



TALCO: Tiling Genome Sequence Alignment using Convergence of Traceback Pointers

Sumit Walia, Cheng Ye, Arkid Bera, Dhruvi Lodhavia and Yatish Turakhia
University of California San Diego

Outline

- Emergence of **Long Genome Sequence Alignment (LGSA)**
- Current LGSA **algorithms, accelerators** and their **limitations**
- **TALCO**: A tiling technique based on convergence of traceback pointers for long genome sequence alignment
- **Key Contributions** and **Results**
- **Conclusion**



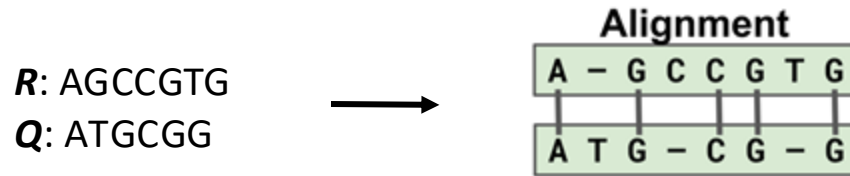
Outline

- Emergence of **Long Genome Sequence Alignment (LGSA)**
- Current LGSA algorithms, accelerators and their limitations
- **TALCO**: A tiling technique based on convergence of traceback pointers for long genome sequence alignment
- **Key Contributions and Results**
- **Conclusion**



Genome Sequence Alignment

- **Sequence alignment:** Given two sequences, Reference (R) and Query (Q), assign gaps (“-”) in R and Q to produce a valid alignment that maximizes the alignment score



Long Genome Sequence Alignment: Applications

- **Third-generation sequencing technologies** (produce reads of length 10 kb - 4 Mb), leads to major breakthroughs in recent past:

METHOD OF THE YEAR: LONG-READ SEQUENCING

To large-scale projects and individual labs, long-read sequencing has delivered new vistas and long wish lists for this technology's future. **By Vivien Marx**

To the delight of scientists across the life sciences, reads, which are the output of sequencing instruments, have been getting longer. Reads might be sequenced DNA or RNA and could one day routinely be entire genomes, transcriptsomes and epigenomes at high throughput and accuracy, and maybe even the amino acid sequences of proteins.

Academics have happy tales about how long-read technologies have empowered their genomic projects. A few companies have facilitated this journey, notably Pacific Biosciences (PacBio) and Oxford Nanopore Technologies (ONT). Of late, other firms presenting long-read approaches include Element Biosciences, Evident and MGI. Ultima Genomics and others have plans in this area.

Long reads have buoyed numerous findings in individual labs. In larger ventures, among the celebrated achievements are those in the Vertebrate Genomes Project (VGP) and the Telomere-to-Telomere Consortium (T2T). A set of papers and news features related to the T2T Consortium can be found as a Nature Collection online. During the T2T project, says Adam Phillippy, a researcher at the National Institutes of Health (NIH) National Human Genome Research Institute (NHGRI) who co-leads the T2T Consortium, the longest read he and his colleagues handled had one million base pairs. Long-read sequencing is being used by the Human Pangenome Reference Consortium (HPRC). The HPRC teams want to assemble the human genome at the T2T level of completion and capture a "better

spectrum of humanity in terms of how they represent allelic diversity," says University of California Santa Cruz researcher Karen Miga, who co-leads the T2T Consortium with Phillippy and is part of the HPRC.

Population-level data from diverse populations are needed, says Heidi Rehm, who, among other appointments, is the chief genomics officer at Massachusetts General Hospital's department of medicine and medical director of the clinical research sequencing platform at the Broad Institute of MIT and Harvard. She and her colleagues have found instances in which Black people received information about risk of a heart condition without sufficient evidence on genetic variants to support it. "Population data had been lacking about these variants, and such data are still limited," says Rehm.

naturemethods

Volume 20 | January 2023 | 6-11 | 6

Oxford Nanopore Technologies (ONT)



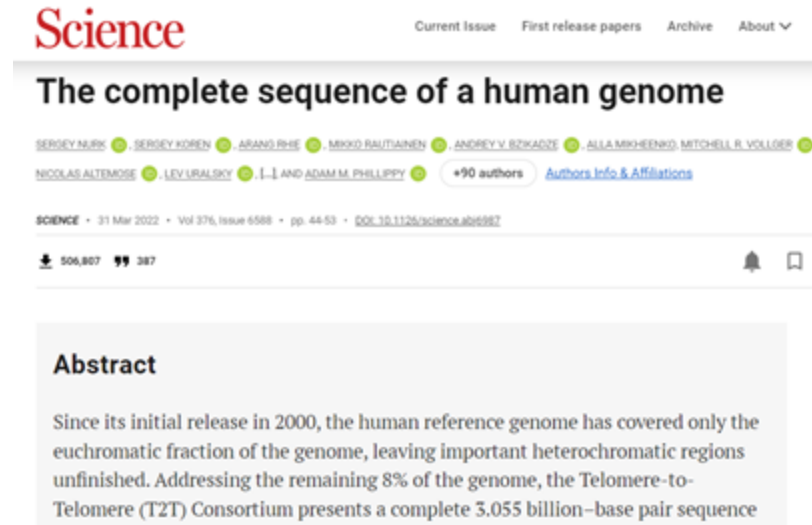
Pacific Biosciences (PacBio)



Marx, Vivien. "Method of the year: long-read sequencing." Nature Methods 20.1 (2023): 6-11.

Long Genome Sequence Alignment: Applications

- Third-generation sequencing technologies (produce reads of length 10 kb - 4 Mb), leads to major breakthroughs in recent past:
 - **Human Genome Assembly**



The screenshot shows the top portion of a Science journal article. The title is "The complete sequence of a human genome". Below the title, the authors are listed: SERGEY NURK, SERGEY KOREN, ARAND RHE, MIKKO RAUTIAINEN, ANDREY V. BZIKADZE, ALLA MIKHAILOV, MITCHELL B. VOLLOS, NICOLAS ALTEMOSE, LEV LURALSKY, I.-I. AND ADAM M. PHILLIPPY, and +90 authors. The article is from Science, 31 Mar 2022, Vol 376, Issue 6588, pp. 44-53, DOI 10.1126/science.aba6987. The abstract text is visible below the author list.

Science Current Issue First release papers Archive About

The complete sequence of a human genome

SERGEY NURK SERGEY KOREN ARAND RHE MIKKO RAUTIAINEN ANDREY V. BZIKADZE ALLA MIKHAILOV MITCHELL B. VOLLOS
NICOLAS ALTEMOSE LEV LURALSKY I.-I. AND ADAM M. PHILLIPPY +90 authors [Authors Info & Affiliations](#)

SCIENCE • 31 Mar 2022 • Vol 376, Issue 6588 • pp. 44-53 • DOI:10.1126/science.aba6987

506,807 387

Abstract

Since its initial release in 2000, the human reference genome has covered only the euchromatic fraction of the genome, leaving important heterochromatic regions unfinished. Addressing the remaining 8% of the genome, the Telomere-to-Telomere (T2T) Consortium presents a complete 3.055 billion-base pair sequence

Nurk, Sergey, et al. "The complete sequence of a human genome." Science 376.6588 (2022): 44-53.



Long Genome Sequence Alignment: Applications

- Third-generation sequencing technologies (produce reads of length 10 kb - 4 Mb), leads to major breakthroughs in recent past:
 - Human Genome Assembly
 - **Rapid genetic diagnosis**

Fastest DNA sequencing technique helps undiagnosed patients find answers in mere hours

A research effort led by Stanford scientists set the first **Guinness World Record for the fastest DNA sequencing technique**, which was used to sequence a human genome in just 5 hours and 2 minutes.

January 12, 2022 - By Hanae Armitage



Researchers were able to quickly determine that Matthew Kunzman's **heart failure** was the result of a **genetic condition** — a finding that cleared the way for him to be placed on a heart transplant list immediately.
Courtesy of Jenny Kunzman

Long Genome Sequence Alignment: Applications

- Third-generation sequencing technologies (produce reads of length 10 kb - 4 Mb), leads to major breakthroughs in recent past:
 - Human Genome Assembly
 - Rapid genetic diagnosis
 - **Characterize structural variations and complex regions**

ARTICLE

Structural Variation of Chromosomes in Autism Spectrum Disorder

Christian R. Marshall,¹ Abdul Noor,² John B. Vincent,² Anath C. Lionel,¹ Lars Feuk,¹ Jennifer Skaug,¹ Mary Shago,³ Rainald Moessner,¹ Dalila Pinto,¹ Yan Ren,¹ Bhooma Thiruvahindrapduram,¹ Andreas Fiebig,⁶ Stefan Schreiber,⁶ Jan Friedman,⁷ Cees E.J. Ketelaars,⁸ Yvonne J. Vos,⁸ Can Ficioglu,⁹ Susan Kirkpatrick,¹⁰ Rob Nicolson,¹¹ Leon Sloman,² Anne Summers,¹² Clare A. Gibbons,¹² Ahmad Teebi,⁴ David Chitayat,⁴ Rosanna Weksberg,⁴ Ann Thompson,¹³ Cathy Vardy,¹⁴ Vicki Crosbie,¹⁴ Sandra Luscombe,¹⁴ Rebecca Baatjes,¹ Lonnie Zwaigenbaum,¹⁵ Wendy Roberts,^{5,16} Bridget Fernandez,¹⁴ Peter Szatmari,¹³ and Stephen W. Scherer^{1,*}

Article

Patterns of somatic structural variation in human cancer genomes

<https://doi.org/10.1038/s41586-019-1913-9>
Received: 22 September 2017
Accepted: 18 November 2019
Published online: 5 February 2020

Yilong Li^{1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,38,39,40,41,42,43,44,45,46,47,48,49,50,51,52,53,54,55,56,57,58,59,60,61,62,63,64,65,66,67,68,69,70,71,72,73,74,75,76,77,78,79,80,81,82,83,84,85,86,87,88,89,90,91,92,93,94,95,96,97,98,99,100,101,102,103,104,105,106,107,108,109,110,111,112,113,114,115,116,117,118,119,120,121,122,123,124,125,126,127,128,129,130,131,132,133,134,135,136,137,138,139,140,141,142,143,144,145,146,147,148,149,150,151,152,153,154,155,156,157,158,159,160,161,162,163,164,165,166,167,168,169,170,171,172,173,174,175,176,177,178,179,180,181,182,183,184,185,186,187,188,189,190,191,192,193,194,195,196,197,198,199,200,201,202,203,204,205,206,207,208,209,210,211,212,213,214,215,216,217,218,219,220,221,222,223,224,225,226,227,228,229,230,231,232,233,234,235,236,237,238,239,240,241,242,243,244,245,246,247,248,249,250,251,252,253,254,255,256,257,258,259,260,261,262,263,264,265,266,267,268,269,270,271,272,273,274,275,276,277,278,279,280,281,282,283,284,285,286,287,288,289,290,291,292,293,294,295,296,297,298,299,300,301,302,303,304,305,306,307,308,309,310,311,312,313,314,315,316,317,318,319,320,321,322,323,324,325,326,327,328,329,330,331,332,333,334,335,336,337,338,339,340,341,342,343,344,345,346,347,348,349,350,351,352,353,354,355,356,357,358,359,360,361,362,363,364,365,366,367,368,369,370,371,372,373,374,375,376,377,378,379,380,381,382,383,384,385,386,387,388,389,390,391,392,393,394,395,396,397,398,399,400,401,402,403,404,405,406,407,408,409,410,411,412,413,414,415,416,417,418,419,420,421,422,423,424,425,426,427,428,429,430,431,432,433,434,435,436,437,438,439,440,441,442,443,444,445,446,447,448,449,450,451,452,453,454,455,456,457,458,459,460,461,462,463,464,465,466,467,468,469,470,471,472,473,474,475,476,477,478,479,480,481,482,483,484,485,486,487,488,489,490,491,492,493,494,495,496,497,498,499,500,501,502,503,504,505,506,507,508,509,510,511,512,513,514,515,516,517,518,519,520,521,522,523,524,525,526,527,528,529,530,531,532,533,534,535,536,537,538,539,540,541,542,543,544,545,546,547,548,549,550,551,552,553,554,555,556,557,558,559,560,561,562,563,564,565,566,567,568,569,570,571,572,573,574,575,576,577,578,579,580,581,582,583,584,585,586,587,588,589,590,591,592,593,594,595,596,597,598,599,600,601,602,603,604,605,606,607,608,609,610,611,612,613,614,615,616,617,618,619,620,621,622,623,624,625,626,627,628,629,630,631,632,633,634,635,636,637,638,639,640,641,642,643,644,645,646,647,648,649,650,651,652,653,654,655,656,657,658,659,660,661,662,663,664,665,666,667,668,669,670,671,672,673,674,675,676,677,678,679,680,681,682,683,684,685,686,687,688,689,690,691,692,693,694,695,696,697,698,699,700,701,702,703,704,705,706,707,708,709,710,711,712,713,714,715,716,717,718,719,720,721,722,723,724,725,726,727,728,729,730,731,732,733,734,735,736,737,738,739,740,741,742,743,744,745,746,747,748,749,750,751,752,753,754,755,756,757,758,759,760,761,762,763,764,765,766,767,768,769,770,771,772,773,774,775,776,777,778,779,780,781,782,783,784,785,786,787,788,789,790,791,792,793,794,795,796,797,798,799,800,801,802,803,804,805,806,807,808,809,810,811,812,813,814,815,816,817,818,819,820,821,822,823,824,825,826,827,828,829,830,831,832,833,834,835,836,837,838,839,840,841,842,843,844,845,846,847,848,849,850,851,852,853,854,855,856,857,858,859,860,861,862,863,864,865,866,867,868,869,870,871,872,873,874,875,876,877,878,879,880,881,882,883,884,885,886,887,888,889,890,891,892,893,894,895,896,897,898,899,900,901,902,903,904,905,906,907,908,909,910,911,912,913,914,915,916,917,918,919,920,921,922,923,924,925,926,927,928,929,930,931,932,933,934,935,936,937,938,939,940,941,942,943,944,945,946,947,948,949,950,951,952,953,954,955,956,957,958,959,960,961,962,963,964,965,966,967,968,969,970,971,972,973,974,975,976,977,978,979,980,981,982,983,984,985,986,987,988,989,990,991,992,993,994,995,996,997,998,999,1000}

Circulation: Genomic and Precision Medicine

Volume 14, Issue 4, August 2021
<https://doi.org/10.1161/CIRCGEN.120.003223>

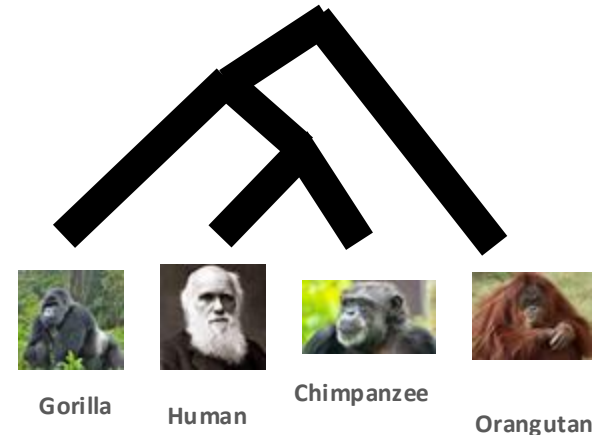


RESEARCH LETTERS

Long-Read Sequence Confirmed a Large Deletion Including *MYH6* and *MYH7* in an Infant of Atrial Septal Defect and Atrial Arrhythmias

Long Genome Sequence Alignment: Applications

- Third-generation sequencing technologies (produce reads of length 10 kb - 4 Mb), leads to major breakthroughs in recent past:
 - Human Genome Assembly
 - Rapid genetic diagnosis
 - Characterize structural variations and complex regions
- **New insights into the evolution of different species through whole genome analysis**



Long Genome Sequence Alignment: Applications

- Third-generation sequencing technologies (produce reads of length 10 kb - 4 Mb), leads to major breakthroughs in recent past:
 - Human Genome Assembly
 - Rapid genetic diagnosis
 - Characterize structural variations and complex regions
- New insights into the evolution of different species through whole genome analysis

**Bottlenecked by
LGSA**

Outline

- Emergence of **Long Genome Sequence Alignment**
- Current LGSA **algorithms, accelerators** and their **limitations**
- **TALCO**: A tiling technique based on convergence of traceback pointers for long genome sequence alignment
- **Key Contributions** and **Results**
- **Conclusion**



Broad Classification of Alignment Algorithms

Classical Dynamic Programming (DP) Algorithms

Ex: *Needleman-Wunsch, Smith-Waterman*

Non-Classical Algorithms

Ex: *WFA, $O(ND)$*

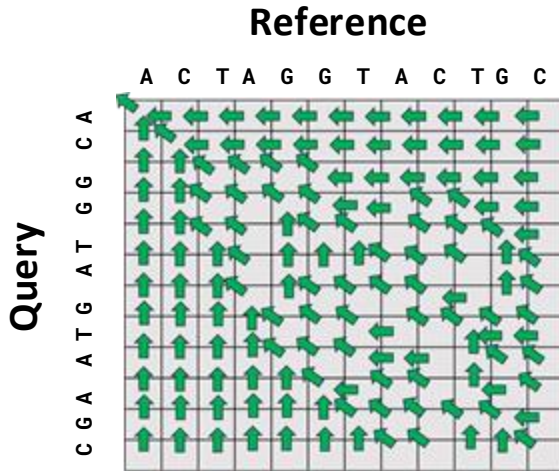
Classical DP based Alignment Algorithms

Classical Dynamic Programming (DP) Algorithms

Ex: *Needleman-Wunsch, Smith-Waterman*

Non-Classical Algorithms

Ex: *WFA, $O(ND)$*



1. Matrix Fill (Store traceback pointers)



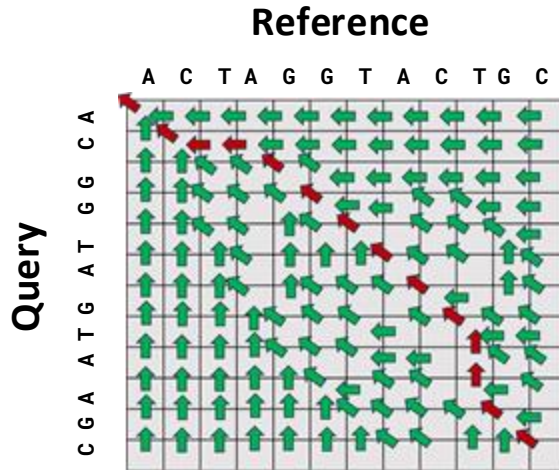
Classical DP based Alignment Algorithms

Classical Dynamic Programming (DP) Algorithms

Ex: Needleman-Wunsch, Smith-Waterman

Non-Classical Algorithms

Ex: WFA, $O(ND)$



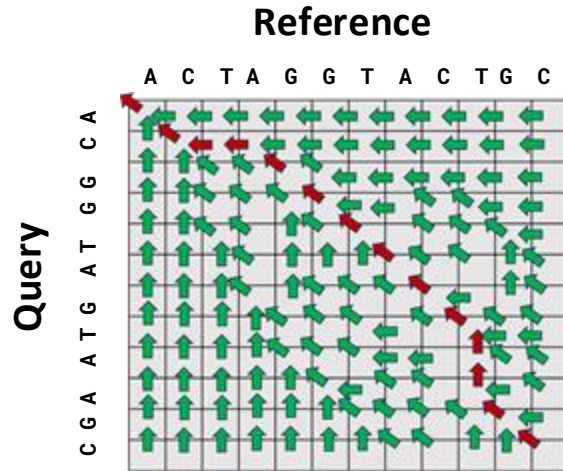
A	C	T	A	G	G	T	A	C	T	-	-	G	C
A	C	-	-	G	G	T	A	G	T	A	A	G	C

1. Matrix Fill (Store traceback pointers)
2. Optimal traceback path

Non-Classical Alignment Algorithms

Classical Dynamic Programming (DP) Algorithms

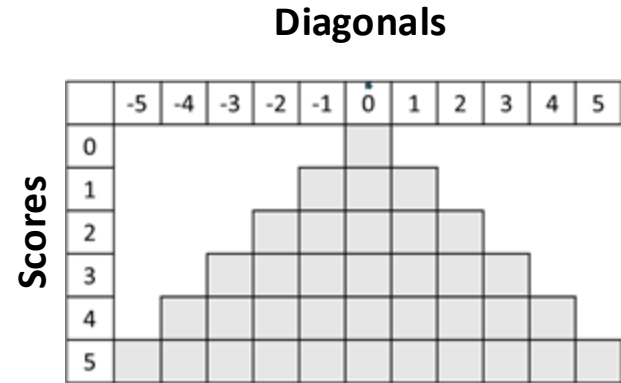
Ex: Needleman-Wunsch, Smith-Waterman



1. Matrix Fill (Store traceback pointers)
2. Optimal traceback path

Non-Classical Algorithms

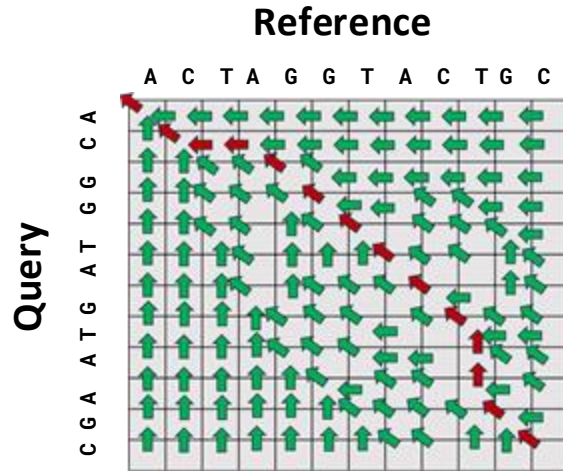
Ex: WFA, $O(ND)$



Non-Classical Alignment Algorithms

Classical Dynamic Programming (DP) Algorithms

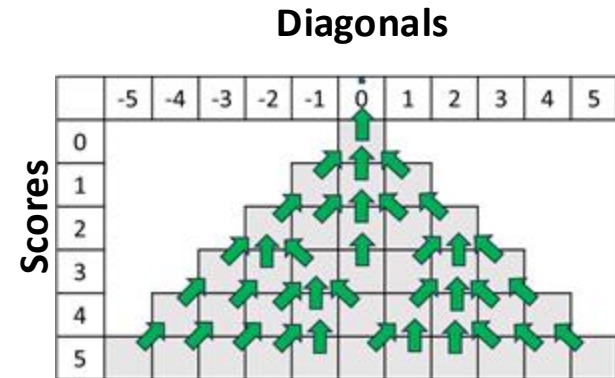
Ex: Needleman-Wunsch, Smith-Waterman



1. Matrix Fill (Store traceback pointers)
2. Optimal traceback path

Non-Classical Algorithms

Ex: WFA, $O(ND)$

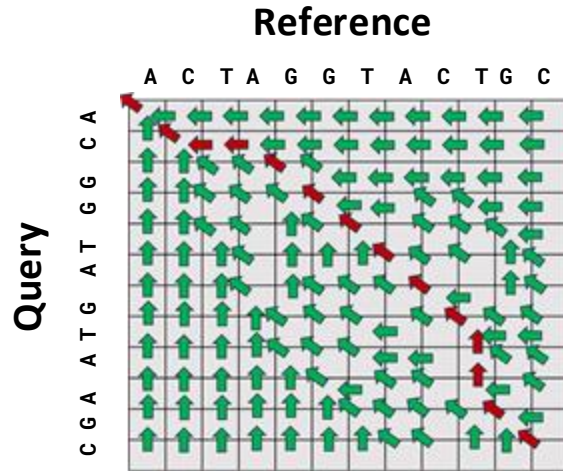


1. Matrix Fill (Store traceback pointers)

Non-Classical Alignment Algorithms

Classical Dynamic Programming (DP) Algorithms

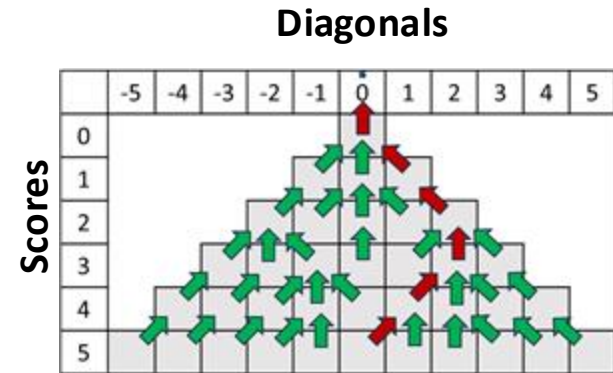
Ex: Needleman-Wunsch, Smith-Waterman



1. Matrix Fill (Store traceback pointers)
2. Optimal traceback path

Non-Classical Algorithms

Ex: WFA, $O(ND)$



1. Matrix Fill (Store traceback pointers)
2. **Optimal traceback path**

Comparison: Classical-DP and Non-Classical

Classical Dynamic Programming (DP) Algorithms

Ex: *Needleman-Wunsch, Smith-Waterman*

Non-Classical Algorithms

Ex: *WFA, $O(ND)$*

Both categories of algorithms produce optimal alignments

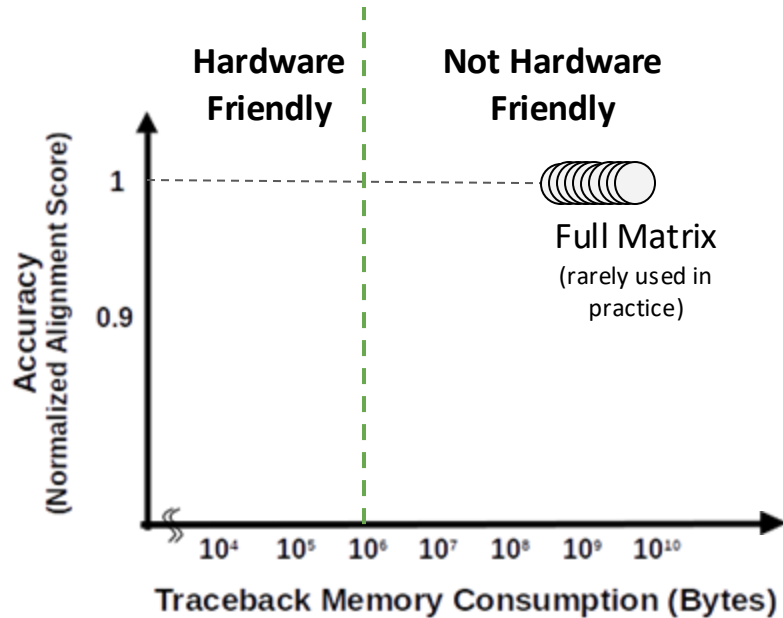
Uniform dependencies
Easier to accelerate

Non-Uniform dependencies
Harder to accelerate

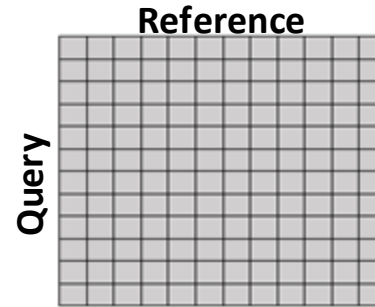
More popular

Very Fast for similar sequences

Full Matrix Sequence Alignment Algorithms



Classical-DP

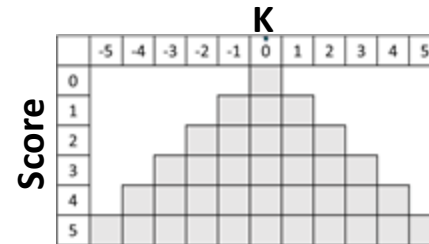


- Computed cells
- Uncomputed cells

Space Complexity $O(NM)$

Ex: Needleman-Wunsch

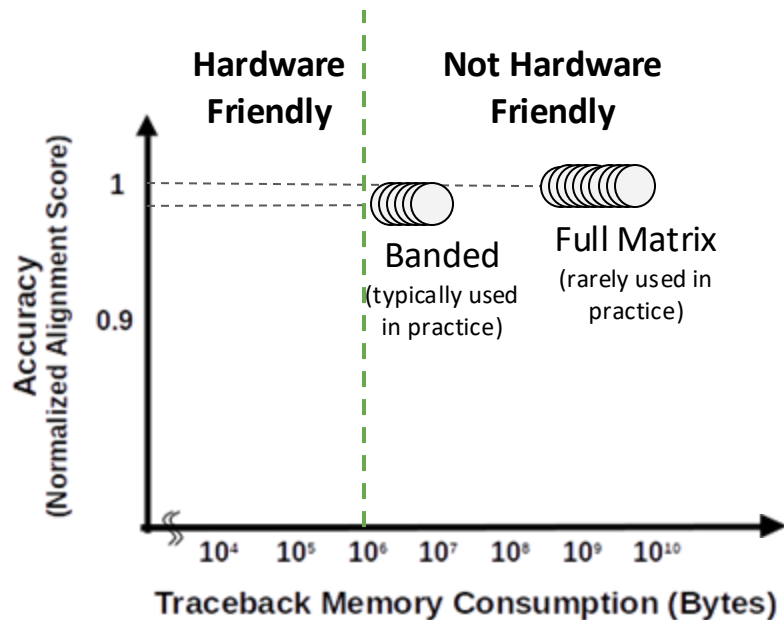
Non Classical



Space Complexity $O(S^2)$

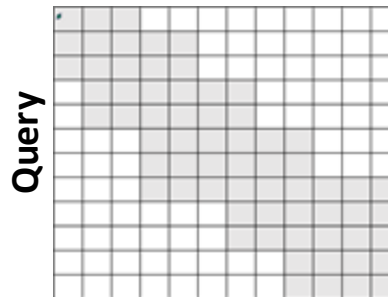
Ex: WFA

Banded Sequence Alignment Algorithms



Classical-DP

Reference

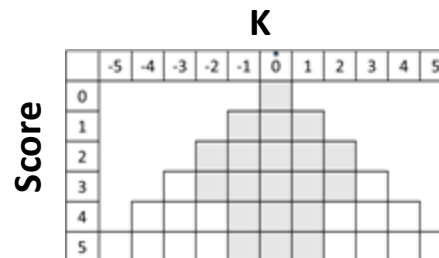


- Computed cells
- Uncomputed cells

Space Complexity
 $O(ND)$

Ex: X-Drop

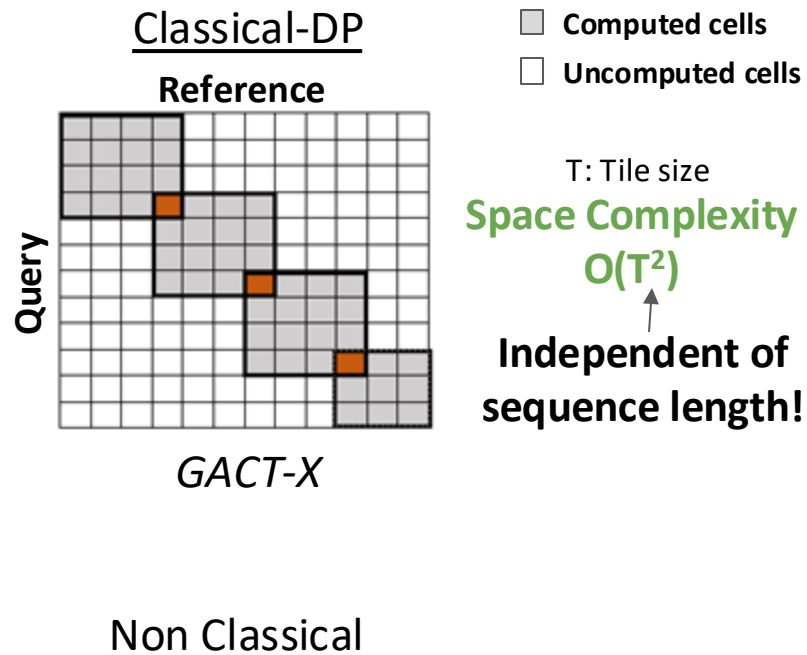
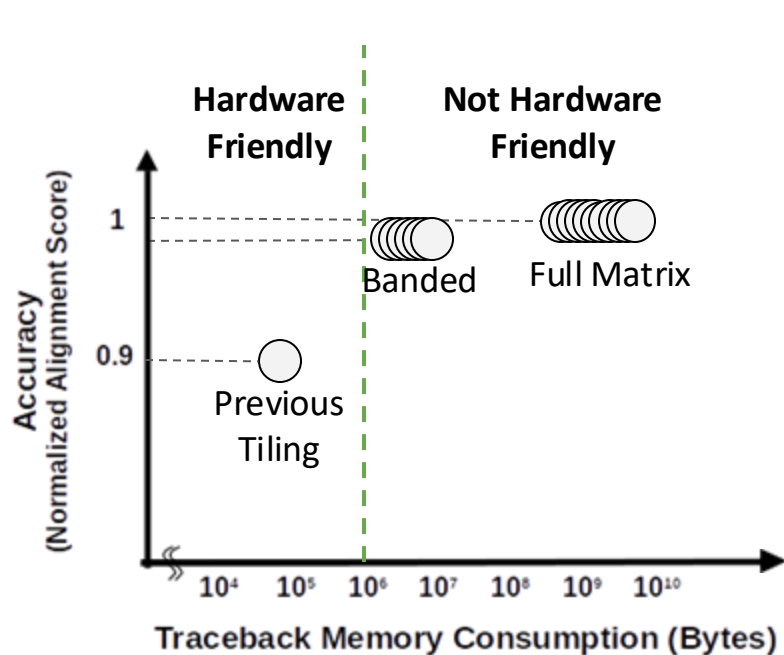
Non Classical



Space Complexity
 $O(SD)$

Ex: WFA-Adapt

Tiling heuristic



Tiling is **never applied to non-classical sequence alignment algorithm**

Architecture Papers Adopting Tiling Heuristics

- **GACT** – Darwin: A Genomics Co-processor Provides up to 15,000X Acceleration on Long Read Assembly (**ASPLOS 2018 Best Paper Award**)
- **GACT-X** – Darwin-WGA: A Co-processor Provides Increased Sensitivity in Whole Genome Alignments with High Speedup

Lower accuracy imposes challenges for tiling-based accelerators to be adopted in critical real-world applications (e.g. medical diagnoses)

Van der Auwera, Geraldine A., et al. "From FastQ data to high-confidence variant calls: the genome analysis toolkit best practices pipeline." Current protocols in bioinformatics 43.1 (2013)

- **RAPIDx**: High-performance ReRAM processing in-memory accelerator for sequence alignment (**TCAD 2023**)
- **GMX**: Instruction Set Extensions for Fast, Scalable, and Efficient Genome Sequence Alignment (**MICRO 2023**)
- **Scrooge**: a fast and memory-frugal genomic sequence aligner for CPUs, GPUs, and ASICs (**Bioinformatics 2023**)



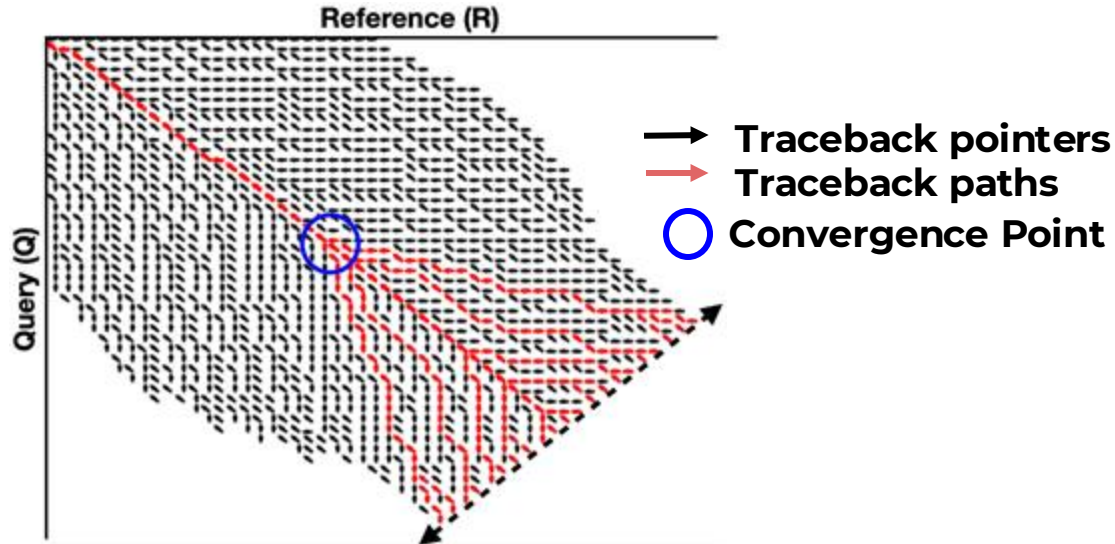
Outline

- Emergence of **Long Genome Sequence Alignment**
- Current LGSA algorithms, accelerators and their limitations
- **TALCO**: A tiling technique based on convergence of traceback pointers for long genome sequence alignment
- **Key Contributions** and **Results**
- **Conclusion**



Key Insight: Convergence of Traceback Paths

TALCO (Tiling Long Genome Alignment using Convergence of Traceback Pointers) is based on the following observation:

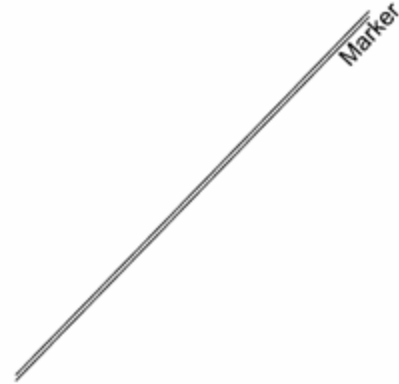


Experiment: Pairwise sequence alignment using Needleman-Wunsch with X-Drop banding

TALCO: Tiling technique for long genome alignment

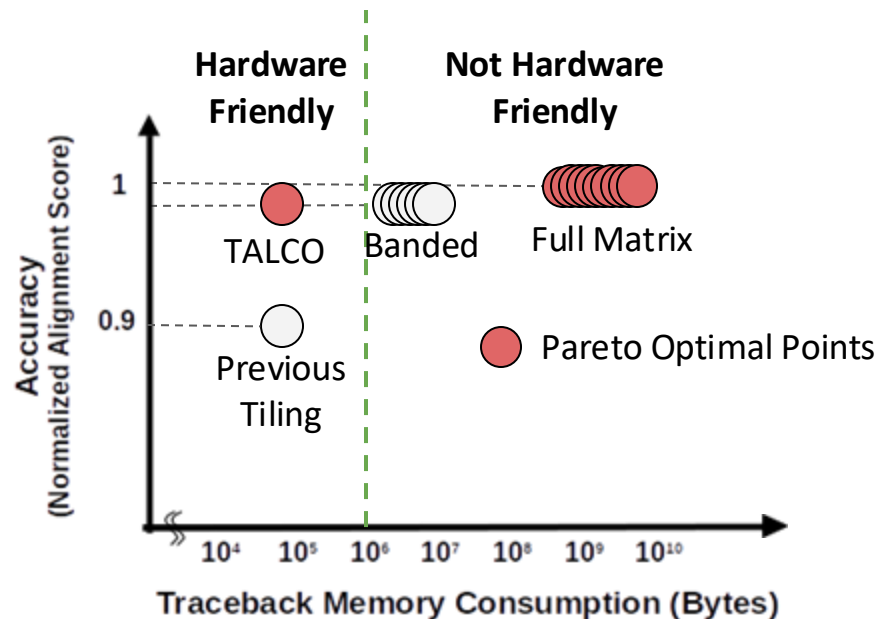
TALCO algorithm has **two** phases:

1. Stores traceback pointers till the **Marker**
1. Find point of convergence of traceback pointers using **pointer-redirection**



TALCO applied to X-Drop Algorithm

TALCO is on the Pareto Optimal Frontier



- Constant Space complexity
- Guarantees optimality under banding constraints

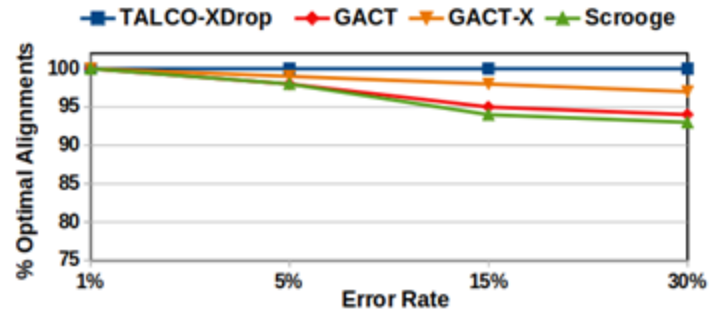
Outline

- Emergence of **Long Genome Sequence Alignment**
- Current LGSA algorithms, accelerators and their limitations
- **TALCO**: A tiling technique based on convergence of traceback pointers for long genome sequence alignment
- **Key Contributions and Results**
- **Conclusion**



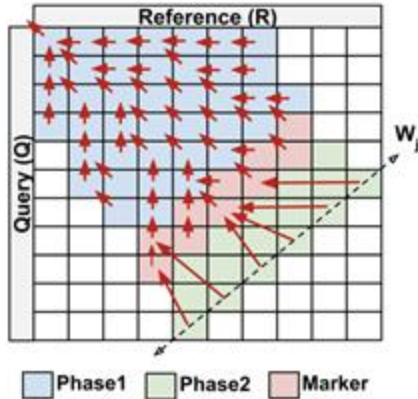
Key Contributions and Results

1. TALCO, **guarantees optimality** under banding constraints

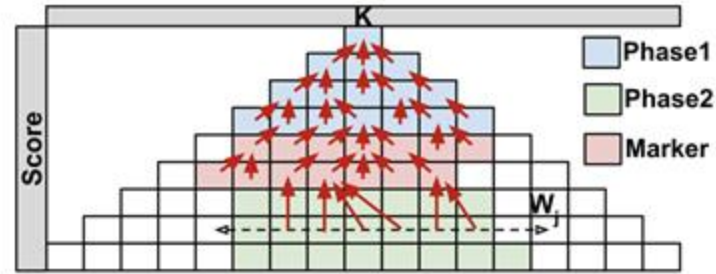


Key Contributions and Results

1. TALCO, **guarantees optimality** under banding constraints
2. We applied TALCO to X-Drop (**TALCO-XDrop**) and WFA-Adapt (**TALCO-WFAA**)



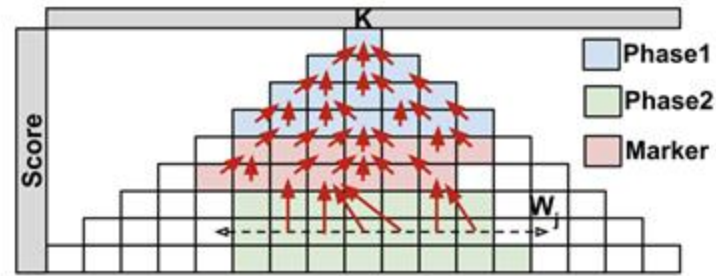
TALCO-XDrop



TALCO-WFAA

Key Contributions and Results

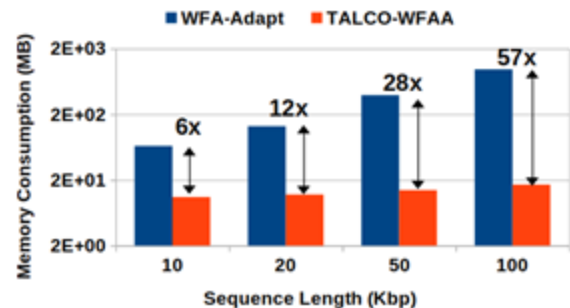
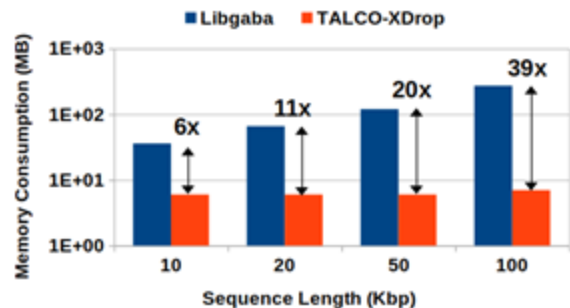
1. TALCO, **guarantees optimality** under banding constraints
2. We applied TALCO to X-Drop (**TALCO- XDrop**) and WFA-Adapt (**TALCO-WFAA**)
3. TALCO-WFAA is the **first accelerator** based on the **WFA-Adapt algorithm** capable of performing arbitrary long sequence alignments



TALCO-WFAA

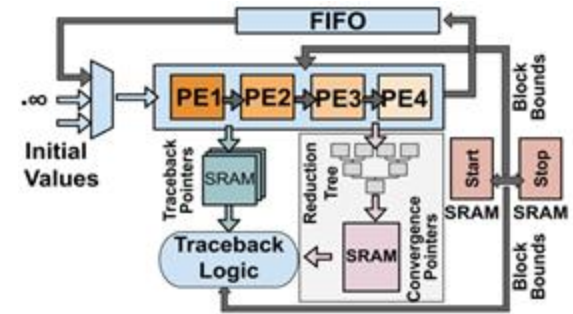
Key Contributions and Results

1. TALCO, **guarantees optimality** under banding constraints
2. We applied TALCO to X-Drop (**TALCO- XDrop**) and WFA-Adapt (**TALCO-WFAA**)
3. TALCO-WFAA is the **first accelerator** based on the **WFA-Adapt algorithm** capable of performing arbitrary long sequence alignments
4. TALCO-XDrop and TALCO-WFAA (**software**) achieves up to **39x** and **57x** improvement in **memory footprint**, respectively, compared to software baselines

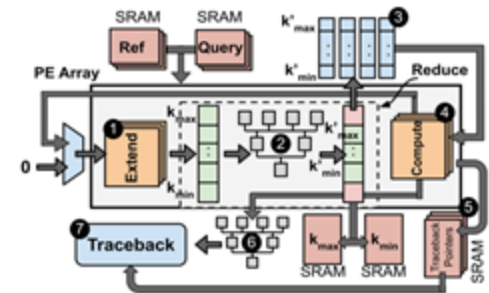


Key Contributions and Results

1. TALCO, **guarantees optimality** under banding constraints
2. We applied TALCO to X-Drop (TALCO- XDrop) and WFA-Adapt (TALCO-WFAA)
3. TALCO-WFAA is the **first accelerator** based on the **WFA-Adapt algorithm** capable of performing arbitrary long sequence alignments
4. TALCO-XDrop and TALCO-WFAA (**software**) achieves up to **39x** and **57x** improvement in memory footprint, respectively, compared to software baselines
5. **Designed hardware accelerators** for TALCO- XDrop and TALCO-WFAA



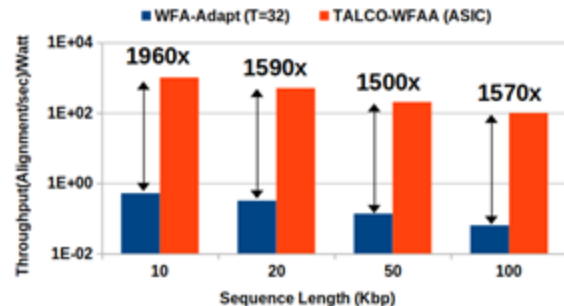
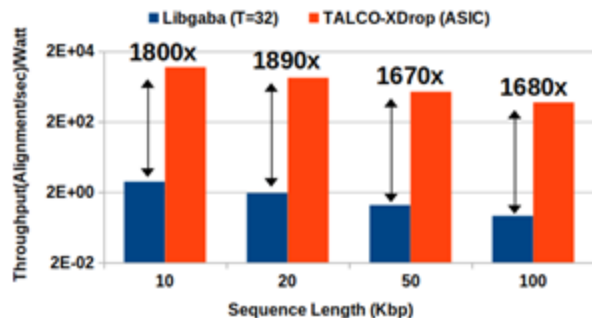
TALCO-XDrop hardware design



TALCO-WFAA hardware design

Key Contributions and Results

1. TALCO, **guarantees optimality** under banding constraints
2. We applied TALCO to X-Drop (**TALCO- XDrop**) and WFA-Adapt (**TALCO-WFAA**)
3. TALCO-WFAA is the **first accelerator** based on the **WFA-Adapt algorithm** capable of performing arbitrary long sequence alignments
4. TALCO-XDrop and TALCO-WFAA (**software**) achieves up to **39x** and **57x** improvement in memory footprint, respectively, compared to software baselines
5. **Designed hardware accelerators** for TALCO- XDrop and TALCO-WFAA
6. TALCO-XDrop and TALCO-WFAA (**ASIC**) achieves up to **~1,900X** and **~2,000X**, respectively, improvement in **alignment throughput/watt** over software baselines



Key Contributions and Results

1. TALCO, **guarantees optimality** under banding constraints
2. We applied TALCO to X-Drop (**TALCO- XDrop**) and WFA-Adapt (**TALCO-WFAA**)
3. TALCO-WFAA is the **first accelerator** based on the **WFA-Adapt algorithm** capable of performing arbitrary long sequence alignments
4. TALCO-XDrop and TALCO-WFAA (**software**) achieves up to **39x** and **57x** improvement in memory footprint, respectively, compared to software baselines
5. **Designed hardware accelerators** for TALCO- XDrop and TALCO-WFAA
6. TALCO-XDrop and TALCO-WFAA (**ASIC**) achieves up to **~1,900X** and **~2,000X**, respectively, improvement in **alignment throughput/watt** over software baselines
7. We **synthesized** TALCO-XDrop and TALCO-WFAA for **FPGAs** available on the **Amazon EC2 FPGA instances**



<https://github.com/TurakhiaLab/TALCO/blob/main/hardware/README.md>

Building on AWS EC2 F1 instance

Follow the below instructions to execute TALCO-XDrop and TALCO-WFAA on the AWS EC2 F1 instance, [f1.2xlarge](#).

- Clone aws-fpga repository

```
git clone https://github.com/aws/aws-fpga
cd aws-fpga
source vitis_runtime_setup.sh
```

- Clone TALCO repository

```
git clone https://github.com/TurakhiaLab/TALCO.git
export TALCO_DIR=$(PWD)/TALCO
cd TALCO/hardware/TALCO-XDrop
```

- Steps for running on the EC2 F1 instance, f1.2xlarge (MODE=hw)

```
source $TALCO_DIR/hardware/scripts/run.sh
$TALCO_DIR/dataset/sequence_fa TALCO_XDrop.awsxc1bin
```



Key Contributions and Results

1. TALCO, **guarantees optimality** under banding constraints
2. We applied TALCO to X-Drop (**TALCO- XDrop**) and WFA-Adapt (**TALCO-WFAA**)
3. TALCO-WFAA is the **first accelerator** based on the **WFA-Adapt algorithm** capable of performing arbitrary long sequence alignments
4. TALCO-XDrop and TALCO-WFAA (**software**) achieves up to **39x** and **57x** improvement in memory footprint, respectively, compared to software baselines
5. **Designed hardware accelerators** for TALCO- XDrop and TALCO-WFAA
6. TALCO-XDrop and TALCO-WFAA (**ASIC**) achieves up to **~1,900X** and **~2,000X**, respectively, improvement in **alignment throughput/watt** over software baselines
7. We **synthesized** TALCO-XDrop and TALCO-WFAA for **FPGAs** available on the **Amazon EC2 FPGA instances**



<https://github.com/TurakhiaLab/TALCO/>



HPCA Artifact Evaluation

Outline

- Emergence of **Long Genome Sequence Alignment**
- Current techniques and their **limitations** to **hardware accelerate** long genome sequence alignment
- **TALCO**: A tiling technique based on convergence of traceback pointers for long genome sequence alignment
- **Key Contributions** and **Results**
- **Conclusion**



Conclusion

- We present TALCO, a novel tiling technique for long genome sequence alignment
 - **Maintains a constant memory footprint**
 - Ensures **optimal alignments** under banding constraints
- We applied TALCO to X-Drop (**TALCO-XDrop**) and WFA-Adapt (**TALCO-WFAA**)
- TALCO-XDrop (TALCO-WFAA) software achieve up to **39X (57X)** improvement in **memory footprint** for long alignments compared to software baselines
- We present **hardware accelerator designs** for **TALCO-XDrop** and **TALCO-WFAA**
- TALCO-XDrop (TALCO-WFAA) ASIC achieves up to **1,900X (2,000X)** improvement in **alignment throughput/watt** over software baselines implementing the same algorithm





TALCO: Tiling Genome Sequence Alignment using Convergence of Traceback Pointers



Sumit Walia
Ph.D student



Cheng Ye
MS student



Arkid Bera
MS student



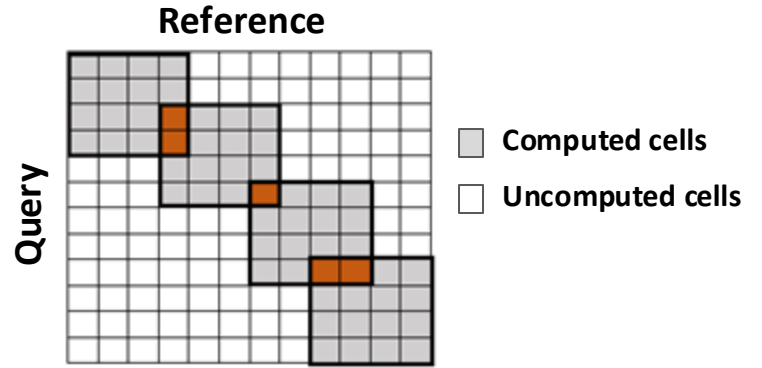
Dhruvi L.
MS student



Yatish Turakhia
Assistant Professor, UCSD

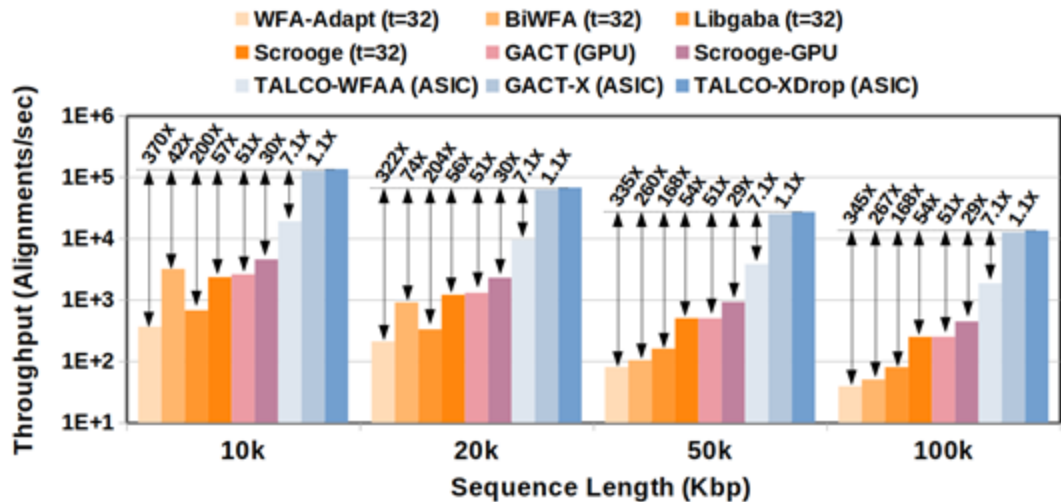
Thank you!

Additional Slides

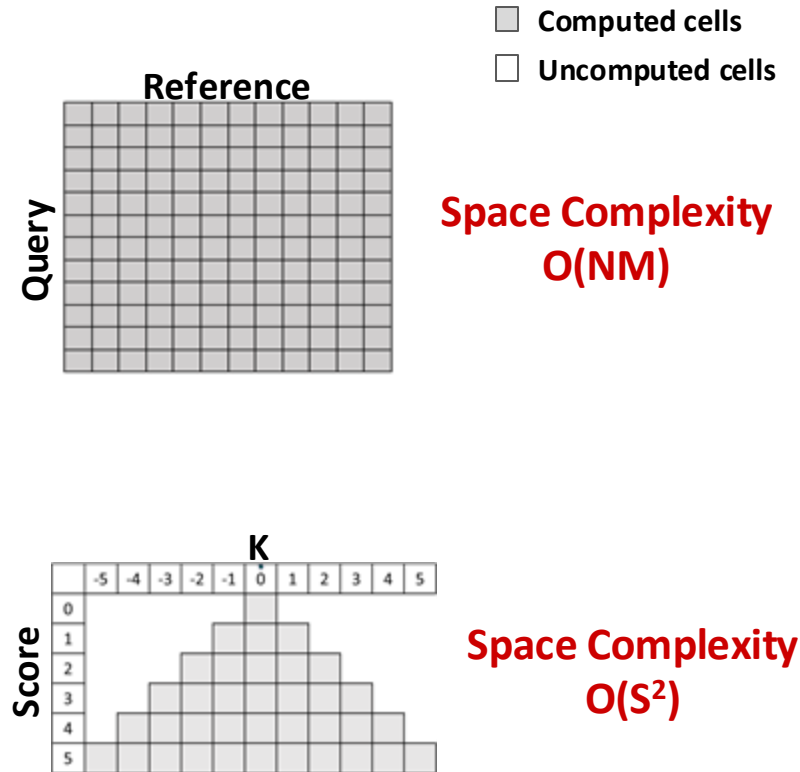
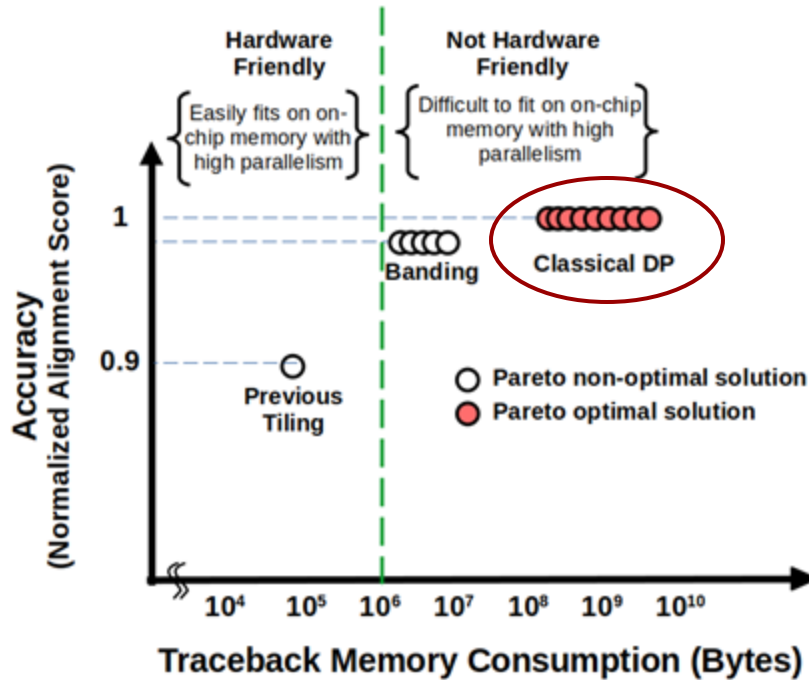


Key Results

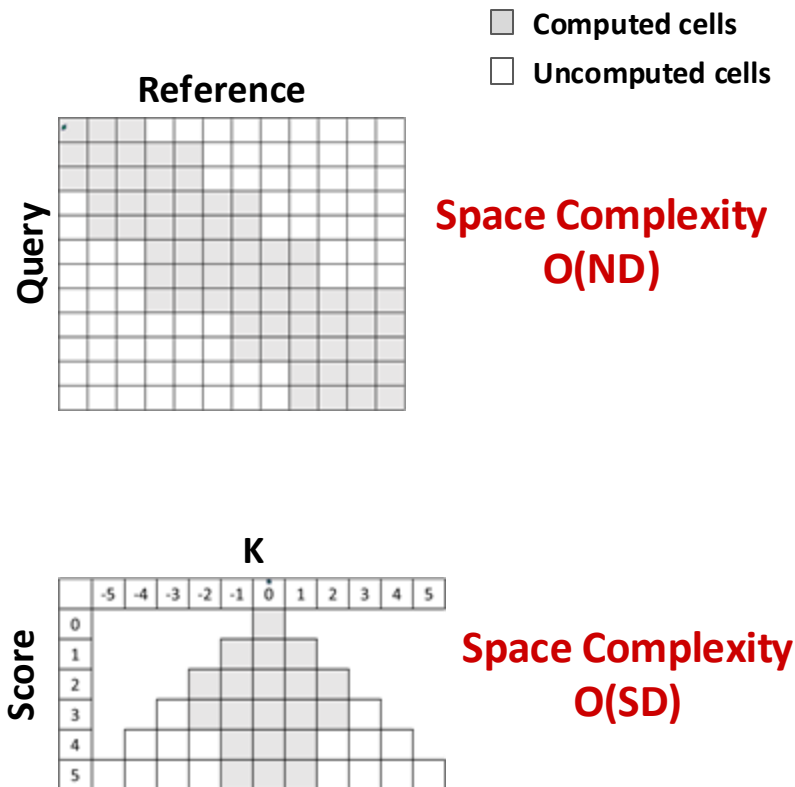
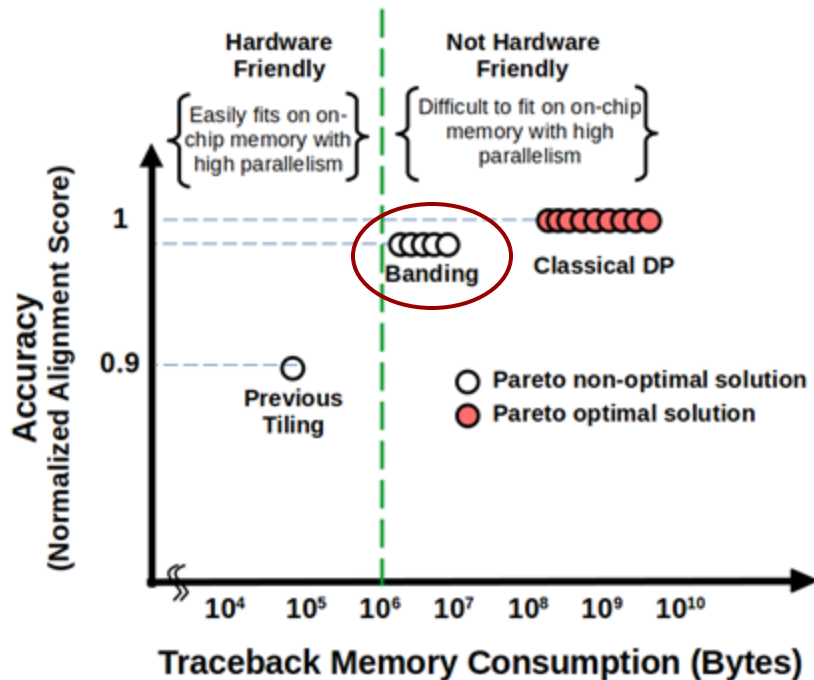
3. TALCO improves the alignment throughput over state-of-the-art GPU and ASIC baselines that implement tiling heuristics by over **50X** and **1.1X**, respectively.



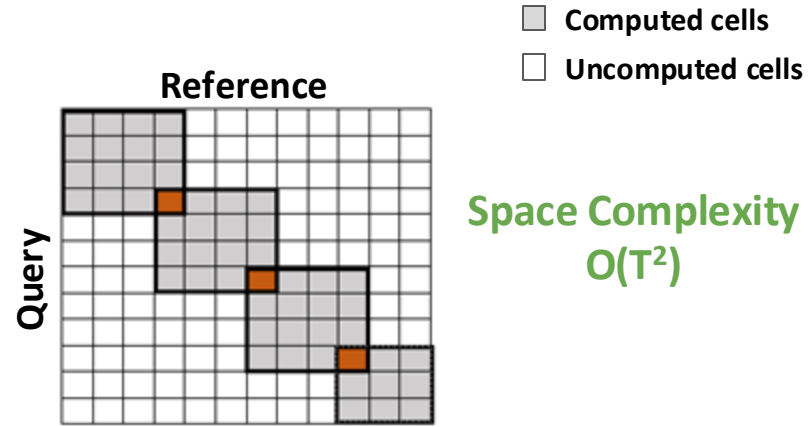
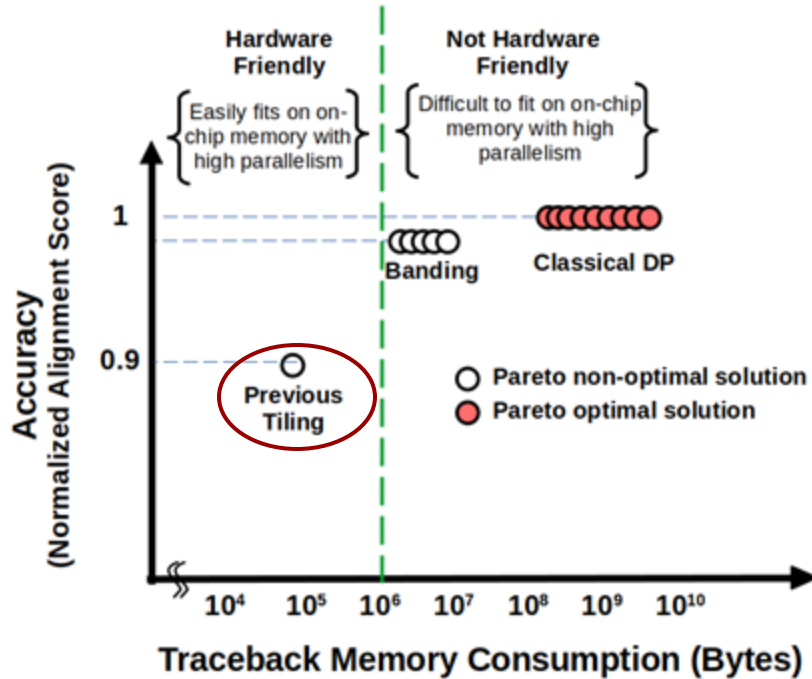
Full Matrix Sequence Alignment Algorithms



Banded Sequence Alignment Algorithms

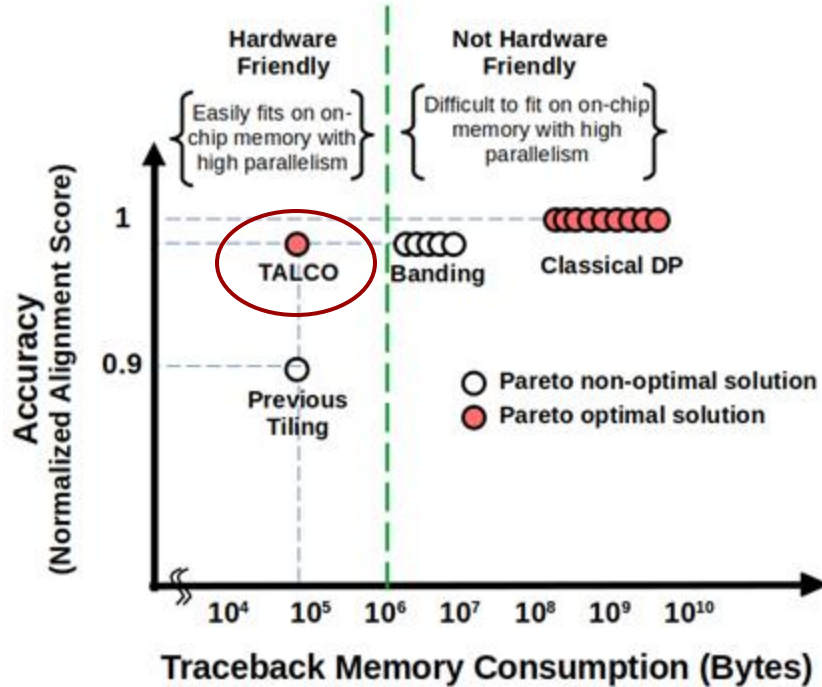


Tiling heuristic



Tiling is never applied to non-classical sequence alignment algorithm

TALCO on Pareto Plot



- Constant Space complexity
- Guarantees optimality under banding constraints (dynamic overlap between consecutive tiles)

Key Contributions

3. We implemented TALCO-XDrop and TALCO-WFAA in software, and achieves up to **39x** and **57x** improvement in **memory footprint**, respectively, compared to software baselines

