## Accelerating Genome Analysis A Primer on an Ongoing Journey

Onur Mutlu <u>omutlu@gmail.com</u> <u>https://people.inf.ethz.ch/omutlu</u>

> March 8, 2018 ETH HAML Seminar







#### Onur Mutlu

- Professor @ ETH Zurich CS, since September 2015 (officially May 2016)
- Strecker Professor @ Carnegie Mellon University ECE/CS, 2009-2016, 2016-...
- PhD from UT-Austin, worked at Google, VMware, Microsoft Research, Intel, AMD
- https://people.inf.ethz.ch/omutlu/
- omutlu@gmail.com (Best way to reach me)
- Office hours: By appointment (email me)

#### Research and Teaching in:

- Computer architecture, computer systems, bioinformatics
- Memory and storage systems
- Hardware security
- Fault tolerance
- Hardware/software cooperation
- ...

#### Current Research Focus Areas

**Research Focus:** Computer architecture, HW/SW, bioinformatics

- Memory and storage (DRAM, flash, emerging), interconnects
- Heterogeneous & parallel systems, GPUs, systems for data analytics
- System/architecture interaction, new execution models, new interfaces
- Hardware security, energy efficiency, fault tolerance, performance
- Genome sequence analysis & assembly algorithms and architectures
- Biologically inspired systems & system design for bio/medicine



**Graphics and Vision Processing** 

#### Four Key Current Directions

Fundamentally Secure/Reliable/Safe Architectures

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

Fundamentally Low-Latency Architectures

Architectures for Genomics, Medicine, Health

#### 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

### 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

#### What Is a Genome Made Of?



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

### The Central Dogma of Molecular Biology



### DNA Under Electron Microscope



### 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







Ion Torrent Proton



Complete Genomics





Illumina NovaSeq 6000

#### **Oxford Nanopore GridION**

... 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.



### The Genomic Era (continued)



14 http://www.economist.com/news/21631808-so-much-genetic-data-so-many-uses-genes-unzipped

### High-Throughput Sequencing (HTS)



### High-Throughput Sequencing (HTS)



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



#### Multiple sequence alignment

PHDHtm				MMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMM	
16082665	Т	acid	10	MASDRKSEGFQSGAGLIRYFEEEEIKGPALDPKLVVYMGIAVAIIVEIAKIFWPP	(55)
13541150	Т	vola	10	MASDKKSEGFQSGAGLIRYFEEEE	(55)
RFAC01077	F	acid	13	- <i>mtsmakdnonen</i> FQSGAG <mark>LIR</mark> YFNEE <mark>E</mark> IKG <mark>PAI</mark> D <mark>EKLIIYIGIAMGVIVELA</mark> KVFWP <mark>V</mark>	(58)
15791336	H	NRC1	10	MSSGQNSGGLMSSAGLVRYFDSEDSNALOIDPRSVVAVGAFFGLVVLLAQFFA	(53)
RAG22196	A	fulg	14	MAKAPKGKAKTPPLMSSAGIMRYFEE-EKTQIKVSPKTILAAGIVTGVLIIILNAYYGLWP-	(68)
RP001000	P	abys	9	MAKEKTTLPPTGAGLMRFFDE-DTRAIKITPKGAVALTLILIIFEIIL	(56)
RPH01741	P	hori	9	MAKEKTTLPPTGAG <mark>LMR</mark> FFD <mark>P-DTRAIKITPKGAIALVLILIIFEILLHVV</mark> GPR <mark>I</mark> FG	(56)
AE000914	М	ther	10	MAKKDKKTLPPSGAGLVRYFEE-ETKGFKLTPEQVVVMSIILAVFCLVLRFSG	(52)
RMJ09857	М	jann	9	MSKREETGLATSAGLIRYMDE-TFSKIRVKPEHVIGVTVAFVIIEAILTYGRFL	(53)
15920503	S	toko	13	-MPSSKKKKETVPLASMAGLIRYYEE-ENEKIKISPKLLIIISIIMVAGVIVASILIPPP	(58)
AE006662	S	solf	11	-MPSSKKKKETVPVMSMAGLIRYYEE-PNEKVKISPKIVIGASLALTIIVIVITKLF	(55)
RPK02491	P	aero	12	MARREKYEGINPFVAAGLIKFSEEGELEKIKLTPRAAVVISLAIIGLLIAINLLLPPL	(58)
RAP00437	A	pern	13	- <i>msv</i> rrrerratPVTAAG <mark>LL</mark> S <mark>FY</mark> EE-YEGK <mark>IKI</mark> SPT <mark>IVVGA</mark> AILVSAVVAAAHIFLP <mark>AV</mark> P-	(59)
5803165	H	sapi	49	SAGTGGMWRFYTE-DSPGLKVGPVLVMSLLFIASVFMLHIWGKYTRS	(96)
13324684	М	musc	49	SAGTGGMWRFYTE-DSPGLKVGPVLVMSLLFIAAVFMLHIWGKYTRS	(96)
6002114	D	mela	53	GAGTGGMWRFYTD-DSPGIKVGPVLVMSLLFIASVFMLHIWGKYNRS	(100)
14574310	C	eleg	32	GGNNGG <mark>LWR</mark> FYT <mark>E-D</mark> STG <mark>LKI</mark> GPVPVLVMSLVFIASVFVLHIWGK <mark>FT</mark> RS	(81)
10697176	Y	lipo	41	GGSSSTMLKLYTD-ESOGLKVDPVVVVLSLGFIFSVVALHILAKVSTK	(91)
6320857	S	cere	40	GGSSSS <mark>ILK</mark> LYTD-PANGFRVDSLVVLFLSVGFIFSVIALHLLTKFTHI	(88)
6320932	s	cere	33	TNSNNSILKIYSD-EATGLRVDPLVVLFLAVGFIFSVVALHVISKVAGK	(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%

Metagenomics, genome assembly, de novo sequencing Question 2: Given a bunch of short sequences, Can you identify the approximate species cluster for genomically unknown organisms (bacteria)?







### The Read Mapping Bottleneck



Illumina HiSeq4000

2 Million bases/minute

300<sup>Million</sup> bases/minute



#### Read Mapping Execution Time Breakdown



### Read Mapping

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



### 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.



Why Is Read Alignment Slow?

 Quadratic-time dynamicprogramming algorithm(s)

 Data dependencies limit the computation parallelism

 Entire matrix computed even though strings may be dissimilar.



#### **Read Alignment**

### 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

### 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]



### 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]



### Advantages of Hash Table Based Mappers

- + Guaranteed to find *all* mappings  $\rightarrow$  sensitive
- + Can tolerate up to errors



http://mrfast.sourceforge.net/

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

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

Alkan+, <u>"Personalized copy number and segmental duplication</u> <u>maps using next-generation sequencing</u>", 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

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



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

#### Adjacency Filtering (AF)



 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



# 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
  - IM simulated reads of human (180 base pairs)
- Evaluation system
  - Intel Core i7 Sandy Bridge machine
  - 16 GB of main memory

### FastHASH Speedup



# 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

Xin et al. BMC Genomics 2013, **14**(Suppl 1):S13 http://www.biomedcentral.com/1471-2164/14/S1/S13



**Open Access** 

#### PROCEEDINGS

### Accelerating read mapping with FastHASH

Hongyi Xin<sup>1</sup>, Donghyuk Lee<sup>1</sup>, Farhad Hormozdiari<sup>2</sup>, Samihan Yedkar<sup>1</sup>, Onur Mutlu<sup>1\*</sup>, Can Alkan<sup>3\*</sup>

*From* The Eleventh Asia Pacific Bioinformatics Conference (APBC 2013) Vancouver, Canada. 21-24 January 2013

### 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

Bioinformatics, 31(10), 2015, 1553–1560 doi: 10.1093/bioinformatics/btu856 Advance Access Publication Date: 10 January 2015 Original Paper

OXFORD

Sequence analysis

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

Hongyi Xin<sup>1,\*</sup>, John Greth<sup>2</sup>, John Emmons<sup>2</sup>, Gennady Pekhimenko<sup>1</sup>, Carl Kingsford<sup>3</sup>, Can Alkan<sup>4,\*</sup> and Onur Mutlu<sup>2,\*</sup>

Xin+, <u>"Shifted Hamming Distance: A Fast and Accurate SIMD-friendly Filter</u> 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

#### Insight: Shifting a String Helps Similarity Search

3 matches 5 mismatches



#### Insight: Shifting a String Helps Similarity Search

7 matches 1 mismatches



# Highly Parallel Matrix Computation



# Key Idea of SHD Filtering

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

#### Amend random zeros: $101 \rightarrow 111 \& 1001 \rightarrow 1111$

Generate 2E+1 masks

Query	:GAGAGAGATATTTAGTGTTGCAGCACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGGGAACATTGTTGGGCCGGA
Reference	<sup>;</sup> GAGAGAGATAGTTAGTGTTGCAGCCACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGAGACATTGTTGGGCCGG

--- Masks after amendment ---

Hamming Mask	: <mark>0000000000100000000000000000000000000</mark>
1-Deletion Mask	:1111111111111111111111111111111111111
2-Deletion Mask	:0000000011111111111111111111111110001111
3-Deletion Mask	:11111111111111111111111111111111111111
1-Insertion Mask	:111111111111111111111111111110001111111
2-Insertion Mask	:00000011111111111111111111000111111111
3-Insertion Mask	:11111111111111111100011111111111111111

Needleman-Wunsch Alignment 

# Alignment vs. Pre-alignment (Filtering)



 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! No data dependencies!

## New Bottleneck: Filtering (Pre-Alignment)

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

 $\rightarrow$ 

 $\rightarrow$ 

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



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

### Ideal Filtering Algorithm



# Alignment vs. Pre-alignment (Filtering)



 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! No data dependencies!

#### Our Solution: GateKeeper



#### GateKeeper Walkthrough

#### AND all masks, ACCEPT iff number of $1' \leq$ Threshold

Amend random zeros:  $101 \rightarrow 111 \& 1001 \rightarrow 1111$ 

Generate 2E+1 masks

Query	$: {\tt GAGAGAGATATTTAGTGTTGCAGCACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGGGAACATTGTTGGGCCGGA$
Reference	<sup>;</sup> GAGAGAGATAGTTAGTGTTGCAGCCACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGAGACATTGTTGGGCCGG

--- Masks after amendment ---

Hamming Mask	: <mark>0000000000100000000000000000000000000</mark>
1-Deletion Mask	:1111111111111111111111111111111111111
2-Deletion Mask	:00000001111111111111111111111111111111
3-Deletion Mask	:11111111111111111111111111111111111111
1-Insertion Mask	:111111111111111111111111111110001111111
2-Insertion Mask	:00000011111111111111111111000111111111
3-Insertion Mask	:11111111111111111100011111111111111111

Needleman-Wunsch Alignment GAGAGAGATATTTAGTGTTGCAG-CACTACAACACAAAAGAGGACCAACTTACGTGTCTAAAAGGGGGGAACATTGTTGGGCCGG

# GateKeeper Walkthrough (cont'd)



#### 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



#### SAFARI

#### GateKeeper vs. SHD

#### GateKeeper

- FPGA (Xilinx VC709)
- Multi-core (parallel)
- Examines a single mapping @ 125 MHz
- Limited to PCIe Gen3(4x) transfer rate (128 bits @ 250MHz)
- Amending requires:
  - (2E+1) 5-input LUT.

#### SHD

- Intel SIMD
- Single-core (sequential)
- Examines a single mapping @ ~2MHz
- Limited to a read length of 128 bp (SSE register size)
- Amending requires:
  - 4(2E+1) bitwise OR.
  - 4(2E+1) packed shuffle.
  - □ 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 <u>https://github.com/BilkentCompGen/GateKeeper</u>

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

#### Sequence analysis

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

Mohammed Alser<sup>1,\*</sup>, Hasan Hassan<sup>2</sup>, Hongyi Xin<sup>3</sup>, Oğuz Ergin<sup>2</sup>, Onur Mutlu<sup>4,\*</sup>, and Can Alkan<sup>1,\*</sup>

#### SAFARI

# Next Talk: MAGNET (AACBB 2018)

- 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.
#### 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

### 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)



Hash Table



#### **Reference Genome**



Х



### 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



### **GRIM-Filter: Bins**

SAFAR

We partition the genome into large sequences (bins).



### **GRIM-Filter: Bitvectors**



### **GRIM-Filter: Bitvectors**



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



### **GRIM-Filter: Checking a Bin**

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



### **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**



### **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 SAFARI

### **3D-Stacked Memory**



significant

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

### **3D-Stacked Memory**

floading



 3D-Stacked DF
 bandwidth as
 Logic Layer e computation f
 Embed GRIMappropriately

SAFARI

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



### **3D-Stacked Memory**

# Micron's HMC

# Micron has working demonstration components

http://images.anandtech.com/doci/9266/HBMCar\_678x452.jpg







floading

### **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

### **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

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

### 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

### **GRIM-Filter Performance**



**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**



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** 

#### 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 <u>BMC Genomics</u>*, 2018. *Proceedings of the <u>16th Asia Pacific Bioinformatics Conference</u> (APBC), Yokohama, Japan, January 2018. <u>arxiv.org Version (pdf)</u>* 

#### GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies

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

#### 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

### 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  $\rightarrow$  Less limit on read length  $\rightarrow$  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 Technology



an emerging and a promising ncing technology read length  $\rightarrow$  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



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

step.

SAFARI

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

#### Nanopore Sequencing Technology and Tools for Genome Assembly: Computational Analysis of the Current State, Bottlenecks, and Future Directions

#### Damla Senol Cali<sup>1,\*</sup>, Jeremie Kim<sup>1,3</sup>, Saugata Ghose<sup>1</sup>, Can Alkan<sup>2\*</sup> and Onur Mutlu<sup>3,1\*</sup>

<sup>1</sup>Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA, USA <sup>2</sup>Department of Computer Engineering, Bilkent University, Bilkent, Ankara, Turkey <sup>3</sup>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

#### 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

#### Acknowledgments

- Prof. Can Alkan, Bilkent University
- Many students
  - Mohammed Alser, Damla Senol Cali, Jeremie Kim
  - Hasan Hassan
  - Hongyi Xin
  - ...
- All papers, source code, and more are at:
  - https://people.inf.ethz.ch/omutlu/projects.htm

## Accelerating Genome Analysis A Primer on an Ongoing Journey

Onur Mutlu <u>omutlu@gmail.com</u> <u>https://people.inf.ethz.ch/omutlu</u>

> March 8, 2018 ETH HAML Seminar




# High-Throughput 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.

## 3- Highly Accurate Filtering Algorithm (cont'd)

#### MAGNET

- Check for substitutions.
- ✓ The longest identical subsequence  $\geq [(m E)/(E + 1)]$ .
  - Extraction & Encapsulation (divide-and-Conquer fashion).



#### **SAFARI**

Substitution

### 3- Highly Accurate Filtering Algorithm (cont'd)

#### MAGNET

- Check for substitutions.
- The longest identical subsequence  $\geq [(m E)/(E + 1)]$ .
- Extraction & Encapsulation (divide-and-Conquer fashion).

Substitution



SA Now divide the problem into two subproblems and repeat

## 3- Highly Accurate Filtering Algorithm (cont'd)

#### MAGNET

- Check for substitutions.
- The longest identical subsequence  $\geq [(m E)/(E + 1)]$ .
- Extraction & Encapsulation (divide-and-Conquer fashion).



SA Counting the encapsulation bits reveals the number of edits

### MAGNET Accelerator

