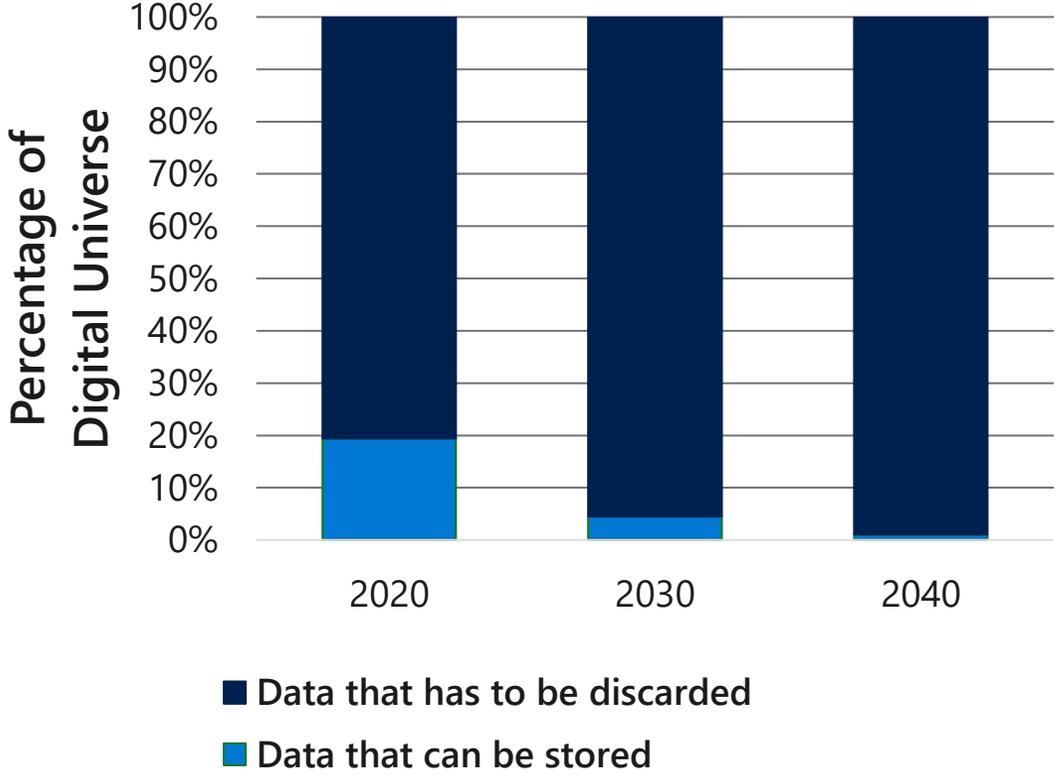
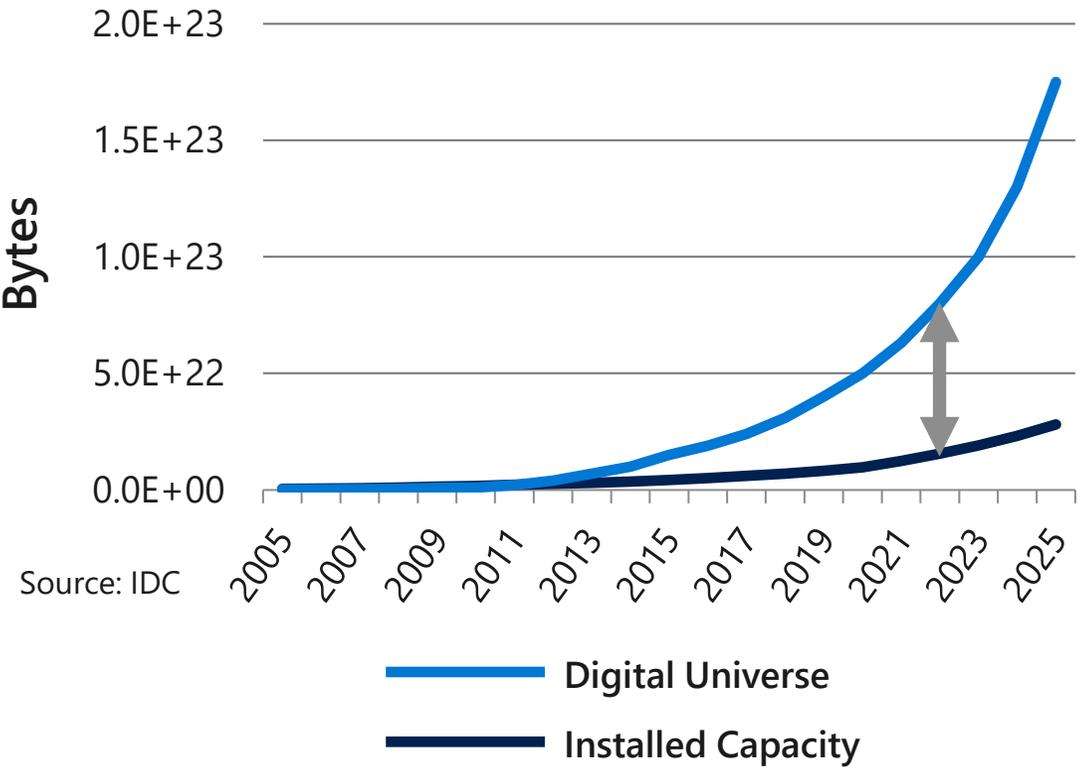


DNA Data Storage and Near-Molecule Processing for the Yottabyte Era

Karin Strauss, Microsoft Research
Luis Ceze, University of Washington



Storage capacity is growing too slowly



When Moore met Feynman

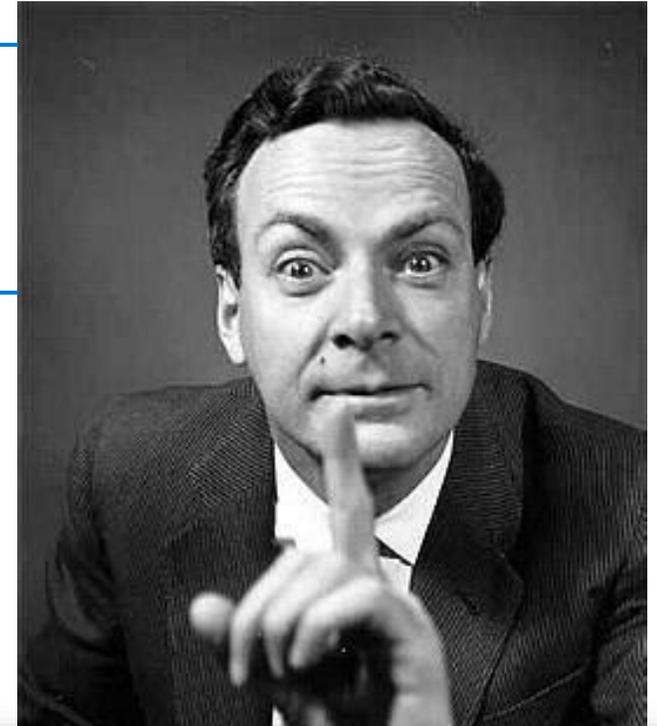


The number of transistors doubles every 18-24 months

The industry roadmaps are based on that continued rate of improvement

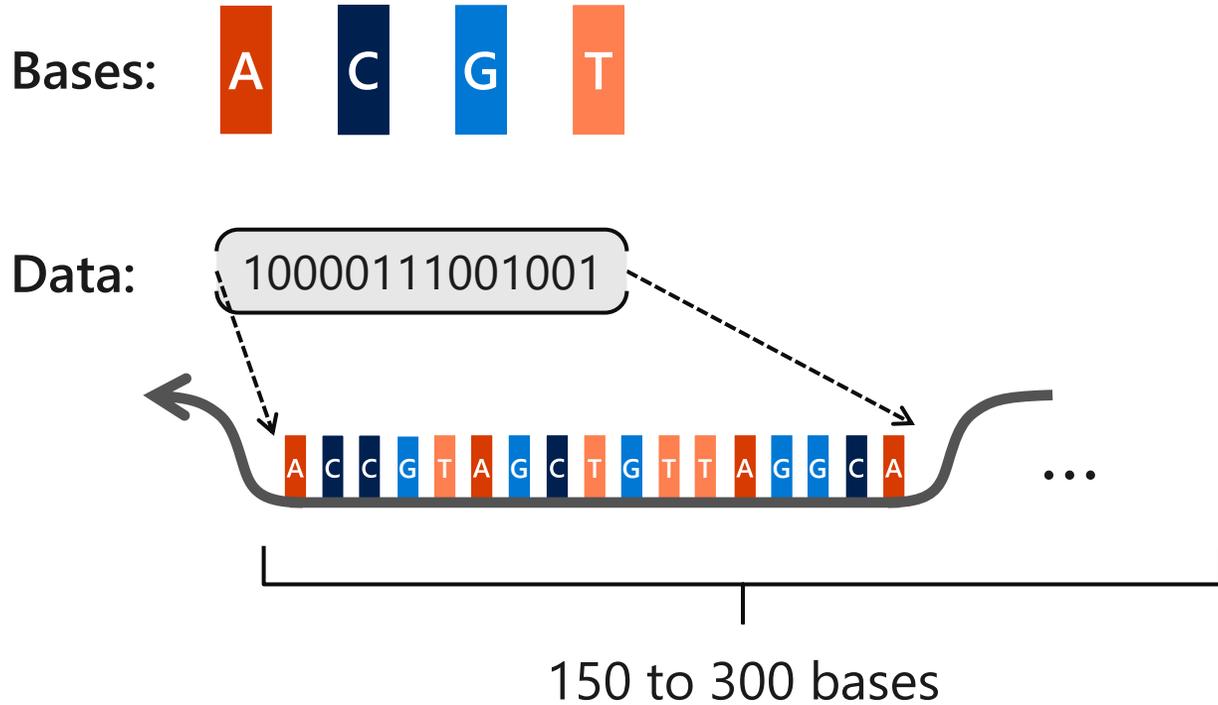
Arrange the atoms the way we want

DNA molecules use approximately 50 atoms for one bit



Let's store data in DNA!

DNA data storage basics



Simple mapping:

Bits	Base
00	A
01	C
10	G
11	T

Store data in synthetic DNA strands



Credit: Tara Brown Photography/University of Washington

Dense, really dense

Cold Storage: 1EB

Size: Two Walmart Supercenters



VS.



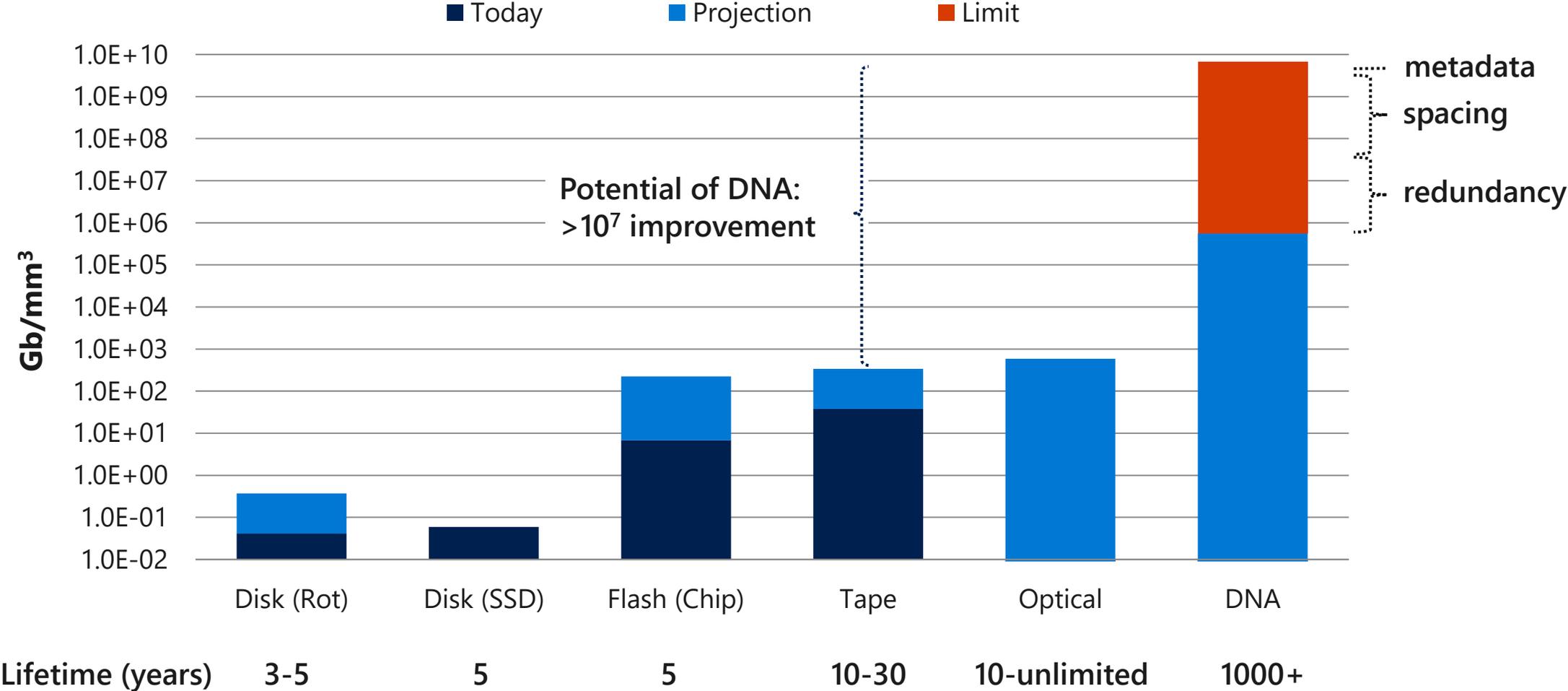
Information durability

DNA "synthetic fossils" last 2,000 – 2,000,000 years



Extreme density makes these conditions cheap and easy to keep

Comparison with other media



No obsolescence



Size of a mainframe
800 Kbases/day



Size of a workstation
80 Gbases/day



Size of a portable SSD
10 Gbases/day

Same medium as read technology improves:



No obsolescence issue, DNA will always be relevant

Same medium as read technology improves:

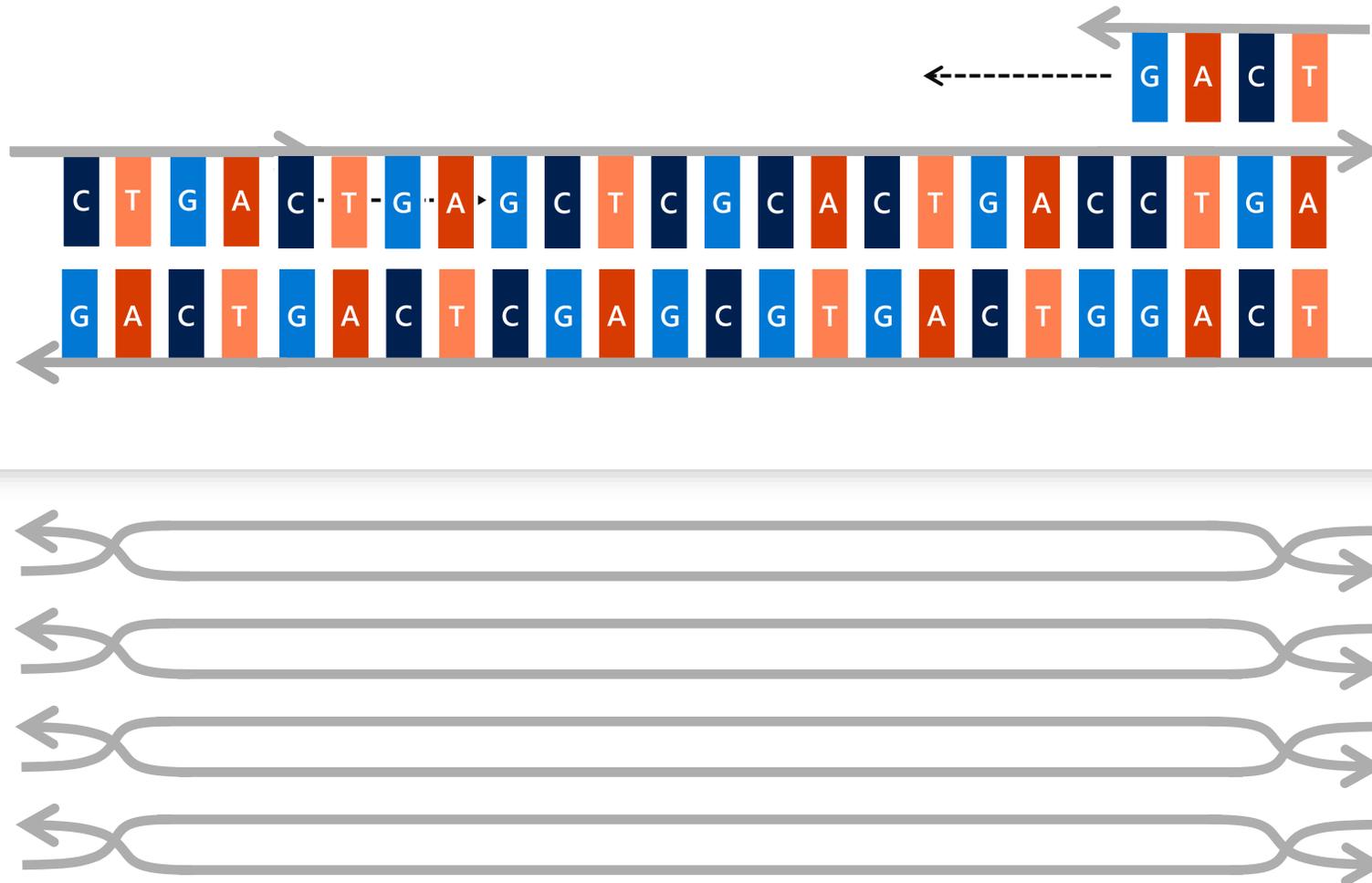


Medium changes as read and write technology improves:

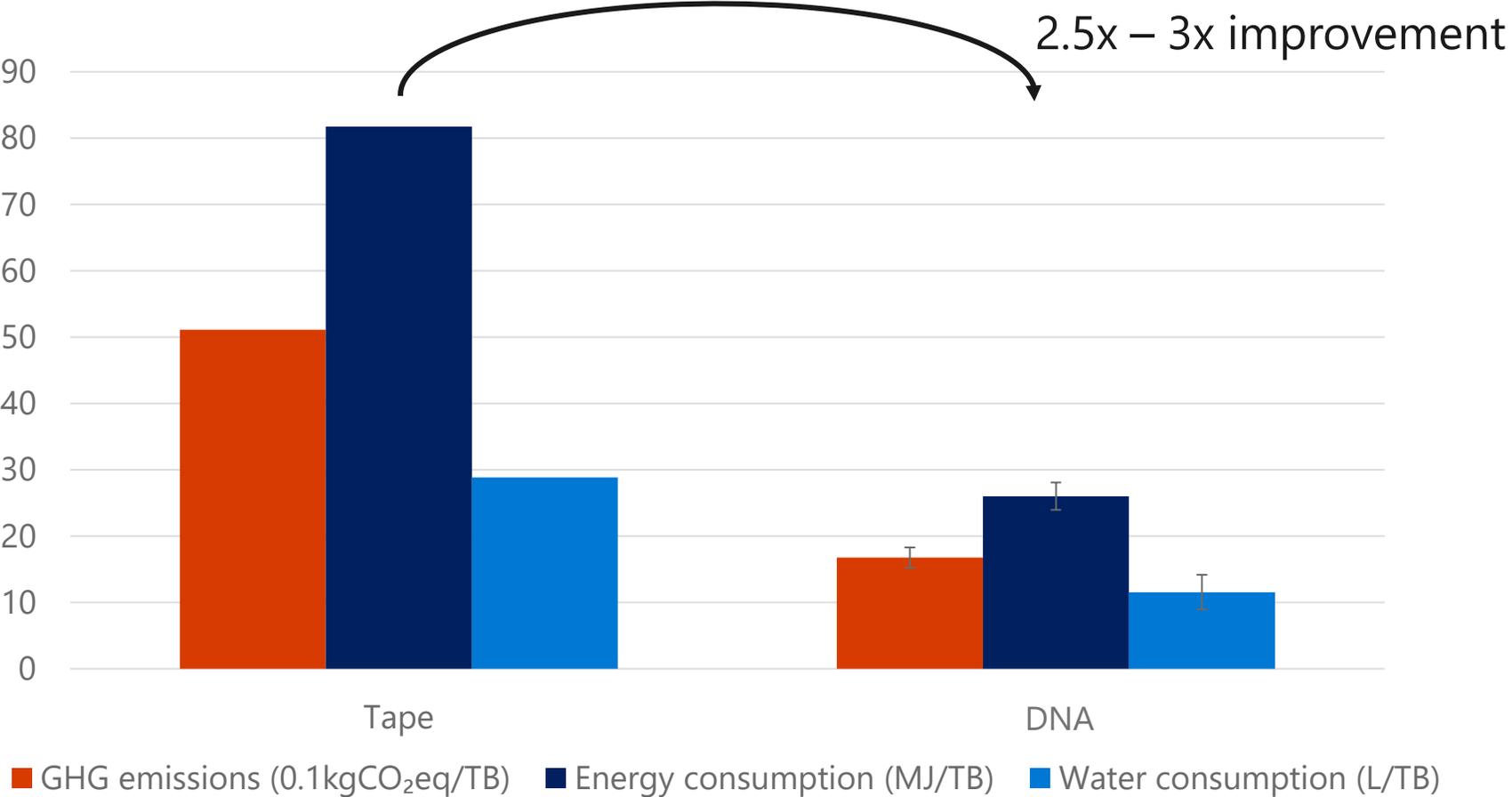


Ability to make copies

Polymerase Chain Reactions (PCR) create copies exponentially



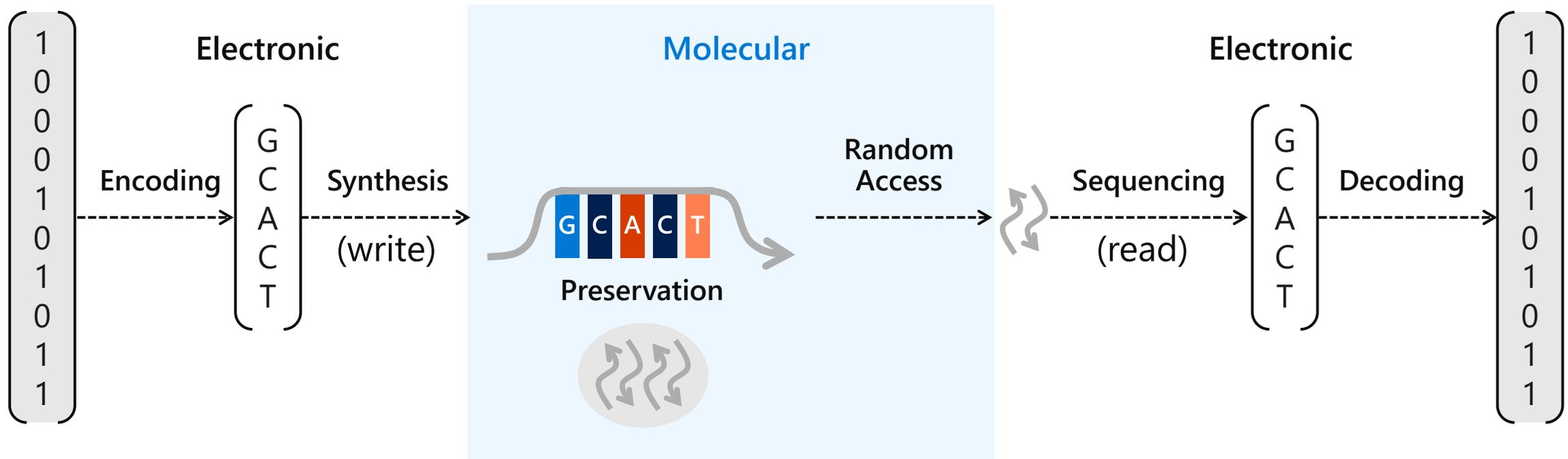
Sustainability



▶ DNA promises to be significantly more sustainable than tape

Nguyen et al., Electronics Goes Green, 2020

DNA storage end-to-end system

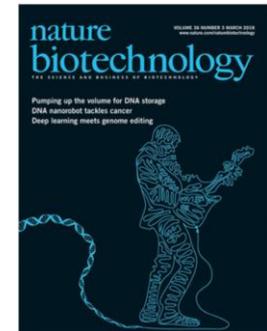


Our results so far

1GB of data stored and fully recovered

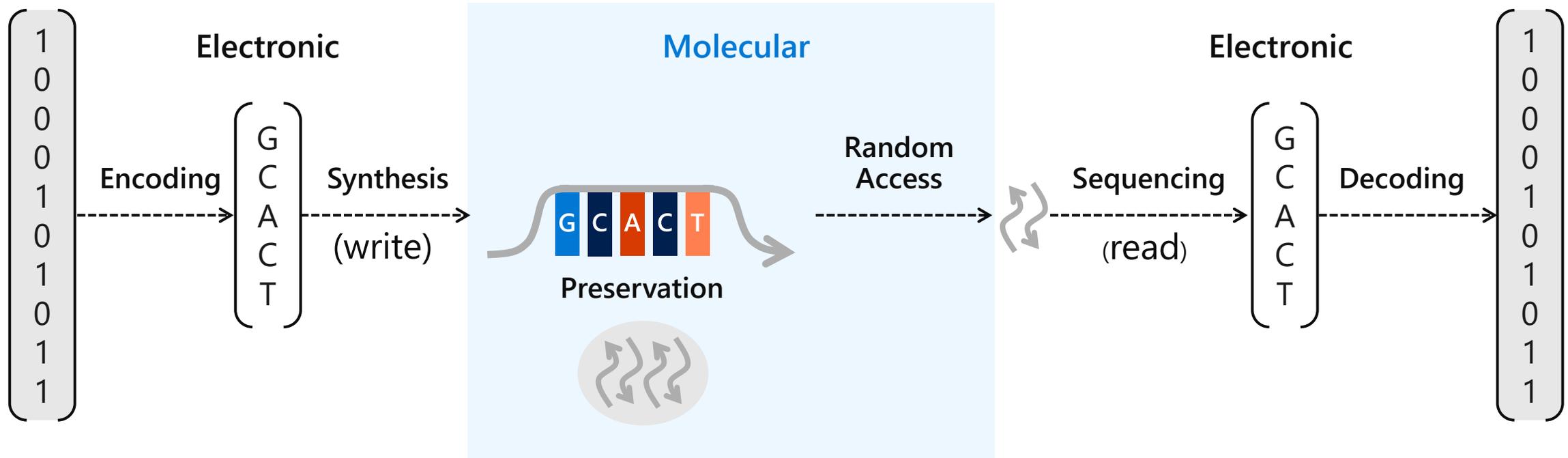


Published



Most data in DNA in peer reviewed publication

DNA storage end-to-end system



DNA encoding



Photo credit: OK Go.

0101000101011100 1101100001111000 0111010111001000 1001010001001001



1 of 4



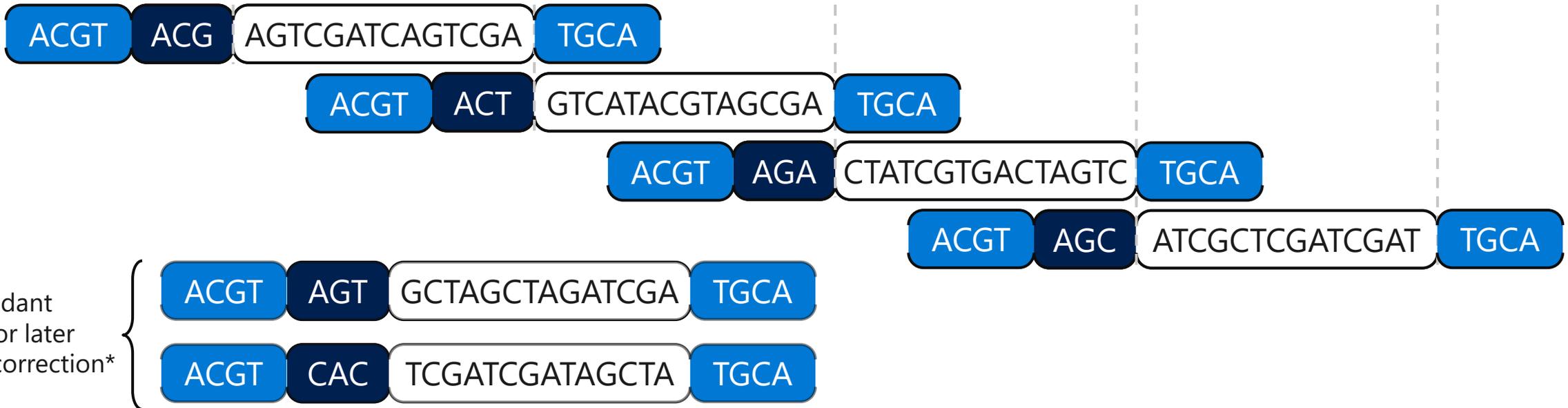
2 of 4



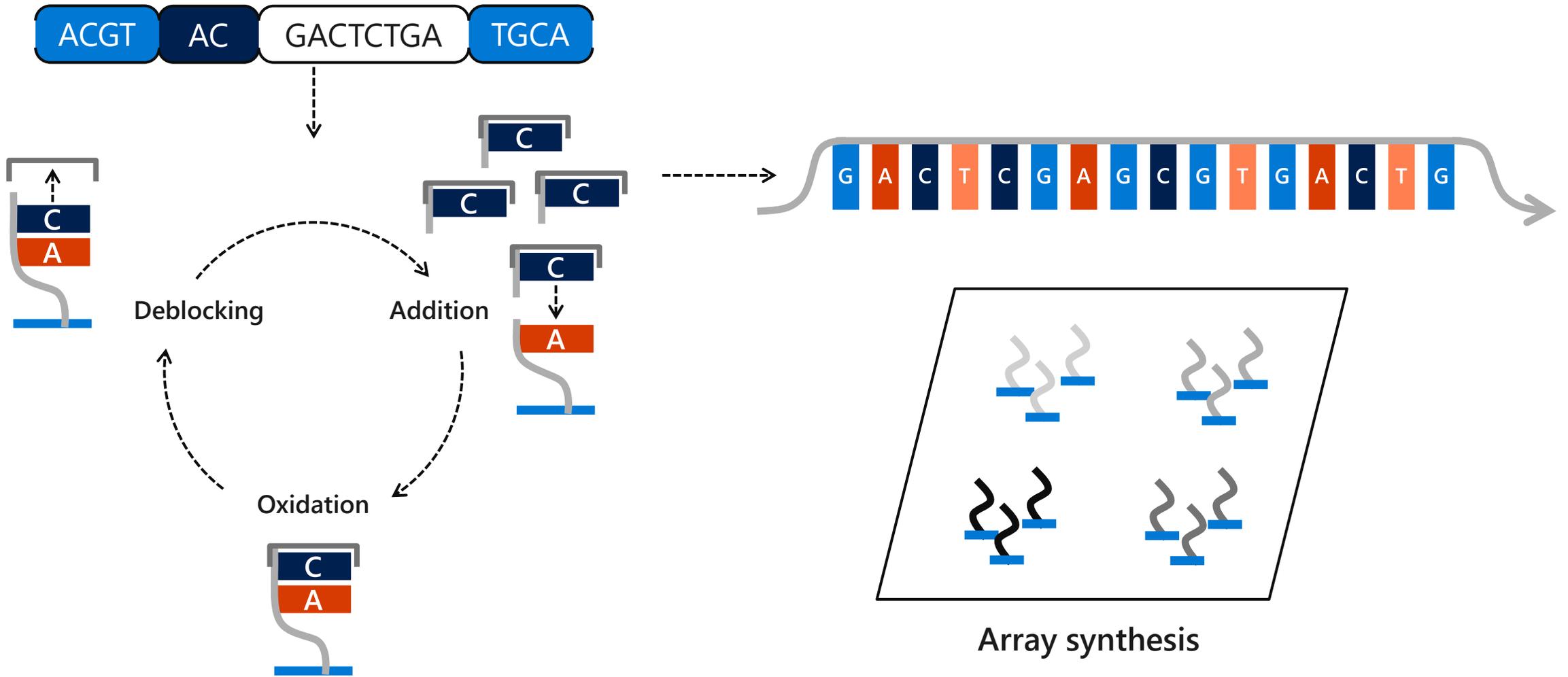
3 of 4



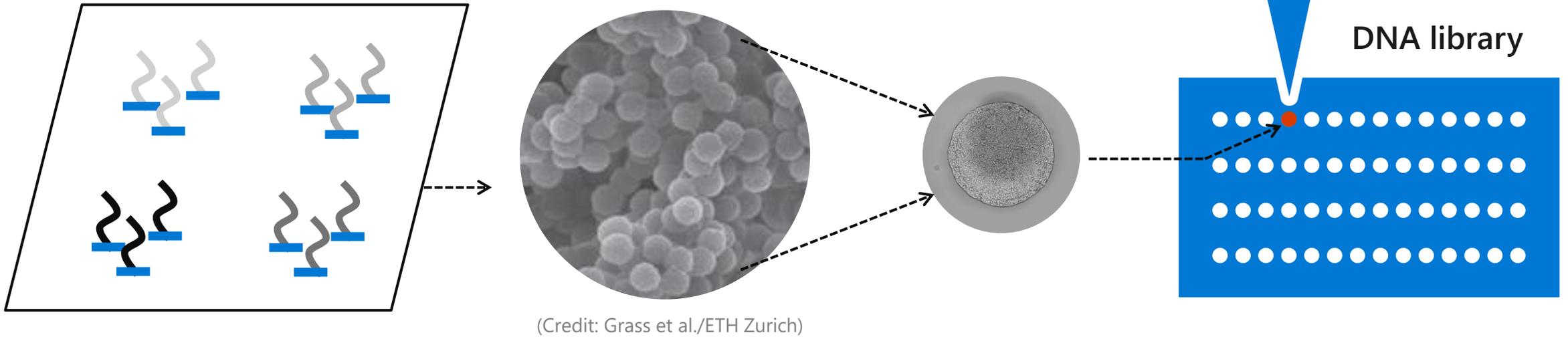
4 of 4



DNA synthesis

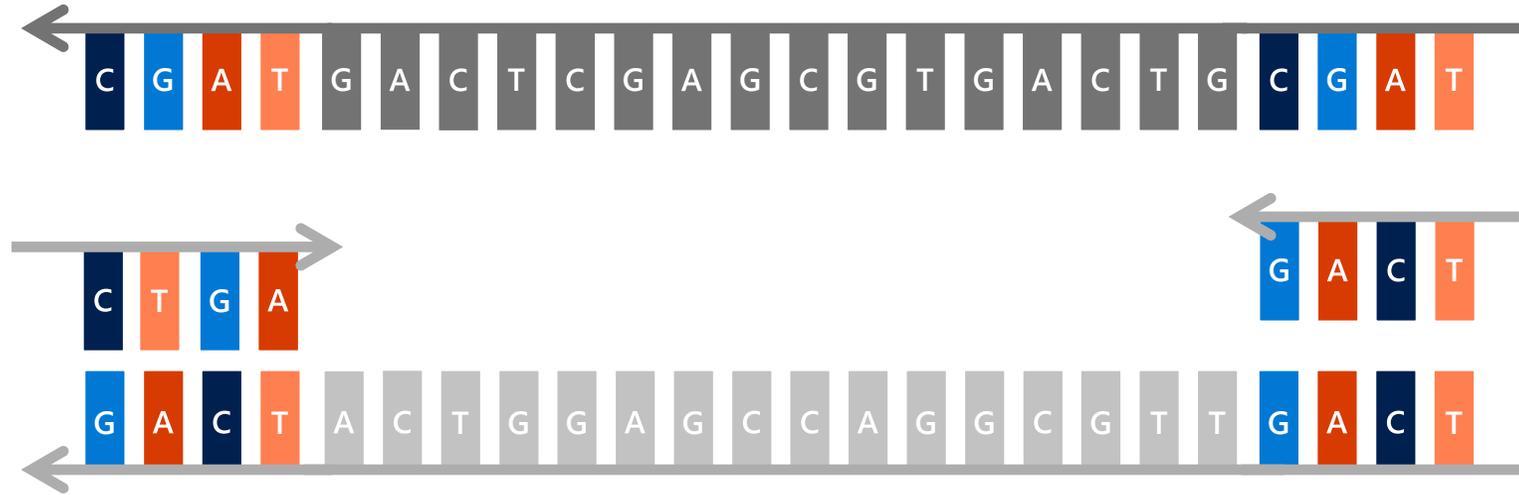


DNA preservation

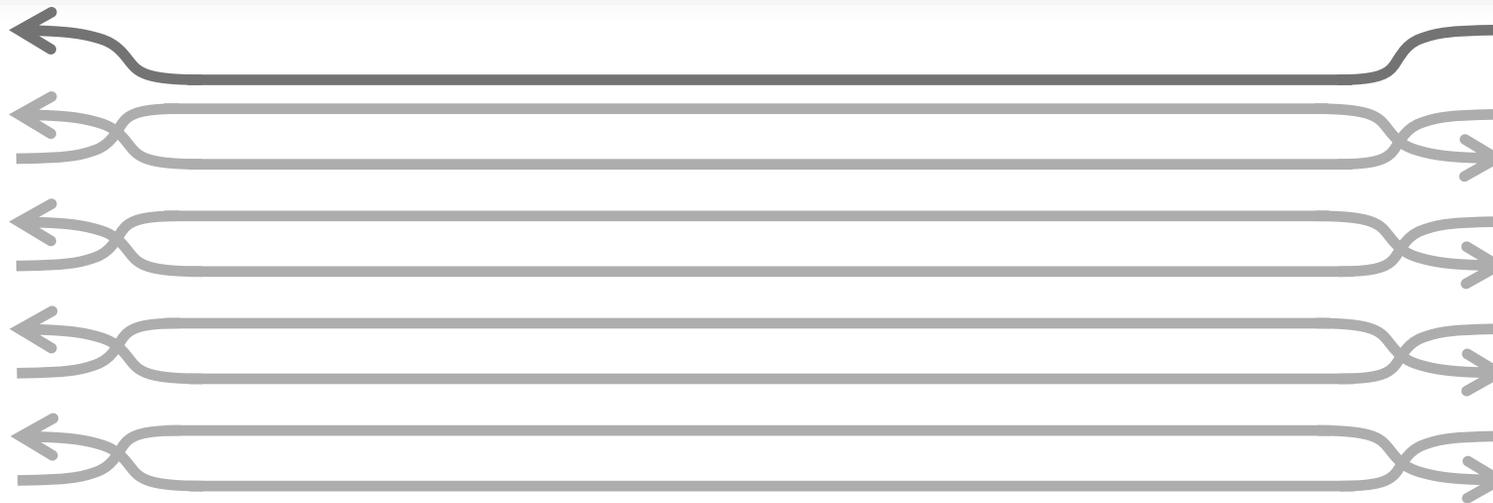


DNA storage random access with PCR

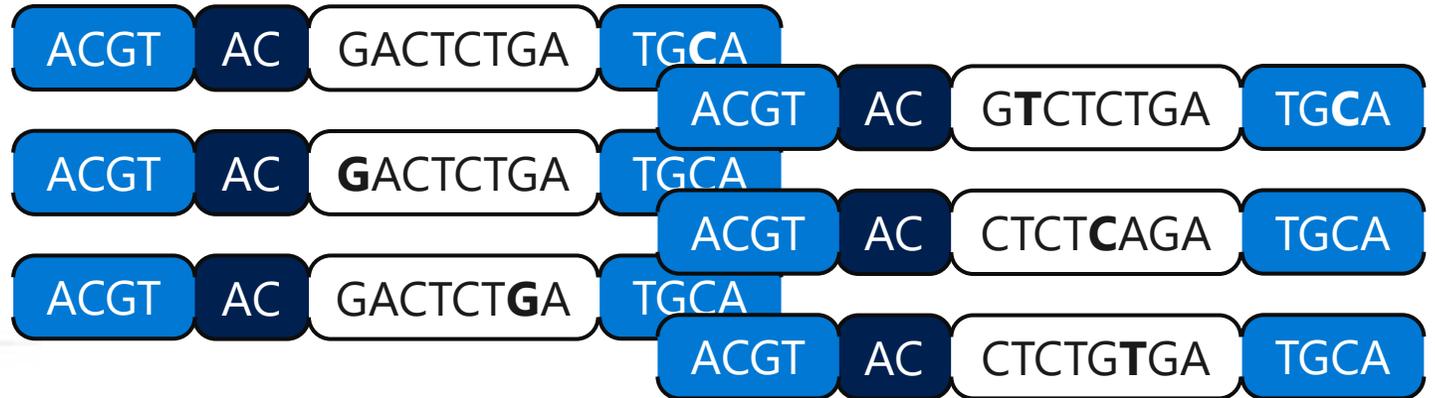
Selecting one item out of two



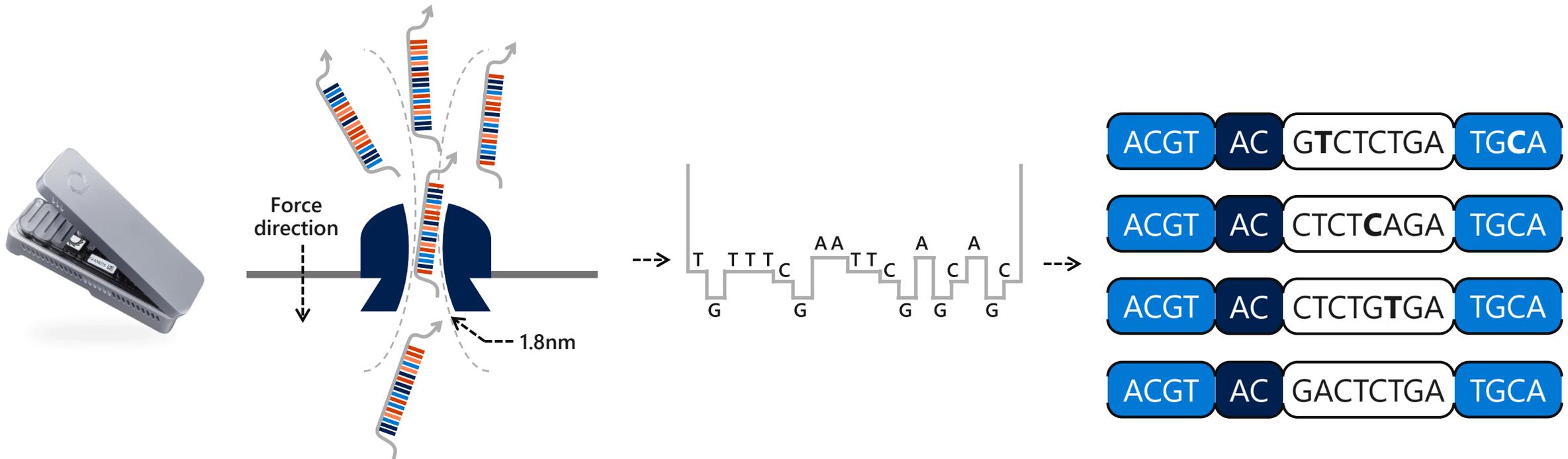
PCR primers are used to access units of storage individually



Reading DNA with sequencing by synthesis

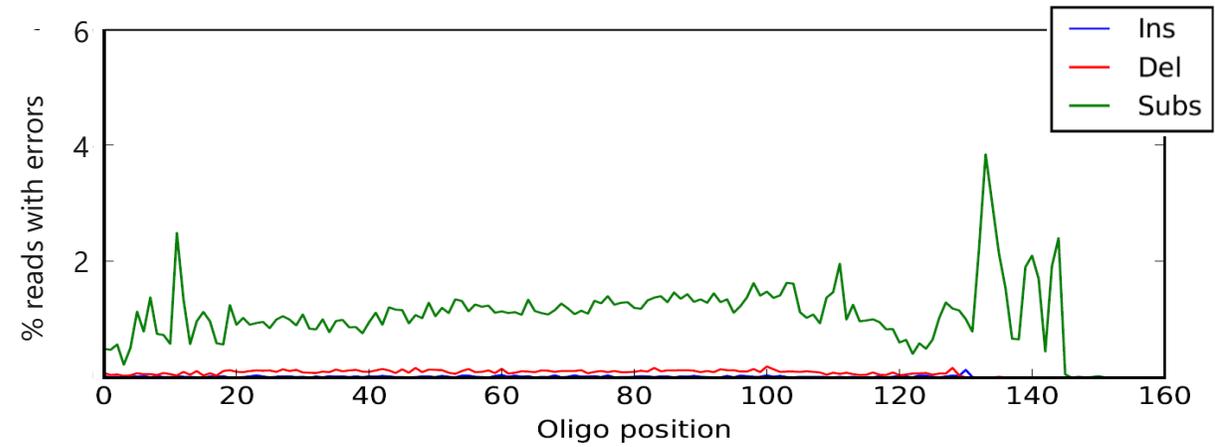


Reading DNA with nanopores

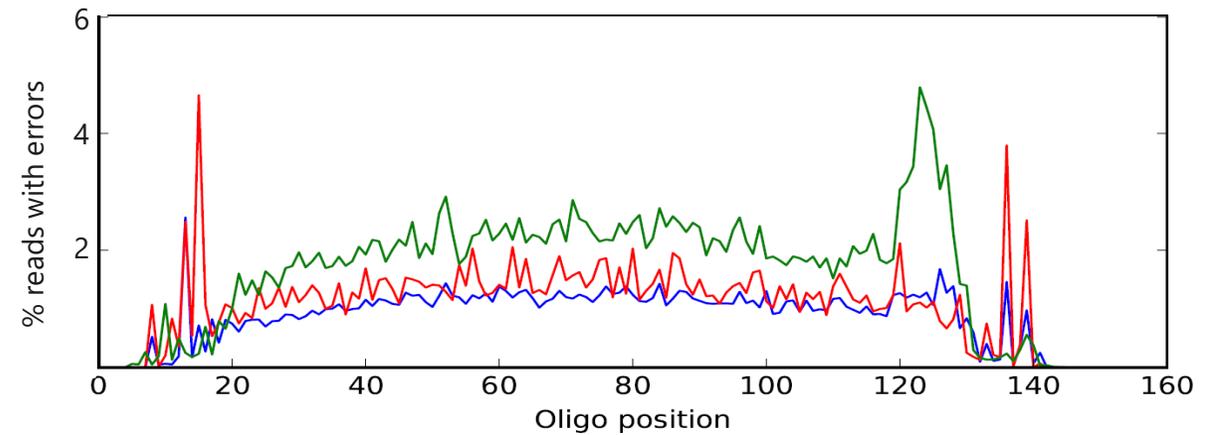


Different error profile across platforms

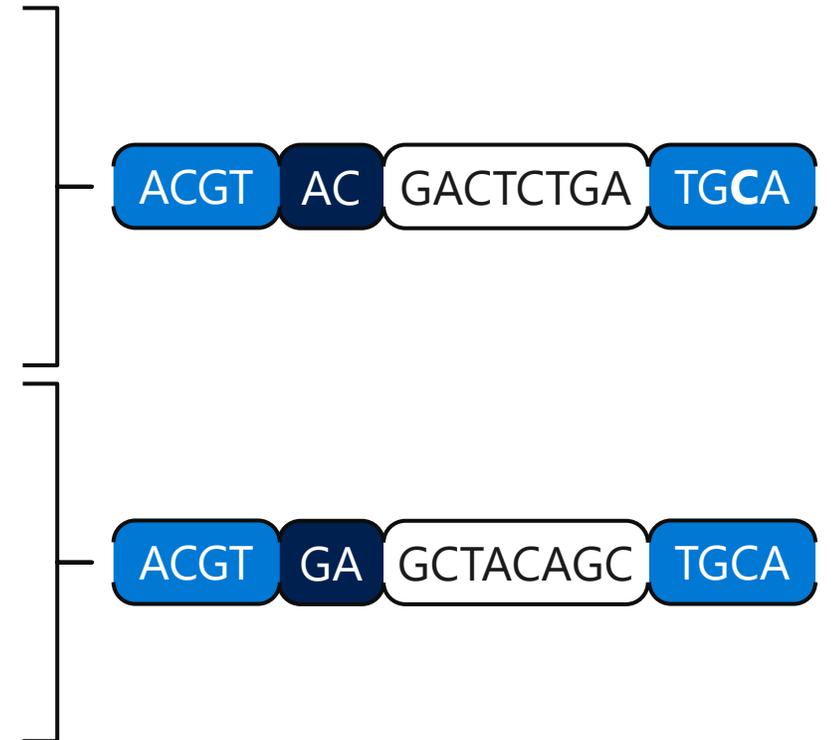
Illumina NextSeq



ONT MinION



DNA decoding and error correction



DNA decoding and error correction

ACGT AC GACTCTGA TGCA



0001 0101000101011100

ACGT GA GCTACAGC TGCA



1100 1001010001001001

0001 0101000101011100

0010 1101100001111000

0011

0100 1001010001001001

1010 0111010111001000

1100 1001010001001001

0001 0101000101011100

0010 1101100001111000

0011 0111010111001000

0100 1001010001001001

DNA decoding and error correction

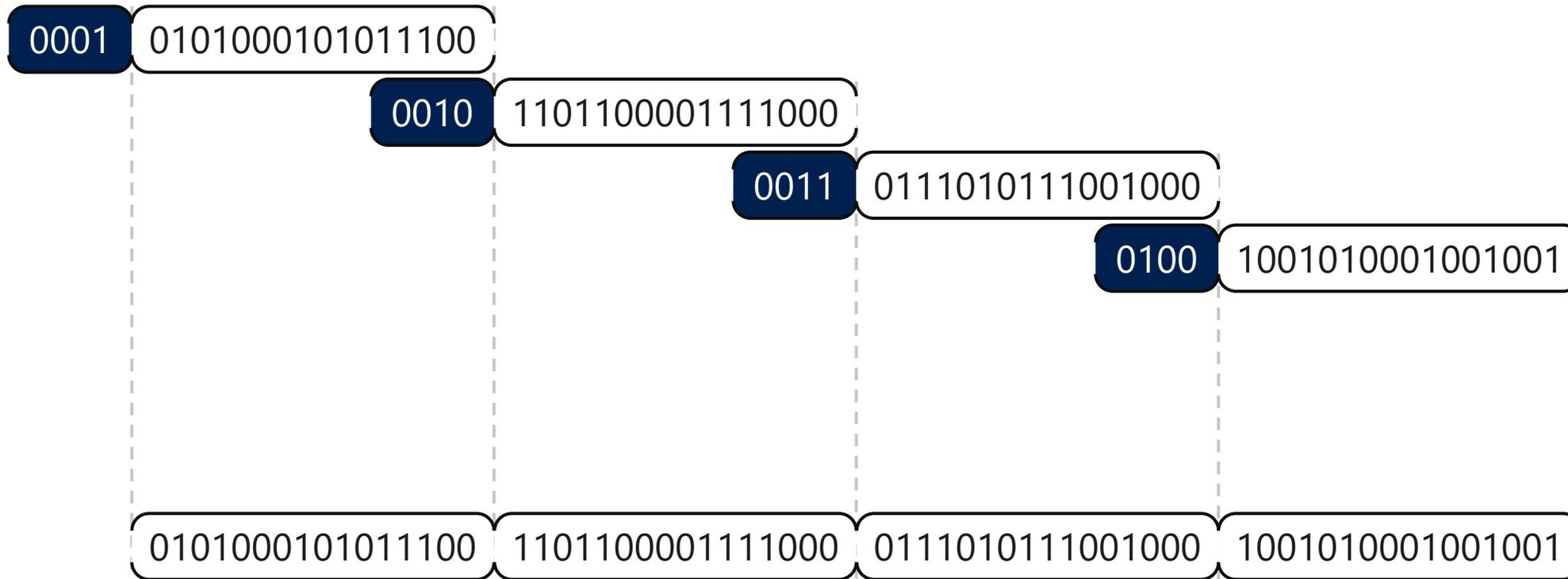
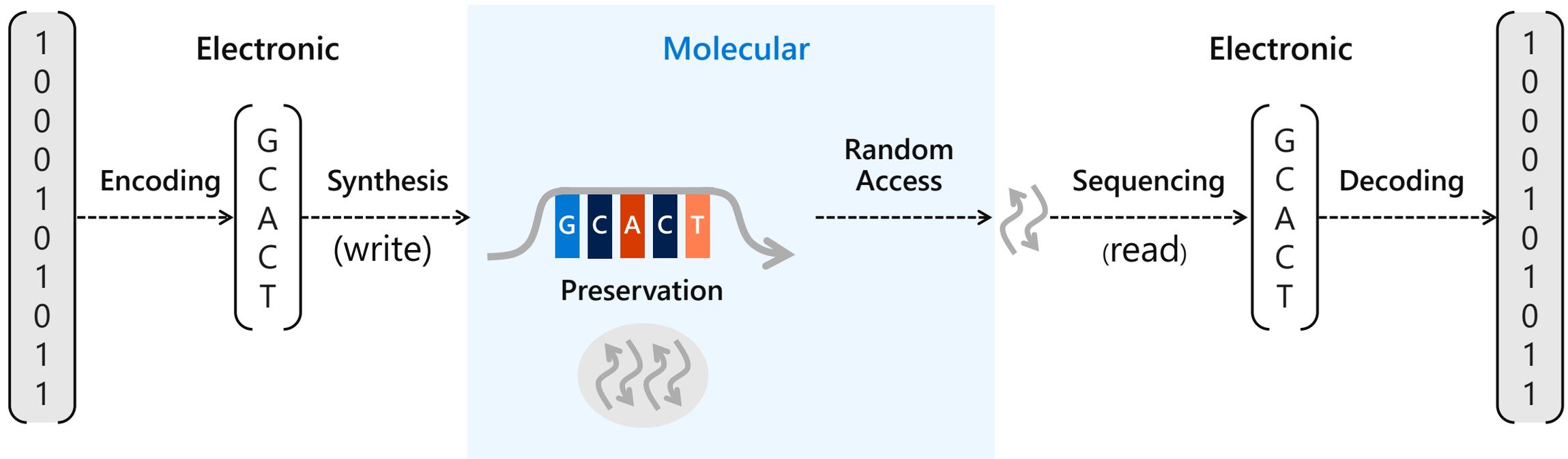
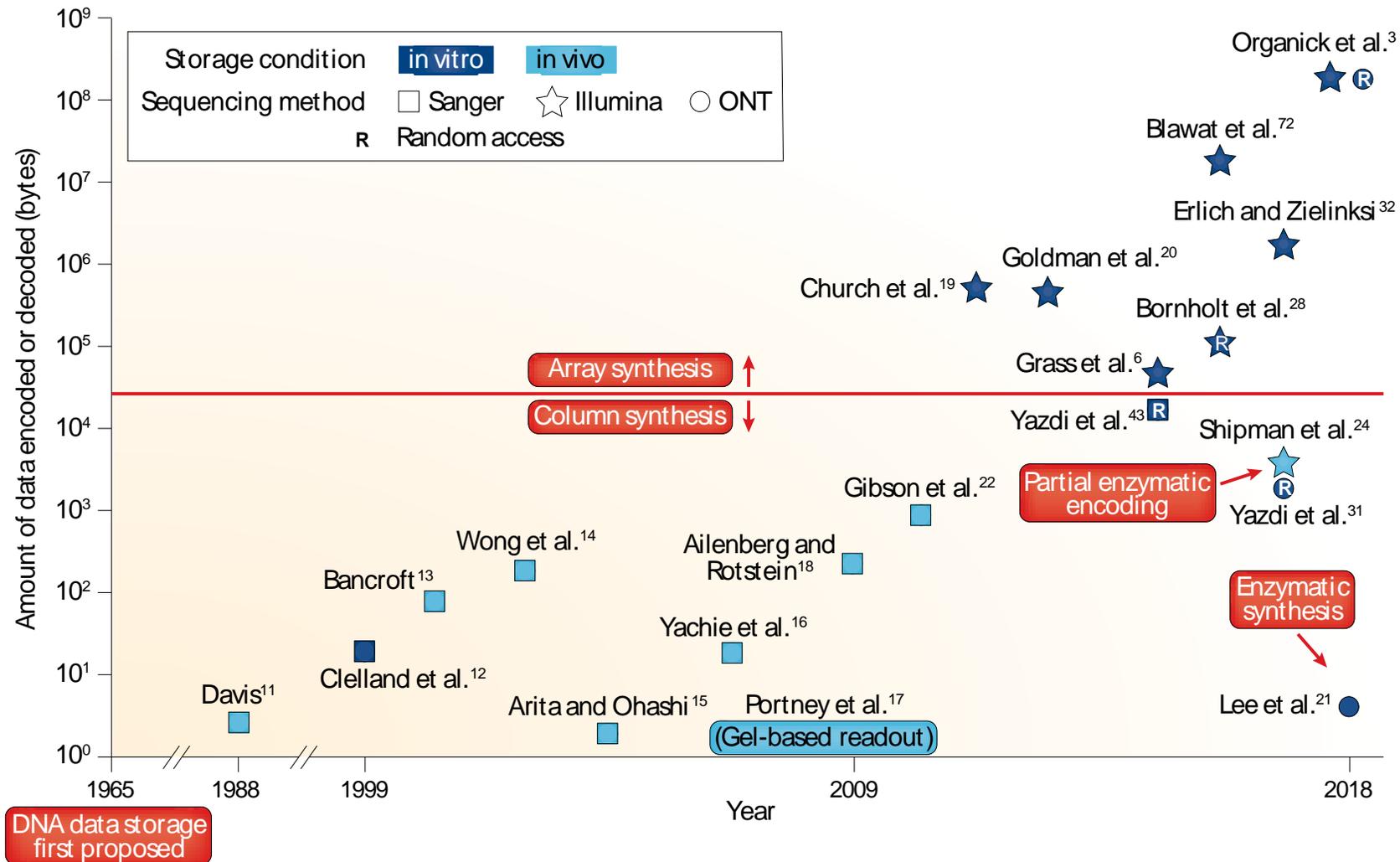


Photo credit: OK Go.

DNA storage end-to-end system

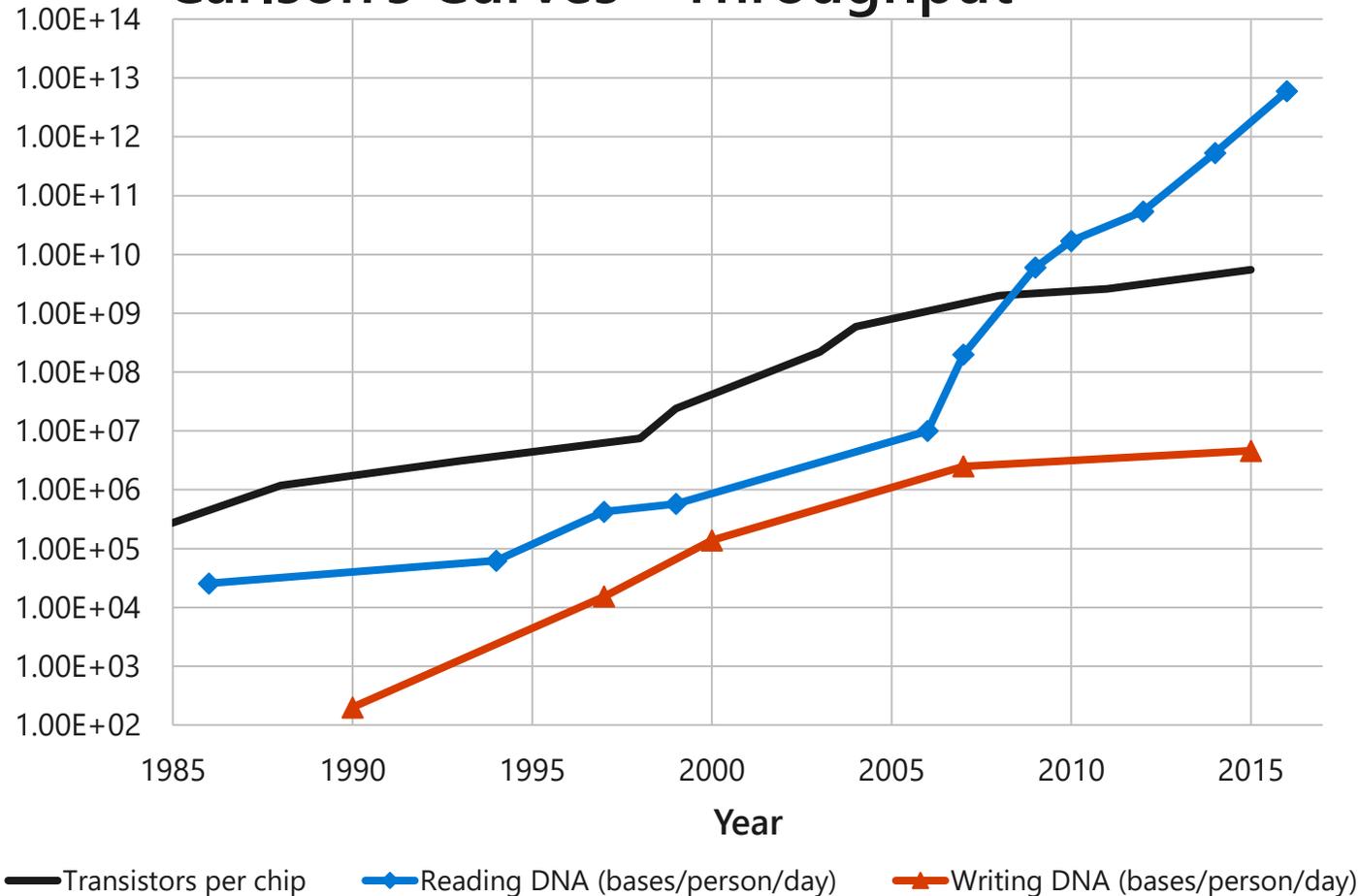


Exponential improvements in DNA data storage



Performance of reading and writing DNA

Carlson's Curves - Throughput



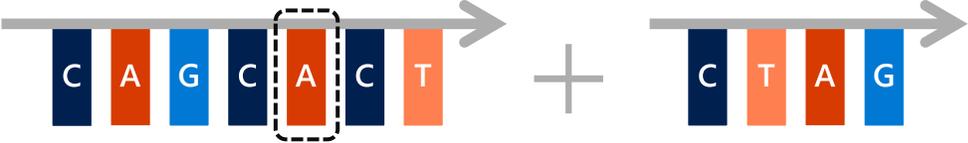
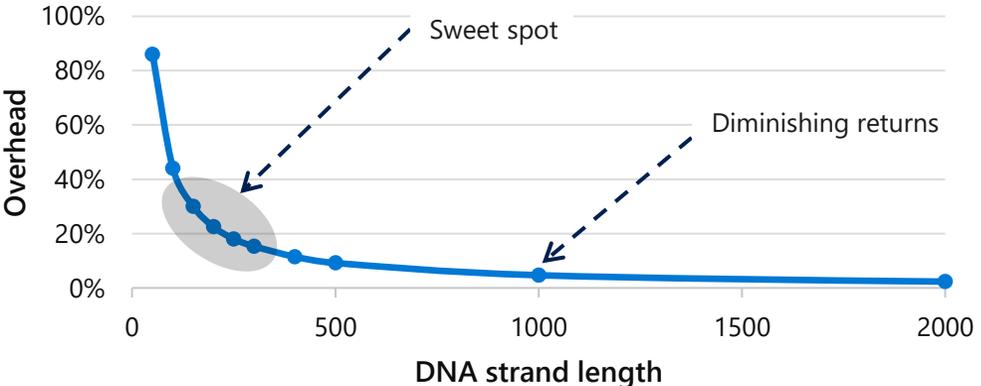
Latency

Synthesis and sequencing are currently batch processes, matches archival storage SLAs (~hours).

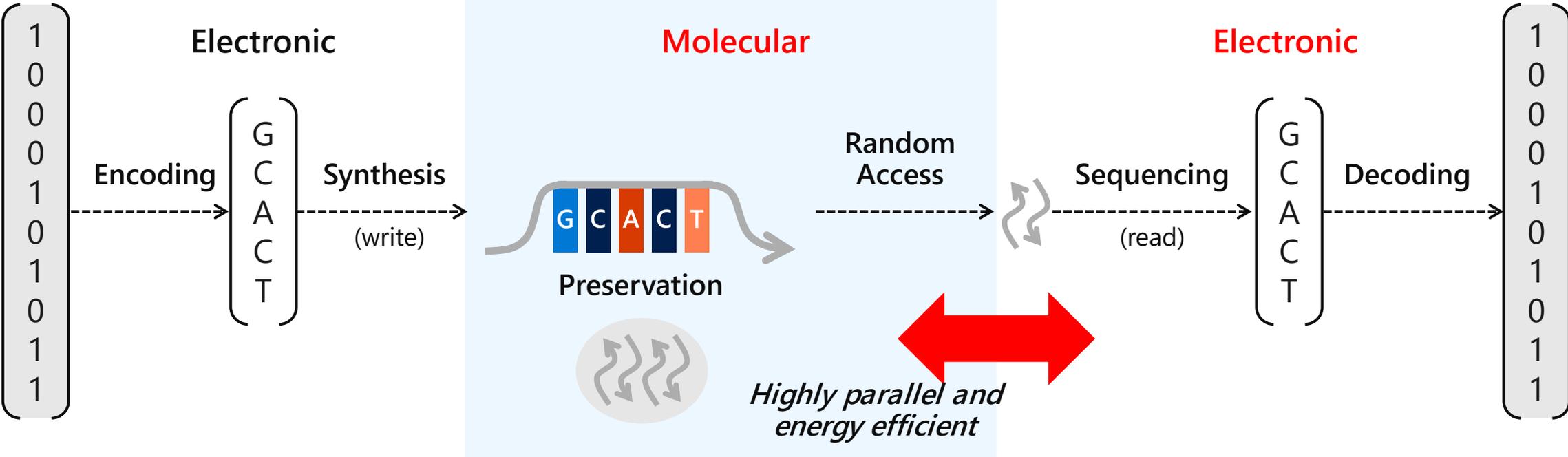
Emerging technologies, like nanopore devices, provide closer to real time latency

Write and read mechanisms

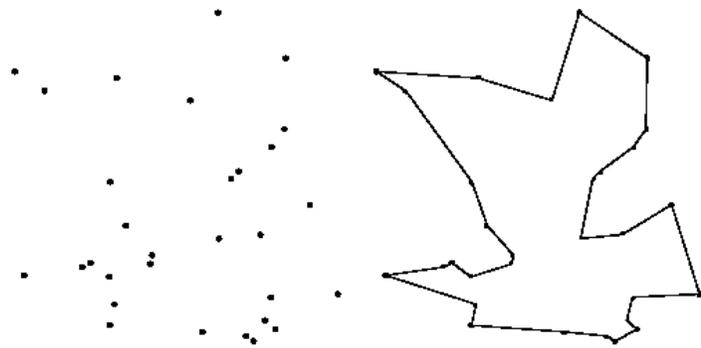
Opportunities for improvements in write and read throughput, latency and cost

Opportunities	Life sciences	Data storage																
Error rate	Single base mutations affect function	<p>Error correcting codes allow data recovery even in the presence of multiple errors</p>  <p>Error types: substitutions, deletions and insertions</p>																
Length ("block size")	Longer sequences have more function	<p>Shorter sequences are faster and easier to make</p>  <table border="1"><caption>Approximate data from the Overhead vs DNA strand length graph</caption><thead><tr><th>DNA strand length</th><th>Overhead (%)</th></tr></thead><tbody><tr><td>100</td><td>85</td></tr><tr><td>200</td><td>45</td></tr><tr><td>300</td><td>30</td></tr><tr><td>400</td><td>20</td></tr><tr><td>500</td><td>15</td></tr><tr><td>1000</td><td>10</td></tr><tr><td>2000</td><td>5</td></tr></tbody></table>	DNA strand length	Overhead (%)	100	85	200	45	300	30	400	20	500	15	1000	10	2000	5
DNA strand length	Overhead (%)																	
100	85																	
200	45																	
300	30																	
400	20																	
500	15																	
1000	10																	
2000	5																	

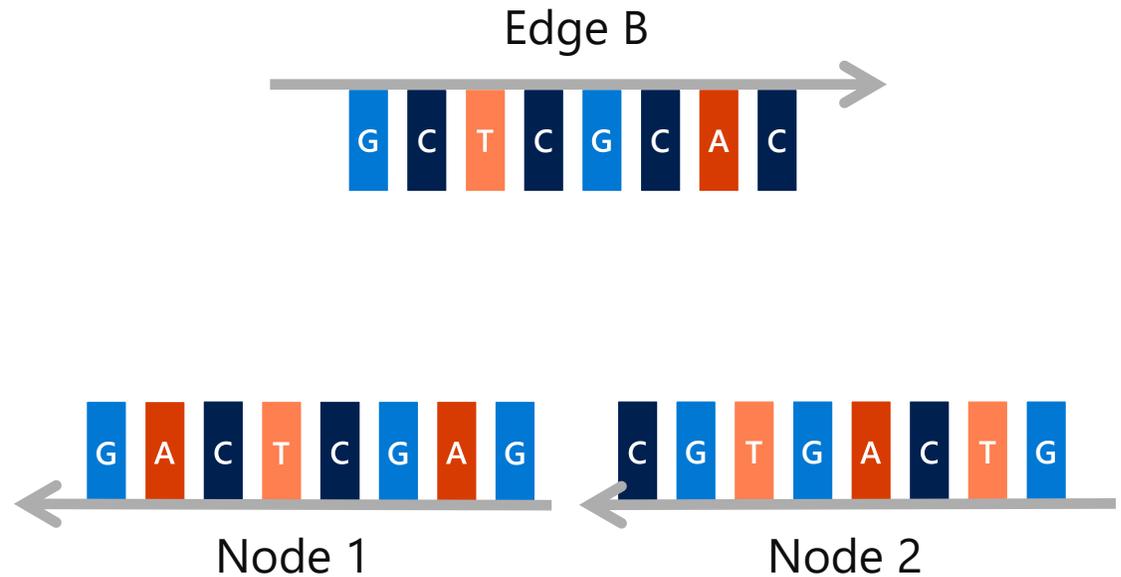
DNA storage end-to-end system w/ integrated computing



DNA computing in the 80s

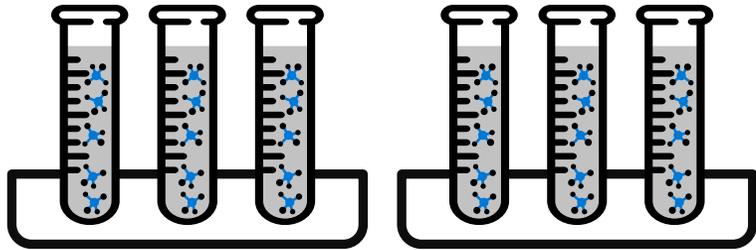


Hamiltonian path problem



Problem: shifts complexity from time to amount of material

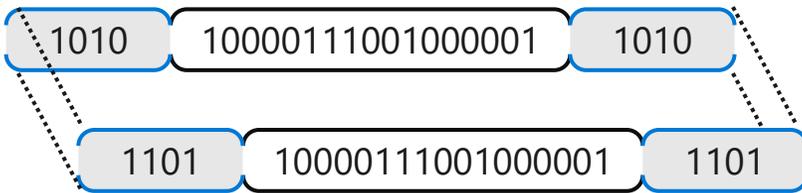
DNA "computing" in the age of big data



Operate over data already stored in DNA

Target polynomial time algorithms

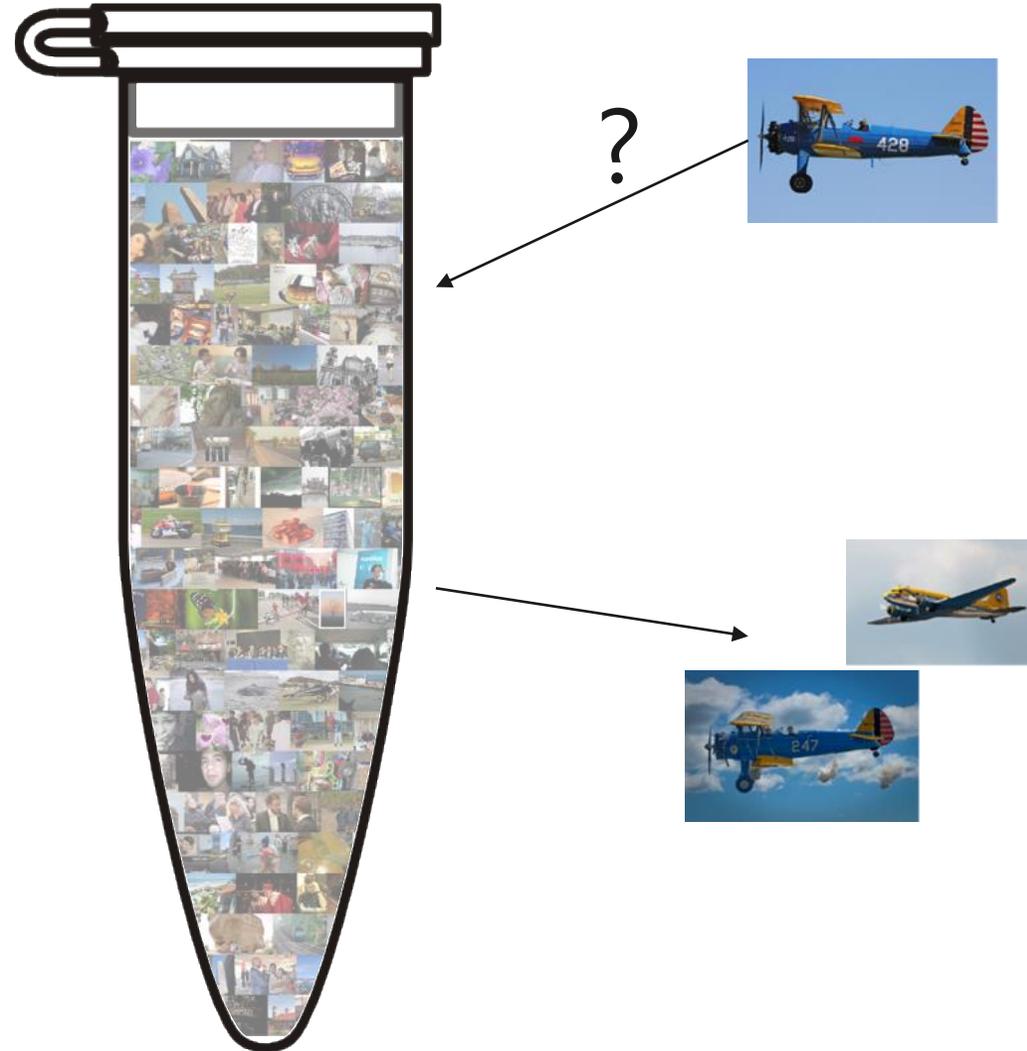
Extremely parallel and energy efficient



Content-based image/video search

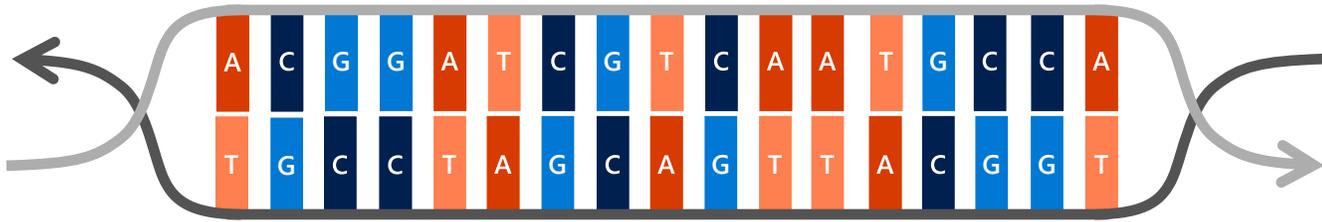


Content-based image/video search in DNA

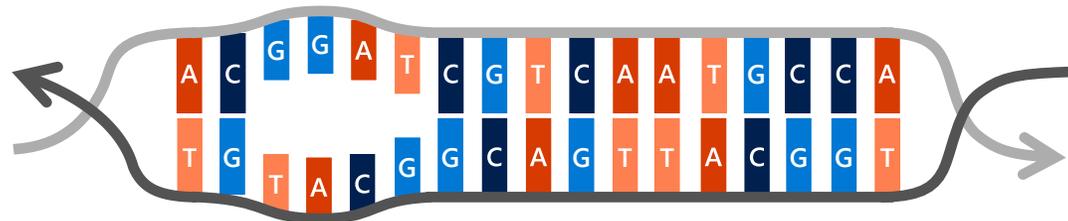


Exploiting matches for exact and approximate search

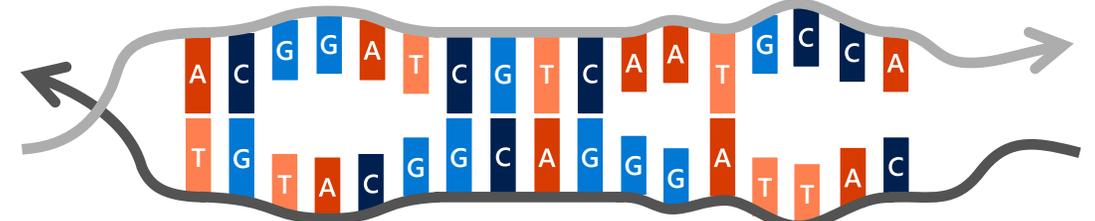
Double helix: complete match



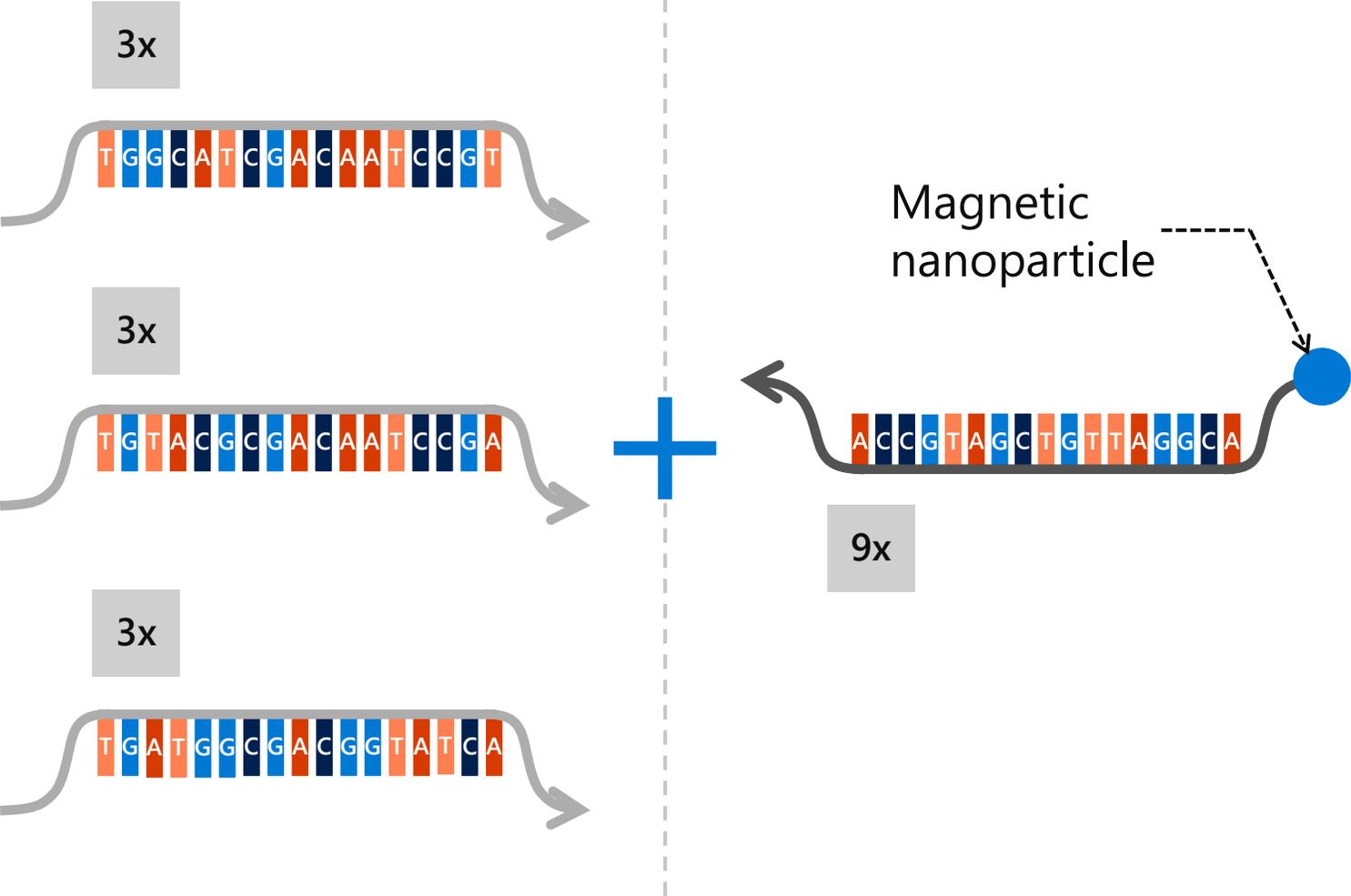
Good partial match



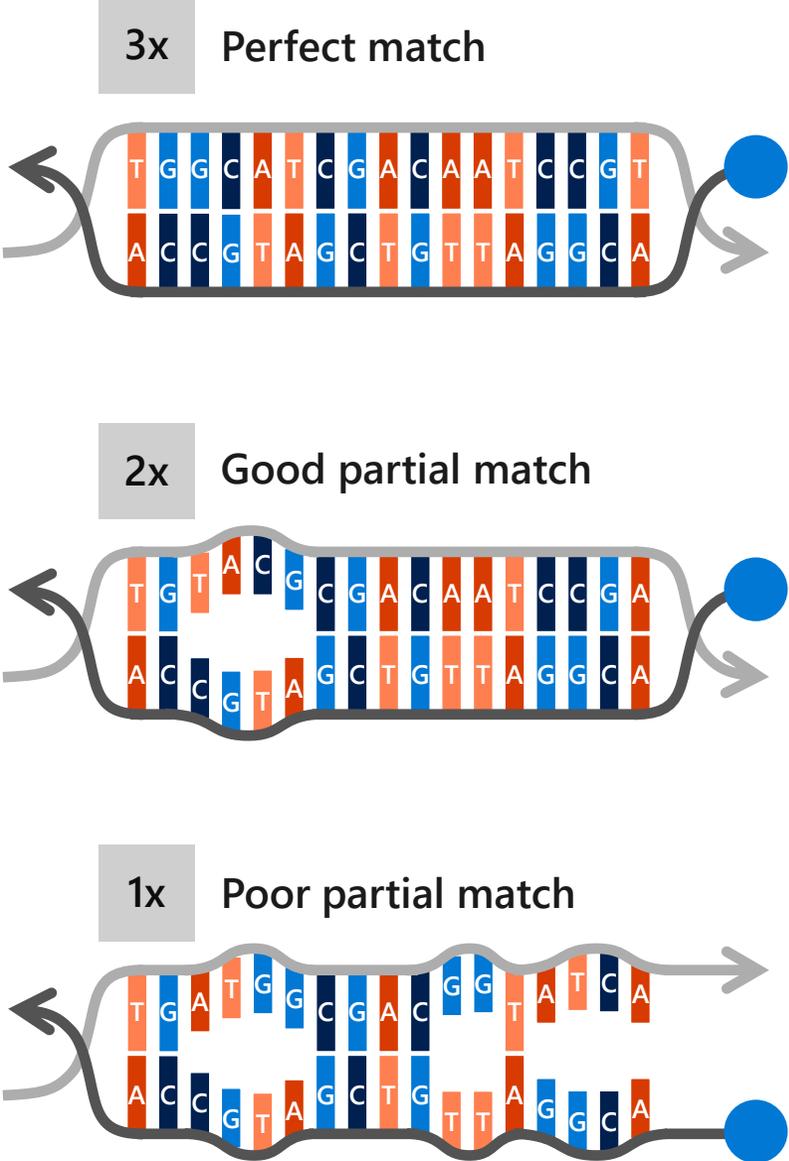
Poor partial match



Searching with DNA

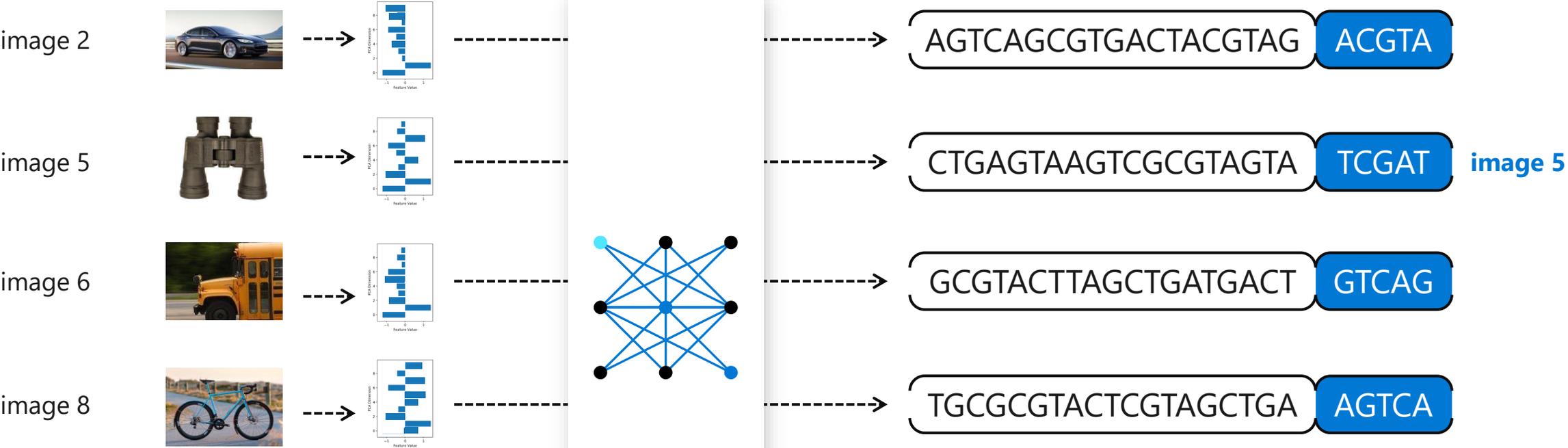


Match-dependent yield



Content based media search

Database/ training

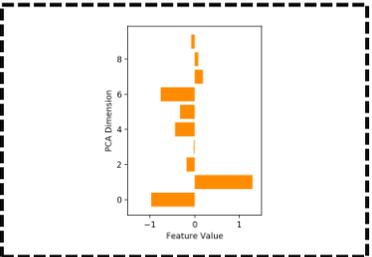


Query/ inference



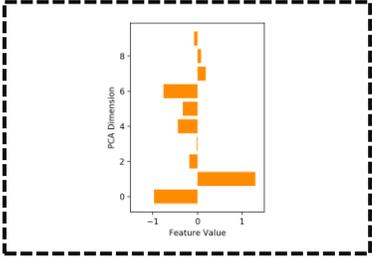
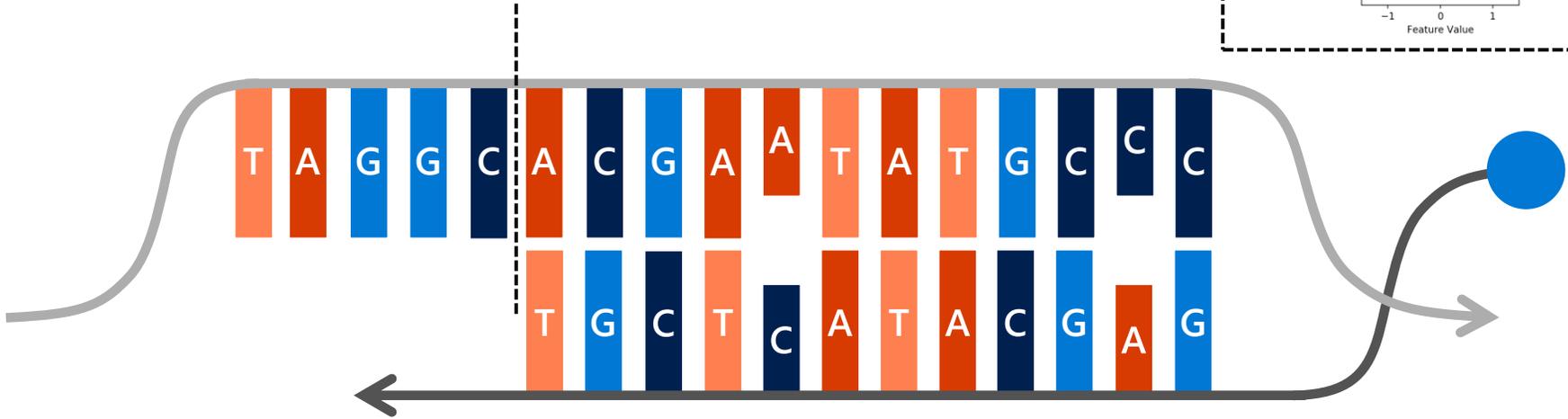
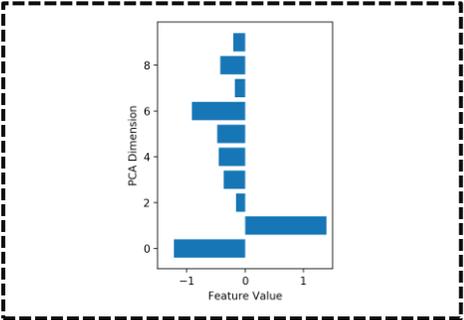
Image 5

TAGGC ACGAATATGCCC

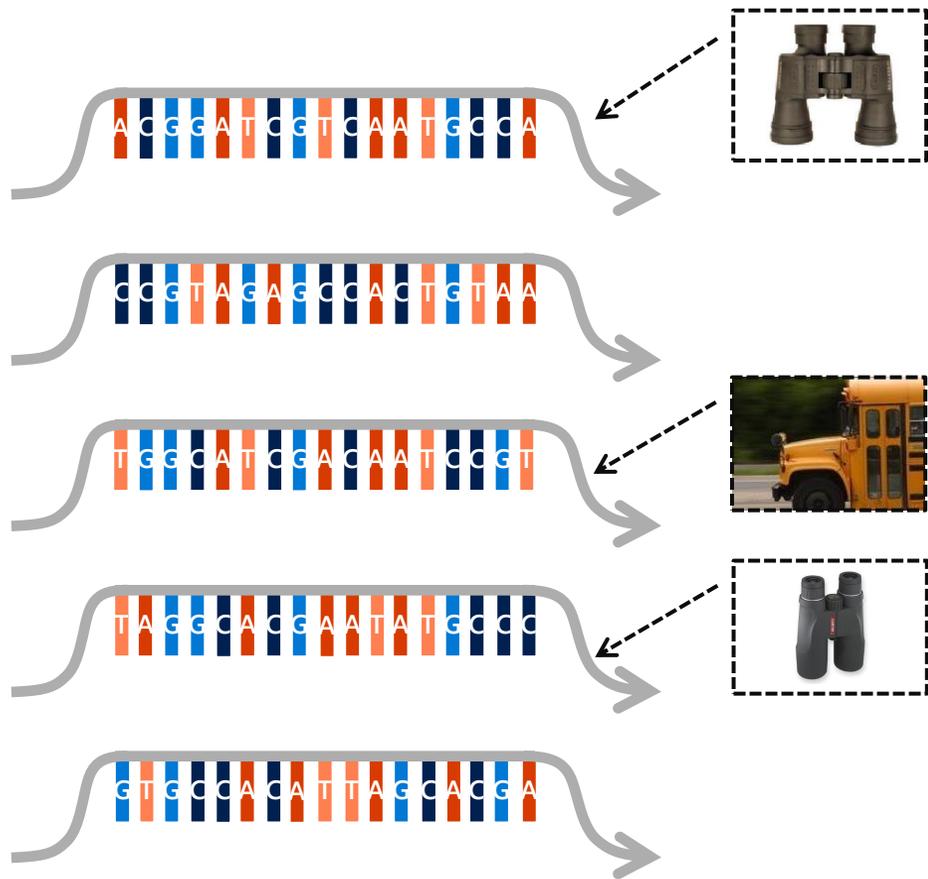


CTCATACGAG

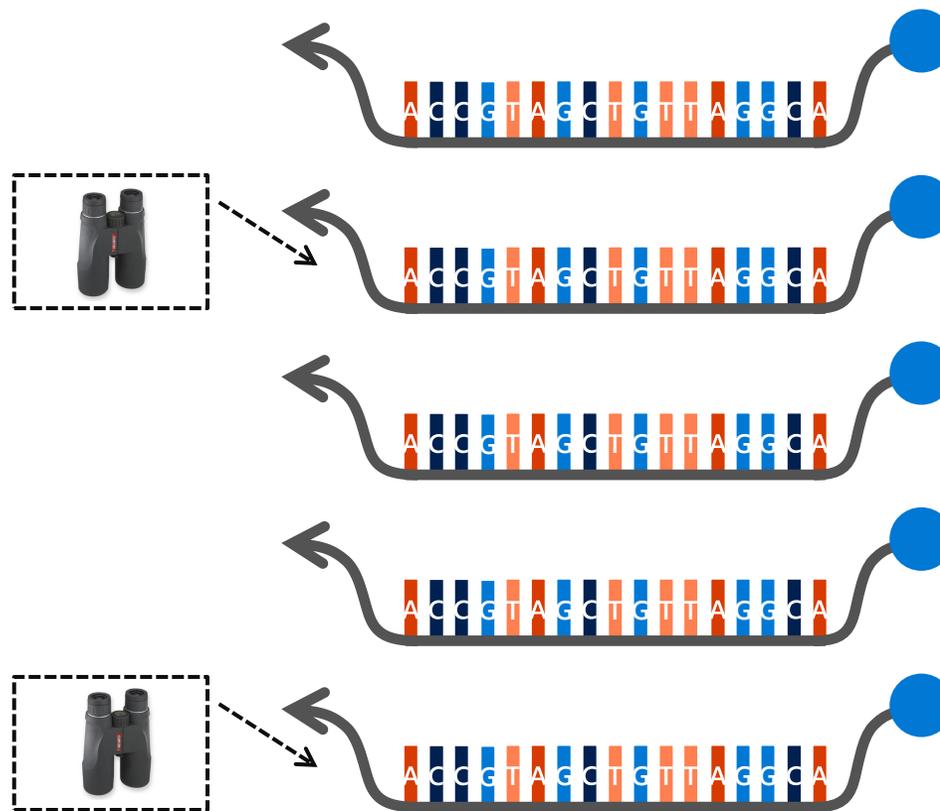
Image 5

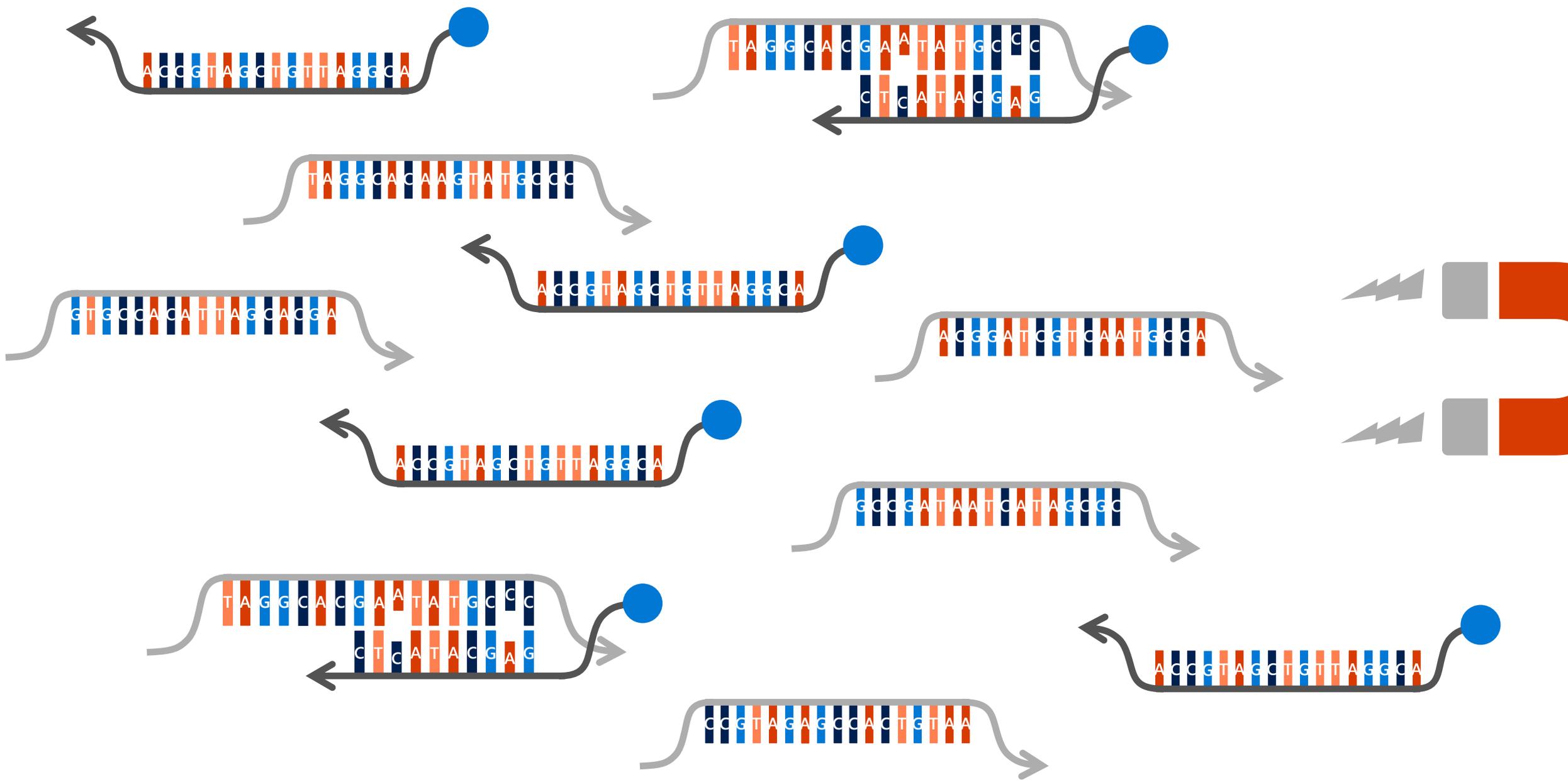


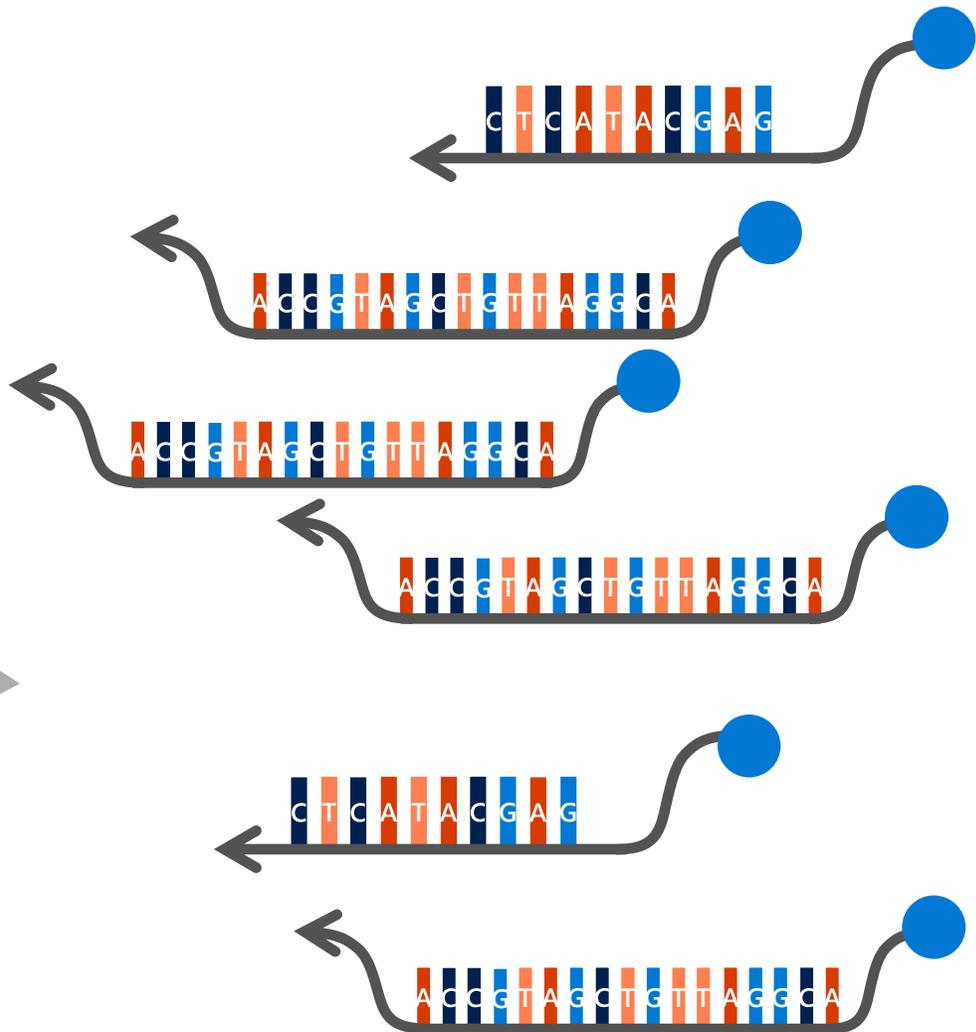
"Database" strands (many)

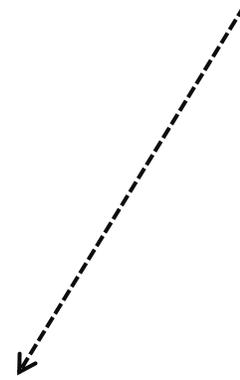
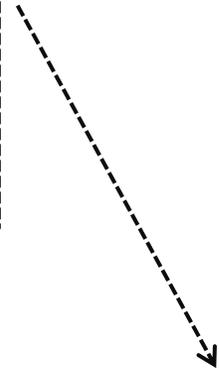


"Query" strands (one type)



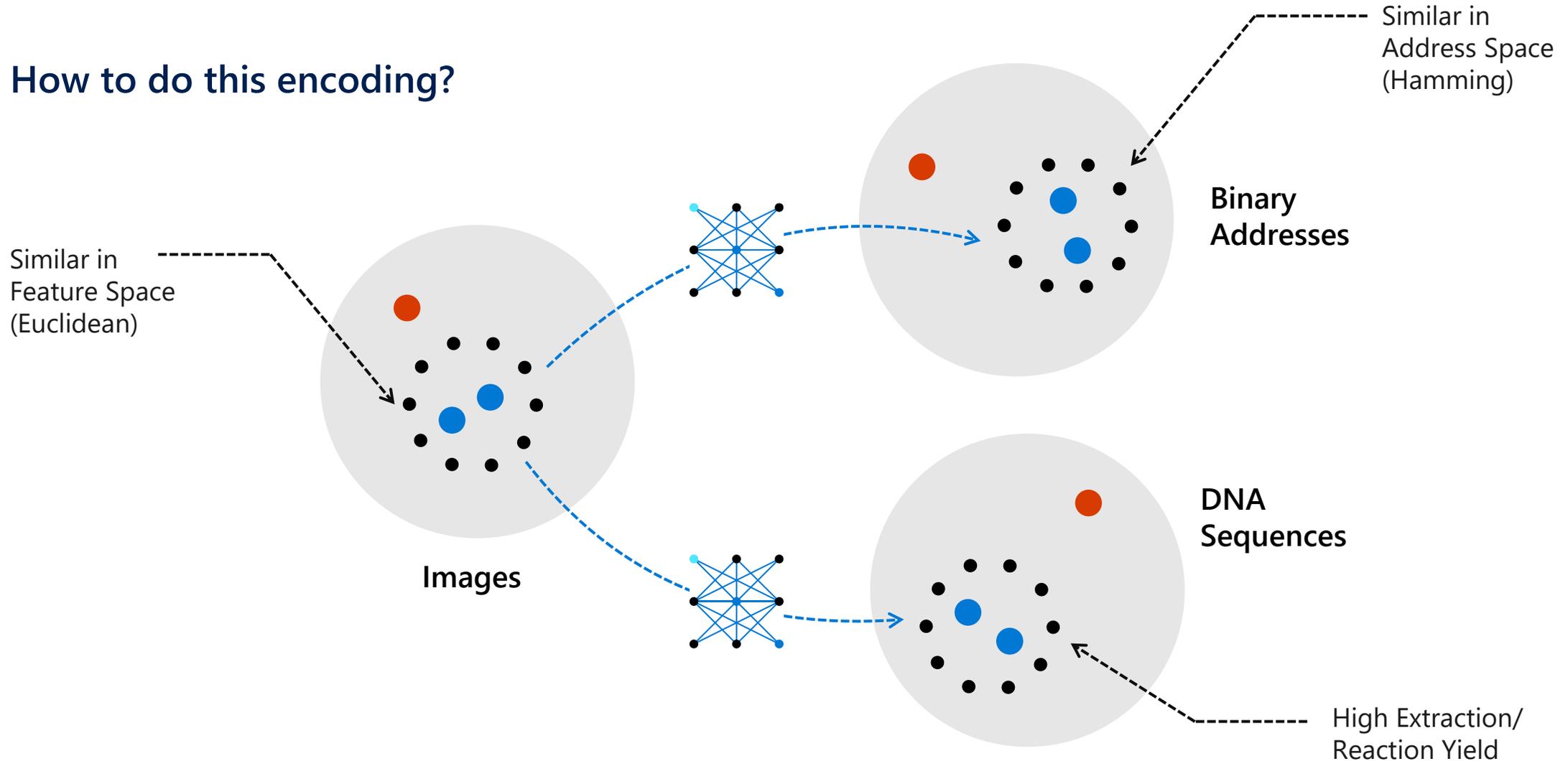




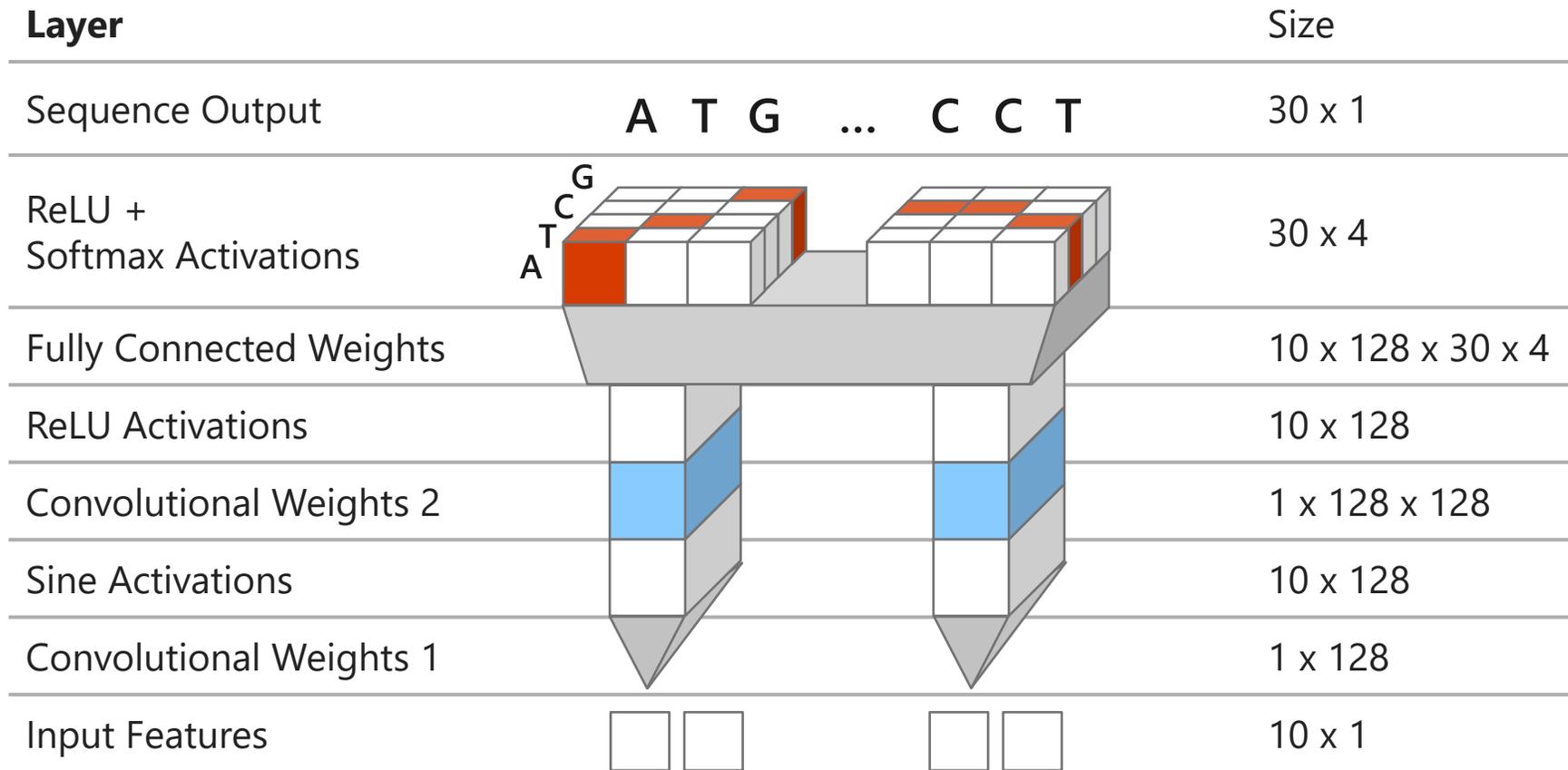


"Semantic" Hashing

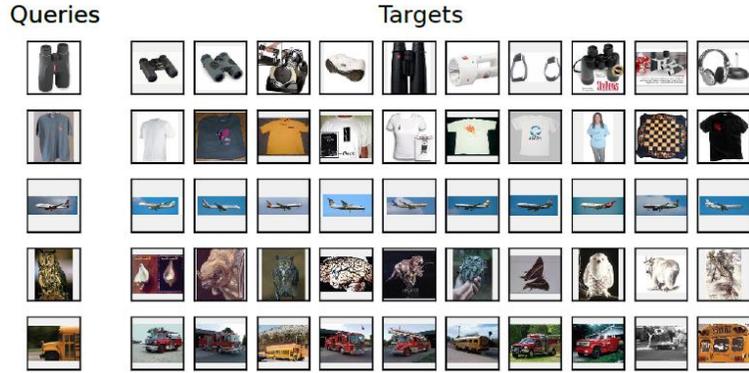
How to do this encoding?



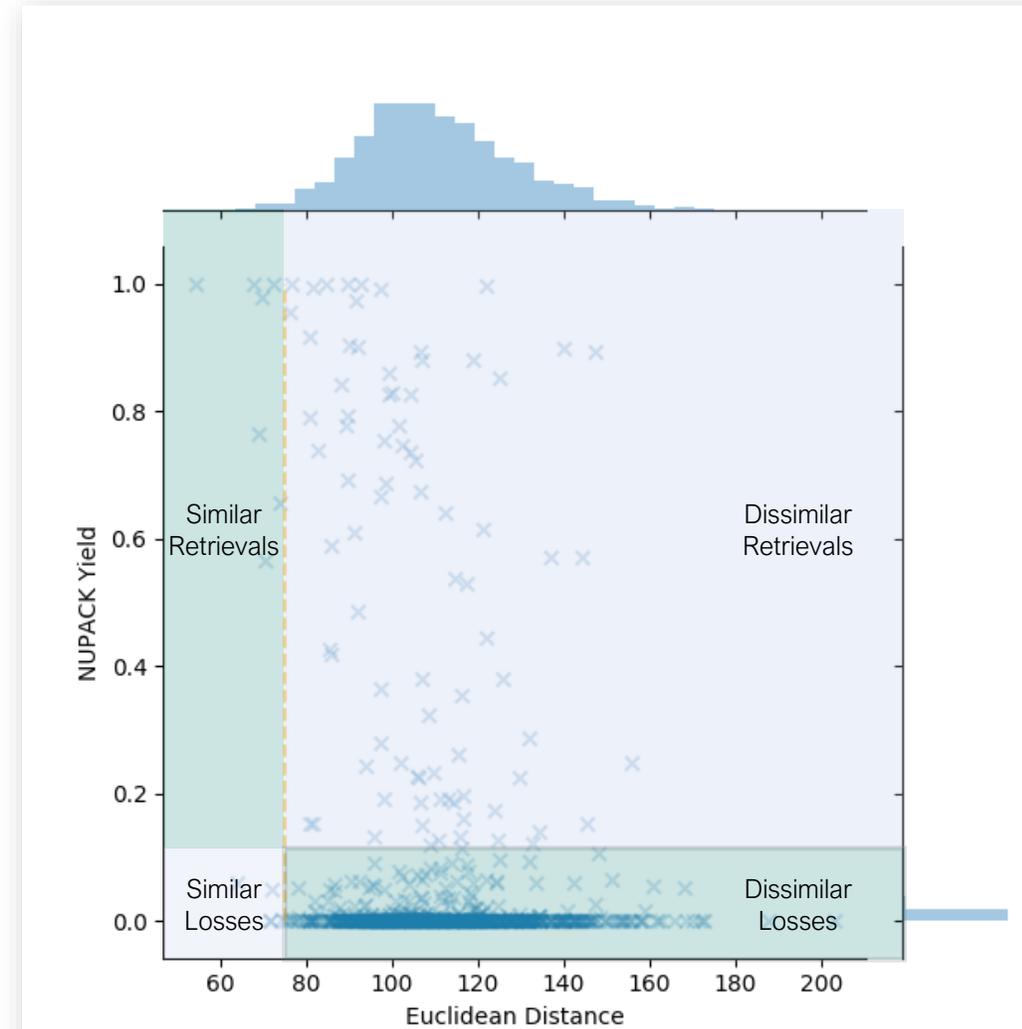
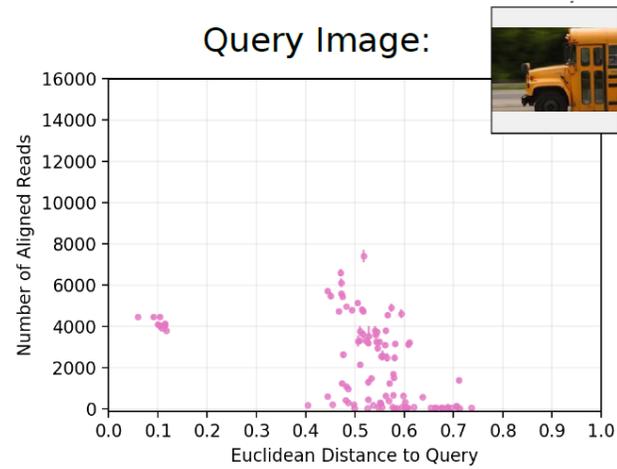
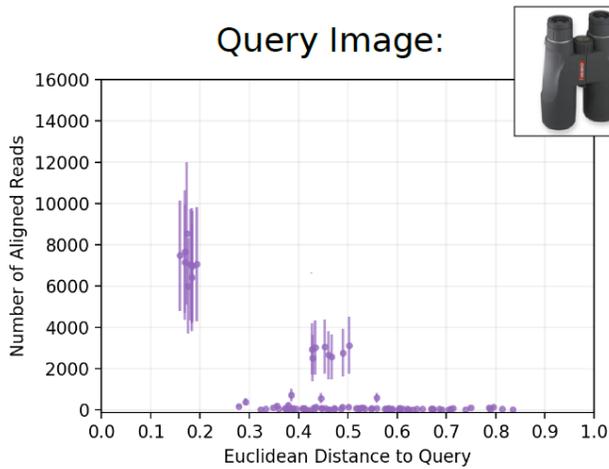
Learning-based encoding



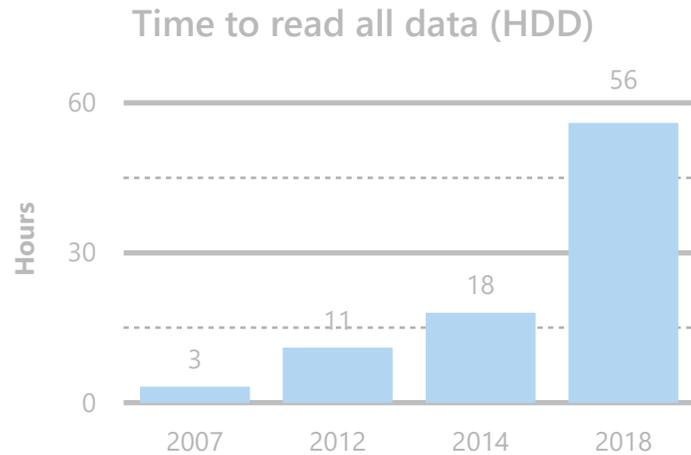
Experiments show encouraging results



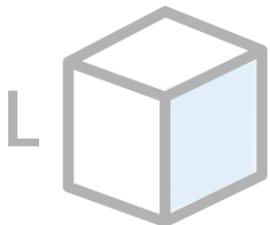
1.6M images, 3 queries
Magnetic bead extraction
Illumina-based sequencing



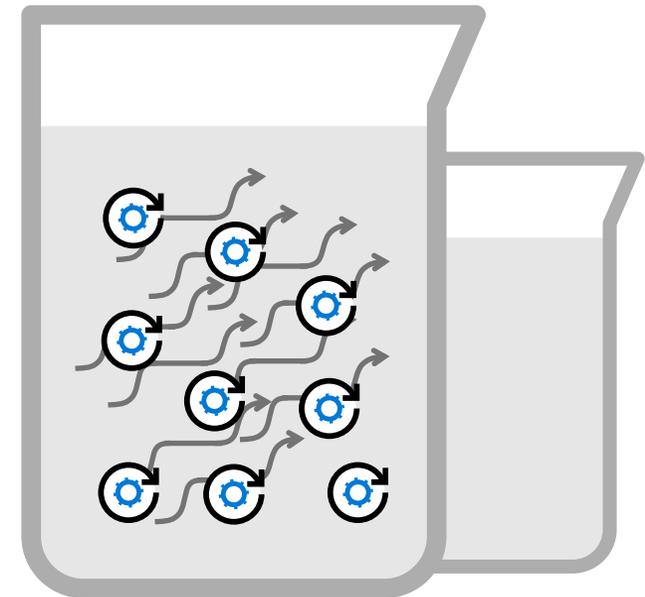
Yottabyte-scale near-molecule computing?



Capacity/bandwidth going up

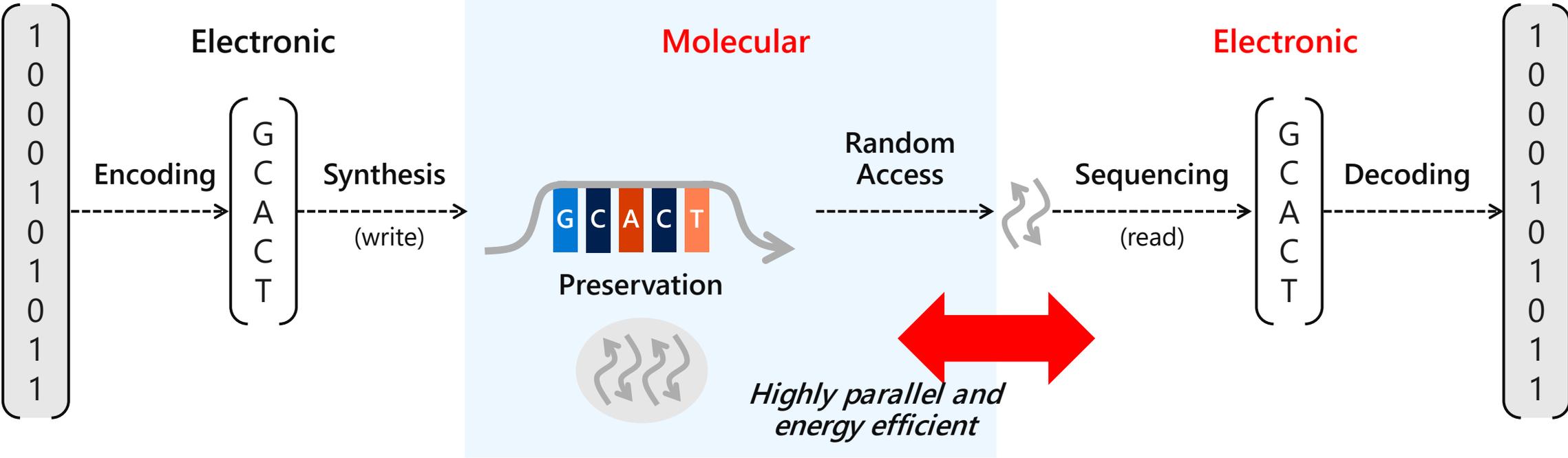


$$\frac{\text{Capacity} \sim L^3}{\text{Bandwidth} \sim L^2}$$

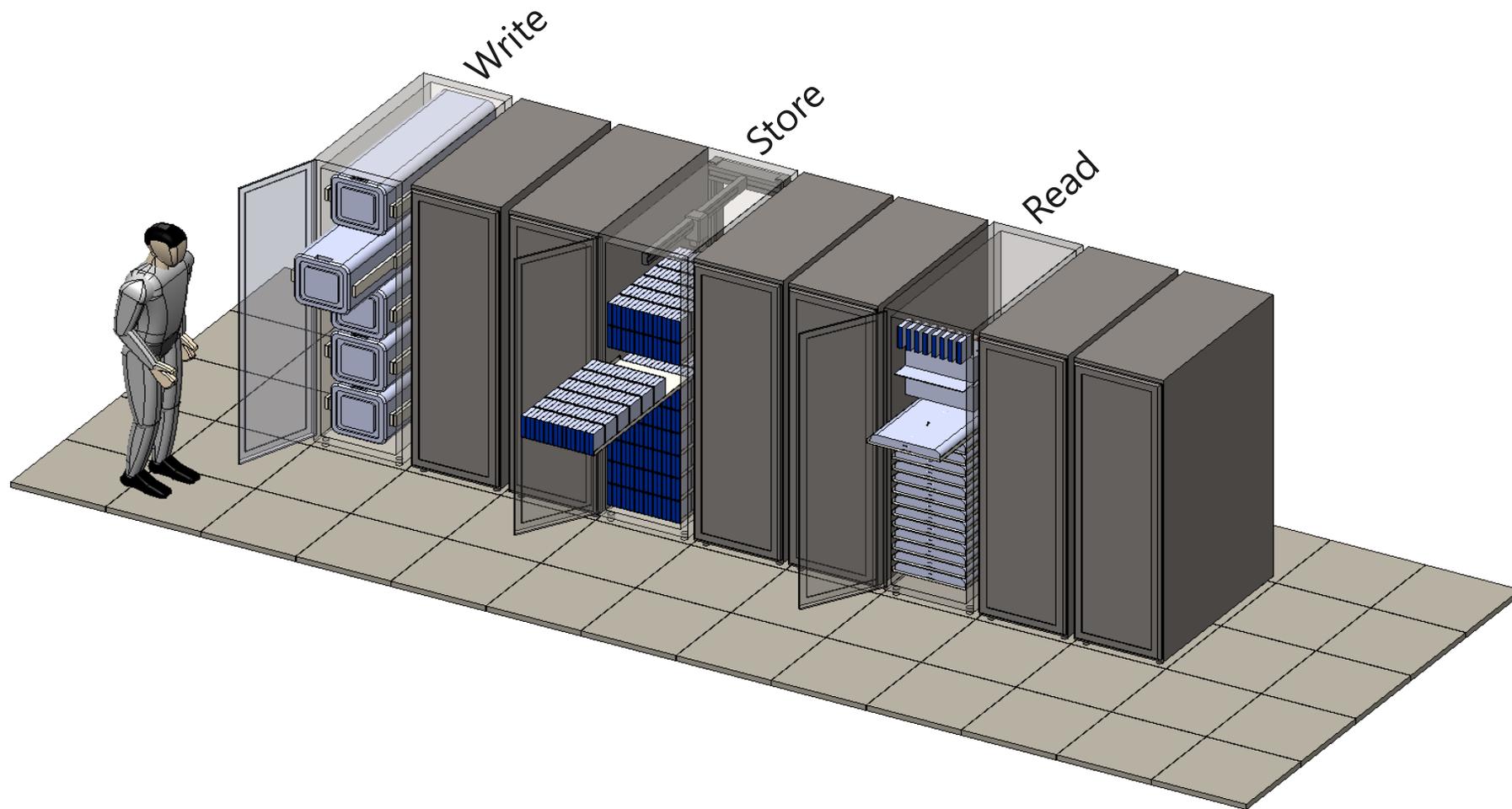


Physically "diffusing" computation through data offers parallelism and virtually unlimited access bandwidth. Yes, at a higher latency.

DNA storage end-to-end system w/ integrated computing

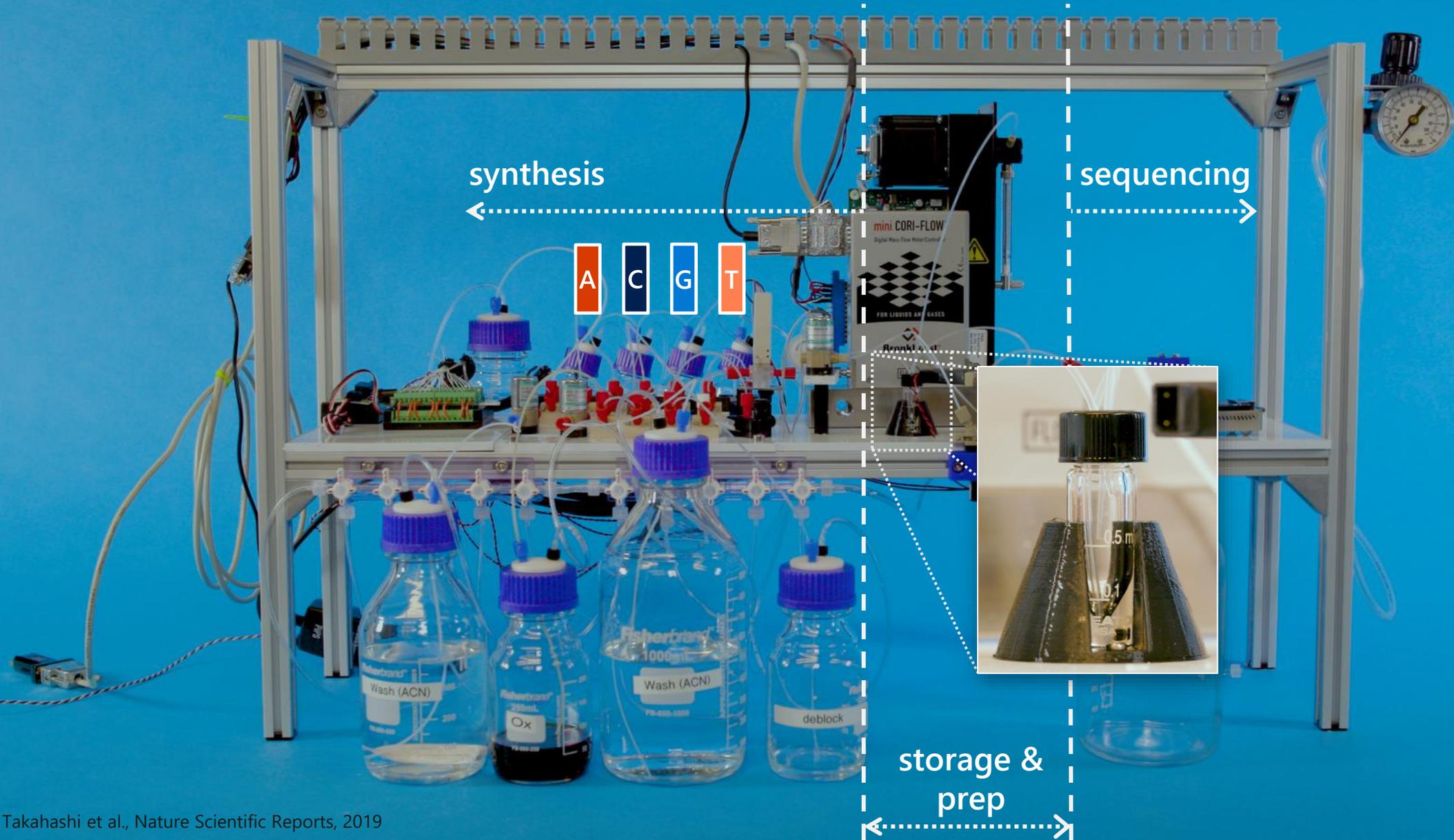


End-to-end system in a datacenter



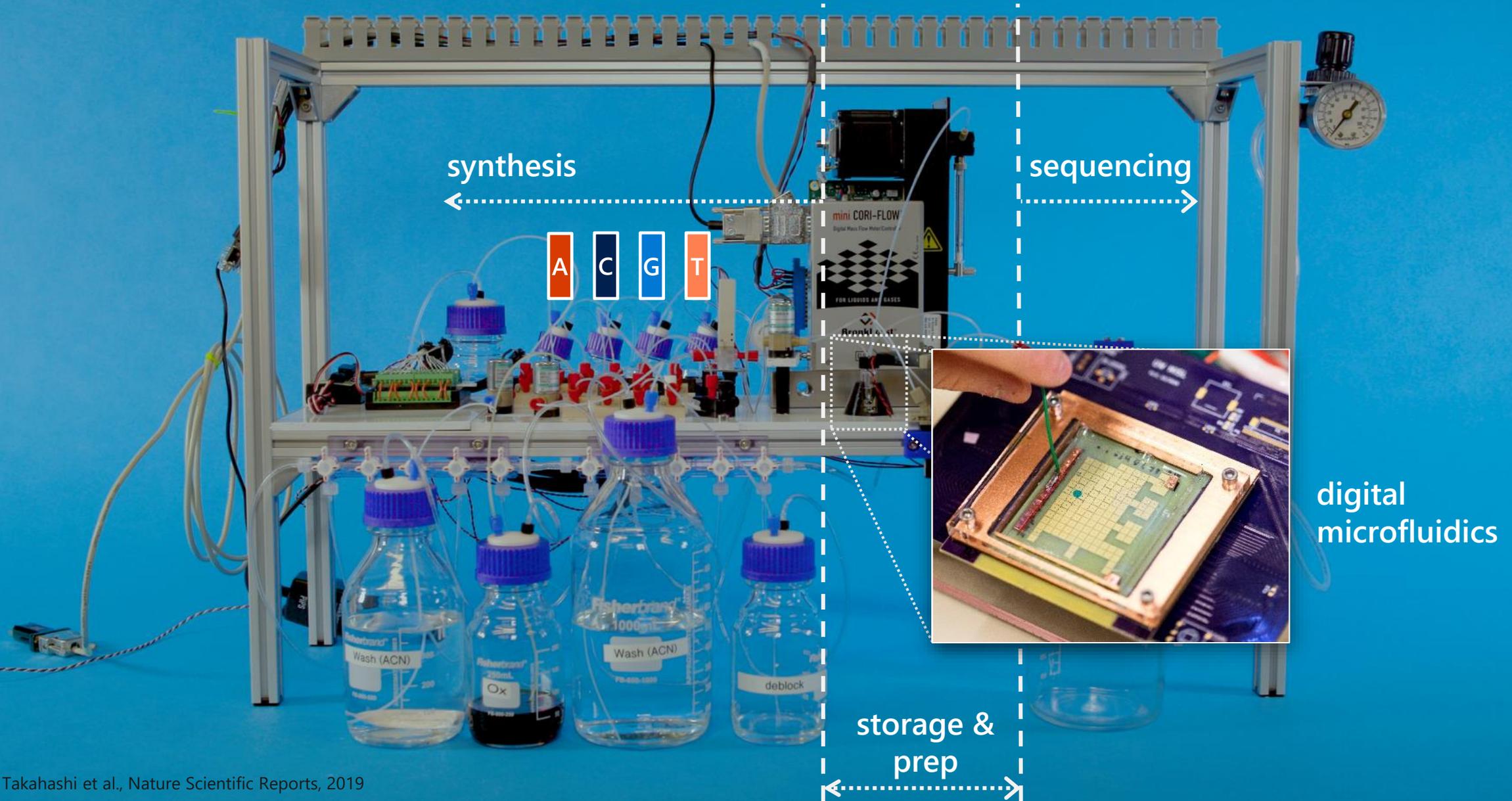


First fully automated DNA data storage system



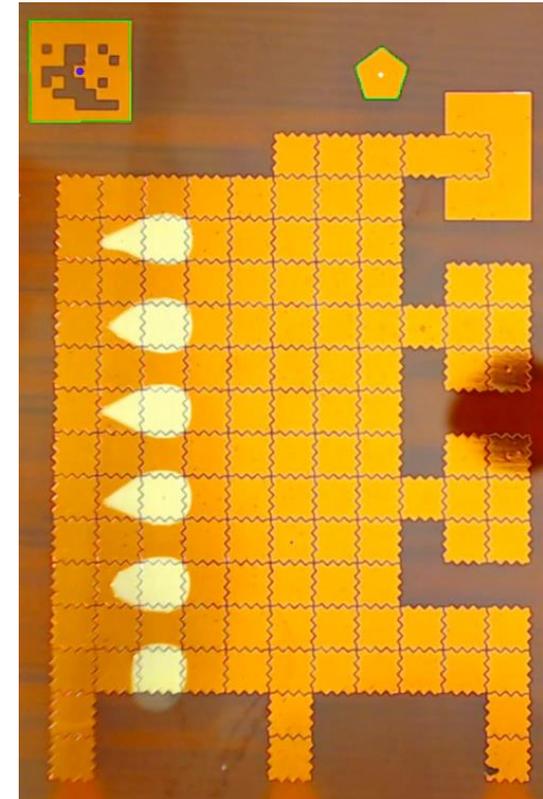
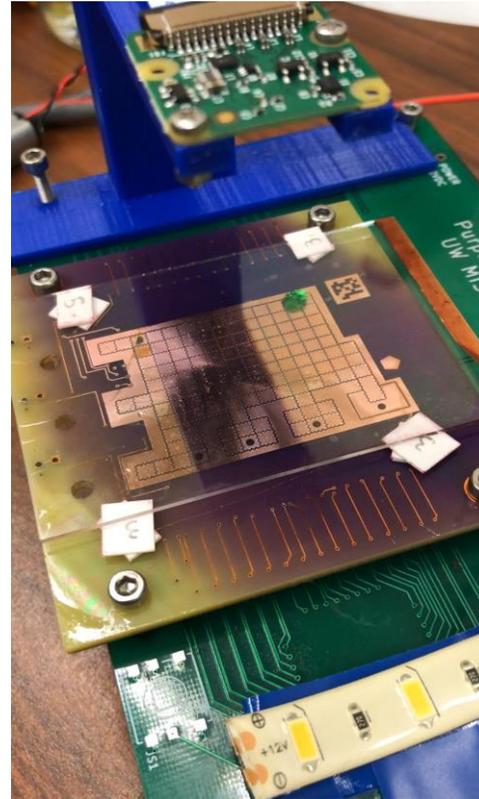
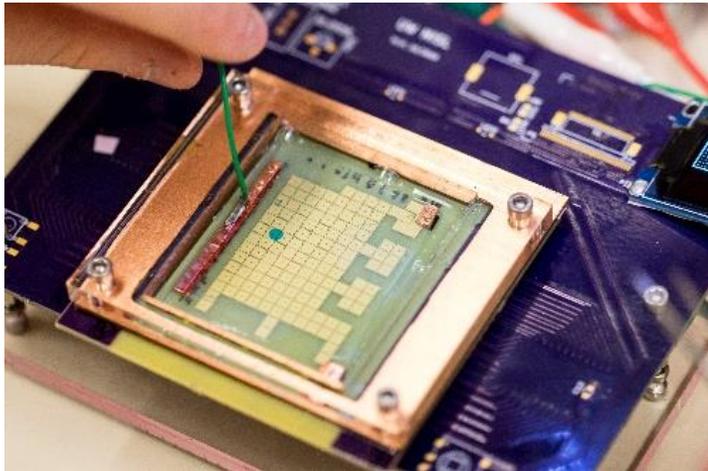
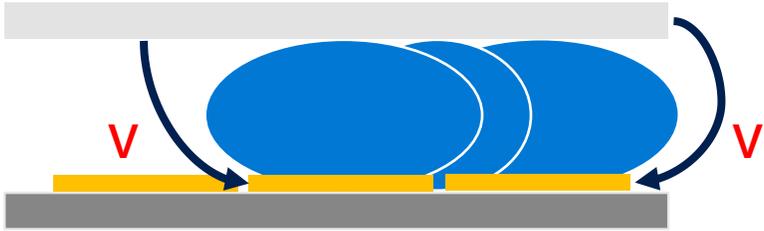


First fully automated DNA data storage system

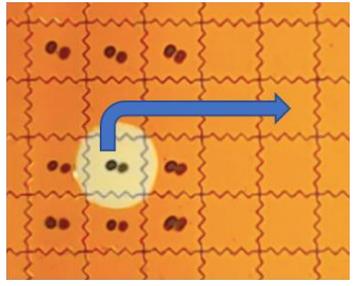
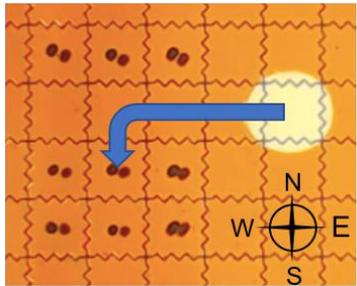
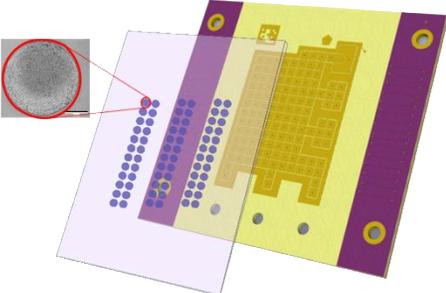


Digital microfluidics

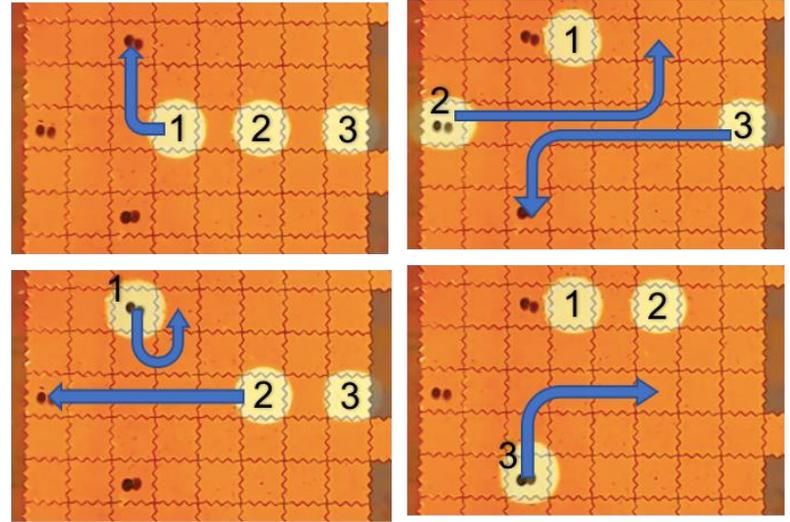
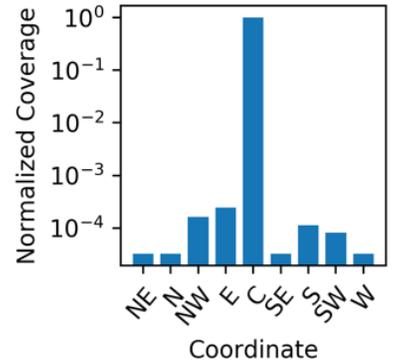
Versatile platform to implement wet lab preparation protocols



Random access with spots+digital microfluidics



60s dwell time
33ng mass



No measurable contamination

Affordable full-stack SW/HW digital microfluidics platform

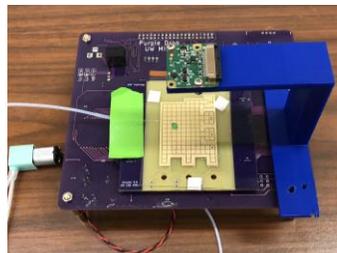
High-level programming with *Puddle*

```
def thermocycle(droplet, temps_and_times):  
    for temp, time in temps_and_times:  
        heat(droplet, temp, time)  
    if droplet.volume < MIN_VOLUME:  
        droplet += input("water", min_volume)  
  
def pcr(droplet, n_iter):  
    thermocycle(droplet, n_iter * [  
        (95, 3 * minutes),  
        (62, 30 * seconds),  
        (72, 20 * seconds),  
    ])
```

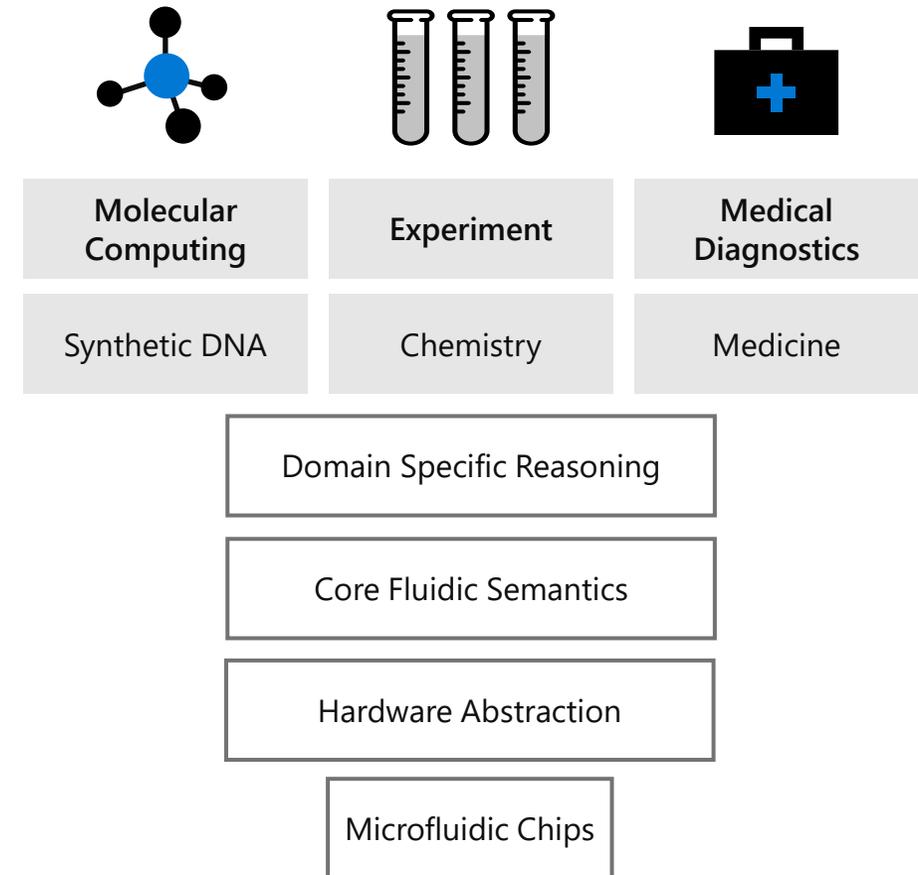
"Assembly code"

```
activate(3,0)  
activate(3,1)  
activate(3,2)  
...
```

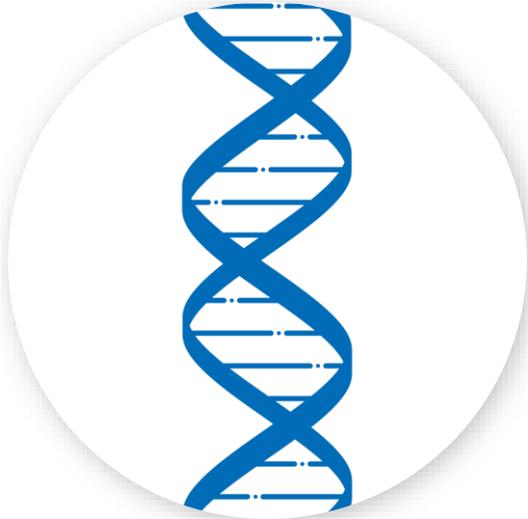
Hardware



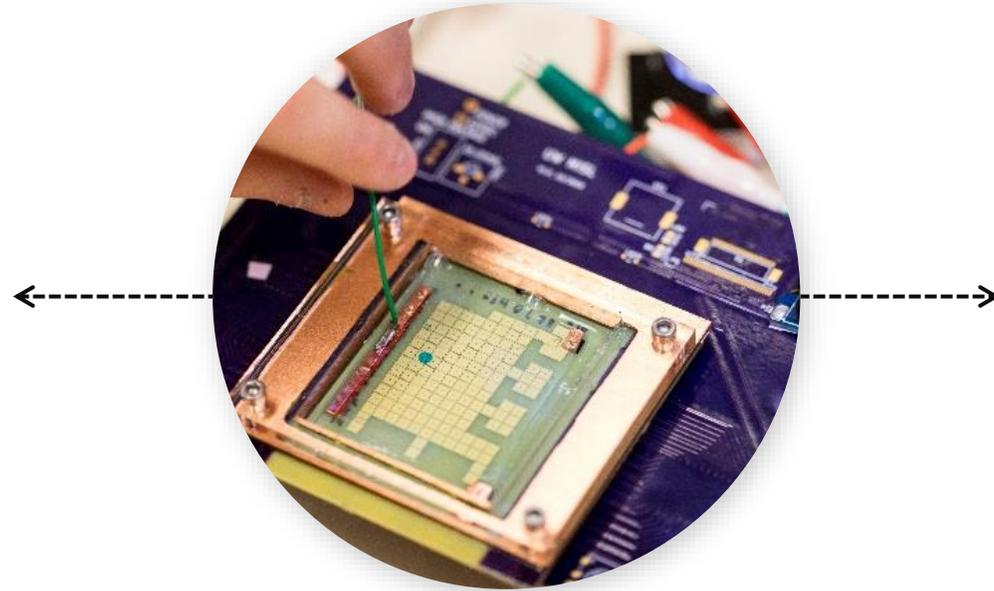
Willsey et al., ASPLOS, 2019; Stephenson et al., IEEE MICRO 2020.



Hardware, software, wetware

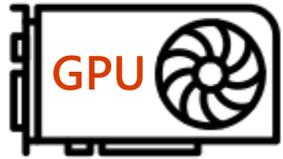


Molecular domain

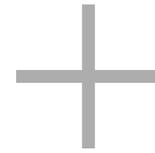
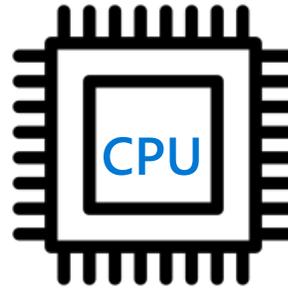


Electronic domain

Future hybrid systems



Special purpose

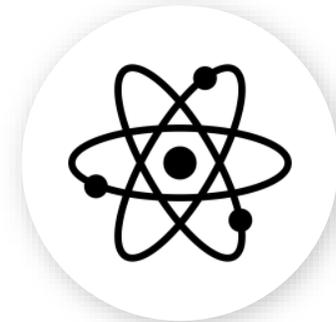


General purpose

Electronics:
Ultra low latency,
engineerable, perfect control



Biomolecules:
Self assembly, massive
data and efficiency



Quantum:
Massive specialized parallel
computing, little data



Questions?

<https://misl.cs.washington.edu>

<https://www.microsoft.com/en-us/research/project/dna-storage/>