# Machine Learning Driven Memory and Storage Systems

Onur Mutlu omutlu@gmail.com

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

19 December 2023

**IBM** Research



**ETH** zürich





# 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

### Data is Key for Modern Workloads



### **In-memory Databases**

[Mao+, EuroSys'12; Clapp+ (**Intel**), IISWC'15]



### **In-Memory Data Analytics**

[Clapp+ (**Intel**), IISWC'15; Awan+, BDCloud'15]



**Graph/Tree Processing** [Xu+, IISWC'12; Umuroglu+, FPL'15]



**Datacenter Workloads** [Kanev+ (**Google**), ISCA'15]

# Exponential Growth of Neural Networks



Source: https://youtu.be/Bh13Idwcb0Q?t=283

SAFAR

### Huge Demand for Performance & Efficiency

5

### Data Overwhelms Modern Machines





#### **In-memory Databases**

#### **Graph/Tree Processing**

### Data → performance & energy bottleneck



### In-Memory Data Analytics

[Clapp+ (**Intel**), IISWC'15; Awan+, BDCloud'15]



**Datacenter Workloads** [Kanev+ (**Google**), ISCA'15]

### Data is Key for Modern Workloads





**Google's web browser** 



### **TensorFlow Mobile**

Google's machine learning framework



**Google's video codec** 



### Data Overwhelms Modern Machines



### Data → performance & energy bottleneck



**Google's video codec** 



### 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 <u>23rd International Conference on Architectural Support for Programming</u> <u>Languages and Operating Systems</u> (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 Boroumand<sup>1</sup>Saugata Ghose<sup>1</sup>Youngsok Kim<sup>2</sup>Rachata Ausavarungnirun<sup>1</sup>Eric Shiu<sup>3</sup>Rahul Thakur<sup>3</sup>Daehyun Kim<sup>4,3</sup>Aki Kuusela<sup>3</sup>Allan Knies<sup>3</sup>Parthasarathy Ranganathan<sup>3</sup>Onur Mutlu<sup>5,1</sup>9

# 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 <u>30th International Conference on Parallel Architectures and Compilation</u> <u>Techniques</u> (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**

Amirali Boroumand<sup>†</sup>Saugata Ghose<sup>‡</sup>Berkin Akin<sup>§</sup>Ravi Narayanaswami<sup>§</sup>Geraldo F. Oliveira<sup>★</sup>Xiaoyu Ma<sup>§</sup>Eric Shiu<sup>§</sup>Onur Mutlu<sup>★†</sup>

<sup>†</sup>Carnegie Mellon Univ. <sup>°</sup>Stanford Univ. <sup>‡</sup>Univ. of Illinois Urbana-Champaign <sup>§</sup>Google <sup>\*</sup>ETH Zürich



# An Intelligent Architecture Handles Data Well



### How to Handle Data Well

#### **Ensure data does not overwhelm** the components

 via intelligent algorithms, architectures & system designs: algorithm-architecture-devices

### **Take advantage of** vast amounts of **data** and metadata

to improve architectural & system-level decisions

Understand and exploit properties of (different) data

to improve algorithms & architectures in various metrics

# 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

### Fundamentally Better Architectures

# **Data-centric**

# **Data-driven**

# **Data-aware**



### We Need to Revisit the Entire Stack

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

### We can get there step by step

### A Blueprint for Fundamentally Better Architectures

#### Onur Mutlu, "Intelligent Architectures for Intelligent Computing Systems" Invited Paper in Proceedings of the Design, Automation, and Test in Europe Conference (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)]

### Intelligent Architectures for Intelligent Computing Systems

Onur Mutlu ETH Zurich omutlu@gmail.com

# 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 (<u>Retrospective (pdf) Full Issue</u>).* 

Self-Optimizing Memory Controllers: A Reinforcement Learning Approach

Engin İpek<sup>1,2</sup> Onur Mutlu<sup>2</sup> José F. Martínez<sup>1</sup> Rich Caruana<sup>1</sup>

<sup>1</sup>Cornell 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 <u>54th International Symposium on Microarchitecture</u> (<i>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 Bera<sup>1</sup> Konstantinos Kanellopoulos<sup>1</sup>

Anant V. Nori<sup>2</sup> 7 Onur Mutlu<sup>1</sup>

Taha Shahroodi<sup>3,1</sup>

<sup>1</sup>ETH Zürich <sup>2</sup>Processor Architecture Research Labs, Intel Labs <sup>3</sup>TU Delft

Sreenivas Subramoney<sup>2</sup>

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

<sup>1</sup>ETH Zürich <sup>2</sup>Intel Processor Architecture Research Lab <sup>3</sup>LIRMM, 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. [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

Self Optimizing Memory 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 (<u>Retrospective (pdf) Full Issue</u>).* 

Self-Optimizing Memory Controllers: A Reinforcement Learning Approach

Engin İpek<sup>1,2</sup> Onur Mutlu<sup>2</sup> José F. Martínez<sup>1</sup> Rich Caruana<sup>1</sup>

<sup>1</sup>Cornell University, Ithaca, NY 14850 USA

 $^2$  Microsoft Research, Redmond, WA 98052 USA

# 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	$\begin{array}{l} \text{ACT} \rightarrow \text{READ} \\ \text{ACT} \rightarrow \text{WRITE} \end{array}$	tRCD	15ns
	$\mathrm{ACT} \to \mathrm{PRE}$	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

# Why So Many Timing Constraints? (I)



**Figure 4.** DRAM bank operation: Steps involved in serving a memory request [17]  $(V_{PP} > V_{DD})$ 

Category	RowCmd↔RowCmd		RowCmd↔ColCmd		ColCmd↔ColCmd			ColCmd→DATA			
Name	tRC	tRAS	tRP	tRCD	tRTP	$tWR^*$	tCCD	$tRTW^{\dagger}$	$tWTR^*$	CL	CWL
Commands	A→A	A→P	P→A	A→R/W	$R \rightarrow P$	$W^*\!\rightarrow\!P$	$R(W) \rightarrow R(W)$	$R \rightarrow W$	$W^* \rightarrow R$	<b>R→DATA</b>	W→DATA
Scope	Bank	Bank	Bank	Bank	Bank	Bank	Channel	Rank	Rank	Bank	Bank
Value (ns)	$\sim$ 50	~35	13-15	13-15	~7.5	15	5-7.5	11-15	~7.5	13-15	10-15

A: ACTIVATE- P: PRECHARGE- R: READ- W: WRITE

+ W: WRITE \* Goes into effect after the last write *data*, not from the WRITE command † Not explicitly specified by the JEDEC DDR3 standard [18]. Defined as a function of other timing constraints.

Table 1. Summary of DDR3-SDRAM timing constraints (derived from Micron's 2Gb DDR3-SDRAM datasheet [33])

#### Kim et al., "A Case for Exploiting Subarray-Level Parallelism (SALP) in DRAM," ISCA 2012.

# Why So Many Timing Constraints? (II)



Figure 6. Charge Flow Between the Cell Capacitor ( $C_C$ ), Bitline Parasitic Capacitor ( $C_B$ ), and the Sense-Amplifier ( $C_B \approx 3.5C_C$  [39])

Table 2. Timing Constraints (DDR3-1066)	[4	3]
---	----	----

Phase	Commands	Name	Value
1	$\begin{array}{l} \mathrm{ACT} \rightarrow \mathrm{READ} \\ \mathrm{ACT} \rightarrow \mathrm{WRITE} \end{array}$	tRCD	15ns
	$ACT \rightarrow PRE$	tRAS	37.5ns
2	$\begin{array}{l} \text{READ} \rightarrow data \\ \text{WRITE} \rightarrow data \end{array}$	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

Lee et al., "Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture," HPCA 2013.

### DRAM Controller Design Is Becoming More Difficult



- Heterogeneous agents: CPUs, GPUs, and HWAs
- Main memory interference between CPUs, GPUs, HWAs
- Many timing constraints for various memory types
- Many goals at the same time: performance, fairness, QoS, energy efficiency, ...

# 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?

# Memory Controller: Performance Function



How to schedule requests to maximize system performance?

# Self-Optimizing DRAM Controllers

- 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.

Ipek+, "Self Optimizing Memory Controllers: A Reinforcement Learning Approach," ISCA 2008.

## Self-Optimizing DRAM Controllers



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

# Self-Optimizing DRAM Controllers

- 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


### Self-Optimizing DRAM Controllers

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

*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

### Self Optimizing DRAM Controllers

+ 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

### 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 (<u>Retrospective (pdf) Full Issue</u>).* 

Self-Optimizing Memory Controllers: A Reinforcement Learning Approach

Engin İpek<sup>1,2</sup> Onur Mutlu<sup>2</sup> José F. Martínez<sup>1</sup> Rich Caruana<sup>1</sup>

<sup>1</sup>Cornell University, Ithaca, NY 14850 USA

 $^2$  Microsoft Research, Redmond, WA 98052 USA

**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 <u>54th International Symposium on Microarchitecture</u> (<i>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 Bera<sup>1</sup> Konstantinos Kanellopoulos<sup>1</sup>

Anant V. Nori<sup>2</sup> Taha Shahroodi<sup>3,1</sup> Onur Mutlu<sup>1</sup>

<sup>1</sup>ETH Zürich <sup>2</sup>Processor Architecture Research Labs, Intel Labs <sup>3</sup>TU Delft

Sreenivas Subramoney<sup>2</sup>

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





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

### **Executive Summary**

- Background: Prefetchers predict addresses of future memory requests by associating memory access patterns with pieces of program and system information (called feature)
- **Problem**: Three key shortcomings of prior prefetchers:
  - Predict mainly using a single program feature
  - Lack inherent system awareness (e.g., memory bandwidth usage)
  - Lack in-silicon customizability
- **Goal**: Design a prefetching framework that:
  - Learns from multiple features and inherent system-level feedback
  - Can be customized in silicon to use different features and/or prefetching objectives
- Contribution: Pythia, which formulates prefetching as reinforcement learning problem
  - Takes adaptive prefetch decisions using multiple features and system-level feedback
  - Can be **customized in silicon** for target workloads via simple configuration registers
  - Proposes a realistic and practical implementation of RL algorithm in hardware
- Key Results:

SAFARI

- Evaluated using a wide range of workloads from SPEC CPU, PARSEC, Ligra, Cloudsuite
- Outperforms best prefetcher by **3.4%**, **7.7%** and **17%** in 1/4/bw-constrained cores
- Customizing Pythia leads to up to 7.8% more performance over basic Pythia across Ligra workloads

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

# **Talk Outline**

**Key Shortcomings of Prior Prefetchers** 

#### Formulating Prefetching as Reinforcement Learning

**Pythia: Overview** 

**Evaluation of Pythia and Key Results** 

Conclusion



# **Prefetching Basics**

- Predicts addresses of long-latency memory requests and fetches data before the program demands it
- Associates access patterns from past memory requests with program context or system information

#### **Program Feature** → Access Pattern

#### • Example program features

- Program counter (PC)
- Page number
- Page offset
- Cacheline delta
- ...
- Or a combination of these attributes

## **Key Shortcomings in Prior Prefetchers**

• We observe three key shortcomings that significantly limit performance benefits of prior prefetchers





### (1) Single-Feature Prefetch Prediction

 Provides good performance gains mainly on workloads where the feature-to-pattern correlation exists



### (1) Single-Feature Prefetch Prediction

• Provides good performance gains mainly on workloads where the feature-to-pattern correlation exists



# (2) Lack of Inherent System Awareness

- Little understanding of **undesirable effects** (e.g., memory bandwidth usage, cache pollution, ...)
  - Performance loss in **resource-constrained** configurations



Similar coverage

SAFARI

Lower overpredictions

Yet, lower performance

# (2) Lack of Inherent System Awareness

- Little understanding of **undesirable effects** (e.g., memory bandwidth usage, cache pollution, ...)
  - Performance loss in **resource-constrained** configurations



# (3) Lack of In-silicon Customizability

• Feature **statically** selected at design time

SAFA

- **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

# **Talk Outline**

**Key Shortcomings of Prior Prefetchers** 

#### Formulating Prefetching as Reinforcement Learning

**Pythia: Overview** 

#### **Evaluation of Pythia and Key Results**

Conclusion



# **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, ...



### What is State?



# 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

Prefetch from addres

A+offset (0)

Features of memory

(e.g., PC)

# What is Reward?

- Defines the **objective** of Pythia
- Encapsulates two metrics:

- Features of memory request to address A (e.g., PC) Processor & Memory subsystem
- 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<sub>AT</sub>)
- Accurate but late (R<sub>AL</sub>)
- Loss of coverage (R<sub>CL</sub>)
- Inaccurate
  - With low memory b/w usage (R<sub>IN</sub>-L)
  - With high memory b/w usage (R<sub>IN</sub>-H)
- No-prefetch
  - With low memory b/w usage (R<sub>NP</sub>-L)
  - With high memory b/w usage(R<sub>NP</sub>-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** 

#### **Steering Pythia's Objective via Reward Values**

Customizing reward values to make Dythic concernative towards p Strict Pythia configuration



# **Talk Outline**

**Key Shortcomings of Prior Prefetchers** 

#### Formulating Prefetching as Reinforcement Learning

**Pythia: Overview** 

#### **Evaluation of Pythia and Key Results**

#### Conclusion



# **Pythia Overview**

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



## **Architecting QVStore**



# **Architecting the QVStore**



# **Organization of QVStore**

- A monolithic two-dimensional table?
  - Indexed by state and action values
- State-space increases **exponentially** with #bits



# **Organization of QVStore**

- We partition QVStore into k vaults [k = number of features in state]
  - Each vault corresponds to one feature and stores the Qvalues of feature-action pairs



#### To retrieve Q(S,A) for each action

- Query each vault in parallel with feature and action
- Retrieve feature-action
  Q-value from each vault
- Compute MAX of all feature-action Q-values

MAX ensures the Q(S,A) is driven by the constituent feature that has highest Q( $\phi$ ,A)
# **Organization of QVStore**

- We further partition each vault into multiple planes
  - Each plane stores a partial Q-value of a feature-action pair

## To retrieve Q(φ,A) for each action

- Query each plane in parallel with hashed feature and action
- Retrieve partial featureaction Q-value from each plane
- Compute SUM of all partial feature-action Q-values



# **Organization of QVStore**

We further partition each vault into multiple planes
 Each plane stores a partial Q-value of a feature-action pair

**1. Enables sharing of partial Q-values between similar feature values, shortens prefetcher training time** 

parallel with hashed feature and action

2. Reduces chances of sharing partial Q-values across widely different feature values

feature-action Q-values

# More in the Paper

- Pipelined search operation for QVStore
- Reward assignment and **QVStore update**
- Automatic design-space exploration
  - Feature types
  - Actions
  - Reward and Hyperparameter values



# More in the Paper

• Pipelined search operation for QVStore

#### Reward assignment and OVStore undate

#### Pythia: A Customizable Hardware Prefetching Framework Using Online Reinforcement Learning

Rahul Bera1Konstantinos Kanellopoulos1Anant V. Nori2Taha Shahroodi3,1Sreenivas Subramoney2Onur Mutlu111ETH Zürich2Processor Architecture Research Labs, Intel Labs3TU Delft

- Reward and Hyperparameter values

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



# **Talk Outline**

**Key Shortcomings of Prior Prefetchers** 

## Formulating Prefetching as Reinforcement Learning

**Pythia: Overview** 

**Evaluation of Pythia and Key Results** 

Conclusion



# **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+, 3<sup>rd</sup> Prefetching Championship, 2019]
- SPP+DSPatch [Bera+, MICRO'19]
- SPP+PPF [Bhatia+, ISCA'20]

# **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:

- $R_{AT} = +20$ ;  $R_{AL} = +12$ ;  $R_{NP}$ -H=-2;  $R_{NP}$ -L=-4;
- $R_{IN}$ -H=-14;  $R_{IN}$ -L=-8;  $R_{CL}$ =-12

# **List of Evaluated Features**

#### Table 3: List of program control-flow and data-flow components used to derive the list of features for exploration

<b>Control-flow Component</b>	Data-flow Component	
<ol> <li>PC of load request</li> <li>PC-path (XOR-ed last-3 PCs)</li> <li>PC XOR-ed branch-PC</li> <li>None</li> </ol>	<ol> <li>Load cacheline address</li> <li>Page number</li> <li>Page offset</li> <li>Load address delta</li> <li>Sequence of last-4 offsets</li> <li>Sequence of last-4 deltas</li> <li>Offset XOR-ed with delta</li> <li>None</li> </ol>	





# **Basic Pythia Configuration**

#### Table 2: Basic Pythia configuration derived from our automated design-space exploration

Features	PC+Delta,Sequence of last-4 deltas		
<b>Prefetch Action List</b>	{-6,-3,-1,0,1,3,4,5,10,11,12,16,22,23,30,32}		
<b>Reward Level Values</b>	$\begin{array}{llllllllllllllllllllllllllllllllllll$		
Hyperparameters	$\alpha = 0.0065, \gamma = 0.556, \epsilon = 0.002$		



# **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**

- Reward value customization
- Strict Pythia configuration
  - Increase the rewards for no prefetching
  - Decrease the rewards for inaccurate prefetching



- Strict Pythia is more conservative in generating prefetch requests than the basic Pythia
- Evaluate on all Ligra graph processing workloads

## **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)



# More in the Paper

- Performance comparison with **unseen traces** 
  - Pythia provides equally high performance benefits
- Comparison against multi-level prefetchers
  - Pythia outperforms prior best multi-level prefetchers
- Understanding Pythia's learning with a case study
  - We reason towards the correctness of Pythia's decision
- Performance sensitivity towards different features and hyperparameter values
- Detailed single-core and four-core performance

## **Performance on Previously-Unseen Workloads**

- Evaluated with 500 traces from value prediction championship
  - No prefetcher has been trained on these traces



Pythia outperforms MLOP and Bingo by 8.3% and 3.5% in single-core

And 9.7% and 5.4% in four-core





# More in the Paper

Performance comparison with unseen traces
 Pythia provides equally high performance benefits

#### Comparison against multi-level prefetchers

#### Pythia: A Customizable Hardware Prefetching Framework Using Online Reinforcement Learning

Rahul Bera <sup>1</sup>	Konst	antinos Kanellopoulos <sup>1</sup>	Anant V. Nori <sup>2</sup>	Taha Shahroodi <sup>3,1</sup>
		Sreenivas Subramoney <sup>2</sup>	Onur Mutlu <sup>1</sup>	
<sup>1</sup> ETH 2	Zürich	<sup>2</sup> Processor Architecture Rese	arch Labs, Intel Labs	<sup>3</sup> TU Delft

- Performance sensitivity towards different features and hyperp <u>https://arxiv.org/pdf/2109.12021.pdf</u>
- Detailed single-core and four-core performance

# **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

SAFAR

CMU-SAFARI/Pythia Public		<ul> <li>Unwat</li> </ul>	ch 👻 3	☆ Star	9	😵 Fork	2
<> Code 🕢 Issues 🕴 Pull reques	ts 🕞 Actions 🛄 Projects 🖽 Wiki 😲 Security	∠ Insights 龄 Set	ttings				
🐉 master 👻 1 branch 🛯 🛇 5 tags	Go to file A	dd file ▼ Code ▼	About				ş
rahulbera Github pages documentatio	n 🗸 diefc65 7 hours	ago 🕚 40 commits	A custo framew learning	mizable h ork using a as descr	ardwar online ibed in	e prefetcl reinforcer the MICF	hing mer RO
branch	Initial commit for MICRO'21 artifact evaluation	2 months ago	2021 pa	aper by Be	ra and		
config	Initial commit for MICRO'21 artifact evaluation	2 months ago	Kanello	poulos et	al.		
docs	Github pages documentation	7 hours ago	ି arxiv	.org/pdf/2	109.120	21.pdf	
experiments	Added chart visualization in Excel template	2 months ago	machin	e-learning			
inc	Updated README	8 days ago	reinford	ement-learn	ing	refetcher	
prefetcher	Initial commit for MICRO'21 artifact evaluation	2 months ago	microar	chitecture	cache	-replaceme	ent
replacement	Initial commit for MICRO'21 artifact evaluation	2 months ago	branch	-predictor	champ	sim-simula	tor
scripts	Added md5 checksum for all artifact traces to verify download	2 months ago	champs	sim-tracer			
src	Initial commit for MICRO'21 artifact evaluation	2 months ago	🛱 Rea	dme			
tracer	Initial commit for MICRO'21 artifact evaluation	2 months ago	ă <u>t</u> ă Viev	v license			
🗅 .gitignore	Initial commit for MICRO'21 artifact evaluation	2 months ago	Cite کہ	this reposi	tory 👻		
CITATION.cff	Added citation file	8 days ago					
	Updated LICENSE	2 months ago	Release	es 5			
LICENSE.champsim	Initial commit for MICRO'21 artifact evaluation	2 months ago	♡ v1.3 21 da	Latest			

93

# **Talk Outline**

**Key Shortcomings of Prior Prefetchers** 

## Formulating Prefetching as Reinforcement Learning

**Pythia: Overview** 

## **Evaluation of Pythia and Key Results**

### Conclusion



# **Executive Summary**

- Background: Prefetchers predict addresses of future memory requests by associating memory access patterns with pieces of program and system information (called feature)
- **Problem**: Three key shortcomings of prior prefetchers:
  - Predict mainly using a single program feature
  - Lack inherent system awareness (e.g., memory bandwidth usage)
  - Lack in-silicon customizability
- Goal: Design a prefetching framework that:
  - Learns from multiple features and inherent system-level feedback
  - Can be customized in silicon to use different features and/or prefetching objectives
- Contribution: Pythia, which formulates prefetching as reinforcement learning problem
  - Takes adaptive prefetch decisions using multiple features and system-level feedback
  - Can be customized in silicon for target workloads via simple configuration registers
  - Proposes a realistic and practical implementation of RL algorithm in hardware
- Key Results:

SAFARI

- Evaluated using a wide range of workloads from SPEC CPU, PARSEC, Ligra, Cloudsuite
- Outperforms best prefetcher by **3.4%**, **7.7%** and **17%** in 1/4/bw-constrained cores
- Customizing Pythia leads to up to 7.8% more performance over basic Pythia across Ligra workloads

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



# 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





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

# **Pythia Discussion**

#### • FAQs

- <u>Why RL?</u>
- What about large page?
- What's the prefetch degree?
- <u>Can customization happen during</u> <u>workload execution?</u>
- Can runtime mixing create problem?

#### Simulation and Methodology

- Basic Pythia configuration
- System parameters
- Configuration of prefetchers
- Evaluated workloads
- Feature selection

- Detailed Design
  - Reward structure
  - Design overview
  - **QVStore Organization**

#### More Results

- <u>Comparison against other adaptive</u> <u>prefetchers</u>
- Comparison against Context prefetcher
- Feature combination sensitivity
- <u>Hyperparameter sensitivity</u>
- Comparison with multi-level prefetchers
- Performance in unseen workloads
- Single-core s-curve
- Four-core s-curve
- Detailed performance analysis
- Benefit of bandwidth awareness
- Case study
- Customizing rewards
- Customizing features

# 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 <u>54th International Symposium on Microarchitecture</u> (<i>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 Bera<sup>1</sup> Konstantinos Kanellopoulos<sup>1</sup>

Anant V. Nori<sup>2</sup> T Onur Mutlu<sup>1</sup>

ori<sup>2</sup> Taha Shahroodi<sup>3,1</sup>

<sup>1</sup>ETH Zürich <sup>2</sup>Processor Architecture Research Labs, Intel Labs <sup>3</sup>TU Delft

Sreenivas Subramoney<sup>2</sup>

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

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

<sup>1</sup>ETH Zürich <sup>2</sup>Intel Processor Architecture Research Lab <sup>3</sup>LIRMM, 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







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

## **The Key Problem**



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



## **Traditional Solutions**



# ၂ Employ sophisticated prefetchers

# Increase size of on-chip caches

# Key Observation 1





*# off-chip loads without any prefetcher* 

## **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

## **Key Contribution**

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



### 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:



#### 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



Low 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





### **Predicting using POPET**

• Uses simple table lookups, addition, and comparison









# **Training POPET**



# **Features Used in Hermes**

### Table 1: The initial set of program features used for automated feature selection. $\oplus$ represents a bitwise XOR operation.

Features without control-flow information	Features with control-flow information	
	8. Load PC	
1. Load virtual address	9. PC $\oplus$ load virtual address	
2. Virtual page number	10. $PC \oplus virtual page number$	
3. Cacheline offset in page	11. $PC \oplus cacheline offset$	
4. First access	12. PC + first access	
5. Cacheline offset + first access	13. PC $\oplus$ byte offset	
6. Byte offset in cacheline	14. $PC \oplus word offset$	
7. Word offset in cacheline	15. Last-4 load PCs	
	16. Last-4 PCs	

#### **Table 2: POPET configuration parameters**

Selected features	<ul> <li>PC ⊕ cacheline offset</li> <li>PC ⊕ byte offset</li> <li>PC + first access</li> <li>Cacheline offset + first access</li> <li>Last-4 load PCs</li> </ul>
Threshold values	$ au_{act} = -18, T_N = -35, T_P = 40$

# **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

#### **Latency Configuration**



#### Cache round-trip latency

- L1-D: 5 cycles
- L2: **15** cycles
- LLC: **55** cycles
- Hermes request issue latency (incurred after address translation)

Depends on

Interconnect between POPET and MC



### **Single-Core Performance Improvement**



Hermes provides nearly 90% 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



# **Effect of Cache Hierarchy Access Latency**



Hermes can provide even higher performance benefit in future processors with bigger and slower on-chip caches

On-chip cache hierarchy access latency (in processor cycles)



### **Effect of ROB Size**





### **Effect of LLC Size**





### Accuracy and Coverage with Different Prefetchers



POPET's accuracy and coverage increases significantly in absence of a data prefetcher



### **Overhead of Hermes**



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

# More in the Paper

- Performance sensitivity to:
  - Cache hierarchy access latency
  - Hermes request issue latency
  - Activation threshold
  - ROB size (in extended version on arXiv)
  - LLC size (in extended version on arXiv)
- Accuracy, coverage, and performance analysis against HMP and TTP
- Understanding usefulness of each program feature
- Effect on stall cycle reduction
- Performance analysis on an eight-core system

# More in the Paper

#### Performance sensitivity to:



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

Rahul Bera<sup>1</sup> Konstantinos Kanellopoulos<sup>1</sup> Shankar Balachandran<sup>2</sup> David Novo<sup>3</sup> Ataberk Olgun<sup>1</sup> Mohammad Sadrosadati<sup>1</sup> Onur Mutlu<sup>1</sup>

<sup>1</sup>ETH Zürich <sup>2</sup>Intel Processor Architecture Research Lab <sup>3</sup>LIRMM, 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

#### \_

# To Summarize...

# Summary

### Hermes enables 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



# Summary

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



# Hermes is Open Sourced





# All workload traces





# 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]

# off-chip predictors

riment f	iles and rollup script	6 davs ago
	Predictor type	Description
	Base	Always NO
	Basic	Simple confidence counter-based threshold
iement	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
S.CSV	HMP-GSkew	Hit-miss predictor with GSkew prediction
nple p	HMP-Ensemble	Hit-miss predictor with all three types combined
	TTP	Tag-tracking based predictor
	Perc	Perceptron-based OCP used in this paper

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

# **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;
        uint64_t seed;
13
         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
26
         virtual void reset_stats();
         virtual void train(ooo model instr *arch instr, uint32 t data index, LSQ ENTRY *lq entry);
27
28
         virtual bool predict(ooo model instr *arch instr, uint32 t data index, LSQ ENTRY *lq entry);
29
    };
30
31
    #endif /* OFFCHIP PRED BASE H */
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...









# 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







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

# **Hermes Discussion**

#### • FAQs

- What are the selected set of program features?
- <u>Can you provide some intuition on why these</u> <u>features work?</u>
- What happens in case of a misprediction?
- <u>What's the performance headroom for off-chip</u> <u>prediction?</u>
- <u>Do you see a variance of different features in final</u> <u>prediction accuracy?</u>

#### Simulation Methodology

- System parameters
- Evaluated workloads

- More Results
  - Percentage of off-chip requests
  - <u>Reduction in stall cycles by reducing the</u> <u>critical path</u>
  - Fraction of off-chip load requests
  - Accuracy and coverage of POPET
  - Effect of different features
  - Are all features required?
  - <u>1C performance</u>
  - <u>1C performance line graph</u>
  - <u>1C performance against prior predictors</u>
  - Effect on stall cycles
  - <u>8C performance</u>
  - Sensitivity:
    - Hermes request issue latency
    - <u>Cache hierarchy access latency</u>
    - Activation threshold
    - <u>ROB size</u>
    - LLC size
  - Power overhead
  - Accuracy without prefetcher
  - <u>Main memory request overhead with</u> <u>different prefetchers</u>

# Hermes Paper [MICRO 2022]

 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

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

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

# 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

144
# **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

#### SAFARI

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

# **Talk Outline**

**Key Shortcomings of Prior Data Placement Techniques** 

### Formulating Data Placement as Reinforcement Learning

Sibyl: Overview

**Evaluation of Sibyl and Key Results** 

Conclusion



# **Hybrid Storage System Basics**

### Address Space (Application/File System View)



# **Hybrid Storage System Basics**

Logical Address Space (Application/File System View)





### **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 Adaptivity (1/2)

### **Workload Changes**

Prior data placement techniques consider only a few workload characteristics that are statically tuned



# Lack of Adaptivity (2/2)

**Changes in Device Types and Configurations** 

Do not consider **underlying storage device characteristics** (e.g., changes in the level asymmetry in read/write latencies, garbage collection)



# 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

# **Talk Outline**

**Key Shortcomings of Prior Data Placement Techniques** 

### Formulating Data Placement as Reinforcement Learning

Sibyl: Overview

### **Evaluation of Sibyl and Key Results**

Conclusion



### **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

• 6-dimensional vector of state features

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

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





### **Selected State Attributes**

### Table 1: State features used by Sibyl

Feature	Description	# of bins	Encoding (bits)
size <sub>t</sub>	Size of the requested page (in pages)	8	8
$type_t$	Type of the current request (read/write)	2	4
intr <sub>t</sub>	Access interval of the requested page	64	8
$cnt_t$	Access count of the requested page	64	8
$cap_t$	Remaining capacity in the fast storage device	8	8
curr <sub>t</sub>	Current placement of the requested page (fast/slow)	2	4

# 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
  SAFARI

# **Reward Function**

**Reward.** After every data placement decision at time-step<sup>4</sup> t, Sibyl gets a reward from the environment at time-step t + 1 that acts as a feedback to Sibyl's previous action. To achieve Sibyl's performance goal, we craft the reward function R as follows:

 $R = \begin{cases} \frac{1}{L_t} & \text{if no eviction of a page from the} \\ fast storage to the slow storage \\ max(0, \frac{1}{L_t} - R_p) & \text{in case of eviction} \end{cases}$ (1)

where  $L_t$  and  $R_p$  represent the last served request latency and eviction penalty, respectively. If the fast storage is running out of free space, there might be evictions in the background from the fast <sup>4</sup>In HSS, a time-step is defined as a new storage request.

storage to the slow storage. Therefore, we add an eviction penalty  $(R_p)$  to guide Sibyl to place only performance-critical pages in the fast storage. We empirically select  $R_p$  to be equal to  $0.001 \times L_e$  ( $L_e$  is the time spent in evicting pages from the fast storage to the slow storage), which prevents the agent from aggressively placing all requests into the fast storage device.

# 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 learns to proactively evict or promote a page

# **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**



# **Training and Inference Networks**

 Training and inference networks allow parallel execution

- Observation vector as the input
- Probablility distribution of the actions (place data in the fast or the slow storage) **Fully-connected** layer (30 neurons) swish activation **Fully-connected** layer (20 neurons) **Observation vector**
- Produces probability distribution of Q-values
- <size<sub>t</sub>, type<sub>t</sub>, intr<sub>t</sub>, cnt<sub>t</sub>, cap<sub>t</sub>, curr<sub>t</sub>>

### **RL-Based Data Placement Algorithm**

### Algorithm 1 Sibyl's reinforcement learning-based data placement algorithm

1:	Intialize: the experience but	ffer EB to capacity $e_{EB}$			
2:	Intialize: the training netw	ork with random weights $ heta$			
3:	3: Intialize: the inference network with random weights $\hat{ heta}$				
4:	4: Intialize: the observation vector $O_t = O(s_1)$ with storage request $s_1 = \{req_t\}$ , and				
	host and storage features				
5:	5: for all storage requests do				
6:	if (rand() $< \epsilon$ ) then	$\blacktriangleright$ with probability $\epsilon$ , perform exploration			
7:	: random action $a_t$				
8:	else	▶ with probability 1- $\epsilon$ , perform exploitation			
9:	$a_t = argmax_a Q_t(a)$	▶ select action with the highest $Q_t$ value from inference network			
10:	execute $a_t$	place the requested page to fast or slow storage			
11:	if no eviction then				
12:	$r_t \leftarrow \frac{1}{L_t}$	▶ reward, given no eviction of a page from fast to slow storage			
13:	else				
14:	$r_t \leftarrow \max(0, \frac{1}{L_t} - R_p)$	▶ reward with an eviction penalty in case of an eviction			
15:	5: store experience $(O_t, a_t, r_t, O(t+1))$ in EB				
16:	if (num requests in $EB == e_{EB}$ ) then $\triangleright$ train training network when EB is full				
17:	sample random batc	hes of experiences from EB, which are in format			
$(O_j, a_j, r_j, O(j+1))$ $\triangleright$ where $O_j$ represents an observation at a time instant j from EB					
18:	Perform stochastic	gradient descent > update the training network weights			
19:	$\hat{ heta} \leftarrow  heta$	▶ copy the training network weights to the inference network			

# **Hyperparameter Tuning**

Table 2: Hyper-parameters considered for tuning



Figure 14: Sensitivity of Sibyl throughput to: (a) the discount factor ( $\gamma$ ), (b) the learning rate ( $\alpha$ ), (c) the exploration rate ( $\epsilon$ ), averaged across 14 workloads (normalized to Fast-Only)

176

SAFARI

# **Talk Outline**

**Key Shortcomings of Prior Data Placement Techniques** 

### Formulating Data Placement as Reinforcement Learning

Sibyl: Overview

### **Evaluation of Sibyl and Key Results**

### Conclusion



# **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

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

#### SAFARI



### **Performance-Oriented HSS Configuration**





### **Performance-Oriented HSS Configuration**



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

#### SAFARI



### **Performance-Oriented HSS Configuration**





### **Performance-Oriented HSS Configuration**



### 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



# 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 using 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

#### SAFARI

## 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



## **Talk Outline**

**Key Shortcomings of Prior Data Placement Techniques** 

### Formulating Data Placement as Reinforcement Learning

Sibyl: Overview

### **Evaluation of Sibyl and Key Results**

#### Conclusion



# 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





# 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

196

196

### 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	YouTube Live	M2: Basics of NAND Flash- Based SSDs m PDF m PPT	Required Recommended	
W3	19.10	You Time Live	M3: NAND Flash Read/Write Operations	Required Recommended	
W4	26.10	You the Live	M4: Processing inside NAND Flash	Required Recommended	
W5	02.11	You Time Live	M5: Advanced NAND Flash Commands & Mapping	Required Recommended	
W6	09.11	You Tube Live	M6: Processing inside Storage	Required Recommended	
W7	23.11	You the Live	M7: Address Mapping & Garbage Collection	Required Recommended	
W8	30.11	Yeu Tute Live	M8: Introduction to MQSim	Required Recommended	
W9	14.12	Ynu Tube Live	M9: Fine-Grained Mapping and Multi-Plane Operation-Aware Block Management	Required Recommended	
W10	04.01.2023	You Tube 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	Required Recommended	
			M10d: DeepSkatch: A New Machine Learning-Based Reference Search Technique for Post-Deduplication Delta Compression m PDF m PPT m Paper	Required Recommended	
W11	11.01	You the Live	M11: FLIN: Enabling Fairness and Enhancing Performance in Modern NVMe Solid State Drives m PDF m PPT	Required	
W12	25.01	You the 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



Computer Architecture - Fall 2021

# Watch on • Volutibe • UtcpOTTmvqOE?te4238=

#### Fall 2021 Lectures & Schedule

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 @(PDF) @(PPT)	Required Mentioned		
W2	07.10 Thu.	You two	L3a: Memory Systems: Challenges and Opportunities (PDF) m(PPT)	Described Suggested		HW 1 Out
			L3b: Course Info & Logistics			
			L3c: Memory Performance Attacks	Described Suggested		
	08.10 Fri.	You Tube Live	L4a: Memory Performance Attacks (PDF)  (PPT)	Described Suggested	Lab 2 Out	
			L4b: Data Retention and Memory Refresh	Described Suggested		
			L4c: RowHammer	Described Suggested		

# Machine Learning Driven Memory and Storage Systems

Onur Mutlu omutlu@gmail.com

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

19 December 2023

**IBM** Research



**ETH** zürich



# **PYTHIA BACKUP**

# **Reward Assignment to EQ Entry**

- Every action gets inserted into EQ
- Reward is assigned to each EQ entry **before or during** the eviction
- During EQ insertion: for actions
  - Not to prefetch
  - Out-of-page prefetch



# **Reward Assignment to EQ Entry**

- Every action gets inserted into EQ
- Reward is assigned to each EQ entry **before or during** the eviction
- **During EQ insertion**: for actions
  - Not to prefetch
  - Out-of-page prefetch
- During EQ residency:



2

State

Vector

Demand

Request

Look up QVStore

3

prefetch

Q-Value Store

Memory

Hierarchv

 In case address of a demand matches with address in EQ (signifies accurate prefetch)

# **Reward Assignment to EQ Entry**

- Every action gets inserted into EQ
- Reward is assigned to each EQ entry **before or during** the eviction
- During EQ insertion: for actions
  - Not to prefetch
  - Out-of-page prefetch
- During EQ residency:
  - In case address of a demand matches with address in EQ (signifies accurate prefetch)
- During EQ eviction:
  - In case no reward is assigned till eviction (signifies inaccurate prefetch)



SAFARI

## **Performance S-curve: Single-core**



## **Performance S-curve: Four-core**



FAQs

# **Pythia Discussion**

#### • FAQs

- <u>Why RL?</u>
- What about large page?
- What's the prefetch degree?
- <u>Can customization happen during</u> <u>workload execution?</u>
- Can runtime mixing create problem?

#### Simulation and Methodology

- Basic Pythia configuration
- System parameters
- Configuration of prefetchers
- Evaluated workloads
- Feature selection

- Detailed Design
  - <u>Reward structure</u>
  - Design overview
  - **QVStore Organization**

#### More Results

- <u>Comparison against other adaptive</u> <u>prefetchers</u>
- Comparison against Context prefetcher
- Feature combination sensitivity
- <u>Hyperparameter sensitivity</u>
- Comparison with multi-level prefetchers
- Performance in unseen workloads
- Single-core s-curve
- Four-core s-curve
- Detailed performance analysis
- Benefit of bandwidth awareness
- Case study
- Customizing rewards
- Customizing features

## Why RL? Why Not Supervised Learning?

- Determining the **benefits of prefetching** (i.e., whether a decision was good for performance or not) is **not easy** 
  - Depends on a complex set of metrics
    - Coverage, accuracy, timeliness
    - Effects on system: b/w usage, pollution, cross-application interference, ...
  - Dynamically-changing environmental conditions change the benefit
  - Delayed feedback due to long latency (might not receive feedback at all for inaccurate prefetches!)
- Differs from classification tasks (e.g., branch prediction)
  - Performance strongly correlates mainly to accuracy
  - Does not depend on environment
  - Bounded feedback delay

#### SAFARI



# What About Large Pages?

- Pythia's framework can be easily extended to incorporate additional prefetch actions (i.e., possible prefetch offsets for the page size)
- To decrease the storage overhead
  - Prune action space via automatic design-space exploration
  - Hash action values to retrieve Q-values





# What is the Prefetch Degree? Is It Managed by the RL Agent?

- Pythia employs a simple degree selector, separate from the RL agent
  - If the agent has selected the same prefetch action (O) multiple times in a row, Pythia increases the degree (A+2O, A+3O, ...)
  - At most degree 4
- Future works on managing degree by the RL agent





# Can the Customization Be Done While the Workload is Running?

- Certainly.
- Pythia, being an **online learning** technique, will autonomously adapt (and optimize) its policy to use the new program features or the modified reward values



### **Can Runtime Workload Mix Create an Issue?**

- We implement the bandwidth usage feedback using a counter in the memory controller. Thus Pythia already has a global view of the memory bandwidth usage that incorporates all workloads running on a multi-core system
- We evaluate a diverse set (300 of each category) of fourcore, eight-core, twelve-core random workload mixes
- Based on our evaluation, we observe that Pythia dynamically adapts itself to varying workload demands





### How does Pythia Compare Against Other Adaptive Prefetching Solutions?

- We compare Pythia against IBM POWER7<sup>[5]</sup> prefetcher
  - Adaptively selects prefetcher degree/configuration by monitoring program IPC





# How Does Pythia Compare Against the Context Prefetcher?

- Pythia widely differs from the Context Prefetcher (CP)<sup>[6]</sup> in all three aspects: state, action, and reward. The key differences are:
  - CP does not consider system-level feedback
  - CP models the agent as a contextual bandit which takes myopic prefetch decisions as compared to Pythia
  - CP requires compiler support to extract software-level features



### Pythia outperforms CP-HW by 5.3% in single-core and 7.6% in four-core system

#### SAFARI

[6] Leeor et al., ISCA'15



### How Pythia's Performance Changes With Various State Definitions You Have Swept?

• In total we evaluate state defined as any-one, any-two, and any-three combinations of 32 features



**Performance** gain ranges from 20.7% to 22.4%

**Coverage ranges from 66.2% to 71.5%** 

**Overprediction** ranges from 26.7% to 32.2%

SAFARI


### Is Pythia Sensitive to Hyperparameters?

Not setting hyperparameters can significantly impact the overall performance improvement



Changing  $\varepsilon$  from 0.002 to 1.0 drops perf. by 16%

Changing  $\alpha$  from 0.0065 to 1.0 drops perf. by 5.4%





### How Does Pythia Compare Against Commercial Multi-level Prefetchers?



Pythia outperforms IPCP [7] by 14.2% on average in 150-MTPS

DRAM MTPS (in log scale)

#### SAFARI

[6] Prakalapati et al., ISCA'20



# **Does Pythia Perform Equally Well for Unseen Workloads?**

- Evaluated with 500 traces from value prediction championship
  - No prefetcher has been trained on these traces



Pythia outperforms MLOP and Bingo by 8.3% and 3.5% in single-core

And 9.7% and 5.4% in four-core





# **Basic Pythia Configuration**

#### Table 2: Basic Pythia configuration derived from our automated design-space exploration

Features	PC+Delta,Sequence of last-4 deltas
<b>Prefetch Action List</b>	{-6,-3,-1,0,1,3,4,5,10,11,12,16,22,23,30,32}
<b>Reward Level Values</b>	$\begin{array}{llllllllllllllllllllllllllllllllllll$
Hyperparameters	$\alpha = 0.0065, \gamma = 0.556, \epsilon = 0.002$



# **System Parameters**

#### **Table 5: Simulated system parameters**

Core	1-12 cores, 4-wide OoO, 256-entry ROB, 72/56-entry LQ/SQ	
Branch Pred.	Perceptron-based [69], 20-cycle misprediction penalty	
L1/L2	Private, 32KB/256KB, 64B line, 8 way, LRU, 16/32 MSHRs, 4-	
Caches	cycle/14-cycle round-trip latency	
LLC	2MB/core, 64B line, 16 way, SHiP [133], 64 MSHRs per LLC Bank,	
	34-cycle round-trip latency	
Main Memory1C: Single channel, 1 rank/channel; 4C: Dual channel, ranks/channel; 8C: Quad channel, 2 ranks/channel; 8 banks/rank, 2400 MTPS, 64b data-bus/channel, 2KB row buf /bank, tRCD=15ns, tRP=15ns, tCAS=12.5ns		





# **Configuration of Prefetchers**

#### **Table 7: Configuration of evaluated prefetchers**

<b>SPP</b> [78]	256-entry ST, 512-entry 4-way PT, 8-entry GHR	6.2 KB
<b>Bingo</b> [27]	2KB region, 64/128/4K-entry FT/AT/PHT	46 KB
<b>MLOP</b> [111]	128-entry AMT, 500-update, 16-degree	8 KB
DSPatch [30]	Same configuration as in [30]	3.6 KB
<b>PPF</b> [32]	Same configuration as in [32]	39.3 KB
Pythia	2 features, 2 vaults, 3 planes, 16 actions	25.5 KB



## **Evaluated Workloads**

#### Table 6: Workloads used for evaluation

Suite	# Workloads	# Traces	Example Workloads
SPEC06	16	28	gcc, mcf, cactusADM, lbm,
SPEC17	12	18	gcc, mcf, pop2, fotonik3d,
PARSEC	5	11	canneal, facesim, raytrace,
Ligra	13	40	BFS, PageRank, Bellman-ford,
Cloudsuite	2 4	53	cassandra, cloud9, nutch,



# **List of Evaluated Features**

#### Table 3: List of program control-flow and data-flow components used to derive the list of features for exploration

<b>Control-flow Component</b>	Data-flow Component
<ol> <li>PC of load request</li> <li>PC-path (XOR-ed last-3 PCs)</li> <li>PC XOR-ed branch-PC</li> <li>None</li> </ol>	<ol> <li>Load cacheline address</li> <li>Page number</li> <li>Page offset</li> <li>Load address delta</li> <li>Sequence of last-4 offsets</li> <li>Sequence of last-4 deltas</li> <li>Offset XOR-ed with delta</li> <li>None</li> </ol>





# **MORE RESULTS**

# **Performance S-curve: Single-core**







## **Performance S-curve: Four-core**







### **Single-core Coverage & Overprediction**





# **Detailed Performance**



#### SAFARI

**1**229

# **Benefit of Bandwidth Awareness**







**Case Study** 



Figure 13: Q-value curves of PC+Delta feature values (a) 0x436a81+0 and (b) 0x4377c5+0 in 459.GemsFDTD-1320B.



# **Customizing Rewards**



Figure 14: Performance and main memory bandwidth usage of prefetchers in Ligra-CC.



Figure 15: Performance of the basic and strict Pythia configurations on the Ligra workload suite.



# **Customizing Features**



Figure 16: Performance of the basic and feature-optimized Pythia on the SPEC CPU2006 suite.



### **Hermes Discussion**

#### • FAQs

- What are the selected set of program features?
- <u>Can you provide some intuition on why these</u> <u>features work?</u>
- What happens in case of a misprediction?
- <u>What's the performance headroom for off-chip</u> <u>prediction?</u>
- <u>Do you see a variance of different features in final</u> prediction accuracy?

#### Simulation Methodology

- System parameters
- Evaluated workloads

- More Results
  - Percentage of off-chip requests
  - <u>Reduction in stall cycles by reducing the</u> <u>critical path</u>
  - Fraction of off-chip load requests
  - Accuracy and coverage of POPET
  - Effect of different features
  - Are all features required?
  - <u>1C performance</u>
  - <u>1C performance line graph</u>
  - <u>1C performance against prior predictors</u>
  - Effect on stall cycles
  - <u>8C performance</u>
  - Sensitivity:
    - Hermes request issue latency
    - <u>Cache hierarchy access latency</u>
    - Activation threshold
    - <u>ROB size</u>
    - LLC size
  - <u>Power overhead</u>
  - Accuracy without prefetcher
  - <u>Main memory request overhead with</u> <u>different prefetchers</u>

# HERMES BACKUP

### **Initial Set of Program Features**

Features without control-flow information	Features with control-flow information
	8. Load PC
1. Load virtual address	9. PC $\oplus$ load virtual address
2. Virtual page number	10. PC $\oplus$ virtual page number
3. Cacheline offset in page	11. PC $\oplus$ cacheline offset
4. First access	12. PC + first access
5. Cacheline offset + first access	13. PC $\oplus$ byte offset
6. Byte offset in cacheline	14. PC $\oplus$ word offset
7. Word offset in cacheline	15. Last-4 load PCs
	16. Last-4 PCs

### **Selected Set of Program Features**

### Five features

- $PC \oplus cacheline offset$
- $PC \oplus byte offset$
- PC + first access
- Last-4 load PCs

### A **binary hint** that

represents whether or not a cacheblock has been recently touched



### When A Feature Works/Does Not Work?



#### Without prefetcher

- PC + first access
- Cacheline offset + first access

#### With a simple stride prefetcher

• Cacheline offset + first access



### What Happens in case of a Misprediction?

- Two cases of mispredictions:
- Predicted on-chip but actually goes off-chip
  - Loss of performance improvement opportunity

No need for misprediction detection and recovery

#### Predicted off-chip but actually is on-chip

 Memory controller forwards the data to LLC if and only if a load to the same address have already missed LLC and arrived at the memory controller

#### No need for misprediction detection and recovery



### **Performance Headroom of Off-Chip Prediction**



### **System Parameters**

SAFARI

#### **Table 4: Simulated system parameters**

Core	1 and 8 cores, 6-wide fetch/execute/commit, 512-entry ROB, 128/72-entry LQ/SQ, Perceptron branch predictor [61] with 17-cycle misprediction penalty
L1/L2 Caches	Private, 48KB/1.25MB, 64B line, 12/20-way, 16/48 MSHRs, LRU, 5/15-cycle round-trip latency [25]
LLC	3MB/core, 64B line, 12 way, 64 MSHRs/slice, SHiP [122], 55-cycle round-trip latency [24, 25], <b>Pythia</b> prefetcher [32]
Main Memory	<b>1C:</b> 1 channel, 1 rank per channel; <b>8C:</b> 4 channels, 2 ranks per channel; 8 banks per rank, DDR4-3200 MTPS, 64b databus per channel, 2KB row buffer per bank, tRCD=12.5ns, tRP=12.5ns, tCAS=12.5ns
Hermes	Hermes-O/P: 6/18-cycle Hermes request issue latency



SAFARI

#### **Table 5: Workloads used for evaluation**

Suite	#Workloads	#Traces	Example Workloads
SPEC06	14	22	gcc, mcf, cactusADM, lbm,
SPEC17	11	23	gcc, mcf, pop2, fotonik3d,
PARSEC	4	12	canneal, facesim, raytrace,
Ligra	11	20	BFS, PageRank, Radii,
CVP	33	33	integer, floating-point, server,



### **Observation: Not All Off-Chip Loads are Prefetched**



Nearly 50% of the loads are still not prefetched

### **Observation: Not All Off-Chip Loads are Prefetched**



70% of these off-chip loads blocks ROB

### **Observation: With Large Cache Comes Longer Latency**

• On-chip cache access latency significantly contributes to the latency of an off-chip load



### **Observation: With Large Cache Comes Longer Latency**

• On-chip cache access latency significantly contributes to the latency of an off-chip load



40% of stall cycles caused by an off-chip load can be eliminated by removing on-chip cache access latency from its critical path





### What Fraction of Load Requests Goes Off-Chip?



### **Off-Chip Prediction Quality:** *Defining Metrics*





### **Off-Chip Prediction Quality:** Analysis



### **Off-Chip Prediction Quality:** Analysis



POPET provides off-chip predictions with high-accuracy and high-coverage



### **Effect of Different Features**



Combination of features provides both higher accuracy and higher coverage than any individual feature



### Are All Features Required? (1)



### No single feature individually provides highest prediction accuracy across *all* workloads


### Are All Features Required? (2)



#### No single feature individually provides highest prediction coverage across *all* workloads



### **Single-Core Performance**

SAFARI



#### Hermes in combination with Pythia outperforms Pythia alone in every workload category



### **Single-Core Performance Line Graph**





### Single-Core Performance Against Prior Predictors



**POPET provides higher performance benefit** than prior predictors

Hermes with POPET achieves nearly 90% performance improvement of the Ideal Hermes





### **Effect on Stall Cycles**



Hermes reduces off-chip load induced stall cycles on average by 16.2% (up-to 51.8%)



# **Eight-Core Performance**



# Hermes in combination with Pythia outperforms Pythia alone by **5.1%** on average



### **Effect of Hermes Request Issue Latency**



Hermes in combination with Pythia outperforms Pythia alone even with a 24-cycle Hermes request issue latency

Hermes request issue latency (in processor cycles)



### **Effect of Cache Hierarchy Access Latency**



Hermes can provide even higher performance benefit in future processors with bigger and slower on-chip caches

On-chip cache hierarchy access latency (in processor cycles)



### **Effect of Activation Threshold**



With increase in activation threshold 1. Accuracy increases 2. Coverage decreases





#### **Power Overhead**





### **Effect of ROB Size**





### **Effect of LLC Size**





#### Accuracy and Coverage with Different Prefetchers



POPET's accuracy and coverage increases significantly in absence of a data prefetcher



### **Increase in Main Memory Requests**





# **SIBYL BACKUP**

### **Performance on Unseen Workloads**



H&M (H&L) HSS configuration, Sibyl outperforms RNN-HSS and Archivist by 46.1% (54.6%) and 8.5% (44.1%), respectively

# **Performance Analysis**

#### **Performance-Oriented HSS Configuration**





## **Performance on Mixed Workloads**



# **Performance on Mixed Workloads**



# **Performance on Mixed Workloads**



#### **Performance With Different Features**



Sibyl autonomously decides which features are important to maximize the performance of the running workload

### **Sensitivity to Fast Storage Capacity**



# **Explainability Analysis**



# **Training and Inference Network**

 Training and inference networks allow parallel execution

• Observation vector as the input



- Produces probability distribution of Q-values
- <size<sub>t</sub>, type<sub>t</sub>, intr<sub>t</sub>, cnt<sub>t</sub>, cap<sub>t</sub>, curr<sub>t</sub>>