# Memory Systems and Memory-Centric Computing Topic 1.2: Memory Fundamentals

Onur Mutlu omutlu@gmail.com https://people.inf.ethz.ch/omutlu 16 July 2024 HiPEAC ACACES Summer School 2024



**ETH** zürich

### What Will You Learn in This Course?

- Memory Systems and Memory-Centric Computing
   July 15-19, 2024
- Topic 1: Memory Trends, Challenges, Opportunities, Basics
- Topic 2: Memory-Centric Computing
- Topic 3: Memory Robustness: RowHammer, RowPress & Beyond
- Topic 4: Machine Learning Driven Memory Systems
- Topic 5 (another course): Architectures for Genomics and ML
- Topic 6 (unlikely): Non-Volatile Memories and Storage
- Topic 7 (unlikely): Memory Latency, Predictability & QoS
- Major Overview Reading:
  - Mutlu et al., "A Modern Primer on Processing in Memory," Book Chapter on Emerging Computing and Devices, 2022.

# Memory Fundamentals

#### Memory (Programmer's View)



#### Abstraction: Virtual vs. Physical Memory

- Programmer sees virtual memory
  - Can assume the memory is "infinite"
- Reality: Physical memory size is much smaller than what the programmer assumes
- The system (system software + hardware, cooperatively) maps virtual memory addresses to physical memory
  - The system automatically manages the physical memory space transparently to the programmer
- + Programmer does not need to know the physical size of memory nor manage it  $\rightarrow$  A small physical memory can appear as a huge one to the programmer  $\rightarrow$  Life is easier for the programmer
- -- More complex system software and architecture

A classic example of the programmer/(micro)architect tradeoff

Requires indirection and mapping between virtual and physical address spaces

### (Physical) Memory System

- You need a larger level of storage to manage a small amount of physical memory automatically
   → Physical memory has a backing store: disk
- We will first start with the physical memory system
- For now, ignore the virtual  $\rightarrow$  physical indirection
  - However, virtual memory needs to be re-examined in modern systems, especially with memory-centric computing and heterogeneous memories and systems

### Ideal Memory

- Zero access time (latency)
- Infinite capacity
- Zero cost
- Infinite bandwidth (to support multiple accesses in parallel)
- Zero energy

#### The Problem

- Ideal memory's requirements oppose each other
- Bigger is slower
  - Bigger  $\rightarrow$  Takes longer to determine the location
- Faster is more expensive
  - Memory technology: SRAM vs. DRAM vs. Disk vs. Tape
- Higher bandwidth is more expensive
  - Need more banks, more ports, higher frequency, or faster technology

### Memory Technology: DRAM

- Dynamic random access memory
- Capacitor charge state indicates stored value
  - Whether the capacitor is charged or discharged indicates storage of 1 or 0
  - 1 capacitor
  - 1 access transistor
- Capacitor leaks through the RC path
  - DRAM cell loses charge over time
  - DRAM cell needs to be refreshed



 Read Liu et al., "RAIDR: Retention-aware Intelligent DRAM Refresh," ISCA 2012.

# Accessing a DRAM Cell



#### [Seshadri+ MICRO'17]



### Memory Technology: SRAM

- Static random access memory
- Two cross coupled inverters store a single bit
  - □ Feedback path enables the stored value to persist in the "cell"
  - 4 transistors for storage
  - 2 transistors for access



### An Aside: Phase Change Memory

- Phase change material (chalcogenide glass) exists in two states:
  - Amorphous: Low optical reflexivity and high electrical resistivity
  - Crystalline: High optical reflexivity and low electrical resistivity



PCM is resistive memory: High resistance (0), Low resistance (1)

Lee, Ipek, Mutlu, Burger, "Architecting Phase Change Memory as a Scalable DRAM Alternative," ISCA 2009.

### Reading: PCM as Main Memory: Idea in 2009

Benjamin C. Lee, Engin Ipek, Onur Mutlu, and Doug Burger,
<u>"Architecting Phase Change Memory as a Scalable DRAM Alternative"</u> *Proceedings of the <u>36th International Symposium on Computer</u>*<u>Architecture</u> (ISCA), pages 2-13, Austin, TX, June 2009. <u>Slides (pdf)</u>
One of the 13 computer architecture papers of 2009 selected as Top
Picks by IEEE Micro. Selected as a CACM Research Highlight.
2022 Persistent Impact Prize.

#### Architecting Phase Change Memory as a Scalable DRAM Alternative

Benjamin C. Lee† Engin Ipek† Onur Mutlu‡ Doug Burger†

†Computer Architecture Group Microsoft Research Redmond, WA {blee, ipek, dburger}@microsoft.com

SAFARI

‡Computer Architecture Laboratory Carnegie Mellon University Pittsburgh, PA onur@cmu.edu

# Reading: More on PCM As Main Memory

 Benjamin C. Lee, Ping Zhou, Jun Yang, Youtao Zhang, Bo Zhao, Engin Ipek, Onur Mutlu, and Doug Burger,
 "Phase Change Technology and the Future of Main Memory" <u>IEEE Micro</u>, Special Issue: Micro's Top Picks from 2009 Computer Architecture Conferences (MICRO TOP PICKS), Vol. 30, No. 1, pages 60-70, January/February 2010.

# Phase-Change Technology and the Future of Main Memory

### Intel Optane Persistent Memory (2019)

- Non-volatile main memory
- Based on 3D-XPoint Technology



#### SAFARI <u>https://www.storagereview.com/intel\_optane\_dc\_persistent\_memory\_module\_pmm</u><sup>16</sup>

### Charge vs. Resistive Memories

#### Charge Memory (e.g., DRAM, Flash)

- Write data by capturing charge Q
- Read data by detecting voltage V

#### Resistive Memory (e.g., PCM, STT-MRAM, memristors)

- Write data by pulsing current dQ/dt
- Read data by detecting resistance R

# Promising Resistive Memory Technologies

#### PCM

- Inject current to change material phase
- Resistance determined by phase

#### STT-MRAM

- Inject current to change magnet polarity
- Resistance determined by polarity
- Memristors/RRAM/ReRAM
  - Inject current to change atomic structure
  - Resistance determined by atom distance

# More on Emerging Memory Technologies



#### https://www.youtube.com/watch?v=AIE1rD9G\_YU&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=28

# More on Emerging Memory Technologies



Onur Mutlu Lectures

#### https://www.youtube.com/watch?v=pmLszWGmMGQ&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=29

EDIT VIDEO

ANALYTICS

### More on Memory Technologies



Onur Mutlu Lectures 16.3K subscribers

#### https://www.youtube.com/watch?v=pmLszWGmMGQ&list=PL5Q2soXY2Zi9xidyIgBxUz7xRPS-wisBN&index=29

EDIT VIDEO

ANALYTICS

### A Bit on Flash Memory & SSDs

Flash memory was a very "doubtful" emerging technology
 for at least two decades



By YU CAI, SAUGATA GHOSE, ERICH F. HARATSCH, YIXIN LUO, AND ONUR MUTLU

ABSTRACT | NAND flash memory is ubiquitous in everyday life today because its capacity has continuously increased and

**KEYWORDS** | Data storage systems; error recovery; fault tolerance; flash memory; reliability; solid-state drives

SAFARI

https://arxiv.org/pdf/1711.11427.pdf

### A Flash Memory SSD Controller



#### **Fig. 1.** (a) SSD system architecture, showing controller (Ctrl) and chips. (b) Detailed view of connections between controller components and chips.

Cai+, "Error Characterization, Mitigation, and Recovery in Flash Memory Based Solid State Drives," Proc. IEEE 2017.

#### https://arxiv.org/pdf/1711.11427.pdf

#### Lecture on Flash Memory & SSDs

Planar Scaling Redu Reliability Scal	ws. 3D NAND Flash Memorywwwwwwwwwwwwwwwwwwwwwwwwwwwwwwwwwww	
© ETH ZÜRICH HAUPTGEBÄUDE Computer Architecture - Lecture 26: Flash Memory and	Solid-State Drives (ETH Zürich, Fall 2020)	
Onur Mutlu Lectures 19.7K subscribers		I 43 √ 0 ↔ SHARE =+ SAVE ···· ANALYTICS EDIT VIDEO

#### SAFARI

### 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	Required Recommended	
W3	19.10	You Tube 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 Tube Live	M5: Advanced NAND Flash Commands & Mapping	Required Recommended	
W6	09.11	Yeu Tute Live	M6: Processing inside Storage	Required Recommended	
W7	23.11	You Time Live	M7: Address Mapping & Garbage Collection	Required Recommended	
W8	30.11	You Tube Live	M8: Introduction to MQSim	Required Recommended	
W9	14.12	Ynu Ture Live	M9: Fine-Grained Mapping and Multi-Plane Operation-Aware Block Management	Required Recommended	
W10	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: DeepSketch: A New Machine Learning-Based Reference Search Technique for Post-Deduplication Delta Compression maPDF miPPT maPaper	Required Recommended		
W11	11.01	You the Live	M11: FLIN: Enabling Fairness and Enhancing Performance in Modern NVMe Solid State Drives im PDF im PPT	Required	
W12	25.01	You the Premiere	M12: Flash Memory and Solid- State Drives	Recommended	

### Lectures on Memory Technologies

- Computer Architecture, Fall 2020, Lecture 15
  - Emerging Memory Technologies (ETH, Fall 2020)
  - https://www.youtube.com/watch?v=AlE1rD9G\_YU&list=PL5Q2soXY2Zi9xidyIgBxUz 7xRPS-wisBN&index=28
- Computer Architecture, Fall 2020, Lecture 16a
  - Opportunities & Challenges of Emerging Memory Tech (ETH, Fall 2020)
  - https://www.youtube.com/watch?v=pmLszWGmMGQ&list=PL5Q2soXY2Zi9xidyIgBx Uz7xRPS-wisBN&index=29
- Computer Architecture, Fall 2020, Lecture 3b
  - Memory Systems: Challenges & Opportunities (ETH, Fall 2020)
  - https://www.youtube.com/watch?v=Q2FbUxD7GHs&list=PL5Q2soXY2Zi9xidyIgBxU z7xRPS-wisBN&index=6

#### https://www.youtube.com/onurmutlulectures

### General Principle: Interleaving (Banking)

#### Interleaving (banking)

- Problem: a single monolithic large memory array takes long to access and does not enable multiple accesses in parallel
- Goal: Reduce the latency of memory array access and enable multiple accesses in parallel
- Idea: Divide a monolithic large array into multiple banks that can be accessed independently (in the same cycle or in consecutive cycles)
  - Each bank is smaller than the entire memory storage
  - Accesses to different banks can be overlapped
- A Key Issue: How do you map data to different banks? (i.e., how do you interleave data across banks?)

#### Memory Bank Organization and Operation



#### Read access sequence:

1. Decode row address & drive word-lines

2. Selected bits drive bit-lines

- Entire row read
- 3. Amplify row data

4. Decode column address & select subset of row

- Send to output
- 5. Precharge bit-lines
  - For next access

### Memory Banking

- Memory is divided into banks that can be accessed independently; banks share address and data buses (to reduce memory chip pins)
- Can start and complete one bank access per cycle
- Can sustain N concurrent accesses if all N go to different banks



DDCA 2023 Lecture 19 https://www.youtube.com/watch?v=gkMaO3yJMz0&list=PL5Q2soXY2Zi-EImKxYYY1SZuGiOAOBKaf&index=24

### Why Memory Hierarchy?

- We want both fast and large
- But we cannot achieve both with a single level of memory
- Idea: Have multiple levels of storage (progressively bigger and slower as the levels are farther from the processor) and ensure most of the data the processor needs is kept in the fast(er) level(s)

### Memory Hierarchy



# Caching Basics: Exploit Temporal Locality

- Idea: Store recently accessed data in automatically managed fast memory (called cache)
- Anticipation: the data will be accessed again soon
- Temporal locality principle
  - Recently accessed data will be again accessed in the near future
  - This is what Maurice Wilkes had in mind:
    - Wilkes, "Slave Memories and Dynamic Storage Allocation," IEEE Trans. On Electronic Computers, 1965.
    - "The use is discussed of a fast core memory of, say 32000 words as a slave to a slower core memory of, say, one million words in such a way that in practical cases the effective access time is nearer that of the fast memory than that of the slow memory."

# Caching Basics: Exploit Spatial Locality

- Idea: Store addresses adjacent to the recently accessed one in automatically managed fast memory
  - Logically divide memory into equal size blocks
  - Fetch to cache the accessed block in its entirety
- Anticipation: nearby data will be accessed soon
- Spatial locality principle
  - Nearby data in memory will be accessed in the near future
    - E.g., sequential instruction access, array traversal
  - □ This is what IBM 360/85 implemented
    - 16 Kbyte cache with 64 byte blocks
    - Liptay, "Structural aspects of the System/360 Model 85 II: the cache," IBM Systems Journal, 1968.

#### A Note on Manual vs. Automatic Management

- Manual: Programmer manages data movement across levels
  - -- too painful for programmers on substantial programs
  - "core" vs "drum" memory in the 50's
  - still done in some embedded processors (on-chip scratch pad SRAM in lieu of a cache)
- Automatic: Hardware manages data movement across levels, transparently to the programmer
  - ++ programmer's life is easier
  - simple heuristic: keep most recently used items in cache
  - the average programmer doesn't need to know about it
    - You don't need to know how big the cache is and how it works to write a "correct" program! (What if you want a "fast" program?)

#### Automatic Management in Memory Hierarchy

 Wilkes, "Slave Memories and Dynamic Storage Allocation," IEEE Trans. On Electronic Computers, 1965.

#### Slave Memories and Dynamic Storage Allocation

#### M. V. WILKES

#### Summary

The use is discussed of a fast core memory of, say, 32 000 words as a slave to a slower core memory of, say, one million words in such a way that in practical cases the effective access time is nearer that of the fast memory than that of the slow memory.

"By a slave memory I mean one which automatically accumulates to itself words that come from a slower main memory, and keeps them available for subsequent use without it being necessary for the penalty of main memory access to be incurred again."

### Historical Aside: Other Cache Papers

- Fotheringham, "Dynamic Storage Allocation in the Atlas Computer, Including an Automatic Use of a Backing Store," CACM 1961.
  - http://dl.acm.org/citation.cfm?id=366800

 Bloom, Cohen, Porter, "Considerations in the Design of a Computer with High Logic-to-Memory Speed Ratio," AIEE Gigacycle Computing Systems Winter Meeting, Jan. 1962.
## Cache in 1962 (Bloom, Cohen, Porter)



FIGURE 1. MEMORY (In this example, store is destructive read).

## A Modern Memory Hierarchy



## The DRAM Subsystem

## DRAM Subsystem Organization

- Channel
- DIMM
- Rank
- Chip
- Bank
- Row/Column



## DRAM Organization



## **DRAM** Operations



ACTIVATE (ACT): Fetch the row's content into the **row buffer** 

2 Column Access (RD/WR): Read/Write the target column and drive to I/O

PRECHARGE (PRE):Prepare the arrayfor a new ACTIVATE

## DRAM Refresh



DRAM Refresh **is the key maintenance operation** to **avoid bit flips** due to charge leakage

DRAM Refresh **activates** a row and **precharges** the bank

**Problem:** DRAM Refresh **blocks** accesses to the **whole bank / rank** 

# Page Mode DRAM

- A DRAM bank is a 2D array of cells: rows x columns
- A "DRAM row" is also called a "DRAM page"
- "Sense amplifiers" also called "row buffer"
- Each address is a <row,column> pair
- Access to a "closed row"
  - Activate command opens row (placed into row buffer)
  - Read/write command reads/writes column in the row buffer
  - Precharge command closes the row and prepares the bank for next access
- Access to an "open row"
  - No need for activate command

### The DRAM Bank Structure



## **DRAM Bank Operation**



# The DRAM Chip

- Consists of multiple banks (8-16 are common today)
- Banks share command/address/data buses
- The chip itself has a narrow interface (4-16 bits per read)
- Changing the number of banks, size of the interface (pins), whether or not command/address/data buses are shared has significant impact on DRAM system cost

### 128M x 8-bit DRAM Chip



## How Multiple Banks Help



Before: No Overlapping Assuming accesses to different DRAM rows



# Address Mapping (Single Channel)

- Single-channel system with 8-byte memory bus
  - □ 2GB memory, 8 banks, 16K rows & 2K columns per bank

### Row interleaving

Consecutive rows of memory in consecutive banks

Row (14 bits)	Bank (3 bits)	Column (11 bits)	Byte in bus (3 bits)
· · · · · · · · · · · · · · · · · · ·	( )		<b>, , ,</b>

Accesses to consecutive cache blocks serviced in a pipelined manner

### Cache block interleaving

- Consecutive cache block addresses in consecutive banks
- 64 byte cache blocks

Row (14 bits)	High Column	Bank (3 bits)	Low Col.	Byte in bus (3 bits)
	8 bits		3 bits	
Accesses to consecutive cache blocks can be serviced in parallel				

## Bank Mapping Randomization

 DRAM controller can randomize the address mapping to banks so that bank conflicts are less likely



• Reading:

Rau, "Pseudo-randomly Interleaved Memory," ISCA 1991.

# Address Mapping (Multiple Channels)

C Row (14 bits)	Bank (3 bits	s) Column (11 bits)		Byte in bus (3 bits)
Row (14 bits)	C Bank (3 bits	s) Column (11 bits)		Byte in bus (3 bits)
Row (14 bits)	Bank (3 bits)	C Column (11 bits)		Byte in bus (3 bits)
Row (14 bits)	Bank (3 bits)	Column (11 bits)	С	Byte in bus (3 bits)

### Where are consecutive cache blocks?

C Row (14 bits)	High Colum	n	Bank (3 bits	s)	Low Col		Byte in bus (3 bits)
	8 bits		3 bits				
Row (14 bits)	C High Colum	n	Bank (3 bits	s)	Low Col		Byte in bus (3 bits)
	8 bits		3 bits				
Row (14 bits)	High Column	С	Bank (3 bits	s)	Low Col		Byte in bus (3 bits)
	8 bits				3 bits		
Row (14 bits)	High Column	В	ank (3 bits)	С	Low Col		Byte in bus (3 bits)
	8 bits	3 bits					
Row (14 bits)	High Column	В	ank (3 bits)	L	.ow Col.	С	Byte in bus (3 bits)
	8 bits				3 bits		

### Interaction with Virtual → Physical Mapping

 Operating System influences where an address maps to in DRAM



- Operating system can influence which bank/channel/rank a virtual page is mapped to.
- It can perform page coloring to
  - Minimize bank conflicts
  - Minimize inter-application interference [Muralidhara+ MICRO'11]
  - Minimize latency in the network [Das+ HPCA'13]

### Memory Channel Partitioning

 Sai Prashanth Muralidhara, Lavanya Subramanian, Onur Mutlu, Mahmut Kandemir, and Thomas Moscibroda,
"Reducing Memory Interference in Multicore Systems via Application-Aware Memory Channel Partitioning"
Proceedings of the <u>44th International Symposium on</u> Microarchitecture (MICRO), Porto Alegre, Brazil, December 2011. <u>Slides (pptx)</u>

### Reducing Memory Interference in Multicore Systems via Application-Aware Memory Channel Partitioning

Sai Prashanth Muralidhara Pennsylvania State University smuralid@cse.psu.edu Lavanya Subramanian Carnegie Mellon University Isubrama@ece.cmu.edu Onur Mutlu Carnegie Mellon University onur@cmu.edu

Mahmut Kandemir Pennsylvania State University kandemir@cse.psu.edu Thomas Moscibroda Microsoft Research Asia moscitho@microsoft.com

# Application-to-Core Mapping

 Reetuparna Das, Rachata Ausavarungnirun, Onur Mutlu, Akhilesh Kumar, and Mani Azimi,
"Application-to-Core Mapping Policies to Reduce Memory System Interference in Multi-Core Systems"
Proceedings of the <u>19th International Symposium on High-Performance</u> Computer Architecture (HPCA), Shenzhen, China, February 2013. Slides (pptx)

#### Application-to-Core Mapping Policies to Reduce Memory System Interference in Multi-Core Systems

Reetuparna Das\* Rachata Ausavarungnirun† Onur Mutlu† Akhilesh Kumar‡ Mani Azimi‡ University of Michigan\* Carnegie Mellon University† Intel Labs‡

## On Reducing Bank Conflicts

Read Sections 1 through 4 of:

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



Figure 1. DRAM bank organization

### Subarray Level Parallelism

 Yoongu Kim, Vivek Seshadri, Donghyuk Lee, Jamie Liu, and Onur Mutlu, <u>"A Case for Exploiting Subarray-Level Parallelism (SALP) in</u> <u>DRAM"</u> *Proceedings of the <u>39th International Symposium on Computer</u> <u>Architecture</u> (ISCA), Portland, OR, June 2012. <u>Slides (pptx)</u>* 

#### A Case for Exploiting Subarray-Level Parallelism (SALP) in DRAM

Yoongu Kim Vivek Seshadri Donghyuk Lee Jamie Liu Onur Mutlu Carnegie Mellon University

### Generalized Memory Structure



### Generalized Memory Structure



Kim+, "A Case for Exploiting Subarray-Level Parallelism in DRAM," ISCA 2012. Lee+, "Decoupled Direct Memory Access," PACT 2015.

## The DRAM Subsystem A Top-Down View

## DRAM Subsystem Organization

- Channel
- DIMM
- Rank
- Chip
- Bank
- Row/Column



## The DRAM Subsystem



**DIMM** (Dual in-line memory module)

## Breaking down a DIMM (module)



## Breaking down a DIMM (module)



## Rank



### Breaking down a Rank



### Breaking down a Chip



### Breaking down a Bank



# Digging Deeper: DRAM Bank Operation



### A DRAM Bank Internally Has Sub-Banks



Figure 1. DRAM bank organization

### Another View of a DRAM Bank



Seshadri+, "In-DRAM Bulk Bitwise Execution Engine," ADCOM 2020.

### More on DRAM Basics & Organization

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

See Section 2 for comprehensive DRAM Background

### In-DRAM Bulk Bitwise Execution Engine

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

https://arxiv.org/pdf/1905.09822.pdf
### DRAM Subsystem Organization

- Channel
- DIMM
- Rank
- Chip
- Bank
- Row/Column















**Physical memory space** 



A 64B cache block takes 8 I/O cycles to transfer.

During the process, 8 columns are read sequentially.

# Simulating Memory

### Ramulator: A Fast and Extensible DRAM Simulator [IEEE Comp Arch Letters'15]

#### Ramulator Motivation

- DRAM and Memory Controller landscape is changing
- Many new and upcoming standards
- Many new controller designs
- A fast and easy-to-extend simulator is very much needed

Segment	DRAM Standards & Architectures
Commodity	DDR3 (2007) [14]; DDR4 (2012) [18]
Low-Power	LPDDR3 (2012) [17]; LPDDR4 (2014) [20]
Graphics	GDDR5 (2009) [15]
Performance	eDRAM [28], [32]; RLDRAM3 (2011) [29]
3D-Stacked	WIO (2011) [16]; WIO2 (2014) [21]; MCDRAM (2015) [13]; HBM (2013) [19]; HMC1.0 (2013) [10]; HMC1.1 (2014) [11]
Academic	SBA/SSA (2010) [38]; Staged Reads (2012) [8]; RAIDR (2012) [27]; SALP (2012) [24]; TL-DRAM (2013) [26]; RowClone (2013) [37]; Half-DRAM (2014) [39]; Row-Buffer Decoupling (2014) [33]; SARP (2014) [6]; AL-DRAM (2015) [25]
	Table 1. Landscape of DRAM-based memory

#### Ramulator

- Provides out-of-the box support for many DRAM standards:
  - DDR3/4, LPDDR3/4, GDDR5, WIO1/2, HBM, plus new proposals (SALP, AL-DRAM, TLDRAM, RowClone, and SARP)
- ~2.5X faster than fastest open-source simulator
- Modular and extensible to different standards

Simulator (clang -O3)	Cycles (10 <sup>6</sup> )		Runtime (sec.)		Req/sec (10 <sup>3</sup> )		Memory	
	Random	Stream	Random	Stream	Random	Stream	( <i>MB</i> )	
Ramulator	652	411	752	249	133	402	2.1	
DRAMSim2	645	413	2,030	876	49	114	1.2	
USIMM	661	409	1,880	750	53	133	4.5	
DrSim	647	406	18,109	12,984	6	8	1.6	
NVMain	666	413	6,881	5,023	15	20	4,230.0	

Table 3. Comparison of five simulators using two traces

#### Case Study: Comparison of DRAM Standards

Standard	d Rate (MT/s)	Timing (CL-RCD-RP)	Data-Bus (Width×Chan.)	Rank-per-Chan	BW (GB/s)
DDR3	1,600	11-11-11	$64$ -bit $\times 1$	1	11.9
DDR4	2,400	16-16-16	$64$ -bit $\times 1$	1	17.9
SALP <sup>†</sup>	1,600	11-11-11	64-bit $\times$ 1	1	11.9
LPDDR	3 1,600	12 - 15 - 15	$64$ -bit $\times 1$	1	11.9
LPDDR	4 2,400	22-22-22	$32$ -bit $\times 2^*$	1	17.9
GDDR5	[12] 6,000	18-18-18	$64$ -bit $\times 1$	1	44.7
HBM	1,000	7-7-7	$128$ -bit $\times 8^*$	1	119.2
WIO	266	7-7-7	$128$ -bit $\times 4^*$	1	15.9
WIO2	1,066	9-10-10	128-bit $\times$ 8*	1	127.2



#### Ramulator 2.0

 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

#### Optional Assignment: Ramulator 2.0

- Review the Ramulator 2.0 paper
  - Email me your review (<u>omutlu@gmail.com</u>)
- Download and run Ramulator 2.0
  - Compare DDR4, LPDDR5, HBM2 for benchmarks of your choice (provided in Ramulator repository)
  - Email me your report (<u>omutlu@gmail.com</u>)

This may help you get into memory systems research quickly

# On Virtual Memory

#### Access Control & Protection Mechanisms

- Are based on virtual memory (VM), invented in 1950s
- VM has not changed much even after decades of technology scaling and memory system improvements
- VM causes large performance problems and is responsible for large complexity, power, energy
- VM is poor for fine-grained security and access control
- VM hinders innovation in heterogeneous (accelerator) systems and new architectures (e.g., processing near data)

#### It is time to rethink virtual memory

#### Virtual Memory: Parting Thoughts

- Virtual Memory is one of the most successful examples of
  - architectural support for programmers
  - how to partition work between hardware and software
  - hardware/software cooperation
  - programmer/architect tradeoff
- Going forward: How does virtual memory fare and scale into the future? Five key trends:
  - Increasing, huge physical memory sizes (local & remote)
  - Hybrid physical memory systems (DRAM + NVM + SSD)
  - Many accelerators in the system accessing physical memory
  - Virtualized systems (hypervisors, software virtualization, local and remote memories)
  - Processing in memory systems near-data accelerators

### Rethinking Virtual Memory

Nastaran Hajinazar, Pratyush Patel, Minesh Patel, Konstantinos Kanellopoulos, Saugata Ghose, Rachata Ausavarungnirun, Geraldo Francisco de Oliveira Jr., Jonathan Appavoo, Vivek Seshadri, and Onur Mutlu, "The Virtual Block Interface: A Flexible Alternative to the Conventional Virtual Memory Framework" Proceedings of the 47th International Symposium on Computer Architecture (ISCA), Virtual, June 2020. [Slides (pptx) (pdf)] [Lightning Talk Slides (pptx) (pdf)] [ARM Research Summit Poster (pptx) (pdf)] [Talk Video (26 minutes)] [Lightning Talk Video (3 minutes)] [Lecture Video (43 minutes)]

#### The Virtual Block Interface: A Flexible Alternative to the Conventional Virtual Memory Framework

Nastaran Hajinazar<sup>\*†</sup> Pratyush Patel<sup>™</sup> Minesh Patel<sup>\*</sup> Konstantinos Kanellopoulos<sup>\*</sup> Saugata Ghose<sup>‡</sup> Rachata Ausavarungnirun<sup> $\odot$ </sup> Geraldo F. Oliveira<sup>\*</sup> Jonathan Appavoo<sup> $\diamond$ </sup> Vivek Seshadri<sup> $\nabla$ </sup> Onur Mutlu<sup>\*‡</sup>

\*ETH Zürich <sup>†</sup>Simon Fraser University <sup>⋈</sup>University of Washington <sup>‡</sup>Carnegie Mellon University  $^{\odot}$ King Mongkut's University of Technology North Bangkok  $^{\diamond}$ Boston University  $^{\nabla}$ Microsoft Research India

#### https://people.inf.ethz.ch/omutlu/pub/VBI-virtual-block-interface\_isca20.pdf

#### Better Virtual Memory (I)

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 Kanellopoulos<sup>1</sup> Hong Chul Nam<sup>1</sup> F. Nisa Bostanci<sup>1</sup> Rahul Bera<sup>1</sup> Mohammad Sadrosadati<sup>1</sup> Rakesh Kumar<sup>2</sup> Davide Basilio Bartolini<sup>3</sup> Onur Mutlu<sup>1</sup> <sup>1</sup>ETH Zürich <sup>2</sup>Norwegian University of Science and Technology <sup>3</sup>Huawei Zurich Research Center

#### SAFARI

https://arxiv.org/pdf/2310.04158

#### Better Virtual Memory (II)

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> (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 Kanellopoulos<sup>1</sup> Rahul Bera<sup>1</sup> Kosta Stojiljkovic<sup>1</sup> Nisa Bostanci<sup>1</sup> Can Firtina<sup>1</sup> Rachata Ausavarungnirun<sup>2</sup> Rakesh Kumar<sup>3</sup> Nastaran Hajinazar<sup>4</sup> Mohammad Sadrosadati<sup>1</sup> Nandita Vijaykumar<sup>5</sup> Onur Mutlu<sup>1</sup>

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

#### https://arxiv.org/abs/2211.12205

## Memory Systems and Memory-Centric Computing Topic 1.2: Memory Fundamentals

Onur Mutlu omutlu@gmail.com https://people.inf.ethz.ch/omutlu 16 July 2024 HiPEAC ACACES Summer School 2024



**ETH** zürich