# A DNA-Based Archival Storage System

Luis Ceze and Karin Strauss
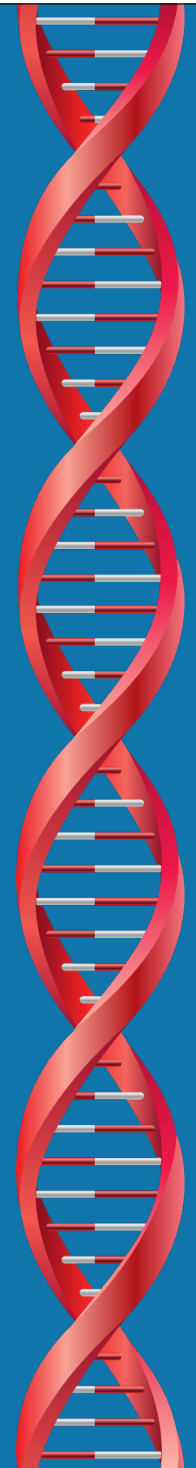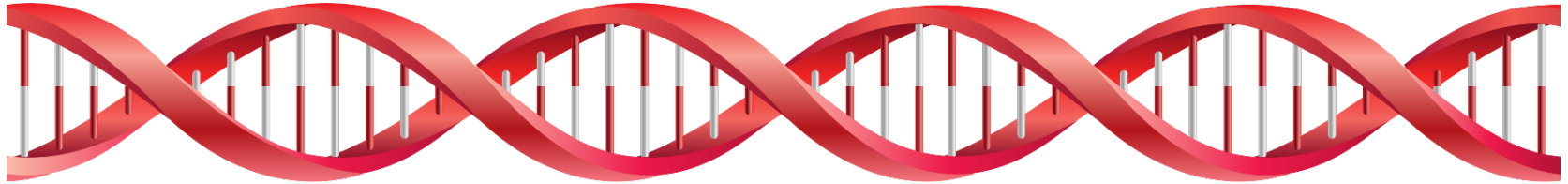
*University of Washington*
*Microsoft Research*

**MISL**

**W** Microsoft Research

*joint work with Doug Carmean, Georg Seelig, James Bornholt, Randolph Lopez, Lee Organick, Rob Carlson, Hsing-Yeh Parker, Yuan Chen, Chris Takahashi, Bichlien Nguyen, Sergey Yekhanin, Siena Dumas Ang, Sharon Newman.*
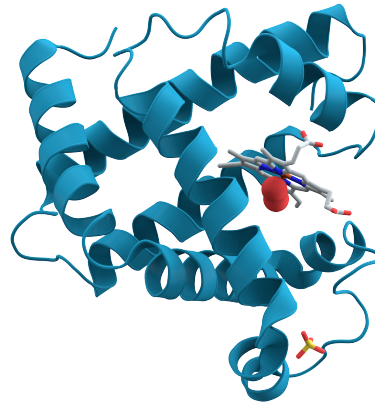
Library of Congress, Sep 2016.
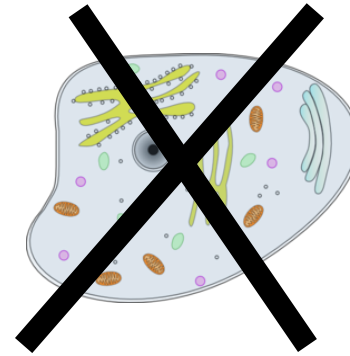
# DNA is the information storage medium for life



Gene

Protein

Function/Characteristic

**MISL UW MSR**

# Using *synthetic* DNA for data storage



*Manufacture* DNA
Dehydrate & store
Read DNA

100101010

But why?

**MISL UW MSR**

# DNA molecules for digital data

**Extremely dense**
1 exabyte in 1 in$^3$
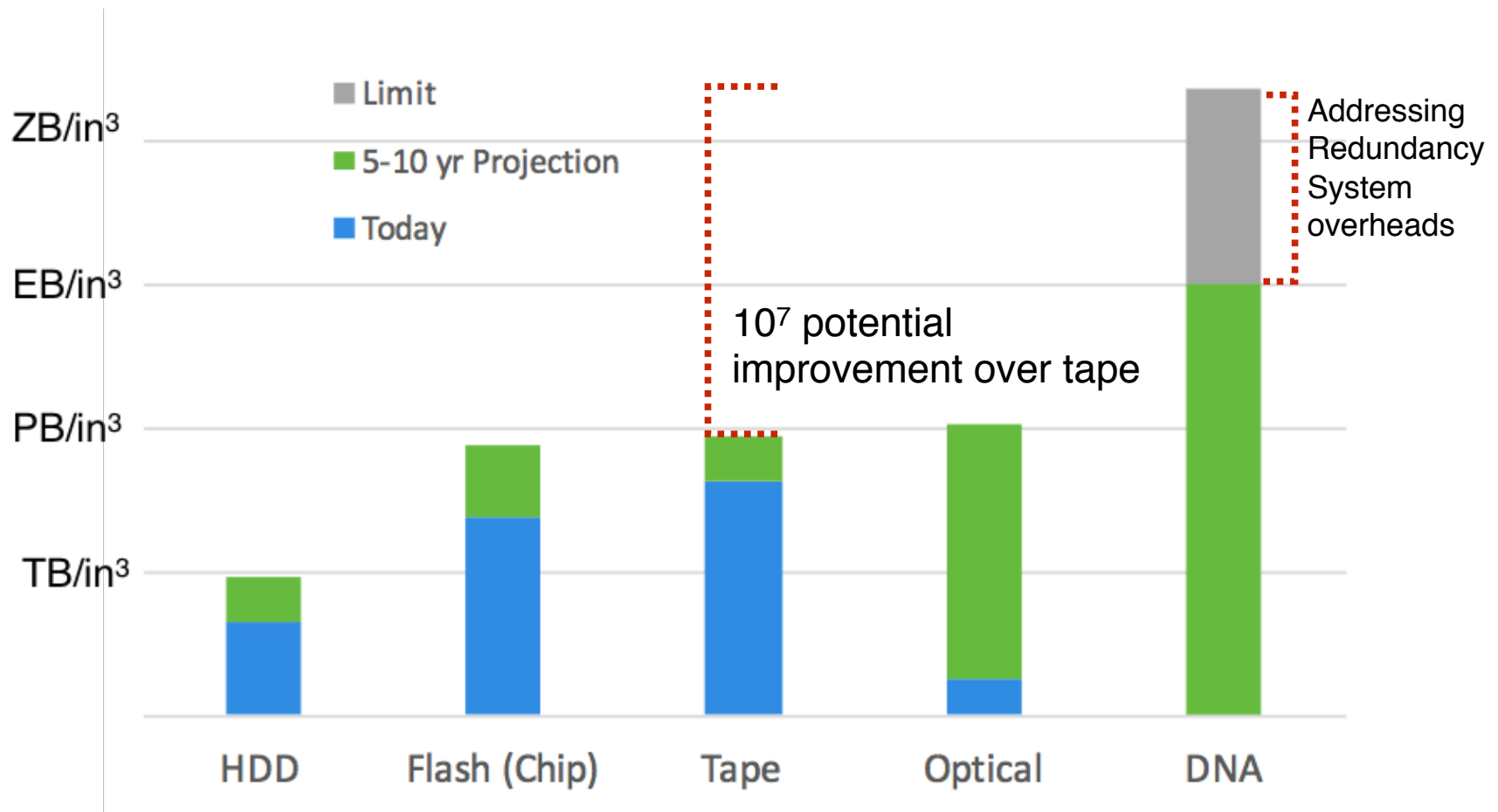
**Extremely durable**
Half life > 500 years

Readers never become obsolete!
(no migration :)

And consumes very little power at rest.

**MISL UW MSR**

# Comparing storage density



Legend:
- **Limit** (grey)
- **5-10 yr Projection** (green)
- **Today** (blue)

Y-axis: ZB/in$^3$, EB/in$^3$, PB/in$^3$, TB/in$^3$

X-axis categories: HDD, Flash (Chip), Tape, Optical, DNA

$10^7$ potential improvement over tape

Addressing Redundancy System overheads

**MISL UW MSR**

# The ultimate storage hierarchy

| | Access Time | Capacity | Durability |
|---|---|---|---|
| Flash | $\mu$s-ms | TBs | ~5 yrs |
| HDD | 10s ms | 100s TBs | ~5 yrs |
| Tape | minutes | PBs | ~10s yrs |
| DNA-based Archival | hours | ZBs | ~100s yrs |

*Our goal: build an integrated DNA storage system.*

**MISL UW MSR**

# DNA molecules

Four nucleotides:

(A) Adenine

(C) Cytosine

(G) Guanine

(T) Thymine

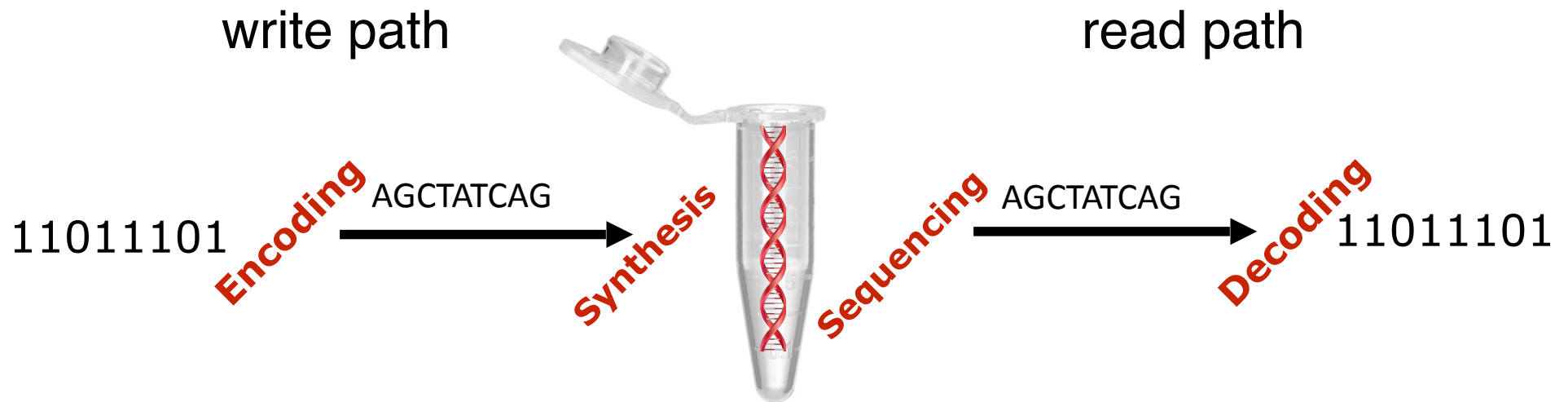DNA strand (oligonucleotide) is a linear sequence of these nucleotides

(G)—(A)—(C)—(A)—(C)—(C)—(T)

Two strands can bind to each other if they are complementary:

(G)—(A)—(C)—(A)—(C)—(T)—(T)

(C)—(T)—(G)—(T)—(G)—(A)—(A)

C, G are complementary

A, T are complementary

**MISL UW MSR**

# DNA data storage at 30,000 feet

write path

read path

11011101 → *Encoding* AGCTATCAG → *Synthesis*

*Sequencing* AGCTATCAG → *Decoding* 11011101

**MISL UW MSR**

# DNA data storage at 30,000 feet

write path

read path

11011101 → *Encoding* AGCTATCAG → *Synthesis*



*Sequencing* AGCTATCAG → *Decoding* 11011101

**MISL UW MSR**

# Encoding digital data in DNA

1010001110010001111001111000101100101001011101...

0 ⇒ A
1 ⇒ C
2 ⇒ G
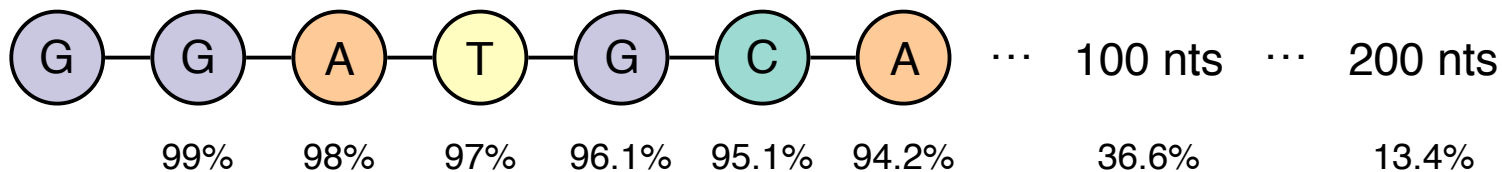3 ⇒ T

⬇

2 2 0 3 2 1 0 1 3 2 1 3 3 0 1 1 2 1 1 0 2 3 3 1

⬇

G G A T G C A C T G C T T A C C G C C A G T T C

Repeated letters are bad: Use base 3 and "rotate" mapping.

G A A A C C T

Synthetic DNA sequences have limited length: Break it up

P[Attach] = 99%

G G A T G C A ⋯ 100 nts ⋯ 200 nts

99%  98%  97%  96.1%  95.1%  94.2%     36.6%     13.4%

**MISL UW MSR**

# Breaking up data into chunks (~150nts)

C G A T  G C A C  T G C T  G A C G  G C T A  G C T C

| A T G T T |
|:---------:|
| A T G T T |
| A T G T T |

File identifiers ("primers")

| A A A A | C A T C C |
|:-------:|:---------:|
| A A A C | C A T C C |
| A A A G | C A T C C |

Addresses within the file

1 of N
2 of N
3 of N

**~ 20 bytes per DNA strand. Many strands per file.**

**MISL UW MSR**

# Errors in writing/reading DNA

MISL UW MSR

# DNA data storage at 30,000 feet

write path

read path

11011101 → *Encoding* AGCTATCAG → *Synthesis*

*Sequencing* AGCTATCAG → *Decoding* 11011101

**MISL UW MSR**
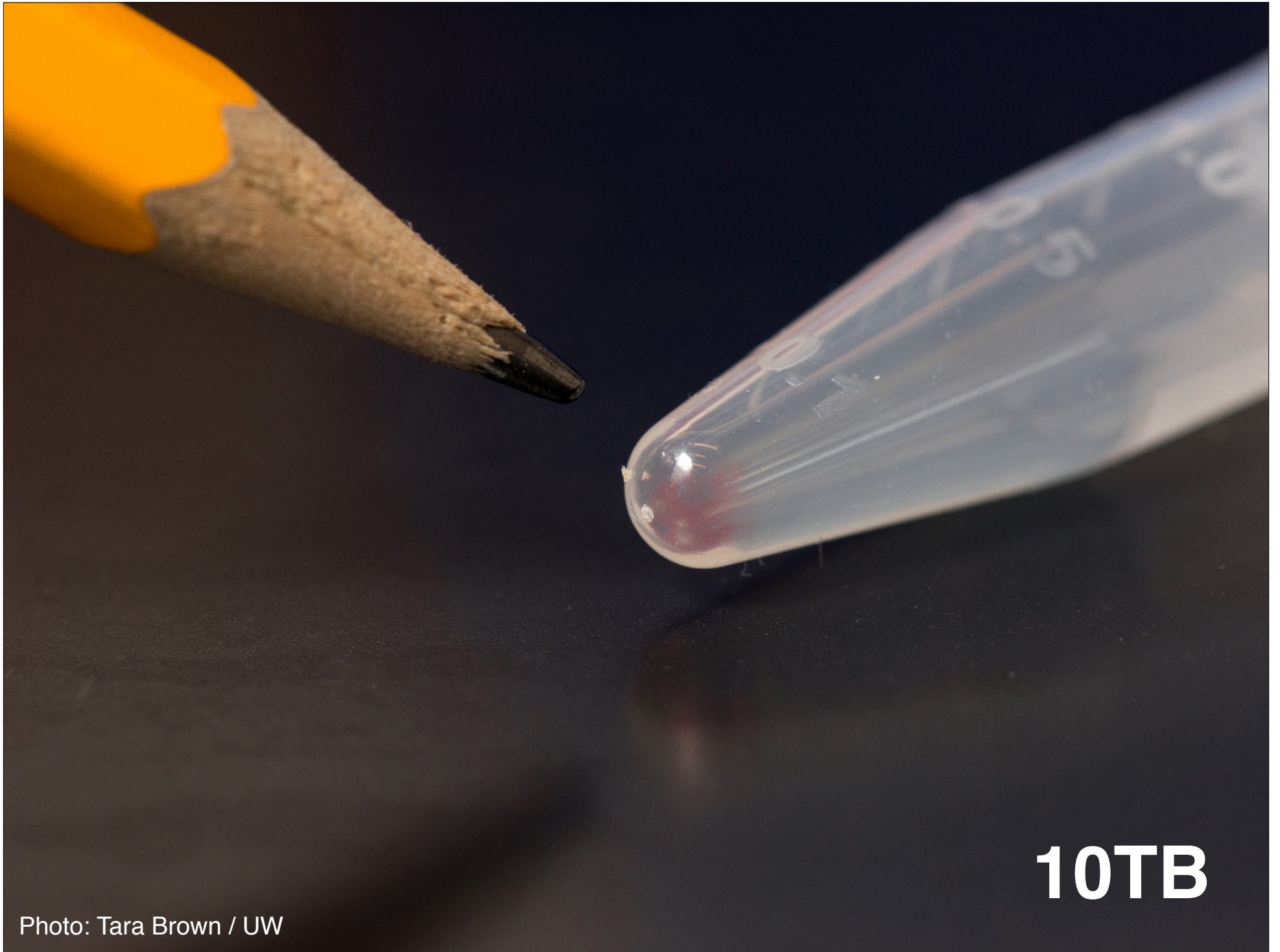
# DNA Synthesis

Manufacturing DNA strands

GACACCT ➜ (G)—(A)—(C)—(A)—(C)—(C)—(T)

- Normally used for life sciences and medicine

- Millions of copies of each sequence

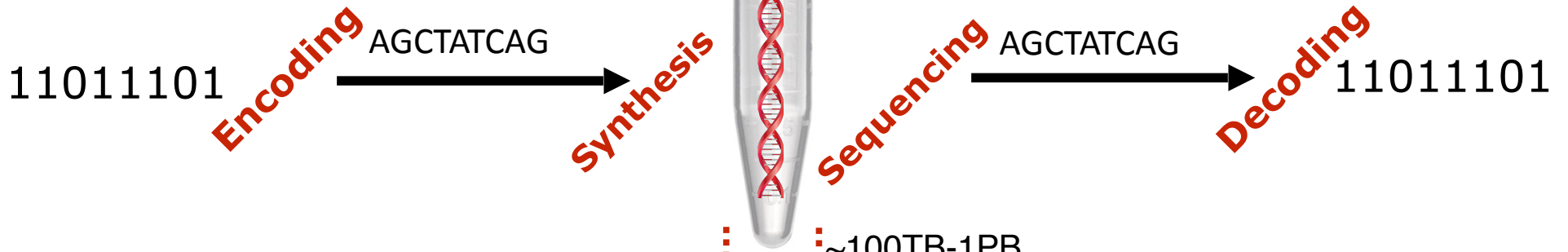- Can make many different sequences in parallel

**Twist Bioscience**

**MISL UW MSR**

10TB

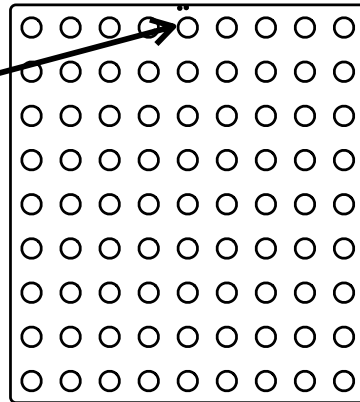Photo: Tara Brown / UW

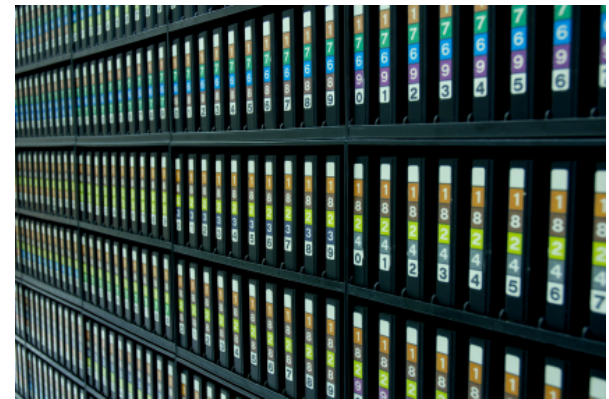# DNA Storage "Library"

write path                                    read path

11011101  **Encoding** → AGCTATCAG → **Synthesis**    **Sequencing** → AGCTATCAG → **Decoding** 11011101

~100TB-1PB

Data address specifies
physical location.

foo.mp4

DNA storage
(physical) library
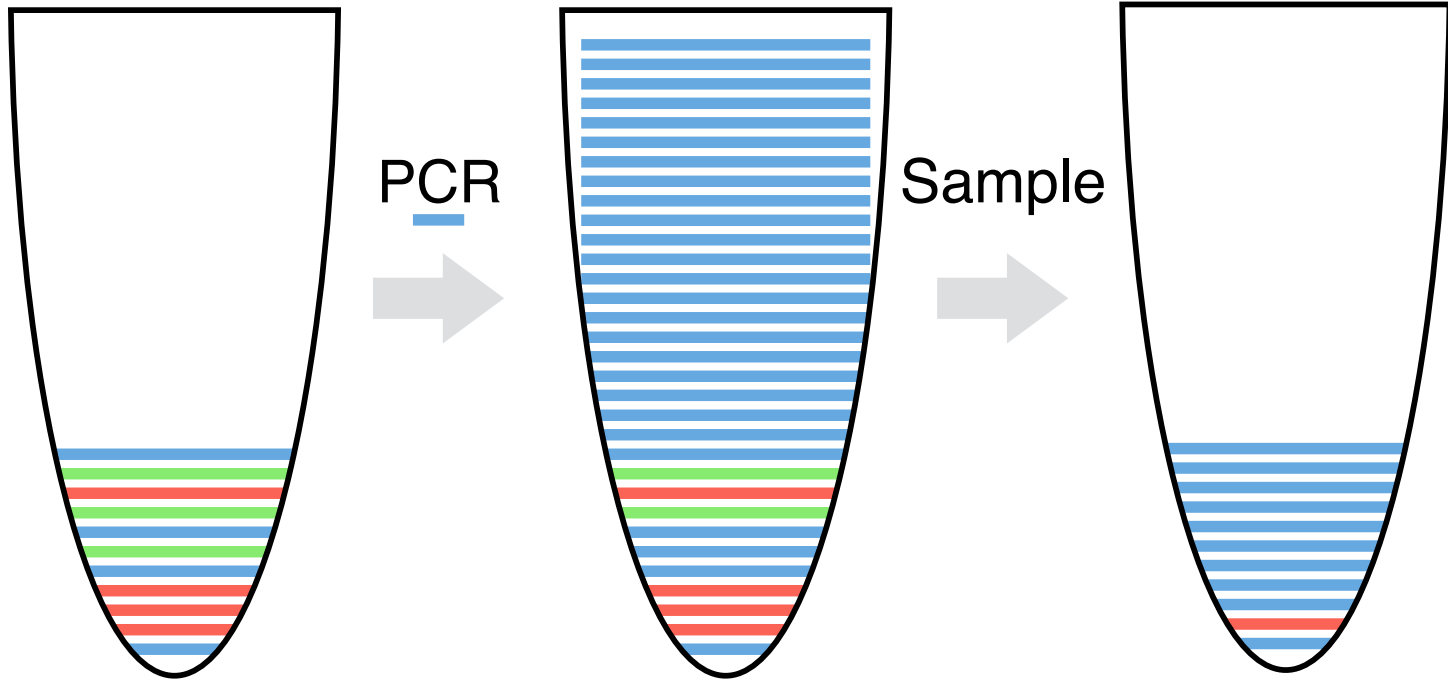
**MISL UW MSR**

# DNA data storage at 30,000 feet

write path

read path

11011101

AGCTATCAG

11011101

AGCTATCAG

*Encoding*

*Synthesis*

*Sequencing*

*Decoding*

**MISL UW MSR**

# Random access?

**MISL UW MSR**

# Random access!

**MISL UW MSR**

# DNA Sequencing

## Reading DNA strands

G—A—C—A—C—C—T  ➡️  GACACCT
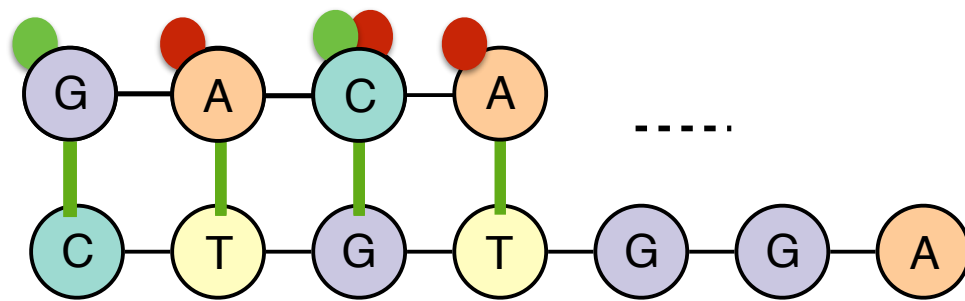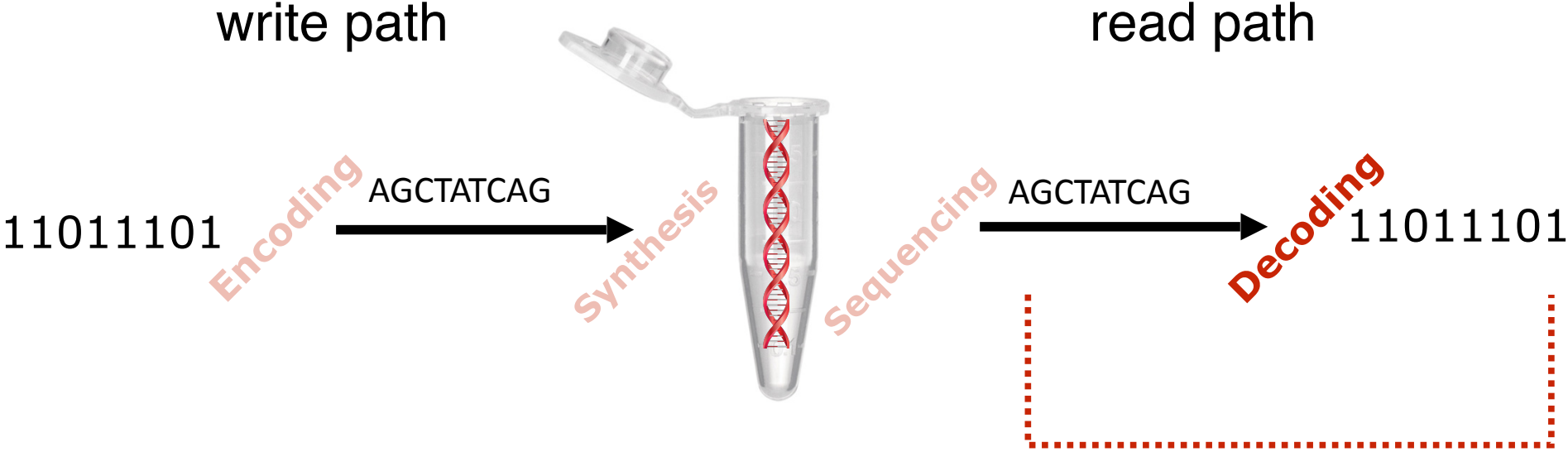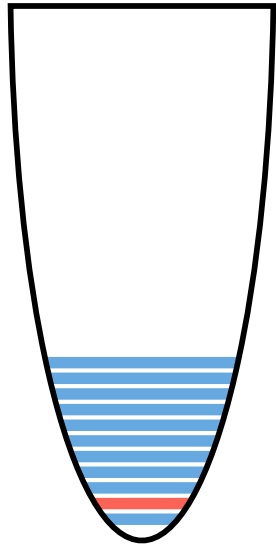
- Normally used for genome sequencing

- Reads many copies of millions of DNA strands at a time

- Currently much higher throughput than synthesis

**MISL UW MSR**

# DNA data storage at 30,000 feet

write path          read path

11011101  *Encoding* →AGCTATCAG→ *Synthesis*  *Sequencing* →AGCTATCAG→ *Decoding* 11011101
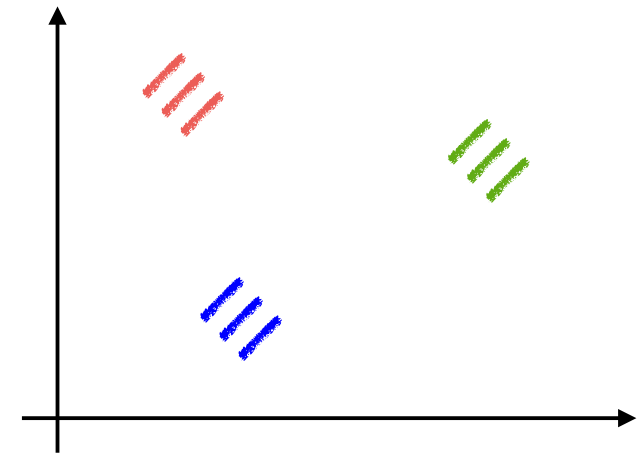
**MISL UW MSR**

# Decoding back to digital data



Sequencing

```
ATGTT GGAT GCAC AAAA CATCC
ATGTT TGCT TACC AAAC CATGC
ATGTT GCCA GTTC ACAGCATCC
ATGTT GGAT GCAC AAGACATCC
ATGTT TGCT TACC CAACCATCC
ATGTT GCCA GTTC AAAGCATCC
ATGTT GGAT GCAC AAGACATCC
ATGTT TGCT TACC CAACCATCC
ATGTT GCCA GTTC AAAGCATCC
........
```

Clustering reads

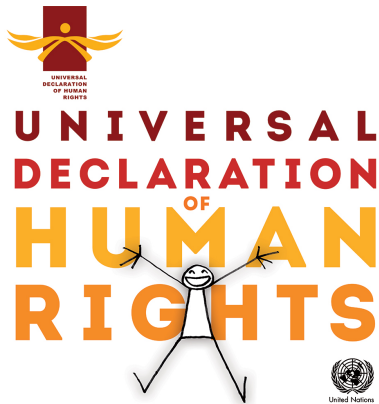Reassemble data

Error correction

…110110101101…

**MISL UW MSR**

# Results

200MB as of July'16.  last year: 1MB
1.5 B nucleotides
10M DNA strands





NEWS

## THIS TOO SHALL PASS RUBE GOLDBERG = DNA

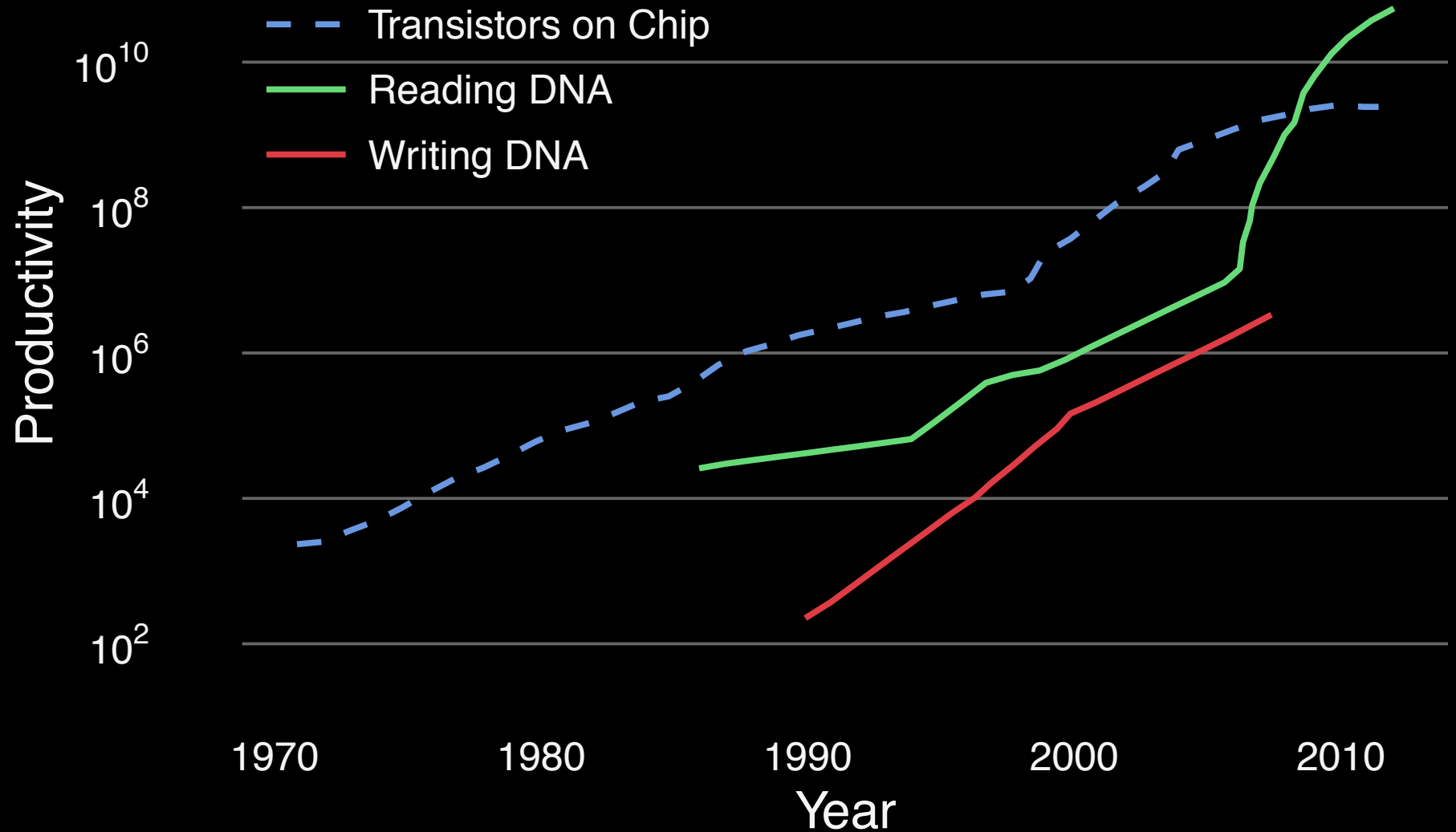JULY 7, 2016

Share → Facebook  Twitter



This Too Shall Pass Rube Goldberg is not just a video anymore! Huge Thanks to Microsoft Research and the University of Washington. Read all about it here.
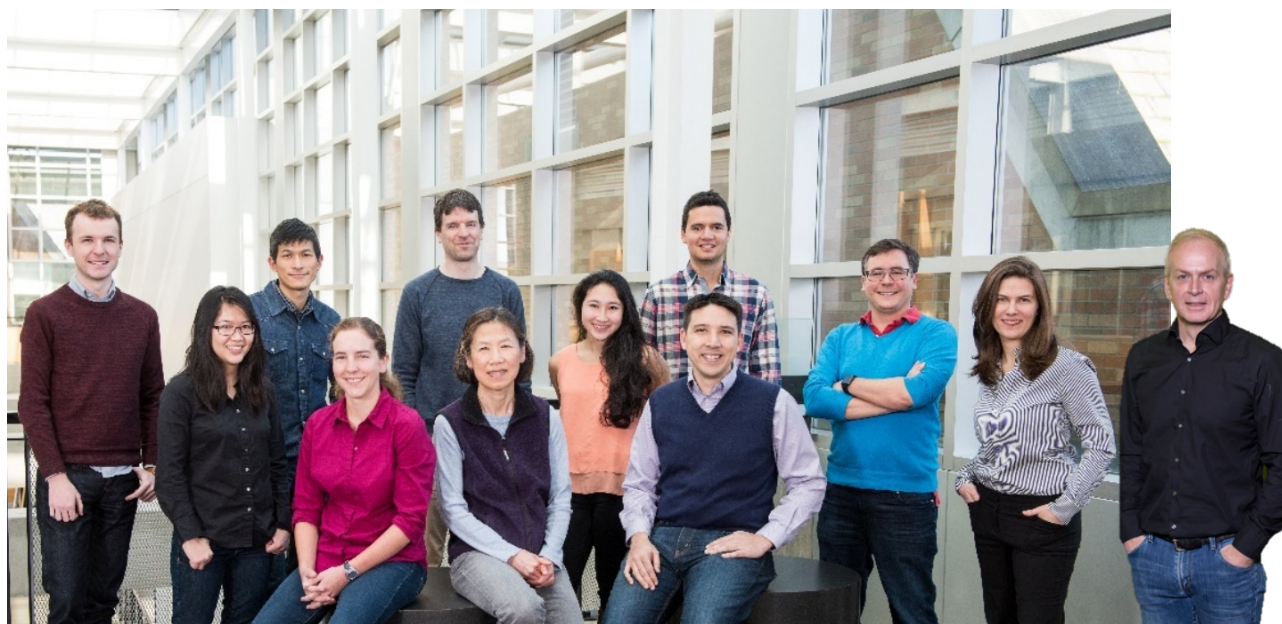
**MISL UW MSR**

10MBs/week ➡ 100GBs/second

# Molecular Information Systems Lab



Computer architects, coding theorists, molecular biologists, fluidics, algorithms, …