ML/AI for Memory System Design & Memory System Design for AI/ML

Onur Mutlu <u>omutlu@gmail.com</u> <u>https://people.inf.ethz.ch/omutlu</u>

23 August 2024 Intel ArchFest







Computing is Bottlenecked by Data



Data is Key for AI, ML, Genomics, ...

Important workloads are all data intensive

 They require rapid and efficient processing of large amounts of data

- Data is increasing
 - We can generate more than we can process
 - We need to perform more sophisticated analyses on more data





Huge Demand for Performance & Efficiency

Exponential Growth of Neural Networks



Huge Demand for Performance & Efficiency



http://www.economist.com/news/21631808-so-much-genetic-data-so-many-uses-genes-unzipped

Data Overwhelms Modern Machines ...

Storage/memory capability

Communication capability

Computation capability

Greatly impacts robustness, energy, performance, cost

Data Movement Overwhelms Modern Machines

 Amirali Boroumand, Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, and Onur Mutlu, "Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks" Proceedings of the 23rd International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), Williamsburg, VA, USA, March 2018.

62.7% of the total system energy is spent on data movement

Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

Amirali Boroumand¹Saugata Ghose¹Youngsok Kim²Rachata Ausavarungnirun¹Eric Shiu³Rahul Thakur³Daehyun Kim^{4,3}Aki Kuusela³Allan Knies³Parthasarathy Ranganathan³Onur Mutlu^{5,1}7

Data Movement Overwhelms Accelerators

Amirali Boroumand, Saugata Ghose, Berkin Akin, Ravi Narayanaswami, Geraldo F. Oliveira, Xiaoyu Ma, Eric Shiu, and Onur Mutlu, "Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks" Proceedings of the 30th International Conference on Parallel Architectures and Compilation Techniques (PACT), Virtual, September 2021. [Slides (pptx) (pdf)] [Talk Video (14 minutes)]

> 90% of the total system energy is spent on memory in large ML models

Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks

Saugata Ghose[‡] Amirali Boroumand[†] Xiaoyu Ma[§] Geraldo F. Oliveira^{*}

Berkin Akin[§] Eric Shiu[§]

Ravi Narayanaswami[§] Onur Mutlu^{*†}

[§]Google [†]Carnegie Mellon Univ. [•]Stanford Univ. [‡]Univ. of Illinois Urbana-Champaign

*ETH Zürich

https://arxiv.org/pdf/2109.14320

Example Energy Breakdowns



In LSTMs and Transducers used by Google, >90% energy spent on off-chip interconnect and DRAM

https://arxiv.org/pdf/2109.14320



An Intelligent Architecture Handles Data Well



Corollaries: Computing Systems Today ...

Are processor-centric vs. data-centric

Make designer-dictated decisions vs. data-driven

Make component-based myopic decisions vs. data-aware

Architectures for Intelligent Machines

Data-centric

Data-driven

Data-aware



A Blueprint for Fundamentally Better Architectures

Onur Mutlu, "Intelligent Architectures for Intelligent Computing Systems" Invited Paper in Proceedings of the <u>Design, Automation, and Test in</u> <u>Europe Conference</u> (DATE), Virtual, February 2021. [Slides (pptx) (pdf)] [IEDM Tutorial Slides (pptx) (pdf)] [Short DATE Talk Video (11 minutes)] [Longer IEDM Tutorial Video (1 hr 51 minutes)] [Virtual Virtual Video (1 hr 51 minutes)] [Virtual Virtual Virtual

Intelligent Architectures for Intelligent Computing Systems

Onur Mutlu ETH Zurich omutlu@gmail.com

1. Memory system design for AI/ML workloads/accelerators

 \rightarrow in-depth exploration of memory system designs for cutting-edge and emerging machine learning accelerators \rightarrow more efficient on-chip and off-chip memory systems

2. AI/ML techniques for improving memory system designs

 \rightarrow comprehensive look at memory system design to make it data driven, i.e., based on machine learning

→ more effective cache/memory/prefetch/thread controllers and data/resource management/mapping/scheduling policies

1. Memory system design for AI/ML workloads/accelerators

2. AI/ML techniques for improving memory system designs

Data-Driven (Self-Optimizing) Architectures

System Architecture Design Today

- Human-driven
 - Humans design the policies (how to do things)
- Many (too) simple, short-sighted policies all over the system
- No automatic data-driven policy learning
- (Almost) no learning: cannot take lessons from past actions

Can we design fundamentally intelligent architectures?

An Intelligent Architecture

- Data-driven
 - Machine learns the "best" policies (how to do things)
- Sophisticated, workload-driven, changing, far-sighted policies
- Automatic data-driven policy learning
- All controllers are intelligent data-driven agents

We need to rethink design (of all controllers)

Self-Optimizing Memory Controllers

 Engin Ipek, Onur Mutlu, José F. Martínez, and Rich Caruana, "Self Optimizing Memory Controllers: A Reinforcement Learning <u>Approach</u>" *Proceedings of the <u>35th International Symposium on Computer Architecture</u> (ISCA), pages 39-50, Beijing, China, June 2008. <i>Selected to the ISCA-50 25-Year Retrospective Issue covering 1996- 2020 in 2023 (Retrospective (pdf) Full Issue).*

Self-Optimizing Memory Controllers: A Reinforcement Learning Approach

Engin İpek^{1,2} Onur Mutlu² José F. Martínez¹ Rich Caruana¹

¹Cornell University, Ithaca, NY 14850 USA

 2 Microsoft Research, Redmond, WA 98052 USA

Self-Optimizing Memory Prefetchers

Rahul Bera, Konstantinos Kanellopoulos, Anant Nori, Taha Shahroodi, Sreenivas Subramoney, and Onur Mutlu, "Pythia: A Customizable Hardware Prefetching Framework Using Online Reinforcement Learning" Proceedings of the 54th International Symposium on Microarchitecture (MICRO), Virtual, October 2021. [Slides (pptx) (pdf)] [Short Talk Slides (pptx) (pdf)] [Lightning Talk Slides (pptx) (pdf)] [Talk Video (20 minutes)] [Lightning Talk Video (1.5 minutes)] [Pythia Source Code (Officially Artifact Evaluated with All Badges)] [arXiv version] Officially artifact evaluated as available, reusable and reproducible.



Pythia: A Customizable Hardware Prefetching Framework Using Online Reinforcement Learning

Rahul Bera¹ Konstantinos Kanellopoulos¹

Anant V. Nori² Taha Shahroodi^{3,1} Onur Mutlu¹

¹ETH Zürich ²Processor Architecture Research Labs, Intel Labs ³TU Delft

Sreenivas Subramoney²

https://arxiv.org/pdf/2109.12021.pdf

Learning-Based Off-Chip Load Predictors

 Rahul Bera, Konstantinos Kanellopoulos, Shankar Balachandran, David Novo, Ataberk Olgun, Mohammad Sadrosadati, and Onur Mutlu,
"Hermes: Accelerating Long-Latency Load Requests via Perceptron-Based Off-Chip Load Prediction"
Proceedings of the <u>55th International Symposium on Microarchitecture</u> (MICRO), Chicago, IL, USA, October 2022.
[Slides (pptx) (pdf)]
[Longer Lecture Slides (pptx) (pdf)]
[Talk Video (12 minutes)]
[Lecture Video (25 minutes)]
[arXiv version]
[Source Code (Officially Artifact Evaluated with All Badges)]
Officially artifact evaluated as available, reusable and reproducible. Best paper award at MICRO 2022.



Hermes: Accelerating Long-Latency Load Requests via Perceptron-Based Off-Chip Load Prediction

Rahul Bera1Konstantinos Kanellopoulos1Shankar Balachandran2David Novo3Ataberk Olgun1Mohammad Sadrosadati1Onur Mutlu1

¹ETH Zürich ²Intel Processor Architecture Research Lab ³LIRMM, Univ. Montpellier, CNRS

https://arxiv.org/pdf/2209.00188.pdf

Self-Optimizing Hybrid SSD Controllers

Gagandeep Singh, Rakesh Nadig, Jisung Park, Rahul Bera, Nastaran Hajinazar, David Novo, Juan Gomez-Luna, Sander Stuijk, Henk Corporaal, and Onur Mutlu, "Sibyl: Adaptive and Extensible Data Placement in Hybrid Storage Systems Using Online Reinforcement Learning" Proceedings of the <u>49th International Symposium on Computer</u> <u>Architecture (ISCA)</u>, New York, June 2022. [Sildes (pptx) (pdf)] [arXiv version] [Sibyl Source Code] [Talk Video (16 minutes)]

Sibyl: Adaptive and Extensible Data Placement in Hybrid Storage Systems Using Online Reinforcement Learning

Gagandeep Singh1Rakesh Nadig1Jisung Park1Rahul Bera1Nastaran Hajinazar1David Novo3Juan Gómez-Luna1Sander Stuijk2Henk Corporaal2Onur Mutlu11ETH Zürich2Eindhoven University of Technology3LIRMM, Univ. Montpellier, CNRS

https://arxiv.org/pdf/2205.07394.pdf

A Blueprint for Fundamentally Better Architectures

Onur Mutlu, "Intelligent Architectures for Intelligent Computing Systems" Invited Paper in Proceedings of the <u>Design, Automation, and Test in</u> <u>Europe Conference</u> (DATE), Virtual, February 2021. [Slides (pptx) (pdf)] [IEDM Tutorial Slides (pptx) (pdf)] [Short DATE Talk Video (11 minutes)] [Longer IEDM Tutorial Video (1 hr 51 minutes)] [Virtual Virtual Video (1 hr 51 minutes)] [Virtual Virtual
Intelligent Architectures for Intelligent Computing Systems

Onur Mutlu ETH Zurich omutlu@gmail.com

Fundamentally Better Architectures

Data-centric

Data-driven

Data-aware



Pythia: Prefetching using Reinforcement Learning

Self-Optimizing Memory Prefetchers

Rahul Bera, Konstantinos Kanellopoulos, Anant Nori, Taha Shahroodi, Sreenivas Subramoney, and Onur Mutlu, "Pythia: A Customizable Hardware Prefetching Framework Using Online Reinforcement Learning" Proceedings of the 54th International Symposium on Microarchitecture (MICRO), Virtual, October 2021. [Slides (pptx) (pdf)] [Short Talk Slides (pptx) (pdf)] [Lightning Talk Slides (pptx) (pdf)] [Talk Video (20 minutes)] [Lightning Talk Video (1.5 minutes)] [Pythia Source Code (Officially Artifact Evaluated with All Badges)] [arXiv version] Officially artifact evaluated as available, reusable and reproducible.



Pythia: A Customizable Hardware Prefetching Framework Using Online Reinforcement Learning

Rahul Bera¹ Konstantinos Kanellopoulos¹

Anant V. Nori² Taha Shahroodi^{3,1} Onur Mutlu¹

¹ETH Zürich ²Processor Architecture Research Labs, Intel Labs ³TU Delft

Sreenivas Subramoney²

https://arxiv.org/pdf/2109.12021.pdf



Pythia

A Customizable Hardware Prefetching Framework Using Online Reinforcement Learning

<u>Rahul Bera</u>, Konstantinos Kanellopoulos, Anant V. Nori, Taha Shahroodi, Sreenivas Subramoney, Onur Mutlu

https://github.com/CMU-SAFARI/Pythia





Mainly use one program context info. for prediction 2 Lack inherent system awareness

Lack in-silicon customizability







Why do prefetchers not perform well?





Lack of In-silicon Customizability

- Feature statically selected at design time
 - **Rigid hardware** designed specifically to exploit that feature
- No way to change program feature and/or change prefetcher's objective in silicon
 - Cannot adapt to a wide range of workload demands



Our Goal

A prefetching framework that can:

1.Learn to prefetch using multiple features and inherent system-level feedback information

2.Be **easily customized in silicon** to use different features and/or change prefetcher's objectives



Our Proposal



Pythia

Formulates prefetching as a reinforcement learning problem



Pythia is named after the oracle of Delphi, who is known for her accurate prophecies https://en.wikipedia.org/wiki/Pythia

Basics of Reinforcement Learning (RL)

• Algorithmic approach to learn to take an **action** in a given **situation** to maximize a numerical **reward**



Environment

- Agent stores Q-values for every state-action pair
 - Expected return for taking an action in a state

Given a state, selects action that provides highest Q-value
SAFARI

Formulating Prefetching as RL

What is State?

k-dimensional vector of features

 $S \equiv \{\phi_S^1, \phi_S^2, \dots, \phi_S^k\}$

• Feature = control-flow + data-flow

Control-flow examples

- PC
- Branch PC
- Last-3 PCs, ...

Data-flow examples

- Cacheline address
- Physical page number
- Delta between two cacheline addresses
- Last 4 deltas, ...

Features of memory request to address A (e.g., PC) Processor & Memory Subsystem

What is Action?

Given a demand access to address A the action is to select prefetch offset "O"

- Action-space: 127 actions in the range [-63, +63]
 - For a machine with 4KB page and 64B cacheline
- Upper and lower limits ensure prefetches do not cross physical page boundary
- A zero offset means no prefetch is generated
- We further **prune** action-space by design-space exploration

SAFARI

Prefetcher

Reward

(e.a., PC)

Prefetch from addre

A+offset (0)

What is Reward?

- Defines the **objective** of Pythia
- Encapsulates two metrics:
 - Prefetch usefulness (e.g., accurate, late, out-of-page, ...)
 - **System-level feedback** (e.g., mem. b/w usage, cache pollution, energy, ...)
- We demonstrate Pythia with memory bandwidth usage as the system-level feedback in the paper


What is Reward?

Seven distinct reward levels

- Accurate and timely (R_{AT})
- Accurate but late (R_{AL})
- Loss of coverage (R_{CL})
- Inaccurate
 - With low memory b/w usage (R_{IN}-L)
 - With high memory b/w usage (R_{IN}-H)
- No-prefetch
 - With low memory b/w usage (R_{NP}-L)
 - With high memory b/w usage(R_{NP}-H)
- Values are set at design time via automatic designspace exploration

- Can be customized further in silicon for higher performance SAFARI



Steering Pythia's Objective via Reward Values

- Example reward configuration for
 - Generating accurate prefetches
 - Making bandwidth-aware prefetch decisions



Highly prefers to generate accurate prefetches

Prefers not to prefetch if memory bandwidth usage is low

Strongly prefers not to prefetch if memory bandwidth usage is high

Steering Pythia's Objective via Reward Values

 Customizing reward values to make Pythia conservative towards prefetching



Highly prefers to generate accurate prefetches

Otherwise prefers not to prefetch



Basic Pythia Configuration

• Derived from automatic design-space exploration

• State: 2 features

- PC+Delta
- Sequence of last-4 deltas

• Actions: 16 prefetch offsets

- Ranging between -6 to +32. Including 0.

• Rewards:

- R_{IN}-H=-14; R_{IN}-L=-8; R_{CL}=-12

More Detailed Pythia Overview

- Q-Value Store: Records Q-values for *all* state-action pairs
- Evaluation Queue: A FIFO queue of recently-taken actions



Simulation Methodology

- Champsim [3] trace-driven simulator
- **150** single-core memory-intensive workload traces
 - SPEC CPU2006 and CPU2017
 - PARSEC 2.1
 - Ligra
 - Cloudsuite
- Homogeneous and heterogeneous multi-core mixes

• Five state-of-the-art prefetchers

- SPP [Kim+, MICRO'16]
- Bingo [Bakhshalipour+, HPCA'19]
- MLOP [Shakerinava+, 3rd Prefetching Championship, 2019]
- SPP+DSPatch [Bera+, MICRO'19]
- SPP+PPF [Bhatia+, ISCA'20]

Performance with Varying Core Count



Performance with Varying Core Count



Performance with Varying DRAM Bandwidth



Performance with Varying DRAM Bandwidth



Pythia outperforms prior best prefetchers for a wide range of DRAM bandwidth configurations



Performance Improvement via Customization



Performance Improvement via Customization



Pythia can extract even higher performance via customization without changing hardware



Pythia's Overhead

• 25.5 KB of total metadata storage per core

- Only simple tables
- We also model functionally-accurate Pythia with full complexity in Chisel [4] HDL



of a desktop-class 4-core Skylake processor (Xeon D2132IT, 60W)



Pythia is Open Source



https://github.com/CMU-SAFARI/Pythia

- MICRO'21 artifact evaluated
- Champsim source code + Chisel modeling code
- All traces used for evaluation

CMU-SAFARI/Pythia Public		 Unwate 	Jnwatch → 3 🖧 Star 9 ♀ Fork 2	
<> Code ⊙ Issues îì Pull reque:	sts 🕑 Actions III Projects III Wiki 🛈 Security	∠ Insights 爺 Set	tings	
్రి master 👻 రి 1 branch 📀 5 tags	Go to file	Add file - Code -	About	٤
rahulbera Github pages documentation	on 🗸 diefc65 7 hours	ago 🕚 40 commits	A customizable h framework using	ardware prefetching online reinforcemen
b ranch	Initial commit for MICRO'21 artifact evaluation	2 months ago	2021 paper by Bera and Kanellopoulos et al.	
Config	Initial commit for MICRO'21 artifact evaluation	2 months ago		
docs	Github pages documentation	7 hours ago		
experiments	Added chart visualization in Excel template	2 months ago	machine-learning	
inc inc	Updated README	8 days ago	reinforcement-learn	ning ure profetebor
prefetcher	Initial commit for MICRO'21 artifact evaluation	2 months ago	microarchitecture	cache-replacement
replacement	Initial commit for MICRO'21 artifact evaluation	2 months ago	branch-predictor	champsim-simulator
scripts	Added md5 checksum for all artifact traces to verify download	2 months ago	champsim-tracer	
src src	Initial commit for MICRO'21 artifact evaluation 2 months ago		🛱 Readme	
tracer	Initial commit for MICRO'21 artifact evaluation 2 months ago		<u>ما</u> ع View license	
🗅 .gitignore	Initial commit for MICRO'21 artifact evaluation	2 months ago	Ç∄ Cite this repository -	
CITATION.cff	Added citation file	8 days ago		
	Updated LICENSE	2 months ago	Releases 5	
LICENSE.champsim	Initial commit for MICRO'21 artifact evaluation	2 months ago	V1.3 Latest	



Pythia Talk Video



Pythia: A Customizable Prefetching Framework Using Reinforcement Learning - MICRO'21 Long Talk



SAFARI https://www.youtube.com/watch?v=6UMFRW3VFPo&list=PL5Q2soXY2Zi--0LrXSQ9sST3N0k0bXp51&index=8

A Lot More in the Pythia Paper

Rahul Bera, Konstantinos Kanellopoulos, Anant Nori, Taha Shahroodi, Sreenivas Subramoney, and Onur Mutlu, "Pythia: A Customizable Hardware Prefetching Framework Using Online Reinforcement Learning" Proceedings of the 54th International Symposium on Microarchitecture (MICRO), Virtual, October 2021. [Slides (pptx) (pdf)] [Short Talk Slides (pptx) (pdf)] [Lightning Talk Slides (pptx) (pdf)] [Talk Video (20 minutes)] [Lightning Talk Video (1.5 minutes)] [Pythia Source Code (Officially Artifact Evaluated with All Badges)] [arXiv version] Officially artifact evaluated as available, reusable and reproducible.



Pythia: A Customizable Hardware Prefetching Framework Using Online Reinforcement Learning

Rahul Bera¹ Konstantinos Kanellopoulos¹

Anant V. Nori² Taha Shahroodi^{3,1} Onur Mutlu¹

¹ETH Zürich ²Processor Architecture Research Labs, Intel Labs ³TU Delft

Sreenivas Subramoney²

https://arxiv.org/pdf/2109.12021.pdf



Pythia

A Customizable Hardware Prefetching Framework Using Online Reinforcement Learning

<u>Rahul Bera</u>, Konstantinos Kanellopoulos, Anant V. Nori, Taha Shahroodi, Sreenivas Subramoney, Onur Mutlu

https://github.com/CMU-SAFARI/Pythia





Hermes: Perceptron-Based Off-Chip Load Prediction

Learning-Based Off-Chip Load Predictors

 Rahul Bera, Konstantinos Kanellopoulos, Shankar Balachandran, David Novo, Ataberk Olgun, Mohammad Sadrosadati, and Onur Mutlu,
 "Hermes: Accelerating Long-Latency Load Requests via Perceptron-Based Off-Chip Load Prediction"
 Proceedings of the <u>55th International Symposium on Microarchitecture</u> (MICRO), Chicago, IL, USA, October 2022.
 [Slides (pptx) (pdf)]
 [Longer Lecture Slides (pptx) (pdf)]
 [Talk Video (12 minutes)]
 [Lecture Video (25 minutes)]
 [arXiv version]
 [Source Code (Officially Artifact Evaluated with All Badges)]
 Officially artifact evaluated as available, reusable and reproducible. Best paper award at MICRO 2022.



Hermes: Accelerating Long-Latency Load Requests via Perceptron-Based Off-Chip Load Prediction

Rahul Bera1Konstantinos Kanellopoulos1Shankar Balachandran2David Novo3Ataberk Olgun1Mohammad Sadrosadati1Onur Mutlu1

¹ETH Zürich ²Intel Processor Architecture Research Lab ³LIRMM, Univ. Montpellier, CNRS

https://arxiv.org/pdf/2209.00188.pdf

Hermes Talk Video



Computer Architecture - Lecture 18: Cutting-Edge Research in Computer Architecture (Fall 2022)



2.4K views Streamed 5 months ago Livestream - Computer Architecture - ETH Zürich (Fall 2022) Computer Architecture, ETH Zürich, Fall 2022 (https://safari.ethz.ch/architecture/f...)

SAFARI

https://www.youtube.com/watch?v=PWWBtrL60dQ&t=3609s







Accelerating Long-Latency Load Requests via Perceptron-Based Off-Chip Load Prediction

Rahul Bera, Konstantinos Kanellopoulos, Shankar Balachandran, David Novo, Ataberk Olgun, Mohammad Sadrosadati, Onur Mutlu

https://github.com/CMU-SAFARI/Hermes





Problem

Long-latency off-chip load requests

Often **stall** processor by **blocking instruction retirement** from Reorder Buffer (ROB)





Traditional Solutions





Increase size of on-chip caches

Key Observation 1





off-chip loads without any prefetcher

Key Observation 2

On-chip cache access latency significantly contributes to off-chip load latency



40% of the stalls can be eliminated by removing on-chip cache access latency from critical path

Caches are Getting Bigger and Slower...



Our Goal

Improve processor performance by **removing on-chip cache access latency** from the **critical path of off-chip loads**



Predicts which load requests are likely to go off-chip

Starts **fetching** data **directly** from **main memory** while concurrently accessing the cache hierarchy

Hermes: Key Contribution



By **learning** from multiple program context information

Hermes Overview





Designing the Off-Chip Load Predictor

History-based prediction

HMP [Yoaz+, ISCA'99] for the L1-D cache Using branch-predictor-like hybrid predictor: Global Gshare and GSkew



POPET provides both higher accuracy and higher performance than predictors inspired from these previous works

- Metadata size increases with cache hierarchy size
- X May need to track **all** cache operations
 - Gets complex depending on the cache hierarchy configuration (e.g., inclusivity, bypassing,...)

Learning from program behavior

Correlate different program features with off-chip loads



.ow storage overhead 🛛 📀



Low design complexity



POPET: Perceptron-Based Off-Chip Predictor

- Multi-feature hashed perceptron model^[1]
 - Each feature has its own weight table
 - Stores correlation between feature value and off-chip prediction





69

Predicting using POPET



Training POPET



Evaluation
Simulation Methodology

- ChampSim trace driven simulator
- **110 single-core** memory-intensive traces
 - SPEC CPU 2006 and 2017
 - PARSEC 2.1
 - Ligra
 - Real-world applications

• **220 eight-core** memory-intensive trace mixes

LLC Prefetchers

- Pythia [Bera+, MICRO'21]
- Bingo [Bakshalipour+, HPCA'19]
- MLOP [Shakerinava+, 3rd Prefetching Championship'19]
- SPP + Perceptron filter [Bhatia+, ISCA'20]
- SMS [Somogyi+, ISCA'06]

Off-Chip Predictors

- History-based: HMP [Yoaz+, ISCA'99]
- Tracking-based: Address Tag-Tracking based Predictor (TTP)
- Ideal Off-chip Predictor

Single-Core Performance Improvement



Hermes provides nearly 90% of performance benefit of Ideal Hermes that has an ideal off-chip load predictor

Increase in Main Memory Requests

Hermes Pythia Pythia + Hermes Pythia + Ideal Hermes



Hermes is more **bandwidth-efficient** than even an efficient prefetcher like Pythia



Performance with Varying Memory Bandwidth



Hermes+Pythia outperforms Pythia across all bandwidth configurations

Performance with Varying Baseline Prefetcher



Overhead of Hermes



*On top of an Intel Alder Lake-like performance-core ^[2] configuration

A Lot More in the Hermes Paper

Performance sensitivity to:



Hermes: Accelerating Long-Latency Load Requests via Perceptron-Based Off-Chip Load Prediction

Rahul Bera¹ Konstantinos Kanellopoulos¹ Shankar Balachandran² David Novo³ Ataberk Olgun¹ Mohammad Sadrosadati¹ Onur Mutlu¹

¹ETH Zürich ²Intel Processor Architecture Research Lab ³LIRMM, Univ. Montpellier, CNRS

Long-latency load requests continue to limit the performance of modern high-performance processors. To increase the latency tolerance of a processor, architects have primarily relied on two key techniques: sophisticated data prefetchers and large on-chip caches. In this work, we show that: (1) even a sophisticated stateof-the-art prefetcher can only predict half of the off-chip load requests on average across a wide range of workloads, and (2) due to the increasing size and complexity of on-chip caches, a large fraction of the latency of an off-chip load request is spent accessing the on-chip cache hierarchy to solely determine that it needs to go off-chip.

The goal of this work is to accelerate off-chip load requests by removing the on-chip cache access latency from their critical path. To this end, we propose a new technique called Hermes, whose key idea is to: (1) accurately predict which load requests off-chip main memory (i.e., an *off-chip load*) often stalls the processor core by blocking the instruction retirement from the reorder buffer (ROB), thus limiting the core's performance [88, 91, 92]. To increase the latency tolerance of a core, computer architects primarily rely on two key techniques. First, they employ increasingly sophisticated hardware prefetchers that can learn complex memory address patterns and fetch data required by future load requests before the core demands them [28, 32, 33, 35, 75]. Second, they significantly scale up the size of the on-chip cache hierarchy with each new generation of processors [10, 11, 16].

Key problem. Despite recent advances in processor core design, we observe two key trends in new processor designs that leave a significant opportunity for performance improvement on the table. First, even a sophisticated state-of-the-art

https://arxiv.org/pdf/2209.00188.pdf

A New Approach to Latency Reduction

Hermes advocates for **off-chip load prediction**, a **different** form of speculation than **load address prediction** employed by prefetchers

Off-chip load prediction can be applied **by itself** or **combined with load address prediction** to provide performance improvement



Hermes: Summary

Hermes employs the first perceptron-based off-chip load predictor



Hermes is Open Source





All workload traces

9eaeca0 5 days ago 🛽 🔁 21 co

long-latency off-chip load requests by



V1.1

SAFARI

13 prefetchers

- Stride [Fu+, MICRO'92]
- Streamer [Chen and Baer, IEEE TC'95]
- SMS [Somogyi+, ISCA'06]
- AMPM [Ishii+, ICS'09]
- Sandbox [Pugsley+, HPCA'14]
- BOP [Michaud, HPCA'16]
- SPP [Kim+, MICRO'16]
- Bingo [Bakshalipour+, HPCA'19]
- SPP+PPF [Bhatia+, ISCA'19]
- DSPatch [Bera+, MICRO'19]
- MLOP [Shakerinava+, DPC-3'19]
- IPCP [Pakalapati+, ISCA'20]
- Pythia [Bera+, MICRO'21]

9 off-chip predictors

Predictor type Description Base Always NO Simple confidence counter-based threshold Basic Random Random Hit-miss predictor with a given positive probability HMP-Local Hit-miss predictor [Yoaz+, ISCA'99] with local prediction HMP-GShare Hit-miss predictor with GShare prediction HMP-GSkew Hit-miss predictor with GSkew prediction **HMP-Ensemble** Hit-miss predictor with all three types combined TTP Tag-tracking based predictor Perceptron-based OCP used in this paper Perc

https://github.com/CMU-SAFARI/Hermes

Easy To Define Your Own Off-Chip Predictor

• Just extend the OffchipPredBase class

```
class OffchipPredBase
 8
    {
 9
    public:
10
        uint32_t cpu;
11
12
        string type;
13
        uint64_t seed;
        uint8 t dram bw; // current DRAM bandwidth bucket
14
15
        OffchipPredBase(uint32_t _cpu, string _type, uint64_t _seed) : cpu(_cpu), type(_type), seed(_seed)
16
         {
17
             srand(seed);
18
            dram_bw = 0;
19
        }
20
        ~OffchipPredBase() {}
21
        void update dram bw(uint8 t dram bw) { dram bw = dram bw; }
22
23
        virtual void print_config();
24
        virtual void dump_stats();
25
        virtual void reset_stats();
26
        virtual void train(ooo_model_instr *arch_instr, uint32_t data_index, LSQ_ENTRY *lq_entry);
27
        virtual bool predict(oco_model_instr *arch_instr, uint32_t data_index, LSQ_ENTRY *lq_entry);
28
    };
29
30
    #endif /* OFFCHIP_PRED_BASE_H */
31
32
```

Easy To Define Your Own Off-Chip Predictor

Define your own train() and predict() functions

```
void OffchipPredBase::train(ooo_model_instr *arch_instr, uint32_t data_index, LSQ_ENTRY *lq_entry)
19
     {
20
        // nothing to train
21
22
    }
23
24
    bool OffchipPredBase::predict(ooo_model_instr *arch_instr, uint32_t data_index, LSQ_ENTRY *lq_entry)
25
    {
        // predict randomly
26
        // return (rand() % 2) ? true : false;
27
        return false;
28
29
   }
```

 Get statistics like accuracy (stat name precision) and coverage (stat name recall) out of the box

> Core_0_offchip_pred_true_pos 2358716 Core_0_offchip_pred_false_pos 276883 Core_0_offchip_pred_false_neg 132145 Core_0_offchip_pred_precision 89.49 Core_0_offchip_pred_recall 94.69

Off-Chip Prediction Can Further Enable...

Prioritizing loads that are likely go off-chip in cache queues and on-chip network routing

Better instruction scheduling of data-dependent instructions

Other ideas to improve **performance** and **fairness** in multi-core system design...

Learning-Based Off-Chip Load Predictors

 Rahul Bera, Konstantinos Kanellopoulos, Shankar Balachandran, David Novo, Ataberk Olgun, Mohammad Sadrosadati, and Onur Mutlu,
 "Hermes: Accelerating Long-Latency Load Requests via Perceptron-Based Off-Chip Load Prediction"
 Proceedings of the <u>55th International Symposium on Microarchitecture</u> (MICRO), Chicago, IL, USA, October 2022.
 [Slides (pptx) (pdf)]
 [Longer Lecture Slides (pptx) (pdf)]
 [Talk Video (12 minutes)]
 [Lecture Video (25 minutes)]
 [arXiv version]
 [Source Code (Officially Artifact Evaluated with All Badges)]
 Officially artifact evaluated as available, reusable and reproducible. Best paper award at MICRO 2022.



Hermes: Accelerating Long-Latency Load Requests via Perceptron-Based Off-Chip Load Prediction

Rahul Bera1Konstantinos Kanellopoulos1Shankar Balachandran2David Novo3Ataberk Olgun1Mohammad Sadrosadati1Onur Mutlu1

¹ETH Zürich ²Intel Processor Architecture Research Lab ³LIRMM, Univ. Montpellier, CNRS

https://arxiv.org/pdf/2209.00188.pdf

Hermes Talk Video



Computer Architecture - Lecture 18: Cutting-Edge Research in Computer Architecture (Fall 2022)



2.4K views Streamed 5 months ago Livestream - Computer Architecture - ETH Zürich (Fall 2022) Computer Architecture, ETH Zürich, Fall 2022 (https://safari.ethz.ch/architecture/f...)

SAFARI

https://www.youtube.com/watch?v=PWWBtrL60dQ&t=3609s







Accelerating Long-Latency Load Requests via Perceptron-Based Off-Chip Load Prediction

Rahul Bera, Konstantinos Kanellopoulos, Shankar Balachandran, David Novo, Ataberk Olgun, Mohammad Sadrosadati, Onur Mutlu

https://github.com/CMU-SAFARI/Hermes





Reinforcement Learning Based DRAM Controllers

DRAM Controller: Functions

- Ensure correct operation of DRAM (refresh and timing)
- Service DRAM requests while obeying timing constraints of DRAM chips
 - Constraints: resource conflicts (bank, bus, channel), minimum write-to-read delays
 - Translate requests to DRAM command sequences
- Buffer and schedule requests for high performance + QoS
 Reordering, row-buffer, bank, rank, bus management
- Manage power consumption and thermals in DRAM
 - Turn on/off DRAM chips, manage power modes

Why Are DRAM Controllers Difficult to Design?

- Need to obey DRAM timing constraints for correctness
 - □ There are many (50+) timing constraints in DRAM
 - tWTR: Minimum number of cycles to wait before issuing a read command after a write command is issued
 - tRC: Minimum number of cycles between the issuing of two consecutive activate commands to the same bank

••••

- Need to keep track of many resources to prevent conflicts
 - Channels, banks, ranks, data bus, address bus, row buffers
- Need to handle DRAM refresh
- Need to manage power consumption
- Need to optimize performance & QoS (in the presence of constraints)
 - Reordering is not simple
 - Fairness and QoS needs complicates the scheduling problem

Many DRAM Timing Constraints

Latency	Symbol	DRAM cycles	Latency	Symbol	DRAM cycles
Precharge	^{t}RP	11	Activate to read/write	^{t}RCD	11
Read column address strobe	CL	11	Write column address strobe	CWL	8
Additive	AL	0	Activate to activate	^{t}RC	39
Activate to precharge	^t RAS	28	Read to precharge	^{t}RTP	6
Burst length	^{t}BL	4	Column address strobe to column address strobe	^{t}CCD	4
Activate to activate (different bank)	^{t}RRD	6	Four activate windows	^{t}FAW	24
Write to read	^{t}WTR	6	Write recovery	^{t}WR	12

Table 4. DDR3 1600 DRAM timing specifications

 From Lee et al., "DRAM-Aware Last-Level Cache Writeback: Reducing Write-Caused Interference in Memory Systems," HPS Technical Report, April 2010.

More on DRAM Operation

- Kim et al., "A Case for Exploiting Subarray-Level Parallelism (SALP) in DRAM," ISCA 2012.
- Lee et al., "Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture," HPCA 2013.



Table 2. Timing Constraints (DDR3-1066) [43]

Phase	Commands	Name	Value	
1	$\begin{array}{c} \text{ACT} \rightarrow \text{READ} \\ \text{ACT} \rightarrow \text{WRITE} \end{array}$	tRCD	15ns	
A	$\mathrm{ACT} \to \mathrm{PRE}$	tRAS	37.5ns	
2	$\begin{array}{l} \text{READ} \rightarrow \textit{data} \\ \text{WRITE} \rightarrow \textit{data} \end{array}$	tCL tCWL	15ns 11.25ns	
	data burst	tBL	7.5ns	
3	$\text{PRE} \rightarrow \text{ACT}$	tRP	15ns	
1&3	$ACT \rightarrow ACT$	tRC (tRAS+tRP)	52.5ns	

DRAM Scheduling Policies (I)

- FCFS (first come first served)
 - Oldest request first
- FR-FCFS (first ready, first come first served)
 - 1. Row-hit first
 - 2. Oldest first
 - Goal: Maximize row buffer hit rate \rightarrow maximize DRAM throughput

DRAM Scheduling Policies (II)

- A scheduling policy is a request prioritization order
- Prioritization can be based on
 - Request age
 - Row buffer hit/miss status
 - Request type (prefetch, read, write)
 - Requestor type (load miss or store miss)
 - Request criticality
 - Oldest miss in the core?
 - How many instructions in core are dependent on it?
 - Will it stall the processor?
 - Interference caused to other cores

••••

Memory Performance Attacks [USENIX SEC'07]

 Thomas Moscibroda and Onur Mutlu, <u>"Memory Performance Attacks: Denial of Memory Service</u> <u>in Multi-Core Systems"</u> *Proceedings of the <u>16th USENIX Security Symposium</u> (USENIX SECURITY), pages 257-274, Boston, MA, August 2007. <u>Slides</u> (ppt)*

Memory Performance Attacks: Denial of Memory Service in Multi-Core Systems

Thomas Moscibroda Onur Mutlu Microsoft Research {moscitho,onur}@microsoft.com

STFM [MICRO'07]

 Onur Mutlu and Thomas Moscibroda,
 "Stall-Time Fair Memory Access Scheduling for Chip Multiprocessors"
 Proceedings of the <u>40th International Symposium on</u> <u>Microarchitecture</u> (MICRO), pages 146-158, Chicago, IL, December 2007. [Summary] [Slides (ppt)]

Stall-Time Fair Memory Access Scheduling for Chip Multiprocessors

Onur Mutlu Thomas Moscibroda

Microsoft Research {onur,moscitho}@microsoft.com

PAR-BS [ISCA'08]

 Onur Mutlu and Thomas Moscibroda, "Parallelism-Aware Batch Scheduling: Enhancing both Performance and Fairness of Shared DRAM Systems" Proceedings of the <u>35th International Symposium on Computer</u> <u>Architecture</u> (ISCA), pages 63-74, Beijing, China, June 2008. [Summary] [Slides (ppt)]

Parallelism-Aware Batch Scheduling: Enhancing both Performance and Fairness of Shared DRAM Systems

Onur Mutlu Thomas Moscibroda Microsoft Research {onur,moscitho}@microsoft.com Variants implemented in Samsung SoC memory controllers

Effective platform level approach and DRAM accesses are crucial to system performance. This paper touches this topics and suggest a superior approach to current known techniques. **Review from ISCA 2008**

ATLAS Memory Scheduler [HPCA'10]

 Yoongu Kim, Dongsu Han, Onur Mutlu, and Mor Harchol-Balter, <u>"ATLAS: A Scalable and High-Performance Scheduling</u> <u>Algorithm for Multiple Memory Controllers"</u> *Proceedings of the <u>16th International Symposium on High-</u> <u>Performance Computer Architecture</u> (HPCA), Bangalore, India, January 2010. <u>Slides (pptx)</u>*

ATLAS: A Scalable and High-Performance Scheduling Algorithm for Multiple Memory Controllers

Yoongu Kim Dongsu Han Onur Mutlu Mor Harchol-Balter Carnegie Mellon University

Thread Cluster Memory Scheduling [MICRO'10]

 Yoongu Kim, Michael Papamichael, Onur Mutlu, and Mor Harchol-Balter,
 "Thread Cluster Memory Scheduling: Exploiting Differences in Memory Access Behavior" Proceedings of the <u>43rd International Symposium on</u>

Microarchitecture (MICRO), pages 65-76, Atlanta, GA, December 2010. <u>Slides (pptx) (pdf)</u>

Thread Cluster Memory Scheduling: Exploiting Differences in Memory Access Behavior

Yoongu Kim yoonguk@ece.cmu.edu

Michael Papamichael papamix@cs.cmu.edu

Onur Mutlu onur@cmu.edu

Mor Harchol-Balter harchol@cs.cmu.edu

Carnegie Mellon University

BLISS [ICCD'14, TPDS'16]

 Lavanya Subramanian, Donghyuk Lee, Vivek Seshadri, Harsha Rastogi, and Onur Mutlu,
 "The Blacklisting Memory Scheduler: Achieving High Performance and Fairness at Low Cost"
 Proceedings of the 32nd IEEE International Conference on Computer Design (ICCD), Seoul, South Korea, October 2014.
 [Slides (pptx) (pdf)]

The Blacklisting Memory Scheduler: Achieving High Performance and Fairness at Low Cost

Lavanya Subramanian, Donghyuk Lee, Vivek Seshadri, Harsha Rastogi, Onur Mutlu Carnegie Mellon University {lsubrama,donghyu1,visesh,harshar,onur}@cmu.edu

Staged Memory Scheduling: CPU-GPU [ISCA'12]

 Rachata Ausavarungnirun, Kevin Chang, Lavanya Subramanian, Gabriel Loh, and Onur Mutlu,
 "Staged Memory Scheduling: Achieving High Performance and Scalability in Heterogeneous Systems" Proceedings of the <u>39th International Symposium on Computer</u> <u>Architecture</u> (ISCA), Portland, OR, June 2012. <u>Slides (pptx)</u>

Staged Memory Scheduling: Achieving High Performance and Scalability in Heterogeneous Systems

Rachata Ausavarungnirun[†] Kevin Kai-Wei Chang[†] Lavanya Subramanian[†] Gabriel H. Loh[‡] Onur Mutlu[†]

[†]Carnegie Mellon University {rachata,kevincha,lsubrama,onur}@cmu.edu

SAFARI

[‡]Advanced Micro Devices, Inc. gabe.loh@amd.com

DASH: Heterogeneous Systems [TACO'16]

- Hiroyuki Usui, Lavanya Subramanian, Kevin Kai-Wei Chang, and Onur Mutlu,
 - **"DASH: Deadline-Aware High-Performance Memory Scheduler for Heterogeneous Systems with Hardware** <u>Accelerators</u>"
 - <u>ACM Transactions on Architecture and Code Optimization</u> (TACO),
 - Vol. 12, January 2016. Presented at the <u>11th HiPEAC Conference</u>, Prague, Czech Republic, January 2016. [<u>Slides (pptx) (pdf)</u>] [<u>Source Code</u>]

DASH: Deadline-Aware High-Performance Memory Scheduler for Heterogeneous Systems with Hardware Accelerators

HIROYUKI USUI, LAVANYA SUBRAMANIAN, KEVIN KAI-WEI CHANG, and ONUR MUTLU, Carnegie Mellon University

MISE: Predictable Performance [HPCA'13]

 Lavanya Subramanian, Vivek Seshadri, Yoongu Kim, Ben Jaiyen, and Onur Mutlu,
 "MISE: Providing Performance Predictability and Improving Fairness in Shared Main Memory Systems"
 Proceedings of the 19th International Symposium on High-Performance Computer Architecture (HPCA), Shenzhen, China, February 2013. Slides (pptx)

MISE: Providing Performance Predictability and Improving Fairness in Shared Main Memory Systems

Lavanya Subramanian Vivek Seshadri

ivek Seshadri Yoongu Kim Ben Jaiyen Onur Mutlu Carnegie Mellon University

ASM: Predictable Performance [MICRO'15]

- Lavanya Subramanian, Vivek Seshadri, Arnab Ghosh, Samira Khan, and Onur Mutlu,
 - "The Application Slowdown Model: Quantifying and Controlling the Impact of Inter-Application Interference at Shared Caches and Main Memory"

Proceedings of the <u>48th International Symposium on Microarchitecture</u> (**MICRO**), Waikiki, Hawaii, USA, December 2015. [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Poster (pptx) (pdf)] [Source Code]

The Application Slowdown Model: Quantifying and Controlling the Impact of Inter-Application Interference at Shared Caches and Main Memory

Lavanya Subramanian^{*}§ Vivek Seshadri^{*} Arnab Ghosh^{*†} Samira Khan^{*‡} Onur Mutlu^{*}

*Carnegie Mellon University §Intel Labs [†]IIT Kanpur [‡]University of Virginia



Memory Controllers are critical to research

They will become even more important

Memory Control is Getting More Complex



- Heterogeneous agents: CPUs, GPUs, and HWAs
- Main memory interference between CPUs, GPUs, HWAs

Many goals, many constraints, many metrics ...
Reality and Dream

- Reality: It is difficult to design a policy that maximizes performance, QoS, energy-efficiency, ...
 - Too many things to think about
 - Continuously changing workload and system behavior



Dream: Wouldn't it be nice if the DRAM controller automatically found a good scheduling policy on its own?

- Problem: DRAM controllers are difficult to design
 - It is difficult for human designers to design a policy that can adapt itself very well to different workloads and different system conditions
- Idea: A memory controller that adapts its scheduling policy to workload behavior and system conditions using machine learning.
- Observation: Reinforcement learning maps nicely to memory control.
- Design: Memory controller is a reinforcement learning agent
 - It dynamically and continuously learns and employs the best scheduling policy to maximize long-term performance.



Figure 2: (a) Intelligent agent based on reinforcement learning principles;

- Dynamically adapt the memory scheduling policy via interaction with the system at runtime
 - Associate system states and actions (commands) with long term reward values: each action at a given state leads to a learned reward
 - Schedule command with highest estimated long-term reward value in each state
 - Continuously update reward values for <state, action> pairs based on feedback from system



 Engin Ipek, Onur Mutlu, José F. Martínez, and Rich Caruana, <u>"Self Optimizing Memory Controllers: A Reinforcement Learning</u> <u>Approach</u>"

Proceedings of the <u>35th International Symposium on Computer Architecture</u> (<i>ISCA), pages 39-50, Beijing, China, June 2008.



Figure 4: High-level overview of an RL-based scheduler.

States, Actions, Rewards

- Reward function
 - +1 for scheduling Read and Write commands
 - 0 at all other times
 - Goal is to maximize long-term data bus utilization

- State attributes
 - Number of reads, writes, and load misses in transaction queue
 - Number of pending writes and ROB heads waiting for referenced row
 - Request's relative ROB order

- Actions
 - Activate
 - Write
 - Read load miss
 - Read store miss
 - Precharge pending
 - Precharge preemptive
 - NOP

Performance Results



Figure 7: Performance comparison of in-order, FR-FCFS, RL-based, and optimistic memory controllers

Large, robust performance improvements over many human-designed policies



Figure 15: Performance comparison of FR-FCFS and RL-based memory controllers on systems with 6.4GB/s and 12.8GB/s peak DRAM bandwidth



+ Continuous learning in the presence of changing environment

+ Reduced designer burden in finding a good scheduling policy. Designer specifies:

1) What system variables might be useful

2) What target to optimize, but not how to optimize it

- -- How to specify different objectives? (e.g., fairness, QoS, ...)
- -- Hardware complexity?
- -- Design **mindset** and flow

More on Self-Optimizing DRAM Controllers (I)

 Engin Ipek, Onur Mutlu, José F. Martínez, and Rich Caruana, <u>"Self Optimizing Memory Controllers: A Reinforcement Learning</u> <u>Approach</u>" *Proceedings of the <u>35th International Symposium on Computer Architecture</u> (ISCA), pages 39-50, Beijing, China, June 2008.*

Self-Optimizing Memory Controllers: A Reinforcement Learning Approach

Engin İpek^{1,2} Onur Mutlu² José F. Martínez¹ Rich Caruana¹

¹Cornell University, Ithaca, NY 14850 USA

² Microsoft Research, Redmond, WA 98052 USA

More on Self-Optimizing DRAM Controllers (II)

 Janani Mukundan and José F. Martinez
<u>MORSE: Multi-Objective Reconfigurable Self-Optimizing Memory Scheduler</u>" Proceedings of the <u>18th International Symposium on High Performance</u> <u>Computer Architecture</u> (HPCA), New Orleans, Louisiana, February 2012.

MORSE: Multi-objective Reconfigurable Self-optimizing Memory Scheduler

Janani Mukundan José F. Martínez

Computer Systems Laboratory Cornell University Ithaca, NY, 14850 USA

http://m3.csl.cornell.edu/



Memory Controllers are critical to research

They will become even more important

Sibyl: Reinforcement Learning based Data Placement in Hybrid SSDs

Self-Optimizing Hybrid SSD Controllers

Gagandeep Singh, Rakesh Nadig, Jisung Park, Rahul Bera, Nastaran Hajinazar, David Novo, Juan Gomez-Luna, Sander Stuijk, Henk Corporaal, and Onur Mutlu, "Sibyl: Adaptive and Extensible Data Placement in Hybrid Storage Systems Using Online Reinforcement Learning" Proceedings of the <u>49th International Symposium on Computer</u> <u>Architecture (ISCA)</u>, New York, June 2022. [Slides (pptx) (pdf)] [arXiv version] [Sibyl Source Code] [Talk Video (16 minutes)]

Sibyl: Adaptive and Extensible Data Placement in Hybrid Storage Systems Using Online Reinforcement Learning

Gagandeep Singh1Rakesh Nadig1Jisung Park1Rahul Bera1Nastaran Hajinazar1David Novo3Juan Gómez-Luna1Sander Stuijk2Henk Corporaal2Onur Mutlu11ETH Zürich2Eindhoven University of Technology3LIRMM, Univ. Montpellier, CNRS

https://arxiv.org/pdf/2205.07394.pdf



Adaptive and Extensible Data Placement in Hybrid Storage Systems Using Online Reinforcement Learning

Gagandeep Singh, Rakesh Nadig, Jisung Park, Rahul Bera, Nastaran Hajinazar, David Novo, Juan Gómez Luna, Sander Stuijk, Henk Corporaal, Onur Mutlu







TU/

8

Executive Summary

- **Background**: A hybrid storage system (HSS) uses multiple different storage devices to provide high and scalable storage capacity at high performance
- **Problem**: Two key shortcomings of prior data placement policies:
 - Lack of adaptivity to:
 - Workload changes
 - Changes in device types and configurations
 - Lack of extensibility to more devices
- **Goal**: Design a data placement technique that provides:
 - Adaptivity, by continuously learning and adapting to the application and underlying device characteristics
 - Easy extensibility to incorporate a wide range of hybrid storage configurations
- Contribution: Sibyl, the first reinforcement learning-based data placement technique in hybrid storage systems that:
 - Provides adaptivity to changing workload demands and underlying device characteristics
 - Can easily extend to any number of storage devices
 - Provides ease of design and implementation that requires only a small computation overhead
- Key Results: Evaluate on real systems using a wide range of workloads
 - Sibyl improves performance by 21.6% compared to the best previous data placement technique in dual-HSS configuration
 - In a tri-HSS configuration, Sibyl outperforms the state-of-the-art-policy policy by 48.2%
 - Sibyl achieves 80% of the performance of an oracle policy with storage overhead of only 124.4 KiB



https://github.com/CMU-SAFARI/Sibyl

Hybrid Storage System Basics

SAFARI

Address Space (Application/File System View)



124

Hybrid Storage System Basics

Logical Address Space (Application/File System View)

ogical Pages.

Performance of a hybrid storage system highly depends on the storage management layer's ability to manage diverse devices and workloads





Key Shortcomings in Prior Techniques

We observe **two key shortcomings** that significantly limit the performance benefits of prior techniques

1. Lack of **adaptivity to**:

- a) Workload changes
- b) Changes in device types and configuration

2. Lack of **extensibility** to more devices



Lack of Extensibility (1/2)

Rigid techniques that require significant effort to accommodate more than two devices

Change in storage configuration







Lack of Extensibility (2/2)

Rigid techniques that require significant effort to accommodate more than two devices

Change in storage configuration



Design a new policy







Our Goal

A data-placement mechanism that can provide:

1.Adaptivity, by continuously learning and adapting to the application and underlying device characteristics

2.Easy extensibility to incorporate a wide range of hybrid storage configurations



Our Proposal



Sibyl Formulates data placement in hybrid storage systems as a **reinforcement learning problem**



Sibyl is an oracle that makes accurate prophecies https://en.wikipedia.org/wiki/Sibyl

Basics of Reinforcement Learning (RL)



Environment

Agent learns to take an **action** in a given **state** to maximize a numerical **reward**



Formulating Data Placement as RL



What is State?

- Limited number of state features:
 - Reduce the implementation overhead
 - RL agent is more sensitive to reward



 $O_t = (size_t, type_t, intr_t, cnt_t, cap_t, curr_t)$

• We **quantize the state representation** into bins to reduce storage overhead





What is Reward?

• Defines the **objective** of Sibyl



- We formulate the reward as a function of the request latency
- Encapsulates three key aspects:
 - Internal state of the device (e.g., read/write latencies, the latency of garbage collection, queuing delays, ...)
 - Throughput
 - Evictions
- More details in the paper

What is Action?

• At every new page request, the action is to select a storage device



- Action can be easily extended to any number of storage devices
- Sibyl evicts a page when the fast device utilization is 100%
- Sibyl promotes a page when there is an update from the application

Talk Outline

Key Shortcomings of Prior Data Placement Techniques

Formulating Data Placement as Reinforcement Learning

Sibyl: Overview

Evaluation of Sibyl and Key Results

Conclusion



Sibyl Execution



Sibyl Design: Overview













RL Training Thread


Periodic Weight Transfer



Evaluation Methodology (1/3)

Real system with various HSS configurations

- Dual-hybrid and tri-hybrid systems



Evaluation Methodology (2/3)

Cost-Oriented HSS Configuration



High-end SSD

SAFARI

Low-end HDD

Performance-Oriented HSS Configuration



Evaluation Methodology (3/3)

• 18 different workloads from:

- MSR Cambridge and Filebench Suites

• Four state-of-the-art data placement baselines:





Cost-Oriented HSS Configuration





Cost-Oriented HSS Configuration



Sibyl consistently outperforms all the baselines for all the workloads



Performance-Oriented HSS Configuration





Performance-Oriented HSS Configuration



Sibyl provides 21.6% performance improvement by dynamically adapting its data placement policy



Performance-Oriented HSS Configuration



Sibyl achieves 80% of the performance of an oracle policy that has complete knowledge of future access patterns

Performance on Tri-HSS



Extending Sibyl for more devices:

- 1. Add a new action
- 2. Add the remaining capacity of the new device as a state feature



Performance on Tri-HSS



Extending Sibyl for more devices:

- 1. Add a new action
- 2. Add the remaining capacity of the new device as a state feature



Performance on Tri-HSS



Extending Sibyl for more devices: 1. Add a new action

Sibyl **outperforms** the state-of-the-art data placement policy by 48.2% in a real tri-hybrid system Sibyl reduces the system architect's burden by providing ease of extensibility

Sibyl's Overhead

• 124.4 KiB of total storage cost

- Experience buffer, inference and training network
- 40-bit metadata overhead per page for state features
- Inference latency of ~10ns
- Training latency of ~2us

Small inference overhead
Satisfies prediction latency

More in the Paper (1/3)

Throughput (IOPS) evaluation

 Sibyl provides high IOPS compared to baseline policies because it indirectly captures throughput (size/latency)

- Evaluation on unseen workloads
 - Sibyl can effectively adapt its policy to highly dynamic workloads

- Evaluation on **mixed workloads**
 - Sibyl provides equally-high performance benefits as in single workloads



More in the Paper (2/3)

- Evaluation on different features
 - Sibyl autonomously decides which features are important to maximize the performance
- Evaluation with different hyperparameter values

- Sensitivity to fast storage capacity
 - Sibyl provides scalability by dynamically adapting its policy to available storage size
- Explainability analysis of Sibyl's decision making
 - Explain Sibyl's actions for different workload characteristics and device configurations

More in the Paper (3/3)

Sibyl: Adaptive and Extensible Data Placement in Hybrid Storage Systems Using Online Reinforcement Learning

Gagandeep Singh1Rakesh Nadig1Jisung Park1Rahul Bera1Nastaran Hajinazar1David Novo3Juan Gómez-Luna1Sander Stuijk2Henk Corporaal2Onur Mutlu11ETH Zürich2Eindhoven University of Technology3LIRMM, Univ. Montpellier, CNRS

https://arxiv.org/pdf/2205.07394.pdf

https://github.com/CMU-SAFARI/Sibyl



Conclusion

- We introduced Sibyl, the first reinforcement learningbased data placement technique in hybrid storage systems that provides
 - Adaptivity
 - Easily extensibility
 - Ease of design and implementation

• We evaluated Sibyl on real systems using many different workloads

- Sibyl **improves performance by 21.6%** compared to the best prior data placement policy in a dual-HSS configuration
- In a tri-HSS configuration, Sibyl outperforms the state-of-the-artdata placement policy by 48.2%
- Sibyl achieves **80% of the performance** of an oracle policy with a storage overhead of only **124.4 KiB**

SAFARI

https://github.com/CMU-SAFARI/Sibyl

Major Directions

- Consider other optimization objectives
 - Energy consumption, endurance of storage devices.....
 - Design better reward structures

• Optimize data migration in hybrid storage systems

- Explore machine learning (ML) techniques to make data migration adaptive and extensible
- How do we coordinate multiple ML techniques?
- How do we improve these policies in other heterogeneous memory systems?
 - DRAM + NVM, CPU Caches + DRAM
 - Design RL models keeping latency constraints in mind

ISCA 2022 Paper, Slides, Videos

 Gagandeep Singh, Rakesh Nadig, Jisung Park, Rahul Bera, Nastaran Hajinazar, David Novo, Juan Gomez-Luna, Sander Stuijk, Henk Corporaal, and Onur Mutlu, "Sibyl: Adaptive and Extensible Data Placement in Hybrid Storage Systems Using Online Reinforcement Learning" Proceedings of the <u>49th International Symposium on Computer</u> <u>Architecture</u> (ISCA), New York, June 2022.
 [Slides (pptx) (pdf)] [arXiv version]
 [Sibyl Source Code] [Talk Video (16 minutes)]

Sibyl: Adaptive and Extensible Data Placement in Hybrid Storage Systems Using Online Reinforcement Learning

Gagandeep Singh1Rakesh Nadig1Jisung Park1Rahul Bera1Nastaran Hajinazar1David Novo3Juan Gómez-Luna1Sander Stuijk2Henk Corporaal2Onur Mutlu11ETH Zürich2Eindhoven University of Technology3LIRMM, Univ. Montpellier, CNRS

https://arxiv.org/pdf/2205.07394.pdf

SSD Course (Spring 2023)

Spring 2023 Edition:

https://safari.ethz.ch/projects_and_seminars/spring2023/ doku.php?id=modern_ssds

Fall 2022 Edition:

https://safari.ethz.ch/projects_and_seminars/fall2022/do ku.php?id=modern_ssds

Youtube Livestream (Spring 2023):

https://www.youtube.com/watch?v=4VTwOMmsnJY&list =PL5Q2soXY2Zi_8qOM5Icpp8hB2SHtm4z57&pp=iAQB

Youtube Livestream (Fall 2022):

- https://www.youtube.com/watch?v=hqLrd-Uj0aU&list=PL5Q2soXY2Zi9BJhenUq4JI5bwhAMpAp13&p p=iAQB
- Project course
 - Taken by Bachelor's/Master's students
 - SSD Basics and Advanced Topics
 - Hands-on research exploration
 - Many research readings

https://www.youtube.com/onurmutlulectures



Fall 2022 Meetings/Schedule

Week	Date	Livestream	Meeting	Learning Materials	Assignments
W1	06.10		M1: P&S Course Presentation	Required Recommended	
W2	12.10	Yeu Ture Live	M2: Basics of NAND Flash- Based SSDs	Required Recommended	
W3	19.10	You Ture Live	M3: NAND Flash Read/Write Operations	Required Recommended	
W4	26.10	Yeu Ture Live	M4: Processing inside NAND Flash	Required Recommended	
W5	02.11	Yeu Tube Live	M5: Advanced NAND Flash Commands & Mapping m PDF m PPT	Required Recommended	
W6	09.11	You Time Live	M6: Processing inside Storage	Required Recommended	
W7	23.11	You Ture Live	M7: Address Mapping & Garbage Collection	Required Recommended	
W8	30.11	You Tune Live	M8: Introduction to MQSim	Required Recommended	
W9	14.12	You Live	M9: Fine-Grained Mapping and Multi-Plane Operation-Aware Block Management amPDF amPPT	Required Recommended	
W10	04.01.2023	Yw me Premiere	M10a: NAND Flash Basics	Required Recommended	
			M10b: Reducing Solid-State Drive Read Latency by Optimizing Read-Retry	Required Recommended	
			M10c: Evanesco: Architectural Support for Efficient Data Sanitization in Modern Flash- Based Storage Systems an PDF im PPT im Paper	Required Recommended	
			M10d: DeepSketch: A New Machine Learning-Based Reference Search Technique for Post-Deduplication Delta Compression mPDF m PPT mPaper	Required Recommended	
W11	11.01	You Dive	M11: FLIN: Enabling Fairness and Enhancing Performance in Modern NVMe Solid State Drives and PDF im PPT	Required	
W12	25.01	You Tube Premiere	M12: Flash Memory and Solid- State Drives	Recommended	

Comp Arch (Fall 2021)

Fall 2021 Edition:

- https://safari.ethz.ch/architecture/fall2021/doku. php?id=schedule
- Fall 2020 Edition:
 - https://safari.ethz.ch/architecture/fall2020/doku. php?id=schedule

Youtube Livestream (2021):

https://www.youtube.com/watch?v=4yfkM_5EFg o&list=PL5Q2soXY2Zi-Mnk1PxjEIG32HAGILkTOF

Youtube Livestream (2020):

- https://www.youtube.com/watch?v=c3mPdZA-Fmc&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN
- Master's level course
 - Taken by Bachelor's/Masters/PhD students
 - Cutting-edge research topics + fundamentals in Computer Architecture
 - 5 Simulator-based Lab Assignments
 - Potential research exploration
 - Many research readings

https://www.youtube.com/onurmutlulectures



Fall 2021 Lectures & Schedule

Watch on

> You

Week	Date	Livestream	Lecture	Readings	Lab	HW
W1	30.09 Thu.	You The Live	L1: Introduction and Basics	Required Mentioned	Lab 1 Out	HW 0 Out
	01.10 Fri.	You Tube Live	L2: Trends, Tradeoffs and Design Fundamentals (mi(PDF) mi(PPT)	Required Mentioned		
W2	07.10 Thu.	You The Live	L3a: Memory Systems: Challenges and Opportunities (PDF) (PDF) (PPT)	Described Suggested		HW 1 Out
			L3b: Course Info & Logistics			
			L3c: Memory Performance Attacks	Described Suggested		
	08.10 Fri.	You The Live	L4a: Memory Performance Attacks	Described Suggested	Lab 2 Out	
			L4b: Data Retention and Memory Refresh	Described Suggested		
			L4c: RowHammer	Described Suggested		

1. Memory system design for AI/ML workloads/accelerators

2. AI/ML techniques for improving memory system designs

Goal: Processing Inside Memory/Storage



Why In-Memory Computation Today?

Huge demand from Applications & Systems

- Data access bottleneck
- Energy & power bottlenecks
- Data movement energy dominates computation energy
- Need all at the same time: performance, energy, sustainability
- We can improve all metrics by minimizing data movement

Huge problems with Memory Technology

- Memory technology scaling is not going well (e.g., RowHammer)
- Scaling issues demand intelligence in memory + new technology

Designs are squeezed in the middle

Processing-in-Memory: Nature of Computation

Two main approaches for Processing-in-Memory:

- 1 Processing-<u>Near</u>-Memory: Design compute logic and memory separately (today) and integrate logic closer to memory
- **2 Processing-<u>Using</u>-Memory**: Use analog operational principles of memory circuitry to perform computation (no compute logic)



A PIM Taxonomy

Nature (of computation)

Using: Use operational properties of memory structures

Near: Add logic close to memory structures

Technology

□ Flash, DRAM, SRAM, RRAM, MRAM, FeRAM, PCM, 3D, ...

Location

 Sensor, Cold Storage, Hard Disk, SSD, Main Memory, Cache, Register File, Memory Controller, Interconnect, ...

• A tuple of the three determines "PIM type"

One can combine multiple "PIM types" in a system

Mindset: Memory as an Accelerator



Memory similar to a "conventional" accelerator

Accelerating Neural Network Inference

Amirali Boroumand, Saugata Ghose, Berkin Akin, Ravi Narayanaswami, Geraldo F. Oliveira, Xiaoyu Ma, Eric Shiu, and Onur Mutlu,
 "Google Neural Network Models for Edge Devices: Analyzing and

 <u>Mitigating Machine Learning Inference Bottlenecks"</u>
 Proceedings of the <u>30th International Conference on Parallel Architectures and</u>

 <u>Compilation Techniques</u> (PACT), Virtual, September 2021.
 [Slides (pptx) (pdf)]
 [Talk Video (14 minutes)]

Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks

Amirali Boroumand**Saugata Ghose*Berkin Akin*Ravi Narayanaswami*Geraldo F. Oliveira*Xiaoyu Ma*Eric Shiu*Onur Mutlu*** Carnegie Mellon Univ.* Stanford Univ.* Univ. of Illinois Urbana-Champaign * Google * ETH Zürich

Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks

Amirali BoroumandSaugata GhoseBerkin AkinRavi NarayanaswamiGeraldo F. OliveiraXiaoyu MaEric ShiuOnur Mutlu

PACT 2021



Executive Summary

<u>Context</u>: We extensively analyze a state-of-the-art edge ML accelerator (Google Edge TPU) using 24 Google edge models

- Wide range of models (CNNs, LSTMs, Transducers, RCNNs)

Problem: The Edge TPU accelerator suffers from three challenges:

- It operates significantly below its peak throughput
- It operates significantly below its theoretical energy efficiency
- It inefficiently handles memory accesses

<u>Key Insight</u>: These shortcomings arise from the monolithic design of the Edge TPU accelerator

- The Edge TPU accelerator design does not account for layer heterogeneity

Key Mechanism: A new framework called Mensa

 Mensa consists of heterogeneous accelerators whose dataflow and hardware are specialized for specific families of layers

Key Results: We design a version of Mensa for Google edge ML models

- Mensa improves performance and energy by 3.0X and 3.1X
- Mensa reduces cost and improves area efficiency

Google Edge Neural Network Models

We analyze inference execution using 24 edge NN models



Diversity Across the Models

Insight I: there is significant variation in terms of layer characteristics across the models



Diversity Within the Models

Insight 2: even within each model, layers exhibit significant variation in terms of layer characteristics

For example, our analysis of edge CNN models shows:



Variation in MAC intensity: up to 200x across layers

Variation in FLOP/Byte: up to 244x across layers



SAFARI

Heterogeneous Accelerators

Identifying Layer Families

Key observation: the majority of layers group into a small number of <u>layer families</u>



Families I & 2: low parameter footprint, high data reuse and MAC intensity \rightarrow <u>compute-centric layers</u>

Families 3, 4 & 5: high parameter footprint, low data reuse and MAC intensity \rightarrow <u>data-centric layers</u>

Mensa: Energy Reduction



Mensa-G reduces energy consumption by 3.0X compared to the baseline Edge TPU
Mensa: Throughput Improvement



Mensa-G improves inference throughput by 3.1X compared to the baseline Edge TPU

Mensa: Highly-Efficient ML Inference

Amirali Boroumand, Saugata Ghose, Berkin Akin, Ravi Narayanaswami, Geraldo F. Oliveira, Xiaoyu Ma, Eric Shiu, and Onur Mutlu,
 "Google Neural Network Models for Edge Devices: Analyzing and

 <u>Mitigating Machine Learning Inference Bottlenecks"</u>
 Proceedings of the <u>30th International Conference on Parallel Architectures and</u>

 <u>Compilation Techniques</u> (PACT), Virtual, September 2021.
 [Slides (pptx) (pdf)]
 [Talk Video (14 minutes)]

Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks

Amirali Boroumand**Saugata Ghose*Berkin Akin*Ravi Narayanaswami*Geraldo F. Oliveira*Xiaoyu Ma*Eric Shiu*Onur Mutlu*** Carnegie Mellon Univ.* Stanford Univ.* Univ. of Illinois Urbana-Champaign * Google * ETH Zürich

Accelerating In-Memory Graph Analytics

Large graphs are everywhere (circa 2015)



Scalable large-scale graph processing is challenging



Key Bottlenecks in Graph Processing



Opportunity: 3D-Stacked Logic+Memory



Hybrid Memory Cube



Tesseract System for Graph Processing

Interconnected set of 3D-stacked memory+logic chips with simple cores



Tesseract System for Graph Processing



Tesseract System for Graph Processing



Evaluated Systems



Tesseract Graph Processing Performance

>13X Performance Improvement



Tesseract Graph Processing System Energy



More on Tesseract

 Junwhan Ahn, Sungpack Hong, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi,

"A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing"

Proceedings of the <u>42nd International Symposium on Computer</u> <u>Architecture</u> (**ISCA**), Portland, OR, June 2015. [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] **Top Picks Honorable Mention by IEEE Micro.** Selected to the ISCA-50 25-Year Retrospective Issue covering 1996-2020 in 2023 (<u>Retrospective (pdf) Full</u> <u>Issue</u>).

A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing

Junwhan Ahn Sungpack Hong[§] Sungjoo Yoo Onur Mutlu[†] Kiyoung Choi junwhan@snu.ac.kr, sungpack.hong@oracle.com, sungjoo.yoo@gmail.com, onur@cmu.edu, kchoi@snu.ac.kr Seoul National University [§]Oracle Labs [†]Carnegie Mellon University

A Short Retrospective (a) 50 Years of ISCA

Retrospective: A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing

 $\begin{array}{ccc} \text{Sungpack Hong}^{\ddagger} & \text{Sungpoo Yoo}^{\nabla} & \text{Onur Mutlu}^{\$} & \text{Kiyoung Choi}^{\nabla} \\ {}^{\ddagger}Oracle \ Labs & {}^{\$}ETH \ Zürich & {}^{\nabla}Seoul \ National \ Universe \end{array}$ Junwhan Ahn[†] [†]Google DeepMind ⁷Seoul National University

Abstract—Our ISCA 2015 paper [1] provides a new pro-grammable processing-in-memory (PIM) architecture and system design that can accelerate key data-intensive applications, with design that can accelerate key data-intensive applications, with a focus on graph processing workloads. Our major idea was to completely rethink the system, including the programming model, data partitioning mechanisms, system support, instruction set architecture, along with near-memory execution units and their communication architecture, such that an important workload can be accelerated at a maximum level using a distributed system of well-connected near-memory accelerators. We built our accelerator system, Tesseract, using 3D-stacked memories with logic layers, where each logic layer contains general-purpose

with logic layers, where each logic fayer contains general-purpose processing cores and cores communicate with each other using a message-passing programming model. Cores could be specialized for graph processing (or any other application to be accelerated). To our knowledge, our paper was the first to completely design a near-memory accelerator system from scratch such that it is both generally programmable and specifically customizable to accelerate important applications, with a case study on major graph processing workloads. Ensuing work in academia and industry showed that similar approaches to system design can greatly benefit both graph processing workloads and other applications, such as machine learning, for which ideas from Tesseract seem to have been influential. This short retrospective provides a brief analysis of our ISCA

This short retrospective provides a brief analysis of our ISCA 2015 paper and its impact. We briefly describe the major ideas and contributions of the work, discuss later works that built on it or were influenced by it, and make some educated guesses on what the future may bring on PIM and accelerator systems.

I. BACKGROUND, APPROACH & MINDSET

We started our research when 3D-stacked memories (e.g., [2-4]) were viable and seemed to have promise for building effective and practical processing-near-memory systems. Such near-memory systems could lead to improvements, but there was little to no research that examined how an accelerator could be completely (re-)designed using such near-memory technology, from its hardware architecture to its programming model and software system, and what the performance and energy benefits could be of such a re-design. We set out to answer these questions in our ISCA 2015 paper [1].

We followed several major principles to design our acceler-ator from the ground up. We believe these principles are still important: a major contribution and influence of our work was in putting all of these together in a cohesive full-system design and demonstrating the large performance and energy benefits that can be obtained from such a design. We see a similar approach in many modern large-scale accelerator systems in machine learning today (e.g., [5-9]). Our principles are:

1. Near-memory execution to enable/exploit the high data access bandwidth modern workloads (e.g., graph processing) need and to reduce data movement and access latency.

2. General programmability so that the system can be easily adopted, extended, and customized for many workloads.

3. Maximal acceleration capability to maximize the performance and energy benefits. We set ourselves free from backward compatibility and cost constraints. We aimed to completely re-design the system stack. Our goal was to explore the maximal performance and energy efficiency benefits we can gain from a near-memory accelerator if we had complete freedom to change things as much as we needed. We contrast this approach to the minimal intrusion approach we also explored in a separate ISCA 2015 paper [10]

4. Customizable to specific workloads, such that we can maximize acceleration benefits. Our focus workload was graph

analytics/processing, a key workload at the time and today. However, our design principles are not limited to graph processing and the system we built is customizable to other workloads as well, e.g., machine learning, genome analysis.

5. Memory-capacity-proportional performance, i.e., processing capability should proportionally grow (i.e., scale) as memory capacity increases and vice versa. This enables scaling of data-intensive workloads that need both memory and compute.

6. Exploit new technology (3D stacking) that enables tight integration of memory and logic and helps multiple above principles (e.g., enables customizable near-memory acceleration capability in the logic layer of a 3D-stacked memory chip).

7. Good communication and scaling capability to support scalability to large dataset sizes and to enable memorycapacity-proportional performance. To this end, we provided scalable communication mechanisms between execution cores and carefully interconnected small accelerator chips to form a large distributed system of accelerator chips.

8. Maximal and efficient use of memory bandwidth to supply the high-bandwidth data access that modern workloads need. To this end, we introduced new, specialized mechanisms for prefetching and a programming model that helps leverage application semantics for hardware optimization.

II. CONTRIBUTIONS AND INFLUENCE

We believe the major contributions of our work were 1) complete rethinking of how an accelerator system should be designed to enable maximal acceleration capability, and 2) the design and analysis of such an accelerator with this mindset and using the aforementioned principles to demonstrate its effectiveness in an important class of workloads.

One can find examples of our approach in modern largescale machine learning (ML) accelerators, which are perhaps the most successful incarnation of scalable near-memory execution architectures. ML infrastructure today (e.g., [5-9]) consists of accelerator chips, each containing compute units and high-bandwidth memory tightly packaged together, and features scale-up capability enabled by connecting thousands of such chips with high-bandwidth interconnection links. The system-wide rethinking that was done to enable such accelerators and many of the principles used in such accelerators resemble our ISCA 2015 paper's approach.

The "memory-capacity-proportional performance" principle we explored in the paper shares similarities with how ML workloads are scaled up today. Similar to how we carefully sharded graphs across our accelerator chips to greatly improve effective memory bandwidth in our paper, today's ML workloads are sharded across a large number of accelerators by leveraging data/model parallelism and optimizing the placement to balance communication overheads and compute scalability [11, 12]. With the advent of large generative models requiring high memory bandwidth for fast training and inference, the scaling behavior where capacity and bandwidth are scaled together has become an essential architectural property to support modern data-intensive workloads.

The "maximal acceleration capability" principle we used in Tesseract provides much larger performance and energy improvements and better customization than the "minimalist" approach that our other ISCA 2015 paper on PIM-Enabled Instructions [10] explored: "minimally change" an existing

system to incorporate (near-memory) acceleration capability to ease programming and keep costs low. So far, the industry has more widely adopted the maximal approach to overcome the pressing scaling bottlenecks of major workloads. The key enabler that bridges the programmability gap between the maximal approach favoring large performance & energy benefits and the minimal approach favoring ease of programming is compilation techniques. These techniques lower well-defined high-level constructs into lower-level primitives [12, 13]; our ISCA 2015 papers [1, 10] and a follow-up work [14] explore them lightly. We believe that a good programming model that enables large benefits coupled with support for it across the entire system stack (including compilers & hardware) will continue to be important for effective near-memory system and accelerator designs [14]. We also believe that the maximal versus minimal approaches that are initially explored in our two ISCA 2015 papers is a useful way of exploring emerging technologies (e.g., near-memory accelerators) to better understand the tradeoffs of system designs that exploit such technologies.

III. INFLUENCE ON LATER WORKS

Our paper was at the beginning of a proliferation of scalable near-memory processing systems designed to accelerate key applications (see [15] for many works on the topic). Tesseract has inspired many near-memory system ideas (e.g., [16-28]) and served as the de facto comparison point for such systems, including near-memory graph processing accelerators that built on Tesseract and improved various aspects of Tesseract. Since machine learning accelerators that use high-bandwidth memory (e.g., [5, 29]) and industrial PIM prototypes (e.g., [30-41]) are now in the market, near-memory processing is no longer an "eccentric" architecture it used to be when Tesseract was originally published.

Graph processing & analytics workloads remain as an important and growing class of applications in various forms, ranging from large-scale industrial graph analysis engines (e.g., [42]) to graph neural networks [43]. Our focus on largescale graph processing in our ISCA 2015 paper increased attention to this domain in the computer architecture community, resulting in subsequent research on efficient hardware architectures for graph processing (e.g., [44-46]).

IV. SUMMARY AND FUTURE OUTLOOK

We believe that our ISCA 2015 paper's principled rethinking of system design to accelerate an important class of data-intensive workloads provided significant value and enabled/influenced a large body of follow-on works and ideas. We expect that such rethinking of system design for key workloads, especially with a focus on "maximal acceleration capability," will continue to be critical as pressing technology and application scaling challenges increasingly require us to think differently to substantially improve performance and energy (as well as other metrics). We believe the principles exploited in Tesseract are fundamental and they will remain useful and likely become even more important as systems become more constrained due to the continuously-increasing memory access and computation demands of future workloads. We also project that as hardware substrates for near-memory acceleration (e.g., 3D stacking, in-DRAM computation, NVMbased PIM, processing using memory [15]) evolve and mature, systems will take advantage of them even more, likely using principles similar to those used in the design of Tesseract.

REFERENCES

- J. Ahn et al., "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing," in ISCA, 2015.
 Hybrid Memory Cube Consortium, "HMC Specification 1.1," 2013.
 J. Jeddeloh and B. Keeth, "Hybrid Memory Cube: New DRAM Archi-tecture Increases Density and Performance" in VLZT, 2012.
 IEDEC, "High Bandwidth Memory (HBM) DRAM," Standard No. JESDD253, 2013.

- [5] N. Jouppi et al., "TPU v4: An Optically Reconfigurable Supercomputer for Machine Learning with Hardware Support for Embedding," in ISCA,
- 2023.
 [6] J. Fowers et al., "A Configurable Cloud-Scale DNN Processor for Real-Time AL," in ISCA, 2018.
 [7] S. Lie, "Cerebras Architecture Deep Dive: First Look Inside the Hard-warr/Software Co-Design for Deep Learning," in *IEEE Micro*, 2023.
 [8] E. Talpes et al., "The Microarchitecture of DOIO, Tesla's Exa-Scale Computer," in *IEEE Micro*, 2023.
 [9] A. Ishin and R., "WILink-Network Switch NVIDIA's Switch Chip for High Communication-Bandwidth SuperPODs," in *Hot Chips*, 2022.
- J. Ahn et al., "PIM-Enabled Instructions: A Low-Overhead, Locality-
- Aware Processing-in-Memory Architecture," in ISCA, 2015. [11] R. Pope et al., "Efficiently Scaling Transformer Inference," in MLSys,
- 2021.
 2023.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.
 2024.

- [15] O. Mutlu et al., "A Modern Primer on Processing in Memory," Emerging Computing: From Devices to Systems, 2021, https://arxiv.org/abs/2012

- Computing: From Devices to Systems, 2021, https://arxiv.org/abs/2012.
 [16] M. Zhang et al., "GraphP: Reducing Communication for PIM-Based Processing with Efficient Data Partition," in *HPCA*, 2018.
 [17] Graph Processing With Efficient Data Partition, "in *HPCA*, 2018.
 [18] Y. Zhuo et al., "GraphI: Scalable PIM-Based Graph Processing," in *MCRO*, 2019.
 [19] G. Dai et al., "GraphI: A Processing-in-Memory Architecture for Large-Scale Graph Processing," *IEE TCAD*, 2018.
 [20] G. Li et al., "GraphI: A Processing-in-Memory Architecture for Large-Scale Graph Processing," *IEE TCAD*, 2018.
 [21] S. Rheind et al., "RDBEYS: Near-Memory Graph Copy Enhanced System-Software," in *MEMSYS*, 2019.
 [22] S. Rheind et al., "GraphI: A Processing Computing," in *JCA*, 2018.
 [23] N. Challapalle et al., "Graph Analytics, "Craph Analytics Accelerator Supporting Sparse Data Representation using Crossbar Architectures," in *ISCA*, 2020.
 [24] M. Zhou et al., "Graph Lacelefficient Accelerator for Carge-Scale Graph Processing," in *ICA*, 2018.

- M. Zhou et al., "Ultra Efficient Acceleration for De Novo Genome Assembly via Near-Memory Computing," in PACT, 2021.
 X. Xie et al., "SpaceA: Sparse Matrix Vector Multiplication on Processing-in-Memory Accelerator," in HPCA, 2021.
 M. Zhou et al., "Hydraph: Accelerator, and the processing with Hybrid Memory-Centric Computing," in DATE, 2021.
 M. Zhou et al., "Hydraph: Accelerator, and the processing with Hybrid Dispatching and Hybrid Partitioning in PIM-based Accelerators," in ISCA. 2002.
- M. Orenes-Vera et al., "Dalorex: A Data-Local Program Execution and Architecture for Memory-Bound Applications," in *HPCA*, 2023.
 J. Choquette, "Nvidia Hopper GPU: Scaling Performance," in *Hot Chips*,
- 2022. [30] F. Devaux, "The True Processing In Memory Accelerator," in *Hot Chips*
- 2019.
 [31] J. Gómez-Luna et al., "Benchmarking a New Paradigm: Experimental Analysis and Characterization of a Real Processing-in-Memory System," *IEEE Access*, 2022.
 [31] J. Gomez-Luna et al., "Evaluating Machine Learning Workloads on

- Analysis and Characterization of a Keal Processing-in-Memory System, *IEEE Accels*, 2022. d. "Evaluating Machine Learning Workloads on Memory-Centric Computing Systems," in *ISPASS*, 2023.
 S. Lee et al., "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology: Industrial Product," in *ISCA*, 2021.
 Y.-C. Kwon et al., "25 A 2 20m 6GB Function-In-Memory DRAM, Based on HBM2 with a 1.2 TFLOPS Programmable Computing Unit Using Bank-Level Parallelism, for Machine Learning Applications," in ISSCC, 2021.
 L. Ke et al., "Near-Memory Processing in Action: Accelerating Person- Bized Recommendation with ADDIMS," *IEEE Micro*, 2021.
 S. Lee et al., "A lymm 1,25V 8Gb, 16Gb/spin GDDR6-based Accelerator-in-Memory Supporting TFLOPS MAC Operations and Var- ious Activation Functions for Deep-Learning Applications," in *ISSCC*, 2022.

- ious Activation Functions for Deep-Learning Applications," in *ISSCC*, 2022.
 [38] D. Niu et al., "184QPS/W 64Mb/mm² 3D Logic-to-DRAM Hybrid Bonding with Process-Near-Memory Engine for Recommendation System," in *ISSCC*, 2022.
 [39] Y. Kwon, "System Architecture and Software Stack for GDDR6-AiM," (40) G Kingh 22 at, "PFGA-based Near-Memory Acceleration of Modern Data-Intensive Applications," *IEEE Micro*, 2021.
 [41] G. Singh et al., "Accelerating Weather Prediction using Near-Memory Reconfigurable Fabric," *ACM TRETS*, 2021.
 [42] S. Hong et al., "PGXD: A Fast Distributed Graph Processing Engine," in *SC*, 2015.
 [43] T. N. Kipf and N. Welling, "Semi Spervised Classification with Graph Learning and CR, 2017.
 [44] L. Nai et al., "GraphWith: Enabling Instruction-Level PIM Offloading in Graph Computing Frameworks," in *HPCA*, 2017.
 [45] M. Besta et al., "Six18. Set-Centric Instruction Set Architecture for Graph Mining on Processing-in-Memory Systems," in *MICRO*, 2021.
 [46] T. J. Ham et al., "Graphicionado: A High-Performance and Energy-Efficient Accelerator for Graph Analytics," in *MICRO*, 2016.

SAFARI

https://arxiv.org/pdf/2306.16093

Accelerating Graph Pattern Mining

Maciej Besta, Raghavendra Kanakagiri, Grzegorz Kwasniewski, Rachata Ausavarungnirun, Jakub Beránek, Konstantinos Kanellopoulos, Kacper Janda, Zur Vonarburg-Shmaria, Lukas Gianinazzi, Ioana Stefan, Juan Gómez-Luna, Marcin Copik, Lukas Kapp-Schwoerer, Salvatore Di Girolamo, Nils Blach, Marek Konieczny, Onur Mutlu, and Torsten Hoefler,
 <u>"SISA: Set-Centric Instruction Set Architecture for Graph Mining on Processing-in-Memory Systems"</u>
 Proceedings of the <u>54th International Symposium on Microarchitecture</u> (<i>MICRO), Virtual, October 2021.
 [Slides (pdf)]
 [Talk Video (22 minutes)]
 [Lightning Talk Video (1.5 minutes)]
 [Full arXiv version]

SISA: Set-Centric Instruction Set Architecture for Graph Mining on Processing-in-Memory Systems

Maciej Besta¹, Raghavendra Kanakagiri², Grzegorz Kwasniewski¹, Rachata Ausavarungnirun³, Jakub Beránek⁴, Konstantinos Kanellopoulos¹, Kacper Janda⁵, Zur Vonarburg-Shmaria¹, Lukas Gianinazzi¹, Ioana Stefan¹, Juan Gómez-Luna¹, Marcin Copik¹, Lukas Kapp-Schwoerer¹, Salvatore Di Girolamo¹, Nils Blach¹, Marek Konieczny⁵, Onur Mutlu¹, Torsten Hoefler¹ ¹ETH Zurich, Switzerland ²IIT Tirupati, India ³King Mongkut's University of Technology North Bangkok, Thailand ⁴Technical University of Ostrava, Czech Republic ⁵AGH-UST, Poland

Processing using DRAM

Background Work: RowClone

 Vivek Seshadri, Yoongu Kim, Chris Fallin, Donghyuk Lee, Rachata Ausavarungnirun, Gennady Pekhimenko, Yixin Luo, Onur Mutlu, Michael A. Kozuch, Phillip B. Gibbons, and Todd C. Mowry,
 "RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization"

Proceedings of the <u>46th International Symposium on Microarchitecture</u> (**MICRO**), Davis, CA, December 2013. [<u>Slides (pptx) (pdf)</u>] [<u>Lightning Session</u> <u>Slides (pptx) (pdf)</u>] [<u>Poster (pptx) (pdf)</u>]

RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization

Vivek Seshadri Yoongu Kim Chris Fallin^{*} Donghyuk Lee vseshadr@cs.cmu.edu yoongukim@cmu.edu cfallin@c1f.net donghyuk1@cmu.edu Rachata Ausavarungnirun Gennady Pekhimenko Yixin Luo rachata@cmu.edu gpekhime@cs.cmu.edu yixinluo@andrew.cmu.edu Onur Mutlu Phillip B. Gibbons[†] Michael A. Kozuch[†] Todd C. Mowry onur@cmu.edu phillip.b.gibbons@intel.com michael.a.kozuch@intel.com tcm@cs.cmu.edu Carnegie Mellon University [†]Intel Pittsburgh

Background Work: PiDRAM

 Ataberk Olgun, Juan Gomez Luna, Konstantinos Kanellopoulos, Behzad Salami, Hasan Hassan, Oguz Ergin, and Onur Mutlu, "PiDRAM: A Holistic End-to-end FPGA-based Framework for Processing-in-DRAM" ACM Transactions on Architecture and Code Optimization (TACO), March 2023.
 [arXiv version] Presented at the 18th HiPEAC Conference, Toulouse, France, January 2023.
 [Slides (pptx) (pdf)]
 [Longer Lecture Slides (pptx) (pdf)]
 [Lecture Video (40 minutes)]
 [PiDRAM Source Code]

PiDRAM: A Holistic End-to-end FPGA-based Framework for <u>Processing-in-DRAM</u>

Ataberk Olgun§

Juan Gómez Luna[§] Konstantinos Kanellopoulos[§] Hasan Hassan[§] Oguz Ergin[†] Onur Mutlu[§]

[§]ETH Zürich [†]TOBB University of Economics and Technology

Behzad Salami[§]

Background Work: In-DRAM Bulk AND/OR

 Vivek Seshadri, Kevin Hsieh, Amirali Boroumand, Donghyuk Lee, Michael A. Kozuch, Onur Mutlu, Phillip B. Gibbons, and Todd C. Mowry,
 <u>"Fast Bulk Bitwise AND and OR in DRAM"</u> <u>IEEE Computer Architecture Letters</u> (CAL), April 2015.

Fast Bulk Bitwise AND and OR in DRAM

Vivek Seshadri*, Kevin Hsieh*, Amirali Boroumand*, Donghyuk Lee*, Michael A. Kozuch[†], Onur Mutlu*, Phillip B. Gibbons[†], Todd C. Mowry* *Carnegie Mellon University [†]Intel Pittsburgh

Background Work: Ambit

 Vivek Seshadri, Donghyuk Lee, Thomas Mullins, Hasan Hassan, Amirali Boroumand, Jeremie Kim, Michael A. Kozuch, Onur Mutlu, Phillip B. Gibbons, and Todd C. Mowry,
 <u>"Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using</u> <u>Commodity DRAM Technology"</u> *Proceedings of the <u>50th International Symposium on</u> <u>Microarchitecture</u> (<i>MICRO*), Boston, MA, USA, October 2017. [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Poster (pptx) (pdf)]

Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology

Vivek Seshadri^{1,5} Donghyuk Lee^{2,5} Thomas Mullins^{3,5} Hasan Hassan⁴ Amirali Boroumand⁵ Jeremie Kim^{4,5} Michael A. Kozuch³ Onur Mutlu^{4,5} Phillip B. Gibbons⁵ Todd C. Mowry⁵

¹Microsoft Research India ²NVIDIA Research ³Intel ⁴ETH Zürich ⁵Carnegie Mellon University

Background: In-DRAM Bulk Bitwise Execution

 Vivek Seshadri and Onur Mutlu,
 "In-DRAM Bulk Bitwise Execution Engine" Invited Book Chapter in Advances in Computers, to appear in 2020.
 [Preliminary arXiv version]

In-DRAM Bulk Bitwise Execution Engine

Vivek Seshadri Microsoft Research India visesha@microsoft.com

Onur Mutlu ETH Zürich onur.mutlu@inf.ethz.ch

Background: SIMDRAM Framework

 Nastaran Hajinazar, Geraldo F. Oliveira, Sven Gregorio, Joao Dinis Ferreira, Nika Mansouri Ghiasi, Minesh Patel, Mohammed Alser, Saugata Ghose, Juan Gomez-Luna, and Onur Mutlu, "SIMDRAM: An End-to-End Framework for Bit-Serial SIMD Computing in DRAM" Proceedings of the 26th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), Virtual, March-April 2021.
 [2-page Extended Abstract]
 [Short Talk Slides (pptx) (pdf)]
 [Talk Slides (pptx) (pdf)]
 [Short Talk Video (5 mins)]
 [Full Talk Video (27 mins)]

SIMDRAM: A Framework for Bit-Serial SIMD Processing using DRAM

*Nastaran Hajinazar^{1,2} Nika Mansouri Ghiasi¹ Juan Gómez-Luna¹

Sven Gregorio¹ Mohammed Alser¹ Onur Mutlu¹ João Dinis Ferreira¹ Saugata Ghose³

¹ETH Zürich ²Simon Fraser University

³University of Illinois at Urbana–Champaign

In-DRAM Lookup-Table Based Execution

João Dinis Ferreira, Gabriel Falcao, Juan Gómez-Luna, Mohammed Alser, Lois Orosa, Mohammad Sadrosadati, Jeremie S. Kim, Geraldo F. Oliveira, Taha Shahroodi, Anant Nori, and Onur Mutlu, "pLUTo: Enabling Massively Parallel Computation in DRAM via Lookup Tables" *Proceedings of the <u>55th International Symposium on Microarchitecture</u> (<i>MICRO*), Chicago, IL, USA, October 2022. [Slides (pptx) (pdf)] [Longer Lecture Slides (pptx) (pdf)] [Lecture Video (26 minutes)] [arXiv version] [Source Code (Officially Artifact Evaluated with All Badges)] *Officially artifact evaluated as available, reusable and reproducible.*



pLUTo: Enabling Massively Parallel Computation in DRAM via Lookup Tables

Mohammed Alser§ João Dinis Ferreira[§] Gabriel Falcao[†] Juan Gómez-Luna§ Mohammad Sadrosadati[§] Lois Orosa[§]∇ Jeremie S. Kim[§] Geraldo F. Oliveira§ Taha Shahroodi[‡] Anant Nori* Onur Mutlu[§] [‡]TU Delft §ETH Zürich [†]IT, University of Coimbra [∇]*Galicia Supercomputing Center* *Intel

SAFARI

https://arxiv.org/pdf/2104.07699.pdf

MIMDRAM: More Flexible Processing using DRAM

Appears at HPCA 2024 <u>https://arxiv.org/pdf/2402.19080.pdf</u>

MIMDRAM: An End-to-End Processing-Using-DRAM System for High-Throughput, Energy-Efficient and Programmer-Transparent Multiple-Instruction Multiple-Data Computing

Geraldo F. Oliveira[†]Ataberk Olgun[†]Abdullah Giray Yağlıkçı[†]F. Nisa Bostancı[†]Juan Gómez-Luna[†]Saugata Ghose[‡]Onur Mutlu[†]

† ETH Zürich ‡ Univ. of Illinois Urbana-Champaign

Our goal is to design a flexible PUD system that overcomes the limitations caused by the large and rigid granularity of PUD. To this end, we propose MIMDRAM, a hardware/software co-designed PUD system that introduces new mechanisms to allocate and control only the necessary resources for a given PUD operation. The key idea of MIMDRAM is to leverage finegrained DRAM (i.e., the ability to independently access smaller segments of a large DRAM row) for PUD computation. MIM-DRAM exploits this key idea to enable a multiple-instruction multiple-data (MIMD) execution model in each DRAM subarray (and SIMD execution within each DRAM row segment).

MIMDRAM: Executive Summary

Problem: Processing-Using-DRAM (PUD) suffers from three issues caused by DRAM's large and rigid access granularity

- <u>Underutilization</u> due to data parallelism variation in (and across) applications
- <u>Limited computation support</u> due to a lack of interconnects
- <u>Challenging programming</u> model due to a lack of compilers

Goal: Design a flexible PUD system that overcomes the three limitations caused by DRAM's large and rigid access granularity

Key Mechanism: MIMDRAM, a hardware/software co-design PUD system

- Key idea: leverage fine-grained DRAM for PUD operation
- **HW**: <u>simple changes</u> to the DRAM array, enabling concurrent PUD operations
 - low-cost interconnects at the DRAM peripherals for data reduction
- SW: <u>compiler</u> and <u>OS</u> support to generate and map PUD instructions

Key Results: MIMDRAM achieves

- **14.3***x*, **30.6***x*, and **6.8***x* the energy efficiency of state-of-the-art PUD systems, a high-end CPU and GPU, respectively
- Small area cost to a DRAM chip (1.11%) and CPU die (0.6%)

SAFARI

https://github.com/CMU-SAFARI/MIMDRAM

Real DRAM Chips Are Already Quite Capable: FC-DRAM & SiMRA

DRAM Chips Are Already (Quite) Capable!

Appears at HPCA 2024 <u>https://arxiv.org/pdf/2402.18736.pdf</u>

Functionally-Complete Boolean Logic in Real DRAM Chips: Experimental Characterization and Analysis

İsmail Emir Yüksel Yahya Can Tuğrul Ataberk Olgun F. Nisa Bostancı A. Giray Yağlıkçı Geraldo F. Oliveira Haocong Luo Juan Gómez-Luna Mohammad Sadrosadati Onur Mutlu

ETH Zürich

We experimentally demonstrate that COTS DRAM chips are capable of performing 1) functionally-complete Boolean operations: NOT, NAND, and NOR and 2) many-input (i.e., more than two-input) AND and OR operations. We present an extensive characterization of new bulk bitwise operations in 256 off-theshelf modern DDR4 DRAM chips. We evaluate the reliability of these operations using a metric called success rate: the fraction of correctly performed bitwise operations. Among our 19 new observations, we highlight four major results. First, we can perform the NOT operation on COTS DRAM chips with 98.37% success rate on average. Second, we can perform up to 16-input NAND, NOR, AND, and OR operations on COTS DRAM chips with high reliability (e.g., 16-input NAND, NOR, AND, and OR with average success rate of 94.94%, 95.87%, 94.94%, and 95.85%, respectively). Third, data pattern only slightly

DRAM Chips Are Already (Quite) Capable!

https://arxiv.org/pdf/2312.02880.pdf

PULSAR: Simultaneous Many-Row Activation for Reliable and High-Performance Computing in Off-the-Shelf DRAM Chips

Ismail Emir Yuksel Yahya Can Tugrul F. Nisa Bostanci Abdullah Giray Yaglikci Ataberk Olgun Geraldo F. Oliveira Melina Soysal Haocong Luo Juan Gomez Luna Mohammad Sadrosadati Onur Mutlu

ETH Zurich

We propose PULSAR, a new technique to enable highsuccess-rate and high-performance PuM operations in off-theshelf DRAM chips. PULSAR leverages our new observation that a carefully-crafted sequence of DRAM commands simultaneously activates up to 32 DRAM rows. PULSAR overcomes the limitations of existing techniques by 1) replicating the input data to improve the success rate and 2) enabling new bulk bitwise operations (e.g., many-input majority, *Multi-RowInit*, and *Bulk-Write*) to improve the performance.

DRAM Chips Are Already (Quite) Capable!

Appears at DSN 2024



Simultaneous Many-Row Activation in Off-the-Shelf DRAM Chips: Experimental Characterization and Analysis

İsmail Emir Yüksel¹ Yahya Can Tuğrul^{1,2} F. Nisa Bostancı¹ Geraldo F. Oliveira¹ A. Giray Yağlıkçı¹ Ataberk Olgun¹ Melina Soysal¹ Haocong Luo¹ Juan Gómez-Luna¹ Mohammad Sadrosadati¹ Onur Mutlu¹ ¹ETH Zürich ²TOBB University of Economics and Technology

The Capability of COTS DRAM Chips

We **demonstrate** that **COTS DRAM chips**:

Can simultaneously activate up to48 rows in two neighboring subarrays

Can perform **NOT operation** with up to **32 output operands**

Can perform up to **16-input** AND, NAND, OR, and NOR operations

2

3

Finding: SiMRA Across Subarrays

Activating two rows in **quick succession** can **simultaneously** activate **multiple rows in neighboring subarrays**



Key Idea: NOT Operation

Connect rows in neighboring subarrays through **a NOT gate** by simultaneously activating rows



Key Idea: NAND, NOR, AND, OR

Manipulate the bitline voltage to express a wide variety of functions using

multiple-row activation in neighboring subarrays







SAFARI *Gao et al., "FracDRAM: Fractional Values in Off-the-Shelf DRAM," in MICRO, 2022. **213**





 $V_{DD} = 1 \& GND = 0$









 $V_{DD}=1 \& GND = 0$







 $V_{DD}=1 \& GND = 0$


Two-Input AND and NAND Operations





 V_{DD} =1 & GND = 0



Two-Input AND and NAND Operations





 $V_{DD} = 1 \& GND = 0$ COM **REF** Х $\left(\right)$ 0 1 $\mathbf{0}$ 1 \mathbf{O} $\left(\right)$ 1 1 1 0 1 AND NAND

Many-Input AND, NAND, OR, and NOR Operations

We can express AND, NAND, OR, and NOR operations by carefully manipulating the **reference voltage**

Functionally-Complete Boolean Logic in Real DRAM Chips: Experimental Characterization and Analysis

İsmail Emir Yüksel Yahya Can Tuğrul Ataberk Olgun F. Nisa Bostancı A. Giray Yağlıkçı Geraldo F. Oliveira Haocong Luo Juan Gómez-Luna Mohammad Sadrosadati Onur Mutlu

ETH Zürich

(More details in the paper)

https://arxiv.org/pdf/2402.18736.pdf

DRAM Testing Infrastructure

- Developed from DRAM Bender [Olgun+, TCAD'23]*
- Fine-grained control over DRAM commands, timings, and temperature



SAFAR *Olgun et al., "<u>DRAM Bender: An Extensible and Versatile FPGA-based Infrastructure</u> to Easily Test State-of-the-art DRAM Chips," TCAD, 2023.

DRAM Chips Tested

- 256 DDR4 chips from two major DRAM manufacturers
- Covers different die revisions and chip densities

Chip Mfr.	#Modules (#Chips)	Die Rev.	Mfr. Date ^a	Chip Density	Chip Org.	Speed Rate
SK Hynix	9 (72)	М	N/A	4Gb	x8	2666MT/s
	5 (40)	А	N/A	4Gb	x8	2133MT/s
	1 (16)	А	N/A	8Gb	x8	2666MT/s
	1 (32)	А	18-14	4Gb	x4	2400MT/s
	1 (32)	А	16-49	8Gb	x4	2400MT/s
	1 (32)	М	16-22	8Gb	x4	2666MT/s
Samsung	1 (8)	F	21-02	4Gb	x8	2666MT/s
	2 (16)	D	21-10	8Gb	x8	2133MT/s
	1 (8)	А	22-12	8Gb	x8	3200MT/s

Performing AND, NAND, OR, and NOR



COTS DRAM chips can perform {2, 4, 8, 16}-input AND, NAND, OR, and NOR operations

Performing AND, NAND, OR, and NOR



COTS DRAM chips can perform 16-input AND, NAND, OR, and NOR operations with very high success rate (>94%)











Data pattern slightly affects

the reliability of AND, NAND, OR, and NOR operations

Available on arXiv

Functionally-Complete Boolean Logic in Real DRAM Chips: Experimental Characterization and Analysis

İsmail Emir Yüksel Yahya Can Tuğrul Ataberk Olgun F. Nisa Bostancı A. Giray Yağlıkçı Geraldo F. Oliveira Haocong Luo Juan Gómez-Luna Mohammad Sadrosadati Onur Mutlu

ETH Zürich

Processing-using-DRAM (PuD) is an emerging paradigm that leverages the analog operational properties of DRAM circuitry to enable massively parallel in-DRAM computation. PuD has the potential to significantly reduce or eliminate costly data movement between processing elements and main memory. A common approach for PuD architectures is to make use of bulk bitwise computation (e.g., AND, OR, NOT). Prior works experimentally demonstrate three-input MAJ (i.e., MAJ3) and two-input AND and OR operations in commercial off-the-shelf (COTS) DRAM chips. Yet, demonstrations on COTS DRAM chips do not provide a functionally complete set of operations (e.g., NAND or AND and NOT).

We experimentally demonstrate that COTS DRAM chips are capable of performing 1) functionally-complete Boolean operations: NOT, NAND, and NOR and 2) many-input (i.e., more than two-input) AND and OR operations. We present an extensive systems and applications [12, 13]. Processing-using-DRAM (PuD) [29–32] is a promising paradigm that can alleviate the data movement bottleneck. PuD uses the analog operational properties of the DRAM circuitry to enable massively parallel in-DRAM computation. Many prior works [29–53] demonstrate that PuD can greatly reduce or eliminate data movement.

A widely used approach for PuD is to perform bulk bitwise operations, i.e., bitwise operations on large bit vectors. To perform bulk bitwise operations using DRAM, prior works propose modifications to the DRAM circuitry [29–31, 33, 35, 36, 43, 44, 46, 48–58]. Recent works [38, 41, 42, 45] experimentally demonstrate the feasibility of executing data copy & initialization [42, 45], i.e., the RowClone operation [49], and a subset of bitwise operations, i.e., three-input bitwise majority (MAJ3) and two-input AND and OR operations in unmodified commercial off-the-shelf (COTS) DRAM chips by operating beyond

https://arxiv.org/pdf/2402.18736.pdf

Summary

- We experimentally demonstrate that commercial off-the-shelf (COTS) DRAM chips can perform:
 - **Functionally-complete** Boolean operations: NOT, NAND, and NOR
 - Up to 16-input AND, NAND, OR, and NOR operations
- We characterize the success rate of these operations on 256 COTS DDR4 chips from two major manufacturers
- We highlight **two key results**:
 - We can perform NOT and
 {2, 4, 8, 16}-input AND, NAND, OR, and NOR operations on COTS DRAM chips with very high success rates (>94%)
 - Data pattern and temperature only slightly affect the reliability of these operations

We believe these empirical results demonstrate the promising potential of using DRAM as a computation substrate

Simultaneous Many-Row Activation in Off-the-Shelf DRAM Chips Experimental Characterization and Analysis





ETH zürich

İsmail Emir Yüksel

Yahya C. Tuğrul F. Nisa Bostancı Geraldo F. Oliveira

A. Giray Yağlıkçı Ataberk Olgun Melina Soysal Haocong Luo

Juan Gómez–Luna Mohammad Sadr Onur Mutlu



In-DRAM Multiple Row Copy (Multi-RowCopy)

Simultaneously activate many rows to copy **one row's content** to **multiple destination rows**

RowClone



Multi-RowCopy



SAFARI

[Seshadri+ MICRO'13]

Key Takeaways from Multi-RowCopy

Key Takeaway 1

COTS DRAM chips are capable of copying one row's data to 1, 3, 7, 15, and 31 other rows at very high success rates

Key Takeaway 2

Multi-RowCopy in COTS DRAM chips is highly resilient to changes in data pattern, temperature, and wordline voltage



Robustness of Multi-RowCopy



COTS DRAM chips can copy one row's content to up to 31 rows with a very high success rate





Data pattern has a small effect

on the success rate of the Multi-RowCopy operation

Also in the Paper: Impact of Temperature & Voltage



Available on arXiv

which we call Multi-RowCopy. Second, storing multiple copies

of MAJX's input operands on all simultaneously activated rows

drastically increases the success rate (i.e., the percentage of

DRAM cells that correctly perform the computation) of the

MAJX operation. For example, MAJ3 with 32-row activation (i.e.,



that bulk bitwise operations are used in a wide variety of important applications, including databases and web search [64, 67, 79, 130, 133–140], data analytics [64, 141–144], graph processing [56, 80, 94, 130, 145], genome analysis [60, 99, 146–149], cryptography [150, 151], set operations [56, 64], and hyperdimensional computing [152–154].

https://arxiv.org/pdf/2405.06081

Our Work is Open Source and Artifact Evaluated

Code Repro	oducible	Dataset Reproducible		
Simra-Dram Public		☆ Edit Pins ▼ ③ Watch 4 ▼	v v Fork 0 v r ★ Starred 6 v	
ያ main 👻 ያግ Branch 🛇 0 Tags	Q Go to file	t Add file - Code -	About ĝ	
🛎 unrealismail Update README.md		a51abfa · last month 🛛 🕙 5 Commits	Source code & scripts for experimental characterization and demonstration of 1)	
DRAM-Bender	initial comit	last month	simultaneous many-row activation, 2) up to nine-input majority operations and 3)	
analysis	initial comit	last month	copying one row's content to up 31 rows	
experimental_data	initial comit	last month	our DSN'24 paper by Yuksel et al. at	
	initial comit	last month	https://arxiv.org/abs/2405.06081	
C README.md	Update README.md	last month	🔟 Readme গ্রুষ্ট View license	
다 README 최 License		Ø∷≣	✓ ActivityE Custom properties	
Simultaneous Ma DRAM Chips: Exp	ny-Row Activation in (erimental Characteriza	Off-the-Shelf ation and Analysis	 ☆ 6 stars ◆ 4 watching ※ 0 forks Report repository 	

https://github.com/CMU-SAFARI/SiMRA-DRAM

A Modern Primer on Processing in Memory

Onur Mutlu^{a,b}, Saugata Ghose^{b,c}, Juan Gómez-Luna^a, Rachata Ausavarungnirun^d

SAFARI Research Group

^aETH Zürich ^bCarnegie Mellon University ^cUniversity of Illinois at Urbana-Champaign ^dKing Mongkut's University of Technology North Bangkok

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna, and Rachata Ausavarungnirun, <u>"A Modern Primer on Processing in Memory"</u> *Invited Book Chapter in <u>Emerging Computing: From Devices to Systems -</u> <u>Looking Beyond Moore and Von Neumann</u>, Springer, to be published in 2021.*

PIM Review and Open Problems (II)

A Workload and Programming Ease Driven Perspective of Processing-in-Memory

Saugata Ghose†Amirali Boroumand†Jeremie S. Kim†§Juan Gómez-Luna§Onur Mutlu§††Carnegie Mellon University§ETH Zürich

Saugata Ghose, Amirali Boroumand, Jeremie S. Kim, Juan Gomez-Luna, and Onur Mutlu, "Processing-in-Memory: A Workload-Driven Perspective" *Invited Article in IBM Journal of Research & Development, Special Issue on Hardware for Artificial Intelligence*, to appear in November 2019. [Preliminary arXiv version]

SAFARI

https://arxiv.org/pdf/1907.12947.pdf

Processing in Memory: Adoption Challenges

Processing using Memory
 Processing near Memory

Eliminating the Adoption Barriers

How to Enable Adoption of Processing in Memory

Potential Barriers to Adoption of PIM

1. Applications & software for PIM

2. Ease of **programming** (interfaces and compiler/HW support)

3. **System** and **security** support: coherence, synchronization, virtual memory, isolation, communication interfaces, ...

4. **Runtime** and **compilation** systems for adaptive scheduling, data mapping, access/sharing control, ...

5. Infrastructures to assess benefits and feasibility

All can be solved with change of mindset

We Need to Revisit the Entire Stack

• With a **memory-centric mindset**

Problem	
Aigorithm	
Program/Language	
System Software	
SW/HW Interface	
Micro-architecture	
Logic	J
Devices	
Electrons	

We can get there step by step

Processing-in-Memory Landscape Today



SAFARI

And, many other experimental chips and startups

Adoption: How to Keep It Simple?

 Junwhan Ahn, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi, "PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture" Proceedings of the <u>42nd International Symposium on</u> <u>Computer Architecture</u> (ISCA), Portland, OR, June 2015. [Slides (pdf)] [Lightning Session Slides (pdf)]

PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture

Junwhan Ahn Sungjoo Yoo Onur Mutlu[†] Kiyoung Choi junwhan@snu.ac.kr, sungjoo.yoo@gmail.com, onur@cmu.edu, kchoi@snu.ac.kr Seoul National University [†]Carnegie Mellon University

Adoption: How to Keep It Simple?

 Junwhan Ahn, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi, "PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture" Proceedings of the <u>42nd International Symposium on</u> <u>Computer Architecture</u> (ISCA), Portland, OR, June 2015. [Slides (pdf)] [Lightning Session Slides (pdf)]

PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture

Junwhan Ahn Sungjoo Yoo Onur Mutlu[†] Kiyoung Choi junwhan@snu.ac.kr, sungjoo.yoo@gmail.com, onur@cmu.edu, kchoi@snu.ac.kr Seoul National University [†]Carnegie Mellon University

Adoption: How to Ease **Programmability**? (I)

 Kevin Hsieh, Eiman Ebrahimi, Gwangsun Kim, Niladrish Chatterjee, Mike O'Connor, Nandita Vijaykumar, Onur Mutlu, and Stephen W. Keckler, "Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems" Proceedings of the <u>43rd International Symposium on Computer</u>

Architecture (ISCA), Seoul, South Korea, June 2016.

[Slides (pptx) (pdf)]

[Lightning Session Slides (pptx) (pdf)]

Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems

Kevin Hsieh[‡] Eiman Ebrahimi[†] Gwangsun Kim^{*} Niladrish Chatterjee[†] Mike O'Connor[†] Nandita Vijaykumar[‡] Onur Mutlu^{§‡} Stephen W. Keckler[†] [‡]Carnegie Mellon University [†]NVIDIA ^{*}KAIST [§]ETH Zürich

Truly Distributed GPU Processing with PIM



Adoption: How to Ease Programmability? (II)

 Geraldo F. Oliveira, Alain Kohli, David Novo, Juan Gómez-Luna, Onur Mutlu,
 "DaPPA: A Data-Parallel Framework for Processingin-Memory Architectures,"
 in PACT SRC Student Competition, Vienna, Austria, October 2023.

DaPPA: A Data-Parallel Framework for Processing-in-Memory Architectures

Geraldo F. Oliveira*

Alain Kohli*

David Novo[‡] Juan Gómez-Luna* Onur Mutlu* [‡]LIRMM, Univ. Montpellier, CNRS

Adoption: How to Ease **Programmability**? (III)

 Jinfan Chen, Juan Gómez-Luna, Izzat El Hajj, YuXin Guo, and Onur Mutlu,
 "SimplePIM: A Software Framework for Productive and Efficient Processing in Memory"
 Proceedings of the <u>32nd International Conference on</u> Parallel Architectures and Compilation Techniques (PACT), Vienna, Austria, October 2023.

SimplePIM: A Software Framework for Productive and Efficient Processing-in-Memory

Jinfan Chen¹ Juan Gómez-Luna¹ Izzat El Hajj² Yuxin Guo¹ Onur Mutlu¹ ¹ETH Zürich ²American University of Beirut

Adoption: How to Ease Programmability? (IV)

Geraldo F. Oliveira, Juan Gomez-Luna, Lois Orosa, Saugata Ghose, Nandita Vijaykumar, Ivan fernandez, Mohammad Sadrosadati, and Onur Mutlu,
 "DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement Bottlenecks"
 IEEE Access, 8 September 2021.
 Preprint in arXiv, 8 May 2021.
 [arXiv preprint]
 [IEEE Access version]
 [DAMOV Suite and Simulator Source Code]
 [SAFARI Live Seminar Video (2 hrs 40 mins)]
 [Short Talk Video (21 minutes)]

DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement Bottlenecks

GERALDO F. OLIVEIRA, ETH Zürich, Switzerland JUAN GÓMEZ-LUNA, ETH Zürich, Switzerland LOIS OROSA, ETH Zürich, Switzerland SAUGATA GHOSE, University of Illinois at Urbana–Champaign, USA NANDITA VIJAYKUMAR, University of Toronto, Canada IVAN FERNANDEZ, University of Malaga, Spain & ETH Zürich, Switzerland MOHAMMAD SADROSADATI, ETH Zürich, Switzerland ONUR MUTLU, ETH Zürich, Switzerland

Adoption: How to Ease **Programmability**? (V)

Appears in IEEE TETC 2023

ALP: Alleviating CPU-Memory Data Movement Overheads in Memory-Centric Systems

Nika Mansouri Ghiasi, Nandita Vijaykumar, Geraldo F. Oliveira, Lois Orosa, Ivan Fernandez, Mohammad Sadrosadati, Konstantinos Kanellopoulos, Nastaran Hajinazar, Juan Gómez Luna, Onur Mutlu

> Abstract—Recent advances in memory technology have enabled near-data processing (NDP) to tackle main memory bottlenecks in modern systems. Prior works partition applications into segments (e.g., instructions, loops, functions) and execute memory-bound segments of the applications on NDP computation units, while mapping the cache-friendly application segments to host CPU cores that access a deeper cache hierarchy. Partitioning applications between NDP and host cores causes inter-segment data movement overhead, which is the overhead from moving data generated from one segment and used in the consecutive segments. This overhead can be large if the segments map to cores in different parts of the system (i.e., host and NDP). Prior works take two approaches to the inter-segment data movement overhead when partitioning applications between NDP and host cores. The first class of works maps segments to NDP or host cores based on the properties of each segment, neglecting the performance impact of the inter-segment data movement. Such partitioning techniques suffer from inter-segment data movement overhead. The second class of works maps segments to host or NDP cores based on the overall memory bandwidth savings of each segment (which depends on the memory bandwidth savings within each segment and the inter-segment data movement overhead between other segments). These works do not offload each segment to the best-fitting core if they incur high inter-segment data movement overhead. Therefore these works miss some of the potential NDP performance benefits. We show that mapping each segment (here basic block) to its best-fitting core based on the properties of each segment, assuming no inter-segment data movement, can provide substantial performance benefits. However, we show that the inter-segment data movement reduces this benefit significantly. To this end, we introduce ALP, a new programmer-transparent technique to leverage the performance benefits of NDP by alleviating the performance impact of inter-segment data movement between host and memory and enabling efficient partitioning of applications

> performance impact of inter-segment data movement between host and memory and enabling efficient partitioning of applications between host and NDP cores. ALP alleviates the inter-segment data movement overhead by *proactively and accurately* transferring the required data between the segments mapped on host and NDP cores. This is based on the key observation that the instructions that generate the inter-segment data stay the same across different executions of a program on different input sets. ALP uses a compiler pass to identify these instructions and uses specialized hardware support to transfer data between the host and NDP cores at runtime. Using both the compiler and runtime information, ALP efficiently maps application segments to either host or NDP cores considering 1) the properties of each segment, 2) the inter-segment data movement overhead between different segments, and 3) whether this inter-segment data movement overhead can be alleviated proactively and in a timely manner. We evaluate ALP across a wide range of workloads and show on average 54.3% and 45.4% speedup compared to executing the application only on the host CPU or only the NDP cores, respectively.

SAFAR

https://arxiv.org/pdf/2212.06292

Adoption: How to Maintain Coherence? (I)

 Amirali Boroumand, Saugata Ghose, Minesh Patel, Hasan Hassan, Brandon Lucia, Kevin Hsieh, Krishna T. Malladi, Hongzhong Zheng, and Onur Mutlu, "LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory" IEEE Computer Architecture Letters (CAL), June 2016.

LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory

Amirali Boroumand[†], Saugata Ghose[†], Minesh Patel[†], Hasan Hassan^{†§}, Brandon Lucia[†], Kevin Hsieh[†], Krishna T. Malladi^{*}, Hongzhong Zheng^{*}, and Onur Mutlu^{‡†} [†]Carnegie Mellon University *Samsung Semiconductor, Inc. [§] TOBB ETÜ [‡]ETH Zürich


Challenge: Coherence for Hybrid CPU-PIM Apps



Adoption: How to Maintain Coherence? (II)

 Amirali Boroumand, Saugata Ghose, Minesh Patel, Hasan Hassan, Brandon Lucia, Kevin Hsieh, Krishna T. Malladi, Hongzhong Zheng, and Onur Mutlu,
"CoNDA: Efficient Cache Coherence Support for Near-Data Accelerators" Proceedings of the <u>46th International Symposium on Computer</u>

Architecture (ISCA), Phoenix, AZ, USA, June 2019.

*Simon Fraser University

CoNDA: Efficient Cache Coherence Support for Near-Data Accelerators

Amirali Boroumand[†]Saugata Ghose[†]Minesh Patel*Hasan Hassan*Brandon Lucia[†]Rachata Ausavarungnirun^{†‡}Kevin Hsieh[†]Nastaran Hajinazar^{◊†}Krishna T. Malladi[§]Hongzhong Zheng[§]Onur Mutlu*[†][†]Carnegie Mellon University*ETH Zürich[‡]KMUTNB

[§]Samsung Semiconductor, Inc.

SAFAR

Adoption: How to Support Synchronization?

 Christina Giannoula, Nandita Vijaykumar, Nikela Papadopoulou, Vasileios Karakostas, Ivan Fernandez, Juan Gómez-Luna, Lois Orosa, Nectarios Koziris, Georgios Goumas, Onur Mutlu, "SynCron: Efficient Synchronization Support for Near-Data-Processing <u>Architectures"</u> *Proceedings of the <u>27th International Symposium on High-Performance Computer</u> <u>Architecture (HPCA)</u>, Virtual, February-March 2021.
[Slides (pptx) (pdf)]
[Short Talk Slides (pptx) (pdf)]
[Talk Video (21 minutes)]
[Short Talk Video (7 minutes)*]

SynCron: Efficient Synchronization Support for Near-Data-Processing Architectures

Christina Giannoula^{†‡} Nandita Vijaykumar^{*‡} Nikela Papadopoulou[†] Vasileios Karakostas[†] Ivan Fernandez^{§‡} Juan Gómez-Luna[‡] Lois Orosa[‡] Nectarios Koziris[†] Georgios Goumas[†] Onur Mutlu[‡] [†]National Technical University of Athens [‡]ETH Zürich ^{*}University of Toronto [§]University of Malaga

SAFARI

Adoption: How to Support Virtual Memory?

 Kevin Hsieh, Samira Khan, Nandita Vijaykumar, Kevin K. Chang, Amirali Boroumand, Saugata Ghose, and Onur Mutlu, <u>"Accelerating Pointer Chasing in 3D-Stacked Memory:</u> <u>Challenges, Mechanisms, Evaluation"</u> *Proceedings of the <u>34th IEEE International Conference on Computer</u> <u>Design</u> (ICCD), Phoenix, AZ, USA, October 2016.*

Accelerating Pointer Chasing in 3D-Stacked Memory: Challenges, Mechanisms, Evaluation

Kevin Hsieh[†] Samira Khan[‡] Nandita Vijaykumar[†] Kevin K. Chang[†] Amirali Boroumand[†] Saugata Ghose[†] Onur Mutlu^{§†} [†]Carnegie Mellon University [‡]University of Virginia [§]ETH Zürich

Adoption: Evaluation Infrastructures

 Haocong Luo, Yahya Can Tugrul, F. Nisa Bostanci, Ataberk Olgun, A. Giray Yaglikci, and Onur Mutlu,
"Ramulator 2.0: A Modern, Modular, and Extensible DRAM Simulator" *Preprint on arxiv*, August 2023.
[arXiv version]
[Ramulator 2.0 Source Code]

Ramulator 2.0: A Modern, Modular, and Extensible DRAM Simulator

Haocong Luo, Yahya Can Tuğrul, F. Nisa Bostancı, Ataberk Olgun, A. Giray Yağlıkçı, and Onur Mutlu

https://arxiv.org/pdf/2308.11030.pdf

SAFARI https://github.com/CMU-SAFARI/ramulator2

Processing-in-Memory: Challenges

To fully support PIM systems, we need to develop:

- **1** Workload characterization methodologies and benchmark suites targeting PIM architectures
- 2 Frameworks that can facilitate the implementation of complex operations and algorithms using PIM primitives
- 3 Compiler support and compiler optimizations targeting PIM architectures
- **4 Operating system support for PIM-aware virtual memory, memory management, data allocation and mapping**
- **5** End-to-End System-on-Chip Design Beyond DRAM

The <u>lack of tools</u> and <u>system support</u> for PIM architectures limit the <u>adoption</u> of PIM systems

SAFARI

An Example: SimplePIM Framework

 Jinfan Chen, Juan Gómez-Luna, Izzat El Hajj, YuXin Guo, and Onur Mutlu,
"SimplePIM: A Software Framework for Productive and Efficient Processing in Memory"
Proceedings of the <u>32nd International Conference on</u> Parallel Architectures and Compilation Techniques (PACT), Vienna, Austria, October 2023.

SimplePIM: A Software Framework for Productive and Efficient Processing-in-Memory

Jinfan Chen¹ Juan Gómez-Luna¹ Izzat El Hajj² Yuxin Guo¹ Onur Mutlu¹ ¹ETH Zürich ²American University of Beirut

SAFARI

Executive Summary

- Real PIM hardware is now available, e.g., UPMEM PIM
- However, programming real PIM hardware is challenging, e.g., need to:
 - Distribute data across PIM memory banks,
 - Manage data transfers between host cores and PIM cores, between PIM cores, and between DRAM bank and PIM scratchpad
 - Launch PIM kernels on the PIM cores, etc.
 - Synchronize properly between threads
- SimplePIM is a high-level programming framework for real PIM hardware
 - Iterators such as map, reduce, and zip
 - Collective communication with broadcast, scatter, and gather
- Implementation on UPMEM and evaluation with six different workloads
 - Reduction, vector add, histogram, linear/logistic regression, K-means
 - 4.4x fewer lines of code compared to hand-optimized code
 - Between 15% and 43% faster than hand-optimized code for three workloads
- Source code: <u>https://github.com/CMU-SAFARI/SimplePIM</u>

Concluding Remarks

Challenge and Opportunity for Future

Fundamentally **Energy-Efficient** (Data-Centric) **Computing Architectures** Challenge and Opportunity for Future

Fundamentally **High-Performance** (Data-Centric) **Computing Architectures** Challenge and Opportunity for Future

Computing Architectures with

Minimal Data Movement



Concluding Remarks

- It is time to design principled system architectures to solve the memory problem
- We must design systems to be balanced, high-performance, and energy-efficient → memory-centric
 - Enable computation capabilities in memory
- This can
 - Lead to orders-of-magnitude improvements
 - Enable new applications & computing platforms
 - **Enable better understanding of nature**
 - ••••

Future of truly memory-centric computing is bright We need to do research & design across the computing stack

Fundamentally Better Architectures

Data-centric

Data-driven

Data-aware



We Need to Revisit the Entire Stack

• With a **memory-centric mindset**

Problem	
Aigorithm	
Program/Language	
System Software	I
SW/HW Interface	I
Micro-architecture	I
Logic	J
Devices	
Electrons	

We can get there step by step

SAFARI

A Modern Primer on Processing in Memory

Onur Mutlu^{a,b}, Saugata Ghose^{b,c}, Juan Gómez-Luna^a, Rachata Ausavarungnirun^d

SAFARI Research Group

^aETH Zürich ^bCarnegie Mellon University ^cUniversity of Illinois at Urbana-Champaign ^dKing Mongkut's University of Technology North Bangkok

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna, and Rachata Ausavarungnirun, "A Modern Primer on Processing in Memory" *Invited Book Chapter in <u>Emerging Computing: From Devices to Systems -</u> <i>Looking Beyond Moore and Von Neumann*, Springer, to be published in 2021.

SAFARI

https://arxiv.org/pdf/2012.03112.pdf

Referenced Papers, Talks, Artifacts

All are available at

https://people.inf.ethz.ch/omutlu/projects.htm

https://www.youtube.com/onurmutlulectures

https://github.com/CMU-SAFARI/

Acknowledgments

SAFARI Research Group safari.ethz.ch



SAFARI Research Group

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

https://safari.ethz.ch/safari-newsletter-january-2021/



https://safari.ethz.ch

SAFARI Research Group: December 2021

<u>https://safari.ethz.ch/safari-newsletter-december-2021/</u>



Think Big, Aim High



f y in 🖸

View in your browser December 2021



SAFARI Newsletter June 2023 Edition

<u>https://safari.ethz.ch/safari-newsletter-june-2023/</u>



Think Big, Aim High



y in D

View in your browser June 2023



SAFARI Newsletter July 2024 Edition

<u>https://safari.ethz.ch/safari-newsletter-july-2024/</u>



SAFARI Introduction & Research

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



Seminar in Computer Architecture - Lecture 5: Potpourri of Research Topics (Spring 2023)

Think BIG, Aim HIGH!



https://www.youtube.com/watch?v=mV2OuB2djEs

Open Source Tools: SAFARI GitHub



https://github.com/CMU-SAFARI/

ML/AI for Memory System Design & Memory System Design for AI/ML

Onur Mutlu <u>omutlu@gmail.com</u> <u>https://people.inf.ethz.ch/omutlu</u>

23 August 2024 Intel ArchFest





Backup Slides

Processing-in-Memory: Challenges

To fully support PIM systems, we need to develop:

- **1** Workload characterization methodologies and benchmark suites targeting PIM architectures
- 2 Frameworks that can facilitate the implementation of complex operations and algorithms using PIM primitives
- 3 Compiler support and compiler optimizations targeting PIM architectures
- **4 Operating system support for PIM-aware virtual memory, memory management, data allocation and mapping**
- **5** End-to-End System-on-Chip Design Beyond DRAM

The <u>lack of tools</u> and <u>system support</u> for PIM architectures limit the <u>adoption</u> of PIM systems

SAFARI

Security Issues in Processing in Memory

- Does PIM make security better or easier?
- Does PIM make security worse?
- Many interesting questions here
- Topic of a separate talk, but we highlight some papers
 - Evaluating Homomorphic Operations on a Real-World Processing-In-Memory System [IISWC 2023]
 - Amplifying Main Memory-Based Timing Covert and Side Channels using Processing-in-Memory Operations [arxiv 2024]

MIMDRAM

MIMDRAM: More Flexible Processing using DRAM

Appears at HPCA 2024 <u>https://arxiv.org/pdf/2402.19080.pdf</u>

MIMDRAM: An End-to-End Processing-Using-DRAM System for High-Throughput, Energy-Efficient and Programmer-Transparent Multiple-Instruction Multiple-Data Computing

Geraldo F. Oliveira[†]Ataberk Olgun[†]Abdullah Giray Yağlıkçı[†]F. Nisa Bostancı[†]Juan Gómez-Luna[†]Saugata Ghose[‡]Onur Mutlu[†]

† ETH Zürich ‡ Univ. of Illinois Urbana-Champaign

Our goal is to design a flexible PUD system that overcomes the limitations caused by the large and rigid granularity of PUD. To this end, we propose MIMDRAM, a hardware/software co-designed PUD system that introduces new mechanisms to allocate and control only the necessary resources for a given PUD operation. The key idea of MIMDRAM is to leverage finegrained DRAM (i.e., the ability to independently access smaller segments of a large DRAM row) for PUD computation. MIM-DRAM exploits this key idea to enable a multiple-instruction multiple-data (MIMD) execution model in each DRAM subarray (and SIMD execution within each DRAM row segment).

SAFARI

Processing-in-Memory: Challenges

To fully support PIM systems, we need to develop:

- **1** Workload characterization methodologies and benchmark suites targeting PIM architectures
- **2** Frameworks that can facilitate the implementation of complex operations and algorithms using PIM primitives
- **3** Compiler support and compiler optimizations targeting PIM architectures
- **4** Operating system support for PIM-aware virtual memory, memory management, data allocation and mapping
- 5 End-to-End System-on-Chip Design Beyond DRAM

The <u>lack of tools</u> and <u>system support</u> for PIM architectures limit the <u>adoption</u> of PIM systems

SAFARI

MIMDRAM

An End-to-End Processing-Using-DRAM System for High-Throughput, Energy-Efficient and Programmer-Transparent Multiple-Instruction Multiple-Data Computing

Geraldo F. Oliveira

Ataberk Olgun

Saugata Ghose

ETH zürich

A. Giray Yağlıkçı

Juan Gómez-Luna

F. Nisa Bostancı

Onur Mutlu





Executive Summary

Problem: Processing-Using-DRAM (PUD) suffers from three issues caused by DRAM's large and rigid access granularity

- <u>Underutilization</u> due to data parallelism variation in (and across) applications
- <u>Limited computation support</u> due to a lack of interconnects
- <u>Challenging programming</u> model due to a lack of compilers

Goal: Design a flexible PUD system that overcomes the three limitations caused by DRAM's large and rigid access granularity

Key Mechanism: MIMDRAM, a hardware/software co-design PUD system

- Key idea: leverage fine-grained DRAM for PUD operation
- **HW**: <u>simple changes</u> to the DRAM array, enabling concurrent PUD operations
 - low-cost interconnects at the DRAM peripherals for data reduction
- SW: <u>compiler</u> and <u>OS</u> support to generate and map PUD instructions

Key Results: MIMDRAM achieves

- **14.3***x*, **30.6***x*, and **6.8***x* the energy efficiency of state-of-the-art PUD systems, a high-end CPU and GPU, respectively
- Small area cost to a DRAM chip (1.11%) and CPU die (0.6%)

SAFARI

https://github.com/CMU-SAFARI/MIMDRAM

Recall: Processing using DRAM

- We can support
 - Bulk bitwise AND, OR, NOT, MAJ
 - Bulk bitwise COPY and INIT/ZERO
 - True Random Number Generation; Physical Unclonable Functions
 - Lookup Table based more complex computation
- At low cost
- Using analog computation capability of DRAM
 - Idea: activating (multiple) rows performs computation
 - Even in commodity off-the-shelf DRAM chips!

30-77X performance and energy improvement

- Seshadri+, "Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology," MICRO 2017.
- Seshadri+"RowClone: Fast and Efficient In-DRAM Copy and Initialization of Bulk Data," MICRO 2013.

Background Work: RowClone

 Vivek Seshadri, Yoongu Kim, Chris Fallin, Donghyuk Lee, Rachata Ausavarungnirun, Gennady Pekhimenko, Yixin Luo, Onur Mutlu, Michael A. Kozuch, Phillip B. Gibbons, and Todd C. Mowry,
"RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization"

Proceedings of the <u>46th International Symposium on Microarchitecture</u> (**MICRO**), Davis, CA, December 2013. [<u>Slides (pptx) (pdf)</u>] [<u>Lightning Session</u> <u>Slides (pptx) (pdf)</u>] [<u>Poster (pptx) (pdf)</u>]

RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization

Vivek Seshadri Yoongu Kim Chris Fallin^{*} Donghyuk Lee vseshadr@cs.cmu.edu yoongukim@cmu.edu cfallin@c1f.net donghyuk1@cmu.edu Rachata Ausavarungnirun Gennady Pekhimenko Yixin Luo rachata@cmu.edu gpekhime@cs.cmu.edu yixinluo@andrew.cmu.edu Onur Mutlu Phillip B. Gibbons[†] Michael A. Kozuch[†] Todd C. Mowry onur@cmu.edu phillip.b.gibbons@intel.com michael.a.kozuch@intel.com tcm@cs.cmu.edu Carnegie Mellon University [†]Intel Pittsburgh

Background Work: PiDRAM

 Ataberk Olgun, Juan Gomez Luna, Konstantinos Kanellopoulos, Behzad Salami, Hasan Hassan, Oguz Ergin, and Onur Mutlu, "PiDRAM: A Holistic End-to-end FPGA-based Framework for Processing-in-DRAM" ACM Transactions on Architecture and Code Optimization (TACO), March 2023.
[arXiv version] Presented at the 18th HiPEAC Conference, Toulouse, France, January 2023.
[Slides (pptx) (pdf)]
[Longer Lecture Slides (pptx) (pdf)]
[Lecture Video (40 minutes)]
[PiDRAM Source Code]

PiDRAM: A Holistic End-to-end FPGA-based Framework for <u>Processing-in-DRAM</u>

Ataberk Olgun§

Juan Gómez Luna
§ Konstantinos Kanellopoulos
§ Hasan Hassan
§ Oğuz Ergin
† Onur Mutlu
§

[§]ETH Zürich [†]TOBB University of Economics and Technology

Behzad Salami[§]
Background Work: In-DRAM Bulk AND/OR

 Vivek Seshadri, Kevin Hsieh, Amirali Boroumand, Donghyuk Lee, Michael A. Kozuch, Onur Mutlu, Phillip B. Gibbons, and Todd C. Mowry,
 <u>"Fast Bulk Bitwise AND and OR in DRAM"</u> <u>IEEE Computer Architecture Letters</u> (CAL), April 2015.

Fast Bulk Bitwise AND and OR in DRAM

Vivek Seshadri*, Kevin Hsieh*, Amirali Boroumand*, Donghyuk Lee*, Michael A. Kozuch[†], Onur Mutlu*, Phillip B. Gibbons[†], Todd C. Mowry* *Carnegie Mellon University [†]Intel Pittsburgh

Background Work: Ambit

 Vivek Seshadri, Donghyuk Lee, Thomas Mullins, Hasan Hassan, Amirali Boroumand, Jeremie Kim, Michael A. Kozuch, Onur Mutlu, Phillip B. Gibbons, and Todd C. Mowry,
 <u>"Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using</u> <u>Commodity DRAM Technology"</u> *Proceedings of the <u>50th International Symposium on</u> <u>Microarchitecture</u> (<i>MICRO*), Boston, MA, USA, October 2017. [Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)] [Poster (pptx) (pdf)]

Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology

Vivek Seshadri^{1,5} Donghyuk Lee^{2,5} Thomas Mullins^{3,5} Hasan Hassan⁴ Amirali Boroumand⁵ Jeremie Kim^{4,5} Michael A. Kozuch³ Onur Mutlu^{4,5} Phillip B. Gibbons⁵ Todd C. Mowry⁵

¹Microsoft Research India ²NVIDIA Research ³Intel ⁴ETH Zürich ⁵Carnegie Mellon University

Background: In-DRAM Bulk Bitwise Execution

 Vivek Seshadri and Onur Mutlu,
 "In-DRAM Bulk Bitwise Execution Engine" Invited Book Chapter in Advances in Computers, to appear in 2020.
 [Preliminary arXiv version]

In-DRAM Bulk Bitwise Execution Engine

Vivek Seshadri Microsoft Research India visesha@microsoft.com Onur Mutlu ETH Zürich onur.mutlu@inf.ethz.ch

Recall: SIMDRAM Framework

 Nastaran Hajinazar, Geraldo F. Oliveira, Sven Gregorio, Joao Dinis Ferreira, Nika Mansouri Ghiasi, Minesh Patel, Mohammed Alser, Saugata Ghose, Juan Gomez-Luna, and Onur Mutlu, "SIMDRAM: An End-to-End Framework for Bit-Serial SIMD Computing in DRAM" Proceedings of the 26th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), Virtual, March-April 2021.
 [2-page Extended Abstract]
 [Short Talk Slides (pptx) (pdf)]
 [Talk Slides (pptx) (pdf)]
 [Short Talk Video (5 mins)]
 [Full Talk Video (27 mins)]

SIMDRAM: A Framework for Bit-Serial SIMD Processing using DRAM

*Nastaran Hajinazar^{1,2} Nika Mansouri Ghiasi¹ Juan Gómez-Luna¹ Sven Gregorio¹ Mohammed Alser¹ Onur Mutlu¹

João Dinis Ferreira¹ Saugata Ghose³

¹ETH Zürich

²Simon Fraser University

³University of Illinois at Urbana–Champaign

Background:

In-DRAM Copy/Init, Majority & NOT Operations

In-DRAM majority is performed by simultaneously activating three DRAM rows



Seshadri, Vivek, et al. " Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology," in MICRO, 2017



Introduction & Background

Limitations of PUD

MIMDRAM

Hardware Overview

Software Support

....

Background: In-DRAM Majority Operations



Processing-Using-DRAM architectures (e.g., SIMDRAM) are very-wide (e.g., 65,536 wide) bit-serial SIMD engines



Oliveira, Geraldo F., et al. " SIMDRAM: An End-to-End Framework for Bit-Serial SIMD Computing in DRAM," in ASPLOS, 2021

Introduction & Background

Limitations of PUD

MIMDRAM

Hardware Overview

Software Support

.

Limitations of PUD Systems: Overview

PUD systems suffer from three sources of inefficiency due to the large and rigid DRAM access granularity

SIMD Underutilization

- due to data parallelism variation within and across applications
- leads to throughput and energy waste

Limited Computation Support

- due to a lack of low-cost interconnects across columns
- limits PUD operations to only parallel map constructs

Challenging Programming Model

- due to a lack of compiler support for PUD systems
- creates a burden on programmers, limiting PUD adoption

Limitations of PUD

Hardware Overvie w

MIMDRAM

. . . .

Programmer's Tasks:	Goal:
Map & align	Just write
data structures	my kernel

High-level code for C[i] = (A[i] > pred[i])? A[i] + B[i] : A[i] - B[i]



Limitations of PUD

MIMDRAM

Hardware Overview

Software Support

•

Programmer's Tasks:		Goal:	
Map & align	Identify		Just write
data structures	array boundaries		my kernel

High-level code for C[i] = (A[i] > pred[i])? A[i] + B[i] : A[i] - B[i]



Limitations of PUD

MIMDRAM

Hardware Overview

Software Support

•

Programmer's Tasks:				Goal:
Map & align	Identify	Manually	Map C to	Just write
data structures	array boundaries	unroll loop	PUD instructions	my kernel

High-level code for C[i] = (A[i] > pred[i])? A[i] + B[i] : A[i] - B[i]

```
for (int i = 0; i < size ; ++ i){
    bool cond = A[i] > pred[i];
    if (cond) C[i] = A[i] + B[i];
    else C[i] = A[i] - B[i];
}
```



Limitations of PUD

MIMDRAM

Hardware Overview

ew Softwa

Software Support Evaluation

Co

Conclusion 298

Programmer's Tasks: Go				Goal:	
Map & align	Identify	Manually	Map C to	Orchestrate	Just write
data structures	array boundaries	unroll loop	PUD instructions	data movement	my kernel

High-level code for C[i] = (A[i] > pred[i])? A[i] + B[i] : A[i] - B[i]



Limitations of PUD

MIMDRAM

Hardware Overview

Software Support

Programmer's Tasks:					Goal:
Map & align data structures	Identify array boundaries	Manually unroll loop	Map C to PUD instructions	Orchestrate data movement	Just write my kernel
	P C[i] = (U D's assem A[i] > pred[i]	bly-like code fo)? A[i] + B[i] : A[i]	o r — B[i]	
	<pre>bbop_trsp_ini bbop_trsp_ini bbop_trsp_ini</pre>	t(A , si t(B , si t(C , si	ze , elm_si ze , elm_si ze , elm_si	ze); ze); ze);	
	<pre>bbop_add(D , bbop_sub(E ,</pre>	A, B, A, B,	size , elm_s size , elm_s	size); size);	

SAFARI

Limitations of PUD • • • • • • • •

MIMDRAM

Hardware Overview

Software Support Evaluation

.

....

300

Problem & Goal



SAFARI

Introduction & Background

Limitations of PUD

MIMDRAM

Hardware Overview

.....

Software Support

Evaluation

Conclusion 301

DRAM's hierarchical organization can enable <u>fine-grained access</u>



Fine-Grained DRAM:

segments the global wordline to access individual DRAM mats

Fine-Grained DRAM:

segments the global wordline to access individual DRAM mats



global sense amplifier

Fine-grained DRAM for energy-efficient DRAM access:

[Cooper-Balis+, 2010]: Fine-Grained Activation for Power Reduction in DRAM [Udipi+, 2010]: Rethinking DRAM Design and Organization for Energy-Constrained Multi-Cores [Zhang+, 2014]: Half-DRAM [Ha+, 2016]: Improving Energy Efficiency of DRAM by Exploiting Half Page Row Access

[Hat, 2010]. Improving Energy Efficiency of DRAW by Exploiting Han

[O'Connor+, 2017]: Fine-Grained DRAM

[Olgun+, 2024]: Sectored DRAM

Sectored DRAM

 Ataberk Olgun, F. Nisa Bostanci, Geraldo F. Oliveira, Yahya Can Tugrul, Rahul Bera, A. Giray Yaglikci, Hasan Hassan, Oguz Ergin, and Onur Mutlu, "Sectored DRAM: A Practical Energy-Efficient and High-Performance Fine-Grained DRAM Architecture" ACM Transactions on Architecture and Code Optimization (TACO), [online] June 2024.
 [arXiv version] [ACM Digital Library version]

Sectored DRAM: A Practical Energy-Efficient and High-Performance Fine-Grained DRAM Architecture

Ataberk Olgun[§] F. Nisa Bostancı^{§†} Geraldo F. Oliveira[§] Yahya Can Tuğrul^{§†} Rahul Bera[§] A. Giray Yağlıkcı[§] Hasan Hassan[§] Oğuz Ergin[†] Onur Mutlu[§]

SAFARI

https://arxiv.org/pdf/2207.13795



Fine-grained DRAM for processing-using-DRAM:

1 Improves SIMD utilization

for a single PUD operation, only access the DRAM mats with target data



Fine-grained DRAM for processing-using-DRAM:

1 Improves SIMD utilization

- for a single PUD operation, only access the DRAM mats with target data
- for multiple PUD operations, execute independent operations concurrently
 → multiple instruction, multiple data (MIMD) execution model

segmented global wordline



global sense amplifier

Fine-grained DRAM for processing-using-DRAM:

1

mproves SIMD utilization

for a single PUD operation, only access the DRAM mats with target data

for multiple PUD operations, execute independent operations concurrently → multiple instruction, multiple data (MIMD) execution model

7 Enables low-cost interconnects for vector reduction

- global and local data buses can be used for inter-/intra-mat communication



Fine-grained DRAM for processing-using-DRAM:

Improves SIMD

- for a single PUD operation, only access the DRAM mats with target data
- for multiple PUD operations, execute independent operations concurrentl
 → multiple instruction, multiple data (MIMD) execution model
- **7** Enables low-cost interconnects for vector reduction
 - global and local data buses can be used for inter-/intra-mat communication

3 Eases programmability

- SIMD parallelism in a DRAM mat is on par with vector ISAs' SIMD width

MIMDRAM: Compiler Support (I)

Transparently: <u>extract</u> SIMD parallelism from an application, and <u>schedule</u> PUD instructions while maximizing <u>utilization</u>

Three new LLVM-based passes targeting PUD execution



MIMDRAM: Compiler Support (II)



Identify SIMD parallelism, generate PUD instructions, and set the appropriate vectorization factor

MIMDRAM: Compiler Support (II)



Identify SIMD parallelism, generate PUD instructions, and set the appropriate vectorization factor

Homology SIMD utilization by allowing the distribution of independent PUD instructions across DRAM mats

MIMDRAM: Compiler Support (III)



Identify SIMD parallelism, generate PUD instructions, and set the appropriate vectorization factor

त्नmprove SIMD utilization by allowing the distribution of independent PUD े instructions across DRAM mats

Generate the appropriate binary for data allocation and PUD instructions

MIMDRAM: System Support

- Instruction set architecture
- Execution & data transposition
- Data coherence
- Address translation
- Data allocation & alignment
- Mat label translation

Evaluation: Methodology Overview

Evaluation Setup

- CPU: Intel Skylake CPU
- GPU: NVIDIA A100 GPU
- PUD: SIMDRAM [Oliveira+, 2021] and DRISA [Li+, 2017]
- PND: Fulcrum [Lenjani+, 2020]
- https://github.com/CMU-SAFARI/MIMDRAM

• Workloads:

- 12 workloads from Polybench, Rodinia, Phoenix, and SPEC2017
- 495 multi-programmed application mixes
- Two-Level Analysis
 - Single application \rightarrow leverages intra-application data parallelism
 - Multi-programmed workload → leverages inter-application

data parallelism

Evaluation:

Single Application Analysis – Energy Efficiency



MIMDRAM significantly improves energy efficiency compared to CPU (30.6x), GPU (6.8x), and SIMDRAM (14.3x)

Evaluation:

Multi-Programmed Workload Analysis



Evaluation: More in the Paper

- MIMDRAM with subarray and bank-level parallelism
 - MIMDRAM provides significant performance gains compared to the baseline CPU (13.2x) and GPU (2x)
- Comparison to DRISA and Fulcrum for multi-programmed workloads
 - MIMDRAM achieves system throughput on par with DRISA and Fulcrum
- MIMDRAM's SIMD utilization versus SIMDRAM
 - MIMDRAM provides **15.6x** the utilization of SIMDRAM
- Area analysis
 - MIMDRAM adds small area cost to a DRAM chip (1.11%) and CPU die (0.6%)



MIMDRAM: Summary

We introduced MIMDRAM,

a hardware/software co-designed processing-using-DRAM system

- Key idea: leverage fine-grained DRAM for processing-using-DRAM operation
- HW: simple changes to DRAM, enabling concurrent instruction execution
 - low-cost interconnects at the DRAM peripherals for data reduction
- SW: <u>compiler</u> and <u>OS</u> support to generate and map instructions

Our evaluation demonstrates that MIMDRAM

- **significantly** improves **performance**, **energy efficiency**, and **throughput** compared to processor-centric (CPU and GPU) and memory-centric (SIMDRAM, DRISA, and Fulcrum) architectures
- incurs small area cost to a DRAM chip and CPU die

https://github.com/CMU-SAFARI/MIMDRAM

SAFARI

Introduction & Background

Limitations of PUD Systems

MIMDRAM Overview

Evaluation

Two Other Works on PIM Programmability

Adoption: How to Ease **Programmability**? (I)

 Geraldo F. Oliveira, Alain Kohli, David Novo, Juan Gómez-Luna, Onur Mutlu,
 "DaPPA: A Data-Parallel Framework for Processingin-Memory Architectures,"
 in PACT SRC Student Competition, Vienna, Austria, October 2023.

DaPPA: A Data-Parallel Framework for Processing-in-Memory Architectures

Geraldo F. Oliveira*

Alain Kohli*

David Novo[‡] Juan Gómez-Luna^{*} Onur Mutlu^{*} [‡]LIRMM, Univ. Montpellier, CNRS

Adoption: How to Ease Programmability? (II)

 Jinfan Chen, Juan Gómez-Luna, Izzat El Hajj, YuXin Guo, and Onur Mutlu,
 "SimplePIM: A Software Framework for Productive and Efficient Processing in Memory"
 Proceedings of the <u>32nd International Conference on</u> Parallel Architectures and Compilation Techniques (PACT), Vienna, Austria, October 2023.

SimplePIM: A Software Framework for Productive and Efficient Processing-in-Memory

Jinfan Chen¹ Juan Gómez-Luna¹ Izzat El Hajj² Yuxin Guo¹ Onur Mutlu¹ ¹ETH Zürich ²American University of Beirut



Adoption: How to Ease Programmability? (II)

 Jinfan Chen, Juan Gómez-Luna, Izzat El Hajj, YuXin Guo, and Onur Mutlu,
 "SimplePIM: A Software Framework for Productive and Efficient Processing in Memory"
 Proceedings of the <u>32nd International Conference on</u> Parallel Architectures and Compilation Techniques (PACT), Vienna, Austria, October 2023.

SimplePIM: A Software Framework for Productive and Efficient Processing-in-Memory

Jinfan Chen¹ Juan Gómez-Luna¹ Izzat El Hajj² Yuxin Guo¹ Onur Mutlu¹ ¹ETH Zürich ²American University of Beirut

SimplePIM: A Software Framework for Productive and Efficient Processing-in-Memory

Jinfan Chen, Juan Gómez Luna, Izzat El Hajj, Yuxin Guo, Onur Mutlu

https://arxiv.org/pdf/2310.01893.pdf https://github.com/CMU-SAFARI/SimplePIM




Executive Summary

- Processing-in-Memory (PIM) promises to alleviate the data movement bottleneck
- Real PIM hardware is now available, e.g., UPMEM PIM
- However, programming real PIM hardware is challenging, e.g.:
 - Distribute data across PIM memory banks,
 - Manage data transfers between host cores and PIM cores, and between PIM cores,
 - Launch PIM kernels on the PIM cores, etc.
- SimplePIM is a high-level programming framework for real PIM hardware
 - Iterators such as map, reduce, and zip
 - Collective communication with broadcast, scatter, and gather
- Implementation on UPMEM and evaluation with six different workloads
 - Reduction, vector add, histogram, linear/logistic regression, K-means
 - 4.4x fewer lines of code compared to hand-optimized code
 - Between 15% and 43% faster than hand-optimized code for three workloads
- Source code: <u>https://github.com/CMU-SAFARI/SimplePIM</u>

Real DRAM Chips Are Already Quite Capable: FC-DRAM & SiMRA

DRAM Chips Are Already (Quite) Capable!

Appears at HPCA 2024 <u>https://arxiv.org/pdf/2402.18736.pdf</u>

Functionally-Complete Boolean Logic in Real DRAM Chips: Experimental Characterization and Analysis

İsmail Emir Yüksel Yahya Can Tuğrul Ataberk Olgun F. Nisa Bostancı A. Giray Yağlıkçı Geraldo F. Oliveira Haocong Luo Juan Gómez-Luna Mohammad Sadrosadati Onur Mutlu

ETH Zürich

We experimentally demonstrate that COTS DRAM chips are capable of performing 1) functionally-complete Boolean operations: NOT, NAND, and NOR and 2) many-input (i.e., more than two-input) AND and OR operations. We present an extensive characterization of new bulk bitwise operations in 256 off-theshelf modern DDR4 DRAM chips. We evaluate the reliability of these operations using a metric called success rate: the fraction of correctly performed bitwise operations. Among our 19 new observations, we highlight four major results. First, we can perform the NOT operation on COTS DRAM chips with 98.37% success rate on average. Second, we can perform up to 16-input NAND, NOR, AND, and OR operations on COTS DRAM chips with high reliability (e.g., 16-input NAND, NOR, AND, and OR with average success rate of 94.94%, 95.87%, 94.94%, and 95.85%, respectively). Third, data pattern only slightly

DRAM Chips Are Already (Quite) Capable!

https://arxiv.org/pdf/2312.02880.pdf

PULSAR: Simultaneous Many-Row Activation for Reliable and High-Performance Computing in Off-the-Shelf DRAM Chips

Ismail Emir Yuksel Yahya Can Tugrul F. Nisa Bostanci Abdullah Giray Yaglikci Ataberk Olgun Geraldo F. Oliveira Melina Soysal Haocong Luo Juan Gomez Luna Mohammad Sadrosadati Onur Mutlu

ETH Zurich

We propose PULSAR, a new technique to enable highsuccess-rate and high-performance PuM operations in off-theshelf DRAM chips. PULSAR leverages our new observation that a carefully-crafted sequence of DRAM commands simultaneously activates up to 32 DRAM rows. PULSAR overcomes the limitations of existing techniques by 1) replicating the input data to improve the success rate and 2) enabling new bulk bitwise operations (e.g., many-input majority, *Multi-RowInit*, and *Bulk-Write*) to improve the performance.

DRAM Chips Are Already (Quite) Capable!

Appears at DSN 2024



Simultaneous Many-Row Activation in Off-the-Shelf DRAM Chips: Experimental Characterization and Analysis

İsmail Emir Yüksel¹ Yahya Can Tuğrul^{1,2} F. Nisa Bostancı¹ Geraldo F. Oliveira¹ A. Giray Yağlıkçı¹ Ataberk Olgun¹ Melina Soysal¹ Haocong Luo¹ Juan Gómez-Luna¹ Mohammad Sadrosadati¹ Onur Mutlu¹ ¹ETH Zürich ²TOBB University of Economics and Technology

The Capability of COTS DRAM Chips

We **demonstrate** that **COTS DRAM chips**:

Can simultaneously activate up to48 rows in two neighboring subarrays

Can perform **NOT operation** with up to **32 output operands**

Can perform up to **16-input** AND, NAND, OR, and NOR operations

2

3

Finding: SiMRA Across Subarrays

Activating two rows in **quick succession** can **simultaneously** activate **multiple rows in neighboring subarrays**



Key Idea: NOT Operation

Connect rows in neighboring subarrays through **a NOT gate** by simultaneously activating rows



Key Idea: NAND, NOR, AND, OR

Manipulate the bitline voltage to express a wide variety of functions using

multiple-row activation in neighboring subarrays







SAFARI *Gao et al., "FracDRAM: Fractional Values in Off-the-Shelf DRAM," in MICRO, 2022. 334





 $V_{DD} = 1 \& GND = 0$









 $V_{DD} = 1 \& GND = 0$







 V_{DD} =1 & GND = 0







 $V_{DD}=1 \& GND = 0$







 $V_{DD} = 1 \& GND = 0$ COM **REF** Х $\left(\right)$ 0 1 $\mathbf{0}$ 1 \mathbf{O} $\left(\right)$ 1 1 1 $\mathbf{0}$ 1 AND NAND

339

Many-Input AND, NAND, OR, and NOR Operations

We can express AND, NAND, OR, and NOR operations by carefully manipulating the **reference voltage**

Functionally-Complete Boolean Logic in Real DRAM Chips: Experimental Characterization and Analysis

İsmail Emir Yüksel Yahya Can Tuğrul Ataberk Olgun F. Nisa Bostancı A. Giray Yağlıkçı Geraldo F. Oliveira Haocong Luo Juan Gómez-Luna Mohammad Sadrosadati Onur Mutlu

ETH Zürich

(More details in the paper)

https://arxiv.org/pdf/2402.18736.pdf

DRAM Testing Infrastructure

- Developed from DRAM Bender [Olgun+, TCAD'23]*
- Fine-grained control over DRAM commands, timings, and temperature



SAFAR *Olgun et al., "<u>DRAM Bender: An Extensible and Versatile FPGA-based Infrastructure</u> to Easily Test State-of-the-art DRAM Chips," TCAD, 2023.

341

DRAM Chips Tested

- 256 DDR4 chips from two major DRAM manufacturers
- Covers different die revisions and chip densities

Chip Mfr.	#Modules (#Chips)	Die Rev.	Mfr. Date ^a	Chip Density	Chip Org.	Speed Rate
SK Hynix	9 (72)	М	N/A	4Gb	x8	2666MT/s
	5 (40)	А	N/A	4Gb	x8	2133MT/s
	1 (16)	А	N/A	8Gb	x8	2666MT/s
	1 (32)	А	18-14	4Gb	x4	2400MT/s
	1 (32)	А	16-49	8Gb	x4	2400MT/s
	1 (32)	М	16-22	8Gb	x4	2666MT/s
Samsung	1 (8)	F	21-02	4Gb	x8	2666MT/s
	2 (16)	D	21-10	8Gb	x8	2133MT/s
	1 (8)	А	22-12	8Gb	x8	3200MT/s

Performing AND, NAND, OR, and NOR



COTS DRAM chips can perform {2, 4, 8, 16}-input AND, NAND, OR, and NOR operations

Performing AND, NAND, OR, and NOR



COTS DRAM chips can perform 16-input AND, NAND, OR, and NOR operations with very high success rate (>94%)











Data pattern slightly affects

the reliability of AND, NAND, OR, and NOR operations

More in the Paper

- Detailed hypotheses & key ideas to perform
 - NOT operation
 - Many-input AND, NAND, OR, and NOR operations
- How the reliability of bitwise operations are affected by
 - The location of activated rows
 - Temperature (for AND, NAND, OR, and NOR)
 - DRAM speed rate
 - Chip density and die revision
- Discussion on the limitations of COTS DRAM chips

Available on arXiv

Functionally-Complete Boolean Logic in Real DRAM Chips: Experimental Characterization and Analysis

İsmail Emir Yüksel Yahya Can Tuğrul Ataberk Olgun F. Nisa Bostancı A. Giray Yağlıkçı Geraldo F. Oliveira Haocong Luo Juan Gómez-Luna Mohammad Sadrosadati Onur Mutlu

ETH Zürich

Processing-using-DRAM (PuD) is an emerging paradigm that leverages the analog operational properties of DRAM circuitry to enable massively parallel in-DRAM computation. PuD has the potential to significantly reduce or eliminate costly data movement between processing elements and main memory. A common approach for PuD architectures is to make use of bulk bitwise computation (e.g., AND, OR, NOT). Prior works experimentally demonstrate three-input MAJ (i.e., MAJ3) and two-input AND and OR operations in commercial off-the-shelf (COTS) DRAM chips. Yet, demonstrations on COTS DRAM chips do not provide a functionally complete set of operations (e.g., NAND or AND and NOT).

We experimentally demonstrate that COTS DRAM chips are capable of performing 1) functionally-complete Boolean operations: NOT, NAND, and NOR and 2) many-input (i.e., more than two-input) AND and OR operations. We present an extensive systems and applications [12, 13]. Processing-using-DRAM (PuD) [29–32] is a promising paradigm that can alleviate the data movement bottleneck. PuD uses the analog operational properties of the DRAM circuitry to enable massively parallel in-DRAM computation. Many prior works [29–53] demonstrate that PuD can greatly reduce or eliminate data movement.

A widely used approach for PuD is to perform bulk bitwise operations, i.e., bitwise operations on large bit vectors. To perform bulk bitwise operations using DRAM, prior works propose modifications to the DRAM circuitry [29–31, 33, 35, 36, 43, 44, 46, 48–58]. Recent works [38, 41, 42, 45] experimentally demonstrate the feasibility of executing data copy & initialization [42, 45], i.e., the RowClone operation [49], and a subset of bitwise operations, i.e., three-input bitwise majority (MAJ3) and two-input AND and OR operations in unmodified commercial off-the-shelf (COTS) DRAM chips by operating beyond

https://arxiv.org/pdf/2402.18736.pdf

Summary

- We experimentally demonstrate that commercial off-the-shelf (COTS) DRAM chips can perform:
 - **Functionally-complete** Boolean operations: NOT, NAND, and NOR
 - Up to 16-input AND, NAND, OR, and NOR operations
- We characterize the success rate of these operations on 256 COTS DDR4 chips from two major manufacturers
- We highlight **two key results**:
 - We can perform NOT and
 {2, 4, 8, 16}-input AND, NAND, OR, and NOR operations on COTS DRAM chips with very high success rates (>94%)
 - Data pattern and temperature only slightly affect the reliability of these operations

We believe these empirical results demonstrate the promising potential of using DRAM as a computation substrate

Simultaneous Many-Row Activation in Off-the-Shelf DRAM Chips Experimental Characterization and Analysis





ETH zürich

İsmail Emir Yüksel

Yahya C. Tuğrul F. Nisa Bostancı Geraldo F. Oliveira

A. Giray Yağlıkçı Ataberk Olgun Melina Soysal Haocong Luo

Juan Gómez–Luna Mohammad Sadr Onur Mutlu



Executive Summary

Motivation:

SAFAR

- Processing-Using-DRAM (PUD) alleviates data movement bottlenecks
- Commercial off-the-shelf (COTS) DRAM chips can perform three-input majority (MAJ3) and in-DRAM copy operations

Goal: To experimentally analyze and understand

- The computational capability of COTS DRAM chips beyond that of prior works
- The robustness of such capability under various operating conditions

Experimental Study: 120 DDR4 chips from two major manufacturers

- COTS DRAM chips can perform MAJ5, MAJ7, and MAJ9 operations and copy one DRAM row to up to 31 different rows at once
- Storing multiple redundant copies of MAJ's input operands (i.e., input replication) drastically increases robustness (>30% higher success rate)

affect the robustness of in-DRAM operations (by up to 11.52% success rate)

https://github.com/CMU-SAFARI/SiMRA-DRAM

Leveraging Simultaneous Many-Row Activation



Copy one row's content to multiple rows



In-DRAM Multiple Row Copy (Multi-RowCopy)

Simultaneously activate many rows to copy **one row's content** to **multiple destination rows**

RowClone



Multi-RowCopy



SAFARI

[Seshadri+ MICRO'13]

Key Takeaways from Multi-RowCopy

Key Takeaway 1

COTS DRAM chips are capable of copying one row's data to 1, 3, 7, 15, and 31 other rows at very high success rates

Key Takeaway 2

Multi-RowCopy in COTS DRAM chips is highly resilient to changes in data pattern, temperature, and wordline voltage



Robustness of Multi-RowCopy



COTS DRAM chips can copy one row's content to up to 31 rows with a very high success rate





Data pattern has a small effect

on the success rate of the Multi-RowCopy operation

Also in the Paper: Impact of Temperature & Voltage



More in the Paper

- Detailed hypotheses and key ideas on
 - Hypothetical row decoder circuitry
 - Input Replication
- More characterization results
 - Power consumption of simultaneous many-row activation
 - Effect of timing delays between ACT-PRE and PRE-ACT commands
 - Effect of temperature and wordline voltage
- Circuit-level (SPICE) experiments for input replication
- Potential performance benefits of enabling new in-DRAM operations
 - Majority-based computation
 - Content destruction-based cold-boot attack prevention
- Discussions on the limitations of tested COTS DRAM chips

Available on arXiv

which we call Multi-RowCopy. Second, storing multiple copies

of MAJX's input operands on all simultaneously activated rows

drastically increases the success rate (i.e., the percentage of

DRAM cells that correctly perform the computation) of the

MAJX operation. For example, MAJ3 with 32-row activation (i.e.,



that bulk bitwise operations are used in a wide variety of important applications, including databases and web search [64, 67, 79, 130, 133–140], data analytics [64, 141–144], graph processing [56, 80, 94, 130, 145], genome analysis [60, 99, 146–149], cryptography [150, 151], set operations [56, 64], and hyperdimensional computing [152–154].

https://arxiv.org/pdf/2405.06081
Our Work is Open Source and Artifact Evaluated

Code Repro	ducible		Dataset Reproducible
		S Edit Pins ▼	▼ ⁹ Fork 0 ▼ ★ Starred 6 ▼
양 main 👻 양 1 Branch 🛇 0 Tags	Q Go to file	t Add file 🔹 <> Code 👻	About - 段
🛃 unrealismail Update README.md		a51abfa · last month 🛛 5 Commits	a51abfa · last month 🕚 5 Commits Source code & scripts for experimental characterization and demonstration of 1)
DRAM-Bender	initial comit	last month	simultaneous many-row activation, 2) up to nine-input majority operations and 3) copying one row's content to up 31 rows in real DDR4 DRAM chips. Described in our DSN'24 paper by Yuksel et al. at https://arxiv.org/abs/2405.06081 Readme View license
analysis	initial comit	last month	
experimental_data	initial comit	last month	
LICENSE	initial comit	last month	
README.md	Update README.md	last month	
C README		∅ :≡	✓r ActivityE Custom properties
Simultaneous Many-Row Activation in Off-the-Shelf DRAM Chips: Experimental Characterization and Analysis			 ☆ 6 stars ◆ 4 watching ♀ 0 forks Report repository

https://github.com/CMU-SAFARI/SiMRA-DRAM

MegIS: Metagenomics In Storage

Background: GenStore

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 Ghiasi¹ Jisung Park¹ Harun Mustafa¹ Jeremie Kim¹ Ataberk Olgun¹ Arvid Gollwitzer¹ Damla Senol Cali² Can Firtina¹ Haiyu Mao¹ Nour Almadhoun Alserr¹ Rachata Ausavarungnirun³ Nandita Vijaykumar⁴ Mohammed Alser¹ Onur Mutlu¹

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

MegIS

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

Nika Mansouri Ghiasi

Mohammad Sadrosadati Harun Mustafa Arvid Gollwitzer Can Firtina

Julien Eudine Haiyu Mao Joël Lindegger Meryem Banu Cavlak

Mohammed Alser Jisung Park Onur Mutlu





What is Metagenomics?

• <u>Metagenomics</u>: Study of genome sequences of diverse organisms within a shared environment (e.g., blood, ocean, soil)



- Overcomes the limitations of traditional genomics
 - Bypasses the need for culturing individual species in isolation







What is Metagenomics?

• <u>Metagenomics</u>: Study of genome sequences of diverse organisms within a shared environment (e.g., blood, ocean, soil)



Has led to groundbreaking advances

- Precision medicine
- Understanding microbial diversity of an environment
- Discovering early warnings of communicable diseases

Metagenomic Analysis



SAFARI (e.g., > 100 TBs in emerging databases)

Motivation

- Case study of the performance of metagenomic analysis tools
- With various state-of-the-art SSD configurations



I/O data movement causes significant performance overhead

Motivation

- Case study on the throughput of metagenomic analysis tools
- With Various state-of-the-art SSD configurations



I/O becomes an even larger overhead (by 2.7x) in systems where other bottlenecks are alleviated



I/O data movement causes significant performance overhead



I/O Overhead is Hard to Avoid

I/O overhead due to accessing large, low-reuse data is hard to avoid

Sampling techniques to shrink database sizes

X Reduce accuracy to levels unacceptable for many use cases

Keeping all data required by metagenomic analysis completely and always resident in main memory

Energy inefficient, costly, unscalable, and unsustainable

- Database sizes **increase rapidly** (doubling every few months)
- Different analyses need **different databases**

Improve metagenomic analysis **performance** by reducing large **data movement overhead** from the storage system in a **cost-effective** manner

Challenges of In-Storage Processing

Existing metagenomic analysis approaches cannot be implemented as an in-storage processing system due to SSD hardware limitations

- Long latency of NAND flash chips
- Limited **DRAM capacity** inside the SSD
- Limited **DRAM bandwidth** inside the SSD



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



MegIS's Steps





Task partitioning and mapping

• Each step executes in its most suitable system







• Enable efficient access patterns to the SSD







• 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

Evaluation: Methodology Overview

Performance, Energy, and Power Analysis

Hardware Components

- Synthesized Verilog model for the in-storage accelerators
- MQSim [Tavakkol+, FAST'18] for SSD's internal operations
- Ramulator [Kim+, CAL'15] for SSD's internal DRAM

Software Components

Measure on a real system:

- AMD[®] EPYC[®] CPU with 128 physical cores
- 1-TB DRAM

Baseline Comparison Points

- Performance-optimized software, Kraken2 [Genome Biology'19]
- Accuracy-optimized software, Metalign [Genome Biology'20]
- PIM hardware-accelerated tool (using processing-in-memory), Sieve [ISCA'21]

SSD Configurations

- **SSD-C:** with SATA3 interface (0.5 GB/s sequential read bandwidth)
- SSD-P: with PCIe Gen4 interface (7 GB/s sequential read bandwidth) SAFARI

Evaluation: Speedup over the Software Baselines



MegIS provides significant speedup over both

Performance-Optimized and Accuracy-Optimized baselines



Evaluation: Speedup over the Software Baselines



MegIS provides significant speedup over both Performance-Optimized and Accuracy-Optimized baselines MegIS improves performance on both

cost-optimized and performance-optimized SSDs

Evaluation: Speedup over the PIM Baseline



MegIS provides significant speedup over the PIM baseline

Evaluation: Reduction in Energy Consumption

• On average across different input sets and SSDs



MegIS provides significant energy reduction over

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

Evaluation: Accuracy, Area, and Power

<u>Accuracy</u>

- Same accuracy as the accuracy-optimized baseline
- Significantly higher accuracy than the performance-optimized and PIM baselines
 - 4.6 5.2× higher F1 scores
 - 3 24% lower L1 norm error

Area and Power

Total for an 8-channel SSD:

- Area: 0.04 mm² (Only 1.7% of the area of three ARM Cortex R4 cores in an SSD controller)
- **Power:** 7.658 mW

Evaluation: System Cost-Efficiency

- Cost-optimized system (\$): With SSD-C and 64-GB DRAM
- Performance-optimized system (\$\$\$): With SSD-P and 1-TB DRAM



MegIS outperforms the baselines even when running on a much less costly system

Evaluation: System Cost-Efficiency

- **Cost-optimized system (\$):** With SSD-C and 64-GB DRAM
- Performance-optimized system (\$\$\$): With SSD-P and 1-TB DRAM



MegIS outperforms the baselines even when running on a much less costly system



More in the Paper

- MegIS's performance when running in-storage processing operations on the **cores existing in the SSD controller**
- MegIS's performance when using the same accelerators outside SSD
- Sensitivity analysis with varying
 - Database sizes
 - Memory capacities
 - #SSDs
 - #Channels
 - #Samples
- MegIS's performance for abundance estimation
 SAFARI

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

Nika Mansouri Ghiasi¹ Mohammad Sadrosadati¹ Harun Mustafa¹ Arvid Gollwitzer¹ Can Firtina¹ Julien Eudine¹ Haiyu Mao¹ Joël Lindegger¹ Meryem Banu Cavlak¹ Mohammed Alser¹ Jisung Park² Onur Mutlu¹ ¹ETH Zürich ²POSTECH

- Database sizes
- Memory capacities
- #SSDs
- #Channels
- #Samples



MegIS's performance for abundance estimation

https://arxiv.org/abs/2406.19113

Conclusion

Metagenomic analysis suffers from significant storage I/O data movement overhead

MegIS

The *first* **in-storage processing** system for *end-to-end* metagenomic analysis Leverages and orchestrates **processing inside** and **outside** the storage system

Improves performance

2.7×-37.2× over performance-optimized software
6.9×-100.2× over accuracy-optimized software
1.5×-5.1× over hardware-accelerated PIM baseline

High accuracy

Same as accuracy-optimized

4.8× higher F1 scores

over performance-optimized/PIM

Reduces energy consumption

5.4× over performance-optimized software
15.2× over accuracy-optimized software
1.9× over hardware-accelerated PIM baseline

Low area overhead

1.7% of the three cores in an SSD controller

Homomorphic Operations on Real PIM Systems

Homomorphic Operations on Real PIM Systems

 Harshita Gupta, Mayank Kabra, Juan Gómez-Luna, Konstantinos Kanellopoulos, and Onur Mutlu,
 <u>"Evaluating Homomorphic Operations on a Real-World Processing-In-Memory System"</u>
 <u>Proceedings of the 2023 IEEE International Symposium on Workload</u>
 <u>Characterization Poster Session (IISWC</u>), Ghent, Belgium, October 2023.
 [arXiv version]
 [Lightning Talk Slides (pptx) (pdf)]
 [Poster (pptx) (pdf)]

Evaluating Homomorphic Operations on a Real-World Processing-In-Memory System

Harshita Gupta* Mayank Kabra* Juan Gómez-Luna Konstantinos Kanellopoulos Onur Mutlu *ETH Zürich*

SAFARI https://arxiv.org/pdf/2309.06545.pdf

Side Channels on PIM Systems

Amplifying Main Memory-Based Timing Covert and Side Channels using Processing-in-Memory Operations

Konstantinos Kanellopoulos^{†*} F. Nisa Bostancı^{†*} Ataberk Olgun[†] A. Giray Yağlıkçı[†] İsmail Emir Yüksel[†] Nika Mansouri Ghiasi[†] Zülal Bingöl^{†‡} Mohammad Sadrosadati[†] Onur Mutlu[†] [†]ETH Zürich [‡]Bilkent University

https://arxiv.org/pdf/2404.11284

Distributed ML Training on Real PIM Systems
Accelerating ML Training on Real PIM Systems

To appear at PACT 2024

Analysis of Distributed Optimization Algorithms on a Real Processing-In-Memory System

Steve Rhyner¹ Haocong Luo¹ Juan Gómez-Luna² Mohammad Sadrosadati¹ Jiawei Jiang³ Ataberk Olgun¹ Harshita Gupta¹ Ce Zhang⁴ Onur Mutlu¹ ¹ETH Zurich ²NVIDIA ³Wuhan University ⁴University of Chicago

Reinforcement Learning on Real PIM Systems

SwiftRL

 Kailash Gogineni, Sai Santosh Dayapule, Juan Gomez-Luna, Karthikeya Gogineni, Peng Wei, Tian Lan, Mohammad Sadrosadati, Onur Mutlu, Guru Venkataramani,
 "SwiftRL: Towards Efficient Reinforcement Learning on Real Processing-In-Memory Systems"

 Proceedings of the 2024 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), Indianapolis, Indiana, May 2024.
 [Slides (pptx) (pdf)]
 [arXiv version]

SwiftRL: Towards Efficient Reinforcement Learning on Real Processing-In-Memory Systems

Kailash Gogineni¹ Sai Santosh Dayapule¹ Juan Gómez-Luna² Karthikeya Gogineni³ Peng Wei¹ Tian Lan¹ Mohammad Sadrosadati² Onur Mutlu² Guru Venkataramani¹ ¹George Washington University, USA ²ETH Zürich, Switzerland ³Independent

SwiftRL: Summary

- Adapted and implemented RL algorithms on a PIM architecture for exploring memory-centric systems in RL training
- Explored **optimization strategies** for enhancing RL workload performance across
 - various data types,
 - sampling strategies (SEQ, RAN, STR)
- Compared PIM-based Q-learning & SARSA on UPMEM PIM (2000 cores) to CPU & GPU
- Achieved near-linear scaling of 15x in performance with a 16x increase in PIM cores (125 to 2000)

MATSA

MATSA

 Ivan Fernandez, Christina Giannoula, Aditya Manglik, Ricardo Quislant, Nika Mansouri Ghiasi, Juan Gomez Luna, Eladio Gutierrez, Oscar Plata and Onur Mutlu, "MATSA: An MRAM-Based Energy-Efficient Accelerator for Time Series Analysis" IEEE Access, March 2024. [arXiv version]
 [IEEE Access version]

Accelerating Time Series Analysis via Processing using Non-Volatile Memories

Ivan Fernandez^{§†¶} *Christina Giannoula^{†‡} *Aditya Manglik[†] Ricardo Quislant[§] Nika Mansouri Ghiasi[†] Juan Gómez-Luna[†] Eladio Gutierrez[§] Oscar Plata[§] Onur Mutlu[†] [§]University of Malaga [†]ETH Zürich [¶]Barcelona Supercomputing Center [‡]National Technical University of Athens

SAFARI

https://arxiv.org/pdf/2211.04369



ApHMM

- Can Firtina, Kamlesh Pillai, Gurpreet S. Kalsi, Bharathwaj Suresh, Damla Senol Cali, Jeremie S. Kim, Taha Shahroodi, Meryem Banu Cavlak, Joel Lindegger, Mohammed Alser, Juan Gomez Luna, Sreenivas Subramoney and Onur Mutlu,
 "ApHMM: Accelerating Profile Hidden Markov Models for Fast and Energyefficient Genome Analysis" <u>ACM Transactions on Architecture and Code Optimization</u> (TACO), February 2024. [arXiv version]
 - [ApHMM Source Code]
 - [ACM Digital Library version]
 - Talk Video HiPEAC

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

Can Firtina¹ Kamlesh Pillai² Gurpreet S. Kalsi² Bharathwaj Suresh² Damla Senol Cali³ Jeremie S. Kim¹ Taha Shahroodi⁴ Meryem Banu Cavlak¹ Joel Lindegger¹ Mohammed Alser¹ Juan Gómez Luna¹ Sreenivas Subramoney² Onur Mutlu¹

Executive Summary

Motivation: Graph structures such as **profile Hidden Markov Models (pHMMs)** are commonly used to accurately analyze biological sequences

Problem: The parameters used in pHMMs are mainly trained and used with a **computationally intensive Baum-Welch algorithm**, causing major performance and energy overhead for many genomics workloads

Goal: Enable rapid, power-efficient, and flexible use of pHMMs for genomics workloads

ApHMM: the first flexible and hardware-software accelerator for pHMMs that can

- 1) Substantially reduce unnecessary data storage, data movement, and computations by effectively co-designing hardware and software together
- 2) Provide a flexible design to support several genomics workloads that use pHMMs

Key Results: Our ASIC implementation compared to CPU, GPU, and FPGA baselines across 3 workloads

- 15.55×-260.03×, 1.83×-5.34×, and 27.97× better performance
- Up to 2622.94× reduction in energy consumption

SAFARI <u>https://github.com/CMU-SAFARI/ApHMM-GPU</u>

RUBICON

RUBICON

 Gagandeep Singh, Mohammed Alser, Kristof Denolf, Can Firtina, Alireza Khodamoradi, Meryem Banu Cavlak, Henk Corporaal and Onur Mutlu, "RUBICON: A Framework for Designing Efficient Deep
 Learning-Based Genomic Basecallers"
 Genome Biology, February 2024.
 [arXiv version]
 [Journal Article]
 [RUBICON Source Code]

RUBICON: A Framework for Designing Efficient Deep Learning-Based Genomic Basecallers

Gagandeep Singh^{a,c} Mohammed Alser^a Kristof Denolf^c Can Firtina^{a,*} Alireza Khodamoradi^c Meryem Banu Cavlak^a Henk Corporaal^b Onur Mutlu^{a,*}

^aDepartment of Information Technology and Electrical Engineering, ETH Zürich, Switzerland
^bDepartment of Electrical Engineering, Eindhoven University of Technology, The Netherlands
^cResearch and Advanced Development, AMD, USA

https://arxiv.org/pdf/2211.03079



Better Virtual Memory: Utopia

Konstantinos Kanellopoulos, Rahul Bera, Kosta Stojiljkovic, Nisa Bostanci, Can Firtina, Rachata Ausavarungnirun, Rakesh Kumar, Nastaran Hajinazar, Mohammad Sadrosadati, Nandita Vijaykumar, and Onur Mutlu, "Utopia: Fast and Efficient Address Translation via Hybrid Restrictive & Flexible Virtual-to-Physical Address Mappings" *Proceedings of the <u>56th International Symposium on Microarchitecture</u> (<i>MICRO*), Toronto, ON, Canada, November 2023. [Slides (pptx) (pdf)] [arXiv version] [Utopia Source Code]

Utopia: Fast and Efficient Address Translation via Hybrid Restrictive & Flexible Virtual-to-Physical Address Mappings

Konstantinos Kanellopoulos¹ Rahul Bera¹ Kosta Stojiljkovic¹ Nisa Bostanci¹ Can Firtina¹ Rachata Ausavarungnirun² Rakesh Kumar³ Nastaran Hajinazar⁴ Mohammad Sadrosadati¹ Nandita Vijaykumar⁵ Onur Mutlu¹

> ¹ETH Zürich ²King Mongkut's University of Technology North Bangkok ³Norwegian University of Science and Technology ⁴Intel Labs ⁵University of Toronto

https://arxiv.org/abs/2211.12205

Utopia: Executive Summary

<u>Problem</u>: Conventional virtual memory (VM) frameworks enable a virtual address to flexibly map to any physical address. This flexibility necessitates large translation structures leading to:

(1) high translation latency and (2) large translation-induced interference in the memory hierarchy

<u>Motivation</u>: Restricting the address mapping leads to compact translation structures and reduces the overheads of address translation. Doing so across the **entire memory** has two major drawbacks:

- (1) Limits core VM functionalities (e.g., data sharing)
- (2) Increases swapping activity in the presence of free physical memory

Key Idea: Utopia is a new hybrid virtual-to-physical address mapping scheme that allows both flexible and restrictive hash-based address mappings to harmoniously co-exist in the system

Utopia manages physical memory using two types of physical memory segments:



Key Results: Outperforms (i) the state-of-the-art contiguity-aware translation scheme by 13%, and (ii) achieves 95% of the performance of an ideal perfect-TLB

https://github.com/CMU-SAFARI/Utopia

Victima

Better Virtual Memory: Victima

Konstantinos Kanellopoulos, Hong Chul Nam, F. Nisa Bostanci, Rahul Bera, Mohammad Sadrosadati, Rakesh Kumar, Davide Basilio Bartolini, and Onur Mutlu,
"Victima: Drastically Increasing Address Translation Reach by Leveraging Underutilized Cache Resources"
Proceedings of the <u>56th International Symposium on Microarchitecture</u> (MICRO), Toronto, ON, Canada, November 2023.
[Slides (pptx) (pdf)]
[arXiv version]
[Victima Source Code (Officially Artifact Evaluated with All Badges)]
Officially artifact evaluated as available, functional, reusable and reproducible.
Distinguished artifact award at MICRO 2023.

Victima: Drastically Increasing Address Translation Reach by Leveraging Underutilized Cache Resources

Konstantinos Kanellopoulos¹ Hong Chul Nam¹ F. Nisa Bostanci¹ Rahul Bera¹ Mohammad Sadrosadati¹ Rakesh Kumar² Davide Basilio Bartolini³ Onur Mutlu¹ ¹ETH Zürich ²Norwegian University of Science and Technology ³Huawei Zurich Research Center

SAFARI

https://arxiv.org/pdf/2310.04158

Executive Summary

Problem: Address translation is a major **performance bottleneck** in data-intensive workloads

Large datasets and irregular memory access patterns lead to **frequent L2TLB misses** (e.g., 20-50 MPKI) and **frequent high-latency** (e.g., 100-150 cycles) page table walks (PTW)

Motivation: Increasing the translation reach (i.e., memory covered by the TLBs) reduces PTWs. However, employing large TLBs leads to increased area, power and latency overheads.

Opportunity: Increase the translation reach of the TLB hierarchy by storing the existing TLB entries within the **existing cache hierarchy**

Victima: New software-transparent scheme that drastically increases the address translation reach of the processor's TLB hierarchy by leveraging the underutilized cache resources

Key Idea:

Transform L₂ cache blocks that store PTEs into blocks that store TLB entries



Key Benefits:

- + Efficient in native/virtualized environments
- + Fully transparent to application/OS software
- + Compatible with huge page schemes

Key Results: Victima (i) outperforms by 5.1% a state-of-the-art large TLB design and (ii) achieves similar performance to an optimistically fast 128K-entry L2 TLB

https://github.com/CMU-SAFARI/Victima

Sectored DRAM

Sectored DRAM

 Ataberk Olgun, F. Nisa Bostanci, Geraldo F. Oliveira, Yahya Can Tugrul, Rahul Bera, A. Giray Yaglikci, Hasan Hassan, Oguz Ergin, and Onur Mutlu, "Sectored DRAM: A Practical Energy-Efficient and High-Performance Fine-Grained DRAM Architecture" ACM Transactions on Architecture and Code Optimization (TACO), [online] June 2024.
 [arXiv version] [ACM Digital Library version]

Sectored DRAM: A Practical Energy-Efficient and High-Performance Fine-Grained DRAM Architecture

Ataberk Olgun[§] F. Nisa Bostancı^{§†} Geraldo F. Oliveira[§] Yahya Can Tuğrul^{§†} Rahul Bera[§] A. Giray Yağlıkcı[§] Hasan Hassan[§] Oğuz Ergin[†] Onur Mutlu[§]

SAFARI

https://arxiv.org/pdf/2207.13795



Adoption: How to Ease Programmability? (II)

 Jinfan Chen, Juan Gómez-Luna, Izzat El Hajj, YuXin Guo, and Onur Mutlu,
 "SimplePIM: A Software Framework for Productive and Efficient Processing in Memory"
 Proceedings of the <u>32nd International Conference on</u> Parallel Architectures and Compilation Techniques (PACT), Vienna, Austria, October 2023.

SimplePIM: A Software Framework for Productive and Efficient Processing-in-Memory

Jinfan Chen¹ Juan Gómez-Luna¹ Izzat El Hajj² Yuxin Guo¹ Onur Mutlu¹ ¹ETH Zürich ²American University of Beirut

SimplePIM: A Software Framework for Productive and Efficient Processing-in-Memory

Jinfan Chen, Juan Gómez Luna, Izzat El Hajj, Yuxin Guo, Onur Mutlu

https://arxiv.org/pdf/2310.01893.pdf https://github.com/CMU-SAFARI/SimplePIM





Executive Summary

- Processing-in-Memory (PIM) promises to alleviate the data movement bottleneck
- Real PIM hardware is now available, e.g., UPMEM PIM
- However, programming real PIM hardware is challenging, e.g.:
 - Distribute data across PIM memory banks,
 - Manage data transfers between host cores and PIM cores, and between PIM cores,
 - Launch PIM kernels on the PIM cores, etc.
- SimplePIM is a high-level programming framework for real PIM hardware
 - Iterators such as map, reduce, and zip
 - Collective communication with broadcast, scatter, and gather
- Implementation on UPMEM and evaluation with six different workloads
 - Reduction, vector add, histogram, linear/logistic regression, K-means
 - 4.4x fewer lines of code compared to hand-optimized code
 - Between 15% and 43% faster than hand-optimized code for three workloads
- Source code: <u>https://github.com/CMU-SAFARI/SimplePIM</u>

A State-of-the-Art PIM System



- In our work, we use the UPMEM PIM architecture
 - General-purpose processing cores called DRAM Processing Units (DPUs)
 - Up to 24 PIM threads, called *tasklets*
 - <u>32-bit integer arithmetic</u>, but multiplication/division are emulated*, as well as floating-point operations
 - 64-MB DRAM bank (MRAM), 64-KB scratchpad (WRAM)

Programming a PIM System (I)

• Example: Hand-optimized histogram with UPMEM SDK

```
... // Initialize global variables and functions for histogram
int main kernel() {
  if (tasklet id == 0)
    mem reset(); // Reset the heap
  ... // Initialize variables and the histogram
  T *input buff A = (T*)mem alloc(2048); // Allocate buffer in scratchpad memory
  for (unsigned int byte index = base tasklet; byte index < input size; byte index += stride) {
    // Boundary checking
    uint32 t l size bytes = (byte index + 2048 >= input size) ? (input size - byte index) : 2048;
    // Load scratchpad with a DRAM block
    mram read((const mram ptr void*) (mram base addr A + byte index), input buff A, l size bytes);
    // Histogram calculation
    histogram(hist, bins, input buff A, 1 size bytes/sizeof(uint32 t));
  barrier wait (&my barrier); // Barrier to synchronize PIM threads
  ... // Merging histograms from different tasklets into one histo dpu
  // Write result from scratchpad to DRAM
  if (tasklet id == 0)
    if (bins * sizeof(uint32 t) <= 2048)
      mram write(histo dpu, ( mram ptr void*)mram base addr histo, bins * sizeof(uint32 t));
    else
      for (unsigned int offset = 0; offset < ((bins * sizeof(uint32 t)) >> 11); offset++) {
        mram write(histo dpu + (offset << 9), ( mram ptr void*)(mram base addr histo +
                  (offset << 11)), 2048);
  return 0;
```

Programming a PIM System (II)

- PIM programming is challenging
 - Manage data movement between host DRAM and PIM DRAM
 - Parallel, serial, broadcast, and gather/scatter transfers
 - Manage data movement between PIM DRAM bank and scratchpad
 - 8-byte aligned and maximum of 2,048 bytes
 - Multithreaded programming model
 - Inter-thread synchronization
 - Barriers, handshakes, mutexes, and semaphores

Our Goal

Design a high-level programming framework that abstracts these hardware-specific complexities and provides a clean yet powerful interface for ease of use and high program performance

The SimplePIM Programming Framework

- SimplePIM provides standard abstractions to build and deploy applications on PIM systems
 - Management interface
 - Metadata for PIM-resident arrays
 - Communication interface
 - Abstractions for host-PIM and PIM-PIM communication
 - Processing interface
 - Iterators (map, reduce, zip) to implement workloads

Management Interface

- Metadata for PIM-resident arrays
 - array_meta_data_t describes a PIM-resident array
 - simple_pim_management_t for managing PIM-resident arrays
- lookup: Retrieves all relevant information of an array

array_meta_data_t* simple_pim_array_lookup(const char* id, simple pim management t* management);

• register: Registers the metadata of an array

void simple_pim_array_register(array_meta_data_t* meta_data, simple pim management t* management);

• free: Removes the metadata of an array

void simple_pim_array_free(const char* id, simple_pim_management_t* management);

The SimplePIM Programming Framework

- SimplePIM provides standard abstractions to build and deploy applications on PIM systems
 - Management interface
 - Metadata for PIM-resident arrays
 - Communication interface
 - Abstractions for host-PIM and PIM-PIM communication
 - Processing interface
 - Iterators (map, reduce, zip) to implement workloads

Host-to-PIM Communication: Broadcast

- SimplePIM Broadcast
 - Transfers a host array to all PIM cores in the system

void simple_pim_array_broadcast(char* const id, void* arr, uint64_t len, uint32_t type_size, simple_pim_management_t* management);



Host-to-PIM Communication: Scatter/Gather

• SimplePIM Scatter

- Distributes an array to PIM DRAM banks

```
void simple_pim_array_scatter(char* const id, void* arr, uint64_t len,
uint32_t type_size, simple_pim_management_t* management);
```

• SimplePIM Gather

SAFARI

- Collects portions of an array from PIM DRAM banks

void* simple_pim_array_gather(char* const id, simple_pim_management_t*
management);



PIM-PIM Communication: AllReduce

- SimplePIM AllReduce
 - Used for algorithm synchronization
 - The programmer specifies an accumulative function

```
void simple_pim_array_allreduce(char* const id, handle_t* handle,
simple_pim_management_t* management);
```





PIM-PIM Communication: AllGather

• SimplePIM AllGather

 Combines array pieces and distributes the complete array to all PIM cores

```
void simple_pim_array_allgather(char* const id, char* new_id,
simple_pim_management_t* management);
```



The SimplePIM Programming Framework

- SimplePIM provides standard abstractions to build and deploy applications on PIM systems
 - Management interface
 - Metadata for PIM-resident arrays
 - Communication interface
 - Abstractions for host-PIM and PIM-PIM communication
 - Processing interface
 - Iterators (map, reduce, zip) to implement workloads

Processing Interface: Map

- Array Map
 - Applies map_func to every element of the data array

void simple_pim_array_map(const char* src_id, const char* dest_id, uint32_t output_type, handle_t* handle, simple_pim_management_t* management);



Processing Interface: Reduction

• Array Reduction

- The map_to_val_func function transforms an input element to an output value and an output index
- The acc_func function accumulates the output values onto the output array

void simple_pim_array_red(const char* src_id, const char* dest_id, uint32_t output_type, uint32_t output_len, handle_t* handle, simple_pim_management_t* management);


Processing Interface: Zip

- Array Zip
 - Takes two input arrays and combines their elements into an output array

void simple_pim_array_zip(const char* src1_id, const char* src2_id, const char* dest_id, simple_pim_management_t* management);





General Code Optimizations

- Strength reduction
- Loop unrolling
- Avoiding boundary checks
- Function inlining
- Adjustment of data transfer sizes

More in the Paper

Strength reduction

Loop unrolling

SimplePIM: A Software Framework for Productive and Efficient Processing-in-Memory

Jinfan Chen¹ Juan Gómez-Luna¹ Izzat El Hajj² Yuxin Guo¹ Onur Mutlu¹ ¹ETH Zürich ²American University of Beirut

Adju <u>https://arxiv.org/pdf/2310.01893.pdf</u>



Evaluation Methodology

- Evaluated system
 - UPMEM PIM system with 2,432 PIM cores with 159 GB of PIM DRAM
- Real-world Benchmarks
 - Vector addition
 - Reduction
 - Histogram
 - K-Means
 - Linear regression
 - Logistic regression
- Comparison to hand-optimized codes in terms of programming productivity and performance

Productivity Improvement (I)

• Example: Hand-optimized histogram with UPMEM SDK

```
... // Initialize global variables and functions for histogram
int main kernel() {
  if (tasklet id == 0)
    mem reset(); // Reset the heap
  ... // Initialize variables and the histogram
  T *input buff A = (T*)mem alloc(2048); // Allocate buffer in scratchpad memory
  for (unsigned int byte index = base tasklet; byte index < input size; byte index += stride) {
    // Boundary checking
    uint32 t l size bytes = (byte index + 2048 >= input size) ? (input size - byte index) : 2048;
    // Load scratchpad with a DRAM block
    mram read((const mram ptr void*) (mram base addr A + byte index), input buff A, l size bytes);
    // Histogram calculation
    histogram(hist, bins, input buff A, 1 size bytes/sizeof(uint32 t));
  barrier wait (&my barrier); // Barrier to synchronize PIM threads
  ... // Merging histograms from different tasklets into one histo dpu
  // Write result from scratchpad to DRAM
  if (tasklet id == 0)
    if (bins * sizeof(uint32 t) <= 2048)
      mram write(histo dpu, ( mram ptr void*)mram base addr histo, bins * sizeof(uint32 t));
    else
      for (unsigned int offset = 0; offset < ((bins * sizeof(uint32 t)) >> 11); offset++) {
        mram write(histo dpu + (offset << 9), ( mram ptr void*)(mram base addr histo +</pre>
                  (offset << 11)), 2048);
  return 0;
```

SAFARI

Productivity Improvement (II)

• Example: SimplePIM histogram

```
// Programmer-defined functions in the file "histo filepath"
void init func (uint32 t size, void* ptr) {
  char* casted value ptr = (char*) ptr;
 for (int i = 0; i < size; i++)</pre>
    casted value ptr[i] = 0;
}
void acc func (void* dest, void* src) {
  *(uint32 t*)dest += *(uint32 t*)src;
void map to val func (void* input, void* output, uint32 t* key) {
  uint32 t \overline{d} = \overline{*}((\text{uint32 t}^*)\text{input});
  *(uint32 t*)output = 1;
  *key = d * bins >> 12;
// Host side handle creation and iterator call
handle t* handle = simple pim create handle ("histo filepath", REDUCE, NULL, 0);
// Transfer (scatter) data to PIM, register as "t1"
simple pim array scatter("t1", src, bins, sizeof(T), management);
// Run histogram on "t1" and produce "t2"
simple pim array red("t1", "t2", sizeof(T), bins, handle, management);
```

Productivity Improvement (III)

• Lines of code (LoC) reduction

	SimplePIM	Hand-optimized	LoC Reduction
Reduction	14	83	5.93×
Vector Addition	14	82	5.86×
Histogram	21	114	5•43×
Linear Regression	48	157	3.27×
Logistic Regression	59	176	2.98×
K-Means	68	206	3.03×

SimplePIM reduces the number of lines of effective code by a factor of 2.98× to 5.93×

Performance Evaluation (I)



SAFARI

Performance Evaluation (II)





SimplePIM scales better than hand-optimized implementations for reduction, histogram, and linear regression

SimplePIM outperforms hand-optimized implementations for vector addition, logistic regression, and k-means by 15%-43%

SAFARI

Discussion

- SimplePIM is devised for PIM architectures with
 - A host processor with access to standard main memory and PIM-enabled memory
 - PIM processing elements (PEs) that communicate via the host processor
 - The number of PIM PEs scales with memory capacity
- SimplePIM emulates the communication between PIM cores via the host processor
- Other parallel patterns can be incorporated in future work
 - Prefix sum and filter can be easily added
 - Stencil and convolution would require fine-grained scattergather for halo cells
 - Random access patterns would be hard to support

SimplePIM: A Software Framework for Productive and Efficient Processing-in-Memory

Jinfan Chen¹ Juan Gómez-Luna¹ Izzat El Hajj² Yuxin Guo¹ Onur Mutlu¹ ¹ETH Zürich ²American University of Beirut

https://arxiv.org/pdf/2310.01893.pdf

Source Code

https://github.com/CMU-SAFARI/SimplePIM

SimplePIM Private		😒 Edit Pins 👻	• Unwatch 3
<mark>ঃ main - </mark> ঃ 1 branch 📀 0 tags		Go to file Add file -	<> Code -
😛 Wangsitu98 interface cleanups, add	ed allreduce and allgather	3421614 2 days ago	37 commits
benchmarks	interface cleanups, added allreduce and	allgather	2 days ago
🖿 lib	interface cleanups, added allreduce and	allgather	2 days ago
🗋 .gitignore	some cleanups		3 weeks ago
B README.md	pushed SimplePIM		last month
E README.md			Ø

SimplePIM @

This project implements SimplePIM, a software framework for easy and efficient in-memory-hardware programming. The code is implemented on UPMEM, an actual, commercially available PIM hardware that combines traditional DRAM memory with general-purpose in-order cores inside the same chip. SimplePIM processes arrays of arbitrary elements on a PIM device by calling iterator functions from the host and provides primitives for communication among PIM cores and between PIM and the host system.

We implement six applications with SimplePIM on UPMEM:

- Vector Addtition
- Reduction
- K-Means Clustering
- Histogram
- Linear Regression
- Logistic Regression

Previous manual UPMEM implementations of the same applications can be found in PrIM benchmark (https://github.com/CMU-SAFARI/prim-benchmarks), dpu_kmeans (https://github.com/upmem/dpu_kmeans) and prim-ml (https://github.com/CMU-SAFARI/pim-ml). These previous implementations can serve as baseline for measuring SimplePIM's performance as well as productivity improvements.

SimplePIM: Summary

- Processing-in-Memory (PIM) promises to alleviate the data movement bottleneck
- Real PIM hardware is now available, e.g., UPMEM PIM
- However, programming real PIM hardware is challenging, e.g.:
 - Distribute data across PIM memory banks,
 - Manage data transfers between host cores and PIM cores, and between PIM cores,
 - Launch PIM kernels on the PIM cores, etc.
- SimplePIM is a high-level programming framework for real PIM hardware
 - Iterators such as map, reduce, and zip
 - Collective communication with broadcast, scatter, and gather
- Implementation on UPMEM and evaluation with six different workloads
 - Reduction, vector add, histogram, linear/logistic regression, K-means
 - 4.4x fewer lines of code compared to hand-optimized code
 - Between 15% and 43% faster than hand-optimized code for three workloads
- Source code: <u>https://github.com/CMU-SAFARI/SimplePIM</u>

SAFARI

Constable



Improving Performance and Power Efficiency by Safely Eliminating Load Instruction Execution

Rahul Bera* Adithya Ranganathan* Joydeep Rakshit Sujit Mahto Anant V. Nori Jayesh Gaur Ataberk Olgun Konstantinos Kanellopoulos Mohammad Sadrosadati Sreenivas Subramoney Onur Mutlu







Key Problem

Load instructions are a key limiter of instruction-level parallelism (ILP)

Data Dependence

Stall load-dependent instructions due to **long load execution latency**

Resource Dependence

Stall other loads due to contention in load execution resources





Prior Works on Tolerating Load Latency

- Load value prediction (LVP) [Lipasti+, ASPLOS'96; Sazeides+, MICRO'96; ...]
- Memory renaming (MRN) [Moshovos+, ISCA'97; Tyson+, MICRO'97; ...]

Mitigate Data Dependence

By **speculatively executing** load-dependent instructions using a predicted load value



Do Not Mitigate Resource Dependence

Predicted load still gets executed to verify speculation, consuming execution resources



Motivation

Safely breaking load data dependency without executing a load instruction may provide additional performance benefits



By finding load instructions that repeatedly produce identical results across dynamic instances



Key Finding I: Global-Stable Loads

- Some loads repeatedly fetch the same data value from same load address across entire workload
 - Both operations, address generation & data fetch, produce identical results across all dynamic instances
 - Prime targets for breaking data dependency without execution

Global-Stable Load



Key Finding I: Global-Stable Loads



Nearly **1 in every 3** dynamic loads is a **global-stable load**

SAFARI

Across a wide range of 90 workloads

In the Paper: Analysis of Global-Stable Loads

- Why do these loads even exist in well-optimized real-world workloads?
 - Accessing global-scope variables
 - Accessing local variables of inline functions
 - Limited set of architectural registers
- Can increasing architectural registers help?
 - Very small change even after doubling x64 registers
- Deeper characterization of global-stable loads
 - Which addressing mode do they use?
 - How far away do they appear in a workload?

SAFARI

In the Paper: Analysis of Global-Stable Loads

Nhu do those loads over exist in well entimize

Constable: Improving Performance and Power Efficiency by Safely Eliminating Load Instruction Execution

*Rahul Bera¹ *Adithya Ranganathan² Joydeep Rakshit² Sujit Mahto² Anant V. Nori² Jayesh Gaur² Ataberk Olgun¹ Konstantinos Kanellopoulos¹ Mohammad Sadrosadati¹ Sreenivas Subramoney² Onur Mutlu¹

¹ETH Zürich ²Intel Processor Architecture Research Lab

Load instructions often limit instruction-level parallelism (ILP) in modern processors due to data and resource dependences they cause. Prior techniques like Load Value Prediction (LVP) and Memory Renaming (MRN) mitigate load data dependence by predicting the data value of a load instruction. However, they fail to mitigate load resource dependence as the predicted load instruction gets executed nonetheless (even on a correct prediction), which consumes hard-to-scale pipeline resources that otherwise could have been used to execute other load instructions.

Our goal in this work is to improve ILP by mitigating both load data dependence and resource dependence. To this end, we propose a purely-microarchitectural technique called Constable, that safely eliminates the execution of load instructions. Constable dynamically identifies load instructions that have restall for multiple cycles, which can limit ILP. On the other hand, a load instruction consumes multiple hard-to-scale hardware resources (e.g., reservation station entry, port to address generation unit, L1-data cache read port) which often causes resource dependence in the pipeline, also limiting ILP.

Prior works propose many latency tolerance techniques to improve ILP by mitigating load data dependence. Load Value Prediction (LVP) [32,42,43,71,98,107,114,139–143,151,153–155, 159,160] and Memory Renaming (MRN) [120,121,147,177,178] are two such widely-studied techniques that mitigate load data dependence via data value speculation. LVP and MRN speculatively execute load-data-dependent instructions using the predicted value of the load instruction, thus improving ILP.

- Hov

SAFAR

https://arxiv.org/pdf/2406.18786

A significant fraction of loads are **global-stable**



But do they limit ILP even when using <u>load value prediction and memory renaming</u>?

Key Finding II: Global-Stable Loads Cause Resource Dependence

All execution cycles where at least one load port is utilized



In an aggressive OoO processor with 6-wide issue, 3 load ports, a **load value predictor** (EVES [Seznec, CVP'18]), and **memory renaming enabled SAFARI**

Key Finding II: Global-Stable Loads Cause Resource Dependence



In an aggressive OoO processor with 6-wide issue, 3 load ports, a load value predictor (EVES [Seznec, CVP'18]), and memory renaming enabled SAFARI

Key Finding II: Global-Stable Loads Cause Resource Dependence

All execution cycles where at least one load port is utilized

Even when using load value prediction and memory renaming, global-stable loads limit ILP due to resource dependence



What's the performance headroom of mitigating the resource dependence?

a **load value predictor (**EVES [Seznec, CVP'18]), and **memory renaming enabled**



Key Finding III: High Performance Headroom



Mitigating both data and resource dependence has more than 2x the performance benefit of mitigating only data dependence of global-stable loads

Ideal elimination of global-stable loads exceeds performance of a processor with **2x wider** load execution

0

Load Execution Resources Lag Behind



Load Execution Resources Lag Behind



Mitigating load resource dependence has high performance potential in recent and future generation processors



Our Goal

To **improve instruction-level parallelism** by mitigating **both load data dependence** and **resource dependence**





A purely-microarchitectural technique

Mitigates both load data dependence and load resource dependence

By **safely** eliminating the **entire execution** of a load instruction

Constable: Key Insight



SAFARI

Constable: Key Insight



Constable: Key Steps



Dynamically identify load instructions that have historically fetched the same data from the same load address (i.e., likely-stable)



Eliminate execution of likely-stable loads by tracking modifications to their source registers and their load addresses



Prior Related Literature

Rich literature on skipping redundant computations by **memoizing previously-computed results**

[Michie, Nature'68; Harbison+, ASPLOS'82; Richardson, SCA'93; Sodani+, ISCA'97; González+, ICPP'99; ...]



Aim to memoize every instruction

including multiple dynamic instances of each instruction



Require large memoization buffer

Often bigger than the size of L1 data cache



Key Improvements over Literature

Rich literature on skipping redundant computations by **memoizing previously-computed results**

[Michie, Nature'68; Harbison+, ASPLOS'82; Richardson, SCA'93; Sodani+, ISCA'97; González+, ICPP'99; ...]

1

Focus only on loads that are likely stable



Lower storage overhead

with high load elimination coverage

SAFARI



Lower design complexity

Fewer port requirements, lower power
Key Improvements over Literature

Rich literature on skipping redundant computations by **memoizing previously-computed results**

[Michie, Nature'68; Harbison+, ASPLOS'82; Richardson, SCA'93; Sodani+, ISCA'97; González+, ICPP'99; ...]



Key Improvements over Literature

Rich literature on skipping redundant computations by **memoizing previously-computed results**

[Michie, Nature'68; Harbison+, ASPLOS'82; Richardson, SCA'93; Sodani+, ISCA'97; González+, ICPP'99; ...]



Design Overview

Constable: Key Steps





Eliminate

by tracking modifications



Identify a Likely-Stable Load



 Using a stability confidence add rax, 0x10 counter per load instruction mov r8, [rbp+0x8] 5 sub rax, r8 Same data & address cmp rsi, rax as last dynamic instance add rax, 0x10 +1mov r8, [rbp+0x8] 6 Stability sub rax, r8 Confidence cmp rsi, rax ret mov r8, [rbp+0x8] 3 Different data or sub rax, r8 different address

cmp rsi, rax

jle 0x40230e

Eliminate a Likely-Stable Load



to handle correct elimination of in-flight loads

Stop Elimination of a Likely-Stable Load



More in the Paper

- Ensuring safe and correct elimination in presence of
 - Out-of-order load issue
 - Multi-threaded & multi-core execution
 - Wrong-path execution
- Integration of Constable into the processor pipeline



- Microarchitecture for breaking data dependence on the eliminated loads
- Microarchitecture of Constable's own structures
 - Read and write port requirements

More in the Paper

• Ensuring safe and correct elimination in presence of

Out of order load iccus

• Ir

- Rea

SAFARI

Constable: Improving Performance and Power Efficiency by Safely Eliminating Load Instruction Execution

*Rahul Bera¹ *Adithya Ranganathan² Joydeep Rakshit² Sujit Mahto² Anant V. Nori² Jayesh Gaur² Ataberk Olgun¹ Konstantinos Kanellopoulos¹ Mohammad Sadrosadati¹ Sreenivas Subramoney² Onur Mutlu¹

¹ETH Zürich ²Intel Processor Architecture Research Lab

Load instructions often limit instruction-level parallelism (ILP) in modern processors due to data and resource dependences they cause. Prior techniques like Load Value Prediction (LVP) and Memory Renaming (MRN) mitigate load data dependence by predicting the data value of a load instruction. However, they fail to mitigate load resource dependence as the predicted load instruction gets executed nonetheless (even on a correct prediction), which consumes hard-to-scale pipeline resources that otherwise could have been used to execute other load instructions.

Our goal in this work is to improve ILP by mitigating both load data dependence and resource dependence. To this end, we propose a purely-microarchitectural technique called Constable, that safely eliminates the execution of load instructions. Constable dynamically identifies load instructions that have restall for multiple cycles, which can limit ILP. On the other hand, a load instruction consumes multiple hard-to-scale hardware resources (e.g., reservation station entry, port to address generation unit, L1-data cache read port) which often causes resource dependence in the pipeline, also limiting ILP.

Prior works propose many latency tolerance techniques to improve ILP by mitigating load data dependence. Load Value Prediction (LVP) [32,42,43,71,98,107,114,139–143,151,153–155, 159,160] and Memory Renaming (MRN) [120,121,147,177,178] are two such widely-studied techniques that mitigate load data dependence via data value speculation. LVP and MRN speculatively execute load-data-dependent instructions using the predicted value of the load instruction, thus improving ILP.

https://arxiv.org/pdf/2406.18786

Evaluation

Methodology

- Industry-grade x86-64 simulator modeling aggressive OoO processor
 - 8-wide fetch, 6-wide issue to 3 load ports, 512-entry ROB
 - With memory renaming, zero/constant/move elimination, branch folding
 - Five prefetchers throughout cache hierarchy
- 90 workloads of wide variety
 - All from SPEC CPU 2017
 - Client (SYSMark, DaCapo, ...)
 - Enterprise (SPECjbb, SPECjEnterprise, ...)
 - Server (BigBench, Hadoop, ...)

Mechanisms compared against

- EVES, the state-of-the-art load value predictor [Seznec, CVP'18]
- Early Load Address Resolution [Bekerman+, ISCA'00]
- Register File Prefetching [Shukla+, ISCA'22]

Configurations

- No simultaneous multi-threading (SMT)
- 2-way SMT

Performance Improvement in noSMT



Constable alone provides similar performance as EVES with only $\frac{1}{2}$ of EVES' storage overhead

Constable on top of EVES outperforms EVES alone

Performance Improvement in 2-way SMT



Performance Improvement in 2-way SMT



Constable provides higher performance benefits in a 2-way SMT processor



Improvement in Resource Efficiency



Improvement in Resource Efficiency



Constable significantly improves resource efficiency by eliminating load instruction execution



Reduction in Dynamic Power

Memory Execution Unit Power

OoO Unit Power





Reduction in Dynamic Power



By eliminating load instruction execution, Constable reduces dynamic power consumption



Area and Power Overhead of Constable's Own Structures



12.4 KB

Storage overhead per core



0.232 mm²

0.0061% area of Intel Alderlake-S processor



SAFARI

Low Energy Up to 10.8 pJ/read and 16.7 pJ/write

More in the Paper

- Load elimination coverage of Constable
 - 23.5% of all dynamic loads are eliminated
- Per-workload performance analysis
 - Up to 31.2% over baseline
 - 60/90 workloads outperforms EVES by more than 5%
- Performance contribution per load category
 - Stack loads contribute the highest
- Performance improvement over prior works
 - 4.7% over Early load address resolution
 - 3.6% over Register file prefetching
- Performance **sensitivity**:
 - Higher performance in every configuration up to 2X load execution width
 - Higher performance in every configuration up to 2X pipeline depth

More in the Paper

• P

• P

SAFAR

High

Load elimination coverage of Constable 23.5% of all dynamic loads are eliminated

Constable: Improving Performance and Power Efficiency by Safely Eliminating Load Instruction Execution

*Rahul Bera¹ *Adithya Ranganathan² Joydeep Rakshit² Sujit Mahto² Anant V. Nori² Jayesh Gaur² Ataberk Olgun¹ Konstantinos Kanellopoulos¹ Mohammad Sadrosadati¹ Sreenivas Subramoney² Onur Mutlu¹

¹ETH Zürich ²Intel Processor Architecture Research Lab

Load instructions often limit instruction-level parallelism (ILP) in modern processors due to data and resource dependences they cause. Prior techniques like Load Value Prediction (LVP) and Memory Renaming (MRN) mitigate load data dependence by predicting the data value of a load instruction. However, they fail to mitigate load resource dependence as the predicted load instruction gets executed nonetheless (even on a correct prediction), which consumes hard-to-scale pipeline resources that otherwise could have been used to execute other load instructions.

Our goal in this work is to improve ILP by mitigating both load data dependence and resource dependence. To this end, we propose a purely-microarchitectural technique called Constable, that safely eliminates the execution of load instructions. Constable dynamically identifies load instructions that have restall for multiple cycles, which can limit ILP. On the other hand, a load instruction consumes multiple hard-to-scale hardware resources (e.g., reservation station entry, port to address generation unit, L1-data cache read port) which often causes resource dependence in the pipeline, also limiting ILP.

Prior works propose many latency tolerance techniques to improve ILP by mitigating load data dependence. Load Value Prediction (LVP) [32,42,43,71,98,107,114,139–143,151,153–155, 159,160] and Memory Renaming (MRN) [120,121,147,177,178] are two such widely-studied techniques that mitigate load data dependence via data value speculation. LVP and MRN speculatively execute load-data-dependent instructions using the predicted value of the load instruction, thus improving ILP.

https://arxiv.org/pdf/2406.18786

oth

To Summarize...

Our Key Findings



A large fraction (34%) of dynamic loads fetch the same data from the same address throughout the entire workload



These global-stable loads cause significant ILP loss due to **resource dependence**



Eliminating global-stable load execution provides more than **2x the performance** benefit of just **breaking** their **load data dependency**



Our Proposal

Constable

Identifies and eliminates loads that repeatedly fetch same data from same address



There's Still Headroom...

Constable successfully eliminates 57% of all global-stable loads at runtime

> 43% of global-stable loads do not get eliminated

We need to understand more

software primitives that generate global-stable loads



Open-Source Tool



A tool to analyze load instructions in any off-the-shelf x86(-64) program

Load-Inspector Public		ⓒ Unwatch 5 \checkmark $\bigcirc \bigcirc \bigcirc \bigcirc$ Fork 0 \checkmark \checkmark Starred 2 \checkmark		
양 main 👻 양 1 Branch 🛇 1 Tags	Q Go to file	t Add file 🔹 <> Code 👻	About ©	
Rahul Bera Updated README		183460c · last week 🚯 5 Commits	A binary instrumentation tool to analyze load instructions in any off-the-shelf	
📄 logo	First commit	last month	x86(-64) program.	
src	Added post-processing script	last month	compiler emulation microarchitecture intel-sde binary-instrumentation	
test/fft	First commit	last month	C Readme	
tools	Added post-processing script	last month	▲IT license	
🗋 .gitignore	Added post-processing script	last month	-∿ Activity	
	First commit	last month	☆ 2 stars	
README.md	Updated README	last week	• 5 watching	
inspector	Added post-processing script	last month	양 0 forks Report repository	
ী install.sh	First commit	last month		

SAFARI

https://github.com/CMU-SAFARI/Load-Inspector

Open-Source Tool



A tool to analyze load instructions in any off-the-shelf x86(-64) program

		°g Fork 0 → 🚖 Starred 2 →					
Study global-stable loads							
	logo	First commit	last month	x86(-64) program.			
Study the effects of increasing architectural registers using APX extension to x64 ISA							
	🗋 README.md	Updated README	last week	⊙ 5 watching			
	inspector	Added post-processing script	last month	Report repository			
	🗅 Install.sh	First commit	last month				

https://github.com/CMU-SAFARI/Load-Inspector



Improving Performance and Power Efficiency by Safely Eliminating Load Instruction Execution







