



A Co-processor Provides Increased Sensitivity in Whole Genome Alignments with High Speedup

Yatish Turakhia*, <u>Sneha D. Goenka</u>*, Prof. Gill Bejerano, Prof. William J. Dally

* Equal contribution

What are whole genome alignments (WGA)?



WGA is correspondence between genomes



Darwin-WGA: A Co-processor for Whole Genome Alignment

Why are whole genome alignments important?



2. WGA help predict functional elements



(Mayor et al. , 2000)



Classical alignment algorithm



Smith-Waterman Algorithm

Inputs

Target sequence (r) GTGTCACTA (L_r = 9)

Query sequence (q) GGCCAACTA (L_q = 9)

Scoring parameters

		A	С	G	A	
	A	2	-1	-1	-1	
N =	C	-1	2	-1	-1	
	G	-1	2	2	-1	
	T	-1	-1	-1	2	
Gap penalty = 1						

Smith-waterman equations

 $V(i, j) = \max \begin{cases} V(i-1, j-1) + W(r_i, q_j) \\ V(i-1, j) + gap \\ V(i, j-1) + gap \\ 0 \end{cases}$

		Target (r)									
		*	G	т	G	т	с	A	с	т	A
	*	0	0	0	0	0	0	0	0	0	0
	G	0	2€	- 1+	- 2 ≁	-1	0	0	0	0	0
	G	0	2	`1	_ 3 €	- 2	-1	0	0	0	0
(d)	с	0	1	`1	2	2	4	-3 🛫	–2≁	-1	0
Ň	с	0	Ó	0	1	1	4	3	`5է	-4∻	-3
e	A	0	0	0	0	þ	3	` 6≁	-5	4	` 6
2	A	0	0	0	0	0	2	5,	5	4	6
-	с	0	0	0	0	0	2	4	7	6,	5
	т	0	0	2 ,	1	2 ,	1	3	6	9,	8
	A	0	0	1	`1	1	`1	3	5	8	11

Alignment GTGTC-A-CTA G-G-CCAACTA



Smith Waterman algorithm intractable on whole genomes

- Smith Waterman algorithm time and space complexity $\sim O(L_r \cdot L_q)$
- Mammalian genomes ~ 10⁹-10¹⁰ base-pairs
- Use heuristics based approaches



Seed-Filter-Extend algorithm (LASTZ)



Seeding finds small matching local patterns



 Seeding finds local matching patterns of fixed length s

• Substrings of query of length s are compared to the target

Substrings start from position 0



Seeding finds small matching local patterns



•Seeding finds local matching patterns of fixed length s

• Substrings of query of length s are compared to the target

Substrings start from position 0



Filtering aligns ~100bp around seed hits



Seed hit

Filter (Ungapped)

•Calculate scores along the seed hit diagonal (match or mismatch)

Track maximum score

•Stop as soon as score falls below (max_score - x)

Does not consider indels



Filtering aligns ~100bp around seed hits



Seed hit

Filter (Ungapped)

•Calculate scores along the seed hit diagonal (match or mismatch)

Track maximum score

•Stop as soon as score falls below (max_score - x)

Does not consider indels



Filtering aligns ~100bp around seed hits



Filter (Ungapped)

•Calculate scores along the seed hit diagonal (match or mismatch)

Track maximum score

•Stop as soon as score falls below (max_score - x)

Does not consider indels



Extension uses Y-drop algorithm



Extend

- \$ \$ \$
- Darwin-WGA: A Co-processor for Whole Genome Alignment

•Start computing score along a row when

score > (max_score - y)

•Stop computing score along a row when

score < (max_score - y)</pre>

Extension provides the final alignments



Anchor

Extend

\$* {*`{*`

Darwin-WGA: A Co-processor for Whole Genome Alignment

Final alignment = right extension + left extension

Why is LASTZ less sensitive?



Darwin-WGA: A Co-processor for Whole Genome Alignment

Increasing indel frequency => increasing need for gapped filtering



Replacing ungapped filtering by gapped filtering slows down the software by 200x



Darwin-WGA algorithm overview



Seeding – D-SOFT

- Target bin and Query chunk determine diagonal band
- •Each seed hit falls in a single diagonal band
- •At most 1 seed hit per diagonal band is extended





Gapped Filtering – Banded Smith Waterman

- •Seed hit extended using banded Smith-Waterman
- Pre-determined band with no traceback
- If (max_score > threshold), the maximum score position (x_{max}) is the anchor or the starting point for alignment extension
- •Gaps considered => better alignments for species further apart





- •Tiled (tile size T, overlap O) implementation inspired by GACT in Darwin*
- •Origin of the next tile lies at the intersection of the current traceback path with the overlap



Outside tile overlap
 Inside tile overlap
 On tile overlap border

* Turakhia et al. , ASPLOS'18

•Extension along a direction continues until a tile is encountered with a non-positive maximum score







- •Y-drop implementation within each tile
- •Adaptive band with traceback
- •Reduces on-chip memory requirement compared to computing whole tile
- Reduces compute time





Workloads in LASTZ v/s Darwin-WGA





Whole genome alignment v/s Read alignment



Darwin-WGA: A Co-processor for Whole Genome Alignment

1. WGA requires aligns less similar sequences

- •Genomes may diverge considerably over evolutionary timescales and have low sequence similarity
- Read alignment deals with highly similar sequences (wellcharacterized sequencing error model)





2. WGA has longer alignments with large indels

- •Whole genome alignments can span millions of basepairs with large indels
- Read alignments span not more than tens of thousands of base-pairs with much shorter indels
- Previous hardware accelerators would require high on-chip memory





Darwin-WGA Framework



- Seeding done in software
- •Banded Smith-Waterman and GACT-X accelerated in hardware as bounded Dynamic Programming with Systolic Arrays



Hardware Acceleration



Darwin-WGA: A Co-processor for Whole Genome Alignment















\$* \$~ \$^



- Banded Smith-Waterman preset band and no traceback
- •GACT-X adaptive band with traceback

Evaluation Framework



Darwin-WGA: A Co-processor for Whole Genome Alignment

Experimental Setup

CPU (baseline)

- AWS c4.8xlarge instance
- 36 vCPUs (18 physical cores)
- LASTZ as software baseline
- Parasail to estimate isosensitive runtime
- \$1.59/hour

FPGA (Darwin-WGA)

- AWS f1.2xlarge instance
- 1 Xilinx Virtex Ultrascale+ FPGA (50 BSW and 2 GACT-X arrays with 32PEs)
- 8 vCPUs
- \$1.65/hour



ASIC

TSMC 40nm DC synthesis (not a chip prototype)

		Configuration	Area (mm ²)	Power (W)
BSW	Logic	64 x (64PE array)	16.6	25.6
GACT-X	Logic	12 x (64PE array)	4.2	6.72
	Traceback SRAM	12 x (64PE x 16KB/PE)	15.1	7.92
DRAM	DDR4-2400R	4 x 32GB	-	3.10
	ΤΟΤΑ	35.9	43.34	





Species and Genome Assembly



- dm6-droSim1 molecular distance comparable to human-monkey
- dm6-dp4 molecular distance comparable to human-chicken



Results



Darwin-WGA: A Co-processor for Whole Genome Alignment

Darwin-WGA finds genes that LASTZ does not

dp4





Indels (shown by arrows) around each seed hit – dropped by ungapped filtering (LASTZ) but retained by gapped filtering (Darwin-WGA)



Darwin-WGA Sensitivity Improvement Versus LASTZ

Species pair	Top-10 Alignment Chain Scores	Matching Base-pairs within Alignments	Number of Aligning Exons (protein-coding genes)	
dm6-droSim1	+0.03%	1.25x	+0.20%	
dm6-droYak2	+0.05%	1.41x	+0.09%	
dm6-dp4	+1.86%	1.42x	+0.41%	
ce11-cb4	+5.73%	3.12x	+2.70%	
	Represent <u>orthologous</u> <u>sequences (</u> derived from "speciation")	Represent <u>paralagous</u> <u>sequences (</u> derived from "duplication")	Represent <u>functionally</u> <u>relevant orthologous</u> <u>sequences</u> , under some selective pressure (at	
False positive rate	species)			



Runtime and Cost Comparison

Species pair	LASTZ runtime (sec)	Iso- sensitive s/w runtime (sec)	Darwin-WGA	runtime (sec)	Darwin-WGA Improvement		
			FPGA	ASIC	FPGA (Perf/\$)	ASIC (Perf/W)	
ce11-cb4	481	64,960	3,823	219	19.1x	1,478x	
dm6- droSim1	643	142,627	5,936	461	23.2x	1,547x	
dm6- droYak2	654	144,454	6,001	469	23.2x	1,539x	
dm6-dp4	557	125,700	4,987	404	24.3x	1,553x	
200x slowdown 200x slowdown							



GACT-X uses 3x less space and time as compared to GACT







- Darwin-WGA replaces ungapped filtering in LASTZ by Banded Smith-Waterman algorithm for higher sensitivity
 - up to 3x matching base-pairs
 - up to 5.7% more orthologs
 - up to 2.1% more exons
- Darwin-WGA outperforms iso-sensitive software
 - FPGA: 24x performance/\$ improvement
 - ASIC: 1,500x performance/Watt improvement
- GACT-X provides 3x improvement in speed and storage efficiency compared to GACT







