# Reshaping DRAM Scaling
# by Enabling System-Memory Cooperation
## *Computer Architecture - Guest Lecture 11(h)*

## Minesh Patel

*SAFARI* **ETH** *zürich*
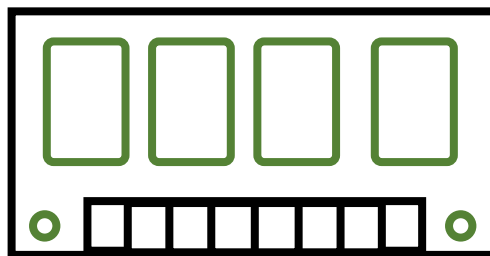
Zürich • 2 November 2023

# Goals of This Talk

1. Reflecting on how we **build and use** memory

2. Considering the **indirect costs** of current design practices

3. Improving how we **address unanticipated problems**

# Executive Summary

- **Problem:** overcoming DRAM scaling challenges requires **new solutions**
- **Observation:** the separation of concerns between DRAM producers and consumers is a barrier to overcoming scaling
  1. **Too rigid** to adapt to unexpected challenges (e.g., Rowhammer)
  2. **Discourages** new solutions based on system-memory cooperation
- **Key idea:** revise the separation to **encourage** new solutions
- **Four case studies**: performance, energy-efficiency, reliability, security
  - We identify **memory testing** as the primary culprit for discouraging new solutions
- **Approach:** two-step plan to revise DRAM standards
  1. **Near-term:** crowdsourcing and publications
  2. **Long-term:** changes to industry-wide DRAM standards
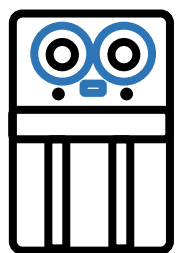
# Talk Outline

- **DRAM Scaling Challenges**
- Addressing Scaling: The Separation of Concerns
  - Problem 1: Inflexibility to Challenges
  - Problem 2: Overly Constraining
- Enabling System-Memory Cooperation
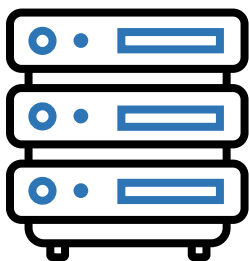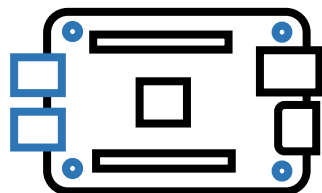- Revising the Separation

# DRAM
*(Dynamic Random Access Memory)*

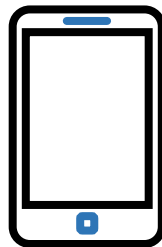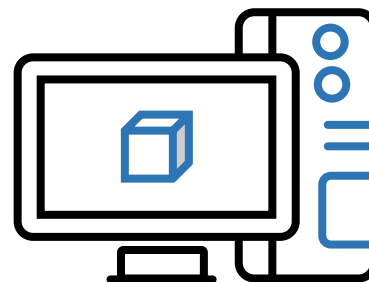Mainframe    Server    Embedded    Mobile    Desktop    Secure    Emerging
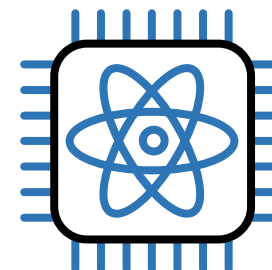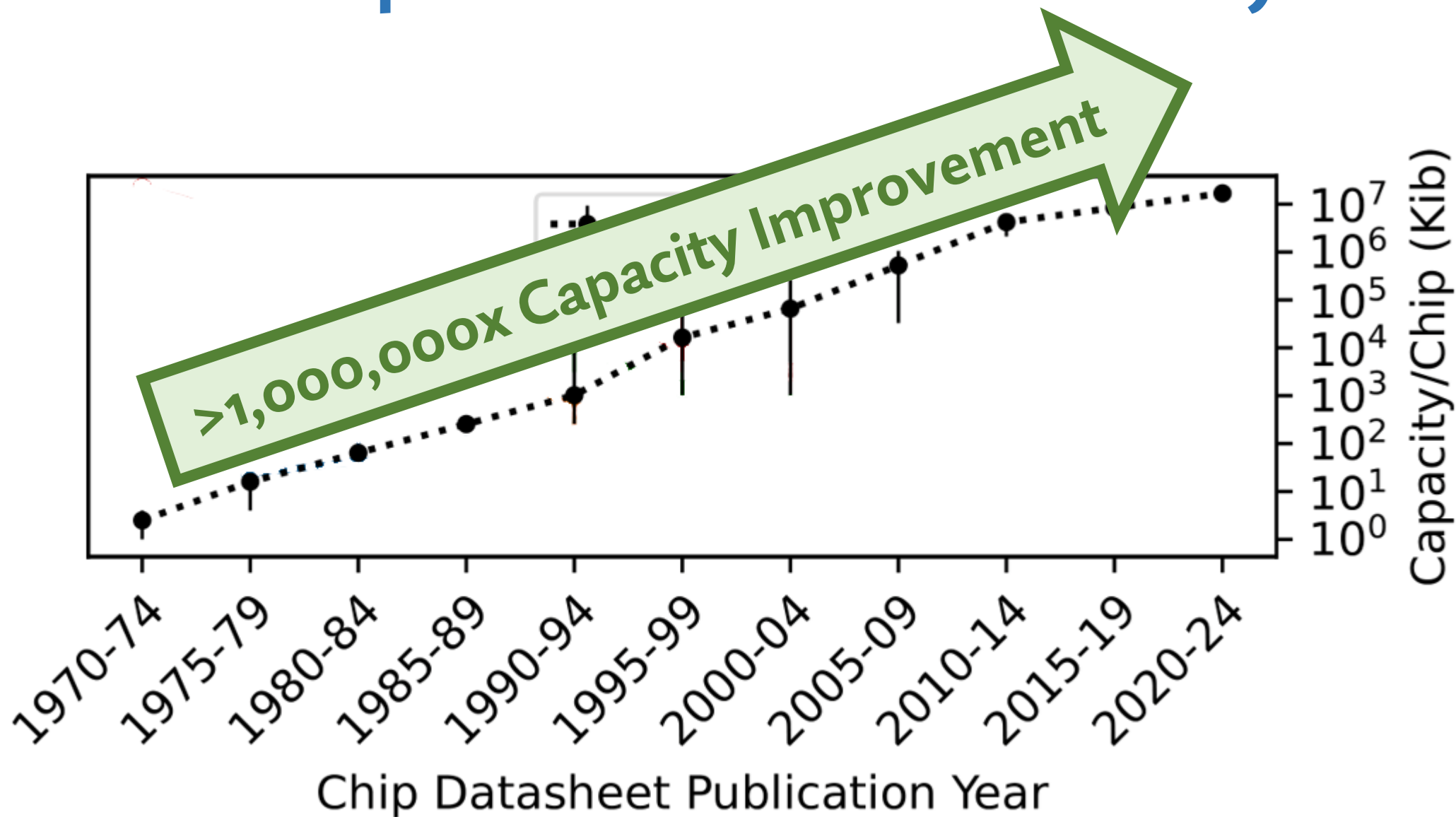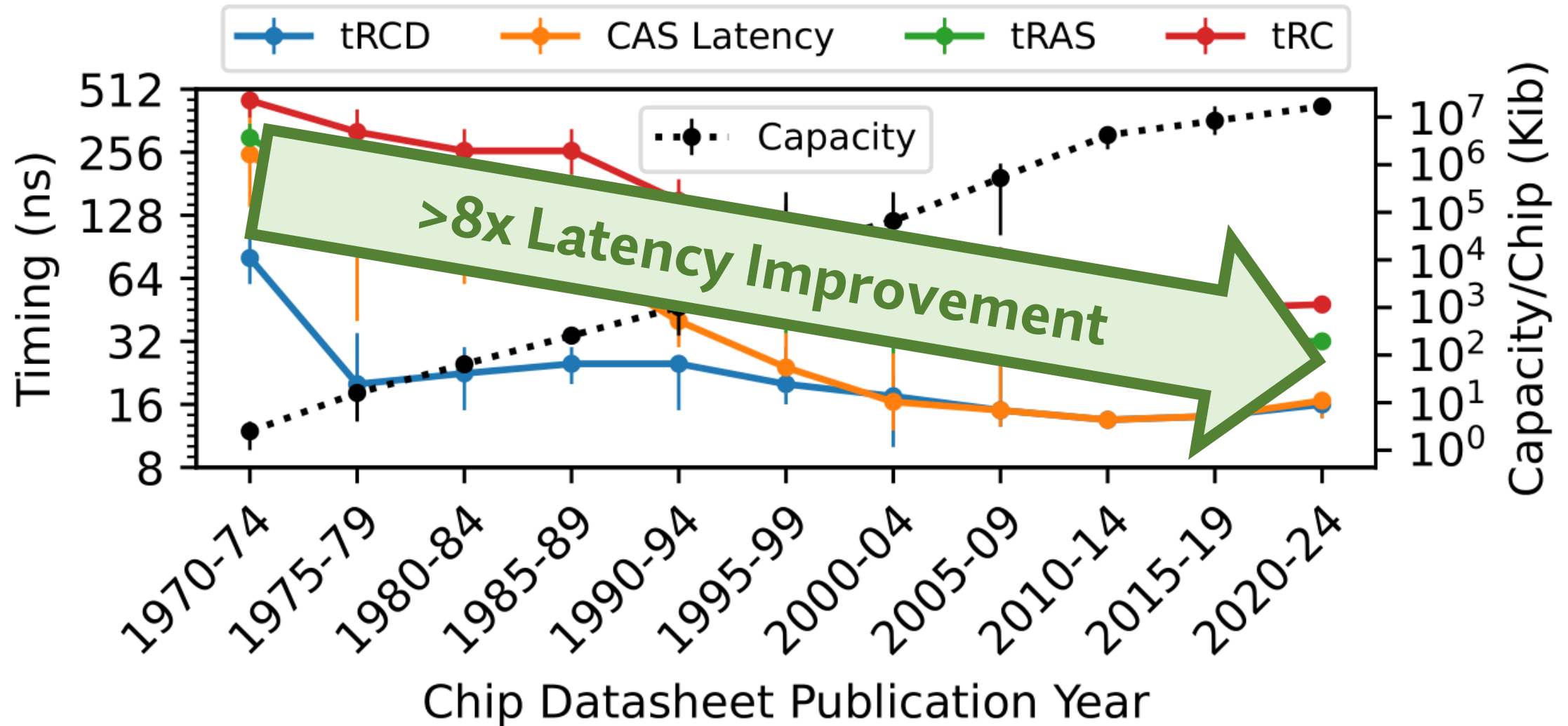
~1970                                                                      2020+

**56.4%** of the global memory market *[Yole Développement, 2022]*

# Generational Improvements to Commodity DRAM



>1,000,000x Capacity Improvement

# Generational Improvements to Commodity DRAM

# Generational Improvements to Commodity DRAM

# DRAM Scaling

- Scaling has traditionally been driven by **DRAM manufacturers**
  - Building efficient DRAM requires **specialized** technologies



**Tightly-Packed DRAM Cell Array**

**1 Mb**
*DIP, 1985*

**512 Mb**
*DDR2 DIMM, 2000*

**32 Gb**
*LPDDR4 FBGA, 2020*

**Base Technology** → *New circuits, materials, fabrication methods, etc.*

# DRAM Errors

Memory Errors

Particle Strike

RowHammer Errors

Electrical Noise

Charge Leakage

**Bit-Level Storage Array**

DRAM suffers from **errors** that cause **data loss** or **system failure** if ignored

# DRAM Scaling Challenges

- Scaling becomes **more difficult** with higher storage density



**1 Mb**
*DIP, 1985*

**512 Mb**
*DDR2 DIMM, 2000*

**32 Gb**
*LPDDR4 FBGA, 2020*

Smaller, denser cells are **less reliable**

Continued scaling is **expensive** due to
the **overheads** of maintaining reliable operation

# Error Mitigation for Further Scaling

**Low**

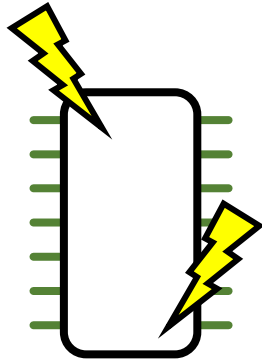No Mitigation Needed

Row/Col Redundancy

Error Correcting Codes

Active (TRR, (D)RFM)

**High**  **???????????????**

**Error Rate**

**$$$$
Is it worth it?**

*spare row*

| data |
|------|

| data | metadata |
|------|----------|

*expanded representation
(resilient to errors)*

# Key DRAM Scaling Challenges



Deeply-Scaled Unreliable Cells

**Performance**
Reducing the long DRAM access latency

**Efficiency**
Improving refresh power and performance

**Reliability**
Mitigating worsening memory errors

**Security**
Addressing worsening RowHammer vulnerability

# Generational Improvements to Commodity DRAM



**Improvements driven by DRAM manufacturers**

**Slow progress**

We **cannot rely** on manufacturers alone to **overcome** DRAM scaling challenges

# Talk Outline

- DRAM Scaling Challenges
- **Addressing Scaling: The Separation of Concerns**
  - Problem 1: Inflexibility to Challenges
  - Problem 2: Overly Constraining
- Enabling System-Memory Cooperation
- Revising the Separation

# The Producer-Consumer Relationship

**DRAM Producers**

**DRAM Consumers**



**DRAM manufacturers**

*e.g., Samsung, SK Hynix, Micron Technologies*

**System Design/Test/Research**

*e.g., Board designers, test engineers, research scientists*

# Producer and Consumer Responsibilities

**DRAM Producers**

**DRAM Consumers**

**Separation of Concerns**

**DRAM Design**

**DRAM Use**

☑ Specialized roles enable **highly-optimized DRAM chips**

☑ Interoperability enables **widespread use** of DRAM

☑ **Preserves trade secrets** among producers and consumers

# Producer and Consumer Responsibilities

**DRAM Producers**

DRAM Design

**DRAM Consumers**

DRAM Use

**Limited solution space for DRAM scaling**

# Producer and Consumer Responsibilities

**DRAM Producers**

**DRAM Consumers**

Separation of Concerns

**DRAM Design**

**DRAM Use**

**Barrier to addressing scaling challenges**

# Observations

- Two major problems **inherent** to the existing separation

**1** The separation is **inflexible** when new challenges occur
*e.g., RowHammer, worsening memory errors*

**2** The separation **constrains the solution space** available to address those challenges

# Talk Outline

- DRAM Scaling Challenges
- Addressing Scaling: The Separation of Concerns
  - **Problem 1: Inflexibility to Challenges**
  - Problem 2: Overly Constraining
- Enabling System-Memory Cooperation
- Revising the Separation

# DRAM Standards

- Separation of concerns is implemented by **DRAM standards**



**DRAM Producers**

**DRAM Standards**

**DRAM Consumers**

**Separation of Concerns**

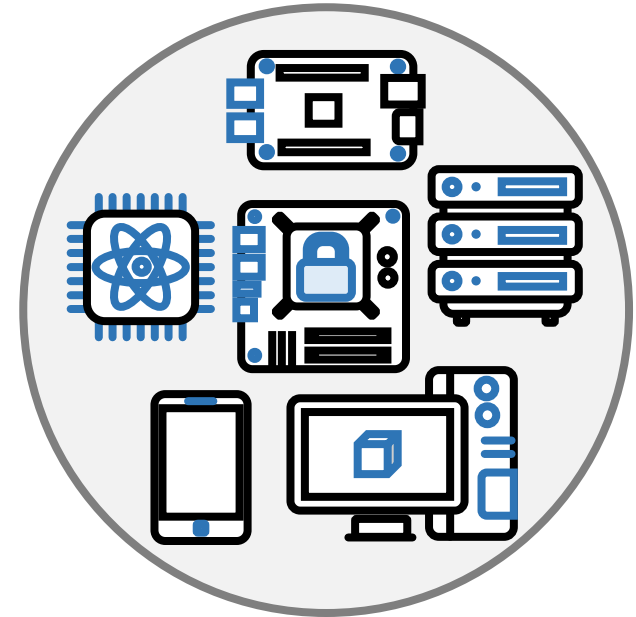# The Evolution of Standards Over Time

- Standards govern **consumer-visible properties** of DRAM chips
  - Interface, configuration, performance characteristics
  - Abstract DRAM design details away from consumers



The exponential fit shown is:
$$0.08e^{0.38(x-1998)} + 143.68$$

# Inflexibility of Standards

- Unfortunately, DRAM standards are **slow to adapt to change**
  - Requires **industry-wide consensus** among producers and consumers

*RowHammer Security Vulnerability*

| | | | |
|---|---|---|---|
| | Target Row Refresh (HBM 2) | RFM (DDR5) | DRFM (DDR5) |

RowHammer Discovered

"RowHammer-Free" Chips Announced

JEDEC RowHammer Task Group Formed

JEDEC RowHammer Recommendations (JEP30{0,1}-1) Released

**2014**  **2017**  **2020**  **2021**  **2022**

**No Complete Solution**

**???**

ISSCC Proposals

Used in LPDDR4

Discussed in HBM3 Standard

ISSCC Papers

Partial Disclosure (DDR5)

*On-Die Error Correcting Codes (ECC)*

# Talk Outline

- DRAM Scaling Challenges
- Addressing Scaling: The Separation of Concerns
  - Problem 1: Inflexibility to Challenges
  - **Problem 2: Overly Constraining**
- Enabling System-Memory Cooperation
- Revising the Separation

# System-Memory Cooperation

**Highly-efficient solutions** based on
a holistic understanding of how scaling impacts the system

↑

**Addressing
Scaling Challenges**



**DRAM**

**Systems**

# Example System-Memory Cooperative Solutions

| To improve… | Consumers can… |
| --- | --- |
| Performance Energy & Power | Exploit slack in operating conditions<br>• Access and refresh timings<br>• Operating voltage and temperature |
| Reliability | Protect against and test for failures<br>• System-level error mitigations<br>• Detailed qualification and validation |
| Security | Implement system-level defenses<br>• RowHammer, cold-boot attacks |

# System-Memory Cooperation

- **Problem:** The separation **discourages** cooperative solutions



**DRAM Producers**

Separation of Concerns

**DRAM Consumers**

*No insight into DRAM design*

*No insight into DRAM use*

# DRAM Operating Space



**Discouraged**

**Possible DRAM Operating Points**

**Standardized Operating Points**

**Cooperative Solutions**

**Encouraged Designs**
*Low TCO*

- **Producer:** fully specified design, proprietary
- **Consumer:** fully supported, predictable behavior

**Discouraged Designs**
*High TCO*

- **Producer:** "out-of-spec", unsupported
- **Consumer:** unknown behavior, out-of-warranty

# Key Idea: Encourage Cooperative Solutions

**1** Enable Cooperative Solutions

Possible DRAM Operating Points

Standardized Operating Points

Cooperative Solutions

Near-Term Solution

**2** Broaden Standards

Possible DRAM Operating Points

Standardized Operating Points

Cooperative Solutions

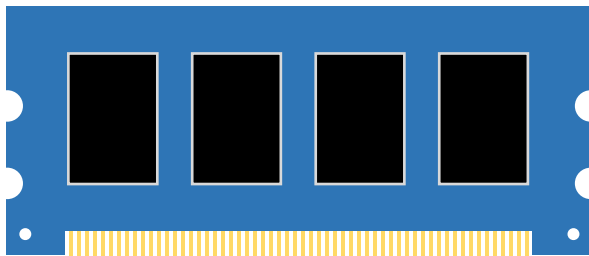**Encouraged (Low TCO)**

Long-Term Solution

# Talk Outline

- DRAM Scaling Challenges
- Addressing Scaling: The Separation of Concerns
  - Problem 1: Inflexibility to Challenges
  - Problem 2: Overly Constraining
- **Enabling System-Memory Cooperation**
- Revising the Separation

# Four Case Studies

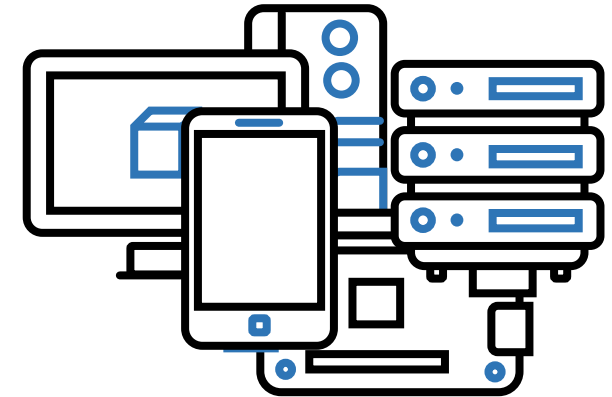- **Goal:** survey system-memory cooperative solutions to understand what holds them back from widespread adoption

**Performance**
Reducing the long DRAM access latency

**Efficiency**
Improving refresh power and performance

**Reliability**
Mitigating worsening memory errors

**Security**
Addressing worsening RowHammer vulnerability

**This Talk**

# Case Study: Mitigating DRAM Refresh Overheads

## DRAM Cell

**Data Encoding**

*access transistor*

*storage capacitor*

*stores one bit of data*

*"charged"*

⬤ = 1 or 0

*design-dependent*

*"discharged"*

◯ = 0 or 1

# Case Study: Mitigating DRAM Refresh Overheads

Every capacitor **leaks charge** over time

**Fully charged**     **Data-retention error**

**DRAM Refresh**

Periodically restores leaked charge **to every cell**
*(default period = 32-64 ms)*

Significant system **performance and energy overhead**

# Case Study: Mitigating DRAM Refresh Overheads

- DRAM refresh performance overheads in a 4-core system *[Patel+, ISCA'17]*

**Average System Performance Overhead** *(Normalized to No-Refresh)*

**More cells to refresh**

| DRAM Chip Size (Gb) | Overhead |
|---|---|
| 8 | 2.9% |
| 16 | 4.7% |
| 32 | 8.1% |
| 64 | 15.9% |

DRAM Chip Size (Gb)

# Case Study: Mitigating DRAM Refresh Overheads

- Fortunately, most DRAM cells **do not fail** at a longer refresh interval



# Failing Cells @ 45°C (y-axis, log scale: 1, 10, 100, 1,000, 10,000, 100,000)

Refresh Period
Scaled to Default (64 ms) (x-axis: 1, 8, 16, 24, 32, 40, 48, 56, 64, 72, 80, 88, 96, 104, 112, 120, 128)

**< 100 failures**

# Problem: Finding Fast-/Slow-Leaking Cells



**fast-leaking**

**slow-leaking**

- Unfortunately, finding those cells requires **memory testing**
  - **Difficult task** that relies on knowing or reverse-engineering DRAM design details
  - E.g., internal cell organization, worst-case testing parameters

**Unsupported** by the separation of concerns

# Finding "Weak" Cells

**Violates the separation**

*Testing*

*Modeling*

**Chip Design Properties**
- Substructure (subarray) size/location
- Internal address mappings
- On-die mechanisms (e.g., ECC, TRR)

**Test Methodology**
- Worst-case data/access patterns
- Chip features to enable/disable

**Test Results**
- Error probabilities and distributions
- Lists of failing cells

**?**

**Modeling**
- Analytical relationships
- Statistical distributions
- Physics-based simulations
- Empirically-measured curves

**"Weak" Cell Locations**

# Example System-Memory Cooperative Solutions

| To Improve... | System designers can... |
|---|---|
| Performance Energy & Power | Exploit expendable operating margins<br>• Access and refresh timings<br>• Operating voltage and temperature |
| Reliability | Protect against various failure modes<br>• System-level error mitigations<br>• Detailed qualification and validation |
| Security | Implement system-level defenses<br>• RowHammer, cold-boot attacks |

All rely on **unstandardized information** about DRAM reliability and testing

# Revision-by-Example: Case Studies

**Performance**
Reducing the long DRAM access latency

**Efficiency**
Improving refresh power and performance
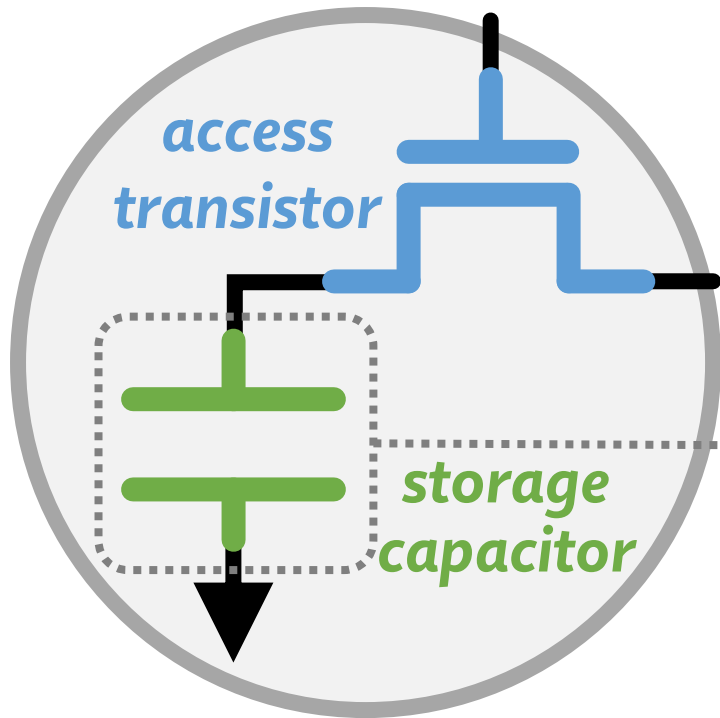
**Reliability**
Mitigating worsening memory errors

**Security**
Addressing worsening RowHammer vulnerability

**All are discouraged by
information that standards abstract away**

# Talk Outline

- DRAM Scaling Challenges
- Addressing Scaling: The Separation of Concerns
  - Problem 1: Inflexibility to Challenges
  - Problem 2: Overly Constraining
- Enabling System-Memory Cooperation
- **Revising the Separation**

# How should we revise the separation of concerns?

- *Near-term* recommendations to the industry should be:
  - **Achievable**: do not rely on specific changes to DRAM hardware
  - **Practical**: preserve the commodity industry (e.g., trade secrets)
  - **Backwards-compatible**: preserve "drop-in use" of DRAM like today
- *Long-term* recommendations can be more encompassing

**Information Transparency**                    **Changes to Future Standards**



**Near Term**                                                        **Long Term**

*New designs enabled by
information transparency can be
adopted in future standards*

# Near-Term Revisions

- We study the **practicality** of **information transparency**



|  | **Black-Box** *Today's DRAM* | **Grey Box** | **White-Box** *"Open-Source" DRAM* |
|---|---|---|---|
| **Producer Burden** | ☑ Low | + Medium | ✘ High |
| **Consumer Burden** | ☑ Low | ☑ Low-Med | ☑ Low-Med |
| **System-Memory Cooperation** | ✘ Discourages | + Enables-Encourages | ☑ Encourages |

**Good middle ground**

# Identifying Information to Release

- Provide transparency in DRAM chip **design and test**
  - **No physical changes** to hardware
  - **Key chip properties** that are useful (or even necessary) for **system-memory cooperation**
- Choose properties that can be **reverse-engineered** (no longer secret)

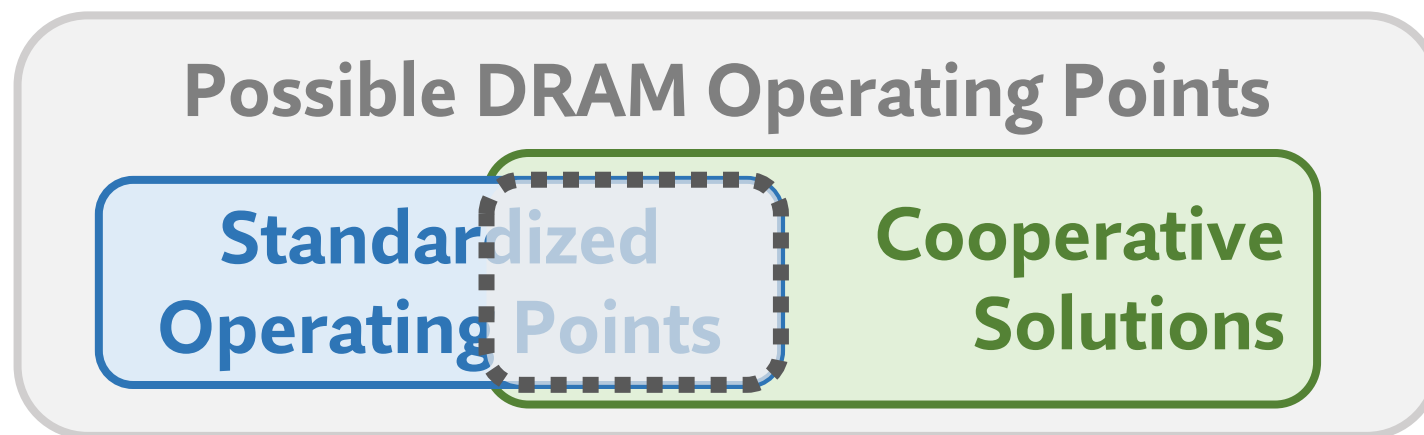| Design Characteristic | Reverse-Engineered By | Use-Case(s) Relying on Knowing the Characteristic |
|---|---|---|
| Cell charge encoding convention (i.e., true- and anti-cell layout) | Testing [78, 95, 98, 189] | Data-retention error modeling and testing for mitigating refresh overheads (e.g., designing worst-case test patterns) [98, 130, 189] |
| On-die ECC details | Modeling and testing [95, 258] | Improving reliability (e.g., designing ECC within the memory controller) [27, 30, 101, 321], mitigating RowHammer [100, 216, 219, 222] |
| Target row refresh (TRR) details | Testing [100, 160] | Modeling and mitigating RowHammer [100, 160, 222] |
| Mapping between internal and external row addresses | Testing [69, 94, 216, 297, 299, 342] | Mitigating RowHammer [87, 94, 216, 297, 298] |
| Row addresses refreshed by each refresh operation | Testing [100] | Mitigating RowHammer [100], improving access timings [70, 211] |
| Substructure organization (e.g., cell array dimensions) | Modeling [69] and testing [39, 67, 69] | Improving DRAM access timings [39, 67, 69] |
| Analytical model parameters (e.g., bitline capacitance) | Modeling and testing [189, 191] | Developing and using error models for improving overall reliability [276], mitigating refresh overheads (e.g., data-retention [191, 271, 275] and VRT [283, 284] models), improving access timings [69], and mitigating RowHammer [270, 343] |

Table 2: Basic DRAM chip design characteristics that are typically assumed or inferred for experimental studies.

44

# Two-Step Plan to Revise DRAM Standards

**1** Enable Cooperative Solutions

Possible DRAM Operating Points

Standardized Operating Points

Cooperative Solutions

**Near-Term Plan**

**2** Broaden Standards

Possible DRAM Operating Points

Standardized Operating Points

Cooperative Solutions

**Long-Term Plan**

45

# Two-Step Plan to Revise DRAM Standards

**1** Enable Cooperative Solutions

**Possible DRAM Operating Points**

**Standardized Operating Points**

**Cooperative Solutions**

**Near-Term Plan**

- Information transparency about DRAM chip design
- Communicate information in two key ways
    1. **Directly provided** by chip manufacturers (e.g., datasheets)
    2. **Crowdsourced database** built by DRAM consumers

# Two-Step Plan to Revise DRAM Standards

- **Broaden DRAM standards** to foster cooperative designs
  - Incorporate **testing and reliability** information that does not exist today
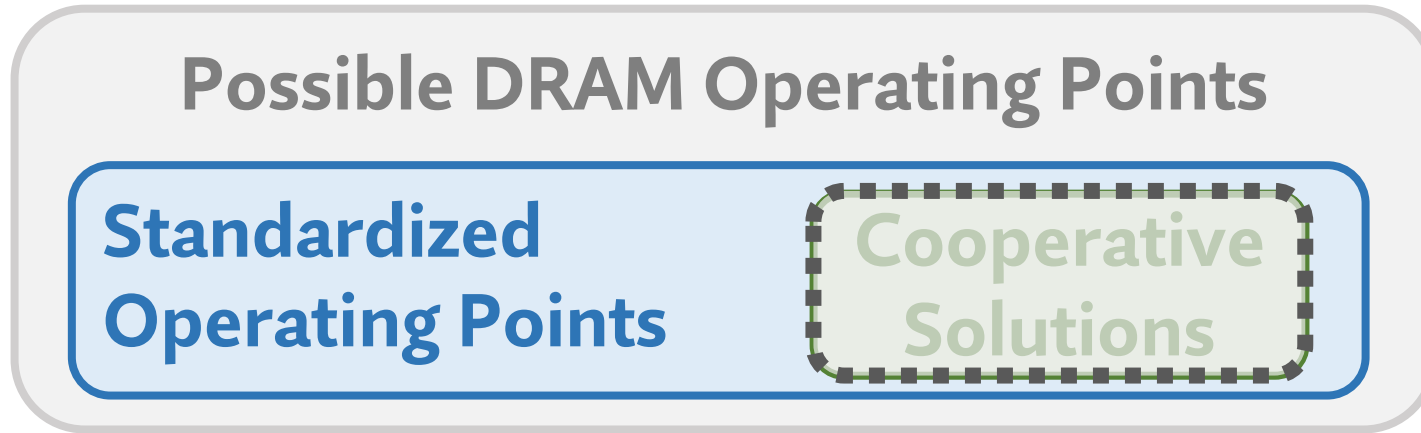  - Specific information will organically grow from:
    - System-memory cooperative solutions driven by the near-term plan
    - Industry-wide need for more efficient scaling solutions

② Broaden Standards

**Possible DRAM Operating Points**

**Standardized Operating Points**

**Cooperative Solutions**

**Long-Term Plan**

# A Case for Transparent Reliability in DRAM Systems

Minesh Patel[†]   Taha Shahroodi[‡†]   Aditya Manglik[†]   A. Giray Yağlıkçı[†]

Ataberk Olgun[†]   Haocong Luo[†]   Onur Mutlu[†]

[†]*ETH Zürich*   [‡]*TU Delft*

Mass-produced commodity DRAM is the preferred choice of main memory for a broad range of computing systems due to its favorable cost-per-bit. However, today's systems have diverse system-specific needs (e.g., performance, energy, reliability) that are difficult to address using one-size-fits-all general-purpose DRAM. Unfortunately, although system designers can theoretically adapt commodity DRAM chips to meet their particular design goals (e.g., by exploiting slack in access timings to improve performance, or implementing system-level RowHammer mitigations), we observe that designers today lack the necessary insight into commodity DRAM chips' reliability characteristics to implement these techniques in practice.

who purchase, test, and/or integrate commodity DRAM chips (e.g., cloud system designers, processor and system-on-a-chip (SoC) architects, memory module designers, test and validation engineers) are free to focus on the particular challenges of the systems they work on instead of dealing with the nuances of building low-cost, high-performance DRAM.

To ensure that system designers can integrate commodity DRAM chips from any manufacturer, the DRAM interface and operating characteristics have long been standardized by the JEDEC consortium [8]. JEDEC maintains a limited set of *DRAM standards* for commodity DRAM chips with different target applications, e.g., general-purpose DDR$n$ [9–11], bandwidth-

Minesh Patel, Taha Shahroodi, Aditya Manglik, A. Giray Yaglikci, Ataberk Olgun, Haocong Luo, Onur Mutlu,
**"A Case for Transparent Reliability in DRAM Systems"**
*arXiv*, April 2022.

# Reshaping DRAM Scaling
# by Enabling System-Memory Cooperation

## Minesh Patel

**SAFARI** **ETH** *zürich*

Zürich • 2 November 2023