

Memory Systems

Fundamentals, Recent Research, Challenges, Opportunities

Lecture 4: Low-Latency Memory

Prof. Onur Mutlu

omutlu@gmail.com

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

9 October 2018

Technion Fast Course 2018

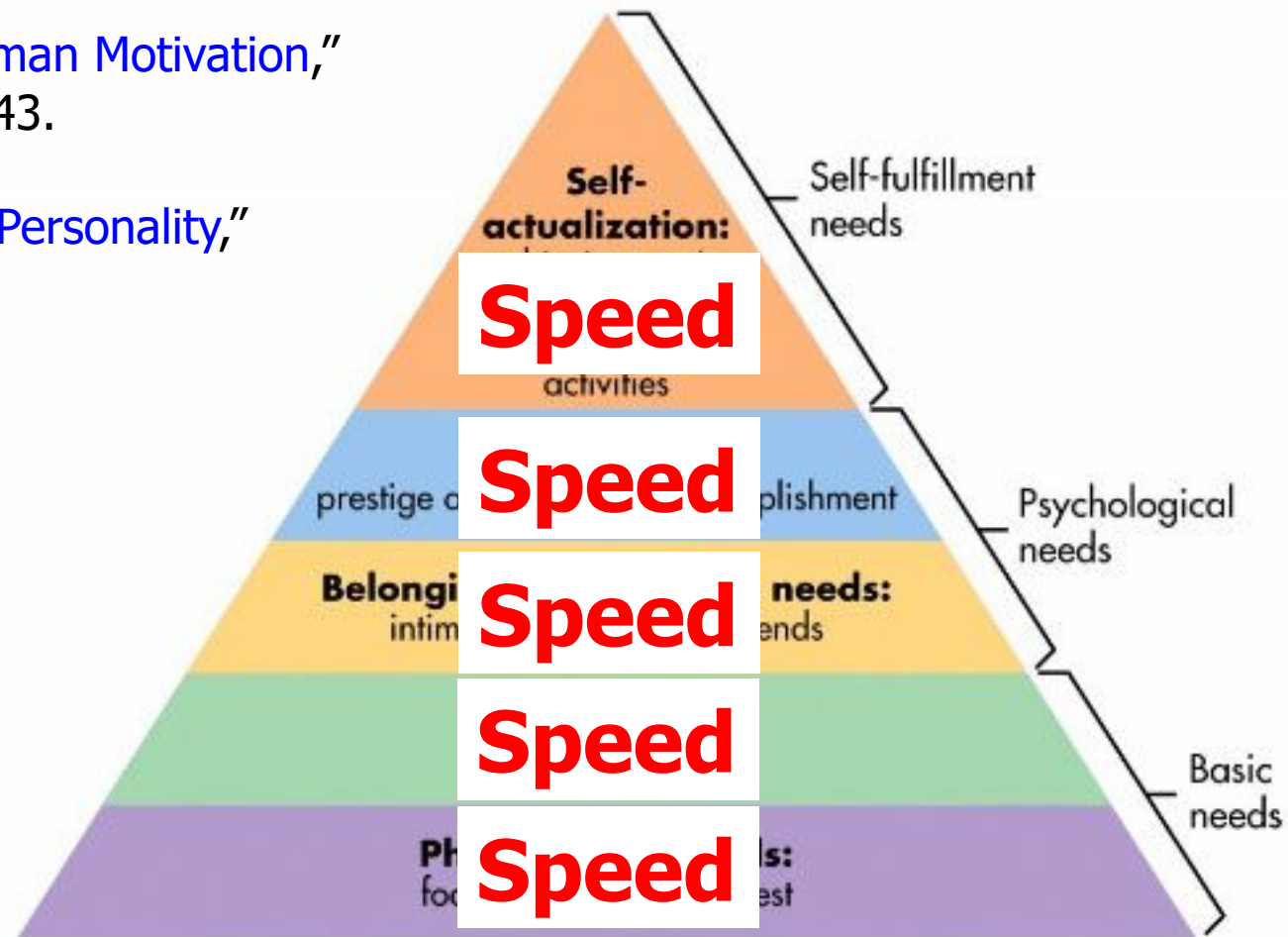
Four Key Directions

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

Maslow's Hierarchy of Needs, A Third Time

Maslow, "A Theory of Human Motivation,"
Psychological Review, 1943.

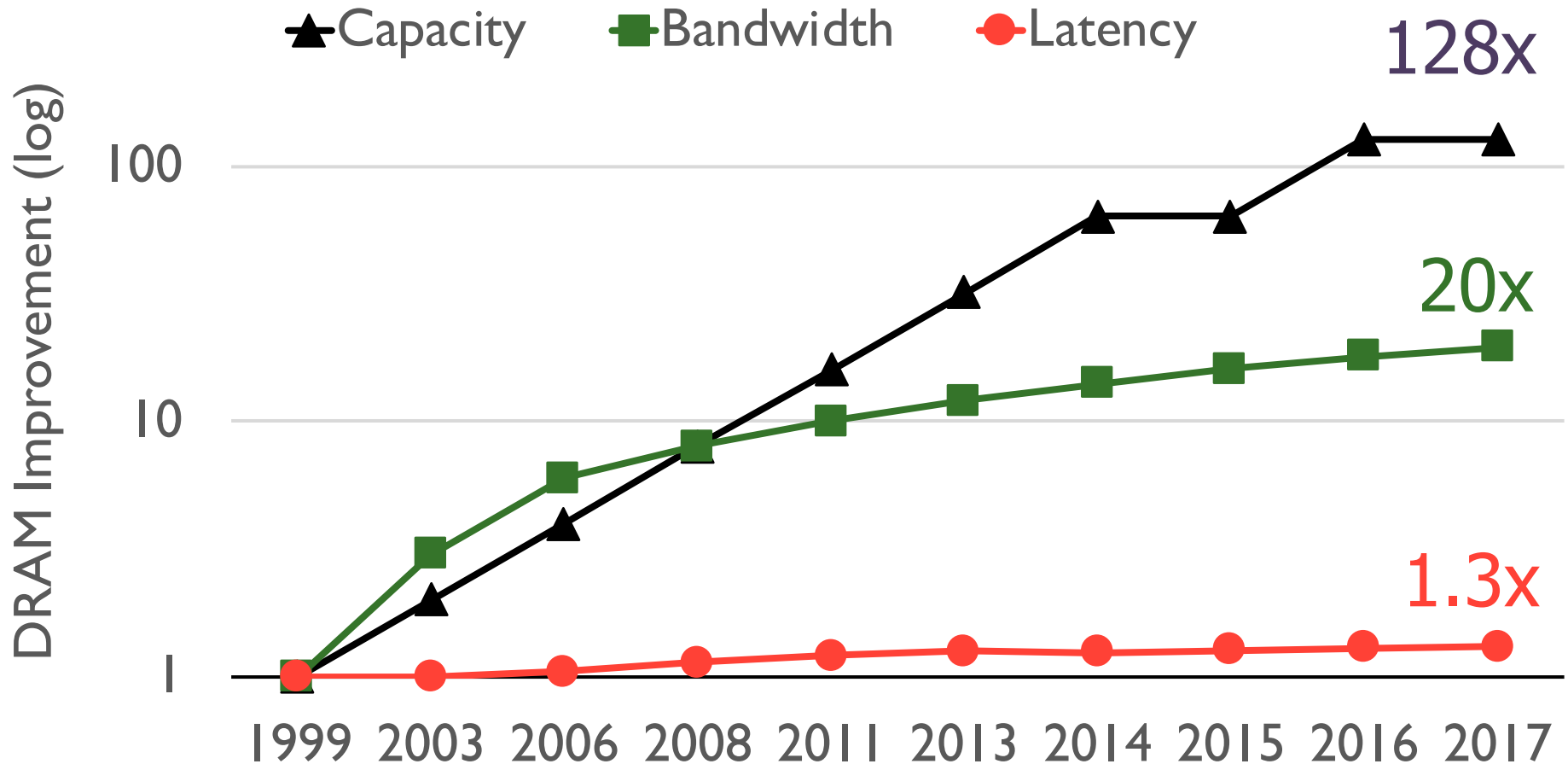
Maslow, "Motivation and Personality,"
Book, 1954-1970.



Fundamentally Low-Latency Computing Architectures

Memory Latency: Fundamental Tradeoffs

Review: Memory Latency Lags Behind



Memory latency remains almost constant

DRAM Latency Is Critical for Performance



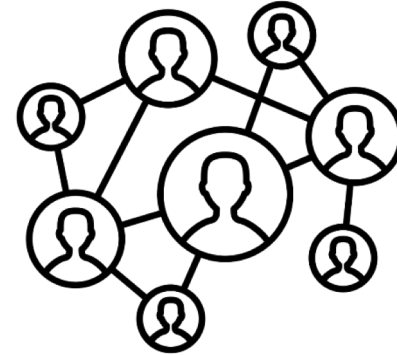
In-memory Databases

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



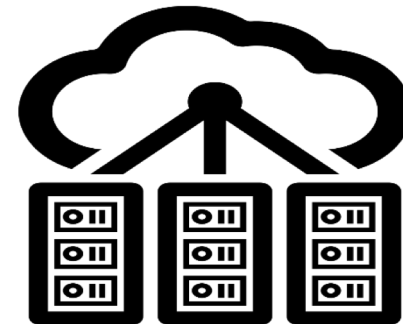
In-Memory Data Analytics

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



Graph/Tree Processing

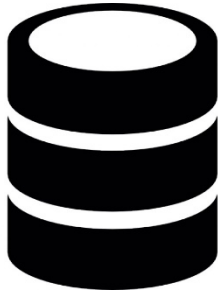
[Xu+, IISWC'12; Umuroglu+, FPL'15]



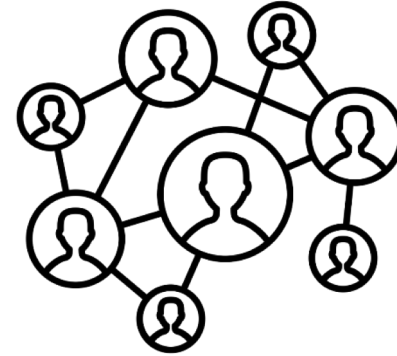
Datacenter Workloads

[Kanev+ (Google), ISCA'15]

DRAM Latency Is Critical for Performance



In-memory Databases



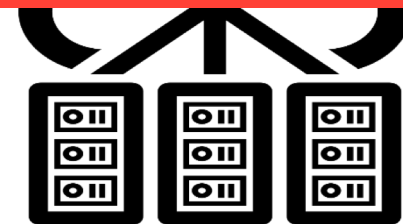
Graph/Tree Processing

Long memory latency → performance bottleneck



In-Memory Data Analytics

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



Datacenter Workloads

[Kanev+ (Google), ISCA'15]

The Memory Latency Problem

- High memory latency is a significant **limiter of system performance and energy-efficiency**
- It is becoming increasingly so with **higher memory contention** in multi-core and heterogeneous architectures
 - Exacerbating the bandwidth need
 - Exacerbating the QoS problem
- It increases **processor design complexity** due to the mechanisms incorporated to tolerate memory latency

Retrospective: Conventional Latency Tolerance Techniques

- Caching [initially by Wilkes, 1965]
 - Widely used, simple, effective, but inefficient, passive
 - Not all applications/phases exhibit temporal or spatial locality
- Prefetching [initially in IBM 360/91, 1967]

**None of These
Fundamentally Reduce
Memory Latency**

ongoing research effort

- Out-of-order execution [initially by Tomasulo, 1967]
 - **Tolerates cache misses that cannot be prefetched**
 - Requires extensive hardware resources for tolerating long latencies

Two Major Sources of Latency Inefficiency

- Modern DRAM is not designed for low latency
 - Main focus is cost-per-bit (capacity)
- Modern DRAM latency is determined by worst case conditions and worst case devices
 - Much of memory latency is unnecessary

**Our Goal: Reduce Memory Latency
at the Source of the Problem**

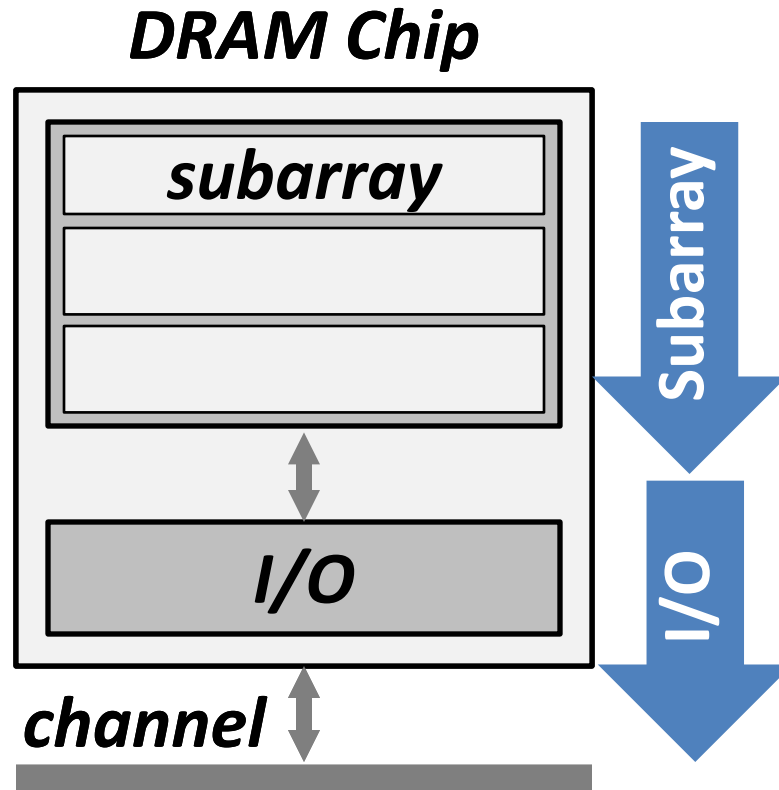
What Causes the Long Memory Latency?

Why the Long Memory Latency?

- Reason 1: Design of DRAM Micro-architecture
 - Goal: Maximize capacity/area, not minimize latency
- Reason 2: “One size fits all” approach to latency specification
 - Same latency parameters for all temperatures
 - Same latency parameters for all DRAM chips (e.g., rows)
 - Same latency parameters for all parts of a DRAM chip
 - Same latency parameters for all supply voltage levels
 - Same latency parameters for all application data
 - ...

Tiered Latency DRAM

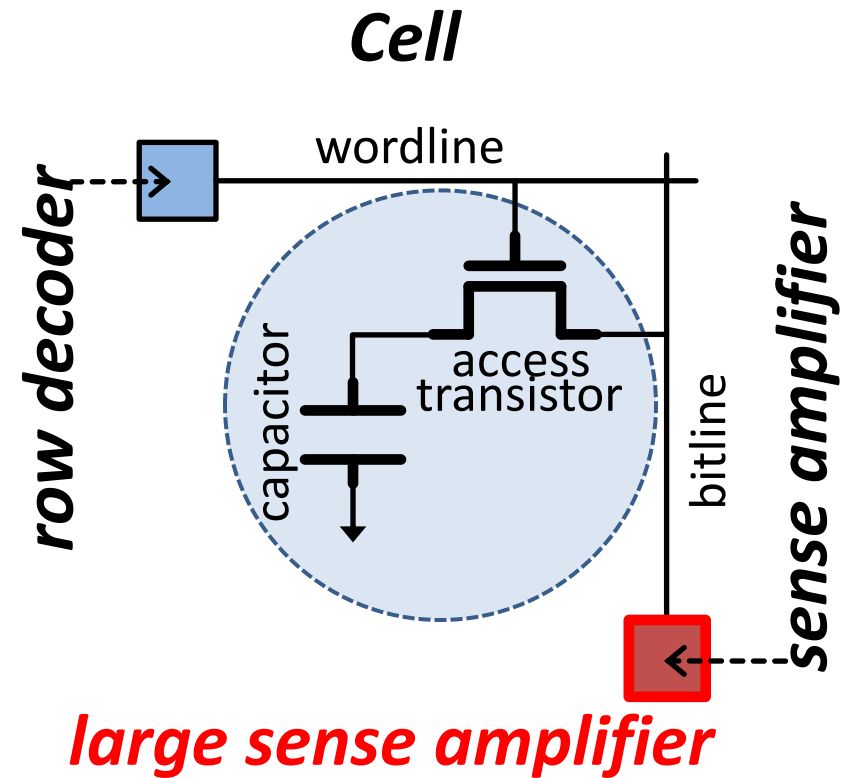
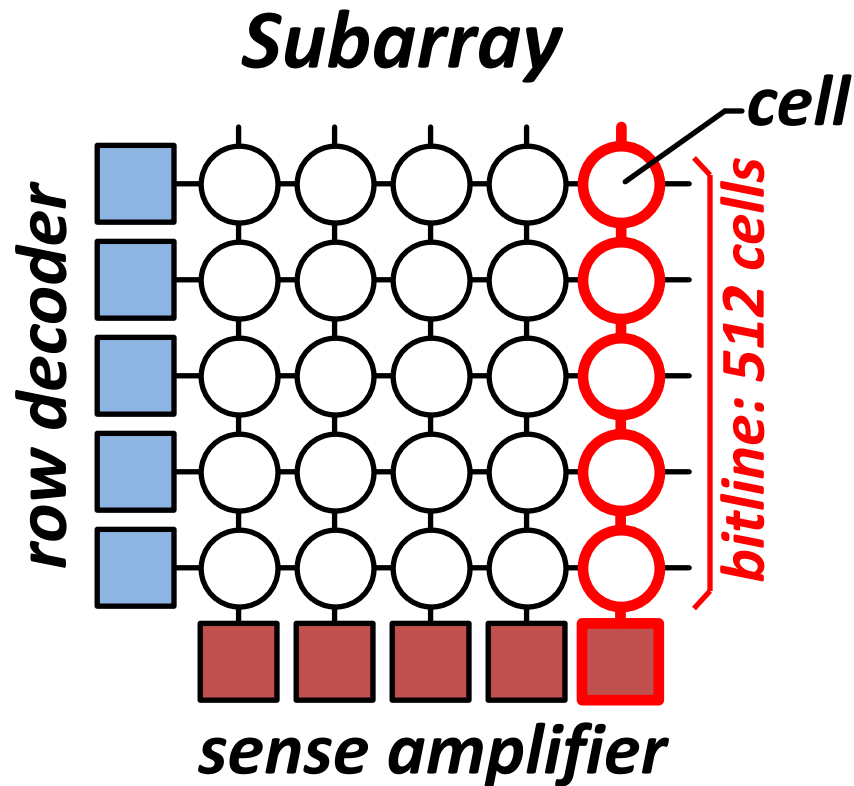
What Causes the Long Latency?



DRAM Latency = Subarray Latency + I/O Latency

Dominant

Why is the Subarray So Slow?

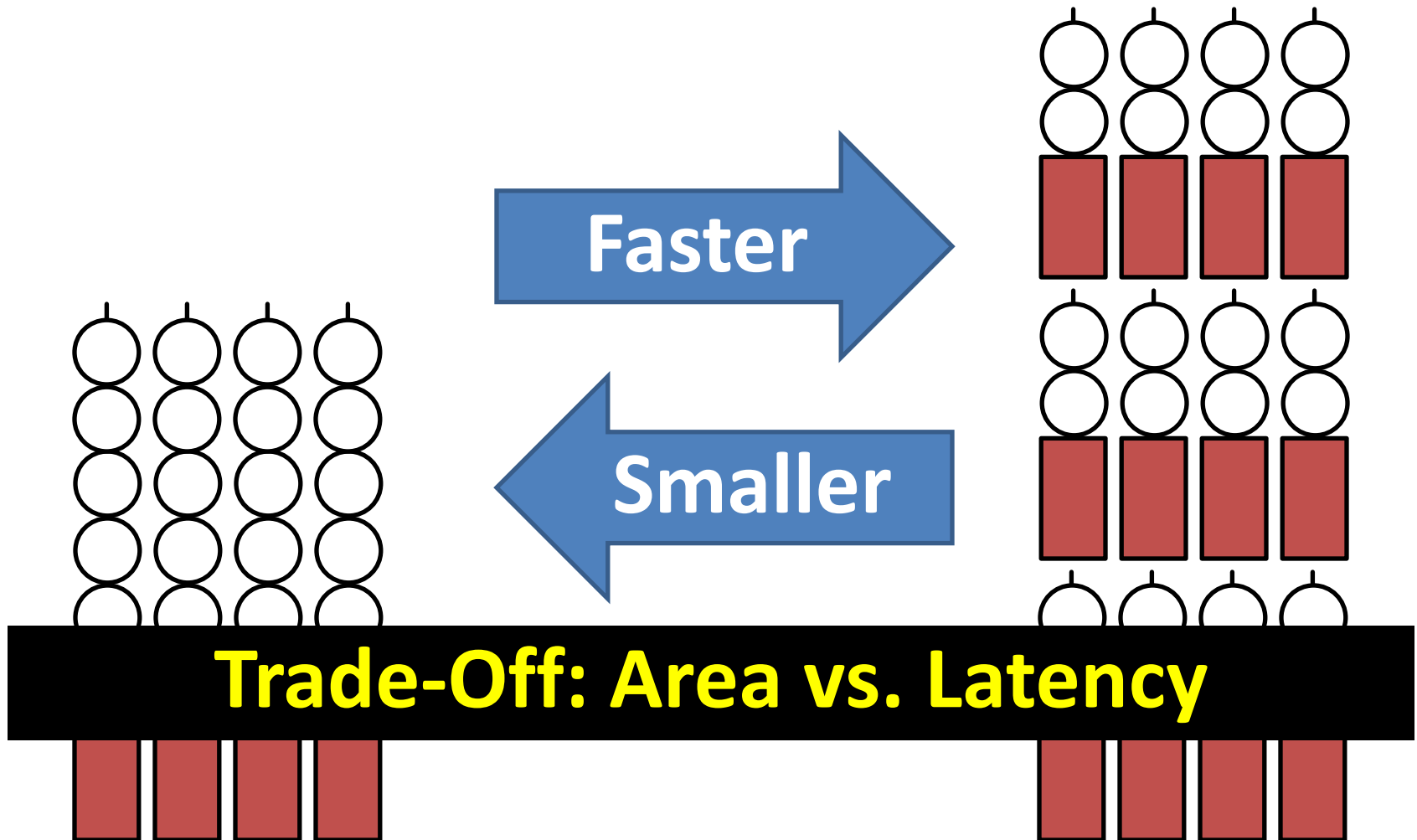


- Long bitline
 - Amortizes sense amplifier cost → Small area
 - Large bitline capacitance → High latency & power

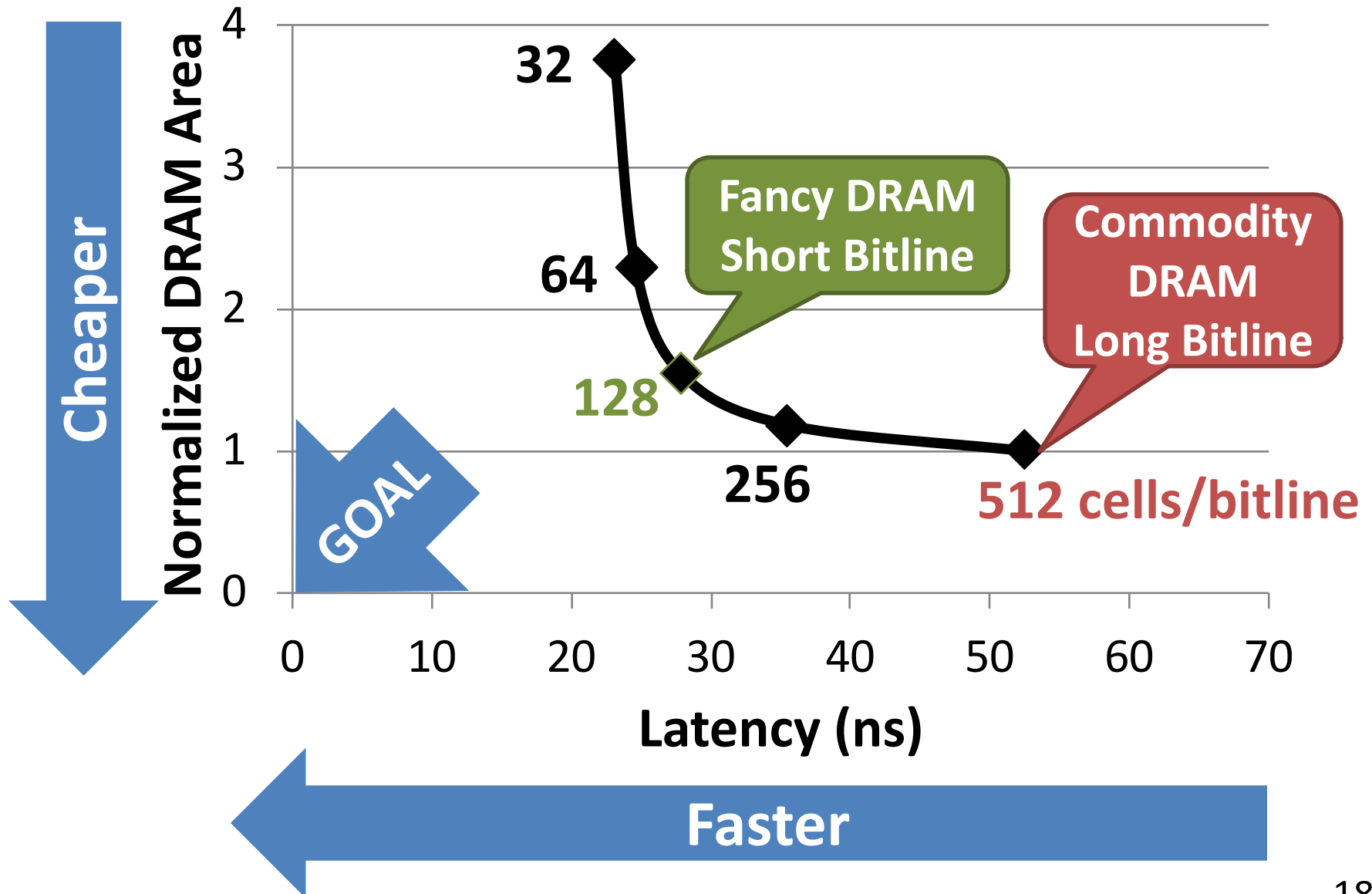
Trade-Off: Area (Die Size) vs. Latency

Long Bitline

Short Bitline



Trade-Off: Area (Die Size) vs. Latency

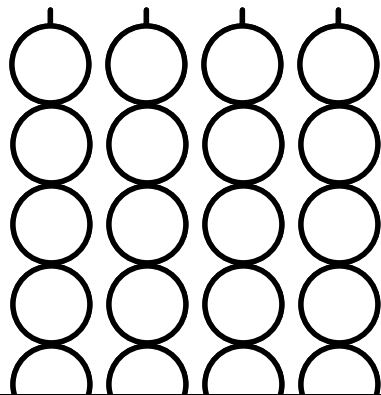


Approximating the Best of Both Worlds

Long Bitline

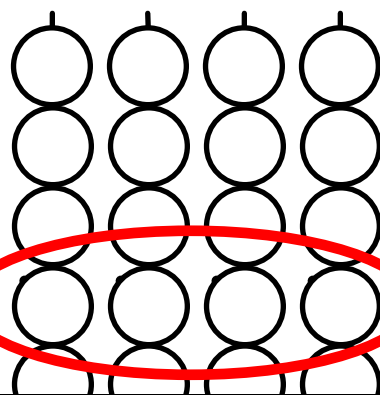
Small Area

~~High Latency~~



Need Isolation

Our Proposal

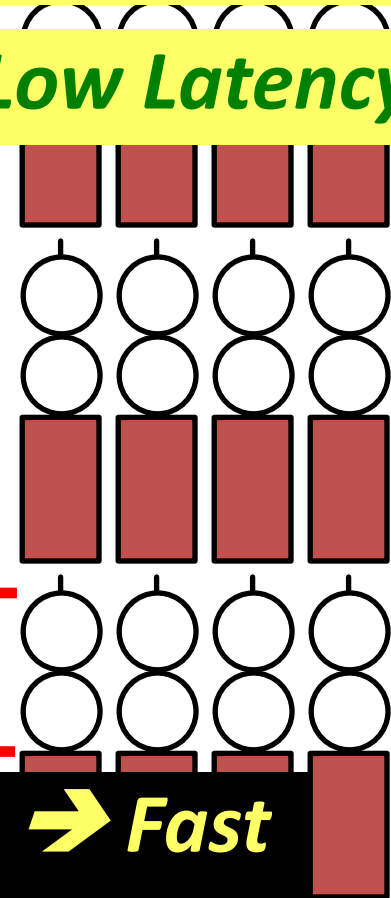


Add Isolation Transistors

Short Bitline

~~Large Area~~

Low Latency



tline → Fast

Approximating the Best of Both Worlds

Long Bitline Tiered-Latency DRAM **Short Bitline**

Small Area

Small Area

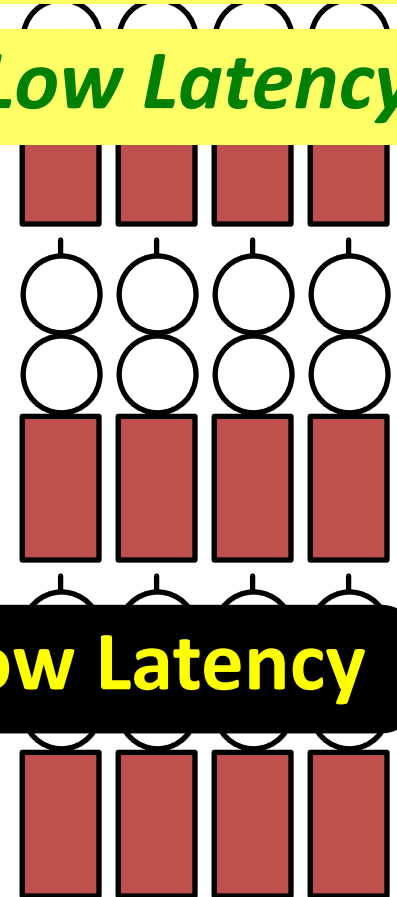
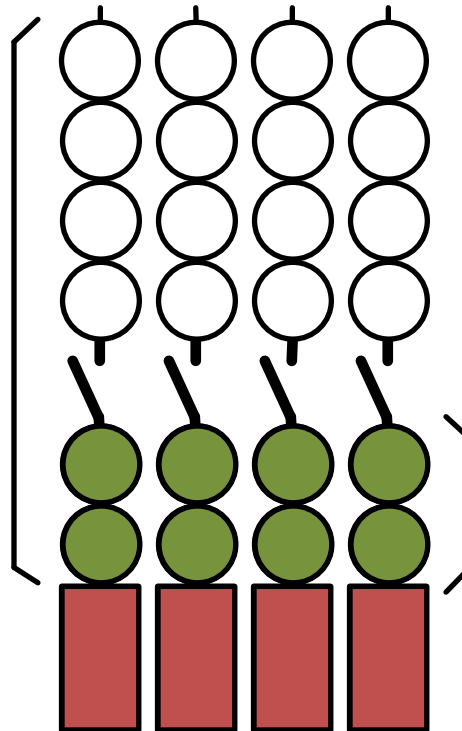
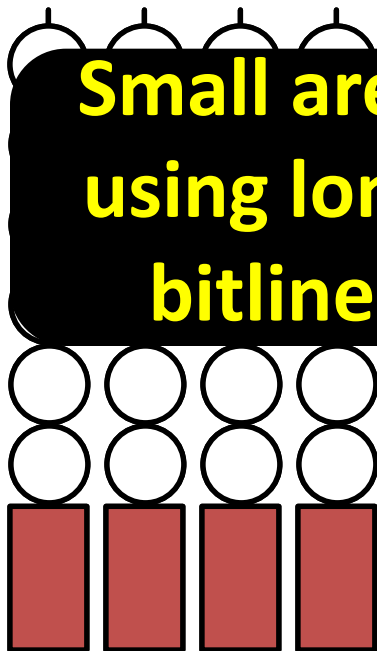
~~*Large Area*~~

~~*High Latency*~~

Low Latency

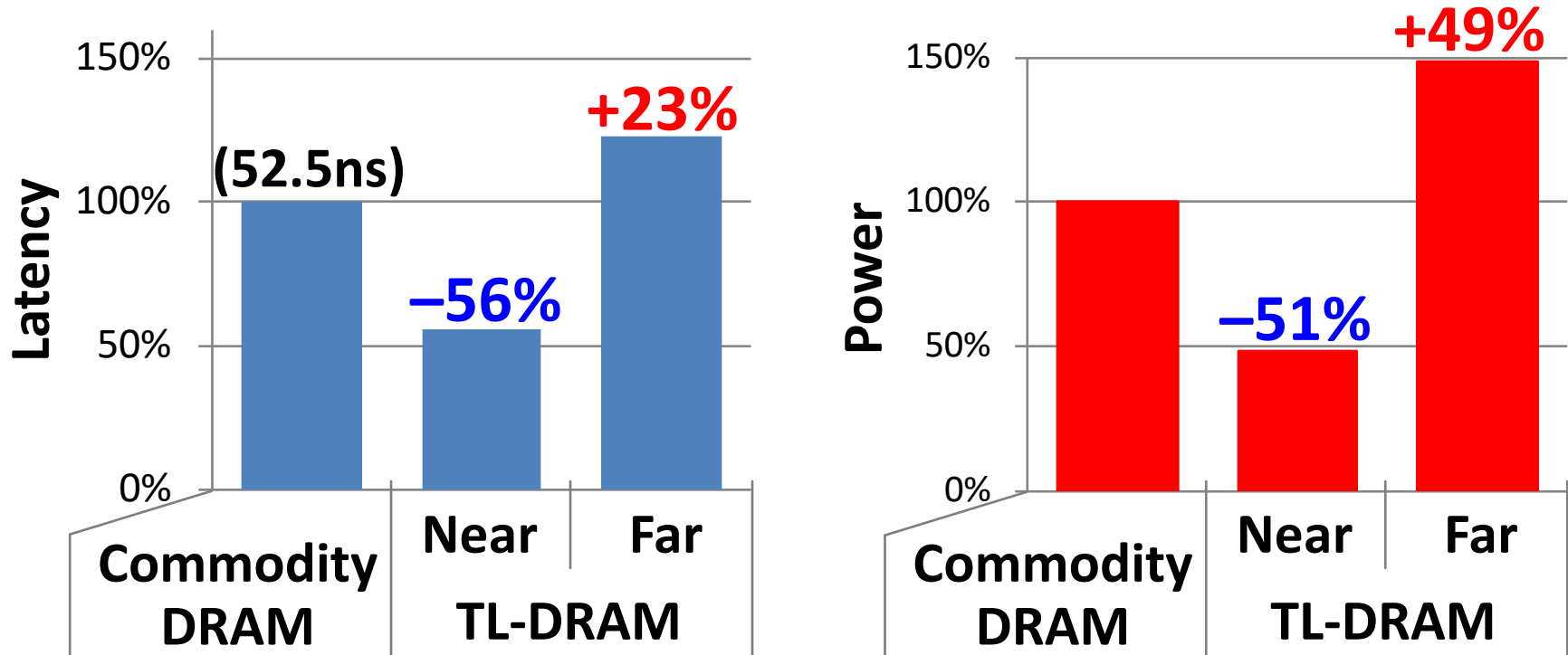
Low Latency

**Small area
using long
bitline**



Commodity DRAM vs. TL-DRAM [HPCA 2013]

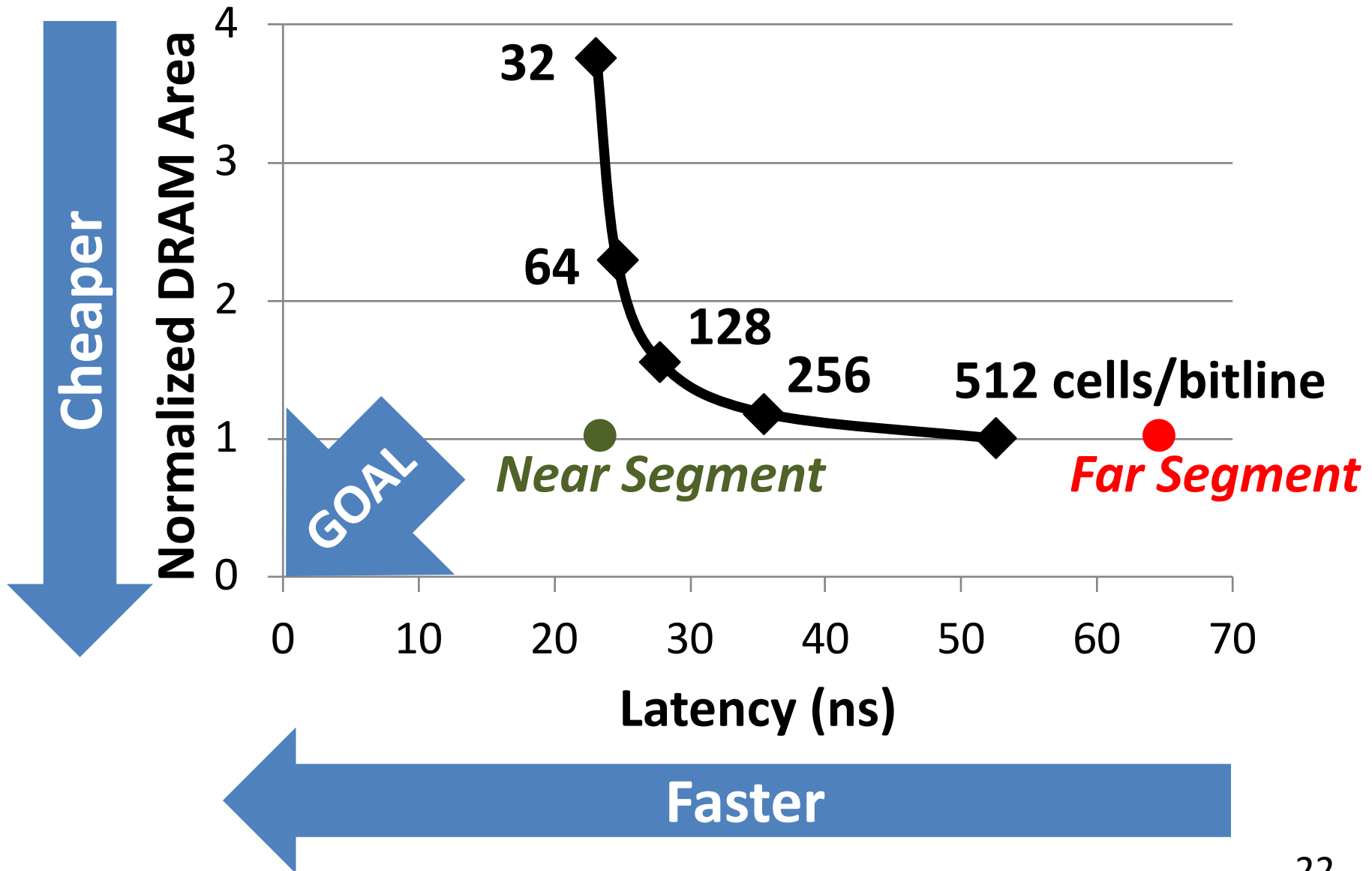
- DRAM Latency (tRC) • DRAM Power



- DRAM Area Overhead

~3%: mainly due to the isolation transistors

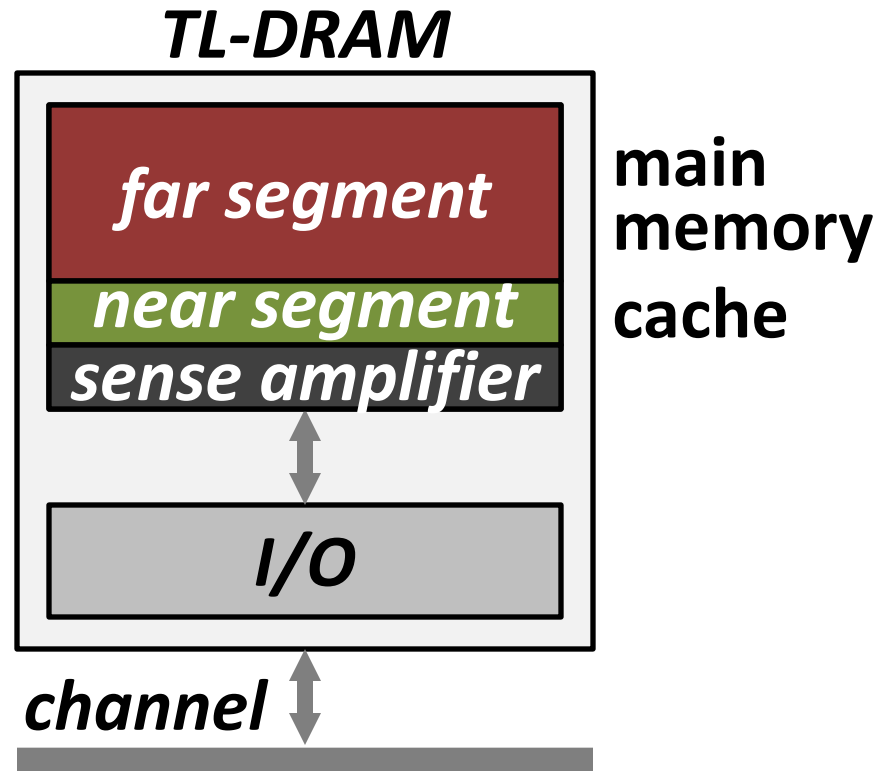
Trade-Off: Area (Die-Area) vs. Latency



Leveraging Tiered-Latency DRAM

- TL-DRAM is a ***substrate*** that can be leveraged by the hardware and/or software
- Many potential uses
 1. Use near segment as hardware-managed ***inclusive*** cache to far segment
 2. Use near segment as hardware-managed ***exclusive*** cache to far segment
 3. Profile-based page mapping by operating system
 4. Simply replace DRAM with TL-DRAM

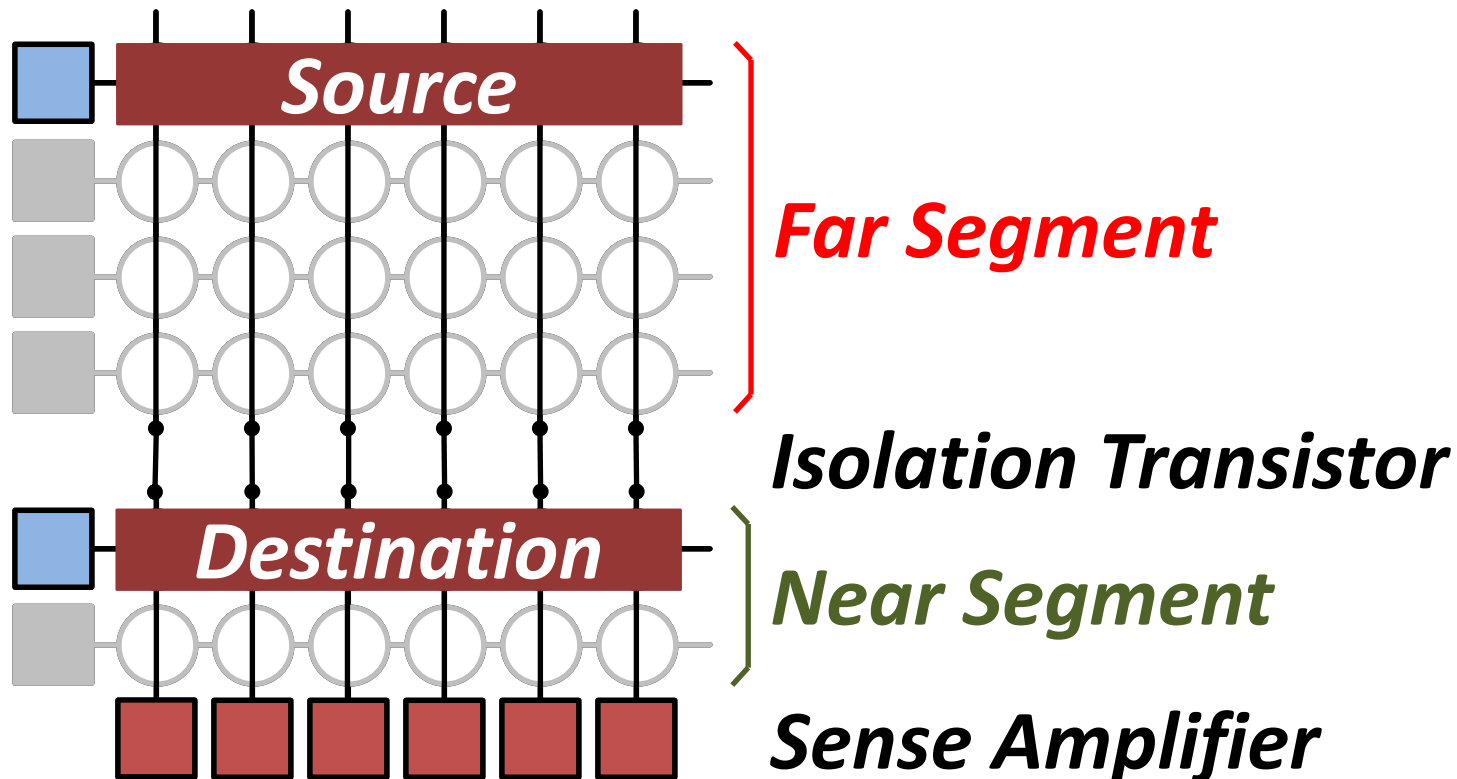
Near Segment as Hardware-Managed Cache



- **Challenge 1:** How to efficiently migrate a row between segments?
- **Challenge 2:** How to efficiently manage the cache?

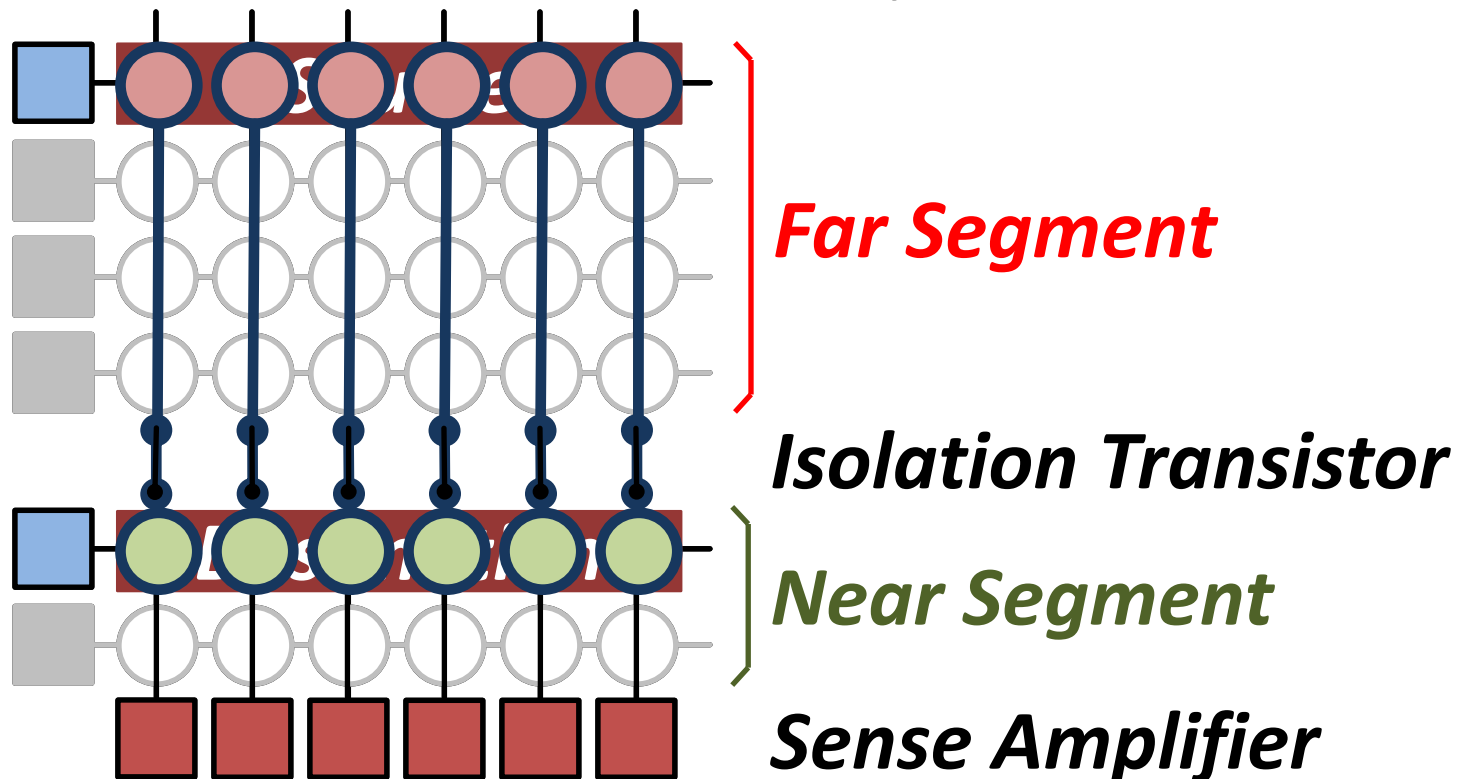
Inter-Segment Migration

- **Goal:** Migrate source row into destination row
- **Naïve way:** Memory controller reads the source row *byte by byte* and writes to destination row *byte by byte*
→ *High latency*



Inter-Segment Migration

- Our way:
 - Source and destination cells *share bitlines*
 - Transfer data from source to destination across *shared bitlines* concurrently



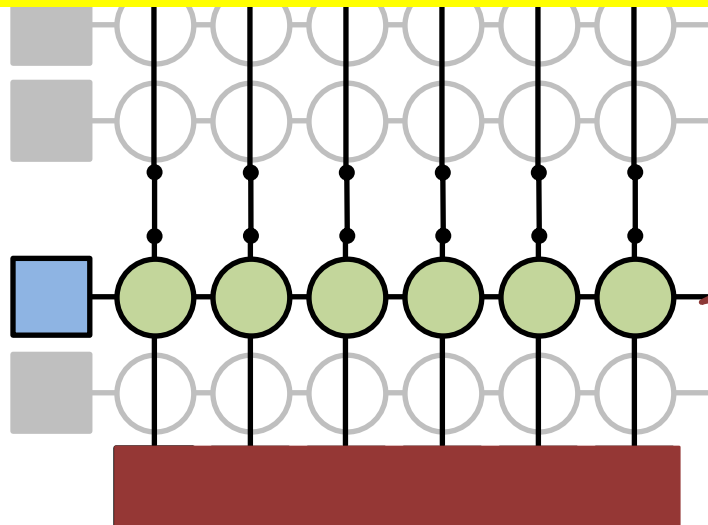
Inter-Segment Migration

- Our way:

- Source and destination cells *share bitlines*
- Transfer data from source cell to destination cell via *shared bitlines* concurrently

Step 1: Activate source row

Migration is overlapped with source row access
Additional ~4ns over row access latency

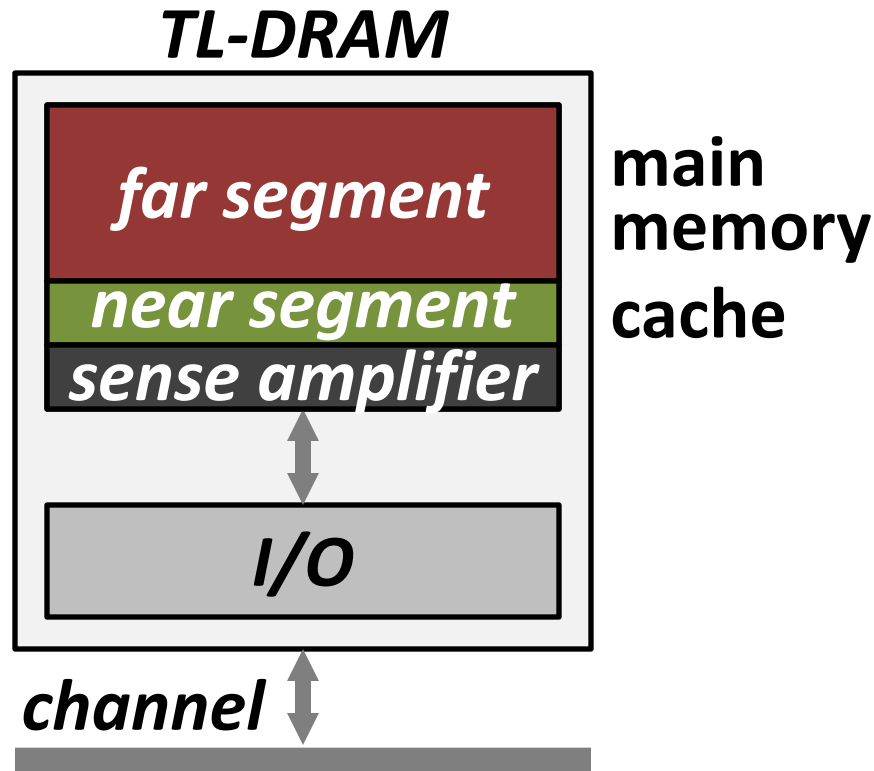


Step 2: Activate destination row to connect cell and bitline

Near Segment

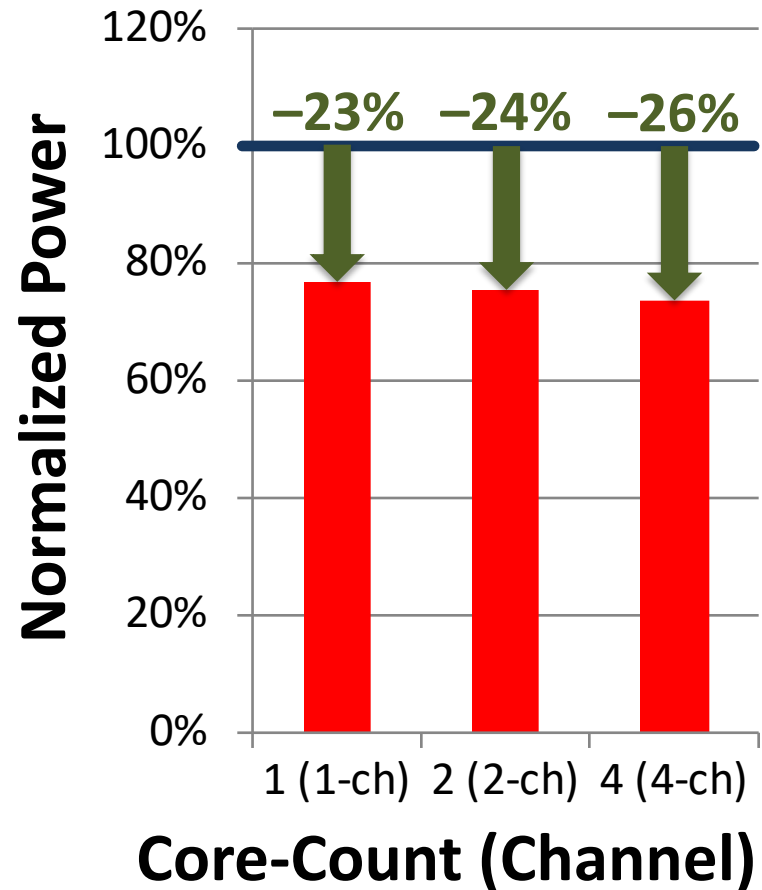
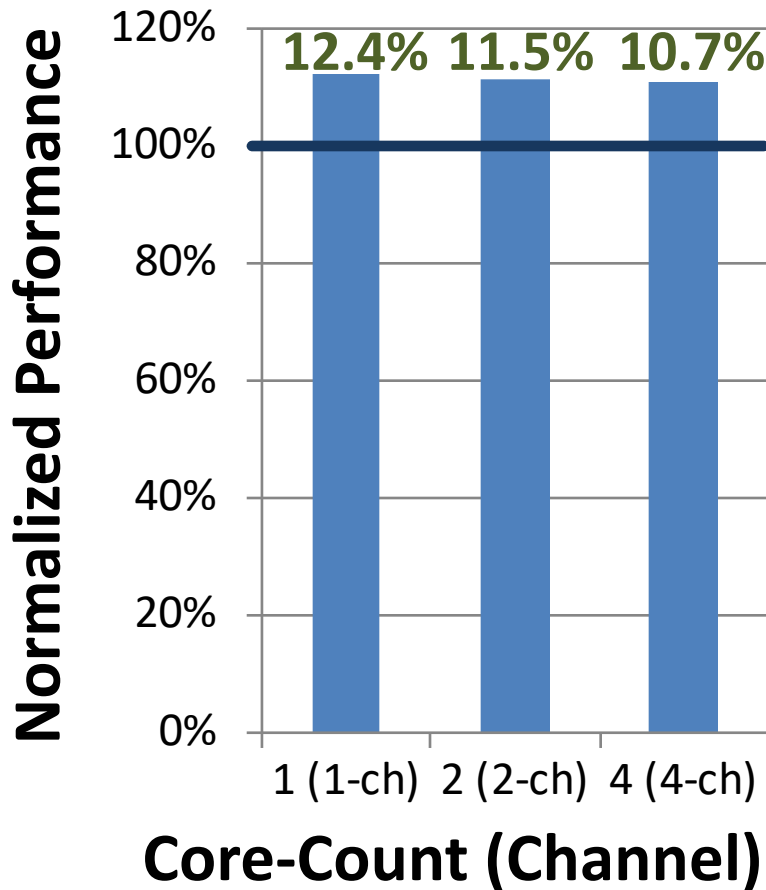
Sense Amplifier

Near Segment as Hardware-Managed Cache



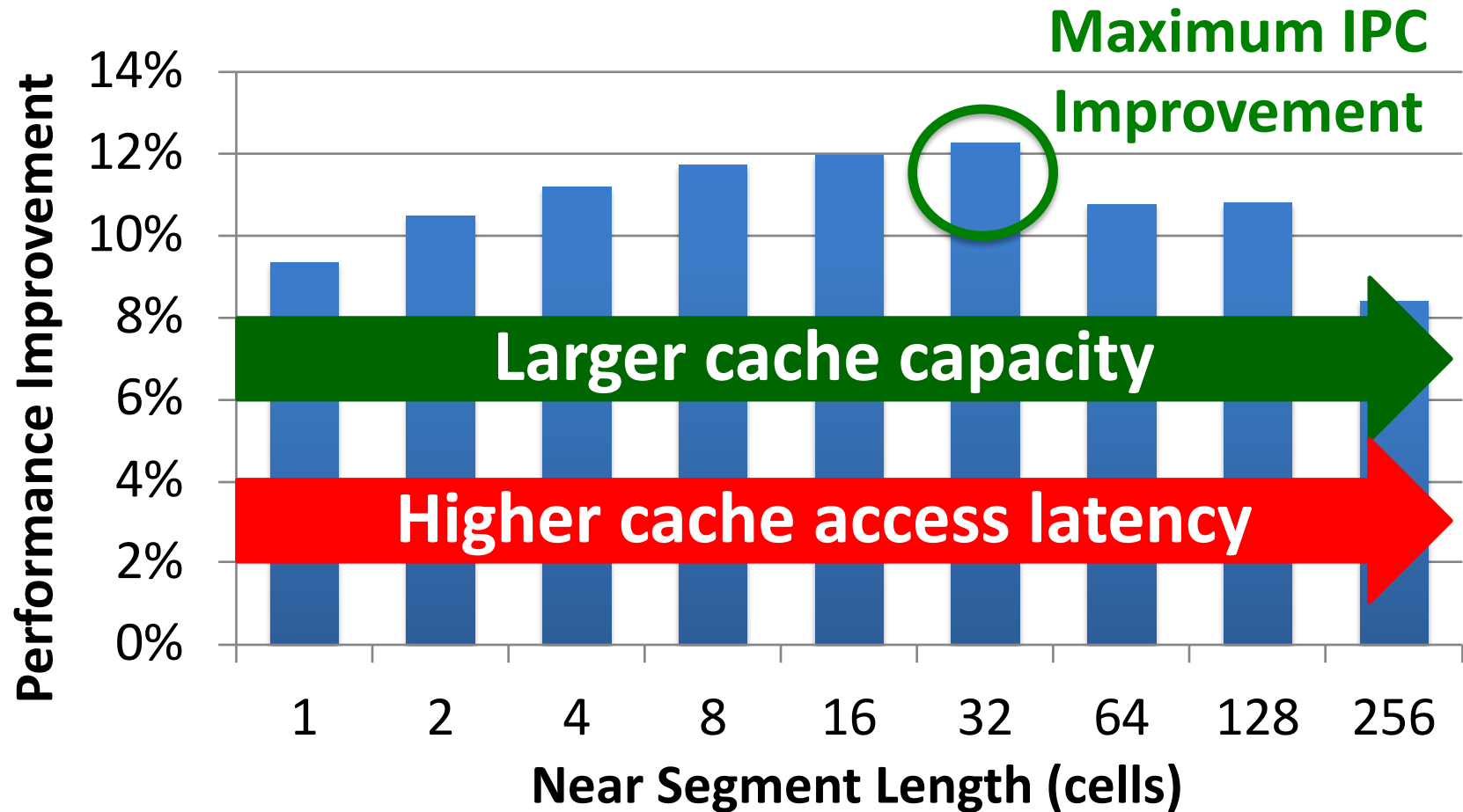
- **Challenge 1:** How to efficiently migrate a row between segments?
- **Challenge 2:** How to efficiently manage the cache?

Performance & Power Consumption



Using near segment as a cache improves performance and reduces power consumption

Single-Core: Varying Near Segment Length



By adjusting the near segment length, we can trade off cache capacity for cache latency

More on TL-DRAM

- Donghyuk Lee, Yoongu Kim, Vivek Seshadri, Jamie Liu, Lavanya Subramanian, and Onur Mutlu,
"Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture"
Proceedings of the 19th International Symposium on High-Performance Computer Architecture (HPCA), Shenzhen, China, February 2013. [Slides \(pptx\)](#)

Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture

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

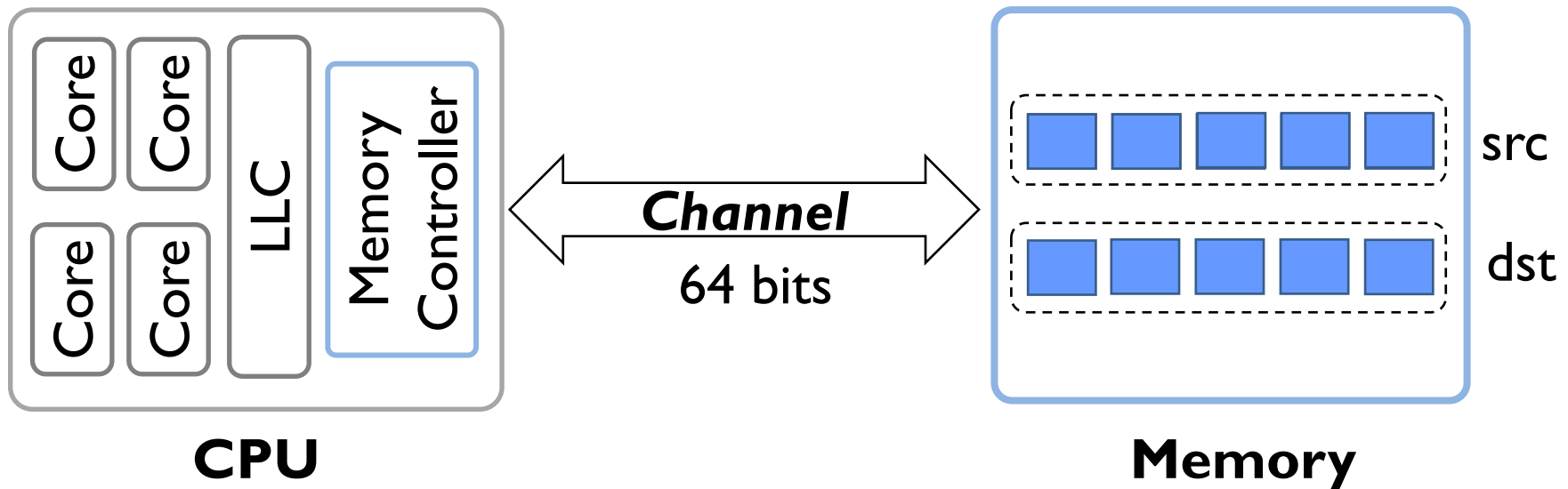
LISA: Low-Cost Inter-Linked Subarrays

[HPCA 2016]

Problem: Inefficient Bulk Data Movement

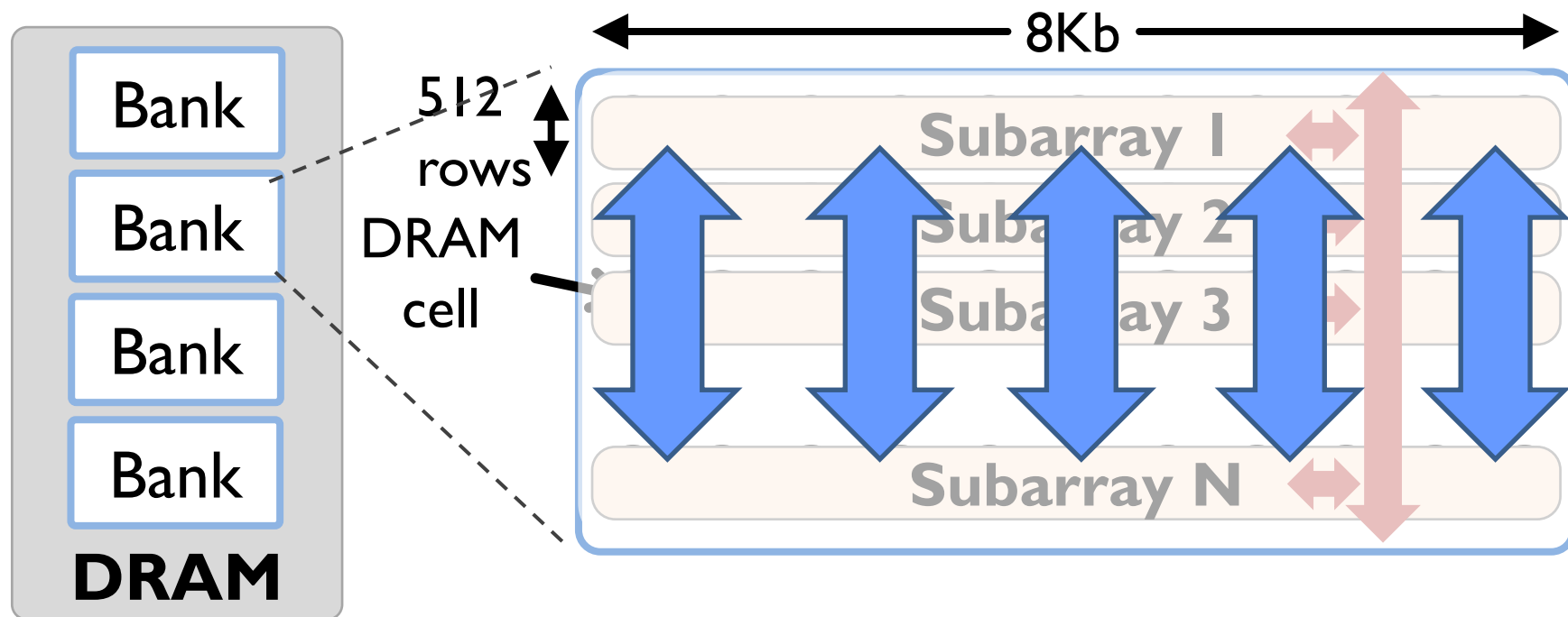
Bulk data movement is a key operation in many applications

– *memmove & memcpy*: 5% cycles in Google's datacenter [Kanev+ ISCA'15]



Long latency and high energy

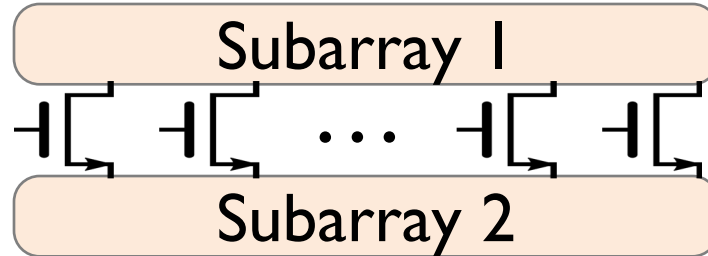
Moving Data Inside DRAM?



Goal: Provide a new substrate to enable wide connectivity between subarrays

Key Idea and Applications

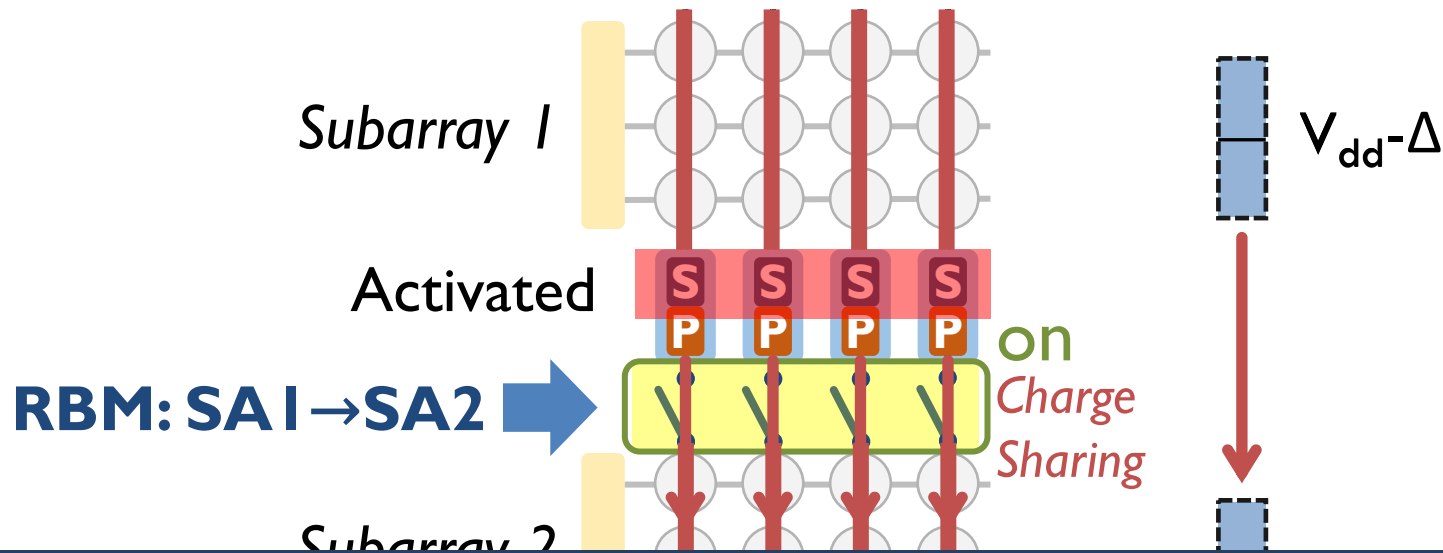
- **Low-cost Inter-linked subarrays (LISA)**
 - Fast bulk data movement between subarrays
 - Wide datapath via isolation transistors: 0.8% DRAM chip area



- LISA is a **versatile substrate** → new applications
 - Fast bulk data copy:** Copy latency 1.363ms→0.148ms (9.2x)
→ 66% speedup, -55% DRAM energy
 - In-DRAM caching:** Hot data access latency 48.7ns→21.5ns (2.2x)
→ 5% speedup
 - Fast precharge:** Precharge latency 13.1ns→5.0ns (2.6x)
→ 8% speedup

New DRAM Command to Use LISA

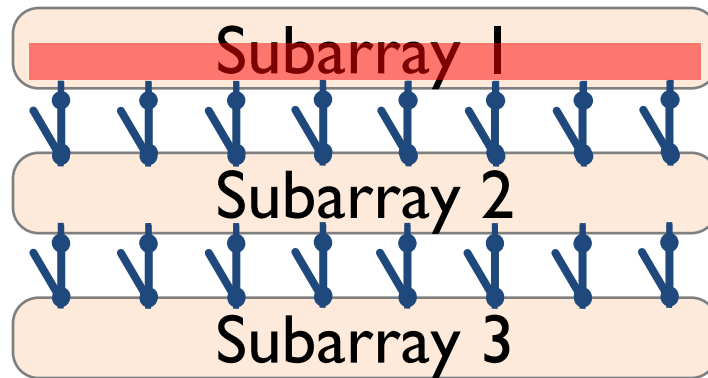
Row Buffer Movement (RBM): Move a row of data in an activated row buffer to a precharged one



RBM transfers an entire row b/w subarrays

RBM Analysis

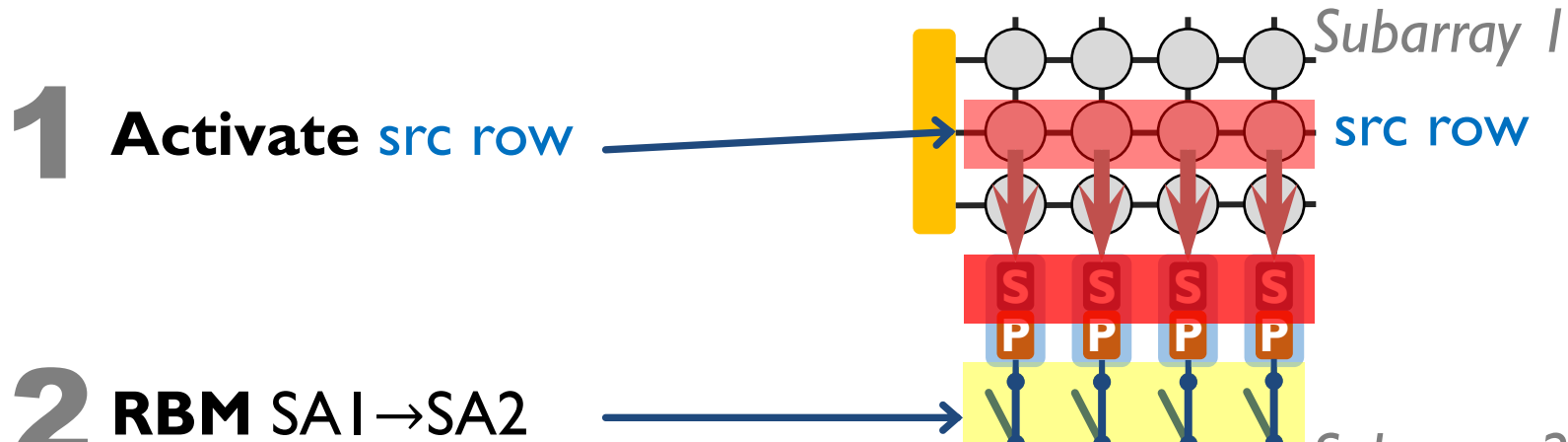
- The range of RBM depends on the DRAM design
 - Multiple RBMs to move data across > 3 subarrays



- Validated with SPICE using worst-case cells
 - NCSU FreePDK 45nm library
- **4KB data in 8ns (w/ 60% guardband)**
→ **500 GB/s, 26x** bandwidth of a DDR4-2400 channel
- **0.8% DRAM chip area overhead [O+ ISCA'14]**

1. Rapid Inter-Subarray Copying (RISC)

- **Goal:** Efficiently copy a row across subarrays
- **Key idea:** Use *RBM* to form a new command sequence

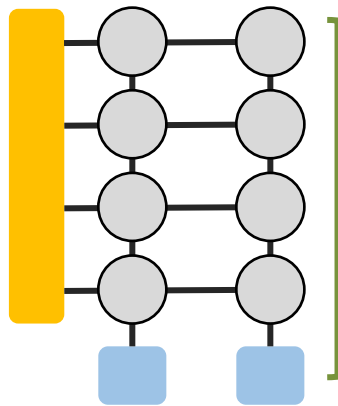


Reduces row-copy latency by 9.2x,
DRAM energy by 48.1x

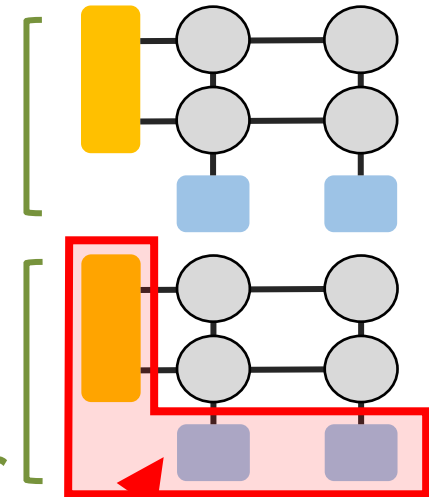
2. Variable Latency DRAM (VILLA)

- **Goal:** Reduce DRAM latency with low area overhead
- **Motivation:** Trade-off between area and latency

**Long Bitline
(DDR_x)**



**Short Bitline
(RLDRAM)**

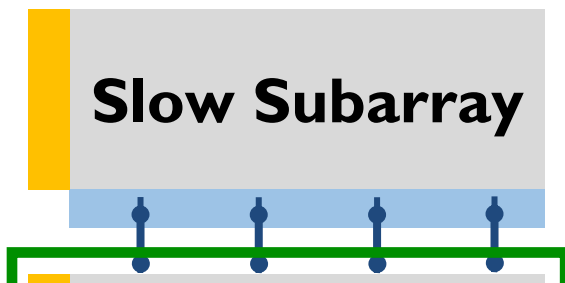


Shorter bitlines → faster
activate and **precharge** time

High area overhead: >40%

2. Variable Latency DRAM (VILLA)

- **Key idea:** Reduce access latency of hot data via a **heterogeneous DRAM** design [Lee+ HPCA'13, Son+ ISCA'13]
- **VILLA:** Add fast subarrays as a **cache** in each bank

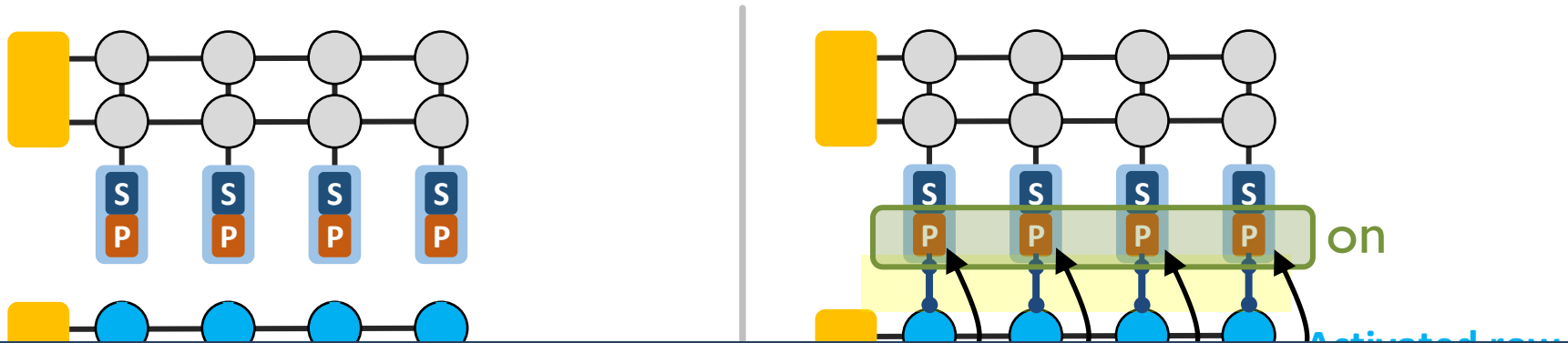


Challenge: VILLA cache requires frequent movement of data rows

Reduces hot data access latency by 2.2x
at only 1.6% area overhead

3. Linked Precharge (LIP)

- **Problem:** The precharge time is limited by the strength of one precharge unit
- **Linked Precharge (LIP):** LISA precharges a subarray using multiple precharge units



Reduces precharge latency by 2.6x
(43% guardband)

More on LISA

- Kevin K. Chang, Prashant J. Nair, Saugata Ghose, Donghyuk Lee, Moinuddin K. Qureshi, and Onur Mutlu,
"Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Movement in DRAM"
Proceedings of the 22nd International Symposium on High-Performance Computer Architecture (HPCA), Barcelona, Spain, March 2016.
[[Slides \(pptx\)](#)] [[pdf](#)]
[[Source Code](#)]

Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Movement in DRAM

Kevin K. Chang[†], Prashant J. Nair^{*}, Donghyuk Lee[†], Saugata Ghose[†], Moinuddin K. Qureshi^{*}, and Onur Mutlu[†]

[†]Carnegie Mellon University ^{*}Georgia Institute of Technology

Why the Long Memory Latency?

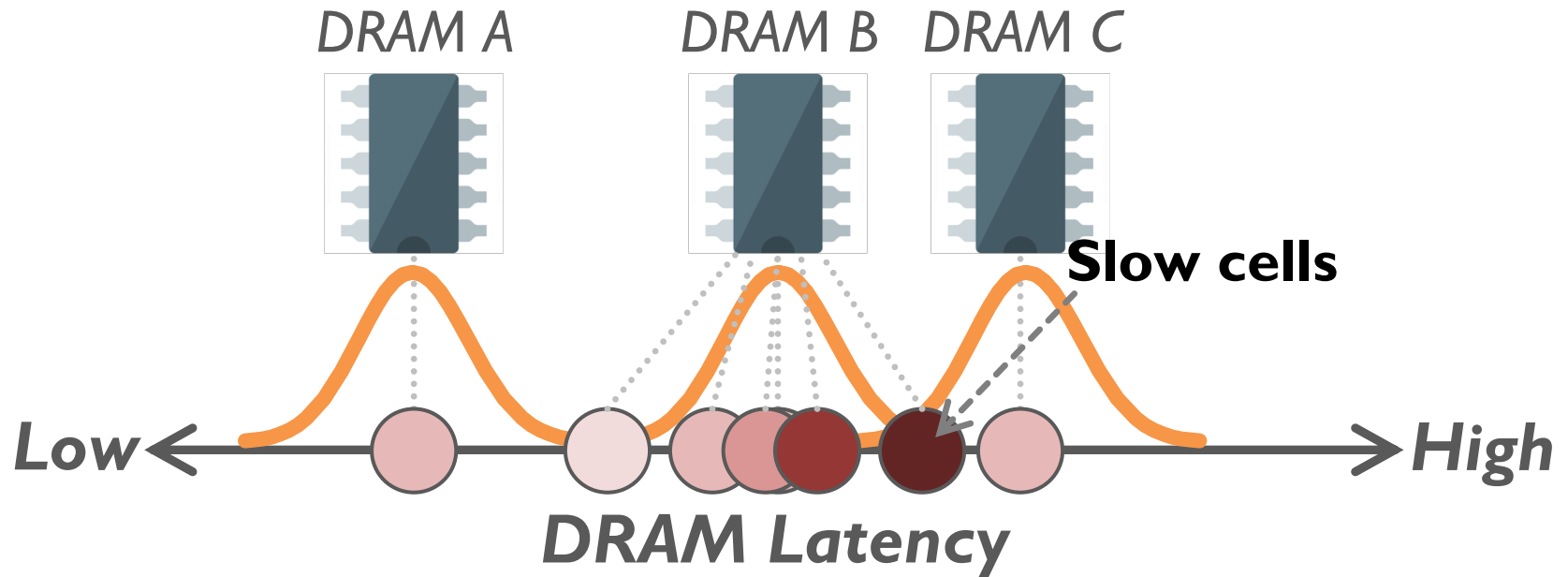
- Reason 1: Design of DRAM Micro-architecture
 - Goal: Maximize capacity/area, not minimize latency
- Reason 2: “One size fits all” approach to latency specification
 - Same latency parameters for all temperatures
 - Same latency parameters for all DRAM chips (e.g., rows)
 - Same latency parameters for all parts of a DRAM chip
 - Same latency parameters for all supply voltage levels
 - Same latency parameters for all application data
 - ...

Tackling the Fixed Latency Mindset

- Reliable operation latency is actually very heterogeneous
 - Across temperatures, chips, parts of a chip, voltage levels, ...
- Idea: Dynamically find out and use the lowest latency one can reliably access a memory location with
 - Adaptive-Latency DRAM [HPCA 2015]
 - Flexible-Latency DRAM [SIGMETRICS 2016]
 - Design-Induced Variation-Aware DRAM [SIGMETRICS 2017]
 - Voltron [SIGMETRICS 2017]
 - DRAM Latency PUF [HPCA 2018]
 - ...
- We would like to find sources of latency heterogeneity and exploit them to minimize latency

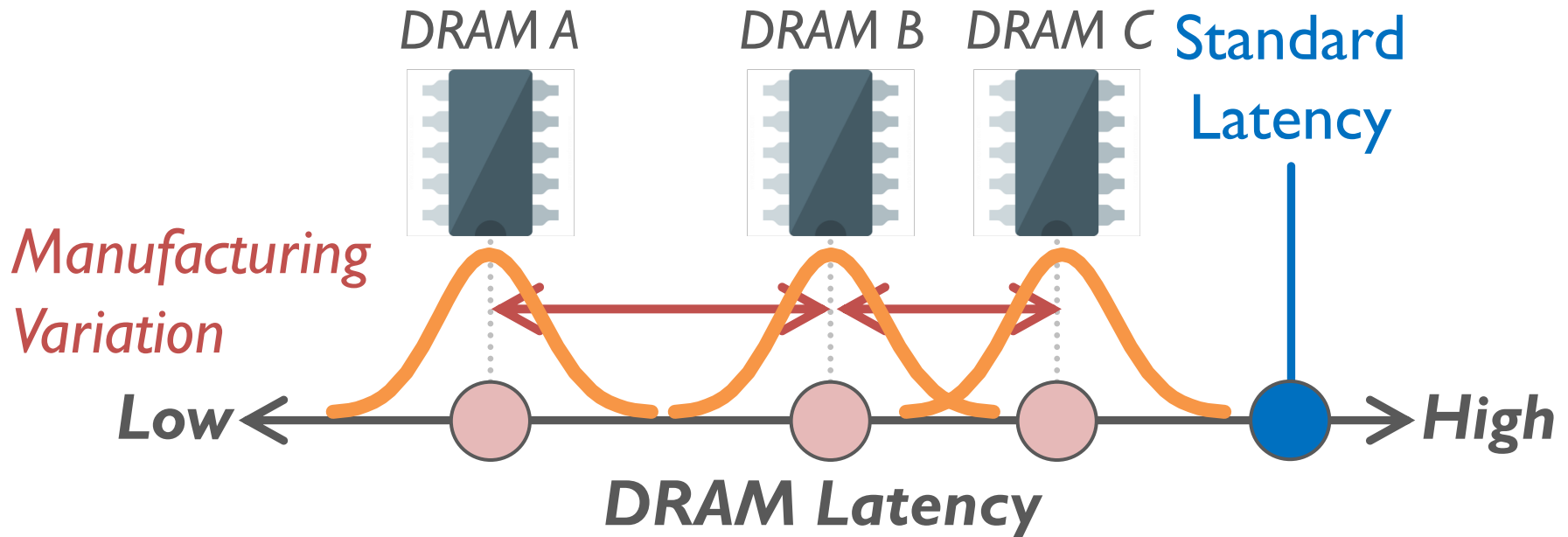
Latency Variation in Memory Chips

Heterogeneous manufacturing & operating conditions →
latency variation in timing parameters



Why is Latency High?

- DRAM latency: Delay as specified in DRAM standards
 - Doesn't reflect true DRAM device latency
- Imperfect manufacturing process → latency variation
- **High standard latency** chosen to increase yield



What Causes the Long Memory Latency?

- **Conservative timing margins!**
- DRAM timing parameters are set to cover the worst case
- **Worst-case temperatures**
 - ❑ 85 degrees vs. common-case
 - ❑ to enable a wide range of operating conditions
- **Worst-case devices**
 - ❑ DRAM cell with smallest charge across any acceptable device
 - ❑ to tolerate process variation at acceptable yield
- This leads to large timing margins for the common case

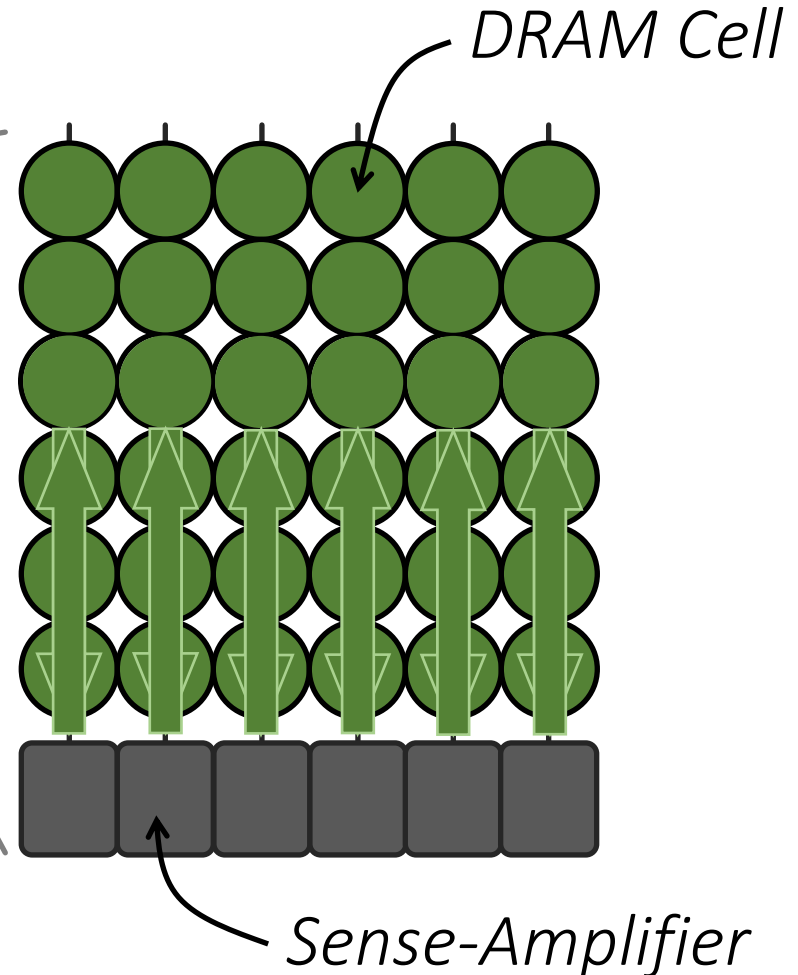
Understanding and Exploiting Variation in DRAM Latency

DRAM Stores Data as Charge

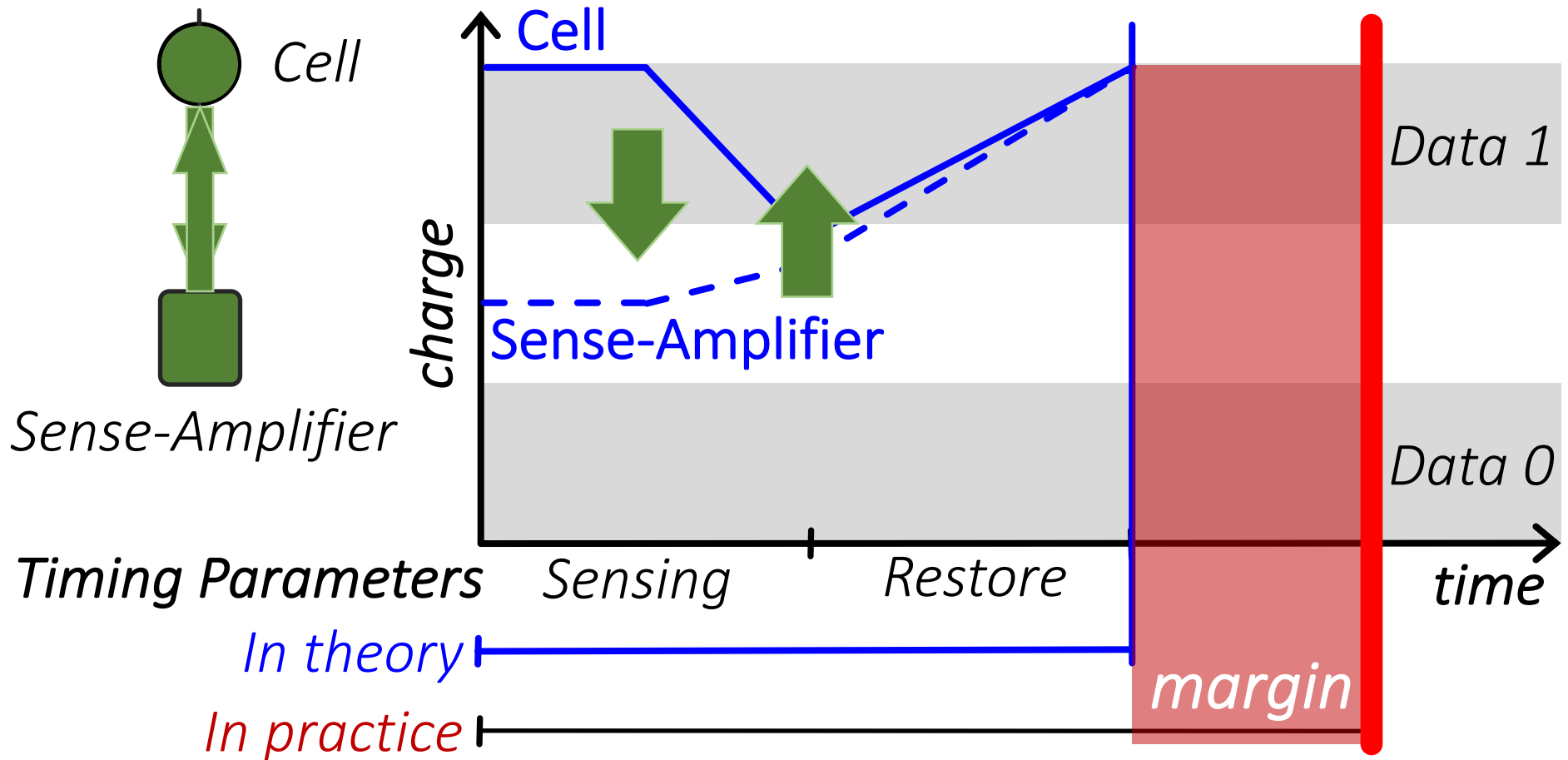


Three steps of
charge movement

1. Sensing
2. Restore
3. Precharge



DRAM Charge over Time



Why does DRAM need the extra timing margin?

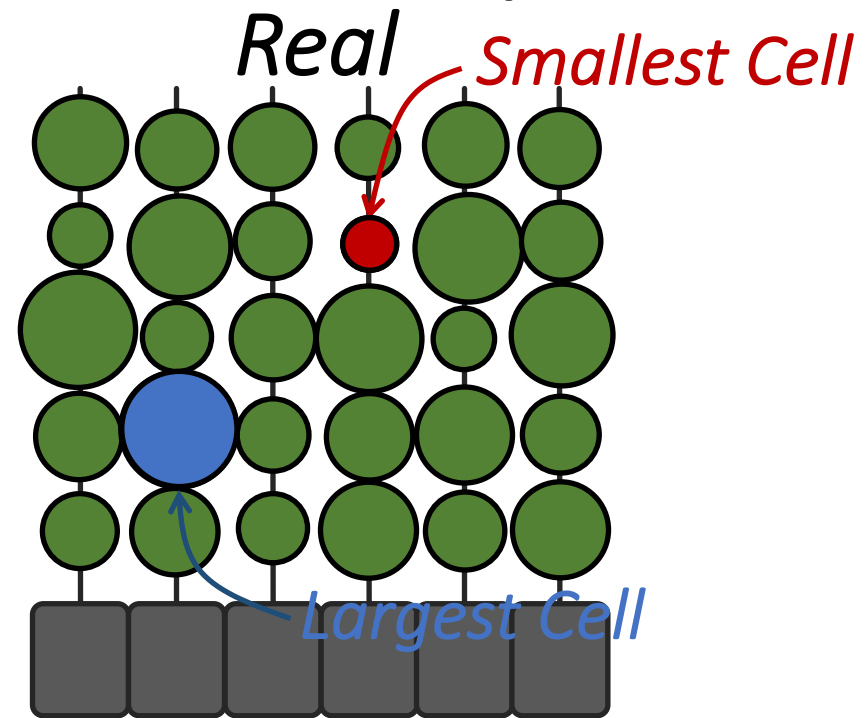
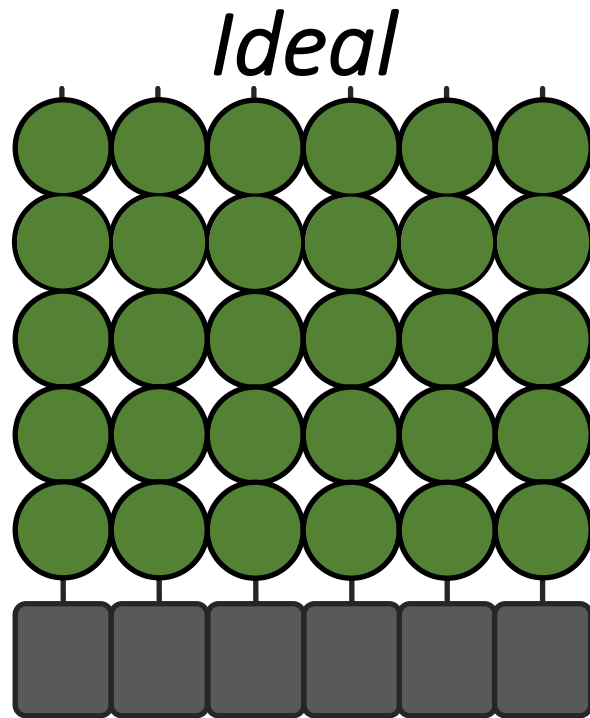
Two Reasons for Timing Margin

1. Process Variation

- DRAM cells are not equal
- Leads to extra timing margin for a cell that can store a large amount of charge

2. Temperature Dependence

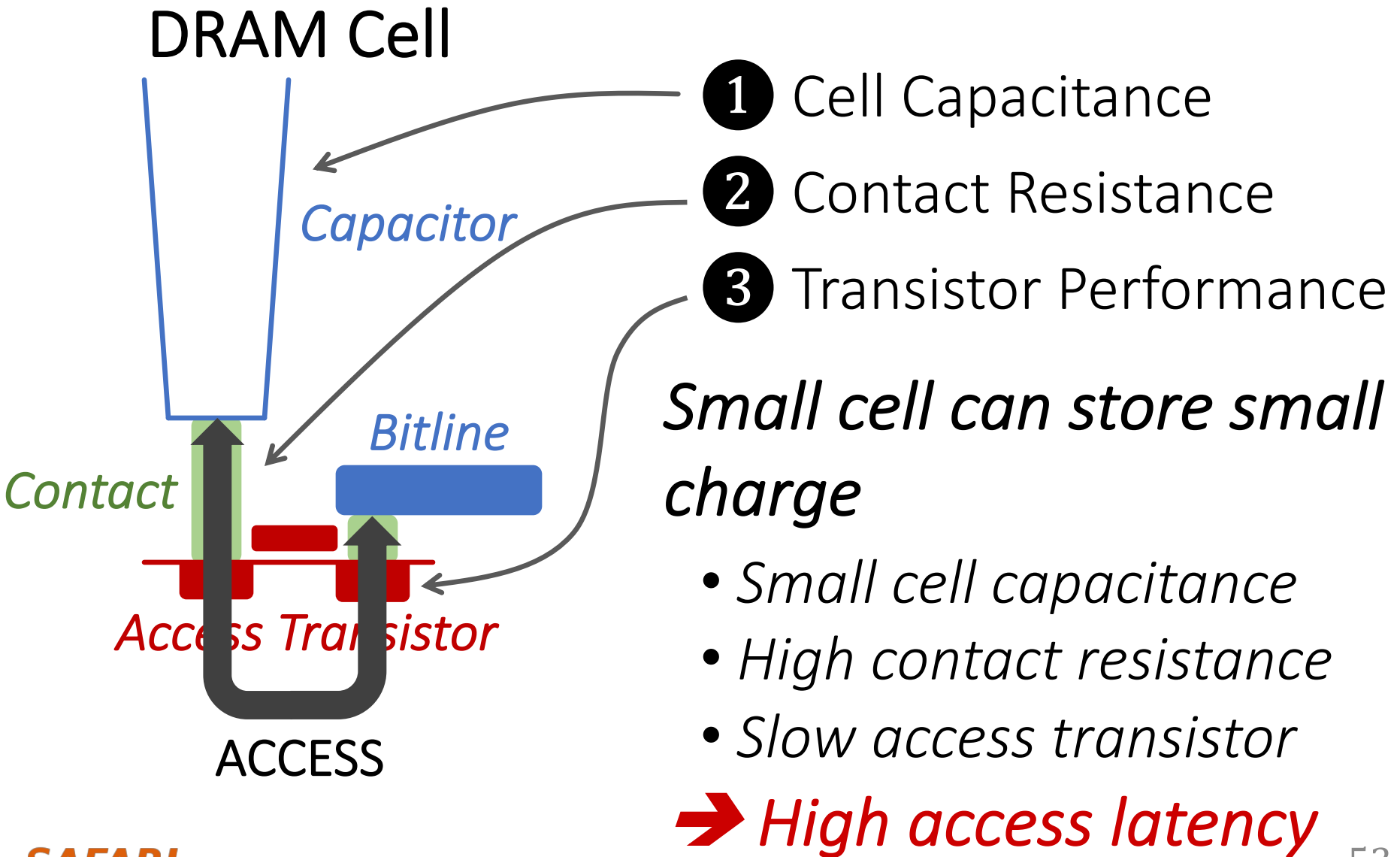
DRAM Cells are Not Equal



Same Size → Large variation in cell size
Same Charge → Large variation in charge
Same Latency → Large variation in access latency

Different Size →
Different Charge →
Different Latency →

Process Variation



Two Reasons for Timing Margin

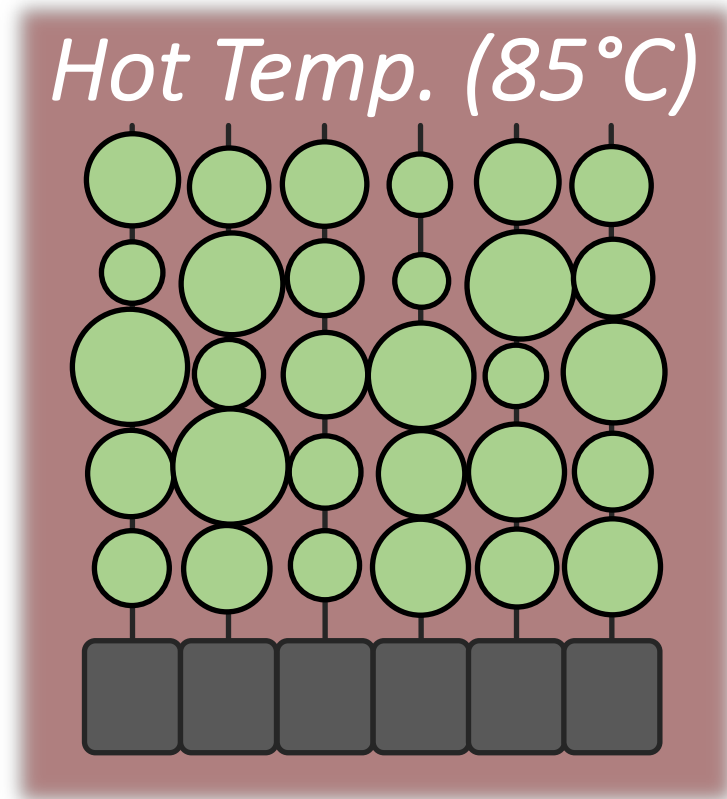
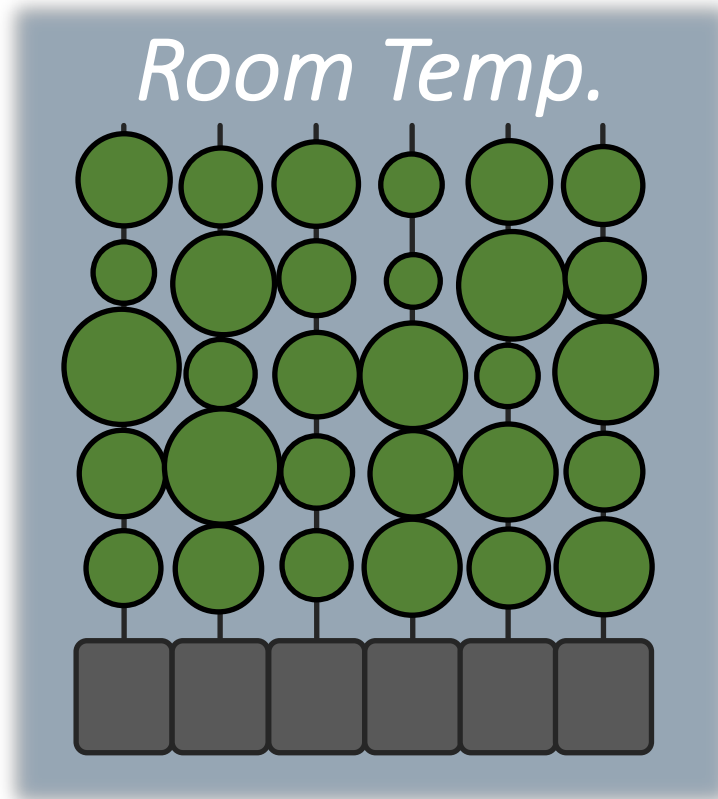
1. Process Variation

- DRAM cells are not equal
- Leads to **extra timing margin** for a cell that can store a large amount of charge

2. Temperature Dependence

- DRAM leaks more charge at higher temperature
- Leads to extra timing margin for cells that operate at low temperature

Charge Leakage Temperature



Cells store small charge at high temperature and large charge at low temperature
→ Large variation in access latency

DRAM Timing Parameters

- *DRAM timing parameters are dictated by the worst-case*
 - The smallest cell with the smallest charge in all DRAM products
 - Operating at the highest temperature
- *Large timing margin for the common-case*

Adaptive-Latency DRAM [HPCA 2015]

- Idea: Optimize DRAM timing for the common case
 - Current temperature
 - Current DRAM module
- Why would this reduce latency?
 - A DRAM cell can store much more charge in the common case (low temperature, strong cell) than in the worst case
 - More charge in a DRAM cell
 - Faster sensing, charge restoration, precharging
 - Faster access (read, write, refresh, ...)

Extra Charge → Reduced Latency

1. Sensing

Sense cells with extra charge faster

→ Lower sensing latency

2. Restore

No need to fully restore cells with extra charge

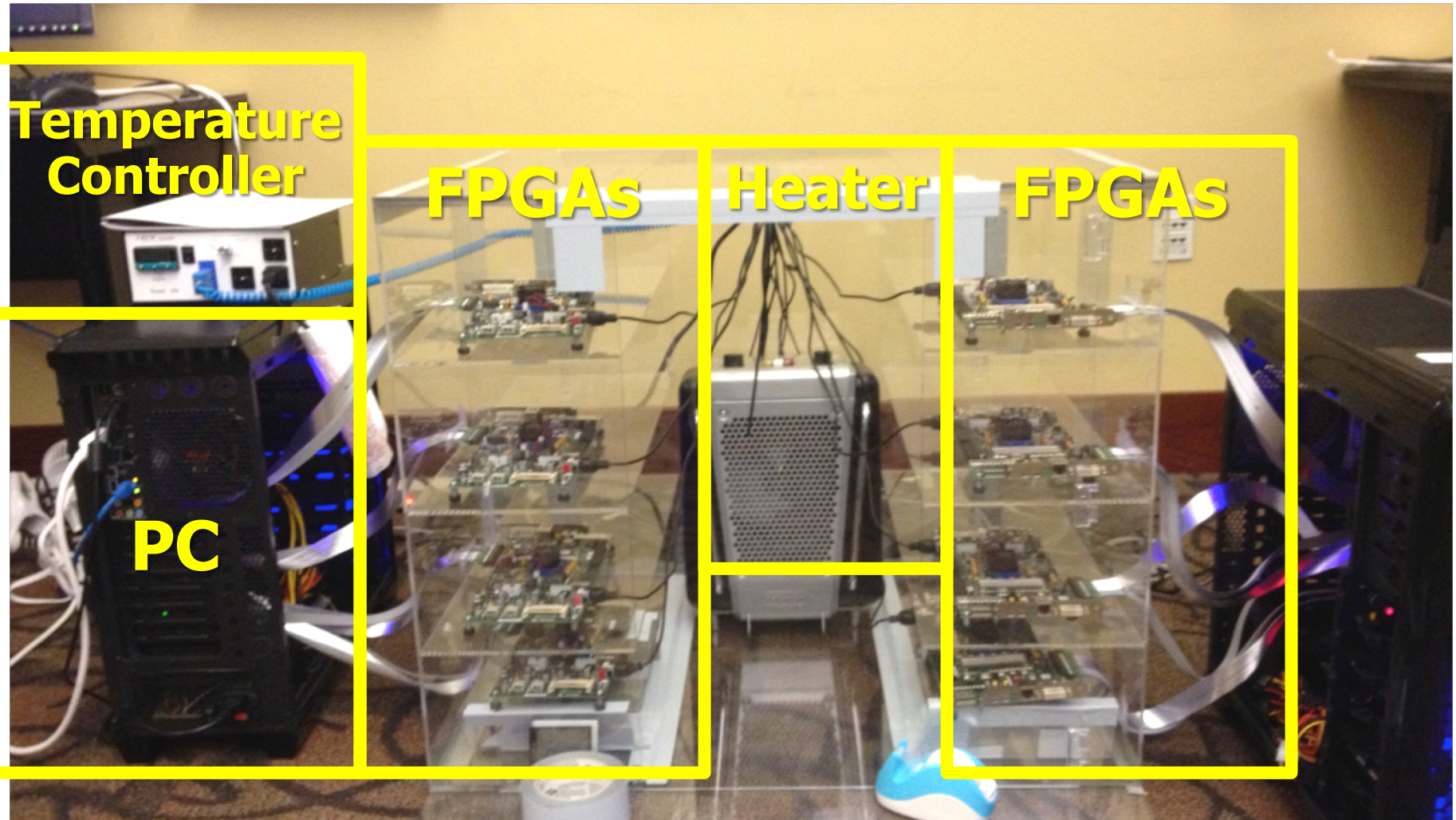
→ Lower restoration latency

3. Precharge

No need to fully precharge bitlines for cells with extra charge

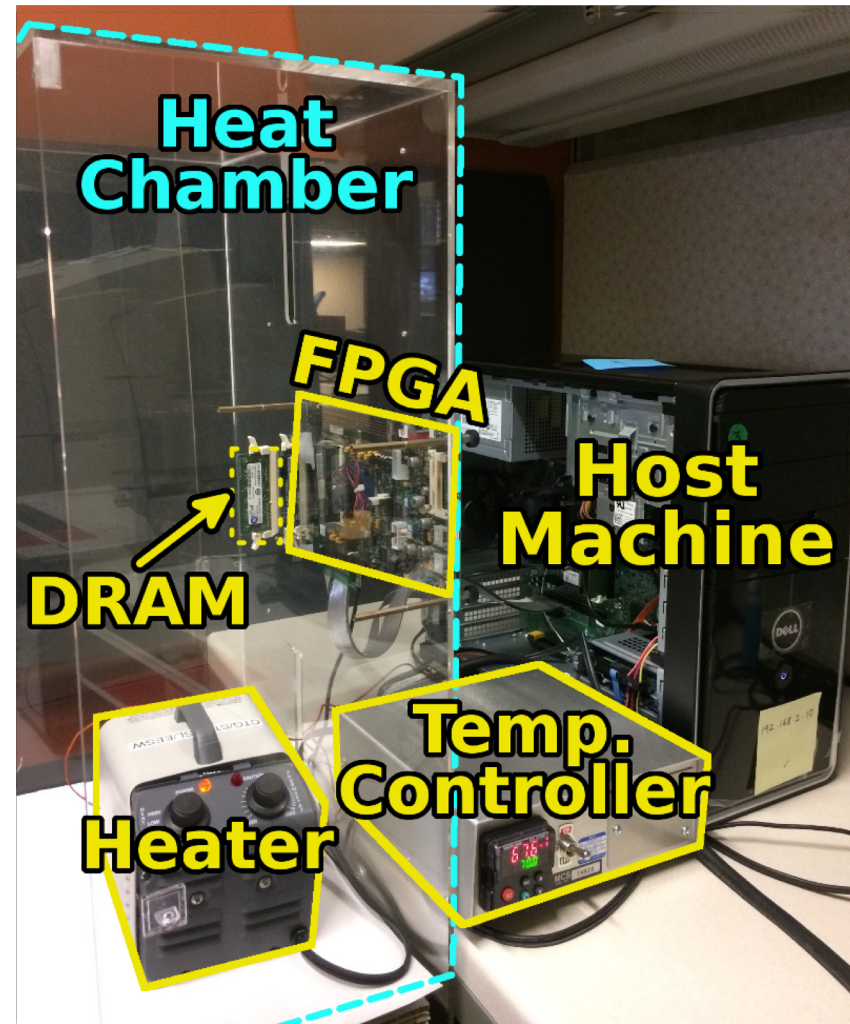
→ Lower precharge latency

DRAM Characterization Infrastructure



DRAM Characterization Infrastructure

- Hasan Hassan et al., **SoftMC: A Flexible and Practical Open-Source Infrastructure for Enabling Experimental DRAM Studies**, HPCA 2017.
- Flexible
- Easy to Use (C++ API)
- Open-source
github.com/CMU-SAFARI/SoftMC



SoftMC: Open Source DRAM Infrastructure

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

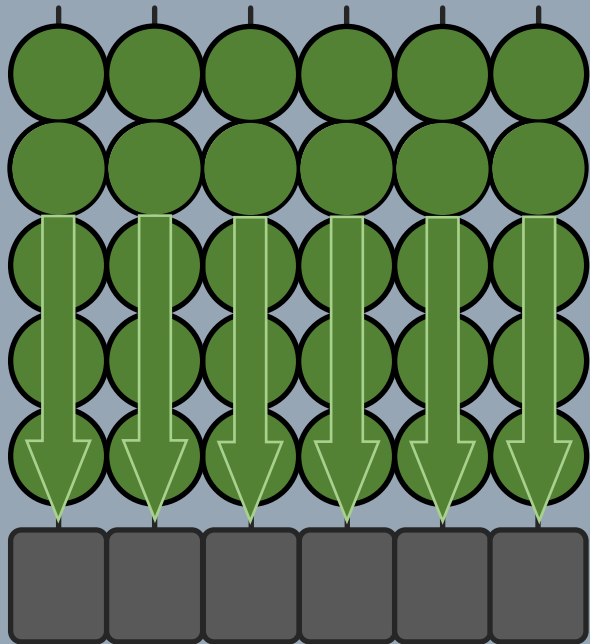
SoftMC: A Flexible and Practical Open-Source Infrastructure for Enabling Experimental DRAM Studies

Hasan Hassan^{1,2,3} Nandita Vijaykumar³ Samira Khan^{4,3} Saugata Ghose³ Kevin Chang³
Gennady Pekhimenko^{5,3} Donghyuk Lee^{6,3} Oguz Ergin² Onur Mutlu^{1,3}

¹*ETH Zürich* ²*TOBB University of Economics & Technology* ³*Carnegie Mellon University*
⁴*University of Virginia* ⁵*Microsoft Research* ⁶*NVIDIA Research*

Observation 1. Faster Sensing

Typical DIMM at Low Temperature



More Charge

Strong Charge
Flow

Faster Sensing

*115 DIMM
Characterization*

Timing
(t_{RCD})

17% ↓

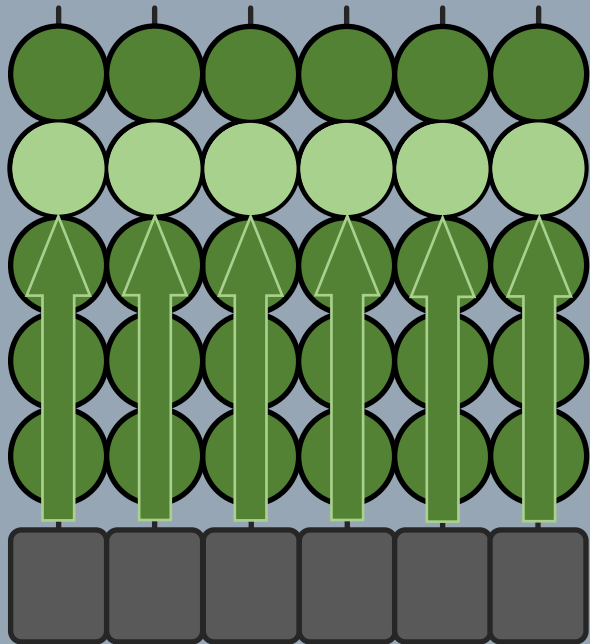
No Errors

Typical DIMM at Low Temperature

➔ *More charge* ➔ *Faster sensing*

Observation 2. Reducing Restore Time

Typical DIMM at Low Temperature



Less Leakage →
Extra Charge

No Need to Fully
Restore Charge

*115 DIMM
Characterization*

Read (t_{RAS})

37% ↓

Write (t_{WR})

54% ↓

No Errors

Typical DIMM at lower temperature

➔ *More charge* ➔ *Restore time reduction*

AL-DRAM

- *Key idea*
 - Optimize DRAM timing parameters online
- *Two components*
 - DRAM manufacturer provides multiple sets of **reliable DRAM timing parameters** at different temperatures for each DIMM
 - System monitors **DRAM temperature** & uses appropriate DRAM timing parameters

DRAM Temperature

- *DRAM temperature measurement*
 - Server cluster: Operates at under 34°C
 - Desktop: Operates at under 50°C
 - *DRAM standard optimized for 85 °C*

DRAM operates at low temperatures
in the common-case

- *Previous works – Maintain low DRAM temperature*
 - David+ ICAC 2011
 - Liu+ ISCA 2007
 - Zhu+ ITherm 2008

Latency Reduction Summary of 115 DIMMs

- *Latency reduction for read & write (55°C)*
 - *Read Latency: 32.7%*
 - *Write Latency: 55.1%*
- *Latency reduction for each timing parameter (55°C)*
 - *Sensing: 17.3%*
 - *Restore: 37.3% (read), 54.8% (write)*
 - *Precharge: 35.2%*

AL-DRAM: Real System Evaluation

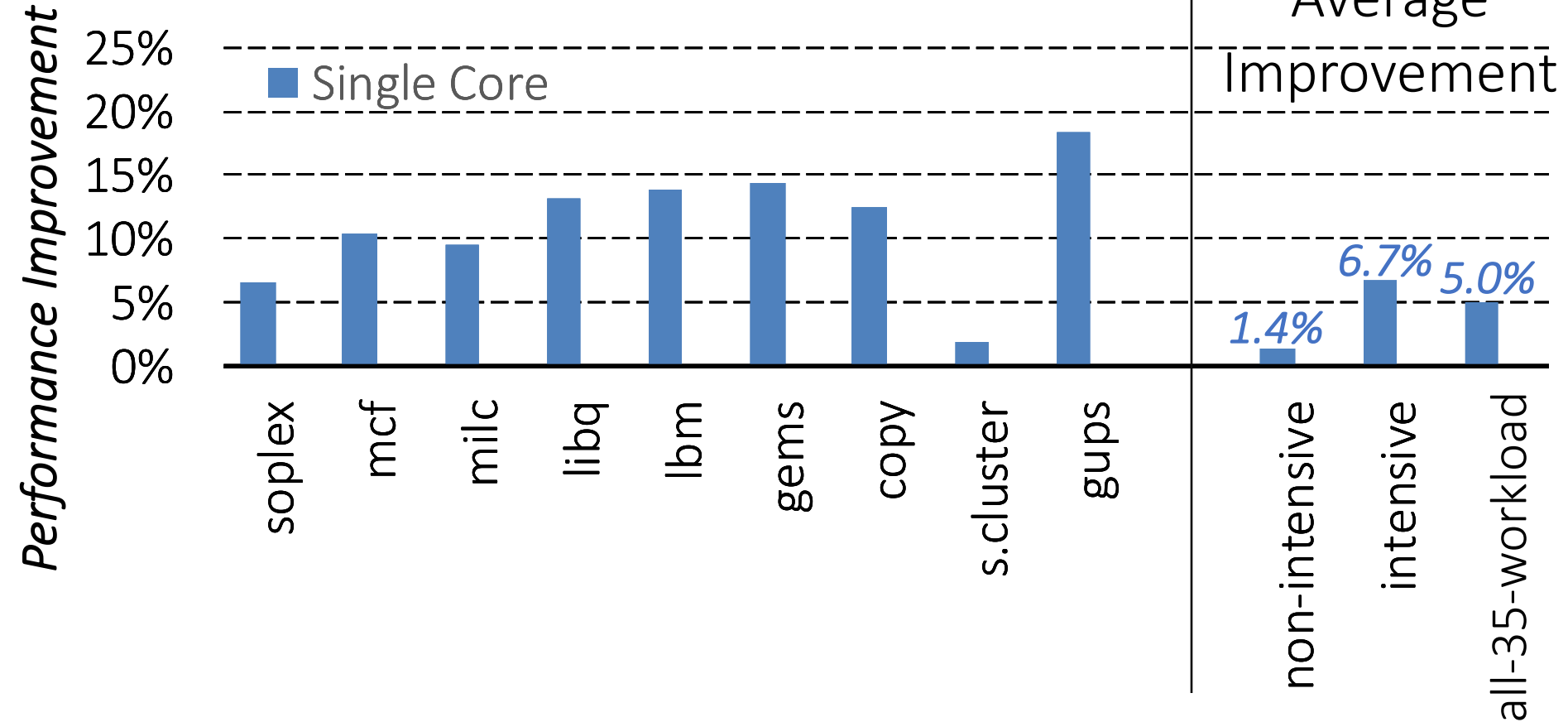
- *System*
 - *CPU: AMD 4386 (8 Cores, 3.1GHz, 8MB LLC)*

D18F2x200_dct[0]_mp[1:0] DDR3 DRAM Timing 0

Reset: 0F05_0505h. See [2.9.3 \[DCT Configuration Registers\]](#).

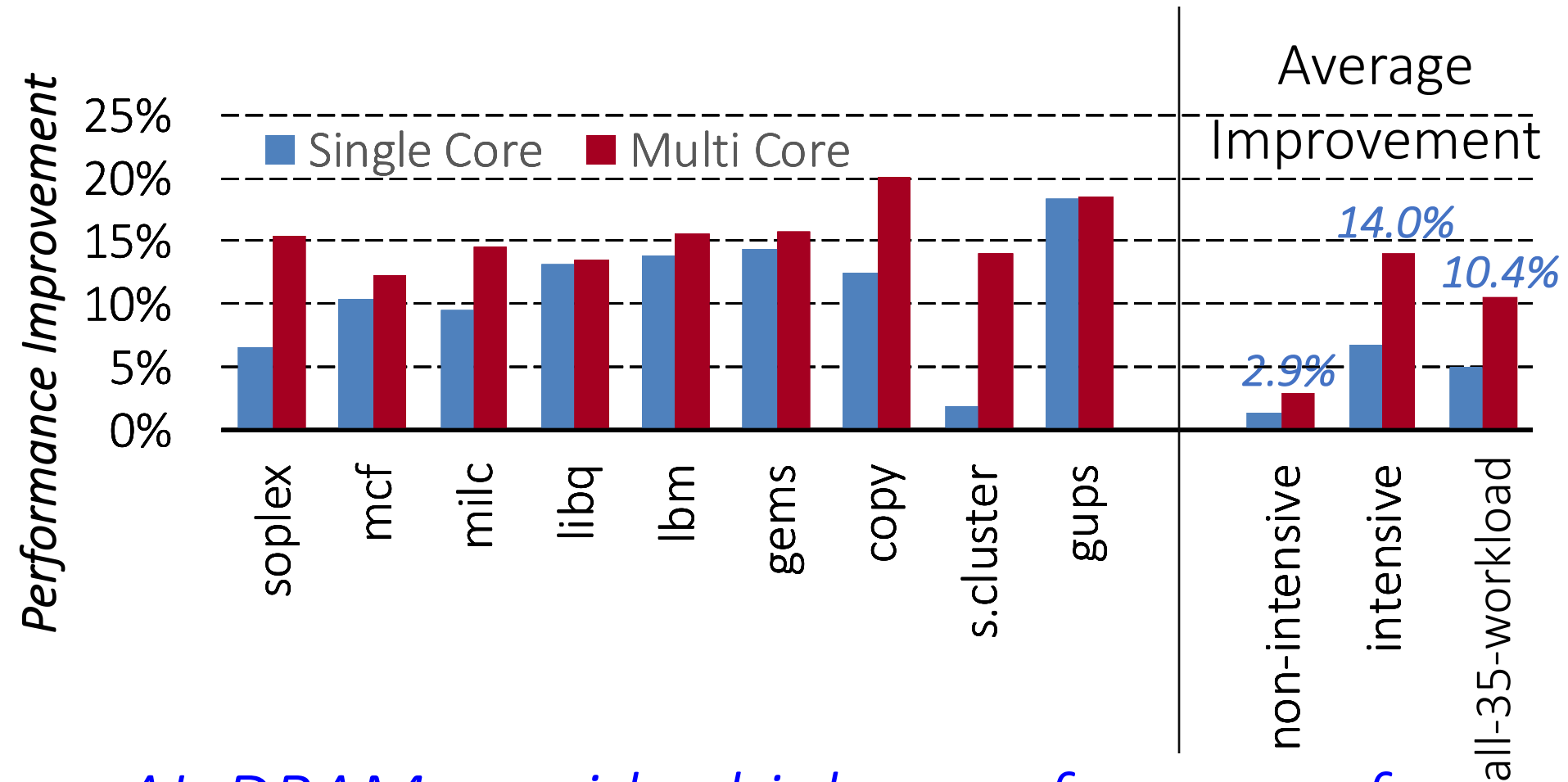
Bits	Description								
31:30	Reserved.								
29:24	Tras: row active strobe. Read-write. BIOS: See 2.9.7.5 [SPD ROM-Based Configuration] . Specifies the minimum time in memory clock cycles from an activate command to a precharge command, both to the same chip select bank. <table><tr><th>Bits</th><th>Description</th></tr><tr><td>07h-00h</td><td>Reserved</td></tr><tr><td>2Ah-08h</td><td><Tras> clocks</td></tr><tr><td>3Fh-2Bh</td><td>Reserved</td></tr></table>	Bits	Description	07h-00h	Reserved	2Ah-08h	<Tras> clocks	3Fh-2Bh	Reserved
Bits	Description								
07h-00h	Reserved								
2Ah-08h	<Tras> clocks								
3Fh-2Bh	Reserved								
23:21	Reserved.								
20:16	Trp: row precharge time. Read-write. BIOS: See 2.9.7.5 [SPD ROM-Based Configuration] . Specifies the minimum time in memory clock cycles from a precharge command to an activate command or auto refresh command, both to the same bank.								

AL-DRAM: Single-Core Evaluation



AL-DRAM improves performance on a real system

AL-DRAM: Multi-Core Evaluation



AL-DRAM provides higher performance for multi-programmed & multi-threaded workloads

Reducing Latency Also Reduces Energy

- AL-DRAM reduces DRAM power consumption by 5.8%
- Major reason: reduction in row activation time

AL-DRAM: Advantages & Disadvantages

■ Advantages

- + Simple mechanism to reduce latency
- + Significant system performance and energy benefits
 - + Benefits higher at low temperature
- + Low cost, low complexity

■ Disadvantages

- Need to determine reliable operating latencies for different temperatures and different DIMMs → higher testing cost
(might not be that difficult for low temperatures)

More on AL-DRAM

- Donghyuk Lee, Yoongu Kim, Gennady Pekhimenko, Samira Khan, Vivek Seshadri, Kevin Chang, and Onur Mutlu,
"Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common-Case"
Proceedings of the 21st International Symposium on High-Performance Computer Architecture (HPCA), Bay Area, CA, February 2015.
[[Slides \(pptx\) \(pdf\)](#)] [[Full data sets](#)]

Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common-Case

Donghyuk Lee Yoongu Kim Gennady Pekhimenko
Samira Khan Vivek Seshadri Kevin Chang Onur Mutlu
Carnegie Mellon University

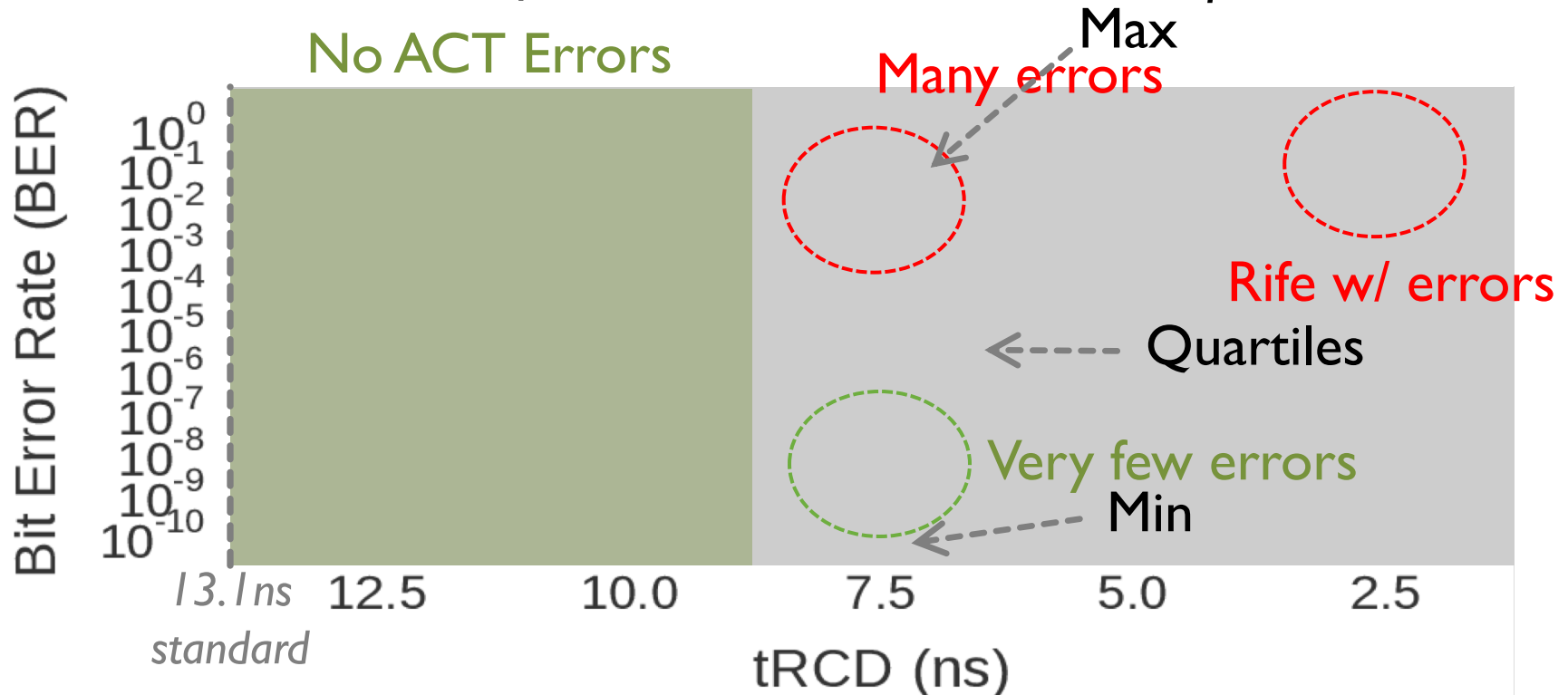
Different Types of Latency Variation

- AL-DRAM exploits latency variation
 - Across time (different temperatures)
 - Across chips

- Is there also latency variation within a chip?
 - Across different parts of a chip

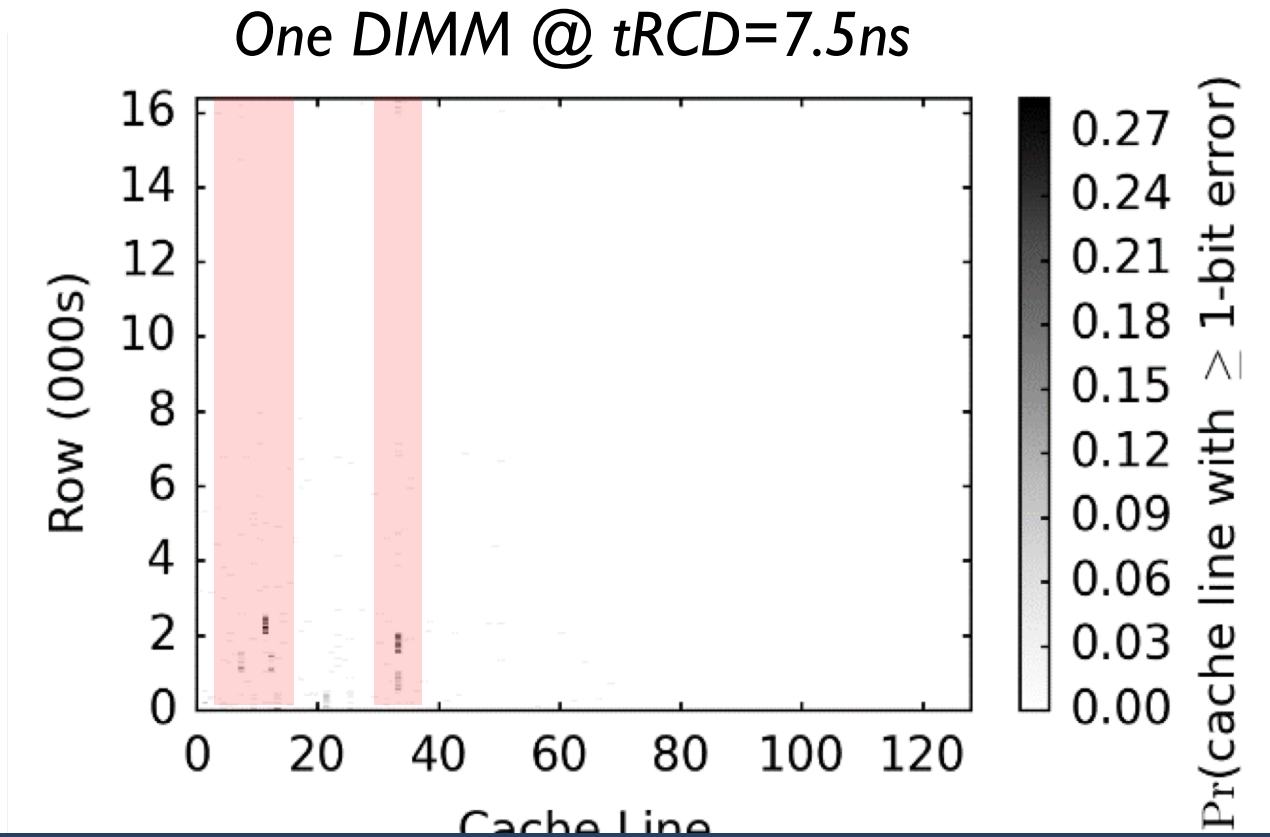
Variation in Activation Errors

Results from 7500 rounds over 240 chips



Modern DRAM chips exhibit significant variation in activation latency

Spatial Locality of Activation Errors

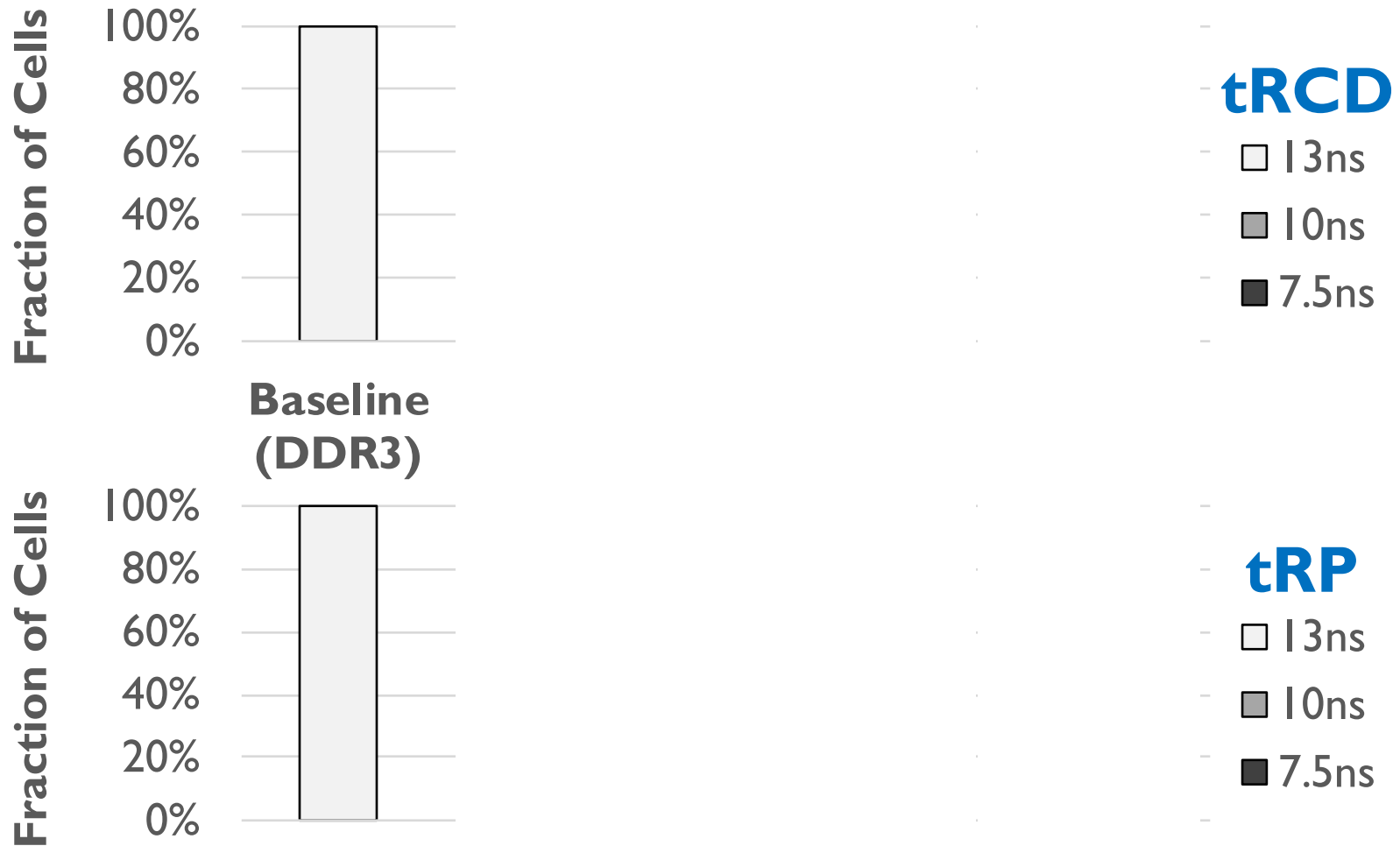


Activation errors are concentrated at certain columns of cells

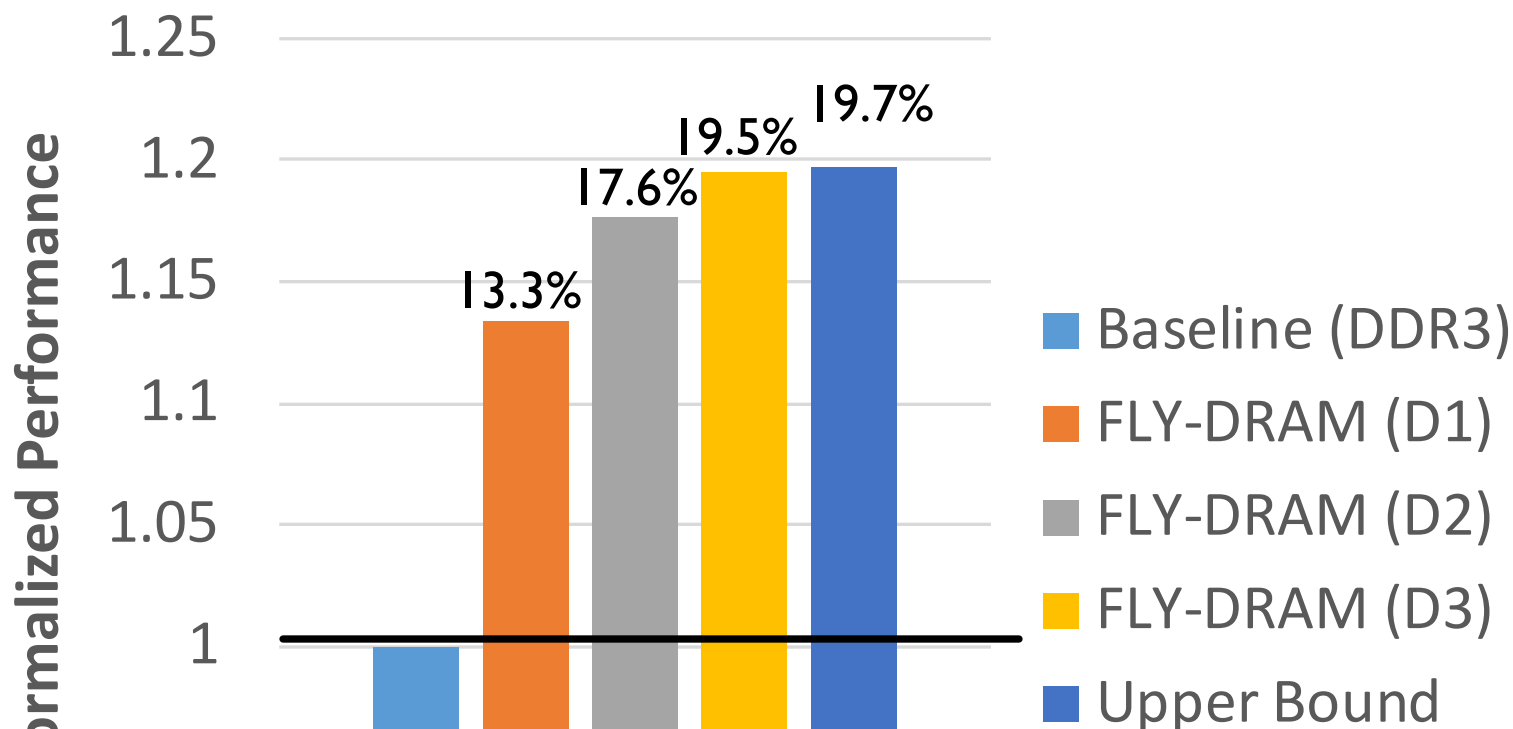
Mechanism to Reduce DRAM Latency

- **Observation:** DRAM timing errors (slow DRAM cells) are concentrated on certain regions
- **Flexible-Latency (FLY) DRAM**
 - A software-transparent design that reduces latency
- **Key idea:**
 - 1) Divide memory into regions of different latencies
 - 2) *Memory controller:* Use lower latency for regions without slow cells; higher latency for other regions

FLY-DRAM Configurations



Results



**FLY-DRAM improves performance
by exploiting spatial latency variation in DRAM**

FLY-DRAM: Advantages & Disadvantages

■ Advantages

- + Reduces latency significantly
- + Exploits significant within-chip latency variation

■ Disadvantages

- Need to determine reliable operating latencies for different parts of a chip → higher testing cost
- Slightly more complicated controller

Analysis of Latency Variation in DRAM Chips

- Kevin Chang, Abhijith Kashyap, Hasan Hassan, Samira Khan, Kevin Hsieh, Donghyuk Lee, Saugata Ghose, Gennady Pekhimenko, Tianshi Li, and Onur Mutlu,

"Understanding Latency Variation in Modern DRAM Chips: Experimental Characterization, Analysis, and Optimization"

*Proceedings of the ACM International Conference on Measurement and Modeling of Computer Systems (**SIGMETRICS**), Antibes Juan-Les-Pins, France, June 2016.*

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

[[Source Code](#)]

Understanding Latency Variation in Modern DRAM Chips: Experimental Characterization, Analysis, and Optimization

Kevin K. Chang¹

Abhijith Kashyap¹

Hasan Hassan^{1,2}

Saugata Ghose¹

Kevin Hsieh¹

Donghyuk Lee¹

Tianshi Li^{1,3}

Gennady Pekhimenko¹

Samira Khan⁴

Onur Mutlu^{5,1}

¹Carnegie Mellon University

²TOBB ETÜ

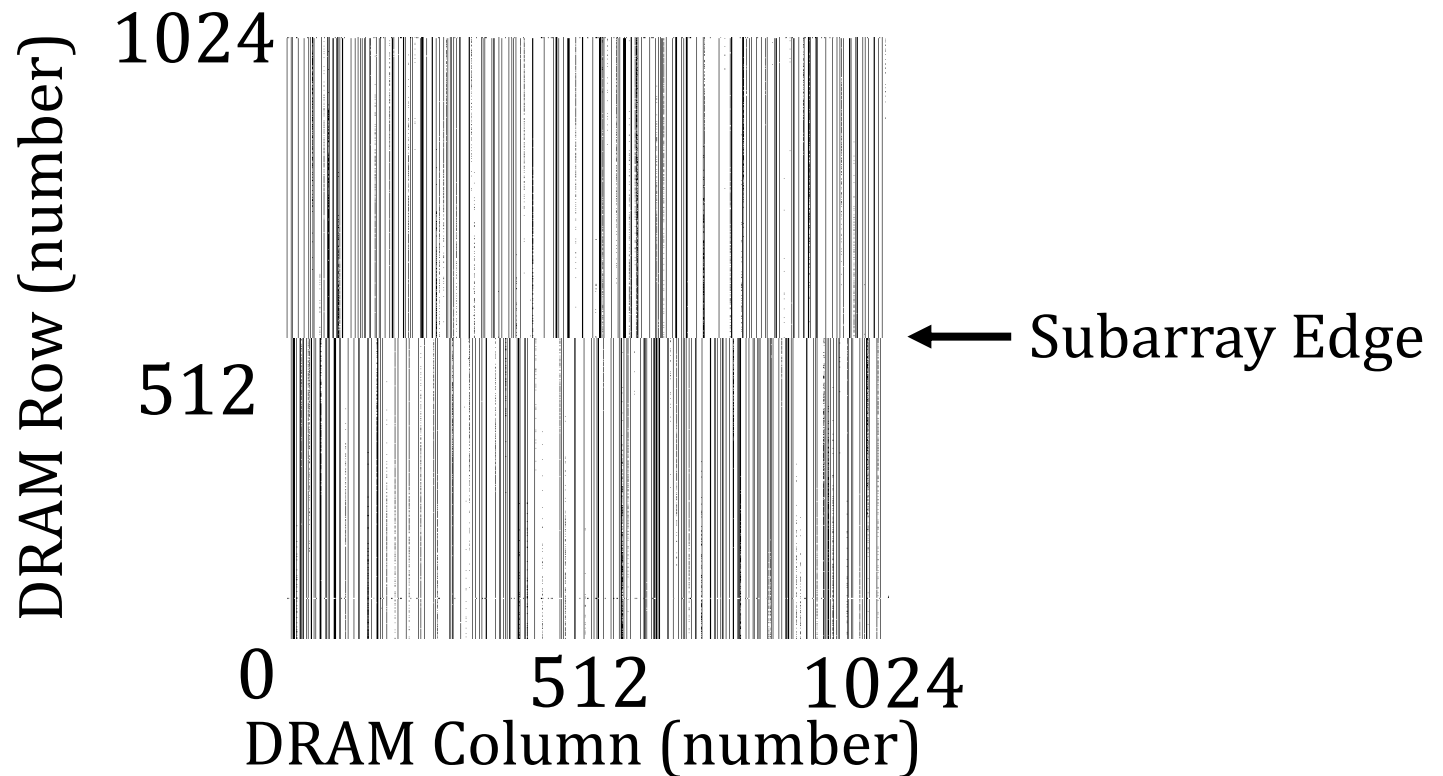
³Peking University

⁴University of Virginia

⁵ETH Zürich

Spatial Distribution of Failures

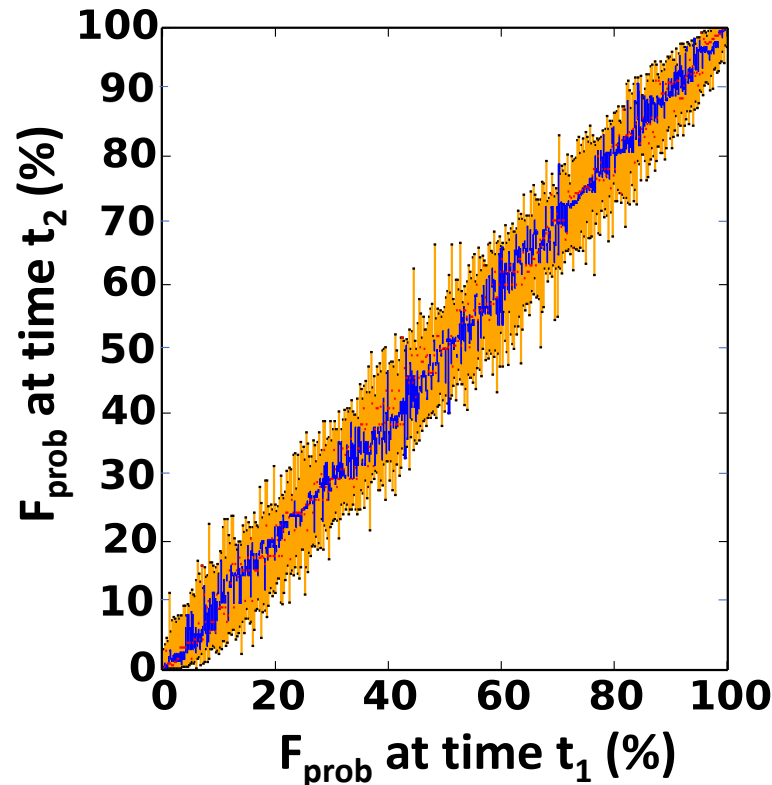
How are activation failures spatially distributed in DRAM?



Activation failures are **highly constrained**
to local bitlines

Short-term Variation

Does a bitline's probability of failure change over time?



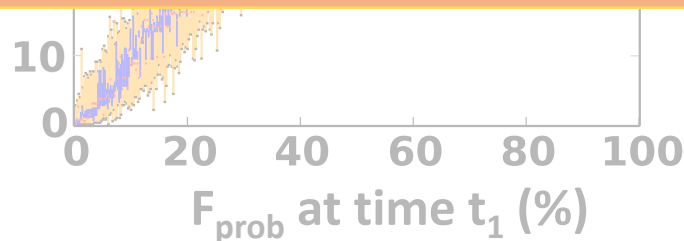
A **weak bitline** is likely to remain **weak** and
a **strong bitline** is likely to remain **strong** over time 82

Short-term Variation

Does a bitline's probability of failure change over time?



We can rely on a **static profile** of weak bitlines to determine whether an access will cause failures

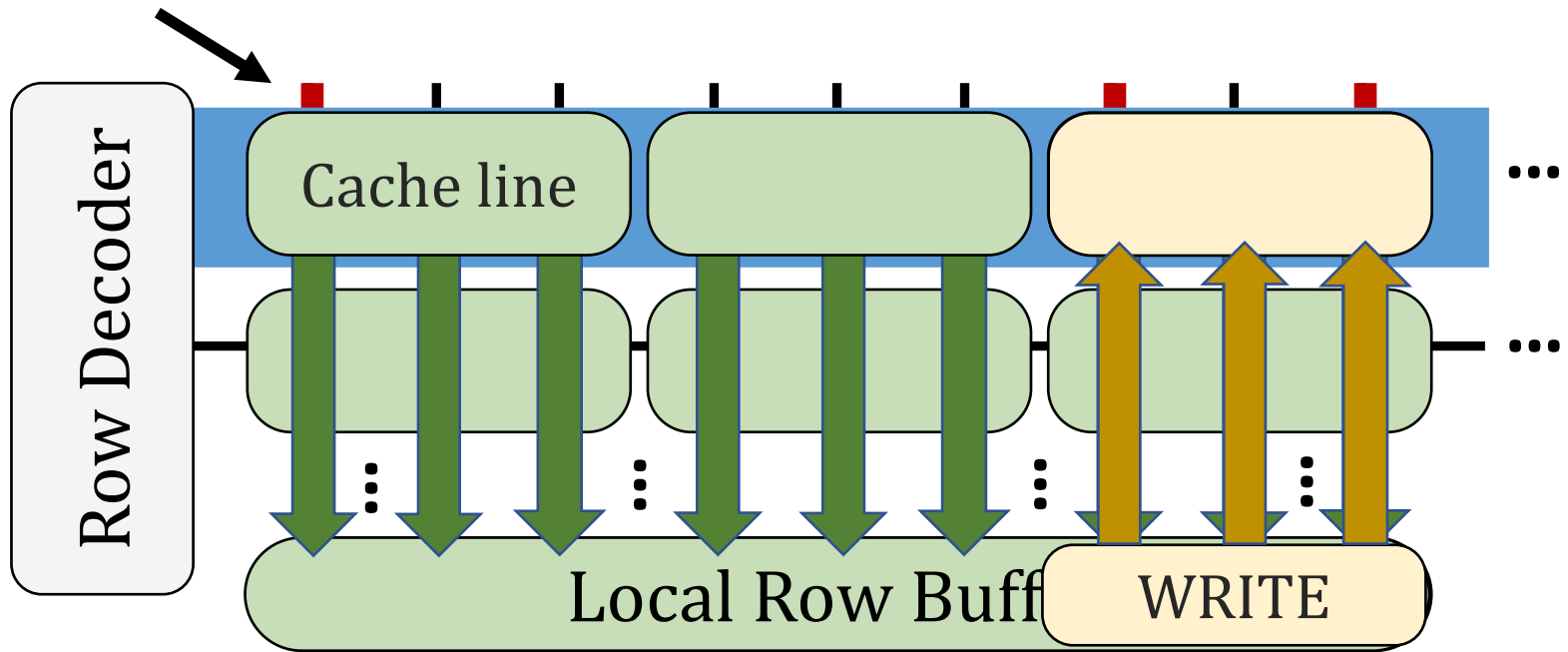


A **weak bitline** is likely to remain **weak** and a **strong bitline** is likely to remain **strong** over time 83

Write Operations

How are write operations affected by reduced t_{RCD} ?

Weak bitline



We can reliably issue write operations
with significantly reduced t_{RCD} (e.g., by 77%)

Solar-DRAM

Uses a **static profile of weak subarray columns**

- Identifies subarray columns as weak or strong
- Obtained in a one-time profiling step

Three Components

1. Variable-latency cache lines (VLC)
2. Reordered subarray columns (RSC)
3. Reduced latency for writes (RLW)

Solar-DRAM

Uses a **static profile of weak subarray columns**

- Identifies subarray columns as weak or strong
- Obtained in a one-time profiling step

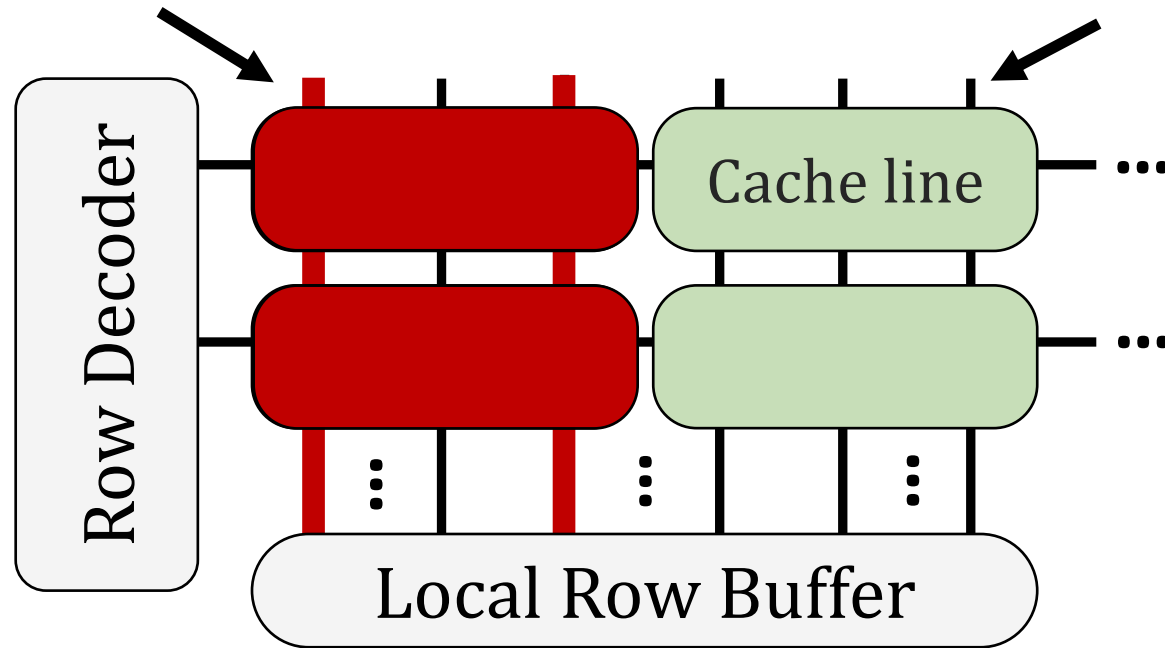
Three Components

1. Variable-latency cache lines (VLC)
2. Reordered subarray columns (RSC)
3. Reduced latency for writes (RLW)

Solar-DRAM: VLC (I)

Weak bitline

Strong bitline



Identify cache lines comprised of **strong bitlines**

Access such cache lines with a **reduced t_{RCD}**

Solar-DRAM

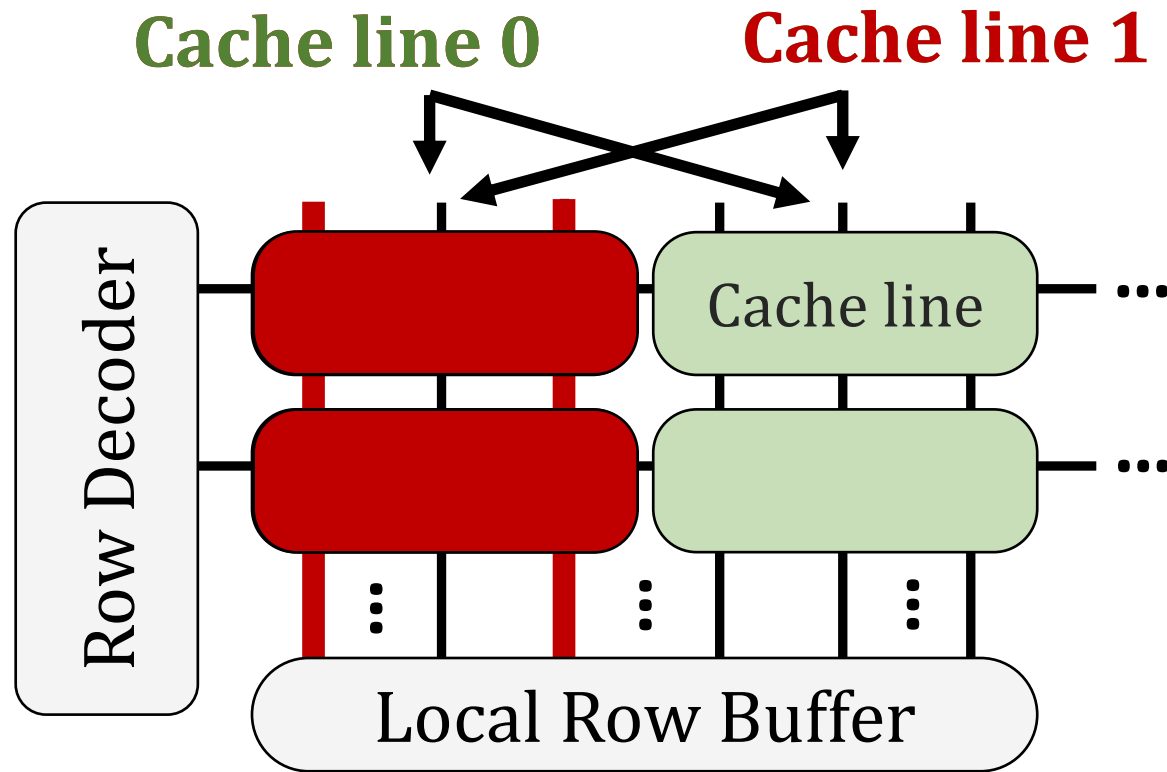
Uses a **static profile of weak subarray columns**

- Identifies subarray columns as weak or strong
- Obtained in a one-time profiling step

Three Components

1. Variable-latency cache lines (VLC)
2. Reordered subarray columns (RSC)
3. Reduced latency for writes (RLW)

Solar-DRAM: RSC (II)



Remap cache lines across DRAM at the memory controller level so cache line 0 will likely map to a **strong** cache line

Solar-DRAM

Uses a **static profile of weak subarray columns**

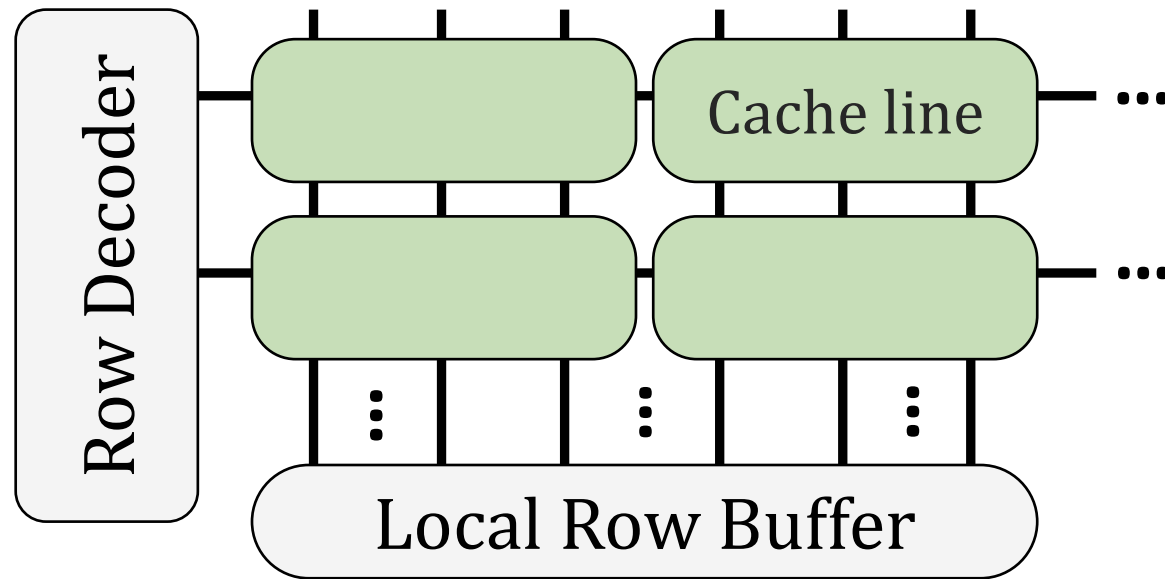
- Identifies subarray columns as weak or strong
- Obtained in a one-time profiling step

Three Components

1. Variable-latency cache lines (VLC)
2. Reordered subarray columns (RSC)
3. Reduced latency for writes (RLW)

Solar-DRAM: RLW (III)

All bitlines are strong when issuing writes



Write to all locations in DRAM with a significantly reduced t_{RCD} (e.g., by 77%)

More on Solar-DRAM

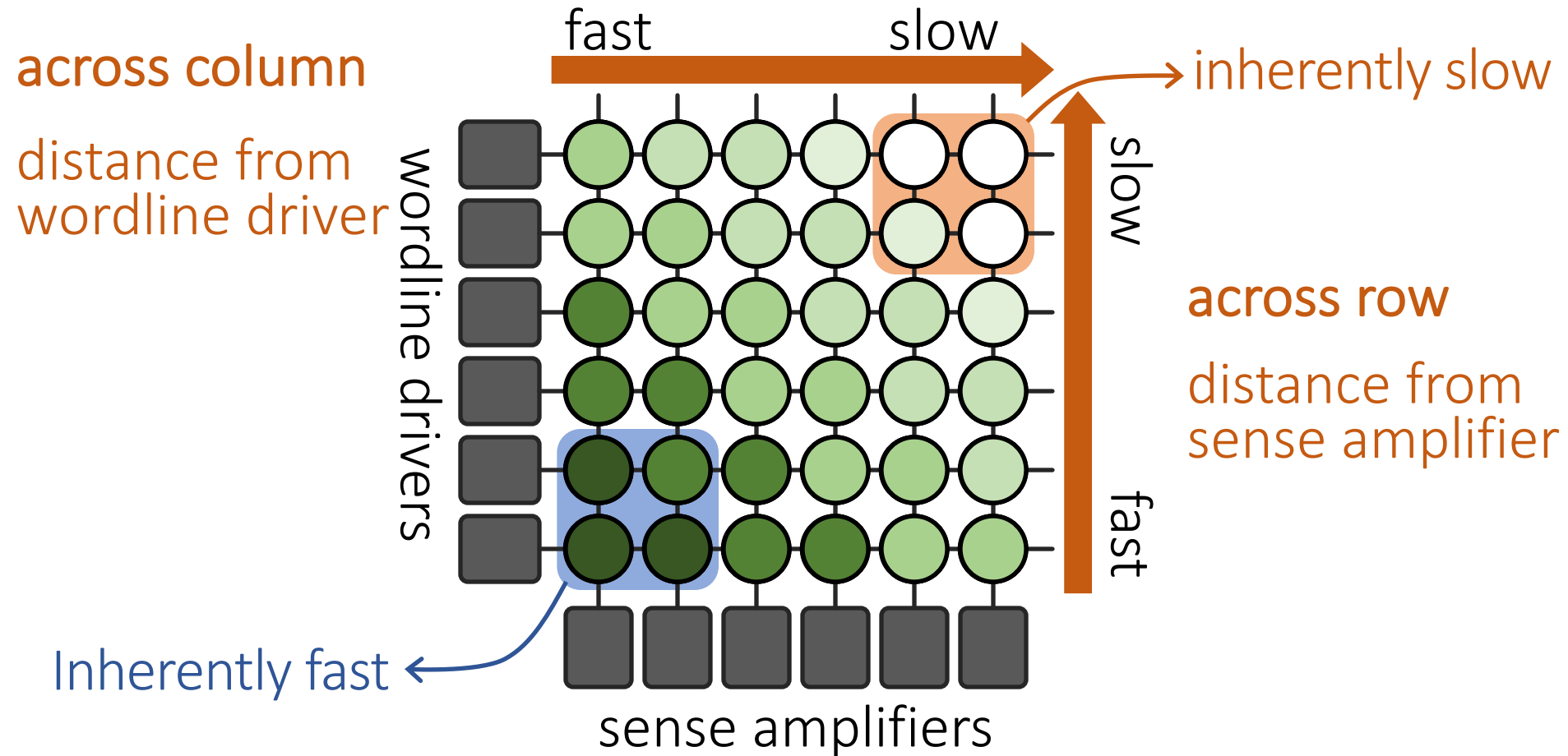
- Jeremie S. Kim, Minesh Patel, Hasan Hassan, and Onur Mutlu,
"Solar-DRAM: Reducing DRAM Access Latency by Exploiting the Variation in Local Bitlines"
Proceedings of the 36th IEEE International Conference on Computer Design (ICCD), Orlando, FL, USA, October 2018.

Solar-DRAM: Reducing DRAM Access Latency by Exploiting the Variation in Local Bitlines

Jeremie S. Kim^{‡§} Minesh Patel[§] Hasan Hassan[§] Onur Mutlu^{§‡}
 ‡Carnegie Mellon University §ETH Zürich

Why Is There Spatial Latency Variation Within a Chip?

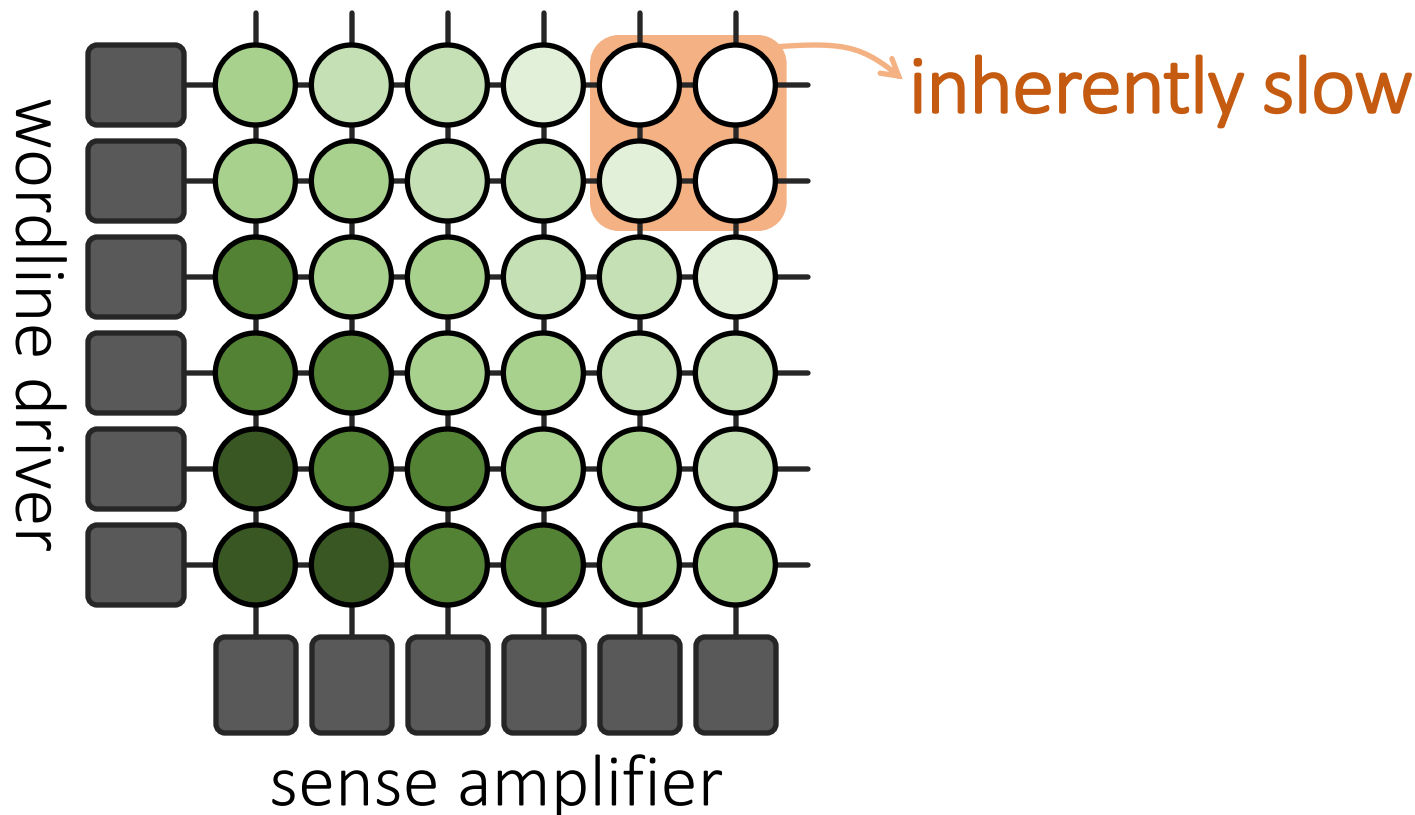
What Is Design-Induced Variation?



Systematic variation in cell access times
caused by the ***physical organization*** of DRAM

DIVA Online Profiling

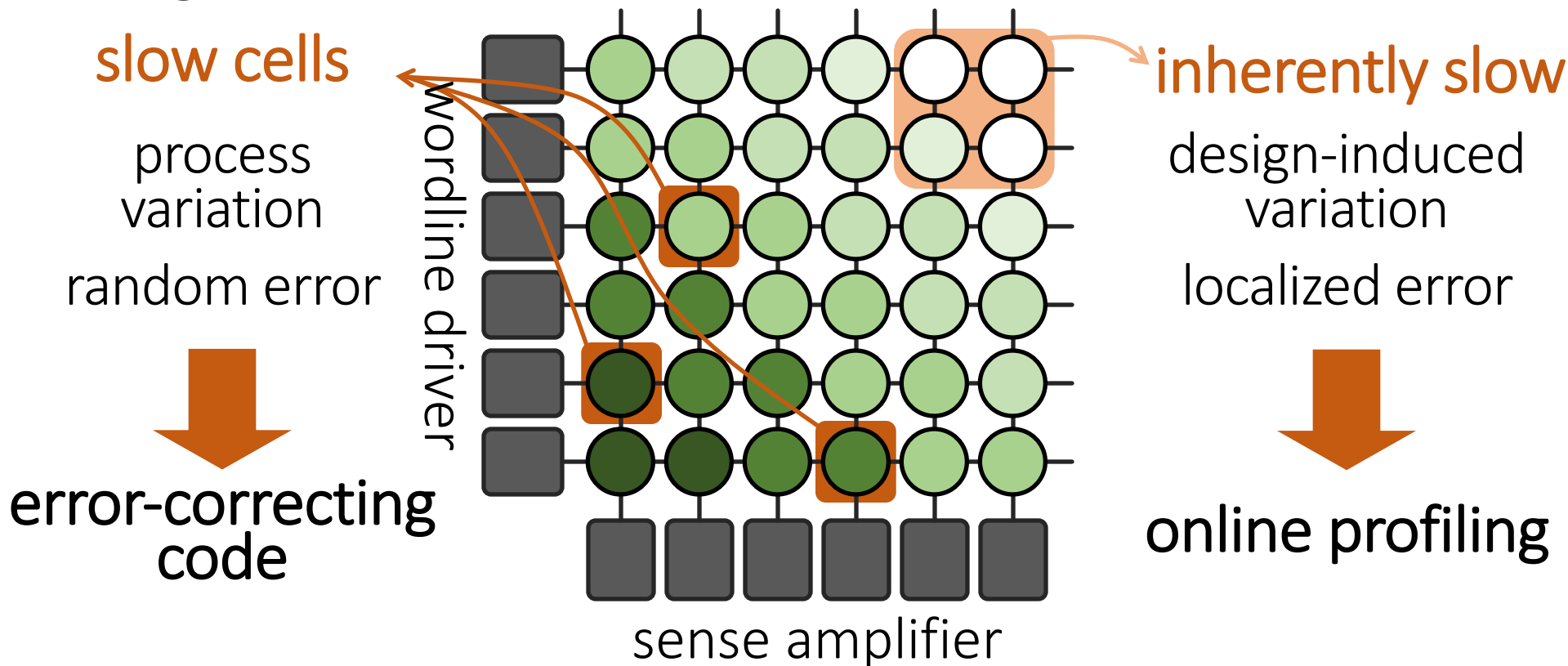
Design-Induced-Variation-Aware



Profile *only slow regions* to determine min. latency
→ *Dynamic* & *low cost* latency optimization

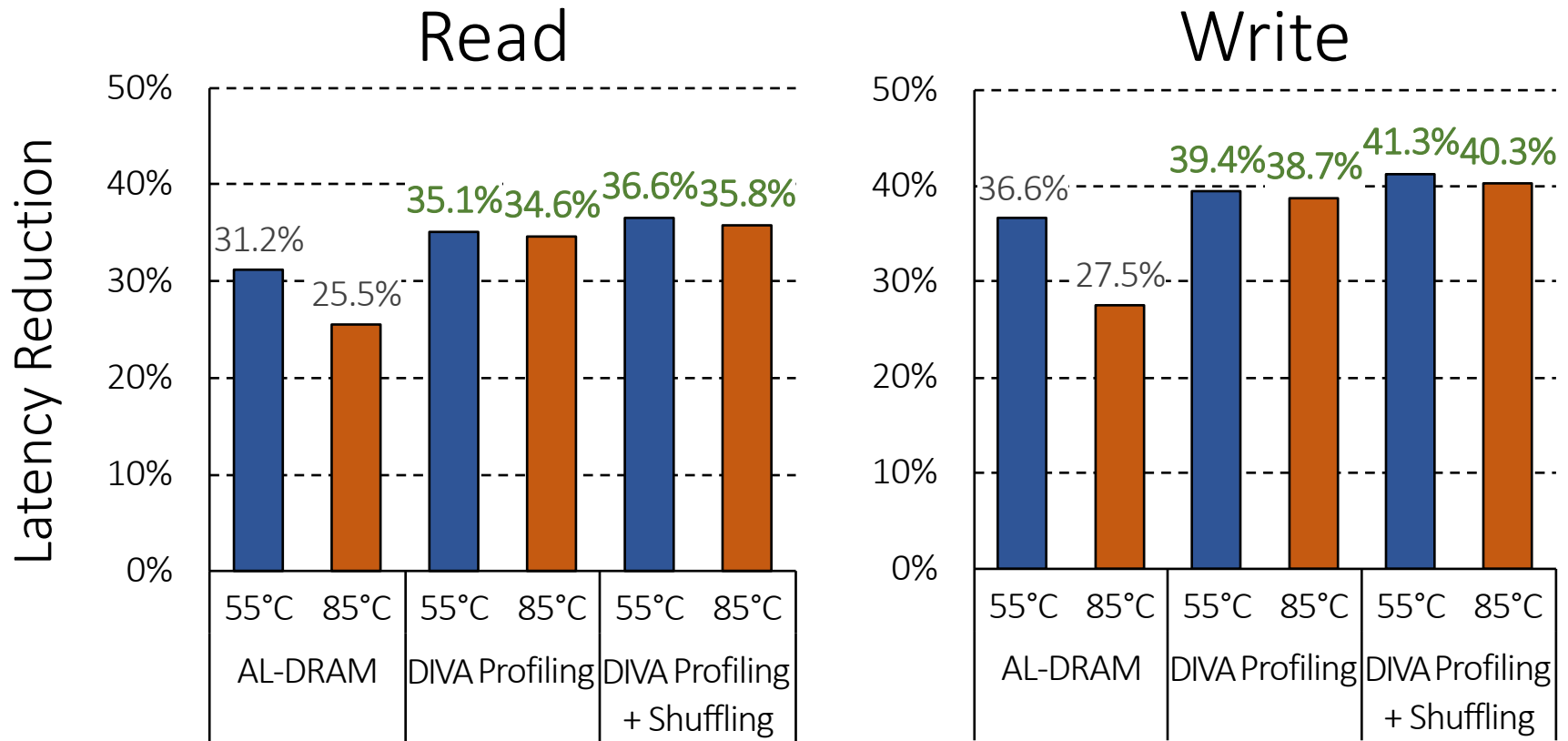
DIVA Online Profiling

Design-Induced-Variation-Aware



Combine **error-correcting codes** & **online profiling**
→ **Reliably** reduce DRAM latency

DIVA-DRAM Reduces Latency



DIVA-DRAM *reduces latency more aggressively*
and uses ECC to correct random slow cells

DIVA-DRAM: Advantages & Disadvantages

■ Advantages

- ++ Automatically finds the lowest reliable operating latency at system runtime (lower production-time testing cost)
- + Reduces latency more than prior methods (w/ ECC)
- + Reduces latency at high temperatures as well

■ Disadvantages

- Requires knowledge of inherently-slow regions
- Requires ECC (Error Correcting Codes)
- Imposes overhead during runtime profiling

Design-Induced Latency Variation in DRAM

- Donghyuk Lee, Samira Khan, Lavanya Subramanian, Saugata Ghose, Rachata Ausavarungnirun, Gennady Pekhimenko, Vivek Seshadri, and Onur Mutlu,
"Design-Induced Latency Variation in Modern DRAM Chips: Characterization, Analysis, and Latency Reduction Mechanisms"
*Proceedings of the ACM International Conference on Measurement and Modeling of Computer Systems (**SIGMETRICS**), Urbana-Champaign, IL, USA, June 2017.*

Design-Induced Latency Variation in Modern DRAM Chips: Characterization, Analysis, and Latency Reduction Mechanisms

Donghyuk Lee, NVIDIA and Carnegie Mellon University

Samira Khan, University of Virginia

Lavanya Subramanian, Saugata Ghose, Rachata Ausavarungnirun, Carnegie Mellon University

Gennady Pekhimenko, Vivek Seshadri, Microsoft Research

Onur Mutlu, ETH Zürich and Carnegie Mellon University

Understanding & Exploiting the Voltage-Latency-Reliability Relationship

High DRAM Power Consumption

- Problem: High DRAM (memory) power in today's systems



>40% in POWER7 (Ware+, HPCA'10)



>40% in GPU (Paul+, ISCA'15)

Low-Voltage Memory

- Existing DRAM designs to help reduce DRAM power by lowering supply voltage conservatively
 - $Power \propto Voltage^2$
- DDR3L (low-voltage) reduces voltage from 1.5V to 1.35V (-10%)
- LPDDR4 (low-power) employs low-power I/O interface with 1.2V (lower bandwidth)

Can we reduce DRAM power and energy by further reducing supply voltage?

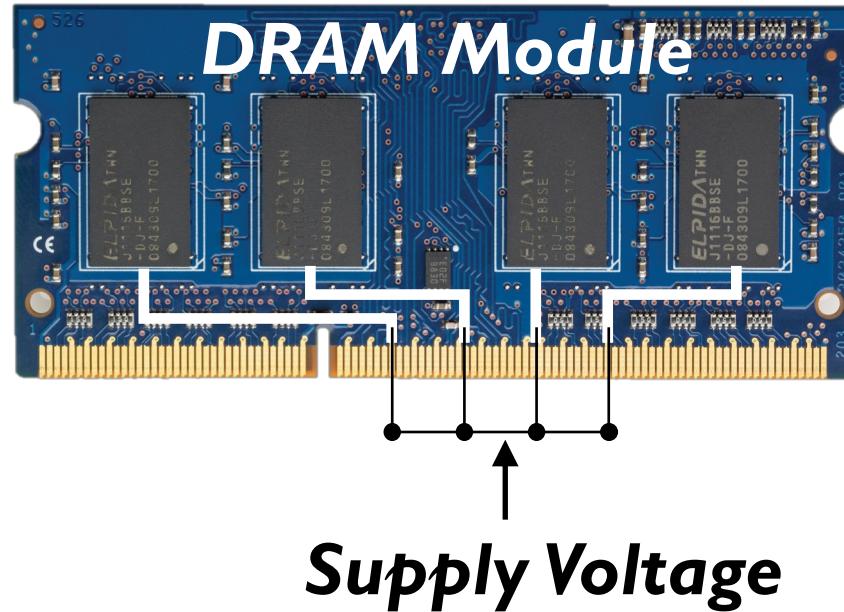
Goals

- 1 Understand and characterize the various characteristics of DRAM under **reduced voltage**
- 2 Develop a mechanism that reduces DRAM energy by **lowering voltage** while keeping performance loss within a target

Key Questions

- How does reducing voltage affect ***reliability*** (errors)?
- How does reducing voltage affect ***DRAM latency***?
- How do we design a new DRAM energy reduction mechanism?

Supply Voltage Control on DRAM



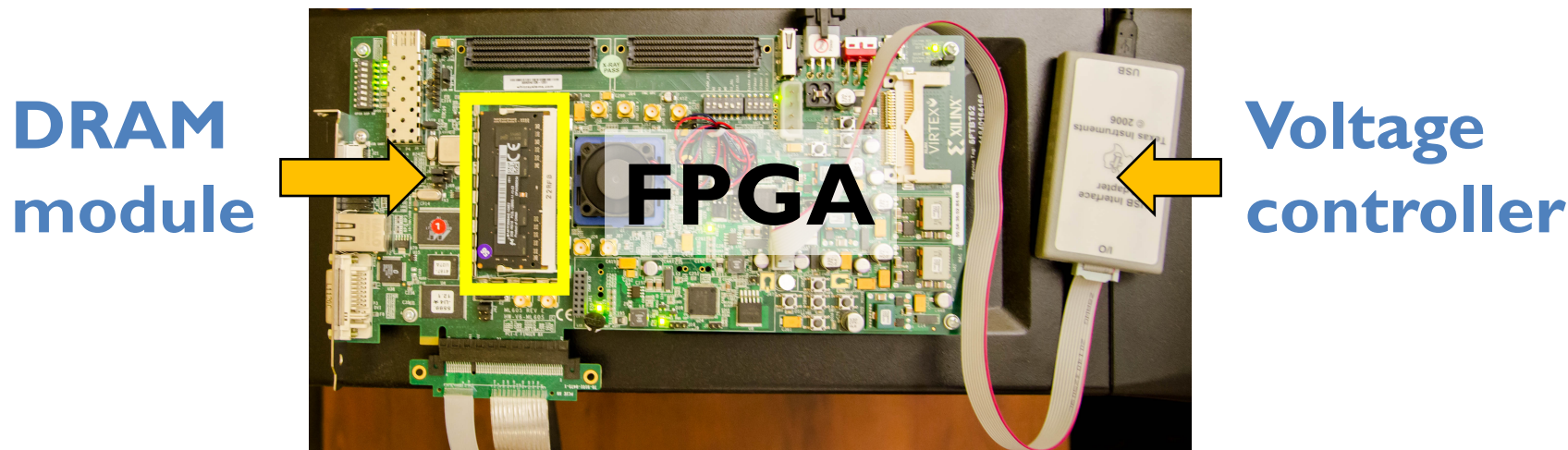
Adjust the *supply voltage* to every chip on the same module

Custom Testing Platform

SoftMC [Hassan+, HPCA'17]: FPGA testing platform to

- 1) Adjust supply voltage to DRAM modules
- 2) Schedule DRAM commands to DRAM modules

Existing systems: DRAM commands not exposed to users

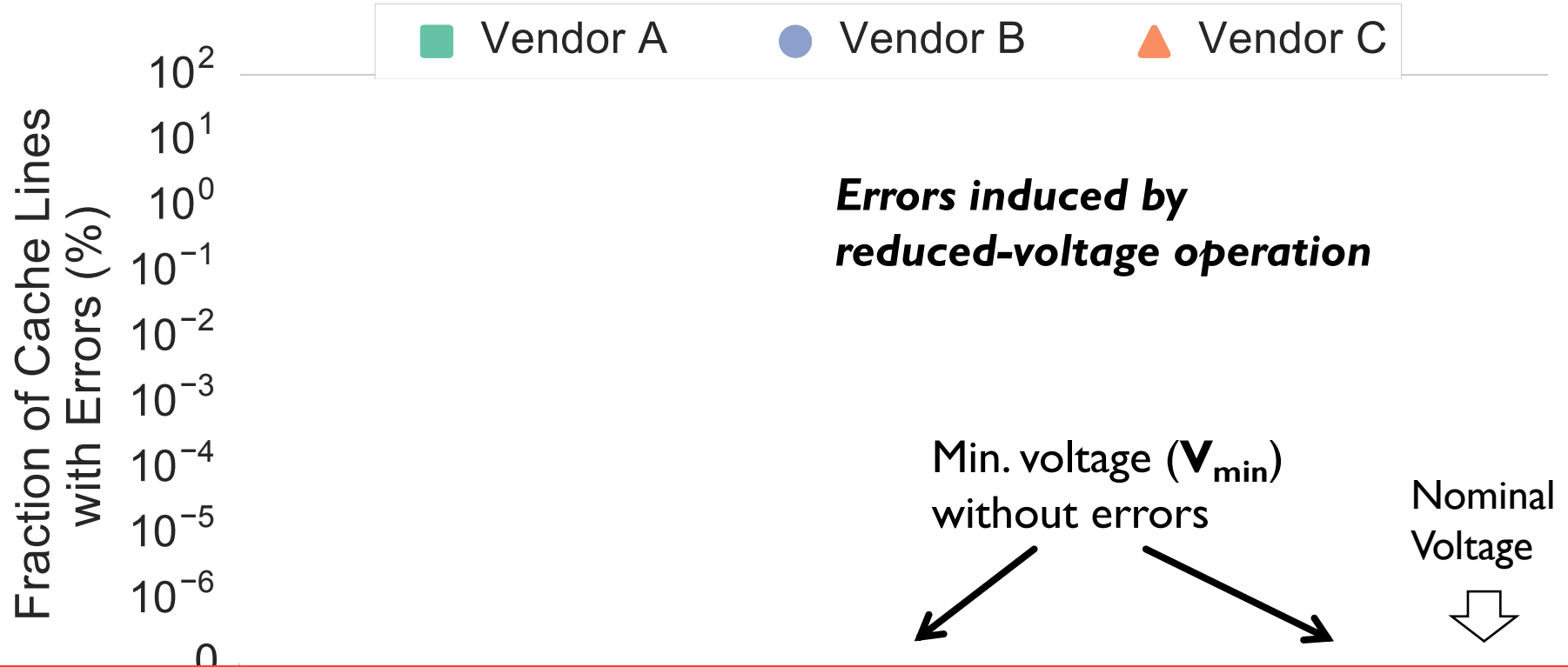


<https://github.com/CMU-SAFARI/DRAM-Voltage-Study>

Tested DRAM Modules

- **124 DDR3L** (low-voltage) DRAM chips
 - **31 SO-DIMMs**
 - **1.35V** (DDR3 uses 1.5V)
 - Density: 4Gb per chip
 - Three major vendors/manufacturers
 - Manufacturing dates: 2014-2016
- Iteratively read every bit in each 4Gb chip under a wide range of supply voltage levels: 1.35V to 1.0V (**-26%**)

Reliability Worsens with Lower Voltage

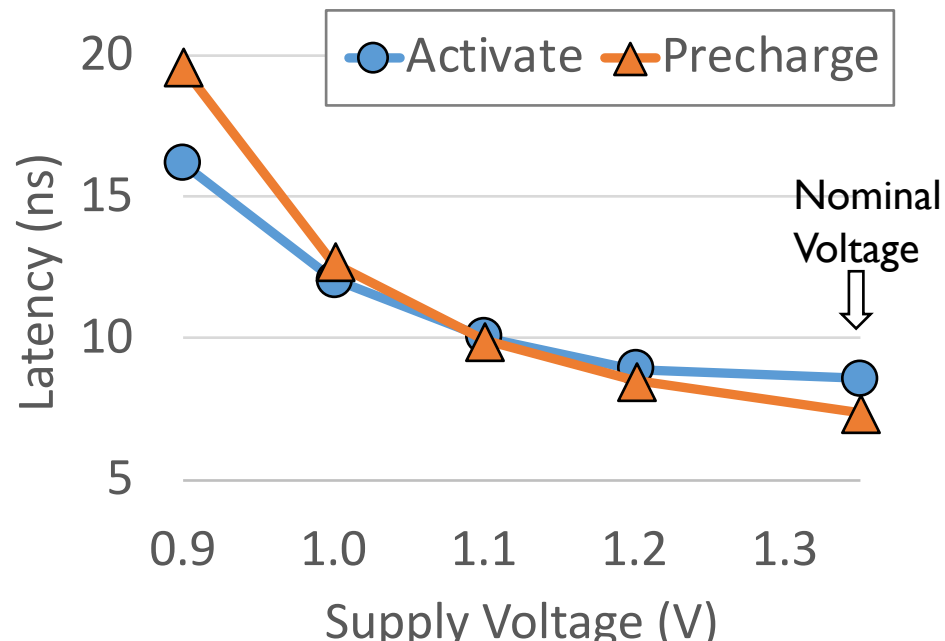
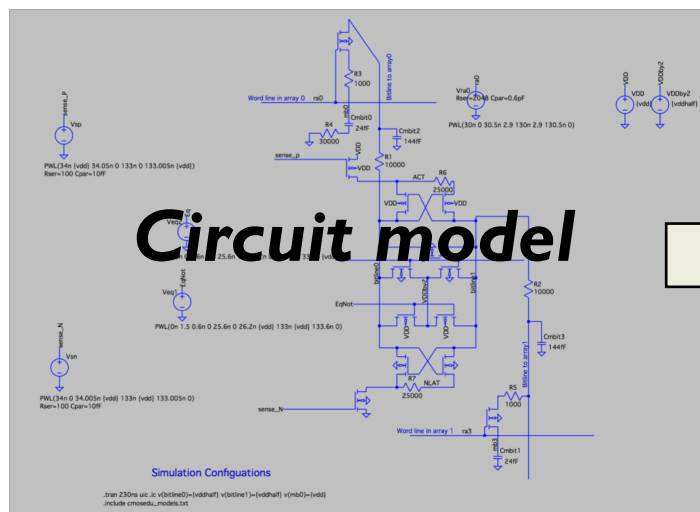


Reducing voltage below V_{min} causes an increasing number of errors

Source of Errors

Detailed circuit simulations (SPICE) of a DRAM cell array to model the behavior of DRAM operations

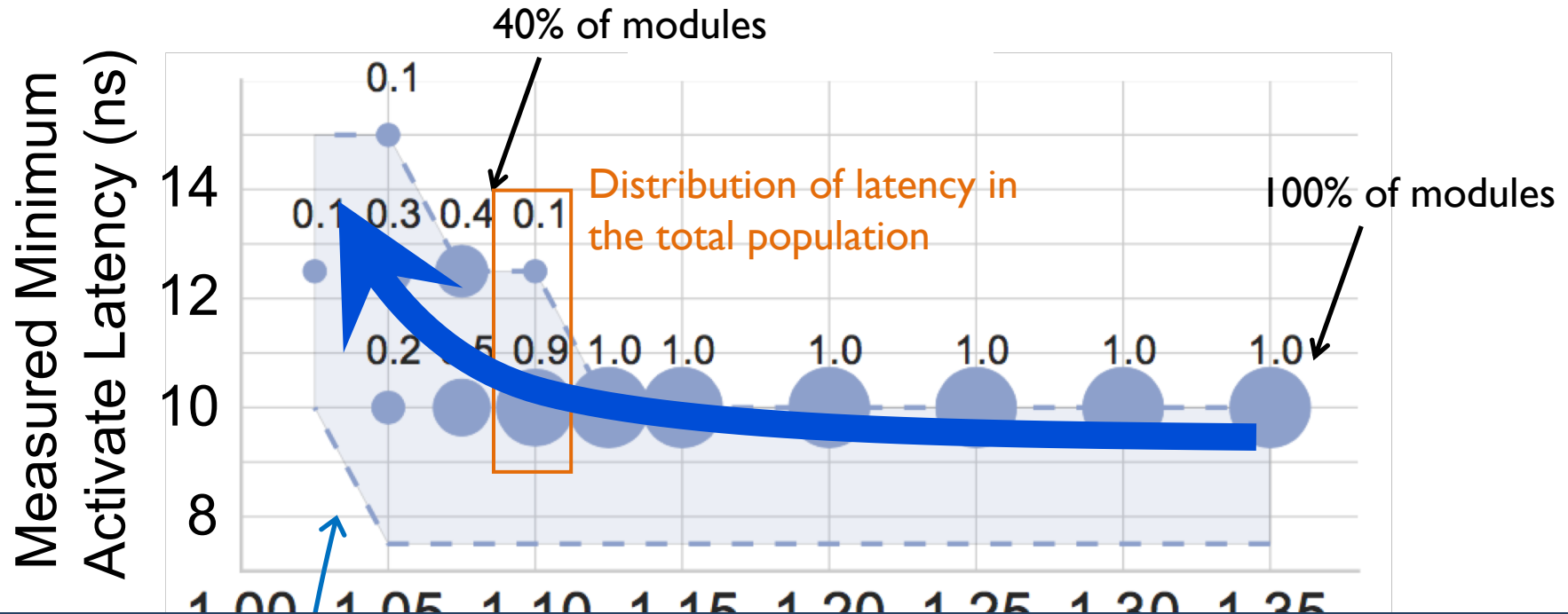
<https://github.com/CMU-SAFARI/DRAM-Voltage-Study>



Reliable low-voltage operation requires higher latency

DIMMs Operating at Higher Latency

Measured minimum latency that *does not* cause errors in DRAM modules



DRAM requires longer latency to access data **without errors** at lower voltage

Spatial Locality of Errors



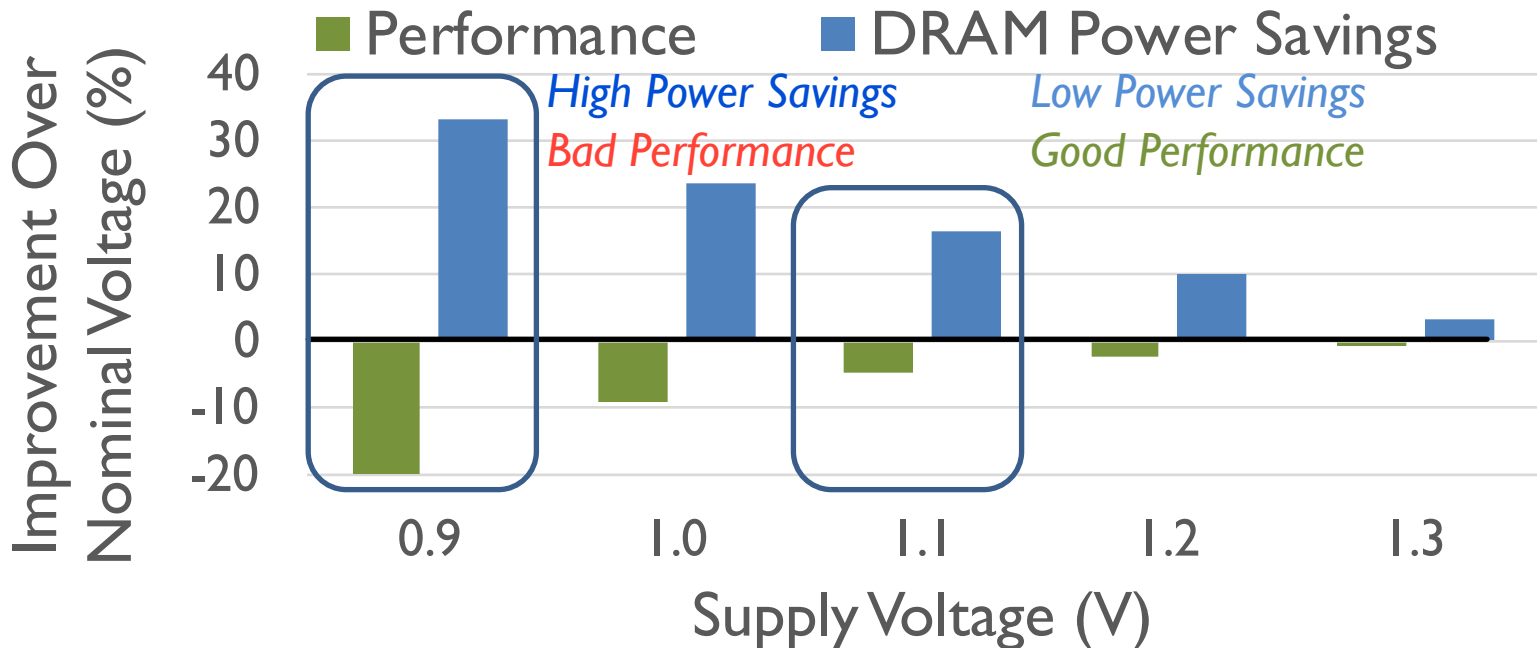
Errors concentrate in certain regions

Summary of Key Experimental Observations

- Voltage-induced errors increase as voltage reduces further below V_{\min}
- Errors exhibit spatial locality
- Increasing the latency of DRAM operations mitigates voltage-induced errors

DRAM Voltage Adjustment to Reduce Energy

- Goal: Exploit the trade-off between voltage and latency to reduce energy consumption
- Approach: Reduce DRAM voltage **reliably**
 - **Performance loss** due to increased latency at lower voltage

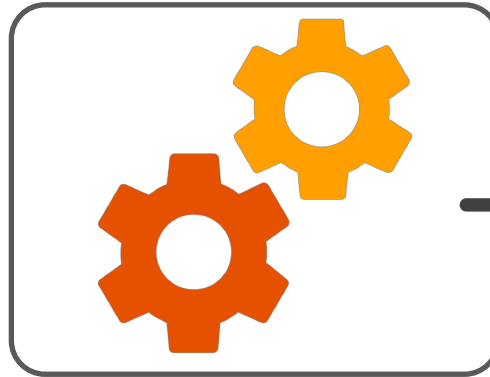


Voltron Overview

Voltron



User specifies the
performance loss target

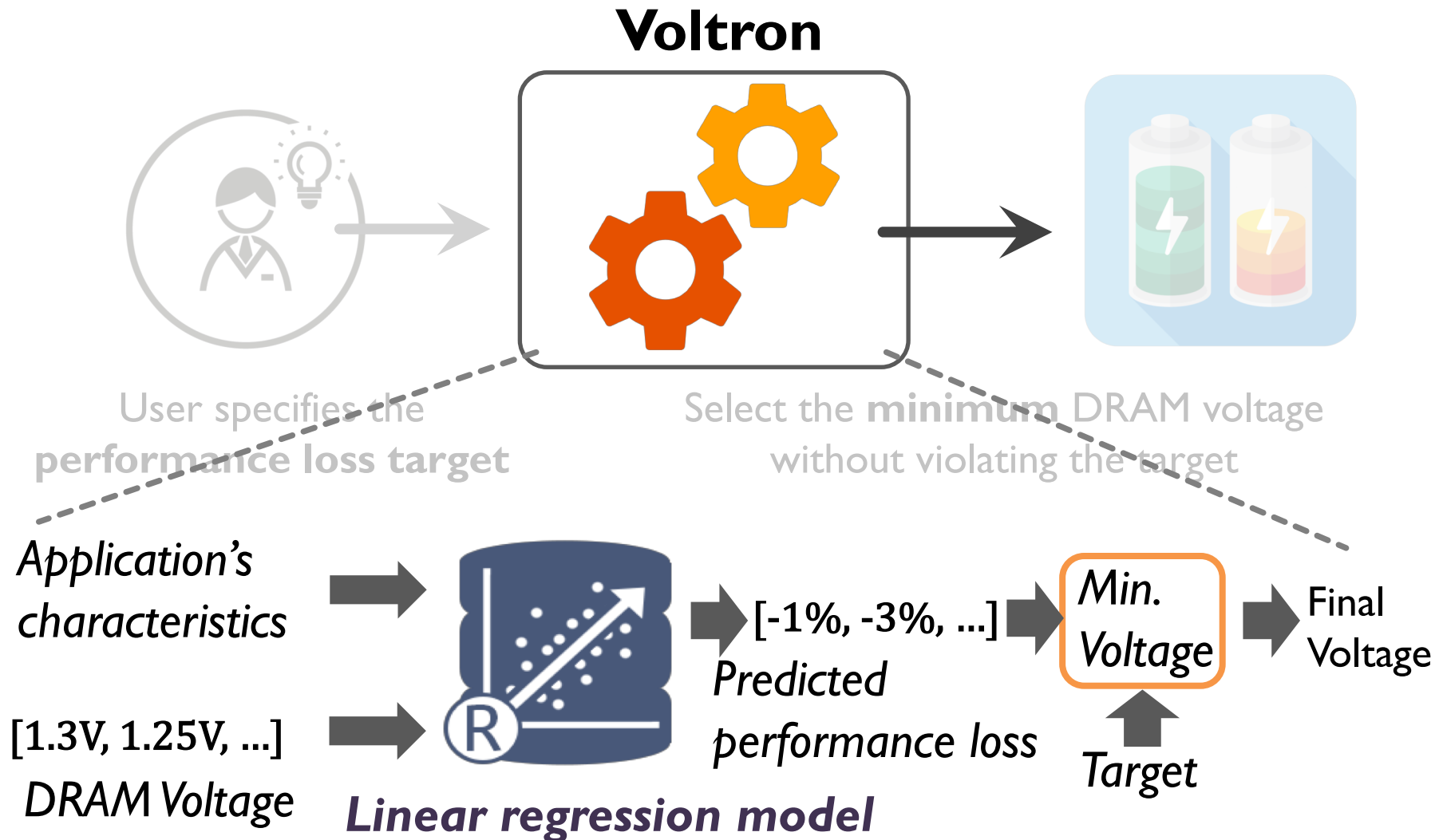


Select the **minimum** DRAM voltage
without violating the target



How do we predict performance loss due to increased latency under low DRAM voltage?

Linear Model to Predict Performance

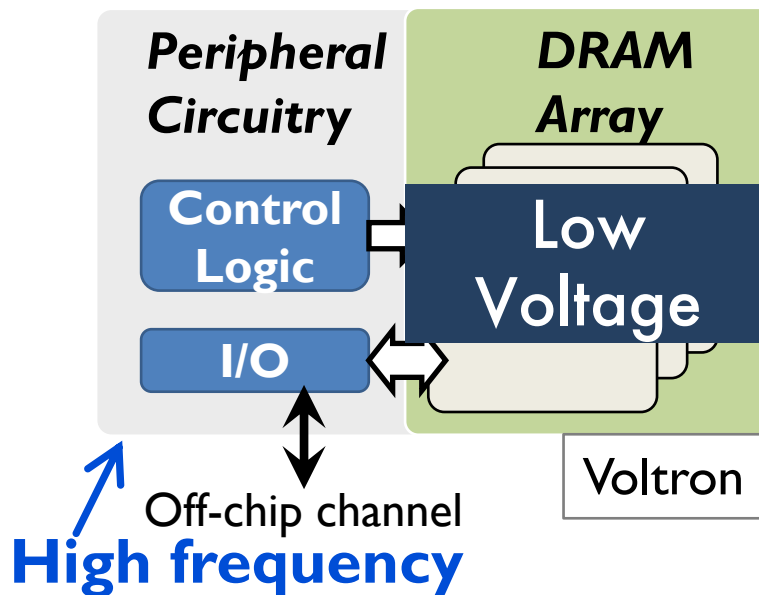
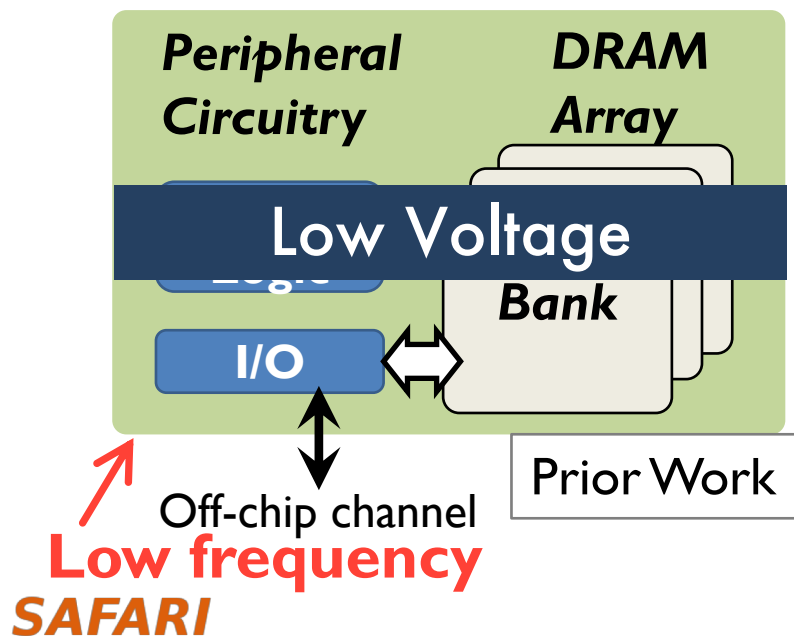


Regression Model to Predict Performance

- Application's characteristics for the model:
 - **Memory intensity**: Frequency of last-level cache misses
 - **Memory stall time**: Amount of time memory requests stall commit inside CPU
- Handling multiple applications:
 - Predict a performance loss for each application
 - Select the minimum voltage that satisfies the performance target for all applications

Comparison to Prior Work

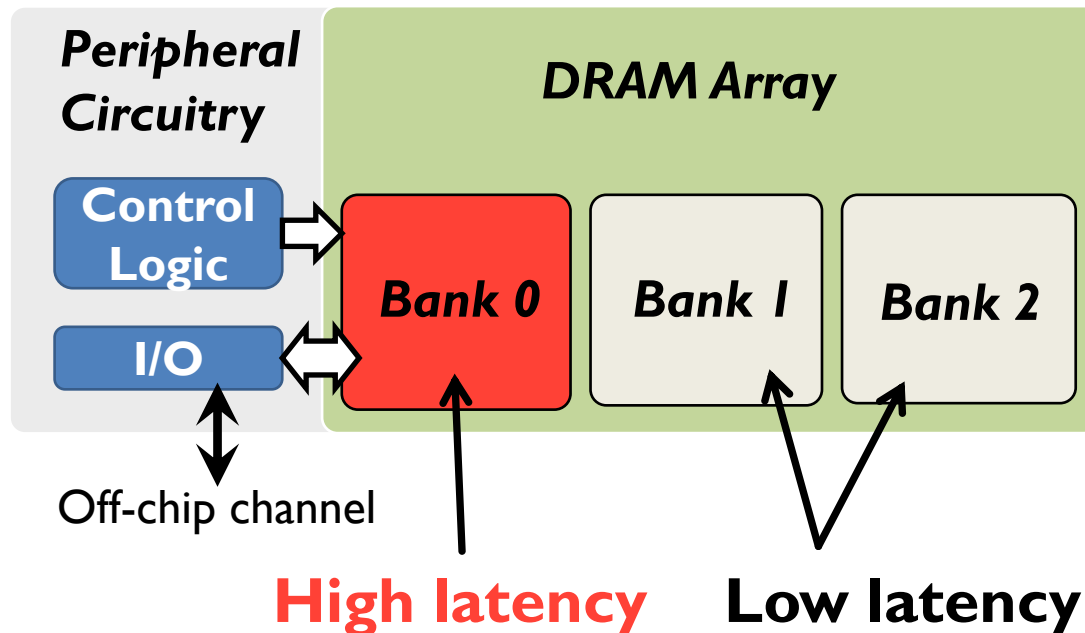
- Prior work: Dynamically scale *frequency and voltage* of the entire DRAM based on bandwidth demand [David+, ICAC'11]
 - Problem: Lowering voltage on the peripheral circuitry decreases channel frequency (memory data throughput)
- Voltron: Reduce voltage to only **DRAM array** without changing the voltage to peripheral circuitry



Exploiting Spatial Locality of Errors

Key idea: Increase the latency only for DRAM banks that observe errors under low voltage

- Benefit: Higher performance

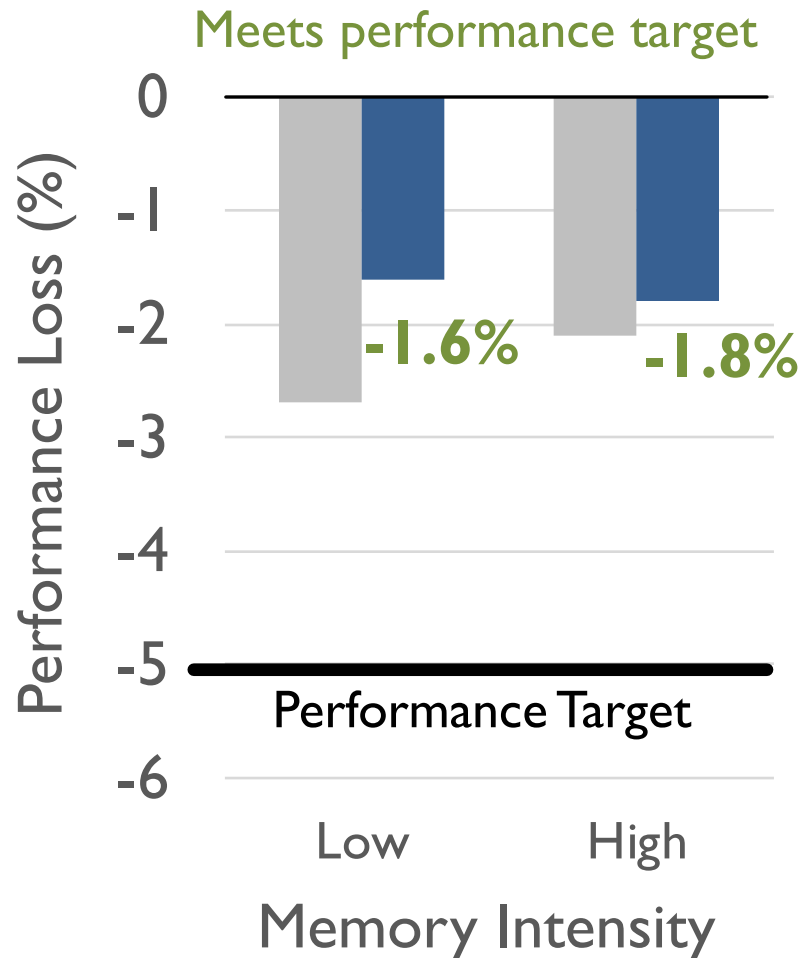
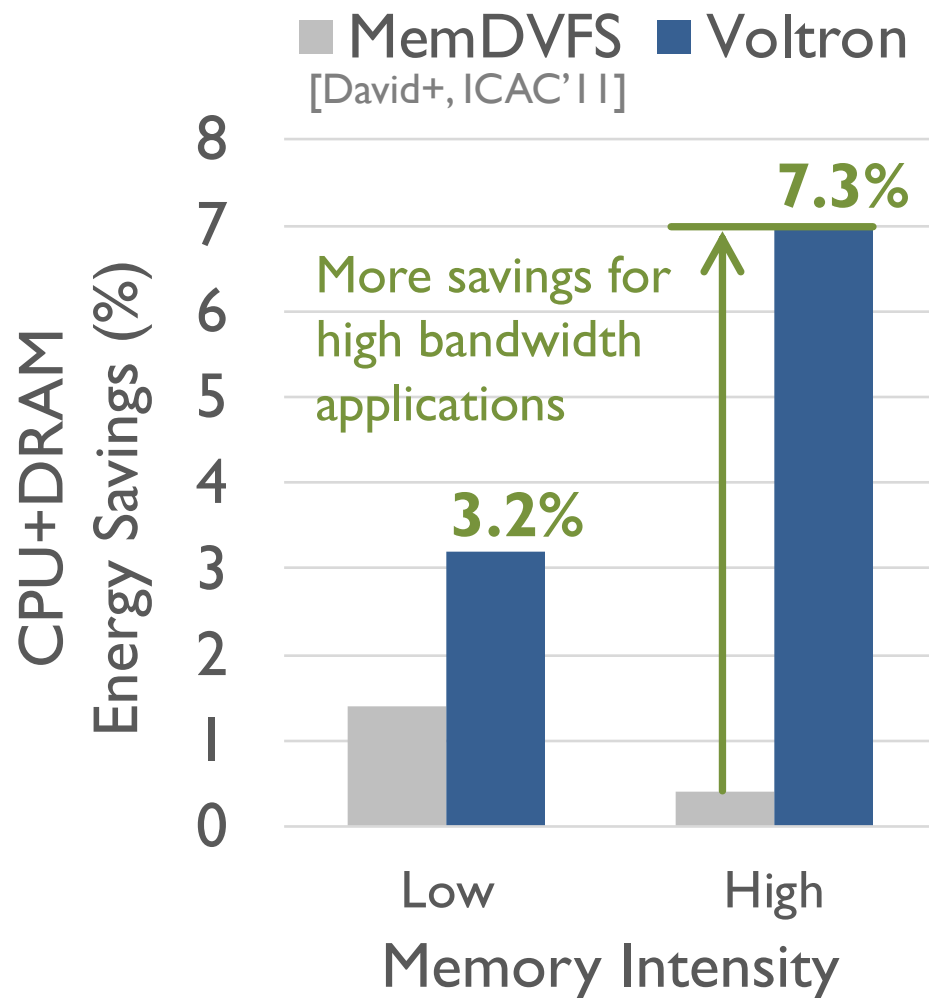


Voltron Evaluation Methodology

- **Cycle-level simulator:** Ramulator [CAL'15]
 - **McPAT** and **DRAMPower** for energy measurement

<https://github.com/CMU-SAFARI/ramulator>
- **4-core** system with DDR3L memory
- **Benchmarks:** SPEC2006, YCSB
- Comparison to prior work: **MemDVFS** [David+, ICAC'11]
 - Dynamic DRAM frequency and voltage scaling
 - Scaling based on the *memory bandwidth consumption*

Energy Savings with Bounded Performance



Voltron: Advantages & Disadvantages

■ Advantages

- + Can trade-off between voltage and latency to improve energy or performance
- + Can exploit the high voltage margin present in DRAM

■ Disadvantages

- Requires finding the reliable operating voltage for each chip → higher testing cost

Analysis of Latency-Voltage in DRAM Chips

- Kevin Chang, A. Giray Yaglikci, Saugata Ghose, Aditya Agrawal, Niladrish Chatterjee, Abhijith Kashyap, Donghyuk Lee, Mike O'Connor, Hasan Hassan, and Onur Mutlu,

"Understanding Reduced-Voltage Operation in Modern DRAM Devices: Experimental Characterization, Analysis, and Mechanisms"

*Proceedings of the ACM International Conference on Measurement and Modeling of Computer Systems (**SIGMETRICS**), Urbana-Champaign, IL, USA, June 2017.*

Understanding Reduced-Voltage Operation in Modern DRAM Chips: Characterization, Analysis, and Mechanisms

Kevin K. Chang[†] Abdullah Giray Yağlıkçı[†] Saugata Ghose[†] Aditya Agrawal[¶] Niladrish Chatterjee[¶]
Abhijith Kashyap[†] Donghyuk Lee[¶] Mike O'Connor^{¶,‡} Hasan Hassan[§] Onur Mutlu^{§,†}

[†]Carnegie Mellon University

[¶]NVIDIA

[‡]The University of Texas at Austin

[§]ETH Zürich

And, What If ...

- ... we can sacrifice reliability of some data to access it with even lower latency?

The DRAM Latency PUF:

Quickly Evaluating Physical Unclonable Functions
by Exploiting the Latency-Reliability Tradeoff
in Modern Commodity DRAM Devices

Jeremie S. Kim Minesh Patel

Hasan Hassan Onur Mutlu



SAFARI

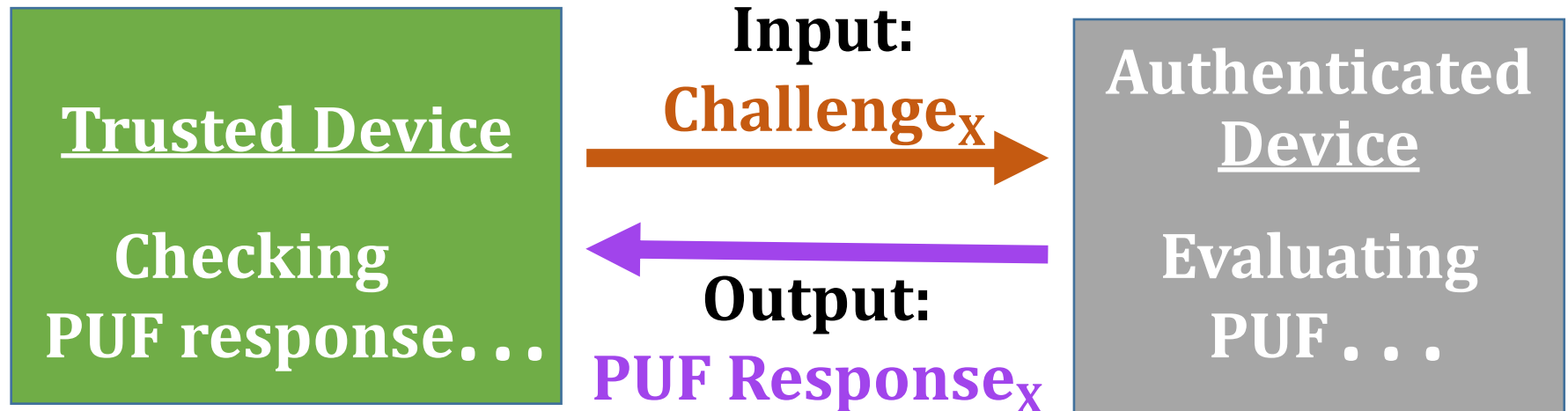
ETH zürich

Carnegie Mellon

Motivation

We want a way to ensure that a system's components are not **compromised**

- **Physical Unclonable Function (PUF)**: a function we **evaluate** on a device to **generate** a **signature** **unique** to the device
- We refer to the unique signature as a **PUF response**
- Often used in a **Challenge-Response Protocol (CRP)**



Motivation

1. We want a **runtime-accessible** PUF
 - Should be evaluated **quickly** with **minimal** impact on concurrent applications
 - Can protect against **attacks that swap system components with malicious parts**
2. DRAM is a **promising substrate** for evaluating PUFs because it is **ubiquitous** in modern systems
 - Unfortunately, current DRAM PUFs are **slow** and get **exponentially slower** at lower temperatures

DRAM Latency Characterization of 223 LPDDR4 DRAM Devices

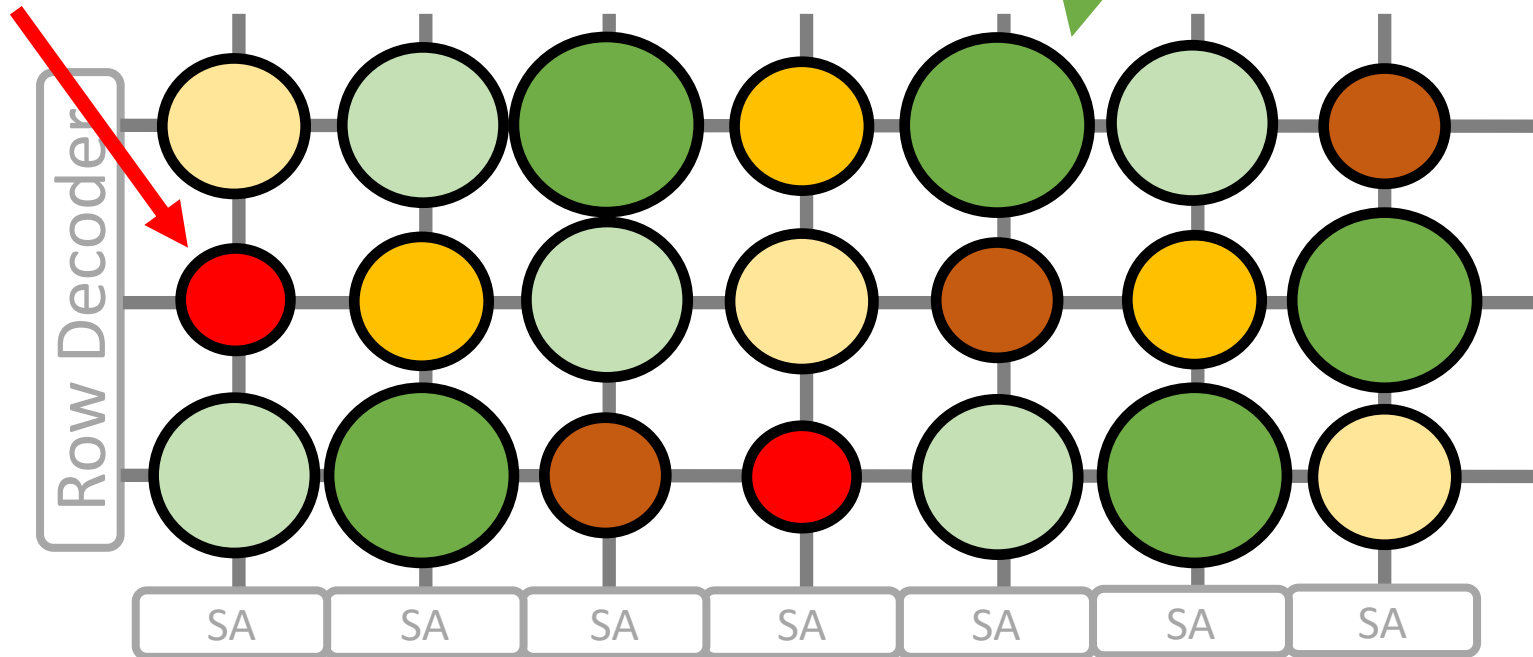
- Latency failures come from accessing DRAM with **reduced** timing parameters.
- **Key Observations:**
 1. A cell's **latency failure** probability is determined by **random process variation**
 2. Latency failure patterns are **repeatable and unique to a device**

DRAM Latency PUF Key Idea

- A cell's latency failure probability is inherently related to **random process variation** from manufacturing
- We can provide **repeatable and unique device signatures** using latency error patterns

High % chance to fail
with reduced t_{RCD}

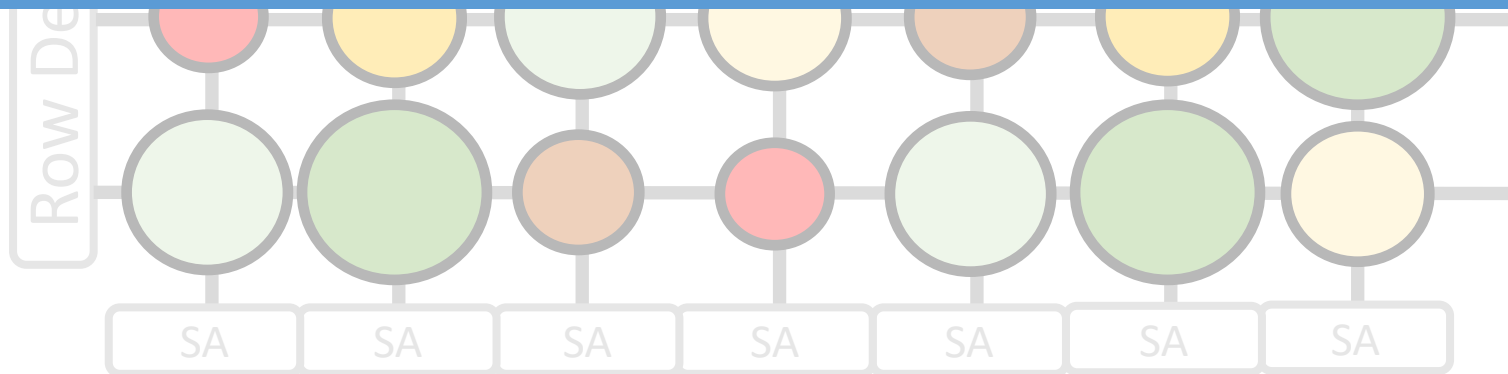
Low % chance to fail
with reduced t_{RCD}



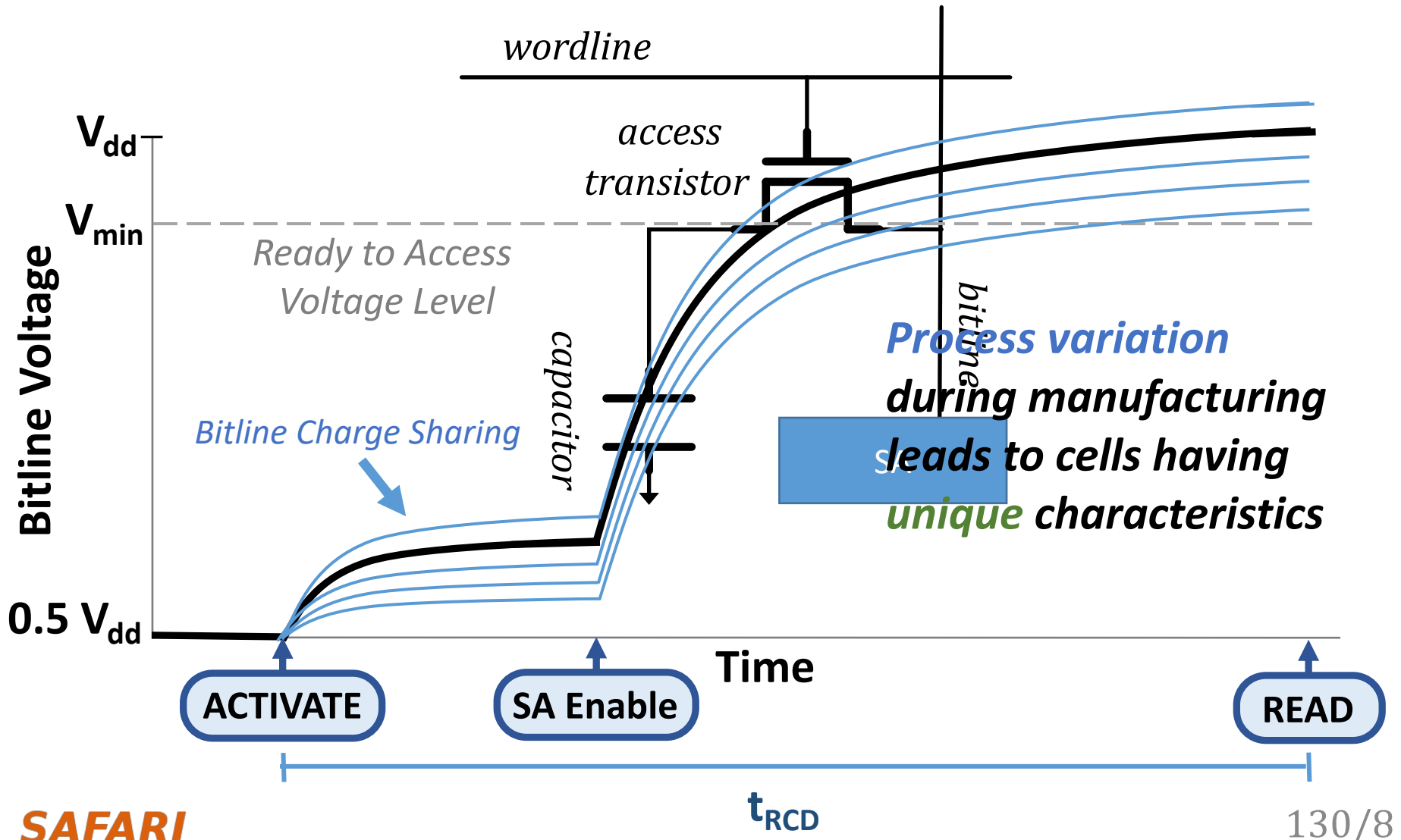
DRAM Latency PUF Key Idea

- A cell's latency failure probability is inherently related to **random process variation** from manufacturing
- We can provide **repeatable and unique device**

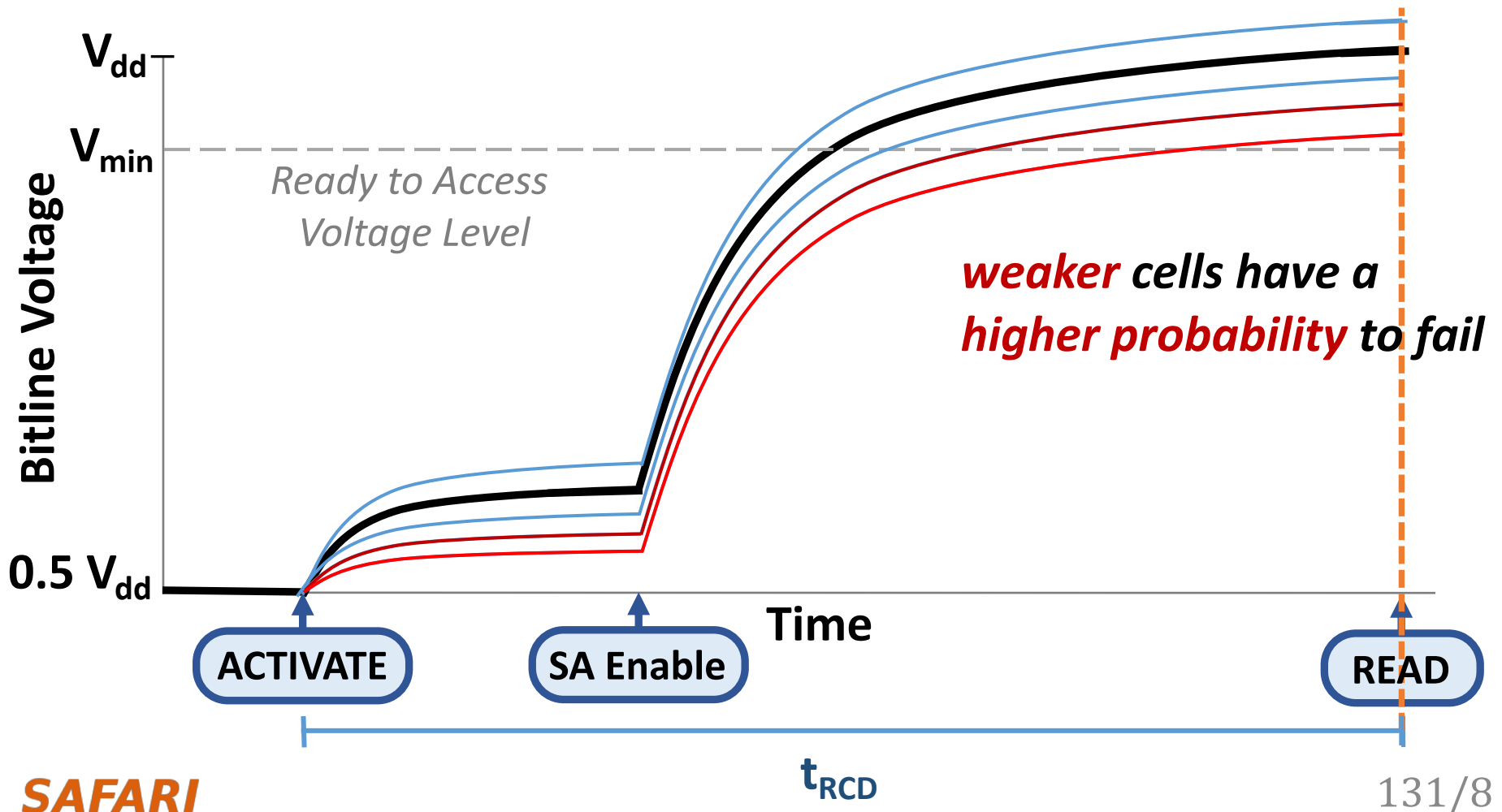
The **key idea** is to compose a PUF response using the DRAM cells that fail with **high probability**



DRAM Accesses and Failures



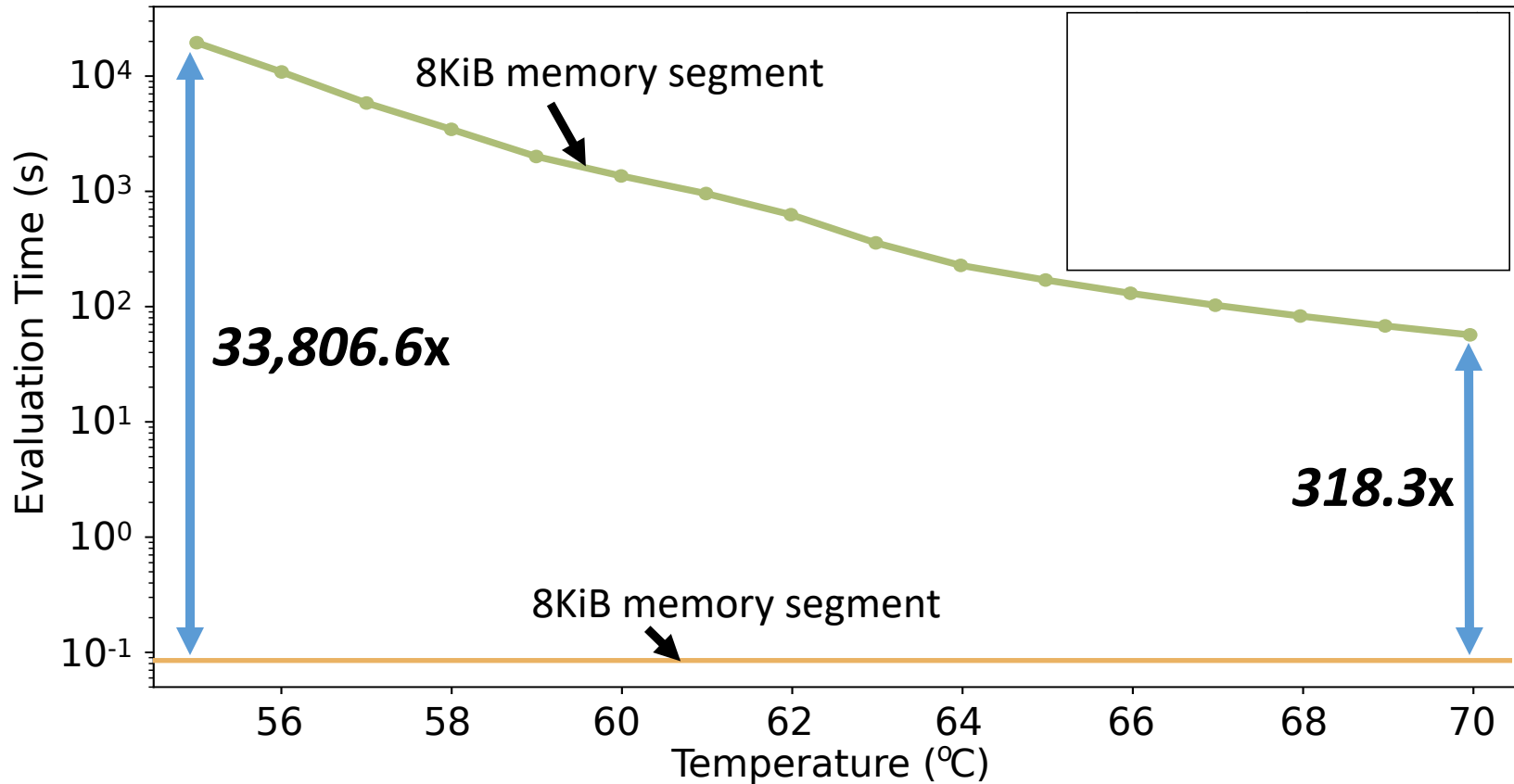
DRAM Accesses and Failures



The DRAM Latency PUF Evaluation

- We generate PUF responses using **latency errors** in a region of DRAM
- The latency error patterns **satisfy PUF requirements**
- The DRAM Latency PUF **generates PUF responses in 88.2ms**

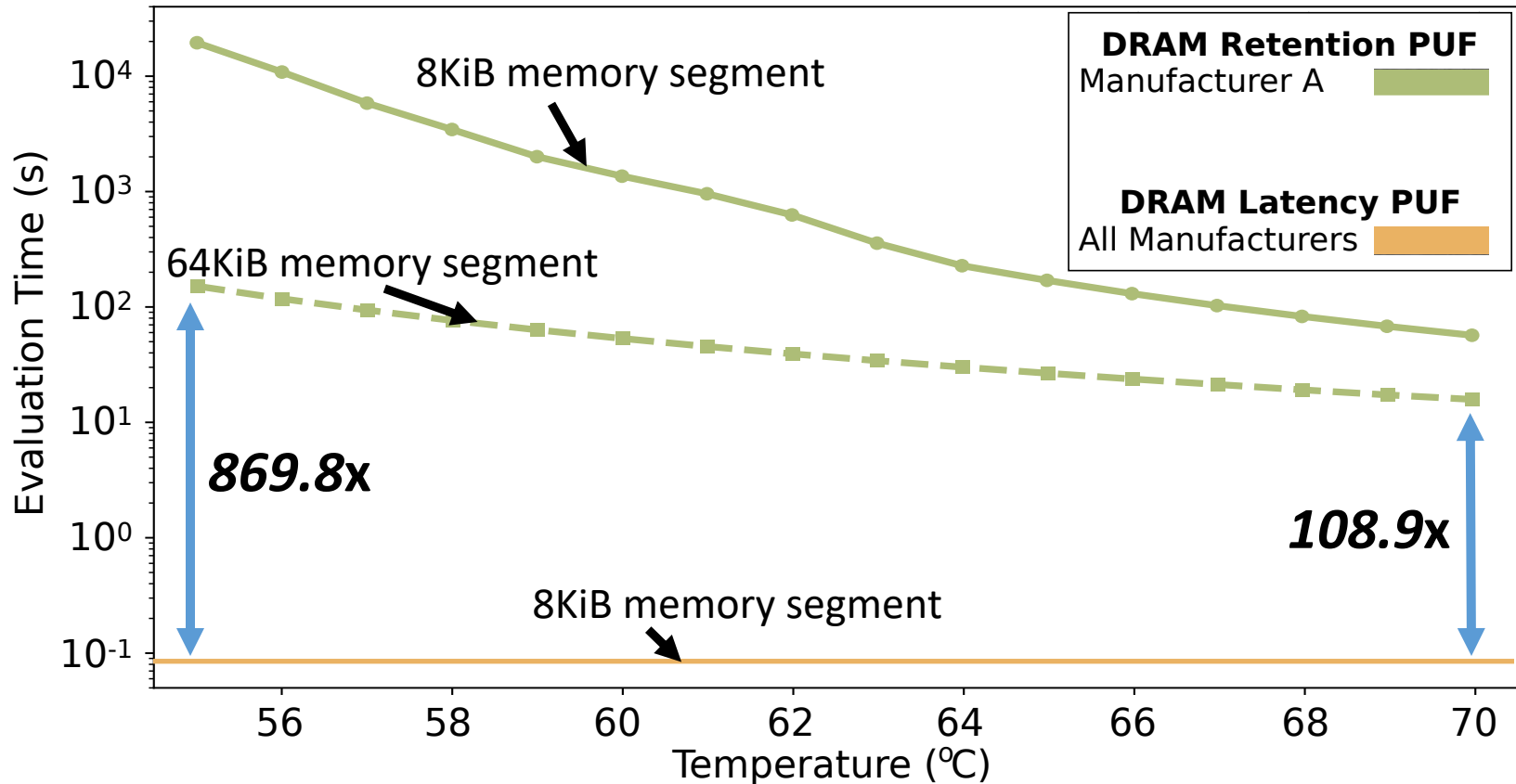
Results – PUF Evaluation Latency



DRAM latency PUF is

1. Fast and constant latency (**88.2ms**)

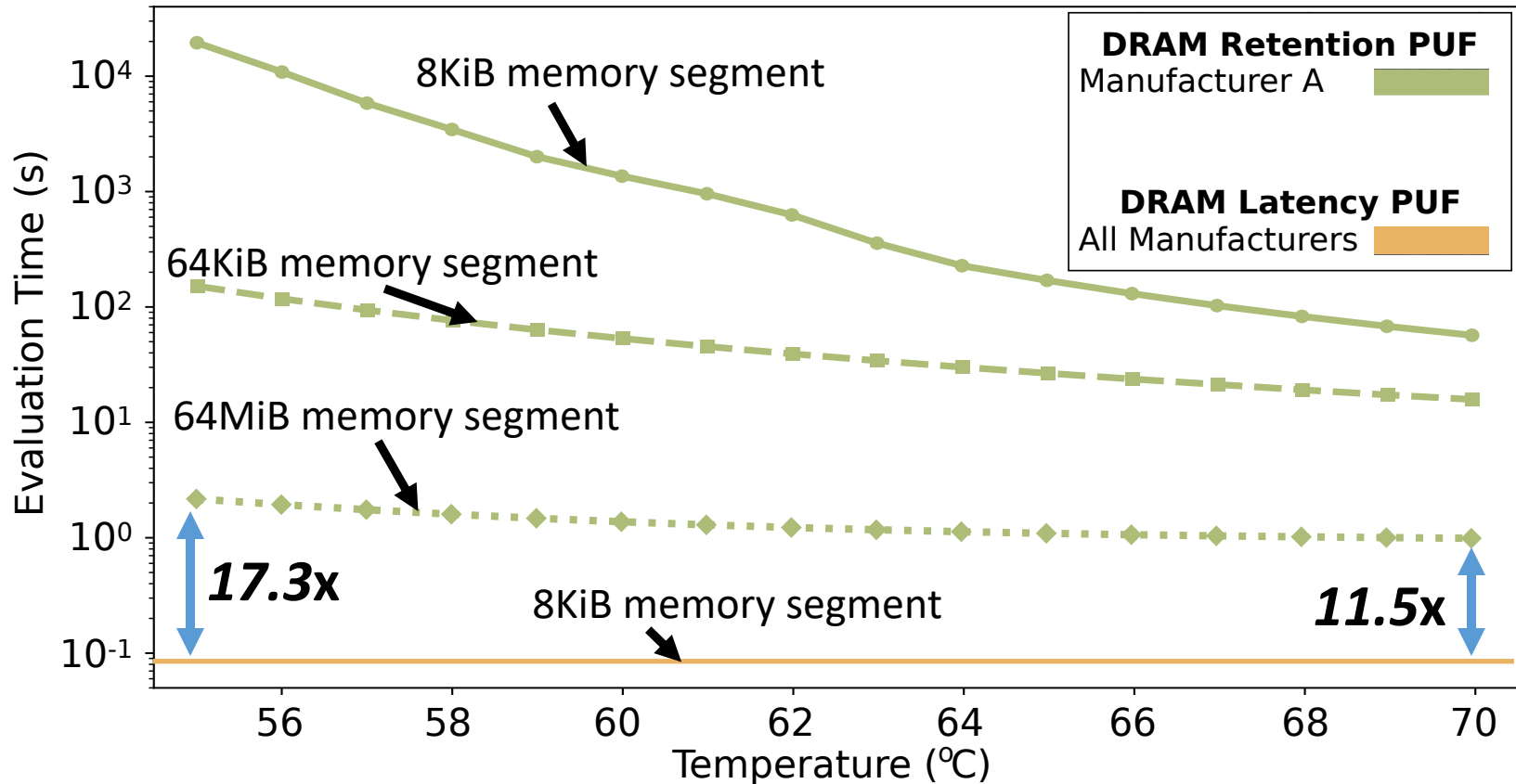
Results – PUF Evaluation Latency



DRAM latency PUF is

1. Fast and constant latency (88.2ms)

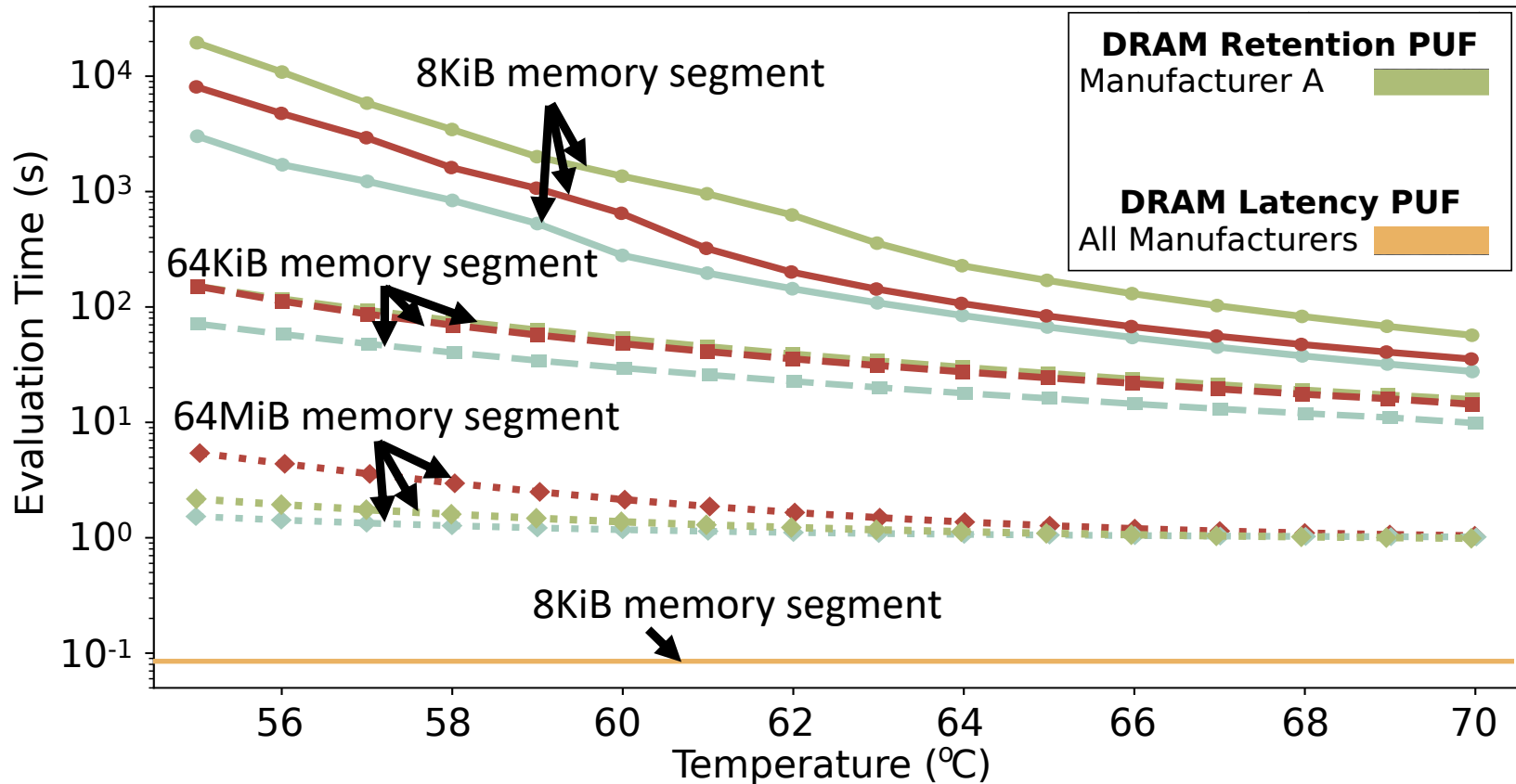
Results – PUF Evaluation Latency



DRAM latency PUF is

1. Fast and constant latency (88.2ms)

Results – PUF Evaluation Latency



DRAM latency PUF is

1. Fast and constant latency (**88.2ms**)
2. On average, **102x/860x** faster than the previous DRAM PUF with the same DRAM capacity overhead (64KiB)

Other Results in the Paper

- How the **DRAM latency PUF** meets the basic requirements for an effective PUF
- A **detailed** analysis on:
 - Devices of **the three major DRAM manufacturers**
 - The **evaluation time** of a PUF
- **Further discussion on:**
 - **Optimizing** retention PUFs
 - **System interference** of DRAM retention and latency PUFs
 - Algorithm to **quickly and reliably** evaluate DRAM latency PUF
 - **Design considerations** for a DRAM latency PUF
 - The DRAM Latency PUF overhead analysis

The DRAM Latency PUF:

Quickly Evaluating Physical Unclonable Functions
by Exploiting the Latency-Reliability Tradeoff
in Modern Commodity DRAM Devices

Jeremie S. Kim Minesh Patel

Hasan Hassan Onur Mutlu



QR Code for the paper

https://people.inf.ethz.ch/omutlu/pub/dram-latency-puf_hpca18.pdf

HPCA 2018

SAFARI



ETH zürich

Carnegie Mellon

DRAM Latency PUFs

- Jeremie S. Kim, Minesh Patel, Hasan Hassan, and Onur Mutlu,
"The DRAM Latency PUF: Quickly Evaluating Physical Unclonable Functions by Exploiting the Latency-Reliability Tradeoff in Modern DRAM Devices"
Proceedings of the 24th International Symposium on High-Performance Computer Architecture (HPCA), Vienna, Austria, February 2018.
[[Lightning Talk Video](#)]
[[Slides \(pptx\)](#)] [[pdf](#)] [[Lightning Session Slides \(pptx\)](#)] [[pdf](#)]

The DRAM Latency PUF:

Quickly Evaluating Physical Unclonable Functions

by Exploiting the Latency-Reliability Tradeoff in Modern Commodity DRAM Devices

Jeremie S. Kim^{†§}

Minesh Patel[§]

Hasan Hassan[§]

Onur Mutlu^{§†}

[†]Carnegie Mellon University

[§]ETH Zürich

Reducing Refresh Latency

On Reducing Refresh Latency

- Anup Das, Hasan Hassan, and Onur Mutlu,
**"VRL-DRAM: Improving DRAM Performance via
Variable Refresh Latency"**
*Proceedings of the 55th Design Automation
Conference (DAC)*, San Francisco, CA, USA, June 2018.

VRL-DRAM: Improving DRAM Performance via Variable Refresh Latency

Anup Das
Drexel University
Philadelphia, PA, USA
anup.das@drexel.edu

Hasan Hassan
ETH Zürich
Zürich, Switzerland
hhasan@ethz.ch

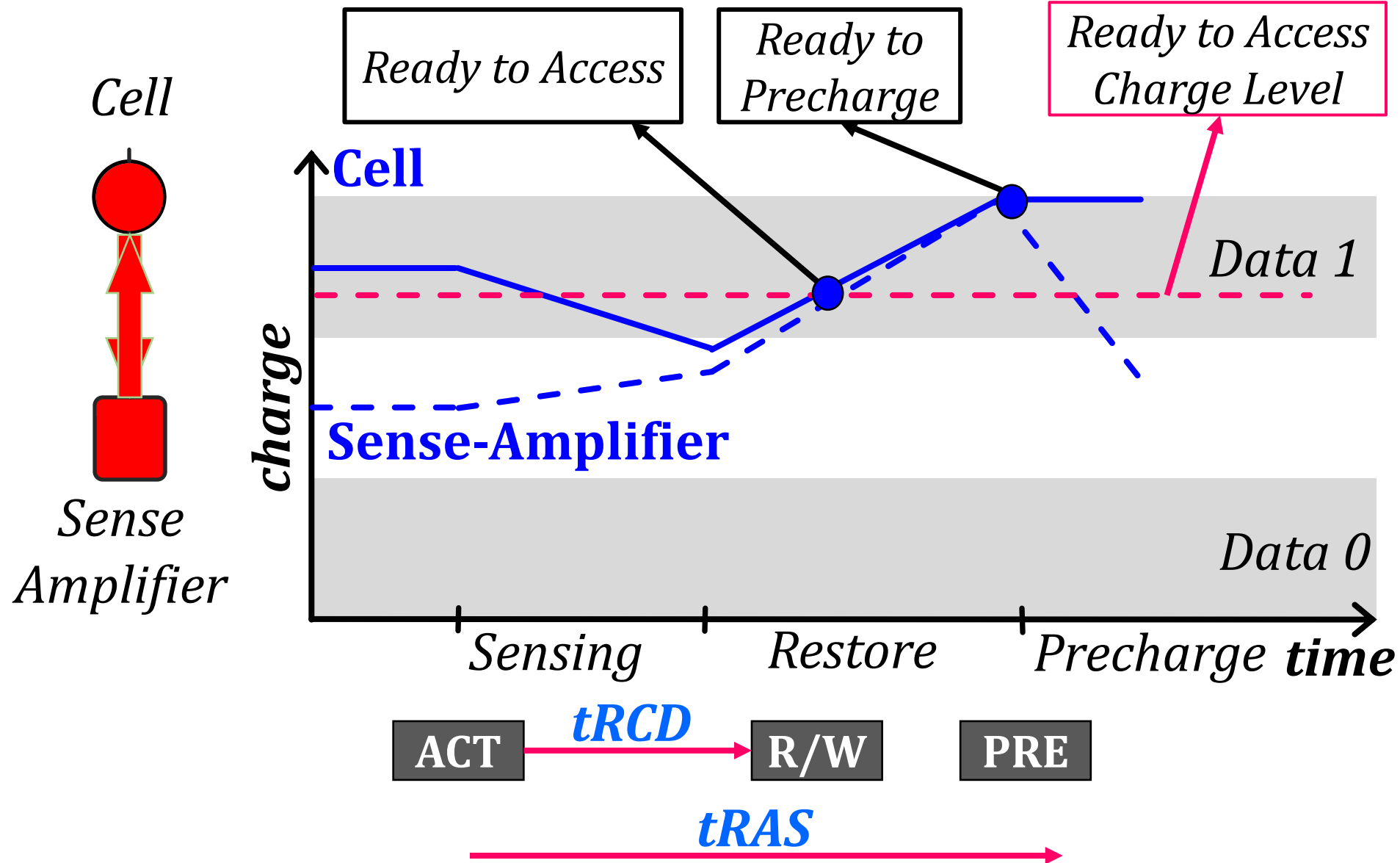
Onur Mutlu
ETH Zürich
Zürich, Switzerland
omutlu@gmail.com

Reducing Memory Latency by Exploiting Memory Access Patterns

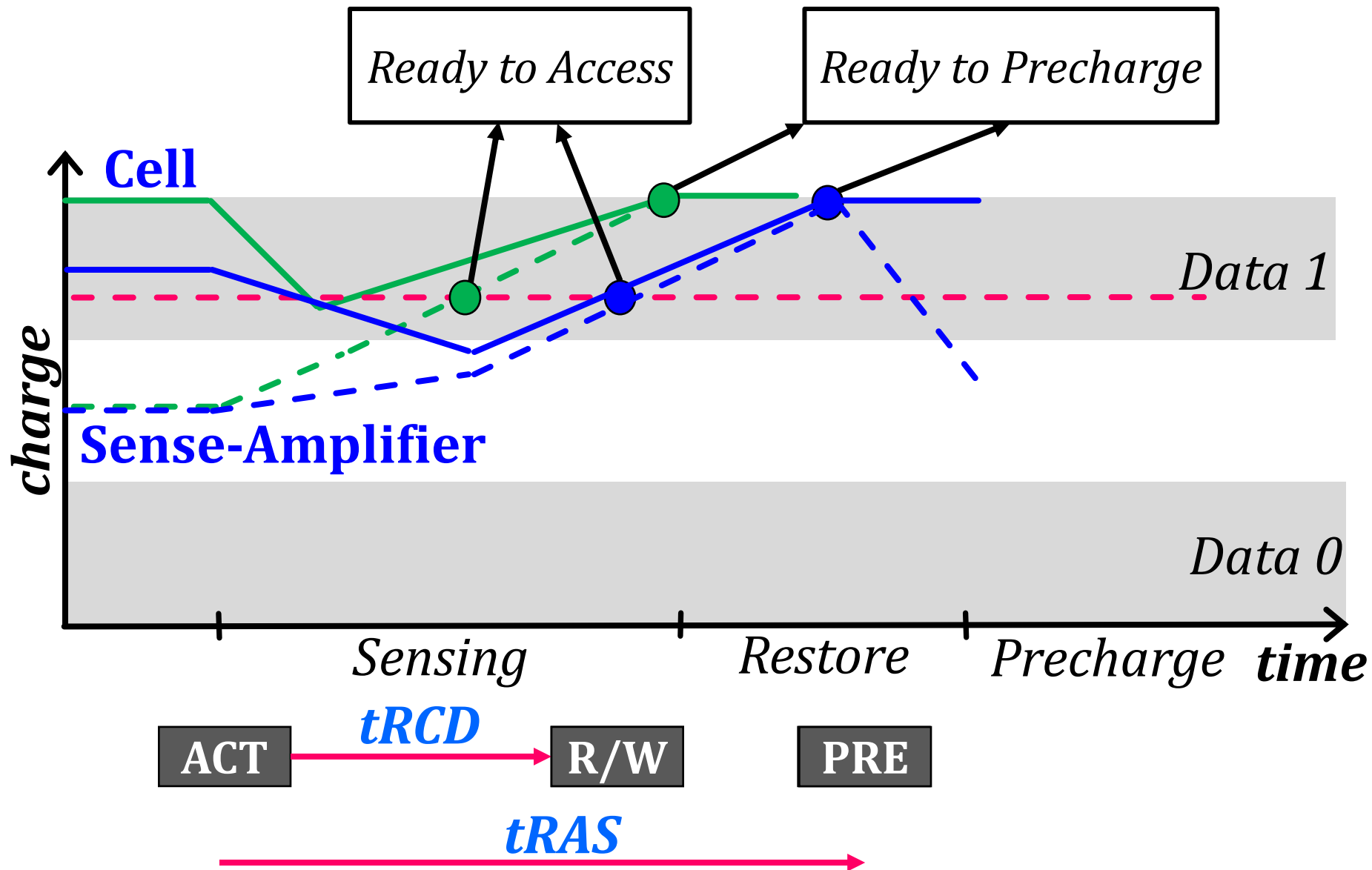
ChargeCache: Executive Summary

- **Goal**: Reduce average DRAM access latency with no modification to the existing DRAM chips
 - **Observations**:
 - 1) A highly-charged DRAM row can be accessed with low latency
 - 2) A row's charge is restored when the row is accessed
 - 3) A recently-accessed row is likely to be accessed again:
- Row Level Temporal Locality (RLTL)**
- **Key Idea**: Track recently-accessed DRAM rows and use lower timing parameters if such rows are accessed again
 - **ChargeCache**:
 - Low cost & no modifications to the DRAM
 - Higher performance (**8.6-10.6%** on average for 8-core)
 - Lower DRAM energy (**7.9%** on average)

DRAM Charge over Time



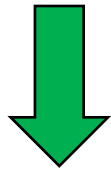
Accessing Highly-charged Rows



Observation 1

A **highly-charged** DRAM row can be accessed with **low latency**

- tRCD: 44%
- tRAS: 37%



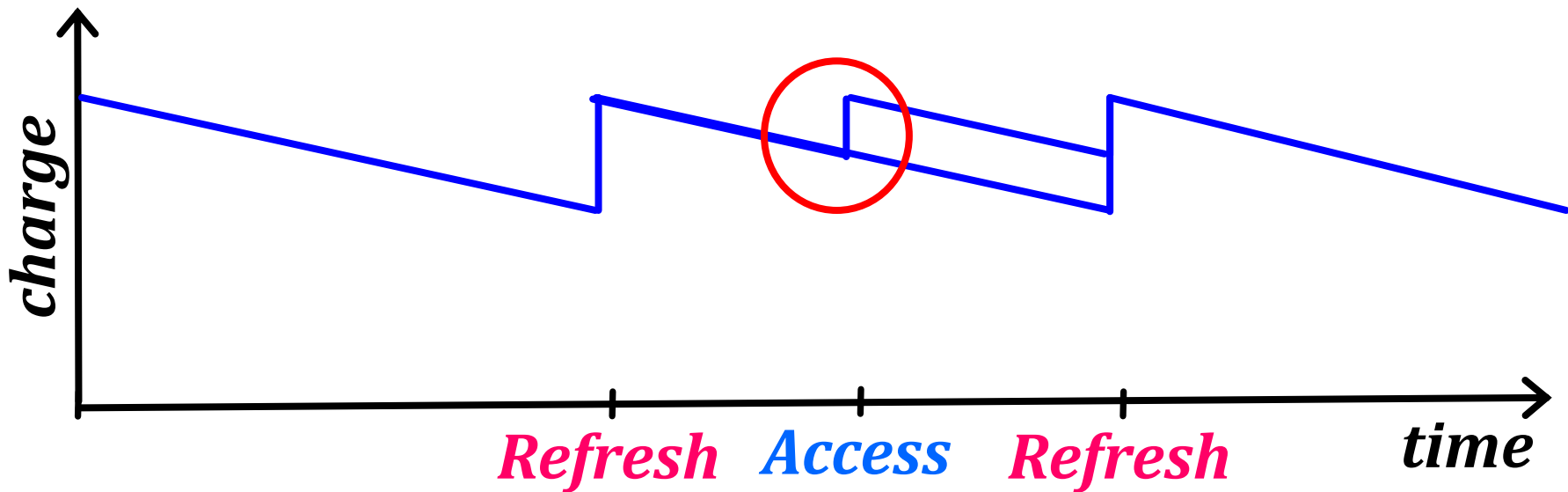
**How does a row become
highly-charged?**

How Does a Row Become Highly-Charged?

DRAM cells **lose charge** over time

Two ways of restoring a row's charge:

- Refresh Operation
- Access



Observation 2

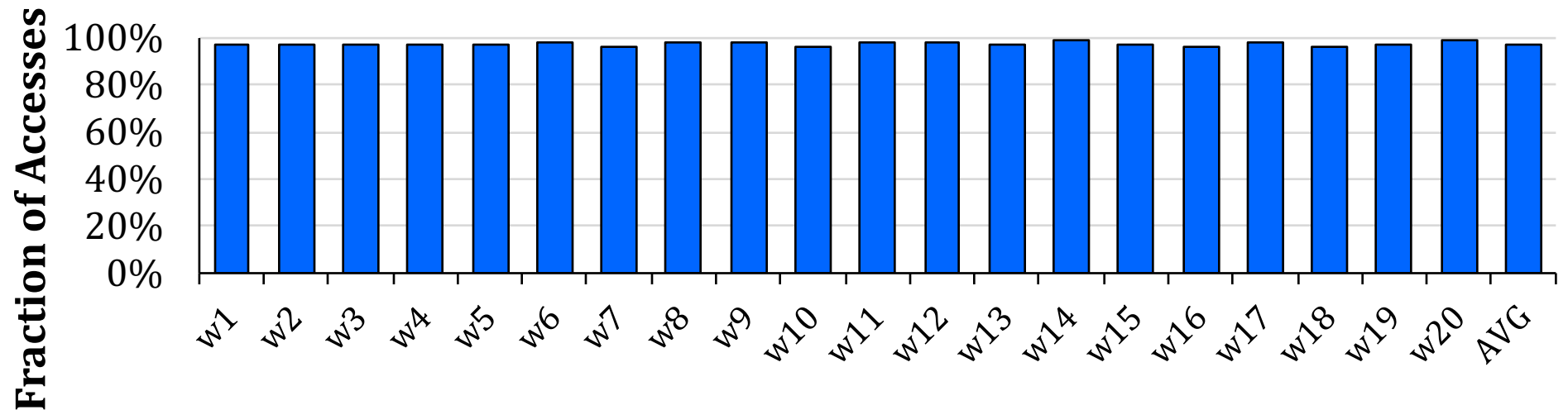
A row's charge is **restored** when the row is **accessed**

How likely is a **recently-accessed row to be accessed again?**

Row Level Temporal Locality (RLTL)

A **recently-accessed** DRAM row is likely to be accessed again.

- t -RLTL: Fraction of rows that are accessed within time t after their previous access

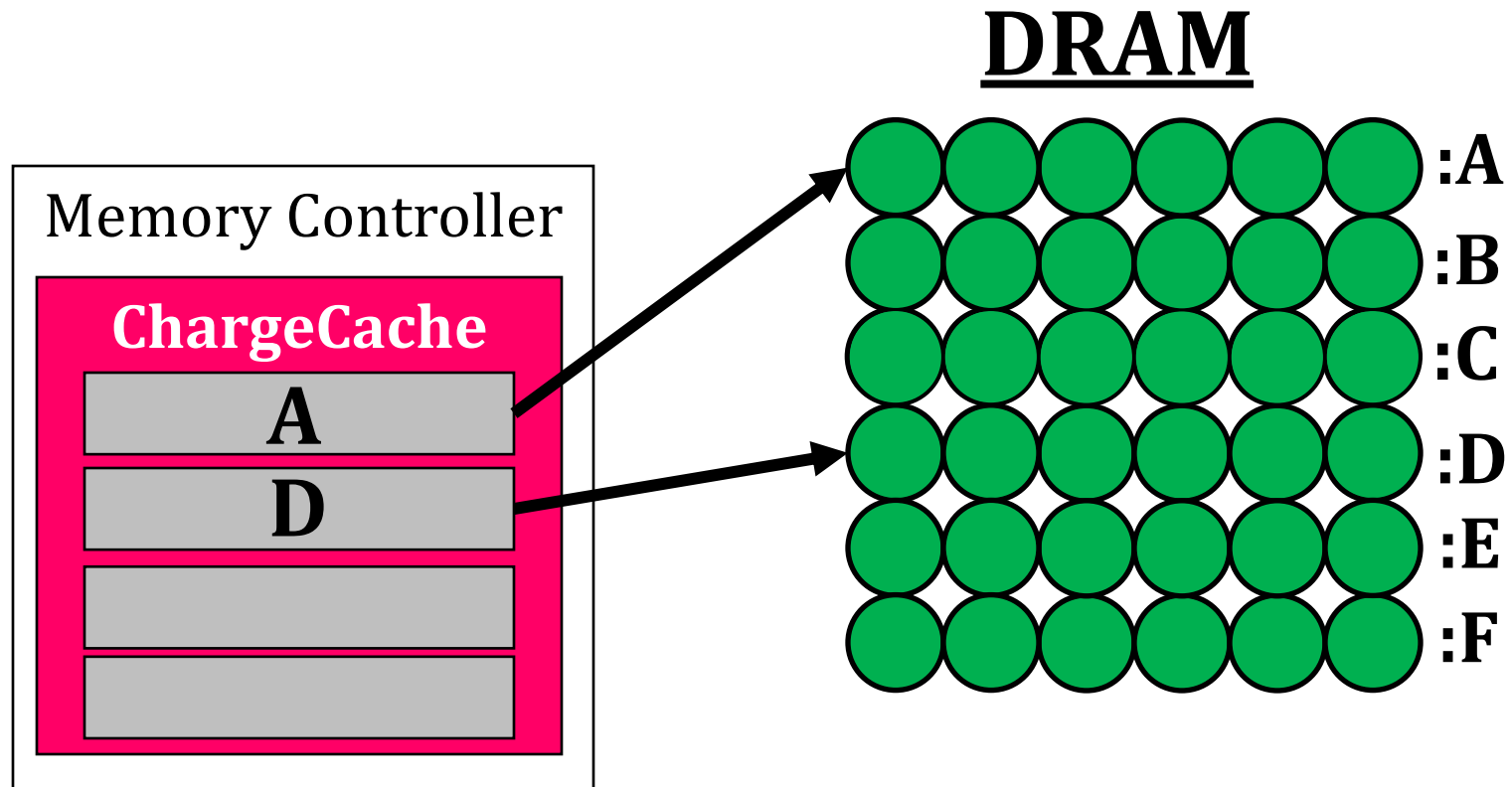


88ns — RLTL for eight-core workloads

Key Idea

Track **recently-accessed** DRAM rows and use **lower timing parameters** if such rows are accessed again

ChargeCache Overview



Requests: A D A 

ChargeCache Hits: Use Default Timings

Area and Power Overhead

- Modeled with CACTI

- Area

- ~5KB for 128-entry ChargeCache
- 0.24% of a 4MB Last Level Cache (LLC) area

- Power Consumption

- 0.15 mW on average (static + dynamic)
- 0.23% of the 4MB LLC power consumption

Methodology

- **Simulator**

- DRAM Simulator (Ramulator *[Kim+, CAL'15]*)
<https://github.com/CMU-SAFARI/ramulator>

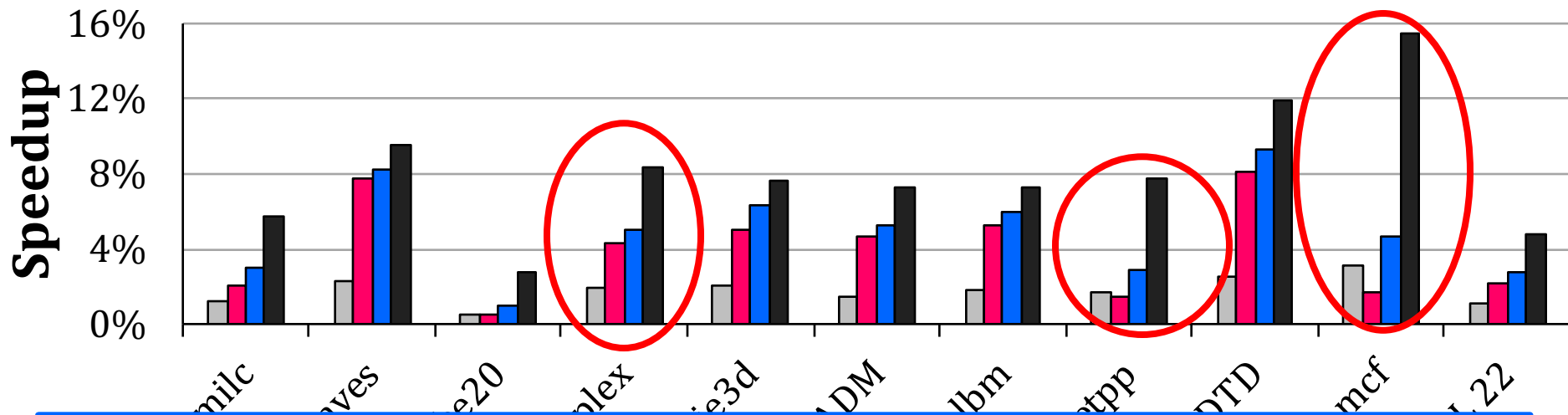
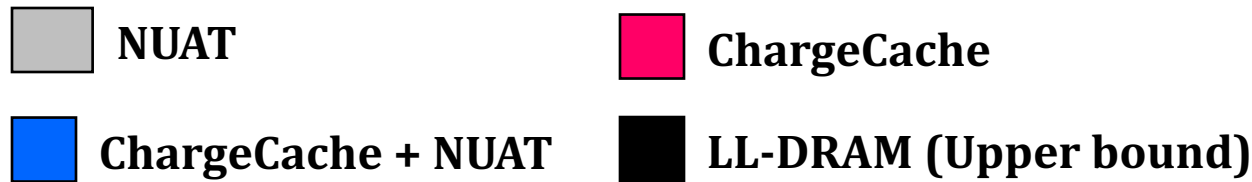
- **Workloads**

- 22 single-core workloads
 - SPEC CPU2006, TPC, STREAM
- 20 multi-programmed 8-core workloads
 - By randomly choosing from single-core workloads
- Execute at least 1 billion representative instructions per core (Pinpoints)

- **System Parameters**

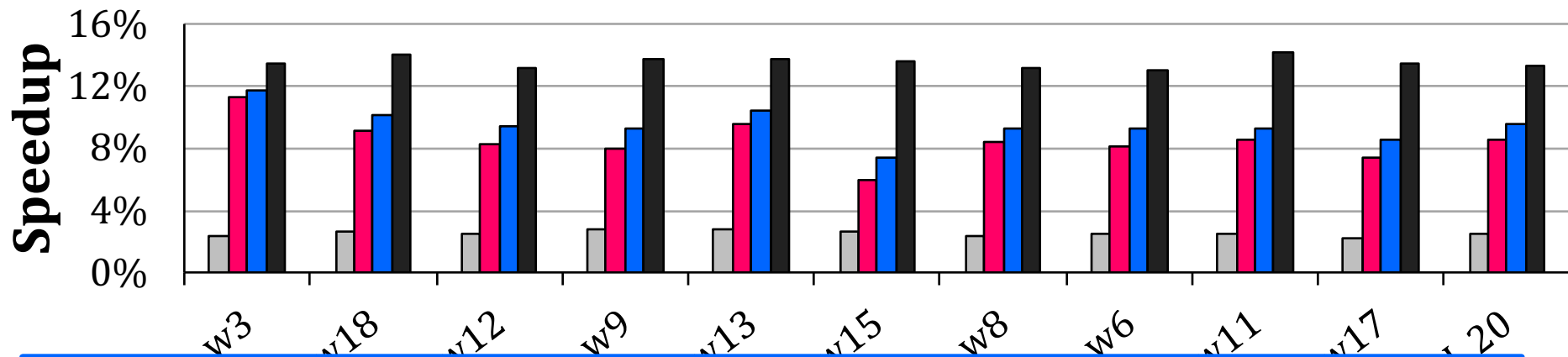
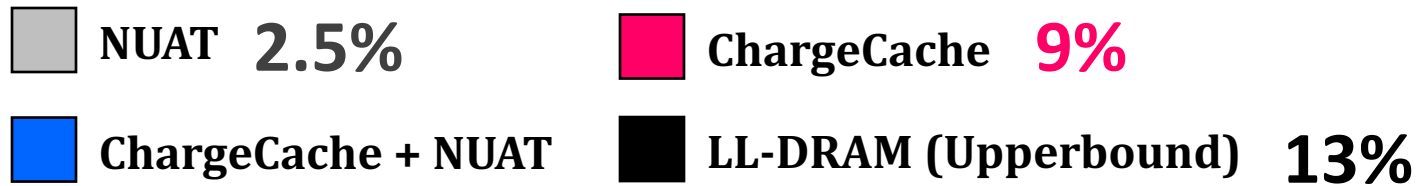
- 1/8 core system with 4MB LLC
- Default tRCD/tRAS of 11/28 cycles

Single-core Performance



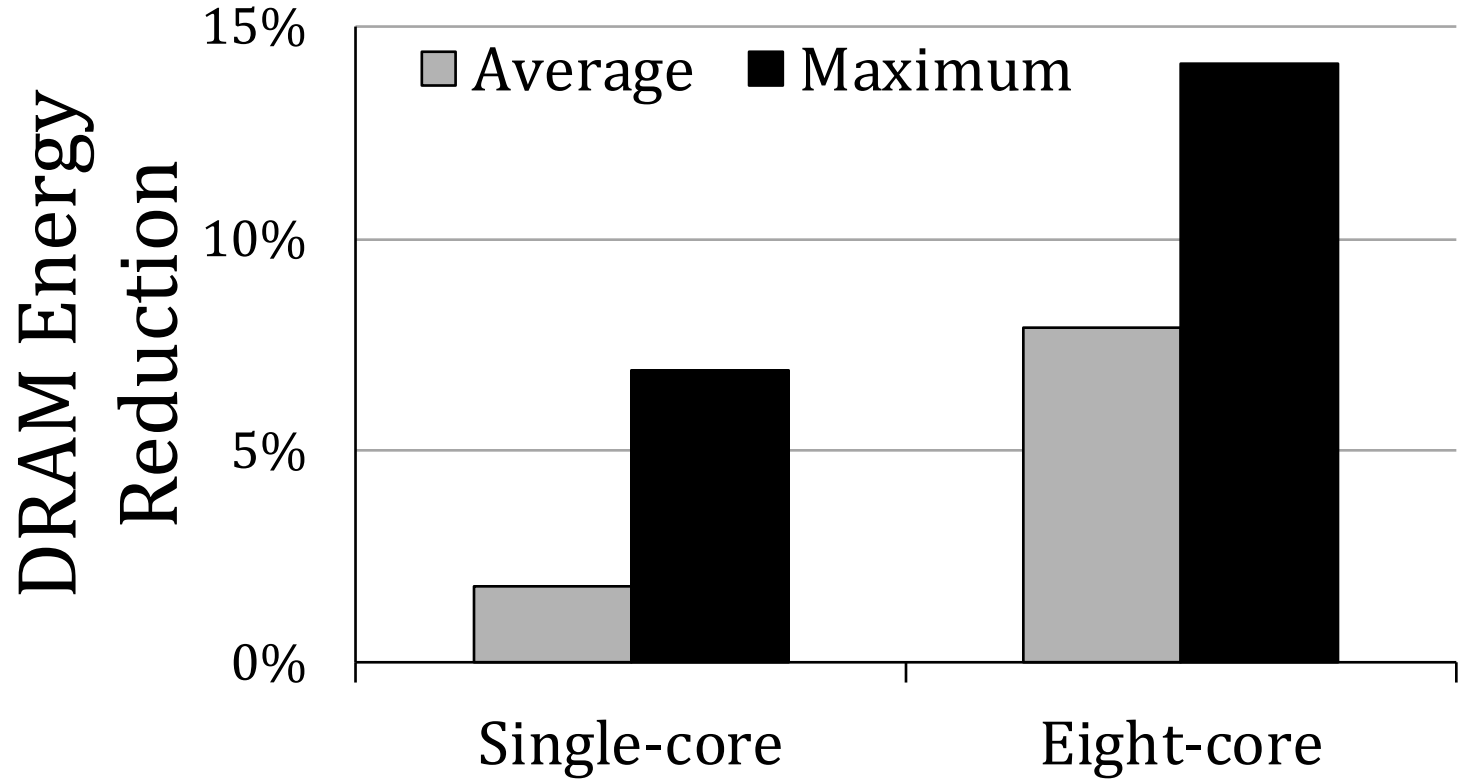
**ChargeCache improves
single-core performance**

Eight-core Performance



ChargeCache significantly improves multi-core performance

DRAM Energy Savings



ChargeCache reduces DRAM energy

More on ChargeCache

- Hasan Hassan, Gennady Pekhimenko, Nandita Vijaykumar, Vivek Seshadri, Donghyuk Lee, Oguz Ergin, and Onur Mutlu,
"ChargeCache: Reducing DRAM Latency by Exploiting Row Access Locality"
Proceedings of the 22nd International Symposium on High-Performance Computer Architecture (HPCA), Barcelona, Spain, March 2016.
[[Slides \(pptx\)](#)] [[pdf](#)]
[[Source Code](#)]

ChargeCache: Reducing DRAM Latency by Exploiting Row Access Locality

Hasan Hassan^{†*}, Gennady Pekhimenko[†], Nandita Vijaykumar[†]
Vivek Seshadri[†], Donghyuk Lee[†], Oguz Ergin^{*}, Onur Mutlu[†]

A Very Recent Work

- Yaohua Wang, Arash Tavakkol, Lois Orosa, Saugata Ghose, Nika Mansouri Ghiasi, Minesh Patel, Jeremie S. Kim, Hasan Hassan, Mohammad Sadrosadati, and Onur Mutlu,

"Reducing DRAM Latency via Charge-Level-Aware Look-Ahead Partial Restoration"

Proceedings of the 51st International Symposium on Microarchitecture (MICRO), Fukuoka, Japan, October 2018.

Reducing DRAM Latency via Charge-Level-Aware Look-Ahead Partial Restoration

Yaohua Wang^{†§} Arash Tavakkol[†] Lois Orosa^{†*} Saugata Ghose[‡] Nika Mansouri Ghiasi[†]
Minesh Patel[†] Jeremie S. Kim^{‡†} Hasan Hassan[†] Mohammad Sadrosadati[†] Onur Mutlu^{‡†}

[†]*ETH Zürich*

[§]*National University of Defense Technology*

[‡]*Carnegie Mellon University*

^{*}*University of Campinas*

On DRAM Power Consumption

Summary: Low-Latency Memory

Summary: Tackling Long Memory Latency

- Reason 1: Design of DRAM Micro-architecture
 - Goal: Maximize capacity/area, not minimize latency
- Reason 2: “One size fits all” approach to latency specification
 - Same latency parameters for all temperatures
 - Same latency parameters for all DRAM chips (e.g., rows)
 - Same latency parameters for all parts of a DRAM chip
 - Same latency parameters for all supply voltage levels
 - Same latency parameters for all application data
 - ...

Fundamentally Low Latency Computing Architectures

On DRAM Power Consumption

VAMPIRE DRAM Power Model

- Saugata Ghose, A. Giray Yaglikci, Raghav Gupta, Donghyuk Lee, Kais Kudrolli, William X. Liu, Hasan Hassan, Kevin K. Chang, Niladrish Chatterjee, Aditya Agrawal, Mike O'Connor, and Onur Mutlu,
"What Your DRAM Power Models Are Not Telling You: Lessons from a Detailed Experimental Study"
*Proceedings of the ACM International Conference on Measurement and Modeling of Computer Systems (**SIGMETRICS**), Irvine, CA, USA, June 2018.*
[[Abstract](#)]

What Your DRAM Power Models Are Not Telling You: Lessons from a Detailed Experimental Study

Saugata Ghose [†]	Abdullah Giray Yağlıkçı ^{‡†}	Raghav Gupta [†]	Donghyuk Lee [§]
Kais Kudrolli [†]	William X. Liu [†]	Hasan Hassan [‡]	Kevin K. Chang [†]
Niladrish Chatterjee [§]	Aditya Agrawal [§]	Mike O'Connor ^{§¶}	Onur Mutlu ^{‡†}

[†]Carnegie Mellon University

[‡]ETH Zürich

[§]NVIDIA

[¶]University of Texas at Austin

Conclusion

Four Key Directions

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

Some Solution Principles (So Far)

- Data-centric system design & intelligence spread around
 - Do not center everything around traditional computation units
- Better cooperation across layers of the system
 - Careful co-design of components and layers: system/arch/device
 - Better, richer, more expressive and flexible interfaces
- Better-than-worst-case design
 - Do not optimize for the worst case
 - Worst case should not determine the common case
- Heterogeneity in design (specialization, asymmetry)
 - Enables a more efficient design (No one size fits all)

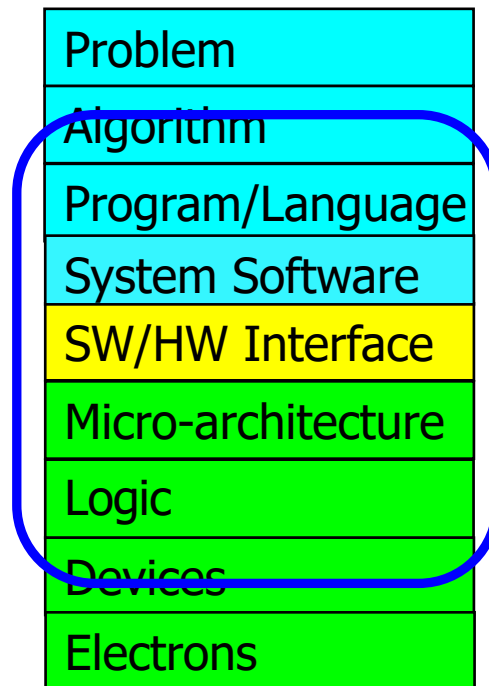
Some Solution Principles (More Compact)

- Data-centric design
- All components intelligent
- Better cross-layer communication, better interfaces
- Better-than-worst-case design
- Heterogeneity
- Flexibility, adaptability

It Is Time to ...

- ... design **principled system architectures** to solve the **memory problem**
- ... design complete systems to be balanced, high-performance, and energy-efficient, i.e., data-centric (or memory-centric)
- ... **make memory a key priority** in system design and optimize it & integrate it better into the system
- This can
 - Lead to **orders-of-magnitude** improvements
 - **Enable new applications & computing platforms**
 - **Enable better understanding of nature**
 - ...

We Need to Revisit the Entire Stack



Memory Systems

Fundamentals, Recent Research, Challenges, Opportunities

Lecture 4: Low-Latency Memory

Prof. Onur Mutlu

omutlu@gmail.com

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

9 October 2018

Technion Fast Course 2018