# Accelerating Genome Analysis
## A Primer on an Ongoing Journey

Onur Mutlu

omutlu@gmail.com

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

January 24, 2018

AACBB Keynote, Vienna

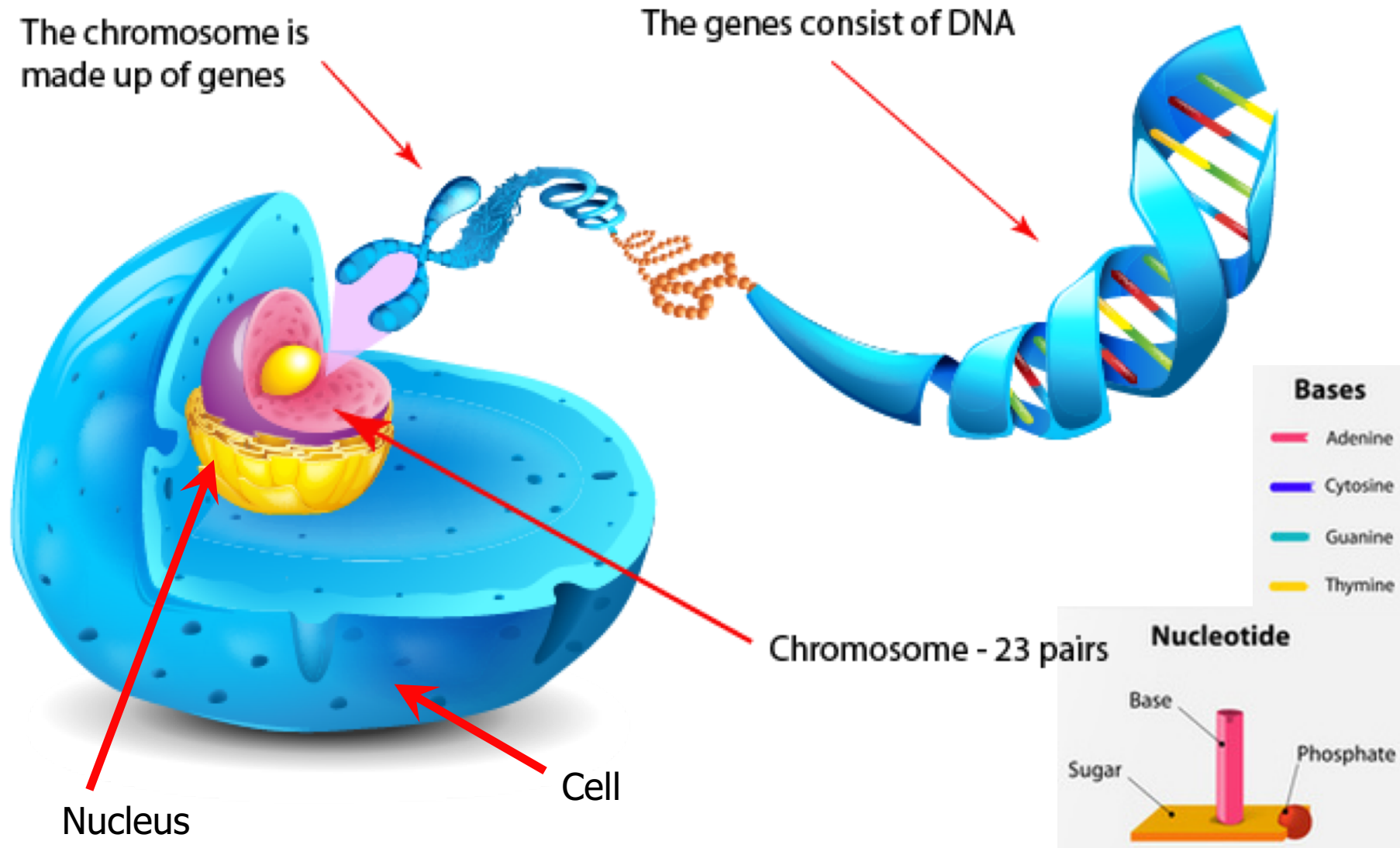Systems@ETH zürich

SAFARI

ETH zürich

# 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

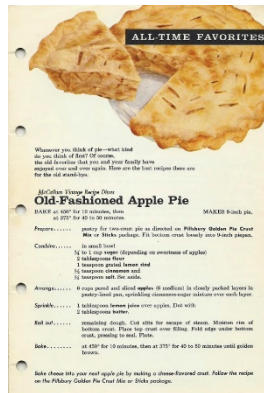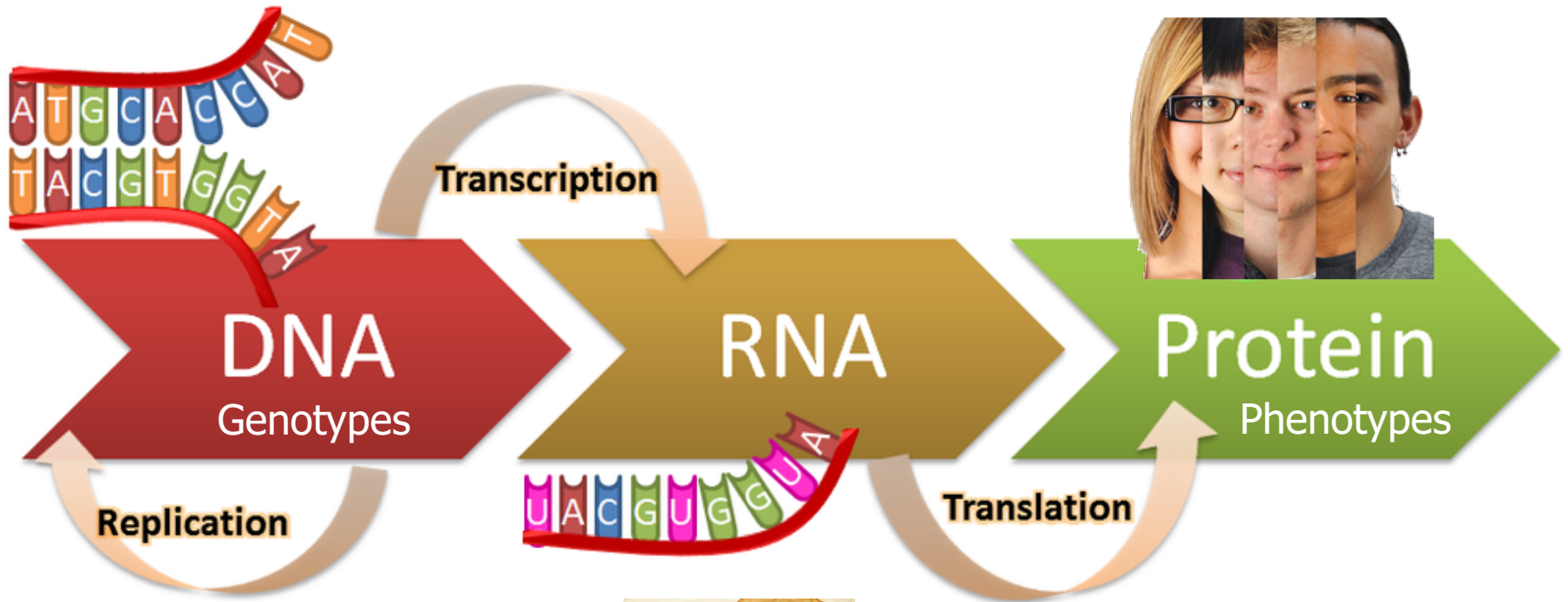**SAFARI**

# 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

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

Chromosome - 23 pairs

**Nucleotide**

Base

Sugar

Phosphate

Nucleus

Cell

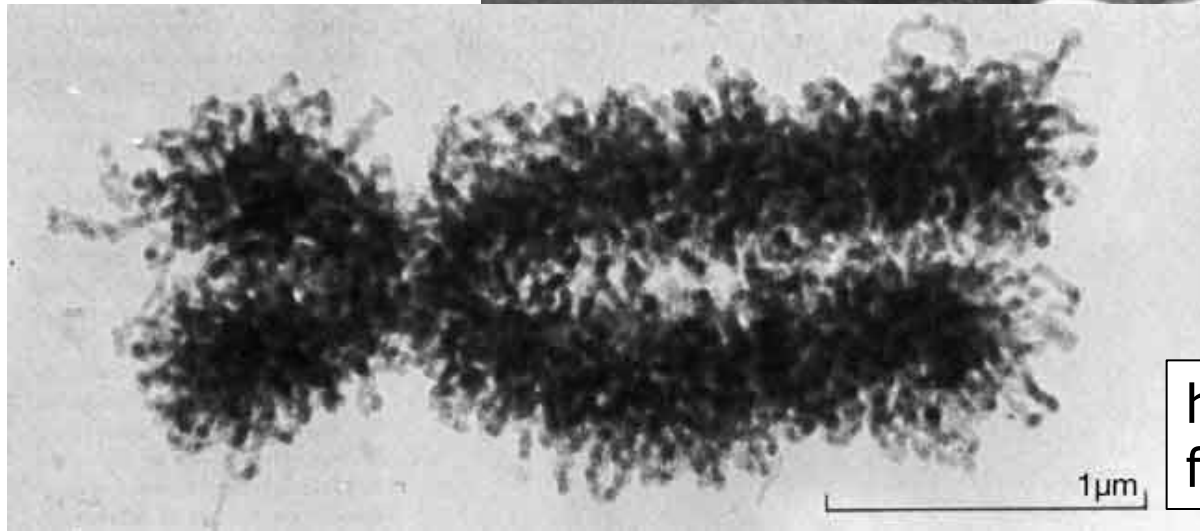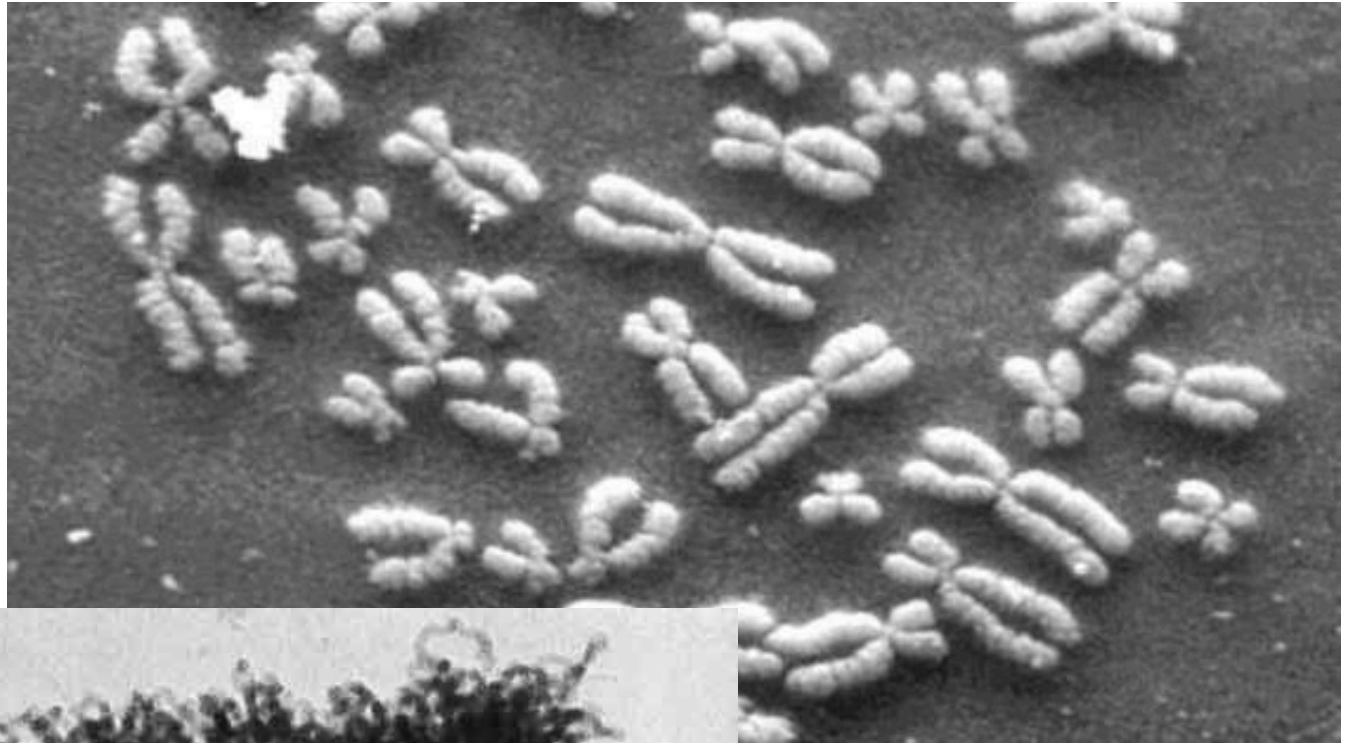The discovery of DNA's double-helical structure (Watson+, 1953)

# The Central Dogma of Molecular Biology

# DNA Under Electron Microscope

human chromosome #12
from HeLa's cell

# 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)

**SAFARI**

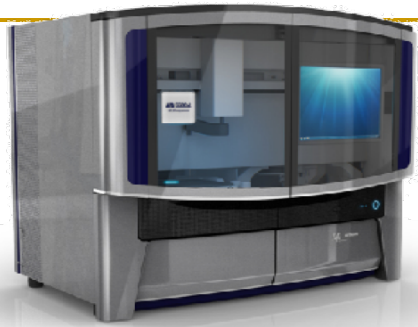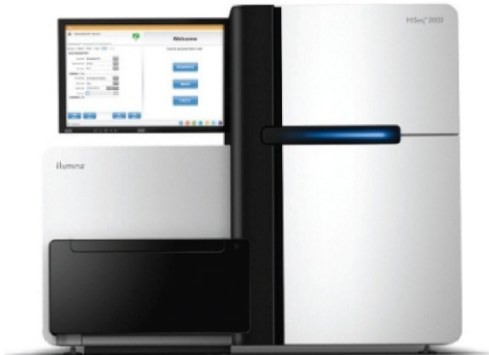# Untangling Yarn Balls & DNA Sequencing

**SAFARI**

# Genome Sequencers

Roche/454

AB SOLiD

Illumina MiSeq

Complete Genomics

Illumina HiSeq2000

Pacific Biosciences RS

Oxford Nanopore MinION

Illumina NovaSeq 6000

Oxford Nanopore GridION

**SAFARI**

Ion Torrent PGM

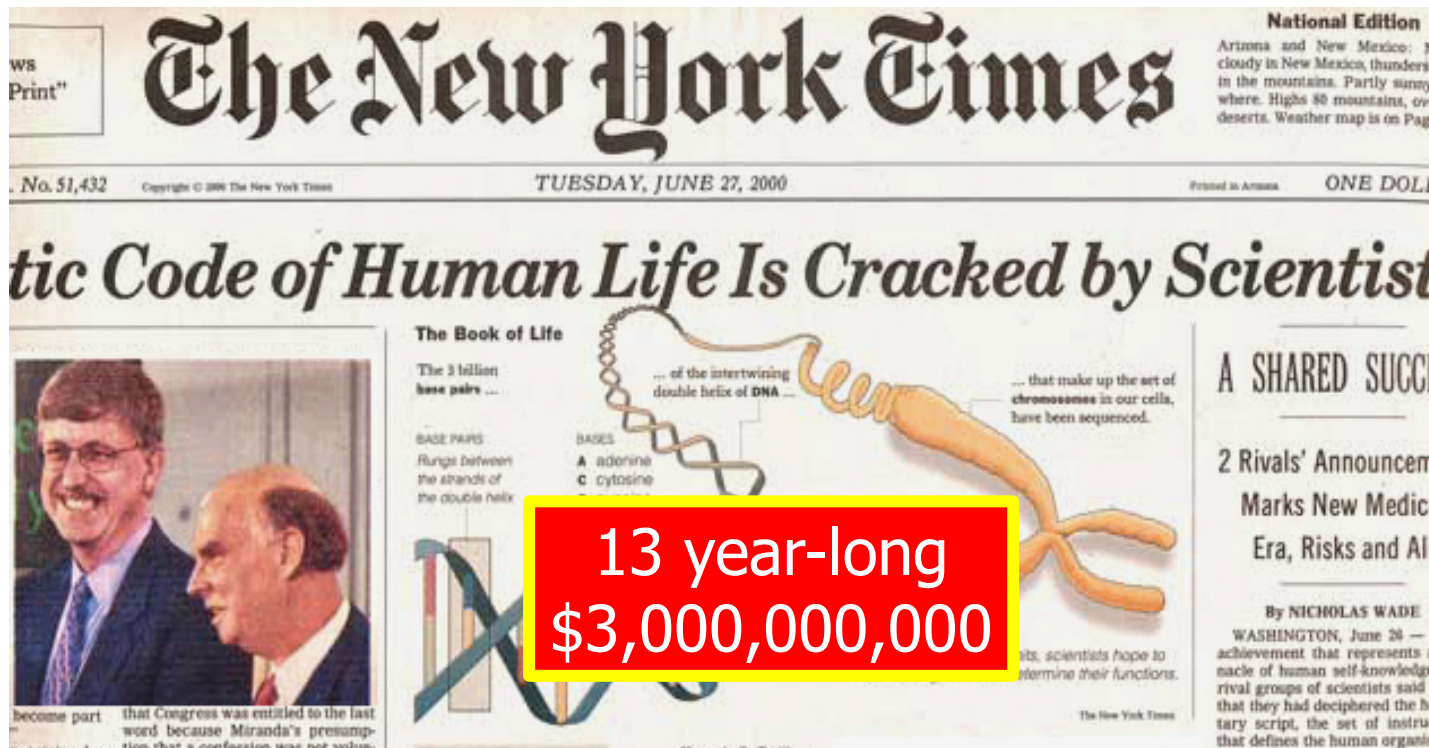Ion Torrent Proton

**... 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

# 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

Cleave fluorescence, wash away

**SAFARI**

# High-Throughput Sequencing (HTS)



Sequence

**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**
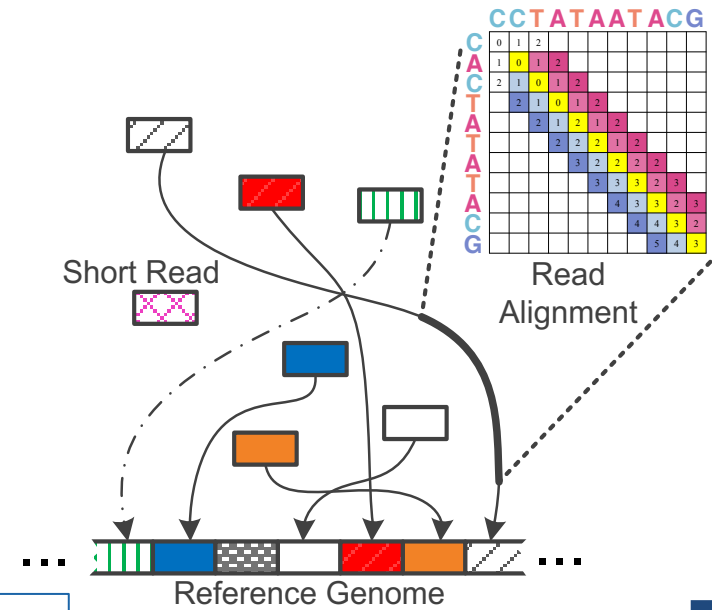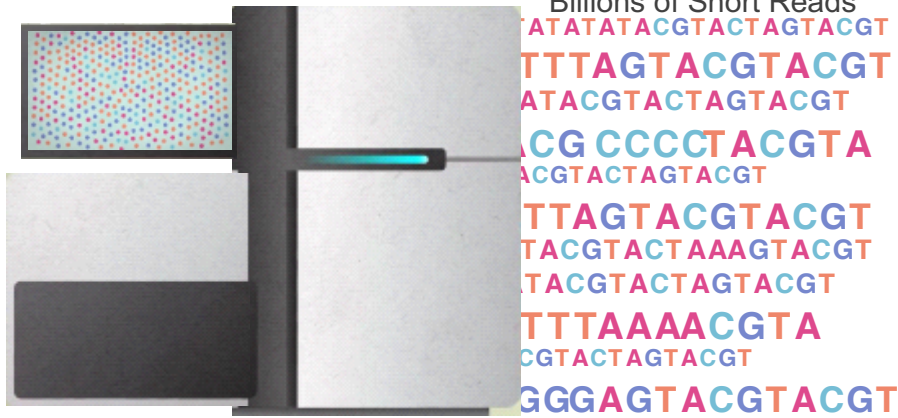  - High throughput
  - High speed
  - Low cost
  - Short reads
    - Amplification step limits the read length since too short or too long fragments are not amplified well.

- **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

**Genome Analysis**

**1  Sequencing**

Billions of Short Reads

**2  Read Mapping**

Short Read
Read Alignment
Reference Genome

**3  Variant Calling**

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

**4  Scientific Discovery**

PRESCRIPTION

**1** **Sequencing**

Billions of Short Reads

**Read Mapping** **2**

Short Read

Read Alignment

Reference Genome

Bottlenecked in Mapping!!

Illumina HiSeq4000

300 M
bases/min

on average

2 M
bases/min
(0.6%)

# The Read Mapping Bottleneck



Illumina HiSeq4000

300 **Million** bases/minute

2 **Million** bases/minute

**150X slower**

# Read Mapping Execution Time Breakdown



SAM printing
3%

candidate alignment
locations (CAL)
4%

Read Verification
93%

SAFARI

# Read Mapping

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

Reference genome

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 → 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.

organization x operation

Ref Read

| o | - | - | r | g | a | n | i | z | a | t | i | o | n |
| o | p | e | r | - | - | - | - | - | a | t | i | o | n |

Ref Read

| o | - | - | r | g | a | n | i | z | a | t | i | o | n |
| o | p | e | r | - | a | - | - | - | - | t | i | o | n |

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

organization x translation

Ref Read

| o | r | g | a | n | i | z | - | a | t | i | o | n |
| t | r | - | a | n | - | s | l | a | t | i | o | n |

Ref Read

| o | r | g | a | n | - | i | z | a | t | i | o | n |
| t | r | - | a | n | s | l | - | a | t | i | o | n |

Ref Read

| o | r | g | a | n | i | z | a | t | i | o | n |
| t | r | - | a | n | s | l | a | t | i | o | n |

# 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, 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 | |

Reference genome

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

| AAAAAAAAAAAT | → | NULL |

| ...... |

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

| ...... |

| ...... |

| ...... |

| TTTTTTTTTTTT | → | 36 | 535 | 123 |

**Once for a reference**

# 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]

AAAAAAAAAAAACCCCCCCCCCCCTTTTTTTTTTTT ← read

CCCCCCCCCCCC ← k-mers
AAAAAAAAAAAA

**Hash Table (HT)**

3224

Reference Genome

AAAAAAAAAAAA → | 12 | 324 | 557 | 940 |

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

TTTTTTTTTTTT → | 36 | 535 | 823 |

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

AAAAAAAAAAAACCCCCCCCCCCCTTTTTTTTTTTT

read

Valid mapping

**Verification/Local Alignment**

# Advantages of Hash Table Based Mappers

- + Guaranteed to find *all* mappings → 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**

# 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**

# Reducing the Cost of Verification

- We observe that most verification (edit distance computation) calculations are unnecessary
  - 1 out of 1000 potential locations passes the verification process

- 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

# 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
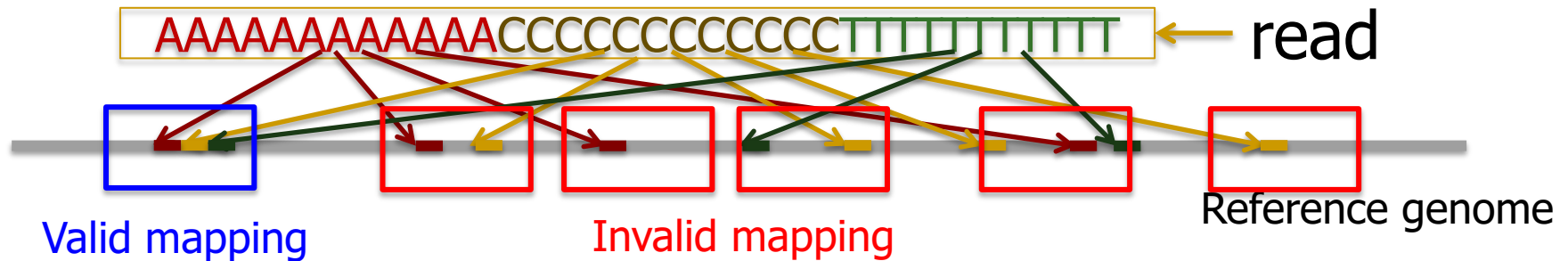  - Read mapper can choose the cheapest $e+1$ k-mers and verify their locations

# 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

# 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)



read

+12   +24

k-mers

**Hash Table (HT)**

Reference Genome

952?

| AAAAAAAAAAAA | | 12 | 324 | 557 | 940 |
| CCCCCCCCCCCC | | 24 | 459 | 744 | 988 | 989 |
| TTTTTTTTTTTT | | 36 | 535 | 123 |

...AAAAAAAAAAAACCCCCCCCCCCCTTTTTTTTTTTT...

AAAAAAAAAAAACCCCCCCCCCCCTTTTTTTTTTTT

# FastHASH Mechanisms [Xin+, BMC Genomics 2013]

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

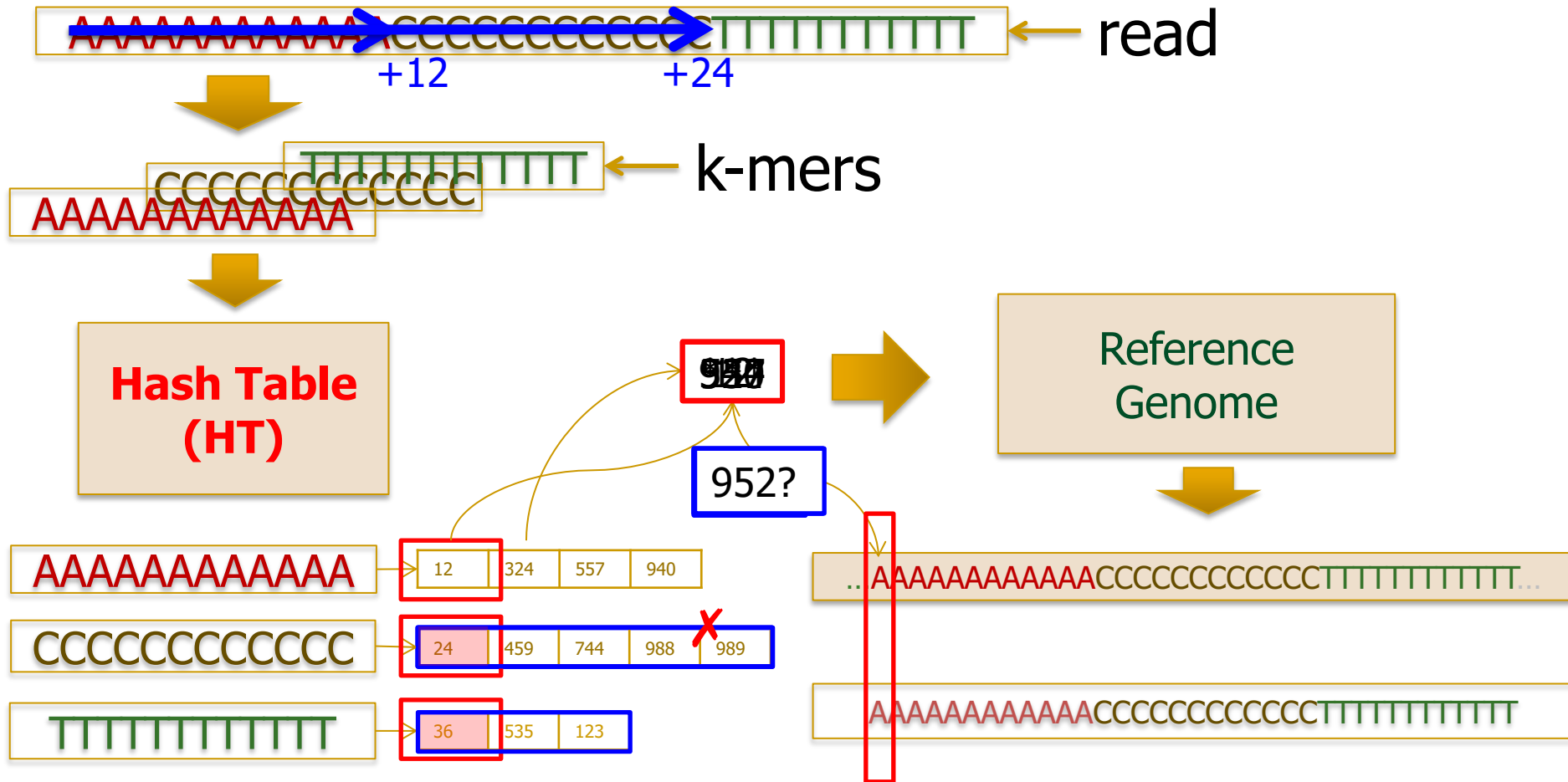- **Cheap K-mer Selection (CKS):** Reduces the absolute number of potential mapping locations
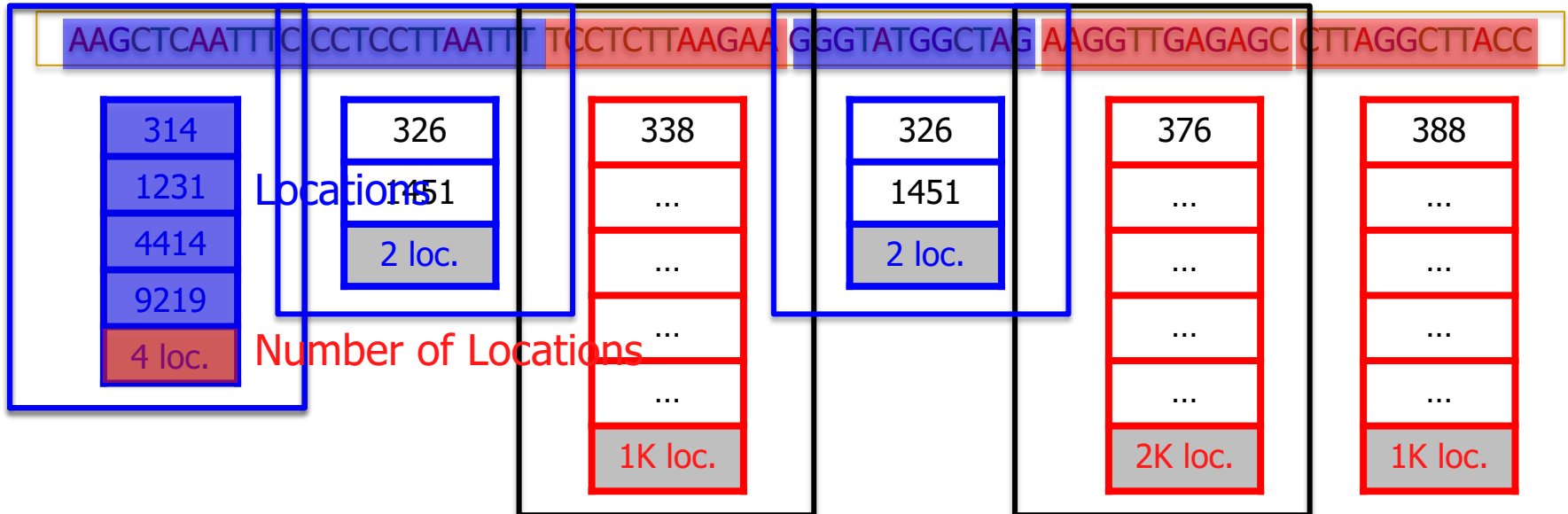
# Cheap K-mer Selection (CKS)

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

- **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 fewest locations to verify

# Cheap K-mer Selection

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

AAGCTCAATTTC CCTCCTTAATTT TCCTCTTAAGAA GGGTATGGCTAG AAGGTTGAGAGC CTTAGGCTTACC

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

Locations

| 326 |
| 1451 |
| 2 loc. |

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

| 326 |
| 1451 |
| 2 loc. |

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

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

Number of Locations

Expensive 3 k-mers

Previous work needs to verify:

3004 locations

FastHASH verifies only:
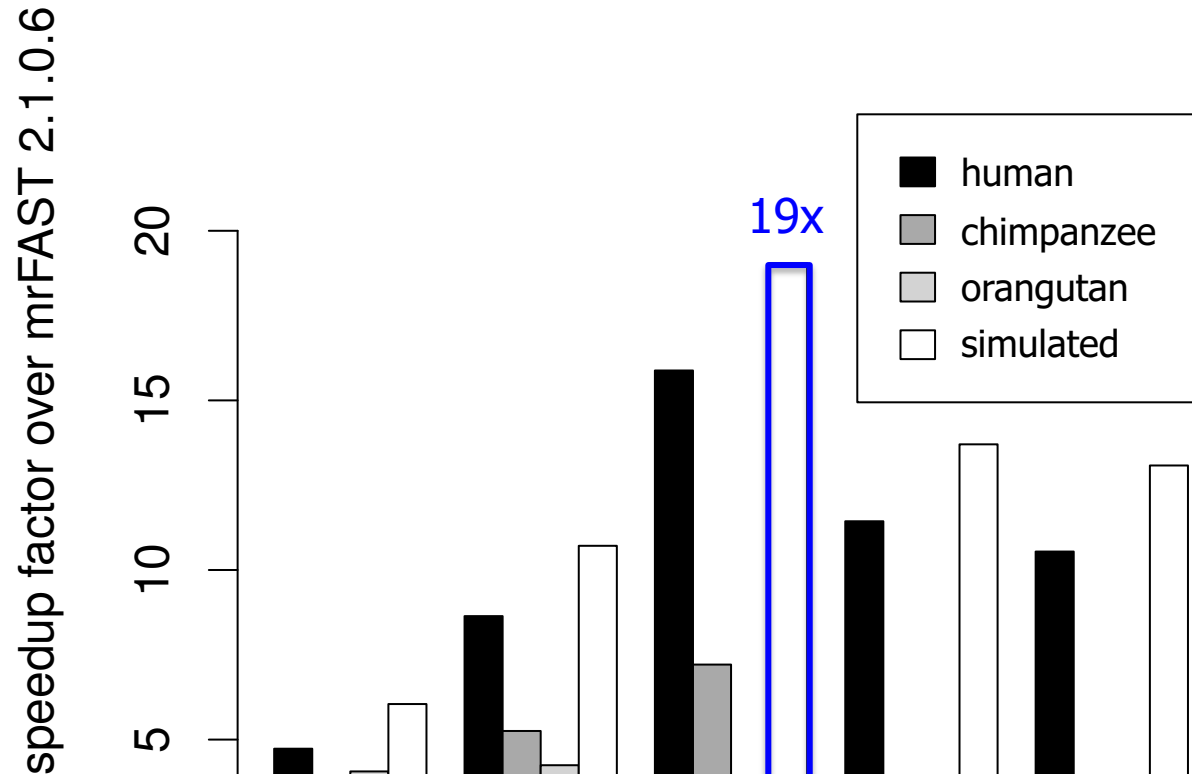
8 locations

# 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
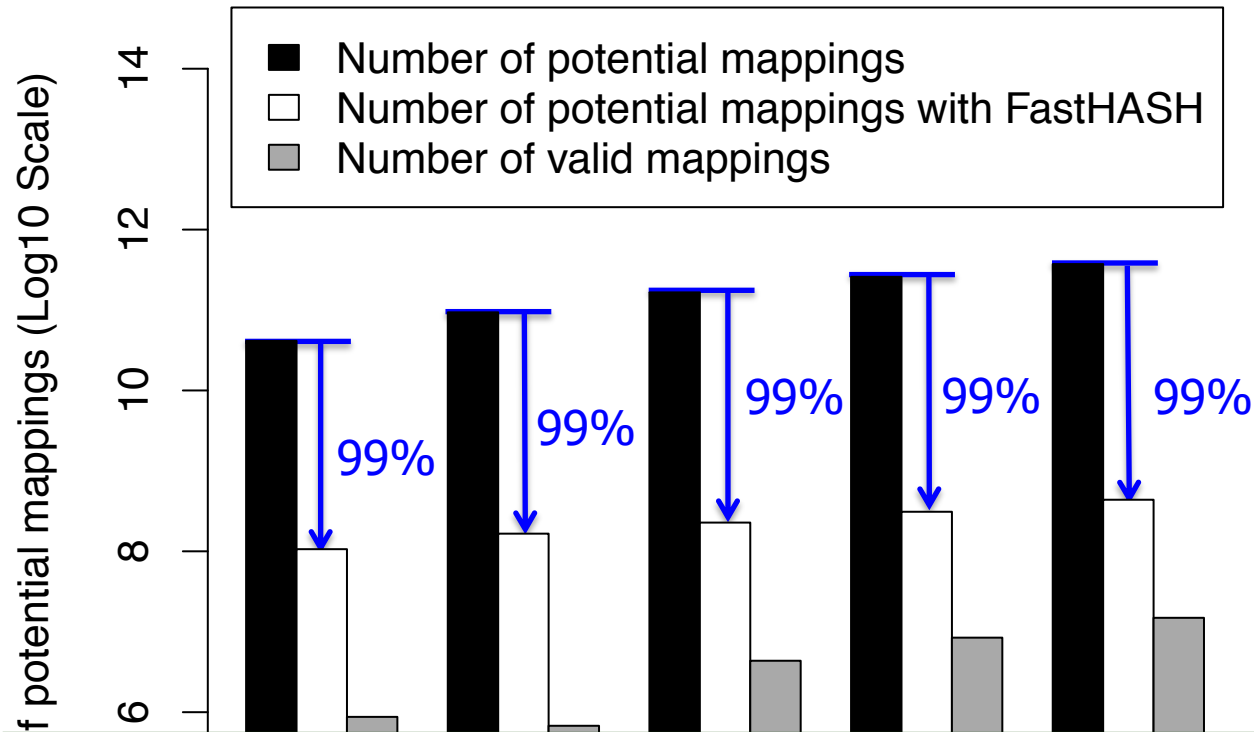  - 16 GB of main memory

# FastHASH Speedup



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

e: edit distance

# Analysis

- Reduction of potential mappings with FastHASH



Reduction of potential mappings with FastHASH

# FastHASH Conclusion

- Problem: Existing read mappers perform poorly in mapping billions 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: To reduce the cost of unnecessary verification
  - 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*]

# 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

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.

- **Key idea:**
  - Compute "Shifted Hamming Distance": AND of 2E Hamming Distances of two strings, to identify 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

# 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 <span style="color:red">expensive</span>
  - We need to align millions to billions of reads

- M                                                    t
f

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

  out mismatches quickly

- Both methods are used by mappers today, but <span style="color:purple">filtering has replaced alignment as the bottleneck</span> **[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)

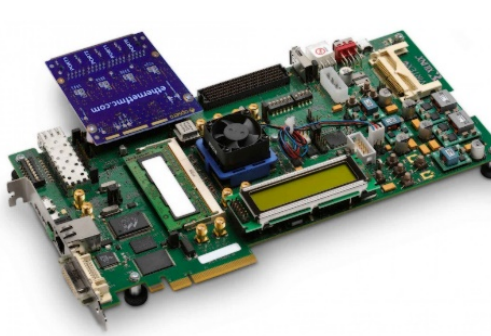## Needleman-Wunsch

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

| 0 | 1 | 2 | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 2 | | | | | | |

**A**

## GateKeeper

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

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | 1 | 1 | 0 | | | | | | |

**A**

- ## Independent vectors can be processed in parallel using hardware technologies



```
                |dp[i][j-1]  // Inser.
dp[i][j]=1+max|dp[i-1][j]   // Del.
                |dp[i-1][j-1]// Subs.
```

Each cell depends on three
pre-computed cells!

```
dp[i][j]=|0 if X[i]=Y[j]
          |1 if X[i]≠Y[j]
```

No data dependencies!

# Our Solution: GateKeeper

Alignment Filter $+$ [FPGA board] $=$ $1^{st}$ FPGA-based Alignment Filter.

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

$x10^{12}$ mappings $\rightarrow$ $x10^{3}$ 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 :00000000001000000000000111111101111000111011010110111111111000100000111101101001 0101
1-Deletion Mask :11111111111001111101111100000000000000000000000000000000000000000011000000000000000
2-Deletion Mask :0000000010110110011111111111111011100011101101011011111111100010001001110110100 1010
3-Deletion Mask :11111111111011101100110111011101100010010011111111111110010110011001011011101 1101111
1-Insertion Mask :111111111110111110111111011101100010010011111111111110010110011000101011101110 111110
2-Insertion Mask :00000010011110011111111100100011010101001101011111111111110111001111111000111 101100
3-Insertion Mask :1111111101110110011000111111111010110111110011001011101111111101110111101011 1001000

                      --- Masks after amendment ---

 Hamming Mask :00000000001000000000000111111111110001111111101111111111111000100000111111111111 1111
1-Deletion Mask :11111111111111111111111100000000000000000000000000000000000000000011000000000000000
2-Deletion Mask :00000000111111111111111111111111111110001111111111111111111111100010001111111 1111110
3-Deletion Mask :11111111111111111111111111111111111000111111111111111111111111111111111111111 1111111
1-Insertion Mask :11111111111111111111111111111111110001111111111111111111111111111000111111111 1111110
2-Insertion Mask :00000001111111111111111111110001111111111111111111111111111111111111100011111 11100
3-Insertion Mask :1111111111111111110001111111111111111111111111111111111111111111111111111111 111000

      AND Mask :000000000010000000000001000000000000000000000000000000000000000000000010000000000 00000
```

```
                GAGAGAGATATTTAGTGTTGCAG-CACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGGAACATTGTTGGGCCGG
Needleman-Wunsch |||||||||| ||||||||||| ||||||||||||||||||||||||||||||||||||||||||::|||||||||||||
 Alignment :    GAGAGAGATAGTTAGTGTTGCAGCCACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGGAGACATTGTTGGGCCGG
```
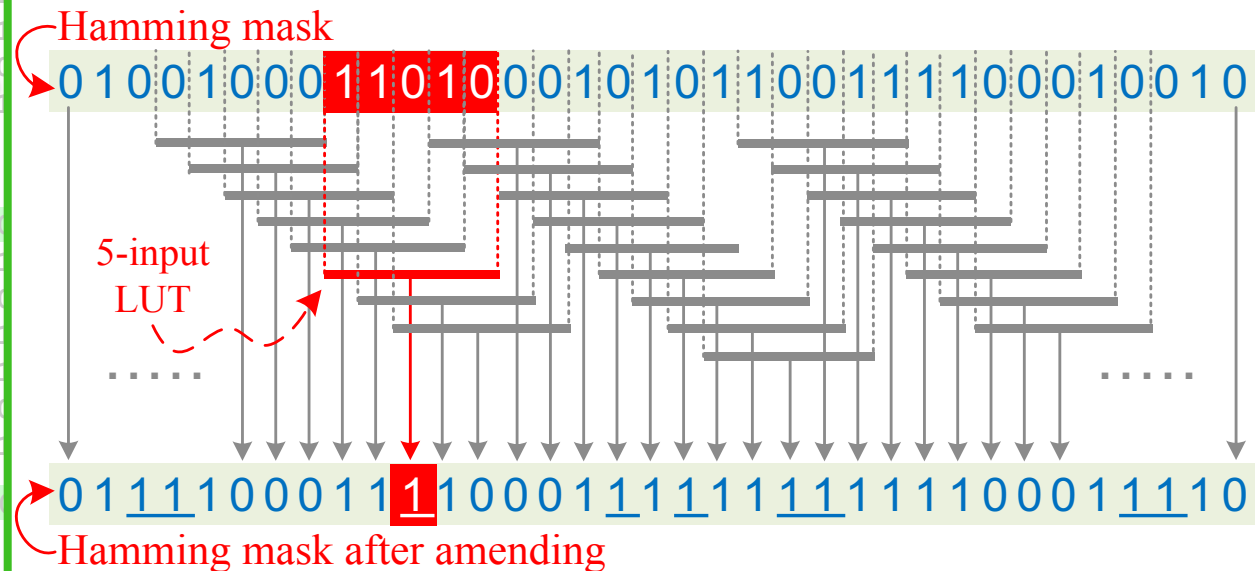
# GateKeeper Walkthrough (cont'd)

**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)
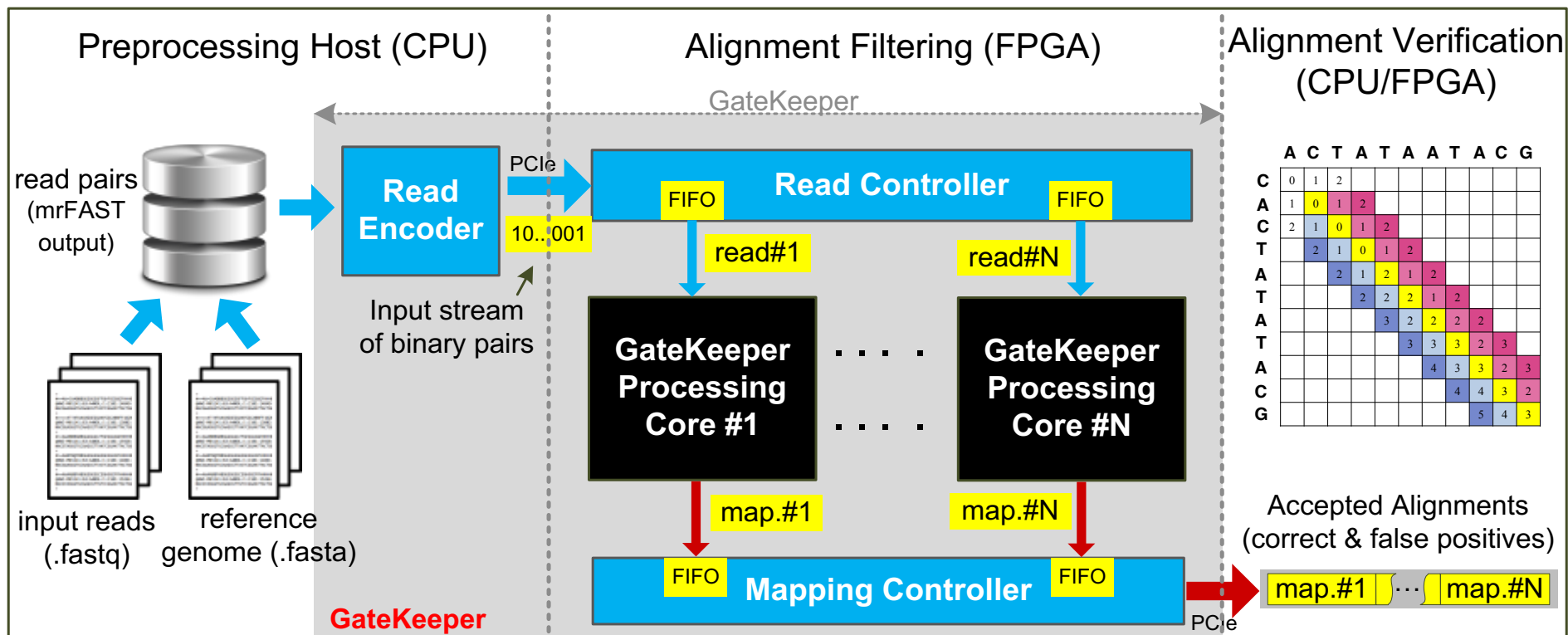- (2E+1) * (ReadLength) 2-XOR operations.

- (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 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 0 0 0 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,*]

# Next Talk: MAGNET

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

- Key Idea: 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**

# 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

# Hash Tables in Read Mapping

**Read Sequence (100 bp)**

A match! Matching...    Mismatch. Mismatch. **False Negative**

Hash Table

**Reference Genome**

**Filter**

37     140
894     1203
1564

**SAFARI**

63

# Read Mapping & Filtering in Memory

We need to design

mapping & filtering algorithms
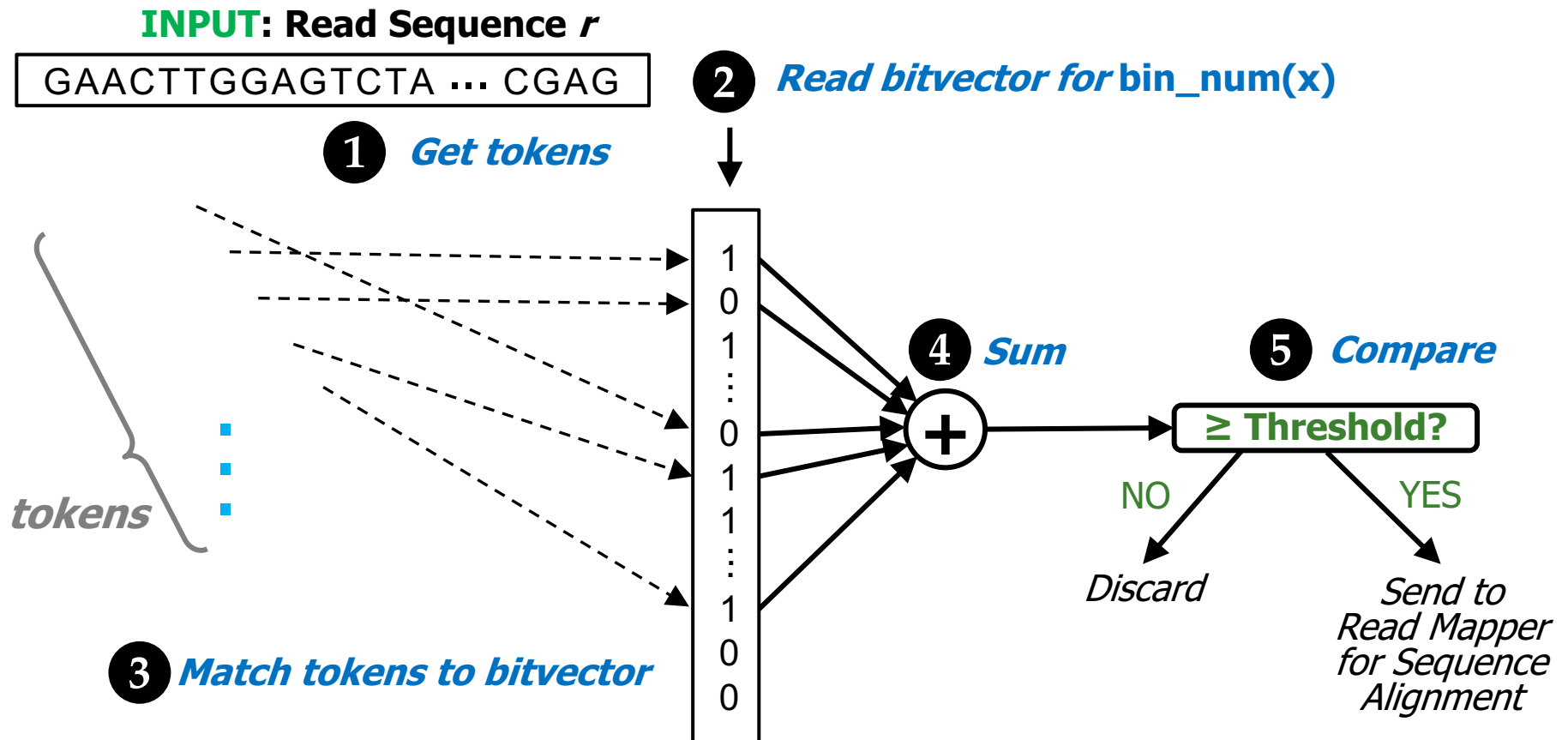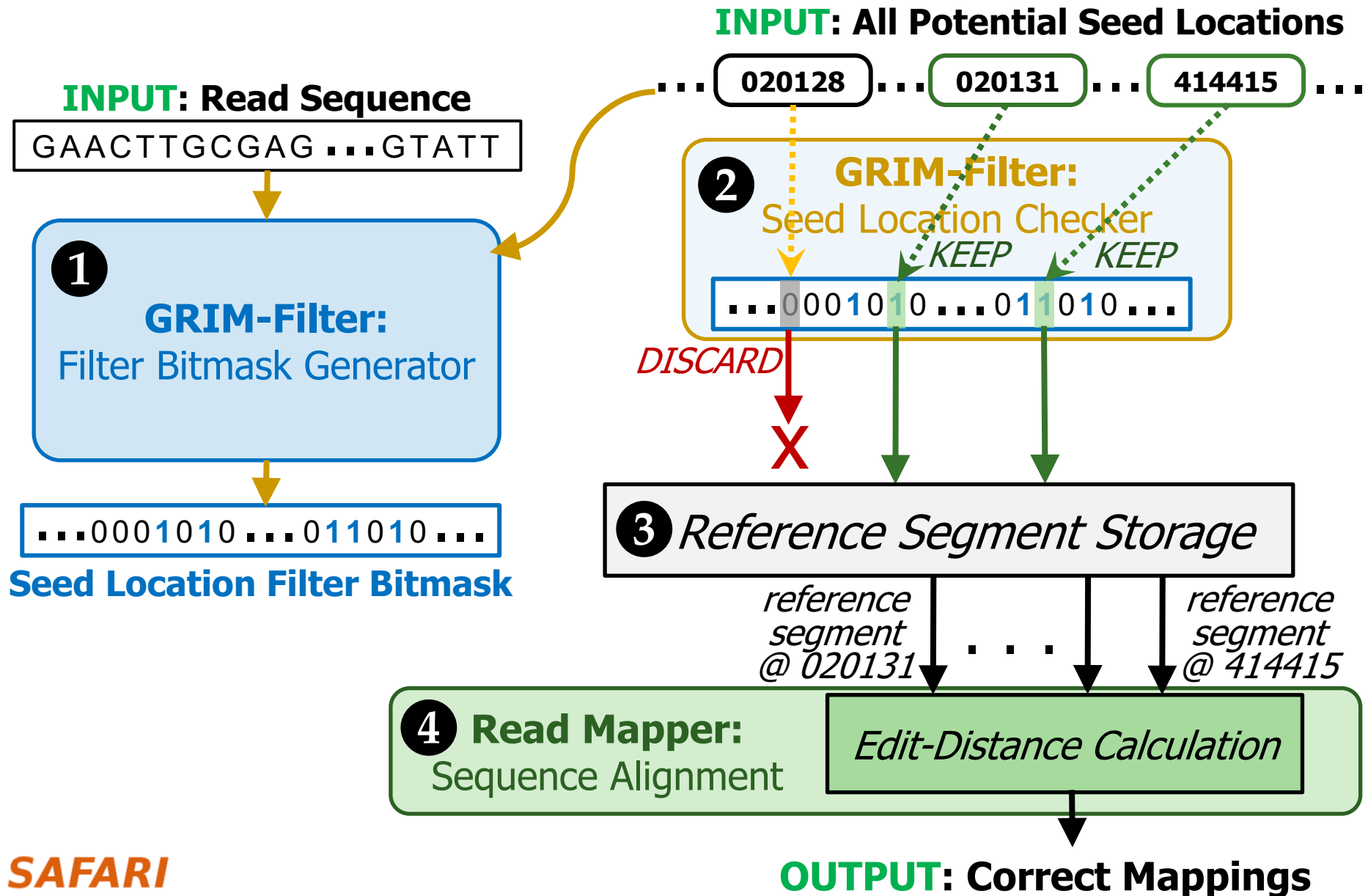
that fit processing-in-memory

# 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 | **AAAAA** **exists** in bin x |
| AAAAC | 0 | |
| AAAAT | 1 | |
| ... | ... | |
| CCCCC | 1 | |
| **CCCCT** | 0 | **CCCCT** **doesn't** **exist** in bin x |
| CCCCG | 0 | |
| ... | ... | |
| GGGGG | 1 | |

**SAFARI**

# GRIM-Filter: Bitvectors

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

*Bin x*

|  |  |
|---|---|
| AAAAA | 0 |
| ... | ... |
| CGTGA | 1 |
| ... | ... |
| TGAGT | 1 |
| ... | ... |
| GAGTC | 1 |
| ... | ... |
| GTGAG | 1 |
| ... | ... |

**Bin x Bitvector**

# GRIM-Filter: Bitvectors



Reference Genome: AAAAACCCCTGCCTTGCATGTAGAAAACTTGACAGGAACTTTTTATCGCA ...

bin₁, bin₂, bin₃, bin₄

**b₁**

| tokens | b₁ |
|--------|----|
| AAAAA | 1 |
| AAAAC | 1 |
| AAAAG | 0 |
| AAAAT | 0 |
| . | . |
| CCCCT | 1 |
| . | . |
| . | . |
| . | . |
| GCATG | 1 |
| . | . |
| TTGCA | 1 |
| . | . |
| TTTTT | 0 |

**b₂**

| | b₂ |
|--|----|
| 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

**SAFARI**

# 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

**SAFARI**

# Our Proposal: GRIM-Filter

1.  Data Structures: Bins & Bitvectors

2.  Checking a Bin

3.  Integrating GRIM-Filter into a Mapper

**SAFARI**

# Our Proposal: GRIM-Filter

1.  Data Structures: Bins & Bitvectors

2.  Checking a Bin

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

**SAFARI**

# Integrating GRIM-Filter into a Read Mapper

**INPUT: Read Sequence**

GAACTTGCGAG...GTATT

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

...0001010...011010...

**Seed Location Filter Bitmask**

**INPUT: All Potential Seed Locations**

...( 020128 )...( 020131 )...( 414415 )...

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

*KEEP*          *KEEP*

...0001010...011010...

*DISCARD*

✗

**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**

# 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

**SAFARI**

# 3D-Stacked Memory

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

**SAFARI**

# 3D-Stacked Memory



Micron's HMC

Micron has working demonstration components

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

http://i1-news.softpedia-static.com/images/news2/Micron-and-Samsung-Join-Force-to-Create-Next-Gen-Hybrid-Memory-2.png

SAFARI

# GRIM-Filter in 3D-Stacked DRAM



- 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



- 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



**Sequence Alignment Error Tolerance ($e$)**

$e = 0.05$

**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"**
*to appear in **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]

# 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
  - No amplification → Less limit on 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

**SAFARI**

# 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

**SAFARI**

# 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.

**SAFARI**

# 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

**SAFARI**

# 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**" 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 [1,*], Jeremie Kim [1,3], Saugata Ghose [1], Can Alkan [2*] and Onur Mutlu [3,1*]

[1] Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA, USA
[2] Department of Computer Engineering, Bilkent University, Bilkent, Ankara, Turkey
[3] Department of Computer Science, Systems Group, ETH Zürich, Zürich, Switzerland

Senol Cali+, "**Nanopore Sequencing Technology and Tools for Genome Assembly: Computational Analysis of the Current State, Bottlenecks and Future Directions**," to appear in 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**

**SAFARI**

# Acknowledgments

- Prof. Can Alkan, Bilkent University

- Many students
  - Mohammed Alser, Damla Senol Cali, Jeremie Kim
  - Hasan Hassan
  - Hongyi Xin
  - ...

# Accelerating Genome Analysis
## A Primer on an Ongoing Journey

Onur Mutlu

omutlu@gmail.com

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

January 24, 2018

AACBB Keynote, Vienna

Systems@**ETH** zürich

**ETH** zürich

**SAFARI**

# High-Throughput Sequencing

**SAFARI**

# Nanopore Sequencing

- **Basecalling** translates the raw signal output of the nanopore sequencer into bases (A, C, G, T) to generate DNA reads.
  - 1) The raw current signal is divided into discrete blocks (events).
  - 2) Each event is decoded into a most-likely set of bases.

- **Deletions** are the dominant error of nanopore sequencing.
  - In the ideal case, each consecutive event should differ by one base. However, in practice, this is not the case because of the non-stable speed of the translocation.
  - Determining the correct length of the homopolymers (*i.e.*, repeating stretches of one kind of base, *e.g.*, AAAAAAA) is challenging.

**SAFARI**

# The Importance of Genome Analysis?

Helps, for example, to answer the following 3 questions:

# 69–92% of the respondents in these studies had positive attitudes towards genomics research and donating their DNA samples.

**Public involvement in pharmacogenomics research: a national survey on public attitudes towards pharmacogenomics research and the willingness to donate DNA samples to a DNA bank in Japan**

Eriko Kobayashi · Nobunori Satoh

**Attitudes and perceptions of patients towards methods of establishing a DNA biobank**

Pulley · Margaret M. Brace · Gordon R. Bernard · Masys

# Genetic research participation in a young adult community sample

Carla L. Storr · Flora Or · William W. Eaton · Nicholas Ialongo

American Journal of Medical Genetics Part A 146A:1696–1706 (2008)

## Miscellaneous

## Genetic research and donation of tissue samples to biobanks. What do potenti sample donors in the Swedish general public think?

# Relationship Between Public Attitudes Toward Genomic Studies Related to Medicine and Their Level of Genomic Literacy in Japan

Izumi Ishiyama,[1] Akiko Nagai,[1] Kaori Muto,[2] Akiko Tamakoshi,[3] Minori Kokado,[4] Kyoko Mimura,[5] Tetsuro Tanzawa,[6] and Zentaro Yamagata[1*]

Åsa Kettis-Lindblad[1], Lena Ring[1,2], Eva Viberth[1], Mats G. Hansson[3]

# Pairwise sequence alignment



BLAST

Distribution of 116 Blast Hits on the Query Sequence

A57075 tensin - chicken (fragment) gi|63805|emb|CAA79215.1| (..S= 492 E=1e-137

Question #1: If I give you a gene sequence, tell me which of the billions of known sequences is most similar to it.

**SAFARI**

# CODIS: Combined DNA Index System

- FBI's program of support for criminal justice.
- CODIS defines 13 human DNA regions (loci) to be stored in the database for personal identification purposes.
- Stored 14.5 million DNA profiles (for offenders, arrestees ..)
- As of September 2016, CODIS has produced over 346,880 hits assisting in more than 332,776 investigations.



https://www.fbi.gov/services/laboratory/biometric-analysis/codis/ndis-statistics

**SAFARI**

# Multiple sequence alignment



Question #2: If I give you a bunch of sequences, tell me where they are the same and where they are different.

**SAFARI**

# Phylogenetic tree

- Reveals the genomic variants that cause diseases.

- Helps understanding the evolutionary relationships among various species.



Manhattan plot

Evolutionary Tree

**SAFARI**

# The genetic similarity between species



Human ~ Chimpanzee
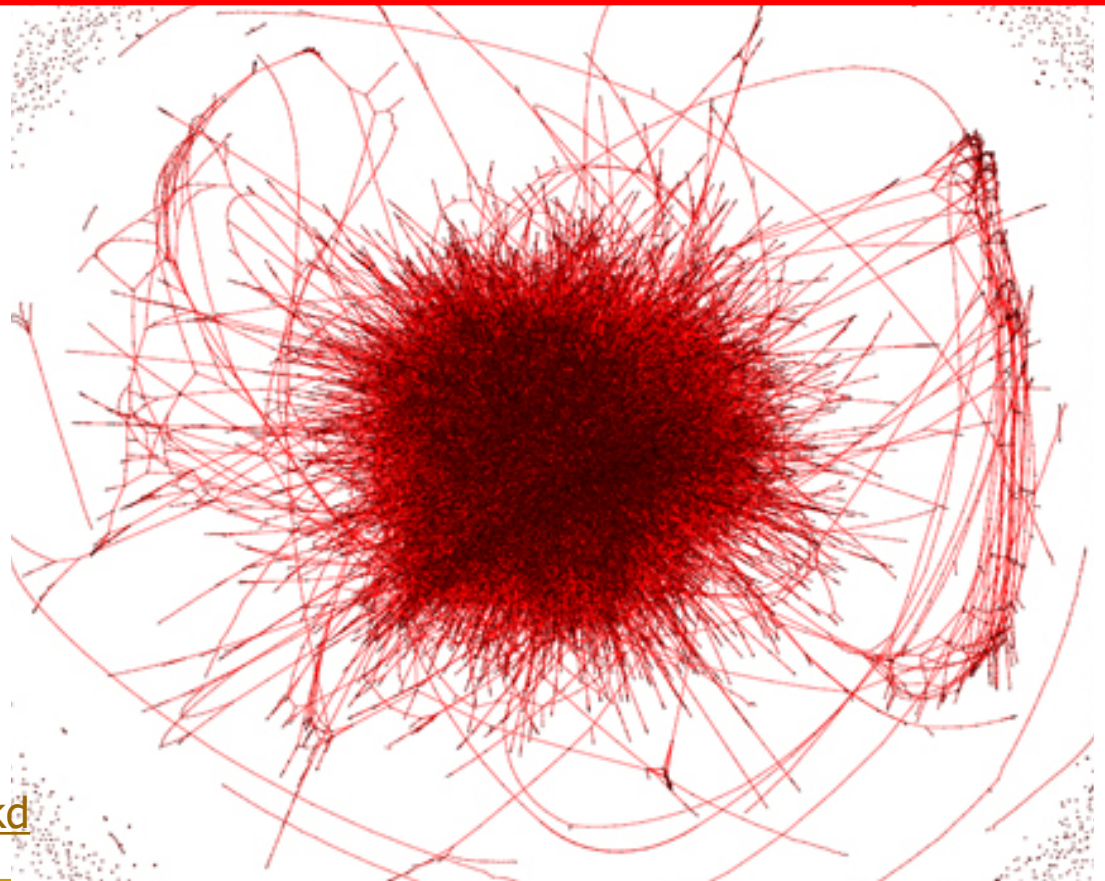96%


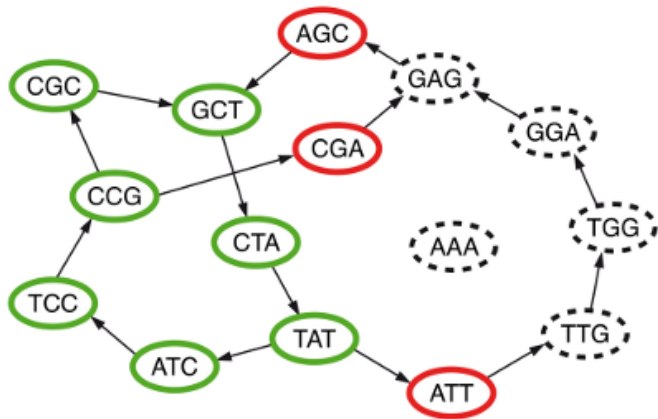
Human ~ Cat
90%



Human ~ Human
99.9%



Human ~ Cow
80%



Human ~ Banana
50-60%

SAFARI

Question #3: 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

**SAFARI**

# ANALYZING THE PROBLEM

**SAFARI**

# Read Mapping

# Key Observations:

- Alignment Verification → **90%** of mapper's execution time.

- **>98%** of candidate locations have high dissimilarity with a given read.

Cheng *et al, BMC bioinformatics* (2015)
Xin *et al, BMC genomics* (2013)

# Read Mappers Timeline



Legend:
- **CPU** (black)
- **GPU** (red)
- **FPGA** (blue)
- **SSE-SIMD** (green)

CUSHAW2 · CUSHAW · CUSHAW2-GPU · SARUMAN · BFAST · BFAST-Olson · BFAST-Yus · RazerS · RazerS 3 · SHRiMP · SHRiMP2 · BWA-Waidyasooriya · BWA-W · BWA · BWA-SW · BWA-MEM · BWA-MEM-FPGA · mrFAST · mrsFAST · +FastHASH · CloudBurst · ProbeMatch · WHAM · Bowtie · Bowtie2 · FHAST · PASS · PASS-bis · Slider · SliderII · SOCS · MAQ · SeqMap · ZOOM · ZOOM Lite · RMAP · SOAP3-FPGA · SOAP · SOAP2 · SOAP3 · SOAP3-dp · GMAP · GSNAP · Exonerate · MUMmer · MUMmer3 · Blat · SSAHA · BLAST · BLASTZ · BLAST+

Years: 1990 ... 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016

**Years**

# Filters, Alignment, and Mappers

| Year | Purpose | Architecture | Platform | alignments#/1sec |
|------|---------|-------------|----------|------------------|
| 2015 | Filter | Shifted Hamming Distance | Intel SSE | **3x faster** 3583 |
| | Alignment | Myers's bit-vector [45] | Intel SSE | 409 |
| | Alignment | Smith-Waterman [40] | Intel SSE | 38 |
| | Mapper | | | 16 |
| 2014 | M | | | 13 |
| | A | | | 1 |
| | | | | 60 |
| 2013 | | | | 17 |
| | | | | 4 |
| | | | | 15 |
| | Mapper | BWT-FM | FPGA(Virtex6) | 1092 |
| | Mapper | BWT-FM | GPU | 17 |
| | Mapper | Hash-Based (BFAST) | FPGA(Virtex6) | 35 |
| 2012 | Alignment | Smith-Waterman | FPGA(Virtex4) | 131 |
| | | | GPU | |
| | | | Cell BE | 16 |
| | | | CPU | 41 |

ideal filter → fast & accurate to compensate the computation overhead

Alignment performance for various state-of-the-art mappers and filters for 100 bp reads with at most 2% mismatch rate.