# 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**

# Overview

- **System design for bioinformatics** is a critical problem
  - It has large scientific, medical, societal, personal implications

- This talk is about accelerating a key step in bioinformatics: genome sequence analysis
  - In particular, read mapping

- Many bottlenecks exist in accessing and manipulating huge amounts of genomic data during analysis

- We will cover various recent ideas to accelerate read mapping
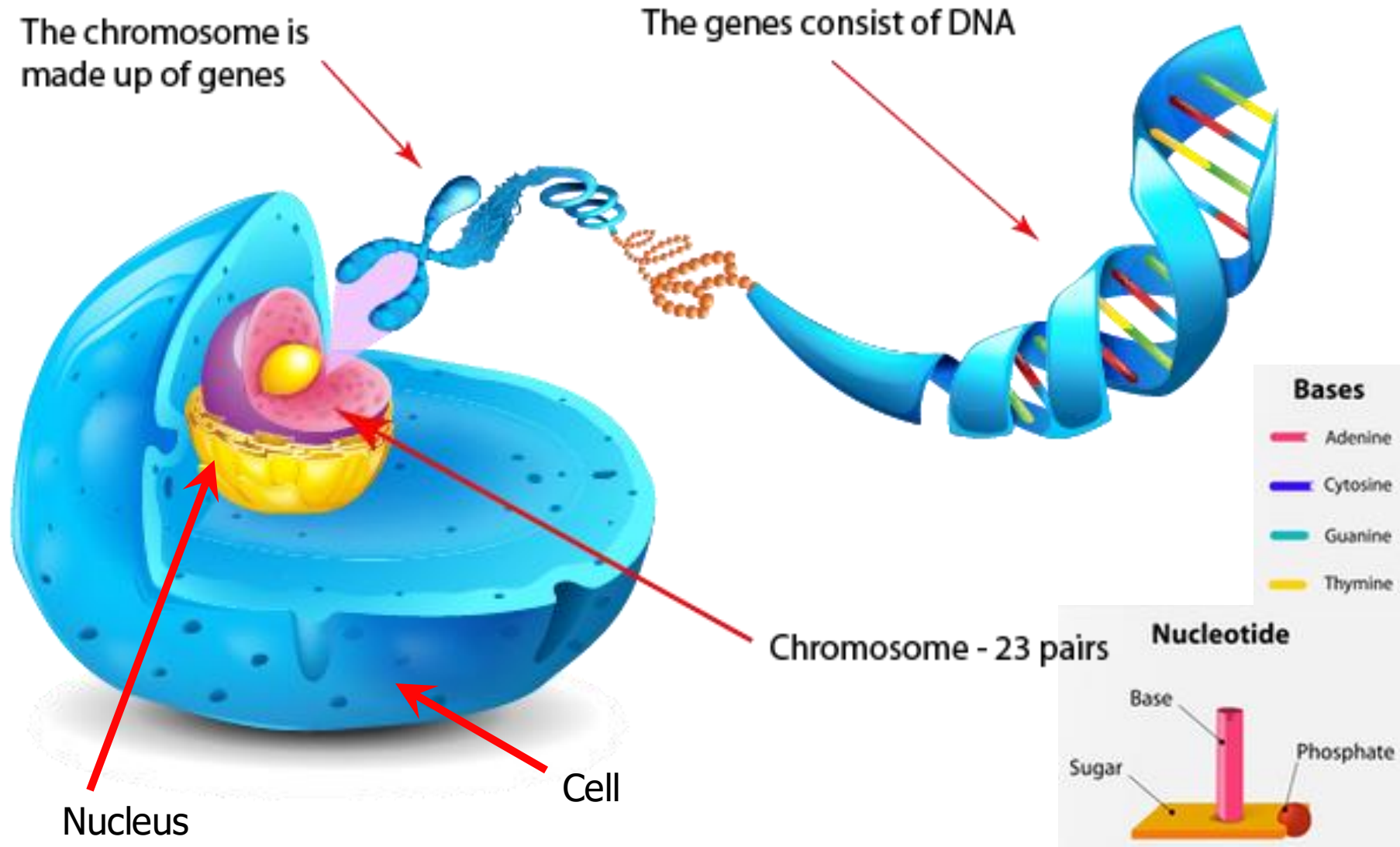  - My personal journey since September 2006

# Our Dream (in 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?
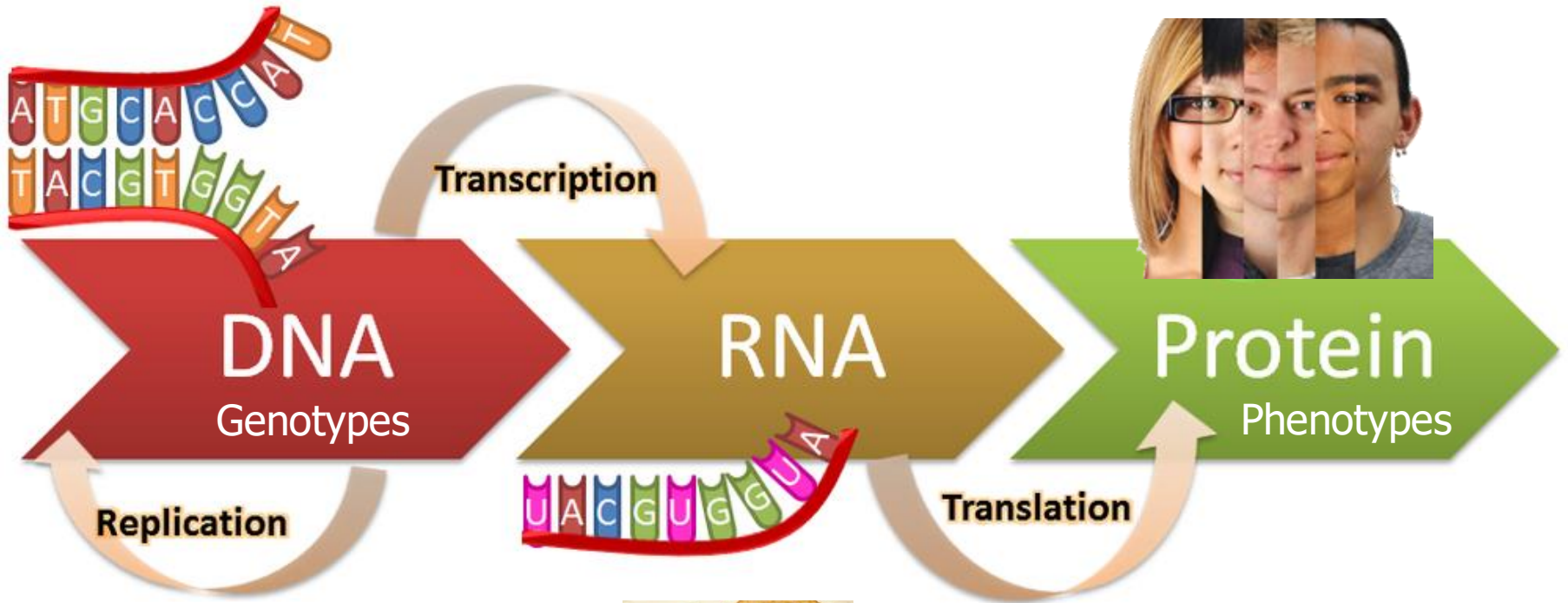  - . . .

# Agenda

- **The Problem: DNA Read Mapping**
  - State-of-the-art Read Mapper Design

- **Algorithmic Acceleration**
  - Exploiting Structure of the Genome
  - Exploiting SIMD Instructions

- **Hardware Acceleration**
  - Specialized Architectures
  - Processing in Memory

- **Future Opportunities: New Sequencing Technologies**
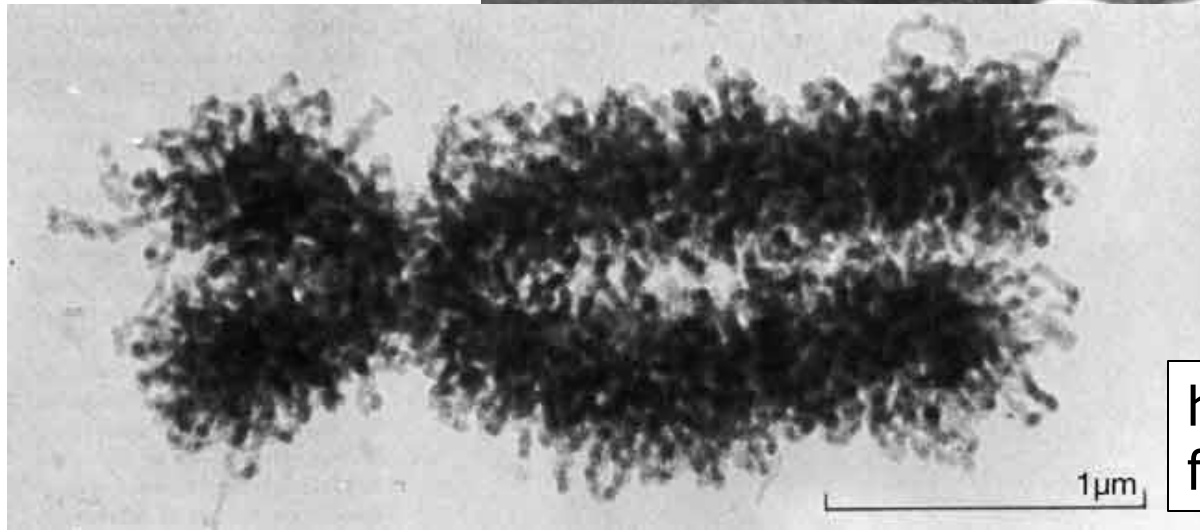
**SAFARI**

# 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
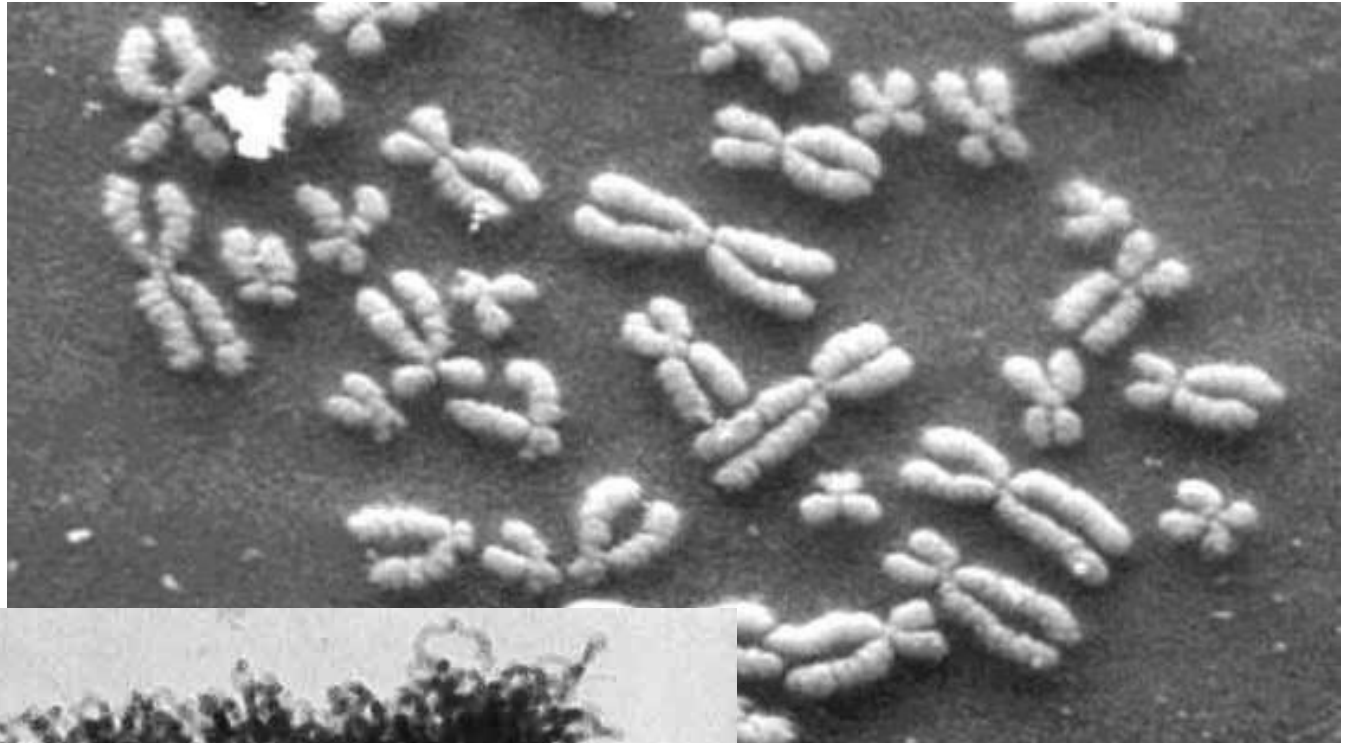- Guanine
- Thymine

**Nucleotide**

Base

Sugar

Phosphate

# The Central Dogma of Molecular Biology

# DNA Under Electron Microscope



human chromosome #12
from HeLa's cell

SAFARI

# DNA Sequencing

- Goal:
  - Find the complete sequence of A, C, G, T's in DNA.

- 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)
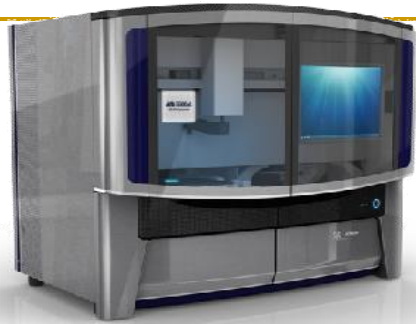
# Untangling Yarn Balls & DNA Sequencing

# Genome Sequencers

Roche/454

AB SOLiD

Illumina MiSeq

Complete Genomics

Illumina HiSeq2000

Pacific Biosciences RS

Oxford Nanopore MinION

Illumina NovaSeq 6000

Oxford Nanopore GridION

Ion Torrent PGM

Ion Torrent Proton

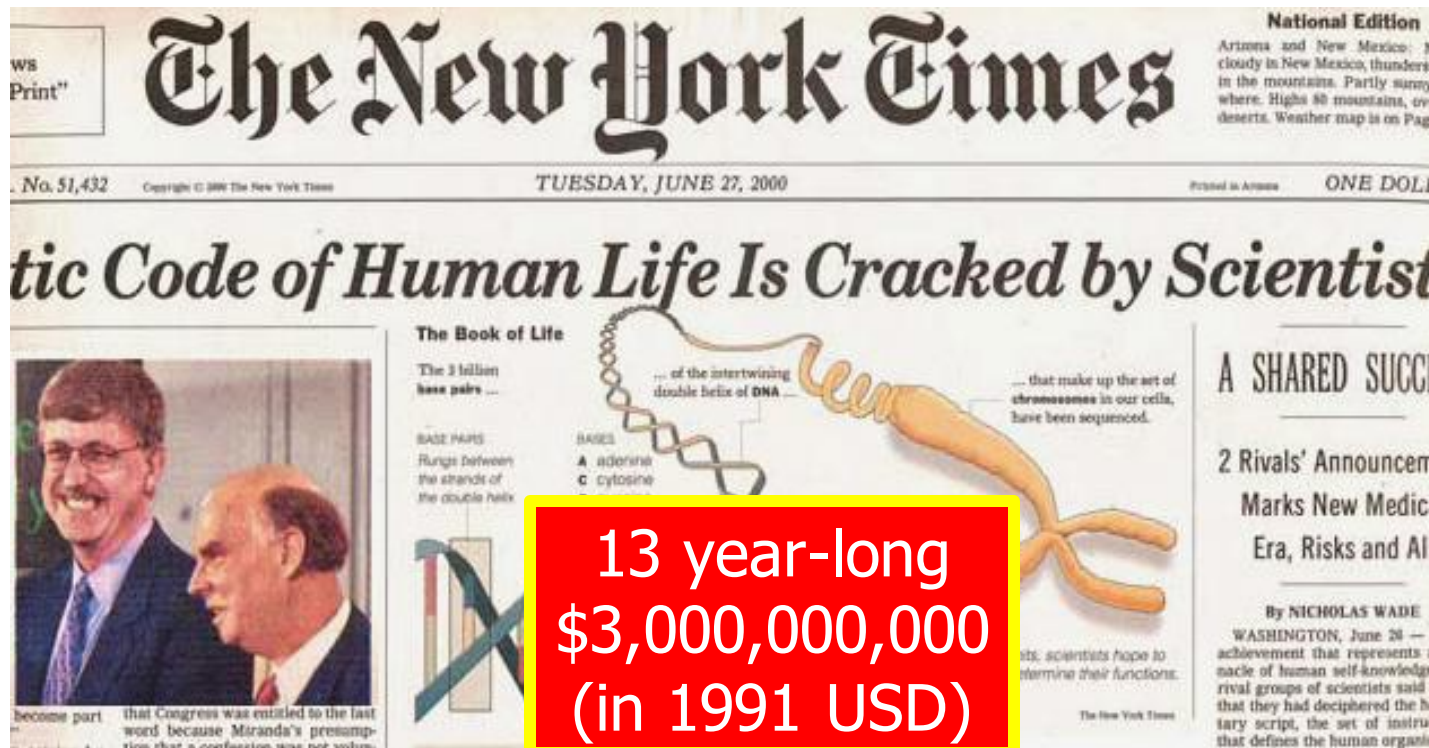**SAFARI**

**… 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
(in 1991 USD)

# The Genomic Era (continued)



development of high-throughput sequencing (HTS) technologies

Cost per Raw Megabase of DNA Sequence

Number of Genomes Sequenced

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

Source: Illumina

# High-Throughput Sequencing (HTS)



flow cell

computer    readout

orange = G    →    AGTG

= Second Generation
= Next Generation
= Massively Parallel Sequencing
= High Throughput Sequencing (HTS)
= Sequencing by Synthesis (Illumina)

Cleave fluorescence, wash away

**SAFARI**

# High-Throughput Sequencing (HTS)

**The sequencer adds the molecule "T" to all bases near the flow cell surface and observes the chemical reaction via a CMOS sensor.**
If a reaction happens then the base is "A"

Glass flow cell surface

As a workaround, HTS technologies sequence random short DNA fragments (75-300 basepairs long) of copies of the original molecule.

# High-Throughput Sequencing

- Massively parallel sequencing technology
  - Illumina, Roche 454, Ion Torrent, SOLID…

- Small DNA fragments are first amplified and then sequenced in parallel, leading to
  - High throughput
  - High speed
  - Low cost
  - Short reads

- Sequencing is done by either reading optical signals as each base is added, or by detecting hydrogen ions instead of light, leading to:
  - Low error rates (relatively)
  - Reads lack information about their order and which part of genome they are originated from

**Genome Analysis**

**1** **Sequencing**

Billions of Short Reads

**2** **Read Mapping**

Short Read

Read Alignment

Reference Genome

**3** **Variant Calling**

reference: TTTATCGCTTCCATGACGCAG
read1:        ATCGC**A**TCC
read2:       TATCGC**A**TC
read3:           C**A**TCCATGA
read4:         CGCTTCCAT
read5:              CCATGACGC
read6:             TTCCATGAC

**4** **Scientific Discovery**

PRESCRIPTION

# Multiple sequence alignment

```
PHDHtm                    ------------------------------------MMMMMMMMMMMMMMMMMMMM------
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 -MPSSKKKSTVPLASMAGLIRYYEE-ENEKIKISPKLLIIISIIMVAGVIVASILIPPP--  (58)
AE006662  S solf   11 -MPSSKKKSTVPVMSMAGLIRYYEE-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 -------------TNSNNSILKIYSD-EATGLRVDPLVVLFLAVGFIFSVVALEVISKVAGK  (82)
```

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

# 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 (bacteria)?



uncleaned de Bruijn graph

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

# Problem

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

**SAFARI**

**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%)

# The Read Mapping Bottleneck



300 Million bases/minute

Read Sequencing**

2 Million bases/minute

Read Mapping*

150x slower

* BWA-MEM
** HiSeqX10, MinION

# Read Mapping Execution Time Breakdown



SAM printing
3%

candidate alignment
locations (CAL)
4%

**Read Verification
93%**

# Read Mapping

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

Reference genome

Reads

DNA, logically physically

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

# Challenges in Read Mapping

- **Need to find many mappings of each read**
  - A short read may map to many locations, especially with High-Throughput DNA Sequencing technologies
  - 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 $\rightarrow$ 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?

# Read Alignment/Verification

- **Edit distance** 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 |

# Why Is Read Alignment Slow?

- **Quadratic-time** dynamic-programming algorithm(s)

- **Data dependencies** limit the computation parallelism

- **Entire matrix** computed even though strings may be dissimilar.



Read Alignment

# Example: Dynamic Programming Table

NETHERLANDS x SWITZERLAND

immediate left,
upper left,
upper entries of its own

|   |    | N | E | T | H | E | R | L | A | N | D  | S  |
|---|----|---|---|---|---|---|---|---|---|---|----|----|
|   |    |   | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| S |    | 1 |   |   |   |   |   |   |   |   |    |    |
| W | 2  |   |   |   |   |   |   |   |   |   |    |    |
| I | 3  |   |   |   |   |   |   |   |   |   |    |    |
| T | 4  |   |   |   |   |   |   |   |   |   |    |    |
| Z | 5  |   |   |   |   |   |   |   |   |   |    |    |
| E | 6  |   |   |   |   |   |   |   |   |   |    |    |
| R | 7  |   |   |   |   |   |   |   |   |   |    |    |
| L | 8  |   |   |   |   |   |   |   |   |   |    |    |
| A | 9  |   |   |   |   |   |   |   |   |   |    |    |
| N | 10 |   |   |   |   |   |   |   |   |   |    |    |
| D | 11 |   |   |   |   |   |   |   |   |   |    |    |

# Example: Dynamic Programming Table

NETHERLANDS x SWITZERLAND

|   |    | N | E | T | H | E | R | L | A | N | D  | S  |
|---|----|---|---|---|---|---|---|---|---|---|----|----|
|   | 0  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| S | 1  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 10 |
| W | 2  | 2 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| I | 3  | 3 | 3 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| T | 4  | 4 | 4 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| Z | 5  | 5 | 5 | 4 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| E | 6  | 6 | 5 | 5 | 5 | 4 | 5 | 6 | 7 | 8 | 9  | 10 |
| R | 7  | 7 | 6 | 6 | 6 | 5 | 4 | 5 | 6 | 7 | 8  | 9  |
| L | 8  | 8 | 7 | 7 | 7 | 6 | 5 | 4 | 5 | 6 | 7  | 8  |
| A | 9  | 9 | 8 | 8 | 8 | 7 | 6 | 5 | 4 | 5 | 6  | 7  |
| N | 10 | 9 | 9 | 9 | 9 | 8 | 7 | 6 | 5 | 4 | 5  | 6  |
| D | 11 | 10 | 10 | 10 | 10 | 9 | 8 | 7 | 6 | 5 | 4  | 5  |

- Matrix-filling is O(mn) time and space.
- Backtrace is O(m + n) time.

# Example: Dynamic Programming

- **Quadratic-time** dynamic-programming algorithm

  **WHY?!**

  NETHERLANDS x SWITZERLAND

  NETHERLANDS x S
- NETHERLANDS x SW
  NETHERLANDS x SWI
  NETERLANDS x SWIT
  NETHERLANDS x SWITZ
  NETHERLANDS x SWITZE
  NETHERLANDS x SWITZER
  NETHERLANDS x SWITZERL
- NETHERLANDS x SWITZERLA
  NETHERLANDS x SWITZERLAN
  NETHERLANDS x SWITZERLAND

|   | N | E | T | H | E | R | L | A | N | D | S |
|---|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| S | 1 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 10 |
| W | 2 | 2 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| I | 3 | 3 | 3 | 4 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| T | 4 | 4 | 4 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| Z | 5 | 5 | 5 | 4 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| E | 6 | 6 | 5 | 5 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| R | 7 | 7 | 6 | 6 | 6 | 5 | 4 | 5 | 6 | 7 | 8 | 9 |
| L | 8 | 8 | 7 | 7 | 7 | 6 | 5 | 4 | 5 | 6 | 7 | 8 |
| A | 9 | 9 | 8 | 8 | 8 | 7 | 6 | 5 | 4 | 5 | 6 | 7 |
| N | 10 | 9 | 9 | 9 | 9 | 8 | 7 | 6 | 5 | 4 | 5 | 6 |
| D | 11 | 10 | 10 | 10 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 5 |

etc etc etc

# Agenda

- **The Problem: DNA Read Mapping**
  - State-of-the-art Read Mapper Design

- **Algorithmic Acceleration**
  - Exploiting Structure of the Genome
  - Exploiting SIMD Instructions

- **Hardware Acceleration**
  - Specialized Architectures
  - Processing in Memory

- **Future Opportunities: New Sequencing Technologies**

**SAFARI**

# Read Mapping Algorithms: Two Styles

- Hash based seed-and-extend (hash table, suffix array, suffix tree)
  - Index the "k-mers" in the genome into a hash table (pre-processing)
  - When searching a read, find the location of a k-mer in the read; then extend through alignment
  - More sensitive (can find all mapping locations), but slow
  - Requires large memory; this can be reduced with cost to run time

- Burrows-Wheeler Transform & Ferragina-Manzini Index based aligners
  - BWT is a compression method used to compress the genome index
  - Perfect matches can be found very quickly, memory lookup costs increase for imperfect matches
  - Reduced sensitivity

# Hash Table Based Read Mappers

- Key Idea
  - Preprocess the reference into a *Hash Table*
  - Use *Hash Table* to map reads

# Hash Table-Based Mappers [Alkan+ Nature Gen'09]



k-mer or 12-mer (string of length k)

Location list—where the k-mer occurs in reference gnome

| AAAAAAAAAAAA | → | 12 | 324 | 577 | 940 |

| AAAAAAAAAAAC | → | 13 | 421 | 412 | 765 | 889 |

| AAAAAAAAAAAT | → | NULL |

| ...... |

| CCCCCCCCCCCC | → | 24 | 459 | 744 | 988 | 989 |

| ...... |

| ...... |

| ...... |

| TTTTTTTTTTTT | → | 36 | 535 | 123 |

Reference genome

Once for a reference

# Hash Table Based Read Mappers

- **Key Idea**
  - Preprocess the reference into a *Hash Table*
  - Use *Hash Table* to map reads

AAAAAAAAAAAACCCCCCCCCCCCTTTTTTTTTTTT ← read

CCCCCCCCCCCC ← k-mers
AAAAAAAAAAAA

**Hash Table (HT)**

| 3124 |

Reference Genome

AAAAAAAAAAAA → | 12 | 324 | 557 | 940 |

CCCCCCCCCCCC → | 24 | 459 | 744 | 988 | 989 |

TTTTTTTTTTTT → | 36 | 535 | 823 |

Valid mapping

...****************************

AAAAAAAAAAAACCCCCCCCCCCC

read

**Verification/Local Alignment**

36

# Advantages of Hash Table Based Mappers

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

nature genetics

http://mrfast.sourceforge.net/

# 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.

# Problem and Goal

- **Poor performance of existing read mappers: Very slow**
  - Verification/alignment takes too long to execute
  - Verification requires a memory access for reference genome + many base-pair-wise comparisons between the reference and the read (edit distance computation)



- **Goal: Speed up the mapper by reducing the cost of verification**

# Overarching Key Idea

**Filter fast** before you align

Minimize costly
edit distance computations

# Agenda

- The Problem: DNA Read Mapping
  - State-of-the-art Read Mapper Design

- <span style="color:blue">Algorithmic Acceleration</span>
  - <span style="color:blue">Exploiting Structure of the Genome</span>
  - Exploiting SIMD Instructions

- Hardware Acceleration
  - Specialized Architectures
  - Processing in Memory

- Future Opportunities: New Sequencing Technologies

**SAFARI**

# Reducing the Cost of Verification

- We observe that <span style="color:red">most verification (edit distance computation) calculations are unnecessary</span>
  - <span style="color:red">1 out of 1000 potential locations passes the verification process</span>

- We observe that we can get rid of unnecessary verification calculations by
  - *Detecting and rejecting early* invalid mappings (filtering)
  - *Reducing* the *number* of potential mappings to examine

# Key Observations [Xin+, BMC Genomics 2013]

- Observation 1
  - Adjacent k-mers in the read should also be adjacent in the reference genome
  - Read mapper can quickly reject mappings that do **not** satisfy this property

- Observation 2
  - Some k-mers are cheaper to verify than others because they have shorter location lists (they occur less frequently in the reference genome)
    - Mapper needs to examine only $e+1$ k-mers' locations to tolerate $e$ errors
  - Read mapper can choose the cheapest $e+1$ k-mers and verify their locations

# FastHASH Mechanisms [Xin+, BMC Genomics 2013]

- **Adjacency Filtering (AF)**: Rejects obviously invalid mapping locations at early stage to avoid unnecessary verifications

- **Cheap K-mer Selection (CKS):** Reduces the absolute number of potential mapping locations to verify

# Adjacency Filtering (AF)

- **Goal:** detect and filter out invalid mappings at early stage
- **Key Insight:** For a valid mapping, adjacent k-mers in the read are also adjacent in the reference genome

AAAAAAAAAAAACCCCCCCCCCCCTTTTTTTTTTTT ← read

Valid mapping          Invalid mapping          Reference genome

- **Key Idea:** search for adjacent locations in the k-mers' location lists
  - If more than $e$ k-mers fail → there must be more than e errors → invalid mapping

# Adjacency Filtering (AF)

# FastHASH Mechanisms [Xin+, BMC Genomics 2013]

- **Adjacency Filtering (AF)**: Rejects obviously invalid mapping locations at early stage to avoid unnecessary verifications

- **Cheap K-mer Selection (CKS):** Reduces the absolute number of potential mapping locations to verify
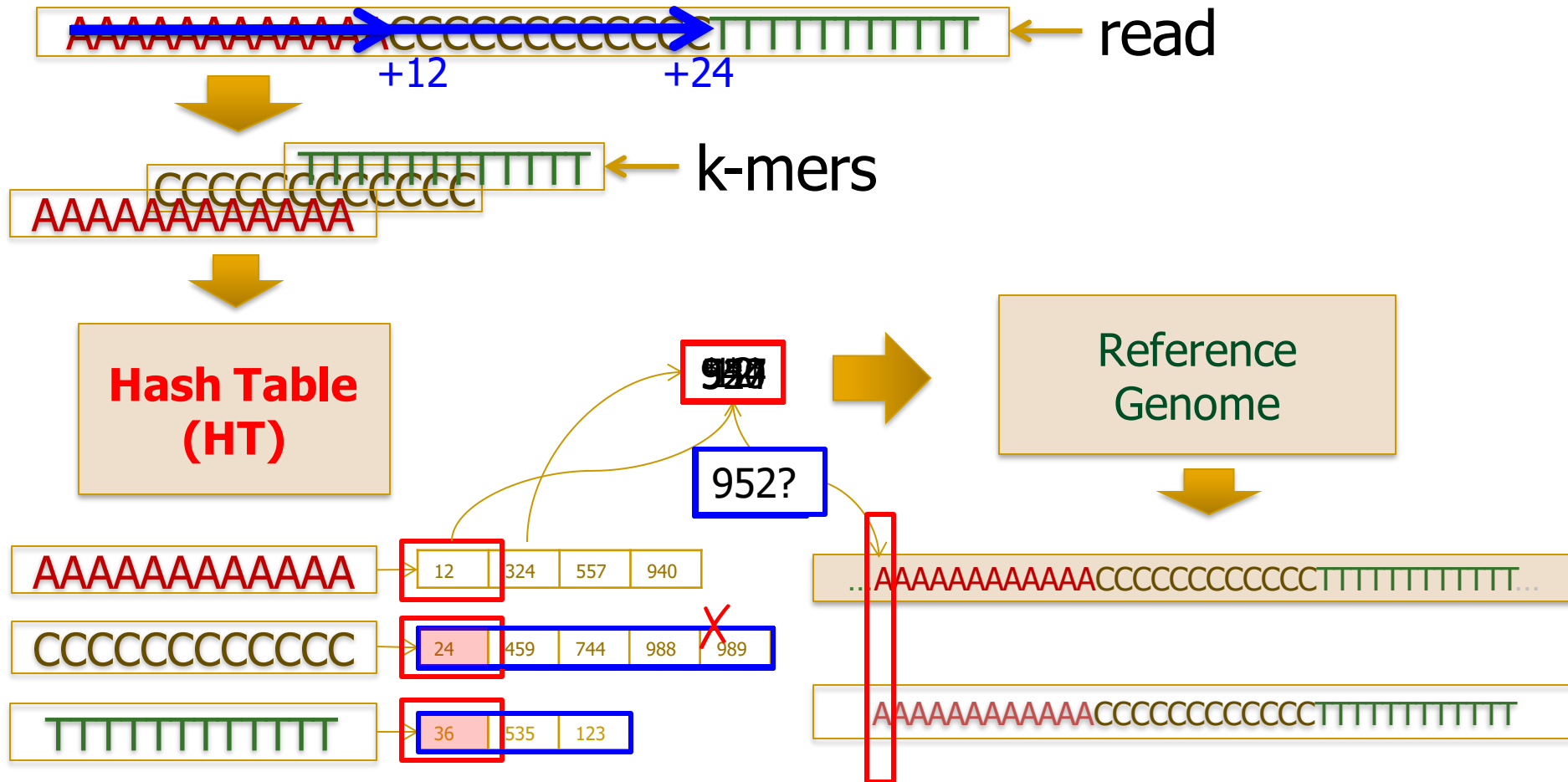
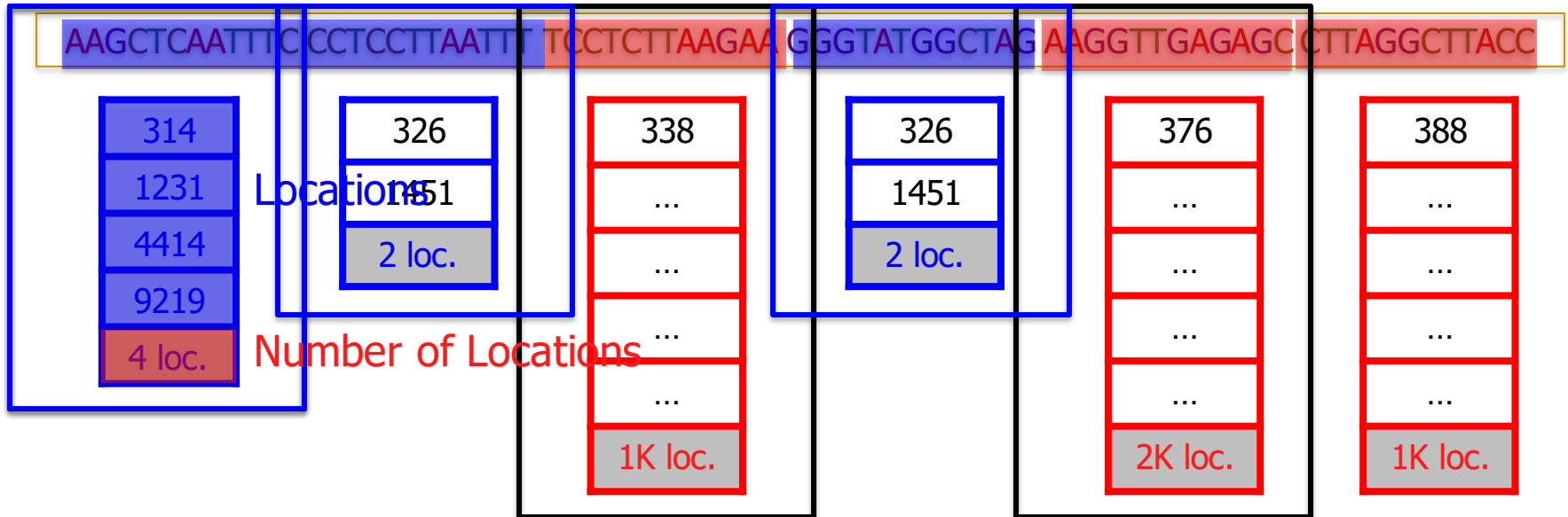# Cheap K-mer Selection (CKS)

- **Goal:** Reduce the number of potential mappings to examine

- **Key insight:**
    - K-mers have different cost to examine: Some k-mers are *cheaper* as they have fewer locations than others (occur less frequently in reference genome)

- **Key idea:**
    - Sort the k-mers based on their number of locations
    - Select the k-mers with the fewest number locations to verify

# Cheap K-mer Selection

- *e*=2 (examine 3 k-mers)        read

AAGCTCAATTTCCCTCCTTAATTTTCCTCTTAAGAAGGGTATGGCTAGAAGGTTGAGAGCCTTAGGCTTACC

| 314 |
| 1231 |
| 4414 |
| 9219 |
| 4 loc. |

Locations

Number of Locations

| 326 |
| 1451 |
| 2 loc. |

| 338 |
| ... |
| ... |
| ... |
| ... |
| ... |
| 1K loc. |

| 326 |
| 1451 |
| 2 loc. |

| 376 |
| ... |
| ... |
| ... |
| ... |
| 2K loc. |

| 388 |
| ... |
| ... |
| ... |
| ... |
| 1K loc. |

**Expensive 3 k-mers**

Previous work needs to verify:

3004 locations

FastHASH verifies only:

8 locations

# Methodology

- Implemented FastHASH on top of state-of-the-art mapper: mrFAST
  - New version mrFAST-2.5.0.0 over mrFAST-2.1.0.6

- Tested with real read sets generated from Illumina platform
  - 1M reads of a human (160 base pairs)
  - 500K reads of a chimpanzee (101 base pairs)
  - 500K reads of a orangutan (70 base pairs)

- Tested with simulated reads generated from reference genome
  - 1M simulated reads of human (180 base pairs)

- Evaluation system
  - Intel Core i7 Sandy Bridge machine
  - 16 GB of main memory

# FastHASH Speedup: Entire Read Mapper



With FastHASH, new mrFAST obtains up to 19x speedup over previous version, without losing valid mappings

# Analysis

- Reduction of potential mappings with FastHASH



FastHASH filters out over 99% of the potential mappings without sacrificing any valid mappings

# FastHASH Conclusion

- Problem: Existing read mappers perform poorly in mapping millions of short reads to the reference genome, in the presence of errors

- Observation: Most of the verification calculations are unnecessary → filter them out

- Key Idea: Exploit the structure of the genome to
  - Reject invalid mappings early (Adjacency Filtering)
  - Reduce the number of possible mappings to examine (Cheap K-mer Selection)

- Key Result: FastHASH obtains up to 19x speedup over the state-of-the-art mapper without losing valid mappings

# More on FastHASH

- Download 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*]

Xin+, **"Accelerating Read Mapping with FastHASH"**, BMC Genomics 2018.   53

# Agenda

- The Problem: DNA Read Mapping
  - State-of-the-art Read Mapper Design

- Algorithmic Acceleration
  - Exploiting Structure of the Genome
  - Exploiting SIMD Instructions

- Hardware Acceleration
  - Specialized Architectures
  - Processing in Memory

- Future Opportunities: New Sequencing Technologies

# An Example: Shifted Hamming Distance

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.**

# Shifted Hamming Distance

- Key observation:
  - If two strings differ by $E$ edits, then every bp match can be aligned in at most $2E$ shifts (of one of the strings).
    - Insight: Shifting a string by one "corrects" for one "error"

- Key idea:
  - Compute "Shifted Hamming Distance": AND of 2E Hamming Distances of two strings, to filter out invalid mappings
    - Uses bit-parallel operations that nicely map to SIMD instructions

- Key result:
  - SHD is 3x faster than SeqAn (the best implementation of Gene Myers' bit-vector algorithm), with only a 7% false positive rate
  - The fastest CPU-based filtering (pre-alignment) mechanism

# Hamming Distance ($\Sigma \oplus$)

3 matches     5 mismatches

*Edit = 1 Deletion*



To cancel the effect of a deletion, we need to shift in the *right* direction

# Insight: Shifting a String Helps Similarity Search

3 matches        5 mismatches



To cancel the effect of the deletion, we need to shift in the right direction

# Insight: Shifting a String Helps Similarity Search

7 matches          1 mismatches

# Shifted Hamming Distance

7 matches    1 mismatches

*Edit = 1 Deletion*

I S T A N B U L

XOR →

0 0 0 1 1 1 1

← XOR

AND

1 1 1 0 0 0 0

# Highly Parallel Matrix Computation

Reference

Query

2 Deletion Hamming masks

**We need to compute 2E+1 vectors, E=edit distance threshold**

$$dp[i][j]= 0 \text{ if } X[i]=Y[j]$$
$$1 \text{ if } X[i] \neq Y[j]$$

No data dependencies!

2 Insertion Hamming masks

# Key Idea of SHD Filtering

| Generate 2E+1 masks | Amend random zeros: 101 → 111 & 1001 → 1111 | AND all masks, ACCEPT iff number of '1' ≤ Threshold |
|---|---|---|

```
         Query :GAGAGAGATATTTAGTGTTGCAGCACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGGAACATTGTTGGGCCGGA
     Reference :GAGAGAGATAGTTAGTGTTGCAGCCACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGAGACATTGTTGGGCCGG

  Hamming Mask :00000000001000000000000111111101111000111011010110111111111000100000111011010010101
1-Deletion Mask :111111111110011111011111000000000000000000000000000000000000000011000000000000000
2-Deletion Mask :000000001011011001111111111111101110001110110101101111111110001001001110110110010 10
3-Deletion Mask :111111111110111011001101110111011000100100111111111111100101100110101101110111 01111
1-Insertion Mask :111111111110111110111111011101100010010011111111111111100101100110000101011101110111110
2-Insertion Mask :000000100111110011111111100100011010101001101011111111111111101110011111100011110 1100
3-Insertion Mask :111111110111011001100011111111010110111110011001011101111111101101111010111001000

              --- Masks after amendment ---

  Hamming Mask :000000000010000000000001111111111110001111111101111111111110001000001111111111111111
1-Deletion Mask :11111111111111111111111100000000000000000000000000000000000000011000000000000000
2-Deletion Mask :000000011111111111111111111111111110001111111111111111111110001000111111111111110
3-Deletion Mask :11111111111111111111111111111110001111111111111111111111111111111111111111111111111
1-Insertion Mask :11111111111111111111111111110001111111111111111111111110001111111111111110
2-Insertion Mask :00000111111111111111111110001111111111111111111111111111111111111110001111111100
3-Insertion Mask :11111111111111110001111111111111111111111111111111111111111111111111111000

       AND Mask :0000000000100000000000010000000000000000000000000000000000000000000001000000000000000

Needleman-Wunsch
      Alignment :GAGAGAGATATTTAGTGTTGCAG-CACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGGAACATTGTTGGGCCGG
                 |||||||||| |||||||||||| ||||||||||||||||||||||||||||||||||||||||||||::|||||||||||||||
                 GAGAGAGATAGTTAGTGTTGCAGCCACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGAGACATTGTTGGGCCGG
```

# Alignment vs. Pre-alignment (Filtering)

Needleman-Wunsch

SHD

C T A T A A T A C G

C T A T A A T A C G



Independent vectors can be processed in parallel using hardware technologies

*DRAM Layers*

*Logic Layer*

# New Bottleneck: Filtering (Pre-Alignment)

Sequencing generates many reads, each of which potentially mapping to many locations

→

Filtering (Pre-alignment) eliminates the need to verify/align read to invalid mapping locations

→

Alignment/verification (costly edit distance computation) is performed **only** on reads that pass the filter

- New bottleneck in read mapping becomes the "filtering (pre-alignment)" step

# Agenda

- **The Problem: DNA Read Mapping**
  - ❑ State-of-the-art Read Mapper Design

- **Algorithmic Acceleration**
  - ❑ Exploiting Structure of the Genome
  - ❑ Exploiting SIMD Instructions

- **Hardware Acceleration**
  - ❑ Specialized Architectures
  - ❑ Processing in Memory

- **Future Opportunities: New Sequencing Technologies**

**SAFARI**

# Location Filtering

- **Alignment** is expensive
  - We need to align millions to billions of reads

- M[...]t
  f[...]

  [...]
  out mismatches quickly

Our goal is to accelerate read mapping by improving the filtering step

- Both methods are used by mappers today, but filtering has replaced alignment as the bottleneck [Xin+, BMC Genomics 2013]

# Ideal Filtering Algorithm

**Minimal False Accept Rate**

**Maximal True Reject Rate**

Filter out all incorrect mappings

**Zero False Reject Rate**

**Faster Than Mapper**

Do not filter out any correct mappings

# Alignment vs. Pre-alignment (Filtering)

## Needleman-Wunsch

C T A T A A T A C G

| | | C | T | A | T | A | A | T | A | C | G |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | | | | | | | | |
| A | 1 | 0 | 1 | 2 | | | | | | | |
| C | 2 | 1 | 0 | 1 | 2 | | | | | | |
| T | | 2 | 1 | 0 | 1 | 2 | | | | | |
| A | | | 2 | 1 | 2 | 1 | 2 | | | | |
| T | | | | 2 | 2 | 2 | 1 | 2 | | | |
| A | | | | | 3 | 2 | 2 | 2 | 2 | | |

## SHD

C T A T A A T A C G

| | | C | T | A | T | A | A | T | A | C | G |
|---|---|---|---|---|---|---|---|---|---|---|---|
| A | | 1 | 1 | 0 | | | | | | | |
| C | | 0 | 1 | 1 | 1 | | | | | | |
| T | | 1 | 0 | 1 | 0 | 1 | | | | | |
| A | | | 1 | 0 | 1 | 0 | 0 | | | | |
| T | | | | 1 | 0 | 1 | 1 | 0 | | | |
| A | | | | | 1 | 0 | 0 | 1 | 0 | | |

Independent vectors can be processed in parallel using hardware technologies

DRAM Layers

Logic Layer

# Our Solution: GateKeeper



Alignment Filter + = 1st 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 Walkthrough

Amend random zeros:
101 → 111 & 1001 → 1111

AND all masks,
ACCEPT iff number of '1' ≤ Threshold

```
       Query :GAGAGAGATATTTAGTGTTGCAGCACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGGAACATTGTTGGGCCGGA
   Reference :GAGAGAGATAGTTAGTGTTGCAGCCACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGGAGACATTGTTGGGCCGG

 Hamming Mask :00000000001000000000000111111101111000111011010110111111111000100000111101101 0010101
1-Deletion Mask :11111111111001111101111100000000000000000000000000000000000000000011000000000000000
2-Deletion Mask :00000000101101110011111111111111101111000111011010110111111111000100010011101101001010
3-Deletion Mask :11111111111011101100110110110110001001001111111111111100101100110010110111011101111
1-Insertion Mask :11111111111011111011111101101100010010011111111111111001011001100010101110110111110
2-Insertion Mask :00000010011111001111111110010001101010100110101011111111111111011100111111100011110 1100
3-Insertion Mask :11111111011011001100111111111010110111110011001011101111111101110111101011100 1000

                         --- Masks after amendment ---

 Hamming Mask :00000000001000000000000111111111110001111111101111111111110001000001111111111111111
1-Deletion Mask :11111111111111111111111100000000000000000000000000000000000000000011000000000000000
2-Deletion Mask :00000000111111111111111111111111111111000111111111111111111111000100011111111111110
3-Deletion Mask :11111111111111111111111111111110001111111111111111111111111111111111111111111111111
1-Insertion Mask :11111111111111111111111111111111000111111111111111111111111111000111111111111111110
2-Insertion Mask :00000111111111111111111110001111111111111111111111111111111111111111000111111100
3-Insertion Mask :11111111111111111000111111111111111111111111111111111111111111111111111111111111000

     AND Mask :00000000001000000000000100000000000000000000000000000000000000000001000000000000000

Needleman-Wunsch   GAGAGAGATATTTAGTGTTGCAG-CACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGGAACATTGTTGGGCCGG
   Alignment :     ||||||||||  |||||||||||    ||||||||||||||||||||||||||||||||||||||||||||||::  |||||||||||
                   GAGAGAGATAGTTAGTGTTGCAGCCACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGGAGACATTGTTGGGCCGG
```

70

# GateKeeper Walkthrough (cont'd)

| Generate 2E+1 masks | Amend random zeros: 101 → 111 & 1001 → 1111 | AND all masks, ACCEPT iff number of '1' ≤ Threshold |
|---|---|---|

- E right-shift registers (length=ReadLength)
- E left-shift registers (length=ReadLength)
- (2E+1) * (ReadLength) 2-XOR operations.

- (2E)*(ReadLength) 2-AND operations.
- (ReadLength/4) 5-input LUT.
- $log_2$ReadLength-bit counter.



Hamming mask

0 1 0 0 1 0 0 0 1 1 0 1 0 0 0 1 0 1 0 1 1 0 0 1 1 1 1 1 0 0 0 1 0 0 1 0

5-input LUT

Hamming mask after amending

0 1 1 1 1 0 0 0 1 1 1 1 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 1 1 1 1 0

- (2E+1)*(ReadLength) 5-input LUT.

# GateKeeper Accelerator Architecture

- **Maximum data throughput** =~13.3 billion bases/sec

- Can examine **8 (300 bp) or 16 (100 bp) mappings concurrently** at 250 MHz

- **Occupies 50%** (100 bp) to **91%** (300 bp) of the FPGA slice LUTs and registers

# GateKeeper vs. SHD

| GateKeeper | SHD |
|---|---|
| ■ FPGA (Xilinx VC709) | ■ Intel SIMD |
| ■ Multi-core (parallel) | ■ Single-core (sequential) |
| ■ Examines a single mapping @ 125 MHz | ■ Examines a single mapping @ ~2MHz |
| ■ Limited to PCIe Gen3(4x) transfer rate (128 bits @ 250MHz) | ■ Limited to a read length of 128 bp (SSE register size) |
| ■ Amending requires:<br>❑ (2E+1) 5-input LUT. | ■ Amending requires:<br>❑ 4(2E+1) bitwise OR.<br>❑ 4(2E+1) packed shuffle.<br>❑ 3(2E+1) shift. |

# GateKeeper: Speed & Accuracy Results

## 90x-130x faster filter

than SHD (Xin et al., 2015) and the Adjacency Filter (Xin et al., 2013)

## 4x lower false accept rate

than the Adjacency Filter (Xin et al., 2013)

## 10x speedup in read mapping

with the addition of GateKeeper to the mrFAST mapper (Alkan et al., 2009)

## Freely available online

github.com/BilkentCompGen/GateKeeper

# Conclusions

- FPGA-based pre-alignment greatly speeds up read mapping
  - 10x speedup of a state-of-the-art mapper (mrFAST)

- FPGA-based pre-alignment can be integrated with the sequencer
  - It can help to hide the complexity and details of the FPGA
  - Enables real-time filtering while sequencing

# More on GateKeeper

- Download and test for yourself
  https://github.com/BilkentCompGen/GateKeeper

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

*Sequence analysis*

# GateKeeper: A New Hardware Architecture for Accelerating Pre-Alignment in DNA Short Read Mapping

Mohammed Alser[1,*], Hasan Hassan[2], Hongyi Xin[3], Oğuz Ergin[2], Onur Mutlu[4,*], and Can Alkan[1,*]

# MAGNET (AACBB 2018, TIR 2017)

- <u>Key observation:</u> the use of **AND operation** to check if a zero (match) exists in a column introduces filtering inaccuracy.

- <u>Key Idea:</u> count the **consecutive zeros** in each mask and select the longest in a divide-and-conquer approach.

- **MAGNET** is **17x to 105x more accurate** than GateKeeper and SHD.

**SAFARI**

# MAGNET Walkthrough

Find the longest segment of consecutive zeros

Exclude the errors from the search space

Divide the problem into two subproblems and repeat

SAFARI

# MAGNET Accelerator

# Agenda

- **The Problem: DNA Read Mapping**
  - State-of-the-art Read Mapper Design

- **Algorithmic Acceleration**
  - Exploiting Structure of the Genome
  - Exploiting SIMD Instructions

- **Hardware Acceleration**
  - Specialized Architectures
  - Processing in Memory

- **Future Opportunities: New Sequencing Technologies**

**SAFARI**

# Read Mapping & Filtering

- Problem: Heavily bottlenecked by Data Movement

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

- Ditto for SHD [Xin+, Bioinformatics 2015]

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

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

**SAFARI**

# Hash Tables in Read Mapping

**Read Sequence (100 bp)**



A Matching... Match!

Mismatch. Missing! **False Negative**

**Hash Table**

37    140
894    1203
1564

**Reference Genome**

**Filter**

**SAFARI**

# Read Mapping & Filtering in Memory

We need to design

mapping & filtering algorithms

that fit processing-in-memory

**SAFARI**

# Our Proposal: GRIM-Filter

1. **Data Structures: Bins & Bitvectors**

2. Checking a Bin

3. Integrating GRIM-Filter into a Mapper

SAFARI

# GRIM-Filter: Bins

- We partition the genome into large sequences (**bins**).

*Bin x - 3*    *Bin x - 1*

... **GGAAATACGTTCAGTCAGTTGGAAATACGTTTTGGGCGTTACTTCTCAGTACGTACAGTACAGTAAAAATGACAGTAAGAC** ...

*Bin x - 2*    *Bin x*

- ❑ Represent each bin with a **bitvector** that holds the occurrence of all permutations of a small string (**token**) in the bin

- ❑ To account for matches that straddle bins, we employ overlapping bins
  - A read will now always completely fall within a single bin

**Bitvector**

| | |
|---|---|
| **AAAAA** | 1 |
| AAAAC | 0 |
| AAAAT | 1 |
| ... | ... |
| CCCCC | 1 |
| **CCCCT** | 0 |
| CCCCG | 0 |
| ... | ... |
| GGGGG | 1 |

**AAAAA exists** in bin x

**CCCCT doesn't exist** in bin x

# GRIM-Filter: Bitvectors

... **C G T G A** G T C ...

Bin x

|  | Bin x Bitvector |
|---|---|
| AAAAA | 0 |
| ... | ... |
| CGTGA | 1 |
| ... | ... |
| TGAGT | 1 |
| ... | ... |
| GAGTC | 0 |
| ... | ... |
| GTGAG | 1 |
| ... | ... |

**SAFARI**

# GRIM-Filter: Bitvectors

*Reference Genome*

bin$_1$
bin$_3$

AAAAACCCCTGCCTTGCATGTAGAAAACTTGACAGGAACTTTTTATCGCA ⬛⬛⬛

bin$_2$
bin$_4$

**tokens**

**b$_1$**

| AAAAA | 1 |
| AAAAC | 1 |
| AAAAG | 0 |
| AAAAT | 0 |
| . | . |
| CCCCT | 1 |
| . | . |
| . | . |
| . | . |
| . | . |
| GCATG | 1 |
| . | . |
| TTGCA | 1 |
| . | . |
| TTTTT | 0 |

**b$_2$**

| AAAAA | 0 |
| AAAAC | 1 |
| AAAAG | 0 |
| . | . |
| AGAAA | 1 |
| . | . |
| GAAAA | 1 |
| . | . |
| GACAG | 1 |
| . | . |
| GCATG | 1 |
| . | . |
| . | . |
| . | . |
| TTTTT | 0 |

• • •

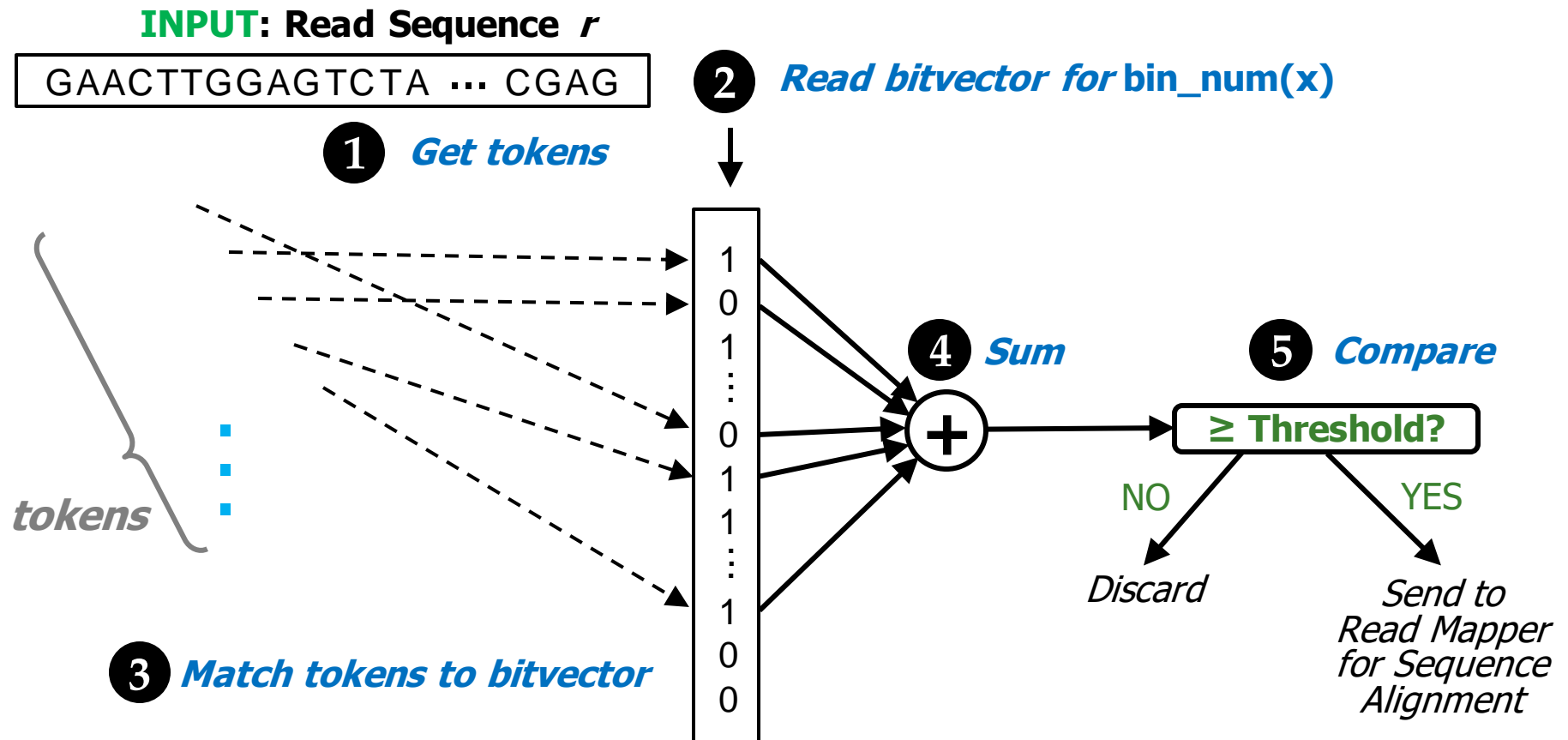Storing all bitvectors requires $4^n * t$ bits in memory, where t = number of bins.

For **bin size** ~200, and **n** = 5, **memory footprint** ~3.8 GB

# Our Proposal: GRIM-Filter

1. Data Structures: Bins & Bitvectors

2. **Checking a Bin**

3. Integrating GRIM-Filter into a Mapper

**SAFARI**

# GRIM-Filter: Checking a Bin

How GRIM-Filter determines whether to **discard** potential match locations in a given bin **prior** to alignment

**INPUT**: Read Sequence *r*

GAACTTGGAGTCTA ··· CGAG

**1** *Get tokens*

**2** *Read bitvector for* **bin_num(x)**

*tokens*

**3** *Match tokens to bitvector*

1
0
1
···
0
1
1
···
1
0
0

**4** *Sum*

(+)

**5** *Compare*

≥ Threshold?

NO        YES

*Discard*     *Send to Read Mapper for Sequence Alignment*
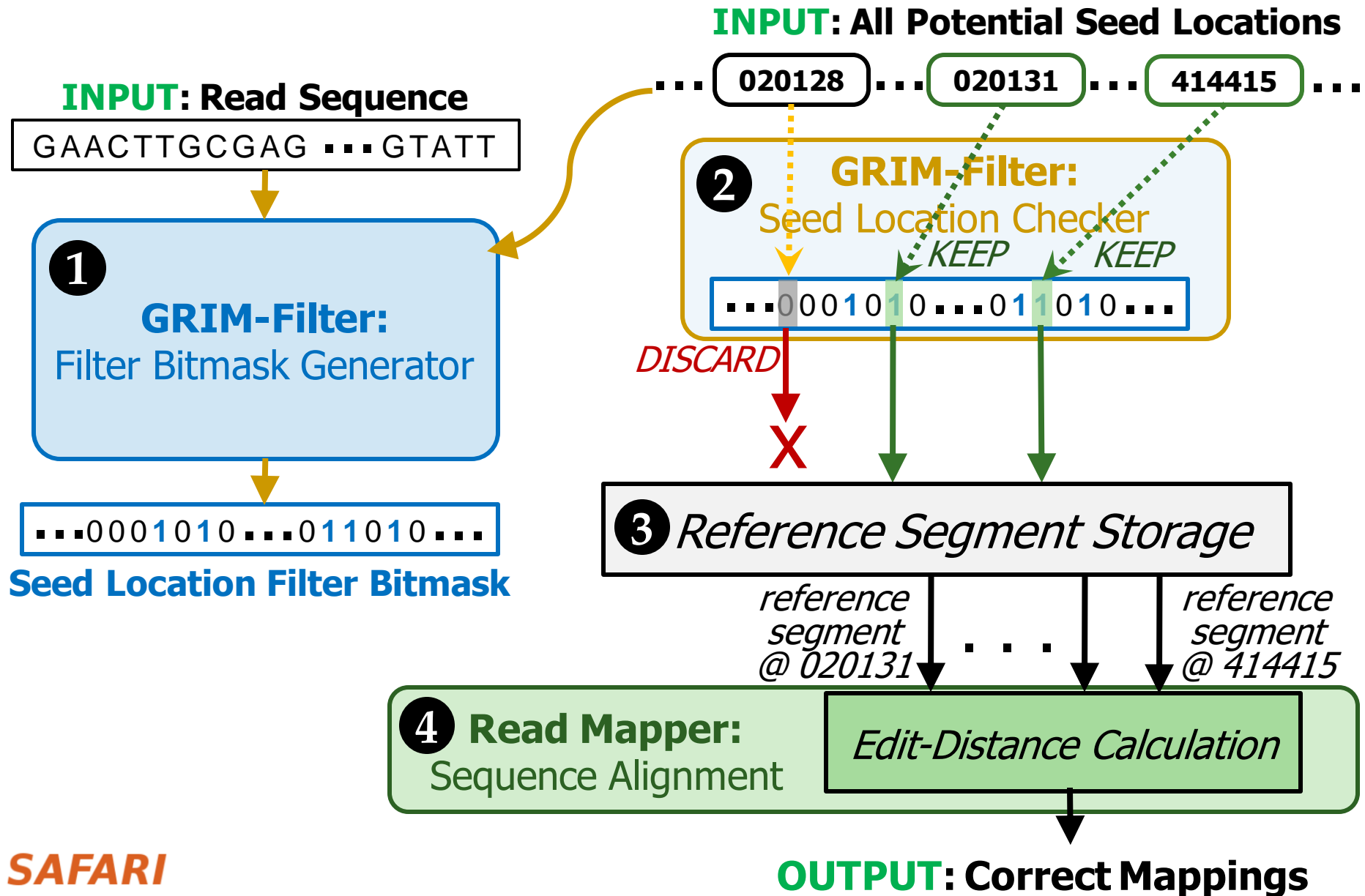
**SAFARI**

# Our Proposal: GRIM-Filter

1. Data Structures: Bins & Bitvectors

2. Checking a Bin

3. Integrating GRIM-Filter into a Mapper

# Our Proposal: GRIM-Filter

1.  Data Structures: Bins & Bitvectors

2.  Checking a Bin

3.  **Integrating GRIM-Filter into a Mapper**

# Integrating GRIM-Filter into a Read Mapper



**INPUT: All Potential Seed Locations**

⋯⋯ 020128 ⋯⋯ 020131 ⋯⋯ 414415 ⋯⋯

**INPUT: Read Sequence**

GAACTTGCGAG ⋯⋯ GTATT

**1 GRIM-Filter:** Filter Bitmask Generator

**2 GRIM-Filter:** Seed Location Checker

KEEP    KEEP

⋯0001010⋯011010⋯

DISCARD

✗

**Seed Location Filter Bitmask**

⋯0001010⋯011010⋯

**3 Reference Segment Storage**

reference segment @ 020131    reference segment @ 414415

**4 Read Mapper:** Sequence Alignment

Edit-Distance Calculation

**OUTPUT: Correct Mappings**

SAFARI

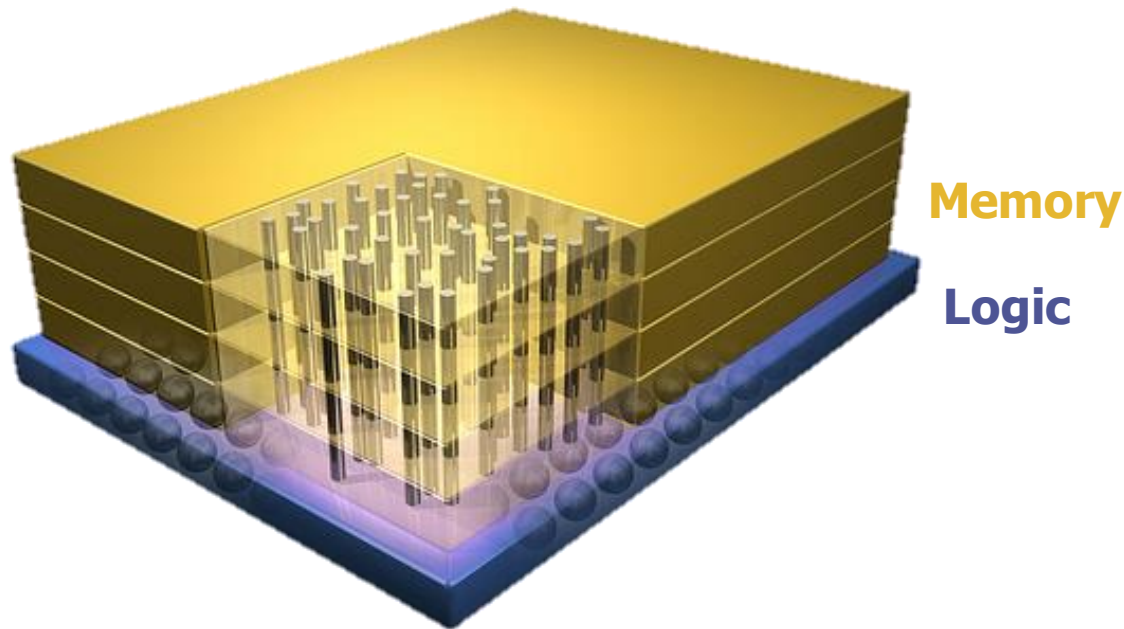# Key Properties of GRIM-Filter

1. **Simple Operations:**
   - ❑ To check a given bin, find the **sum** of all bits corresponding to each token in the read
   - ❑ **Compare** against threshold to determine whether to align

2. **Highly Parallel:** Each bin is operated on independently and there are many many bins

3. **Memory Bound:** Given the frequent accesses to the large bitvectors, we find that GRIM-Filter is memory bound

**These properties together make GRIM-Filter a good algorithm to be run in 3D-Stacked DRAM**

# Opportunity: 3D-Stacked Logic+Memory

**Memory**

**Logic**

Other "True 3D" technologies under development

# DRAM Landscape (circa 2015)

| Segment | DRAM Standards & Architectures |
|---|---|
| Commodity | DDR3 (2007) [14]; DDR4 (2012) [18] |
| Low-Power | LPDDR3 (2012) [17]; LPDDR4 (2014) [20] |
| Graphics | GDDR5 (2009) [15] |
| Performance | eDRAM [28], [32]; RLDRAM3 (2011) [29] |
| 3D-Stacked | WIO (2011) [16]; WIO2 (2014) [21]; MCDRAM (2015) [13]; HBM (2013) [19]; HMC1.0 (2013) [10]; HMC1.1 (2014) [11] |
| Academic | SBA/SSA (2010) [38]; Staged Reads (2012) [8]; RAIDR (2012) [27]; SALP (2012) [24]; TL-DRAM (2013) [26]; RowClone (2013) [37]; Half-DRAM (2014) [39]; Row-Buffer Decoupling (2014) [33]; SARP (2014) [6]; AL-DRAM (2015) [25] |

Table 1. Landscape of DRAM-based memory

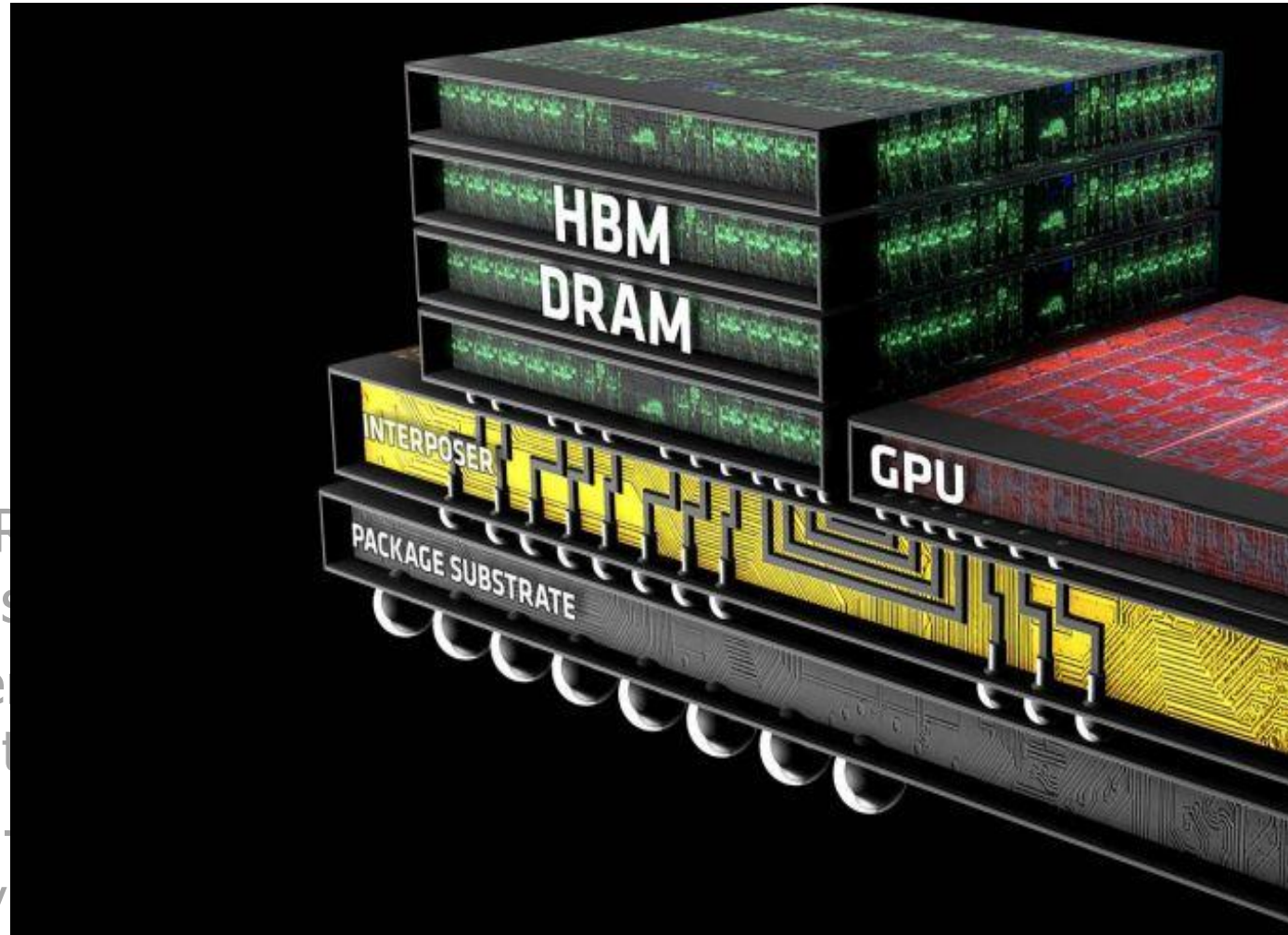Kim+, "Ramulator: A Flexible and Extensible DRAM Simulator", IEEE CAL 2015.

# 3D-Stacked Memory



*DRAM Layers*

*TSVs*

*Logic Layer*

- 3D-Stacked DRAM architecture has **extremely high bandwidth** as well as a stacked customizable logic layer
  - Logic Layer enables **Processing-in-Memory**, via high-bandwidth low-latency access to DRAM layers
  - Embed GRIM-Filter operations into **DRAM logic layer** and appropriately distribute bitvectors throughout memory
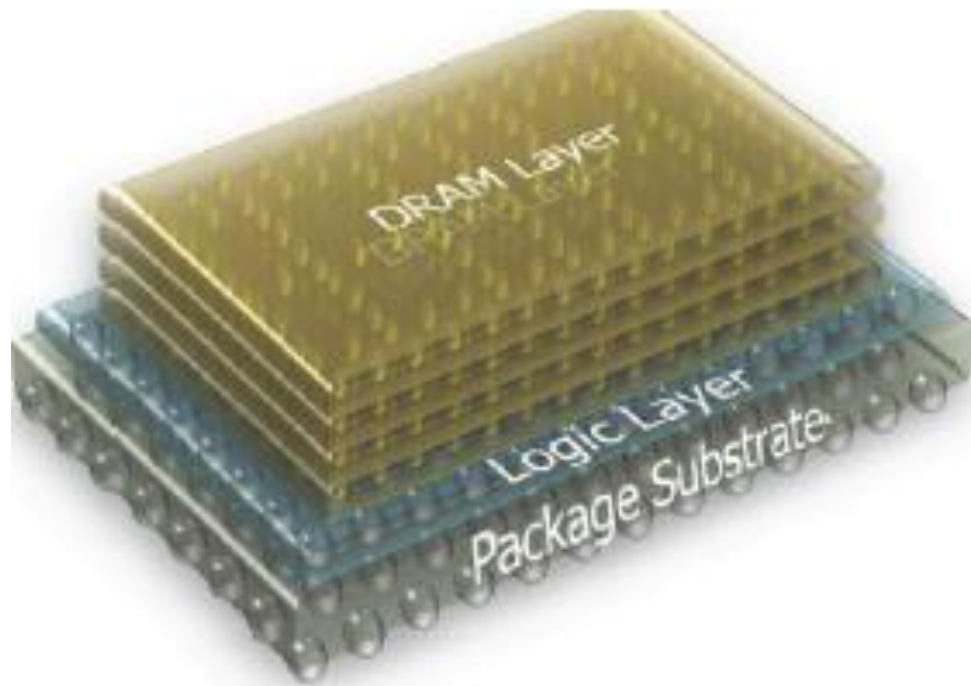
# 3D-Stacked Memory

- 3D-Stacked DR
  **bandwidth** as
  - Logic Layer e
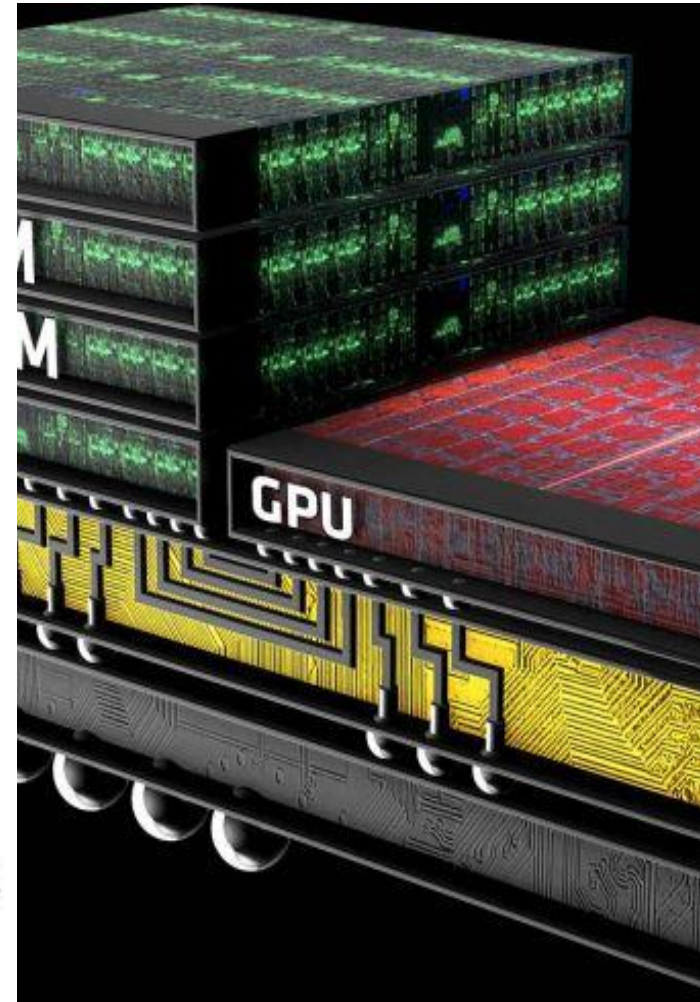    computation t
  - Embed GRIM-
    appropriately

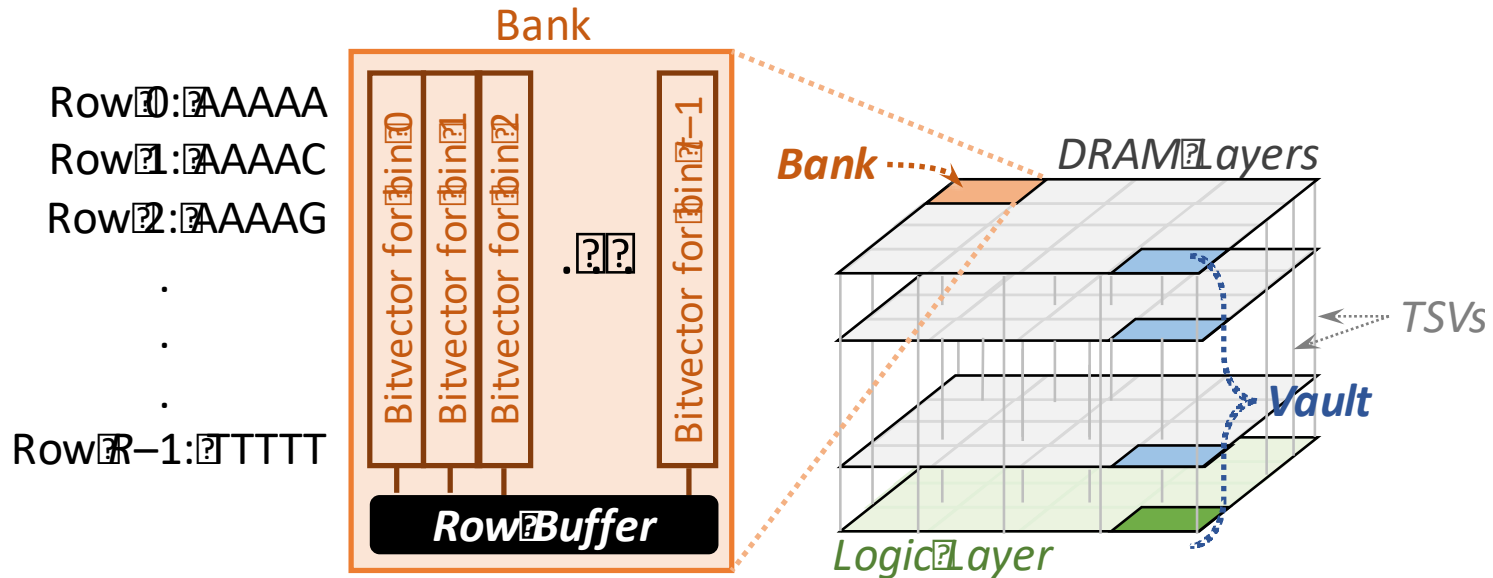**SAFARI**

# 3D-Stacked Memory



Micron's HMC

Micron has working demonstration components

http://images.anandtech.com/doci/9266/HBMCar_678x452.jpg

SAFARI

# GRIM-Filter in 3D-Stacked DRAM

Bank

Row 0: AAAAA
Row 1: AAAAC
Row 2: AAAAG
.
.
.
Row $R$–1: TTTTT

Bitvector for bin 0
Bitvector for bin 1
Bitvector for bin 2
. . .
Bitvector for bin $t$–1

**Row Buffer**

*Bank*
*DRAM Layers*
*TSVs*
*Vault*
*Logic Layer*

- Each DRAM layer is organized as an array of **banks**
  - A **bank** is an array of cells with a row buffer to transfer data

- The layout of bitvectors in a bank enables filtering many bins in parallel

**SAFARI**

# GRIM-Filter in 3D-Stacked DRAM



Per-Vault
Custom GRIM-Filter Logic

Seed Location Filter Bitmask

- Customized logic for accumulation and comparison per genome segment
  - Low area overhead, simple implementation
  - For HBM2, we use 4096 incrementer LUTs, 7-bit counters, and comparators in logic layer

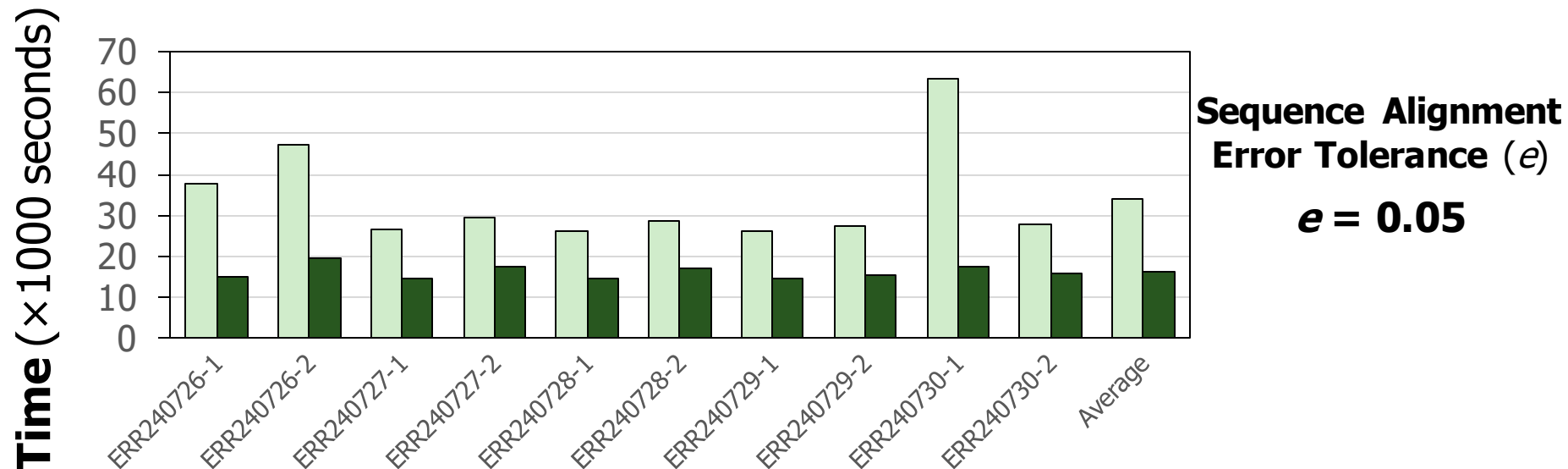**Details are in [Kim+, BMC Genomics 2018]**

*SAFARI*

# Methodology

- Performance simulated using an in-house 3D-Stacked DRAM simulator

- Evaluate 10 real read data sets (From the 1000 Genomes Project)
  - Each data set consists of 4 million reads of length 100

- Evaluate two key metrics
  - Performance
  - False negative rate
    - The fraction of locations that pass the filter but result in a mismatch

- Compare against a state-of-the-art filter, FastHASH **[Xin+, BMC Genomics 2013]** when using mrFAST, but **GRIM-Filter can be used with ANY read mapper**

**SAFARI**

# GRIM-Filter Performance

Benchmarks and their Execution Times



**1.8x-3.7x performance benefit across real data sets**
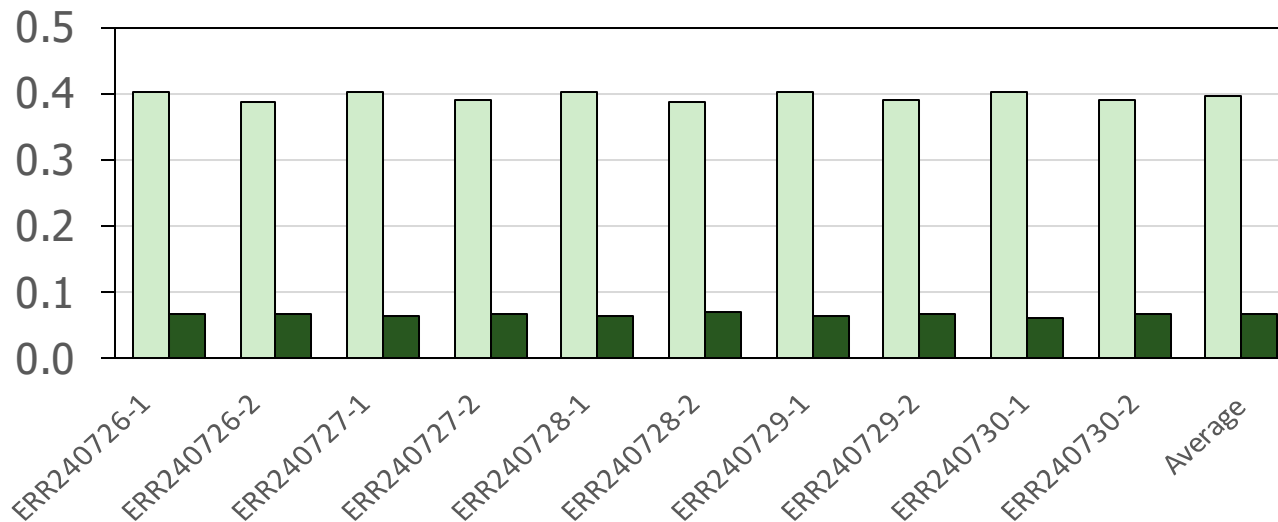
**2.1x average performance benefit**

**GRIM-Filter gets performance due to its hardware-software co-design**

# GRIM-Filter False Negative Rate

## Benchmarks and their False Negative Rates



**5.6x-6.4x False Negative reduction across real data sets**

**6.0x average reduction in False Negative Rate**

**GRIM-Filter utilizes more information available in the read to filter**

SAFARI

# More on GRIM-Filter

- 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*]

# Aside: In-Memory Graph Processing

- Large graphs are everywhere (circa 2015)

| | | | |
|---|---|---|---|
| 36 Million Wikipedia Pages | 1.4 Billion Facebook Users | 300 Million Twitter Users | 30 Billion Instagram Photos |

- Scalable large-scale graph processing is challenging

Bar chart:
- 32 Cores: ~1 speedup
- 128...: ~1.4 speedup, +42%

X-axis: Speedup (0 to 4)

# Key Bottlenecks in Graph Processing

```
for (v: graph.vertices) {
    for (w: v.successors) {
        w.next_rank += weight * v.rank;
    }
}
```

**1. Frequent random memory accesses**

w.rank

w.next_rank

w.edges

…

v

&w

w

**weight * v.rank**

**2. Little amount of computation**

**SAFARI**

# Tesseract System for Graph Processing

Interconnected set of 3D-stacked memory+logic chips with simple cores



Host Processor

Memory-Mapped
Accelerator Interface
(Noncacheable, Physically Addressed)

Memory

Logic

Crossbar Network

In-Order Core

DRAM Controller

LP

PF Buffer

MTP

Message Queue

NI

Ahn+, "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing" ISCA 2015.

# Tesseract System for Graph Processing

Host Processor

Memory-Mapped
Accelerator Interface
(Noncacheable, Physically Addressed)

**Memory**

**Logic**

In-Order Core

DRAM

Crossbar Network

...
...
...
...

Communications via
Remote Function Calls

Message Queue

NI

```
for (v: graph.vertices) {
    for (w: v.successors) {
        w.next_rank += weight * v.rank;
    }
}
```

# Communications In Tesseract (II)

```
for (v: graph.vertices) {
    for (w: v.successors) {
        w.next_rank += weight * v.rank;
    }
}
```

**SAFARI**

# Communications In Tesseract (III)

```
for (v: graph.vertices) {
    for (w: v.successors) {
        put(w.id, function() { w.next_rank += weight * v.rank; });
    }
}
barrier();
```

**Non-blocking Remote Function Call**

Can be **delayed**
until the nearest barrier



Vault #1

Vault #2

**SAFARI**

# Remote Function Call (Non-Blocking)

1. Send function address & args to the remote core
2. Store the incoming message to the message queue
3. Flush the message queue when it is full or a synchronization barrier is reached



put(w.id, function() { w.next_rank += value; })

# Tesseract System for Graph Processing

Host Processor

Memory-Mapped
Accelerator Interface
(Noncacheable, Physically Addressed)

**Memory**

**Logic**

Prefetching

LP

PF Buffer

MTP

DRAM Controller

Message Queue

NI

...

Crossbar Network

...

# Evaluated Systems



| DDR3-OoO | HMC-OoO | HMC-MC | **Tesseract** |
|:---:|:---:|:---:|:---:|
| 8 OoO 4GHz | 8 OoO 4GHz | 128 In-Order 2GHz | 32 Tesseract Cores |
| 102.4GB/s | 640GB/s | 640GB/s | **8TB/s** |

# Tesseract Graph Processing Performance

**>13X Performance Improvement**

On five graph processing algorithms

Ahn+, "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing" ISCA 2015.

# Tesseract Graph Processing Performance



**Memory Bandwidth Consumption**

Memory Bandwidth (TB/s)

- DDR3-OoO: 80GB/s
- HMC-OoO: 190GB/s
- HMC-MC: 243GB/s
- Tesseract: 1.3TB/s
- Tesseract-LP: 2.2TB/s
- Tesseract-LP-MTP: 2.9TB/s

*SAFARI*

# Effect of Bandwidth & Programming Model

**SAFARI**

# Tesseract Graph Processing System Energy



Legend: Memory Layers · Logic Layers · Cores

Categories: HMC-OoO, Tesseract with Prefetching

> 8X Energy Reduction

# More on Tesseract

- Junwhan Ahn, Sungpack Hong, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi,
**"A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing"**
*Proceedings of the 42nd International Symposium on Computer Architecture* (**ISCA**), Portland, OR, June 2015.
[Slides (pdf)] [Lightning Session Slides (pdf)]

## A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing

Junwhan Ahn    Sungpack Hong[§]    Sungjoo Yoo    Onur Mutlu[†]    Kiyoung Choi

junwhan@snu.ac.kr, sungpack.hong@oracle.com, sungjoo.yoo@gmail.com, onur@cmu.edu, kchoi@snu.ac.kr

Seoul National University    [§]Oracle Labs    [†]Carnegie Mellon University

# Agenda

- The Problem: DNA Read Mapping
  - State-of-the-art Read Mapper Design

- Algorithmic Acceleration
  - Exploiting Structure of the Genome
  - Exploiting SIMD Instructions

- Hardware Acceleration
  - Specialized Architectures
  - Processing in Memory

- Future Opportunities: New Sequencing Technologies

**SAFARI**

# Recall: High-Throughput Sequencing

- **Massively parallel sequencing technology**
  - Illumina, Roche 454, Ion Torrent, SOLID…

- **Small DNA fragments are first amplified and then sequenced in parallel, leading to**
  - High throughput
  - High speed
  - Low cost
  - Short reads
    - Amplification step limits the read length since too short or too long fragments are not amplified well.

- **Sequencing is done by either reading optical signals as each base is added, or by detecting hydrogen ions instead of light, leading to:**
  - Low error rates (relatively)
  - Reads lack information about their order and which part of genome they are originated from

# Nanopore Sequencing Technology

- **Nanopore sequencing** is an emerging and a promising single-molecule DNA sequencing technology

- First nanopore sequencing device, **MinION**, made commercially available by **Oxford Nanopore Technologies** (ONT) in **May 2014.**
  - Inexpensive
  - Long read length (> 882K bp)
  - Portable: Pocket-sized
  - Produces data in real-time

# Nanopore Sequencing Technology



... an emerging and a promising ... ncing technology

read length → Longer read length

- First nanopore sequencing device, **MinION**, made commercially available by **Oxford Nanopore Technologies** (ONT) in **May 2014.**
  - ❑ Inexpensive
  - ❑ Long read length (> 882K bp)
  - ❑ Portable: Pocket-sized
  - ❑ Produces data in real-time

# Nanopore Sequencing



- **Nanopore** is a nano-scale hole
- In nanopore sequencers, an **ionic current** passes through the nanopores
- When the DNA strand passes through the nanopore, the sequencer measures the the **change in current**
- This change is used to identify the bases in the strand with the help of **different electrochemical structures** of the different bases

# Advantages of Nanopore Sequencing

Nanopores:

- Do *not* require any labeling of the DNA or nucleotide for detection during sequencing

- Rely on the electronic or chemical structure of the different nucleotides for identification

- Allow sequencing very long reads, and

- Provide portability, low cost, and high throughput.

# Challenges of Nanopore Sequencing

- One major drawback: high error rates

- Nanopore sequence analysis tools have a critical role to:
  - overcome high error rates
  - take better advantage of the technology

- Faster tools are critically needed to:
  - Take better advantage of the real-time data production capability of MinION
  - Enable fast, real-time data analysis

# 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.
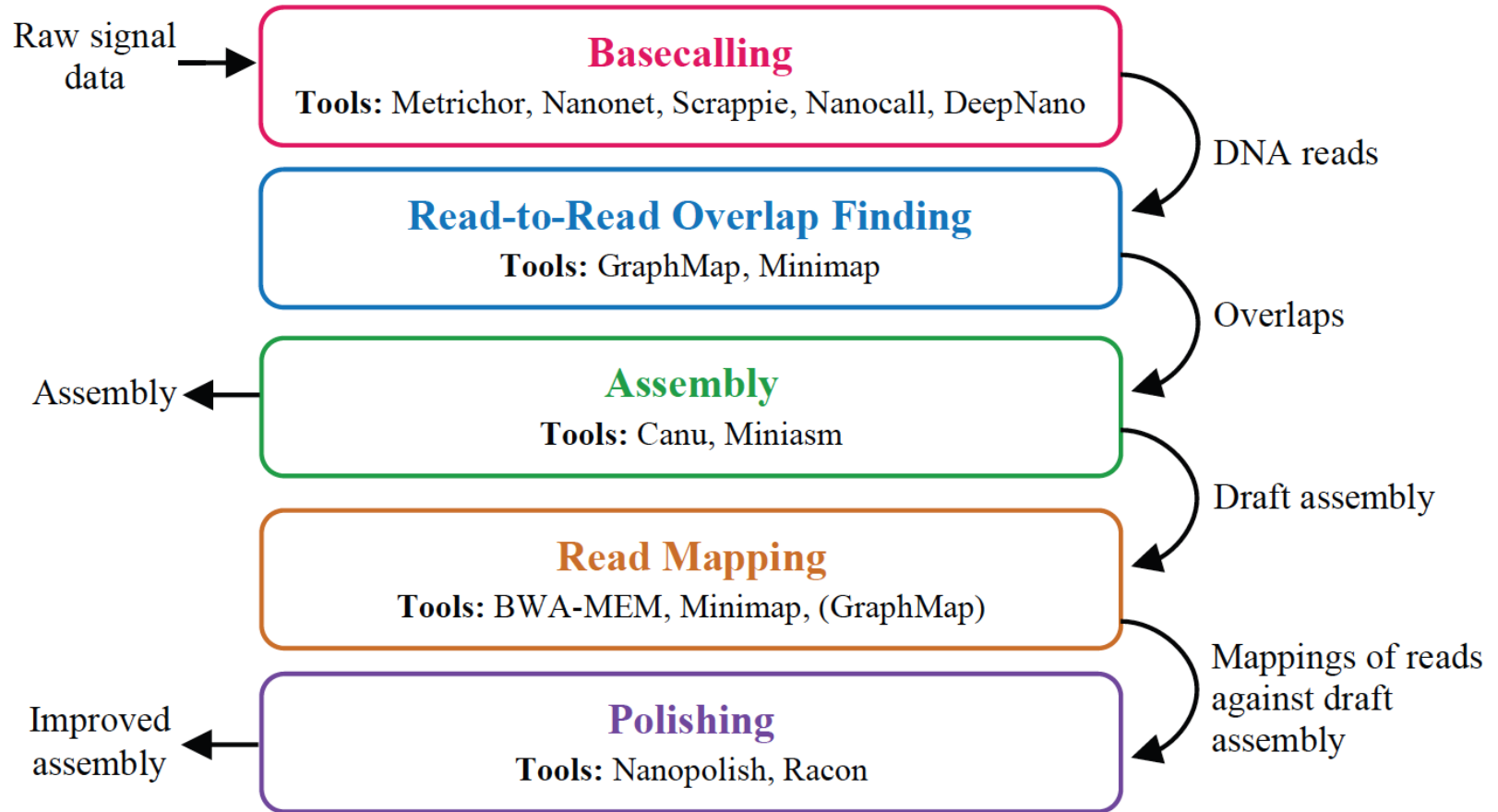
# Nanopore Genome Assembly Tools (I)

**Table 12. Accuracy analysis results for the full pipeline with a focus on the last two steps.**

| | | | | | | | | | Number of Bases | Number of Contigs | Identity (%) | Coverage (%) | Number of Mismatches | Number of Indels |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Metrichor | + | — | + | Canu | + | BWA-MEM | + | Nanopolish | 4,683,072 | 1 | 99.48 | 99.93 | 8,198 | 15,581 |
| 2 | Metrichor | + | Minimap | + | Miniasm+ | BWA-MEM | + | Nanopolish | 4,540,352 | 1 | 92.33 | 96.31 | 162,884 | 182,965 |
| 3 | Metrichor | + | GraphMap | + | Miniasm+ | BWA-MEM | + | Nanopolish | 4,637,916 | 2 | 92.38 | 95.80 | 159,206 | 180,603 |
| 4 | Metrichor | + | — | + | Canu | + | BWA-MEM | + | Racon | 4,650,502 | 1 | 98.46 | 100.00 | 18,036 | 51,842 |
| 5 | Metrichor | + | — | + | Canu | + | Minimap | + | Racon | 4,648,710 | 1 | 98.45 | 100.00 | 17,906 | 52,168 |
| 6 | Metrichor | + | Minimap | + | Miniasm+ | BWA-MEM | + | Racon | 4,598,267 | 1 | 97.70 | 99.91 | 24,014 | 82,906 |
| 7 | Metrichor | + | Minimap | + | Miniasm+ | Minimap | + | Racon | 4,600,109 | 1 | 97.78 | 100.00 | 23,339 | 79,721 |
| 8 | Nanonet | + | — | + | Canu | + | BWA-MEM | + | Racon | 4,622,285 | 1 | 98.48 | 100.00 | 16,872 | 52,509 |
| 9 | Nanonet | + | — | + | Canu | + | Minimap | + | Racon | 4,620,597 | 1 | 98.49 | 100.00 | 16,874 | 52,232 |
| 10 | Nanonet | + | Minimap | + | Miniasm+ | BWA-MEM | + | Racon | 4,593,402 | 1 | 98.01 | 99.97 | 20,322 | 72,284 |
| 11 | Nanonet | + | Minimap | + | Miniasm+ | Minimap | + | Racon | 4,592,907 | 1 | 98.04 | 100.00 | 20,170 | 70,705 |
| 12 | Scrappie | + | — | + | Canu | + | BWA-MEM | + | Racon | 4,673,871 | 1 | 98.40 | 99.98 | 13,583 | 60,612 |
| 13 | Scrappie | + | — | + | Canu | + | Minimap | + | Racon | 4,673,606 | 1 | 98.40 | 99.98 | 13,798 | 60,423 |
| 14 | Scrappie | + | Minimap | + | Miniasm+ | BWA-MEM | + | Racon | 5,157,041 | 8 | 97.87 | 99.80 | 18,085 | 78,492 |
| 15 | Scrappie | + | Minimap | + | Miniasm+ | Minimap | + | Racon | 5,156,375 | 8 | 97.87 | 99.94 | 17,922 | 77,807 |
| 16 | Nanocall | + | — | + | Canu | + | BWA-MEM | + | Racon | 1,383,851 | 86 | 93.49 | 28.82 | 19,057 | 65,244 |
| 17 | Nanocall | + | — | + | Canu | + | Minimap | + | Racon | 1,367,834 | 86 | 94.43 | 28.74 | 15,610 | 55,275 |
| 18 | Nanocall | + | Minimap | + | Miniasm+ | BWA-MEM | + | Racon | 4,707,961 | 5 | 90.75 | 97.11 | 91,502 | 347,005 |
| 19 | Nanocall | + | Minimap | + | Miniasm+ | Minimap | + | Racon | 4,673,069 | 5 | 92.23 | 97.10 | 72,646 | 291,918 |
| 20 | DeepNano | + | — | + | Canu | + | BWA-MEM | + | Racon | 7,429,290 | 106 | 96.46 | 99.24 | 27,811 | 102,682 |
| 21 | DeepNano | + | — | + | Canu | + | Minimap | + | Racon | 7,404,454 | 106 | 96.03 | 99.21 | 34,023 | 110,640 |
| 22 | DeepNano | + | Minimap | + | Miniasm+ | BWA-MEM | + | Racon | 4,566,253 | 1 | 96.76 | 99.86 | 25,791 | 125,386 |
| 23 | DeepNano | + | Minimap | + | Miniasm+ | Minimap | + | Racon | 4,571,810 | 1 | 96.90 | 99.97 | 24,994 | 119,519 |

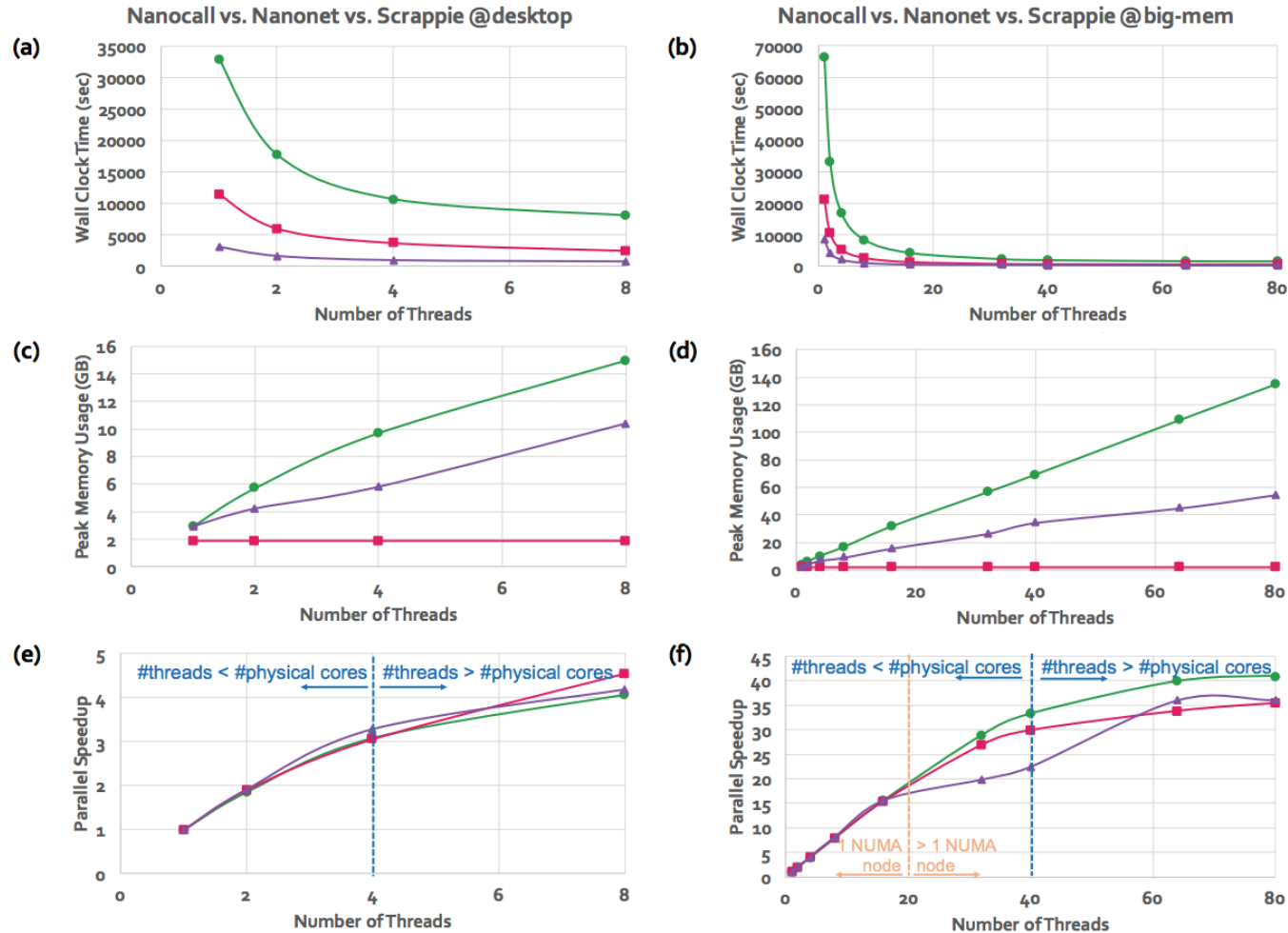Senol Cali+, "**Nanopore Sequencing Technology and Tools for Genome Assembly**" Briefings in Bioinformatics, 2018.

# Nanopore Genome Assembly Tools (II)

**Table 13. Performance analysis results for the full pipeline with a focus on the last two steps.**

| | | | | | | | Step 4: Read Mapper | | | Step 5: Polisher | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Wall Clock Time (h:m:s) | CPU Time (h:m:s) | Memory Usage (GB) | Wall Clock Time (h:m:s) | CPU Time (h:m:s) | Memory Usage (GB) |
| 1 | Metrichor | + — | + Canu | + BWA-MEM | + Nanopolish | | 24:43 | 15:47:21 | 5.26 | 5:51:00 | 191:18:52 | 13.38 |
| 2 | Metrichor | + Minimap | + Miniasm | + BWA-MEM | + Nanopolish | | 12:33 | 7:50:54 | 3.75 | 122:52:00 | 4458:36:10 | 31.36 |
| 3 | Metrichor | + GraphMap | + Miniasm | + BWA-MEM | + Nanopolish | | 12:47 | 7:57:58 | 3.60 | 129:46:00 | 4799:03:51 | 31.31 |
| 4 | Metrichor | + — | + Canu | + BWA-MEM | + Racon | | 24:20 | 15:43:40 | 6.60 | 14:44 | 9:09:22 | 8.11 |
| 5 | Metrichor | + — | + Canu | + Minimap | + Racon | | 3 | 1:35 | 0.26 | 15:12 | 9:45:33 | 14.55 |
| 6 | Metrichor | + Minimap | + Miniasm | + BWA-MEM | + Racon | | 12:10 | 7:48:10 | 5.19 | 15:43 | 9:33:39 | 9.98 |
| 7 | Metrichor | + Minimap | + Miniasm | + Minimap | + Racon | | 3 | 1:24 | 0.26 | 20:28 | 8:57:40 | 18.24 |
| 8 | Nanonet | + — | + Canu | + BWA-MEM | + Racon | | 9:08 | 5:53:18 | 4.84 | 6:33 | 4:02:10 | 4.47 |
| 9 | Nanonet | + — | + Canu | + Minimap | + Racon | | 2 | 54 | 0.26 | 6:45 | 4:17:26 | 7.93 |
| 10 | Nanonet | + Minimap | + Miniasm | + BWA-MEM | + Racon | | 4:40 | 2:58:02 | 3.88 | 7:08 | 4:19:30 | 5.35 |
| 11 | Nanonet | + Minimap | + Miniasm | + Minimap | + Racon | | 2 | 46 | 0.26 | 7:01 | 4:18:48 | 9.53 |
| 12 | Scrappie | + — | + Canu | + BWA-MEM | + Racon | | 33:41 | 21:11:06 | 8.66 | 13:32 | 8:24:44 | 7.58 |
| 13 | Scrappie | + — | + Canu | + Minimap | + Racon | | 3 | 1:39 | 0.27 | 18:45 | 7:43:17 | 13.20 |
| 14 | Scrappie | + Minimap | + Miniasm | + BWA-MEM | + Racon | | 22:41 | 14:31:00 | 6.08 | 14:37 | 8:53:59 | 9.50 |
| 15 | Scrappie | + Minimap | + Miniasm | + Minimap | + Racon | | 3 | 1:27 | 0.27 | 15:10 | 9:02:45 | 12.72 |
| 16 | Nanocall | + — | + Canu | + BWA-MEM | + Racon | | 4:52 | 3:01:15 | 3.80 | 11:07 | 3:26:52 | 5.63 |
| 17 | Nanocall | + — | + Canu | + Minimap | + Racon | | 3 | 1:16 | 0.22 | 7:28 | 2:50:35 | 3.62 |
| 18 | Nanocall | + Minimap | + Miniasm | + BWA-MEM | + Racon | | 16:06 | 10:27:20 | 5.06 | 18:56 | 11:32:45 | 11.47 |
| 19 | Nanocall | + Minimap | + Miniasm | + Minimap | + Racon | | 4 | 1:18 | 0.26 | 11:49 | 7:08:59 | 10.98 |
| 20 | DeepNano | + — | + Canu | + BWA-MEM | + Racon | | 17:36 | 11:30:20 | 4.43 | 12:48 | 7:13:04 | 8.88 |
| 21 | DeepNano | + — | + Canu | + Minimap | + Racon | | 3 | 1:24 | 0.28 | 11:39 | 6:55:01 | 3.73 |
| 22 | DeepNano | + Minimap | + Miniasm | + BWA-MEM | + Racon | | 8:15 | 5:22:29 | 4.11 | 14:16 | 8:34:32 | 10.30 |
| 23 | DeepNano | + Minimap | + Miniasm | + Minimap | + Racon | | 3 | 1:10 | 0.26 | 12:29 | 7:55:32 | 17.11 |

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

SAFARI

# Nanopore Genome Assembly Tools (III)

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

# More on Nanopore Sequencing & Tools

## 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

BiB          arXiv

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]

# Agenda

- The Problem: DNA Read Mapping
  - State-of-the-art Read Mapper Design

- Algorithmic Acceleration
  - Exploiting Structure of the Genome
  - Exploiting SIMD Instructions

- Hardware Acceleration
  - Specialized Architectures
  - Processing in Memory

- Future Opportunities: New Sequencing Technologies

**SAFARI**

# Conclusion

- System design for bioinformatics is a critical problem
  - It has large scientific, medical, societal, personal implications

- This talk is about accelerating a key step in bioinformatics: genome sequence analysis
  - In particular, read mapping

- We covered various recent ideas to accelerate read mapping
  - My personal journey since September 2006

- **Many future opportunities exist**
  - **Especially with new sequencing technologies**
  - **Especially with new applications and use cases**

# Acknowledgments

- Can Alkan, Bilkent University

- Many students at ETH, CMU, Bilkent
  - Mohammed Alser, Damla Senol Cali, Jeremie Kim, Hasan Hassan, Donghyuk Lee, Hongyi Xin, …

- All papers, source code, and more are at:
  - https://people.inf.ethz.ch/omutlu/projects.htm

**SAFARI**

# 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**