# *HeatWatch*

## Improving 3D NAND Flash Memory Device Reliability by Exploiting Self-Recovery and Temperature Awareness

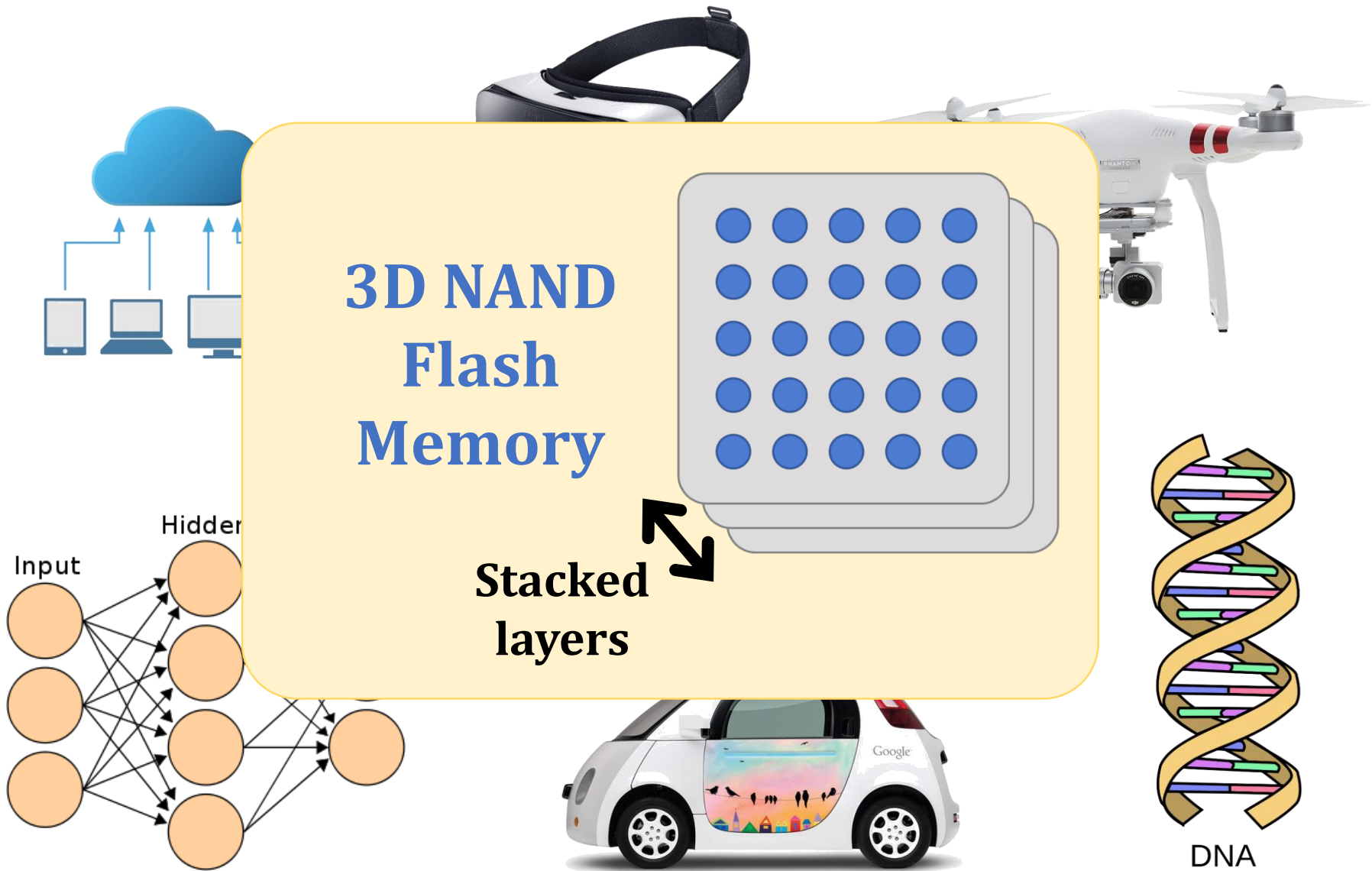**Yixin Luo**    Saugata Ghose    Yu Cai    Erich F. Haratsch    Onur Mutlu

Carnegie Mellon    SK hynix    ETH Zürich

SAFARI    SEAGATE

# Storage Technology Drivers - 2018

**3D NAND Flash Memory**

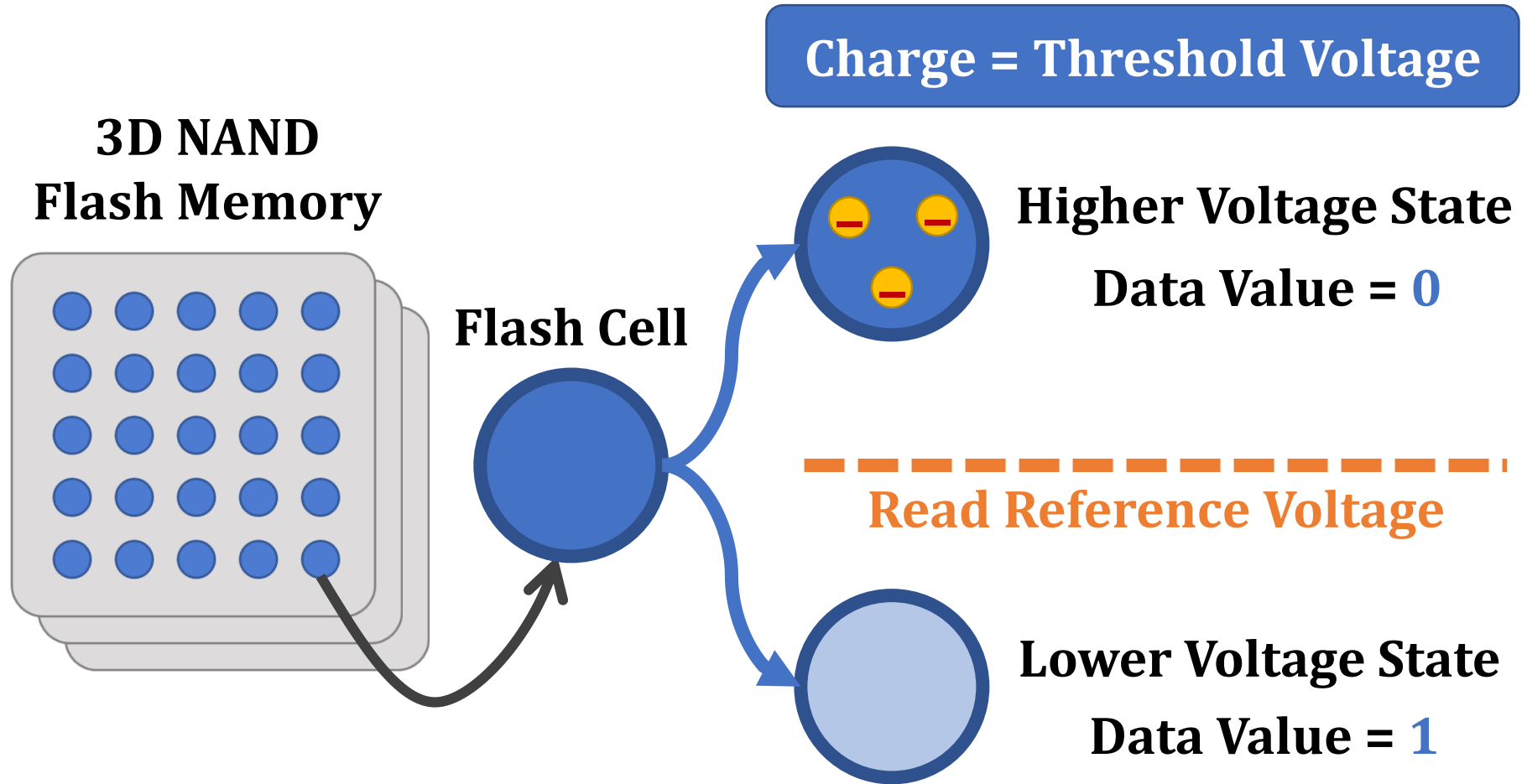Stacked layers

Input

Hidden

DNA

# Executive Summary

- 3D NAND susceptible to **early retention errors**
  - Charge leaks out of flash cell quickly after programming
  - Two unreported factors: ***self-recovery* and *temperature***

- We study *self-recovery* and *temperature* effects

> - **Experimental characterization** of *real* 3D NAND chips

> - **Unified Self-Recovery and Temperature (URT) Model**
>   - Predicts impact of retention loss, wearout, self-recovery, temperature on **flash cell voltage**
>   - **Low prediction error rate: 4.9%**

- We develop a new technique to improve flash reliability

> - **HeatWatch**
>   - Uses URT model to find optimal read voltages for 3D NAND flash
>   - **Improves flash lifetime by 3.85x**
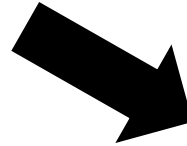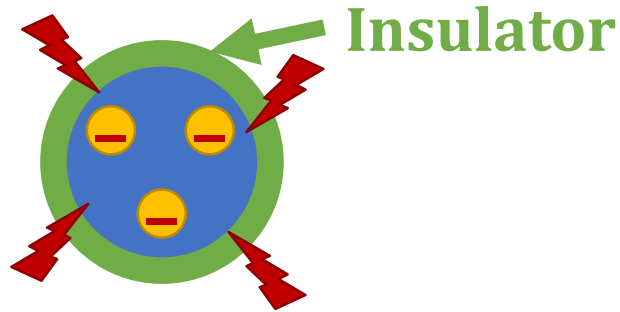
# Outline

- Executive Summary

- **Background on NAND Flash Reliability**

- Characterization of Self-Recovery and Temperature Effect on Real 3D NAND Flash Memory Chips

- URT: Unified Self-Recovery and Temperature Model

- HeatWatch Mechanism

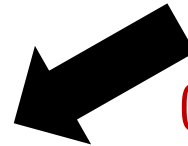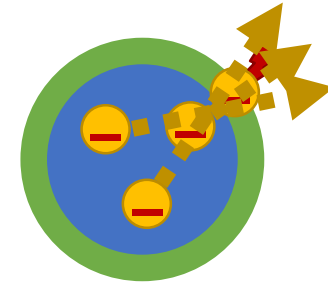- Conclusion

# 3D NAND Flash Memory Background

**3D NAND Flash Memory**

**Flash Cell**

**Charge = Threshold Voltage**

**Higher Voltage State**

**Data Value = 0**

**Read Reference Voltage**

**Lower Voltage State**

**Data Value = 1**

# Flash Wearout

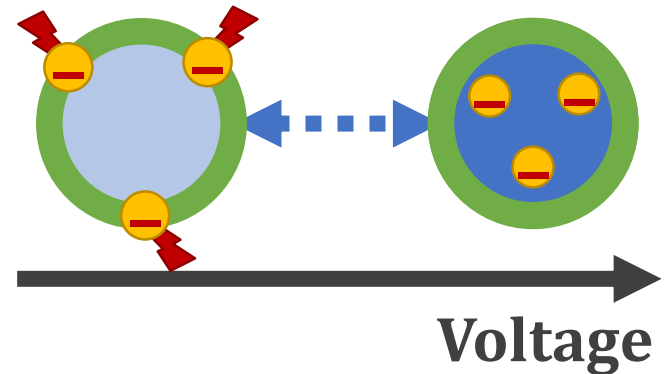**Program/Erase (P/E) → Wearout**

**Insulator**

**Wearout Effects:**

**1. Retention Loss**
(voltage shift over time)

**2. Program Variation**
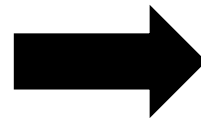(init. voltage difference b/w states)

**Wearout Introduces Errors**

**Voltage**

# Improving Flash Lifetime

**Errors introduced by wearout**
**<span style="color:red">limit flash lifetime</span>**
(measured in P/E cycles)

**Two Ways to Improve**
**Flash Lifetime**

➡

**Exploiting the**
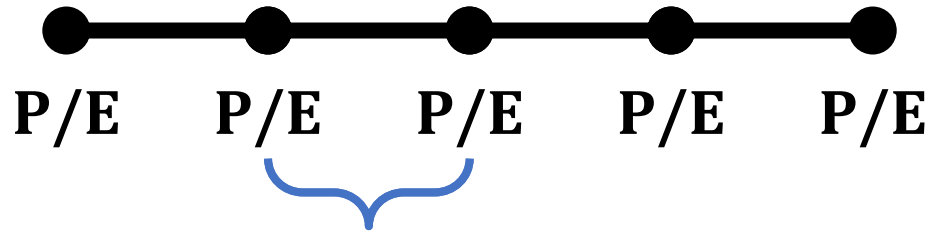**<span style="color:blue">Self-Recovery</span> Effect**
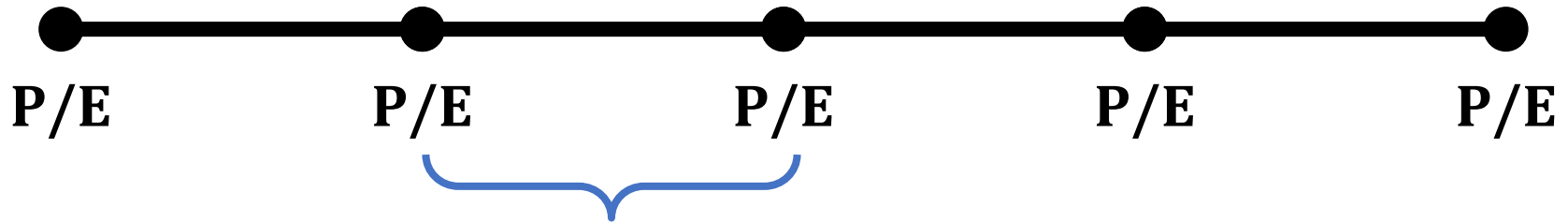
**Exploiting the**
**<span style="color:blue">Temperature</span> Effect**

# Exploiting the Self-Recovery Effect

**Partially repairs damage due to wearout**



**Dwell Time: Idle Time Between P/E Cycles**



**Longer Dwell Time: More Self-Recovery**

**Reduces Retention Loss**
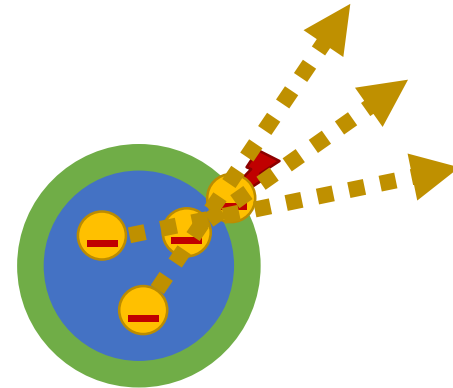
# Exploiting the Temperature Effect

**High Program Temperature**

**Voltage**

**Increases Program Variation**

**High Storage Temperature**

**Accelerates Retention Loss**

# Prior Studies of Self-Recovery/Temperature

|  | **Planar (2D) NAND** | **3D NAND** |
|---|---|---|
| **Self-Recovery Effect** | ✓ Mielke 2006 | ✗ |
| **Temperature Effect** | ✓ JEDEC 2010 (no characterization) | ✗ |

# Outline

- Executive Summary

- Background on NAND Flash Reliability

- **Characterization of Self-Recovery and Temperature Effect on Real 3D NAND Flash Memory Chips**

- URT: Unified Self-Recovery and Temperature Model

- HeatWatch Mechanism

- Conclusion

# Characterization Methodology

- Modified firmware version in the flash controller
  - Control the read reference voltage of the flash chip
  - Bypass ECC to get raw NAND data (with raw bit errors)
- Control temperature with a heat chamber



**Server**

**Heat Chamber**

**SSD**

# Characterized Devices

## Real 30-39 Layer 3D MLC NAND Flash Chips



**2-bit MLC**

**01**

**30- to 39-layer**

# MLC Threshold Voltage Distribution Background

# Characterization Goal



**Characterized Metrics**

**Retention Loss Speed**
(how fast voltage shifts over time)

**Program Variation**
(initial voltage difference between states)

**Characterized Phenomena**

**Self-Recovery Effect**

**Temperature Effect**

# Self-Recovery Effect Characterization Results



Dwell time: Idle time between P/E cycles

**Increasing dwell time from 1 minute to 2.3 hours slows down retention loss speed by 40%**

# Program Temperature Effect Characterization Results



**Increasing program temperature from 0°C to 70°C improves program variation by 21%**

# Storage Temperature Effect Characterization Results



**Lowering storage temperature from 70°C to 0°C slows down retention loss speed by 58%**

# Characterization Summary

**Major Results:**

- *Self-recovery* affects retention loss speed
- Program *temperature* affects program variation
- *Storage temperature* affects retention loss speed

**Unified Model**

**Other Characterizations Methods in the Paper:**

- More detailed results on self-recovery and temperature
  - Effects on error rate
  - Effects on threshold voltage distribution
- Effects of recovery cycle (P/E cycles with long dwell time) on retention loss speed

# Outline

- Executive Summary

- Background on NAND Flash Reliability

- Characterization of Self-Recovery and Temperature Effect on Real 3D NAND Flash Memory Chips

- **URT: Unified Self-Recovery and Temperature Model**

- HeatWatch Mechanism

- Conclusion

# Minimizing 3D NAND Errors

# Predicting the Mean Threshold Voltage

## Our URT Model:

$$V = V_0 + \Delta V$$

**Mean Threshold Voltage**

**Initial Voltage Before Retention (Program Variation)**

**Voltage Shift Due to Retention Loss**

# URT Model Overview

# 1. Program Variation Component

P/E Cycle

Program Temperature

PEC

$T_p$

$V_0$

Initial Voltage

$$V_0 = A \cdot T_p \cdot PEC + B \cdot T_p + C \cdot PEC + D$$

**Validation: $R^2$ = 91.7%**

# 2. Self-Recovery and Retention Component

**Retention Time**  **P/E Cycle**  **Dwell Time**

$t_r$  PEC  $T_d$

**ΔV**

**Retention Shift**

$$\Delta V\left(t_{er}, t_{ed}, PEC\right) = b \cdot (PEC + c) \cdot \ln\left(1 + \frac{t_{er}}{t_0 + a \cdot t_{ed}}\right)$$

**Validation: 3x more accurate than state-of-the-art model**

# 3. Temperature Scaling Component

**Actual Retention Time** — $t_r$

**Storage Temp.** — $T_r$

$t_{r,eff}$

**Effective Retention Time**

**Actual Dwell Time** — $t_d$

**Dwell Temp.** — $T_d$

$t_{d,eff}$

**Effective Dwell Time**

*Arrhenius Equation:*

$$AF = \frac{t_{real}}{t_{room}} = \exp\left(\frac{E_a}{k_B} \cdot \left(\frac{1}{T_{real}} - \frac{1}{T_{room}}\right)\right)$$

**Validation: Adjust an important parameter, $E_a$, from 1.1 eV to 1.04 eV**

# URT Model Summary



**1. Program Variation Component**

**3. Temperature Scaling Component**

**2. Self-Recovery and Retention Component**

$$V = V_0 + \Delta V$$

**Validation: Prediction Error Rate = 4.9%**

# Outline

- Executive Summary

- Background on NAND Flash Reliability

- Characterization of Self-Recovery and Temperature Effect on Real 3D NAND Flash Memory Chips

- URT: Unified Self-Recovery and Temperature Model
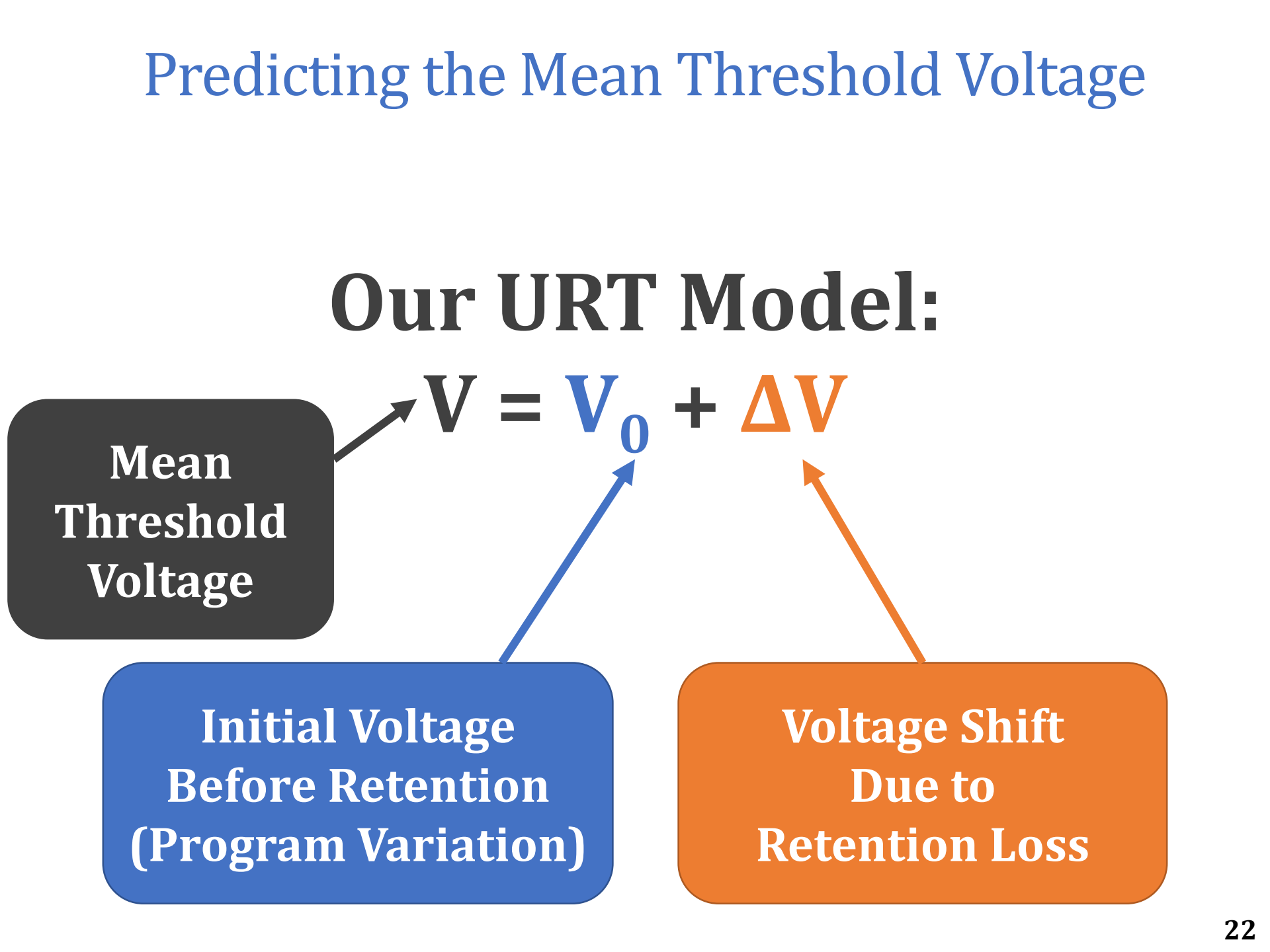
- **HeatWatch Mechanism**

- Conclusion

# HeatWatch Mechanism

- **Key Idea**

  - **Predict change in threshold voltage distribution** by using the URT model

  - **Adapt read reference voltage to near-optimal** $(V_{opt})$ based on predicted change in voltage distribution

# HeatWatch Mechanism Overview

# Tracking SSD Temperature

**Tracking Components**

**SSD Temperature**

**Dwell Time**

**P/E Cycles & Retention Time**

- Use existing sensors in the SSD
- **Precompute** temperature scaling factor at **logarithmic time intervals**

**Prediction Components**

$V_{opt}$ **Prediction**

**Fine-Tuning URT Parameters**

# Tracking Dwell Time

**Tracking Components**

| SSD Temperature | **Dwell Time** | P/E Cycles & Retention Time |

- Only need to log the timestamps of **last 20 full drive writes**
  - Self-recovery effect diminishes after 20 P/E cycles

**Prediction Components**

$V_{opt}$ Prediction

Fine-Tuning URT Parameters

# Tracking P/E Cycles and Retention Time

**Tracking Components**

SSD Temperature

Dwell Time

**P/E Cycles & Retention Time**

- P/E cycle count **already recorded** by SSD
- **Log write timestamp** for each block
- Retention time = read timestamp – write timestamp

**Prediction Components**

$V_{opt}$ Prediction

Fine-Tuning URT Parameters

# Predicting Optimal Read Reference Voltage

## Tracking Components

| SSD Temperature | Dwell Time | P/E Cycles & Retention Time |
|---|---|---|

- **Calculate URT** using tracked information
- Modeling error: 4.9%

## Prediction Components

$V_{opt}$ **Prediction**

**Fine-Tuning URT Parameters**

# Fine-Tuning URT Parameters Online

**Tracking Components**

| SSD Temperature | Dwell Time | P/E Cycles & Retention Time |
|---|---|---|

- Accommodates **chip-to-chip variation**
- Uses **periodic sampling**

**Prediction Components**

$V_{opt}$ Prediction

**Fine-Tuning URT Parameters**

# HeatWatch Mechanism Summary

**Tracking Components**

| SSD Temperature | Dwell Time | P/E Cycles & Retention Time |
|---|---|---|

**Storage Overhead: 0.16% of DRAM in 1TB SSD**

URT

**Prediction Components**

$V_{opt}$ Prediction

Fine-Tuning URT Parameters

**Latency Overhead: < 1% of flash read latency**

# HeatWatch Evaluation Methodology

- **28 real workload storage traces**
  - MSR-Cambridge
  - We use **real dwell time, retention time values** obtained from traces

- **Temperature Model:**
  Trigonometric function + Gaussian noise
  - Represents periodic temperature variation in each day
  - Includes small transient temperature variation

# HeatWatch Greatly Improves Flash Lifetime



**HeatWatch improves lifetime by capturing the effect of retention, wearout, self-recovery, temperature**

# Outline

- Executive Summary

- Background on NAND Flash Reliability

- Characterization of Self-Recovery and Temperature Effect on Real 3D NAND Flash Memory Chips

- URT: Unified Self-Recovery and Temperature Model

- HeatWatch Mechanism

- **Conclusion**

# Conclusion

- 3D NAND susceptible to **early retention errors**
  - Charge leaks out of flash cell quickly after programming
  - Two unreported factors: *self-recovery* **and** *temperature*

- We study *self-recovery* and *temperature* effects
  - **Experimental characterization** of *real* 3D NAND chips

  - **Unified Self-Recovery and Temperature (URT) Model**
    - Predicts impact of retention loss, wearout, self-recovery, temperature on **flash cell voltage**
    - **Low prediction error rate: 4.9%**

- We develop a new technique to improve flash reliability
  - **HeatWatch**
    - Uses URT model to find optimal read voltages for 3D NAND flash
    - **Improves flash lifetime by 3.85x**

# *HeatWatch*

## Improving 3D NAND Flash Memory Device Reliability by Exploiting Self-Recovery and Temperature Awareness

**Yixin Luo**     **Saugata Ghose**     **Yu Cai**     **Erich F. Haratsch**     **Onur Mutlu**

# References to Papers and Talks

# Our FMS Talks and Posters

- FMS 2018
  - Yixin Luo, HeatWatch: Exploiting 3D NAND Self-Recovery and Temperature Effects
  - Saugata Ghose, Enabling Realistic Studies of Modern Multi-Queue SSD Devices
- FMS 2017
  - Aya Fukami, Improving Chip-Off Forensic Analysis for NAND Flash
  - Saugata Ghose, Vulnerabilities in MLC NAND Flash Memory Programming
- FMS 2016
  - Onur Mutlu, ThyNVM: Software-Transparent Crash Consistency for Persistent Memory
  - Onur Mutlu, Large-Scale Study of In-the-Field Flash Failures
  - Yixin Luo, Practical Threshold Voltage Distribution Modeling
  - Saugata Ghose, Write-hotness Aware Retention Management
- FMS 2015
  - Onur Mutlu, Read Disturb Errors in MLC NAND Flash Memory
  - Yixin Luo, Data Retention in MLC NAND Flash Memory
- FMS 2014
  - Onur Mutlu, Error Analysis and Management for MLC NAND Flash Memory

# Our Flash Memory Works (I)

- Summary of our work in NAND flash memory
  - Yu Cai, Saugata Ghose, Erich F. Haratsch, Yixin Luo, and Onur Mutlu, [Error Characterization, Mitigation, and Recovery in Flash Memory Based Solid-State Drives](), Proceedings of the IEEE, Sept. 2017.

- Overall flash error analysis
  - Yu Cai, Erich F. Haratsch, Onur Mutlu, and Ken Mai, [Error Patterns in MLC NAND Flash Memory: Measurement, Characterization, and Analysis](), DATE 2012.
  - Yu Cai, Gulay Yalcin, Onur Mutlu, Erich F. Haratsch, Adrian Cristal, Osman Unsal, and Ken Mai, [Error Analysis and Retention-Aware Error Management for NAND Flash Memory](), ITJ 2013.
  - Yixin Luo, Saugata Ghose, Yu Cai, Erich F. Haratsch, and Onur Mutlu, [Enabling Accurate and Practical Online Flash Channel Modeling for Modern MLC NAND Flash Memory](), IEEE JSAC, Sept. 2016.

# Our Flash Memory Works (II)

- 3D NAND flash memory error analysis
  - Yixin Luo, Saugata Ghose, Yu Cai, Erich F. Haratsch, and Onur Mutlu, Improving 3D NAND Flash Memory Lifetime by Tolerating Early Retention Loss and Process Variation, SIGMETRICS 2018.
  - Yixin Luo, Saugata Ghose, Yu Cai, Erich F. Haratsch, and Onur Mutlu, HeatWatch: Improving 3D NAND Flash Memory Device Reliability by Exploiting Self-Recovery and Temperature-Awareness, HPCA 2018.

- Multi-queue SSDs
  - Arash Tavakkol, Juan Gomez-Luna, Mohammad Sadrosadati, Saugata Ghose, and Onur Mutlu, MQSim: A Framework for Enabling Realistic Studies of Modern Multi-Queue SSD Devices, FAST 2018.
  - Arash Tavakkol, Mohammad Sadrosadati, Saugata Ghose, Jeremie Kim, Yixin Luo, Yaohua Wang, Nika Mansouri Ghiasi, Lois Orosa, Juan G. Luna and Onur Mutlu, FLIN: Enabling Fairness and Enhancing Performance in Modern NVMe Solid State Drives, ISCA 2018.

# Our Flash Memory Works (III)

- Flash-based SSD prototyping and testing platform
  - Yu Cai, Erich F. Haratsh, Mark McCartney, Ken Mai, <u>FPGA-based solid-state drive prototyping platform</u>, FCCM 2011.

- Retention noise study and management
  - Yu Cai, Gulay Yalcin, Onur Mutlu, Erich F. Haratsch, Adrian Cristal, Osman Unsal, and Ken Mai, <u>Flash Correct-and-Refresh: Retention-Aware Error Management for Increased Flash Memory Lifetime</u>, ICCD 2012.
  - Yu Cai, Yixin Luo, Erich F. Haratsch, Ken Mai, and Onur Mutlu, <u>Data Retention in MLC NAND Flash Memory: Characterization, Optimization and Recovery</u>, HPCA 2015.
  - Yixin Luo, Yu Cai, Saugata Ghose, Jongmoo Choi, and Onur Mutlu, <u>WARM: Improving NAND Flash Memory Lifetime with Write-hotness Aware Retention Management</u>, MSST 2015.
  - Aya Fukami, Saugata Ghose, Yixin Luo, Yu Cai, and Onur Mutlu, <u>Improving the Reliability of Chip-Off Forensic Analysis of NAND Flash Memory Devices</u>, Digital Investigation, Mar. 2017.

# Our Flash Memory Works (IV)

- Program and erase noise study
  - Yu Cai, Erich F. Haratsch, Onur Mutlu, and Ken Mai, [Threshold Voltage Distribution in MLC NAND Flash Memory: Characterization, Analysis and Modeling](), DATE 2013.
  - Y. Cai, S. Ghose, Y. Luo, K. Mai, O. Mutlu, and E. F. Haratsch, [Vulnerabilities in MLC NAND Flash Memory Programming: Experimental Analysis, Exploits, and Mitigation Techniques](), HPCA 2017.

- Cell-to-cell interference characterization and tolerance
  - Yu Cai, Onur Mutlu, Erich F. Haratsch, and Ken Mai, [Program Interference in MLC NAND Flash Memory: Characterization, Modeling, and Mitigation](), ICCD 2013.
  - Yu Cai, Gulay Yalcin, Onur Mutlu, Erich F. Haratsch, Osman Unsal, Adrian Cristal, and Ken Mai, [Neighbor-Cell Assisted Error Correction for MLC NAND Flash Memories](), SIGMETRICS 2014.

# Our Flash Memory Works (V)

- Read disturb noise study
  - Yu Cai, Yixin Luo, Saugata Ghose, Erich F. Haratsch, Ken Mai, and Onur Mutlu, [Read Disturb Errors in MLC NAND Flash Memory: Characterization and Mitigation](), DSN 2015.

- Flash errors in the field
  - Justin Meza, Qiang Wu, Sanjeev Kumar, and Onur Mutlu, [A Large-Scale Study of Flash Memory Errors in the Field](), SIGMETRICS 2015.

- Persistent memory
  - Jinglei Ren, Jishen Zhao, Samira Khan, Jongmoo Choi, Yongwei Wu, and Onur Mutlu, [ThyNVM: Enabling Software-Transparent Crash Consistency in Persistent Memory Systems](), MICRO 2015.

# Referenced Papers and Talks

- All are available at
  - https://www.ece.cmu.edu/~safari/pubs.html
  - https://www.ece.cmu.edu/~safari/talks.html

- And, many other previous works on
  - Challenges and opportunities in memory
  - NAND flash memory errors and management
  - Phase change memory as DRAM replacement
  - STT-MRAM as DRAM replacement
  - Taking advantage of persistence in memory
  - Hybrid DRAM + NVM systems
  - NVM design and architecture