

# Computer Architecture

## Lecture 10b: Memory Latency

Prof. Onur Mutlu

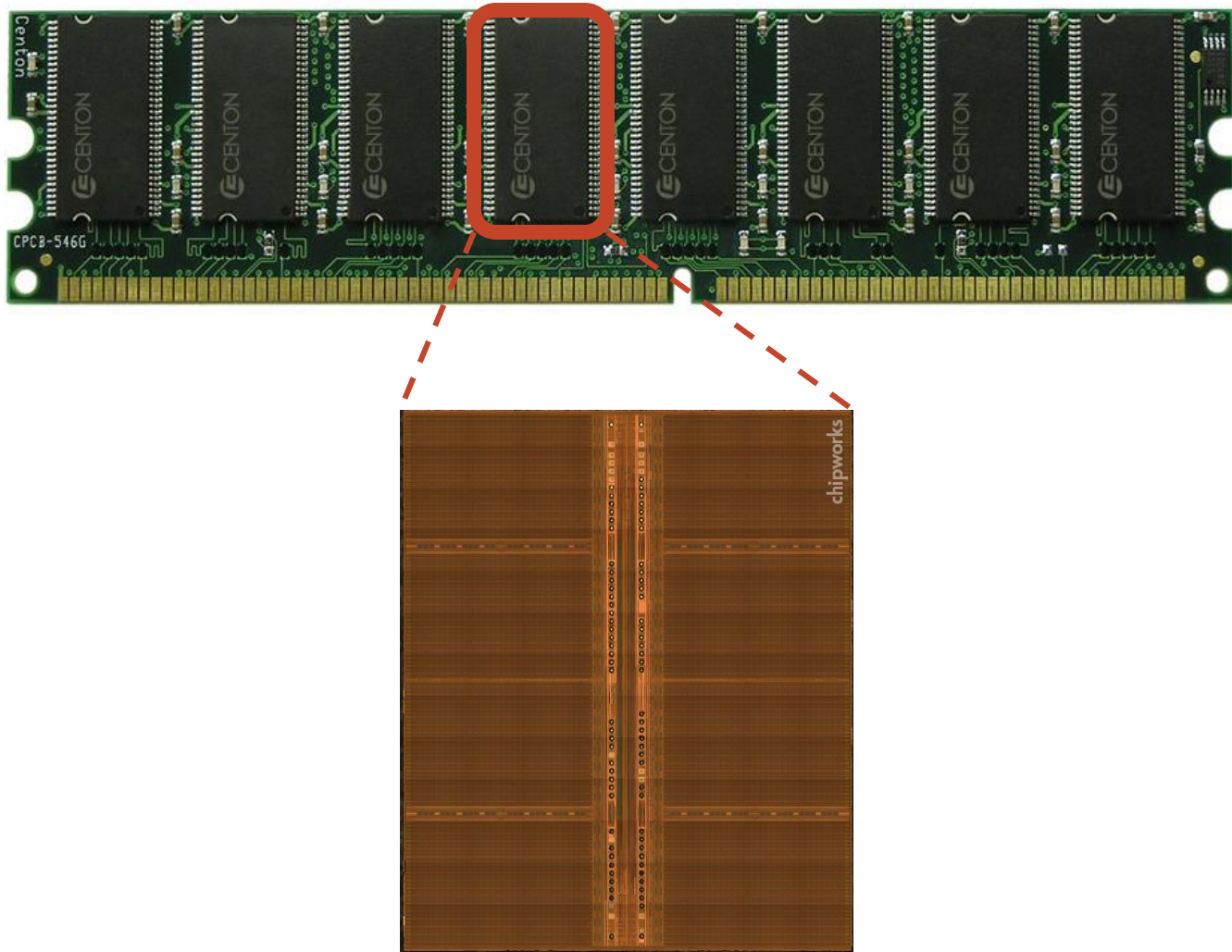
ETH Zürich

Fall 2018

18 October 2018

# DRAM Memory: A Low-Level Perspective

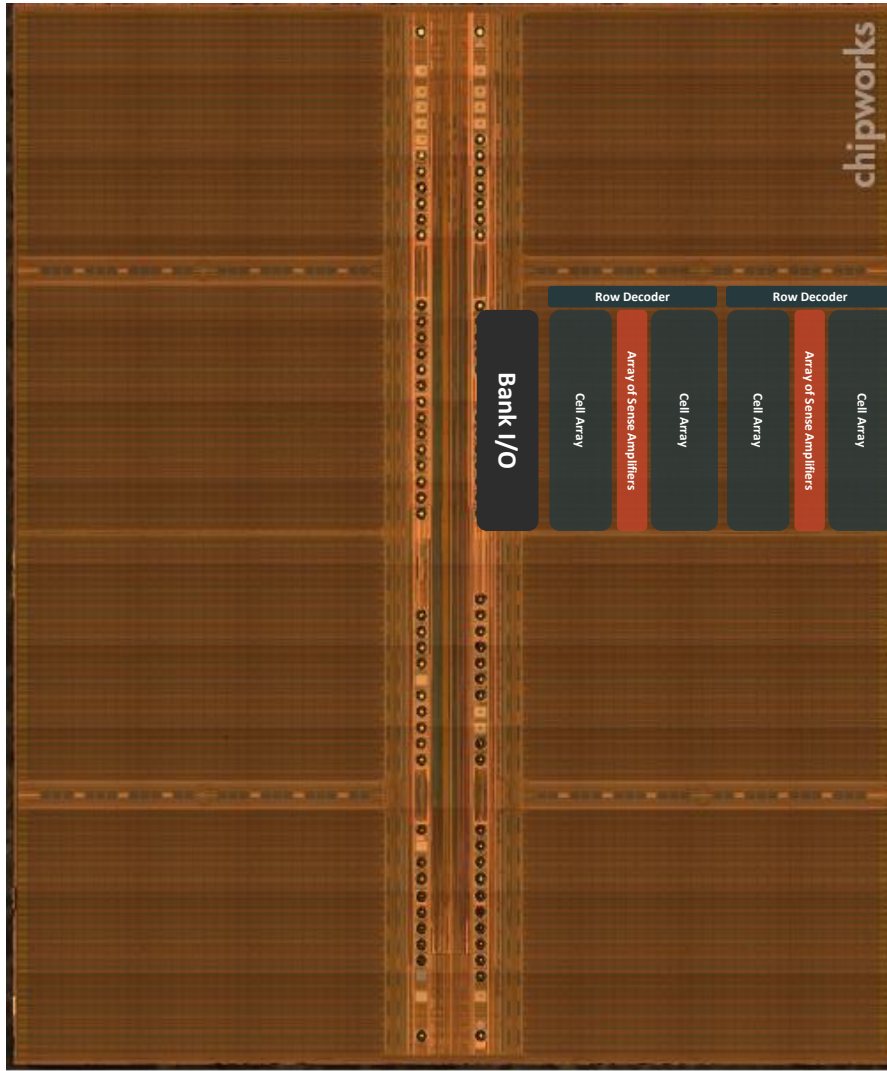
# DRAM Module and Chip



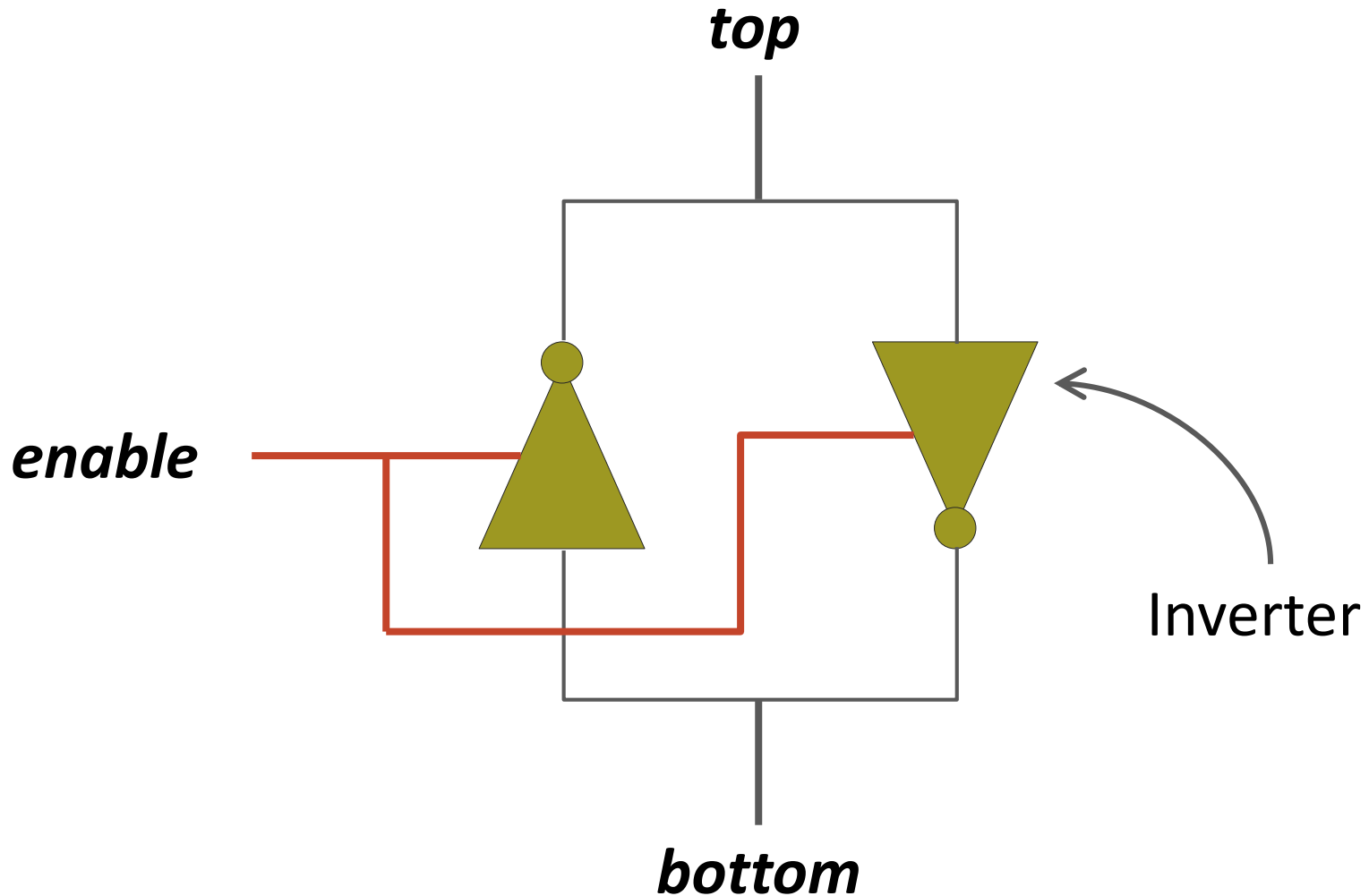
# Goals

- Cost
- Latency
- Bandwidth
- Parallelism
- Power
- Energy
- Reliability
- ...

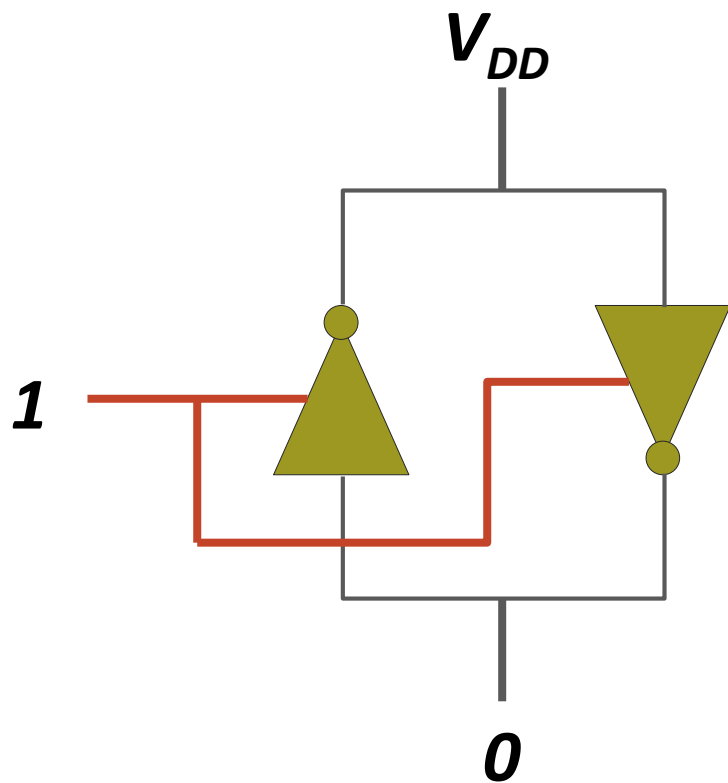
# DRAM Chip



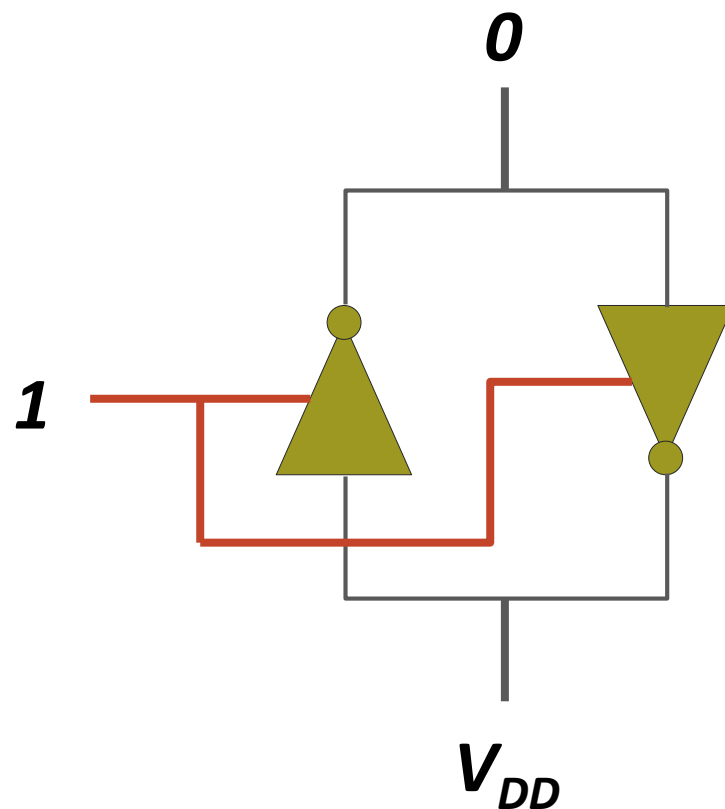
# Sense Amplifier



# Sense Amplifier – Two Stable States

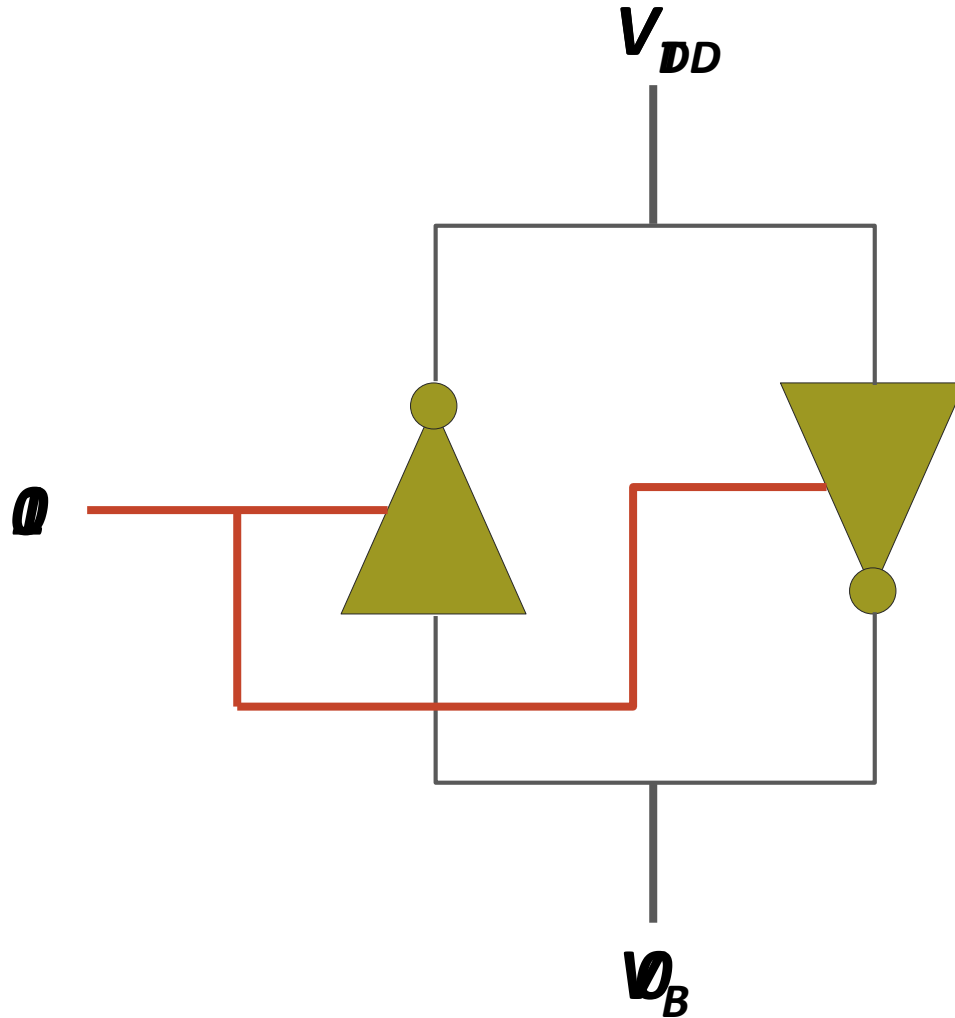


Logical "1"



Logical "0"

# Sense Amplifier Operation



$$V_T > V_B$$



# DRAM Cell – Capacitor



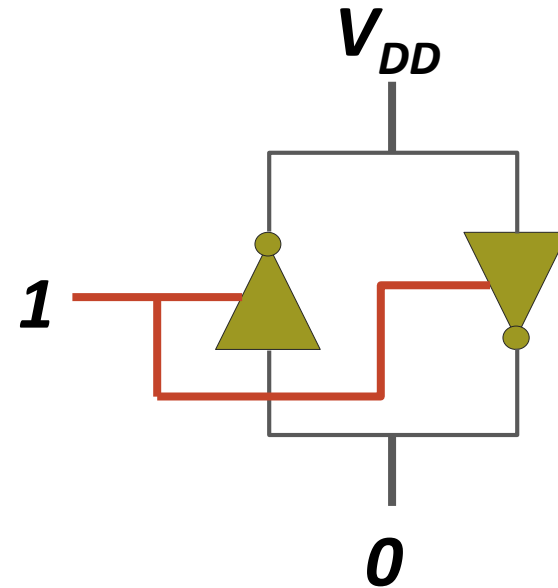
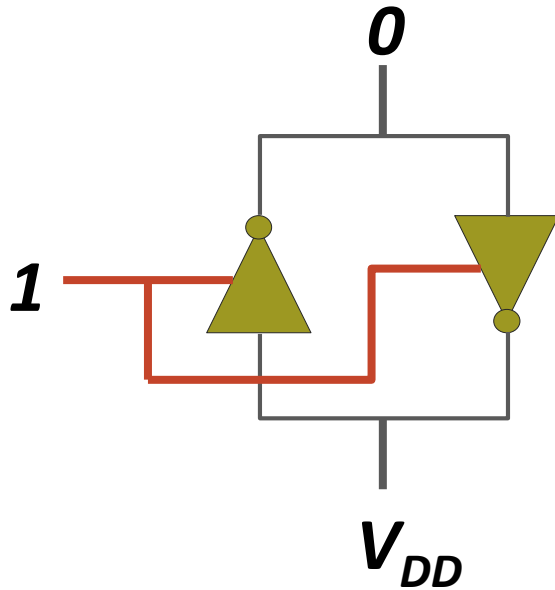
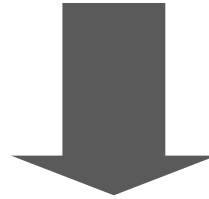
Empty State  
**Logical “0”**



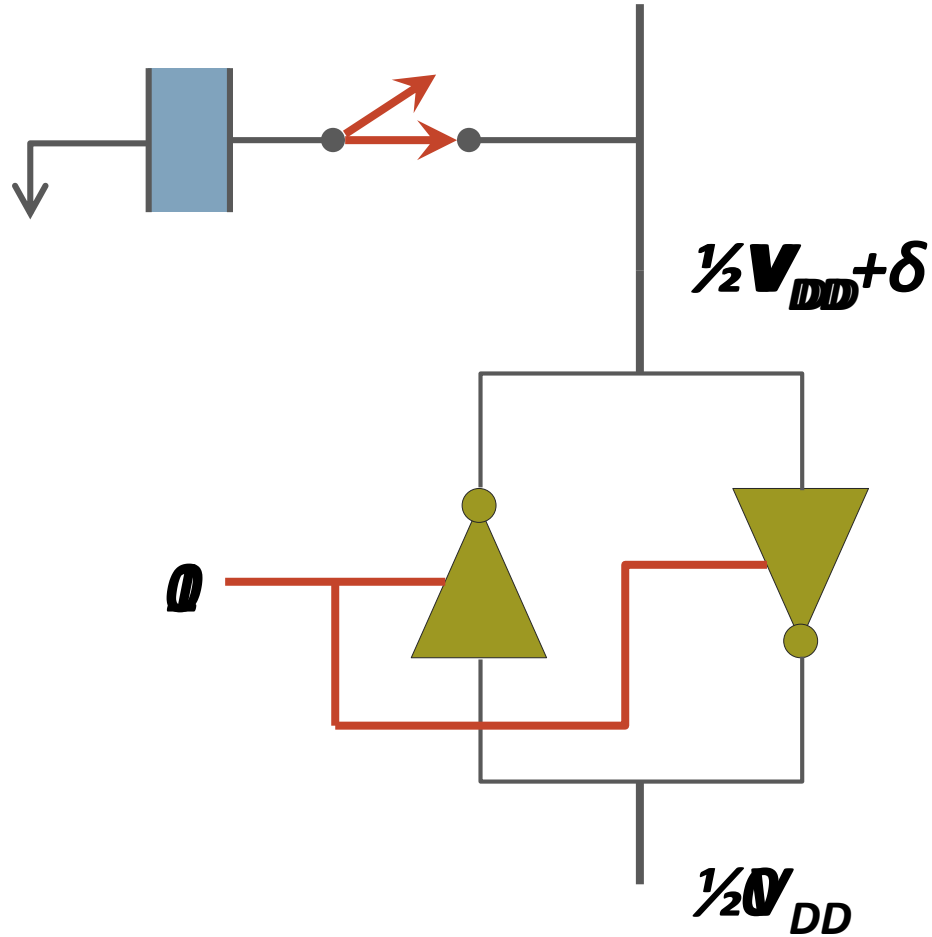
Fully Charged State  
**Logical “1”**

- 1 Small – Cannot drive circuits
- 2 Reading destroys the state

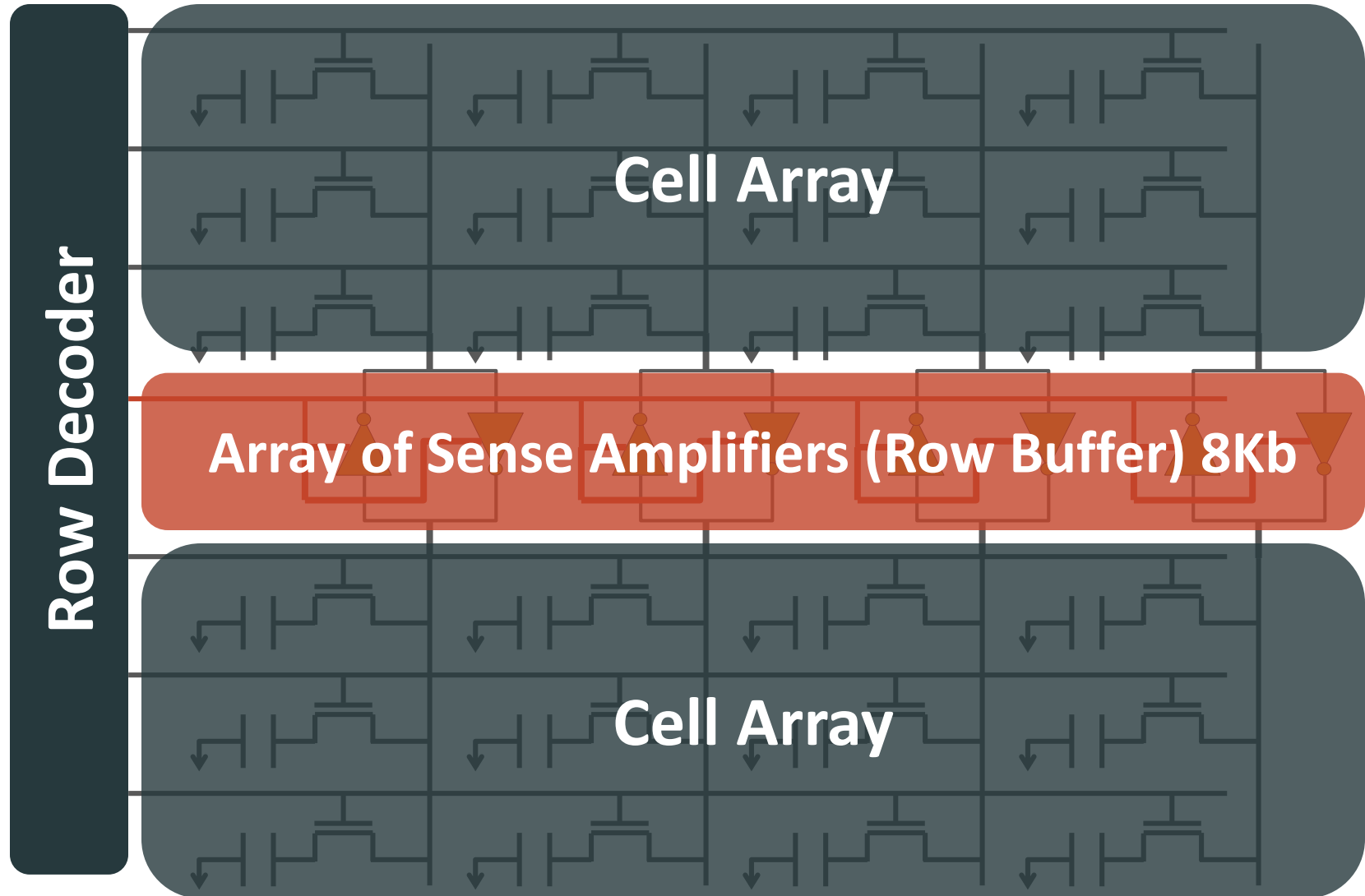
# Capacitor to Sense Amplifier



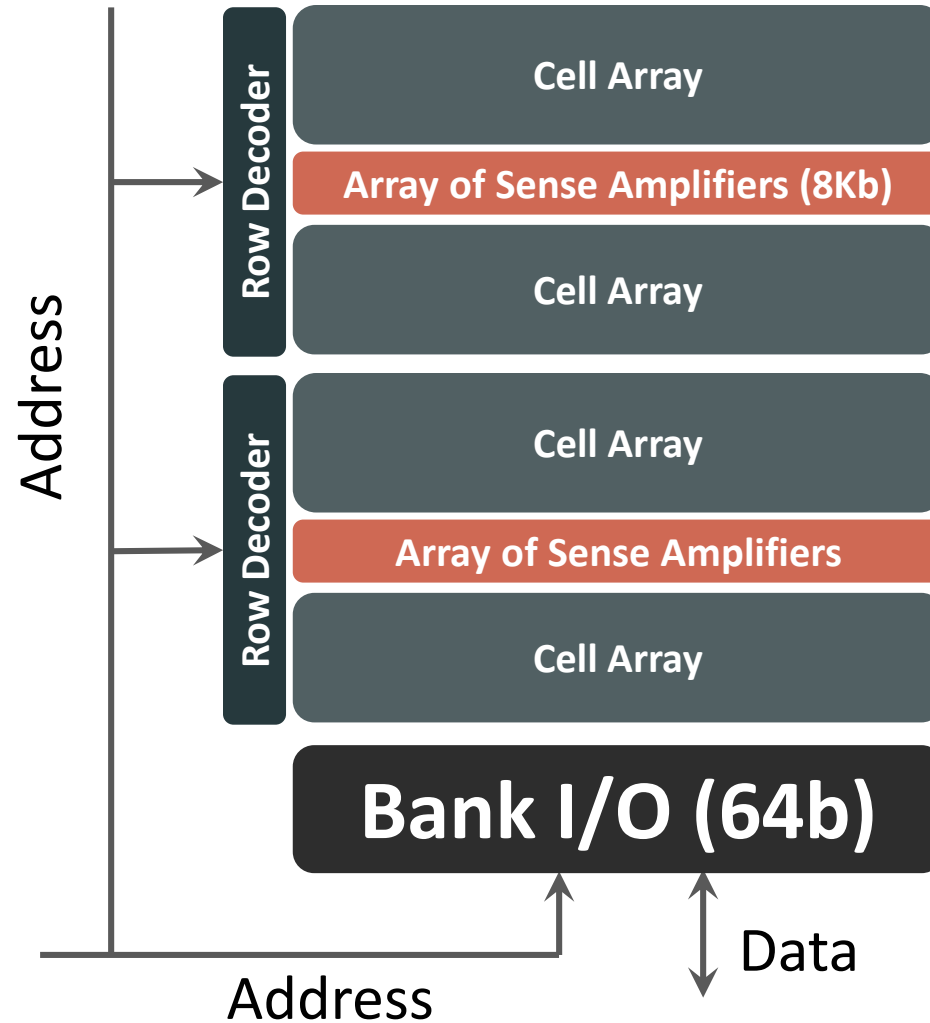
# DRAM Cell Operation



# DRAM Subarray – Building Block for DRAM Chip



# DRAM Bank



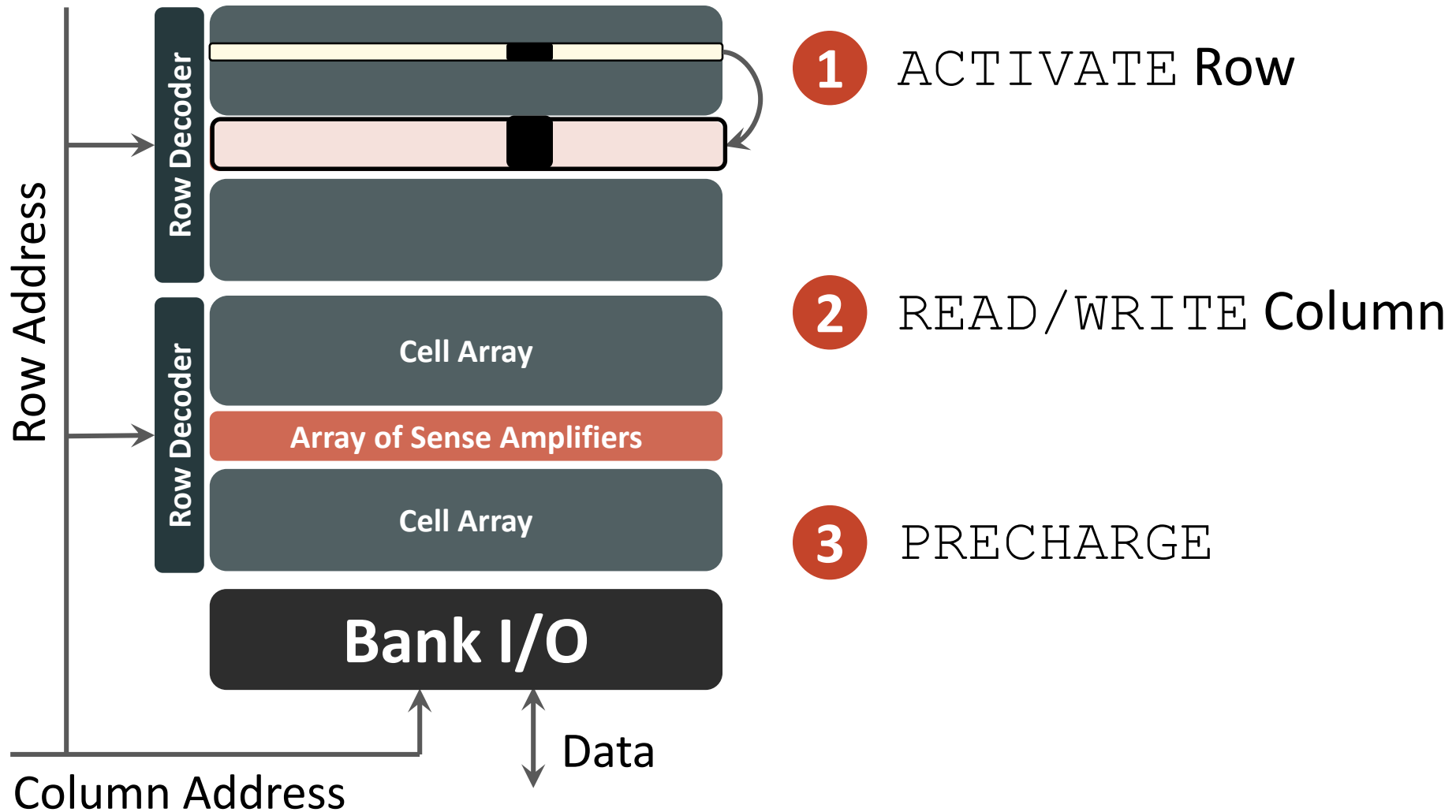
# DRAM Chip

Shared internal bus



Memory channel - 8bits

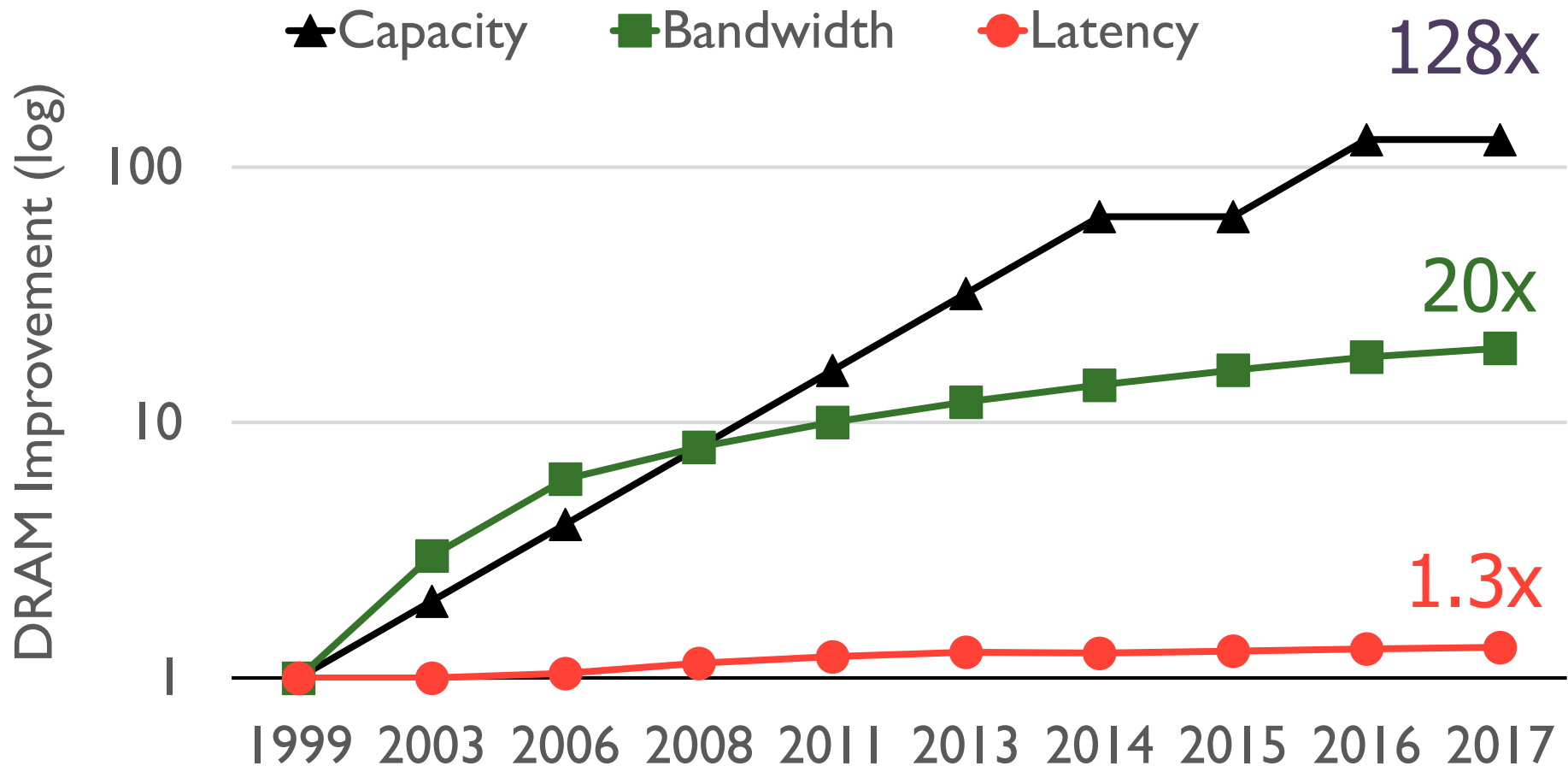
# DRAM Operation



# Memory Latency: Fundamental Tradeoffs



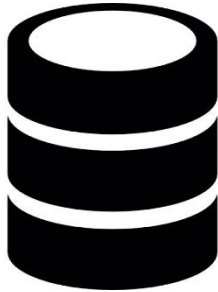
# Review: Memory Latency Lags Behind



Memory latency remains almost constant

# DRAM Latency Is Critical for Performance

---



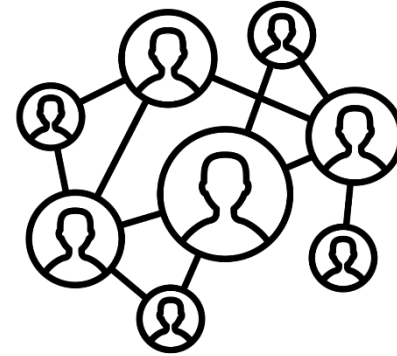
## In-memory Databases

[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]

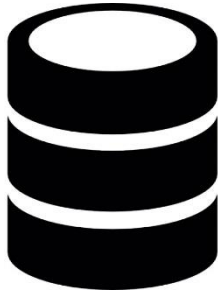


## Datacenter Workloads

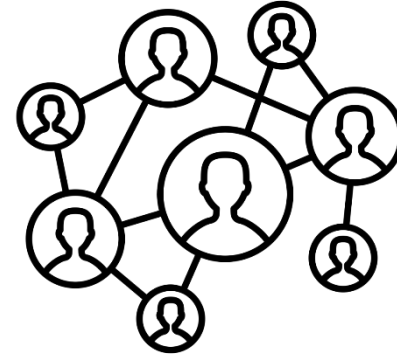
[Kanev+ (Google), ISCA'15]

# DRAM Latency Is Critical for Performance

---



**In-memory Databases**



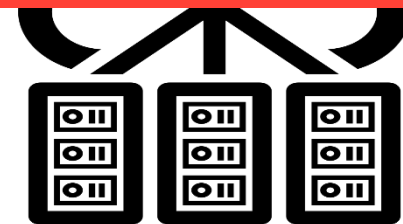
**Graph/Tree Processing**

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?

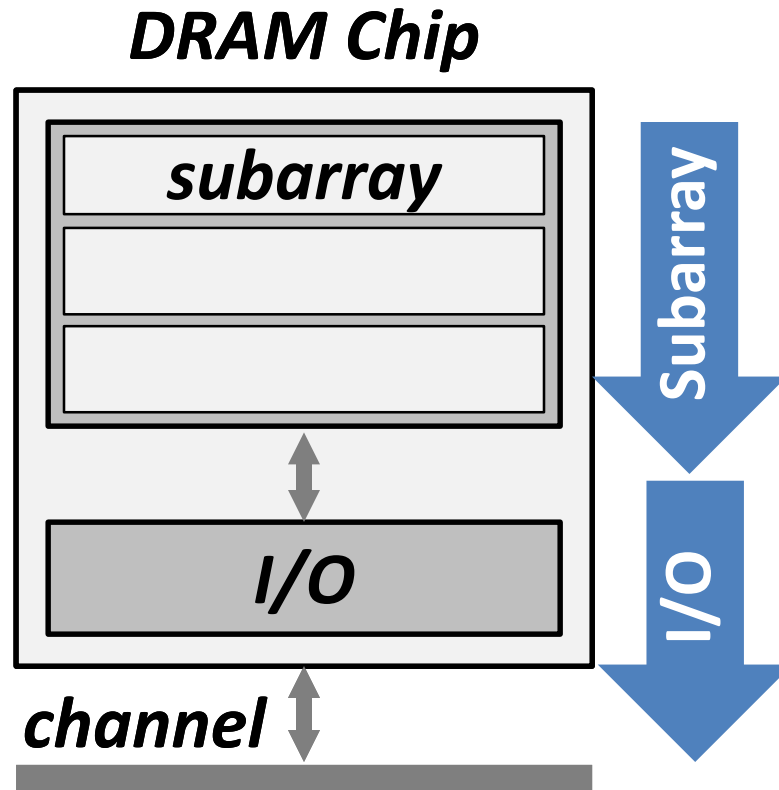
---

- 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
  - 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

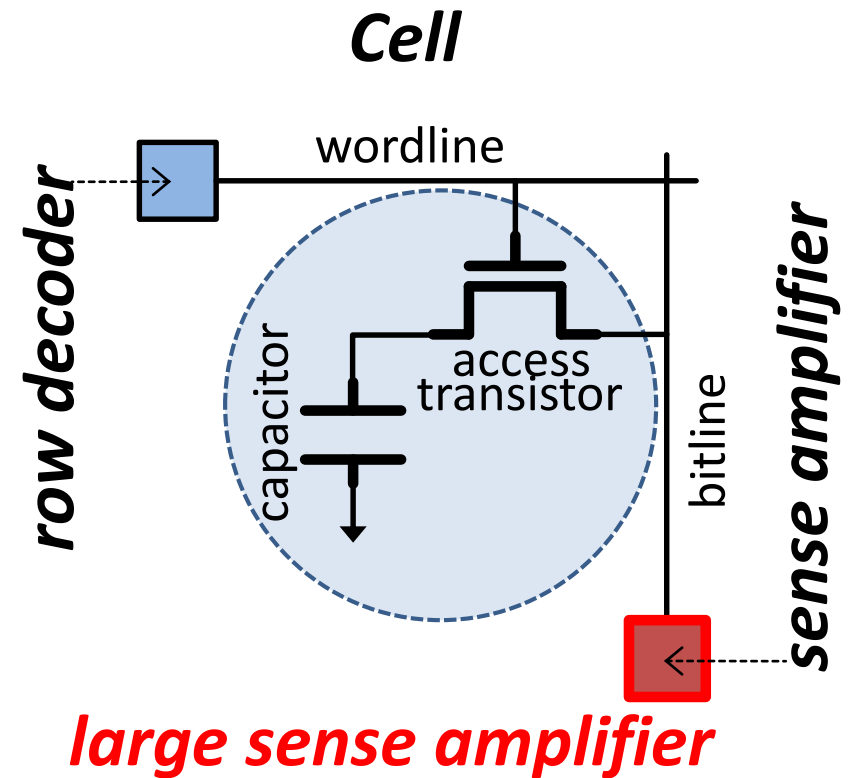
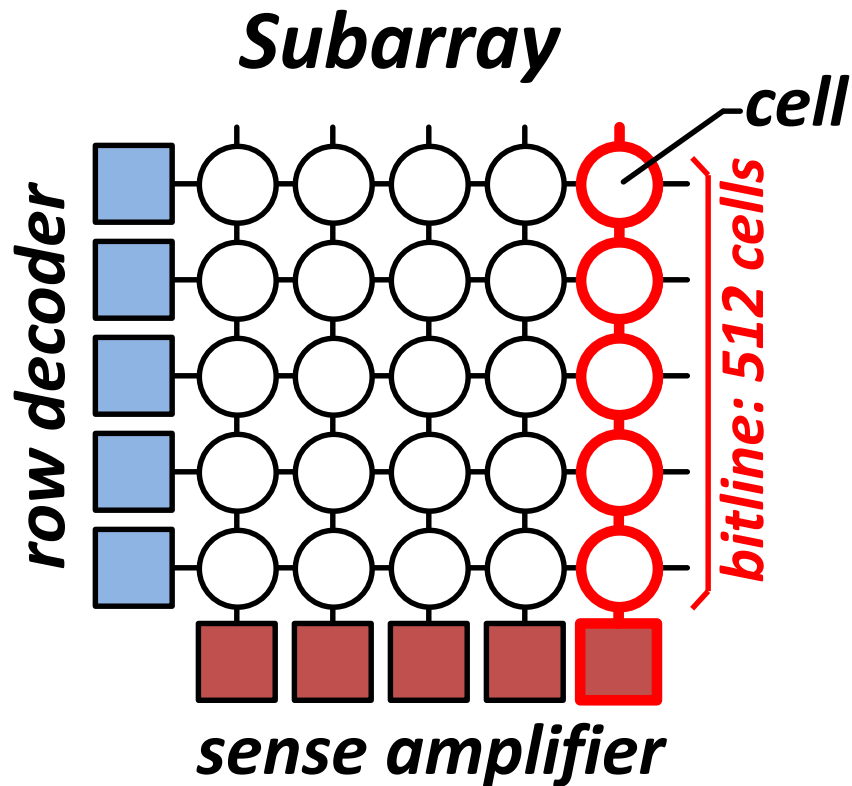
# What Causes the Long Latency?



$DRAM\ Latency = \text{Subarray Latency} + I/O\ Latency$

***Dominant***

# Why is the Subarray So Slow?

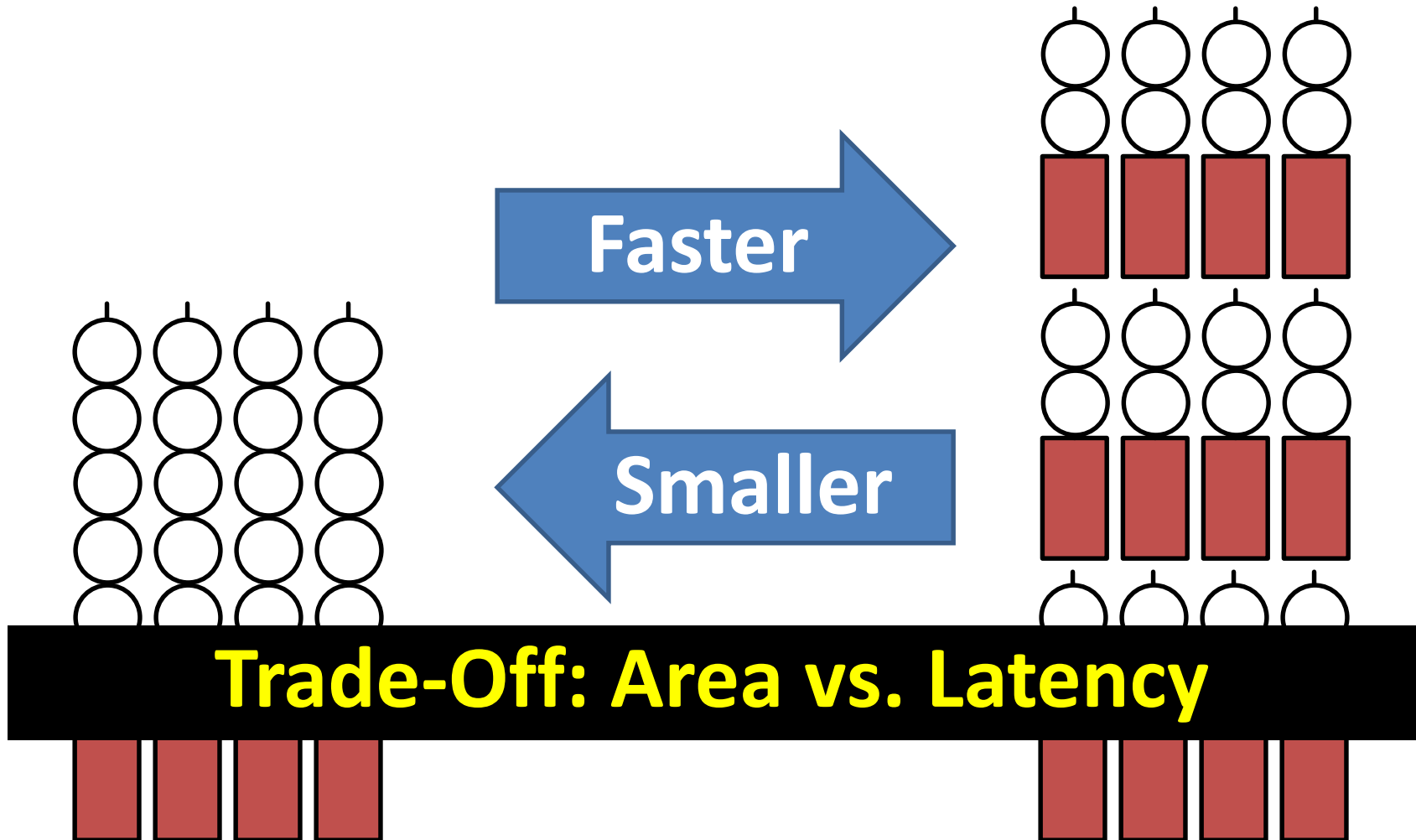


- Long bitline
  - Amortizes sense amplifier cost → Small area
  - Large bitline capacitance → High latency & power

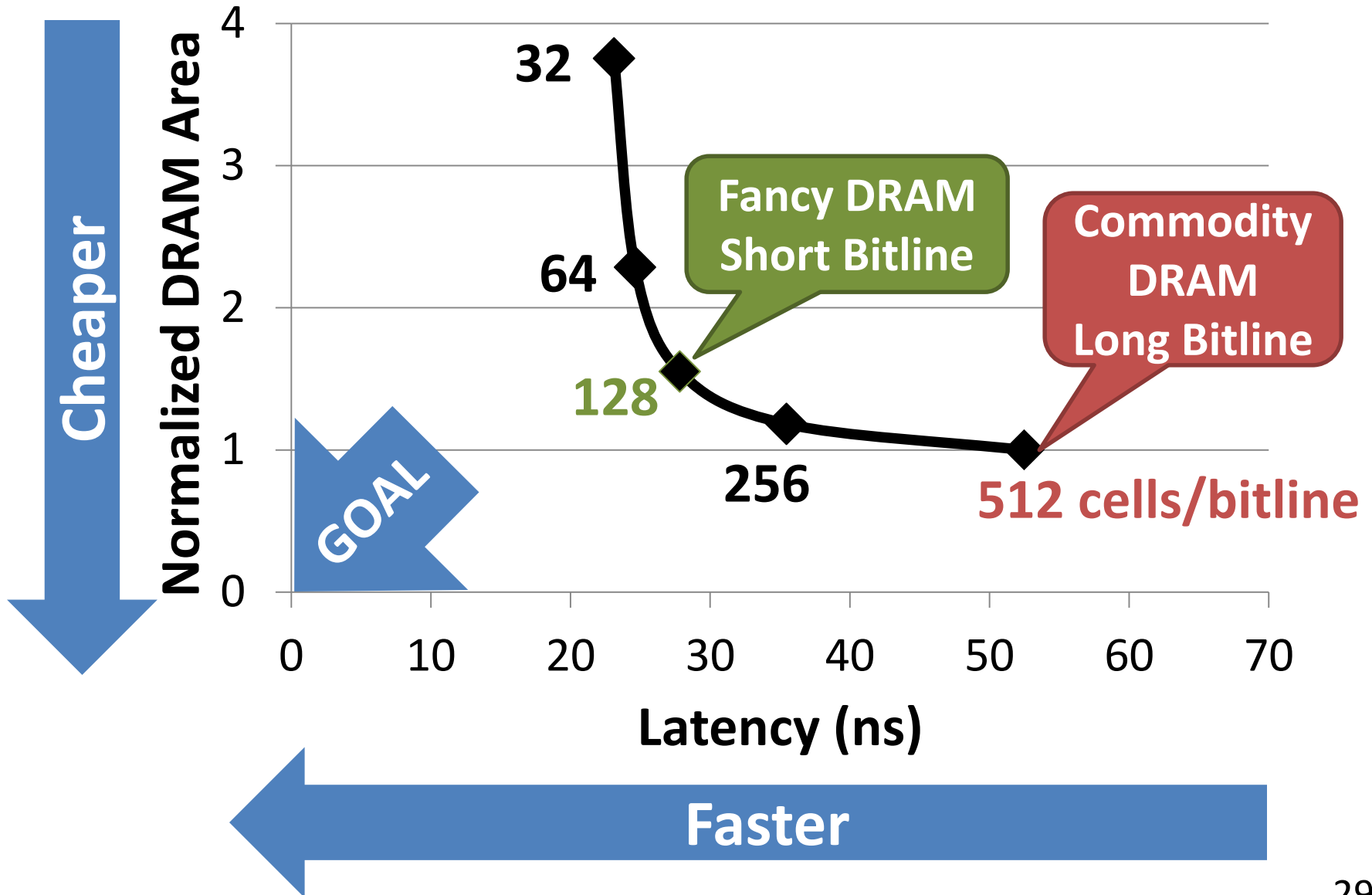
# Trade-Off: Area (Die Size) vs. Latency

Long Bitline

Short Bitline



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

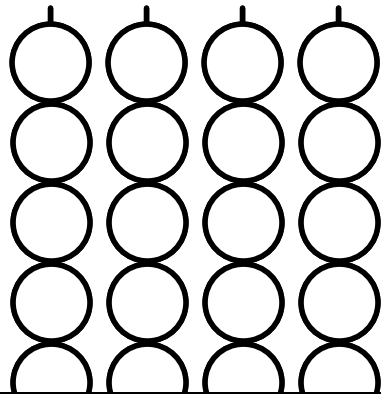


# Approximating the Best of Both Worlds

**Long Bitline**

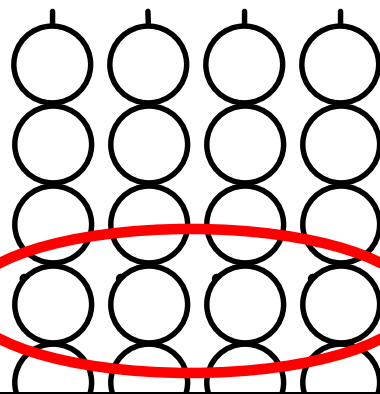
*Small Area*

~~High Latency~~



*Need Isolation*

**Our Proposal**

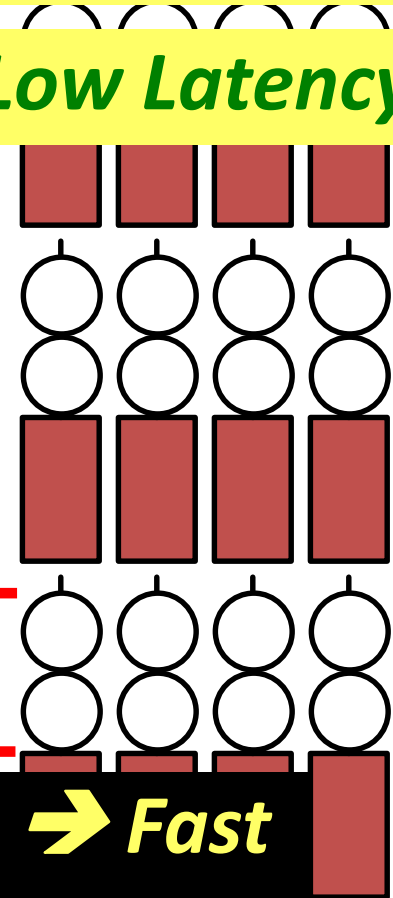


*Add Isolation Transistors*

**Short Bitline**

~~Large Area~~

*Low Latency*



*tline → Fast*

# Approximating the Best of Both Worlds

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

*Small Area*

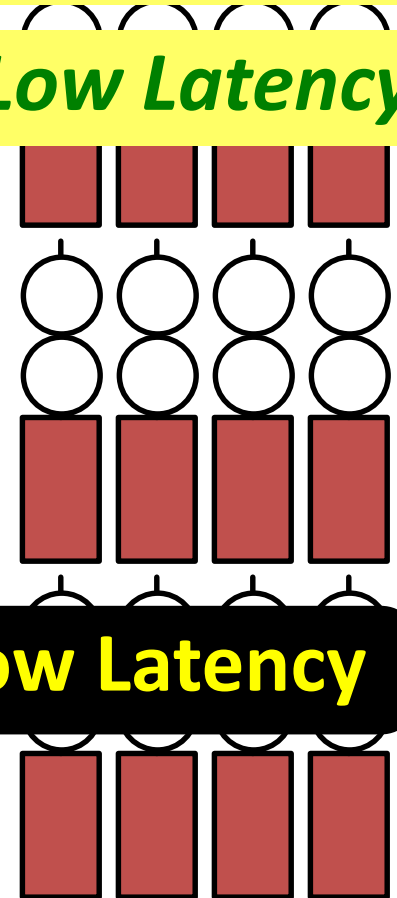
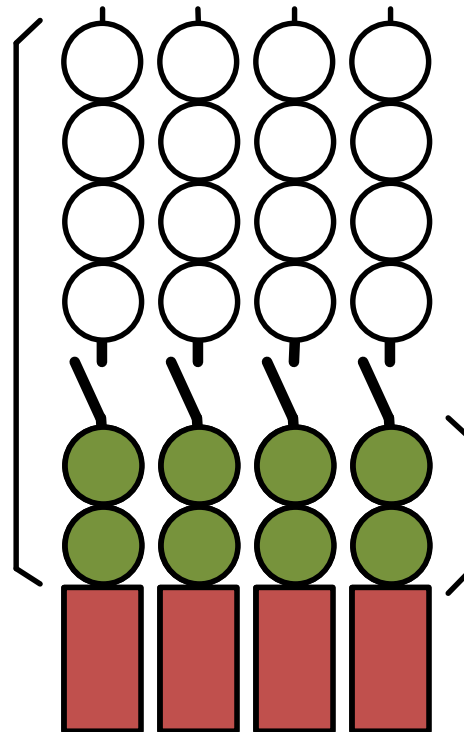
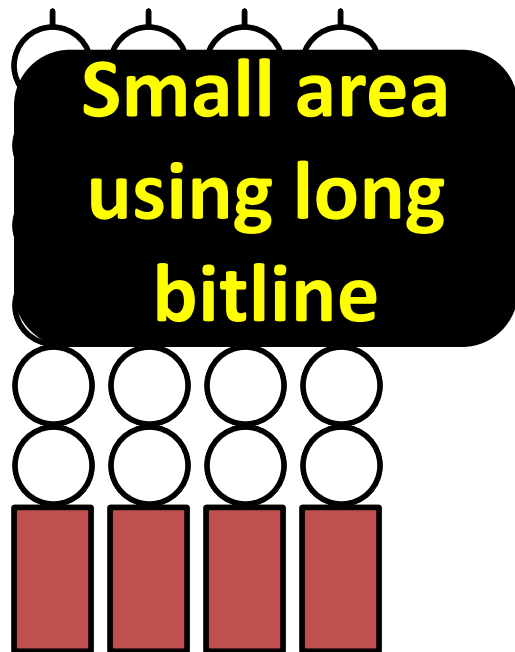
*Small Area*

~~*Large Area*~~

~~*High Latency*~~

*Low Latency*

*Low Latency*



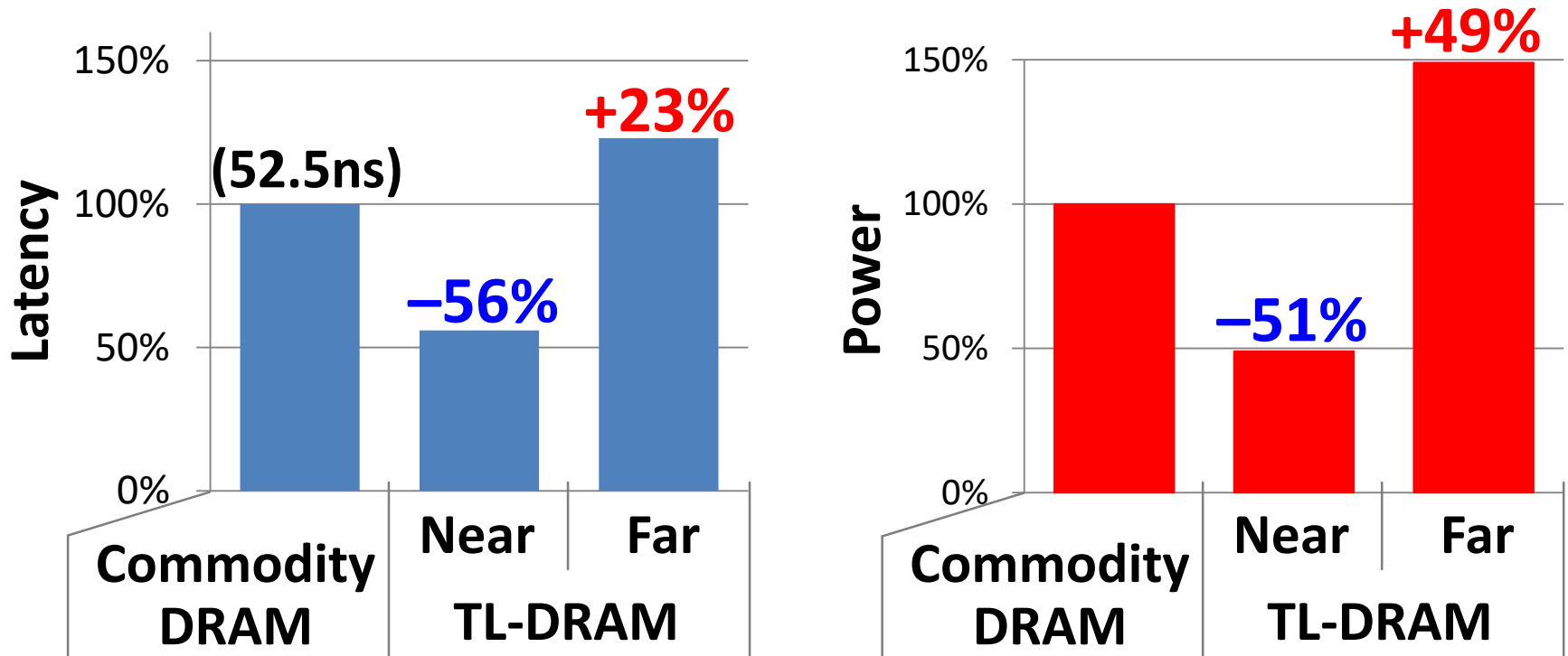
# Latency, Power, and Area Evaluation

- **Commodity DRAM:** 512 cells/bitline
- **TL-DRAM:** 512 cells/bitline
  - Near segment: 32 cells
  - Far segment: 480 cells
- **Latency Evaluation**
  - SPICE simulation using circuit-level DRAM model
- **Power and Area Evaluation**
  - DRAM area/power simulator from Rambus
  - DDR3 energy calculator from Micron



# Commodity DRAM vs. TL-DRAM [HPCA 2013]

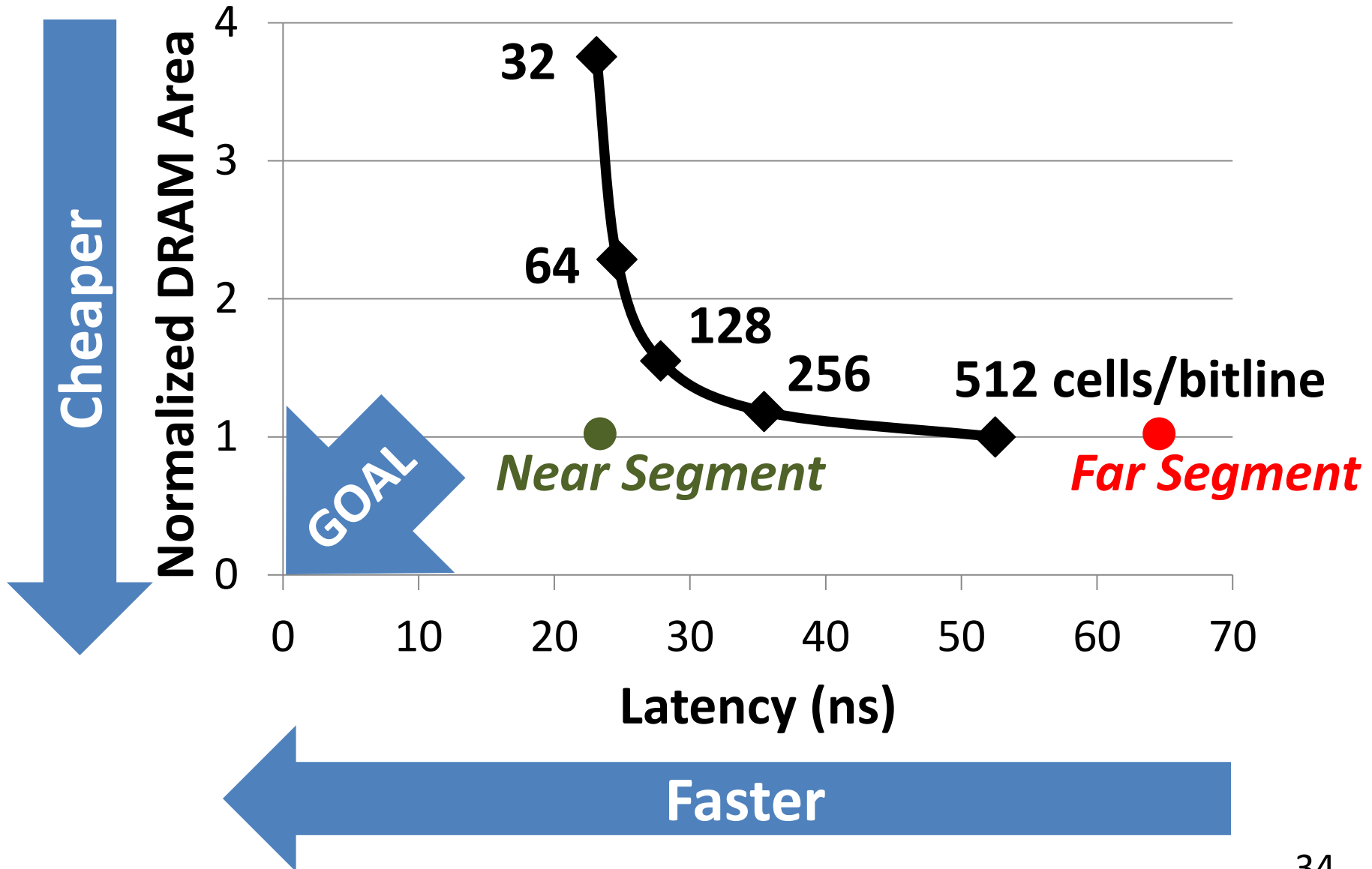
- DRAM Latency (tRC) • DRAM Power



- DRAM Area Overhead

~3%: mainly due to the isolation transistors

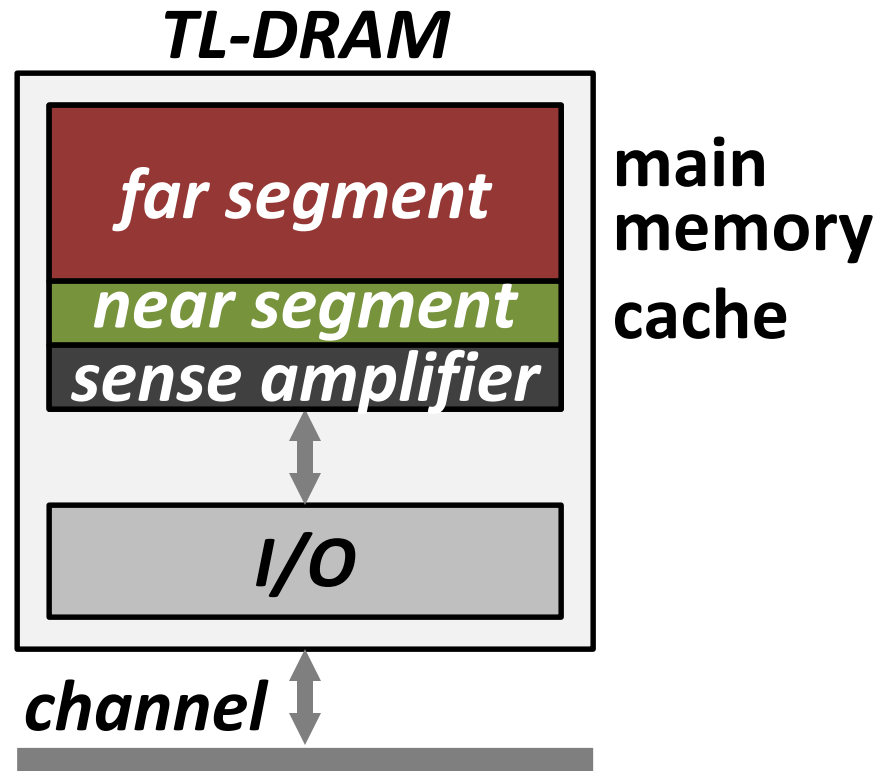
# Trade-Off: Area (Die-Area) vs. Latency



# 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

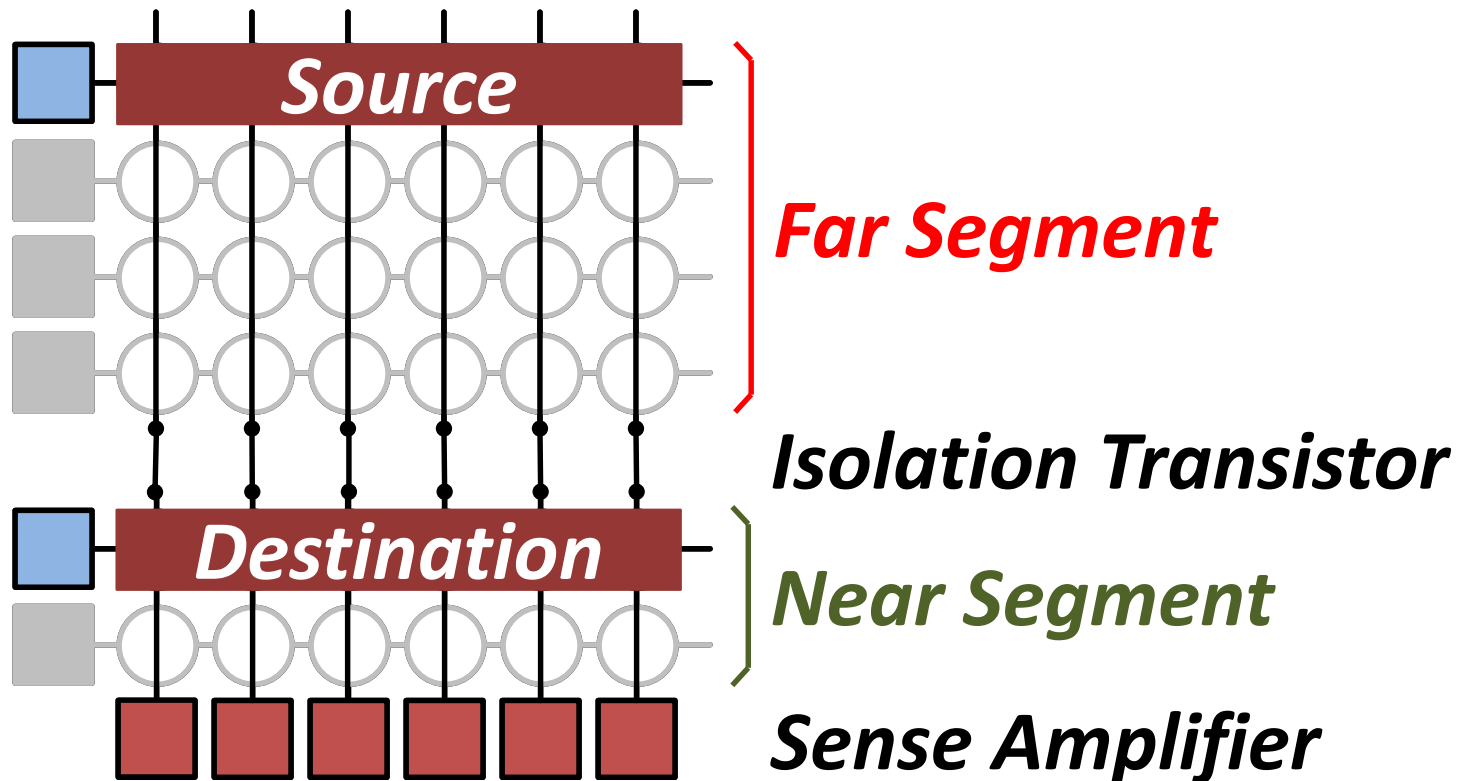
# Near Segment as Hardware-Managed Cache



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

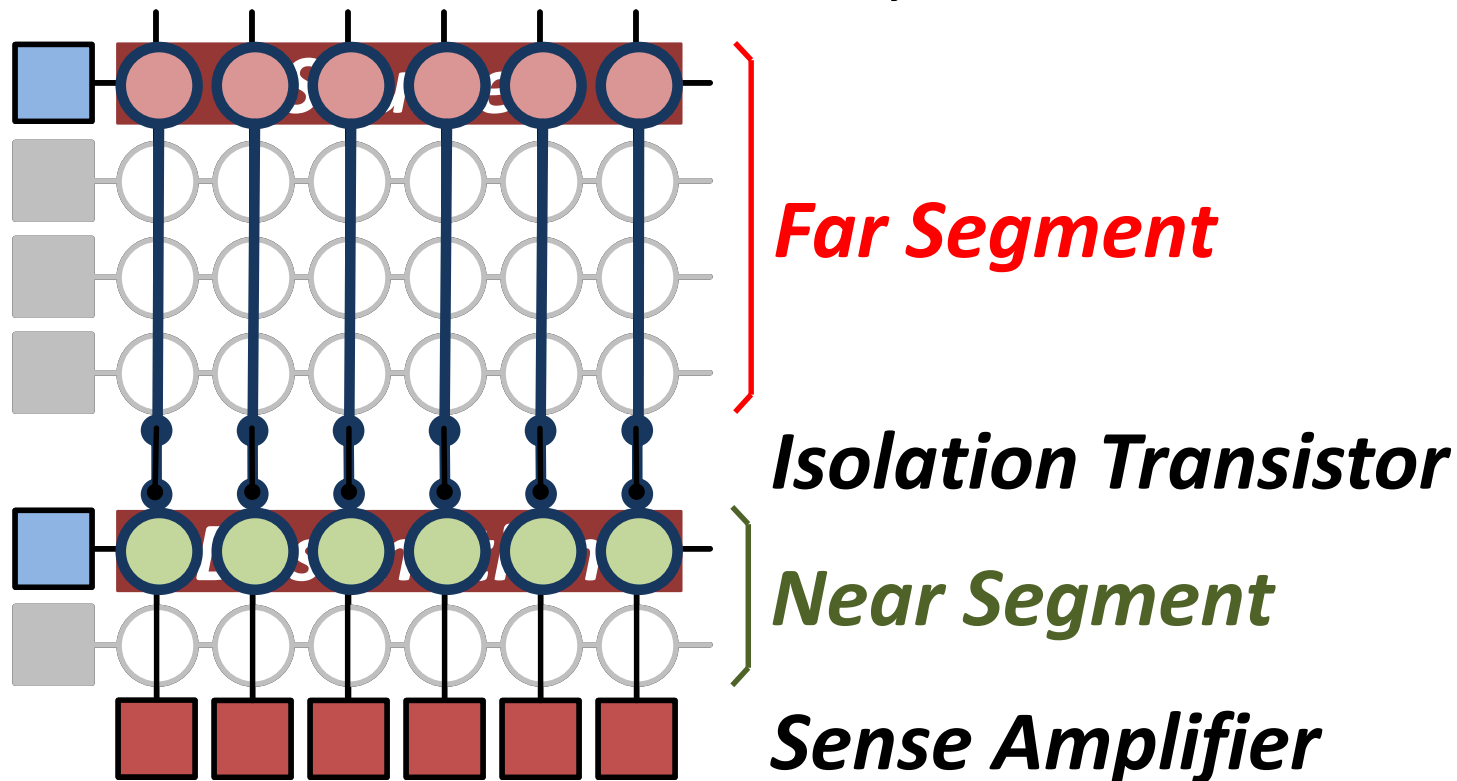
# Inter-Segment Migration

- **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



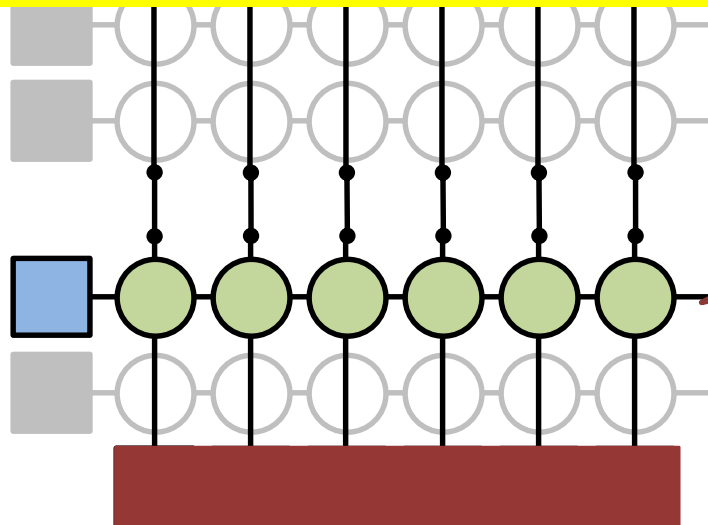
# Inter-Segment Migration

- 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**



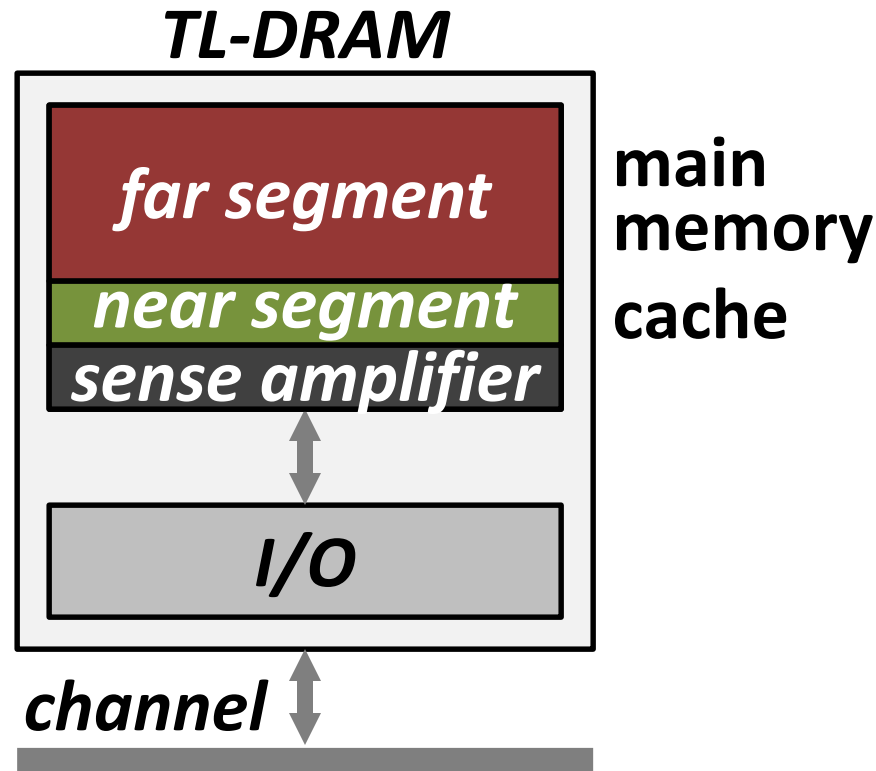
Step 2: Activate destination row to connect cell and bitline

150 pA Transistor

*Near Segment*

*Sense Amplifier*

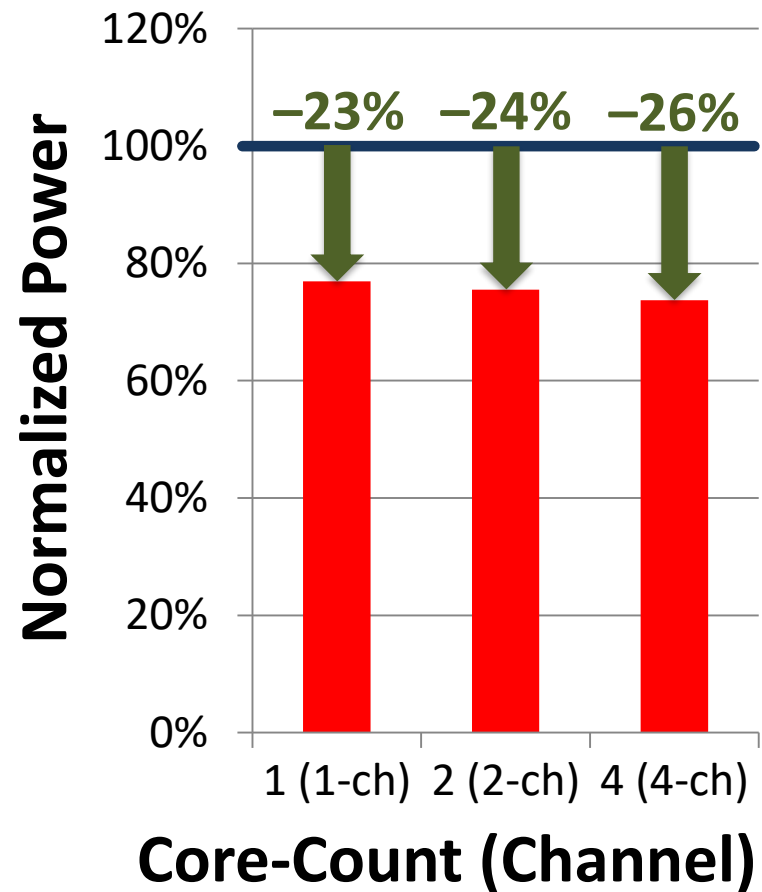
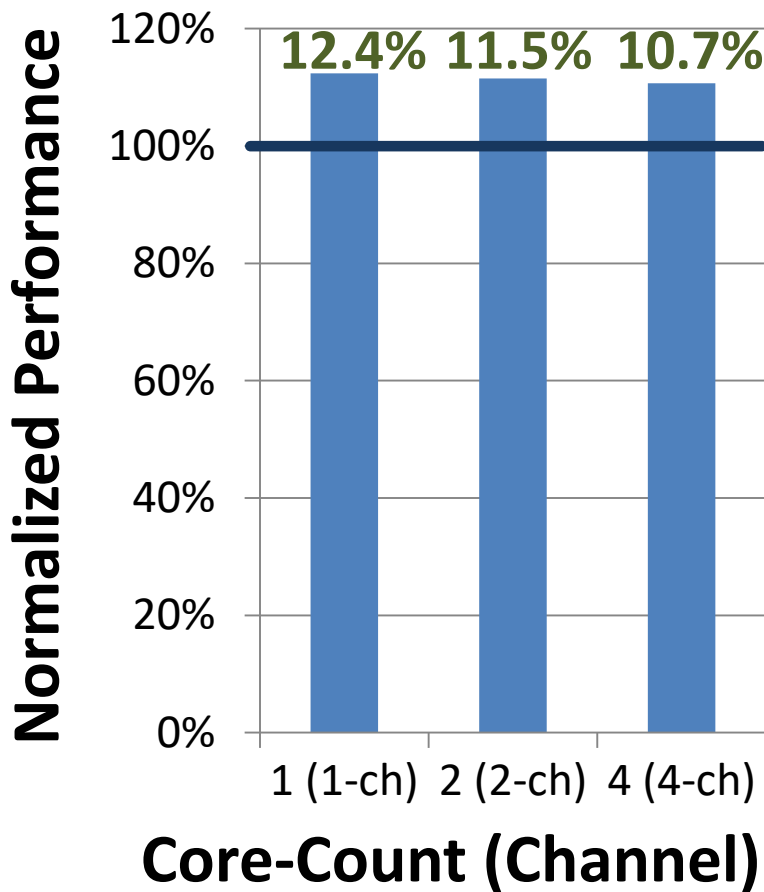
# Near Segment as Hardware-Managed Cache



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

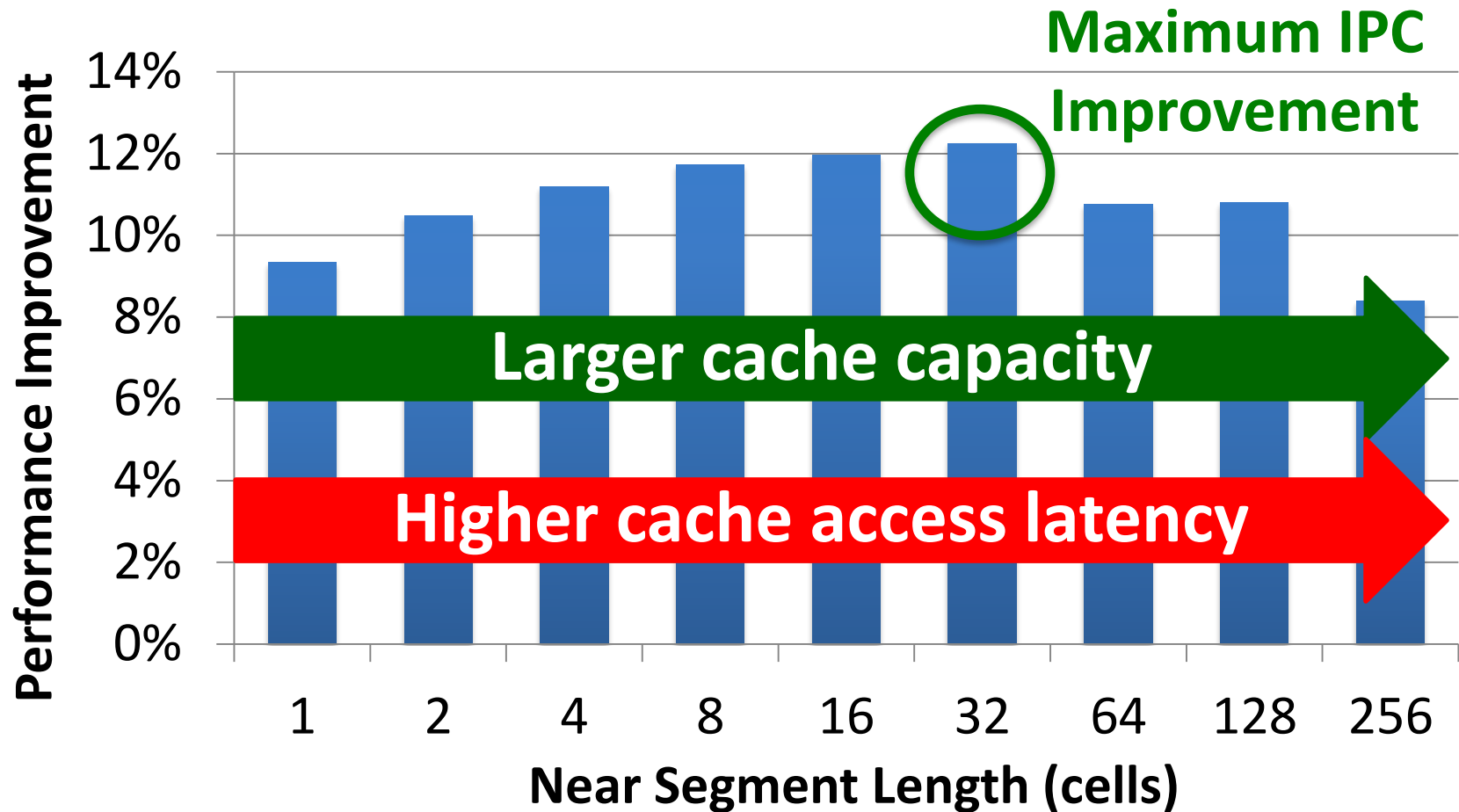


# Performance & Power Consumption



*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

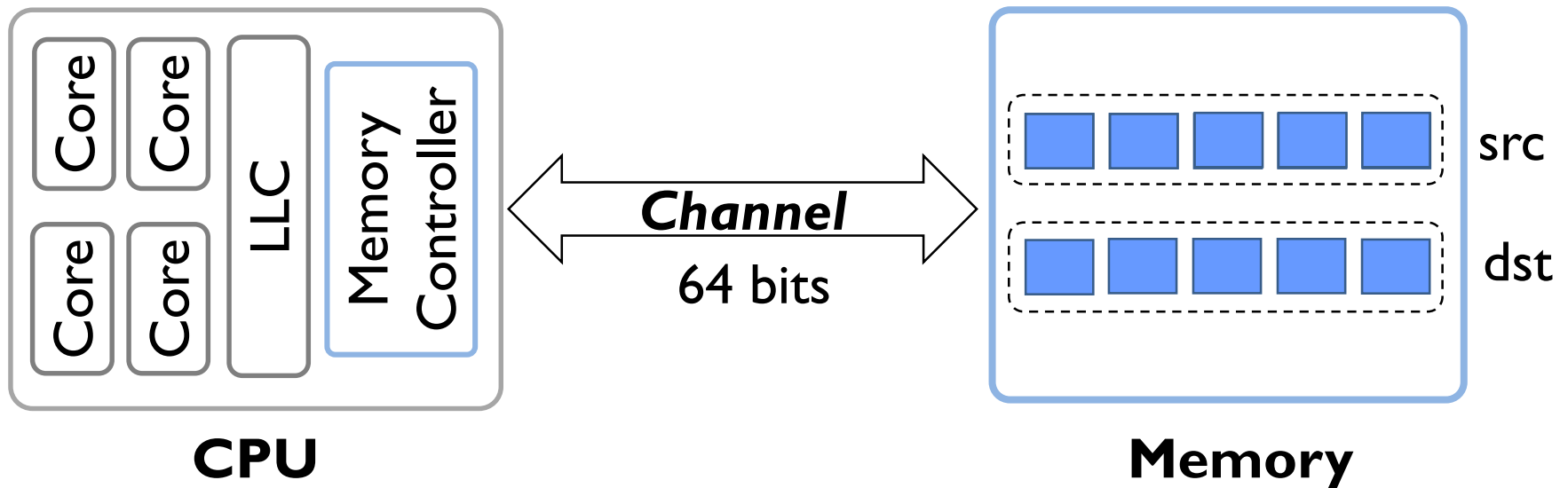
# LISA: Low-Cost Inter-Linked Subarrays

## [HPCA 2016]

# Problem: Inefficient Bulk Data Movement

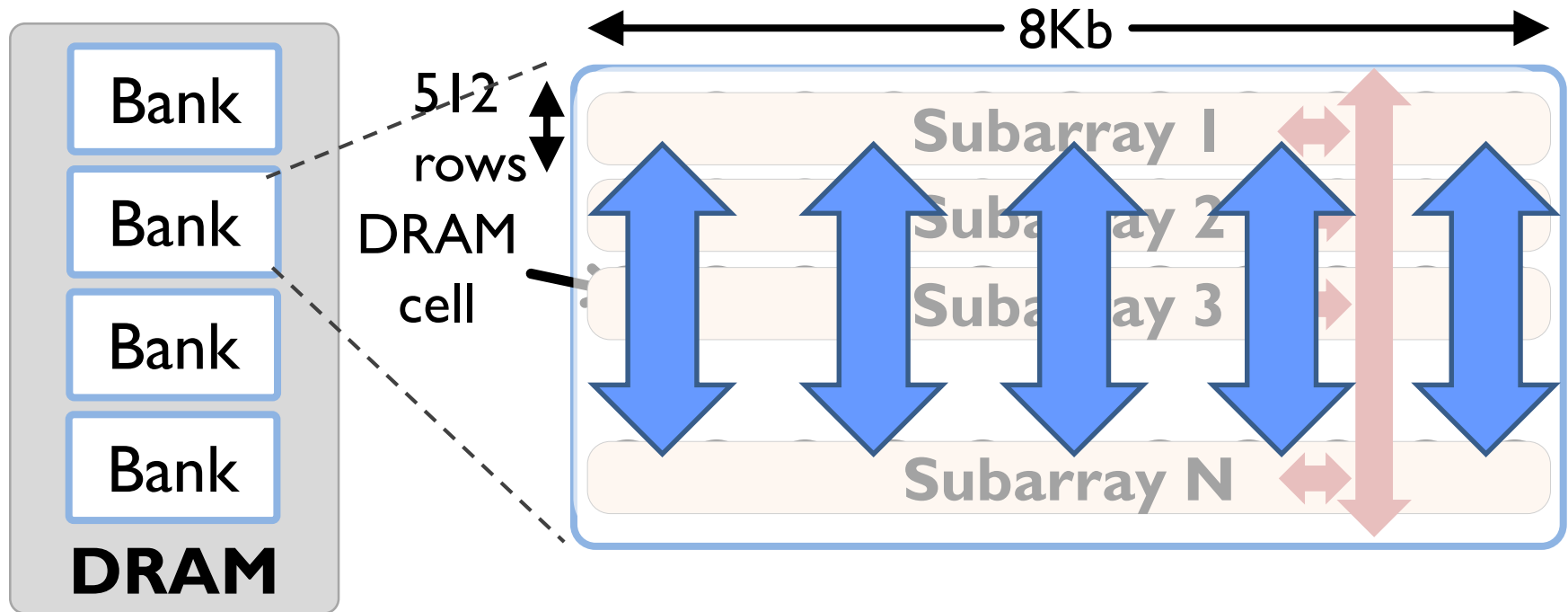
*Bulk data movement is a key operation in many applications*

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



**Long latency and high energy**

# Moving Data Inside DRAM?

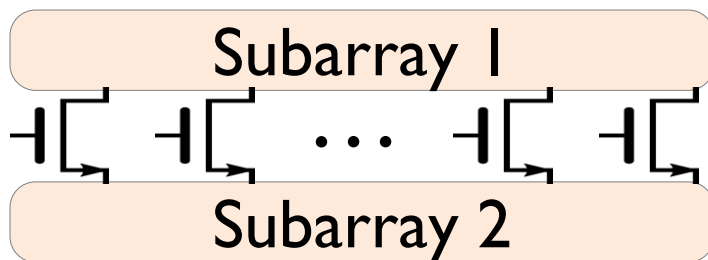


**Goal: Provide a new substrate to enable wide connectivity between subarrays**

# 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)

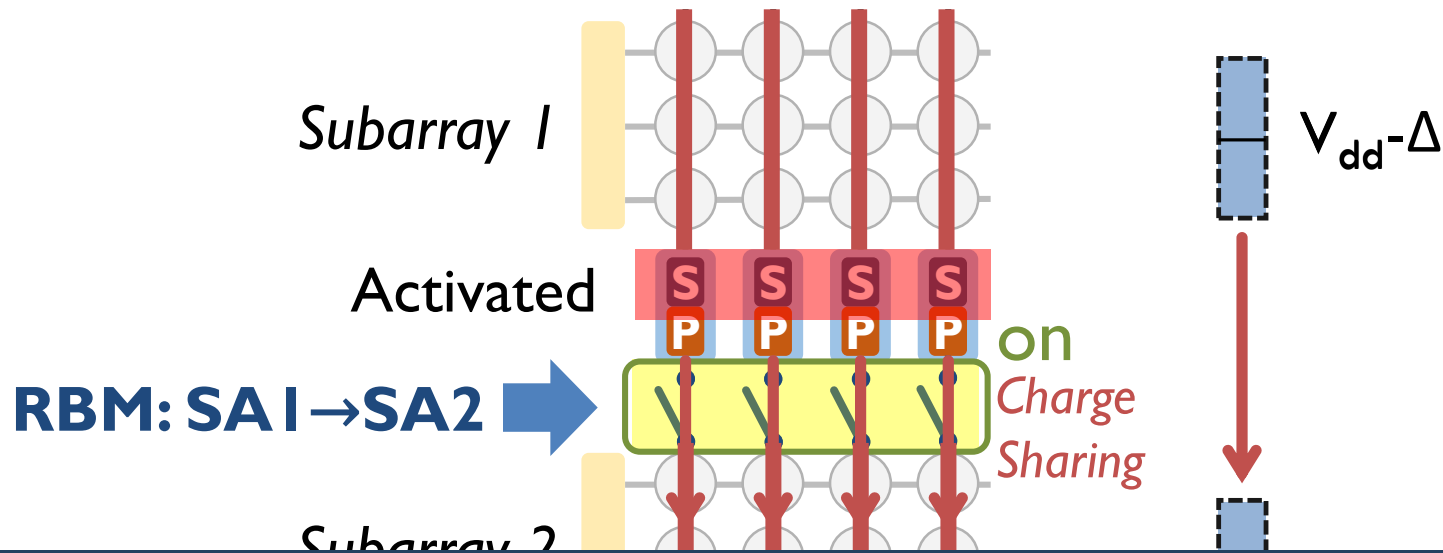
→ 5% speedup

**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



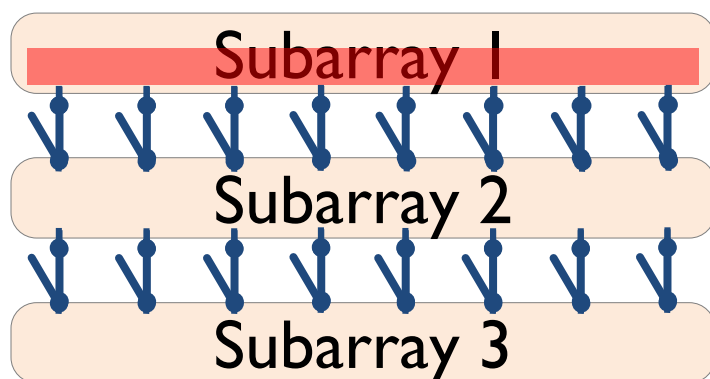
**RBM transfers an entire row b/w subarrays**



# RBM Analysis

---

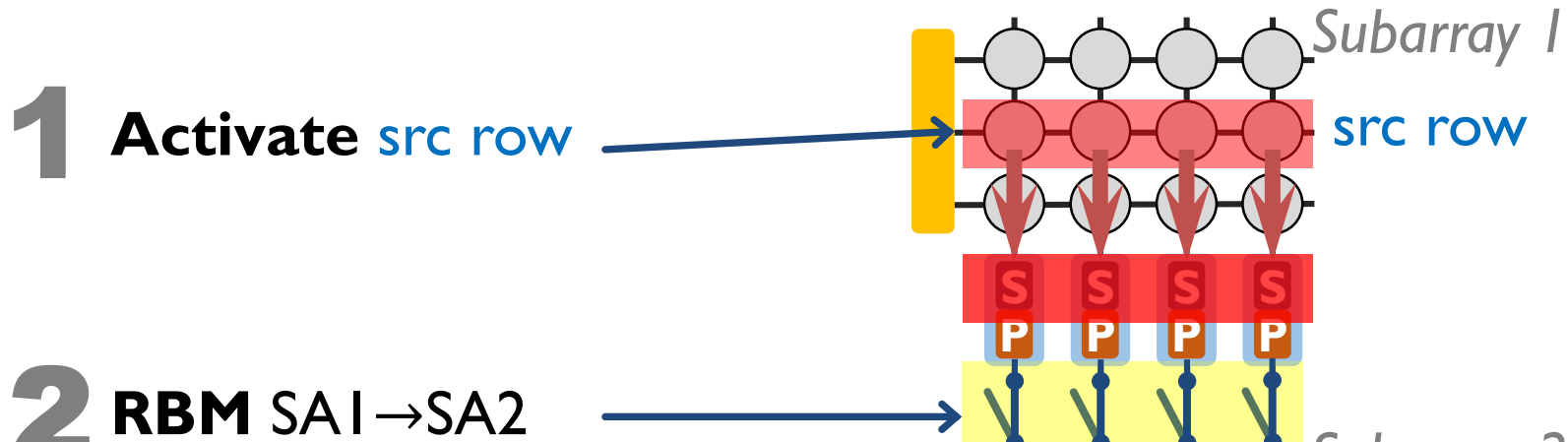
- The range of RBM depends on the DRAM design
  - Multiple RBMs to move data across  $> 3$  subarrays



- 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

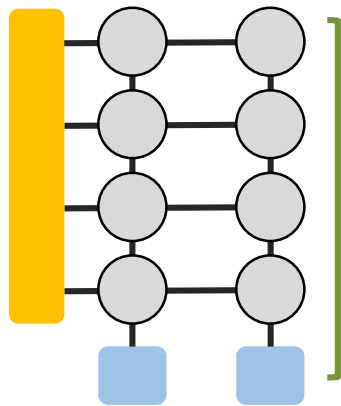


Reduces row-copy latency by 9.2x,  
DRAM energy by 48.1x

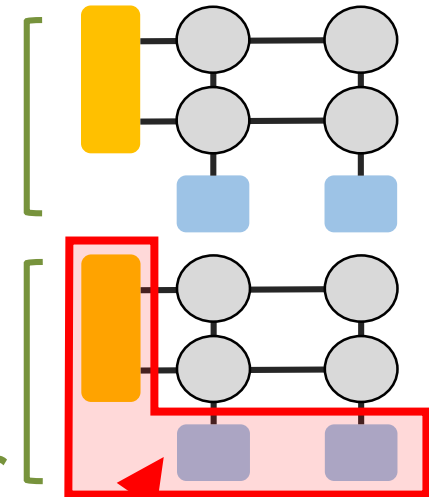
## 2. Variable Latency DRAM (VILLA)

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

**Long Bitline  
(DDR<sub>x</sub>)**



**Short Bitline  
(RLDRAM)**

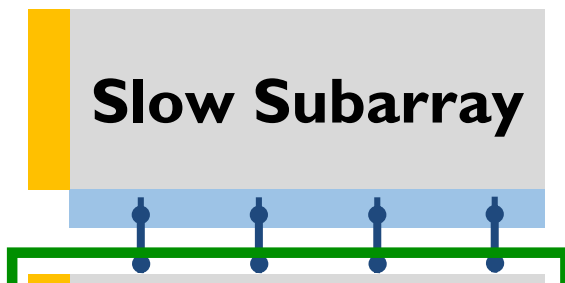


Shorter bitlines → faster  
**activate** and **precharge** time

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

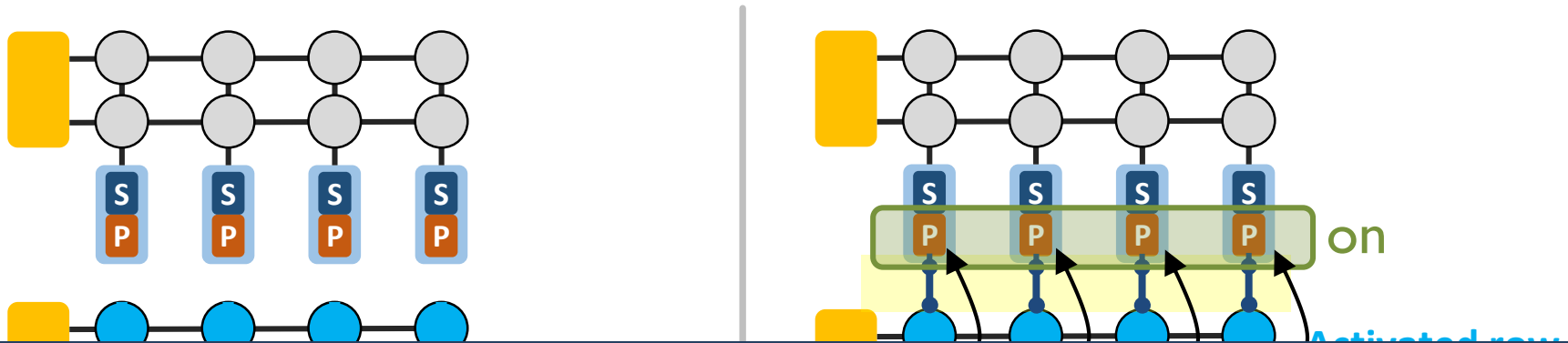


**Challenge:** VILLA cache requires frequent movement of data rows

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<sup>†</sup>, Prashant J. Nair<sup>\*</sup>, Donghyuk Lee<sup>†</sup>, Saugata Ghose<sup>†</sup>, Moinuddin K. Qureshi<sup>\*</sup>, and Onur Mutlu<sup>†</sup>

<sup>†</sup>Carnegie Mellon University    <sup>\*</sup>Georgia Institute of Technology

# What Causes the Long DRAM 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
  - 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
  - ...



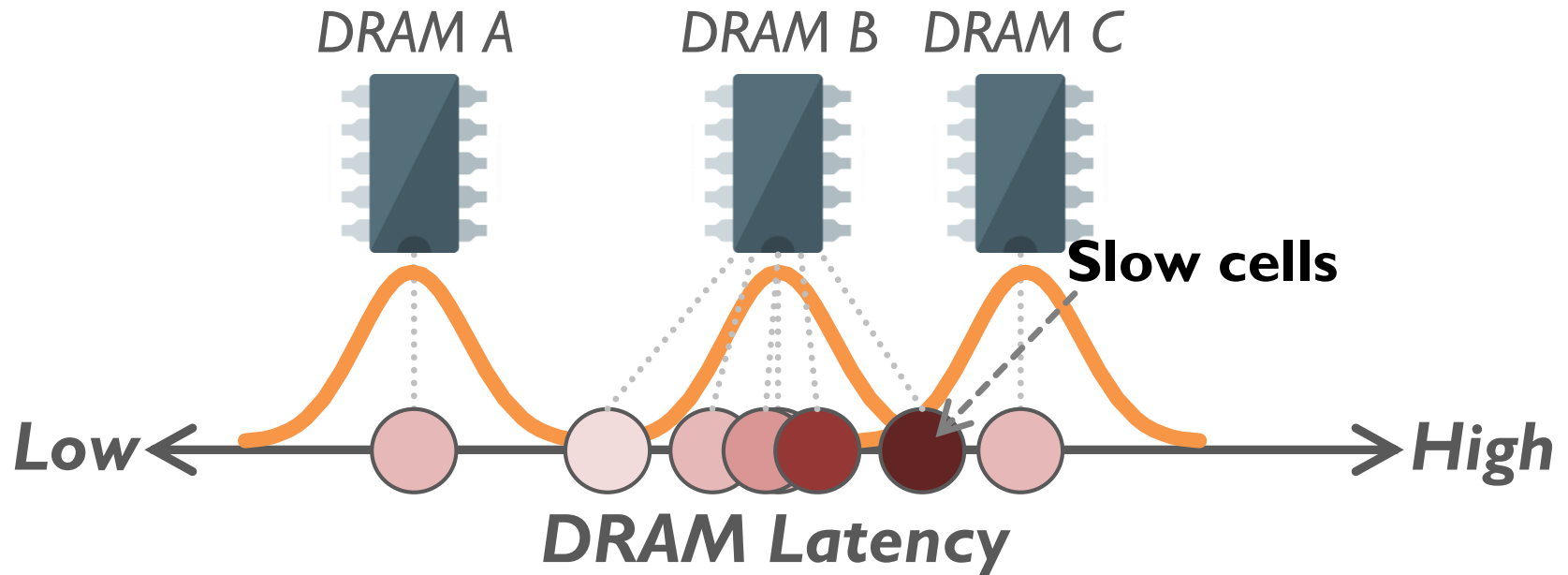
# Tackling the Fixed Latency Mindset

---

- 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

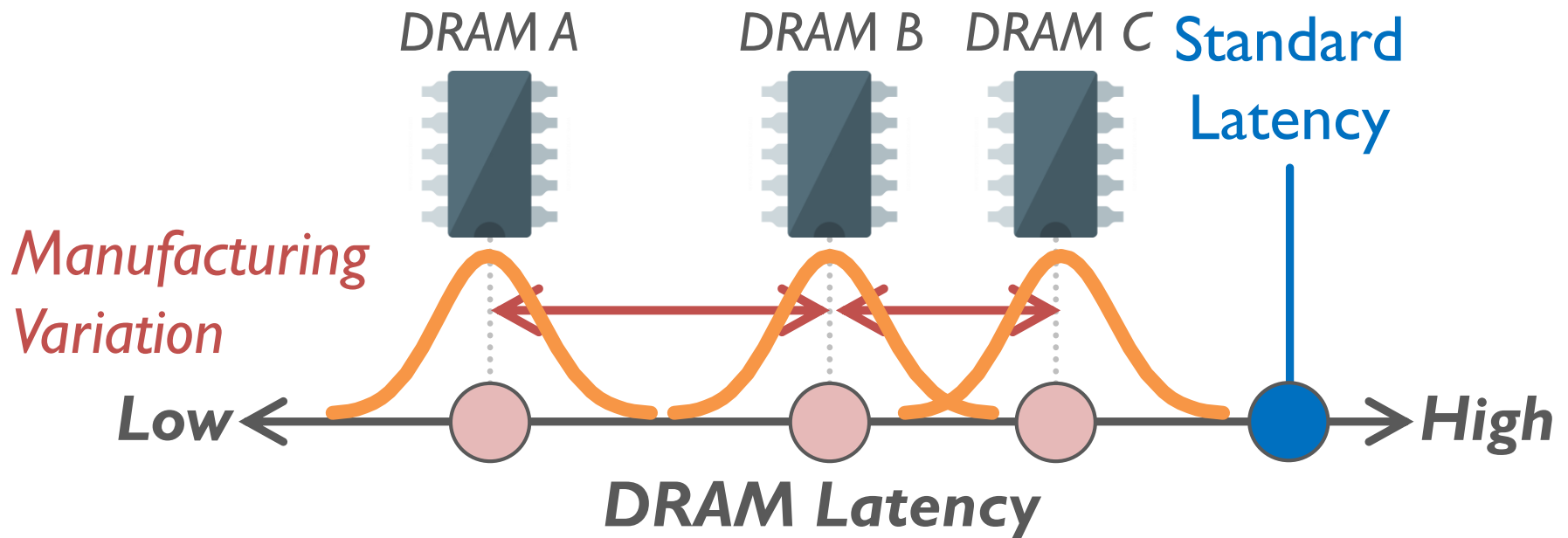
# Latency Variation in Memory Chips

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

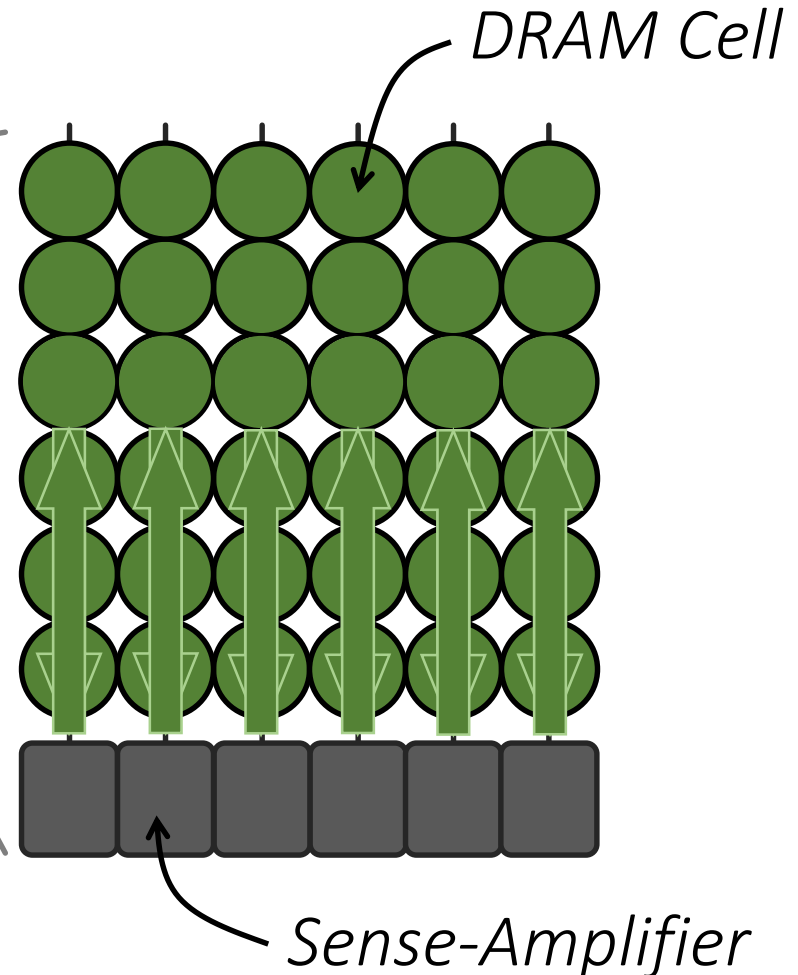
# Understanding and Exploiting Variation in DRAM Latency

# DRAM Stores Data as Charge

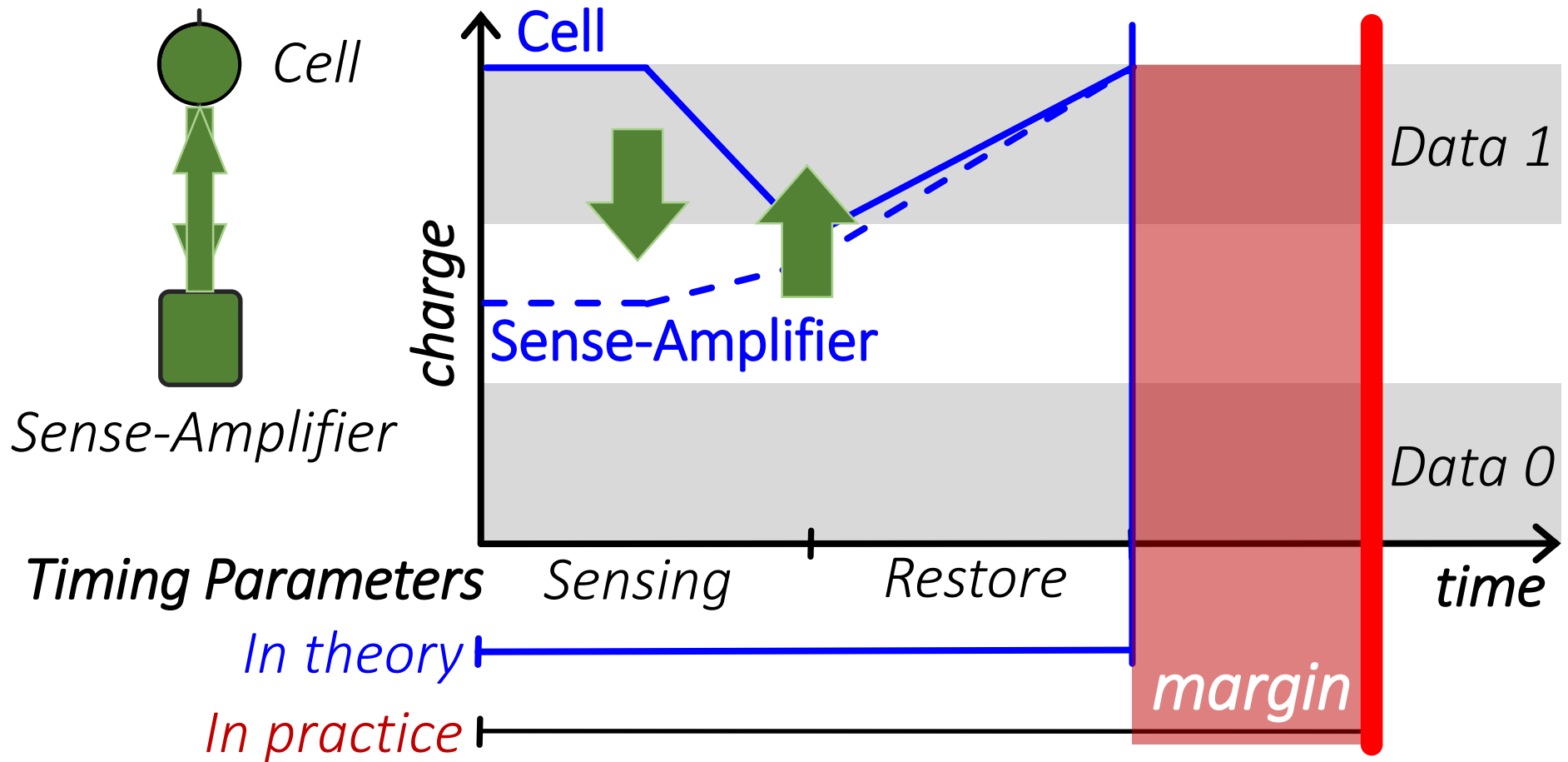


Three steps of  
charge movement

1. Sensing
2. Restore
3. Precharge



# DRAM Charge over Time



*Why does DRAM need the extra timing margin?*

# Two Reasons for Timing Margin

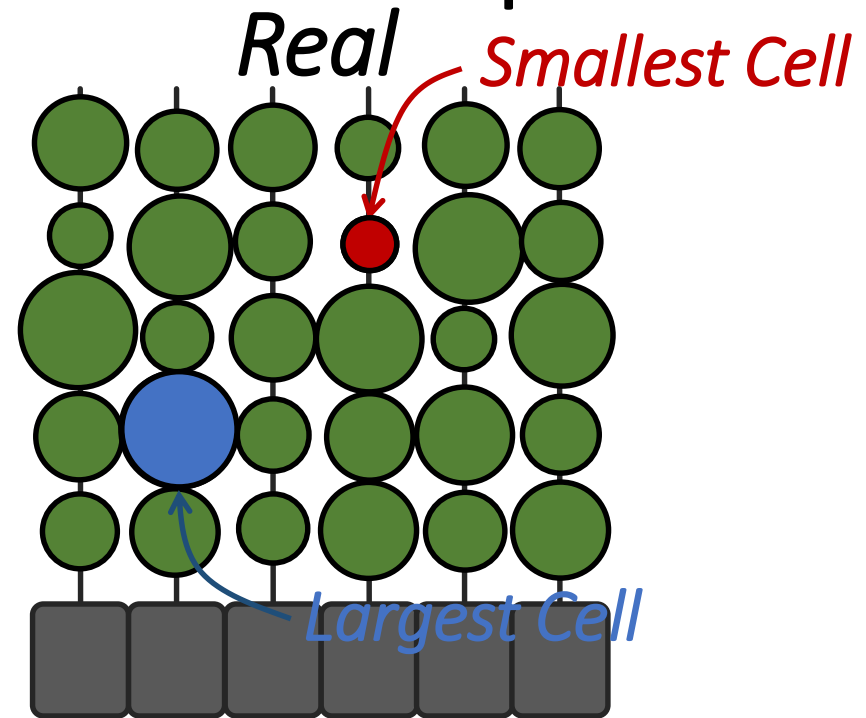
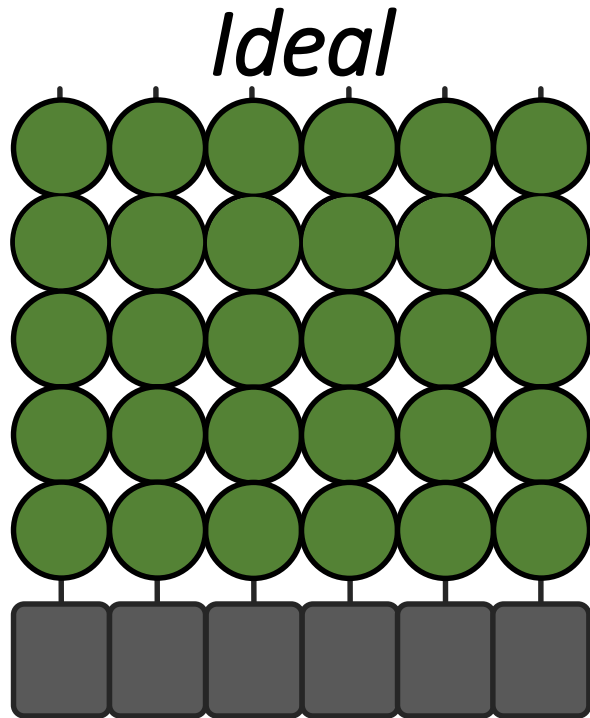
## *1. Process Variation*

- 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



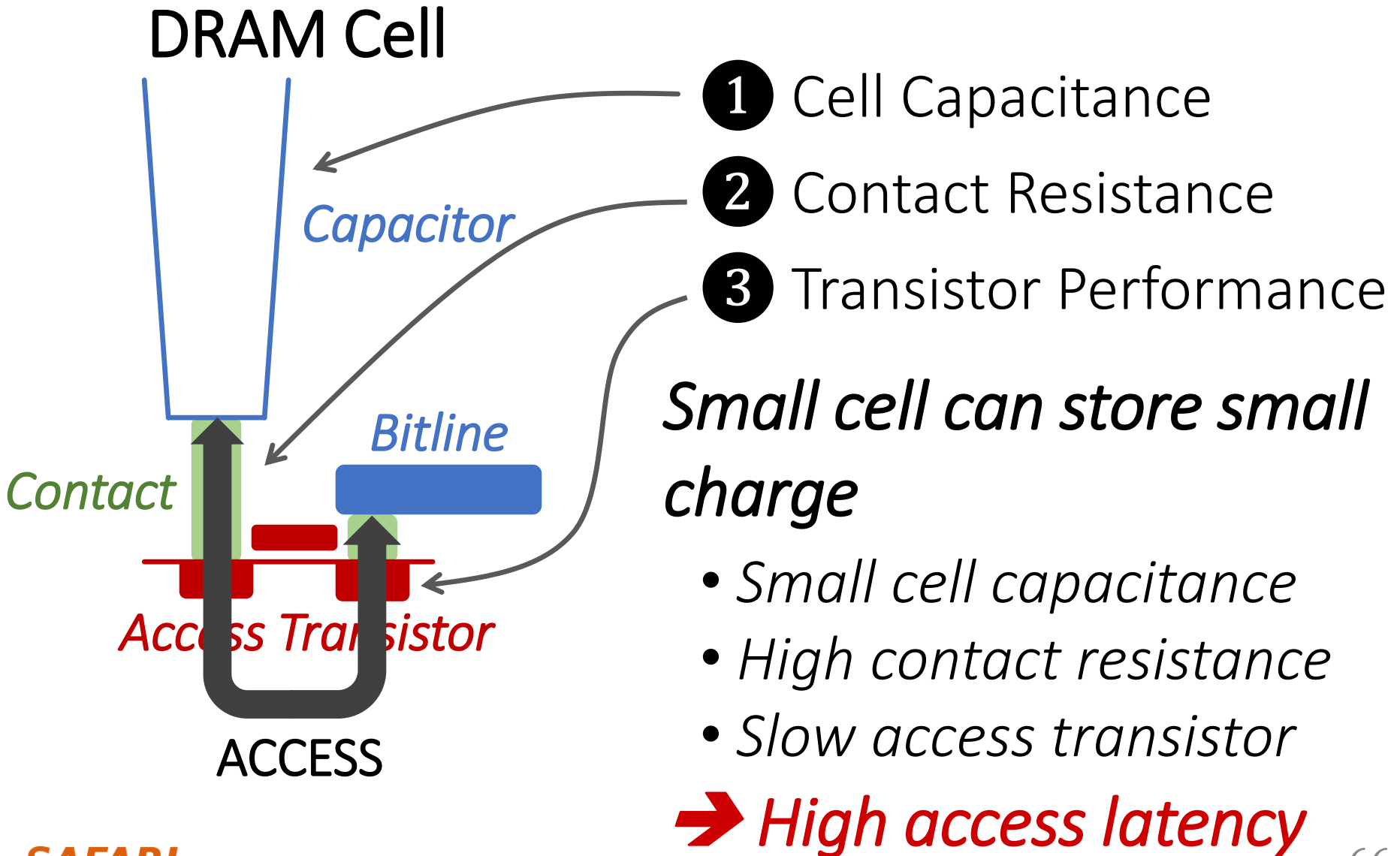
Same Size →  
Same Charge →  
Same Latency →

Large variation in cell size  
Large variation in charge  
Large variation in latency

Different Size →  
Different Charge →  
Different Latency →

Large variation in access latency

# Process Variation



# Two Reasons for Timing Margin

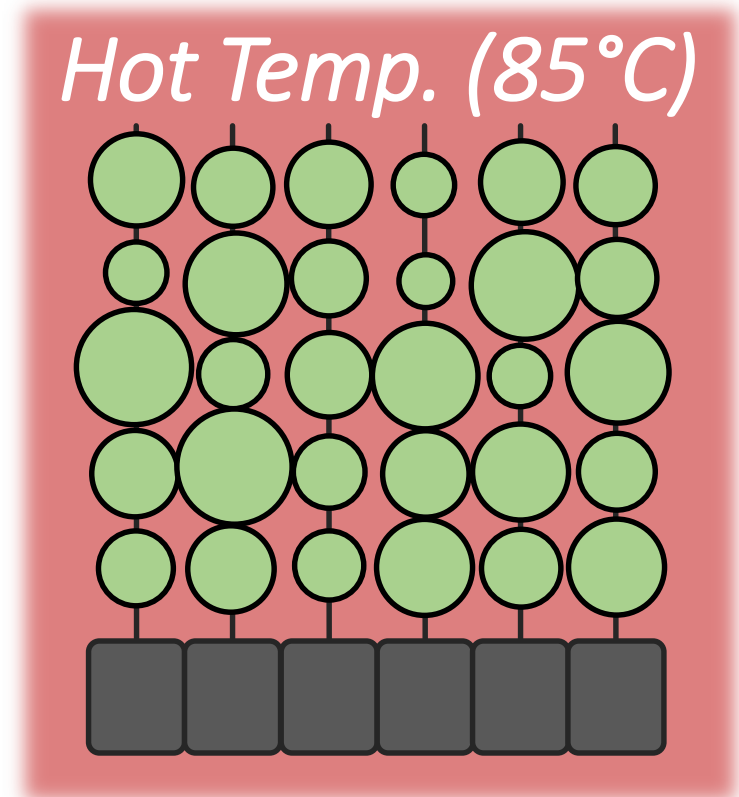
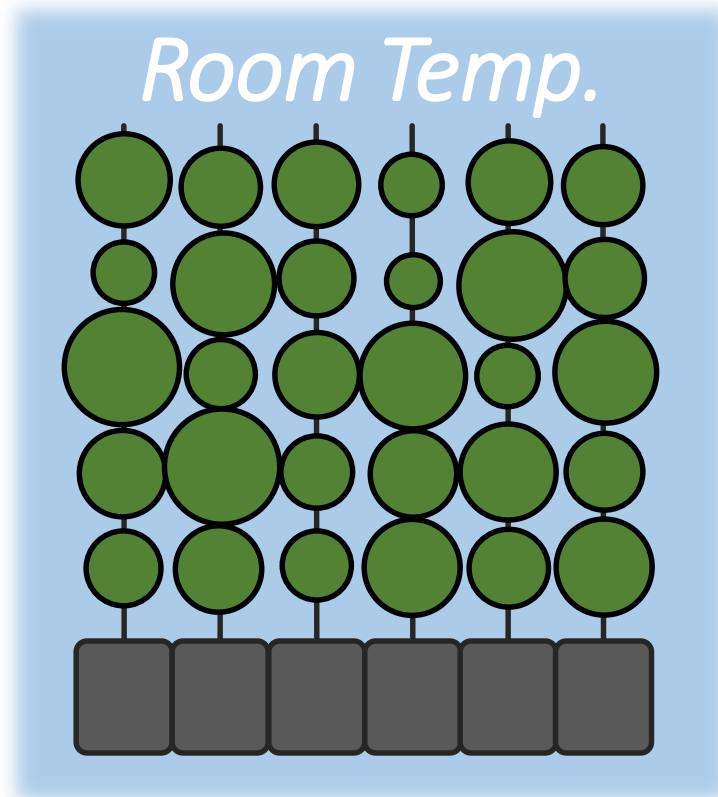
## *1. Process Variation*

- 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

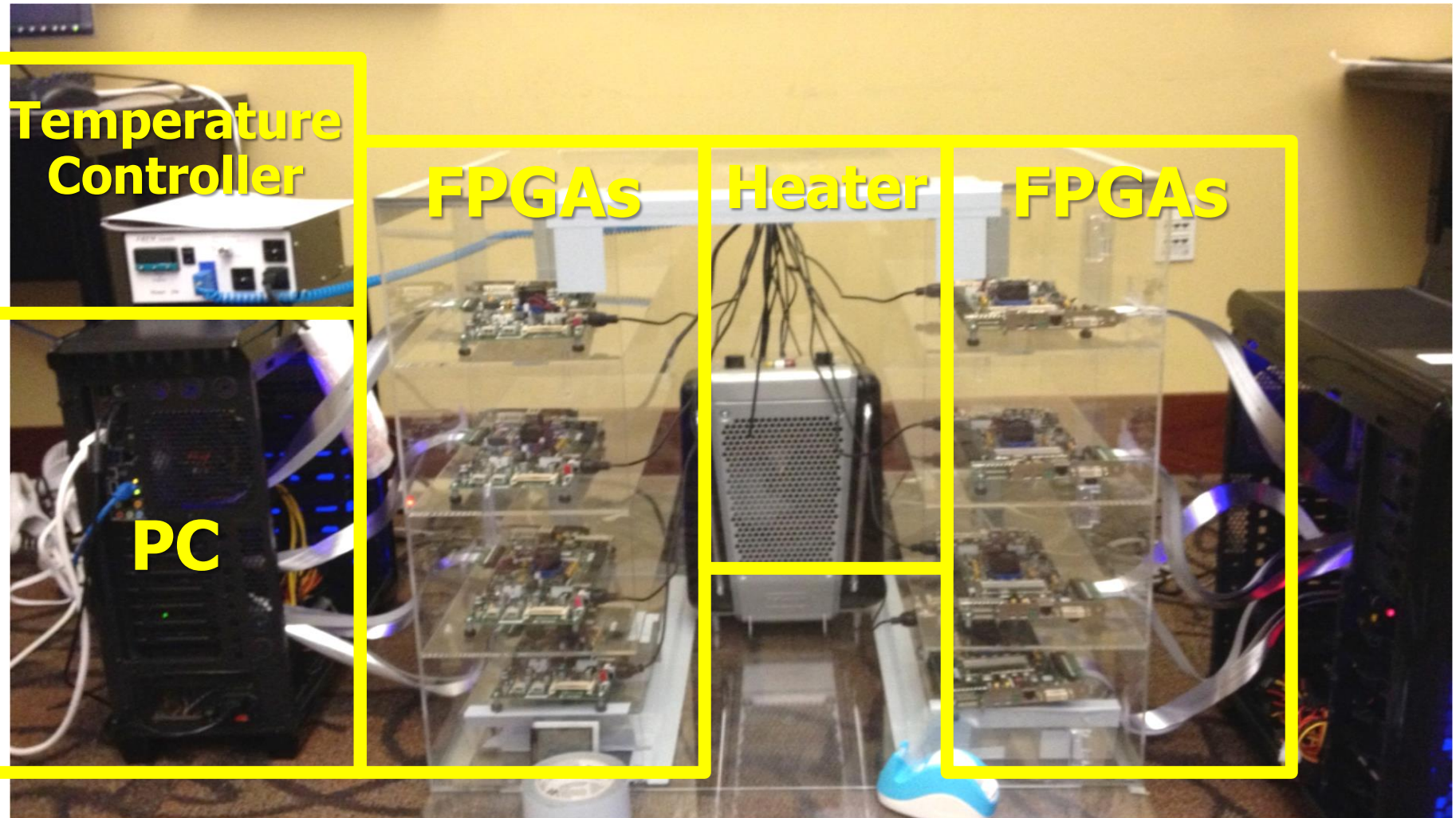
→ Lower restoration latency

## 3. Precharge

No need to fully precharge bitlines for cells with extra charge

→ Lower precharge latency

# DRAM Characterization Infrastructure

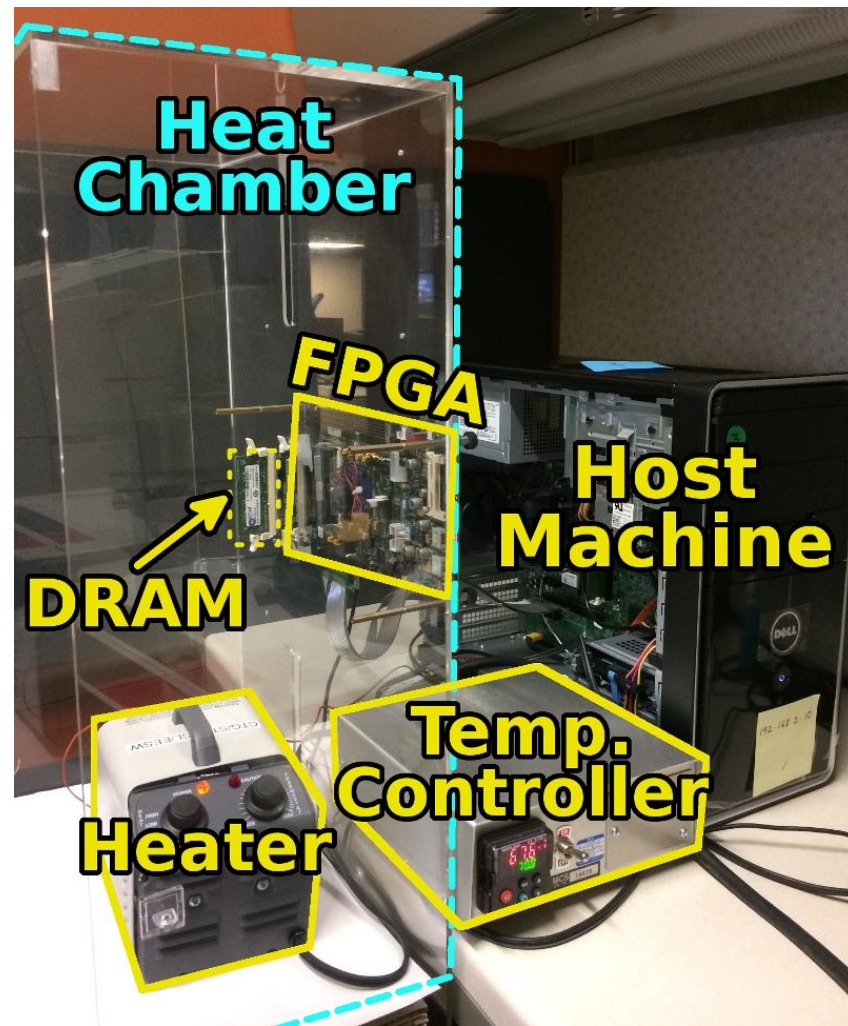




# DRAM Characterization Infrastructure

- 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](https://github.com/CMU-SAFARI/SoftMC)



# SoftMC: Open Source DRAM Infrastructure

---

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

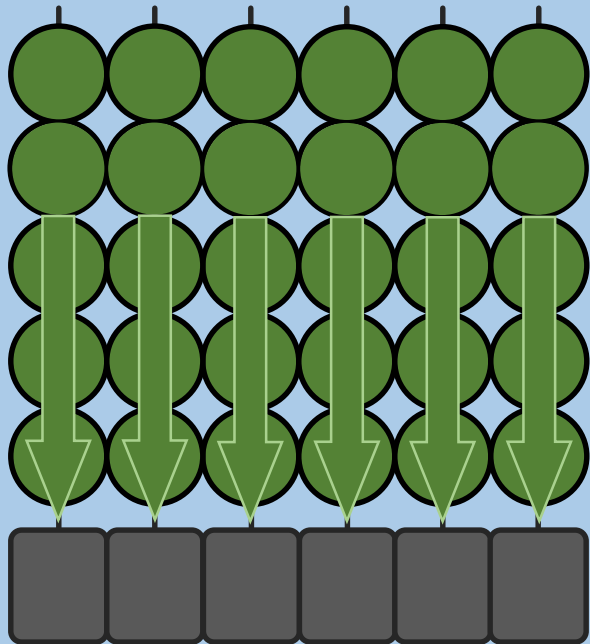
## **SoftMC: A Flexible and Practical Open-Source Infrastructure for Enabling Experimental DRAM Studies**

Hasan Hassan<sup>1,2,3</sup> Nandita Vijaykumar<sup>3</sup> Samira Khan<sup>4,3</sup> Saugata Ghose<sup>3</sup> Kevin Chang<sup>3</sup>  
Gennady Pekhimenko<sup>5,3</sup> Donghyuk Lee<sup>6,3</sup> Oguz Ergin<sup>2</sup> Onur Mutlu<sup>1,3</sup>

<sup>1</sup>*ETH Zürich*   <sup>2</sup>*TOBB University of Economics & Technology*   <sup>3</sup>*Carnegie Mellon University*  
<sup>4</sup>*University of Virginia*   <sup>5</sup>*Microsoft Research*   <sup>6</sup>*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% ↓

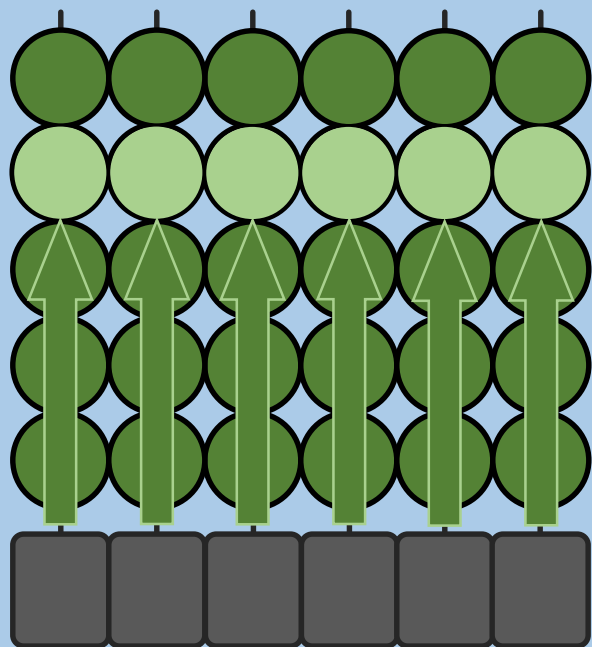
No Errors

*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)*

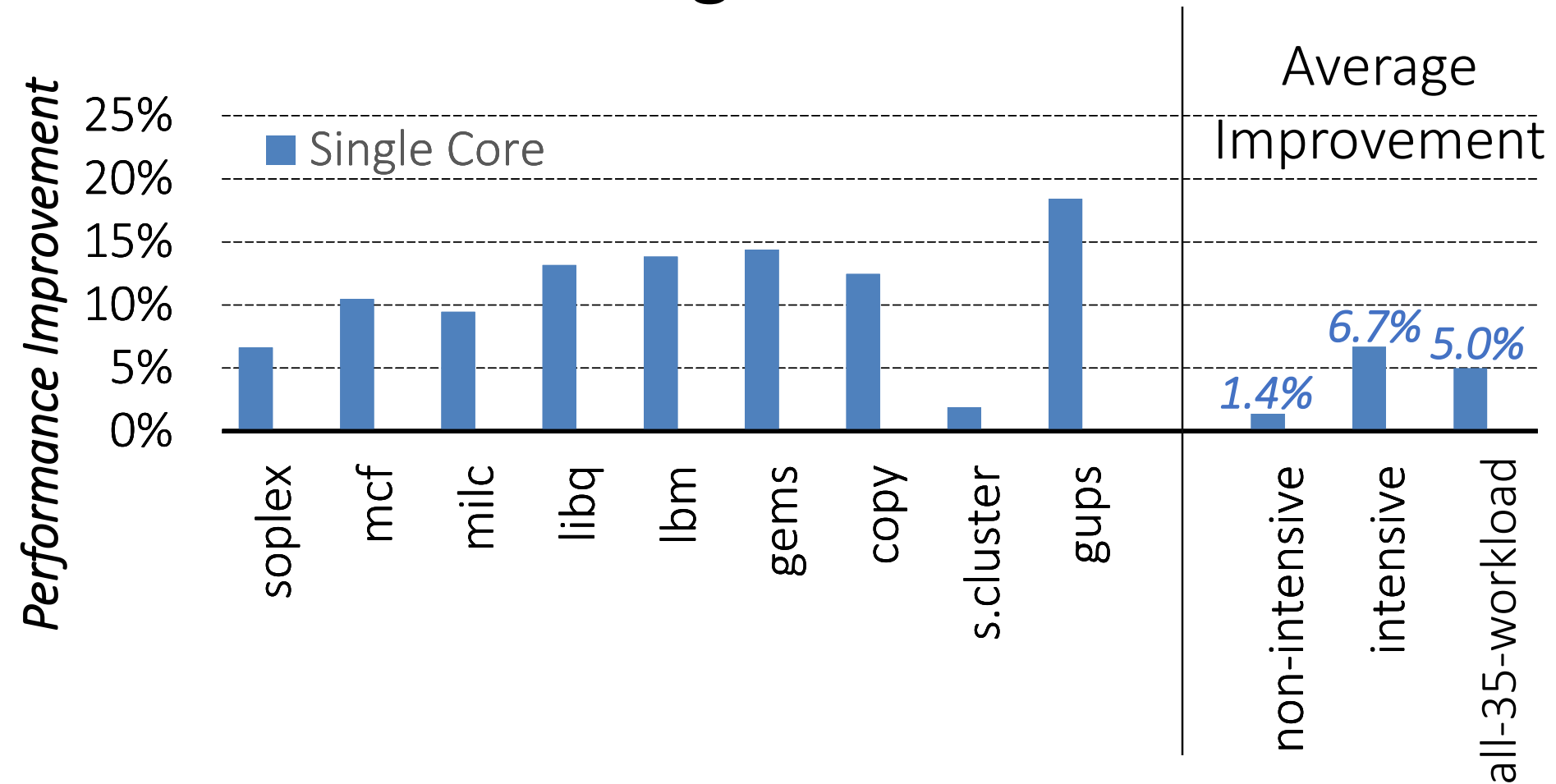
## D18F2x200\_dct[0]\_mp[1:0] DDR3 DRAM Timing 0

Reset: 0F05\_0505h. See [2.9.3 \[DCT Configuration Registers\]](#).

Bits	Description								
31:30	Reserved.								
29:24	<b>Tras: row active strobe.</b> Read-write. BIOS: See <a href="#">2.9.7.5 [SPD ROM-Based Configuration]</a> . 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>&lt;Tras&gt; 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	<b>Trp: row precharge time.</b> Read-write. BIOS: See <a href="#">2.9.7.5 [SPD ROM-Based Configuration]</a> . 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.								

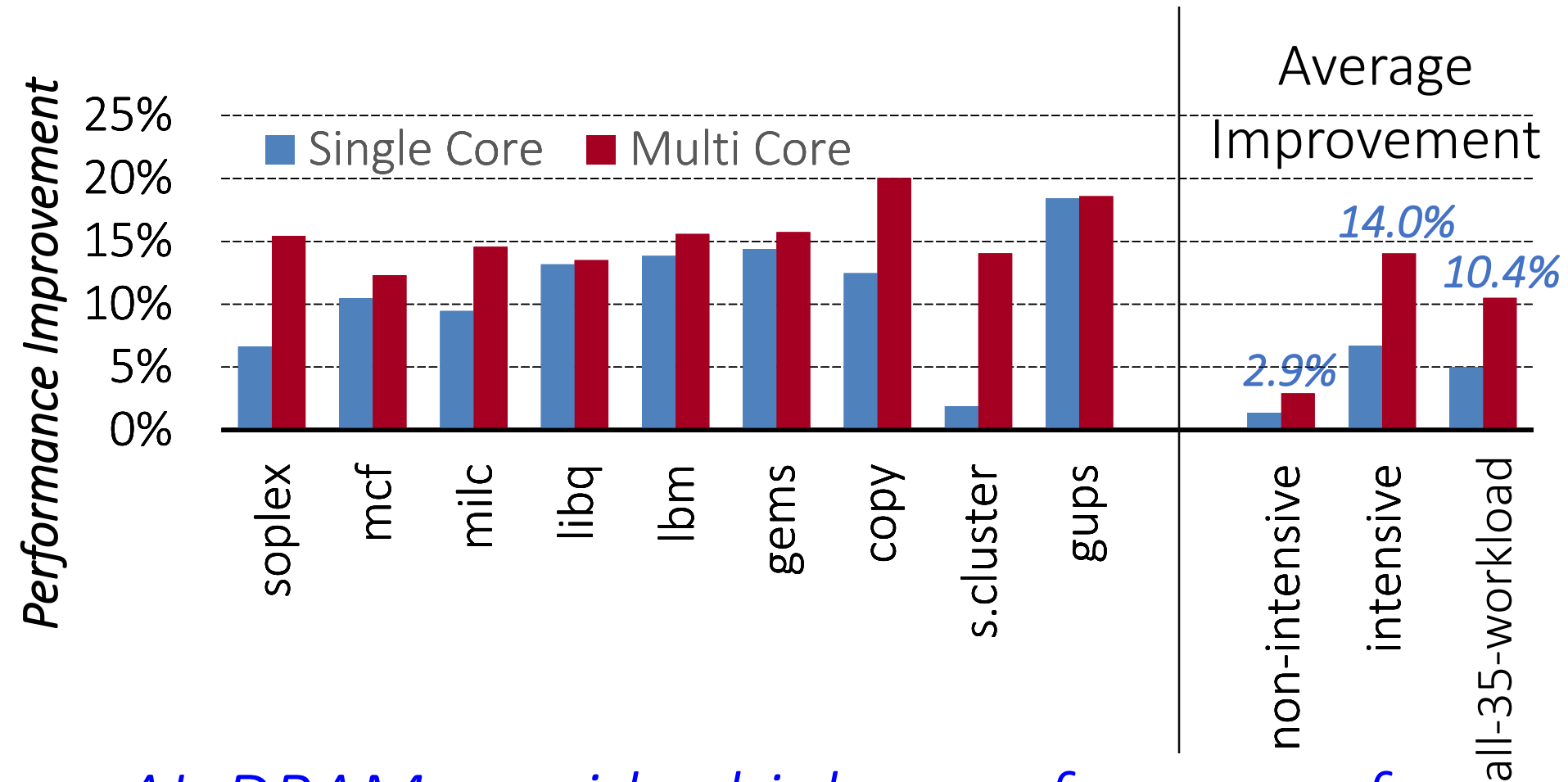


# AL-DRAM: Single-Core Evaluation



*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

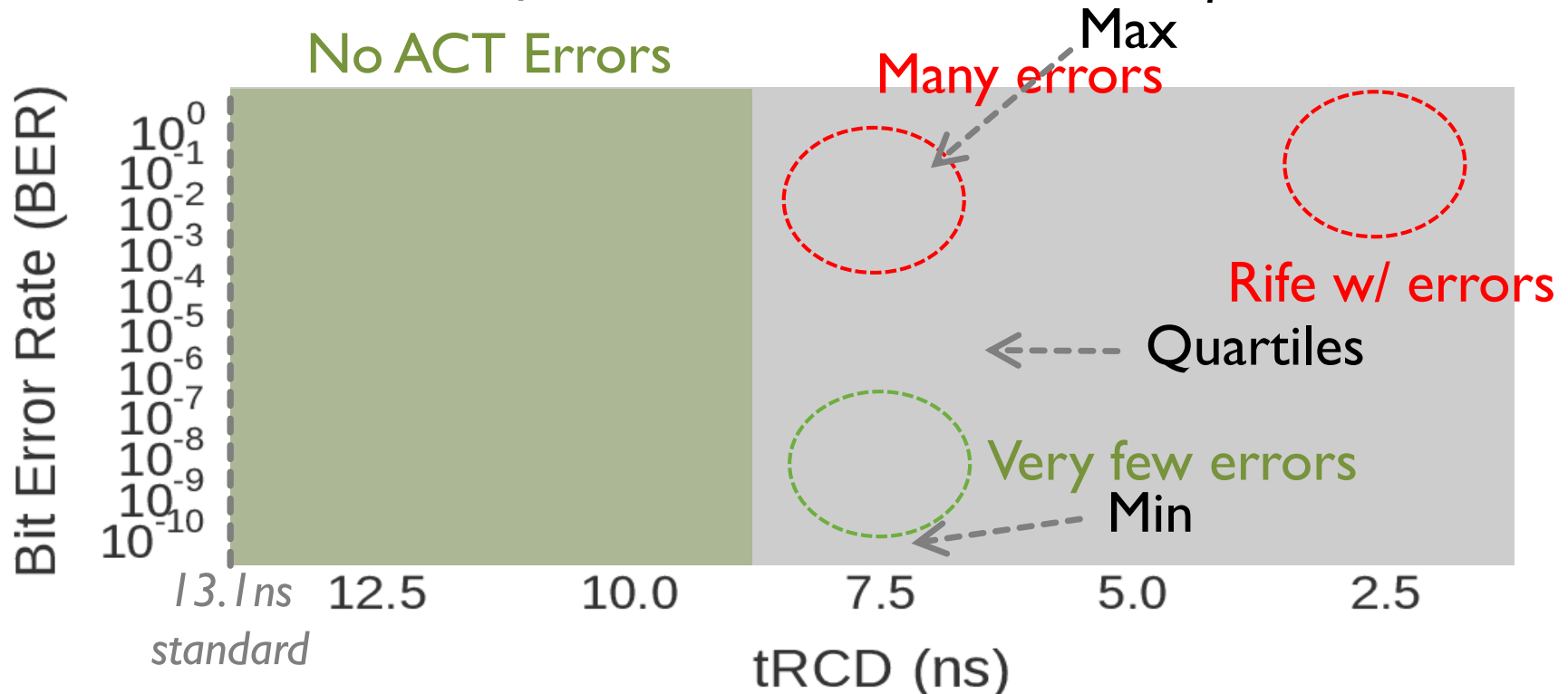
# Different Types of Latency Variation

---

- 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

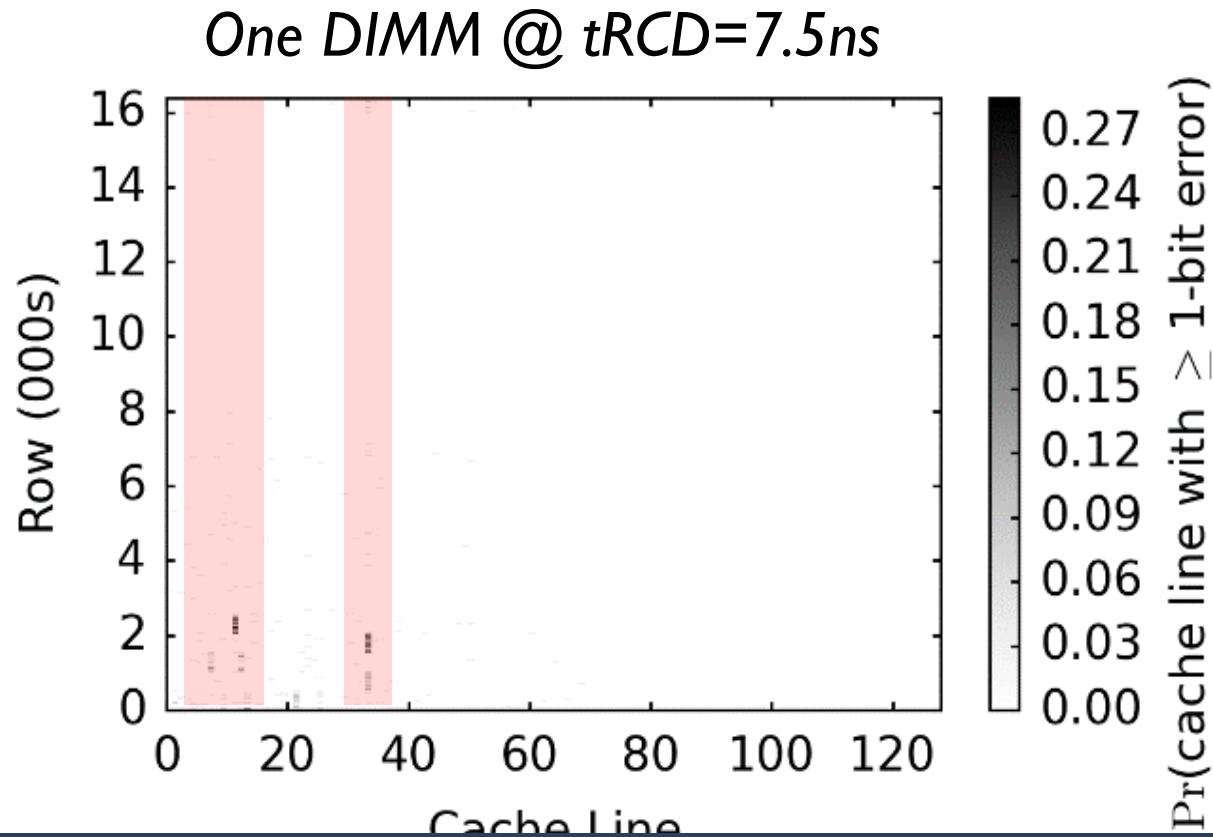
# Variation in Activation Errors

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**

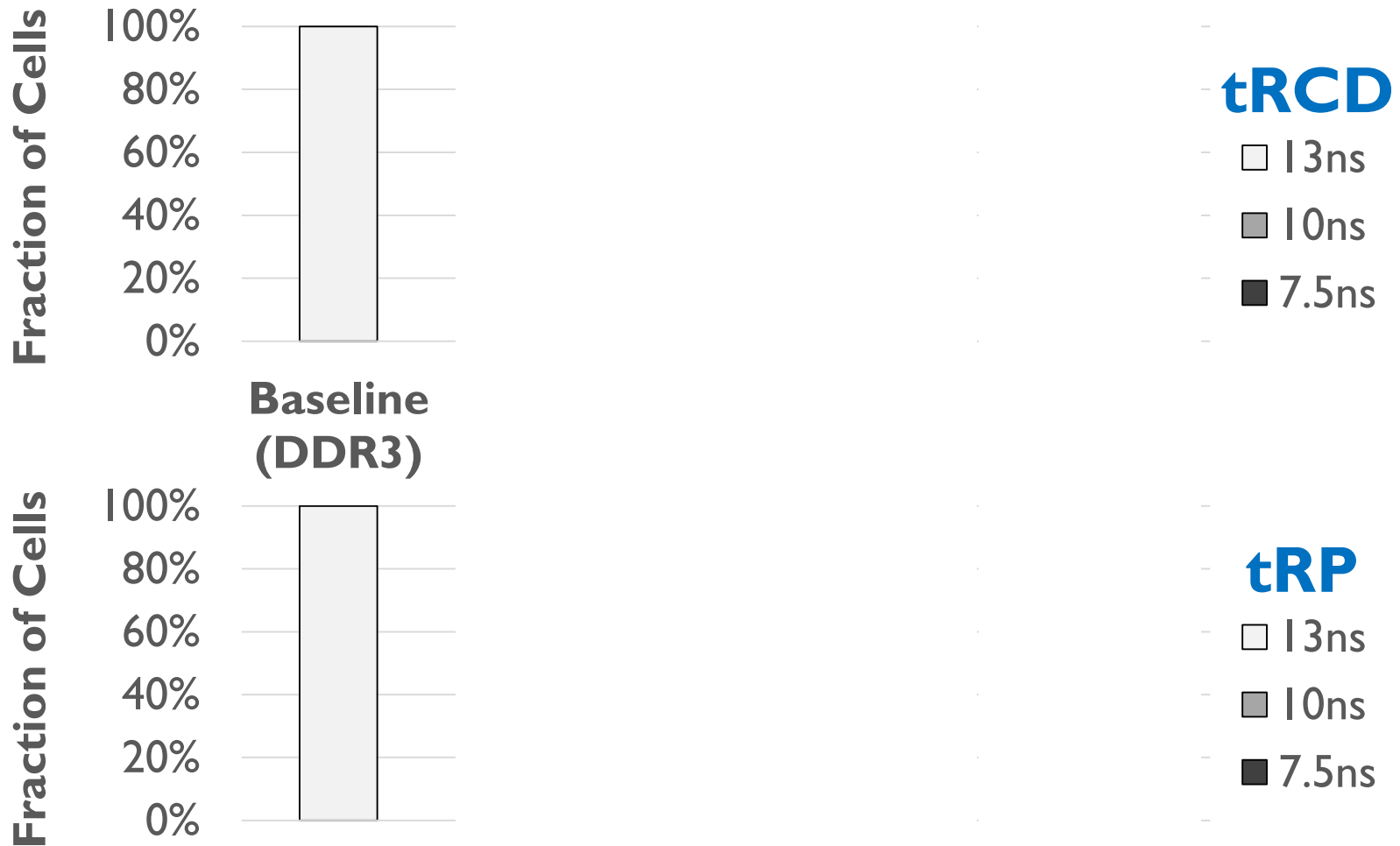


# Mechanism to Reduce DRAM Latency

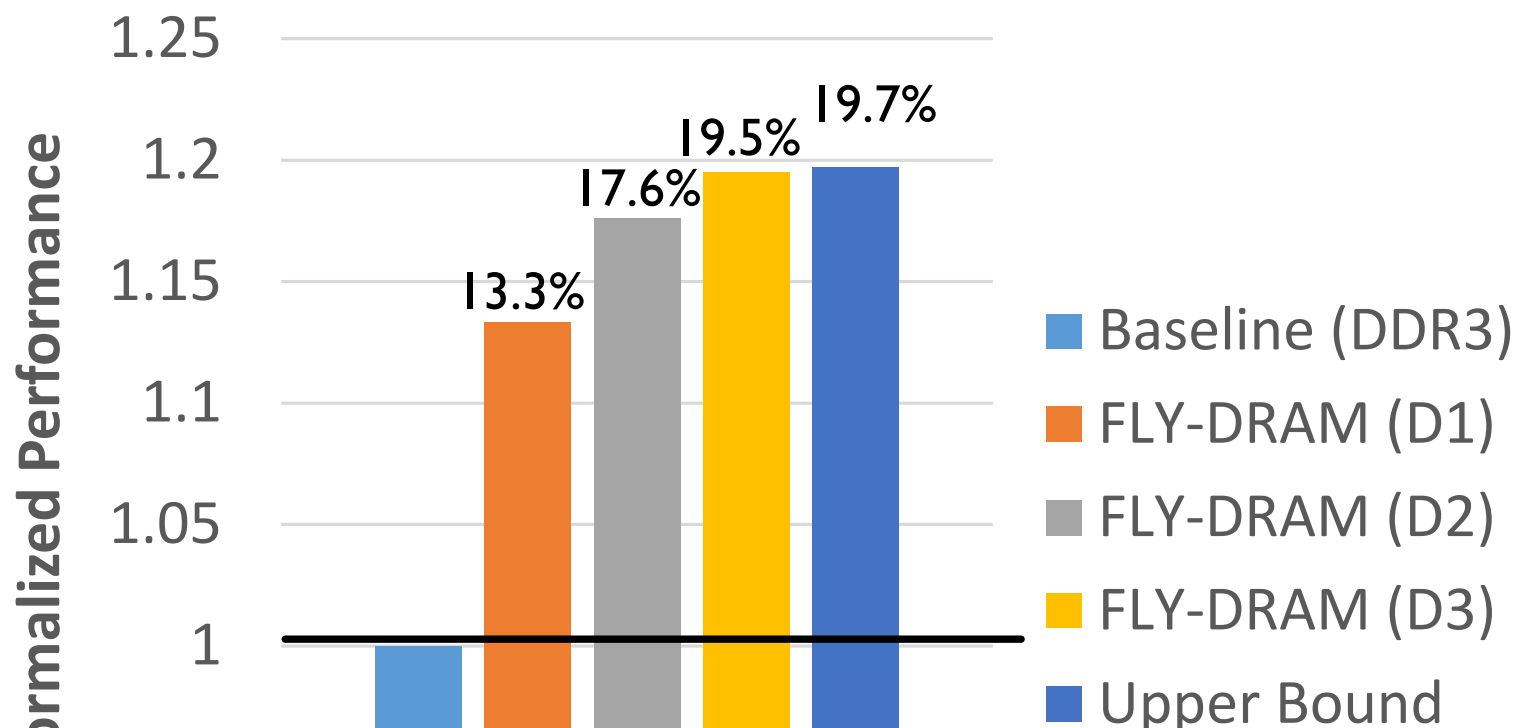
---

- **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. Chang<sup>1</sup>

Abhijith Kashyap<sup>1</sup>

Hasan Hassan<sup>1,2</sup>

Saugata Ghose<sup>1</sup>

Kevin Hsieh<sup>1</sup>

Donghyuk Lee<sup>1</sup>

Tianshi Li<sup>1,3</sup>

Gennady Pekhimenko<sup>1</sup>

Samira Khan<sup>4</sup>

Onur Mutlu<sup>5,1</sup>

<sup>1</sup>Carnegie Mellon University   <sup>2</sup>TOBB ETÜ   <sup>3</sup>Peking University   <sup>4</sup>University of Virginia   <sup>5</sup>ETH Zürich

# Computer Architecture

## Lecture 10b: Memory Latency

Prof. Onur Mutlu

ETH Zürich

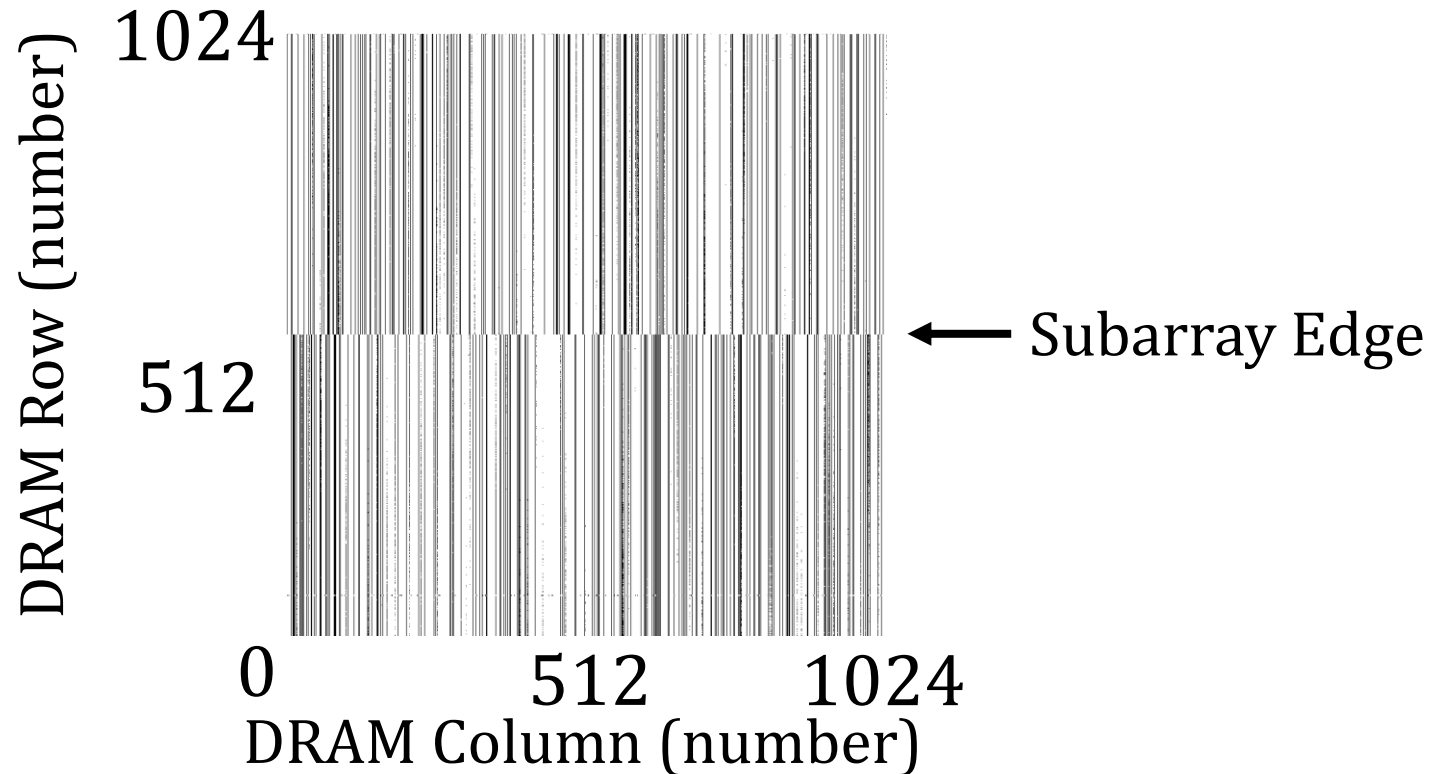
Fall 2018

18 October 2018

We did not cover the following slides in lecture.  
These are for your benefit.

# Spatial Distribution of Failures

How are activation failures spatially distributed in DRAM?

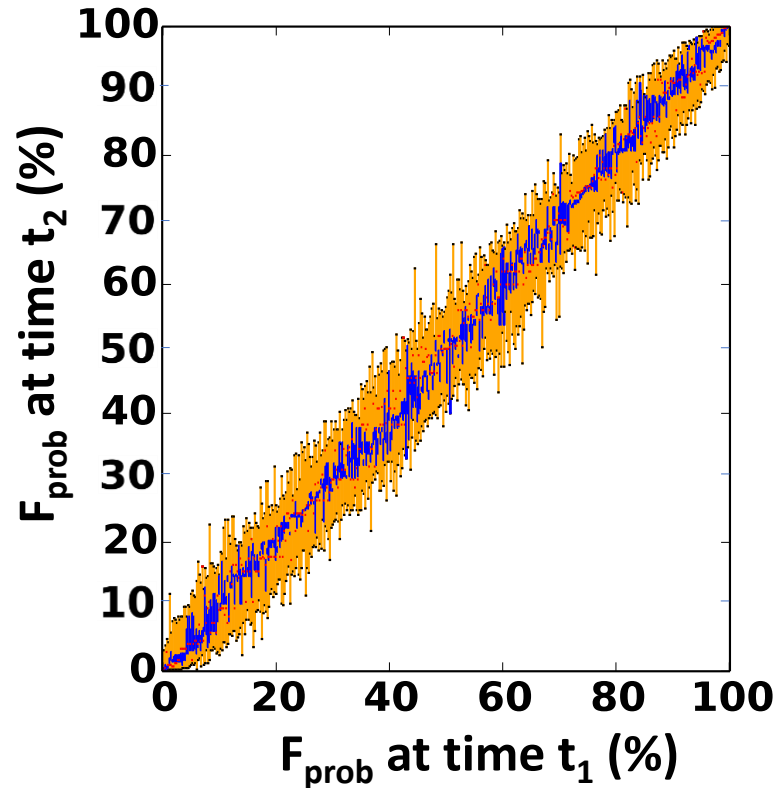


Activation failures are **highly constrained**  
to local bitlines



# Short-term Variation

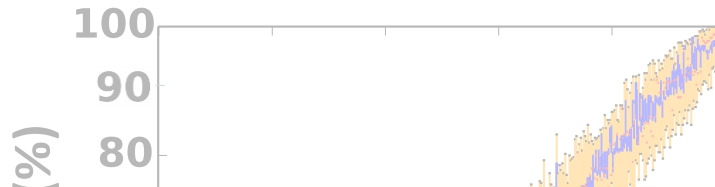
Does a bitline's probability of failure change over time?



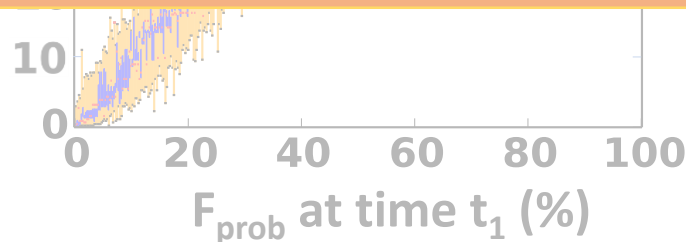
A **weak bitline** is likely to remain **weak** and  
a **strong bitline** is likely to remain **strong** over time 97

# 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

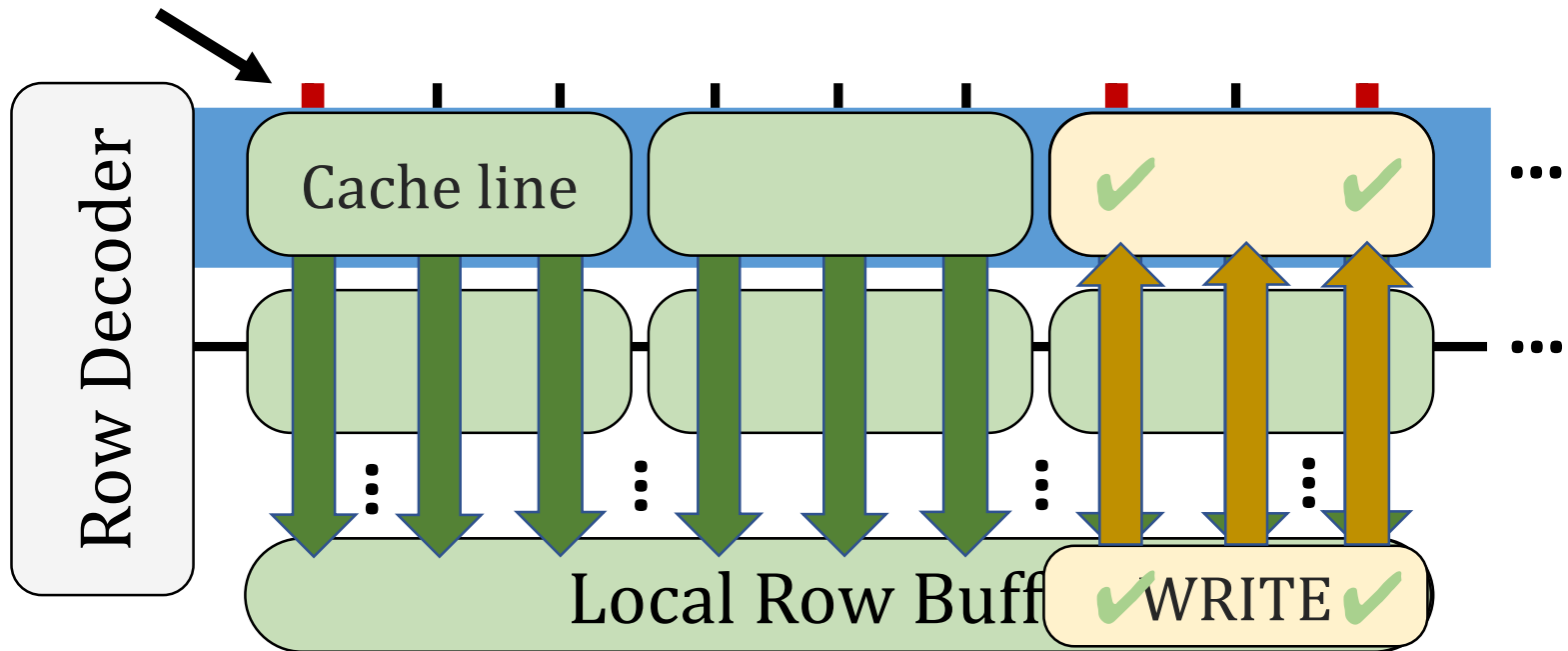


A **weak bitline** is likely to remain **weak** and a **strong bitline** is likely to remain **strong** over time 98

# Write Operations

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

**Weak bitline**



We can reliably issue write operations  
with significantly reduced  $t_{\text{RCD}}$  (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

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

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

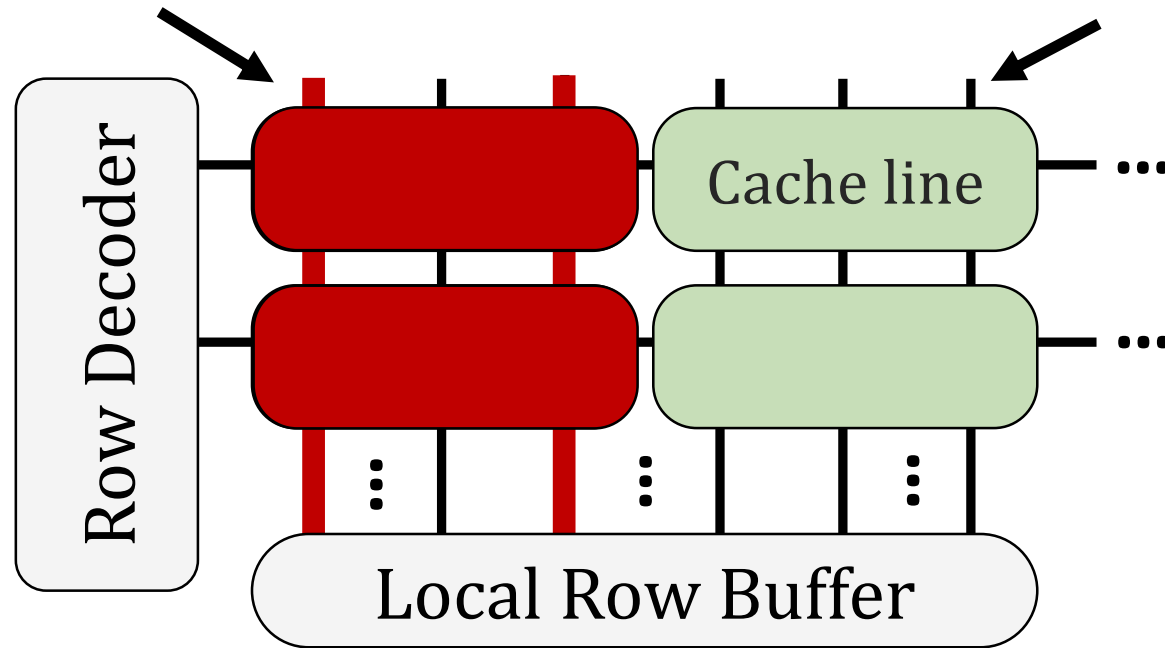
## Three Components

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**

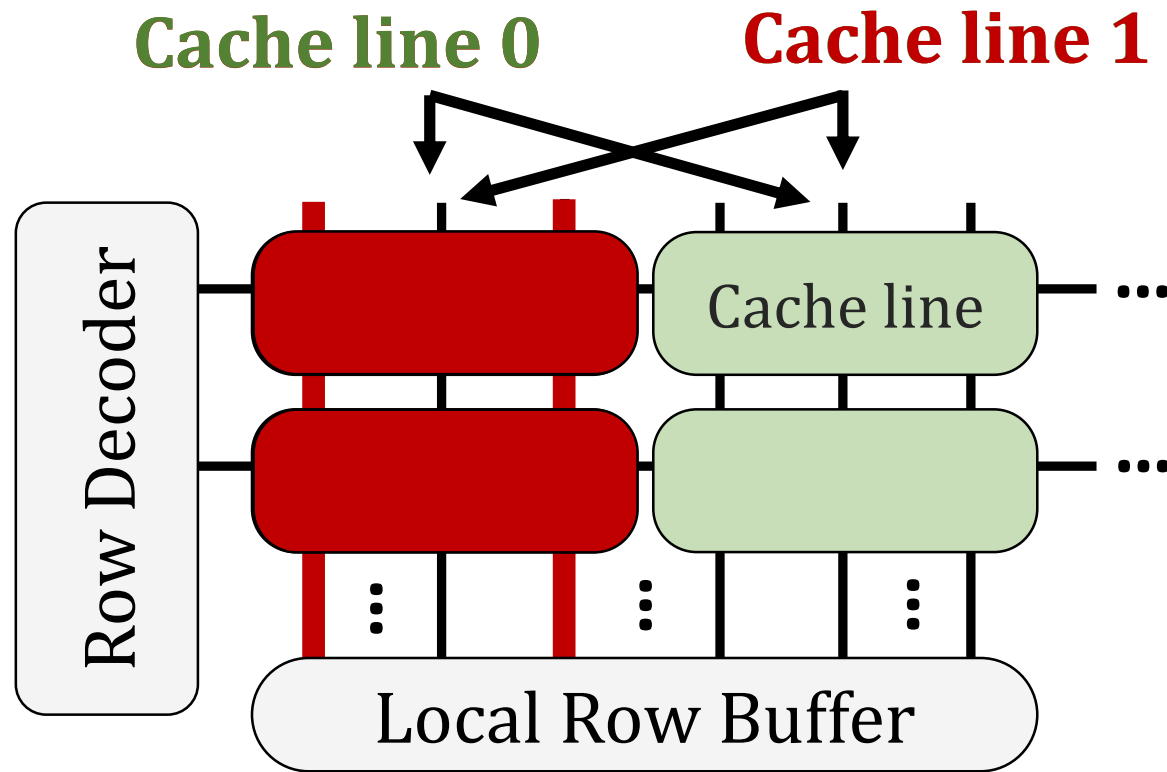
- Identifies subarray columns as weak or strong
- Obtained in a one-time profiling step

## 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**

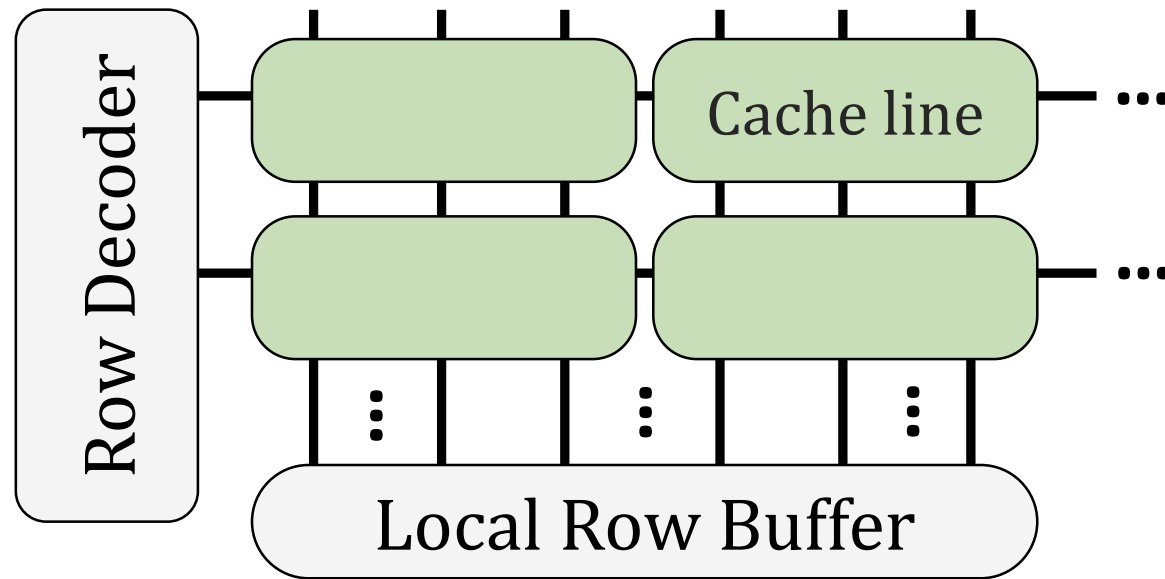
- Identifies subarray columns as weak or strong
- Obtained in a one-time profiling step

## 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  $t_{\text{BCP}}$  (e.g., by 77%)

# More on Solar-DRAM

---

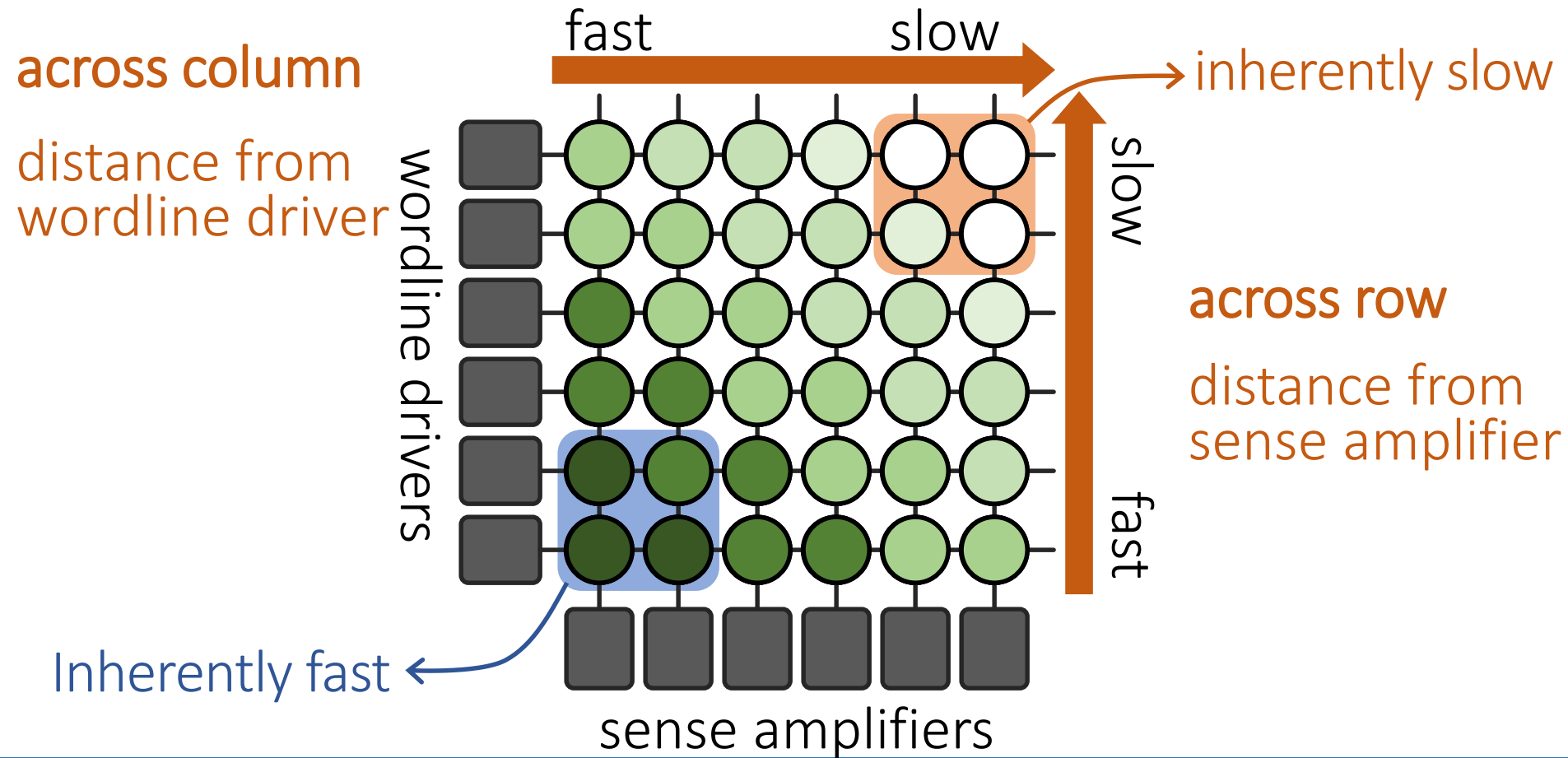
- Jeremie S. Kim, Minesh Patel, Hasan Hassan, and Onur Mutlu,  
**"Solar-DRAM: Reducing DRAM Access Latency by Exploiting the Variation in Local Bitlines"**  
*Proceedings of the 36th IEEE International Conference on Computer Design (ICCD)*, Orlando, FL, USA, October 2018.

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

Jeremie S. Kim<sup>‡§</sup>      Minesh Patel<sup>§</sup>      Hasan Hassan<sup>§</sup>      Onur Mutlu<sup>§‡</sup>  
                         ‡Carnegie Mellon University                           §ETH Zürich

# Why Is There Spatial Latency Variation Within a Chip?

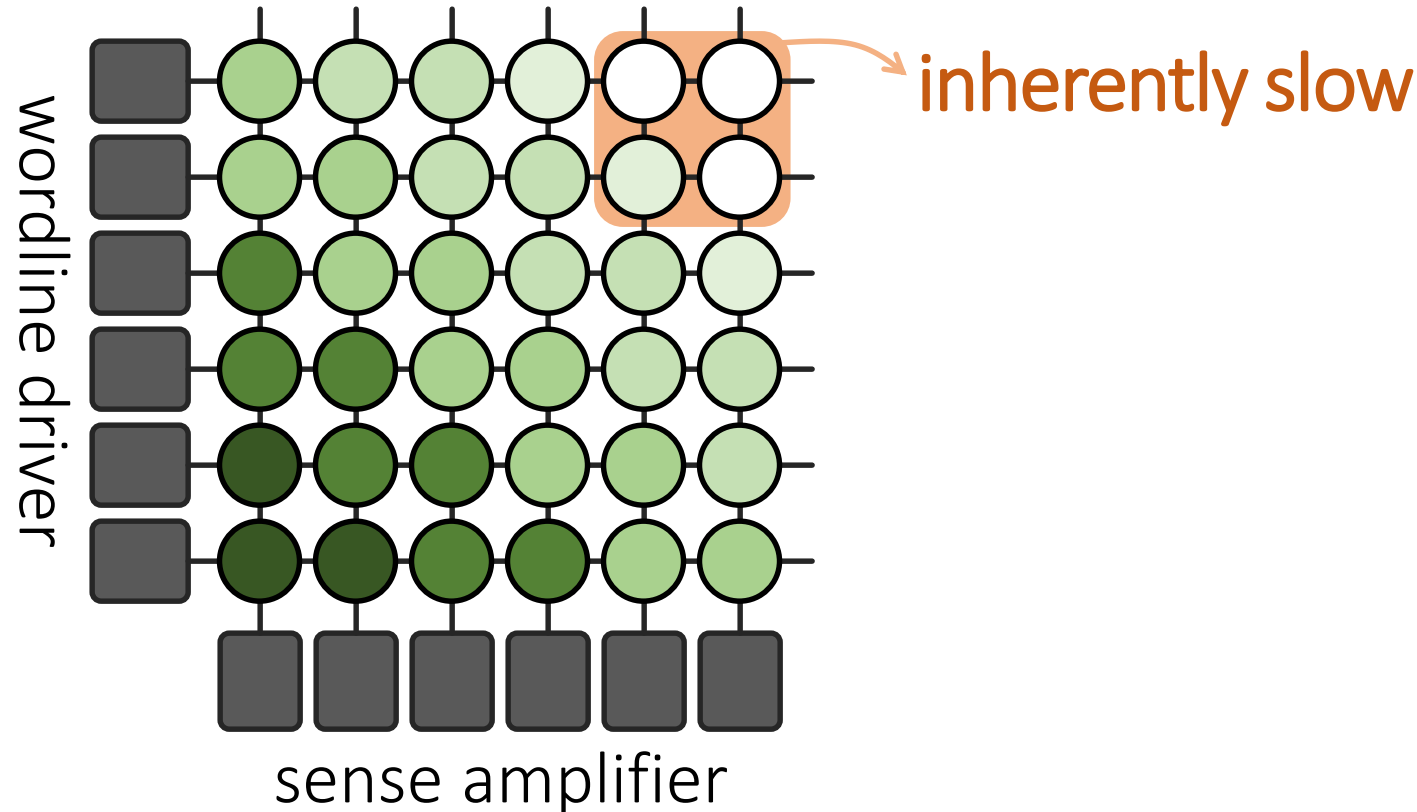
# What Is Design-Induced Variation?



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

# DIVA Online Profiling

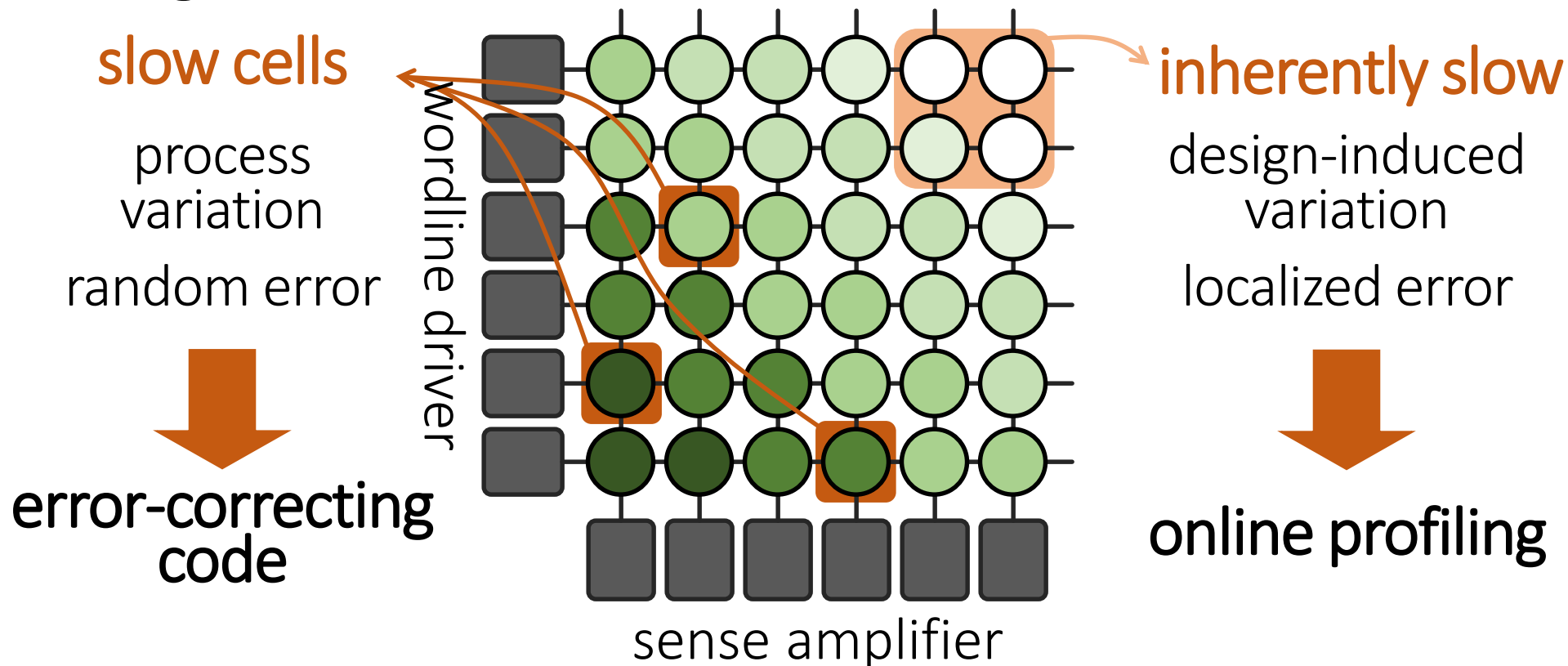
Design-Induced-Variation-Aware



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

# DIVA Online Profiling

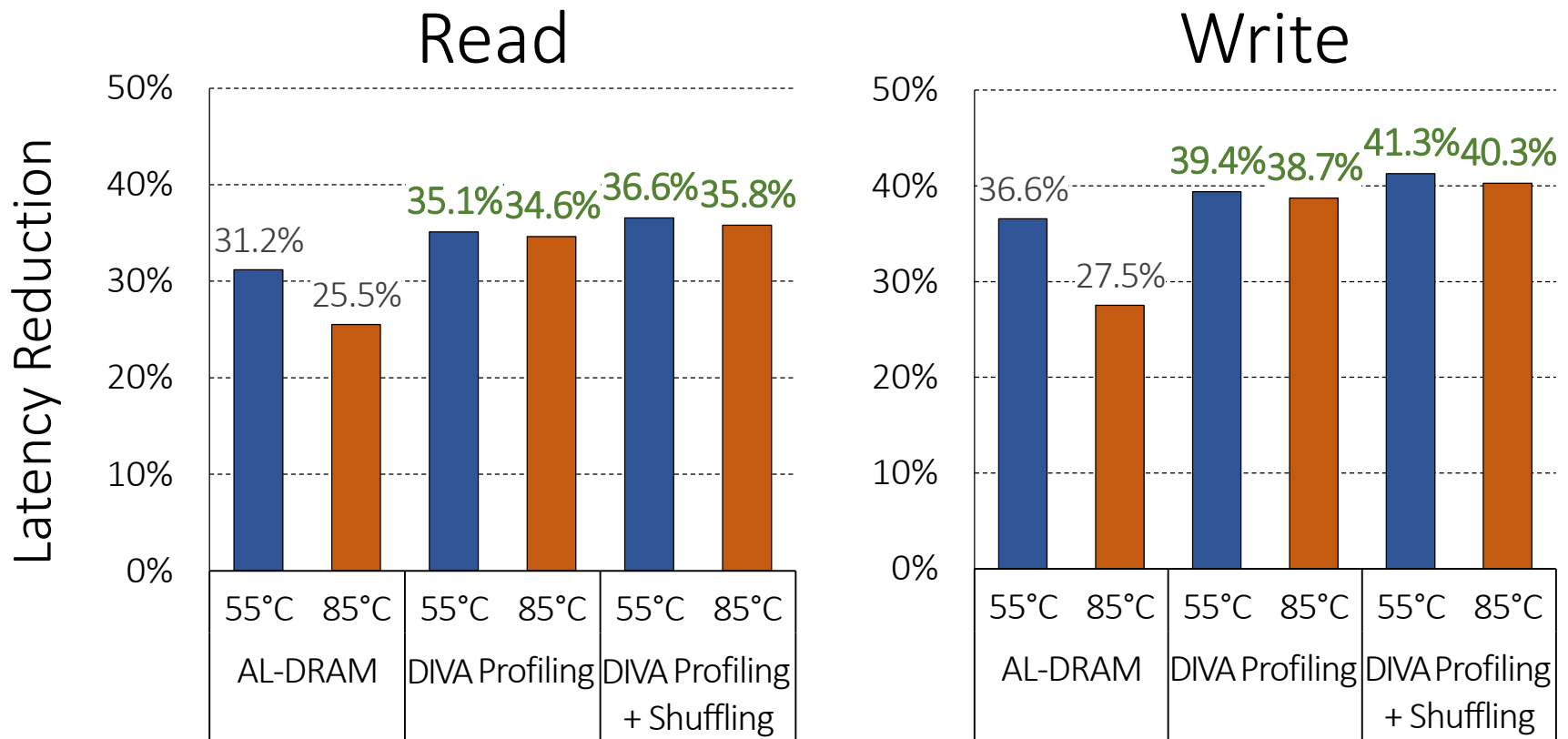
Design-Induced-Variation-Aware



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