

## **RawHash and RawHash2:**

Enabling Fast & Accurate Real-time Analysis  
of Raw Nanopore Signals for Large Genomes  
using a Hash-based Seeding Mechanism

Can Firtina

[canfirtina@gmail.com](mailto:canfirtina@gmail.com)

<https://cfirtina.com>

17 January 2024

Leibniz Institute for Immunotherapy (LIT)

**SAFARI**

**ETH** zürich

# Brief Self Introduction

---



- **Can Firtina**

- Ph.D. Student in [SAFARI Research Group](#) led by [Prof. Onur Mutlu](#)

- **Research interests:** Bioinformatics & Computer Architecture

- Real-time genome analysis
- Similarity search in a large space of genomic data
- Hardware-Algorithm co-design to accelerate genome analysis
- Genome editing
- Error correction

- Get to know **our group and our research**

- **Group website:** <https://safari.ethz.ch/>
- **Contact me:** [canfirtina@gmail.com](mailto:canfirtina@gmail.com)
- **Website:** <https://cfirtina.com>
- **Twitter (aka X):** <https://twitter.com/FirtinaC>

# Professor Mutlu

---



## ■ Onur Mutlu

- ❑ Full Professor @ ETH Zurich ITET (INFK), since September 2015
- ❑ 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](mailto:omutlu@gmail.com) (Best way to reach)
- ❑ <https://people.inf.ethz.ch/omutlu/projects.htm>

## ■ Research and Teaching in:

- ❑ Computer architecture, computer systems, hardware security, bioinformatics
- ❑ Memory and storage systems
- ❑ Hardware security, safety, predictability
- ❑ Fault tolerance
- ❑ Hardware/software cooperation
- ❑ Architectures for bioinformatics, health, medicine
- ❑ ...

# SAFARI Research Group

*Computer architecture, HW/SW, systems, bioinformatics, security, memory*



40+ Researchers

**SAFARI**  
SAFARI Research Group  
safari.ethz.ch

Think BIG, Aim HIGH!

**SAFARI**

<https://safari.ethz.ch>

# Four Key Current Directions

---

- Fundamentally **Secure/Reliable/Safe** Architectures
- Fundamentally **Energy-Efficient** Architectures
  - **Memory-centric** (Data-centric) Architectures
- Fundamentally **Low-Latency and Predictable** Architectures
- Architectures for **AI/ML, Genomics, Medicine, Health**

# Agenda for Today

---

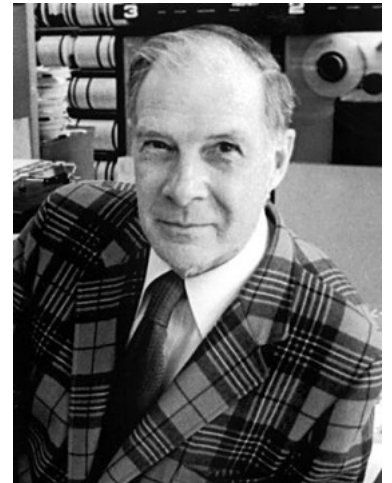
- Cutting-edge in Accelerating Genome Analysis
  - Intelligent genome analysis
  
- Enabling Fast and Accurate Real-time Analysis
  - RawHash and RawHash2
  
- Conclusion

# The Goal of Computing: Beyond Numbers

---

“The purpose of **computing** is [to gain] **insight**, not numbers”

Richard Hamming



---

We need to gain insights  
and observations  
much more efficiently  
than ever before

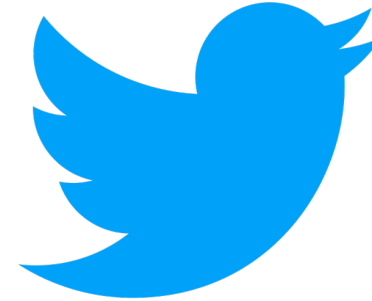


# Big Data is Everywhere

---



Astronomy  
25 zetta-bytes/year



Twitter (now X)  
0.5-15 billion tweets/year



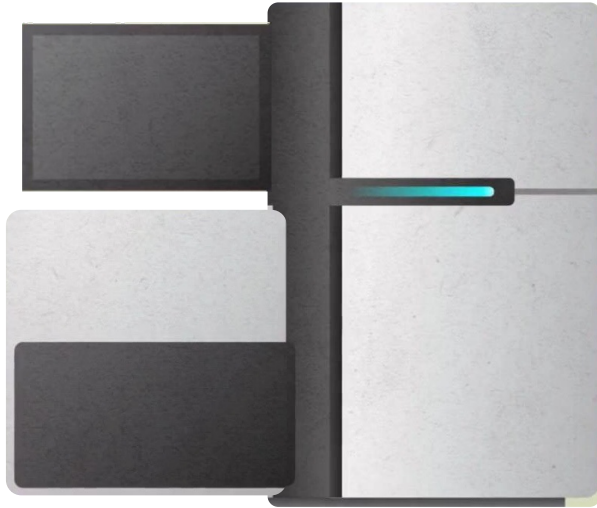
YouTube  
500-900 million hours/year



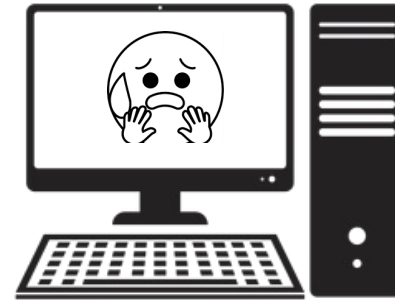
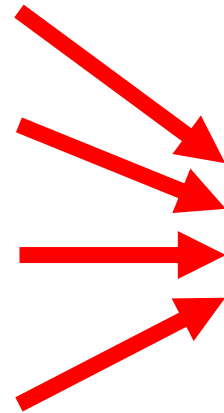
Genomics  
1 zetta-bases/year

# Problems with Data Analysis Today

---



**Special-Purpose** Machine  
for **Data Generation**



**General-Purpose** Machine  
for **Data Analysis**

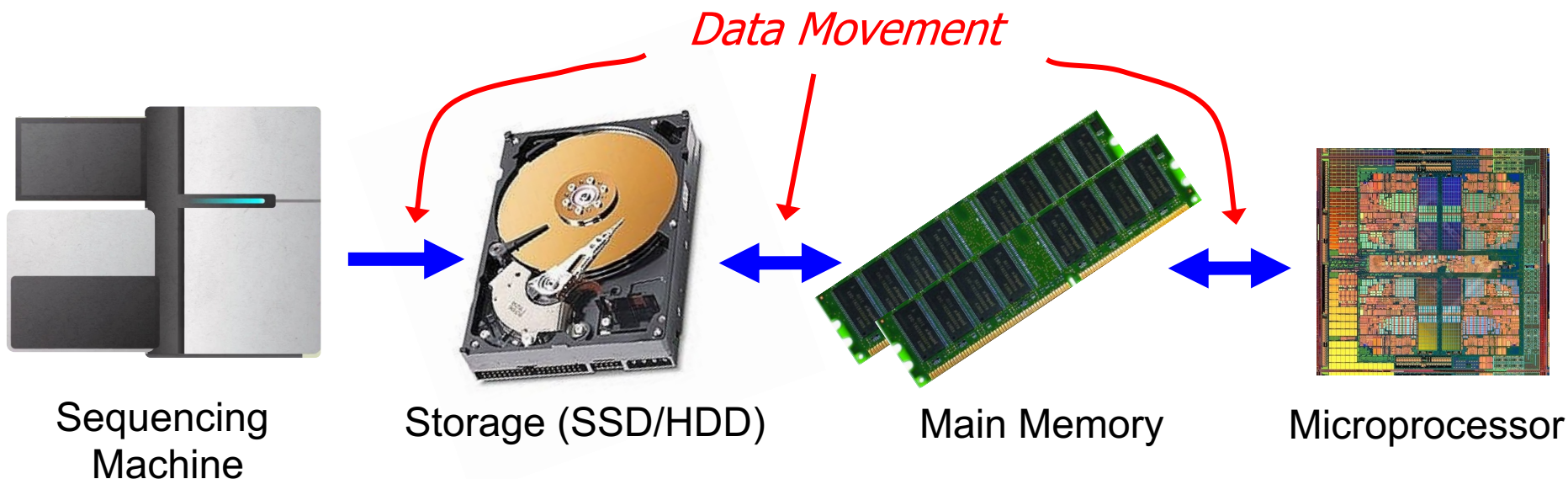
**FAST**

**SLOW**

**Slow and inefficient processing capability**  
**Large amounts of data movement**

# Data Movement Dominates Performance

- **Data movement** dominates performance and is a **major** system **energy bottleneck** (accounting for 40%-62%)



Single **memory** request **consumes** >160x-800x **more** **energy** compared to performing an **addition** operation

\* Boroumand et al., "Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks," ASPLOS 2018

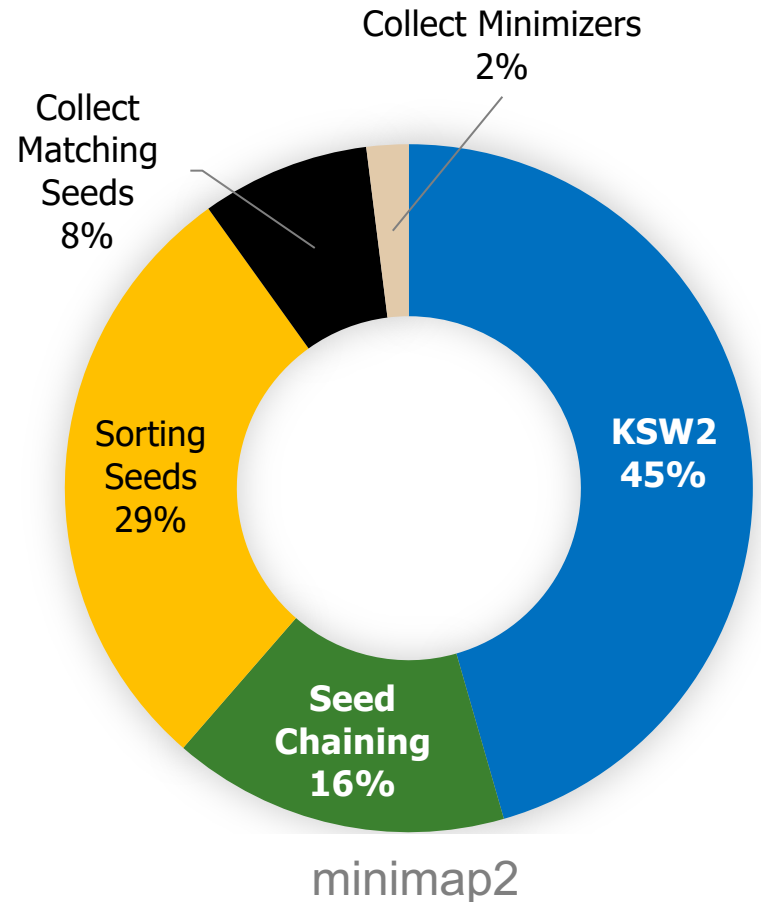
\* Kestor et al., "Quantifying the Energy Cost of Data Movement in Scientific Applications," IISWC 2013

\* Pandiyan and Wu, "Quantifying the energy cost of data movement for emerging smart phone workloads on mobile platforms," IISWC 2014

# Read Mapping Execution Time

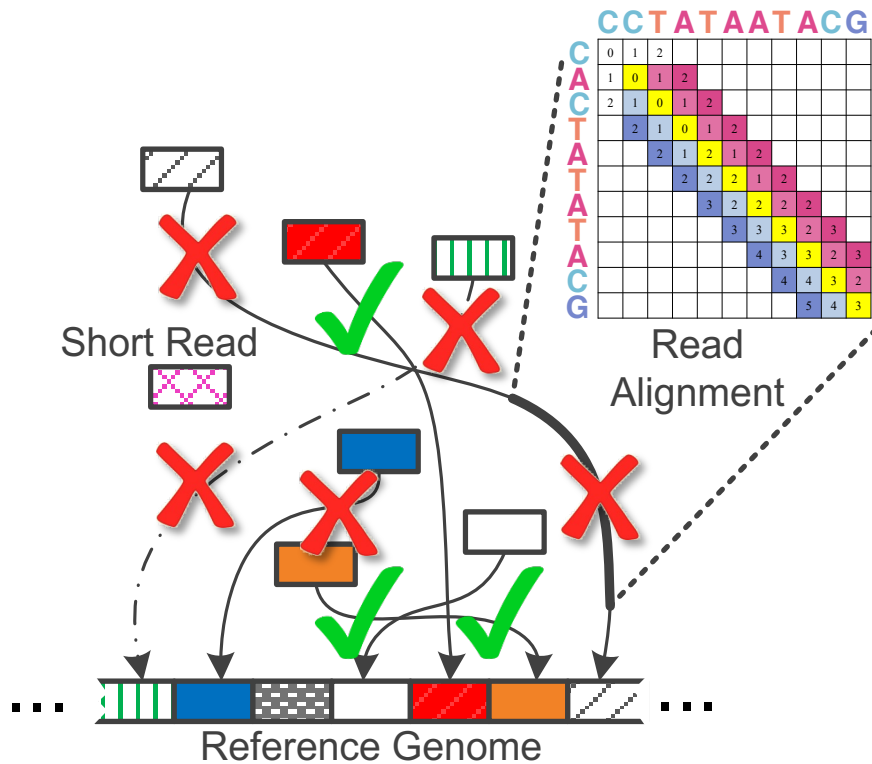
**> 60%**

**of the read mapper's  
execution time is spent  
in sequence alignment**



ONT FASTQ size: 103MB (151 reads), Mean length: 356,403 bp, std: 173,168 bp, longest length: 817,917 bp

# Large Search Space for Mapping Location



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

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

# New Genome Sequencing Technologies

---

## Nanopore sequencing technology and tools for genome assembly: computational analysis of the current state, bottlenecks and future directions

Damla Senol Cali ✉, Jeremie S Kim, Saugata Ghose, Can Alkan, Onur Mutlu

*Briefings in Bioinformatics*, bby017, <https://doi.org/10.1093/bib/bby017>

**Published:** 02 April 2018    **Article history** ▼



Oxford Nanopore MinION

Senol Cali+, “**Nanopore Sequencing Technology and Tools for Genome Assembly: Computational Analysis of the Current State, Bottlenecks and Future Directions**,” *Briefings in Bioinformatics*, 2018.

[[Open arxiv.org version](#)]

# New Genome Sequencing Technologies

---

## Nanopore sequencing technology and tools for genome assembly: computational analysis of the current state, bottlenecks and future directions

Damla Senol Cali ✉, Jeremie S Kim, Saugata Ghose, Can Alkan, Onur Mutlu

*Briefings in Bioinformatics*, bby017, <https://doi.org/10.1093/bib/bby017>

**Published:** 02 April 2018    **Article history** ▼



Oxford Nanopore MinION

Data → performance & energy bottleneck

---

We need intelligent algorithms  
and intelligent architectures  
that handle data well



# Intelligent Genome Analysis

Mohammed Alser, Joel Lindegger, Can Firtina, Nour Almadhoun, Haiyu Mao, Gagandeep Singh, Juan Gomez-Luna, Onur Mutlu

["From Molecules to Genomic Variations: Intelligent Algorithms and Architectures for Intelligent Genome Analysis"](#)

Computational and Structural Biotechnology Journal, 2022

[\[Source code\]](#)



ELSEVIER

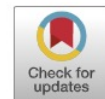


journal homepage: [www.elsevier.com/locate/csbj](http://www.elsevier.com/locate/csbj)



Review

From molecules to genomic variations: Accelerating genome analysis via intelligent algorithms and architectures



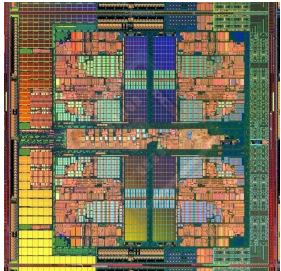
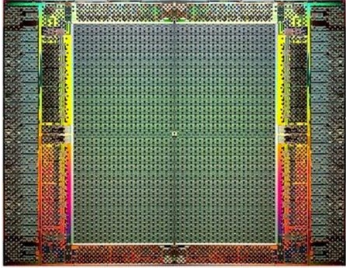
Mohammed Alser\*, Joel Lindegger, Can Firtina, Nour Almadhoun, Haiyu Mao, Gagandeep Singh, Juan Gomez-Luna, Onur Mutlu\*

ETH Zurich, Gloriastrasse 35, 8092 Zürich, Switzerland

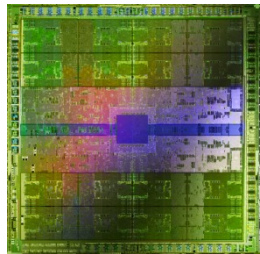
# Pushing Towards New Architectures

Modern systems

FPGAs



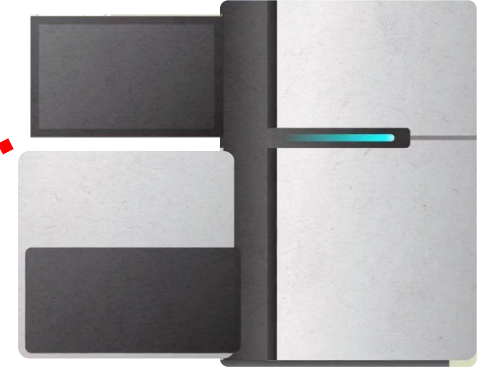
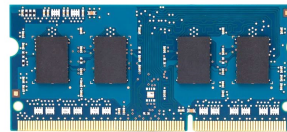
Heterogeneous Processors and Accelerators



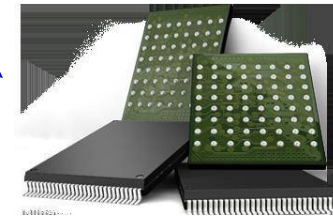
(General Purpose) GPUs



Hybrid Main Memory



Sequencing Machine

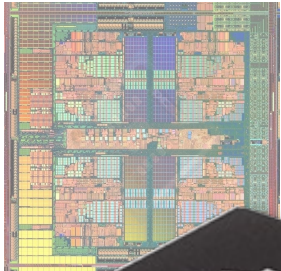
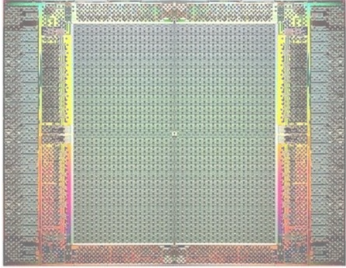


Persistent Memory/Storage

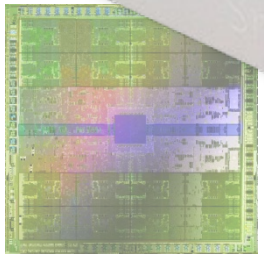
# Pushing Towards New Architectures

Modern systems

FPGAs

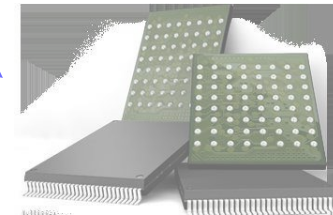


Hetero  
Pro  
Ac



(General Purpose) GPUs

Sequencing  
Machine



Persistent Memory/Storage

# Fast Genome Analysis

---

**Fast** genome analysis

in mere **seconds**

using **limited** computational resources

(e.g., a mobile device).

# Accurate Genome Analysis

---

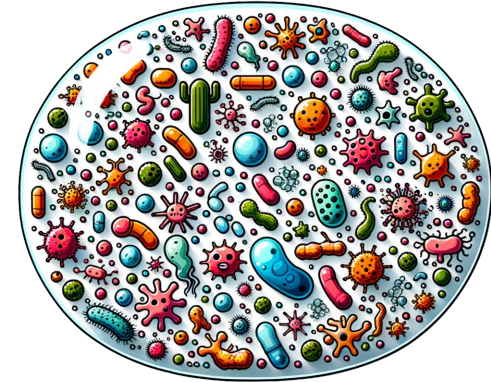
**Accurate** genome analysis  
to **make life-critical decisions**  
and **improving the quality of life**

# Faster, Scalable & Accurate Genome Analysis

---



Understanding **genetic variations, species, and evolution**



Predicting the **presence of pathogens** in an environment



Surveillance of **disease outbreaks**



**Personalized medicine**

# Personalized Medicine in UK

---

“From 2019, **all seriously ill children** in UK will be offered **whole genome sequencing** as part of their care”



# Rapid Surveillance of Disease Outbreaks

Real-time, portable genome sequencing for Ebola surveillance

Figure 1: Deployment of the portable genome surveillance system in Guinea.



University spinout's portable DNA sequencer has proved invaluable in tracking the global spread of coronavirus



Subscribe

Sign In





# Scalable SARS-CoV-2 Testing

## nature biomedical engineering

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾

[nature](#) > [nature biomedical engineering](#) > [articles](#) > [article](#)

Article | [Published: 01 July 2021](#)

## Massively scaled-up testing for SARS-CoV-2 RNA via next-generation sequencing of pooled and barcoded nasal and saliva samples

[Joshua S. Bloom](#) , [Laila Sathe](#), [...] [Valerie A. Arboleda](#) 

[Nature Biomedical Engineering](#) **5**, 657–665 (2021) | [Cite this article](#)

**4675** Accesses | **110** Altmetric | [Metrics](#)

Bloom+, "[Swab-Seq: A high-throughput platform for massively scaled up SARS-CoV-2 testing](#)", *Nature Biomedical Engineering*, 2021

# Large Scale Analysis



---

Applications  
are **only limited**  
by our imagination

# Genome Editing

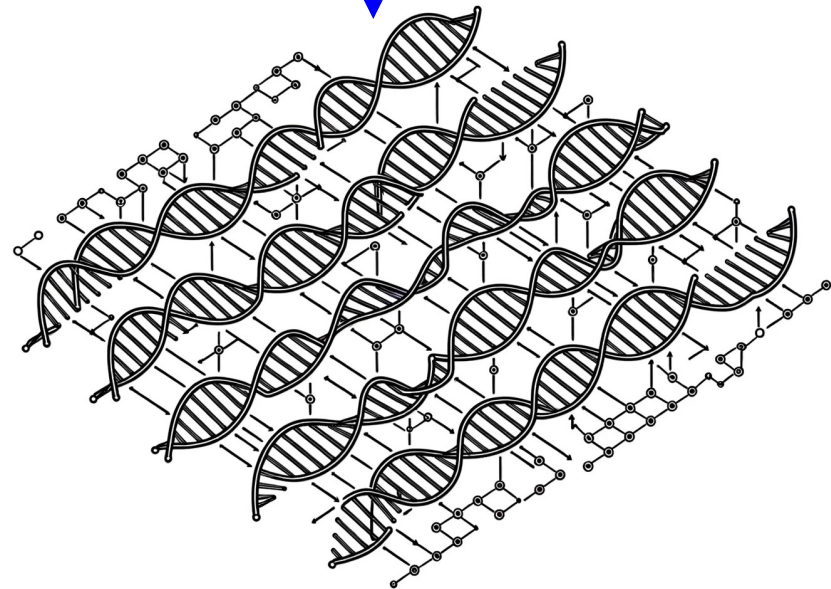
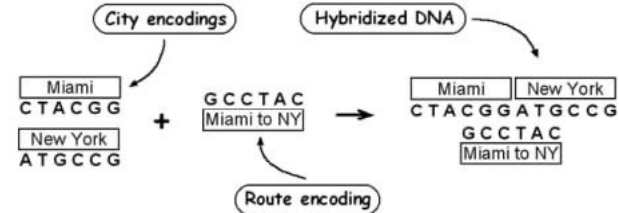
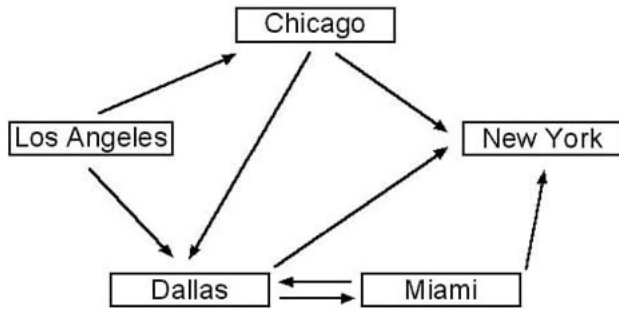


## The Nobel Prize in Chemistry 2020

awarded "for the development of a method of genome editing"



# DNA Computing



**Massive parallelism to solve (hard) problems!**

# Accelerating Genome Analysis [DAC 2023]

---

- Onur Mutlu and Can Firtina,  
**"Accelerating Genome Analysis via Algorithm-Architecture Co-Design"**  
*Invited Special Session Paper in Proceedings of the 60th Design Automation Conference (DAC), San Francisco, CA, USA, July 2023.*  
[[Slides \(pptx\)](#)] [[pdf](#)]  
[[Talk Video](#)] (38 minutes, including Q&A)  
[[Related Invited Paper](#)]  
[[arXiv version](#)]

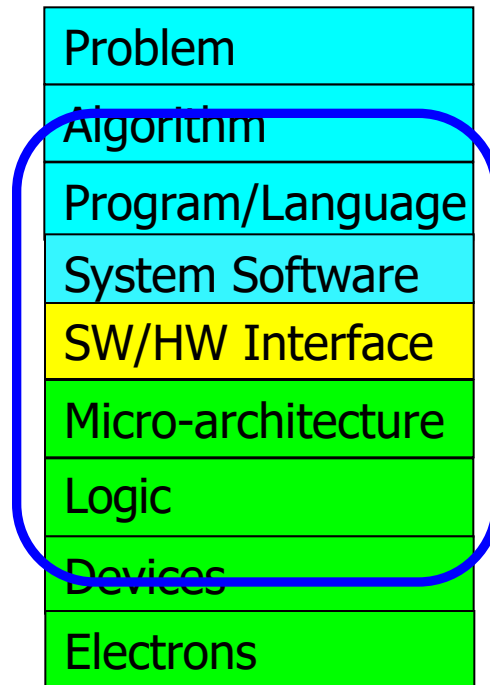
## Accelerating Genome Analysis via Algorithm-Architecture Co-Design

Onur Mutlu    Can Firtina  
*ETH Zürich*

# Algorithm-Arch-Device Co-Design is Critical

---

**Computer Architecture  
(expanded view)**



---

We need intelligent algorithms  
and intelligent architectures  
that handle data well



# New Frontiers: Raw Signal Analysis [ISMB 2023]

- [Can Firtina](#), [Nika Mansouri Ghiasi](#), [Joel Lindegger](#), [Gagandeep Singh](#), [Meryem Banu Cavlak](#), [Haiyu Mao](#), and [Onur Mutlu](#),  
**["RawHash: Enabling Fast and Accurate Real-Time Analysis of Raw Nanopore Signals for Large Genomes"](#)**  
*Proceedings of the [31st Annual Conference on Intelligent Systems for Molecular Biology \(ISMB\)](#) and the [22nd European Conference on Computational Biology \(ECCB\)](#), Jul 2023*  
[\[Bioinformatics Journal version\]](#)  
[\[Slides \(pptx\) \(pdf\)\]](#)  
[\[RawHash Source Code\]](#)

*Bioinformatics*, 2023, **39**, i297–i307  
<https://doi.org/10.1093/bioinformatics/btad272>

ISMB/ECCB 2023

OXFORD

## RawHash: enabling fast and accurate real-time analysis of raw nanopore signals for large genomes

**Can Firtina** <sup>1,\*</sup>, **Nika Mansouri Ghiasi** <sup>1</sup>, **Joel Lindegger** <sup>1</sup>, **Gagandeep Singh** <sup>1</sup>,  
**Meryem Banu Cavlak** <sup>1</sup>, **Haiyu Mao** <sup>1</sup>, **Onur Mutlu** <sup>1,\*</sup>

<sup>1</sup>Department of Information Technology and Electrical Engineering, ETH Zurich, 8092 Zurich, Switzerland

\*Corresponding author. Department of Information Technology and Electrical Engineering, ETH Zurich, Gloriastrasse 35, 8092 Zurich, Switzerland.  
E-mail: [firtinac@ethz.ch](mailto:firtinac@ethz.ch) (C.F.), [omutlu@ethz.ch](mailto:omutlu@ethz.ch) (O.M.)

# Fast and Accurate Real-Time Genome Analysis

---

- Can Firtina, Melina Soysal, Joel Lindegger, and Onur Mutlu,  
**"RawHash2: Accurate and Fast Mapping of Raw Nanopore Signals using a Hash-based Seeding Mechanism"**  
*Preprint on **arxiv**, September 2023.*  
[\[arXiv version\]](#)  
[\[RawHash2 Source Code\]](#)

## **RawHash2: Accurate and Fast Mapping of Raw Nanopore Signals using a Hash-based Seeding Mechanism**














Can Firtina   Melina Soysal   Joel Lindegger   Onur Mutlu  
*ETH Zürich*

# Accelerating ML & Genome Graphs [ACM TACO '23]

- Can Firtina, Kamlesh Pillai, Gurpreet S. Kalsi, Bharathwaj Suresh, Damla Senol Cali, Jeremie S. Kim, Taha Shahroodi, Meryem Banu Cavlak, Joël Lindegger, Mohammed Alser, Juan Gómez Luna, Sreenivas Subramoney, and Onur Mutlu, **"ApHMM: Accelerating Profile Hidden Markov Models for Fast and Energy-Efficient Genome Analysis"** **ACM TACO**, Dec 2023.  
[[Online link at ACM TACO](#)]  
[[arXiv preprint](#)]  
[[ApHMM Source Code](#)]

## ApHMM: Accelerating Profile Hidden Markov Models for Fast and Energy-Efficient Genome Analysis

Just Accepted

**Authors:**  [Can Firtina](#),  [Kamlesh Pillai](#),  [Gurpreet S. Kalsi](#),  [Bharathwaj Suresh](#),  [Damla Senol Cali](#),  
 [Jeremie S. Kim](#),  [Taha Shahroodi](#),  [Meryem Banu Cavlak](#),  [Joël Lindegger](#),  [Mohammed Alser](#),  
 [Juan Gómez Luna](#),  [Sreenivas Subramoney](#),  [Onur Mutlu](#) ([Less](#)) [Authors Info & Claims](#)

ACM Transactions on Architecture and Code Optimization • Accepted on October 2023 • <https://doi.org/10.1145/3632950>

**Published:** 28 December 2023 [Publication History](#)



# Genome Similarity Identification [NARGAB 2023]

- Can Firtina, Jisung Park, Mohammed Alser, Jeremie S. Kim, Damla Senol Cali, Taha Shahroodi, Nika Mansouri Ghiasi, Gagandeep Singh, Konstantinos Kanellopoulos, Can Alkan, and Onur Mutlu,

## **"BLEND: A Fast, Memory-Efficient, and Accurate Mechanism to Find Fuzzy Seed Matches in Genome Analysis"**

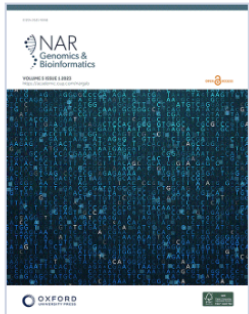
**NAR Genomics and Bioinformatics**, March 2023.

[[Online link at NAR Genomics and Bioinformatics Journal](#)]

[[arXiv preprint](#)]

[[biorXiv preprint](#)]

[[BLEND Source Code](#)]



Volume 5, Issue 1

March 2023

### JOURNAL ARTICLE

## **BLEND: a fast, memory-efficient and accurate mechanism to find fuzzy seed matches in genome analysis**

Can Firtina , Jisung Park, Mohammed Alser, Jeremie S Kim, Damla Senol Cali, Taha Shahroodi, Nika Mansouri Ghiasi, Gagandeep Singh, Konstantinos Kanellopoulos, Can Alkan, Onur Mutlu 

*NAR Genomics and Bioinformatics*, Volume 5, Issue 1, March 2023, lqad004,

# New Applications: Frequent Database Updates

- Jeremie S. Kim\*, Can Firtina\*, M. Banu Cavlak, Damla Senol Cali, Nastaran Hajinazar, Mohammed Alser, Can Alkan, and Onur Mutlu,  
**"AirLift: A Fast and Comprehensive Technique for Remapping Alignments between Reference Genomes"**  
*Proceedings of the 21st Asia Pacific Bioinformatics Conference (APBC)*,  
Changsha, China, April 2023.  
[[AirLift Source Code](#)]  
[[arxiv.org Version \(pdf\)](#)]  
[[Talk Video at BIO-Arch 2023 Workshop](#)]

## METHOD

# AirLift: A Fast and Comprehensive Technique for Remapping Alignments between Reference Genomes

Jeremie S. Kim<sup>1†</sup>, Can Firtina<sup>1†</sup>, Meryem Banu Cavlak<sup>2</sup>, Damla Senol Cali<sup>3</sup>, Nastaran Hajinazar<sup>1,4</sup>, Mohammed Alser<sup>1</sup>, Can Alkan<sup>2</sup> and Onur Mutlu<sup>1,2,3\*</sup>

# Error Correction using ML [Bioinform. 2020]

---

- [Can Firtina](#), [Jeremie S. Kim](#), [Mohammed Alser](#), [Damla Senol Cali](#), [A. Ercument Cicek](#), [Can Alkan](#), and [Onur Mutlu](#),  
**"Apollo: A Sequencing-Technology-Independent, Scalable, and Accurate Assembly Polishing Algorithm"**  
***Bioinformatics***, June 2020.  
[[Source Code](#)]  
[[Online link at Bioinformatics Journal](#)]

## Apollo: a sequencing-technology-independent, scalable and accurate assembly polishing algorithm

FREE

[Can Firtina](#), [Jeremie S Kim](#), [Mohammed Alser](#), [Damla Senol Cali](#), [A Ercument Cicek](#),  
[Can Alkan](#) ✉, [Onur Mutlu](#) ✉

*Bioinformatics*, Volume 36, Issue 12, 15 June 2020, Pages 3669–3679,  
<https://doi.org/10.1093/bioinformatics/btaa179>

**Published:** 13 March 2020    **Article history** ▼

# Accelerating String Matching [MICRO 2020]

- Damla Senol Cali, Gurpreet S. Kalsi, Zülal Bingöl, Can Firtina, Lavanya Subramanian, Jeremie S. Kim, Rachata Ausavarungnirun, Mohammed Alser, Juan Gomez-Luna, Amirali Boroumand, Anant Nori, Allison Scibisz, Sreenivas Subramoney, Can Alkan, Saugata Ghose, and Onur Mutlu, "[GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis](#)"  
*Proceedings of the 53rd International Symposium on Microarchitecture (MICRO)*, Virtual, October 2020.  
[[Lightning Talk Video](#) (1.5 minutes)]  
[[Lightning Talk Slides \(pptx\)](#) ([pdf](#))]  
[[Talk Video](#) (18 minutes)]  
[[Slides \(pptx\)](#) ([pdf](#))]

## GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis

Damla Senol Cali<sup>†</sup><sup>✕</sup> Gurpreet S. Kalsi<sup>✕</sup> Zülal Bingöl<sup>∇</sup> Can Firtina<sup>◇</sup> Lavanya Subramanian<sup>‡</sup> Jeremie S. Kim<sup>◇</sup><sup>†</sup>  
Rachata Ausavarungnirun<sup>○</sup> Mohammed Alser<sup>◇</sup> Juan Gomez-Luna<sup>◇</sup> Amirali Boroumand<sup>†</sup> Anant Nori<sup>✕</sup>  
Allison Scibisz<sup>†</sup> Sreenivas Subramoney<sup>✕</sup> Can Alkan<sup>∇</sup> Saugata Ghose<sup>\*†</sup> Onur Mutlu<sup>◇</sup><sup>†</sup><sup>∇</sup>  
<sup>†</sup>Carnegie Mellon University <sup>✕</sup>Processor Architecture Research Lab, Intel Labs <sup>∇</sup>Bilkent University <sup>◇</sup>ETH Zürich  
<sup>‡</sup>Facebook <sup>○</sup>King Mongkut's University of Technology North Bangkok <sup>\*</sup>University of Illinois at Urbana-Champaign

# Accelerating Genome Graphs [ISCA 2022]

---

- Damla Senol Cali, Konstantinos Kanellopoulos, Joel Lindegger, Zulal Bingol, Gurpreet S. Kalsi, Ziyi Zuo, Can Firtina, Meryem Banu Cavlak, Jeremie Kim, Nika MansouriGhiasi, Gagandeep Singh, Juan Gomez-Luna, Nour Almadhoun Alserr, Mohammed Alser, Sreenivas Subramoney, Can Alkan, Saugata Ghose, and Onur Mutlu,  
**"SeGraM: A Universal Hardware Accelerator for Genomic Sequence-to-Graph and Sequence-to-Sequence Mapping"**  
*Proceedings of the 49th International Symposium on Computer Architecture (ISCA)*, New York, June 2022.  
[\[arXiv version\]](#)

## **SeGraM: A Universal Hardware Accelerator for Genomic Sequence-to-Graph and Sequence-to-Sequence Mapping**

Damla Senol Cali<sup>1</sup> Konstantinos Kanellopoulos<sup>2</sup> Joël Lindegger<sup>2</sup> Zülal Bingöl<sup>3</sup>  
Gurpreet S. Kalsi<sup>4</sup> Ziyi Zuo<sup>5</sup> Can Firtina<sup>2</sup> Meryem Banu Cavlak<sup>2</sup> Jeremie Kim<sup>2</sup>  
Nika Mansouri Ghiasi<sup>2</sup> Gagandeep Singh<sup>2</sup> Juan Gómez-Luna<sup>2</sup> Nour Almadhoun Alserr<sup>2</sup>  
Mohammed Alser<sup>2</sup> Sreenivas Subramoney<sup>4</sup> Can Alkan<sup>3</sup> Saugata Ghose<sup>6</sup> Onur Mutlu<sup>2</sup>

<sup>1</sup>Bionano Genomics   <sup>2</sup>ETH Zürich   <sup>3</sup>Bilkent University   <sup>4</sup>Intel Labs  
<sup>5</sup>Carnegie Mellon University   <sup>6</sup>University of Illinois Urbana-Champaign



# In-Storage Genome Filtering [ASPLOS 2022]

---

- Nika Mansouri Ghiasi, Jisung Park, Harun Mustafa, Jeremie Kim, Ataberk Olgun, Arvid Gollwitzer, Damla Senol Cali, Can Firtina, Haiyu Mao, Nour Almadhoun Alserr, Rachata Ausavarungnirun, Nandita Vijaykumar, Mohammed Alser, and Onur Mutlu, **["GenStore: A High-Performance and Energy-Efficient In-Storage Computing System for Genome Sequence Analysis"](#)**  
*Proceedings of the 27th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, Virtual, February-March 2022.  
[[Lightning Talk Slides \(pptx\) \(pdf\)](#)]  
[[Lightning Talk Video](#) (90 seconds)]

## GenStore: A High-Performance In-Storage Processing System for Genome Sequence Analysis

Nika Mansouri Ghiasi<sup>1</sup> Jisung Park<sup>1</sup> Harun Mustafa<sup>1</sup> Jeremie Kim<sup>1</sup> Ataberk Olgun<sup>1</sup>  
Arvid Gollwitzer<sup>1</sup> Damla Senol Cali<sup>2</sup> Can Firtina<sup>1</sup> Haiyu Mao<sup>1</sup> Nour Almadhoun Alserr<sup>1</sup>  
Rachata Ausavarungnirun<sup>3</sup> Nandita Vijaykumar<sup>4</sup> Mohammed Alser<sup>1</sup> Onur Mutlu<sup>1</sup>

<sup>1</sup>ETH Zürich <sup>2</sup>Bionano Genomics <sup>3</sup>KMUTNB <sup>4</sup>University of Toronto

# Genome Analysis via PIM [MICRO 2022]

---

- Haiyu Mao, Mohammed Alser, Mohammad Sadrosadati, Can Firtina, Akanksha Baranwal, Damla Senol Cali, Aditya Manglik, Nour Almadhoun Alserr, and Onur Mutlu, **"GenPIP: In-Memory Acceleration of Genome Analysis via Tight Integration of Basecalling and Read Mapping"**  
*Proceedings of the 55th International Symposium on Microarchitecture (MICRO)*, Chicago, IL, USA, October 2022.  
[[Slides \(pptx\)](#)] [[pdf](#)]  
[[Longer Lecture Slides \(pptx\)](#)] [[pdf](#)]  
[[Lecture Video](#)] (25 minutes)  
[[arXiv version](#)]

## **GenPIP: In-Memory Acceleration of Genome Analysis via Tight Integration of Basecalling and Read Mapping**

Haiyu Mao<sup>1</sup> Mohammed Alser<sup>1</sup> Mohammad Sadrosadati<sup>1</sup> Can Firtina<sup>1</sup> Akanksha Baranwal<sup>1</sup>  
Damla Senol Cali<sup>2</sup> Aditya Manglik<sup>1</sup> Nour Almadhoun Alserr<sup>1</sup> Onur Mutlu<sup>1</sup>  
<sup>1</sup>*ETH Zürich*      <sup>2</sup>*Bionano Genomics*

# Food Microbiome Profiling using PIM

Taha Shahroodi, Mahdi Zahedi, Can Firtina, Mohammed Alser, Stephan Wong, Onur Mutlu, Said Hamdioui

[“Demeter: A Fast and Energy-Efficient Food Profiler using Hyperdimensional Computing in Memory”](#)

IEEE Access, 2022

**IEEE Access**  
Multidisciplinary | Rapid Review | Open Access Journal

**RESEARCH ARTICLE**

## Demeter: A Fast and Energy-Efficient Food Profiler Using Hyperdimensional Computing in Memory

**TAHA SHAHROODI<sup>ID1</sup>, MAHDI ZAHEDI<sup>ID1</sup>, CAN FIRTINA<sup>2</sup>, MOHAMMED ALSER<sup>ID2</sup>,  
STEPHAN WONG<sup>1</sup>, (Senior Member, IEEE), ONUR MUTLU<sup>ID2</sup>, (Fellow, IEEE),  
AND SAID HAMDIOUI<sup>ID1</sup>, (Senior Member, IEEE)**

<sup>1</sup>Q&CE Department, EEMCS Faculty, Delft University of Technology (TU Delft), 2628 CD Delft, The Netherlands

<sup>2</sup>SAFARI Research Group, D-ITET, ETH Zürich, 8092 Zürich, Switzerland

# Fast and Accurate Real-Time Genome Analysis

---

- Joel Lindegger, Can Firtina, Nika Mansouri Ghiasi, Mohammad Sadrosadati, Mohammed Alser, and Onur Mutlu,  
**"RawAlign: Accurate, Fast, and Scalable Raw Nanopore Signal Mapping via Combining Seeding and Alignment"**  
*Preprint on **arxiv**, October 2023.*  
[\[arXiv version\]](#)  
[\[RawAlign Source Code\]](#)

## **RawAlign: Accurate, Fast, and Scalable Raw Nanopore Signal Mapping via Combining Seeding and Alignment**

Joël Lindegger<sup>§</sup>      Can Firtina<sup>§</sup>      Nika Mansouri Ghiasi<sup>§</sup>  
Mohammad Sadrosadati<sup>§</sup>      Mohammed Alser<sup>§</sup>      Onur Mutlu<sup>§</sup>  
*§ETH Zürich*

# Machine Learning in Genomics

---

- M. Banu Cavlak, Gagandeep Singh, Mohammed Alser, Can Firtina, Joel Lindegger, Mohammad Sadrosadati, Nika Mansouri Ghiasi, Can Alkan, and Onur Mutlu, **"TargetCall: Eliminating the Wasted Computation in Basecalling via Pre-Basecalling Filtering"**  
*Proceedings of the 21st Asia Pacific Bioinformatics Conference (APBC), Changsha, China, April 2023.*  
[[TargetCall Source Code](#)]  
[[arxiv.org Version](#)]  
[[Talk Video at BIO-Arch 2023 Workshop](#)]

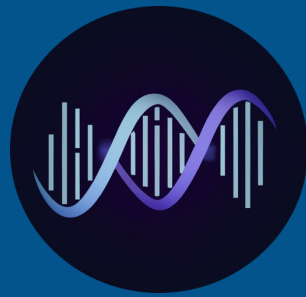
## **TargetCall: Eliminating the Wasted Computation in Basecalling via Pre-Basecalling Filtering**

Meryem Banu Cavlak<sup>1</sup> Gagandeep Singh<sup>1</sup> Mohammed Alser<sup>1</sup> Can Firtina<sup>1</sup> Joël Lindegger<sup>1</sup>  
Mohammad Sadrosadati<sup>1</sup> Nika Mansouri Ghiasi<sup>1</sup> Can Alkan<sup>2</sup> Onur Mutlu<sup>1</sup>  
<sup>1</sup>*ETH Zürich*                      <sup>2</sup>*Bilkent University*

# Agenda for Today

---

- Cutting-edge in Accelerating Genome Analysis
  - Intelligent genome analysis
  
- Enabling Fast and Accurate Real-time Analysis
  - RawHash and RawHash2
  
- Conclusion



# RawHash

Enabling Fast and Accurate Real-Time Analysis  
of Raw Nanopore Signals for Large Genomes

**Can Firtina**

Nika Mansouri Ghiasi

Joel Lindegger

Gagandeep Singh

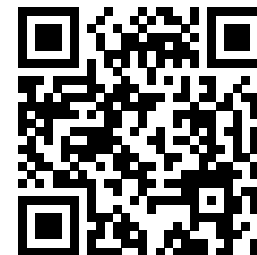
Meryem Banu Cavlak

Haiyu Mao

Onur Mutlu



[Paper](#)



[Code](#)

# Nanopore Sequencing

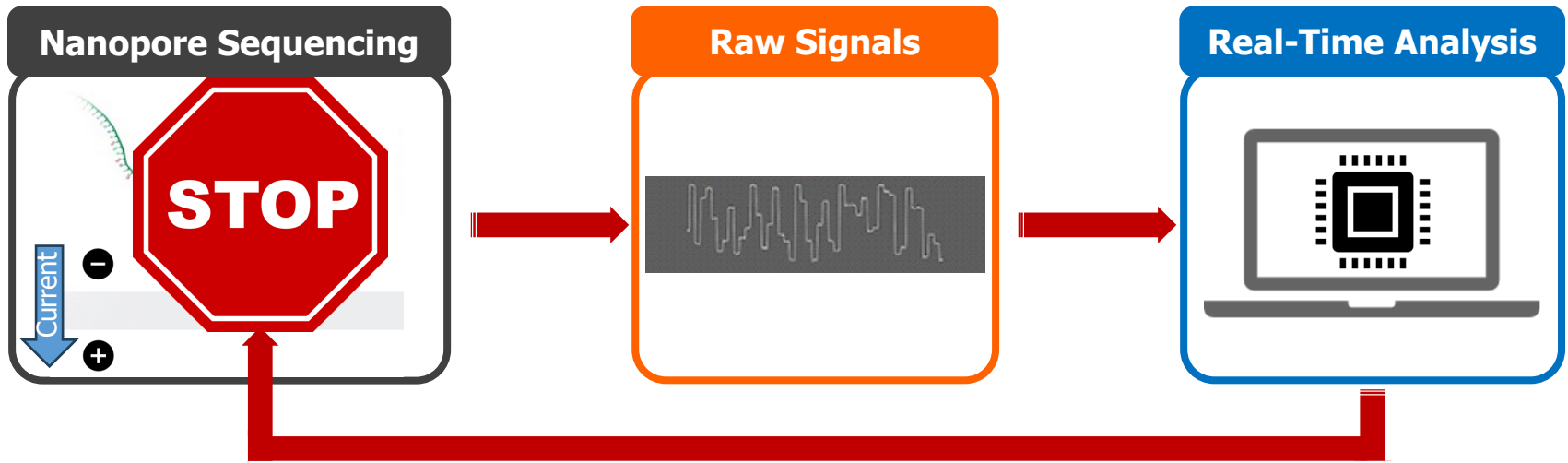
**Nanopore Sequencing:** a widely used sequencing technology

- Can sequence large fragments of nucleic acid molecules (up to >2Mbp)
- Offers high throughput
- Cost-effective
- Enables **real-time genome analysis**





# Real-Time Analysis with Nanopore Sequencing



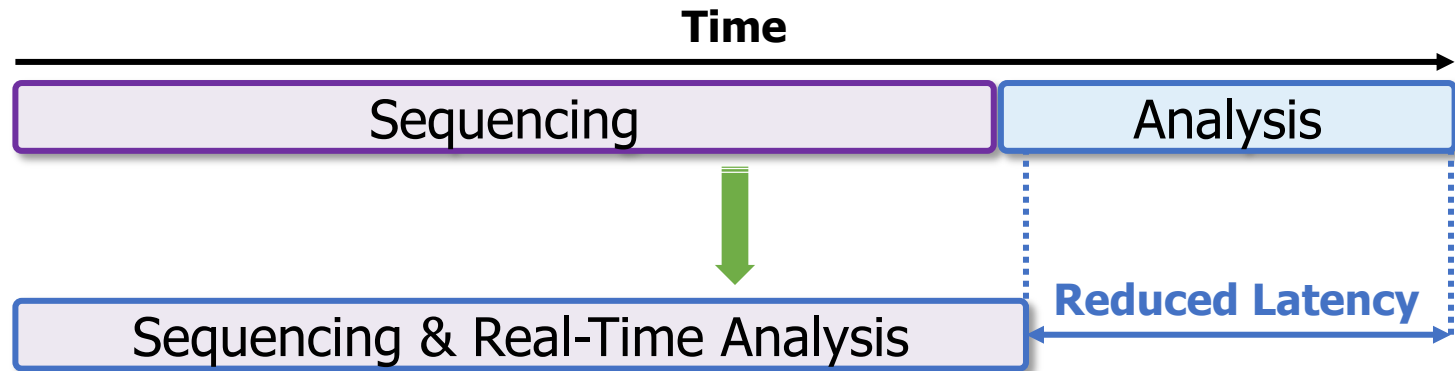
**Raw Signals:** Ionic current measurements generated at a certain **throughput**

**Real-Time Analysis:** Analyzing all raw signals by **matching the throughput**

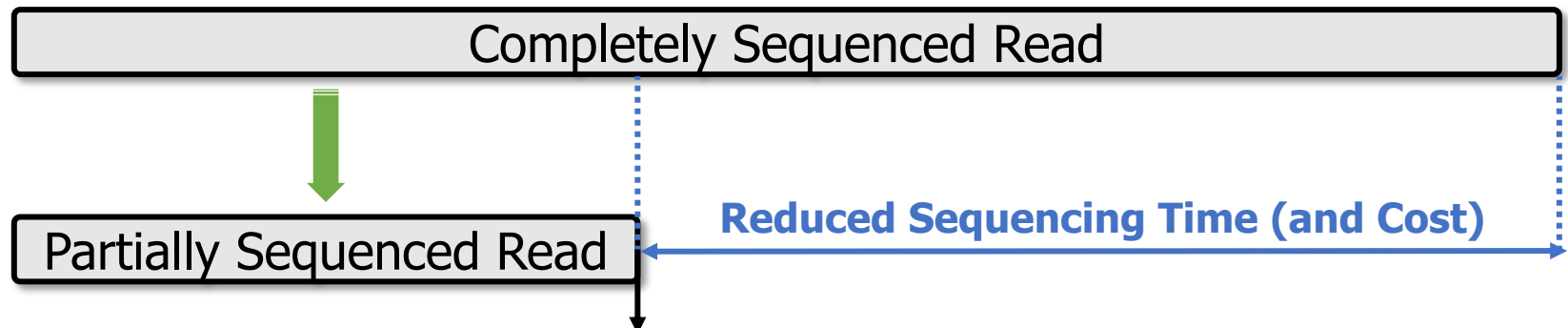
**Real-Time Decisions:** Stopping sequencing **early** based on real-time analysis

# Benefits of Real-Time Genome Analysis

- ✓ **Reducing latency** by overlapping the sequencing and analysis steps



- ✓ **Reducing sequencing time and cost** by stopping sequencing early



Sequencing is stopped early with a real-time decision

# Challenges in Real-Time Genome Analysis

 **Rapid analysis** to match the nanopore sequencer throughput

 **Timely decisions** to stop sequencing as early as possible

 **Accurate analysis** from noisy raw signal data

 **Power-efficient** computation for scalability and portability

# Executive Summary

**Problem:** Real-time analysis of nanopore raw signals is **inaccurate** and **inefficient for large genomes**

**Goal:** Enable **fast** and **accurate** real-time analysis of raw signals for **large genomes**

## Key Contributions:

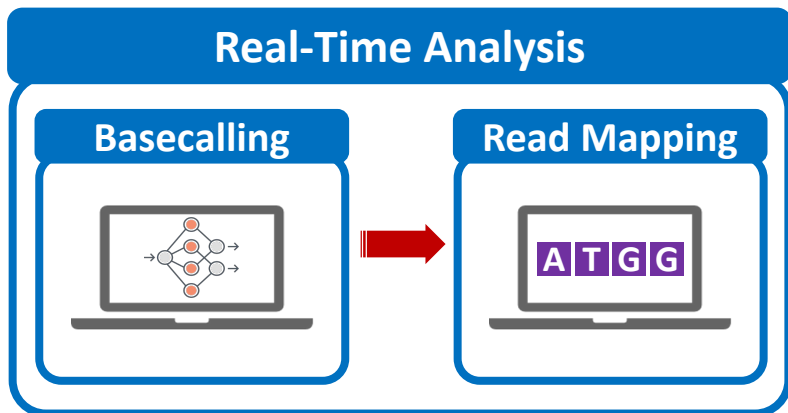
- 1) The **first hash-based mechanism** that can quickly and accurately analyze raw nanopore signals for **large genomes**
- 2) The novel **Sequence Until** technique can accurately and **dynamically stop the entire sequencing of all reads at once** if further sequencing is not necessary

**Key Results:** Across 3 use cases and 5 genomes of varying sizes, RawHash provides

- **25.8× and 3.4× better average throughput** compared to two state-of-the-art works
- **1.14× – 2.13× more accurate mapping results for large genomes**
- Sequence Until **reduces the sequencing time and cost by 15×**

# Existing Solutions

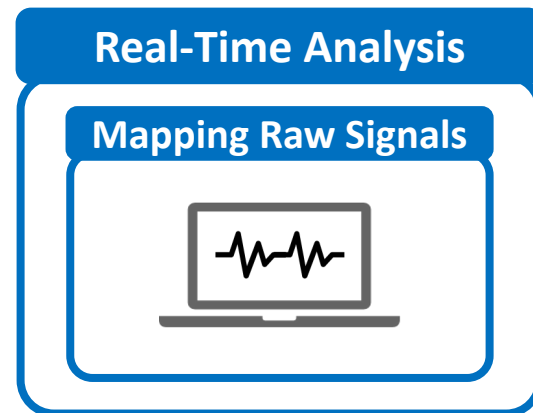
1. Deep neural networks (**DNNs**) for translating **signals** to **bases**



Less noisy analysis from basecalled sequences

**Costly and power-hungry** computational requirements

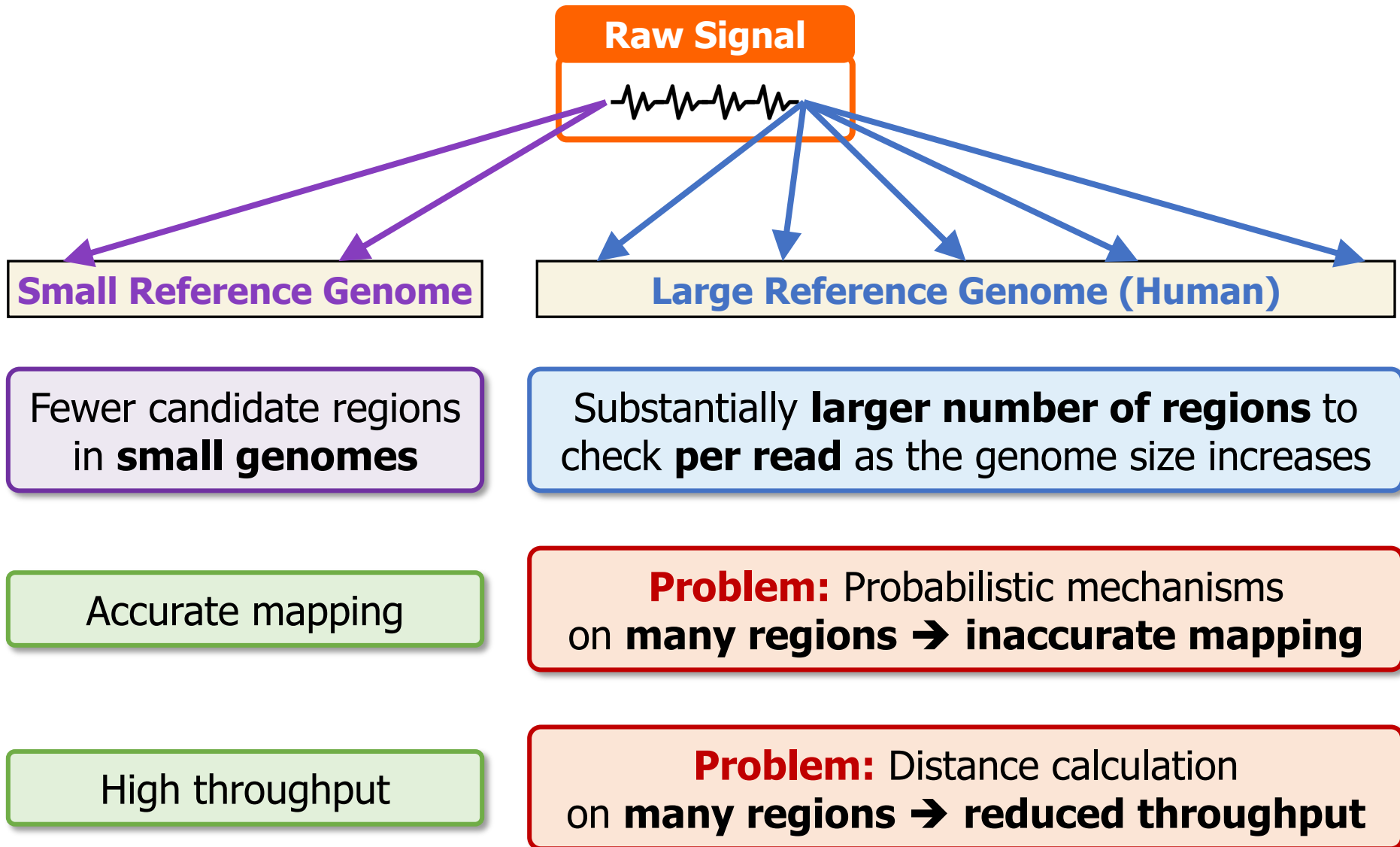
2. Mapping **signals** to reference genomes **without** basecalling



Raw signals contain richer information than bases

Efficient analysis with better scalability and portability

# The Problem – Mapping Raw Signals



# The Problem – Mapping Raw Signals



Existing solutions are  
**inaccurate or inefficient**  
**for large genomes**

Accurate mapping

on many regions → inaccurate mapping

High throughput

**Problem:** Distance calculation  
on many regions → reduced throughput

# Outline

Background

RawHash

Evaluation

Conclusion



# Goal

Enable **fast and accurate real-time analysis**  
of raw nanopore signals **for large genomes**



# RawHash

The **first hash-based search mechanism** to quickly and accurately map raw nanopore signals to reference genomes

**Sequence Until** can accurately and **dynamically stop the entire sequencing run at once** if further sequencing is unnecessary



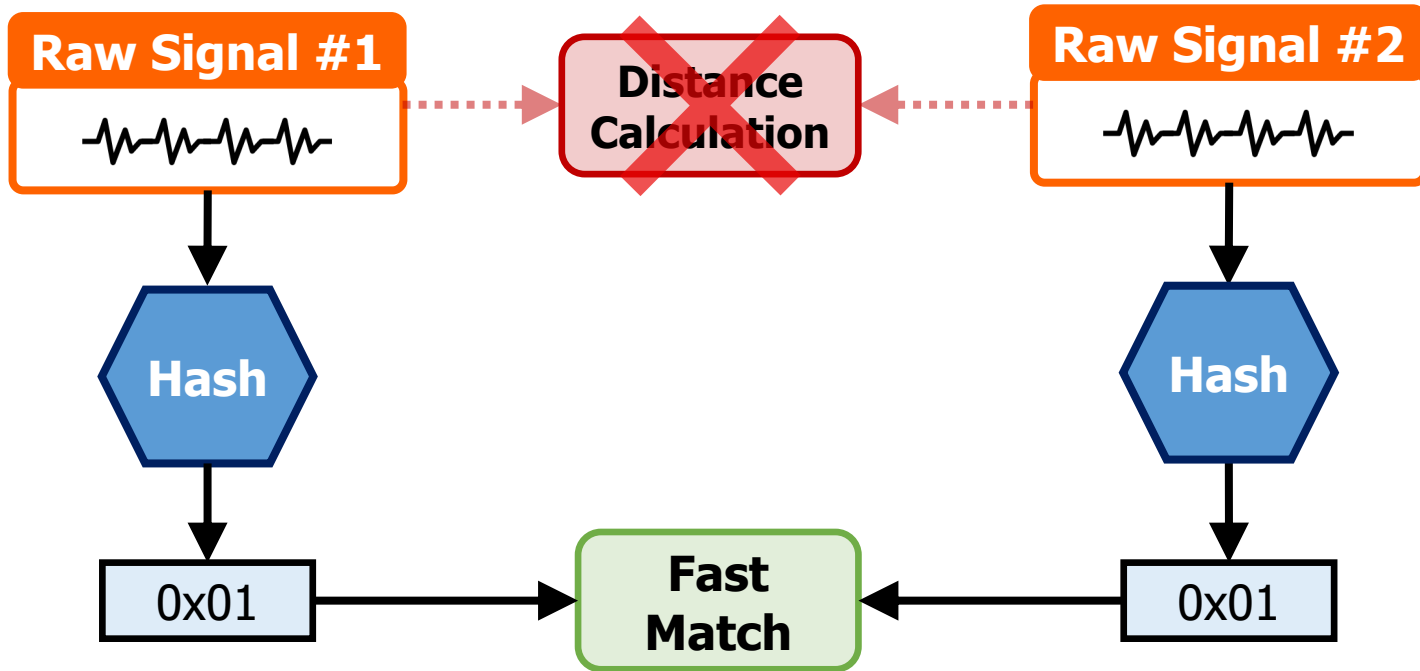
# RawHash

The **first hash-based search mechanism** to quickly and accurately map raw nanopore signals to reference genomes

**Sequence Until** can accurately and **dynamically stop** the entire sequencing run at once if further sequencing is unnecessary

# RawHash – Key Idea

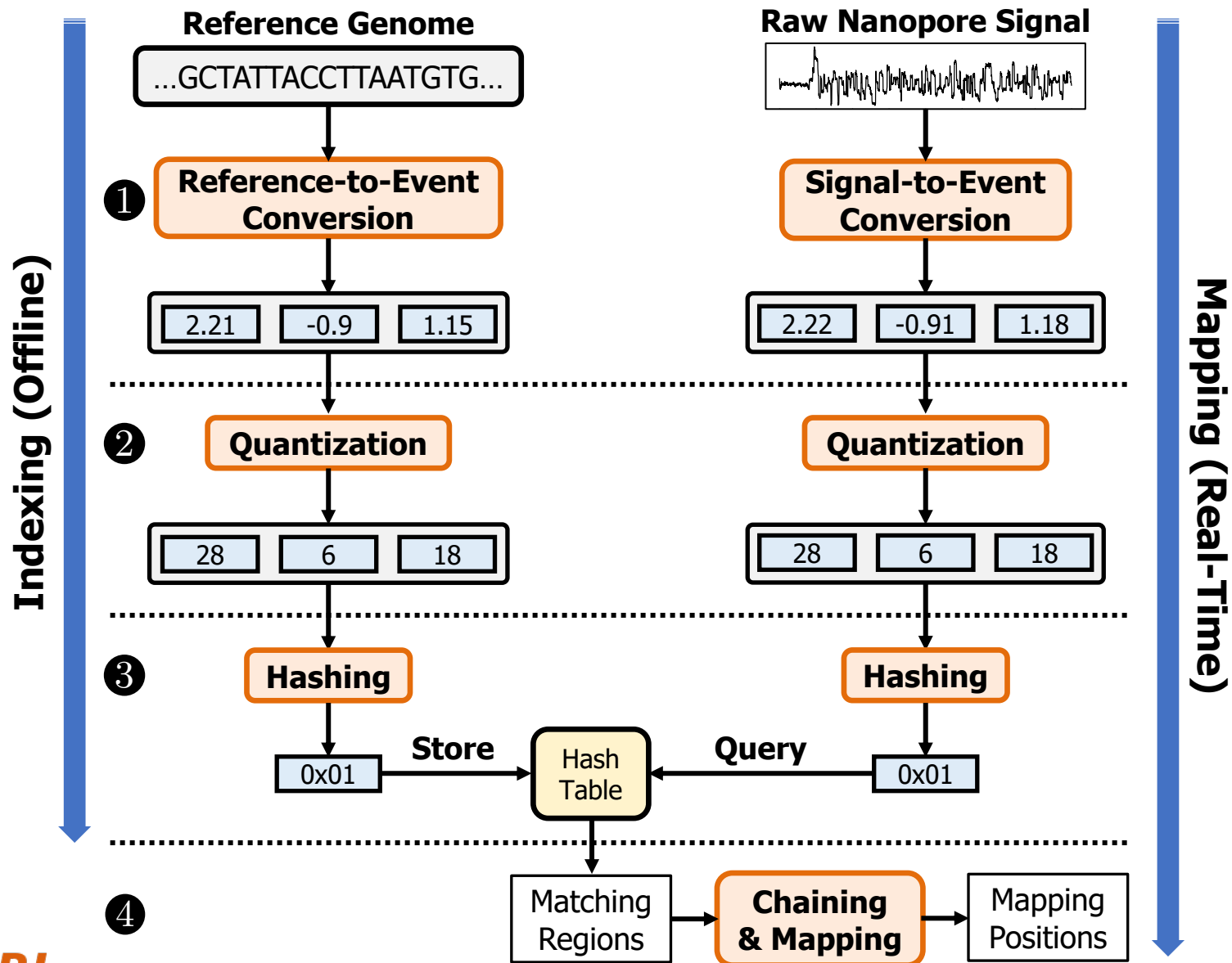
**Key Observation:** **Identical** nucleotides generate **similar** raw signals



**Challenge #1:** Generating the **same** hash value for **similar enough** signals

**Challenge #2:** **Accurately** finding similar regions **as few as possible**

# RawHash Overview

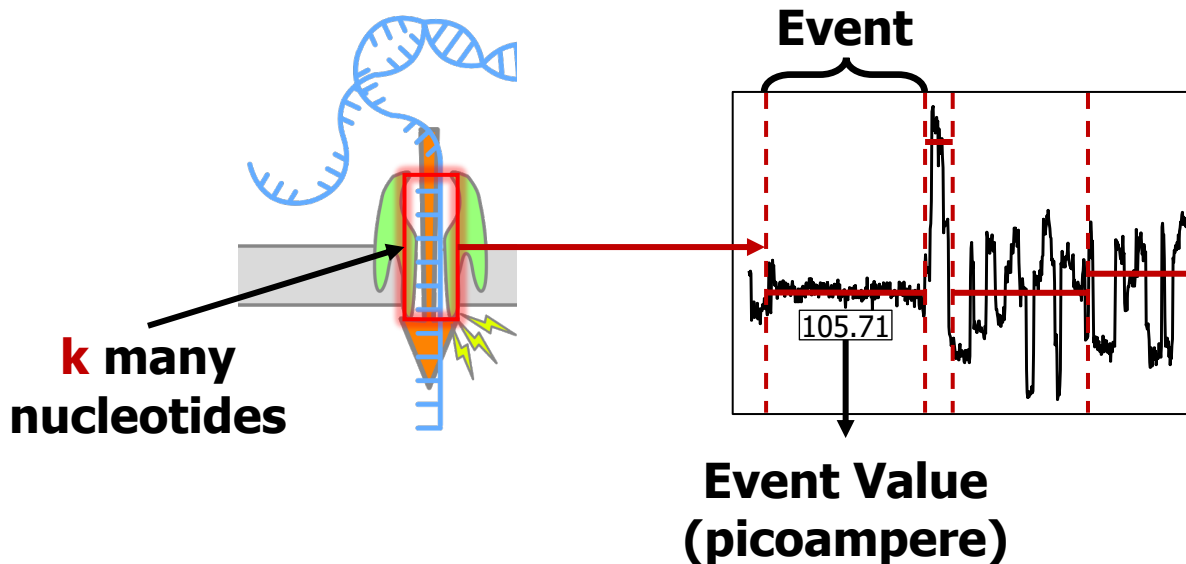


# RawHash Overview



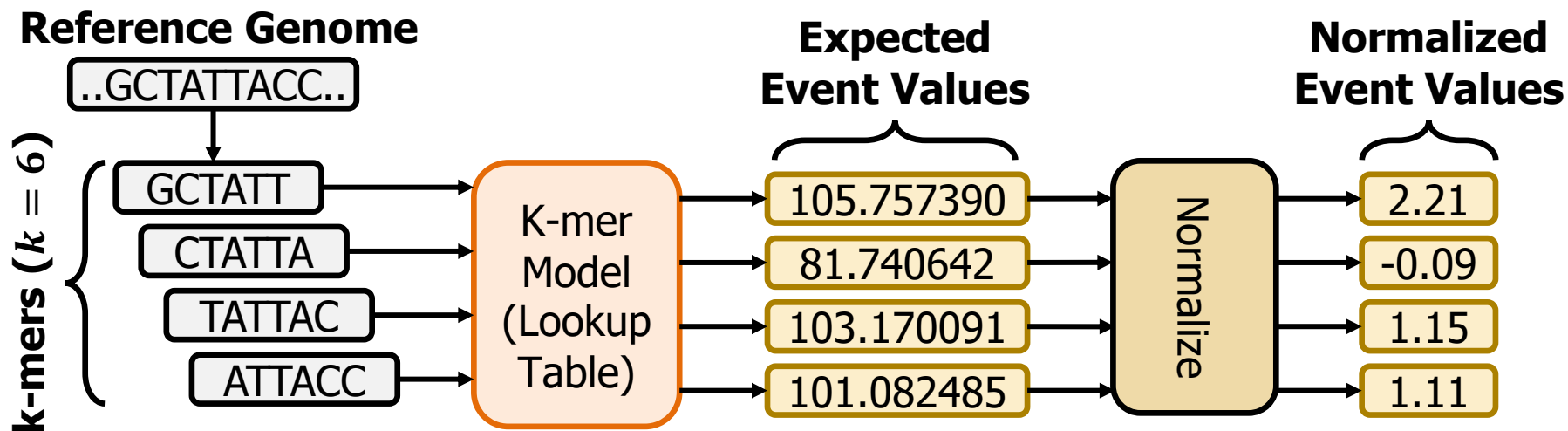
# Events in Raw Nanopore Signals

- **Event:** A **segment** of the raw signal
  - Corresponds to a **particular k-mer**
- **Event detection** finds these segments to identify **k-mers**
  - Start and end positions are marked by abrupt signal changes
  - Statistical methods identify these abrupt changes
  - **Event value:** average of signals **within an event**



# Reference-to-Event Conversion

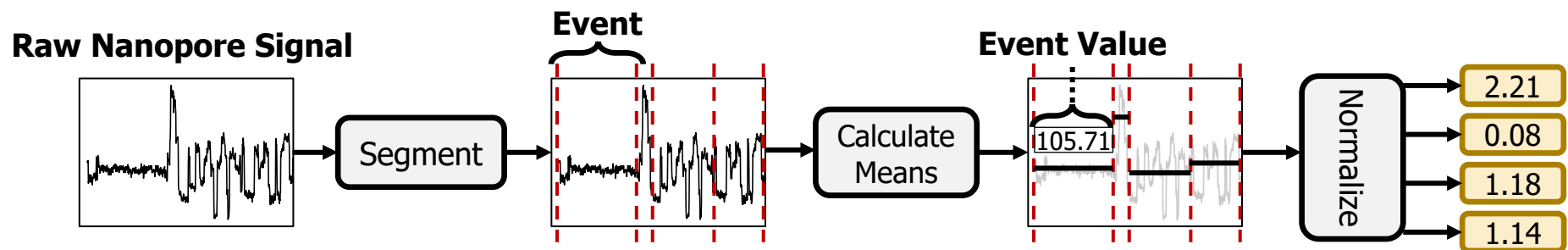
- **K-mer model:** Provides **expected** event values **for each k-mer**
  - Preconstructed based on nanopore sequencer characteristics
- Use the **k-mer model** to convert **all k-mers** of a reference genome to their **expected** event values





# Signal-to-Event Conversion

- **Event detection:** Identifies signal regions corresponding to specific k-mers
  - Uses statistical test (**segmentation**) to spot abrupt signal changes



- Consecutive events → consecutive k-mers

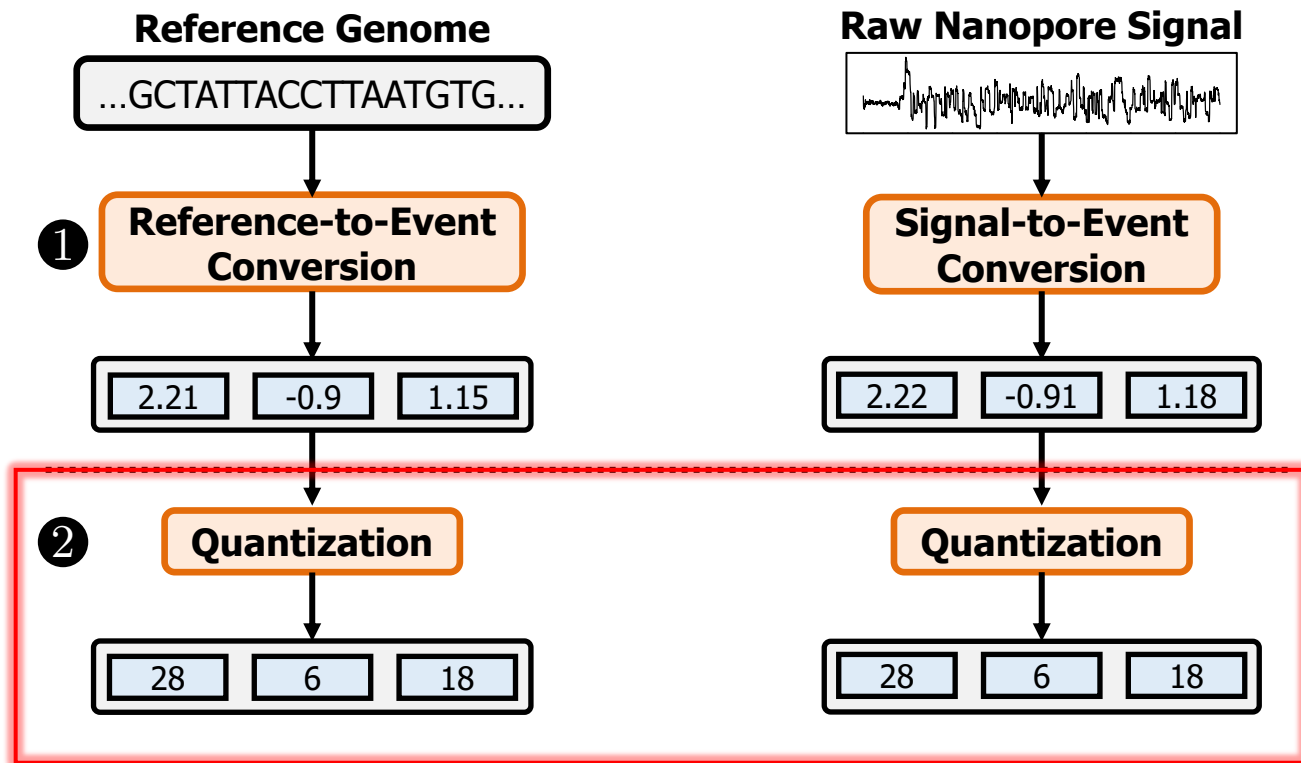
# Signal-to-Event Conversion

- **Event detection:** Identifies signal regions corresponding to specific k-mers
  - Uses statistical test (**segmentation**) to spot abrupt signal changes

Can we directly match signals to each other?

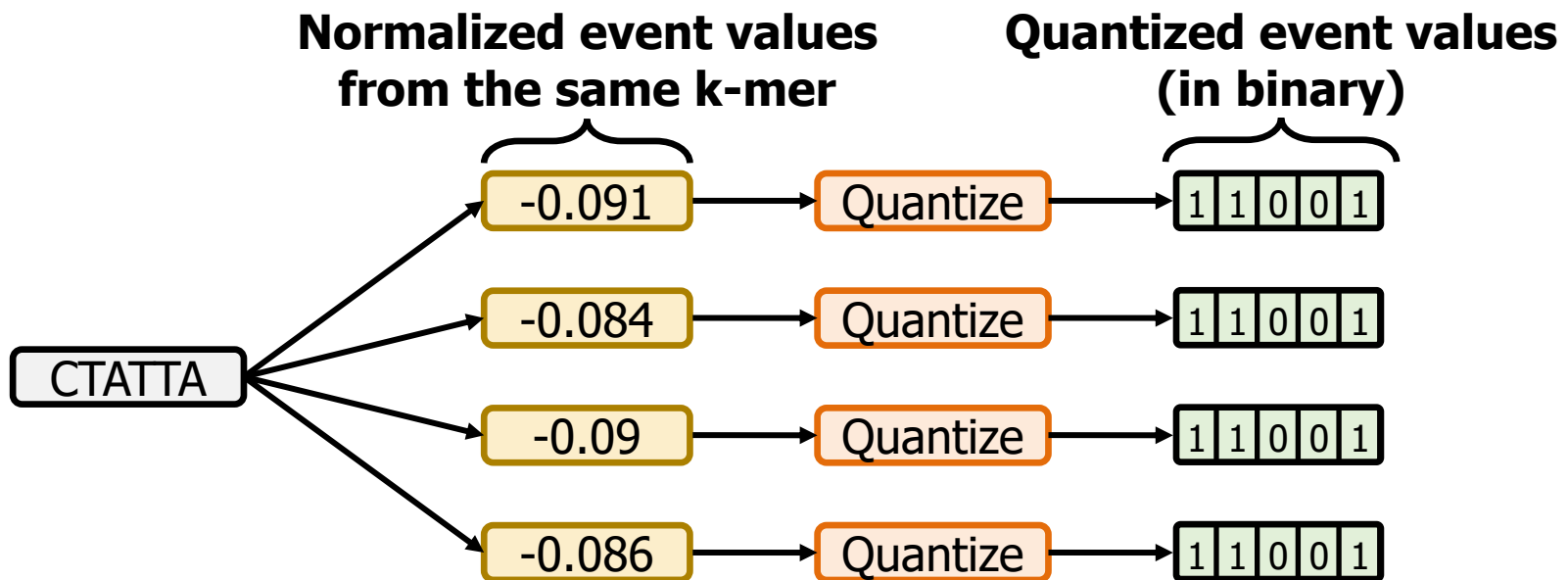
- Consecutive events → consecutive k-mers

# RawHash Overview

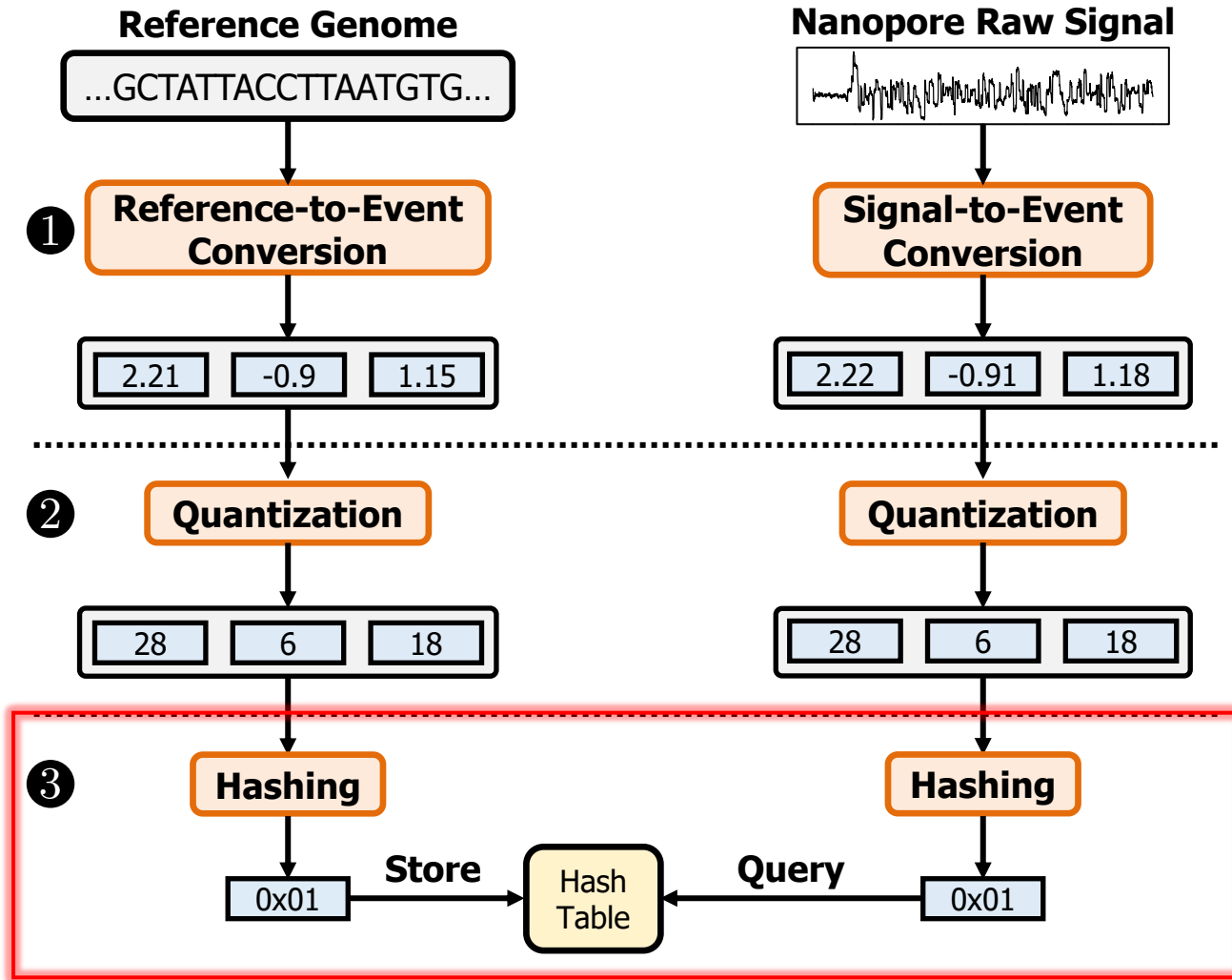


# Quantizing the Event Values

- **Observation:** Slight differences in raw signals from identical k-mers
  - **Challenge:** Direct event value matching is not feasible and accurate
- **Key Idea:** Quantize the event values
  - Enables assigning **identical quantized values** to **similar event values**

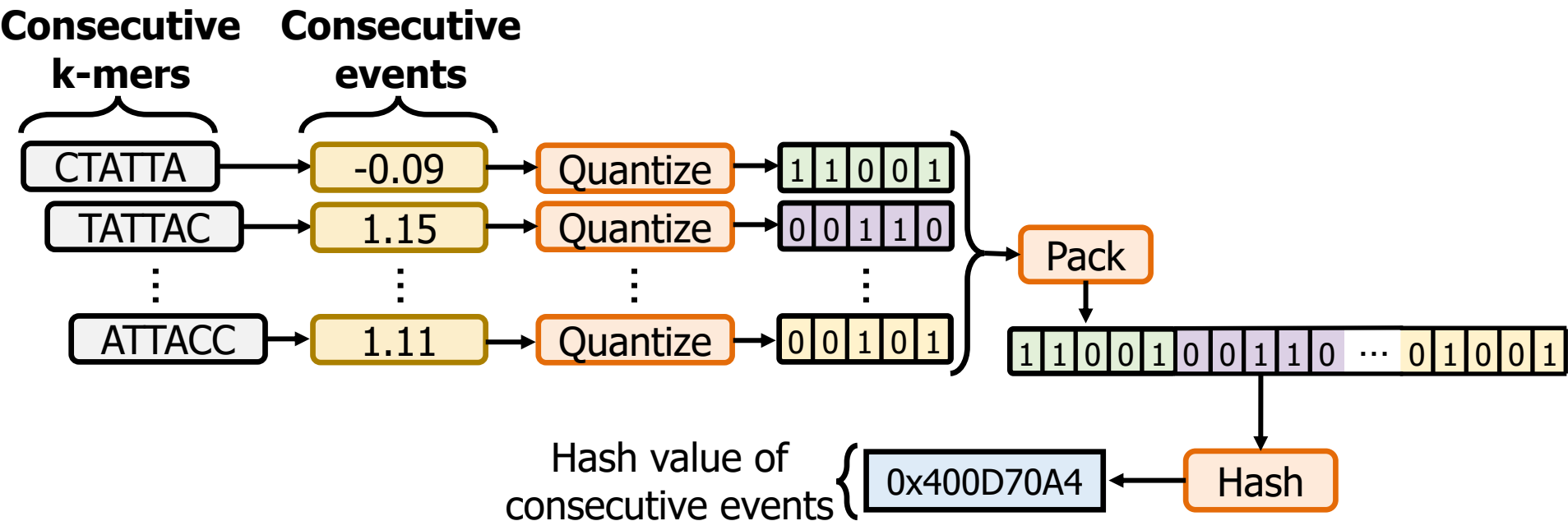


# RawHash Overview

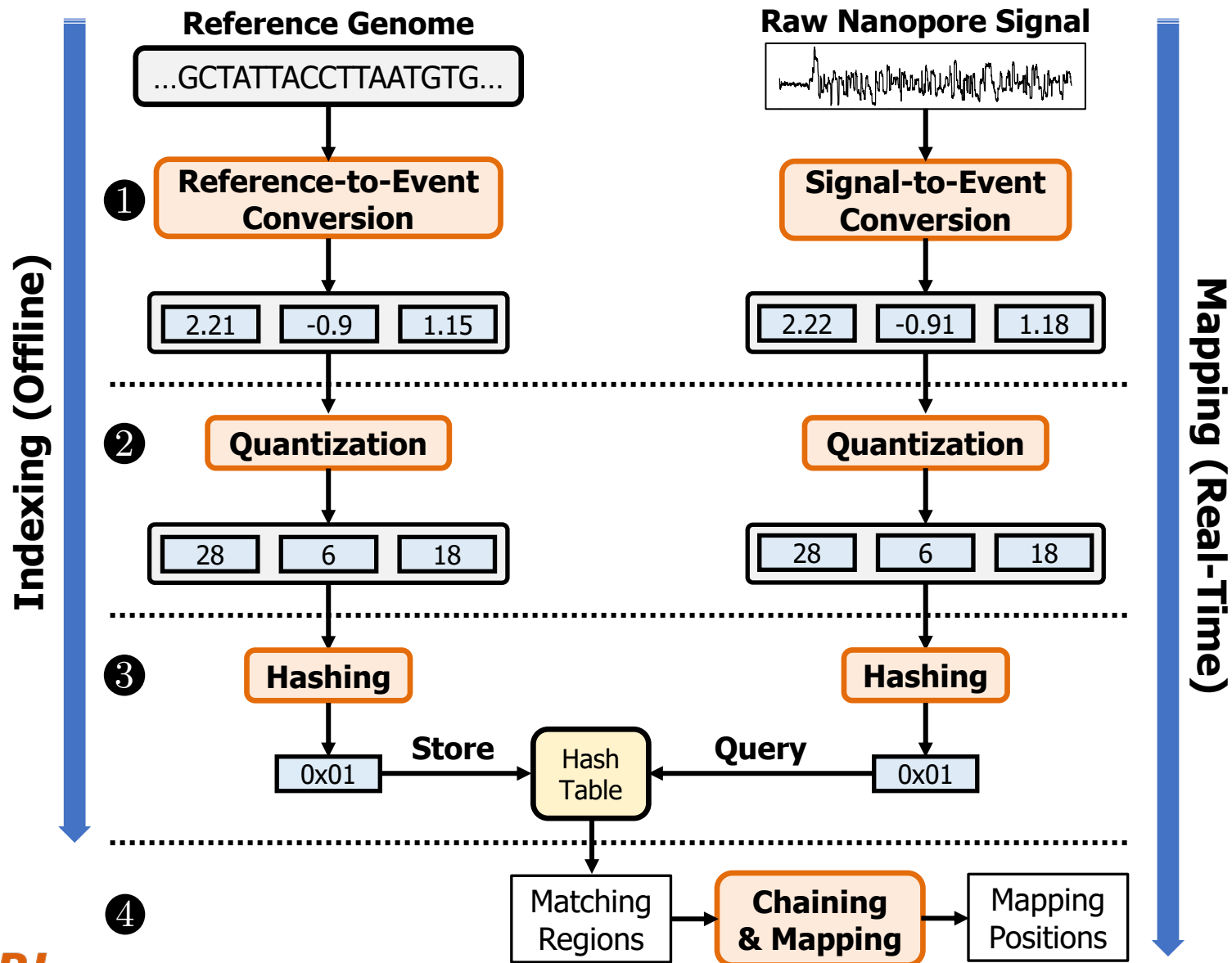


# Hashing for Fast Similarity Search

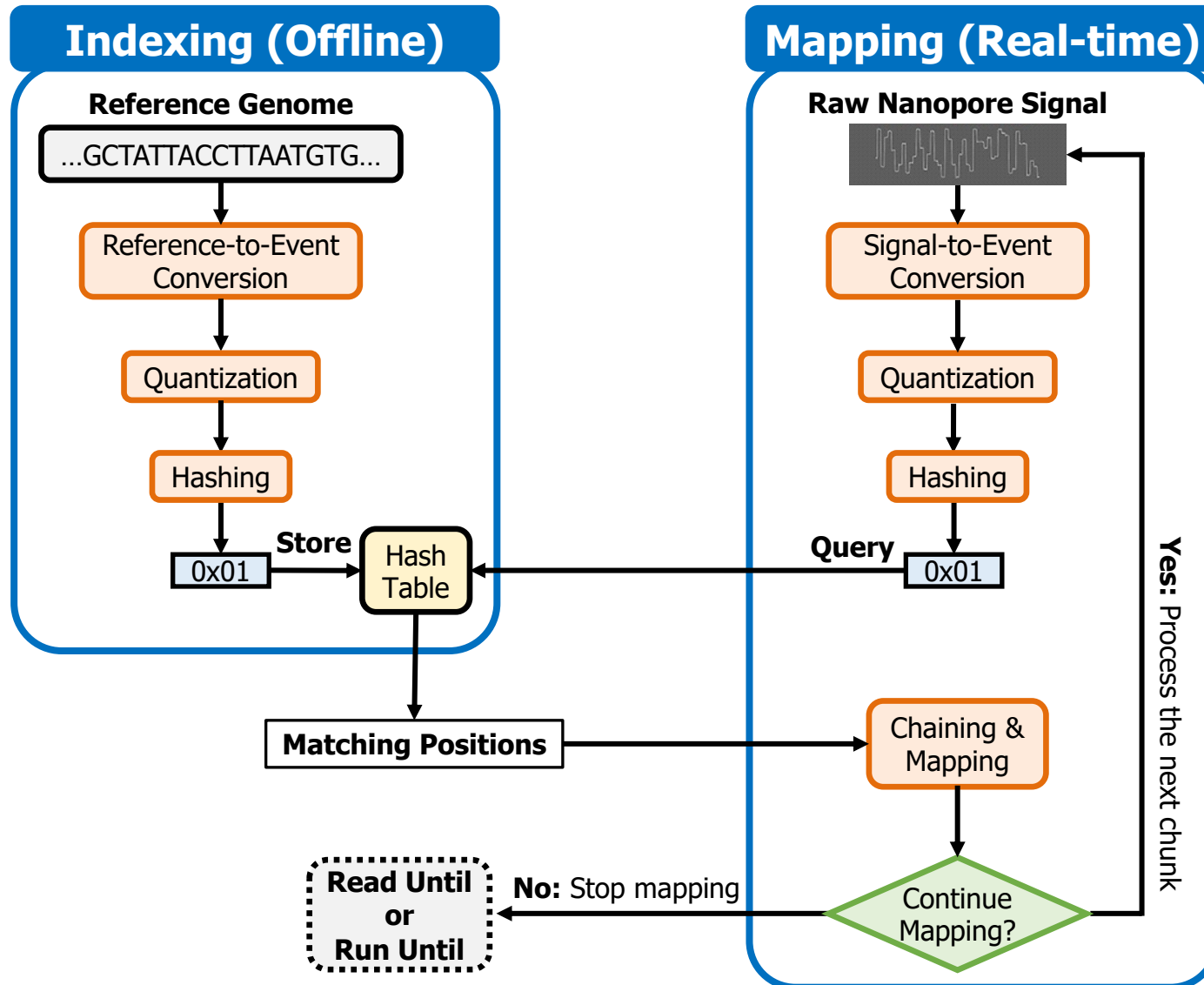
- Each event usually represents a very small k-mer (6 to 9 characters)
  - **Challenge:** Short k-mers are likely to appear in many locations
- **Key Idea:** Create longer k-mers from many **consecutive events**
- **Key Benefit:** Directly match hash values to quickly identify similarities



# RawHash Overview



# Real-Time Mapping using Hash-based Indexing







# RawHash

The **first hash-based search mechanism** to quickly and accurately map raw nanopore signals to reference genomes

**Sequence Until** can accurately and **dynamically stop** the entire sequencing run at once if further sequencing is unnecessary



# RawHash

The **first hash-based search mechanism** to quickly and accurately map raw nanopore signals to reference genomes

**Sequence Until** can accurately and **dynamically stop the entire sequencing run at once** if further sequencing is unnecessary

# The Sequence Until Mechanism

- **Problem:**

- Unnecessary sequencing waste time, power and money

- **Key Idea:**

- **Dynamically** decide if further sequencing of the entire sample is necessary to achieve high accuracy
- Stop sequencing early without sacrificing accuracy

- **Potential Benefits:**

- Significant **reduction in sequencing time and cost**

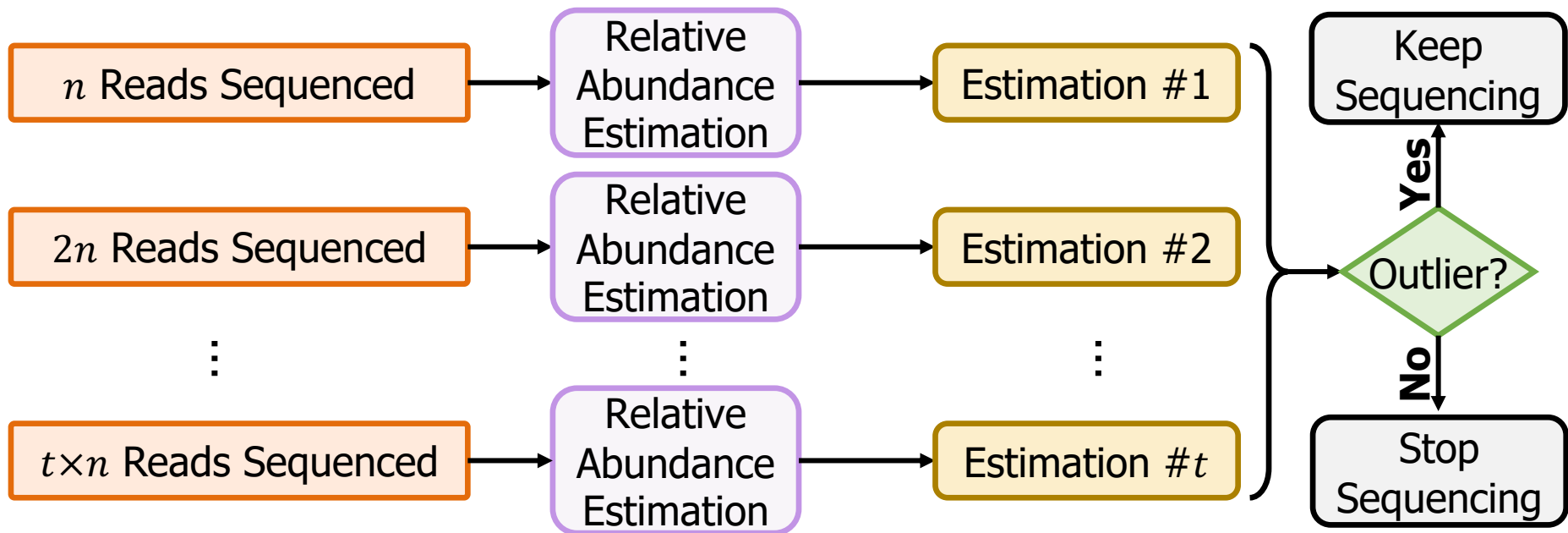
- Example real-time genome analysis use case:

- **Relative abundance estimation**

# The Sequence Until Mechanism

- **Key Steps:**

1. Continuously generate relative abundance estimation after every  $n$  reads
2. Keep the last  $t$  estimation results
3. **Detect outliers** in the results via **cross-correlation** of the recent  $t$  results
4. Absence of outliers indicates **consistent results**
  - Further sequencing **is likely** to generate consistent results → Stop the sequencing



# Outline

Background

RawHash

Evaluation

Conclusion

# Evaluation Methodology

- Compared to **UNCALLED** [Kovaka+, Nat. Biotech. 2021] and **Sigmap** [Zhang+, ISMB/ECCB 2021]
  - **CPU baseline:** AMD EPYC 7742 @2.26GHz
  - **32 threads** for each tool
  
- **Use cases** for real-time genome analysis:
  1. Read mapping
  2. Relative abundance estimation
    - **Benefits of Sequence Until**
  3. Contamination analysis

# Evaluation Methodology

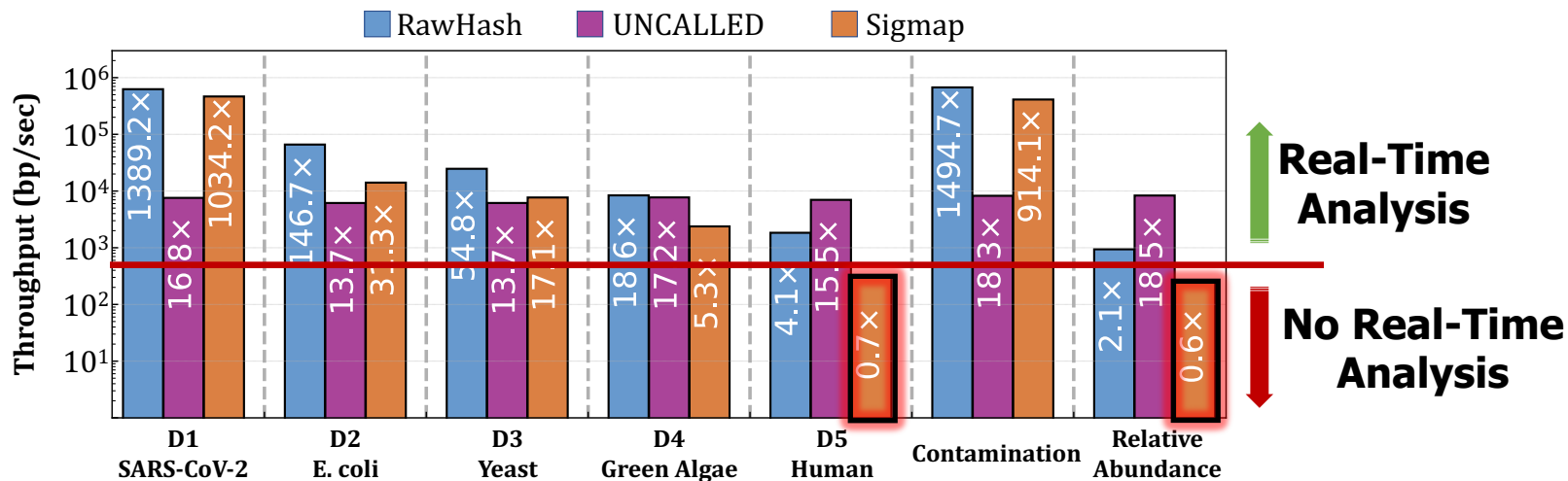
- Evaluation metrics:
  - **Throughput** (bases processed per second)
  - Potential reduction in **sequencing time and cost**
  - **Accuracy**
    - **Baseline:** Mapping basecalled reads using minimap2
    - Precision, recall, and F1 scores
    - Relative abundance estimation distance to ground truth

- **Datasets:**

	Organism	Reads (#)	Bases (#)	Genome Size
<b>Read Mapping</b>				
D1	<i>SARS-CoV-2</i>	1,382,016	594M	29,903
D2	<i>E. coli</i>	353,317	2,365M	5M
D3	<i>Yeast</i>	49,989	380M	12M
D4	<i>Green Algae</i>	29,933	609M	111M
D5	<i>Human HG001</i>	269,507	1,584M	3,117M
<b>Relative Abundance Estimation</b>				
	D1-D5	2,084,762	5,531M	3,246M
<b>Contamination Analysis</b>				
	D1 and D5	1,651,523	2,178M	29,903

# Throughput

- **Real-time analysis requires** faster throughput than sequencer
  - Throughput of a nanopore sequencer: **~450 bp/sec (data generation speed)**



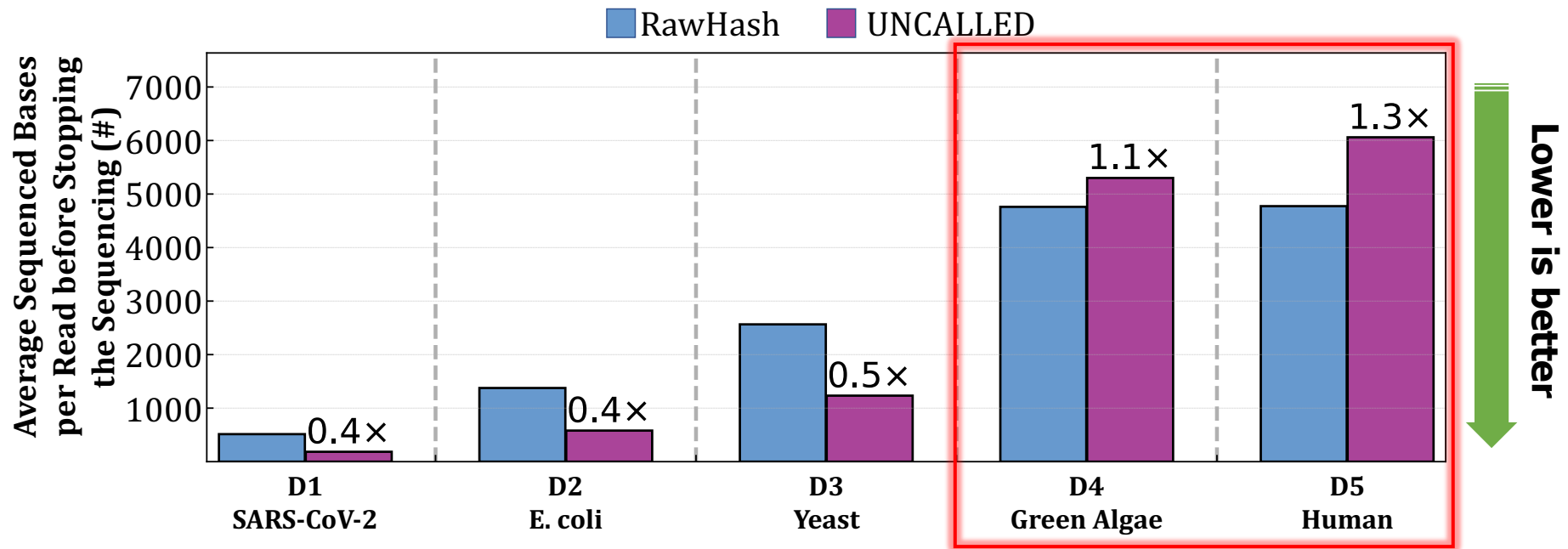
**25.8x and 3.4x** better average throughput compared to **UNCALLED and Sigmap**, respectively

Sigmap **cannot** perform real-time analysis **for large genomes**



# Sequencing Time

- Fewer bases to sequence →
  - Reduction in sequencing time and cost



RawHash **reduces sequencing time and cost**

**for large genomes up to 1.3x compared to UNCALLED**

# Mapping Accuracy

- Read mapping accuracy of each tool and each use case

Dataset		UNCALLED	Sigmap	RawHash
Read Mapping				
D1 <i>SARS-CoV-2</i>	Precision	0.9547	<b>0.9929</b>	0.9868
	Recall	<b>0.9910</b>	0.5540	0.8735
	$F_1$	<b>0.9725</b>	0.7112	0.9267
D2 <i>E. coli</i>	Precision	0.9816	<b>0.9842</b>	0.9573
	Recall	<b>0.9647</b>	0.9504	0.9009
	$F_1$	<b>0.9731</b>	0.9670	0.9282
D3 <i>Yeast</i>	Precision	0.9459	0.9856	<b>0.9862</b>
	Recall	<b>0.9366</b>	0.9123	0.8412
	$F_1$	0.9412	<b>0.9475</b>	0.9079
D4 <i>Green Algae</i>	Precision	0.8836	<b>0.9741</b>	0.9691
	Recall	0.7778	<b>0.8987</b>	0.7015
	$F_1$	0.8273	<b>0.9349</b>	0.8139
D5 <i>Human HG001</i>	Precision	0.4867	0.4287	<b>0.8959</b>
	Recall	0.2379	0.2641	<b>0.4054</b>
	$F_1$	0.3196	0.3268	<b>0.5582</b>

Dataset		UNCALLED	Sigmap	RawHash
Relative Abundance Estimation				
D1-D5	Precision	0.7683	0.7928	<b>0.9484</b>
	Recall	0.1273	0.2739	<b>0.3076</b>
	$F_1$	0.2184	0.4072	<b>0.4645</b>
Contamination Analysis				
D1, D5	Precision	<b>0.9378</b>	0.7856	0.8733
	Recall	<b>0.9910</b>	0.5540	0.8735
	$F_1$	<b>0.9637</b>	0.6498	0.8734

**For Large Genomes:** RawHash provides the **best accuracy**

in all metrics, resulting in **1.14× - 2.13×** improvement in  $F_1$  score

# Relative Abundance Estimation Accuracy

- Estimating the ratio of genomes in a sample in real-time
  - **Distance:** Euclidean distance compared to the ground truth distance
  - The dataset includes a large reference genome

Tool	Estimated Relative Abundance Ratios					Distance
	<i>SARS-CoV-2</i>	<i>E. coli</i>	<i>Yeast</i>	<i>Green Algae</i>	<i>Human</i>	
Ground Truth	0.0929	0.4365	0.0698	0.1179	0.2828	N/A
UNCALLED	0.0026	0.5884	0.0615	0.1313	0.2161	0.1895
Sigmap	0.0419	0.4191	0.1038	0.0962	0.3390	0.0877
RawHash	0.1249	0.4701	0.0957	0.0629	0.2464	<b>0.0847</b>

RawHash provides the **best relative abundance estimation** closest to the ground truth estimation

# Real Implementation of Sequence Until

- Running RawHash by using
  - **RawHash (100%)**: The entire sample **without Sequence Until**
  - **RawHash (7%)**: RawHash **with Sequence Until** where Sequence Until dynamically stops the entire sequencing after sequencing **7% of the sample**

Tool	Estimated Relative Abundance Ratios in 50,000 Random Reads					
	<i>SARS-CoV-2</i>	<i>E. coli</i>	<i>Yeast</i>	<i>Green Algae</i>	<i>Human</i>	Distance
RawHash (100%)	0.0270	0.3636	0.3062	0.1951	0.1081	N/A
RawHash + Sequence Until (7%)	0.0283	0.3539	0.3100	0.1946	0.1133	0.0118

Sequence Until enables sequencing **only 7% (~1/15)**  
of the entire sample **with high accuracy**

# Simulating Sequence Until

- Real relative abundance results using the entire set of reads

Tool	Estimated Relative Abundance Ratios					
	<i>SARS-CoV-2</i>	<i>E. coli</i>	<i>Yeast</i>	<i>Green Algae</i>	<i>Human</i>	Distance
Ground Truth	0.0929	0.4365	0.0698	0.1179	0.2828	N/A
UNCALLED	0.0026	0.5884	0.0615	0.1313	0.2161	0.1895
Sigmap	0.0419	0.4191	0.1038	0.0962	0.3390	0.0877
RawHash	0.1249	0.4701	0.0957	0.0629	0.2464	<b>0.0847</b>

- Simulating the benefits of Sequence Until by
  - Using **a random portion** (25%, 10%, 1%, ...) of the sample

Tool	Estimated Relative Abundance Ratios					
	<i>SARS-CoV-2</i>	<i>E. coli</i>	<i>Yeast</i>	<i>Green Algae</i>	<i>Human</i>	Distance
Ground Truth	0.0929	0.4365	0.0698	0.1179	0.2828	N/A
UNCALLED (25%)	0.0026	0.5890	0.0613	0.1332	0.2139	0.1910
RawHash (25%)	0.0271	0.4853	0.0920	0.0786	0.3170	<b>0.0995</b>
UNCALLED (10%)	0.0026	0.5906	0.0611	0.1316	0.2141	0.1920
RawHash (10%)	0.0273	0.4869	0.0963	0.0772	0.3124	<b>0.1004</b>
UNCALLED (1%)	0.0026	0.5750	0.0616	0.1506	0.2103	0.1836
RawHash (1%)	0.0259	0.4783	0.0987	0.0882	0.3088	<b>0.0928</b>
UNCALLED (0.1%)	0.0040	0.4565	0.0380	0.1910	0.3105	0.1242
RawHash (0.1%)	0.0212	0.5045	0.1120	0.0810	0.2814	<b>0.1136</b>
UNCALLED (0.01%)	0.0000	0.5551	0.0000	0.0000	0.4449	0.2602
RawHash (0.01%)	0.0906	0.6122	0.0000	0.0000	0.2972	<b>0.2232</b>

# Simulating Sequence Until

- Real relative abundance results using the entire set of reads

Tool	Estimated Relative Abundance Ratios					Distance
	<i>SARS-CoV-2</i>	<i>E. coli</i>	<i>Yeast</i>	<i>Green Algae</i>	<i>Human</i>	
Ground Truth	0.0929	0.4365	0.0698	0.1179	0.2828	N/A
UNCALLED	0.0026	0.5884	0.0615	0.1313	0.2161	0.1895
Sigmap	0.0419	0.4191	0.1038	0.0962	0.3390	0.0877

## UNCALLED and RawHash benefit from Sequence Until

significantly **by up to 100×** reductions in sequencing time and costs

Tool	<i>SARS-CoV-2</i>	<i>E. coli</i>	<i>Yeast</i>	<i>Green Algae</i>	<i>Human</i>	Distance
Ground Truth	0.0929	0.4365	0.0698	0.1179	0.2828	N/A
UNCALLED (25%)	0.0026	0.5890	0.0613	0.1332	0.2139	0.1910
RawHash (25%)	0.0271	0.4853	0.0920	0.0786	0.3170	<b>0.0995</b>
UNCALLED (10%)	0.0026	0.5906	0.0611	0.1316	0.2141	0.1920
RawHash (10%)	0.0273	0.4869	0.0963	0.0772	0.3124	<b>0.1004</b>
UNCALLED (1%)	0.0026	0.5750	0.0616	0.1506	0.2103	0.1836
RawHash (1%)	0.0259	0.4783	0.0987	0.0882	0.3088	<b>0.0928</b>
UNCALLED (0.1%)	0.0040	0.4565	0.0380	0.1910	0.3105	0.1242
RawHash (0.1%)	0.0212	0.5045	0.1120	0.0810	0.2814	<b>0.1136</b>
UNCALLED (0.01%)	0.0000	0.5551	0.0000	0.0000	0.4449	0.2602
RawHash (0.01%)	0.0906	0.6122	0.0000	0.0000	0.2972	<b>0.2232</b>

# More in the Paper

- **More Results**

- **Mapping time** per read
- Overall **computational resources** required by each tool
  - Peak memory usage, CPU time and real time in the indexing and mapping steps
- **Performance breakdown** of the steps in RawHash

- **Details of all mechanisms and configurations**

- Details of the **quantization** and **hashing** mechanism
- Details of the **parameter configurations**
- Trade-offs between the **DNN-based approaches** and raw signal mapping approaches

# RawHash

- [Can Firtina](#), [Nika Mansouri Ghiasi](#), [Joel Lindegger](#), [Gagandeep Singh](#), [Meryem Banu Cavlak](#), [Haiyu Mao](#), and [Onur Mutlu](#),

## **"RawHash: Enabling Fast and Accurate Real-Time Analysis of Raw Nanopore Signals for Large Genomes"**

*Proceedings of the [31st Annual Conference on Intelligent Systems for Molecular Biology \(ISMB\)](#) and the [22nd European Conference on Computational Biology \(ECCB\)](#), Jul 2023*

[[arXiv preprint](#)]

[[Source Code](#)]

*Bioinformatics*, 2023, **39**, i297–i307

<https://doi.org/10.1093/bioinformatics/btad272>

ISMB/ECCB 2023



OXFORD

---

## RawHash: enabling fast and accurate real-time analysis of raw nanopore signals for large genomes

[Can Firtina](#) <sup>1,\*</sup>, [Nika Mansouri Ghiasi](#) <sup>1</sup>, [Joel Lindegger](#) <sup>1</sup>, [Gagandeep Singh](#) <sup>1</sup>,  
[Meryem Banu Cavlak](#) <sup>1</sup>, [Haiyu Mao](#) <sup>1</sup>, [Onur Mutlu](#) <sup>1,\*</sup>

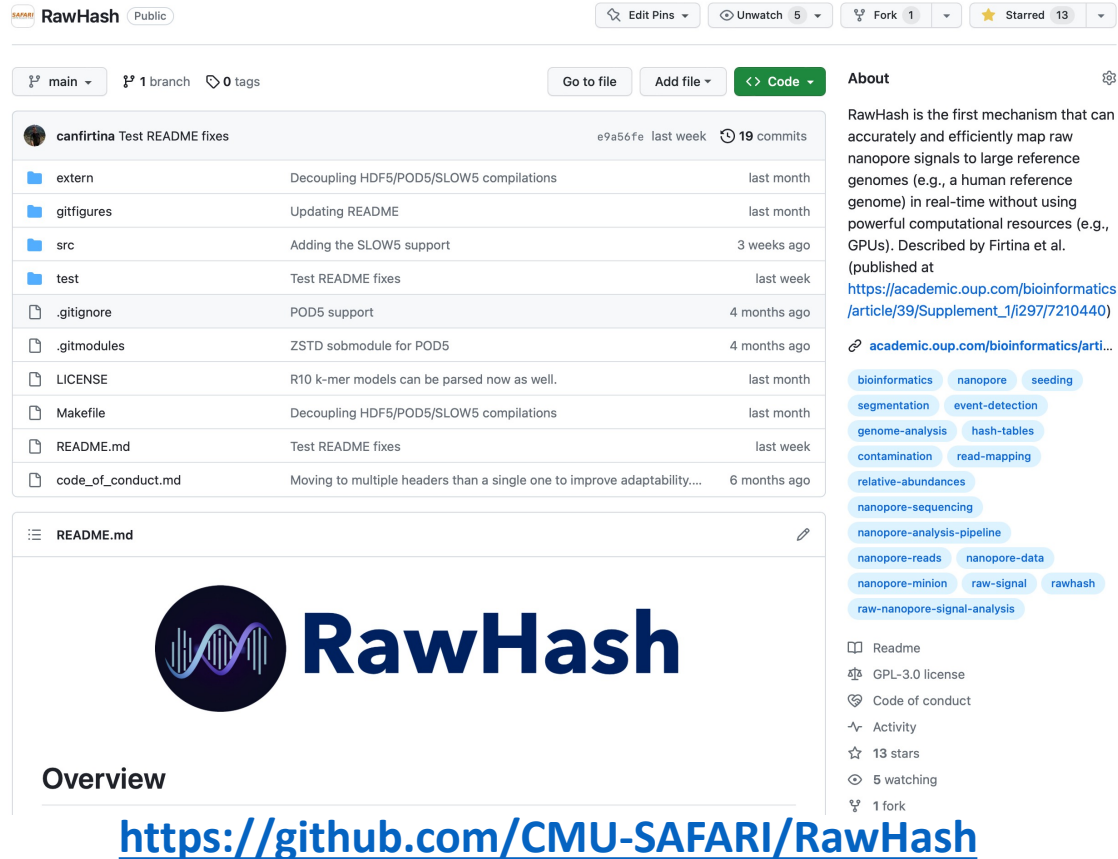
<sup>1</sup>Department of Information Technology and Electrical Engineering, ETH Zurich, 8092 Zurich, Switzerland

\*Corresponding author. Department of Information Technology and Electrical Engineering, ETH Zurich, Gloriastrasse 35, 8092 Zurich, Switzerland.  
E-mail: [firtinac@ethz.ch](mailto:firtinac@ethz.ch) (C.F.), [omutlu@ethz.ch](mailto:omutlu@ethz.ch) (O.M.)



# RawHash Source Code

- Supports **all major raw signal file formats and flow cell versions**
  - FAST5, POD5, S/BLOW5 file formats
- Easy-to-use scripts
  - To download all the datasets
  - To reproduce all of our results
- You can write your outlier function for Sequence Until
  - Easily integrate Sequence Until
- Upcoming Feature:
  - Integrating the MinKNOW API



RawHash

Public

Edit Pins Unwatch 5 Fork 1 Starred 13


main 1 branch 0 tags

Go to file Add file Code

canfirtina Test README fixes e9a56fe last week 19 commits

extern	Decoupling HDF5/POD5/SLOW5 compilations	last month
gitfigures	Updating README	last month
src	Adding the SLOW5 support	3 weeks ago
test	Test README fixes	last week
.gitignore	POD5 support	4 months ago
.gitmodules	ZSTD submodule for POD5	4 months ago
LICENSE	R10 k-mer models can be parsed now as well.	last month
Makefile	Decoupling HDF5/POD5/SLOW5 compilations	last month
README.md	Test README fixes	last week
code_of_conduct.md	Moving to multiple headers than a single one to improve adaptability...	6 months ago

README.md



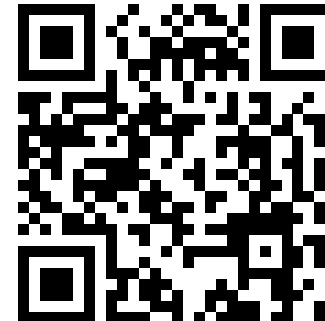
## RawHash

### Overview

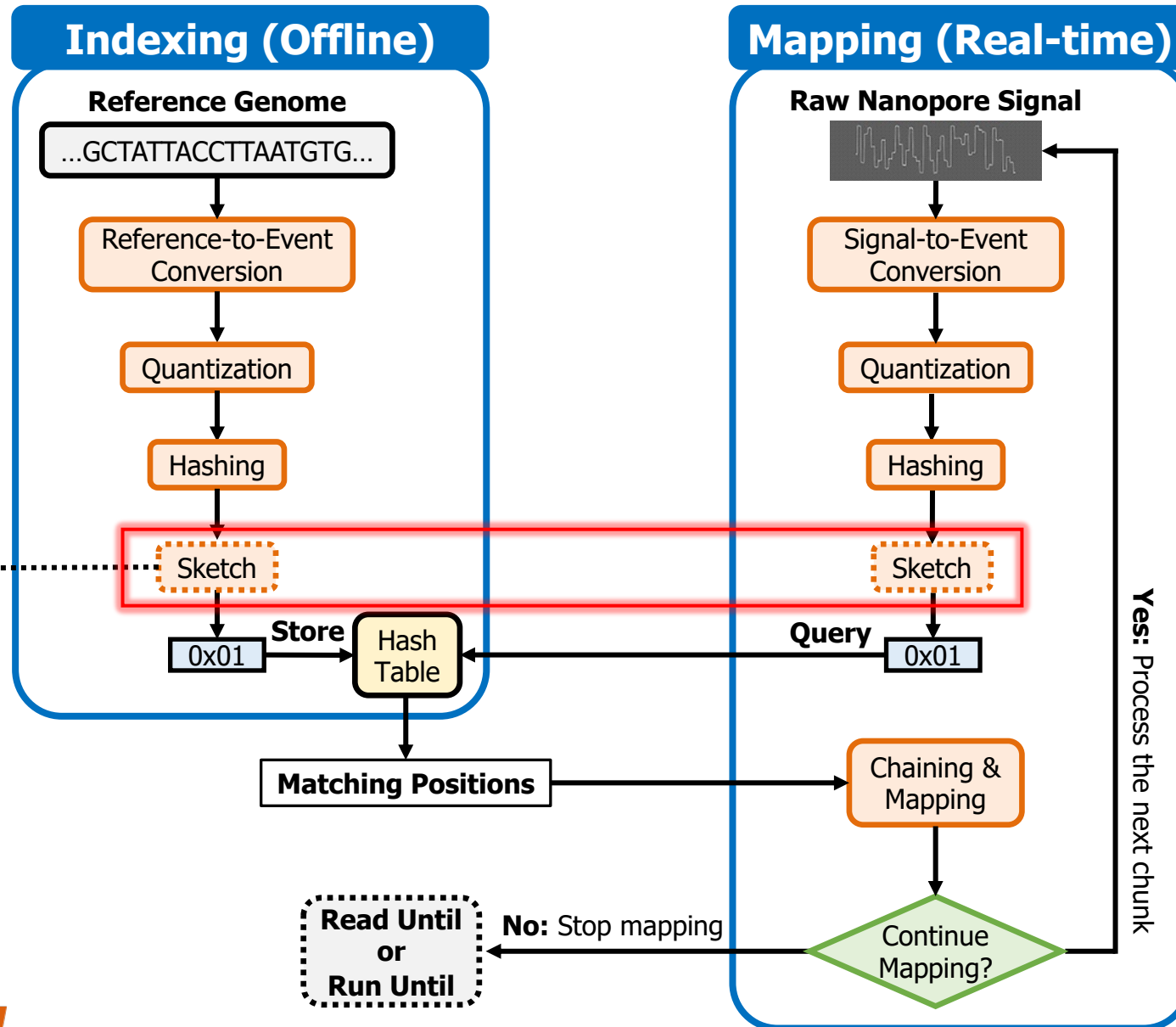
<https://github.com/CMU-SAFARI/RawHash>

bioinformatics nanopore seeding segmentation event-detection genome-analysis hash-tables contamination read-mapping relative-abundances nanopore-sequencing nanopore-analysis-pipeline nanopore-reads nanopore-data nanopore-minion raw-signal rawhash raw-nanopore-signal-analysis

Readme GPL-3.0 license Code of conduct Activity 13 stars 5 watching 1 fork



# Sketching with Hash-based Indexing



# Outline

Background

RawHash

Evaluation

Conclusion

# Conclusion

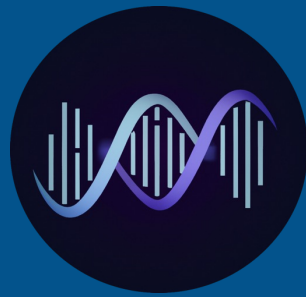
## Key Contributions:

- 1) The **first hash-based mechanism** that can quickly and accurately analyze raw nanopore signals for **large genomes**
- 2) The novel **Sequence Until** technique can accurately and **dynamically stop the entire sequencing of all reads at once** if further sequencing is not necessary

- Key Results:** Across 3 use cases and 5 genomes of varying sizes, RawHash provides
- **25.8× and 3.4× better average throughput** compared to two state-of-the-art works
  - **1.14× – 2.13× more accurate mapping results for large genomes**
  - Sequence Until **reduces the sequencing time and cost by 15×**

## Many opportunities for analyzing raw nanopore signals in real-time:

- Many hash-based **sketching techniques** can now be used for raw signals
- **Indexing is very cheap:** Many future use cases with the on-the-fly index construction
- We should rethink the algorithms to perform downstream analysis fully using raw signals



# RawHash

Enabling Fast and Accurate Real-Time Analysis  
of Raw Nanopore Signals for Large Genomes

**Can Firtina**

Nika Mansouri Ghiasi

Joel Lindegger

Gagandeep Singh

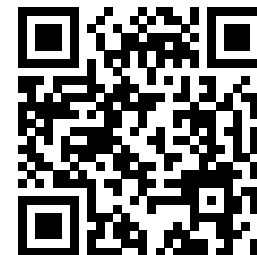
Meryem Banu Cavlak

Haiyu Mao

Onur Mutlu



[Paper](#)



[Code](#)

# Fast and Accurate Real-Time Genome Analysis

---

- Can Firtina, Melina Soysal, Joel Lindegger, and Onur Mutlu,  
**"RawHash2: Accurate and Fast Mapping of Raw Nanopore Signals using a Hash-based Seeding Mechanism"**  
*Preprint on **arxiv**, September 2023.*  
[\[arXiv version\]](#)  
[\[RawHash2 Source Code\]](#)

## **RawHash2: Accurate and Fast Mapping of Raw Nanopore Signals using a Hash-based Seeding Mechanism**

Can Firtina   Melina Soysal   Joel Lindegger   Onur Mutlu  
*ETH Zürich*

# Optimizations in RawHash2 (1)

---

- **More sensitive** chaining implementation with penalty scores
  - **Benefits:** Enables filtering dissimilar regions quickly
  - **Downside:** Additional computations with costly log operations
  
- **Weighted mapping decisions**
  - **Benefit #1:** 'Learned' mapping decisions based on the weights chosen from empirical analysis
  - **Benefit #2:** Faster and more accurate decisions
  
- Frequency **filters**
  - Filters the seeds that frequently appear before chaining
  - **Benefits:** Reduced workload on chaining without significantly affecting accuracy
  - **Downside:** Less sensitive mapping due to removed seeds

# Optimizations in RawHash2 (2)

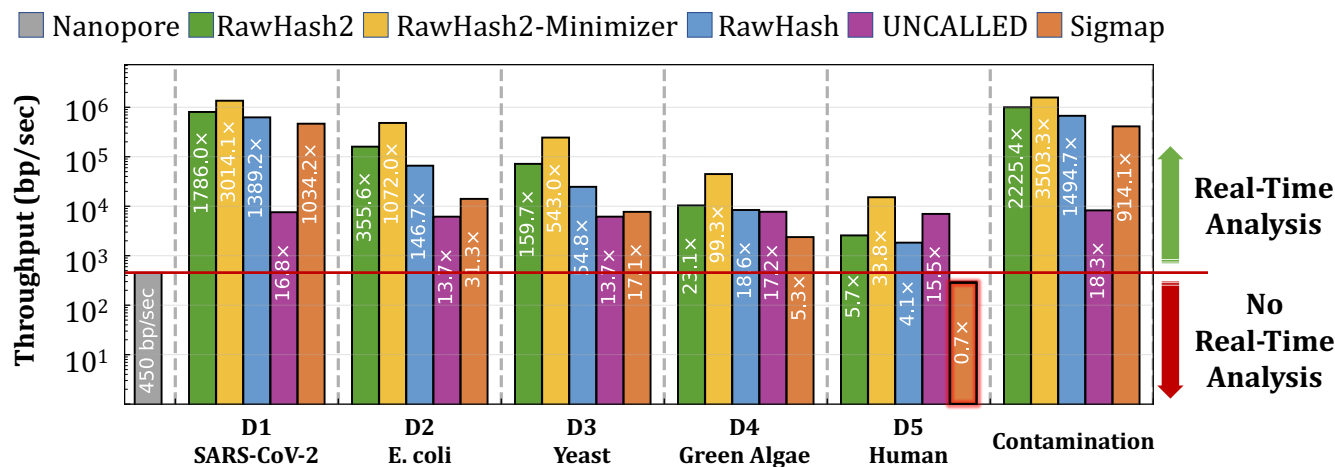
---

- New sketching techniques such as **minimizers** and **BLEND**
  - Enables integration of widely studied sketching techniques
  - **Benefits:** Can take advantage of these techniques (e.g., reduced storage requirements)
- Support for the recent improvements in the technology
  - Support for **new data formats:** POD5 and S/BLOW5
  - Support for **newer nanopore chemistry** versions: R10.4



# Results – Throughput

- **Real-time analysis requires** faster throughput than sequencer
  - Throughput of a nanopore sequencer: **~450 bp/sec (data generation speed)**



**2.3x** better average throughput RawHash

# Results – Accuracy

Dataset		UNCALLED	Sigmap	RawHash	RawHash2	RawHash2-Minimizer
Read Mapping						
D1 <i>SARS-CoV-2</i>	Precision	0.9547	<b>0.9929</b>	0.9868	0.9857	0.9602
	Recall	<b>0.9910</b>	0.5540	0.8735	0.8842	0.7080
	$F_1$	<b>0.9725</b>	0.7112	0.9267	0.9322	0.8150
D2 <i>E. coli</i>	Precision	0.9816	0.9842	0.9573	<b>0.9864</b>	0.9761
	Recall	<b>0.9647</b>	0.9504	0.9009	0.8934	0.7805
	$F_1$	<b>0.9731</b>	0.9670	0.9282	0.9376	0.8674
D3 <i>Yeast</i>	Precision	0.9459	0.9856	<b>0.9862</b>	0.9567	0.9547
	Recall	<b>0.9366</b>	0.9123	0.8412	0.8942	0.7792
	$F_1$	0.9412	<b>0.9475</b>	0.9079	0.9244	0.8581
D4 <i>Green Algae</i>	Precision	0.8836	<b>0.9741</b>	0.9691	0.9264	0.9198
	Recall	0.7778	<b>0.8987</b>	0.7015	0.8659	0.6711
	$F_1$	0.8273	<b>0.9349</b>	0.8139	0.8951	0.7760
D5 <i>Human HG001</i>	Precision	0.4867	0.4287	<b>0.8959</b>	0.8830	0.8111
	Recall	0.2379	0.2641	0.4054	<b>0.4317</b>	0.1862
	$F_1$	0.3196	0.3268	0.5582	<b>0.5799</b>	0.3028
Contamination						
D1 and D5	Precision	0.9378	0.7856	0.8733	<b>0.9393</b>	0.9330

RawHash2 is more accurate than RawHash **in all cases**

# Results – Average Sequencing Length

Tool	<i>SARS-CoV-2</i>	<i>E. coli</i>	<i>Yeast</i>	<i>Green Algae</i>	<i>Human</i>	<i>Contamination</i>
Average sequenced base length per read						
UNCALLED	<b>184.51</b>	<b>580.52</b>	<b>1,233.20</b>	5,300.15	6,060.23	1,582.63
RawHash	513.95	1,376.14	2,565.09	4,760.59	4,773.58	742.56
RawHash2	488.46	1,234.39	1,715.31	<b>2,077.39</b>	<b>3,441.43</b>	<b>681.94</b>
RawHash2-Minimizer	566.42	1,763.76	2,339.41	2,891.55	4,090.68	787.82
Average sequenced number of chunks per read						
Sigmap	<b>1.01</b>	<b>2.11</b>	<b>4.14</b>	5.76	10.40	2.06
RawHash	1.24	3.20	5.83	10.72	10.70	2.41
RawHash2	1.18	2.93	4.02	<b>4.84</b>	<b>7.78</b>	<b>1.68</b>
RawHash2-Minimizer	1.39	4.16	5.45	6.66	9.17	1.89

RawHash2 uses fewer bases to sequence than RawHash in all cases

RawHash2 uses the smallest number of bases to sequence for larger genomes

# Fast and Accurate Real-Time Genome Analysis

---

- Can Firtina, Melina Soysal, Joel Lindegger, and Onur Mutlu,  
**"RawHash2: Accurate and Fast Mapping of Raw Nanopore Signals using a Hash-based Seeding Mechanism"**  
*Preprint on **arxiv**, September 2023.*  
[\[arXiv version\]](#)  
[\[RawHash2 Source Code\]](#)

## **RawHash2: Accurate and Fast Mapping of Raw Nanopore Signals using a Hash-based Seeding Mechanism**

Can Firtina   Melina Soysal   Joel Lindegger   Onur Mutlu  
*ETH Zürich*

# Agenda for Today

---

- Cutting-edge in Accelerating Genome Analysis
  - Intelligent genome analysis
  
- Enabling Fast and Accurate Real-time Analysis
  - RawHash and RawHash2
  
- Conclusion

Things Are Happening In Industry

# Illumina DRAGEN Bio-IT Platform (2018)

- Processes whole genome at 30x coverage in ~25 minutes with hardware support for data compression



[emea.illumina.com/products/by-type/informatics-products/dragen-bio-it-platform.html](https://emea.illumina.com/products/by-type/informatics-products/dragen-bio-it-platform.html)  
[emea.illumina.com/company/news-center/press-releases/2018/2349147.html](https://emea.illumina.com/company/news-center/press-releases/2018/2349147.html)

# NextSeq 2000 with Analysis Capability

## NextSeq 1000/2000 Integrates DRAGEN Bio-IT Platform On-Board

### DRAGEN Bio-IT platform:

- Fast
- Accurate
- Industry standard pipelines
- For both novice and expert users

### Pipelines available on-board:

- DRAGEN Enrichment pipeline
- DRAGEN RNA pipeline
- DRAGEN Germline
- DRAGEN Single Cell RNA
- Generate FASTQ via BCL Convert
- *Additional pipelines available in BaseSpace Sequence Hub*

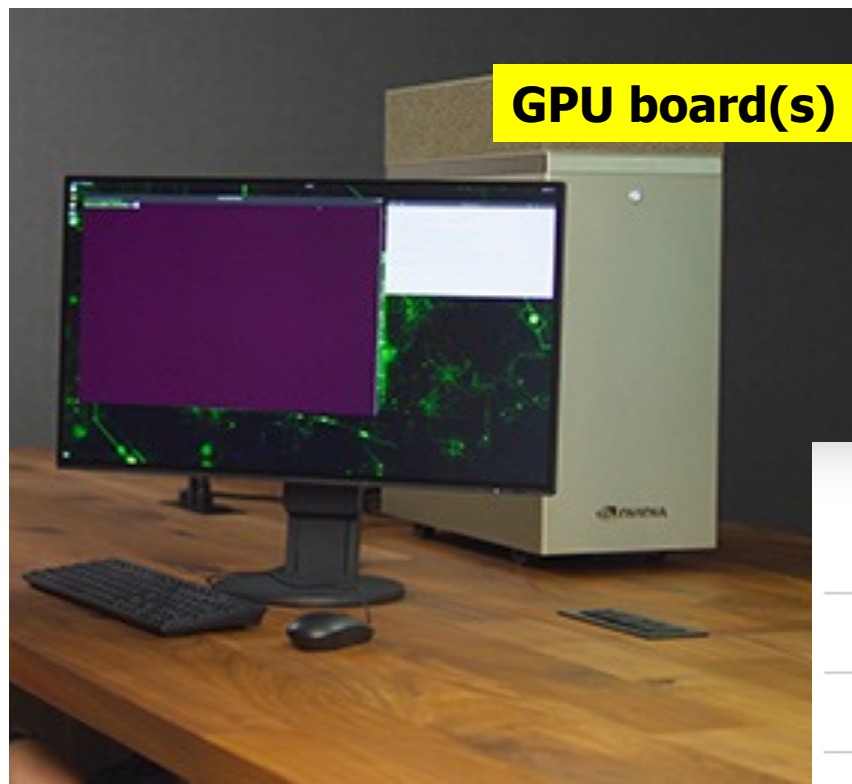


illumina®

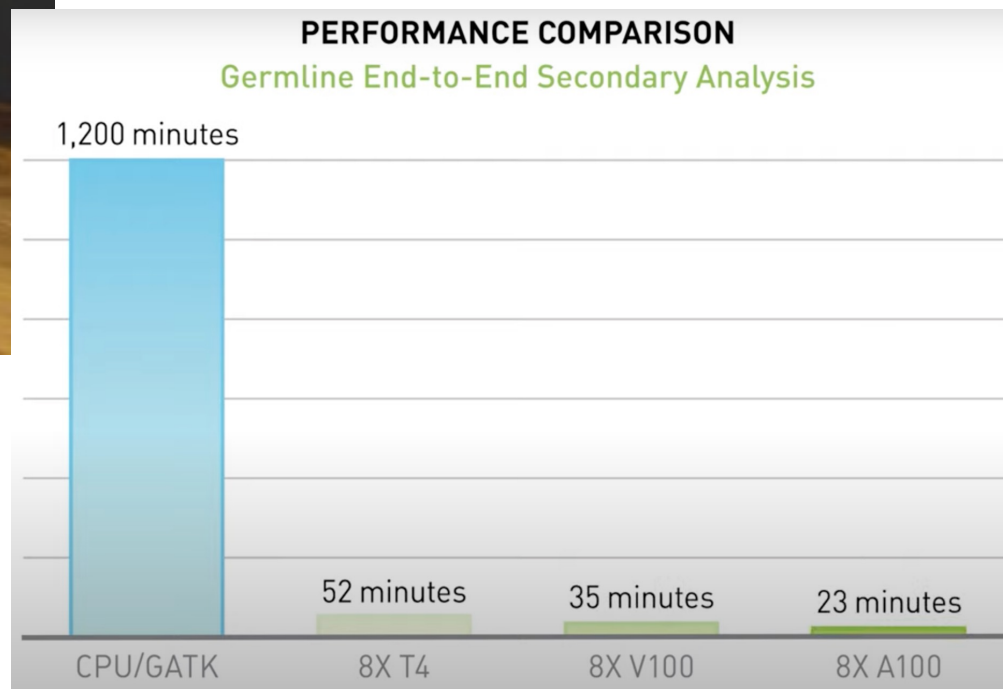
For Research Use Only.  
Not for use in diagnostic procedures.



# NVIDIA Clara Parabricks (2020)



**A University of Michigan startup in 2018 joined NVIDIA in 2020**

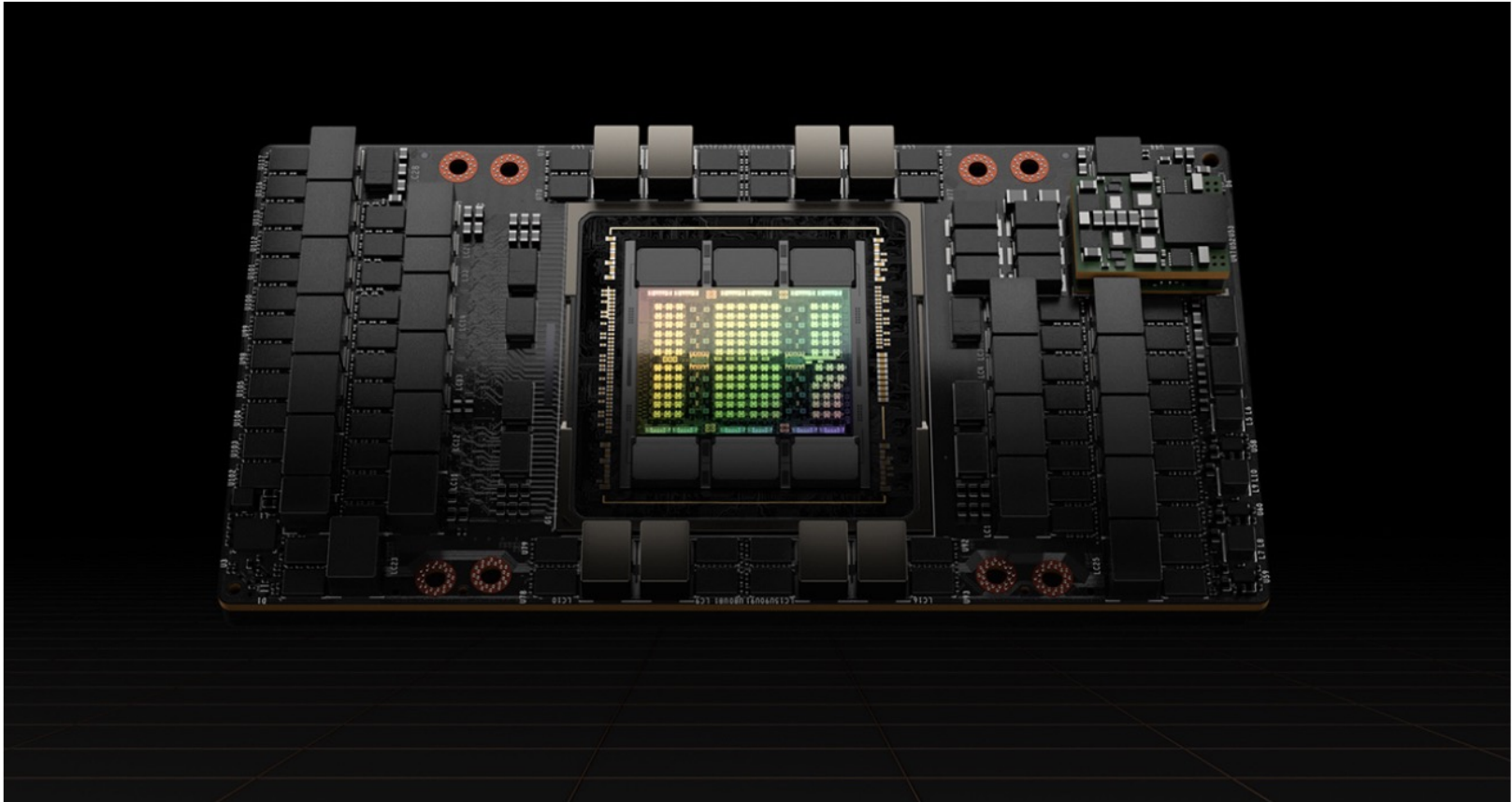


# NVIDIA Hopper DPX Instructions (2022)

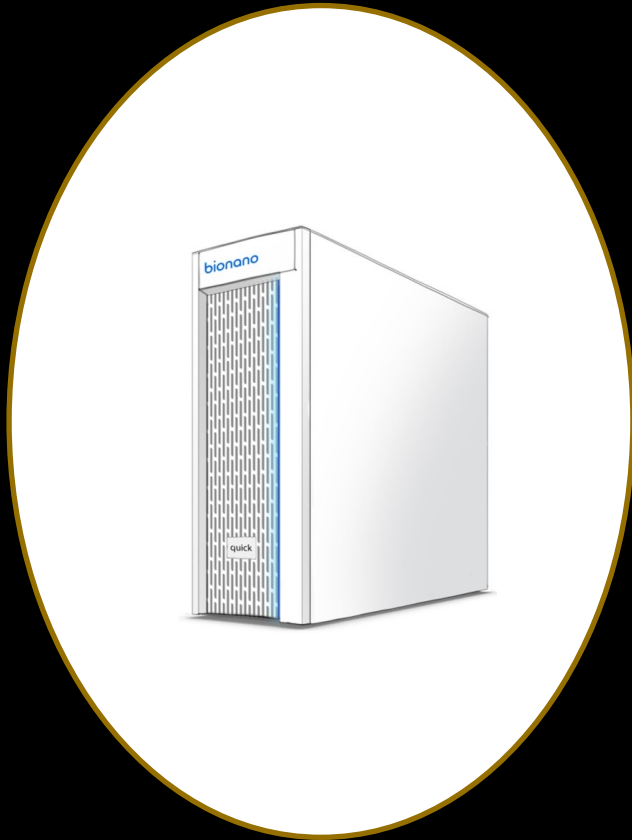
## NVIDIA Hopper GPU Architecture Accelerates Dynamic Programming Up to 40x Using New DPX Instructions

Dynamic programming algorithms are used in healthcare, robotics, quantum computing, data science and more.

March 22, 2022 by DION HARRIS



- We are accelerating the transformation in how we analyze the human genome!



## Bionano & NVIDIA:

*Accelerating Analysis for Fast Time to Results*



Technological solution to **support higher throughput**



**New high-performance algorithms** from Bionano



**Powered by NVIDIA RTX™ 6000 Ada Generation GPUs**

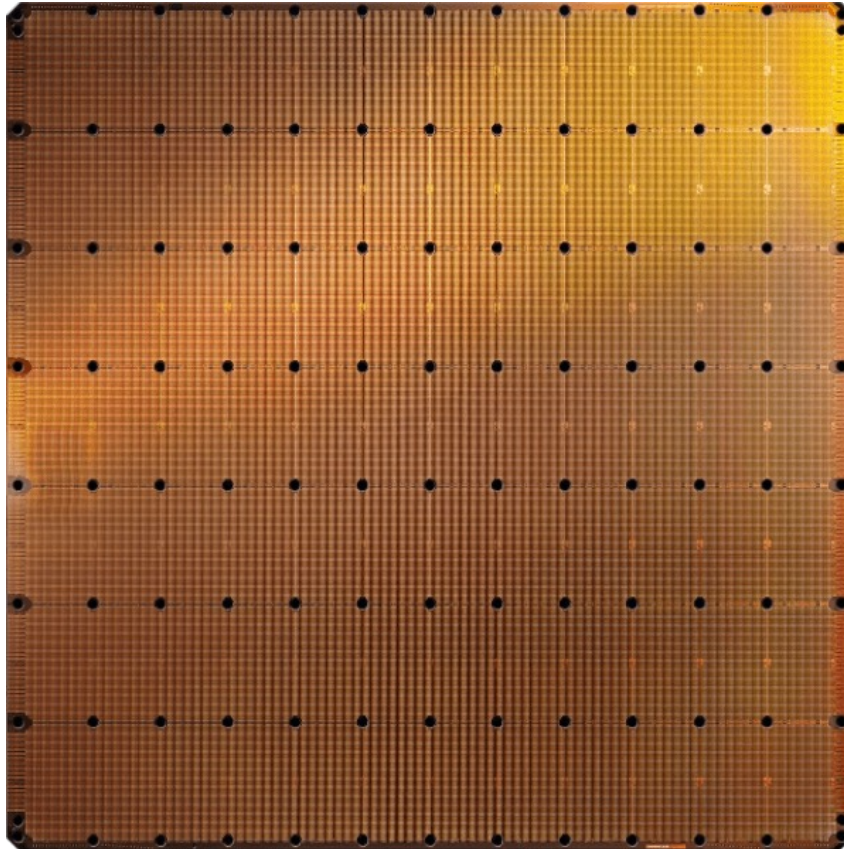


Analysis of highly complex cancer whole genomes in **less than 2 hours**



Workflow tailored for a **small lab and IT footprint**

# Cerebras's Wafer Scale Engine (2021)



**Cerebras WSE-2**  
2.6 Trillion transistors  
46,225 mm<sup>2</sup>

- The largest ML accelerator chip (2021)
- 850,000 cores



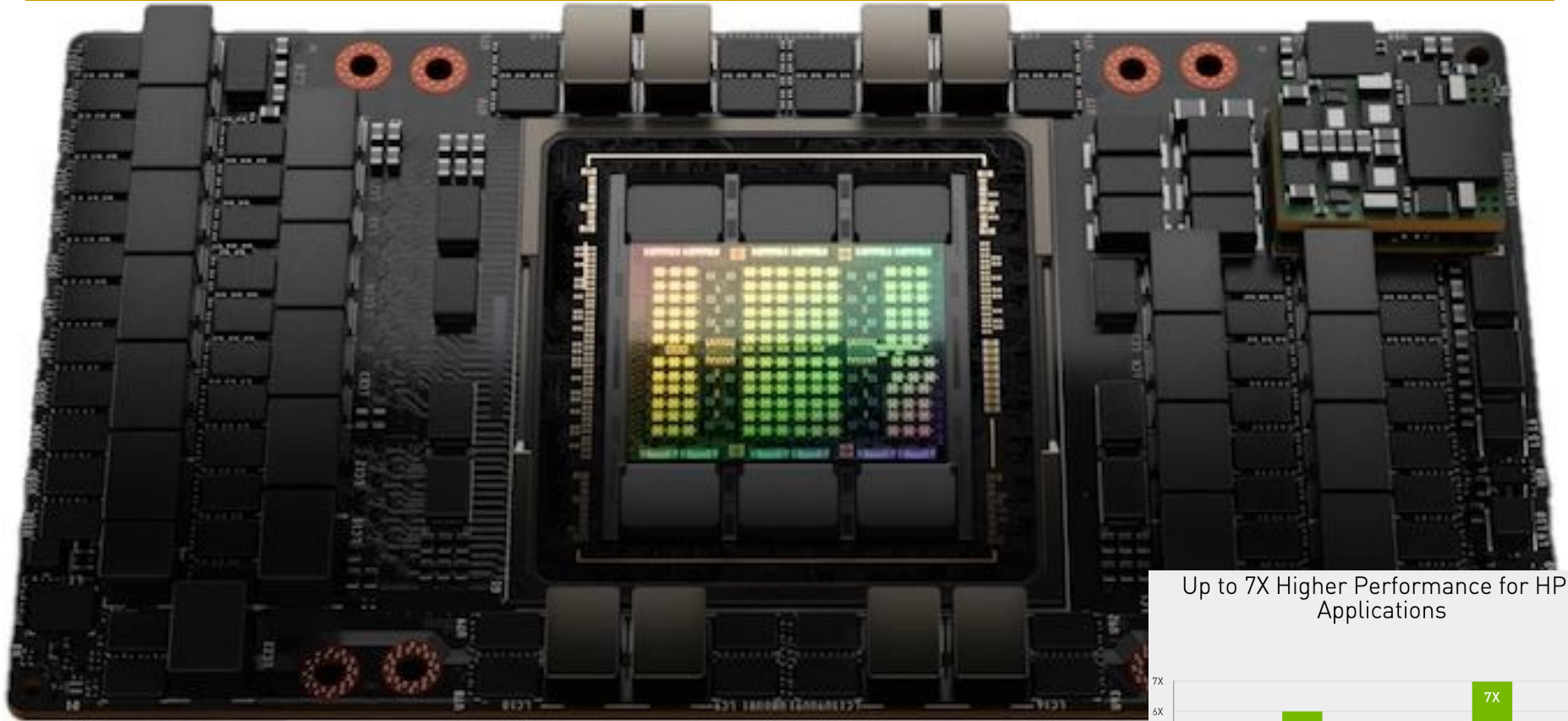
**Largest GPU**  
54.2 Billion transistors  
826 mm<sup>2</sup>  
NVIDIA Ampere GA100

<https://www.anandtech.com/show/14758/hot-chips-31-live-blogs-cerebras-wafer-scale-deep-learning>

**SAFARI**

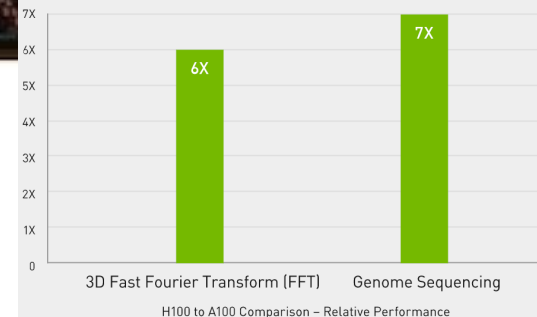
<https://www.cerebras.net/cerebras-wafer-scale-engine-why-we-need-big-chips-for-deep-learning/>

# NVIDIA H100 (2022)



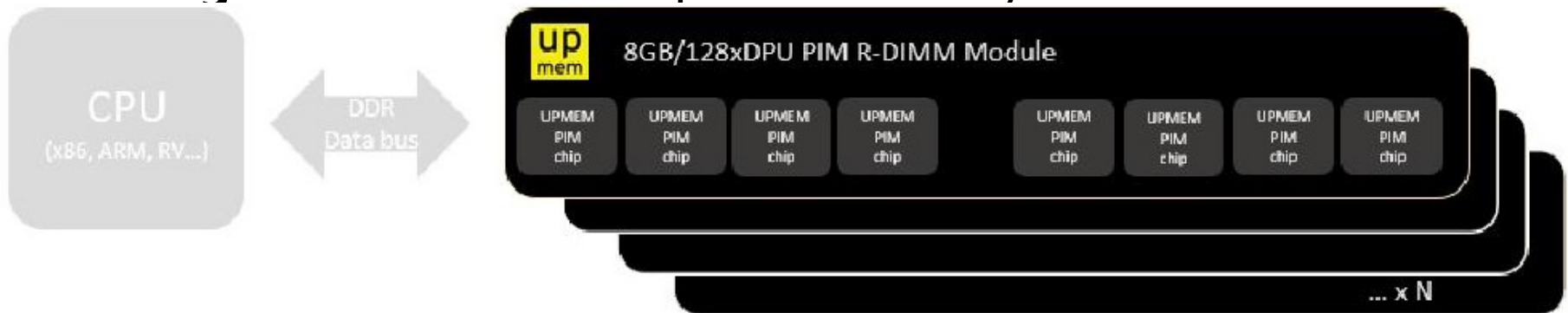
NVIDIA is claiming a **7x improvement** in dynamic programming algorithm (**DPX instructions**) performance on a single H100 versus naïve execution on an A100.

Up to 7X Higher Performance for HPC Applications



# UPMEM Processing-in-DRAM Engine (2019)

- **Processing in DRAM Engine**
- Includes **standard DIMM modules**, with a **large number of DPU processors** combined with DRAM chips.
- Replaces **standard DIMMs**
  - DDR4 R-DIMM modules
    - 8GB+128 DPUs (16 PIM chips)
    - Standard 2x-nm DRAM process
  - **Large amounts of** compute & memory bandwidth

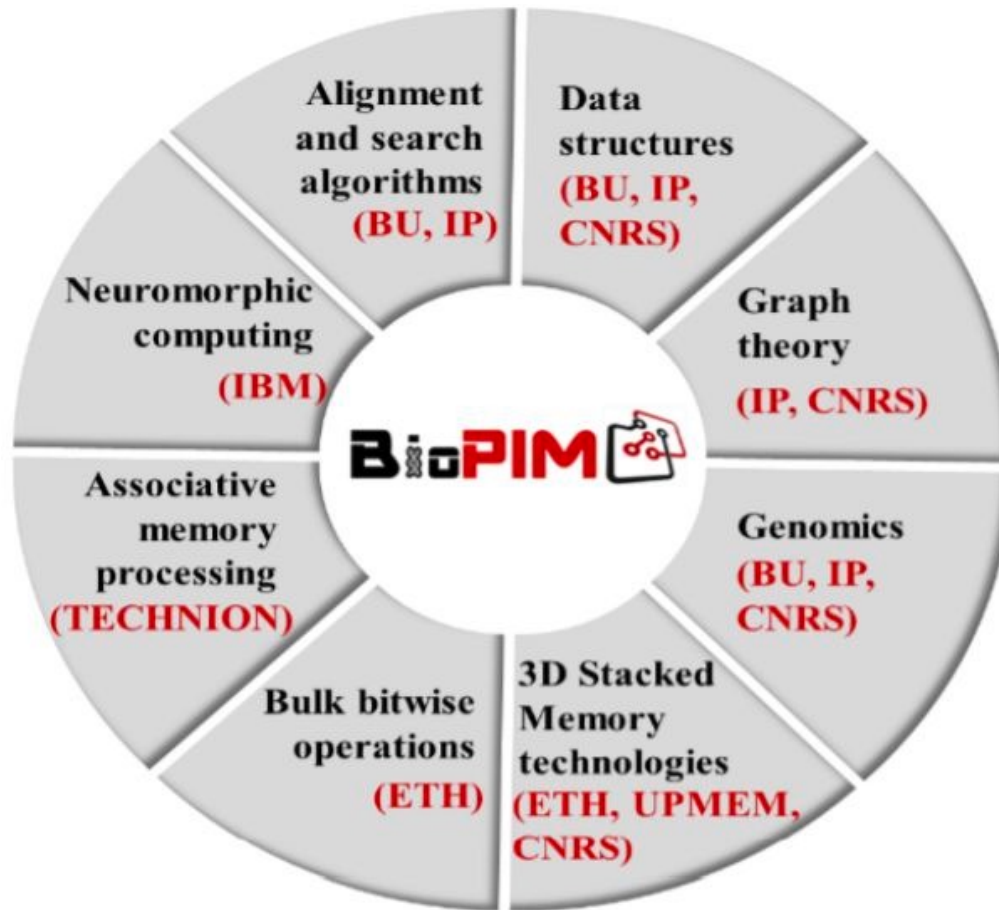


<https://www.anandtech.com/show/14750/hot-chips-31-analysis-inmemory-processing-by-upmem>

<https://www.upmem.com/video-upmem-presenting-its-true-processing-in-memory-solution-hot-chips-2019/>

# BioPIM (2022)

---



The vision of **BioPIM** is the realization of **cheap, ultra-fast and ultra-low energy mobile genomics** that eliminates the current dependence of sequence analysis on large and power-hungry computing clusters/data-centers.

# Fast Genome Analysis...

- Onur Mutlu,  
**"Accelerating Genome Analysis: A Primer on an Ongoing Journey"**  
*Invited Lecture at [Technion](#), Virtual, 26 January 2021.*  
[[Slides \(pptx\)](#) ([pdf](#))]  
[[Talk Video](#) (1 hour 37 minutes, including Q&A)]  
[[Related Invited Paper \(at IEEE Micro, 2020\)](#)]

Insight: Shifting a String Helps Similarity Search

7 matches    1 mismatch

ISTANBUL

ISTNBUL

ISTNBUL

81

Onur Mutlu - Invited Lecture @Technion: Accelerating Genome Analysis: A Primer on an Ongoing Journey

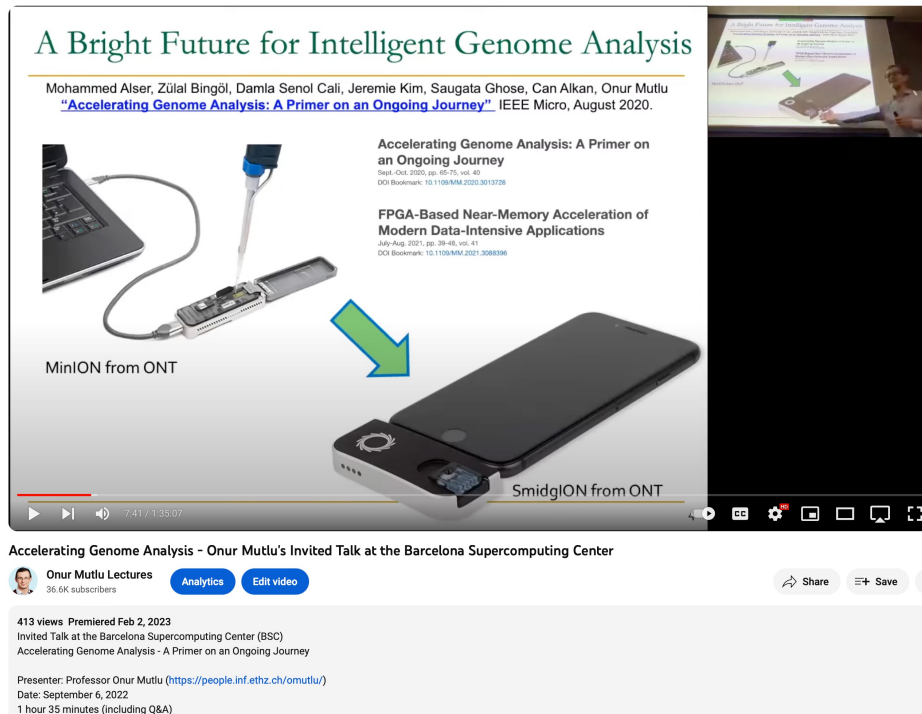
566 views · Premiered Feb 6, 2021

👍 31    🗨️ 0    ➦ SHARE    📌 SAVE    ⋮



# More on Fast Genome Analysis...

- Onur Mutlu,  
**"Accelerating Genome Analysis"**  
*Invited Talk at the Barcelona Supercomputing Center (BSC), Barcelona, Spain, 6 September 2022.*  
[[Slides \(pptx\)](#)] [[pdf](#)]  
[[Talk Video](#) (1 hour 35 minutes, including Q&A)]  
[[Related Invited Paper \(at IEEE Micro, 2020\)](#)]  
[[Related Invited Paper \(at Computational and Structural Biology Journal, 2022\)](#)]



**A Bright Future for Intelligent Genome Analysis**

Mohammed Alser, Zülal Bingöl, Damla Senol Cali, Jeremie Kim, Saugata Ghose, Can Alkan, Onur Mutlu  
"Accelerating Genome Analysis: A Primer on an Ongoing Journey" IEEE Micro, August 2020.

**Accelerating Genome Analysis: A Primer on an Ongoing Journey**  
Sept-Oct 2020, pp. 46-75, vol. 40  
DOI Bookmark: 10.1109/MM.2020.3013726

**FPGA-Based Near-Memory Acceleration of Modern Data-Intensive Applications**  
July-Aug. 2021, pp. 38-48, vol. 41  
DOI Bookmark: 10.1109/MM.2021.3068396

MinION from ONT

SmidgION from ONT

Accelerating Genome Analysis - Onur Mutlu's Invited Talk at the Barcelona Supercomputing Center

**Onur Mutlu Lectures**  
36.6K subscribers

Analytics Edit video

Share Save

413 views Premiered Feb 2, 2023  
Invited Talk at the Barcelona Supercomputing Center (BSC)  
Accelerating Genome Analysis - A Primer on an Ongoing Journey

Presenter: Professor Onur Mutlu (<https://people.inf.ethz.ch/omutlu/>)  
Date: September 6, 2022  
1 hour 35 minutes (including Q&A)

# More on Accelerating Genome Analysis

- Can Firtina,  
**"Enabling Accurate, Fast, and Memory-Efficient Genome Analysis via Efficient and Intelligent Algorithms"**  
*Talk at UC Berkeley, Berkeley, CA, United States, May 27, 2022.*  
[\[Slides \(pptx\) \(pdf\)\]](#)  
[\[Talk Video \(1 hour 6 minutes\)\]](#)



Enabling Accurate, Fast, and Memory-Efficient Genome Analysis - Can Firtina (Talk at UC Berkeley)

# More on Real-Time Genome Analysis

- Can Firtina, ["RawHash: Enabling Fast and Accurate Real-Time Analysis of Raw Nanopore Signals for Large Genomes"](#)

*Proceedings Talk at ISMB-ECCB, Lyon, France, 25 July 2023.*

[[Slides \(pptx\)](#) ([pdf](#))]

[[Talk Video](#) (18 minutes)]

**RawHash – Key Idea**

**Key Observation:** Identical nucleotides generate similar raw signals

The diagram illustrates the RawHash process. It shows two raw signals, 'Raw Signal #1' and 'Raw Signal #2', each represented by a waveform. A dashed arrow labeled 'Distance Calculation' connects them, but this arrow is crossed out with a red 'X', indicating that direct distance calculation is not used. Instead, each signal is processed by a 'Hash' function (represented by a blue hexagon) to produce a '0x01' value. These two '0x01' values are then compared in a 'Fast Match' step (represented by a green rounded rectangle).

**Challenge #1:** Generating the **same** hash value for **similar enough** signals

**Challenge #2:** **Accurately** finding similar regions **as few as possible**

**SAFARI** 14

RawHash: Enabling Fast and Accurate Real-Time Analysis of Raw Nanopore Signals | ISMB-ECCB 2023

Onur Mutlu Lectures  
36.1K subscribers

Analytics Edit video

Share Save

294 views Premiered Aug 15, 2023  
Talk of "RawHash: Enabling Fast and Accurate Real-Time Analysis of Raw Nanopore Signals for Large Genomes" at ISMB-ECCB 2023  
Presenter: Can Firtina  
Duration: 18:58 minutes

# Accelerating Genome Analysis [DAC 2023]

---

- Onur Mutlu and Can Firtina,  
**"Accelerating Genome Analysis via Algorithm-Architecture Co-Design"**  
*Invited Special Session Paper in Proceedings of the 60th Design Automation Conference (DAC), San Francisco, CA, USA, July 2023.*  
[[Slides \(pptx\)](#)] [[pdf](#)]  
[[Talk Video](#)] (38 minutes, including Q&A)  
[[Related Invited Paper](#)]  
[[arXiv version](#)]

## Accelerating Genome Analysis via Algorithm-Architecture Co-Design

Onur Mutlu   Can Firtina  
*ETH Zürich*

# BIO-Arch Workshop at RECOMB 2023

■ April 14, 2023

## BIO-Arch: Workshop on Hardware Acceleration of Bioinformatics Workloads

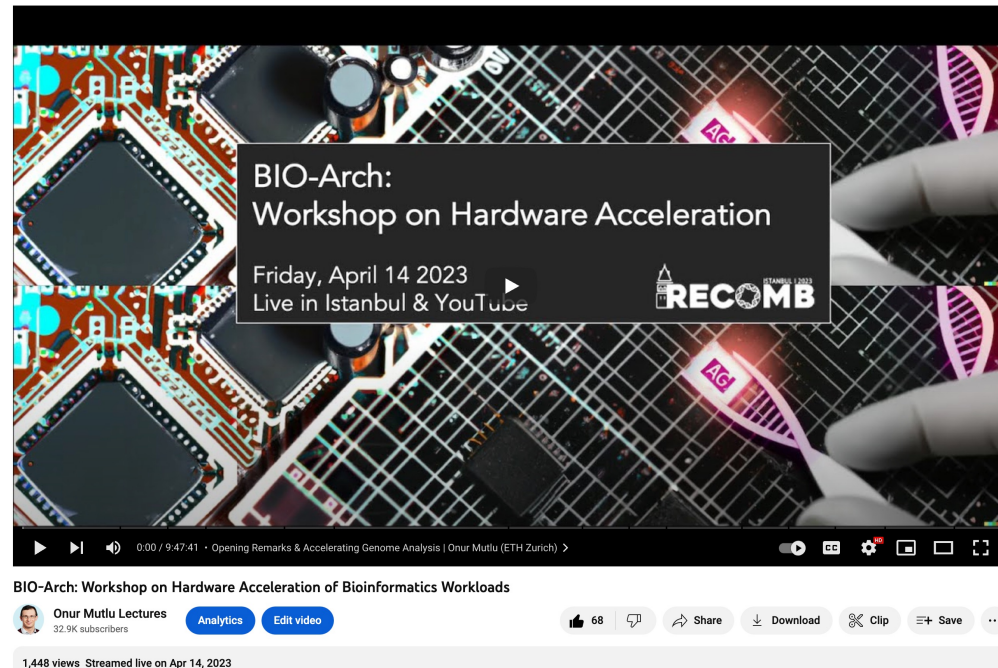
### About

BIO-Arch is a new forum for presenting and discussing new ideas in accelerating bioinformatics workloads with the co-design of hardware & software and the use of new computer architectures. Our goal is to discuss new system designs tailored for bioinformatics. BIO-Arch aims to bring together researchers in the bioinformatics, computational biology, and computer architecture communities to strengthen the progress in accelerating bioinformatics analysis (e.g., genome analysis) with efficient system designs that include hardware acceleration and software systems tailored for new hardware technologies.

### Venue

BIO-Arch will be held in [The Social Facilities of Istanbul Technical University](#) on **April 14**. Detailed information about how to arrive at the venue location with various transportation options can be found on [the RECOMB website](#).

Our panel discussion will be held in conjunction with the main RECOMB conference. The panel discussion will be held in [Marriott Şişli](#) on **April 17 at 17:00**. You can find



<https://www.youtube.com/watch?v=2rCsb4-nLmg>

**SAFARI**

<https://safari.ethz.ch/recomb23-arch-workshop/>

# Genomics Course (Fall 2023)

- **Fall 2023 Edition:**

- [https://safari.ethz.ch/projects\\_and\\_seminars/fall2023/doku.php?id=bioinformatics](https://safari.ethz.ch/projects_and_seminars/fall2023/doku.php?id=bioinformatics)

- **Spring 2023 Edition:**

- [https://safari.ethz.ch/projects\\_and\\_seminars/spring2023/doku.php?id=bioinformatics](https://safari.ethz.ch/projects_and_seminars/spring2023/doku.php?id=bioinformatics)

- **Youtube Livestream (Fall 2023):**

- [https://youtube.com/playlist?list=PL5Q2soXY2Zi\\_00wyOjiMShG4t2QPZoeE3](https://youtube.com/playlist?list=PL5Q2soXY2Zi_00wyOjiMShG4t2QPZoeE3)

- Project course

- Taken by Bachelor's/Master's students
  - Genomics lectures
  - Hands-on research exploration
  - Many research readings

<https://www.youtube.com/onurmutlectures>

Complete Lecture Playlist (Spring 2023):

Fall 2023 Schedule

Week	Date	Livestream	Meeting
W0	05.10 Thu.		<b>L0: Project Introductions and Q&amp;A</b>
W1	11.10 Wed.	<b>You Tube Live</b>	<b>L1: P&amp;S Course Introduction &amp; Scope</b> <a href="#">PDF</a> <a href="#">PPT</a>
W2	25.10 Wed.		<b>L2: Introduction to Genome Analysis</b> <a href="#">PDF</a> <a href="#">PPT</a>
W3	01.11 Wed.		<b>L3: From Molecules to Data: An Overview of DNA Sequencing Technologies</b> <a href="#">PDF</a> <a href="#">PPT</a>
W4	08.11 Wed.		<b>L4a: Fundamentals of Sequence Alignment: Algorithms and Applications</b> <a href="#">PDF</a> <a href="#">PPT</a> <b>L4b: Optimizing Sequence Search: Hashing, Indexing, and Filtering Techniques</b> <a href="#">PDF</a> <a href="#">PPT</a>

# The Future is Bright for Genome Analysis

---

- We covered various recent ideas to
  - Accelerate genome analysis
  - Analyze genomes in ways that were not possible before
- Enabling cost-effective, portable, fast, and accurate genome analysis has many implications
  - What are the new applications to enable with these unique benefits?
- Can we do even better?
  - Understanding and modifying the sequencing process for analyzing other types of biological data
- **Many future opportunities exist**
  - **Especially with new sequencing technologies**
  - **Especially with new applications and use cases**

## **RawHash and RawHash2:**

Enabling Fast & Accurate Real-time Analysis  
of Raw Nanopore Signals for Large Genomes  
using a Hash-based Seeding Mechanism

Can Firtina

[canfirtina@gmail.com](mailto:canfirtina@gmail.com)

<https://cfirtina.com>

17 January 2024

Leibniz Institute for Immunotherapy (LIT)

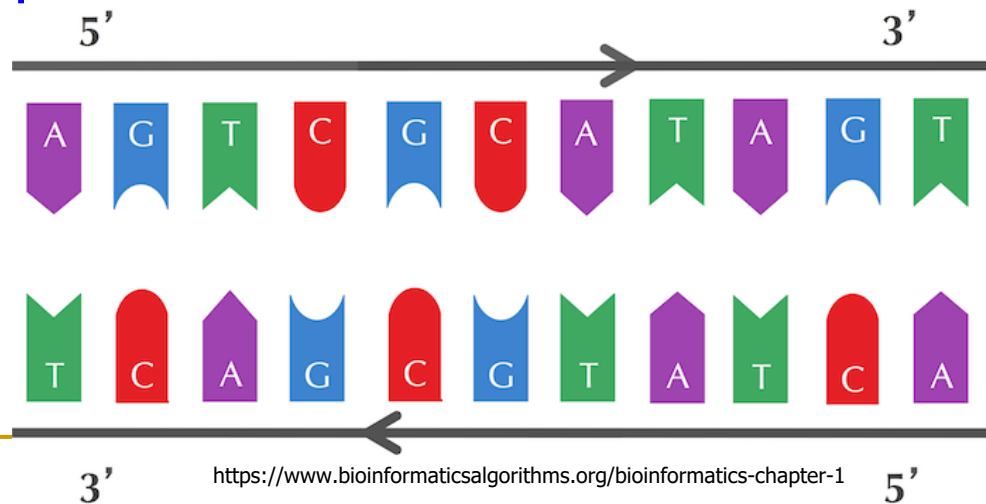
**SAFARI**

**ETH** zürich



# Challenges in Read Mapping

- Need to find many **mappings** of **each read**
- Need to **tolerate variances/sequencing errors** in each read
- Need to **map** each read **very fast** (i.e., performance is important, life critical in some cases)
- Need to **map** reads to both **forward and reverse strands**





# A Tsunami of Sequencing Data

A Tera-scale increase in sequencing production in the past 25 years		
Genes & Operons	1990	<b>Kilo</b> = 1,000
Bacterial genomes	1995	<b>Mega</b> = 1,000,000
Human genome	2000	<b>Giga</b> = 1,000,000,000
Human microbiome	2005	<b>Tera</b> = 1,000,000,000,000
50K Microbiomes	2015	<b>Peta</b> = 1,000,000,000,000,000
what is expected for the next 15 years ? (a Giga?)		
200K Microbiomes	2020	<b>Exa</b> = 1,000,000,000,000,000,000
1M Microbiomes	2025	<b>Zetta</b> = 1,000,000,000,000,000,000,000
Earth Microbiome	2030	<b>Yotta</b> = 1,000,000,000,000,000,000,000,000

Source:  
[@kyrpides](#)

# Solving the Puzzle

---

.FASTA file



Reference genome



.FASTQ file



Reads



<https://www.pacb.com/smrt-science/smrt-sequencing/hifi-reads-for-highly-accurate-long-read-sequencing/>



# Obtaining .FASTQ Files

- <https://www.ncbi.nlm.nih.gov/sra/ERR240727>



Full ▾

Send to: ▾

## [ERX215261](#): Whole Genome Sequencing of human TSI NA20754

1 ILLUMINA (Illumina HiSeq 2000) run: 4.1M spots, 818.7M bases, 387.2Mb downloads

**Design:** Illumina sequencing of library 6511095, constructed from sample accession SRS001721 for study accession SRP000540. This is part of an Illumina multiplexed sequencing run (9340\_1). This submission includes reads tagged with the sequence TTAGGCAT.

**Submitted by:** The Wellcome Trust Sanger Institute (SC)

**Study:** Whole genome sequencing of (TSI) Toscani in Italia HapMap population

[PRJNA33847](#) • [SRP000540](#) • [All experiments](#) • [All runs](#)

**Sample:** Coriell GM20754

[SAMN00001273](#) • SRS001721 • [All experiments](#) • [All runs](#)

*Organism:* [Homo sapiens](#)

### Library:

*Name:* 6511095

*Instrument:* Illumina HiSeq 2000

*Strategy:* WGS

*Source:* GENOMIC

*Selection:* RANDOM

*Layout:* PAIRED

*Construction protocol:* Standard

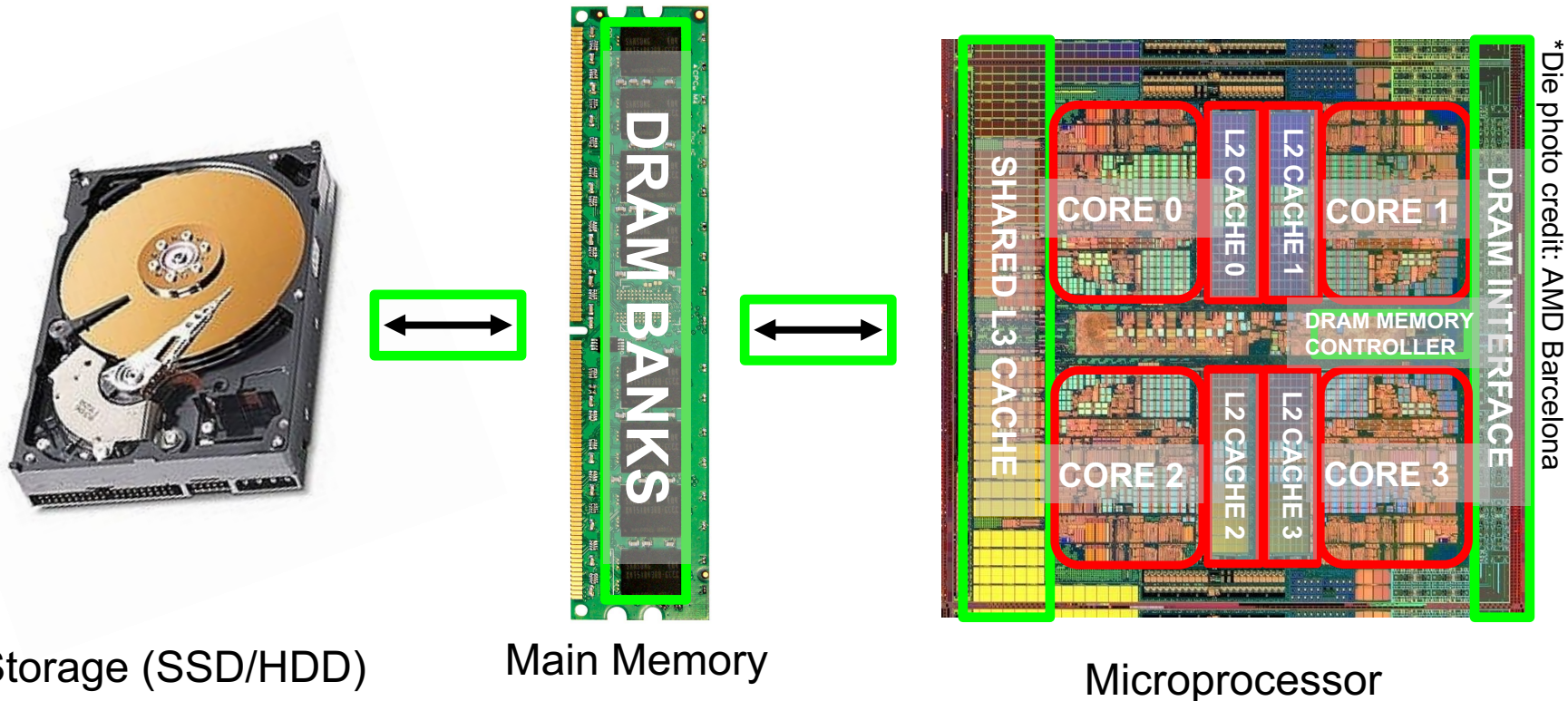
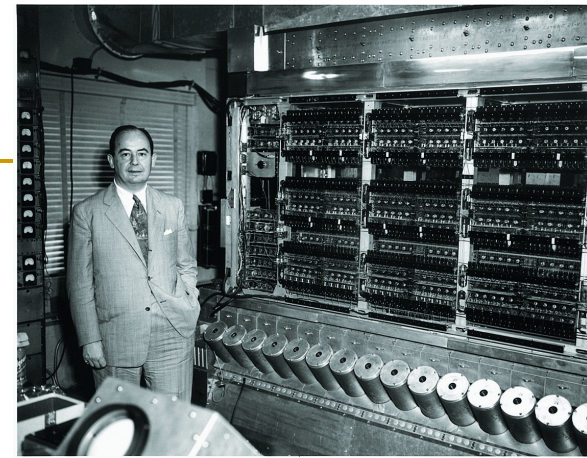
**Runs:** 1 run, 4.1M spots, 818.7M bases, [387.2Mb](#)

Run	# of Spots	# of Bases	Size	Published
<a href="#">ERR240727</a>	4,093,747	818.7M	387.2Mb	2013-03-22

# Today's Computing Systems

von Neumann model, 1945

where the **CPU** can **access data** stored in an off-chip main memory only through **power-hungry bus**



# The Problem

---

Data analysis  
is performed  
far away from the data



# Read Mapping

---

Map **reads** to a known reference genome with some minor differences allowed



DNA Sample  
"chemical format"



Reads  
"text format"



Subject genome  
"text format"

# 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 (can find all mapping locations), 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

# An Example of Hash Table Based Mappers

---

- + Guaranteed to find *all* mappings → very sensitive
- + Can tolerate up to *e* errors

nature  
genetics

<https://github.com/BilkentCompGen/mrfast>

---

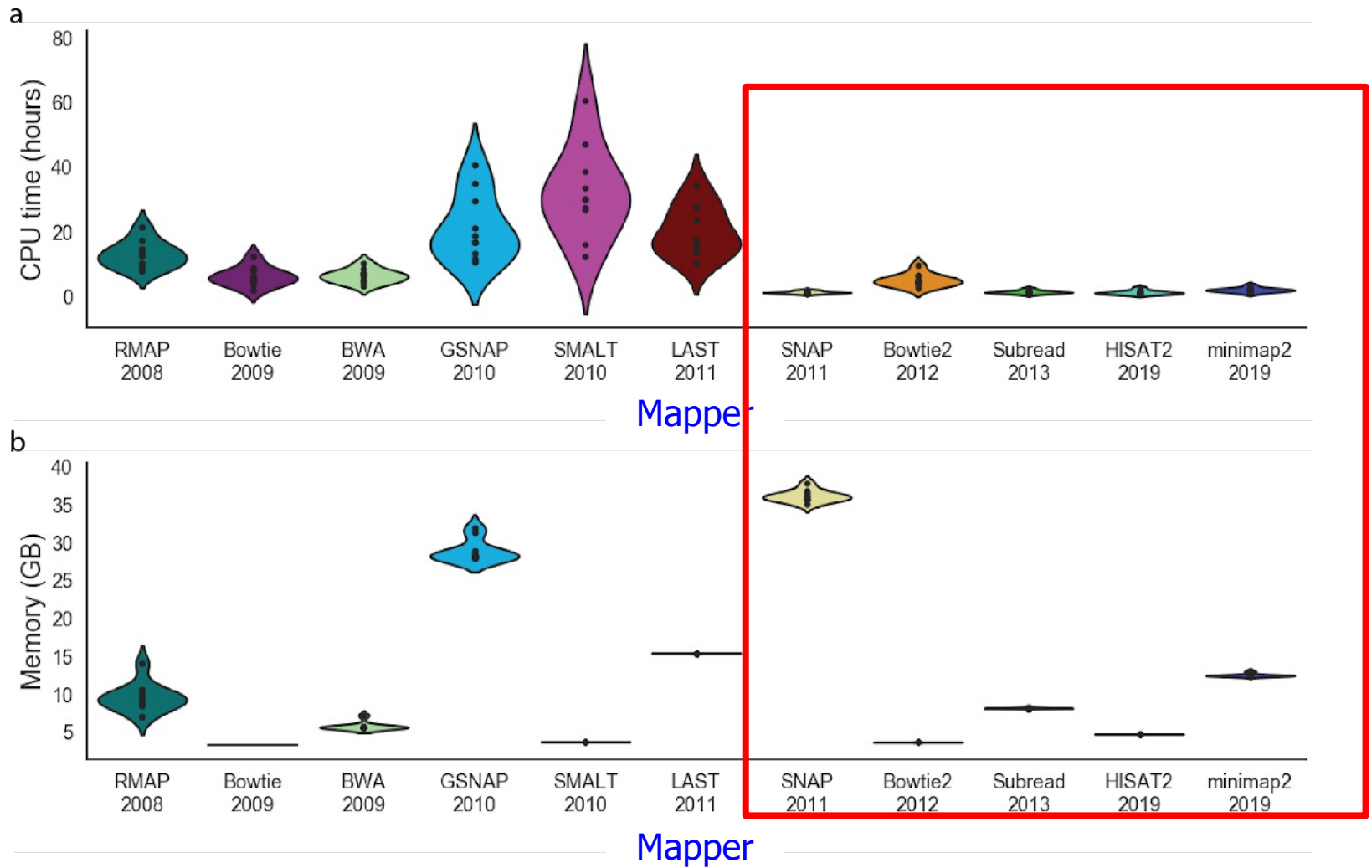
## Personalized copy number and segmental duplication maps using next-generation sequencing

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+, "[Personalized copy number and segmental duplication maps using next-generation sequencing](#)", Nature Genetics 2009.

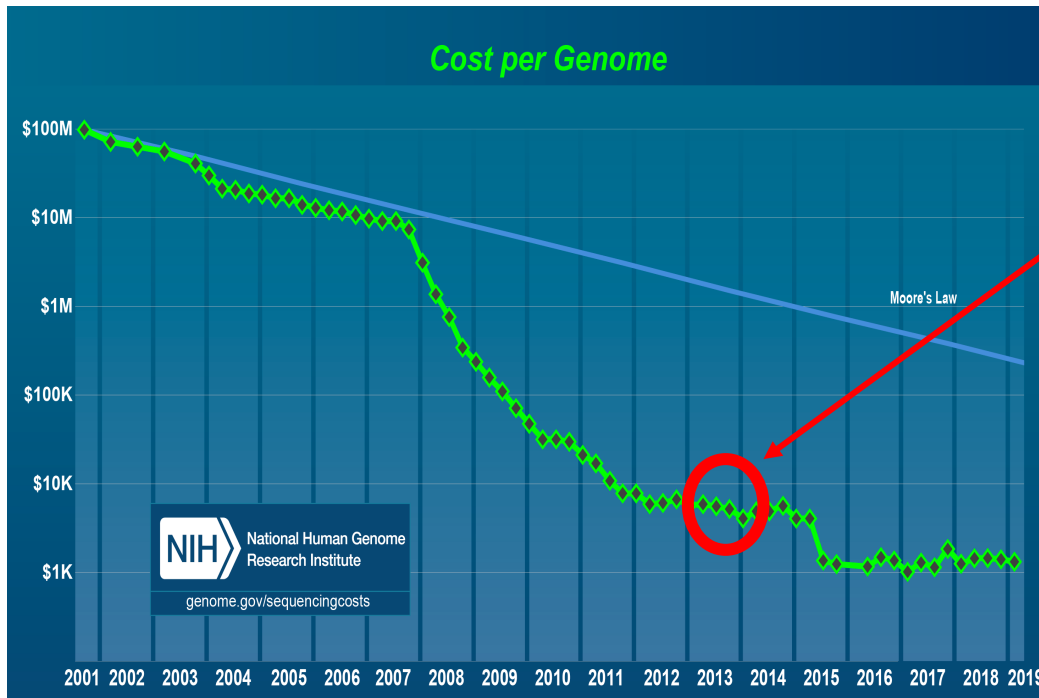
---

# Performance of Read Mapping

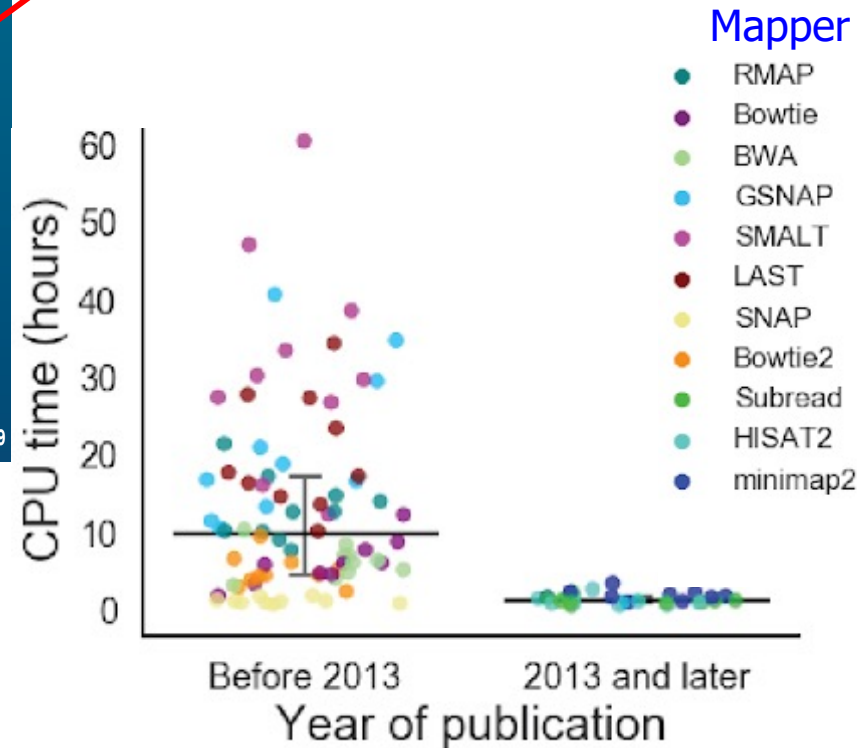


Alser+, "[Technology dictates algorithms: Recent developments in read alignment](#)",  
Genome Biology, 2021

# The Need for Speed



Did we realize the **need** for **faster** genome analysis?



Alser+, "[Technology dictates algorithms: Recent developments in read alignment](#)",  
Genome Biology, 2021

# Sequence Alignment in Unavoidable

- **Quadratic-time** dynamic-programming algorithm **WHY?!**

Enumerating all possible prefixes

- NETHERLANDS x SWITZERLAND
- NETHERLANDS x S
- NETHERLANDS x SW
- NETHERLANDS x SWI
- NETHERLANDS x SWIT
- NETHERLANDS x SWITZ
- NETHERLANDS x SWITZE
- NETHERLANDS x SWITZER
- NETHERLANDS x SWITZERL
- NETHERLANDS x SWITZERLA
- NETHERLANDS x SWITZERLAN
- NETHERLANDS x SWITZERLAND

	N	E	T	H	E	R	L	A	N	D	S	
	0	1	2	3	4	5	6	7	8	9	10	11
S	1	2	3	4	5	6	7	8	9	10	10	
W	2	3	4	5	6	7	8	9	10	11		
I	3	3	4	5	6	7	8	9	10	11		
T	4	4	4	3	4	5	6	7	8	9	10	11
Z	5	5	5	4	4	5	6	7	8	9	10	11
E	6	6	5	5	5	4	5	6	7	8	9	10
R	7	7	6	6	6	5	4	5	6	7	8	9
L	8	8	7	7	7	6	5	4	5	6	7	8
A	9	9	8	8	8	7	6	5	4	5	6	7
N	10	9	9	9	9	8	7	6	5	4	5	6
D	11	10	10	10	10	9	8	7	6	5	4	5

# Sequence Alignment in Unavoidable

- **Quadratic-time** dynamic-programming algorithm

Enumerating all possible prefixes

- **Data dependencies** limit the computation parallelism

Processing row (or column) after another

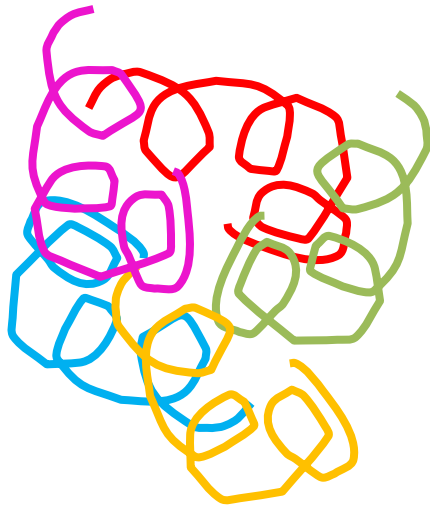
- **Entire matrix** is computed even though strings can be dissimilar.

Number of differences is computed only at the backtraking step.

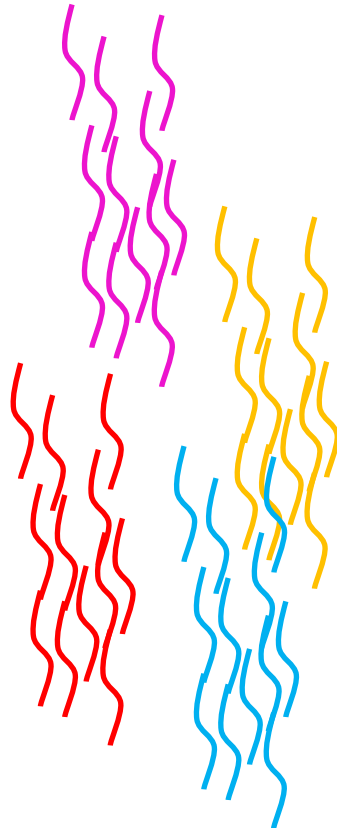
		N	E	T	H	E	R	L	A	N	D	S
	0	1	2	3	4	5	6	7	8	9	10	11
S	1	1	2	3	4	5	6	7	8	9	10	10
W	2	2	2	3	4	5	6	7	8	9	10	11
I	3	3	3	3	4	5	6	7	8	9	10	11
T	4	4	4	3	4	5	6	7	8	9	10	11
Z	5	5	5	4	4	5	6	7	8	9	10	11
E	6	6	5	5	5	4	5	6	7	8	9	10
R	7	7	6	6	6	5	4	5	6	7	8	9
L	8	8	7	7	7	6	5	4	5	6	7	8
A	9	9	8	8	8	7	6	5	4	5	6	7
N	10	9	9	9	9	8	7	6	5	4	5	6
D	11	10	10	10	10	9	8	7	6	5	4	5

# Metagenomics Analysis

Reads from different **unknown** donors at sequencing time are mapped to **many known reference** genomes

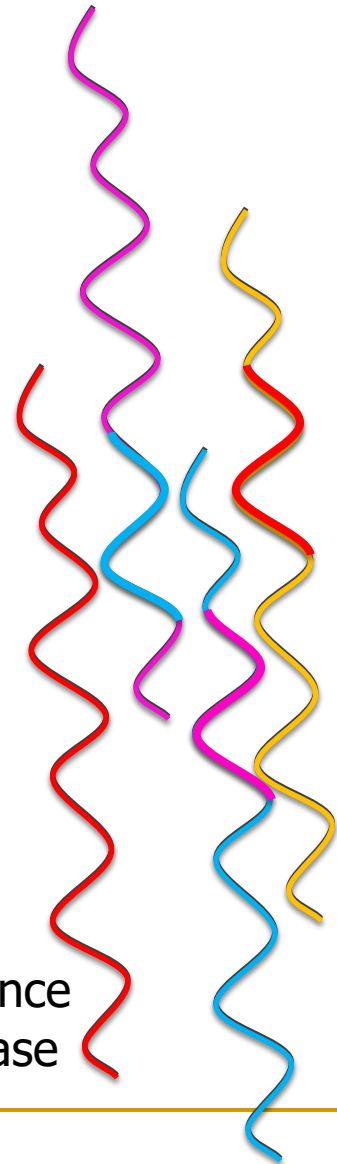


genetic material recovered directly from environmental samples



Reads "text format"

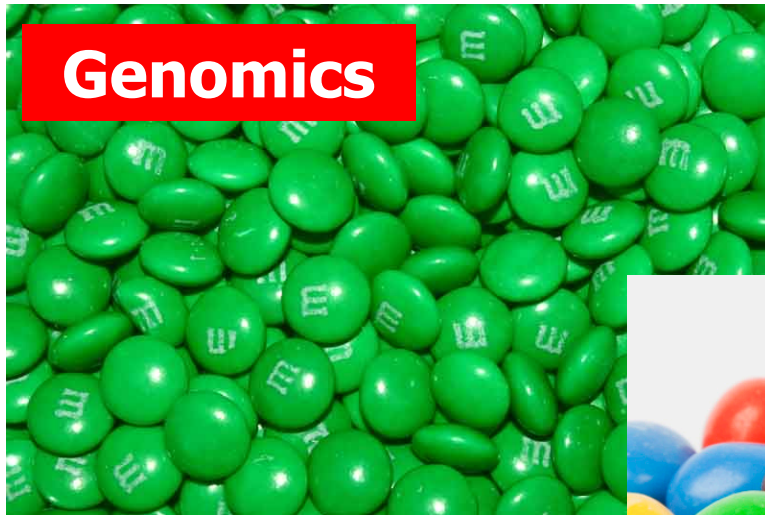
Reference Database





# Genomics vs. Metagenomics

---



# Practical Similarity Identification Seeds

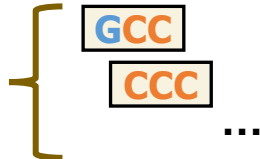
Reference



Read



K-mers



K-mers Locations

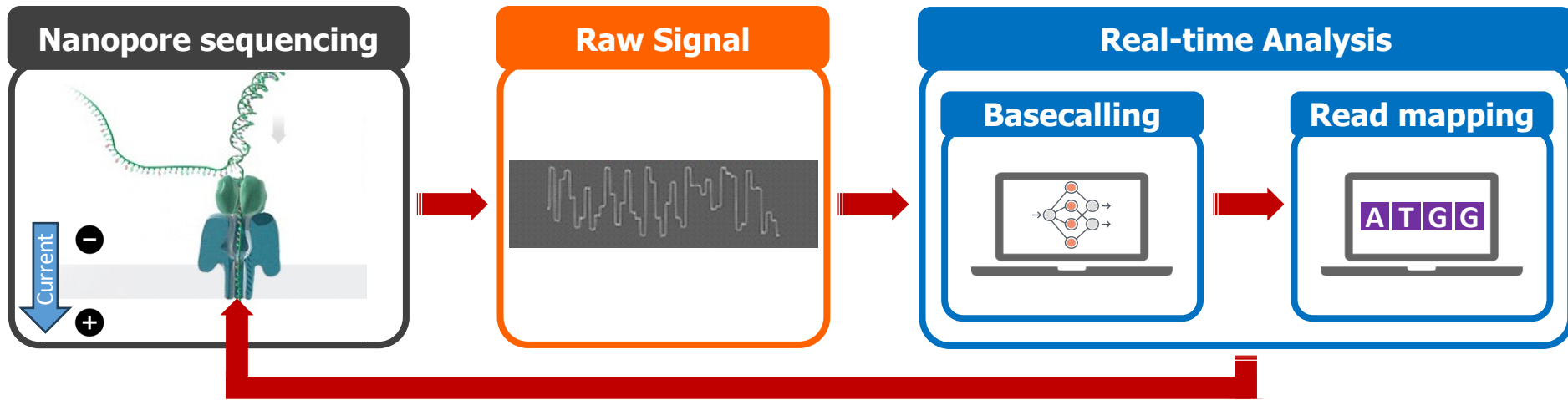
GCC	7		
CCC	8		
CAA	1		
AAA	31	101	
CCA	25	230	400
...	...	...	...

Index (Hash Table)

<b>Seeding</b>	Determine <b>potential matching regions (seeds)</b> in the reference genome
<b>Seed Filtering (e.g., Chaining)</b>	<b>Prune</b> some seeds in the reference genome
<b>Alignment</b>	Determine the <b>exact differences</b> between the read and the reference genome

# Existing Solutions – Real-time Basecalling

Deep neural networks (**DNNs**) for translating **signals** to **bases**

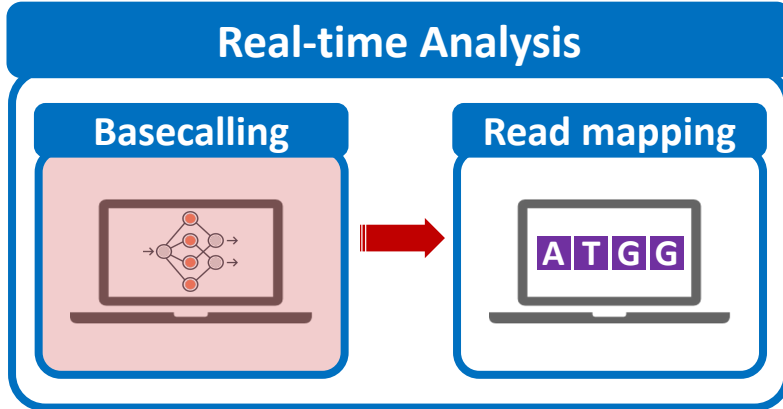


DNNs provide **less noisy analysis** from basecalled sequences

**Costly and power-hungry** computational requirements

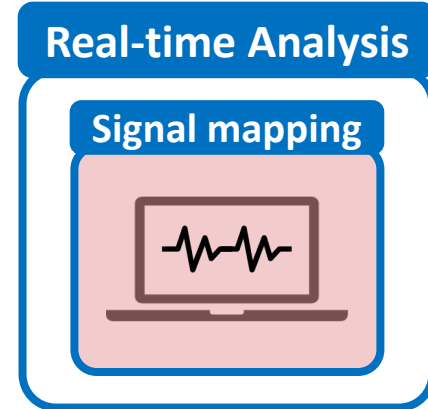
# The Problem

The existing solutions are **ineffective for large genomes**



**Costly and energy-hungry computations to basecall each read:**

Portable sequencing becomes challenging with resource-constrained devices



Larger number of reference regions **cannot be handled accurately or quickly**, rendering existing solutions **ineffective for large genomes**

# Applications of Read Until

**Depletion:** Reads mapping to a particular reference genome is ejected

- Removing contaminated reads from a sample
- Relative abundance estimation
- Controlling low/high-abundance genomes in a sample
- Controlling the sequencing of depth of a genome

**Enrichment:** Reads **not** mapping to a particular reference genome is ejected

- Purifying the sample to ensure it contains only the selected genomes
- Removing the host genome (e.g., human) in contamination analysis

# Applications of Run Until and Sequence Until

**Run Until:** Stopping the sequencing without informative decision from analysis

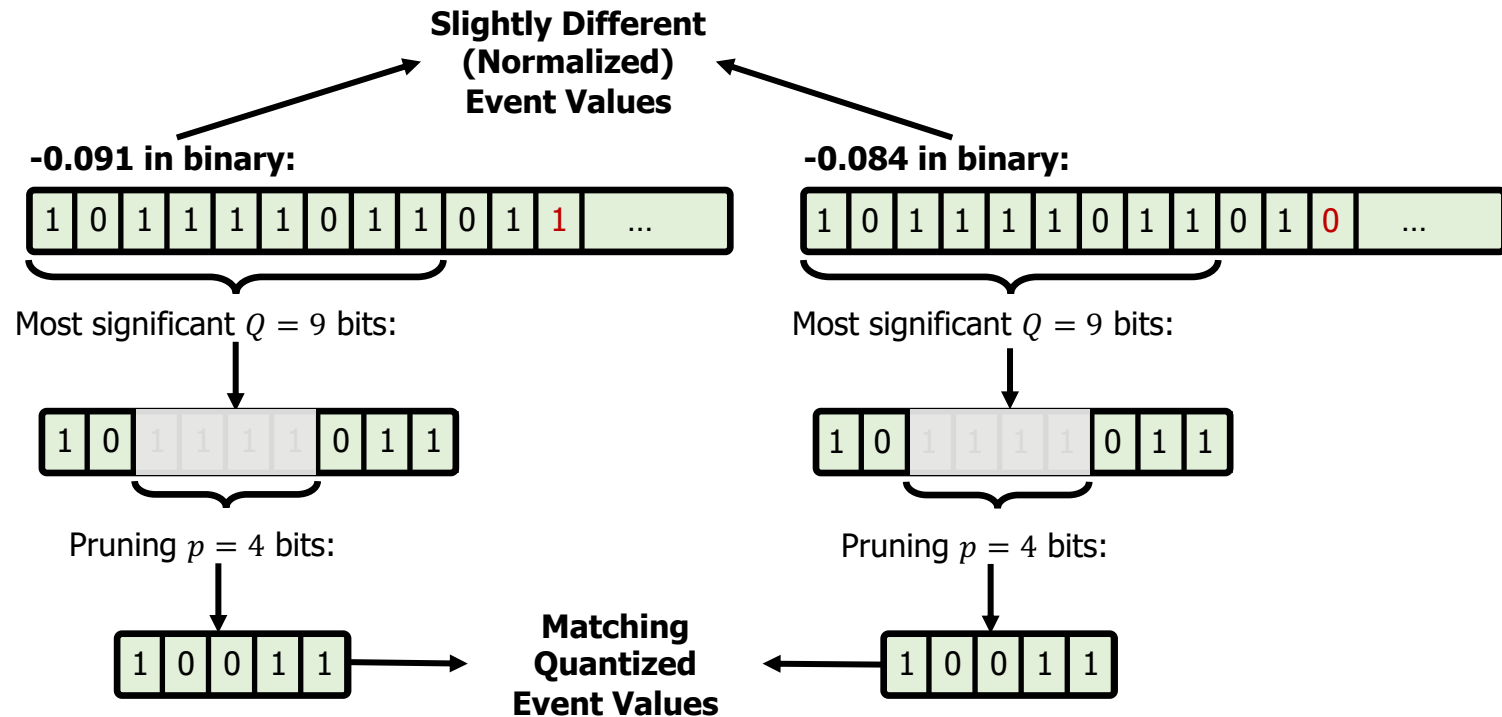
- Stopping when reads reach to a particular depth of coverage
- Stopping when the abundance of all genomes reach a particular threshold

**Sequence Until:** Stopping the sequencing based on information decision

- Stopping when relative abundance estimations do not change substantially (for high-abundance genomes)
- Stopping when finding that the sample is contaminated with a particular set of genomes
- ...

# Details: Quantizing the Event Values

- **Observation:** Identical k-mers generate similar raw signals
  - **Challenge:** Their corresponding event values can be slightly different
- **Key Idea:** Quantize the event values
  - To enable assigning the **same quantized value** to the **similar event values**



# Average Sequenced Bases and Chunks

Tool	<i>SARS-CoV-2</i>	<i>E. coli</i>	<i>Yeast</i>	<i>Green Algae</i>	<i>Human</i>
Average sequenced base length per read					
UNCALLED	<b>184.51</b>	<b>580.52</b>	<b>1,233.20</b>	5,300.15	6,060.23
RawHash	513.95	1,376.14	2,565.09	<b>4,760.59</b>	<b>4,773.58</b>
Average sequenced number of chunks per read					
Sigmap	<b>1.01</b>	<b>2.11</b>	<b>4.14</b>	<b>5.76</b>	<b>10.40</b>
RawHash	1.24	3.20	5.83	10.72	10.70

RawHash **reduces sequencing time and cost for large genomes**

up to **1.3**× compared to UNCALLED

Although Sigmap processes less number of chunks than RawHash, it fails to provide real-time analysis capabilities for large genomes



# Breakdown Analysis of the RawHash Steps

Tool	Fraction of entire runtime (%)				
	<i>SARS-CoV-2</i>	<i>E. coli</i>	<i>Yeast</i>	<i>Green Algae</i>	<i>Human</i>
File I/O	0.00	0.00	0.00	0.00	0.00
Signal-to-Event	21.75	1.86	1.01	0.53	0.02
Sketching	0.74	0.06	0.04	0.03	0.00
Seeding	3.86	4.14	3.52	6.70	5.39
Chaining	73.50	93.92	95.42	92.43	94.46
Seeding + Chaining	77.36	98.06	98.94	99.14	99.86

The entire runtime is **bottlenecked by the chaining step**

# Required Computation Resources in Indexing

Tool	<i>Contamination</i>	<i>SARS-CoV-2</i>	<i>E. coli</i>	<i>Yeast</i>	<i>Green Algae</i>	<i>Human</i>	<i>Relative Abundance</i>
CPU Time (sec)							
UNCALLED	8.72	9.00	11.08	18.62	285.88	4,148.10	4,382.38
Sigmap	0.02	0.04	8.66	24.57	449.29	36,765.24	40,926.76
RawHash	0.18	0.13	2.62	4.48	34.18	1,184.42	788.88
Real time (sec)							
UNCALLED	1.01	1.04	2.67	7.79	280.27	4,190.00	4,471.82
Sigmap	0.13	0.25	9.31	25.86	458.46	37,136.61	41,340.16
RawHash	0.14	0.10	1.70	2.06	15.82	278.69	154.68
Peak memory (GB)							
UNCALLED	0.07	0.07	0.13	0.31	11.96	48.44	47.81
Sigmap	0.01	0.01	0.40	1.04	8.63	227.77	238.32
RawHash	0.01	0.01	0.35	0.76	5.33	83.09	152.80

The indexing step of RawHash is **orders of magnitude faster** than the indexing steps of UNCALLED and Sigmap, especially **for large genomes**

RawHash requires **larger memory space** than UNCALLED

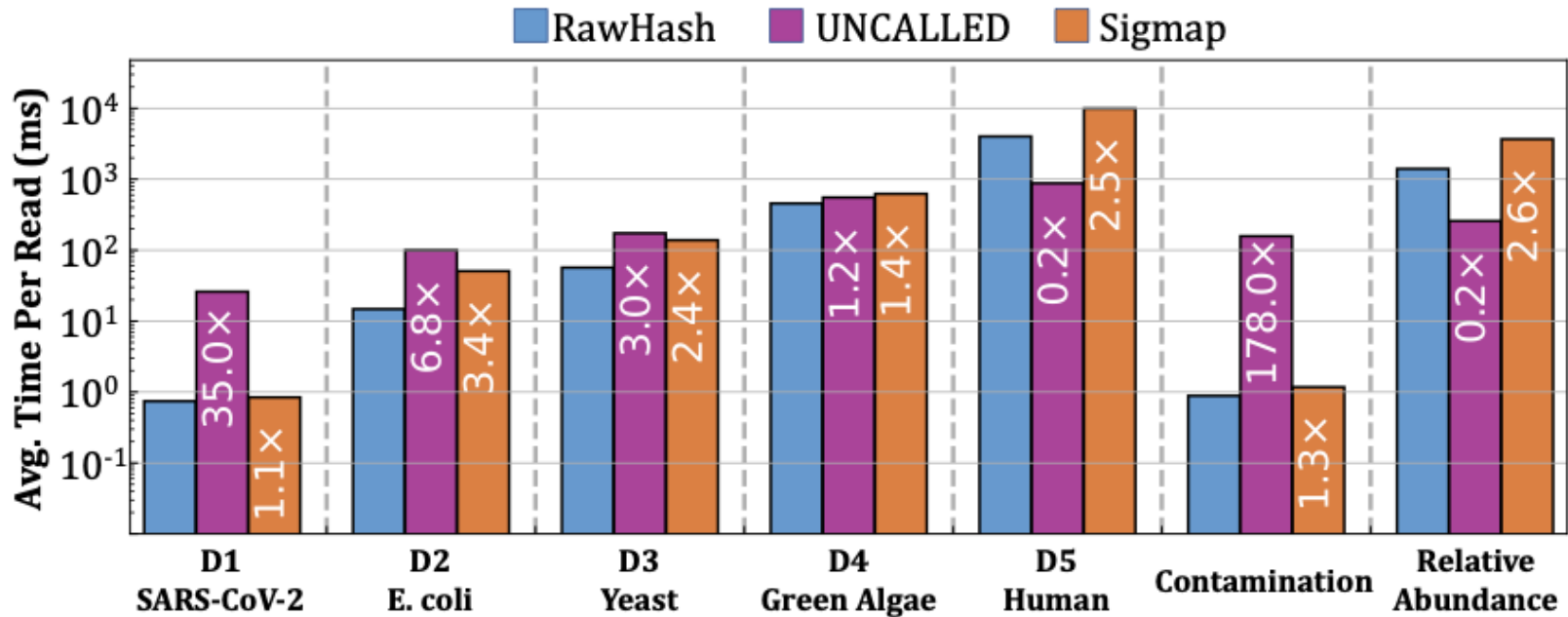
# Required Computation Resources in Mapping

Tool	Contamination	SARS-CoV-2	<i>E. coli</i>	Yeast	Green Algae	Human	Relative Abundance
CPU Time (sec)							
UNCALLED	265,902.26	36,667.26	35,821.14	8,933.52	16,769.09	262,597.83	586,561.54
Sigmap	4,573.18	1,997.84	23,894.70	11,168.96	31,544.55	4,837,058.90	11,027,652.91
RawHash	3,721.62	1,832.56	8,212.17	4,906.70	25,215.23	2,022,521.48	4,738,961.77
Real time (sec)							
UNCALLED	20,628.57	2,794.76	1,544.68	285.42	2,138.91	8,794.30	19,409.71
Sigmap	6,725.26	3,222.32	2,067.02	1,167.08	2,398.83	158,904.69	361,443.88
RawHash	3,917.49	1,949.53	957.13	215.68	1,804.96	65,411.43	152,280.26
Peak memory (GB)							
UNCALLED	0.65	0.19	0.52	0.37	0.81	9.46	9.10
Sigmap	111.69	28.26	111.11	14.65	29.18	311.89	489.89
RawHash	4.13	4.20	4.16	4.37	11.75	52.21	55.31

The mapping step of RawHash is **significantly faster than Sigmap** for all genomes, and **faster than UNCALLED for small genomes**

RawHash requires **larger memory space** than UNCALLED

# Average Mapping Time per Read



The mapping step of RawHash is **significantly faster than Sigmap** for all genomes, and **faster than UNCALLED for small genomes**

# Parameter Configurations

<b>Tool</b>	<b><i>Contamination</i></b>	<b><i>SARS-CoV-2</i></b>	<b><i>E. coli</i></b>	<b><i>Yeast</i></b>	<b><i>Green Algae</i></b>	<b><i>Human</i></b>	<b><i>Relative Abundance</i></b>
RawHash	-x viral -t 32	-x viral -t 32	-x sensitive -t 32	-x sensitive -t 32	-x fast -t 32	-x fast -t 32	-x fast -t 32
UNCALLED				map -t 32			
Sigmap				-m -t 32			
Minimap2				-x map-ont -t 32			

<b>Preset (-x)</b>	<b>Corresponding parameters</b>	<b>Usage</b>
viral	-e 5 -q 9 -l 3	Viral genomes
sensitive	-e 6 -q 9 -l 3	Small genomes (i.e., < 50M bases)
fast	-e 7 -q 9 -l 3	Large genomes (i.e., > 50M bases)

# Versions

<b>Tool</b>	<b>Version</b>	<b>Link to the Source Code</b>
RawHash	0.9	<a href="https://github.com/CMU-SAFARI/RawHash/tree/8042b1728e352a28fcc79c2efd80c8b631fe7bac">https://github.com/CMU-SAFARI/RawHash/tree/8042b1728e352a28fcc79c2efd80c8b631fe7bac</a>
UNCALLED	2.2	<a href="https://github.com/skovaka/UNCALLED/tree/74a5d4e5b5d02fb31d6e88926e8a0896dc3475cb">https://github.com/skovaka/UNCALLED/tree/74a5d4e5b5d02fb31d6e88926e8a0896dc3475cb</a>
Sigmap	0.1	<a href="https://github.com/haowenz/sigmap/tree/c9a40483264c9514587a36555b5af48d3f054f6f">https://github.com/haowenz/sigmap/tree/c9a40483264c9514587a36555b5af48d3f054f6f</a>
Minimap2	2.24	<a href="https://github.com/lh3/minimap2/releases/tag/v2.24">https://github.com/lh3/minimap2/releases/tag/v2.24</a>