# PIM Architectures for Bioinformatics

Dr. Konstantina Koliogeorgi

kkoliogeorgi@ethz.ch

https://people.inf.ethz.ch/

ICS 2025

08 June 2025

**SAFARI**

**ETH** *zürich*

# Brief Self Introduction

- **Konstantina Koliogeorgi**
  - ❑ Senior Researcher and Lecturer @ SAFARI
  - ❑ PhD, National Technical University of Athens, 2023
  - ❑ kkoliogeorgi@safari.ethz.ch

- **Research & Teaching Areas**
  - ❑ Hardware/Software Co-Design
  - ❑ Heterogeneous System Architecture
  - ❑ Reconfigurable Computing and Architectures
  - ❑ Hardware Acceleration
  - ❑ Optimized Architectures for Genome analysis
  - ❑ High Level Synthesis Tools
  - ❑ Design Space Exploration

# Agenda

- **Brief Introduction** to Genomics

- **Data Movement Bottlenecks** during analysis

- Designing **algorithms** and **architectures** that tackle data movement overhead
  - Target Multiple Steps of Pipeline
  - Leverage Processing-In-Memory
  - Leverage In-Storage Processing

**SAFARI**

# Faster, Scalable & Accurate Genome Analysis

**Uncovering and treating diseases**
linked to genomic variations
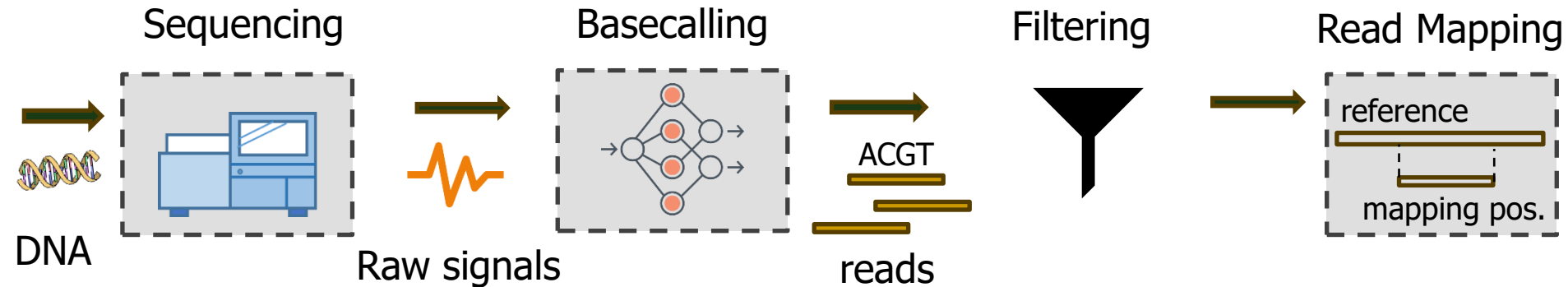
**Altering genomes** to solve
fundamental challenges of life

Detecting **pathogens**
in the environment

Rapid surveillance of
**disease outbreaks**

**And, many, many other applications ...**

# Typical Genome Sequence Analysis



Sequencing

Basecalling

Filtering

Read Mapping

DNA

Raw signals

ACGT

reads

reference

mapping pos.

## Sequencing
extract small fragments of the original DNA sequence

## Basecalling
convert raw signals to DNA bases

**Deep Neural Networks**

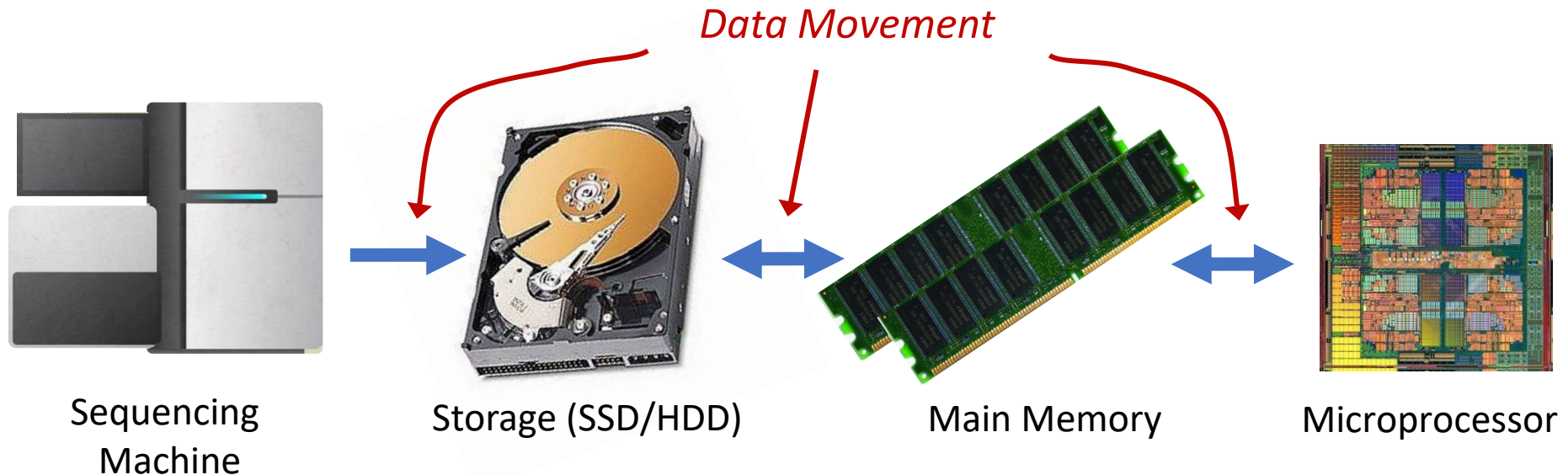## Read mapping
aligns reads to potential matching locations in the reference genome

**SAFARI**

# Significant barrier
# to genome analyses

# Data Movement Dominates Performance

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

*Data Movement*

Sequencing Machine → Storage (SSD/HDD) ↔ Main Memory ↔ Microprocessor

Single memory request consumes **>160× - 800× 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

**SAFARI**

# Data analysis is performed far away from the data

**SAFARI**

We need to orchestrate algorithms and architectures to handle data well

**SAFARI**

# Genomic Analysis Steps in Memory

We need to design algorithms
that fit processing-in-memory

# Processing Using Memory

**SAFARI**

https://www.youtube.com/watch?v=HNd4skQrt6I

# Processing Near Memory



Computer Architecture - Lecture 8: Processing near Memory (Fall 2021)

759 views • Streamed live on Oct 22, 2021

# Using Real PIM System



Computer Architecture - Lecture 9: Real PIM Systems: UPMEM Case Study (Fall 2021)

https://www.youtube.com/watch?v=TuVw_SKaTCo

# Raw Signal Translation using PIM [MICRO '23]

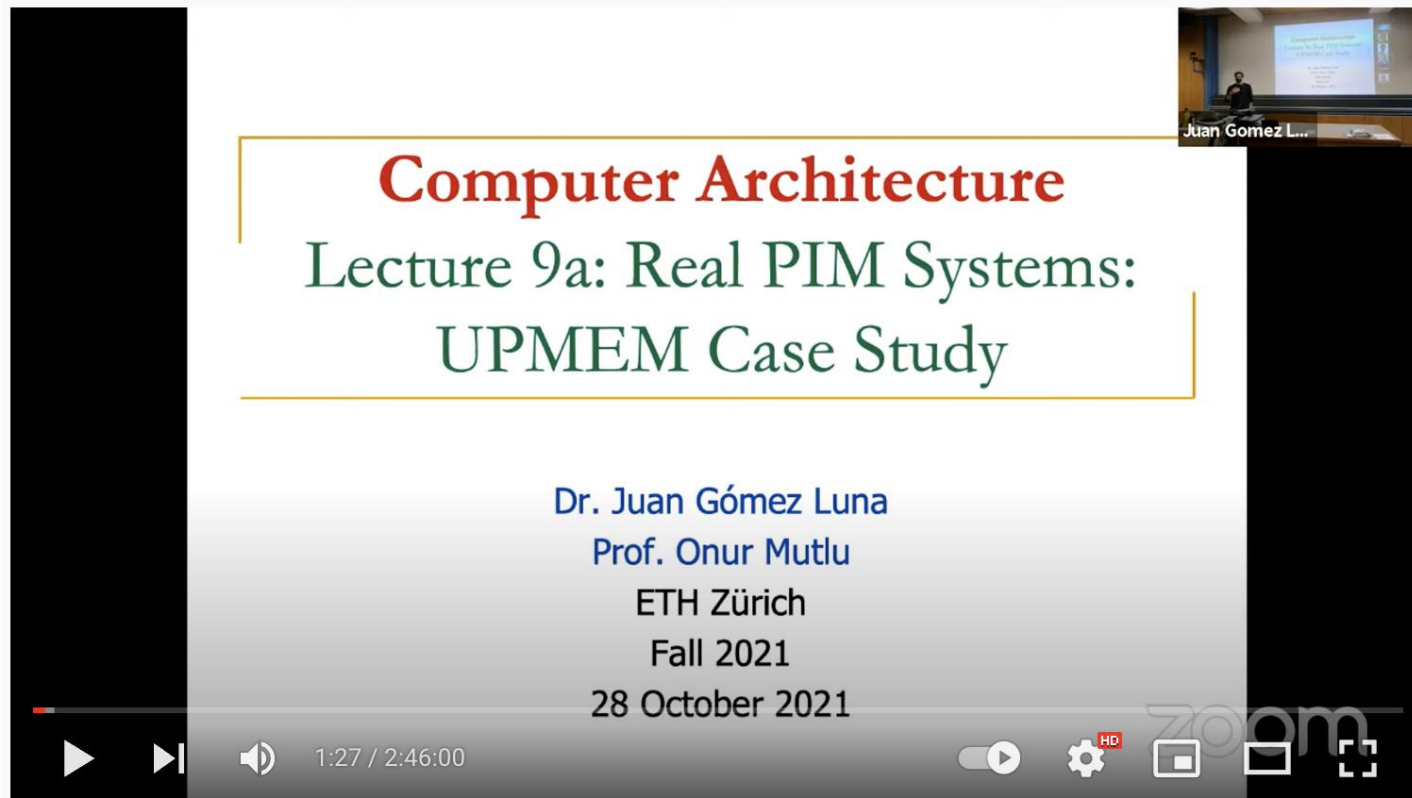- Taha Shahroodi, Gagandeep Singh, Mahdi Zahedi, Haiyu Mao, Joel Lindegger, Can Firtina, Stephan Wong, Onur Mutlu, and Said Hamdioui, **"Swordfish: A Framework for Evaluating Deep Neural Network-based Basecalling using Computation-In-Memory with Non-Ideal Memristors"** *Proceedings of the 56th International Symposium on Microarchitecture (MICRO), Toronto, ON, Canada, November 2023.* [Slides (pptx) (pdf)] [arXiv version]

## Swordfish: A Framework for Evaluating Deep Neural Network-based Basecalling using Computation-In-Memory with Non-Ideal Memristors

Taha Shahroodi[1]    Gagandeep Singh[2,3]    Mahdi Zahedi[1]    Haiyu Mao[3]    Joel Lindegger[3]    Can Firtina[3]
Stephan Wong[1]    Onur Mutlu[3]    Said Hamdioui[1]

[1]TU Delft    [2]AMD Research    [3]ETH Zürich

# Using PIM for filtering

- Jeremie S. Kim, Damla Senol Cali, Hongyi Xin, Donghyuk Lee, Saugata Ghose, Mohammed Alser, Hasan Hassan, Oguz Ergin, Can Alkan, and Onur Mutlu,
**"GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies"**
*BMC Genomics*, 2018.
*Proceedings of the 16th Asia Pacific Bioinformatics Conference* (**APBC**), Yokohama, Japan, January 2018.
arxiv.org Version (pdf)

## GRIM-Filter: Fast seed location filtering in DNA read mapping using processing-in-memory technologies

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

*From* The Sixteenth Asia Pacific Bioinformatics Conference 2018
Yokohama, Japan. 15-17 January 2018

# GRIM-Filter in 3D-Stacked DRAM

1. **Highly Parallel mechanism**

2. **Memory Bound:** Given the frequent accesses to memory, we find that GRIM-Filter is memory bound

**These properties together make GRIM-Filter a good algorithm to be run in 3D-Stacked DRAM**

SAFARI

# AIM (PIM Sequence Alignment Framework)

Safaa Diab, Amir Nassereldine, Mohammed Alser, Juan Gómez-Luna,
Onur Mutlu, Izzat El Hajj
"[A Framework for High-throughput Sequence Alignment using Real Processing-in-Memory Systems](#)"
arXiv, 2022
[[Source code](#)]

# Accelerating Sequence-to-Graph Mapping

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

[1]Bionano Genomics   [2]ETH Zürich   [3]Bilkent University   [4]Intel Labs
[5]Carnegie Mellon University   [6]University of Illinois Urbana-Champaign

# Genome Sequence Analysis

- Mapping the reads to a reference genome (i.e., **read mapping**) is a c*ritical step* in genome sequence analysis

**Linear Reference:** ACG**T**ACGT

**Read:**  ACG**G**

Alternative Sequence: ACG**G**ACGT

Alternative Sequence: ACG**TT**ACGT

Alternative Sequence: ACG–ACGT

*Sequence-to-Sequence (S2S) Mapping*

**Graph-based Reference:**



**Read:** ACGG

*Sequence-to-Graph (S2G) Mapping*

*Sequence-to-graph mapping* results in **notable quality improvements.**
However, it is a **more difficult** computational problem,
with **no prior hardware design.**

# Analysis of State-of-the-Art Tools

Based on our analysis with **GraphAligner** and **vg:** SW

**Observation 1:** Alignment step is the bottleneck

**Observation 2:** **Alignment suffers from high cache miss rates**

**Observation 3:** **Seeding suffers from the DRAM latency bottleneck**

**Observation 4:** Baseline tools scale sublinearly

HW

**Observation 5:** Existing S2S mapping accelerators are unsuitable for the S2G mapping problem

**Observation 6:** Existing graph accelerators are unable to handle S2G alignment

# SeGraM: First Graph Mapping Accelerator

**Our Goal:**

**Specialized, high-performance, scalable, and low-cost** algorithm/hardware co-design that alleviates bottlenecks in **multiple steps** of sequence-to-graph mapping

**SeGraM:** *First universal algorithm/hardware co-designed genomic mapping accelerator* that can effectively and efficiently support:

- Sequence-to-graph mapping
- Sequence-to-sequence mapping
- Both short and long reads

# Overall System Design of SeGraM

**High-Bandwidth Memory (HBM):** Enables low-latency and highly-parallel memory access

# Accelerating Basecalling + Read Mapping via PIM

- 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[1]  Mohammed Alser[1]  Mohammad Sadrosadati[1]  Can Firtina[1]  Akanksha Baranwal[1]
Damla Senol Cali[2]  Aditya Manglik[1]  Nour Almadhoun Alserr[1]  Onur Mutlu[1]

[1]ETH Zürich     [2]Bionano Genomics

# Overview: Two Limitations

**Multiple steps in genome analysis**

**Large data movement** between multiple steps

A lot of **wasted computation** done on data that is later discovered to be **useless**

# Limitation 1: Large Data Movement

❑ Using a human dataset in [NC'19] as an example:

| Raw Signals | Basecalling | Reads | Read quality control | High-quality reads | Read mapping | Mapped reads |

**3913 GB**        **546 GB**        **437 GB**              **382 GB**

**Large data movement** between genome analysis steps

[NC'19] Rory Bowden, Robert W Davies, Andreas Heger, Alistair T Pagnamenta, Mariateresa de Cesare, Laura E Oikkonen, Duncan Parkes, Colin Freeman, Fatima Dhalla, Smita Y Patel, et al. Sequencing of human genomes with nanopore technology. Nature Communications, 2019.

# Limitation 2: Wasted Computation

❑ Using a human dataset in [NC'19] as an example:

**100%**          **79.5%**        **69.5%**

| Raw Signals | Basecalling | Reads | Read quality control | High-quality reads | Read mapping | Mapped reads |

**✖ Low-quality reads**
**20.5%**

**✖ Unmapped reads**
**10%**

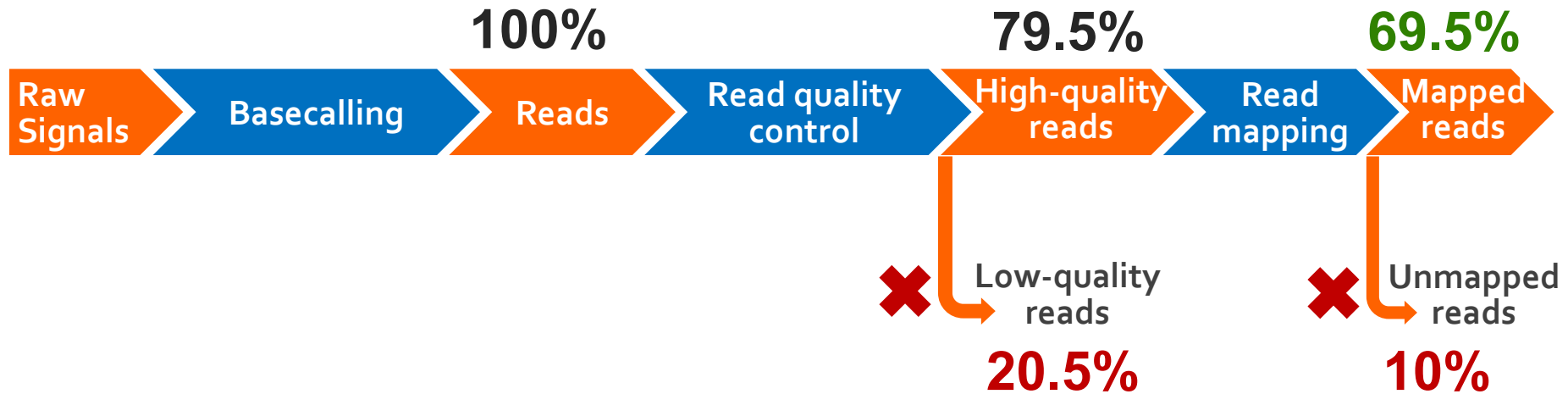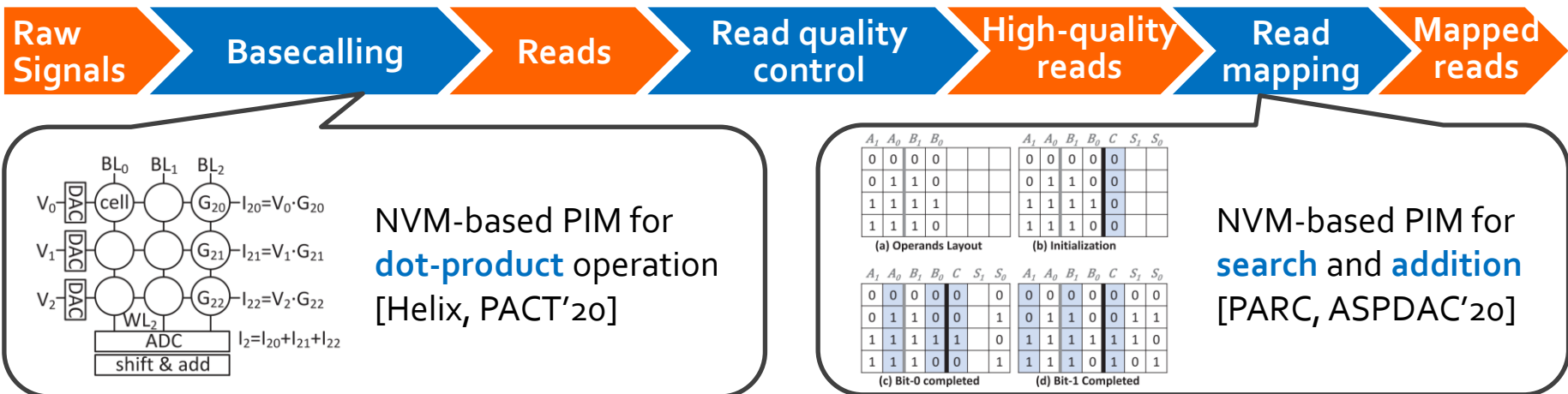> **A considerable amount of computation on useless data due to**
> - **Low-quality reads**
> - **Unmapped reads**

[NC'19] Rory Bowden, Robert W Davies, Andreas Heger, Alistair T Pagnamenta, Mariateresa de Cesare, Laura E Oikkonen, Duncan Parkes, Colin Freeman, Fatima Dhalla, Smita Y Patel, et al. Sequencing of human genomes with nanopore technology. Nature Communications, 2019.

# State-of-the-art Works

❑ NVM-based PIM is an efficient technique to reduce data movement by processing data using or near memory

| Raw Signals | Basecalling | Reads | Read quality control | High-quality reads | Read mapping | Mapped reads |

NVM-based PIM for **dot-product** operation [Helix, PACT'20]

NVM-based PIM for **search** and **addition** [PARC, ASPDAC'20]

o Reduce the data movement in a single genome analysis step
o Exacerbate the data movement overhead between analysis steps

**No prior work tackles data movement between analysis steps and reduces useless computation**

# Goal and Opportunities

> **Goal:** Efficiently accelerate the entire genome analysis pipeline while **minimizing data movement and useless computation**

❑ We perform a study to quantify potential performance benefits

○ Results are normalized to the performance of GPU

# Overview: GenPIP

❑ **GenPIP:** A fast and energy-efficient **in-memory** acceleration system for the <u>Gen</u>ome analysis <u>PIP</u>eline via **tight integration of genome analysis steps**

❑ **GenPIP** has two key techniques

- ○ **Chunk-based pipeline (CP)**
  - ▪ **Provides fine-grained collaboration** of genome analysis steps

- ○ **Early rejection (ER)**
  - ▪ **Timely stops the execution on useless data** by predicting which reads will not be useful

❑ **GenPIP** outperforms state-of-the-art software & hardware solutions using **CPU**, **GPU**, and **optimistic PIM** by **41.6×**, **8.4x**, and **1.4x**, respectively.

# Key Results – Performance



GenPIP provides **41.6x**, **8.4x**, and **1.4x** speedup over CPU, GPU, and optimistic PIM

**Both CP and ER are critical** to the speedup

# Key Results – Energy Efficiency



GenPIP provides **32.8x**, **20.8x**, and **1.37x** energy savings over CPU, GPU, and optimistic PIM

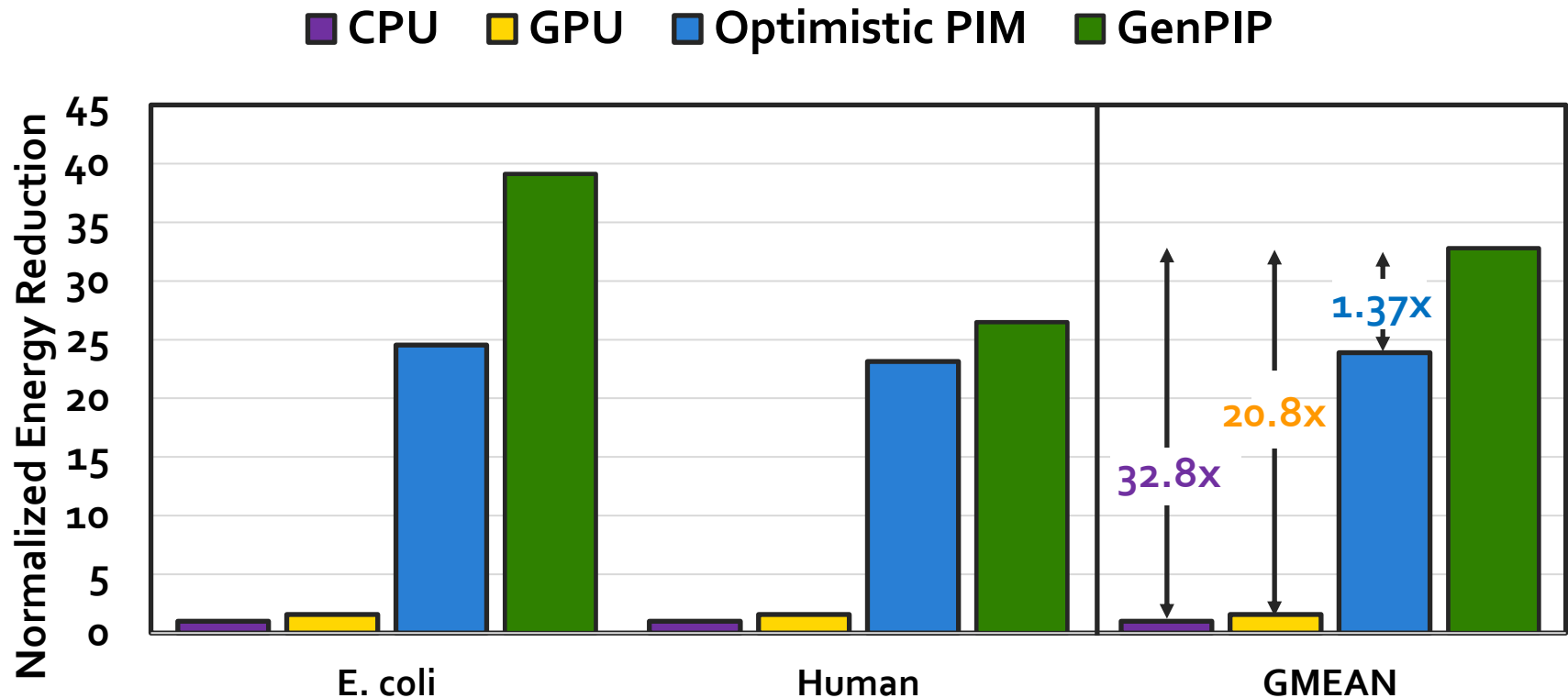**ER is especially critical** to the energy efficiency

# Accelerating Basecalling + Read Mapping via PIM

- 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[1]    Mohammed Alser[1]    Mohammad Sadrosadati[1]    Can Firtina[1]    Akanksha Baranwal[1]
Damla Senol Cali[2]    Aditya Manglik[1]    Nour Almadhoun Alserr[1]    Onur Mutlu[1]
[1]ETH Zürich        [2]Bionano Genomics

# Can we **process data closer** to where it is **stored**?

# 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[1]    Jisung Park[1]    Harun Mustafa[1]    Jeremie Kim[1]    Ataberk Olgun[1]
Arvid Gollwitzer[1]    Damla Senol Cali[2]    Can Firtina[1]    Haiyu Mao[1]    Nour Almadhoun Alserr[1]
Rachata Ausavarungnirun[3]    Nandita Vijaykumar[4]    Mohammed Alser[1]    Onur Mutlu[1]

[1]ETH Zürich   [2]Bionano Genomics   [3]KMUTNB   [4]University of Toronto

# Genome Sequence Analysis

**Data Movement from Storage** →

**Alignment**

| Storage System | | Main Memory | Cache | Computation Unit (CPU or Accelerator) |

✕ **Computation overhead**

✕ **Data movement overhead**

# Accelerating Genome Sequence Analysis

**Heuristics**   **Accelerators**   **Filters**

**Storage System**

**Main Memory**   **Cache**   **Computation Unit** (CPU or Accelerator)

✓ **Computation overhead**

✗ **Data movement overhead**

**SAFARI**

# Key Idea

*Filter* reads that do *not* require alignment *inside the storage system*

AAGCGTTCCTTGGCA
GGGGCCAGAATG
AACCTTTGGGTCCA
GAATGGGGCCA
TTTTCCCCGGGGCCA
GCTTCCAGAATG

**Filtered Reads**

**Main Memory**

**Cache**

**Computation Unit (CPU or Accelerator)**

**Exactly-matching reads**
Do not need expensive approximate string matching during alignment

**Non-matching reads**
Do not have potential matching locations and can skip alignment

**SAFARI**

# GenStore

*Filter* reads that do *not* require alignment *inside the storage system*

| GenStore-Enabled Storage System | | Main Memory | | Cache | Computation Unit (CPU or Accelerator) |
|---|---|---|---|---|---|

✓ **Computation overhead**

✓ **Data movement overhead**

**GenStore provides significant speedup (1.4x - 33.6x) and energy reduction (3.9x – 29.2x) at low cost**

# 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[1]   Jisung Park[1]   Harun Mustafa[1]   Jeremie Kim[1]   Ataberk Olgun[1]
Arvid Gollwitzer[1]   Damla Senol Cali[2]   Can Firtina[1]   Haiyu Mao[1]   Nour Almadhoun Alserr[1]
Rachata Ausavarungnirun[3]   Nandita Vijaykumar[4]   Mohammed Alser[1]   Onur Mutlu[1]

[1]ETH Zürich   [2]Bionano Genomics   [3]KMUTNB   [4]University of Toronto

# In-Storage Metagenomics [ISCA 2024]

- Nika Mansouri Ghiasi, Mohammad Sadrosadati, Harun Mustafa, Arvid Gollwitzer, Can Firtina, Julien Eudine, Haiyu Mao, Joel Lindegger, Meryem Banu Cavlak, Mohammed Alser, Jisung Park, and Onur Mutlu,
  **"MegIS: High-Performance and Low-Cost Metagenomic Analysis with In-Storage Processing"**
  Proceedings of the *51st Annual International Symposium on Computer Architecture* (**ISCA**), Buenos Aires, Argentina, July 2024.
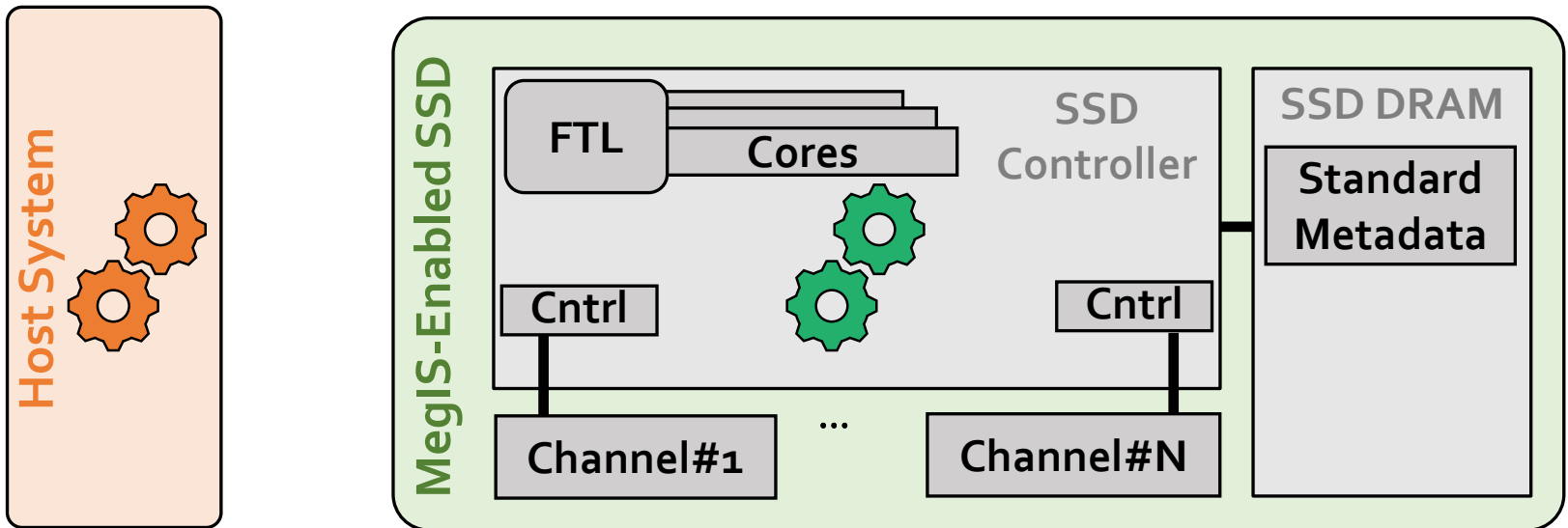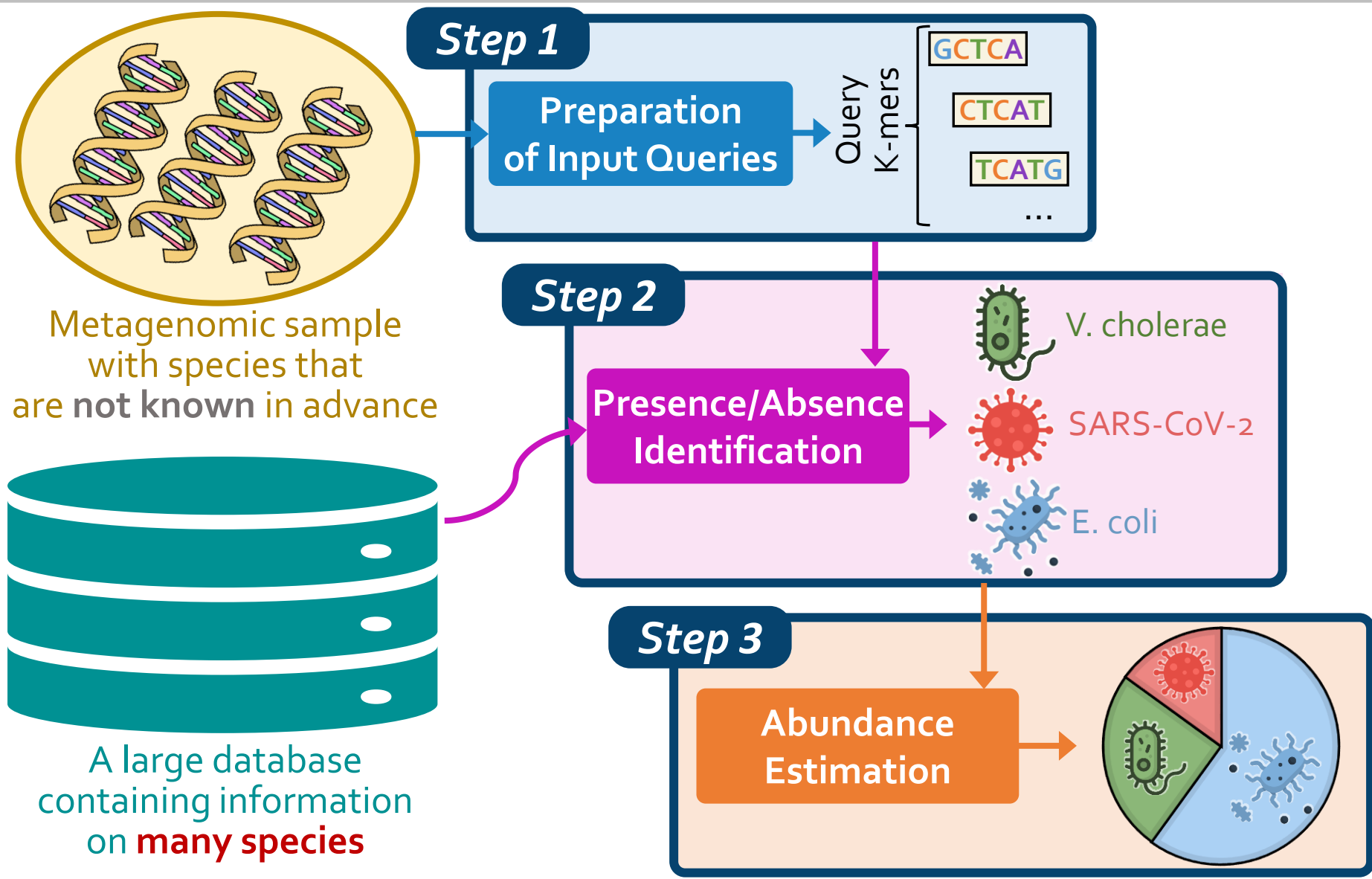  [Slides (pptx) (pdf)]
  [arXiv version]

## MegIS: High-Performance, Energy-Efficient, and Low-Cost Metagenomic Analysis with In-Storage Processing

Nika Mansouri Ghiasi[1]     Mohammad Sadrosadati[1]     Harun Mustafa[1]     Arvid Gollwitzer[1]
Can Firtina[1]     Julien Eudine[1]     Haiyu Mao[1]     Joël Lindegger[1]     Meryem Banu Cavlak[1]
Mohammed Alser[1]     Jisung Park[2]     Onur Mutlu[1]
[1]ETH Zürich     [2]POSTECH

https://arxiv.org/pdf/2406.19113
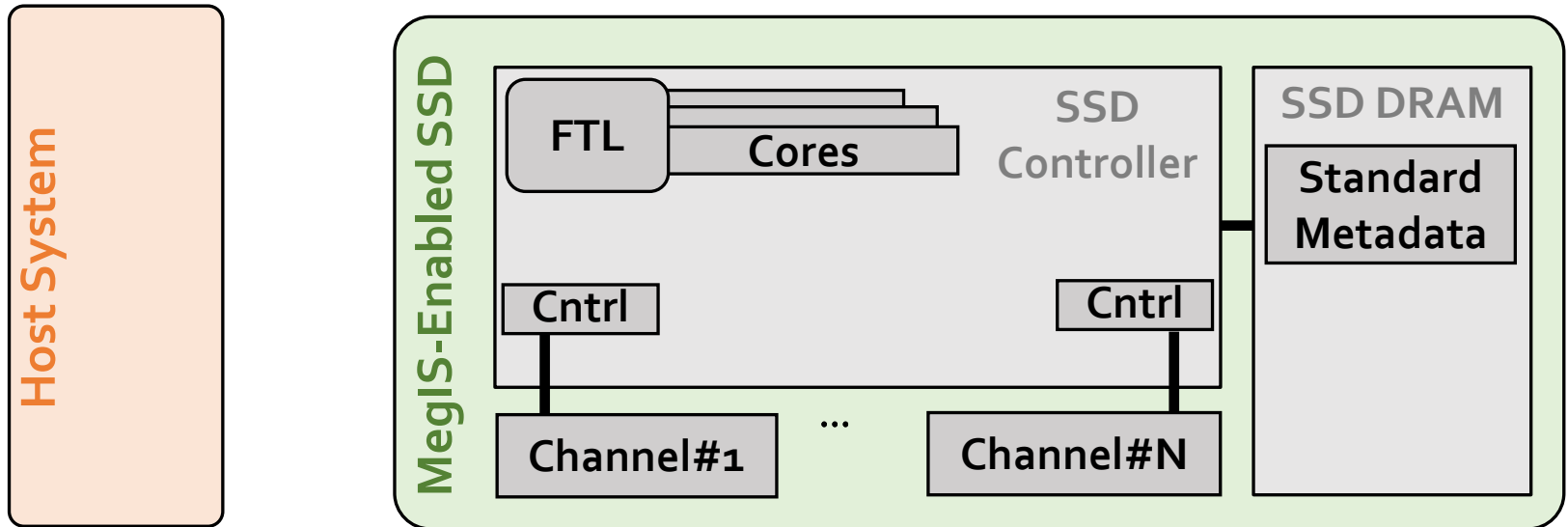
# MegIS: Metagenomics In-Storage

- First in-storage system for *end-to-end* metagenomic analysis

- **Idea:** Cooperative in-storage processing for metagenomic analysis

  - Hardware/software co-design between the **storage system** and **host system**

**SAFARI**

# MegIS's Steps



Metagenomic sample
with species that
are **not known** in advance

A large database
containing information
on **many species**

**Step 1**
Preparation
of Input Queries

Query K-mers
GCTCA
CTCAT
TCATG
...

**Step 2**
Presence/Absence
Identification

V. cholerae
SARS-CoV-2
E. coli

**Step 3**
Abundance
Estimation
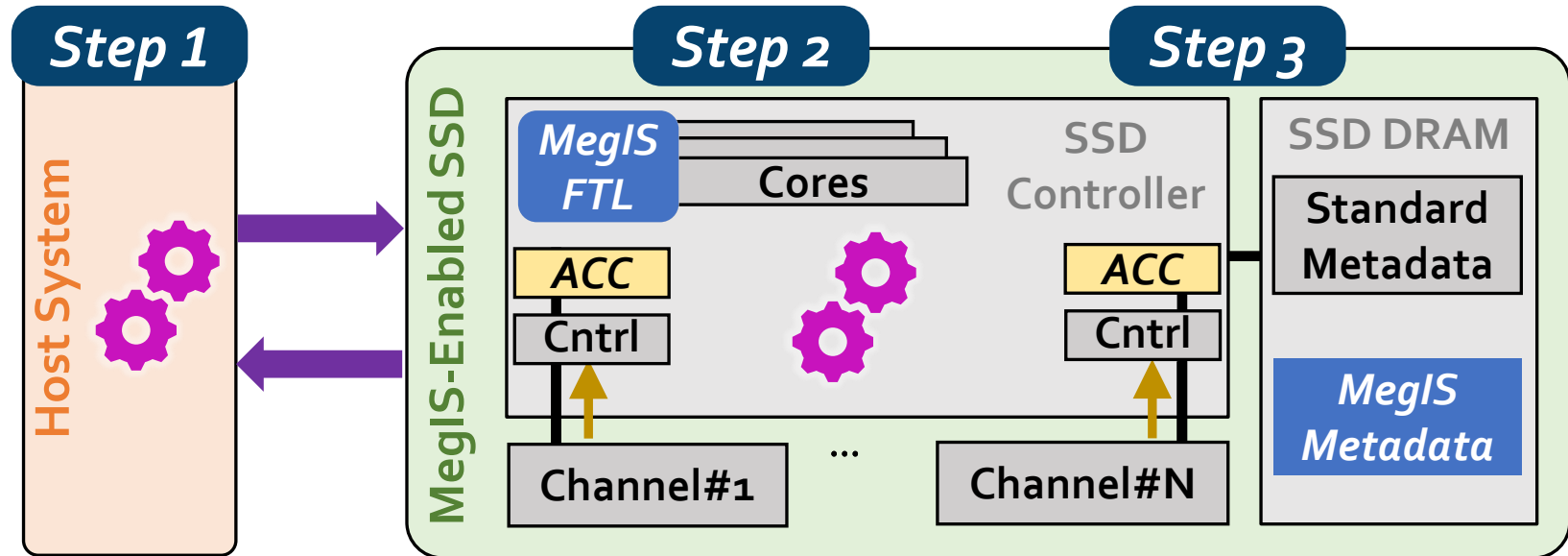
SAFARI

# MegIS Hardware-Software Co-Design

# MegIS Hardware-Software Co-Design

**Task partitioning and mapping**
- *Each step executes in its most suitable system*

**Data/computation flow coordination**
- *Reduce communication overhead*
- *Reduce #writes to flash chips*



**Storage-aware algorithms**
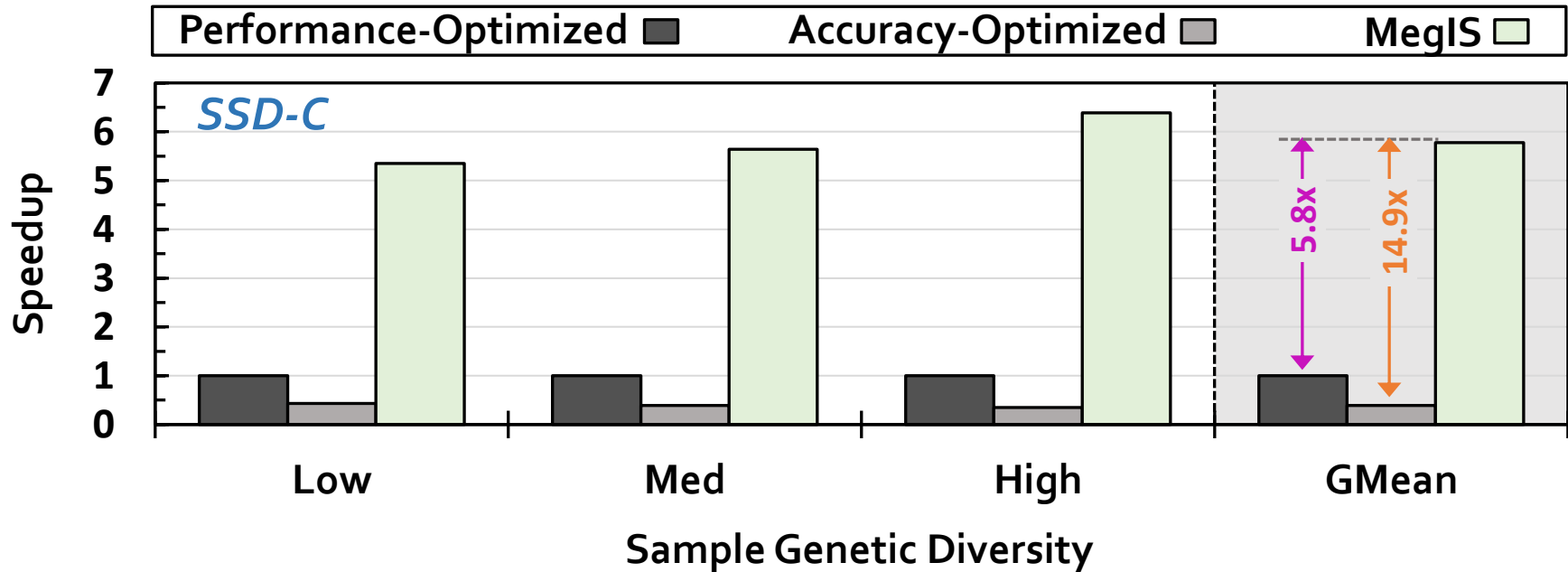- *Enable efficient access patterns to the SSD*

**Lightweight in-storage accelerators**
- *Minimize SRAM/DRAM buffer spaces needed inside the SSD*

**Data mapping scheme and Flash Translation Layer (FTL)**
- *Specialize to the characteristics of metagenomic analysis*
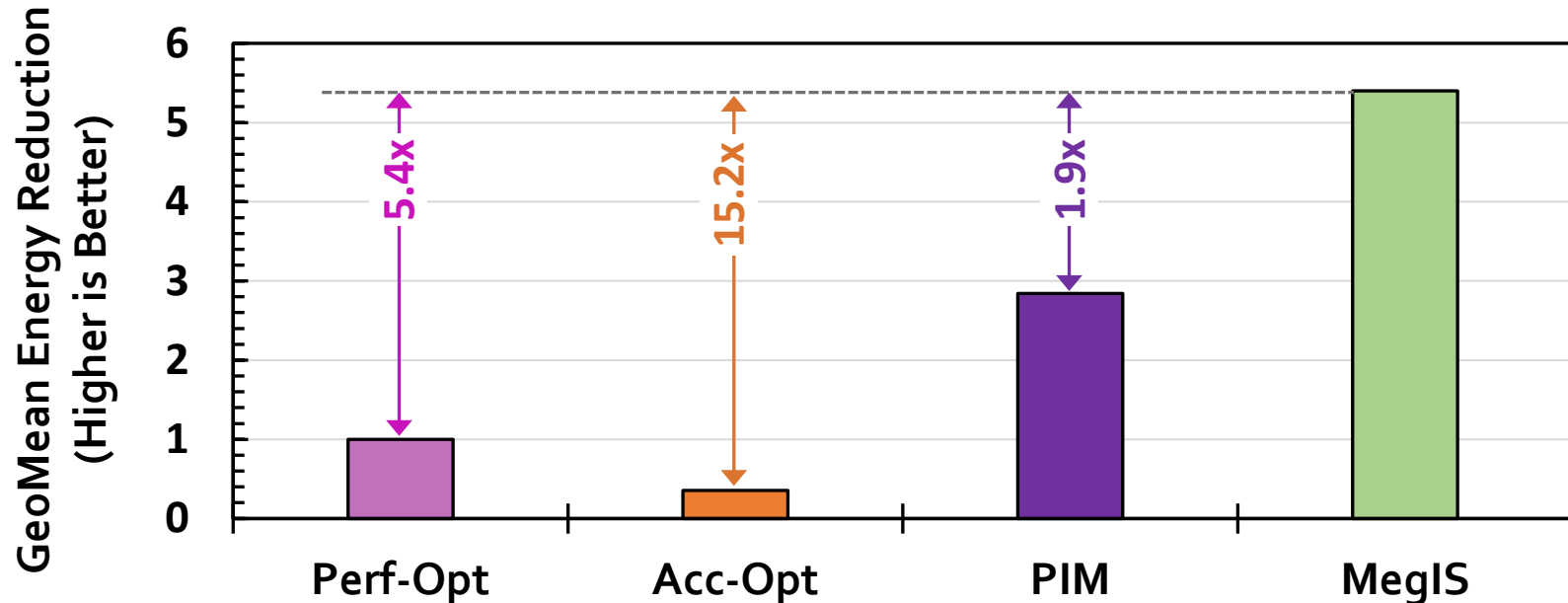- *Leverage the SSD's full internal bandwidth*

SAFARI

# Evaluation: Speedup over the Software Baselines



**MegIS provides significant speedup over both Performance-Optimized and Accuracy-Optimized baselines**

- On average across different input sets and SSDs



MegIS provides significant energy reduction over

the **Performance-Optimized**, **Accuracy-Optimized**, and **PIM** baselines

# In-Storage Metagenomics [ISCA 2024]

- Nika Mansouri Ghiasi, Mohammad Sadrosadati, Harun Mustafa, Arvid Gollwitzer, Can Firtina, Julien Eudine, Haiyu Mao, Joel Lindegger, Meryem Banu Cavlak, Mohammed Alser, Jisung Park, and Onur Mutlu,
  **"MegIS: High-Performance and Low-Cost Metagenomic Analysis with In-Storage Processing"**
  *Proceedings of the 51st Annual International Symposium on Computer Architecture (**ISCA**)*, Buenos Aires, Argentina, July 2024.
  [Slides (pptx) (pdf)]
  [arXiv version]

## MegIS: High-Performance, Energy-Efficient, and Low-Cost Metagenomic Analysis with In-Storage Processing

Nika Mansouri Ghiasi[1]    Mohammad Sadrosadati[1]    Harun Mustafa[1]    Arvid Gollwitzer[1]
Can Firtina[1]    Julien Eudine[1]    Haiyu Mao[1]    Joël Lindegger[1]    Meryem Banu Cavlak[1]
Mohammed Alser[1]    Jisung Park[2]    Onur Mutlu[1]
[1]ETH Zürich    [2]POSTECH

**SAFARI**

https://arxiv.org/pdf/2406.19113

# Conclusion

- System design for bioinformatics is a critical problem
  - It has large scientific, medical, societal, personal implications

- This talk is about accelerating genomics by alleviating data movement bottleneck

- We covered various recent works on individual algorithms and pipelines
  - **PnM, PuM, ISP**

- **Many future opportunities exist**
  - **Especially with new sequencing technologies**
  - **Especially with new applications and use cases**

**SAFARI**

# PIM Architectures
# for Bioinformatics

Dr. Konstantina Koliogeorgi

kkoliogeorgi@ethz.ch

https://people.inf.ethz.ch/

ICS 2025

08 June 2025

**SAFARI**

**ETH** zürich