# Computer Architecture
## Lecture 3a: Introduction to Genome Sequence Analysis

Prof. Onur Mutlu

ETH Zürich

Fall 2020
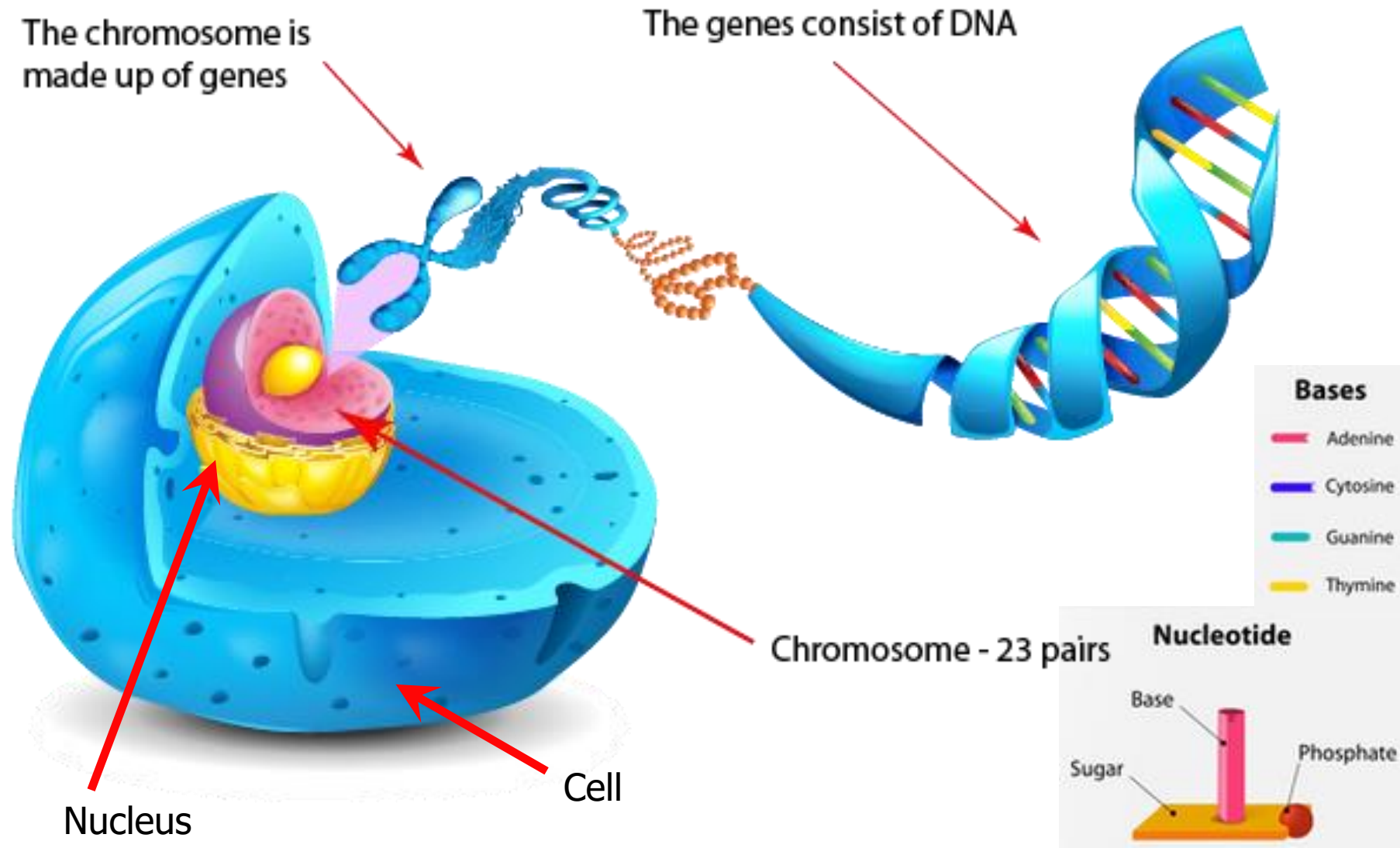
24 September 2020

# Four Key Problems + Directions

- Fundamentally Secure/Reliable/Safe Architectures

- Fundamentally Energy-Efficient Architectures
  - Memory-centric (Data-centric) Architectures

- Fundamentally Low-Latency and Predictable Architectures

- Architectures for AI/ML, Genomics, Medicine, Health

**SAFARI**

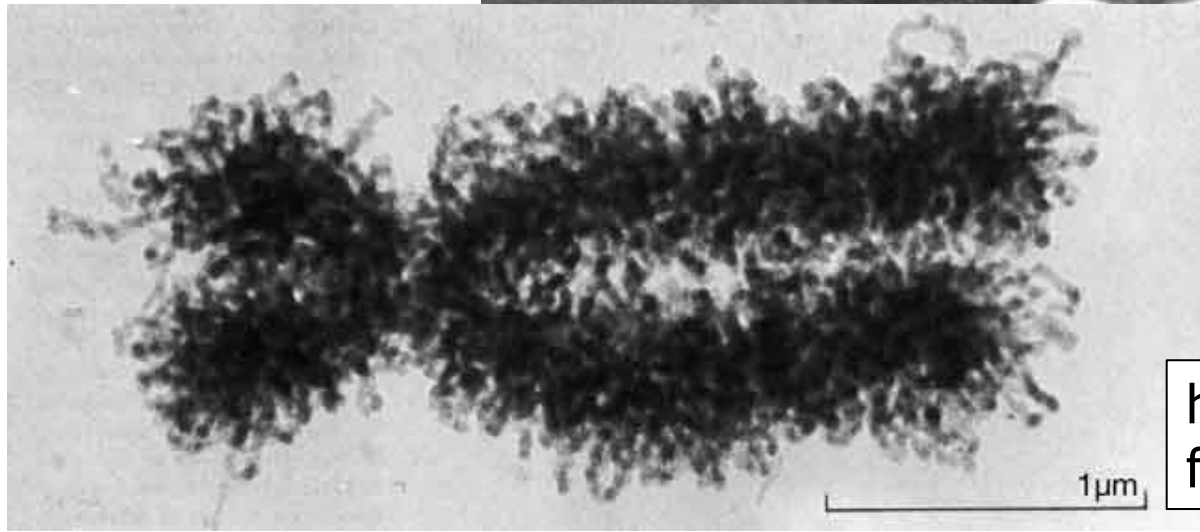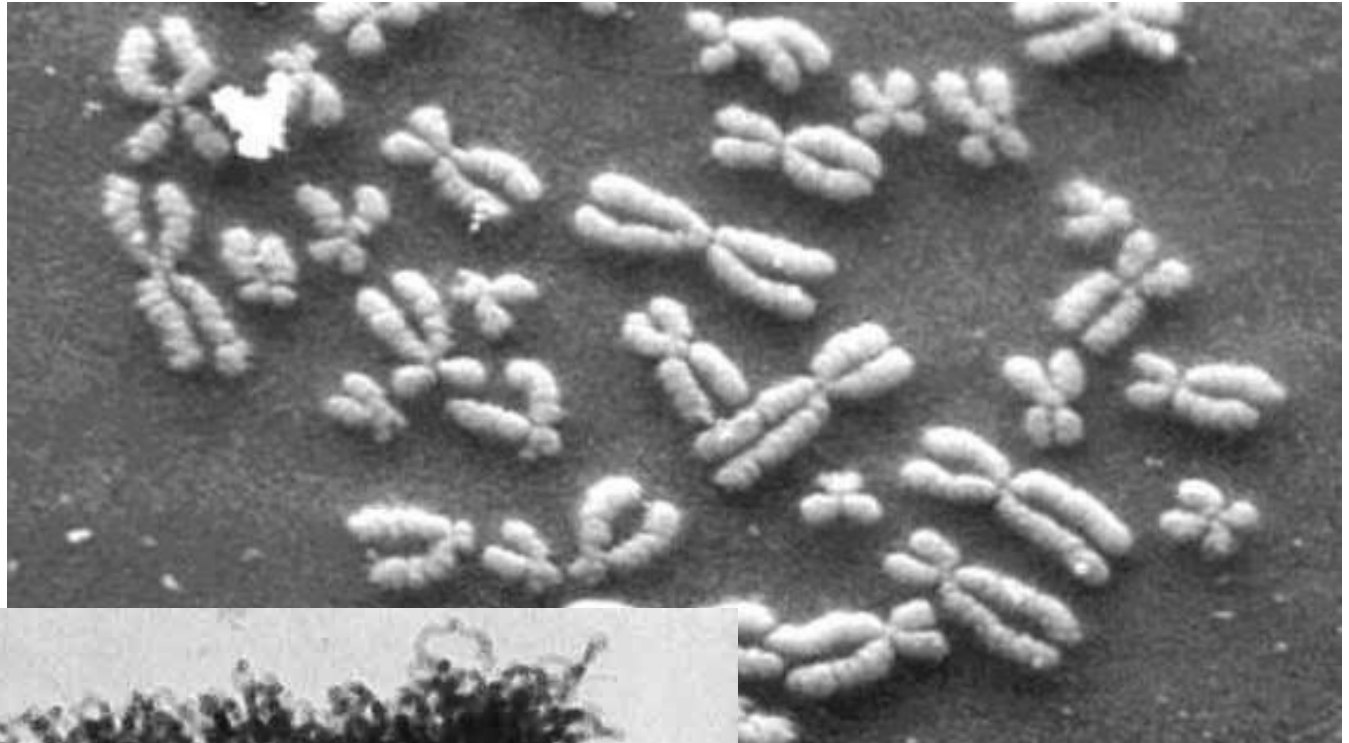# A Motivating Detour: Genome Sequence Analysis

# Our Dream (circa 2007)

- An embedded device that can perform comprehensive genome analysis in real time (within a minute)
    - Which of these DNAs does this DNA segment match with?
    - What is the likely genetic disposition of this patient to this drug?
    - What disease/condition might this particular DNA/RNA piece associated with?
    - . . .

# What Is a Genome Made Of?

The chromosome is made up of genes

The genes consist of DNA

Nucleus

Cell

Chromosome - 23 pairs

**Bases**
- Adenine
- Cytosine
- Guanine
- Thymine

**Nucleotide**

Base

Sugar

Phosphate

The discovery of DNA's double-helical structure (Watson+, 1953)

# DNA Under Electron Microscope



human chromosome #12 from HeLa's cell
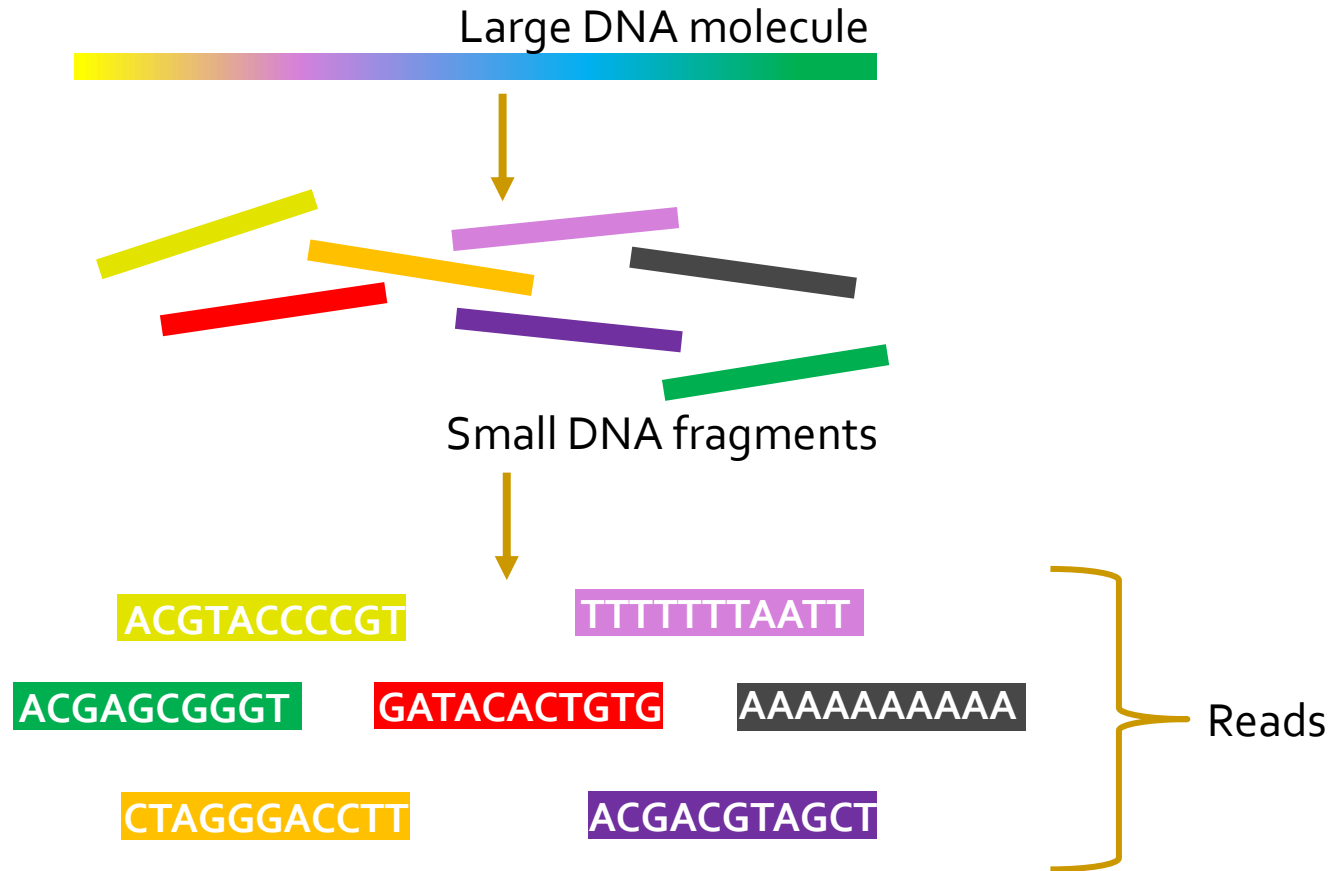
1μm

# Genome Sequencing

- **Goal:**
  - Find the complete sequence of A, C, G, T's in DNA (or RNA).

- **Challenge:**
  - There is no machine that takes long DNA as an input, and gives the complete sequence as output
  - All sequencing machines chop DNA into pieces and identify relatively small pieces (but not how they fit together)

# Genome Sequencing



Large DNA molecule

Small DNA fragments

ACGTACCCCGT
TTTTTTTAATT

ACGAGCGGGT
GATACACTGTG
AAAAAAAAAA

CTAGGGACCTT
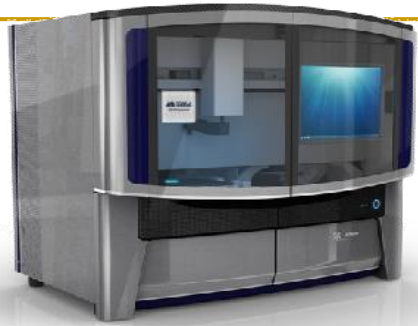ACGACGTAGCT

Reads

# Untangling Yarn Balls & DNA Sequencing

# Genome Sequencers


Roche/454


AB SOLiD


Illumina MiSeq


Complete Genomics


Illumina HiSeq2000


Pacific Biosciences RS


Oxford Nanopore MinION


Illumina NovaSeq 6000


Ion Torrent PGM


Ion Torrent Proton
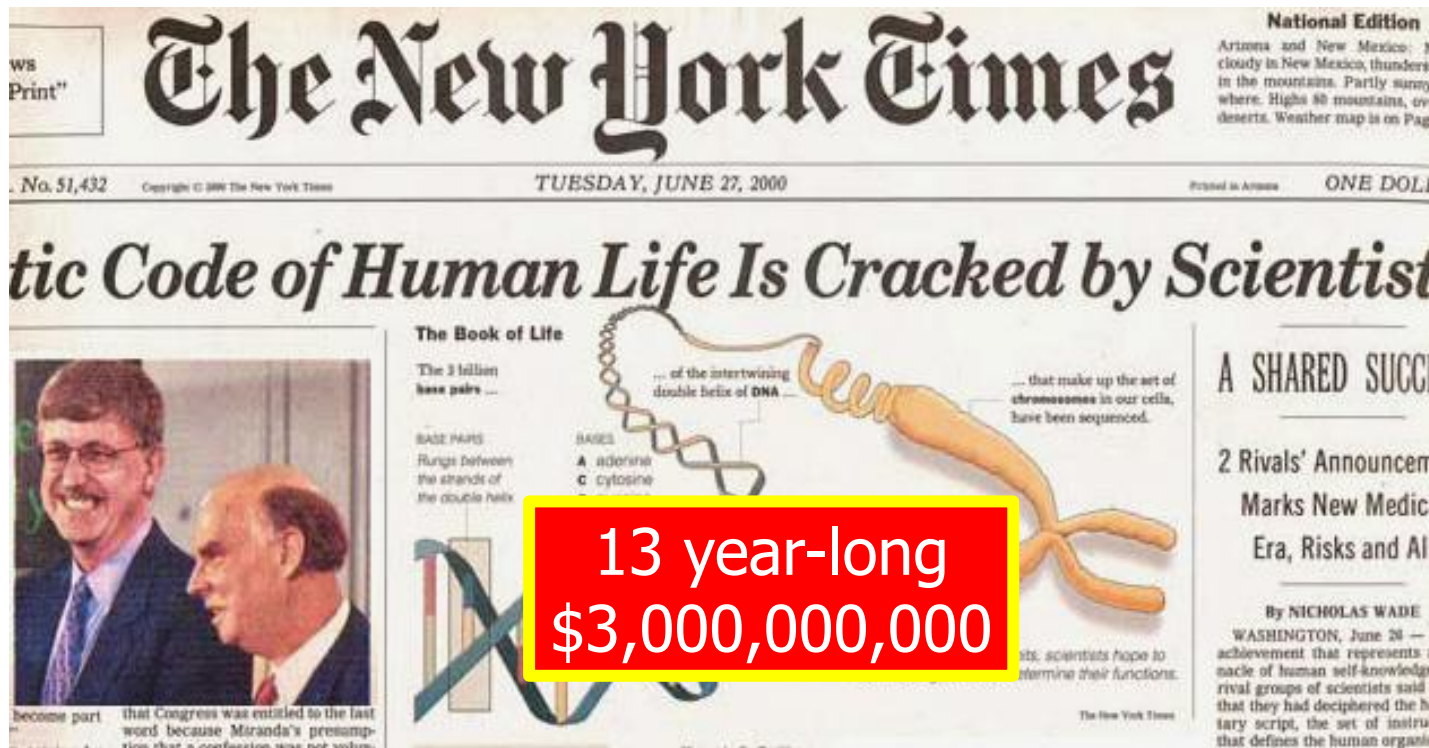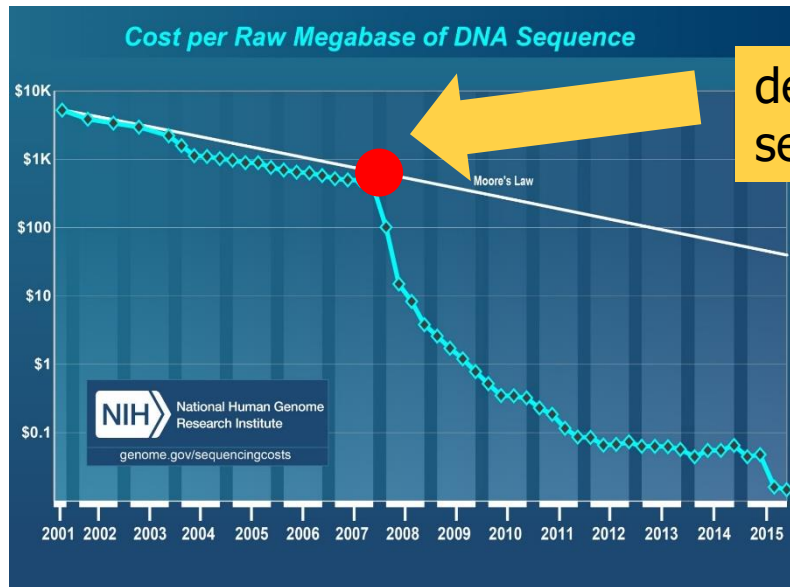

Oxford Nanopore GridION

**SAFARI**

**… and more! All produce data with different properties.**

# The Genomic Era

- 1990-2003: The Human Genome Project (HGP) provides a complete and accurate sequence of all **DNA base pairs** that make up the human genome and finds 20,000 to 25,000 human genes.
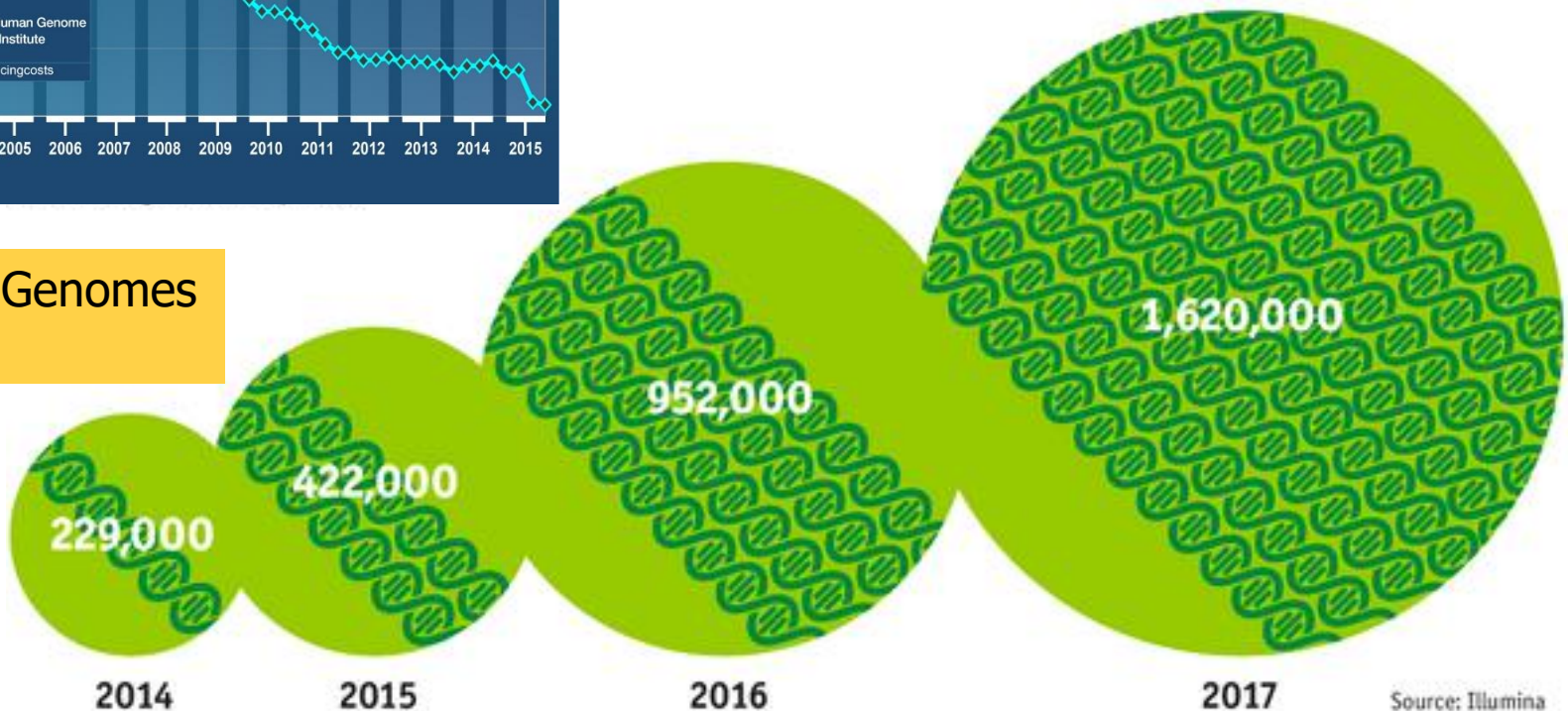


13 year-long
$3,000,000,000
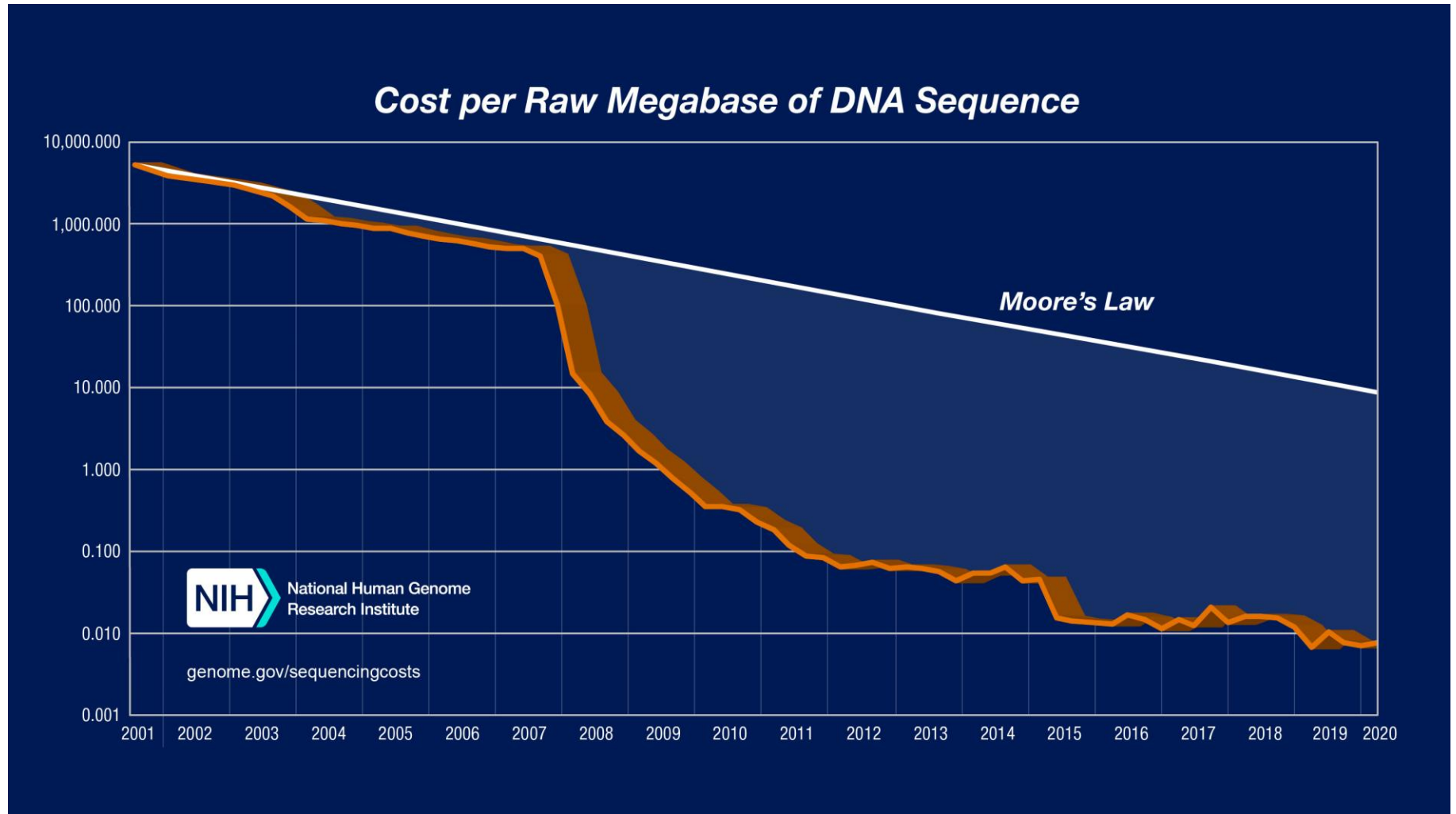
# The Genomic Era (continued)



Cost per Raw Megabase of DNA Sequence

development of high-throughput sequencing (HTS) technologies

Number of Genomes Sequenced

229,000 — 2014
422,000 — 2015
952,000 — 2016
1,620,000 — 2017

Source: Illumina

# Cost of Sequencing



**Cost per Raw Megabase of DNA Sequence**

Moore's Law

National Human Genome Research Institute

genome.gov/sequencingcosts

*From NIH (https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data)

**SAFARI**

# Cost of Sequencing (cont.)



*From NIH (https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data)

**SAFARI**

**1 Sequencing**

Billions of Short Reads

**2 Read Mapping**

Short Read

Read Alignment

Reference Genome

**Genome Analysis**

**3 Variant Calling**

reference: TTTATCGCTTCCATGACGCAG
read1:        ATCGCATCC
read2:       TATCGCATC
read3:          CATCCATGA
read4:         CGCTTCCAT
read5:              CCATGACGC
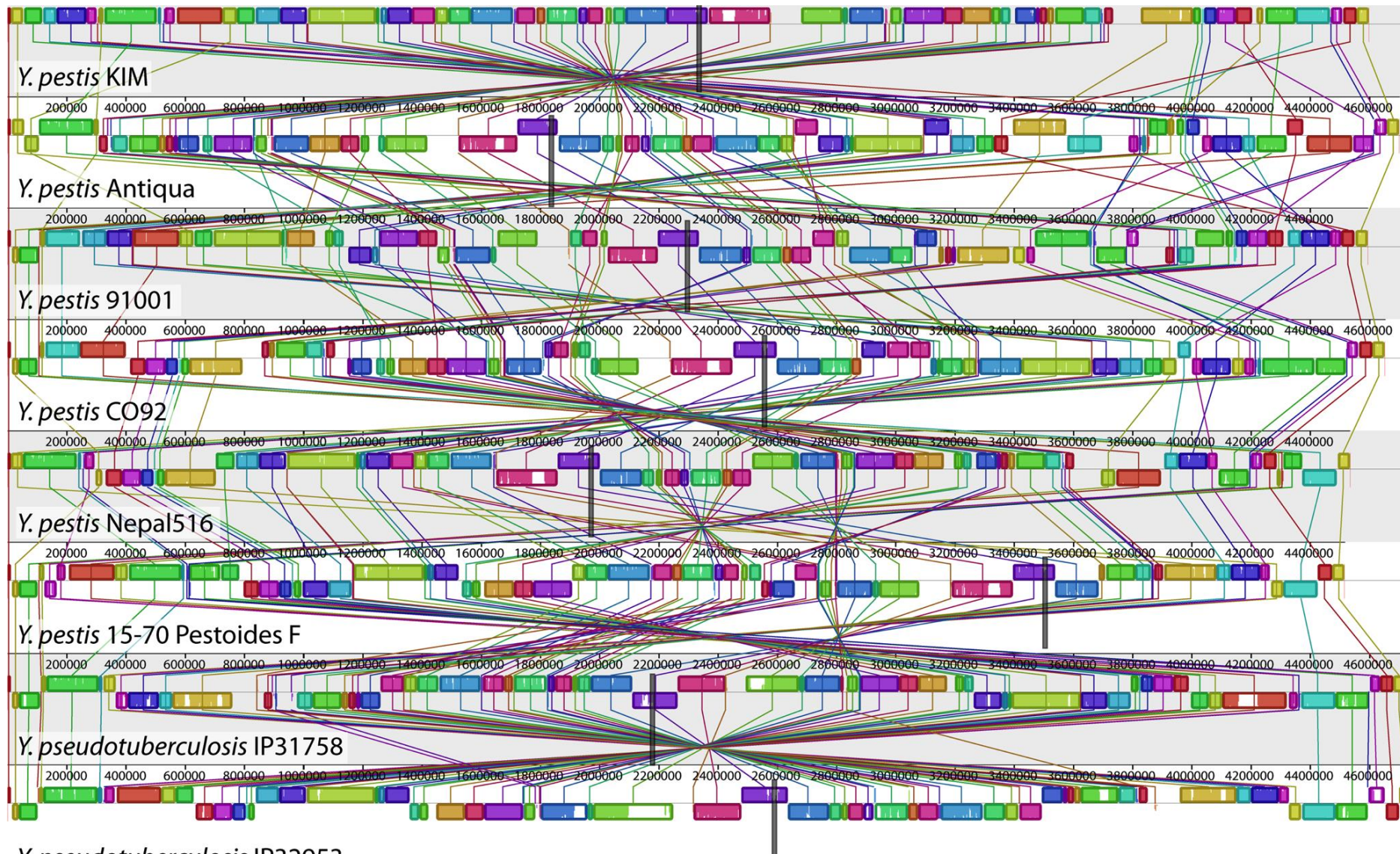read6:             TTCCATGAC

**4 Scientific Discovery**

PRESCRIPTION

# Multiple sequence alignment

```
PHDHtm                   ---------------------------------MMMMMMMMMMMMMMMMMM------
16082665   T acid   10 ----MASDRKSEGFQSGAGLIRYFEEEEIKGPALDPKLVVYMGIAVAIIVEIAKIFWPP---   (55)
13541150   T volc   10 ----MASDKKSEGFQSGAGLIRYFEEEEIKGPALDPKLVVYIGIAVAIMVELAKIFWPP---   (55)
RFAC01077  F acid   13 -MTSMAKDNQNENFQSGAGLIRYFNEEEIKGPAIDPKLIIYIGIAMGVIVELAKVFWPV---   (58)
15791336   H NRC1   10 ----MSSGQNSGGLMSSAGLVRYFDSEDSNALQIDPRSVVAVGAFFGLVVLLAQFFA-----   (53)
RAG22196   A fulg   14 MAKAPKGKAKTPPLMSSAGIMRYFEE-EKTQIKVSPKTILAAGIVTGVLIIILNAYYGLWP-   (68)
RPO01000   P abys    9 -----MAKEKTTLPPTGAGLMRFFDE-DTRAIKITPKGAVALTLILIIFEIILEVVGPRIFG   (56)
RPH01741   P hori    9 -----MAKEKTTLPPTGAGLMRFFDE-DTRAIKITPKGAIALVLILIIFEILLEVVGPRIFG   (56)
AE000914   M ther   10 ----MAKKDKKTLPPSGAGLVRYFEE-ETKGFKLTPEQVVVMSIILAVFCLVLRFSG-----   (52)
RMJ09857   M jann    9 -----MSKRESTGLATSAGLIRYMDE-TFSKIRVKPEHVIGVTVAFVIIEAILTYGRFL---   (53)
15920503   S toko   13 -MPSSKKKETVPLASMAGLIRYYEE-ENEKIKISPKLLIIISIIMVAGVIVASILIPPP--   (58)
AE006662   S solf   11 -MPSSKKKETVPVMSMAGLIRYYEE-ENEKVKISPKIVIGASLALTIIVIVITKLF-----   (55)
RPK02491   P aero   12 --MARRRKYEGLNPFVAAGLIKFSEEGELEKIKLTPRAAVVISLAIIGLLIAINLLLPPL--   (58)
RAP00437   A pern   13 -MSVRRRRERRATPVTAAGLLSFYEE-YEGKIKISPTIVVGAAILVSAVVAAAEIFLPAVP-   (59)

5803165    H sapi   49 ------------SAGTGGMWRFYTE-DSPGLKVGPVPVLVMSLLFIASVFMLEIWGKYTRS   (96)
13324684   M musc   49 ------------SAGTGGMWRFYTE-DSPGLKVGPVPVLVMSLLFIAAVFMLEIWGKYTRS   (96)
6002114    D mela   53 ------------GAGTGGMWRFYTD-DSPGIKVGPVPVLVMSLLFIASVFMLEIWGKYNRS   (100)
14574310   C eleg   32 ------------GGNNGGLWRFYTE-DSTGLKIGPVPVLVMSLVFIASVFVLEIWGKFTRS   (81)
10697176   Y lipo   41 ------------GGSSSTMLKLYTD-ESQGLKVDPVVVMVLSLGFIFSVVALEILAKVSTK   (91)
6320857    S cere   40 ------------GGSSSSIILKLYTD-EANGFRVDSLVVLFLSVGFIFSVIALELLTKFTHI   (88)
6320932    S cere   33 ------------TNSNNSIILKIYSD-EATGLRVDPLVVLFLAVGFIFSVVALEVISKVAGK   (82)
```

Example Question: If I give you a bunch of sequences, tell me where they are the same and where they are different.

**SAFARI**

# Genome Sequence Alignment: Example

17

# The Genetic Similarity Between Species



Human ~ Chimpanzee
96%

Human ~ Cat
90%

Human ~ Human
99.9%

Human ~ Cow
80%

Human ~ Banana
50-60%

Question 2: Given a bunch of short sequences, Can you identify the approximate species cluster for genomically unknown organisms?



uncleaned de Bruijn graph

http://math.oregonstate.edu/~koslickd

**SAFARI**

19

Billions of Short Reads

**1 Sequencing**

Short Read

Read Alignment

Reference Genome

**Read Mapping 2**

Bottlenecked in Mapping!!

Illumina HiSeq4000

300 M

bases/min

GAGTCAGAATTTGAC

on average

2 M

bases/min

(0.6%)

# Problem

**Need to construct
the entire genome
from many reads**

# Genome Sequencing



Large DNA molecule

Small DNA fragments

ACGTACCCCGT          TTTTTTTAATT

ACGAGCGGGT    GATACACTGTG    AAAAAAAAAA          Reads

CTAGGGACCTT          ACGACGTAGCT

# Genome Sequence Analysis

ACGTACCCGT          TTTTTTTAATT

ACGAGCGGGT    GATACACTGTG    AAAAAAAAAA    Reads

CTAGGGACCTT          ACGACGTAGCT

**Read Mapping,** method of aligning the reads against a known reference genome to **detect matches and variations.**

*De novo* **Assembly,** method of merging the reads in order to **construct** the original sequence.

Reference Genome

Original Sequence

# Read Mapping

- Map many short DNA fragments (reads) to a known reference genome with some differences allowed

Reference genome

DNA, logically physically

Reads

Mapping short reads to reference genome is challenging (billions of 50-300 base pair reads)

# Read Alignment/Verification

- **<u>Edit distance</u>** is defined as the minimum number of edits (i.e. insertions, deletions, or substitutions) needed to make the read exactly match the reference segment.

NETHERLANDS x SWITZERLAND

| N | E | - | T | H | E | R | L | A | N | D | S |
|---|---|---|---|---|---|---|---|---|---|---|---|
| S | W | I | T | Z | E | R | L | A | N | D | - |

| match |
|---|
| deletion |
| insertion |
| mismatch |

# Challenges in Read Mapping

- Need to find many mappings of each read
  - How can we find all mappings efficiently?

- Need to tolerate small variances/errors in each read
  - Each individual is different: Subject's DNA may slightly differ from the reference (Mismatches, insertions, deletions)
  - How can we efficiently map each read with up to $e$ errors present?

- Need to map each read very fast (i.e., performance is important)
  - Human DNA is 3.2 billion base pairs long → Millions to billions of reads (State-of-the-art mappers take weeks to map a human's DNA)
  - How can we design a much higher performance read mapper?

# Our First Step: Comprehensive Mapping

- **+ Guaranteed to find *all* mappings → sensitive**
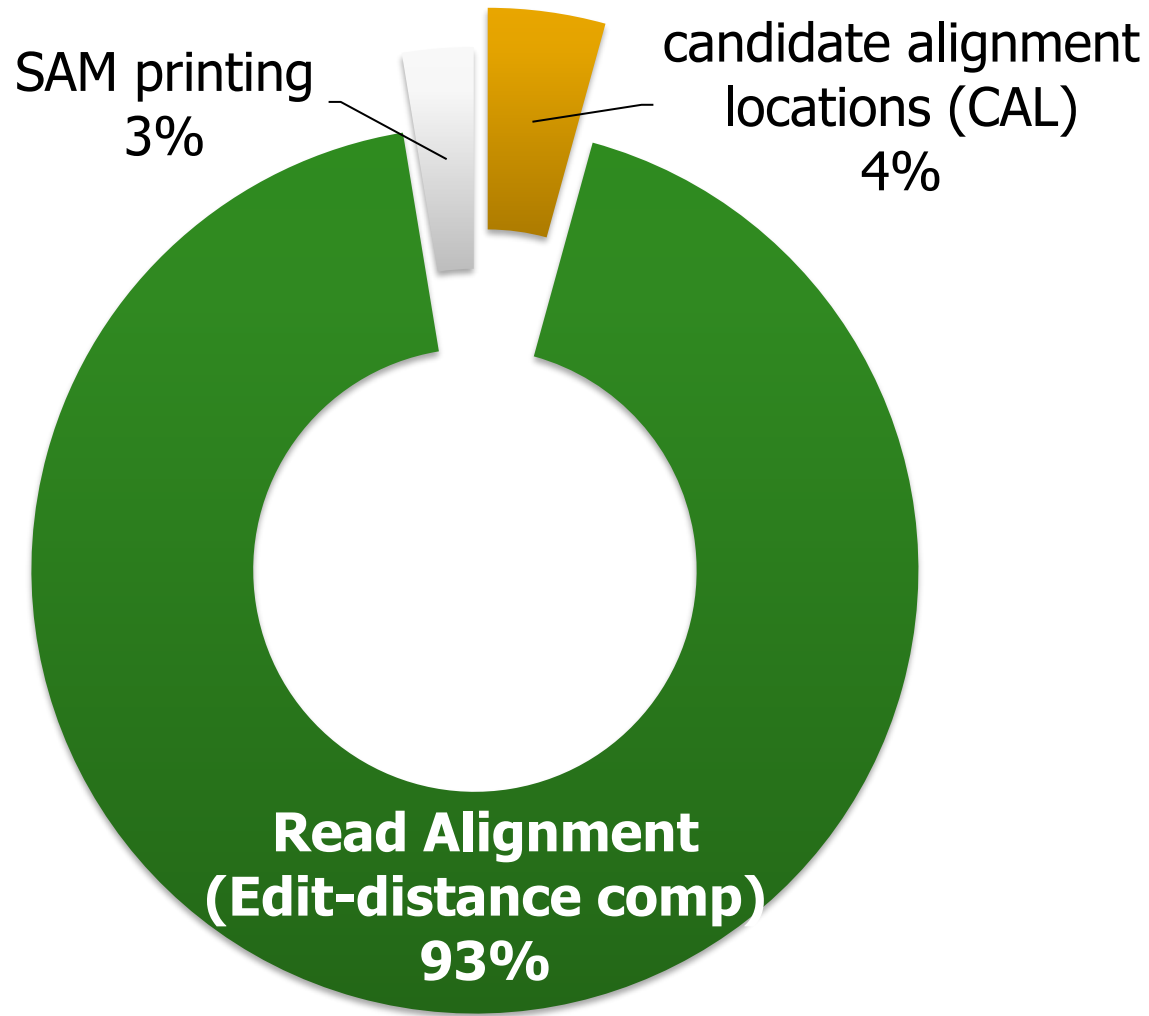- **+ Can tolerate up to *e* errors**

**nature genetics**

# Personalized copy number and segmental duplication maps using next-generation sequencing

Can Alkan[1,2], Jeffrey M Kidd[1], Tomas Marques-Bonet[1,3], Gozde Aksay[1], Francesca Antonacci[1], Fereydoun Hormozdiari[4], Jacob O Kitzman[1], Carl Baker[1], Maika Malig[1], Onur Mutlu[5], S Cenk Sahinalp[4], Richard A Gibbs[6] & Evan E Eichler[1,2]

Alkan+, **"Personalized copy number and segmental duplication maps using next-generation sequencing"**, Nature Genetics 2009.

# Read Mapping Execution Time Breakdown



SAM printing
3%

candidate alignment
locations (CAL)
4%

**Read Alignment
(Edit-distance comp)
93%**

# The Read Mapping Bottleneck



Illumina HiSeq4000

300 Million bases/minute

2 Million bases/minute

150X slower

**SAFARI**

# Idea

**Filter fast** before you align

Minimize costly
"approximate string comparisons"

# Our First Filter: Pure Software Approach

- Download the source code and try for yourself
  - [Download link to FastHASH](#)

**BMC Genomics**

**PROCEEDINGS**                                    **Open Access**

# Accelerating read mapping with FastHASH

Hongyi Xin[1], Donghyuk Lee[1], Farhad Hormozdiari[2], Samihan Yedkar[1], Onur Mutlu[1*], Can Alkan[3*]

31

# Shifted Hamming Distance: SIMD Acceleration

https://github.com/CMU-SAFARI/Shifted-Hamming-Distance

Sequence analysis

## Shifted Hamming distance: a fast and accurate SIMD-friendly filter to accelerate alignment verification in read mapping

Hongyi Xin[1,*], John Greth[2], John Emmons[2], Gennady Pekhimenko[1], Carl Kingsford[3], Can Alkan[4,*] and Onur Mutlu[2,*]

Xin+, **"Shifted Hamming Distance: A Fast and Accurate SIMD-friendly Filter to Accelerate Alignment Verification in Read Mapping"**, **Bioinformatics 2015.**

# GateKeeper: FPGA-Based Alignment Filtering

**Alignment Filter** $+$  $=$ **1**st FPGA-based Alignment Filter.

Low Speed & High Accuracy
Medium Speed, Medium Accuracy
High Speed, Low Accuracy

$x10^{12}$ mappings

$x10^{3}$ mappings

Billions of Short Reads

**1** High throughput DNA sequencing (HTS) technologies

**2** Read Pre-Alignment Filtering
Fast & Low False Positive Rate

**3** Read Alignment
Slow & Zero False Positives

# GateKeeper: FPGA-Based Alignment Filtering

- Mohammed Alser, Hasan Hassan, Hongyi Xin, Oguz Ergin, Onur Mutlu, and Can Alkan
  **"GateKeeper: A New Hardware Architecture for Accelerating Pre-Alignment in DNA Short Read Mapping"**
  ***Bioinformatics***, [published online, May 31], 2017.
  [Source Code]
  [Online link at Bioinformatics Journal]

## GateKeeper: a new hardware architecture for accelerating pre-alignment in DNA short read mapping

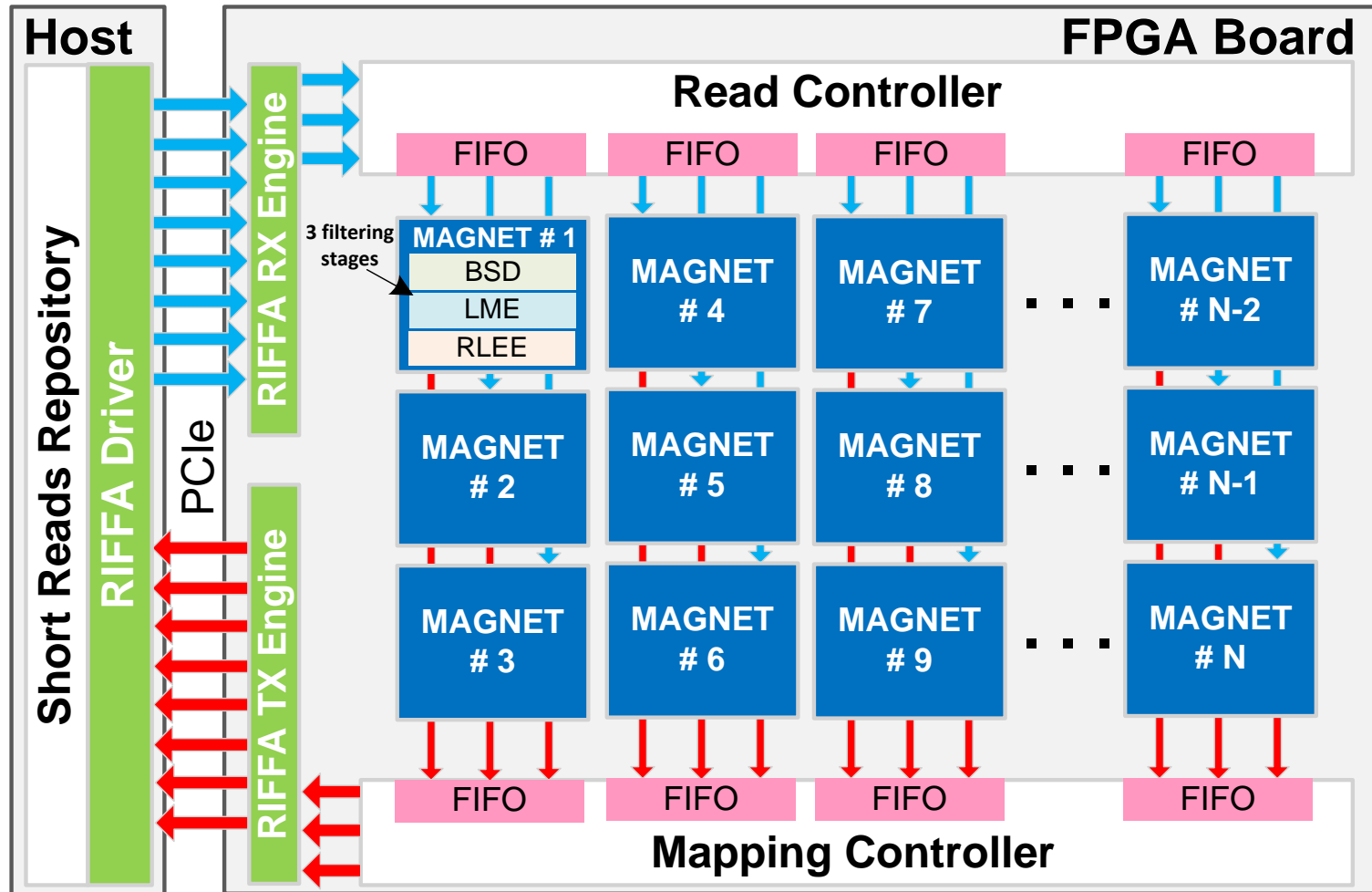Mohammed Alser ✉, Hasan Hassan, Hongyi Xin, Oğuz Ergin, Onur Mutlu ✉, Can Alkan ✉

# MAGNET Accelerator [Alser+, TIR 2017]

**SAFARI**

# Newest Work: Shouji [Alser+, Bioinformatics 2019]

Mohammed Alser, Hasan Hassan, Akash Kumar, Onur Mutlu, and Can Alkan,
**"Shouji: A Fast and Efficient Pre-Alignment Filter for Sequence Alignment"**
***Bioinformatics***, [published online, March 28], 2019.
[Source Code]
[Online link at Bioinformatics Journal]

Sequence alignment

# Shouji: a fast and efficient pre-alignment filter for sequence alignment

Mohammed Alser[1,2,3,*], Hasan Hassan[1], Akash Kumar[2], Onur Mutlu[1,3,*] and Can Alkan[3,*]

[1]Computer Science Department, ETH Zürich, Zürich 8092, Switzerland, [2]Chair for Processor Design, Center For Advancing Electronics Dresden, Institute of Computer Engineering, Technische Universität Dresden, 01062 Dresden, Germany and [3]Computer Engineering Department, Bilkent University, 06800 Ankara, Turkey

*To whom correspondence should be addressed.

SAFARI

36

# DNA Read Mapping & Filtering

- Problem: **Heavily bottlenecked by Data Movement**

- GateKeeper FPGA performance limited by DRAM bandwidth [Alser+, Bioinformatics 2017]

- Ditto for SHD on SIMD [Xin+, Bioinformatics 2015]

- Solution: Processing-in-memory can alleviate the bottleneck

- However, we need to design mapping & filtering algorithms to fit processing-in-memory

# In-Memory DNA Sequence Analysis

- Jeremie S. Kim, Damla Senol Cali, Hongyi Xin, Donghyuk Lee, Saugata Ghose, Mohammed Alser, Hasan Hassan, Oguz Ergin, Can Alkan, and Onur Mutlu,
  **"GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies"**
  *BMC Genomics*, 2018.
  *Proceedings of the 16th Asia Pacific Bioinformatics Conference* (**APBC**), Yokohama, Japan, January 2018.
  arxiv.org Version (pdf)

# GRIM-Filter: Fast seed location filtering in DNA read mapping using processing-in-memory technologies

Jeremie S. Kim[1,6*], Damla Senol Cali[1], Hongyi Xin[2], Donghyuk Lee[3], Saugata Ghose[1], Mohammed Alser[4], Hasan Hassan[6], Oguz Ergin[5], Can Alkan[4*] and Onur Mutlu[6,1*]

# Quick Note: Key Principles and Results

- **Two key principles:**
  - Exploit the structure of the genome to minimize computation
  - Morph and exploit the structure of the underlying hardware to maximize performance and efficiency

- **Algorithm-architecture co-design** for DNA read mapping
  - **Speeds up** read mapping by **~300X (sometimes more)**
  - **Improves accuracy** of read mapping in the presence of errors

Xin et al., "Accelerating Read Mapping with FastHASH," BMC Genomics 2013.

Xin et al., "Shifted Hamming Distance: A Fast and Accurate SIMD-friendly Filter to Accelerate Alignment Verification in Read Mapping," Bioinformatics 2015.

Alser et al., "GateKeeper: A New Hardware Architecture for Accelerating Pre-Alignment in DNA Short Read Mapping," Bioinformatics 2017.

Kim et al., "Genome Read In-Memory (GRIM) Filter," BMC Genomics 2018.

# New Genome Sequencing Technologies

## Nanopore sequencing technology and tools for genome assembly: computational analysis of the current state, bottlenecks and future directions

Damla Senol Cali ✉, Jeremie S Kim, Saugata Ghose, Can Alkan, Onur Mutlu

Oxford Nanopore MinION

Senol Cali+, "**Nanopore Sequencing Technology and Tools for Genome Assembly: Computational Analysis of the Current State, Bottlenecks and Future Directions**," Briefings in Bioinformatics, 2018.
[Preliminary arxiv.org version]

# Future of Genome Sequencing & Analysis



MinION from ONT
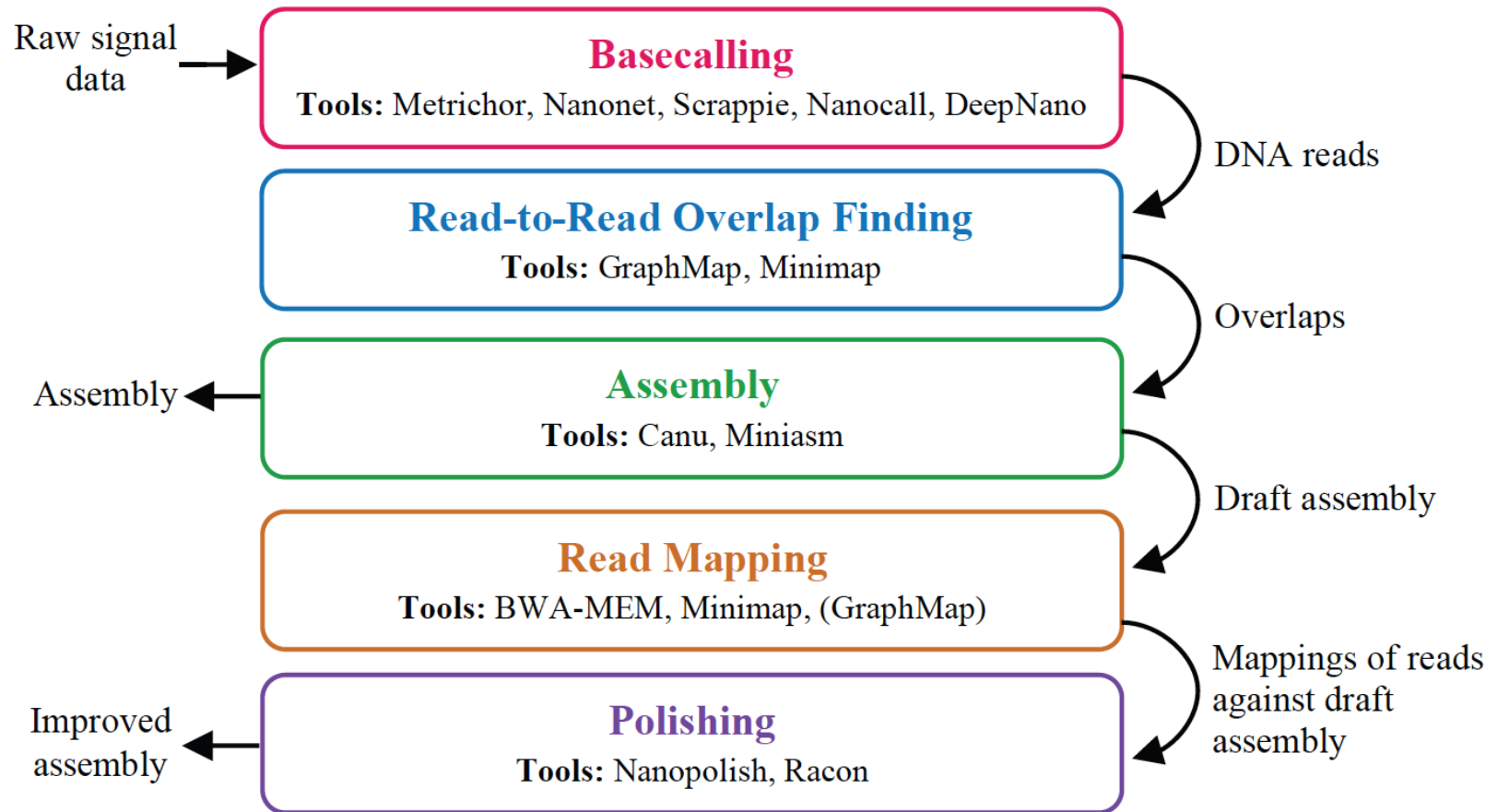
SmidgION from ONT

# Nanopore Genome Assembly Pipeline

Raw signal data →

**Basecalling**
**Tools:** Metrichor, Nanonet, Scrappie, Nanocall, DeepNano

→ DNA reads

**Read-to-Read Overlap Finding**
**Tools:** GraphMap, Minimap

→ Overlaps

Assembly ←

**Assembly**
**Tools:** Canu, Miniasm

→ Draft assembly

**Read Mapping**
**Tools:** BWA-MEM, Minimap, (GraphMap)

→ Mappings of reads against draft assembly

Improved assembly ←

**Polishing**
**Tools:** Nanopolish, Racon

**Figure 1. The analyzed genome assembly pipeline using nanopore sequence data, with its five steps and the associated tools for each step.**

Senol Cali+, "**Nanopore Sequencing Technology and Tools for Genome Assembly,**" Briefings in Bioinformatics, 2018.

**SAFARI**

42

# Recall Our Dream (from 2007)

- An embedded device that can perform comprehensive genome analysis in real time (within a minute)

- Still a long ways to go
    - Energy efficiency
    - Performance (latency)
    - Security
    - **Huge memory bottleneck**

# Why Do We Care? An Example from 2020

## 200 Oxford Nanopore sequencers have left UK for China, to support rapid, near-sample coronavirus sequencing for outbreak surveillance

Fri 31st January 2020

Following extensive support of, and collaboration with, public health professionals in China, Oxford Nanopore has shipped an additional 200 MinION sequencers and related consumables to China. These will be used to support the ongoing surveillance of the current coronavirus outbreak, adding to a large number of the devices already installed in the country.



Each MinION sequencer is approximately the size of a stapler, and can provide rapid sequence information about the coronavirus.



700Kg of Oxford Nanopore sequencers and consumables are on their way for use by Chinese scientists in understanding the current coronavirus outbreak.



**SAFARI**

44

# Sequencing of COVID-19
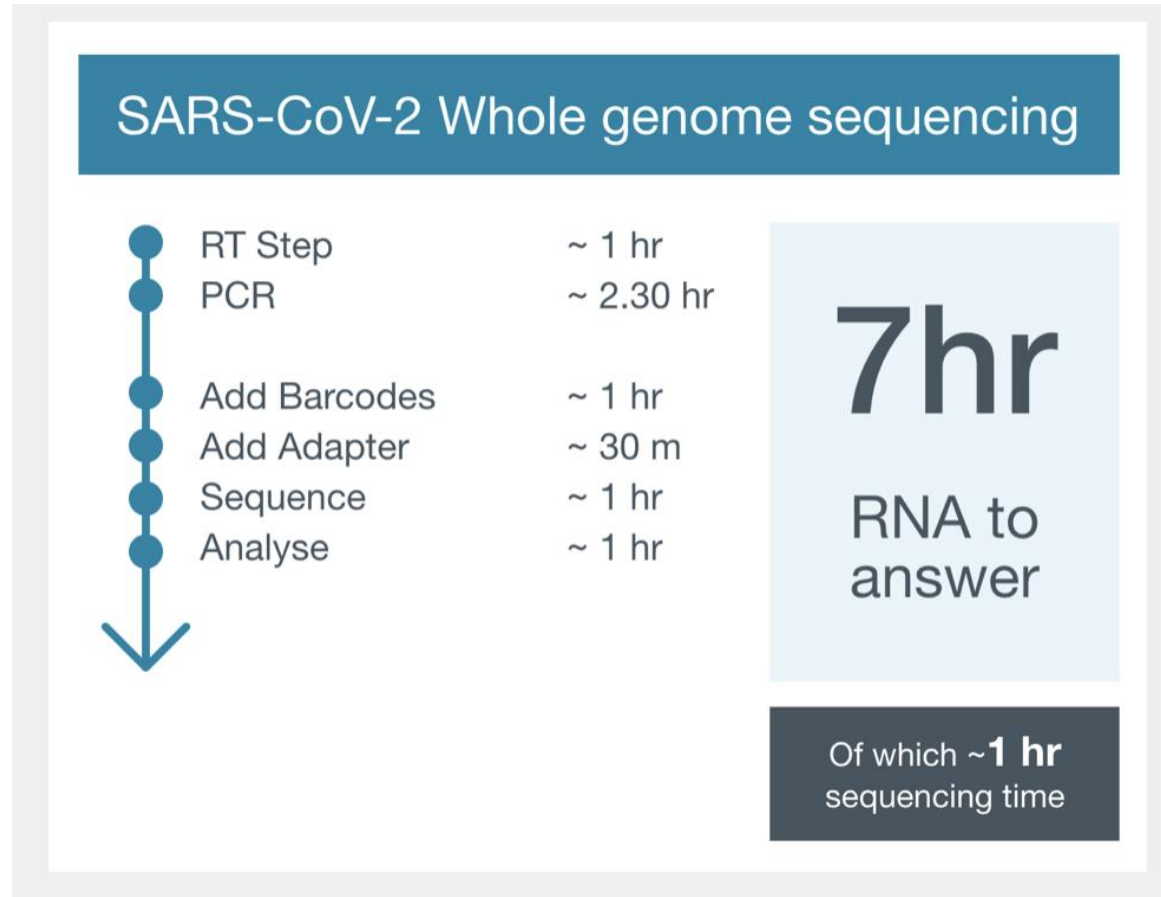
- **Whole genome sequencing (WGS) and sequence data analysis are important**
  - To detect the virus from a human sample such as saliva, Bronchoalveolar fluid etc.
  - To understand the sources and modes of transmission of the virus
  - To discover the genomic characteristics of the virus, and compare with better-known viruses (e.g., 02-03 SARS epidemic)
  - To design and evaluate the diagnostic tests

- **Two key areas of COVID-19 genomic research**
  - To sequence the genome of the virus itself, COVID-19, in order to track the mutations in the virus.
  - To explore the genes of infected patients. This analysis can be used to understand why some people get more severe symptoms than others, as well as, help with the development of new treatments in the future.

# COVID-19 Nanopore Sequencing (I)

**SAFARI**

# COVID-19 Nanopore Sequencing (II)

- From ONT (https://nanoporetech.com/covid-19/overview)

# Future of Genome Sequencing & Analysis



MinION from ONT

SmidgION from ONT

# More on Genome Analysis: Another Talk

- Onur Mutlu,
  **"Accelerating Genome Analysis: A Primer on an Ongoing Journey"**
  *Keynote talk at 2nd Workshop on Accelerator Architecture in Computational Biology and Bioinformatics* (**AACBB**)*, Washington, DC, USA, February 2019.
  [Slides (pptx)(pdf)]
  [Video]

## Accelerating Genome Analysis

### A Primer on an Ongoing Journey

Onur Mutlu

omutlu@gmail.com

https://people.inf.ethz.ch/omutlu

16 February 2019

AACBB Keynote Talk

**SAFARI**        **ETH**zürich        **Carnegie Mellon**

https://www.youtube.com/watch?v=hPnSmfwu2-A

# Recall Our Axiom

To achieve the highest energy efficiency and performance:

## we must take the expanded view
of computer architecture

| Problem |
|---|
| Algorithm |
| Program/Language |
| System Software |
| SW/HW Interface |
| Micro-architecture |
| Logic |
| Devices |
| Electrons |

**Co-design across the hierarchy:**
**Algorithms to devices**

**Specialize as much as possible**
**within the design goals**

# Our Axiom Applies Well to Genome Analysis

**Computer Architecture (expanded view)**

| Problem |
| Algorithm |
| Program/Language |
| System Software |
| SW/HW Interface |
| Micro-architecture |
| Logic |
| Devices |
| Electrons |

# Algorithm-Arch-Device Co-Design is Critical

**Computer Architecture (expanded view)**

| Problem |
| --- |
| Algorithm |
| Program/Language |
| System Software |
| SW/HW Interface |
| Micro-architecture |
| Logic |
| Devices |
| Electrons |

# GenASM [MICRO 2020]

- Damla Senol Cali, Gurpreet S. Kalsi, Zulal Bingol, Can Firtina, Lavanya Subramanian, Jeremie S. Kim, Rachata Ausavarungnirun, Mohammed Alser, Juan Gomez-Luna, Amirali Boroumand, Anant Nori, Allison Scibisz, Sreenivas Subramoney, Can Alkan, Saugata Ghose, and Onur Mutlu,
  **"GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis"**
  *Proceedings of the 53rd International Symposium on Microarchitecture* (**MICRO**), Virtual, October 2020.
  [ARM Research Summit Talk Video (21 minutes)]
  [ARM Research Summit Short Talk Video (15 minutes)]
  [ARM Research Summit Short Talk Video and Q&A (31 minutes)]
  [ARM Research Summit Talk Slides (pptx) (pdf)]
  [ARM Research Summit Short Talk Slides (pptx) (pdf)]

## GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis

Damla Senol Cali[†⋈]   Gurpreet S. Kalsi[⋈]   Zülal Bingöl[▽]   Can Firtina[◇]   Lavanya Subramanian[‡]   Jeremie S. Kim[◇†]

Rachata Ausavarungnirun[⊙]   Mohammed Alser[◇]   Juan Gomez-Luna[◇]   Amirali Boroumand[†]   Anant Nori[⋈]

Allison Scibisz[†]   Sreenivas Subramoney[⋈]   Can Alkan[▽]   Saugata Ghose[⋆†]   Onur Mutlu[◇†▽]

[†]*Carnegie Mellon University*   [⋈]*Processor Architecture Research Lab, Intel Labs*   [▽]*Bilkent University*   [◇]*ETH Zürich*
[‡]*Facebook*   [⊙]*King Mongkut's University of Technology North Bangkok*   [⋆]*University of Illinois at Urbana–Champaign*
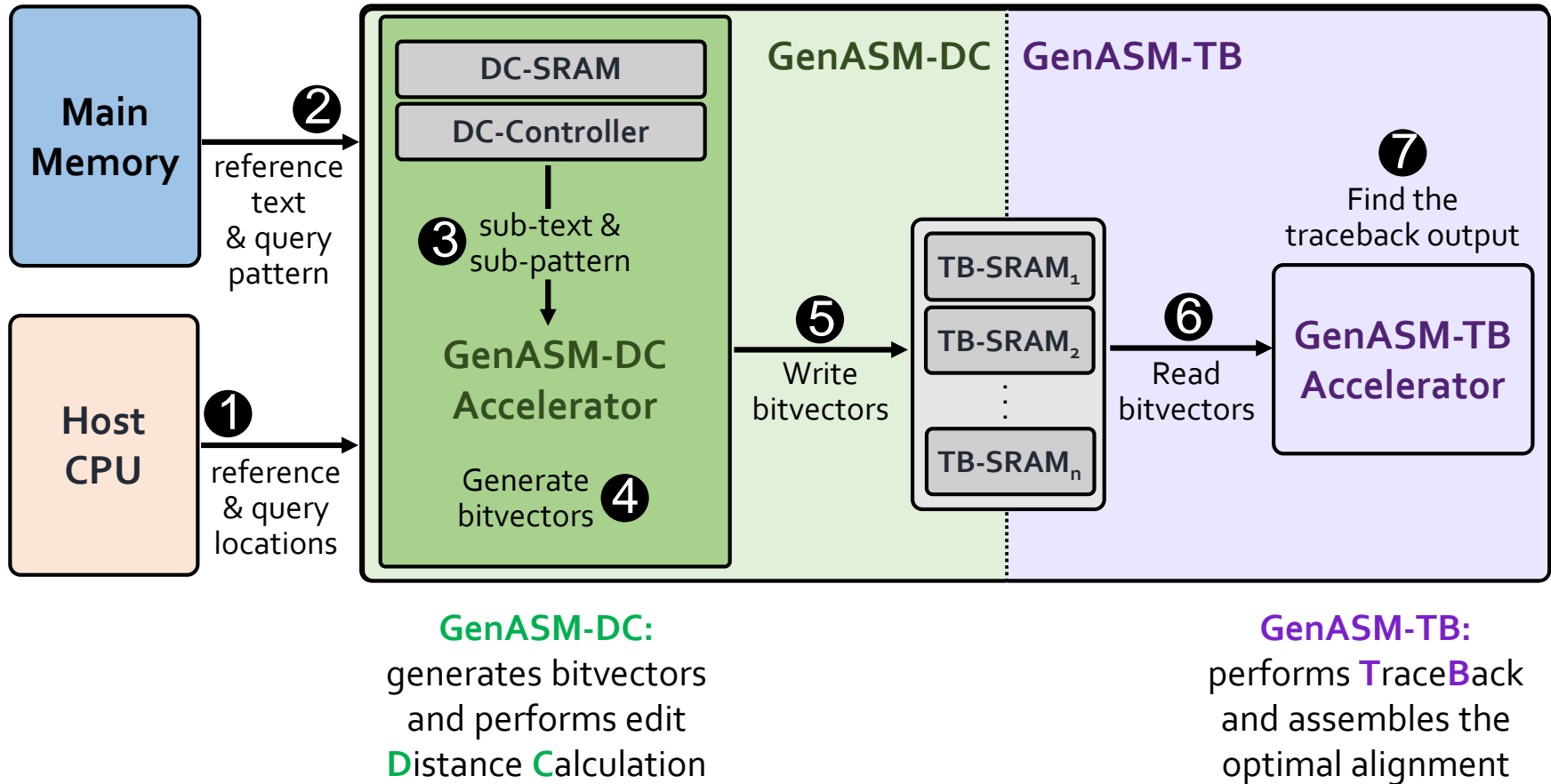
**SAFARI**

53

# Problem & Our Goal

❑ Multiple steps of read mapping require *approximate string matching*
  o ASM enables read mapping to account for sequencing errors and genetic variations in the reads

❑ ASM makes up a significant portion of read mapping (more than 70%)

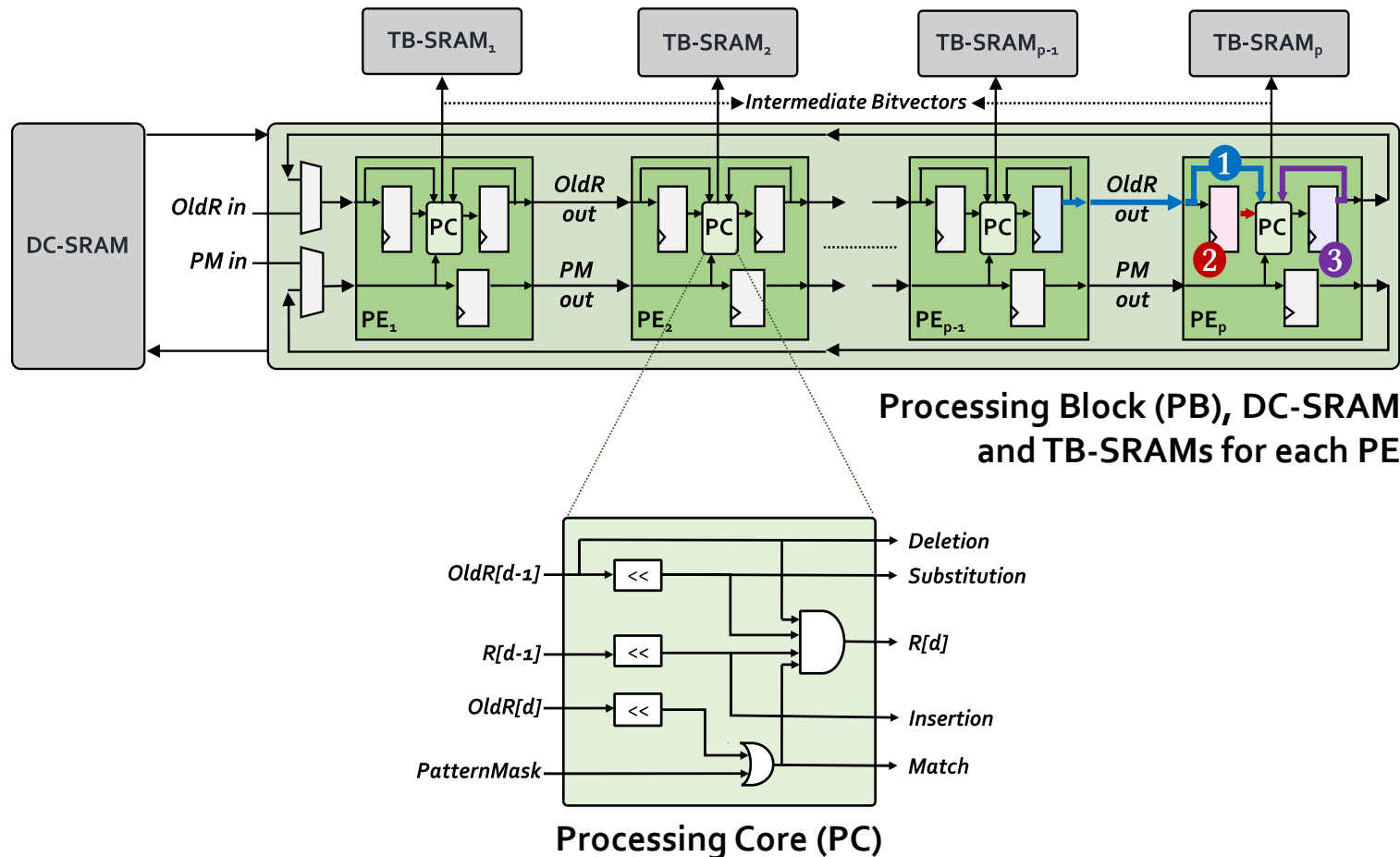❑ One of the major bottlenecks of genome sequence analysis

## Our Goal:
Accelerate approximate string matching by designing a fast and flexible framework, which can be used to accelerate *multiple steps* of the genome sequence analysis pipeline

# GenASM: Hardware Design



GenASM-DC:
generates bitvectors
and performs edit
Distance Calculation

GenASM-TB:
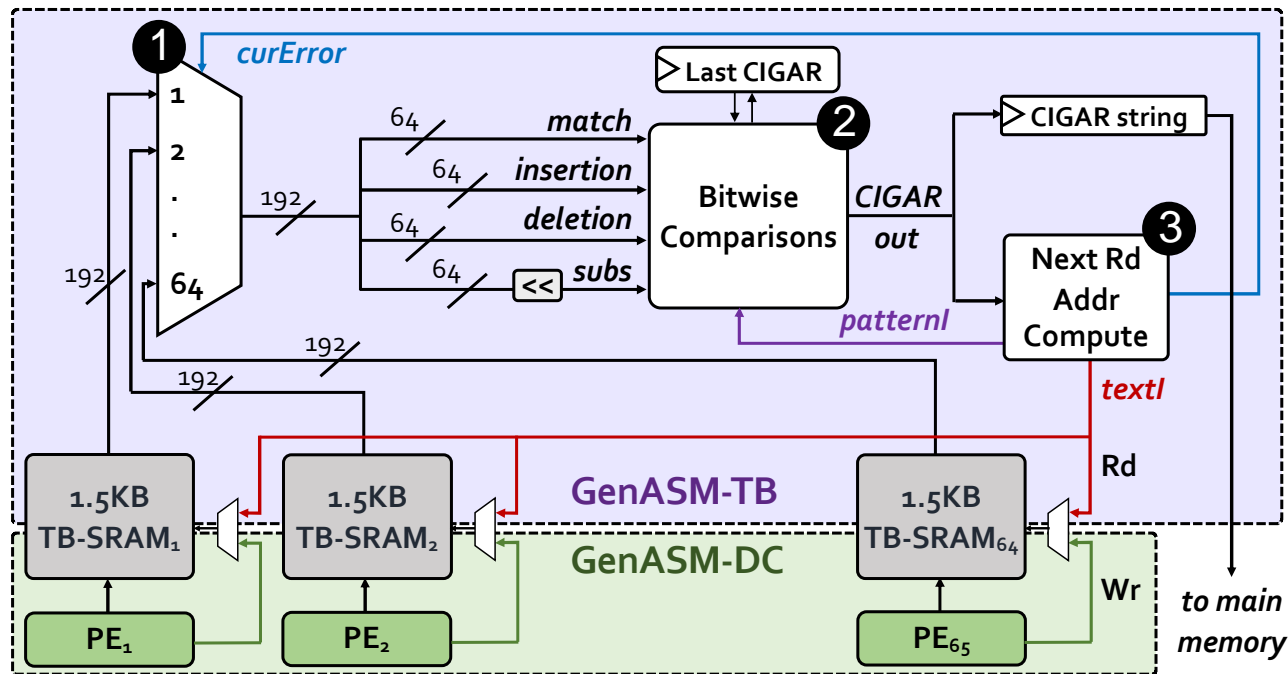performs TraceBack
and assembles the
optimal alignment

# GenASM-DC: Hardware Design

❏ **Linear cyclic systolic array** based accelerator

  o Optimized to reduce memory bandwidth and memory footprint



**Processing Block (PB), DC-SRAM and TB-SRAMs for each PE**

**Processing Core (PC)**

# GenASM-TB: Hardware Design



❑ Very simple logic:

1) Reads the bitvectors from one of the TB-SRAMs using the computed address

2) Performs the required bitwise comparisons to find the traceback output for the current position

3) Computes the next TB-SRAM address to read the new set of bitvectors

# GenASM [MICRO 2020]

Damla Senol Cali, Gurpreet S. Kalsi, Zulal Bingol, Can Firtina, Lavanya Subramanian, Jeremie S. Kim, Rachata Ausavarungnirun, Mohammed Alser, Juan Gomez-Luna, Amirali Boroumand, Anant Nori, Allison Scibisz, Sreenivas Subramoney, Can Alkan, Saugata Ghose, and Onur Mutlu,

**"GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis"**

*Proceedings of the 53rd International Symposium on Microarchitecture* (**MICRO**), Virtual, October 2020.

## GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis

Damla Senol Cali[†⋈]    Gurpreet S. Kalsi[⋈]    Zülal Bingöl[▽]    Can Firtina[◇]    Lavanya Subramanian[‡]    Jeremie S. Kim[◇†]
Rachata Ausavarungnirun[⊙]    Mohammed Alser[◇]    Juan Gomez-Luna[◇]    Amirali Boroumand[†]    Anant Nori[⋈]
Allison Scibisz[†]    Sreenivas Subramoney[⋈]    Can Alkan[▽]    Saugata Ghose[⋆†]    Onur Mutlu[◇†▽]

[†]*Carnegie Mellon University*    [⋈]*Processor Architecture Research Lab, Intel Labs*    [▽]*Bilkent University*    [◇]*ETH Zürich*
[‡]*Facebook*    [⊙]*King Mongkut's University of Technology North Bangkok*    [⋆]*University of Illinois at Urbana–Champaign*

# Recall Our Dream (from 2007)

- An embedded device that can perform comprehensive genome analysis in real time (within a minute)

- Still a long ways to go
  - Energy efficiency
  - Performance (latency)
  - Security
  - **Huge memory bottleneck**

# Four Key Directions

- Fundamentally Secure/Reliable/Safe Architectures

- Fundamentally Energy-Efficient Architectures
  - Memory-centric (Data-centric) Architectures

- Fundamentally Low-Latency and Predictable Architectures

- Architectures for AI/ML, Genomics, Medicine, Health

SAFARI

# Memory & Storage

# **Computer Architecture**
# Lecture 3a: Introduction to Genome Sequence Analysis

Prof. Onur Mutlu

ETH Zürich

Fall 2020

24 September 2020