Hardware Accelerated Next Generation Sequencing Analysis

Arun Subramaniyan

Staff Bioinformatics Scientist, Illumina



Workshop on Secure and Adaptive Computing from Genomics to Healthcare The 33rd Annual Symposium on Field-Programmable Custom Computing Machines May 4th 2025

*All views are own and do not reflect those of employer

Sequencing data is increasing in volume and diversity





Illumina Genome Analyzer, 2005 1 Gbases per day

Illumina NovaSeq X, 2025 4 Tbases per day



Sequencing data is increasing in volume and diversity



Illumina Genome Analyzer, 2005 1 Gbases per day



Illumina NovaSeq X, 2025 4 Tbases per day



illumina®

Read length: 100-350bp Per base inaccuracy: 0.1%



Oxford NANOPORE Technologies

> Read length: 1kb-1Mbp Per base inaccuracy: 1-15%



1000x increase in sequencing fragment length 10 - 100x increase in sequencing error rate

Credits: Illumina, DataBase Center for Life Science (DBCLS), https://www.ecseq.com/support/ngs/do_you_have_two_colors_or_four_colors_in_lllumina

GPUs and custom hardware are being used for sequence analysis









GPU



Microsoft

Stanford University

illumina[®] DRAGEN Bio-IT Platform



FPGA

Modern Sequence Analysis Pipelines



Sequencing Basecalling **Read Alignment** Variant Calling A T C G T G C A G T T T Reference genome (R) Ref. (R) **CGTGAAG** Squiggle CGTGAAG (a) Reference-GAAGTTT $\mathcal{N}\mathcal{N}$ GAAGTTT Aligned Guided Reads ATCGTG Reads Read (Q) h.fast5 Assembly (~millions-billions) TGAAGT Genome sizes: **₽.**bcl hastq .bam vcf ~Mega -- Gigabases Short Reads BWA-MEM2 (~30-40 %) Bustard GATK Haplotype Caller, Platypus (~40%) Software tools: Long Reads Clair Bonito Minimap2 **Error Correction and Polishing** Overlap, Layout, Consensus GA GT CTG Overlap **G**+**T**+**C** Region GT CTG TC ***G+T+C** G-A GAGTCTG (b) De-Novo GAA TG TC GTCTGTC Consensus Assembly GAATGTC **GA GTC TGTC GAGTCTGTC** GTC TGTC Reads Long Reads Consensus GAA TGTC GAGTCTGTC In .fasta Assembly h.bam Layout Flye, Minimap2 Racon, Nanopolish

Modern Sequence Analysis Pipelines

Sequencing Basecalling **Read Alignment** Variant Calling A T C G T G G A G T T T Reference genome (R) Ref. (R) **CGTGAAG** Squiggle CGTGAAG (a) Reference-GAAGTTT $\mathcal{N}\mathcal{N}$ GAAGTTT Aligned Guided Reads ATCGTG Reads Read (Q) h.fast5 Assembly (~millions-billions) TGAAGT Genome sizes: bcl. h.fastq .bam vcf ، ~Mega -- Gigabases BWA-MEM2 (~30-40 %) Short Reads Bustard GATK Haplotype Caller, Platypus (~40%) Software tools: Long Reads Clair Bonito Minimap2 **Error Correction and Polishing** Overlap, Layout, Consensus GA GT CTG ₲**₊**₸₊₢ Region GT CTG TC ***G+T+C** (G)→(A GAGTCTG (b) De-Novo GAA TG TC GTCTGTC Consensus Assembly GAATGTC **GA GTC TGTC** GAGTCTGTC GTC TGTC Reads Long Reads Consensus GAA TGTC GAGTCTGTC _____.fasta Assembly .bam Layout Flye, Minimap2 Racon, Nanopolish (c) Abundance Estimation **Metagenomics** Classification Bacterial + Viral + Human Reference genomes - Species 1 %Viral %Bacterial Species 2 Long Reads Species 3 Reads %Human **Read Alignment** bam. **Phylogenetic Tree** Centrifuge, Minimap2

Modern Sequence Analysis Pipelines

⁴



https://genomicsbench.eecs.umich.edu

- Computationally intensive kernels drawn from well maintained software tools
- Covers the major steps of modern sequence analysis pipelines
- Includes both short and long read analysis algorithms
- Small/large input datasets

Benchmark	Input Datatype	Applications	Chosen Tool	% Time Spent in Tool	Parallelism Motif
				(single-thread)	
fmi	Short reads	Read Alignment	BWA-MEM2	38%	Tree Traversal
		Metagenomics Classification			
bsw	Short reads	Read Alignment	BWA-MEM2	31%	Dynamic Programming
		De-Novo Assembly			
dbg	Short reads	Variant Calling	Platypus	65%	Graph Construction
		De-Novo Assembly			Hash Table
phmm	Short reads	Variant Calling	GATK Haplotype Caller	70%	Dynamic Programming
_		Error Correction			
chain	Long reads	De-Novo Assembly	Minimap2	47.4 %	Dynamic Programming (1D)
		Read Alignment			
spoa	Long reads	Error Correction	Racon	75 %	Dynamic Programming
					Graph Construction
abea	Long reads	Basecalling	Nanopolish	71.4%	Dynamic Programming
	-	Variant Calling	_		
grm	NA	Population Genomics	PLINK2	92.8 %	Dense Matrix Multiplication
nn-base	Long reads	Basecalling	Bonito	95 %	FP Matrix Multiplication
nn-variant	Long reads	Variant Calling	Clair	57.2 %	FP Matrix Multiplication

Benchmark	Input Datatype	Applications	Chosen Tool	% Time Spent in Tool	Parallelism Motif	
				(single-thread)		
fmi	Short reads	Read Alignment	BWA-MEM2	38%	Tree Traversal	
		Metagenomics Classification				
bsw	Short reads	Read Alignment	BWA-MEM2	31%	Dynamic Programming	
		De-Novo Assembly				
dbg	Short reads	Variant Calling	Platypus	65%	Graph Construction	
		De-Novo Assembly			Hash Table	
phmm	Short reads	Variant Calling	GATK Haplotype Caller	70%	Dynamic Programming	
-		Error Correction				
chain	Long reads	De-Novo Assembly	Minimap2	47.4 %	Dynamic Programming (1D)	
		Read Alignment				
spoa	Long reads	Error Correction	Racon	75 %	Dynamic Programming	
-					Graph Construction	
abea	Long reads	Basecalling	Nanopolish	71.4%	Dynamic Programming	
		Variant Calling				
grm	NA	Population Genomics	PLINK2	92.8 %	Dense Matrix Multiplication	
nn-base	Long reads	Basecalling	Bonito	95 %	FP Matrix Multiplication	
nn-variant	Long reads	Variant Calling	Clair	57.2 %	FP Matrix Multiplication	

- 5 dynamic programming kernels in GenomicsBench
- Contribute 30-75% execution time in their respective software tools

- 5 dynamic programming kernels in GenomicsBench
- Contribute 30-75% execution time in their respective software tools
 - Chaining
 - Banded Smith-Waterman
 - Pairwise Hidden Markov Model
 - Partial Order Alignment
 - Adaptive Banded Event Alignment

- 5 dynamic programming kernels in GenomicsBench
- Contribute 30-75% execution time in their respective software tools
 - Chaining
 - Banded Smith-Waterman
 - Pairwise Hidden Markov Model
 - Partial Order Alignment
 - Adaptive Banded Event Alignment

Diverse compute requirements

Data dependencies: 1-D, 2-D

Inputs: Sequence, Graph

Computation type: Integer, Floating point

Memory-Intensive Index Lookup Benchmarks: e.g., FM-index search

Reference : G T A T A T C \$



Memory-Intensive Index Lookup Benchmarks: e.g., FM-index search

Reference : G T A T A T C \$



Irregular memory accesses in FM-index search

Reference : G T A T A T C \$



Read: TAT

Machine learning algorithms are also gaining traction



nn-base (Basecalling)

Machine learning algorithms are also gaining traction





nn-base (Basecalling)

nn-variant (Variant calling)

Abundant data parallelism in sequence analysis applications

Kernel	Parallelism granularity	Source of parallelism
FM-index search (fmi)	Read	# OCC Table Lookups
Banded Smith-Waterman (bsw)	Seed	# Cell Updates
Pairwise Hidden Markov Model (phmm)	Genome Region	# Cell Updates
Chaining (chain)	Read	# Input Anchors
Partial Order Alignment (poa)	Read Chunk	# Cell Updates

Abundant data parallelism available

~billions of independent reads ~millions of independent genome regions

Case Study – Accelerating FM-index Search

FMD-index seeding trades off memory bandwidth for space

Single Character Lookup



FMD-index seeding trades off memory bandwidth for space

Single Character Lookup



FMD-index seeding trades off memory bandwidth for space



BWA-MEM:

102 kB / read









Also, present is a 2nd level index table for 15-mers having >= 256 occurrences

Instead trade-off memory capacity for reducing memory bandwidth



https://github.com/bwa-mem2/bwa-mem2/tree/ert

"Accelerated seeding for genome sequence alignment with enumerated radix trees." In 2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA), pp. 388-401. IEEE, 2021.

• Task: calculate the scores in the entire table.

- Task: calculate the scores in the entire table.
- Initialization: the scores in first row and first column.

- Task: calculate the scores in the entire table.
- Initialization: the scores in first row and first column.
- Objective Function: calculate the score of a cell based on its upper, left and diagonal neighbors.



- Task: calculate the scores in the entire table.
- Initialization: the scores in first row and first column.
- Objective Function: calculate the score of a cell based on its upper, left and diagonal neighbors.



- Task: calculate the scores in the entire table.
- Initialization: the scores in first row and first column.
- Objective Function: calculate the score of a cell based on its upper, left and diagonal neighbors.



1D dynamic programming



Chaining

L. Guo, J. Lau, Z. Ruan, P. Wei, and J. Cong, "Hardware Acceleration of Long Read Pairwise Overlapping in Genome Sequencing: A Race Between FPGA and GPU," in 2019 IEEE 27th International Symposium On Field-Programmable Custom Computing Machines (FCCM), April 2019.

1D dynamic programming

2D dynamic programming



Chaining

Banded Smith-Waterman

L. Guo, J. Lau, Z. Ruan, P. Wei, and J. Cong, "Hardware Acceleration of Long Read Pairwise Overlapping in Genome Sequencing: A Race Between FPGA and GPU," in 2019 IEEE 27th International Symposium On Field-Programmable Custom Computing Machines (FCCM), April 2019.

Dynamic programming benchmarks: Floating point 2D

Pairwise Hidden Markov Model



j i ▼
Dynamic programming benchmarks: Floating point 2D

Pairwise Hidden Markov Model





j i ▼ **Dynamic programming benchmarks: Floating point 2D**

Pairwise Hidden Markov Model



Floating Point

Dynamic programming benchmarks: Sequence to Graph

Partial Order Alignment



Sequence to Graph Alignment

Dynamic programming benchmarks: Sequence to Graph

Partial Order Alignment



Sequence to Graph Alignment

Dynamic programming benchmarks: Sequence to Graph

Partial Order Alignment



Sequence to Graph Alignment

Case Study – Accelerating Pairwise Hidden Markov Model (pairHMM)

Xiao Wu, UMich

Pruning algorithm for pairHMM with same output guarantees

- Compute upper bound probability by replacing floating-point multiplication with log-domain addition and rounding up (high_{res})
- Prune matrix and identify promising paths
- Compute lower bound probability with single-precision floating-point only for small unpruned area (low_{res})
- Recompute matrix in floating point if bound checks fail



Xiao Wu, UMich

Pruning algorithm for pairHMM with same output guarantees

- Compute upper bound probability by replacing floating-point ٠ multiplication with log-domain addition and rounding up (high_{res})
- Prune matrix and identify promising paths ٠
- Compute lower bound probability with single-precision floating-point ٠ only for small unpruned area (low_{res})
- Recompute matrix in floating point if bound checks fail

Original algorithm





Pruned version



Xiao Wu, UMich

Haplotype

Т

alignment

Т

Т Т

Squares processed

by pruning

machine

ΑΤ

Result = sum of last row **Dominated by this**

G

Т

С

Squares need to be

point machine

processed by floating

Pruning algorithm for pairHMM with same output guarantees

- Compute upper bound probability by replacing floating-point multiplication with log-domain addition and rounding up (high_{res})
- Prune matrix and identify promising paths
- Compute lower bound probability with single-precision floating-point only for small unpruned area (low_{res})
- Recompute matrix in floating point if bound checks fail



Results:

43x fewer cells computed in precise floating point

8.3x higher throughput (GCUPS) than floating-point ASIC of the same area

Only need to

compute this

Read

accurately A

A T

G

С

Α



Case Study – Generalized Dynamic Programming Acceleration



ISA Extension





ISA Extension



https://www.nvidia.com/en-us/data-center/technologies/hopper-architecture/



ISA Extension





ISA Extension





Weights used from BWA. Data: <u>HG002 (NA24385)</u> paired-end protocol using Illumina Sequencers.

Diversity of dynamic programming in genomic kernels

Kernel	Application	Dimension and Size	Dependency	Data Type
bsw	Read Alignment	2D~120×60	Last 2 Wave-fronts	Int 8/16
pairHMM	Variant Calling	$2D \sim 100 \times 60$	Last 2 Wave-fronts	Floating Point
роа	Error Correction	$2D \sim 1000 \times 500$	Graph Structure	Int 32
chain	Read Alignment	1D ~20000	Last N (\sim 25) Anchors	Int 32
-	•			

Different objective functions computed in each cell

GenDP: Support for intra-cell objective functions



(a) Reduction Data-Flow in PairHMM



(b) Reduction Data-Flow in BSW





Compute unit – 2-level reduction tree

"GenDP: A framework of dynamic programming acceleration for genome sequencing analysis." In *Proceedings of the 50th Annual International Symposium on Computer Architecture*, pp. 1-15. 2023. 25 Communications of the ACM, Volume 68, Issue 05 Pages 81 - 90 https://doi.org/10.1145/3712168





Integer and floating point PE array





point PE array Local depe



Local dependency - 1-Dimension systolic PE array with FIFO









Software managed scratchpad memory for long range dependencies

GenDP performance

• Metrics: Throughput/Area – Million Cell Updates per Second/*mm*2 (MCUPS/*mm*2)

GenDP performance

- Metrics: Throughput/Area Million Cell Updates per Second/mm2 (MCUPS/mm2)
- GenDP has 2.8x slowdown when compared to custom accelerators



Conclusion

- Genomic sequence data is rapidly increasing in volume and diversity
- Modern sequence analysis pipelines present unique computation challenges that multi-core CPUs and GPUs have struggled to meet
- The GenomicsBench benchmark suite is an initial attempt to capture this diversity with the goal of enabling the development of custom hardware architectures
- GenomicsBench has already helped to develop several hardware accelerators
 - Enumerated Radix Tree accelerator for memory-intensive exact match search
 - GenDP for dynamic programming acceleration
 - Pruning-based accelerator for PairHMM

Acknowledgement

University of Michigan

<u>Reetuparna Das (Advisor)</u> <u>Satish Narayanasamy</u> <u>David Blaauw</u> Jack Wadden Kush Goliya Xiao Wu Nathan Ozog

Intel Labs

Sanchit Misra Md. Vasimuddin Somnath Paul

illumina®

Microsoft Research/UC Berkeley AMPLab

PRECISION HEALTH UNIVERSITY OF MICHIGAN Applications Driving Architectures



<u>Bill Bolosky</u> Ravi Pandya Matei Zaharia David Patterson



- **DPAx**: programmable dynamic programming (DP) accelerator.
- **DPMap**: map the objective function of DP algorithm to DPAx accelerator.

- **DPAx**: programmable dynamic programming (DP) accelerator.
- **DPMap**: map the objective function of DP algorithm to DPAx accelerator.



- **DPAx**: programmable dynamic programming (DP) accelerator.
- **DPMap**: map the objective function of DP algorithm to DPAx accelerator.



- **DPAx**: programmable dynamic programming (DP) accelerator.
- **DPMap**: map the objective function of DP algorithm to DPAx accelerator.



• **DPMap**: map the objective function of DP algorithm to programmable compute units in DPAx.

• **DPMap**: map the objective function of DP algorithm to programmable compute units in DPAx.



• **DPMap**: map the objective function of DP algorithm to programmable compute units in DPAx.



- **DPMap**: map the objective function of DP algorithm to programmable compute units in DPAx.
 - **Partitioning**: Break the data-flow graph with 4-input ALU and Multiplier
 - **Seeding**: Look for vertices that could be mapped to the 2nd level
 - **Refinement**: Break the single-strand structure



- **DPMap**: map the objective function of DP algorithm to programmable compute units in DPAx.
 - **Partitioning**: Break the data-flow graph with 4-input ALU and Multiplier
 - **Seeding**: Look for vertices that could be mapped to the 2nd level
 - **Refinement**: Break the single-strand structure





GenDP performance

• Metrics: Throughput/Area – Million Cell Updates per Second/*mm*2 (MCUPS/*mm*2)
GenDP performance

- Metrics: Throughput/Area Million Cell Updates per Second/mm2 (MCUPS/mm2)
- GenDP achieves 157.8× throughput/mm² over GPU



GenDP performance

- Metrics: Throughput/Area Million Cell Updates per Second/mm2 (MCUPS/mm2)
- GenDP achieves 157.8× throughput/mm² over GPU
- GenDP has 2.8x slowdown when compared to custom accelerators

