A Large-Scale Study of
# Flash Memory Errors in the Field

**Justin Meza**
Qiang Wu
Sanjeev Kumar
Onur Mutlu

facebook

Carnegie Mellon University

# Overview

*First study of flash reliability:*

- at a large scale

- in the field

# Overview

SSD lifecycle

Access pattern
dependence

New
reliability
trends

Read
disturbance

Temperature

# Overview

## SSD lifecycle

**Early detection** lifecycle period distinct from hard disk drive lifecycle.

# Overview

*SSD lifecycle*

We **do not** observe the effects of **read disturbance** errors in the field.

**Read disturbance**

*Temperature*

# Overview

SSD lifecycle

New

Access
depe                                            ance

**Throttling SSD usage** helps mitigate temperature-induced errors.

*Temperature*

# Overview

SSD lifecycle

Access pattern dependence

We quantify the effects of the **page cache** and **write amplification** in the field.

Temperature

# Outline

- background and motivation
- server SSD architecture
- error collection/analysis methodology
- SSD reliability trends
- summary

# Background and motivation

# Flash memory

- persistent
- high performance
- hard disk alternative
- used in solid-state drives (SSDs)

# Flash memory

- persistent
- high performance
- hard disk alternative
- used in solid-state drives (SSDs)
- **prone to a variety of errors**
  - wearout, disturbance, retention

# Our goal

## *Understand SSD reliability:*

- at a large scale
  - millions of device-days, across four years
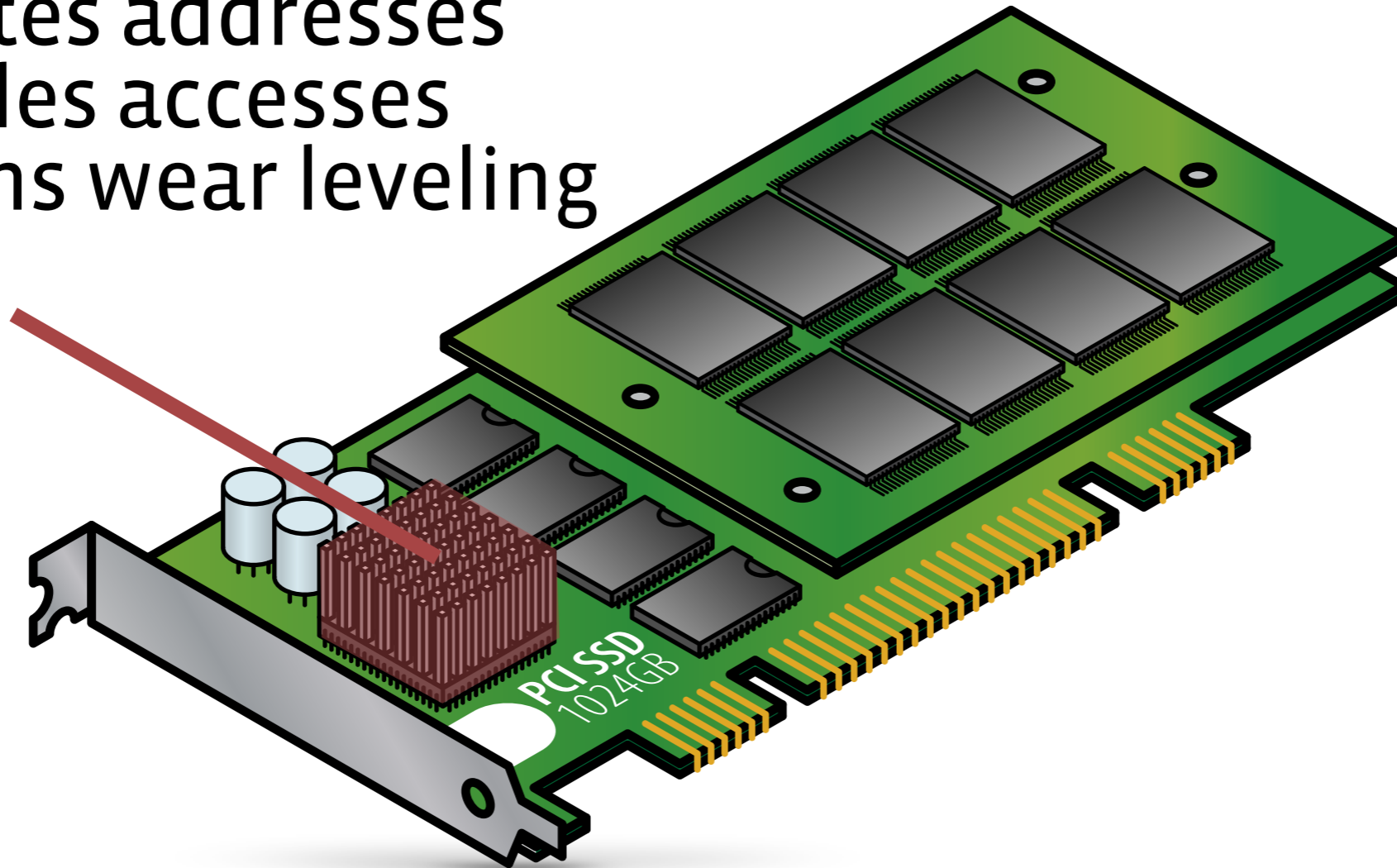
- in the field
  - realistic workloads and systems

# Server SSD architecture

PCI SSD
1024GB

PCIe

Flash chips

# SSD controller

- translates addresses
- schedules accesses
- performs wear leveling

User data
```
10011111 11001111 11000011 00001101
10101110 11100101 11111001 01111011
00011001 11011101 11100011 11111000
11011111 01001101 11110000 10111111
00000001 11011110 00000101 01010110
00001011 10000010 11111110 00011100
...
```

ECC metadata
```
01001100 01001101 11010010 01000000
10011100 10111111 10101111 11000101
```

# Types of errors

## Small errors

- 10's of flipped bits per KB
- silently corrected by SSD controller

## Large errors

- 100's of flipped bits per KB
- corrected by host using driver
- referred to as SSD failure

# Types of errors

## Small errors

We examine *large errors (SSD failures)* in this study.

## Large errors

- ~100's of flipped bits per KB
- corrected by host using driver
- refer to as SSD failure

# Error collection/analysis methodology

# SSD data measurement

- metrics stored on SSDs
- measured across SSD lifetime

# SSD characteristics

- 6 different system configurations
  - 720GB, 1.2TB, and 3.2TB SSDs
  - servers have 1 or 2 SSDs
  - this talk: representative systems
- 6 months to 4 years of operation
- 15TB to 50TB read and written

# Bit error rates (BER)

- BER = bit errors per bits transmitted
- *1 error per 385M bits* transmitted to *1 error per 19.6B bits* transmitted
  - averaged across all SSDs in each system type
- *10x to 1000x lower* than prior studies
  - large errors, SSD performs wear leveling

# A few SSDs cause most errors

# A few SSDs cause most errors

# A few SSDs cause most errors



*What factors contribute to SSD failures in the field?*

# Analytical methodology

- not feasible to log every error
- instead, analyze **lifetime counters**
- **snapshot-based** analysis

|          |         |     |     |     |
| -------- | ------- | --- | --- | --- |
| Errors   | 54,326  | 0   | 2   | 10  |
| Data written | 10TB | 2TB | 5TB | 6TB |

| | | | | |
|---|---|---|---|---|
| Errors | 54,326 | 0 | 2 | 10 |
| Data written | 10TB | 2TB | 5TB | 6TB |

*2014-11-1*

| Errors | 54,326 | 0 | 2 | 10 |
| Data written | 10TB | 2TB | 5TB | 6TB |

*2014-11-1*

Errors

Data written

| | | | | |
|---|---|---|---|---|
| Errors | 54,326 | 0 | 2 | 10 |
| Data written | 10TB | 2TB | 5TB | 6TB |

2014-11-1

| Errors | 54,326 | 0 | 2 | 10 |
|---|---|---|---|---|
| Data written | 10TB | 2TB | 5TB | 6TB |

2014-11-1

Errors

Data written

| | | | |
|---|---|---|---|
| Errors | 54,326 | 0 | 2 | 10 |
| Data written | 10TB | 2TB | 5TB | 6TB |

*2014-11-1*

# SSD reliability trends

# SSD lifecycle

# Access pattern dependence

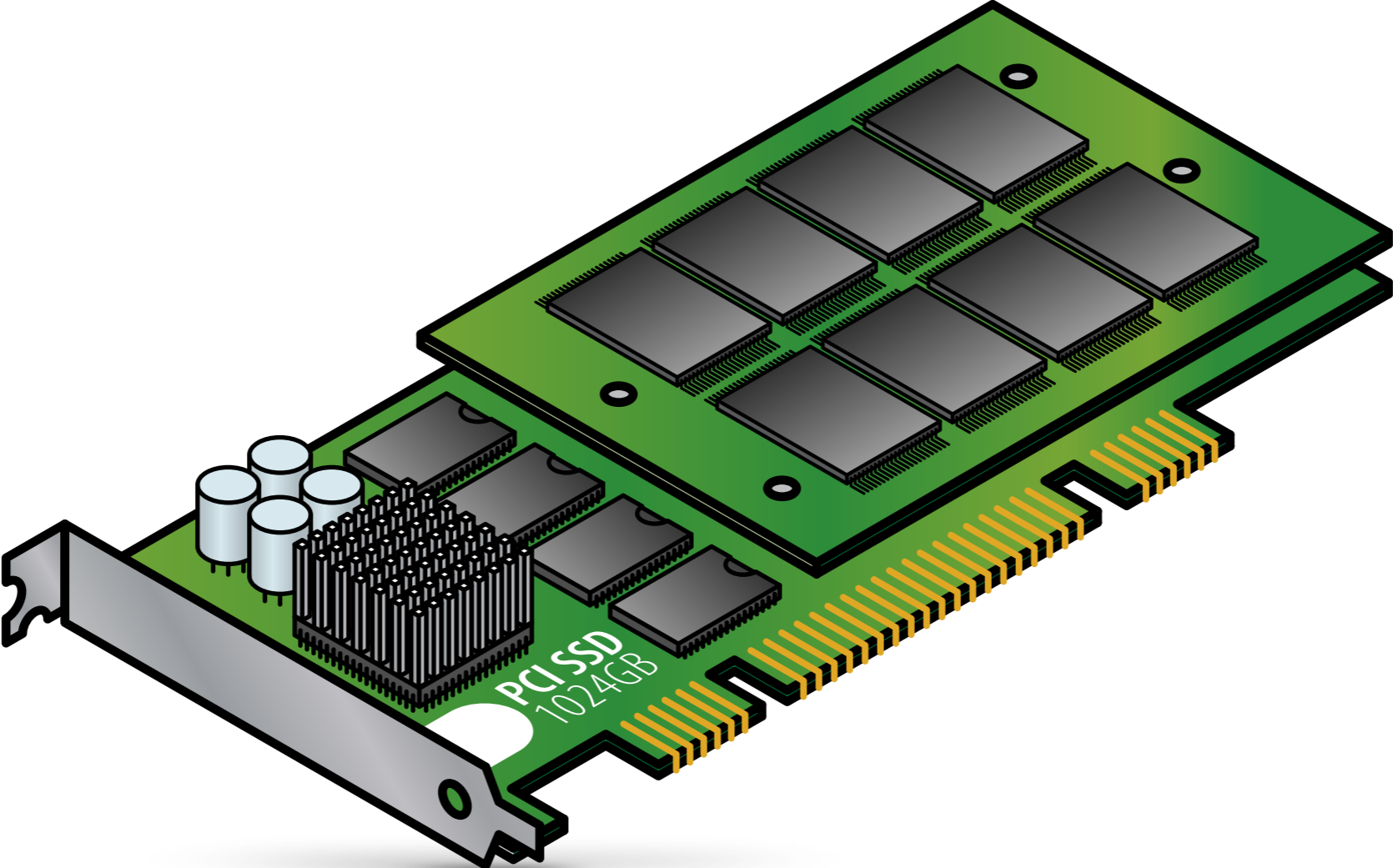# New reliability trends

# Read disturbance

# Temperature

*SSD lifecycle*

**New reliability trends**

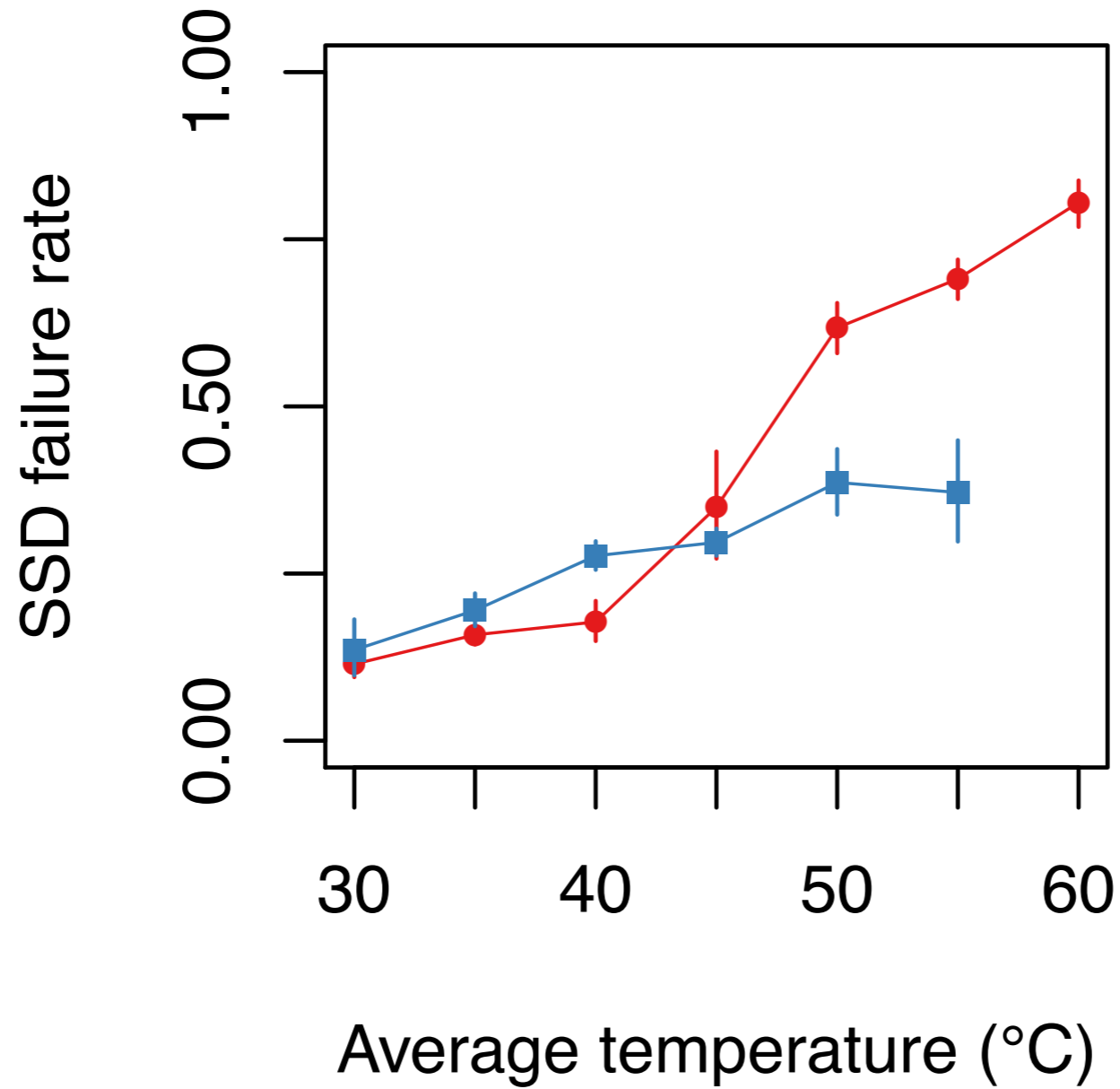*Access pattern dependence*

*Read disturbance*

*Temperature*

# Storage lifecycle background:
## *the bathtub curve* for disk drives



**Failure rate**

**Usage**

[Schroeder+,FAST'07]

# Storage lifecycle background:
## *the **bathtub curve** for disk drives*

**Early failure period**

**Wearout period**

**Failure rate**

**Useful life period**

**Usage**

[Schroeder+,FAST'07]

Storage lifecycle background:
*the* **bathtub curve** for disk drives

**Early failure**

**Wearout**

Fai...

*period*

Usage

**Do SSDs display similar lifecycle periods?**

[Schroeder+,FAST'07]

*Use* **data written to flash**
**to** *examine SSD lifecycle*

*(time-independent utilization metric)*

# SSD lifecycle

**Early detection** lifecycle period distinct from hard disk drive lifecycle.

Temperature

SSD lifecycle

Access pattern
dependence

**New reliability trends**

*Read disturbance*

Temperature

# Read disturbance

- reading data can disturb contents
- failure mode identified in *lab setting*
- under *adversarial workloads*

# Read disturbance

- 
- 
- under adversarial workloads

**_Does read disturbance affect SSDs in the field?_**

*Examine SSDs with* high **flash R/W** *ratios* *and* **most data read** *to understand read effects*

*(isolate effects of read vs. write errors)*

**3.2TB, 1 SSD (average R/W = 2.14)**

**1.2TB, 1 SSD (average R/W = 1.15)**

SSD lifecycle

We **do not** observe the effects of **read disturbance** errors in the field.

**Read disturbance**

Temperature

SSD lifecycle

Access pattern
dependence

New
reliability
trends

Read
disturbance

Temperature

Temperature sensor

**720GB, 1 SSD** **720GB, 2 SSDs**

# High temperature: may throttle or shut down

PCI SSD
1024GB

SSD lifecycle

New

Access

depe

ance

**Throttling SSD usage** helps mitigate temperature-induced errors.

*Temperature*

# Access pattern effects

***System buffering***

- data served from OS caches
- decreases SSD usage

***Write amplification***

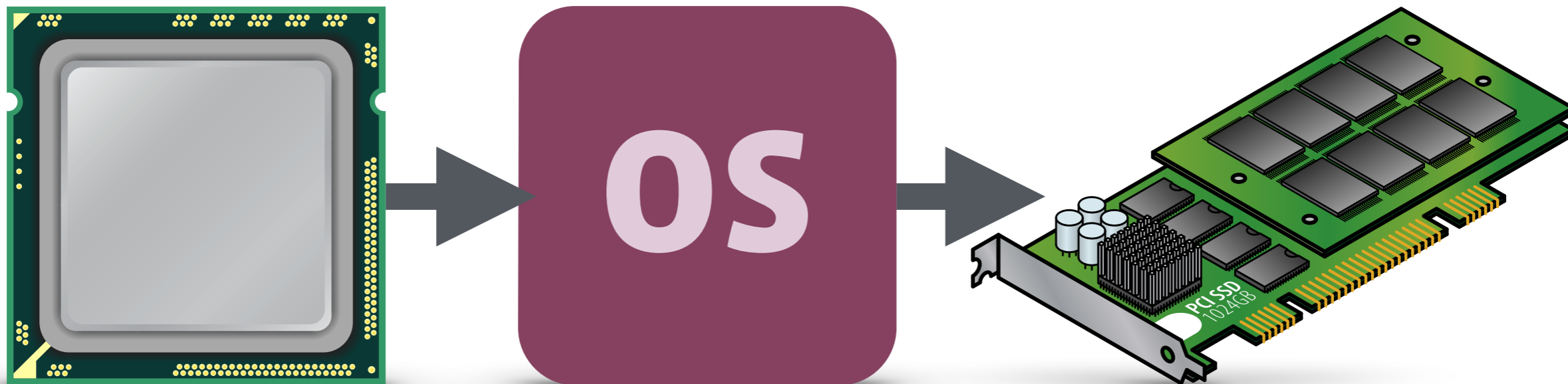- updates to small amounts of data
- increases erasing and copying

# Access pattern effects

## *System buffering*
- data served from OS caches
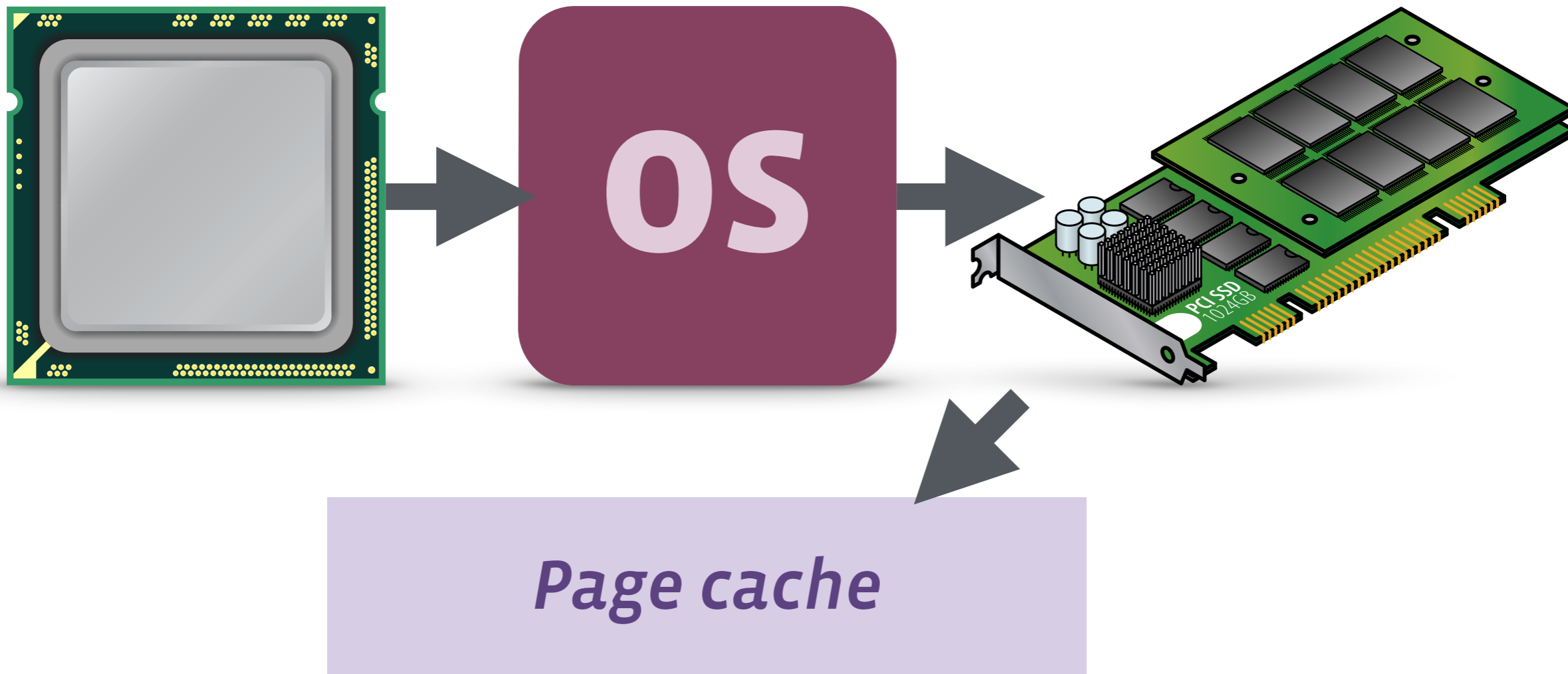- decreases SSD usage
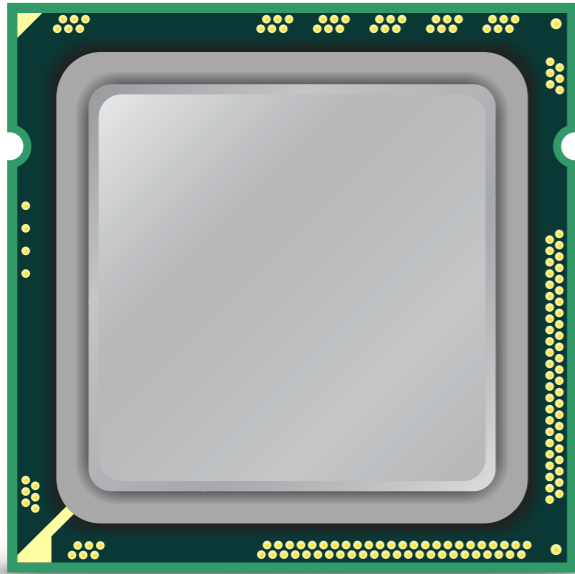
## *Write amplification*
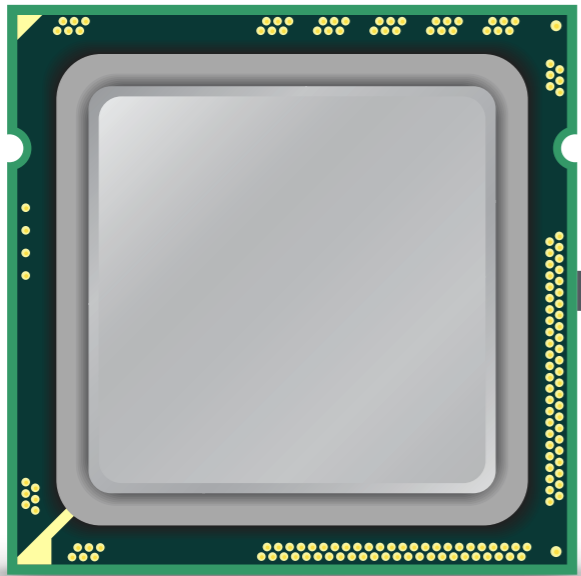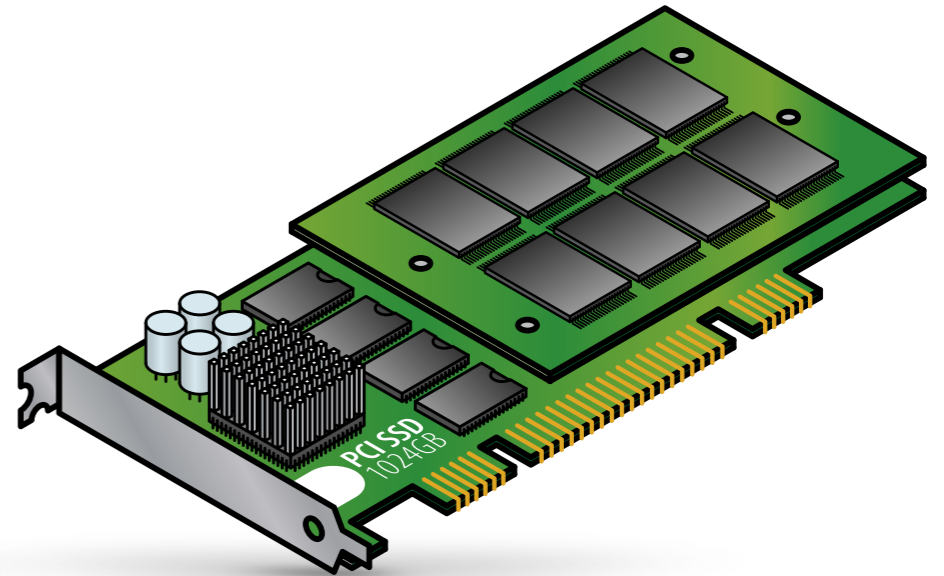- updates to small amounts of data
- increases erasing and copying
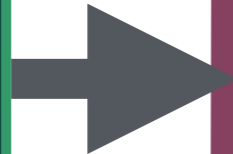
Page cache

**Page cache**

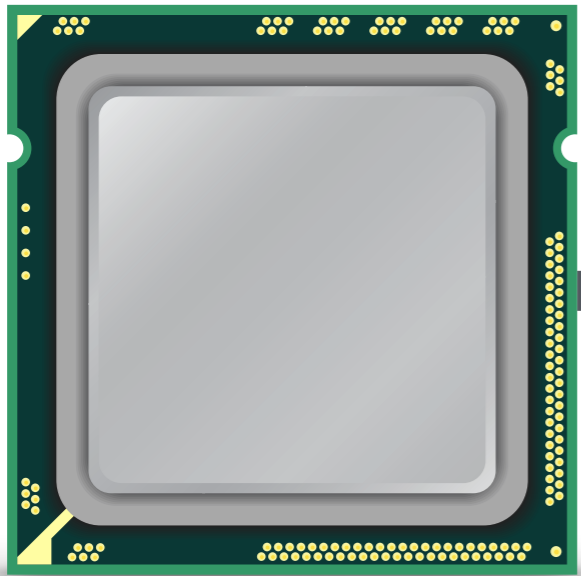Page cache

Page cache

**OS**

**Page cache**

Page cache

Page cache

# System caching *reduces* the impact of SSD writes

**OS**

**Page cache**

# 720GB, 2 SSDs

# Access pattern effects

## *System buffering*
- data served from OS caches
- decreases SSD usage

## ***Write amplification***
- updates to small amounts of data
- increases erasing and copying

# Flash devices use a
## **translation layer**
### to locate data
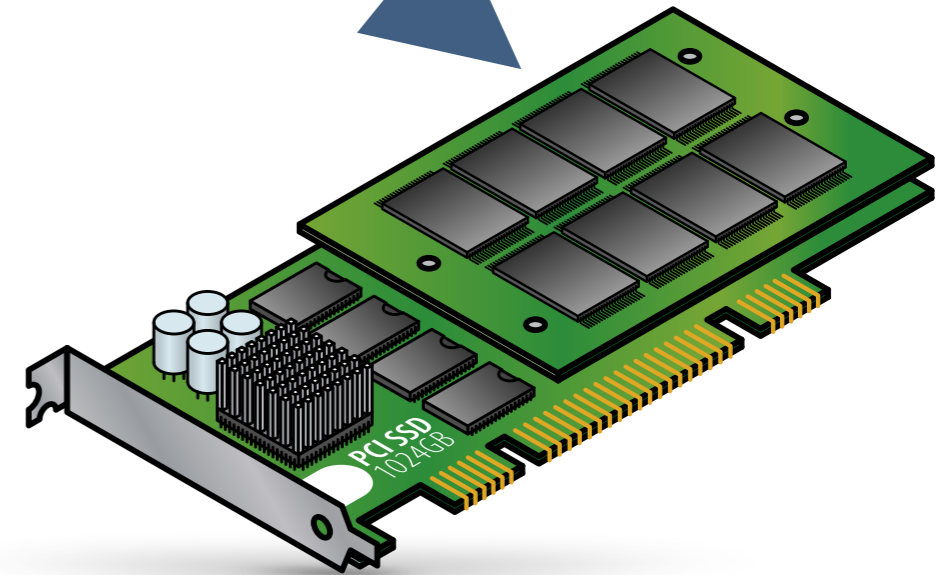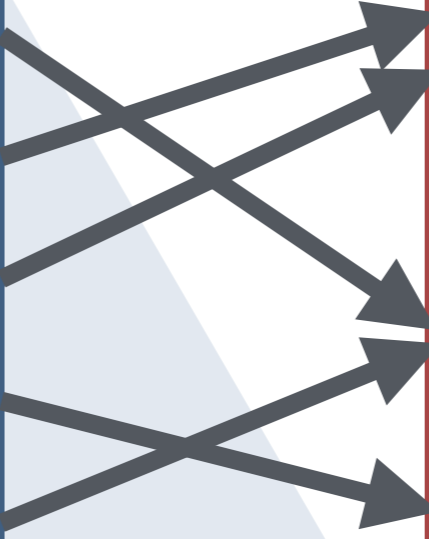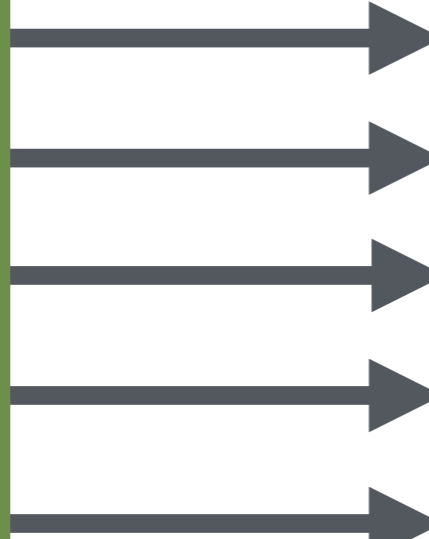
**OS**

# Translation layer

Logical address space

Physical address space

$\langle offset_1, size_1 \rangle$
$\langle offset_2, size_2 \rangle$
...

OS

PCI SSD
1024GB

# *Sparse* data layout

**more** *translation metadata*

potential for *higher* **write amplification**

# *Dense* data layout

## less *translation metadata*
## potential for *lower* write amplification

# *Use* **translation data size**
## *to examine effects of data layout*

*(relates to application access patterns)*

**720GB, 1 SSD**

Denser — Sparser

SSD failure rate vs. Translation data (GB)

# Write amplification in the field
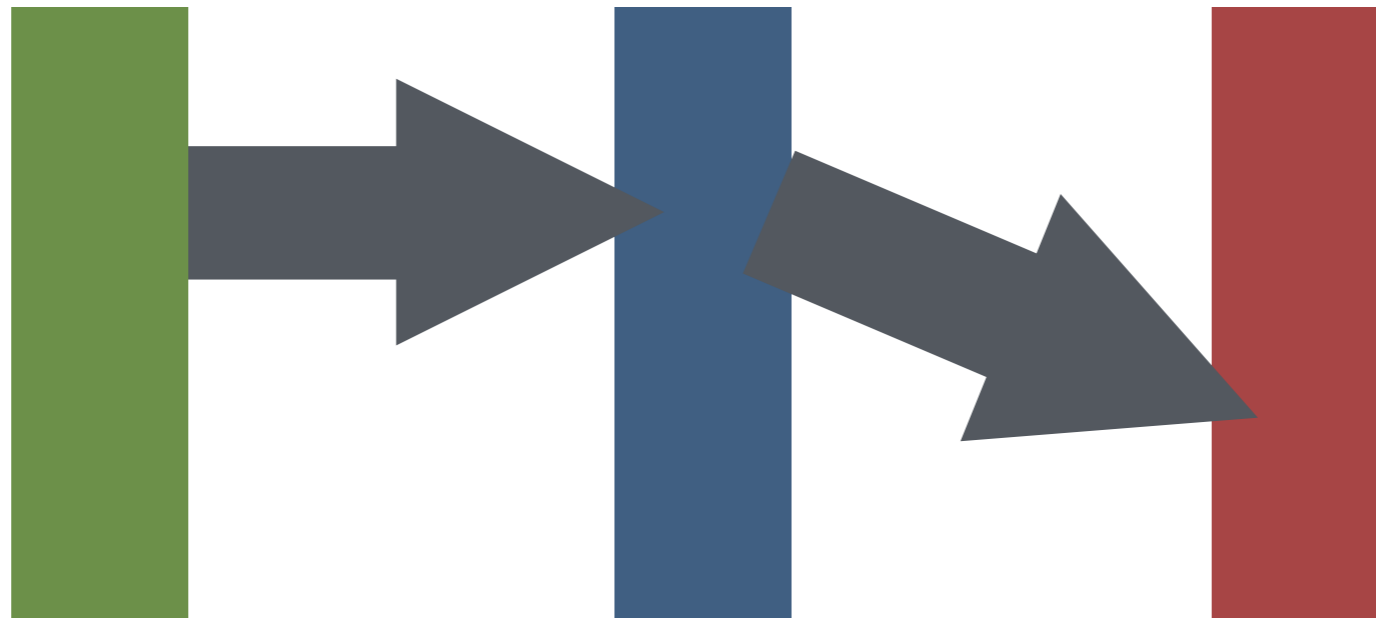
## Graph search

SSD failure rate

1.00
0.50
0.00

0.25    0.45

*Translation data (GB)*

## Key-value store

1.00
0.50
0.00

0.25    0.45

*Translation data (GB)*

AM buffer usa

DRAM buffer

SSD lifecycle

**Access pattern dependence**

We quantify the effects of the **page cache** and **write amplification** in the field.

Temperature

# More results in paper

- *Block erasures and discards*
- *Page copies*
- *Bus power consumption*

# *Summary*

- Large scale
- In the field

# Summary



SSD lifecycle

Access pattern dependence

New reliability trends

Read disturbance

Temperature

# Summary

*SSD lifecycle*

**Early detection** lifecycle period distinct from hard disk drive lifecycle.

*Access*

*depe*

*ance*

*trends*

*Temperature*

# Summary

SSD lifecycle

Read
disturbance

Temperature

We **do not** observe the effects of **read disturbance** errors in the field.

# Summary

SSD lifecycle

New

Access
depe    ance

**Throttling SSD usage** helps mitigate temperature-induced errors.

*Temperature*

# Summary

SSD lifecycle

**Access pattern dependence**

We quantify the effects of the **page cache** and **write amplification** in the field.

Temperature

A Large-Scale Study of
# Flash Memory Errors in the Field

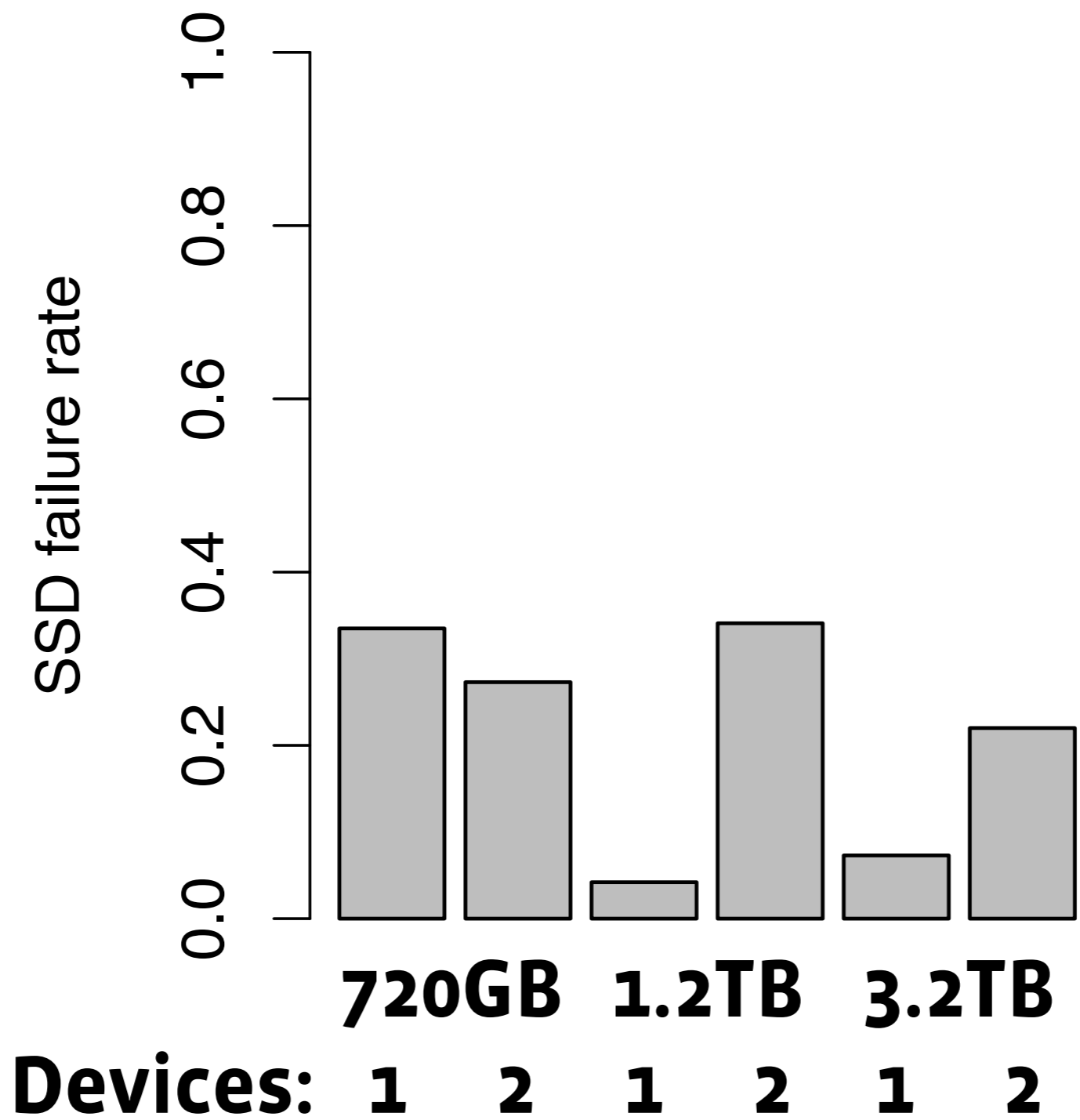**Justin Meza**
Qiang Wu
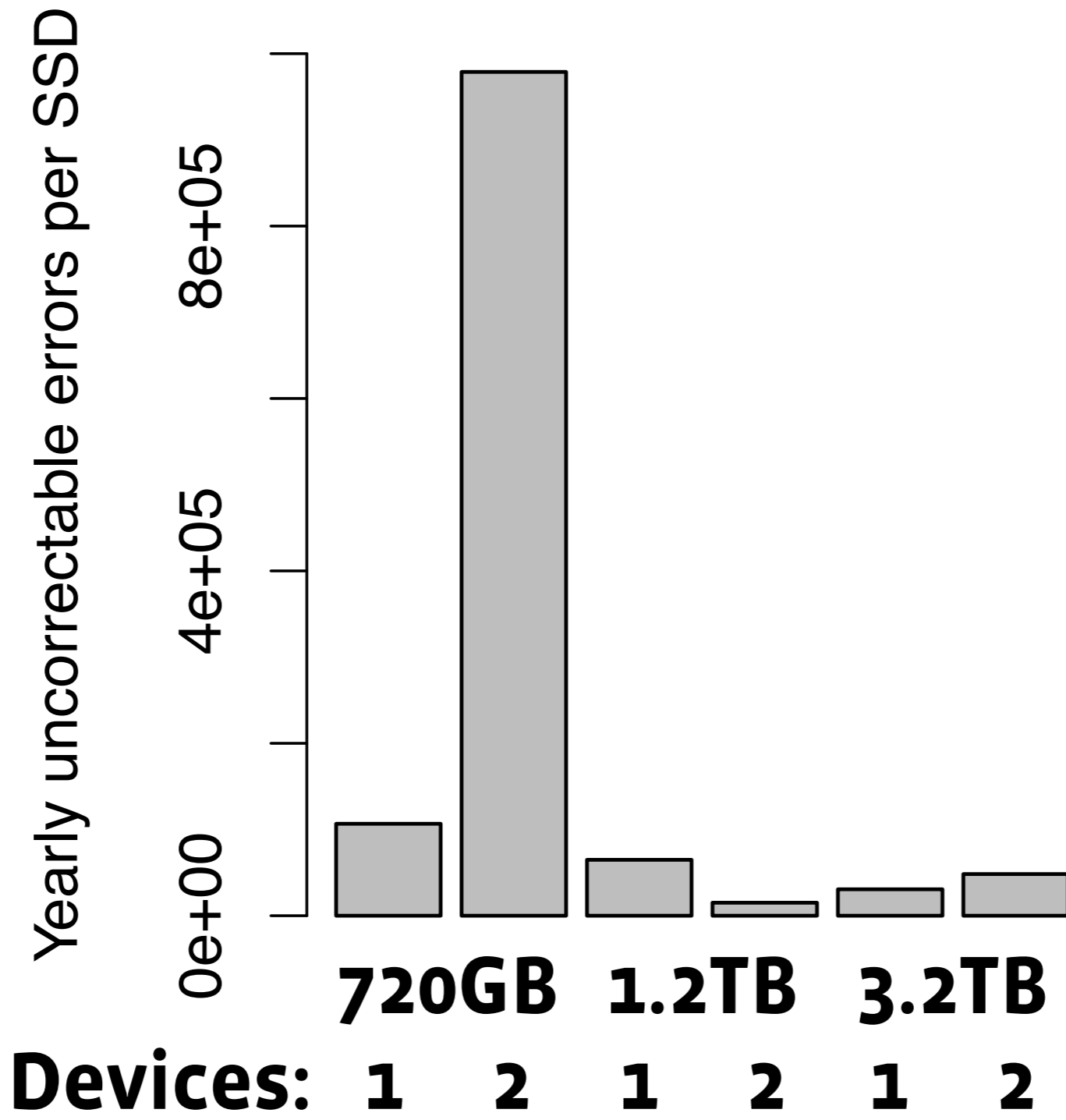Sanjeev Kumar
Onur Mutlu

facebook

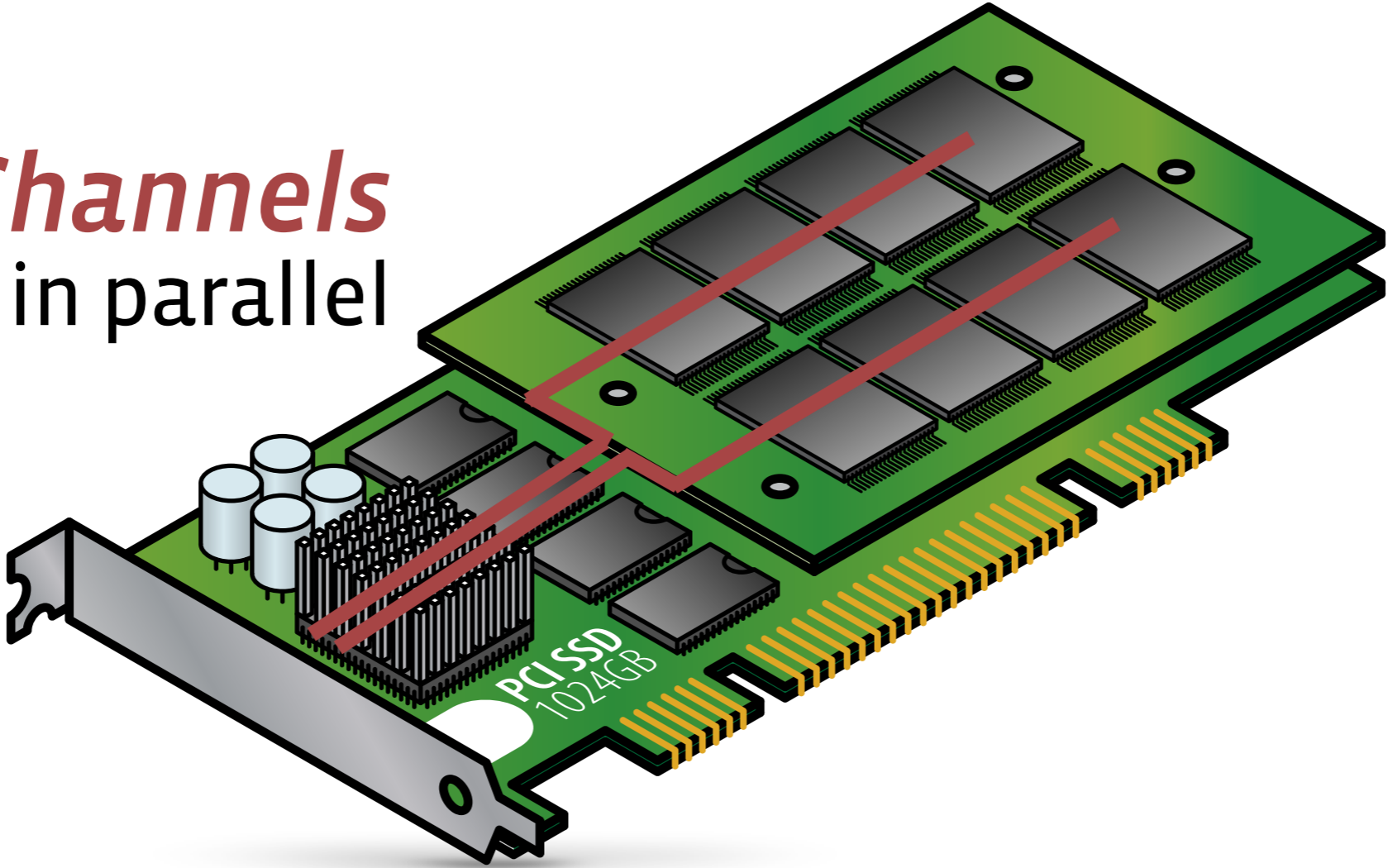Carnegie Mellon University

# Backup slides

# System characteristics

| SSD capacity | PCIe | Average age (years) | SSDs per server | Average written (TB) | Average read (TB) |
|---|---|---|---|---|---|
| 720GB | V1, X4 | 2.4 | 1 | 27.2 | 23.8 |
| | | | 2 | 48.5 | 45.1 |
| 1.2TB | V2, X4 | 1.6 | 1 | 37.8 | 43.4 |
| | | | 2 | 18.9 | 30.6 |
| 3.2TB | V2, X4 | 0.5 | 1 | 23.9 | 51.1 |
| | | | 2 | 14.8 | 18.2 |

*Channels*
operate in parallel

PCI SSD
1024GB

# DRAM buffer

- stores address translations
- may buffer writes