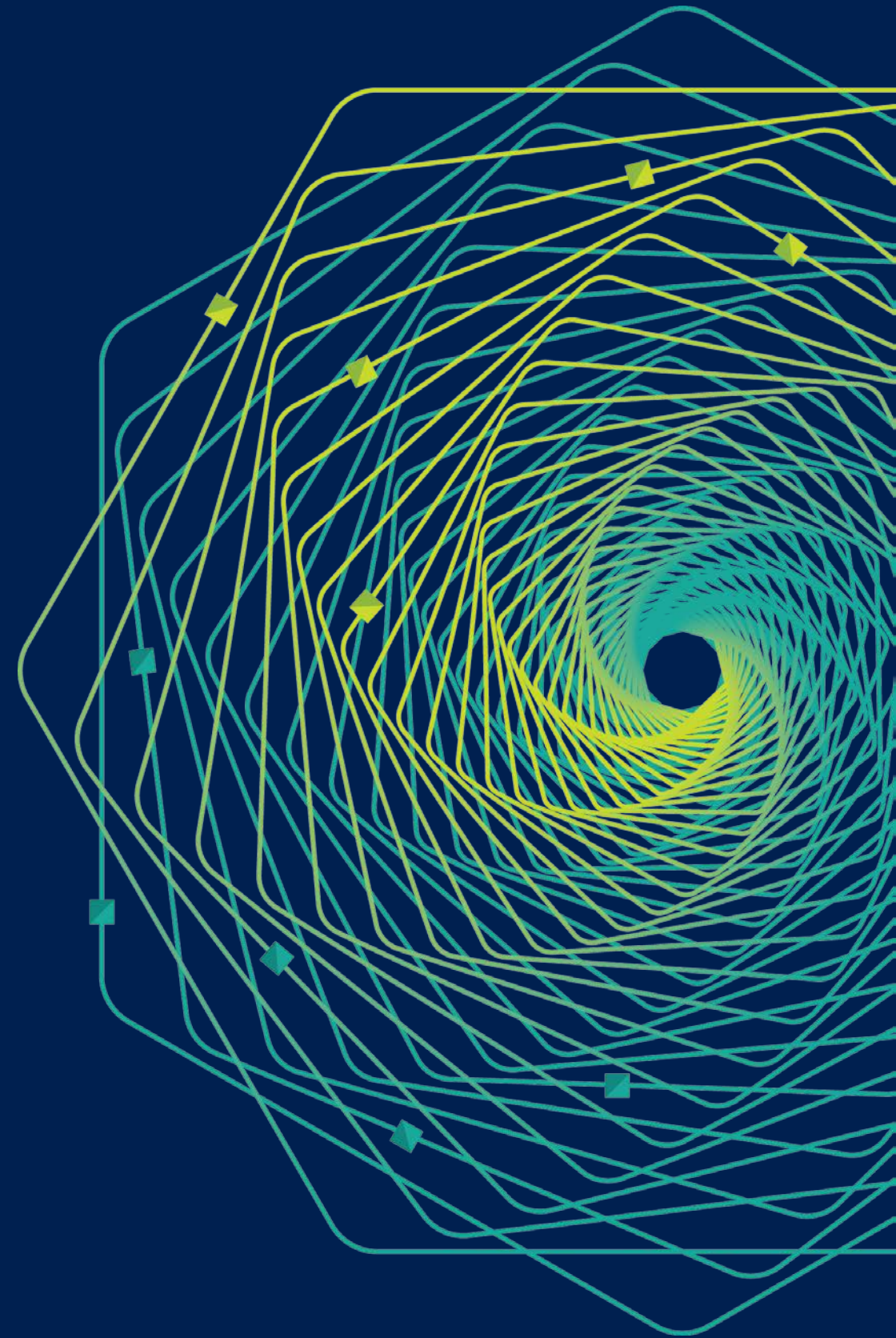


# Research Faculty Summit 2018

Systems | Fueling future disruptions



# Getting Polymers to Tell a Story: DNA Data Storage and Processing-in-Molecules

Luis Ceze

Molecular Information Systems Lab  
Sampa Lab for Hardware/Software Co-Design  
Paul G. Allen School of Computer Science & Engineering  
University of Washington

joint work with **Karin Strauss**, Georg Seelig, Doug Carmean, Sergey Yekhanin, Lee Organick, Yuan-Jyue Chen, Bichlien Nguyen, Chris Takahashi, Ashley Stephenson, Pranav Vaid, Sharon Newmann, Cyrus Rashtchian, Miklos Racz, Siena Ang, David Ward, Randolph Lopez, Max Willsey, Kendall Stewart, James Bornholt, Rob Carlson, Hsing-Yeh Parker.

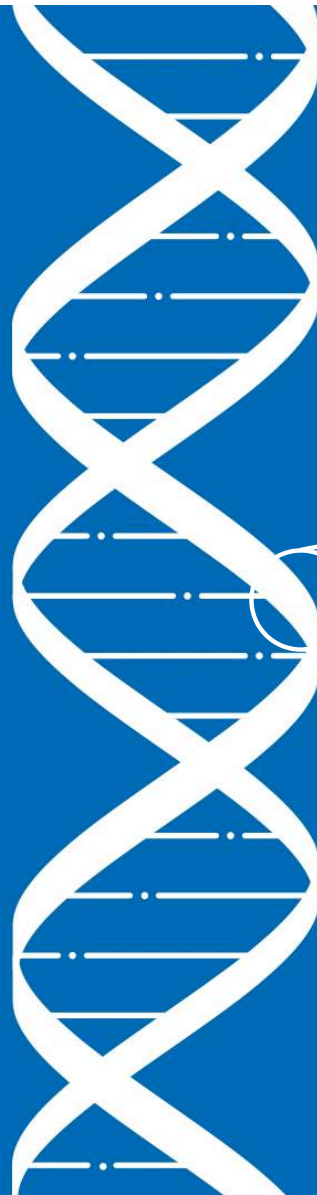
Microsoft Faculty Summit 2018



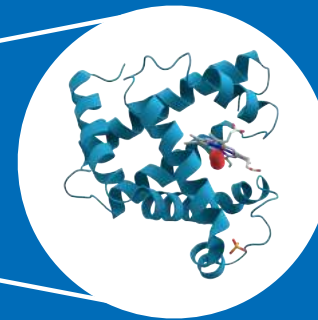
**W** PAUL G. ALLEN SCHOOL  
OF COMPUTER SCIENCE & ENGINEERING



**DNA** is Nature's  
information  
storage medium



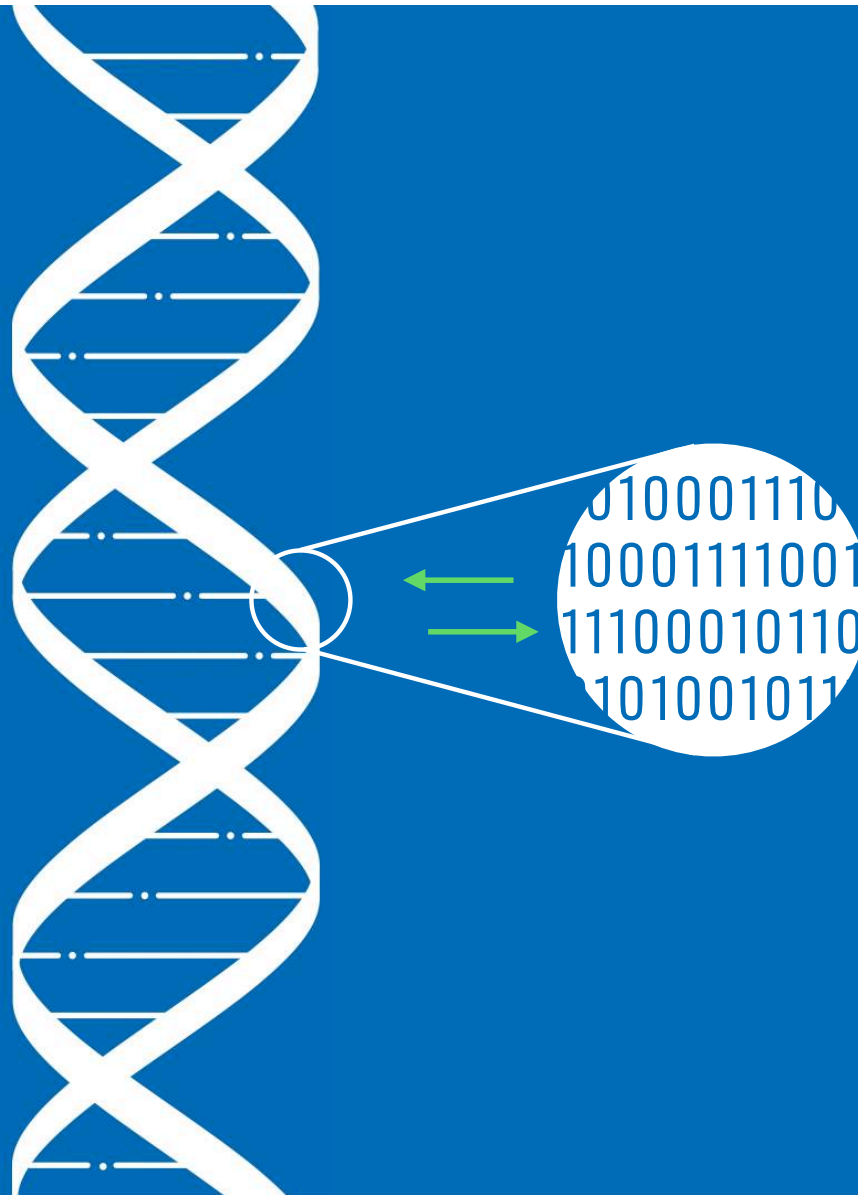
Gene →



Protein

Function/  
Characteristic

Using Synthetic  
**DNA** for Data  
Storage  
Manufactured DNA



# DNA Molecules for Digital Data



Extremely Dense

1 Exabyte in 1 in<sup>3</sup>

Extremely Durable

1,000s of Years

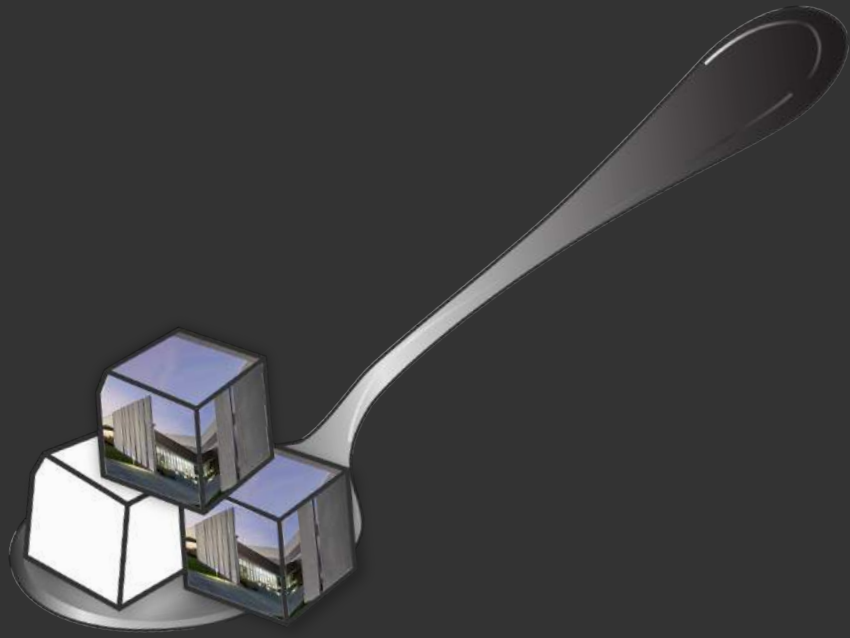
Making Copies Is Nearly Free

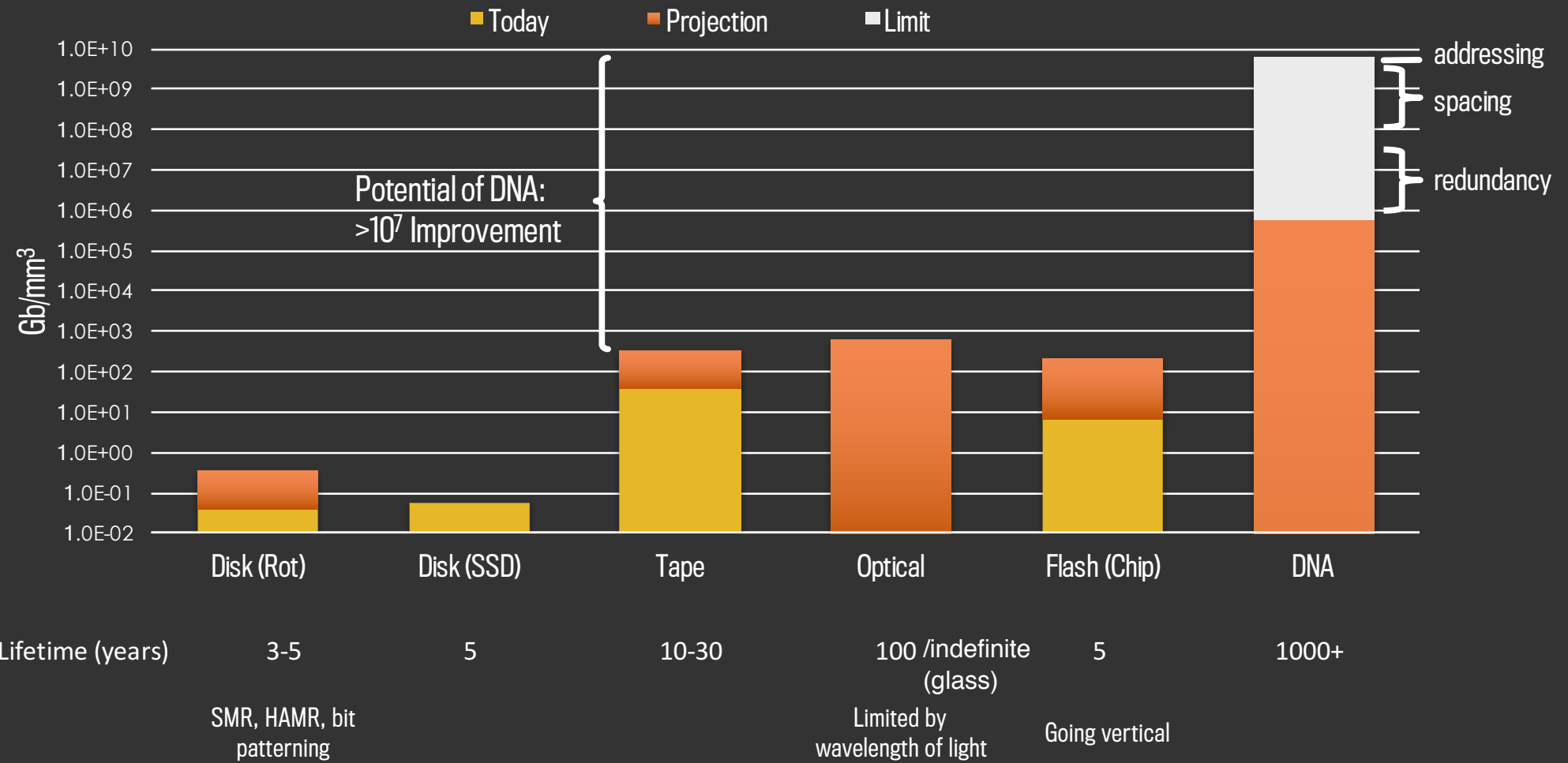
Never Gets Obsolete



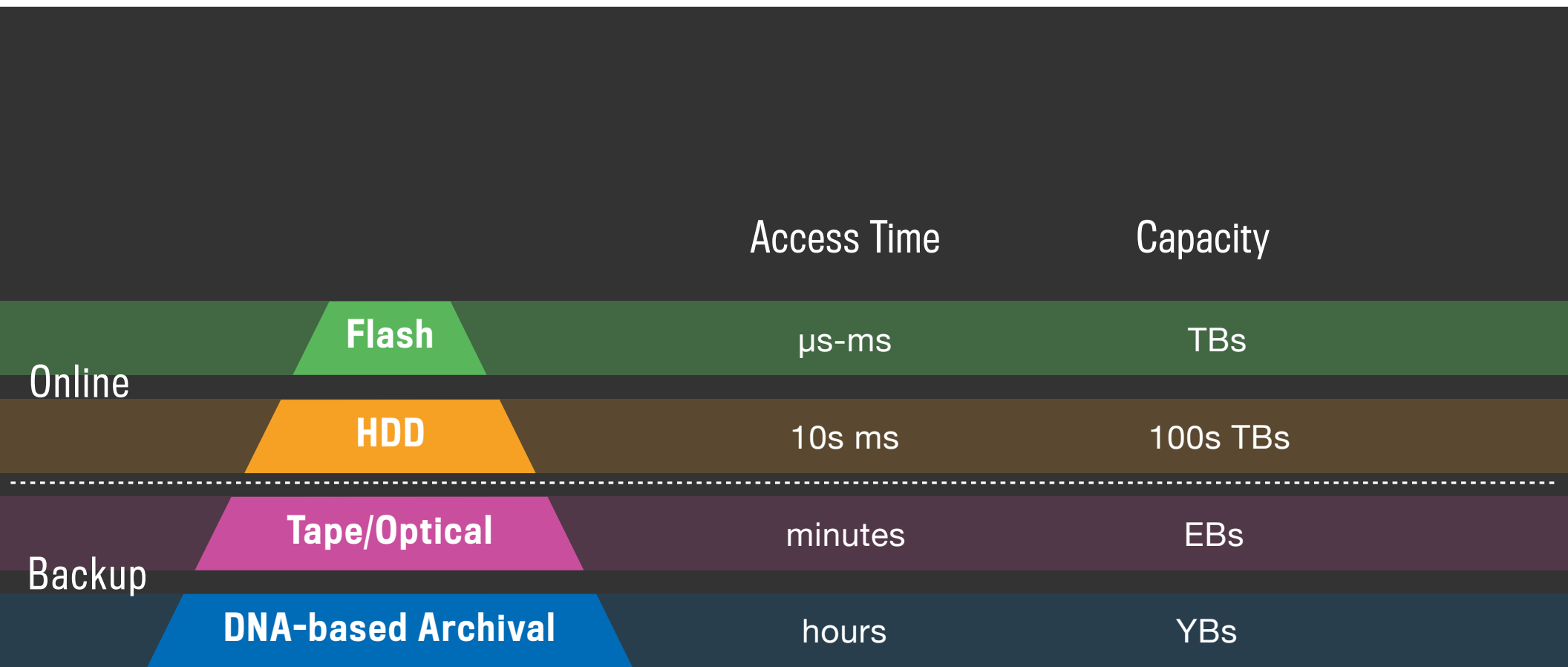


**A few  
exabytes**





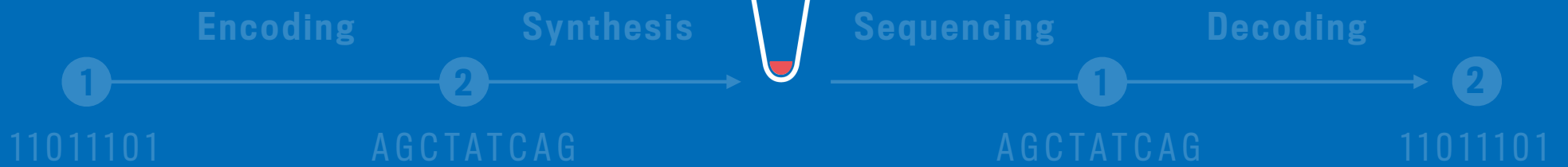




[in ASPLOS'16]

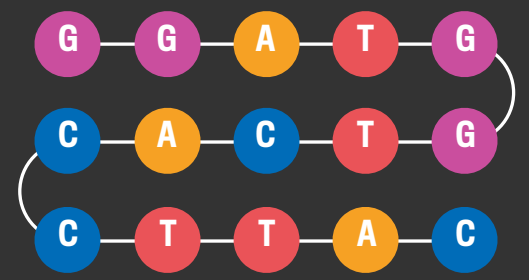
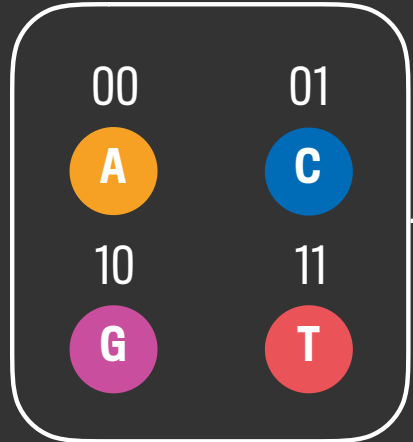
# Write Path

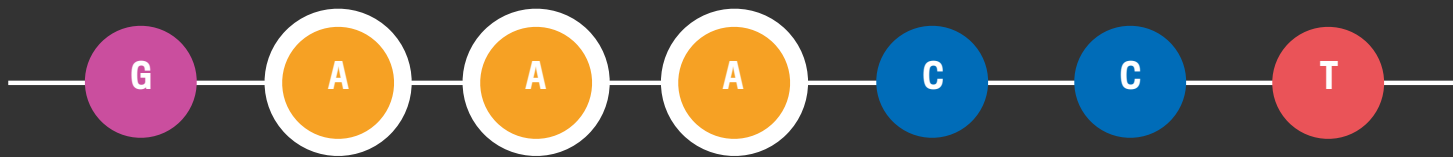
# Read Path



# Encoding Digital Data in DNA

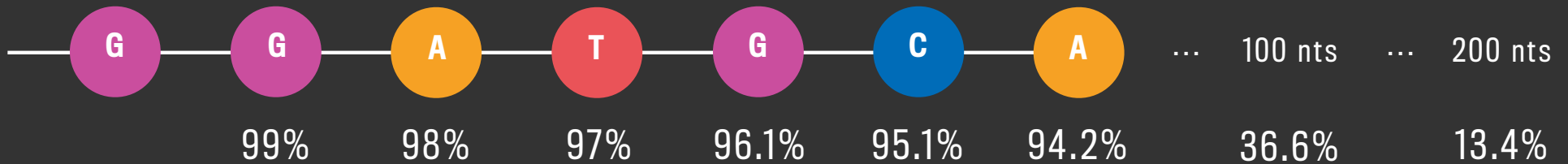
1010001110  
0100011110  
0111110001  
0110010100





Repeated Letters Are Bad: **Avoid them with Randomization & Rotation.**

$P[\text{Attach}] = 99\%$



Synthetic DNA has limited length: **Break it into chunks.**

# Break into chunks and add redundancy

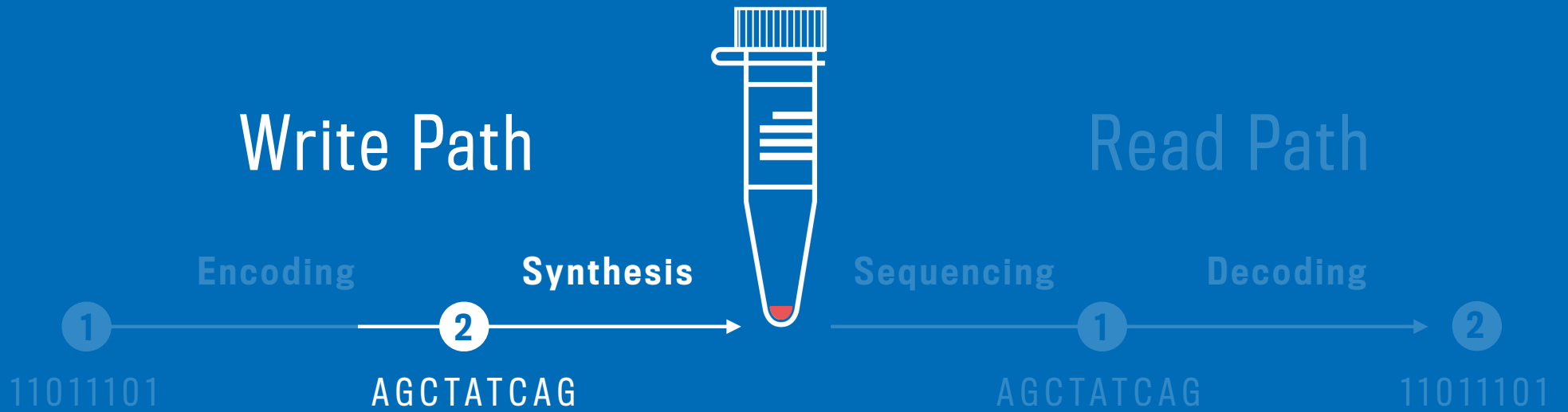


Redundant Data in  
Additional DNA  
Strands. Many  
Possibilities:  
Parity, Reed  
Solomon, LDPC, ...

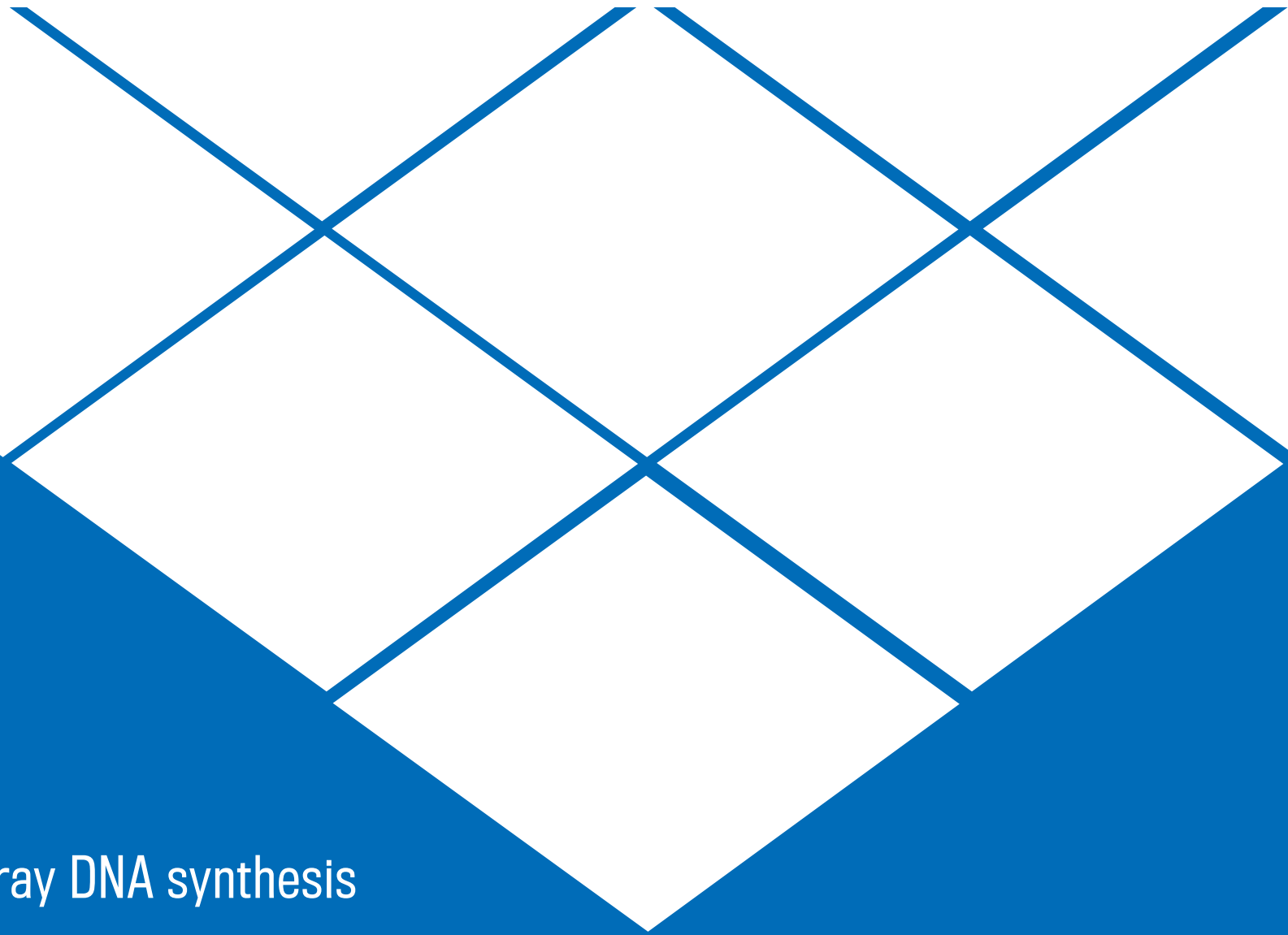
~20 Bytes per 150nt DNA sequence. Many sequences per file. ~1% (in, del, sub) error per base.

# Write Path

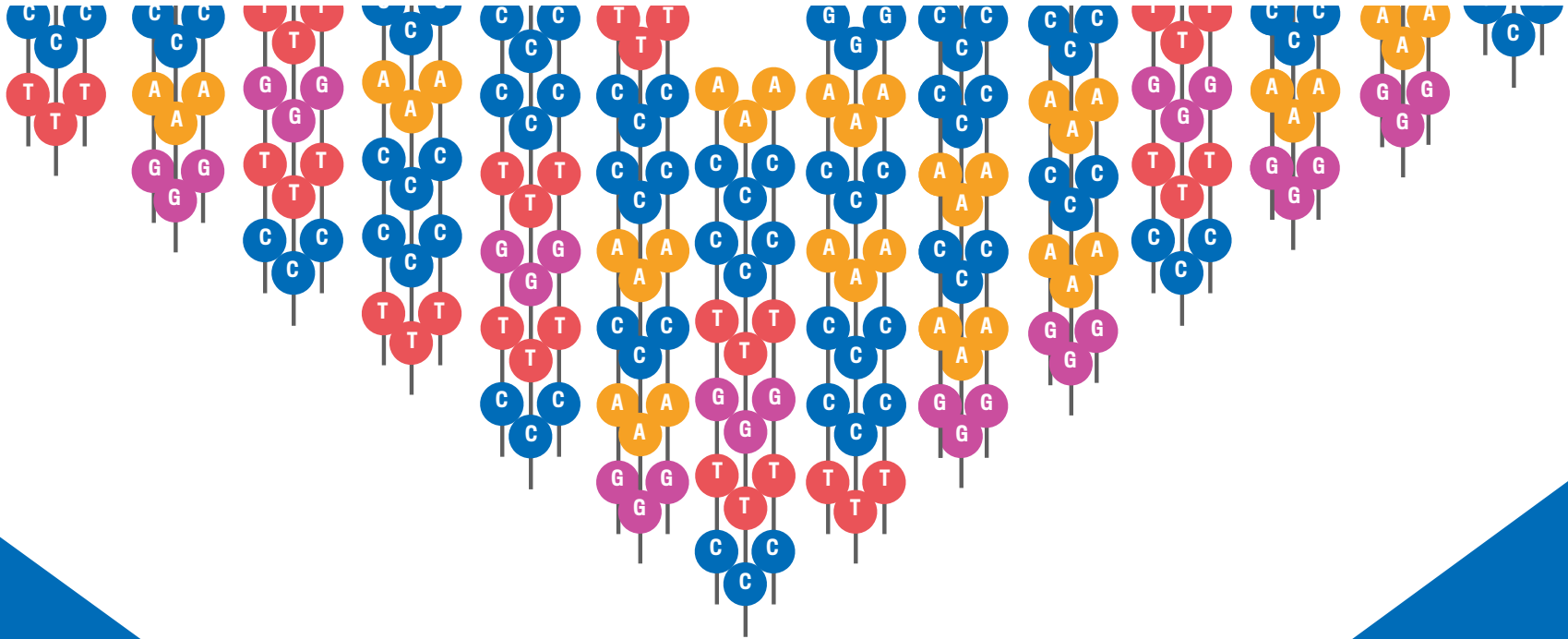
# Read Path







Large array DNA synthesis

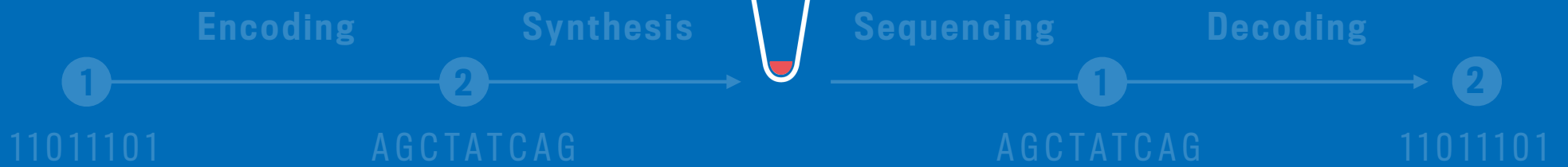


Large array DNA synthesis

Each spot grows many copies  
of a given sequence

# Write Path

# Read Path



# DNA Sequencing

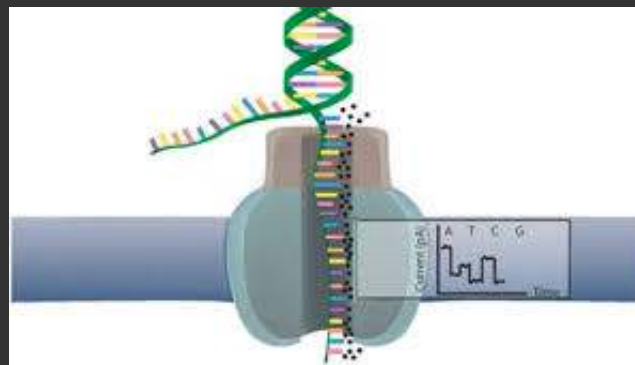
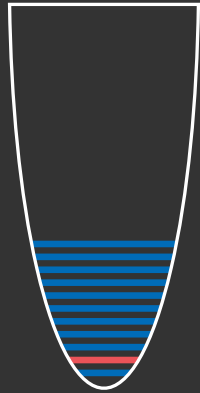
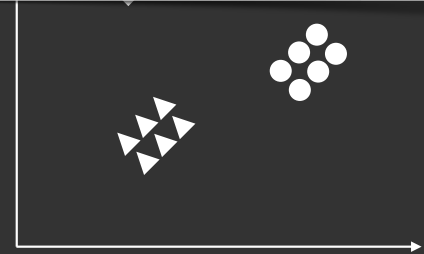
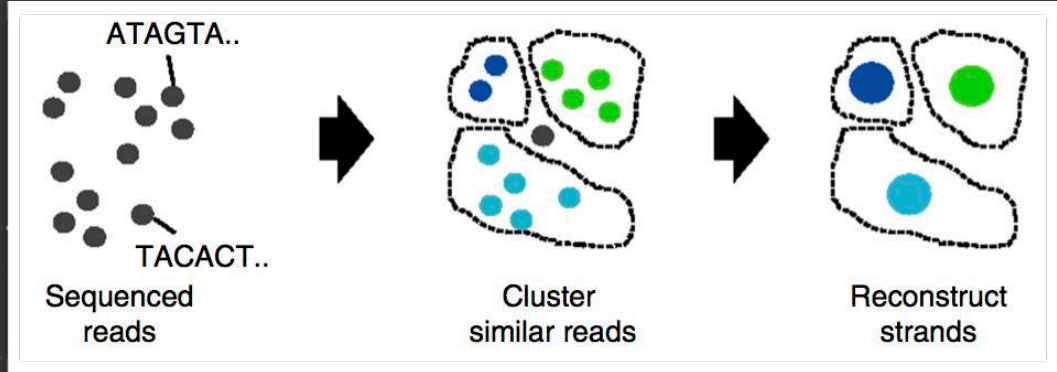


Image credit: Oxford Nanopore



ATGTT TGC  
ATGTT GCC  
ATGTT GGA  
ATGTT TGCT TACC CAAC CATCC  
ATGTT GCCA GTTC AAAG CATGC  
ATGTT GGAT GCAC AAGA CATCC  
ATGTT TGCT TACC CAAC CATCC  
ATGTT GCCA GTTC AAAG CATCC

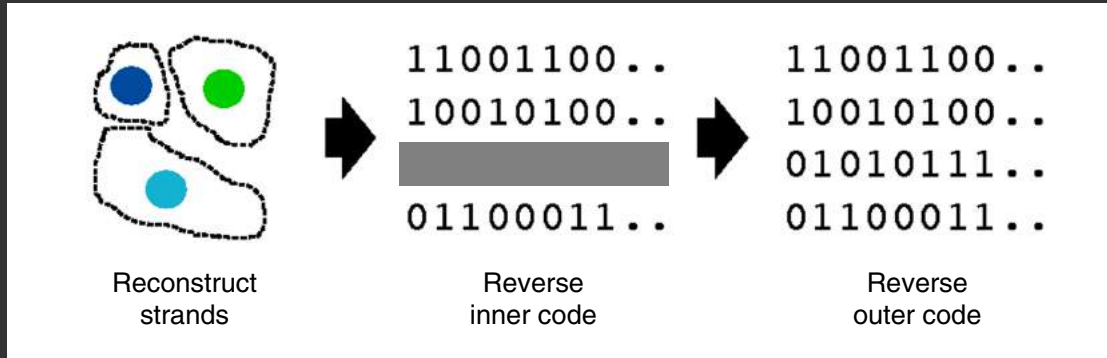


Selected DNA strands

Sequencing (Noisy reads)

Clustering (No reference)

[in NIPS'17]



*ithm. (a. i) We designed a primer library for our PCR-based reader.*

...  
11011  
01011  
01...



Reassemble  
Data

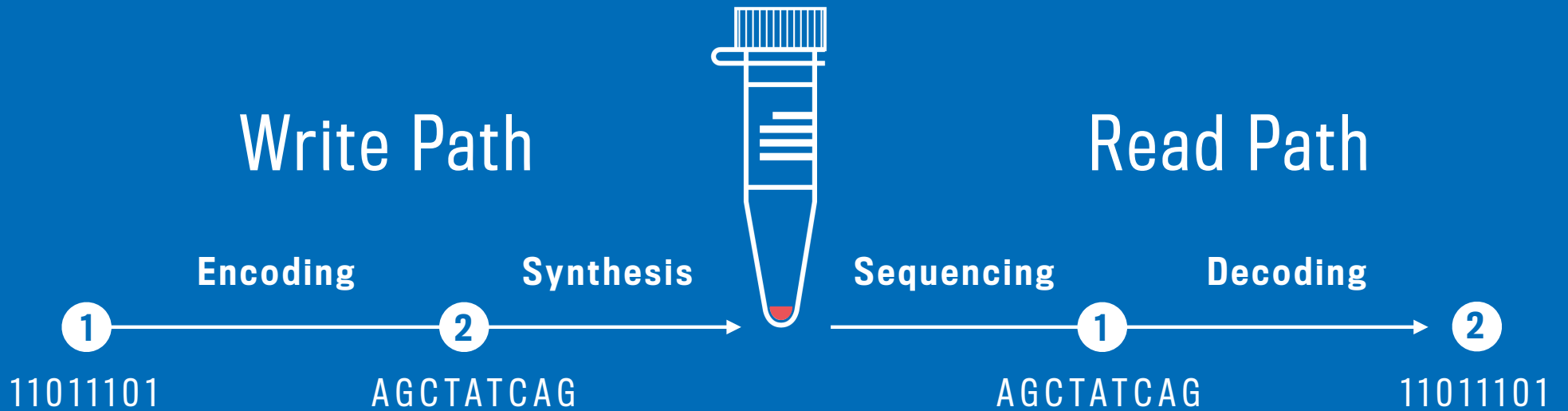
Error  
Correction

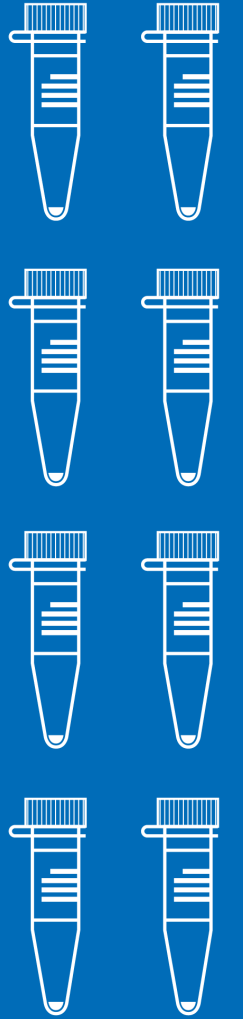
Error-Free  
Data

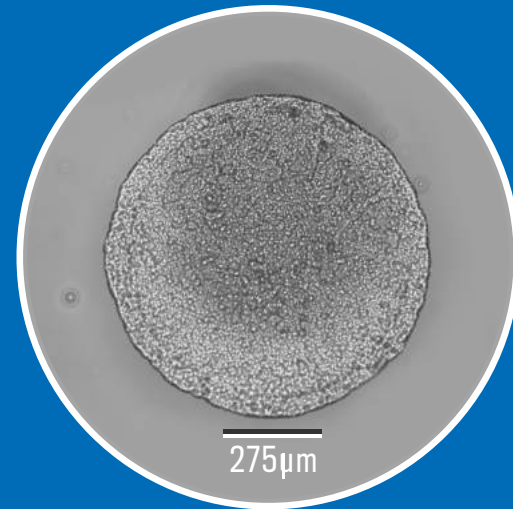
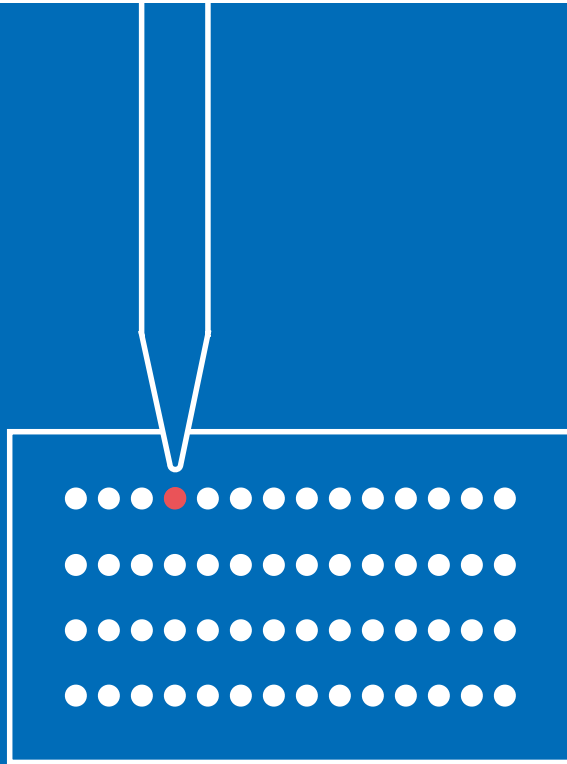


# Write Path

# Read Path

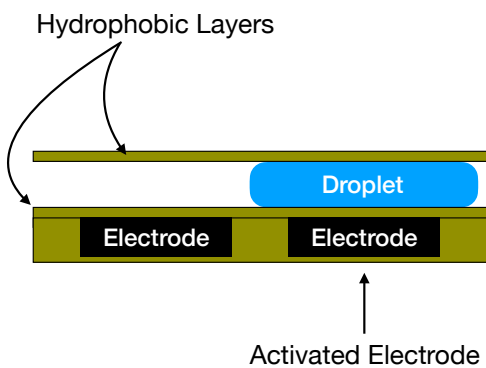




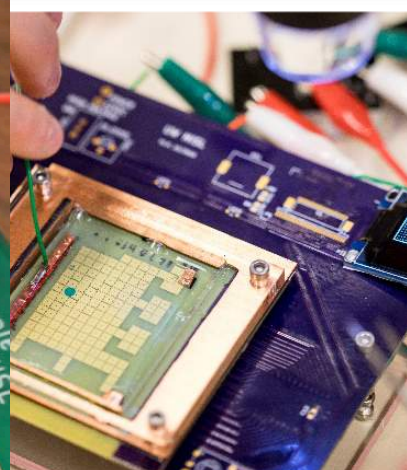
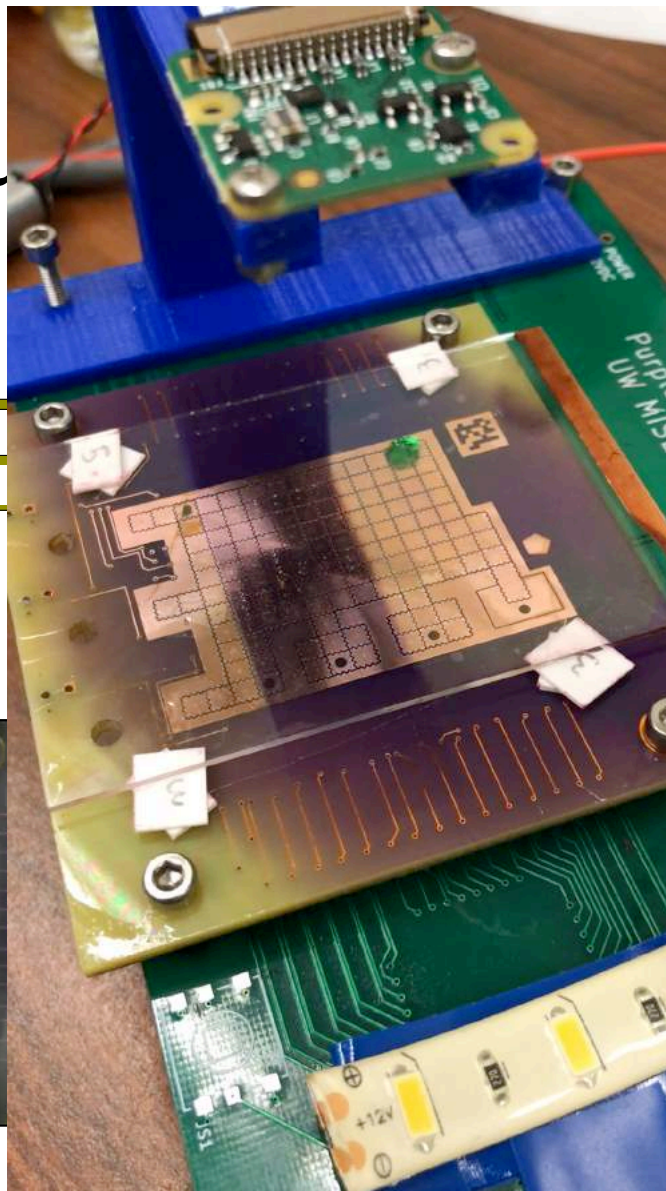


~100TB per spot

# Digital microfluidics



Molecular domain



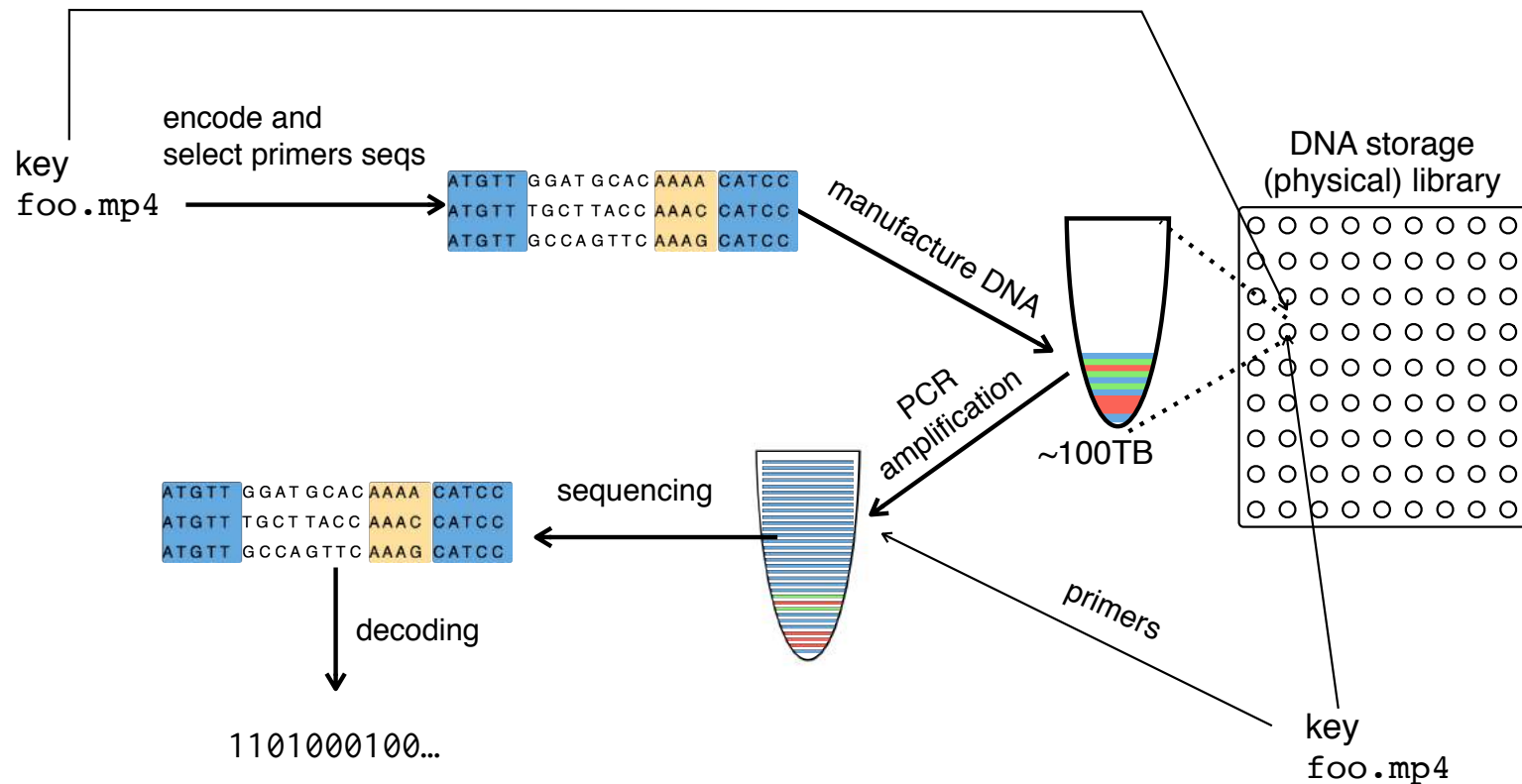
Electronic domain

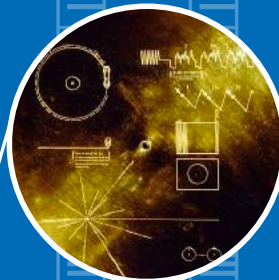




# Putting it all together as (key, value) store

Data address specifies physical location and primer for random access.





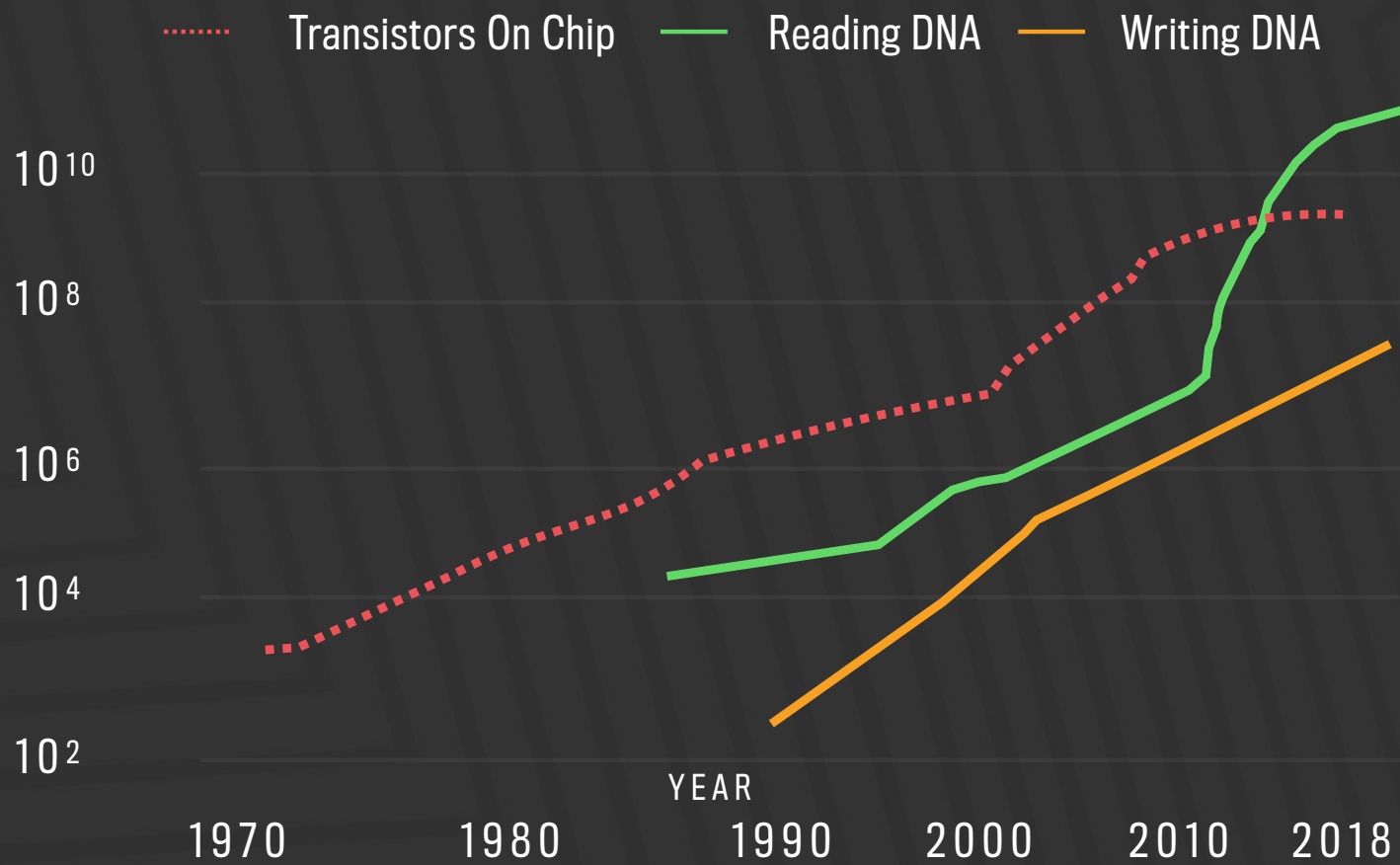
Over 700MB. 50M+ sequences. 9B+ Nucleotides,  
4B+ reads. Demonstrated random access w/ 40+ objects.  
Illumina and Nanopore sequencing readout.

[Nature Biotechnology'18]



A photograph of several blue network cables plugged into a switch or patch panel. The background is dark with many out-of-focus yellow and green lights, creating a bokeh effect. The cables have white labels with black text. One label clearly shows 'fsw003.p038'. Another label shows 'CB-FB-208004'. A third label shows '50 981/940ms'. A fourth label shows 'S/N 1K-1405320-019M'.

10MBs/day → 100GBs/second



Quality	
Life sciences	perfect
Data storage	arbitrarily tolerant

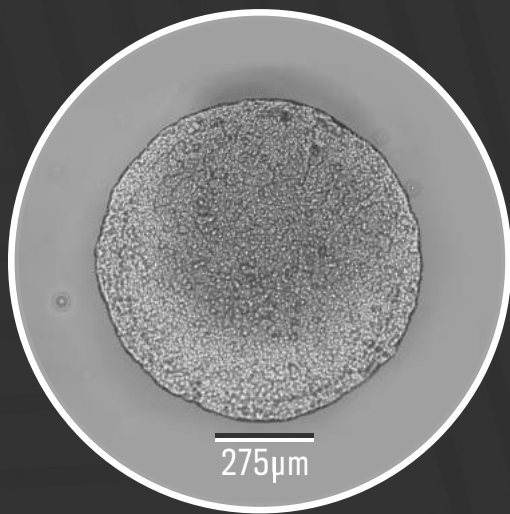
Source: Robert Carlson

# Scalability

- Data centers offer perfect abstraction for “exotic technologies” (Carmean)
- Large-scale fluidics for synthesis, manipulation and sequencing
- Throughput of ~1TB/s at the data-center level
- Computational cost significant
  - Today: ~2.8KB/s encode, ~1KB/s decode on 16 core Xeon.

# Beyond DNA Data Storage

# DNA “computing” in the age of big data

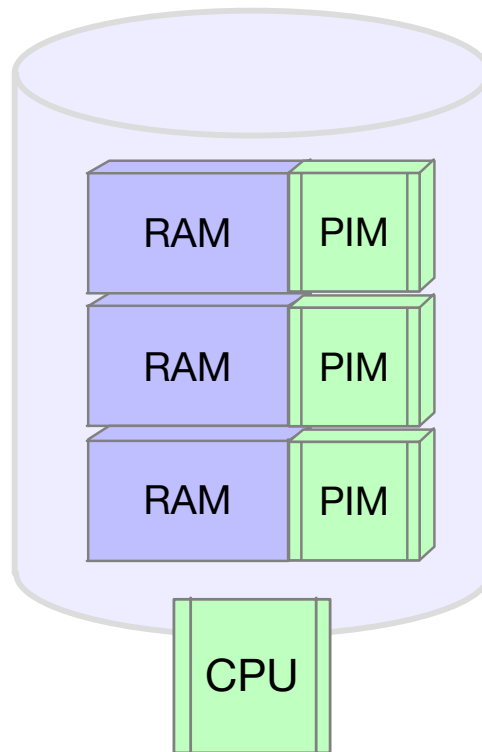


~100TB per spot

If DNA data storage succeeds, what if we could process data directly in DNA?

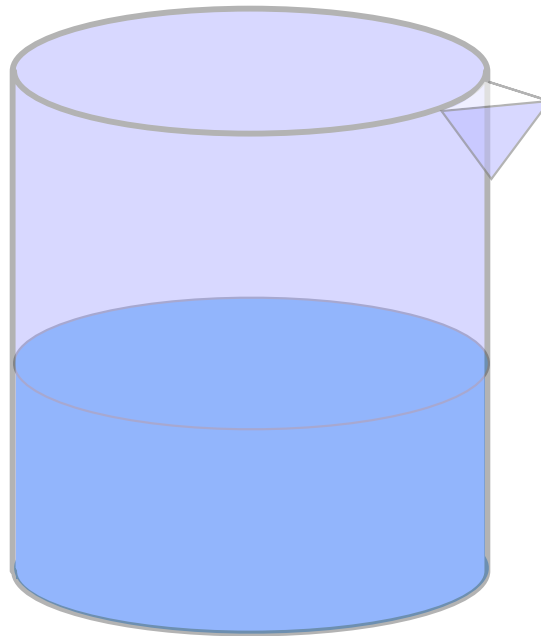
Extremely parallel and energy efficient

# Processing-in-Memory

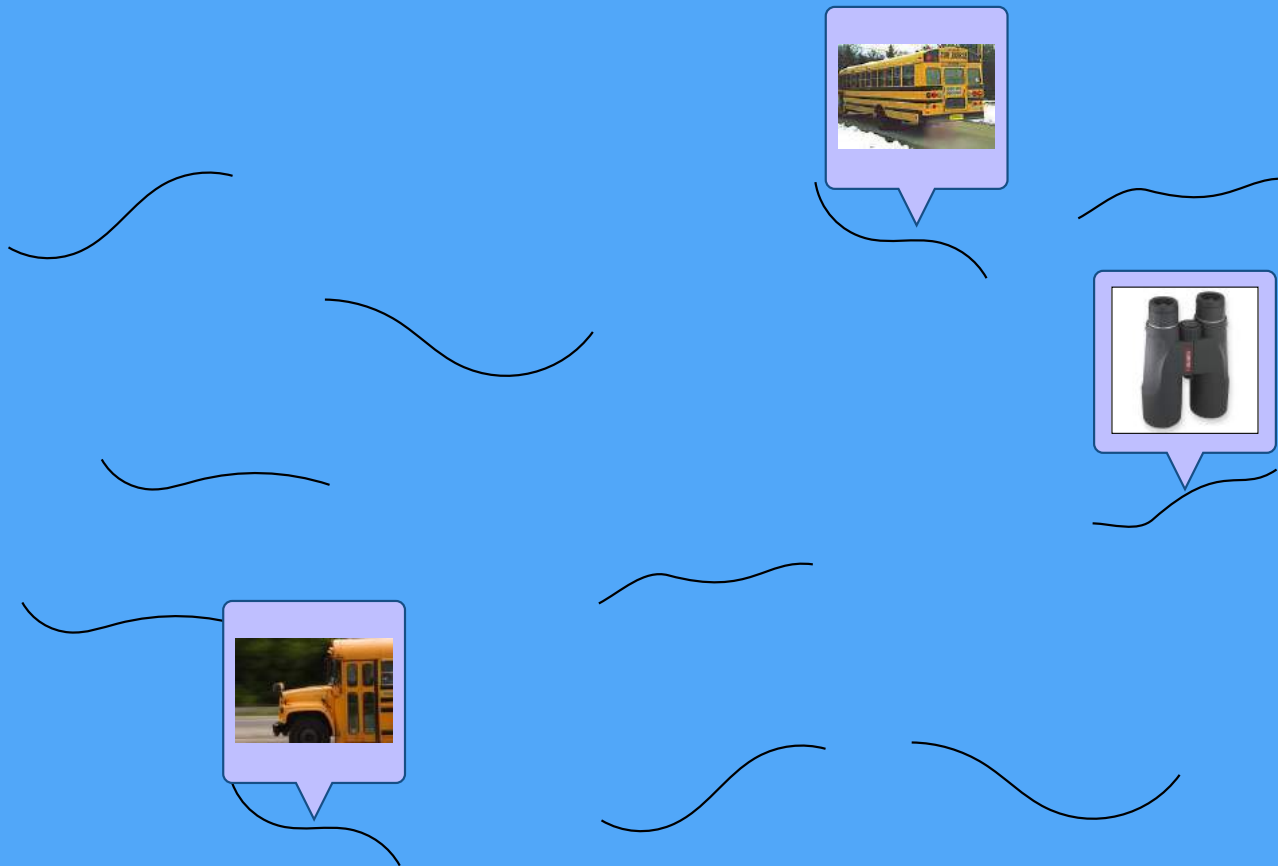


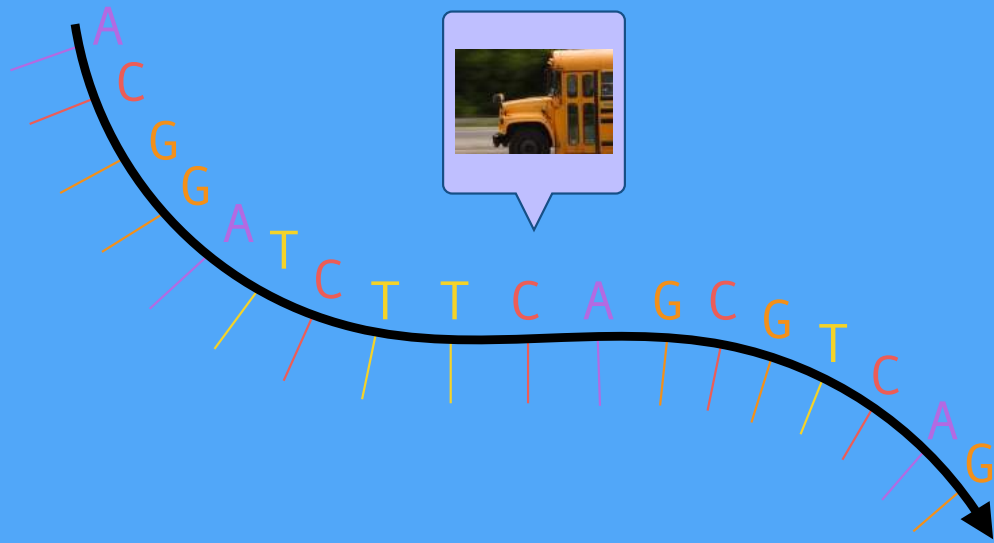


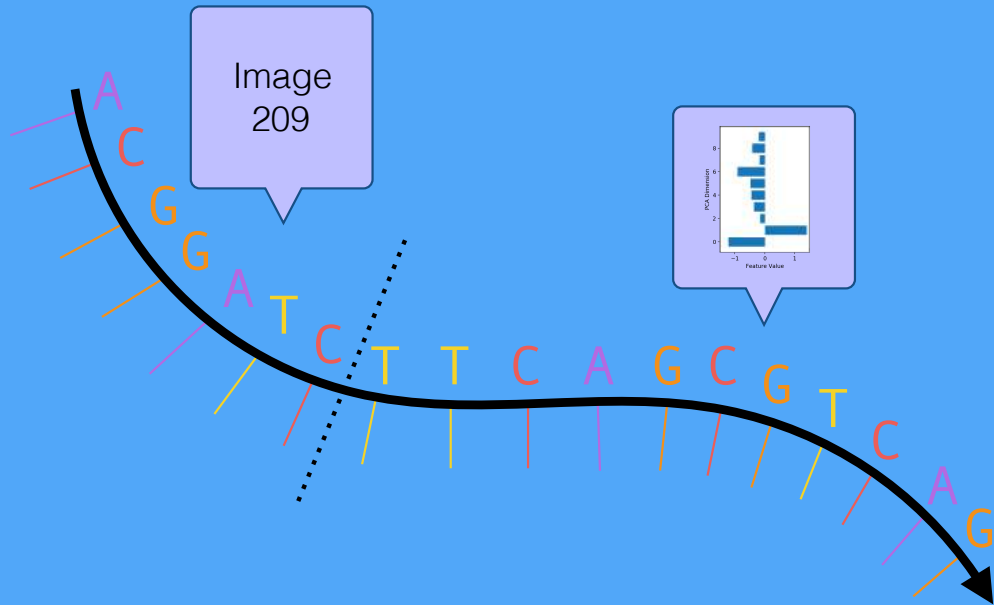
# Processing-in-Molecules



# Storage *and* Processing-in-Molecules (DNA)







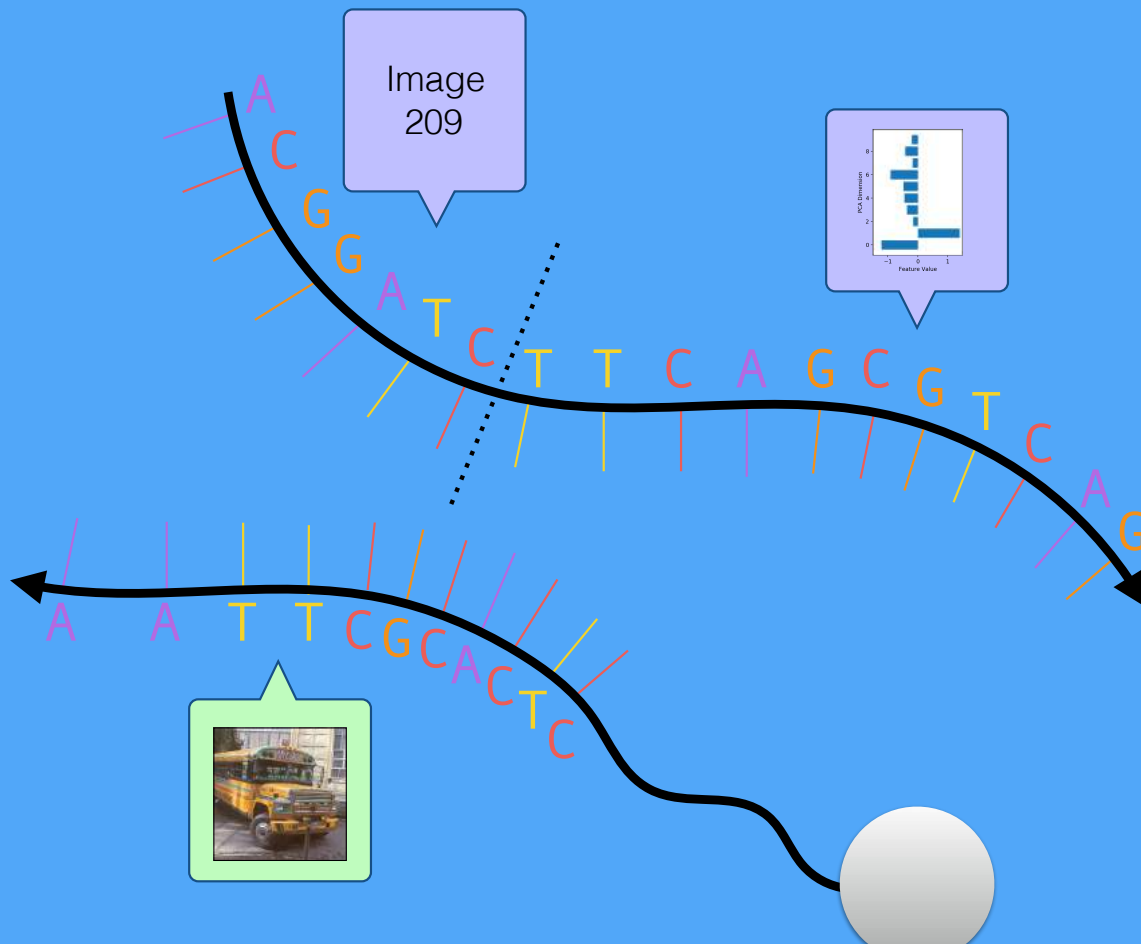
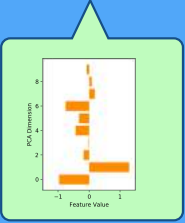
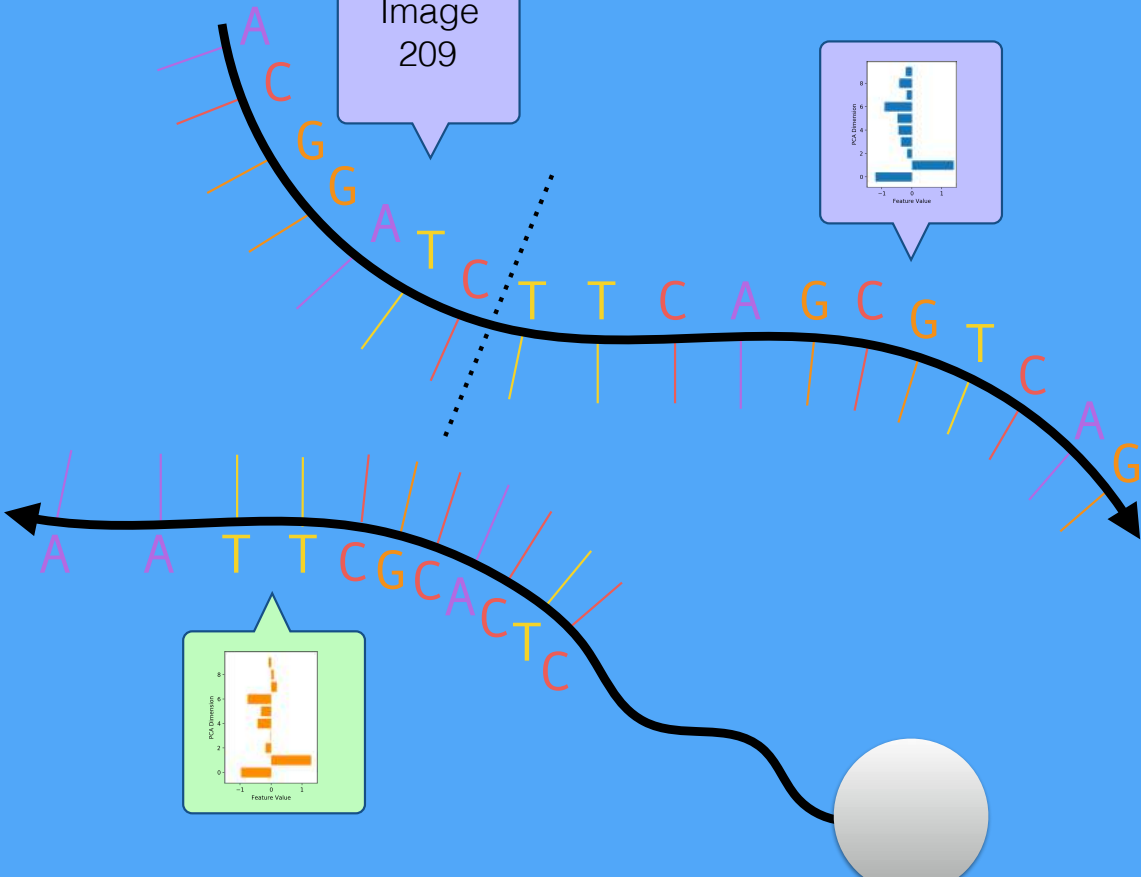
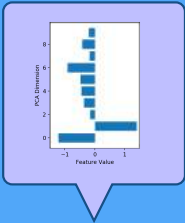
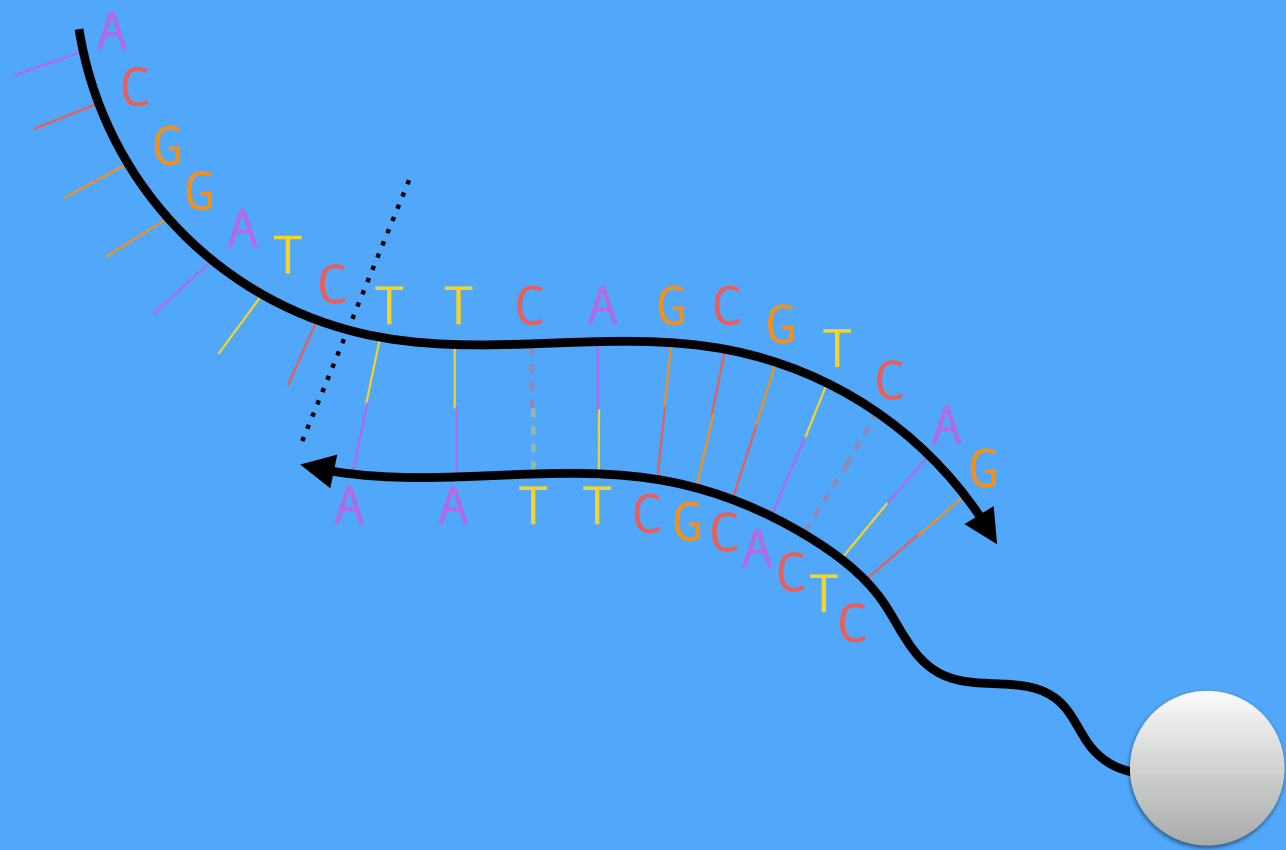
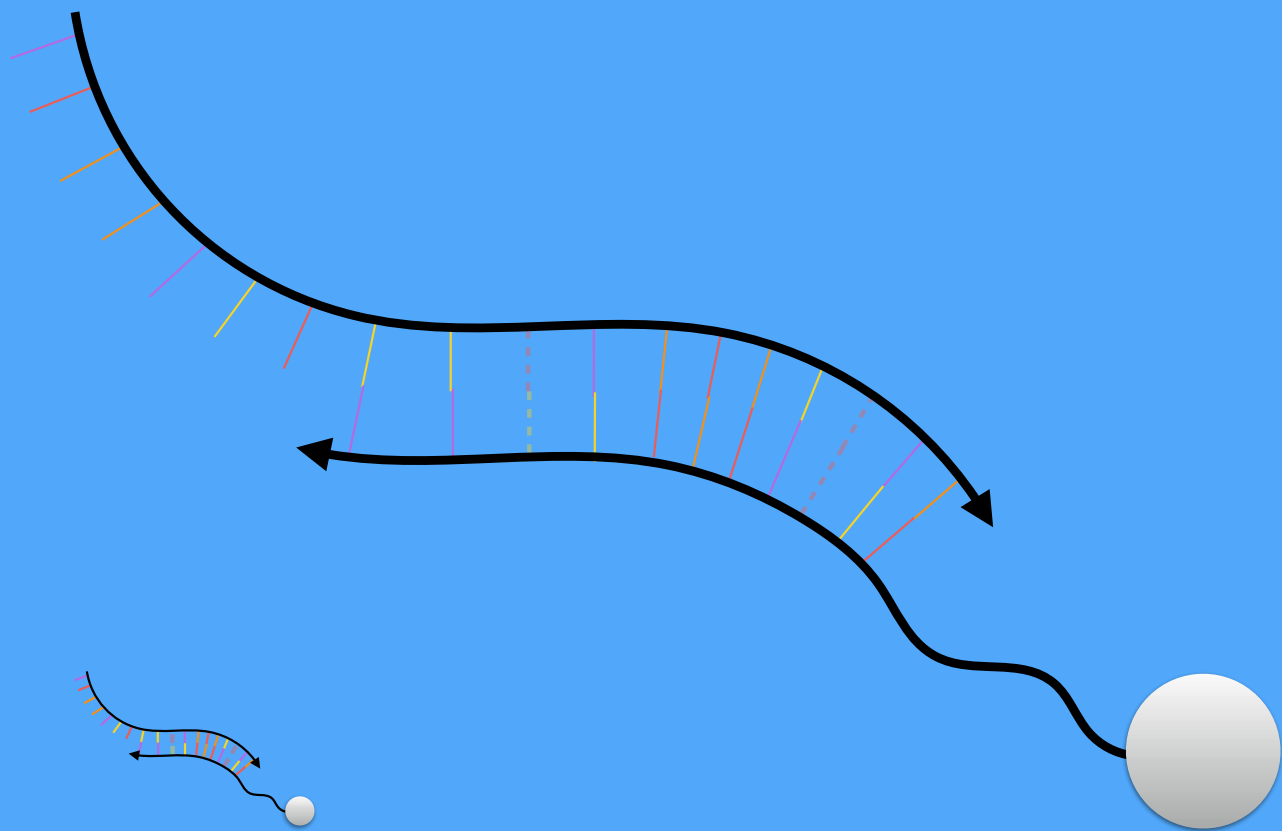


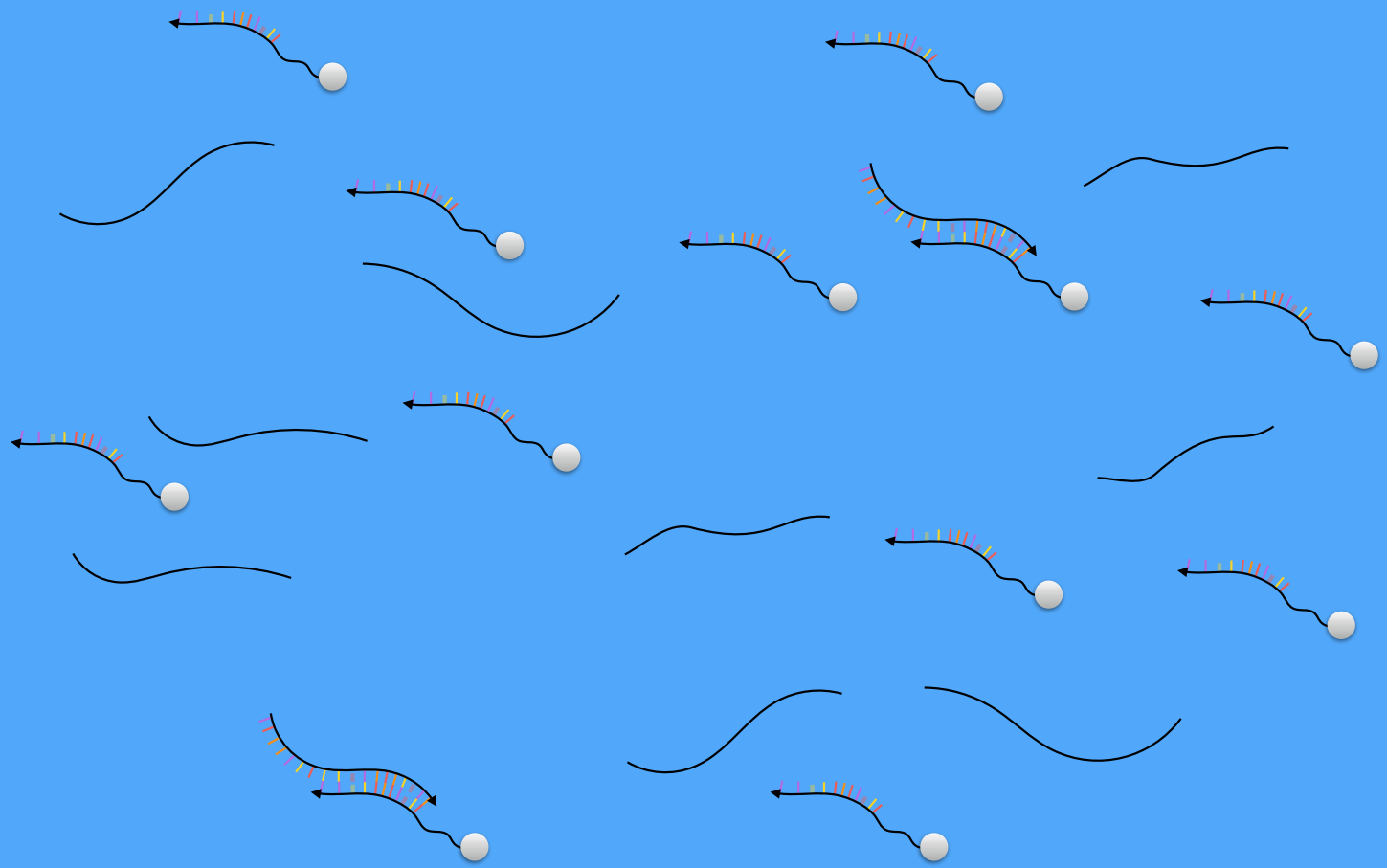
Image  
209

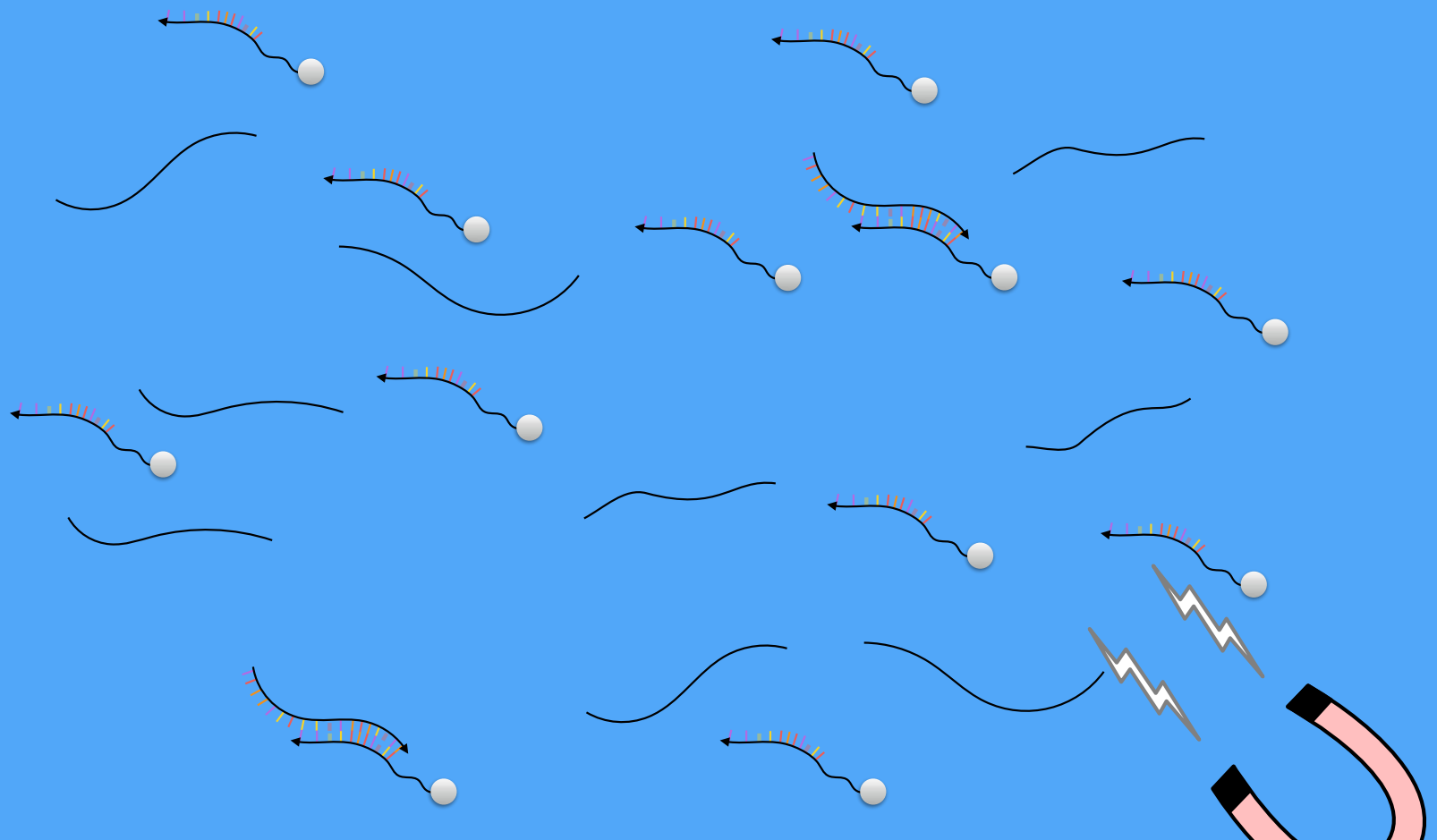


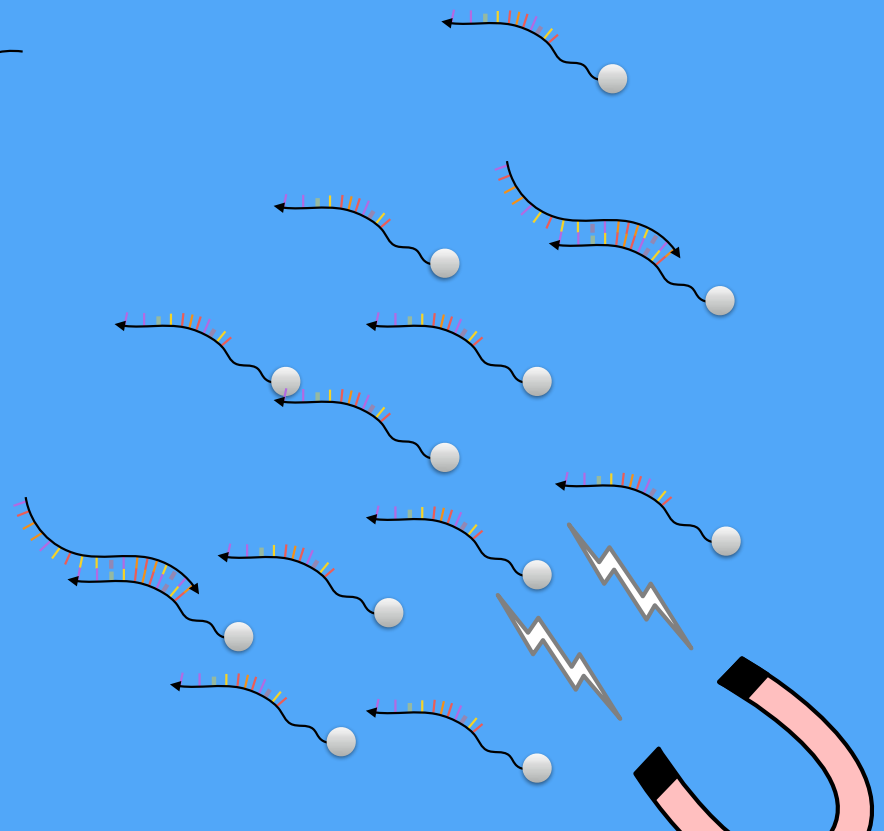
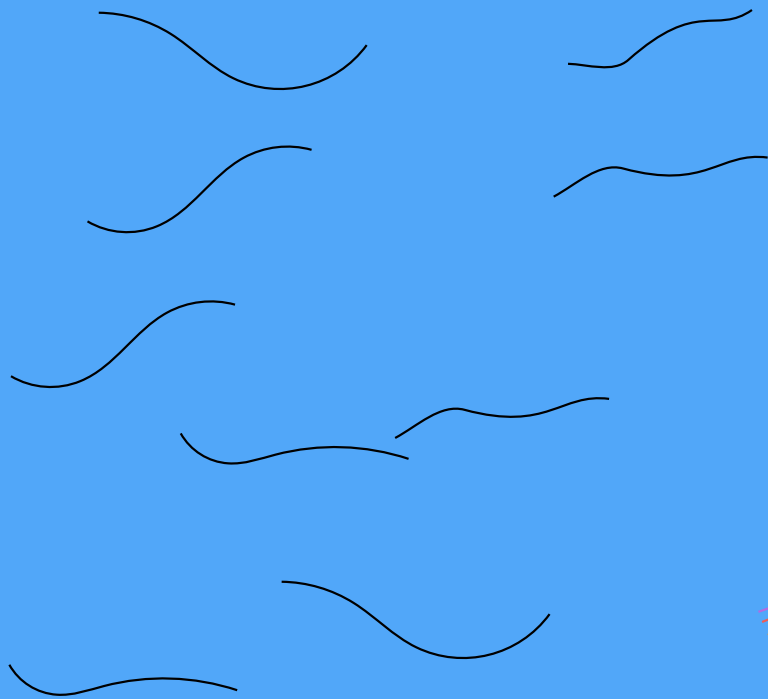


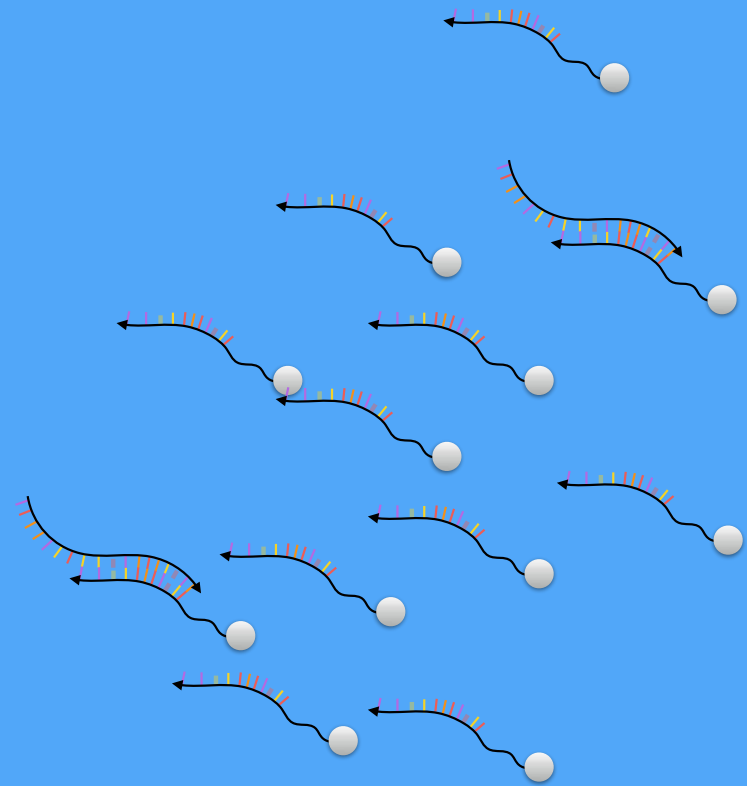


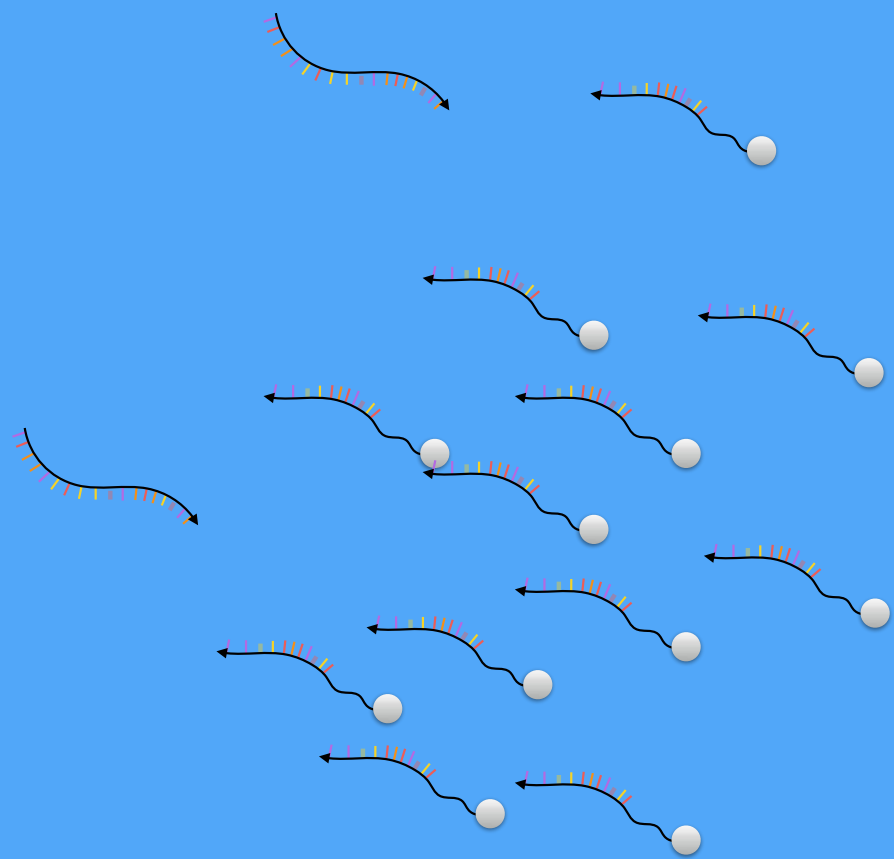


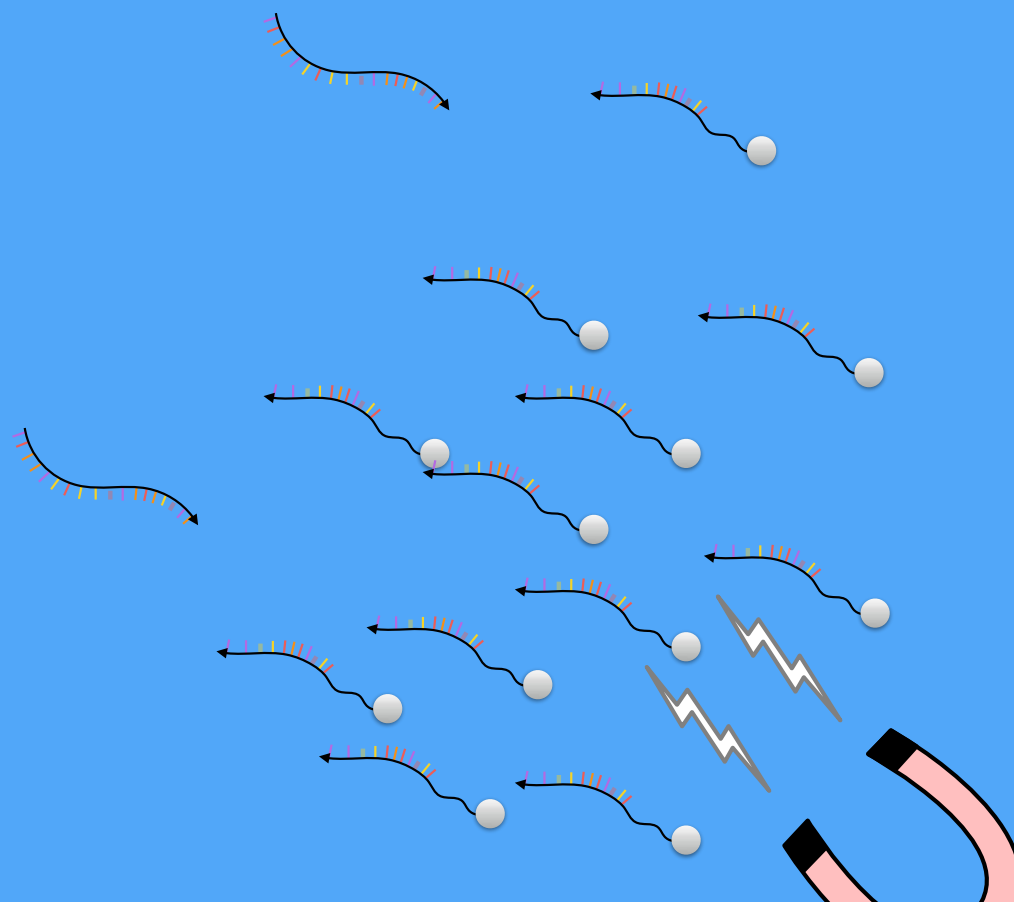




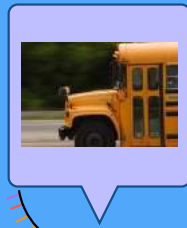






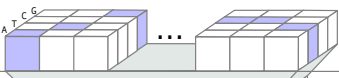






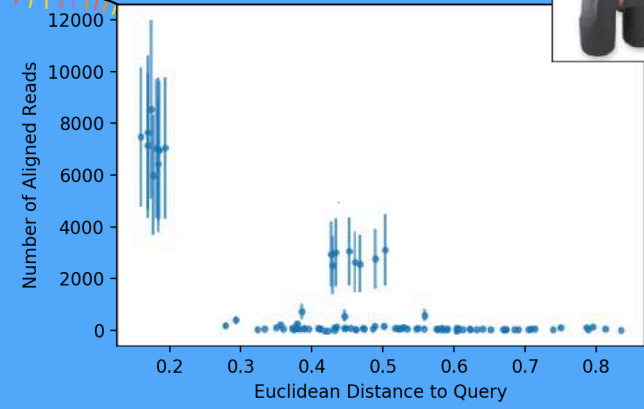
# Yottabyte-scale Associate Memories?

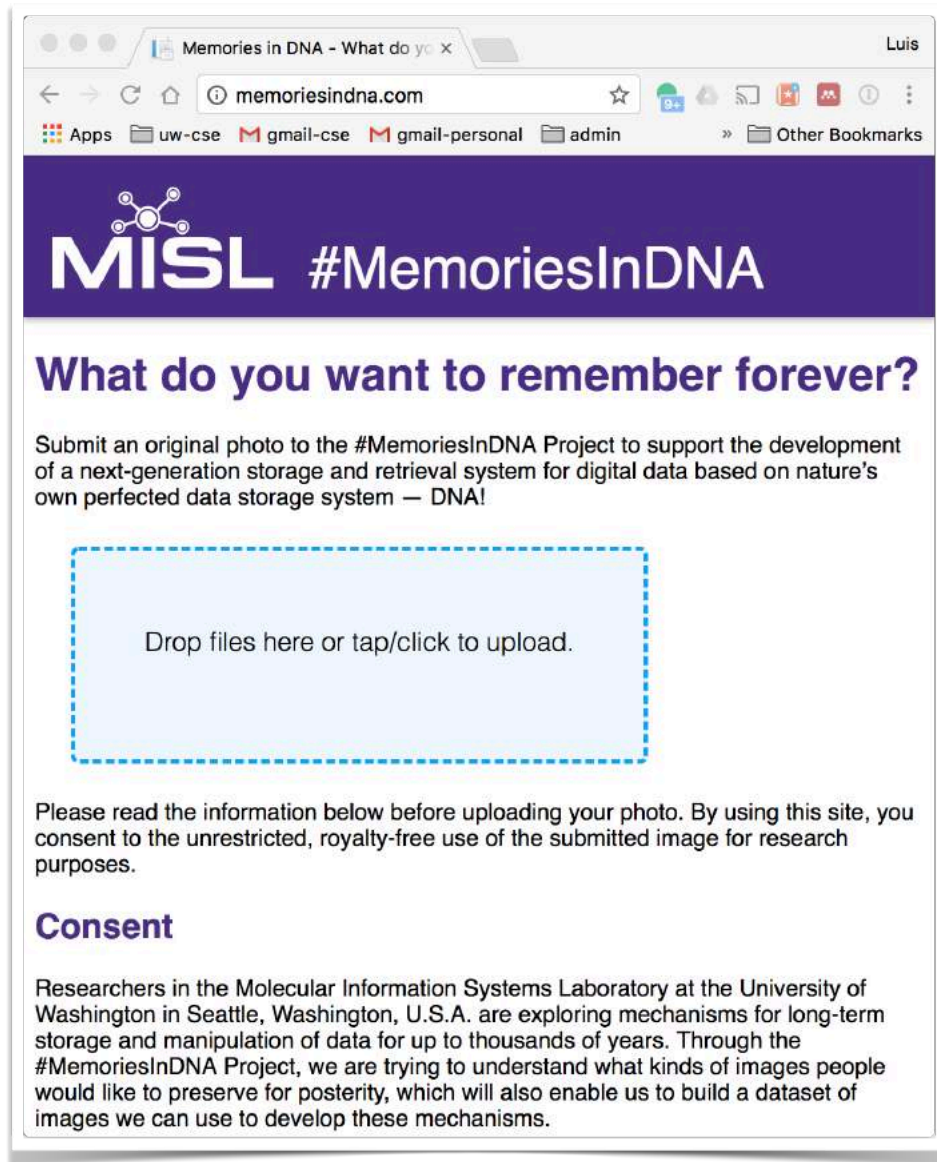


Query Image:



Layer		Size
Sequence Output	A T G ... C C T	30 x 1
ReLU + Softmax Activations		30 x 4
Fully Connected Weights		10 x 128 x 30 x 4
ReLU Activations	...	10 x 128
Convolutional Weights 2		1 x 128 x 128
Sine Activations	...	10 x 128
Convolutional Weights 1		1 x 128
Input Features		10 x 1





The image shows a screenshot of a web browser displaying the website 'memoriesindna.com'. The browser's address bar shows the URL and several tabs are open. The website has a dark blue header with the 'MISL #MemoriesInDNA' logo. Below the header, the main heading asks 'What do you want to remember forever?'. A paragraph explains the project's goal: to support the development of a next-generation storage and retrieval system for digital data based on DNA. A light blue dashed box contains the text 'Drop files here or tap/click to upload.' Below this, a 'Consent' section provides information about the Molecular Information Systems Laboratory at the University of Washington, stating that researchers are exploring mechanisms for long-term storage and manipulation of data for up to thousands of years. The consent text mentions that through the #MemoriesInDNA Project, they aim to understand what kinds of images people would like to preserve for posterity, which will enable them to build a dataset of images for developing these mechanisms.

Memories in DNA - What do you want to remember forever?

memoriesindna.com

MISL #MemoriesInDNA

## What do you want to remember forever?

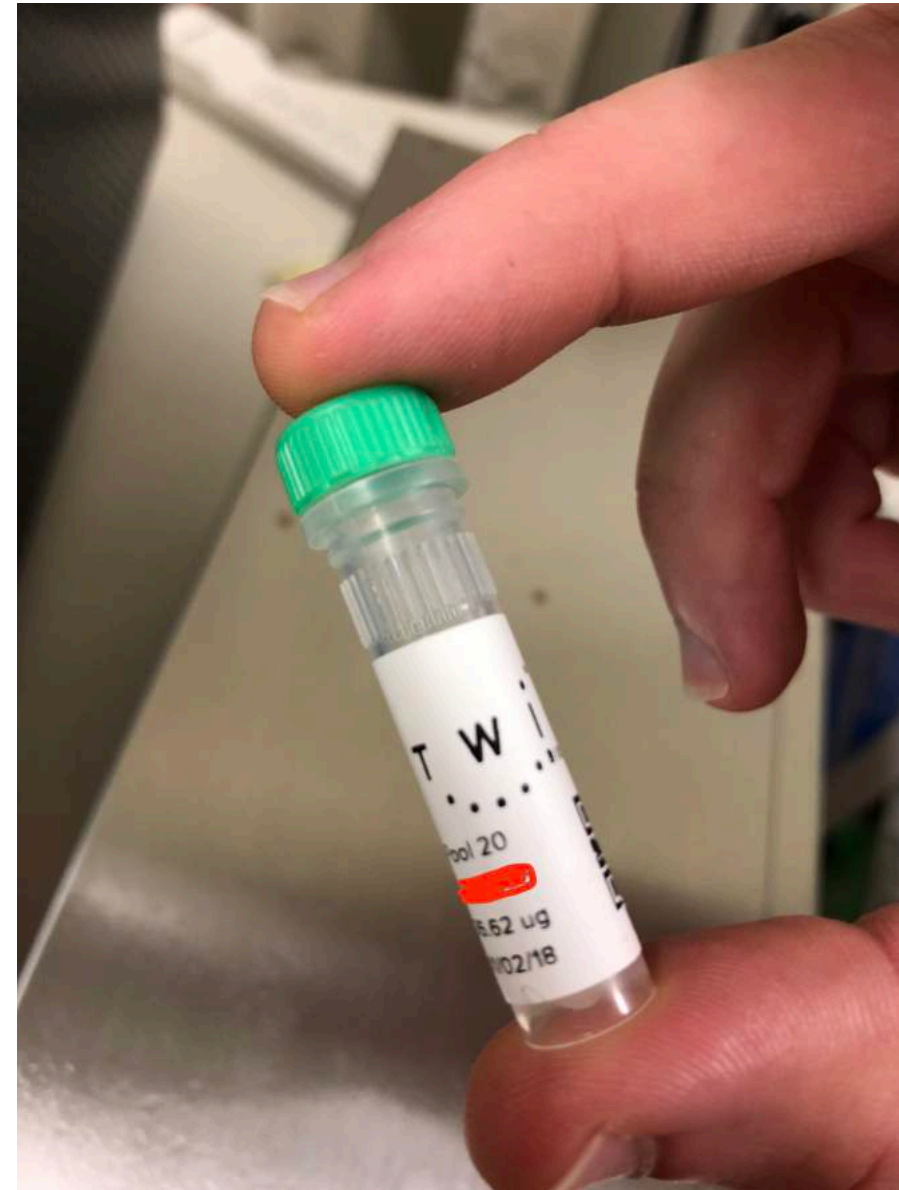
Submit an original photo to the #MemoriesInDNA Project to support the development of a next-generation storage and retrieval system for digital data based on nature's own perfected data storage system — DNA!

Drop files here or tap/click to upload.

Please read the information below before uploading your photo. By using this site, you consent to the unrestricted, royalty-free use of the submitted image for research purposes.

### Consent

Researchers in the Molecular Information Systems Laboratory at the University of Washington in Seattle, Washington, U.S.A. are exploring mechanisms for long-term storage and manipulation of data for up to thousands of years. Through the #MemoriesInDNA Project, we are trying to understand what kinds of images people would like to preserve for posterity, which will also enable us to build a dataset of images we can use to develop these mechanisms.



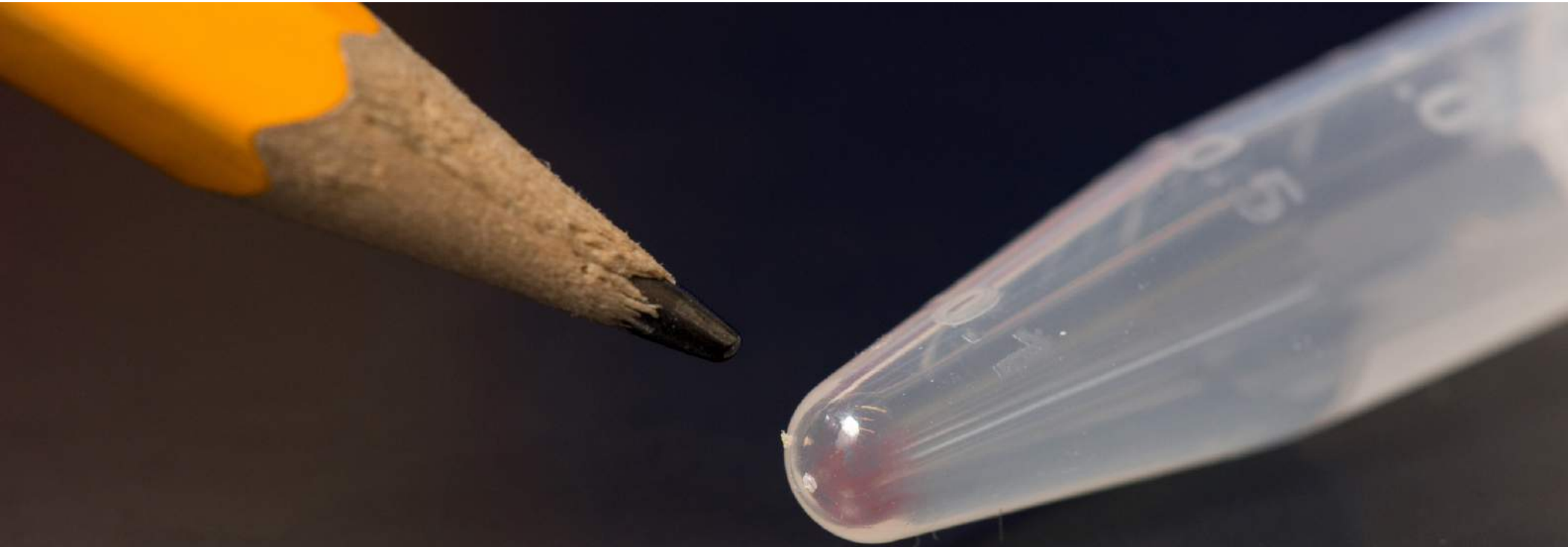






HD Video of OK Go's This Too Shall Pass — watched 100M+ times





~10 million copies of the HD movie

Photo: Tara Brown / UW





Thank you!

