# Understanding and Improving Device Access Complexity

Asim Kadav
(with Prof. Michael M. Swift)
University of Wisconsin-Madison

# Devices enrich computers



★ **Keyboard**
★ **Sound**
★ **Printer**
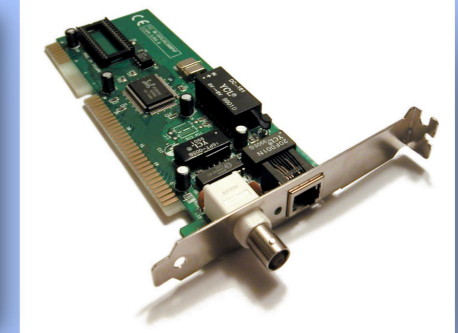★ **Network**
★ **Storage**

# Devices enrich computers

★ **Keyboard**
★ **Sound**
★ **Printer**
★ **Network**
★ **Storage**

★ **Keyboard**
★ **Flash storage**
★ **Graphics**
★ **WIFI**
★ **Headphones**
★ **SD card**
★ **Camera**
★ **Accelerometers**
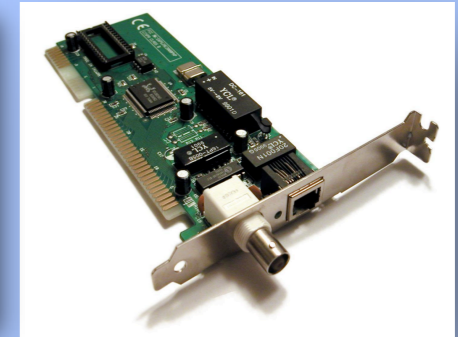★ **GPS**
★ **Touch display**
★ **NFC**

# Huge growth in number of devices

**New I/O devices: accelerometers, GPUS, GPS, touch**

# Huge growth in number of devices

**New I/O devices: accelerometers, GPUS, GPS, touch**
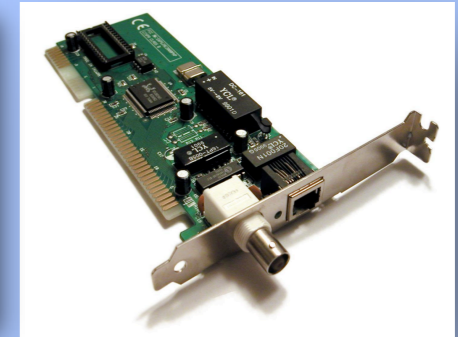
**Many buses: USB, PCI-e, thunderbolt**

# Huge growth in number of devices

**New I/O devices: accelerometers, GPUS, GPS, touch**
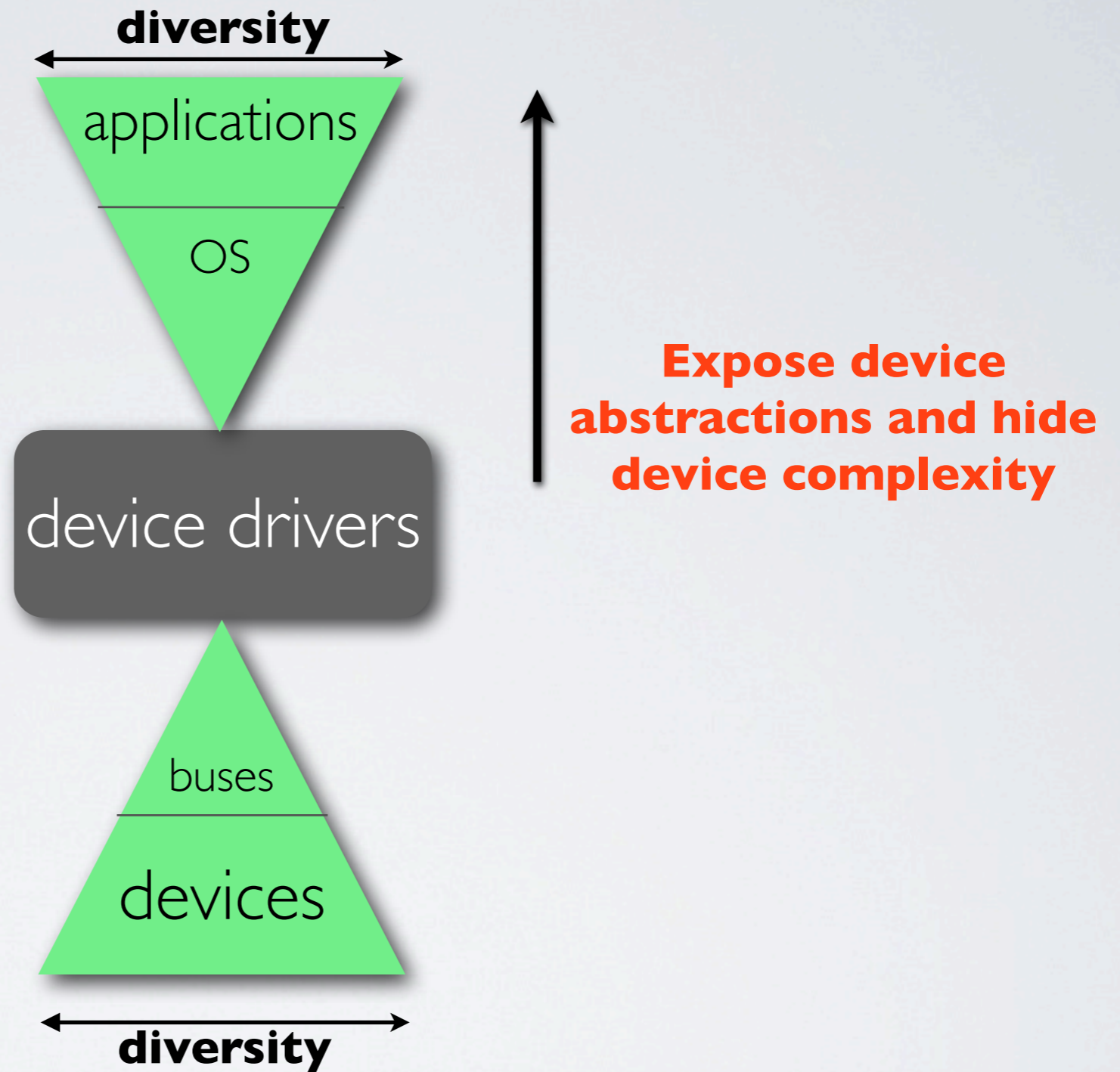
**Many buses: USB, PCI-e, thunderbolt**

**Heterogeneous OS support: 10G ethernet vs card readers**

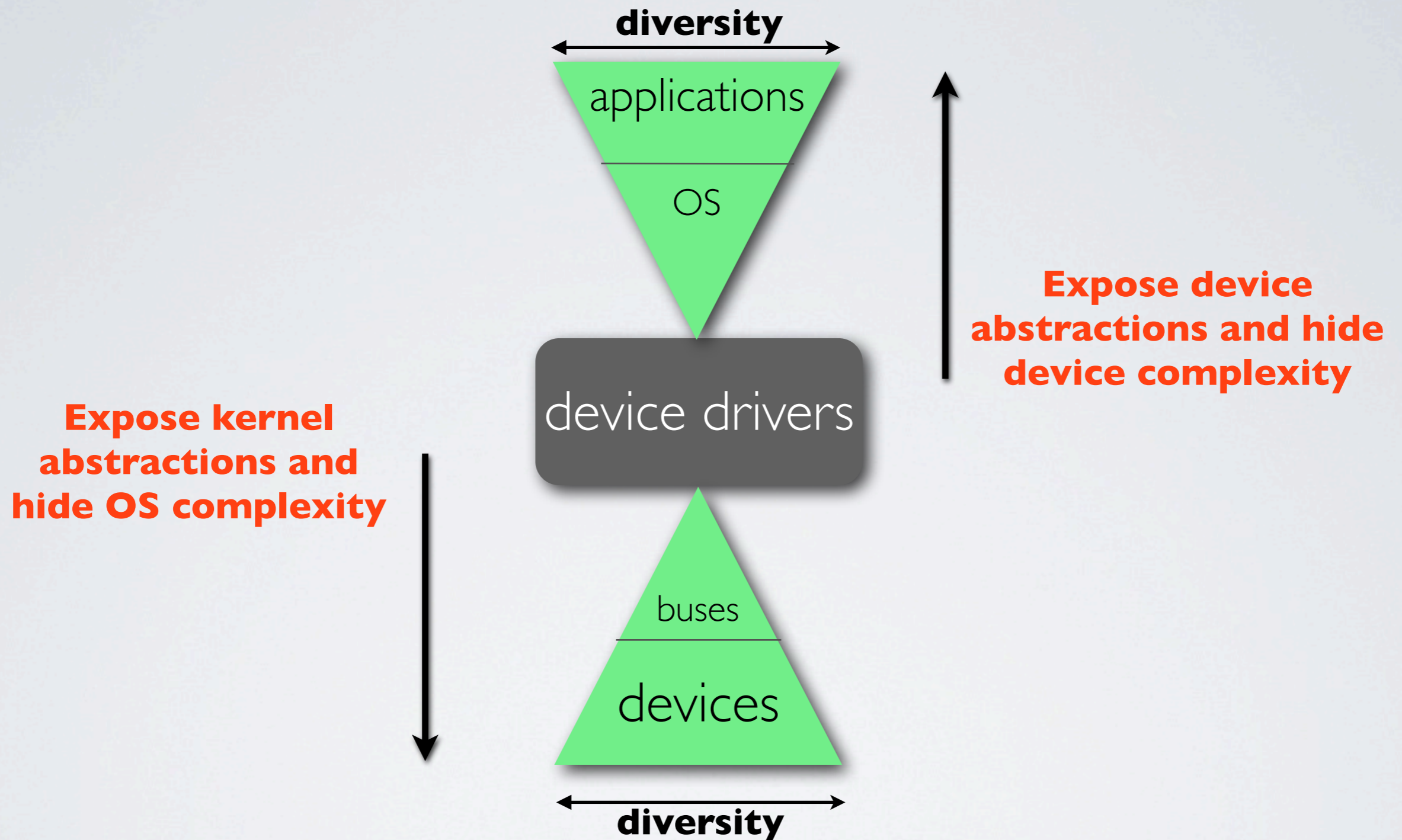# Device drivers: OS interface to devices

# Device drivers: OS interface to devices

**diversity**

applications

OS

device drivers

buses

devices

**diversity**

**Expose device abstractions and hide device complexity**

# Device drivers: OS interface to devices

**diversity**

applications

OS

device drivers

**Expose device abstractions and hide device complexity**

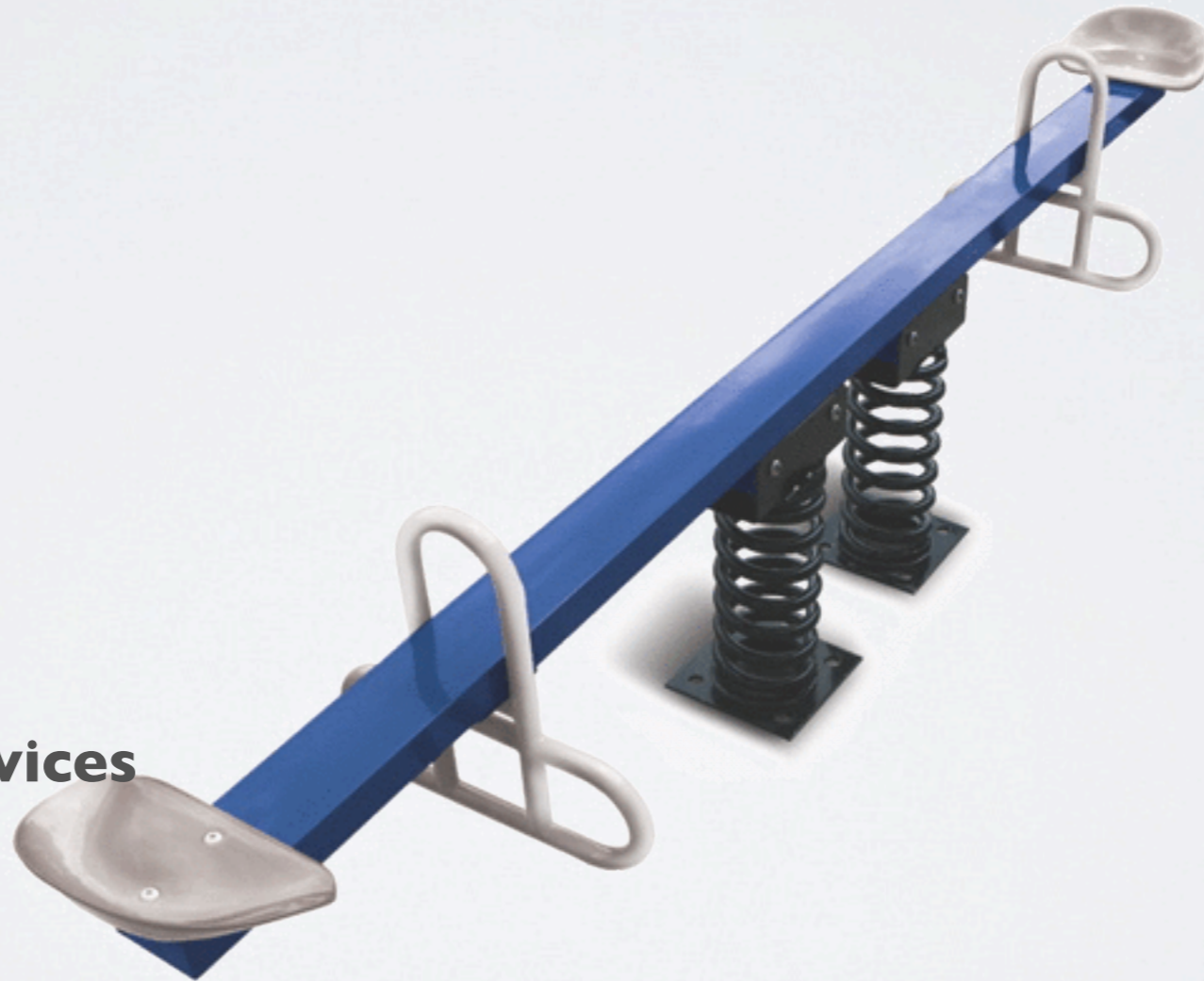**Expose kernel abstractions and hide OS complexity**

buses

devices

**diversity**

**Allow diverse set of applications and OS services to access diverse set of devices**

# Evolution of devices hurts device access

**Efficient device support in OS**

**Evolution of devices**

# Evolution of devices hurts device access

Simplicity
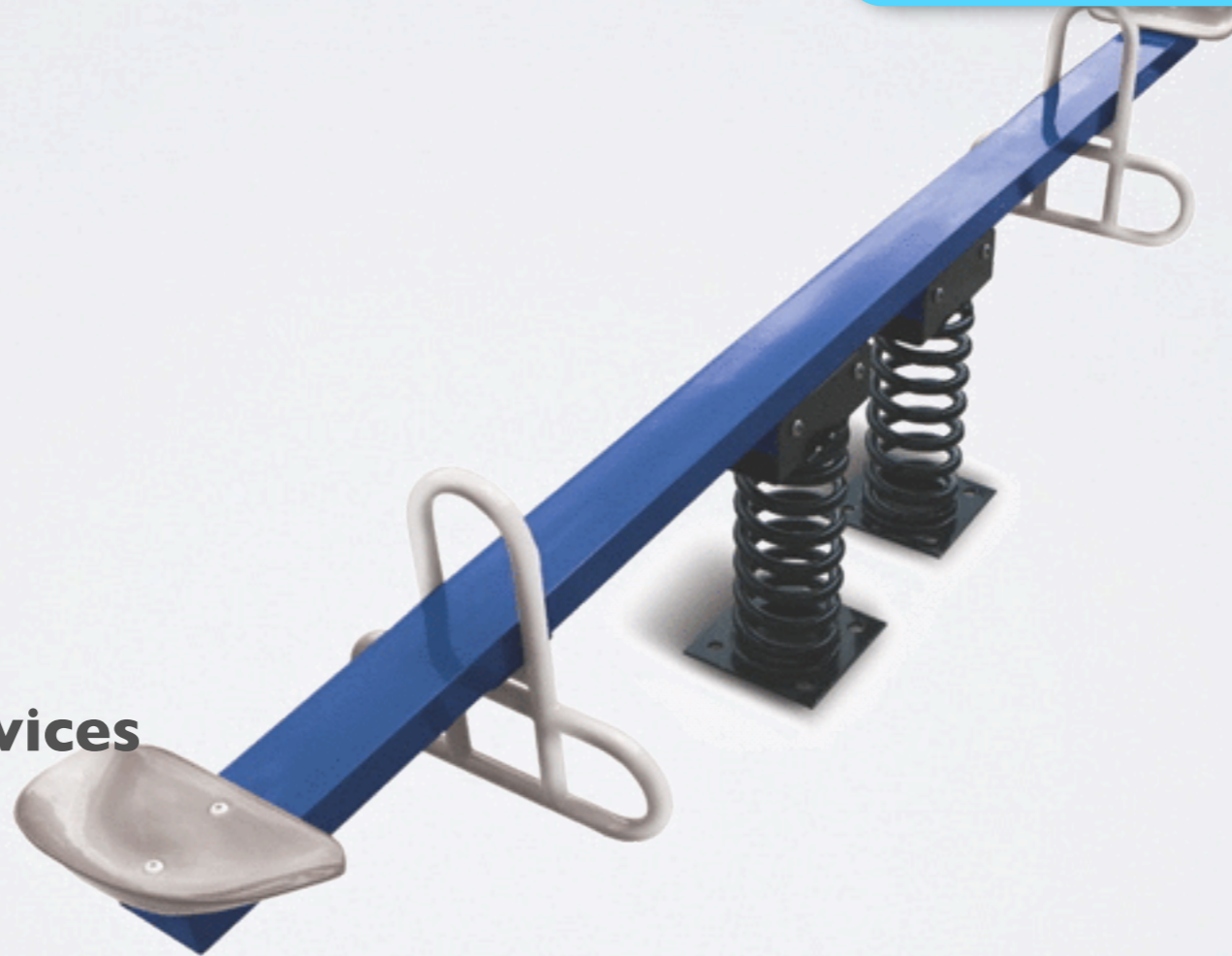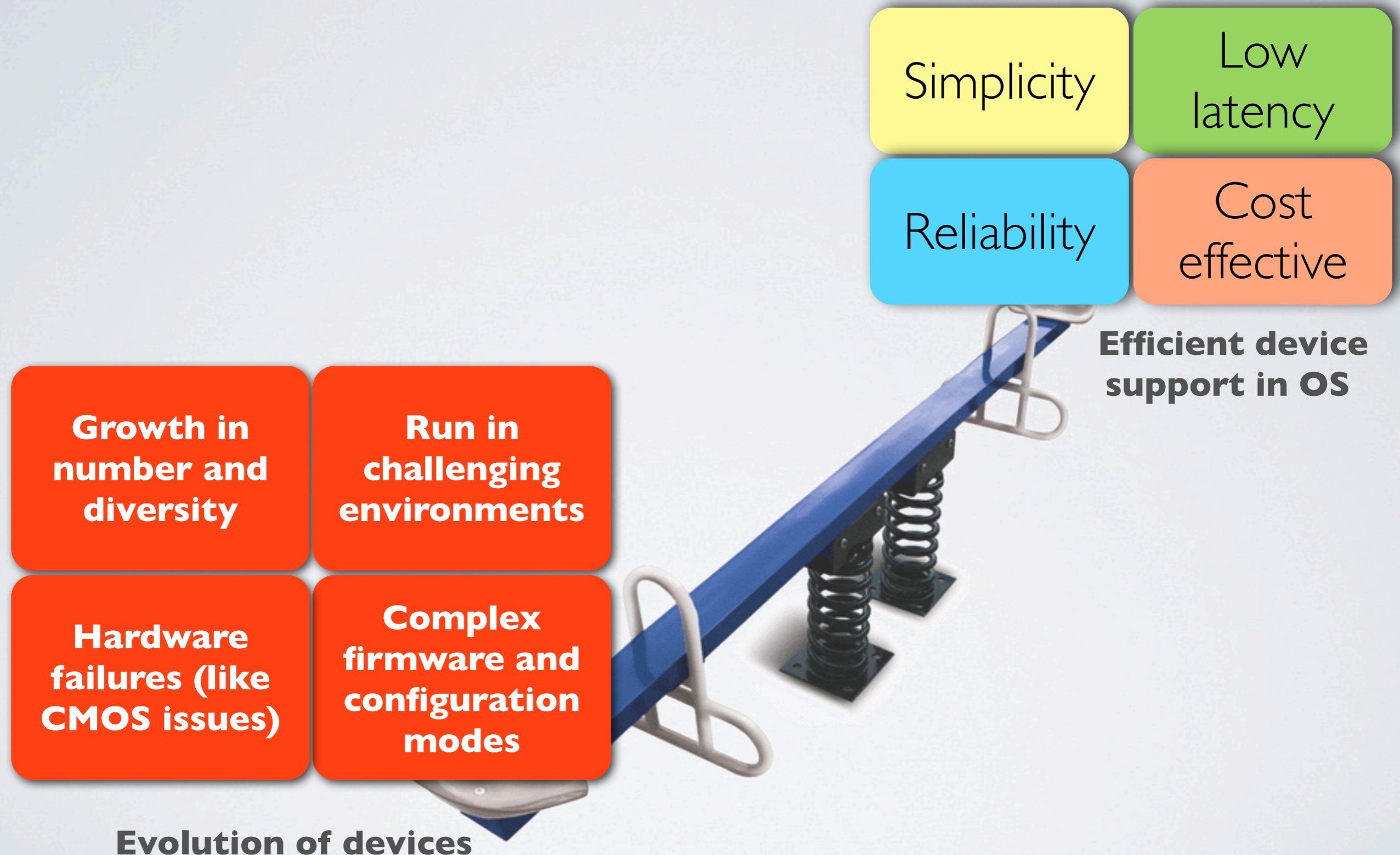
Low latency

Reliability

Cost effective

**Efficient device support in OS**

**Evolution of devices**

# Evolution of devices hurts device access

**Simplicity**

**Low latency**

**Reliability**

**Cost effective**

**Efficient device support in OS**

**Growth in number and diversity**

**Run in challenging environments**

**Hardware failures (like CMOS issues)**

**Complex firmware and configuration modes**

**Evolution of devices**

# Evolution of devices hurts device access

**Tools and mechanisms to address increasing device complexity**

Simplicity

Low latency

Reliability

Cost effective

**Efficient device support in OS**

**Growth in number and diversity**

**Run in challenging environments**

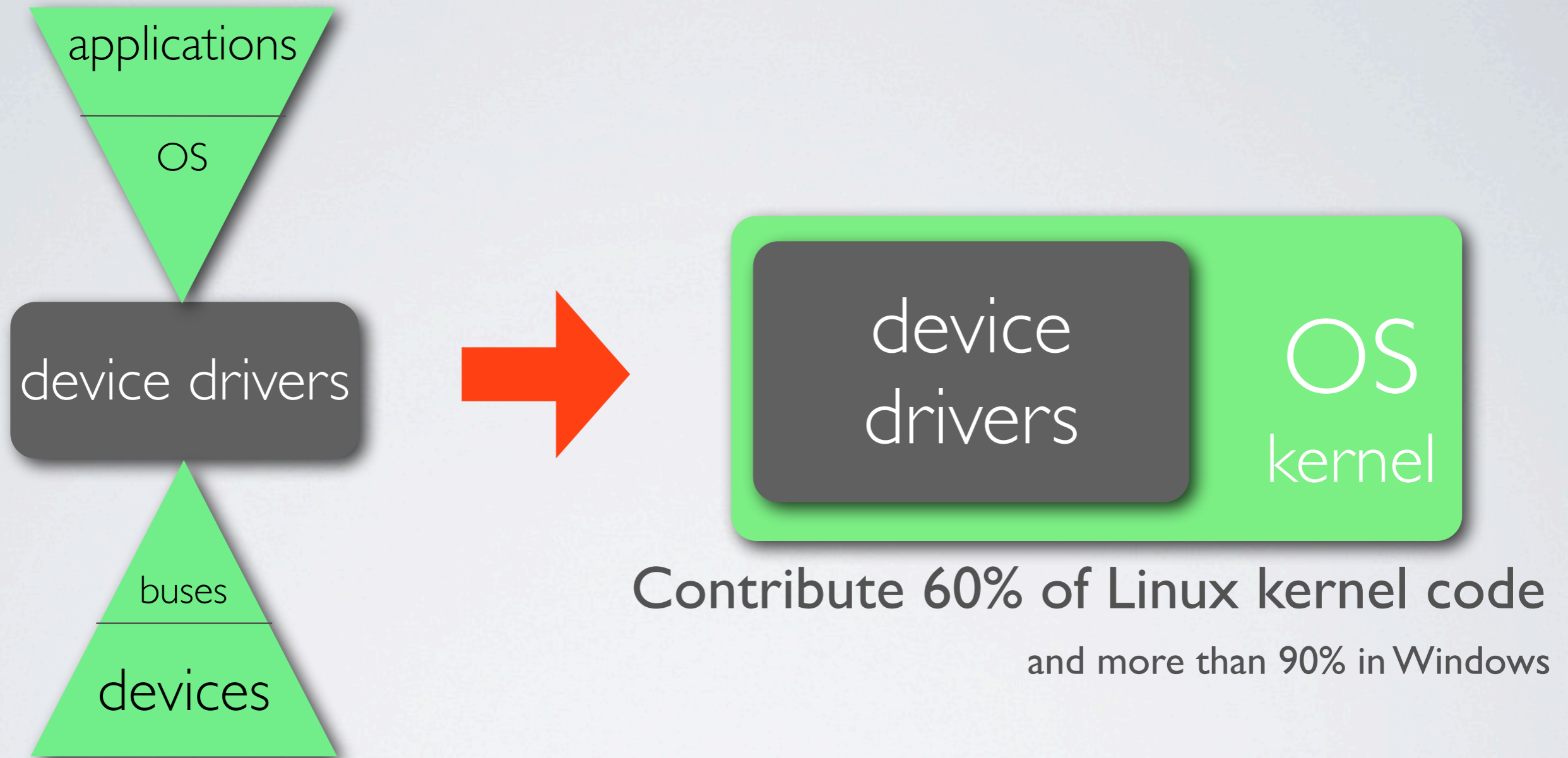**Hardware failures (like CMOS issues)**

**Complex firmware and configuration modes**

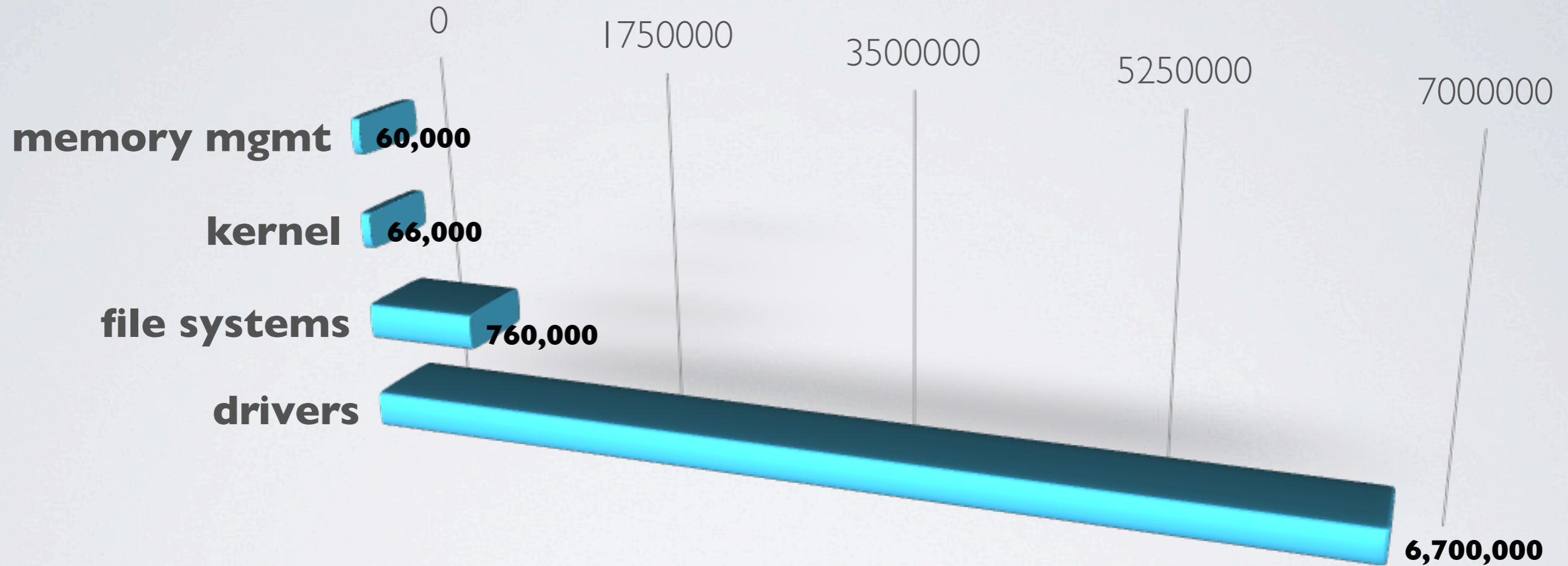**Evolution of devices**

# Growth in drivers hurts understanding of drivers

# Growth in drivers hurts understanding of drivers



Contribute 60% of Linux kernel code

and more than 90% in Windows

# Growth in drivers hurts understanding of drivers

**Lines of code in Linux 3.8**

0      1750000      3500000      5250000      7000000

**memory mgmt** — 60,000

**kernel** — 66,000

**file systems** — 760,000

**drivers** — 6,700,000

# Growth in drivers hurts understanding of drivers

**Lines of code in Linux 3.8**

| | |
|---|---|
| memory mgmt | 60,000 |
| kernel | 66,000 |
| file systems | 760,000 |
| drivers | 6,700,000 |

Axis values: 0, 1750000, 3500000, 5250000, 7000000

**Understand the software complexity and improve driver code**

# Last decade: Focus on the driver-kernel interface

+

device drivers

OS kernel

3rd party developers

# Last decade: Focus on the driver-kernel interface



3rd party developers

+

device drivers

OS kernel

# Last decade: Focus on the driver-kernel interface



3rd party developers

\+

device drivers

OS kernel

Recipe for disaster

# Re-use lessons from existing driver research

| Improvement | System | Validation | | |
|---|---|---|---|---|
| | | Drivers | Bus | Classes |
| New functionality | Shadow driver migration [OSR09] | 1 | 1 | 1 |
| | RevNIC [Eurosys 10] | 1 | 1 | 1 |
| Reliability | Nooks [SOSP 03] | 6 | 1 | 2 |
| | XFI [ OSDI 06] | 2 | 1 | 1 |
| | CuriOS [OSDI 08] | 2 | 1 | 2 |
| Type Safety | SafeDrive [OSDI 06] | 6 | 2 | 3 |
| | Singularity [Eurosys 06] | 1 | 1 | 1 |
| Specification | Nexus [OSDI 08] | 2 | 1 | 2 |
| | Termite [SOSP 09] | 2 | 1 | 2 |
| Static analysis tools | Windows SDV [Eurosys 06] | Many | Many | Many |
| | Coverity [CACM 10] | All | All | All |
| | Cocinelle [Eurosys 08] | All | All | All |

# Re-use lessons from existing driver research

| Improvement | System | Validation | | |
|---|---|---|---|---|
| | | Drivers | Bus | Classes |
| New functionality | Shadow driver migration [OSR09] | 1 | 1 | 1 |
| | RevNIC [Eurosys 10] | 1 | 1 | 1 |
| Reliability | Nooks [SOSP 03] | 6 | 1 | 2 |
| | XFI [OSDI 06] | 2 | 1 | 1 |
| | CuriOS [OSDI 08] | 2 | 1 | 2 |
| Type Safety | SafeDrive [OSDI 06] | 6 | 2 | 3 |
| | Singularity [Eurosys 06] | 1 | 1 | 1 |
| Specification | Nexus [OSDI 08] | 2 | 1 | 2 |
| | Termite [SOSP 09] | 2 | 1 | 2 |

**Large kernel subsystems and validity of few device types result in limited adoption of research solutions**

# Re-use lessons from existing driver research

| Improvement | System | Validation | | |
|---|---|---|---|---|
| | | Drivers | Bus | Classes |
| New functionality | Shadow driver migration [OSR09] | 1 | 1 | 1 |
| | RevNIC [Eurosys 10] | 1 | 1 | 1 |
| Reliability | Nooks [SOSP 03] | 6 | 1 | 2 |
| | XFI [ OSDI 06] | 2 | 1 | 1 |
| | CuriOS [OSDI 08] | 2 | 1 | 2 |
| Type Safety | SafeDrive [OSDI 06] | 6 | 2 | 3 |
| Static analysis tools | Windows SDV [Eurosys 06] | Many | Many | Many |
| | Coverity [CACM 10] | All | All | All |
| | Cocinelle [Eurosys 08] | All | All | All |

**Limited kernel changes + Applicable to lots of drivers => Real Impact**

# Re-use lessons from existing driver research

| Improvement | System | Validation | | |
|---|---|---|---|---|
| | | Drivers | Bus | Classes |
| New functionality | Shadow driver migration [OSR09] | 1 | 1 | 1 |
| | RevNIC [Eurosys 10] | 1 | 1 | 1 |
| Reliability | Nooks [SOSP 03] | 6 | 1 | 2 |
| | XFI [ OSDI 06] | 2 | 1 | 1 |
| | CuriOS [OSDI 08] | 2 | 1 | 2 |
| Type Safety | SafeDrive [OSDI 06] | 6 | 2 | 3 |
| Static analysis tools | Windows SDV [Eurosys 06] | Many | Many | Many |

**Limited kernel changes + Applicable to lots of drivers => Real Impact**

**Design goal: Complete solution that limits kernel changes and applies to all drivers**

# Goal: Address software and hardware complexity

★ **Understand and improve device access in the face of rising hardware and software complexity**

# Goal: Address software and hardware complexity

★ **Understand and improve device access in the face of rising hardware and software complexity**

Increasing hardware complexity

Reliability against hardware failures

# Goal: Address software and hardware complexity

★ **Understand and improve device access in the face of rising hardware and software complexity**

Increasing hardware complexity

Reliability against hardware failures

1

Increasing hardware complexity

Low latency device availability

2

# Goal: Address software and hardware complexity

★ **Understand and improve device access in the face of rising hardware and software complexity**

Increasing hardware complexity

Reliability against hardware failures

1

Increasing hardware complexity

Low latency device availability

2

Increasing software complexity

Better understanding of driver code

3

# Contributions/Outline

First research consideration of hardware failures in drivers

**SOSP '09**

Tolerate device failures

Largest study of drivers to understand their behavior and verify research assumptions

**ASPLOS '12**

Understand drivers and potential opportunities

Introduce checkpoint/restore in drivers for low latency fault tolerance

**ASPLOS '13**

Transactional approach for low latency recovery

# What happens when devices misbehave?

# What happens when devices misbehave?

★ **Drivers make it better**

# What happens when devices misbehave?

★ **Drivers make it better**
★ **Drivers make it worse**
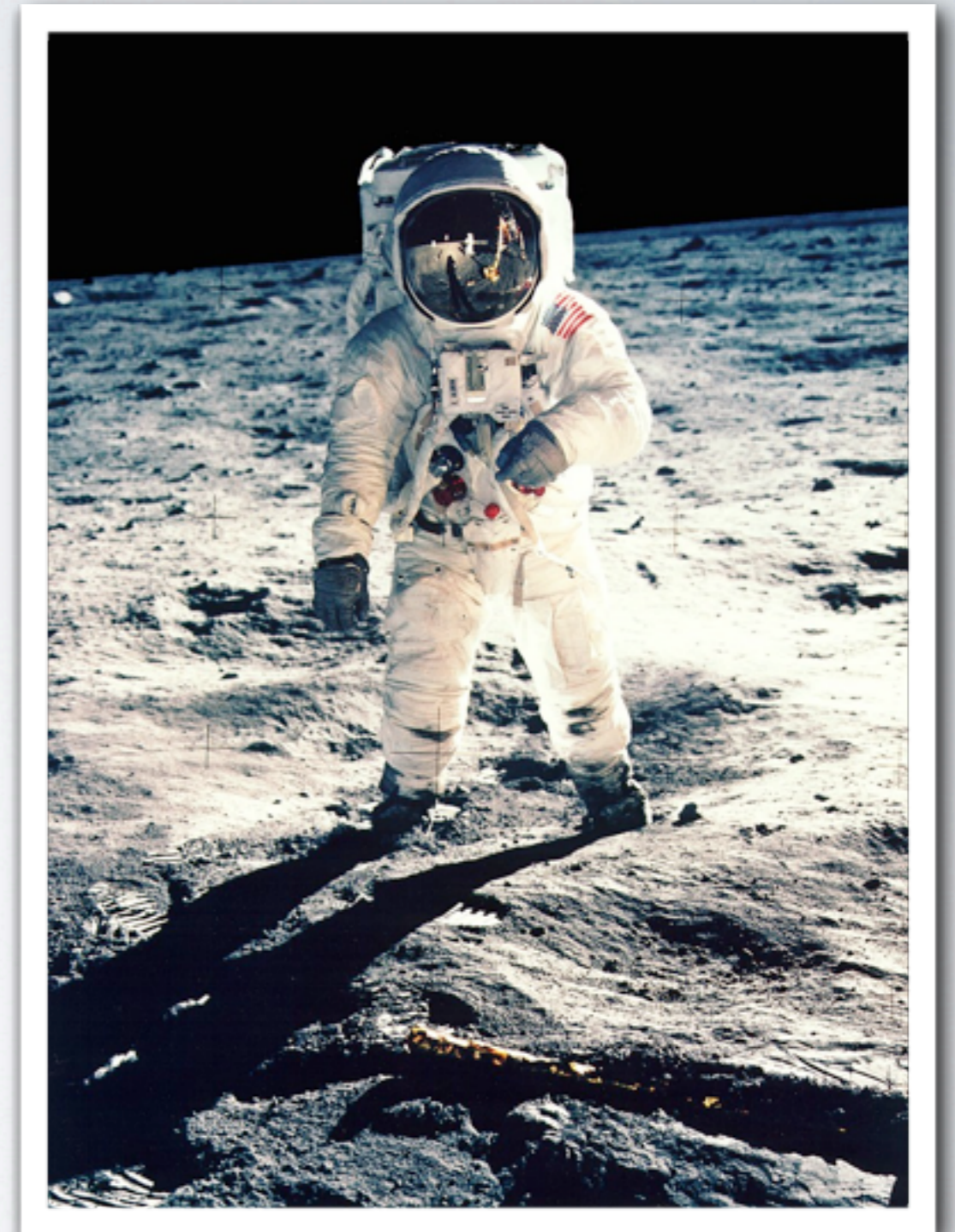
# What happens when devices misbehave?

* **Drivers make it better**
* **Drivers make it worse**

**Early example: Apollo 11 1969**

* **Hardware design bug almost aborted the landing**

* **Assumptions about antenna in driver led to extra CPU**

* **Scientists on-board had to manually prioritize critical tasks**

# Current state of OS-hardware interaction
## **2013**

# Current state of OS-hardware interaction
## 2013

★ **Many device drivers often assume device perfection**
  **- Common Linux network driver: 3c59x.c**

# Current state of OS-hardware interaction
# **2013**

★ **Many device drivers often assume device perfection**
 **- Common Linux network driver: 3c59x.c**

```
while (ioread16(ioaddr + Wn7_MasterStatus))
              & 0x8000);
```

# Current state of OS-hardware interaction
# 2013

★ **Many device drivers often assume device perfection**
  **- Common Linux network driver: 3c59x.c**

```
while (ioread16(ioaddr + Wn7_MasterStatus))
        & 0x8000);
```

HANG!

# Current state of OS-hardware interaction
# 2013

★ **Many device drivers often assume device perfection**
  **- Common Linux network driver: 3c59x.c**

```
while (ioread16(ioaddr + Wn7_MasterStatus))
        & 0x8000);
```
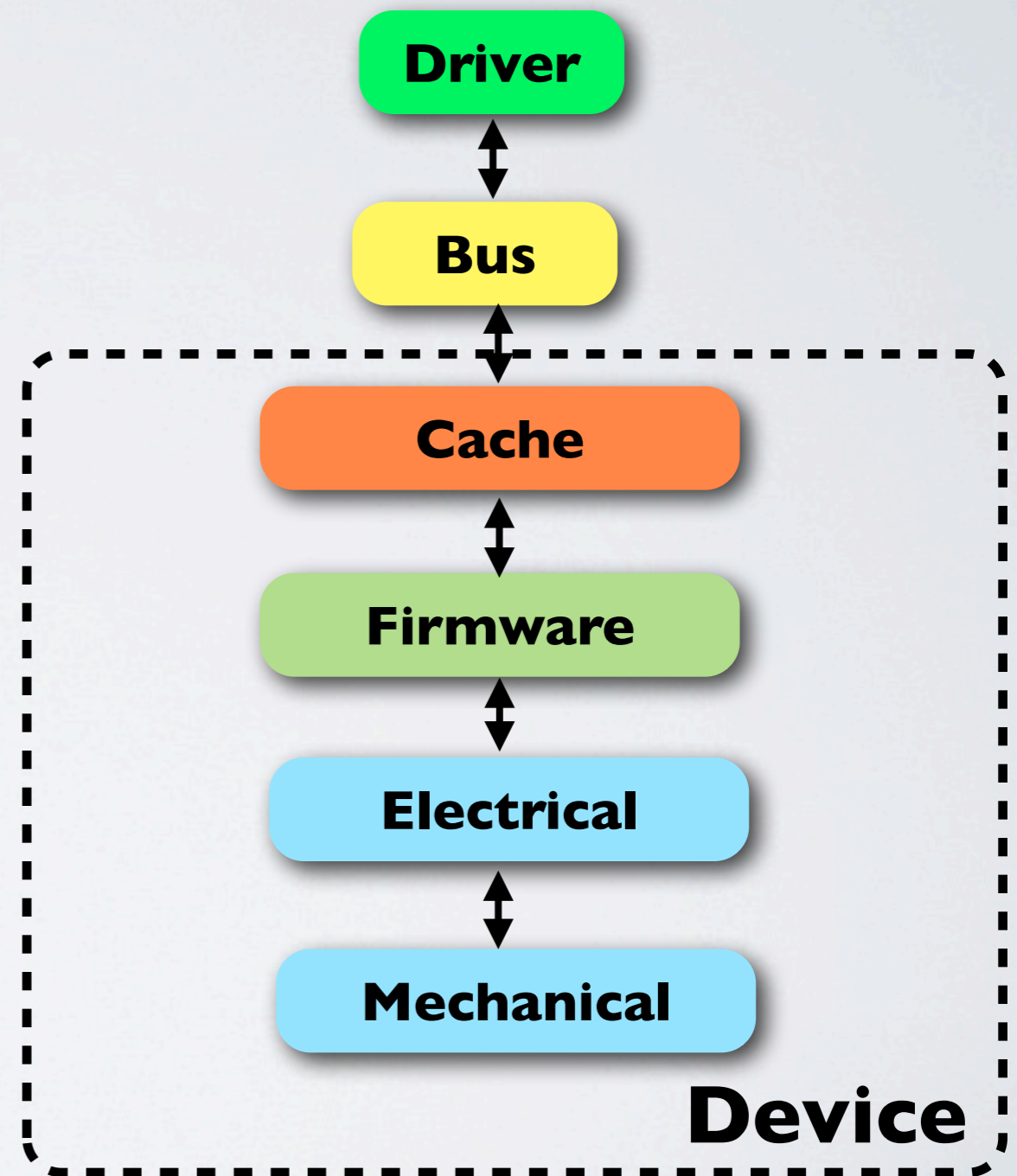
HANG!

**Hardware dependence bug: Device malfunction can crash the system**

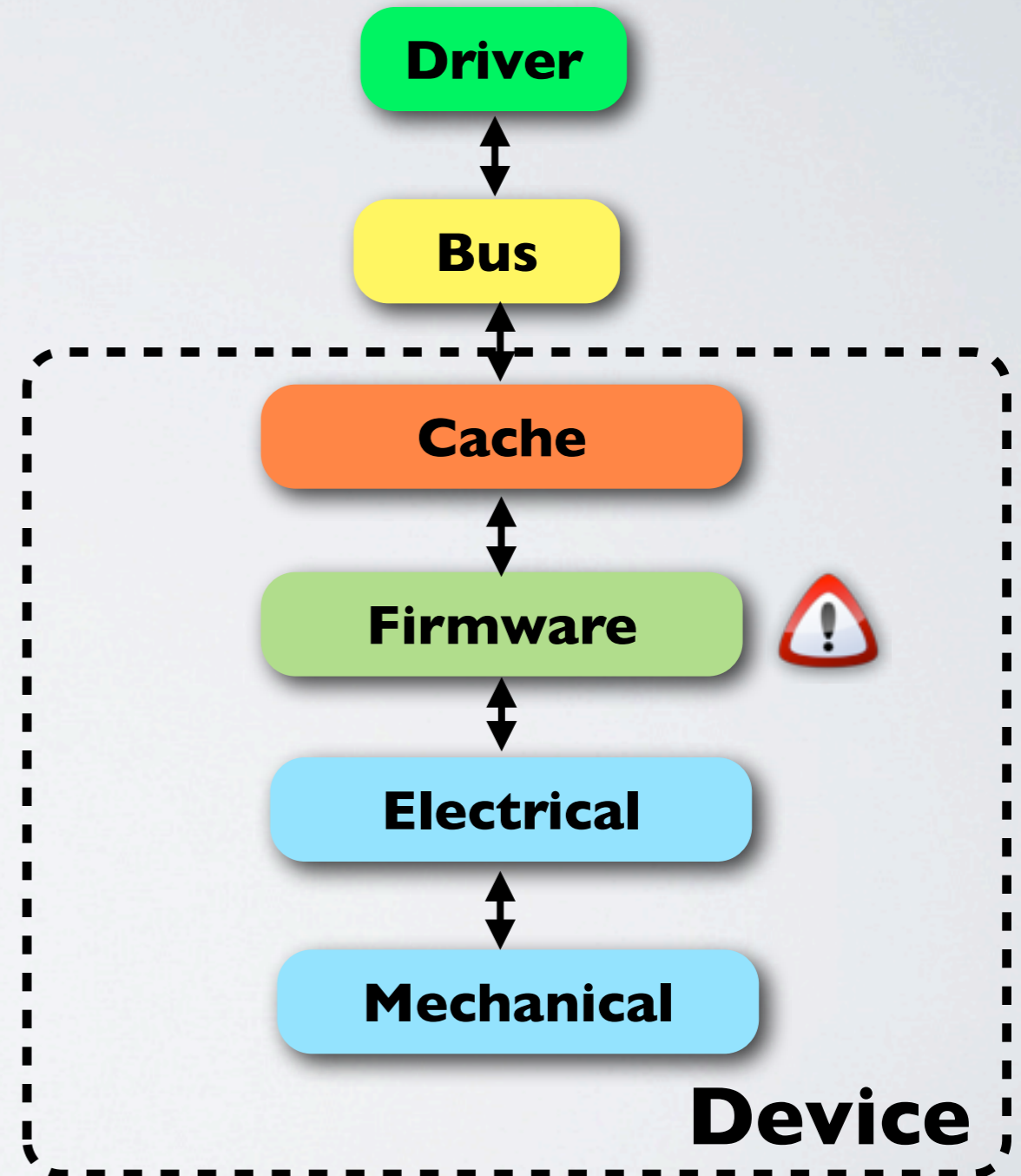# Sources of hardware misbehavior

★ **Sources of hardware misbehavior**

# Sources of hardware misbehavior

★ **Sources of hardware misbehavior**

★ **Firmware/Design bugs**

# Sources of hardware misbehavior

★ **Sources of hardware misbehavior**

★ **Firmware/Design bugs**

★ **Device wear-out, insufficient burn-in**

★ **Bridging faults**

# Sources of hardware misbehavior

★ **Sources of hardware misbehavior**

- ★ **Firmware/Design bugs**
- ★ **Device wear-out, insufficient burn-in**
- ★ **Bridging faults**
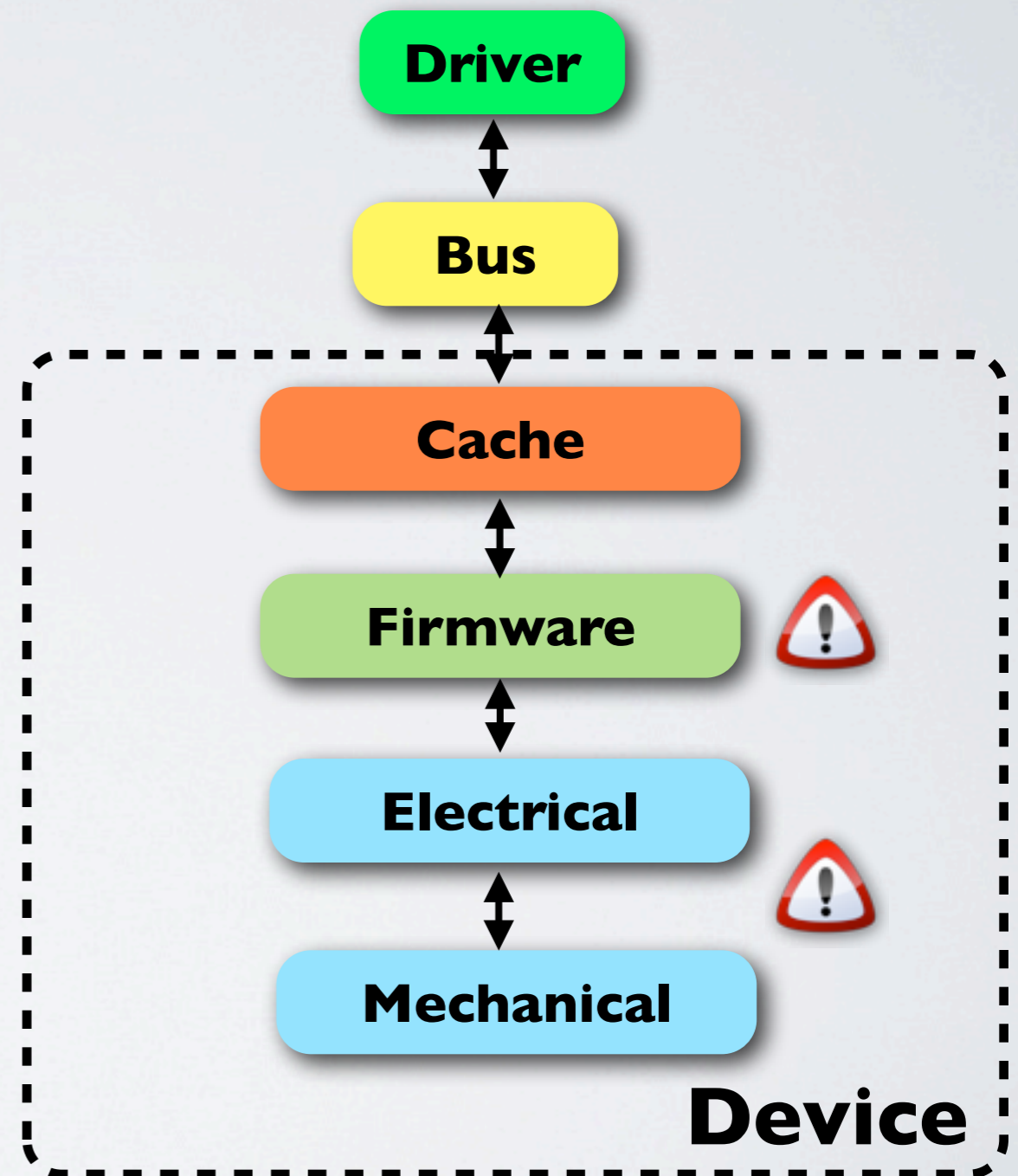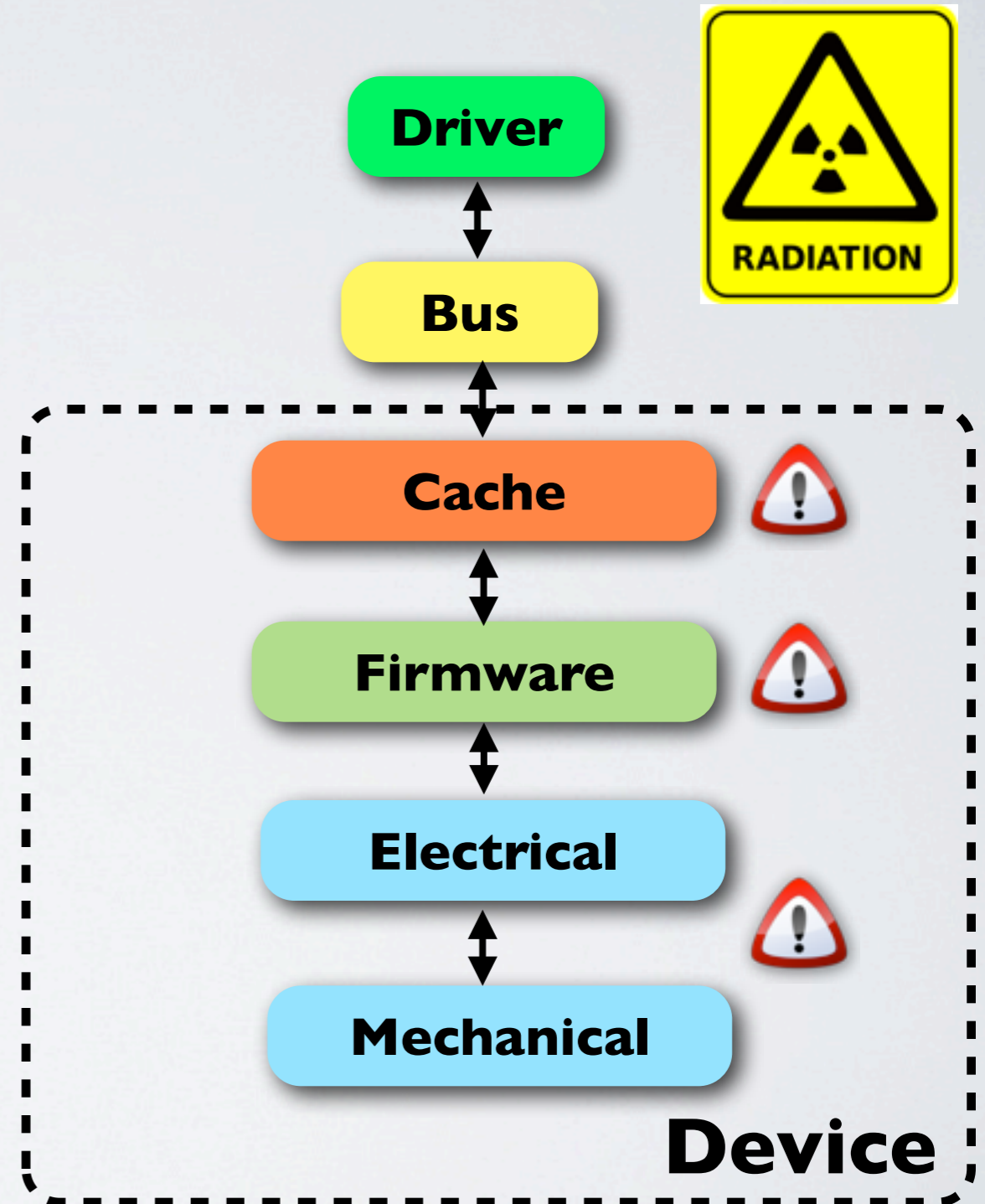- ★ **Electromagnetic interference, radiation, heat**

# Sources of hardware misbehavior

★ **Sources of hardware misbehavior**

  ★ **Firmware/Design bugs**

  ★ **Device wear-out, insufficient burn-in**

  ★ **Bridging faults**

  ★ **Electromagnetic interference, radiation, heat**

★ **Results of misbehavior**

  ★ **Corrupted/stuck-at inputs**

  ★ **Timing errors**

  ★ **Interrupt storms/missing interrupts**

  ★ **Incorrect memory access**

# An evidence:

# An evidence:



[1] Fault resilient drivers for Longhorn server, May 2004. Microsoft Corp.

# An evidence:

Transient hardware failures caused **8%** of all crashes and **9%** of all unplanned reboots [1]

[1] Fault resilient drivers for Longhorn server, May 2004. Microsoft Corp.

# An evidence:

**Windows Server**

Transient hardware failures caused **8%** of all crashes and

**9%** of all unplanned reboots [1]
* ★ Systems work fine after reboots
* ★ Vendors report returned device was faultless

[1] Fault resilient drivers for Longhorn server, May 2004. Microsoft Corp.

# An evidence:

**Windows Server**

Transient hardware failures caused **8%** of all crashes and

**9%** of all unplanned reboots [1]
- ★ Systems work fine after reboots
- ★ Vendors report returned device was faultless

Existing solution is **hand-coded** hardened drivers

- ★ Crashes reduce from **8%** to **3%**

[1] Fault resilient drivers for Longhorn server, May 2004. Microsoft Corp.

# How do hardware dependence bugs manifest?

# How do hardware dependence bugs manifest?

1

Drivers use device
data in critical control and data paths

```
printk("%s",msg[inb(regA)]);
```

# How do hardware dependence bugs manifest?

**1** Drivers use device
data in critical control and data paths

```
printk("%s",msg[inb(regA)]);
```

**2** Drivers do not report device
malfunction to system log

```
if (inb(regA)!= 5)  {
    return; //do nothing
}
```

# How do hardware dependence bugs manifest?

**1** Drivers use device
data in critical control and data paths

```
printk("%s",msg[inb(regA)]);
```

**2** Drivers do not report device
malfunction to system log

```
if (inb(regA)!= 5)  {
    return; //do nothing
}
```

**3** Drivers do not detect or recover from
device failures

```
if (inb(regA)!= 5) {
    panic();
}
```

# Vendor recommendations for driver developers

| Recommendation | Summary | Recommended by | | | |
|---|---|---|---|---|---|
| | | Intel | Sun | MS | Linux |
| Validation | Input validation | ● | ● | ● | |
| | Read once& CRC data | ● | ● | | ● |
| | DMA protection | ● | ● | | |
| Timing | Infinite polling | ● | ● | ● | |
| | Stuck interrupt | | ● | | |
| | Lost request | | | ● | |
| | Avoid excess delay in OS | | | ● | |
| | Unexpected events | ● | | ● | |
| Reporting | Report all failures | ● | ● | ● | |
| Recovery | Handle all failures | | ● | ● | |
| | Cleanup correctly | ● | ● | | |
| | Do not crash on failure | ● | | ● | ● |
| | Wrap I/O memory access | ● | ● | ● | ● |

# Vendor recommendations for driver developers

| Recommendation | Summary | Recommended by | | | |
|---|---|---|---|---|---|
| | | Intel | Sun | MS | Linux |
| Validation | Input validation | ● | ● | ● | |
| | Read once& CRC data | ● | ● | | ● |
| | DMA protection | ● | ● | | |
| Timing | Infinite polling | ● | ● | ● | |
| | Stuck interrupt | | ● | | |
| | Lost request | | | ● | |
| | Avoid excess delay in OS | | | ● | |
| | Unexpected events | ● | | ● | |
| Reporting | Report all failures | ● | ● | ● | |
| Recovery | Handle all failures | | ● | ● | |
| | Cleanup correctly | ● | ● | | |
| | Do not crash on failure | ● | | ● | ● |
| | Wrap I/O memory access | ● | ● | ● | ● |

# Vendor recommendations for driver developers

| Recommendation | Summary | Recommended by | | | |
|---|---|---|---|---|---|
| | | Intel | Sun | MS | Linux |
| Validation | Input validation | ● | ● | ● | |
| | Read once& CRC data | ● | ● | | ● |
| | DMA protection | ● | ● | | |
| Timing | Infinite polling | ● | ● | ● | |
| | Stuck interrupt | | ● | | |
| | Lost request | | | ● | |
| | Avoid excess delay in OS | | | ● | |
| | Unexpected events | ● | | ● | |
| Reporting | Report all failures | ● | ● | ● | |
| Recovery | Handle all failures | | ● | ● | |
| | Cleanup correctly | ● | ● | | |
| | Do not crash on failure | ● | | ● | ● |
| | Wrap I/O memory access | ● | ● | ● | ● |

# Vendor recommendations for driver developers

| Recommendation | Summary | Recommended by | | | |
|---|---|---|---|---|---|
| | | Intel | Sun | MS | Linux |
| Validation | Input validation | ● | ● | ● | |
| | Read once& CRC data | ● | ● | | ● |
| | DMA protection | ● | ● | | |
| Timing | Infinite polling | ● | ● | ● | |
| | Stuck interrupt | | ● | | |
| | Lost request | | | ● | |
| | Avoid excess delay in OS | | | ● | |
| | Unexpected events | ● | | ● | |
| Reporting | Report all failures | ● | ● | ● | |
| Recovery | Handle all failures | | ● | ● | |
| | Cleanup correctly | ● | ● | | |
| | Do not crash on failure | ● | | ● | ● |
| | Wrap I/O memory access | ● | ● | ● | ● |

# Vendor recommendations for driver developers

| Recommendation | Summary | Recommended by | | | |
|---|---|---|---|---|---|
| | | Intel | Sun | MS | Linux |
| Validation | Input validation | ● | ● | ● | |
| | Read once& CRC data | ● | ● | | ● |
| | DMA protection | ● | ● | | |
| Timing | Infinite polling | ● | ● | ● | |
| | Stuck interrupt | | ● | | |
| | Lost request | | | ● | |
| | Avoid excess delay in OS | | | ● | |
| | Unexpected events | ● | | ● | |
| Reporting | Report all failures | ● | ● | ● | |
| Recovery | Handle all failures | | ● | ● | |
| | Cleanup correctly | ● | ● | | |
| | Do not crash on failure | ● | | ● | ● |
| | Wrap I/O memory access | ● | ● | ● | ● |

# Vendor recommendations for driver developers

| Recommendation | Summary | Recommended by | | | |
|---|---|---|---|---|---|
| | | Intel | Sun | MS | Linux |
| Validation | Input validation | ● | ● | ● | |
| | Read once& CRC data | ● | ● | | ● |
| | DMA protection | ● | ● | | |
| Timing | Infinite polling | ● | ● | ● | |
| | Stuck interrupt | | ● | | |
| | Lost request | | | ● | |
| | Avoid excess delay in OS | | | ● | |
| | Unexpected events | ● | | ● | |
| Reporting | Report all failures | ● | ● | ● | |
| Recovery | Handle all failures | | ● | ● | |

**Goal: Automatically implement as many recommendations as possible in commodity drivers**

# Carburizer [SOSP '09]

Goal: Tolerate hardware device failures in software through hardware failure detection and recovery

# Carburizer [SOSP '09]

Goal: Tolerate hardware device failures in software through hardware failure detection and recovery

**Static analysis component**

★ **Detect and fix hardware dependence bugs**

★ **Detect and generate missing error reporting information**

# Carburizer [SOSP '09]

Goal: Tolerate hardware device failures in software through hardware failure detection and recovery

**Static analysis component**

★ **Detect and fix hardware dependence bugs**

★ **Detect and generate missing error reporting information**

**Runtime component**

★ **Detect interrupt failures**

★ **Provide automatic recovery**

18

# Carburizer architecture

# Carburizer architecture

# Carburizer architecture

# Carburizer architecture

**Bug detection and automatic fix generation**

**Recovery and interrupt watchdog**

**OS Kernel**

Carburizer

Compiler

```
If (c==0) {
.
print ("Driver
init");
}
.
.
```

Driver

```
If (c==0) {
.
print ("Driver init");
}
.
.
```

Hardened Driver Binary

# Carburizer architecture



**Bug detection and automatic fix generation**

**Recovery and interrupt watchdog**

**OS Kernel**

Carburizer

```
If (c==0) {
.
print ("Driver
init");
}
.
.
```

Driver

Compiler

```
If (c==0) {
.
print ("Driver init");
}
.
.
```

Kernel Interface

Carburizer Runtime

Hardened Driver Binary

Faulty Hardware

# Hardening drivers

# Hardening drivers

- **Goal: Remove hardware dependence bugs**
  - ★ **Find driver code that uses data from device**
  - ★ **Ensure driver performs validity checks**

# Hardening drivers

- **Goal: Remove hardware dependence bugs**
  - ★ **Find driver code that uses data from device**
  - ★ **Ensure driver performs validity checks**

- **Carburizer detects and fixes hardware bugs :**

# Hardening drivers

- **Goal: Remove hardware dependence bugs**
  - ★ **Find driver code that uses data from device**
  - ★ **Ensure driver performs validity checks**

- **Carburizer detects and fixes hardware bugs :**

**Infinite polling**

# Hardening drivers

- **Goal: Remove hardware dependence bugs**
  - ★ **Find driver code that uses data from device**
  - ★ **Ensure driver performs validity checks**

- **Carburizer detects and fixes hardware bugs :**

**Infinite polling**

**Unsafe array reference**

# Hardening drivers

- **Goal: Remove hardware dependence bugs**
  - ★ **Find driver code that uses data from device**
  - ★ **Ensure driver performs validity checks**

- **Carburizer detects and fixes hardware bugs :**

**Infinite polling**

**Unsafe array reference**

**Unsafe pointer reference**

# Hardening drivers

- **Goal: Remove hardware dependence bugs**
  - ★ **Find driver code that uses data from device**
  - ★ **Ensure driver performs validity checks**

- **Carburizer detects and fixes hardware bugs :**

**Infinite polling**

**Unsafe array reference**

**Unsafe pointer reference**

**System panic calls**

20

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

**Types of device I/O**

★ **Port I/O :** `inb/inw`
★ **Memory-mapped I/O :** `readl/readw`
★ **DMA buffers**
★ **Data from USB packets**

**OS**

**network card**

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test ()  {
```

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test () {
        a = read1();
```

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test ()  {
        a = read1();
```

Tainted Variables

a

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test ()  {
        a = readl();
        b = inb();
```

Tainted Variables

a

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test () {
    a = readl();
    b = inb();
```

Tainted Variables
_____

a
b

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test ()  {
        a = readl();
        b = inb();
        c = b;
```

Tainted Variables
___

a
b

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test () {
        a = readl();
        b = inb();
        c = b;
```

Tainted Variables

a
b
c

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test () {
    a = readl();
    b = inb();
    c = b;
    d = c + 2;
```

Tainted Variables
_____

a
b
c

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test ()  {
        a = readl();
        b = inb();
        c = b;
        d = c + 2;
```

Tainted Variables

a
b
c
d

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test ()  {
     a = readl();
     b = inb();
     c = b;
     d = c + 2;
     return d;
```

Tainted Variables

a
b
c
d

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test () {
    a = readl();
    b = inb();
    c = b;
    d = c + 2;
    return d;
```

Tainted Variables

a
b
c
d
test()

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test ()  {
        a = readl();
        b = inb();
        c = b;
        d = c + 2;
        return d;
}
```

**Tainted Variables**

a
b
c
d
test()

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test ()   {
      a = readl();
      b = inb();
      c = b;
      d = c + 2;
      return d;
}
int set()       {
```

Tainted Variables

a
b
c
d
test()

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test ()  {
        a = readl();
        b = inb();
        c = b;
        d = c + 2;
        return d;
}
int set()      {
        e = test();
```

Tainted Variables

a
b
c
d
test()

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test ()  {
      a = readl();
      b = inb();
      c = b;
      d = c + 2;
      return d;
}
int set()      {
      e = test();
```

Tainted Variables

a
b
c
d
test()
e

# Finding sensitive code

★ **First pass: Identify tainted variables that contain data from device**

```
int test ()   {
      a = readl();
      b = inb();
      c = b;
      d = c + 2;
      return d;
}
int set()      {
      e = test();
}
```

Tainted Variables
_____

a
b
c
d
test()
e

# Detecting risky uses of tainted variables

★ **Second pass: Identify risky uses of tainted variables**

★ **Example: Infinite polling**

  ★ **Driver waiting for device to enter particular state**

  ★ **Solution: Detect loops where all terminating conditions depend on tainted variables**

  ★ **Extra analyses to existing timeouts**

# Infinite polling

★ **Infinite polling of devices can cause system lockups**

```
static int amd8111e_read_phy(………)
{
 ...
  reg_val = readl(mmio + PHY_ACCESS);
  while (reg_val & PHY_CMD_ACTIVE)
        reg_val = readl(mmio + PHY_ACCESS);
  ...
}
```

AMD 8111e network driver(amd8111e.c)

# Hardware data used in array reference

★ **Tainted variables used as array indexes**
★ **Detect existing range/not NULL checks**

```
static void __init attach_pas_card(...)
{
    if ((pas_model = pas_read(0xFF88)))
    {
        ...
        sprintf(temp, "%s rev %d",
          pas_model_names[(int) pas_model], pas_read(0x2789));
        ...
}
```

Pro Audio Sound driver (pas2_card.c)

# Analysis results over the Linux kernel

| Driver class | Infinite polling | Static array | Dynamic array | Panic calls |
|---|---|---|---|---|
| net | 117 | 2 | 21 | 2 |
| scsi | 298 | 31 | 22 | 121 |
| sound | 64 | 1 | 0 | 2 |
| video | 174 | 0 | 22 | 22 |
| other | 381 | 9 | 57 | 32 |
| Total | 860 | 43 | 89 | 179 |

* ★ Analyzed/Built 6300 driver files (2.8 million LOC) in 37 min
* ★ Found **992** hardware dependence bugs in driver code
* ★ False positive rate: 7.4% (manual sampling of 190 bugs)

# Analysis results over the Linux kernel

| Driver class | Infinite polling | Static array | Dynamic array | Panic calls |
|---|---|---|---|---|
| net | 117 | 2 | 21 | 2 |
| scsi | 298 | 31 | 22 | 121 |
| sound | | | | |
| video | | | | 2 |
| other | 381 | 9 | 57 | 32 |
| Total | 860 | 43 | 89 | 179 |

**Lightweight and usable technique to find hardware dependence bugs**

- ★ Analyzed/Built 6300 driver files (2.8 million LOC) in 37 min
- ★ Found **992** hardware dependence bugs in driver code
- ★ False positive rate: 7.4%  (manual sampling of 190 bugs)

# Repairing drivers

* **Carburizer automatically generates repair code**
  * **Inserts failure detection and recovery service callout**

# Repairing drivers

* **Carburizer automatically generates repair code**
  * **Inserts failure detection and recovery service callout**

**Infinite polling**

**Unsafe array reference**

**Unsafe pointer reference**

**System panic calls**

# Repairing drivers

* **Carburizer automatically generates repair code**
  * **Inserts failure detection and recovery service callout**

**Timeout checks**

**Infinite polling**

**Unsafe array reference**

**Unsafe pointer reference**

**System panic calls**

# Repairing drivers

* **Carburizer automatically generates repair code**
  * **Inserts failure detection and recovery service callout**

**Timeout checks**

**Array bounds check**

**Infinite polling**

**Unsafe array reference**

**Unsafe pointer reference**

**System panic calls**

# Repairing drivers

* **Carburizer automatically generates repair code**
  * **Inserts failure detection and recovery service callout**

**Timeout checks**

**Array bounds check**

**Not null checks**

**Infinite polling**

**Unsafe array reference**

**Unsafe pointer reference**

**System panic calls**

# Runtime fault recovery : Shadow drivers

- **Carburizer calls generic recovery service if check fails**
- **Low cost transparent recovery**
  - ★ **Based on shadow drivers**
  - ★ **Records state of driver at all times**
  - ★ **Transparently restarts and replays recorded state on failure**
- **No isolation required (like Nooks)**

Driver-Kernel Interface

Shadow Driver

Taps

Device Driver

Device

**Swift [OSDI '04]**

# Carburizer automatically fixes infinite loops

```
timeout = rdtscll(start) + (cpu/khz/HZ)*2;
reg_val = readl(mmio + PHY_ACCESS);
while (reg_val & PHY_CMD_ACTIVE)        {
        reg_val = readl(mmio + PHY_ACCESS);

        if (_cur < timeout)
            rdtscll(_cur);
        else
            __recover_driver();



}
```

**Timeout code added**

`AMD 8111e network driver(amd8111e.c)`

*Code simplified for presentation purposes

# Carburizer automatically adds bounds checks

```
static void __init attach_pas_card(...)
{

  if ((pas_model = pas_read(0xFF88)))
  {
    ...
    if ((pas_model< 0)) || (pas_model>= 5))
        __recover_driver();
    ...
    sprintf(temp, "%s rev %d",
      pas_model_names[(int) pas_model], pas_read(0x2789));

}
```

**Array bounds detected and check added**

Pro Audio Sound driver (pas2_card.c)

*Code simplified for presentation purposes

# Fault injection and performance

- ★ **Synthetic fault injection on network drivers**

# Fault injection and performance

★ **Synthetic fault injection on network drivers**

| Device/ Driver | Original Driver | | Carburizer | | |
|---|---|---|---|---|---|
| | Behavior | Detection | Behavior | Detection | Recovery |
| 3COM 3C905 | CRASH | None | RUNNING | Yes | Yes |
| DEC DC 21x4x | CRASH | None | RUNNING | Yes | Yes |

# Fault injection and performance

★ **Synthetic fault injection on network drivers**

| Device/ Driver | Original Driver | | Carburizer | | |
|---|---|---|---|---|---|
| | Behavior | Detection | Behavior | Detection | Recovery |
| 3COM 3C905 | CRASH | None | RUNNING | Yes | Yes |
| DEC DC 21x4x | CRASH | None | RUNNING | Yes | Yes |

★ **< 0.5% throughput overhead and no CPU overhead with network drivers**

# Fault injection and performance

★ **Synthetic fault injection on network drivers**

| Device/ Driver | Original Driver | | Carburizer | | |
|---|---|---|---|---|---|
| | Behavior | Detection | Behavior | Detection | Recovery |
| 3COM 3C905 | CRASH | None | RUNNING | Yes | Yes |
| DEC DC 21x4x | CRASH | None | RUNNING | Yes | Yes |

★ **< 0.5% throughput overhead and no CPU overhead with network drivers**

**Carburizer failure detection and transparent recovery works and has very low overhead**

# Summary

| Recommendation | Summary | Recommended by | | | | Carburizer Ensures |
|---|---|---|---|---|---|---|
| | | Intel | Sun | MS | Linux | |
| Validation | Input validation | ● | ● | ● | | ● |
| | Read once& CRC data | ● | ● | | ● | |
| | DMA protection | ● | ● | | | |
| Timing | Infinite polling | ● | ● | ● | | ● |
| | Stuck interrupt | | ● | | | ● |
| | Lost request | | | ● | | ● |
| | Avoid excess delay in OS | | | ● | | |
| | Unexpected events | ● | | ● | | |
| Reporting | Report all failures | ● | ● | ● | | ● |
| Recovery | Handle all failures | | ● | ● | | ● |
| | Cleanup correctly | ● | ● | | | ● |
| | Do not crash on failure | ● | | ● | ● | ● |
| | Wrap I/O memory access | ● | ● | ● | ● | |

# Summary

| Recommendation | Summary | Recommended by | | | | Carburizer Ensures |
|---|---|---|---|---|---|---|
| | | Intel | Sun | MS | Linux | |
| Validation | Input validation | ● | ● | ● | | ● |
| | Read once& CRC data | ● | ● | | ● | |
| | DMA protection | ● | ● | | | |
| Timing | Infinite polling | ● | ● | ● | | ● |
| | Stuck interrupt | | ● | | | ● |
| | Lost request | | | ● | | ● |
| | Avoid excess delay in OS | | | ● | | |
| | Unexpected events | ● | | ● | | |
| Reporting | Report all failures | ● | ● | ● | | ● |
| | Wrap I/O memory access | ● | ● | ● | ● | |

**Carburizer improves system reliability by automatically ensuring that hardware failures are tolerated in software**

# Contributions beyond research

★ **Linux Plumbers Conference [Sep '11]**

★ **LWN Article with paper & list of bugs [Feb '12]**

★ **Released patches to the Linux kernel**

★ **Tool + source available for download at:**

  `http://bit.ly/carburizer`

# Recovery performance: device initialization is slow

★ **Multi-second device probe**

   ★ **Identify device**

   ★ **Cold boot device**

   ★ **Setup device/driver structures**

   ★ **Configuration/Self-test**

Module registration

Allocate device structures

Map BAR and I/O ports

Register device operations

Detect chipset capabilities

Cold boot the device

Verify EEPROM checksum

Self test?

Self test on boot

Set chipset specific ops

Allocate driver structures

Configure device

Device ready for requests

# Recovery performance: device initialization is slow

★ **Multi-second device probe**

   ★ **Identify device**

   ★ **Cold boot device**

   ★ **Setup device/driver structures**

   ★ **Configuration/Self-test**

Module registration

Allocate device structures

Map BAR and I/O ports

Register device operations

Detect chipset capabilities

Cold boot the device

Verify EEPROM checksum

Self test?

Self test on boot

Set chipset specific ops

Allocate driver structures

Configure device

**Device ready for requests**

★ **What does slow device re-initialization hurt?**

   ★ **Fault tolerance: Driver recovery**

   ★ **Virtualization: Live migration, cloning**

   ★ **OS functions: Boot, upgrade**

# Recovery functionality: assumes drivers follow class behavior

Driver-Kernel entrypoints

Shadow Driver

Taps

Device Driver

Device

- ★ **Kernel exports standard entry points for every class (like "packet send" for network class)**

- ★ **Shadow drivers records state by interposing class defined entry points**

- ★ **Recovery = Restart and replay of captured state**

- ★ **Do drivers have additional state?**

# Recovery functionality: assumes drivers follow class behavior

Driver-Kernel entrypoints

Taps

Shadow Driver

Device Driver

Device

★ **Kernel exports standard entry points for every class (like "packet send" for network class)**

★ **Shadow drivers records state by interposing class defined entry points**

★ **Recovery = Restart and replay of captured state**

★ **Do drivers have additional state?**

**How many drivers obey class behavior?**

# Outline

Tolerate device failures

Understand drivers and potential opportunities

**Overview**
**Recovery specific results**

Transactional approach for cheap recovery

# Our view of drivers is narrow

**Drivers
6.7 million LOC in
Linux**

# Our view of drivers is narrow



Drivers
6.7 million LOC in
Linux

Driver
Research
(avg. 2.2
drivers/
system)

# Our view of drivers is narrow



Drivers
6.7 million LOC in
Linux

Driver Research (avg. 2.2 drivers/ system)

Bugs

# Our view of drivers is narrow

# Understanding Modern Device Drivers[ASPLOS 2012]

# Understanding Modern Device Drivers[ASPLOS 2012]

**Study source of all Linux drivers for x86 (~3200  drivers)**

# Understanding Modern Device Drivers[ASPLOS 2012]

**Study source of all Linux drivers for x86 (~3200 drivers)**

**Driver properties**

★ **Code properties**
★ **Verify research assumptions**

# Understanding Modern Device Drivers[ASPLOS 2012]

**Study source of all Linux drivers for x86 (~3200 drivers)**

**Driver properties**

**Driver interaction**

★ **Code properties**
★ **Verify research assumptions**

★ **Driver kernel & device interaction**
★ **Driver architecture**

# Understanding Modern Device Drivers[ASPLOS 2012]

**Study source of all Linux drivers for x86 (~3200 drivers)**

**Driver properties**

**Driver interaction**

**Driver similarity**

★ **Code properties**
★ **Verify research assumptions**

★ **Driver kernel & device interaction**
★ **Driver architecture**

★ **7 million lines of code needed?**

# Study methodology

★ **Static source analysis of 3200 drivers in Linux 2.6.37.6 (May 2011)**

# Study methodology

★ **Static source analysis of 3200 drivers in Linux 2.6.37.6 (May 2011)**

**Driver properties**

★ **Identify driver entry points, kernel and bus callouts**

# Study methodology

★ **Static source analysis of 3200 drivers in Linux 2.6.37.6 (May 2011)**

**Driver properties**

★ **Identify driver entry points, kernel and bus callouts**

★ **Device class, sub-class, chipsets**

# Study methodology

★ **Static source analysis of 3200 drivers in Linux 2.6.37.6 (May 2011)**

**Driver properties**

★ **Identify driver entry points, kernel and bus callouts**
  ★ **Device class, sub-class, chipsets**
  ★ **Bus properties & other properties (like module params)**

# Study methodology

★ **Static source analysis of 3200 drivers in Linux 2.6.37.6 (May 2011)**

**Driver properties**

★ **Identify driver entry points, kernel and bus callouts**
- ★ **Device class, sub-class, chipsets**
- ★ **Bus properties & other properties (like module params)**
- ★ **Driver functions registered as entry points (purpose)**

# Study methodology

★ **Static source analysis of 3200 drivers in Linux 2.6.37.6 (May 2011)**

**Driver properties**

★ **Identify driver entry points, kernel and bus callouts**
   ★ **Device class, sub-class, chipsets**
   ★ **Bus properties & other properties (like module params)**
   ★ **Driver functions registered as entry points (purpose)**

```
#include <nothing>
unsigned main()
{
write : Hello all;
write : I know !;
write : not real;
write : +p ;
return all;
}
```

**For every driver**

**Driver entry points**

**xmit**

**close**

**open**

**probe**

# Study methodology

★ **Static source analysis of 3200 drivers in Linux 2.6.37.6 (May 2011)**

**Driver properties**

★ **Identify driver entry points, kernel and bus callouts**
  ★ **Device class, sub-class, chipsets**
  ★ **Bus properties & other properties (like module params)**
  ★ **Driver functions registered as entry points (purpose)**

**For every driver**

**Driver entry points**

open

xmit    close    probe

# Study methodology

**Driver properties**

★ Identify driver entry points, kernel and bus callouts

**Driver interactions**

★ **Reverse propagate information to aggregate bus, device and kernel behavior**

**Driver entry points**

**xmit**

**open**

**close**

**probe**

**kmalloc**

39

# Study methodology

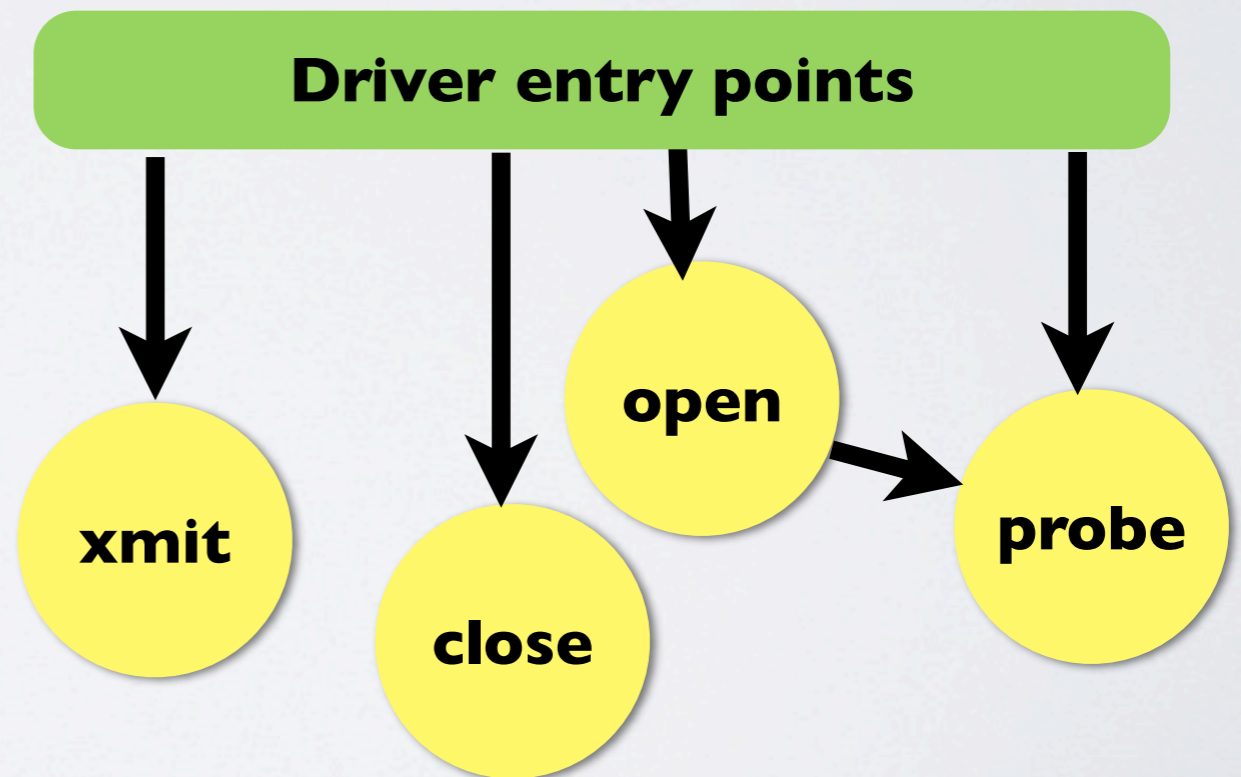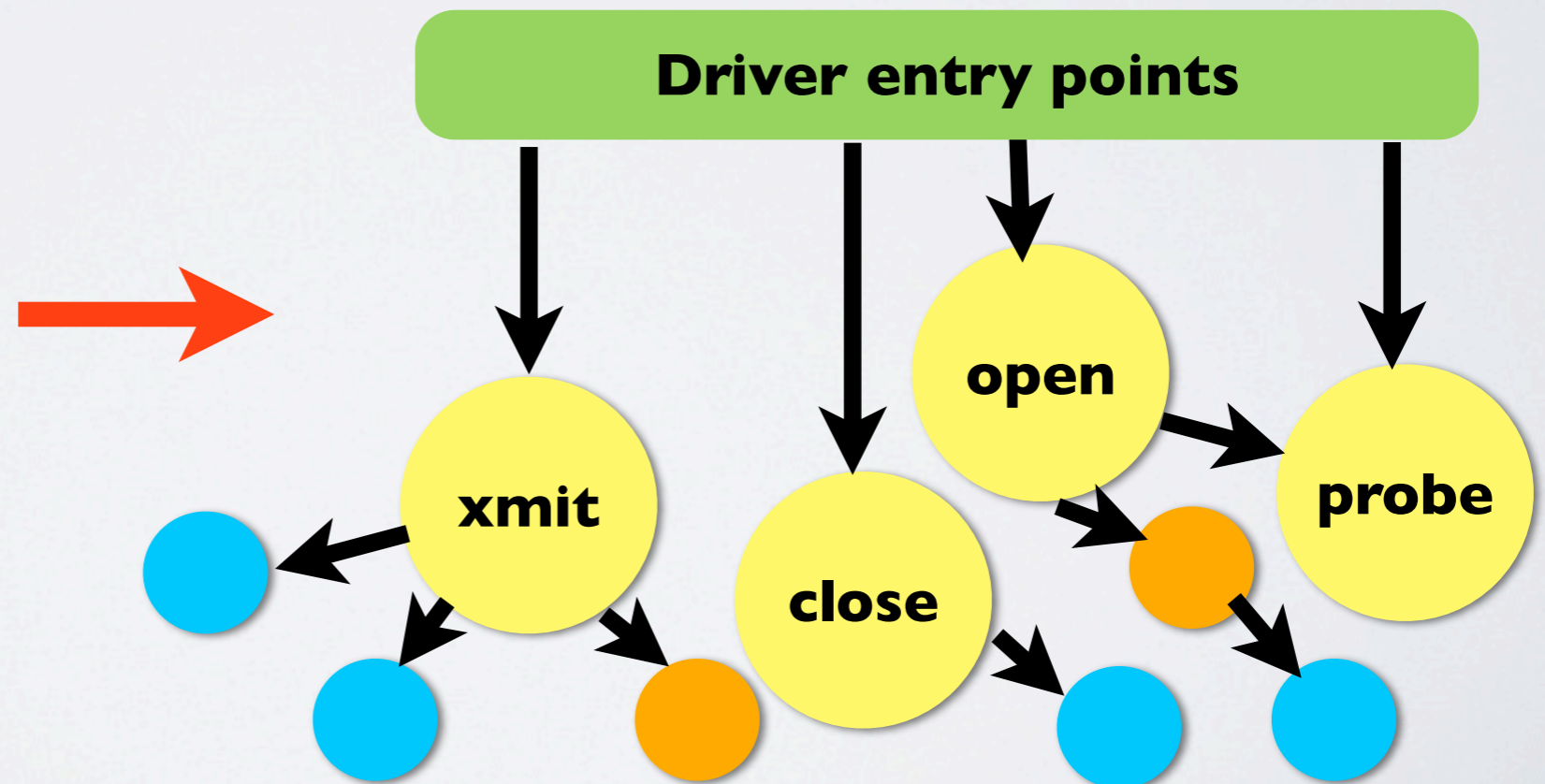★ **Static source analysis of 3200 drivers in Linux 2.6.37.6 (May 2011)**

**Driver properties**

★ Identify driver wide and function specific properties of all drivers

**Driver interactions**

★ Reverse propagate information to aggregate bus, device and kernel behavior

**Driver similarity**

★ **Use statistical clustering techniques and static analysis to identify similar code**

# Contributions/Outline

Tolerate device failures

Understand drivers and potential opportunities

**Overview**
**Recovery specific results**

Transactional approach for cheap recovery

# Driver Code Characteristics

Driver Code Characteristics

42

# Driver Code Characteristics

★ **Initialization/cleanup – 36%**

★ **Core I/O & interrupts – 23%**

★ **Device configuration – 15%**

★ **Power management – 7.4%**

★ **Device ioctl – 6.2%**

Percent-age of LOC

42

# Driver Code Characteristics

★ **Initialization/cleanup – 36%**

★ **Core I/O & interrupts – 23%**

★ **Device configuration – 15%**

★ **Power management – 7.4%**

★ **Device ioctl – 6.2%**

Percent-age of LOC

0

10

**Initialization code dominates driver LOC and adds to complexity**

char drivers: acpi, bluetooth, cdrom, char, crypto, edac, firewire, gpio, gpu, hid, hwmon, input, isdn, leds, media, message, misc, parport, platform, pnp, serial, sound, tty, video, watchdog

block drivers: ata, block, ide, md, mtd, scsi

net drivers: atm, infiniband, net, uwb

Columns: init, cleanup, ioctl, config, power, error, proc, core, intr

# Problem 2: Shadow drivers assume drivers follow class behavior



- ★ **Class definition includes:**
  - ★ **Callbacks registered with the bus, device and kernel subsystem**

# Problem 2: Shadow drivers assume drivers follow class behavior



* **Class definition includes:**
  * **Callbacks registered with the bus, device and kernel subsystem**

**How many drivers follow class behavior and how much code does this add?**

# Problem 2(a): Drivers do behave outside class definitions

★ **Non-class behavior in device drivers:**

- **module parameters, unique ioctls, procfs/sysfs interactions**

DW1520 Wireless-N WLAN Half-Mini Card Properties

General | Advanced | Driver | Details | Power Management

The following properties are available for this network adapter. Click the property you want to change on the left, and then select its value on the right.

Property:

- Disable Upon Wired Connect
- Fragmentation Threshold
- IBSS 54g(tm) Protection Mode
- IBSS Mode
- Locally Administered MAC Address
- Location
- Minimum Power Consumption
- PLCP Header
- Priority & VLAN
- Rate (802.11a)
- Rate (802.11b/g)

Value:

USA

```
$ echo 1 > /sys/class/sound/mixer/
               device/enable
```

**Windows WLAN card
config via private ioctls**

**Linux sound card config via sysfs**

# Problem 2(a): Drivers do behave outside class definitions

★ **Non-class behavior in device drivers:**

**- module parameters, unique ioctls, procfs/sysfs interactions**



```
$ echo 1 > /sys/class/sound/mixer/
              device/enable
```

**Windows WLAN card
config via private ioctls**

**Linux sound card config via sysfs**

**Overall 44% of drivers have non-class behavior and research making this assumption will not apply**

# Problem 2(b): Too many classes



★ **"Understanding Modern Device Drivers"** ASPLOS 2012

# Problem 2(b): Too many classes



**Class-specific driver recovery leads to a large kernel recovery subsystem**

★ **"Understanding Modern Device Drivers"** ASPLOS 2012

# Few other results

**Driver properties**

★ **Many assumptions made by driver research does not hold:**
  ★ **44% of drivers do not obey class behavior**
  ★ **15% drivers perform significant processing**
  ★ **28% drivers support multiple chipsets**

**Driver interactions**

★ **USB bus offers efficient access (as compared to PCI, Xen)**
  ★ **Supports high # devices/driver (standardized code)**
  ★ **Coarse-grained access**

**Driver similarity**

★ **400, 000 lines of code similar to code elsewhere and ripe for improvement via:**
  ★ **Procedural abstractions**
  ★ **Better multiple chipset support**
  ★ **Table driver programming**

★ **More results in "Understanding Modern Device Drivers"** **ASPLOS 2012**

# Outline

Tolerate device failures

Understand drivers and potential opportunities

Transactional approach for cheap recovery

**Checkpoint/restore**
**FGFT**
**Future work and conclude**

# Limitations of restart/replay recovery

★ **Device save/restore limited to restart/replay**

- ★ **Slow: Device initialization is complex (multiple seconds)**
- ★ **Incomplete: Unique device semantics not captured**
- ★ **Hard: Need to be written for every class of drivers**
- ★ **Large changes: Introduces new, large kernel subsystem**

Driver-Kernel Interface

Shadow Driver

Taps

Device Driver

Device

# Limitations of restart/replay recovery

★ **Device save/restore limited to restart/replay**

- ★ **Slow: Device initialization is complex (multiple seconds)**
- ★ **Incomplete: Unique device semantics not captured**
- ★ **Hard: Need to be written for every class of drivers**
- ★ **Large changes: Introduces new, large kernel subsystem**

Driver-Kernel Interface

Shadow Driver

Taps

Device Driver

Device

**Checkpoint/restore of device and driver state removes the need to reboot device and replay state**

# Checkpointing drivers is hard

★ **Easy to capture memory state**



network driver ↔ network card

# Checkpointing drivers is hard

★ **Easy to capture memory state**



checkpoint

**network driver**

**network card**

# Checkpointing drivers is hard

★ **Easy to capture memory state**



★ **Device state is not captured**

   ★ **Device configuration space**

# Checkpointing drivers is hard

★ **Easy to capture memory state**



★ **Device state is not captured**
  - ★ **Device configuration space**
  - ★ **Internal device registers and counters**

# Checkpointing drivers is hard

★**Easy to capture memory state**



checkpoint

network
driver

network
card

★ **Device state is not captured**

★ **Device configuration space**

★ **Internal device registers and counters**

★ **Memory buffer addresses used for DMA**

# Checkpointing drivers is hard

★**Easy to capture memory state**



★ **Device state is not captured**

  ★ **Device configuration space**

  ★ **Internal device registers and counters**

  ★ **Memory buffer addresses used for DMA**

★ **Unique for every device**

# Checkpointing drivers is hard

★ **Easy to capture memory state**

**checkpoint**

**Intuition: Operating systems already capture device state during power management**

**card**

★ **Device state is not captured**

  ★ **Device configuration space**

  ★ **Internal device registers and counters**

  ★ **Memory buffer addresses used for DMA**

★ **Unique for every device**

# Intuition with power management



★ **Refactor power management code for device checkpoints**

  ★ Correct: Developer captures unique device semantics

  ★ Fast: Avoids probe and latency critical for applications

★ **Ask developers to export checkpoint/restore in their drivers**

# Device checkpoint/restore from PM code

**Suspend**

| |
|---|
| Save config state |
| Save register state |
| Disable device |
| Save DMA state |
| Suspend device |

**Resume**

| |
|---|
| Restore config state |
| Restore register state |
| Restore or reset DMA state |
| Re-attach/Enable device |
| Device Ready |

# Device checkpoint/restore from PM code

**Suspend**

- Save config state
- Save register state

<br>

- Save DMA state
- Suspend device

**Resume**

- Restore config state
- Restore register state
- Restore or reset DMA state
- Re-attach/Enable device
- Device Ready

# Device checkpoint/restore from PM code

**Suspend**

**Resume**

| Save config state |
| Save register state |

| Save DMA state |

| Restore config state |
| Restore register state |
| Restore or reset DMA state |
| Re-attach/Enable device |
| Device Ready |

# Device checkpoint/restore from PM code

**Suspend**

**Resume**

| Save config state |
| Save register state |
| Save DMA state |

| Restore config state |
| Restore register state |
| Restore or reset DMA state |
| Re-attach/Enable device |
| Device Ready |

# Device checkpoint/restore from PM code

**Checkpoint**

**Resume**

| Checkpoint |
|:---:|
| Save config state |
| Save register state |
| Save DMA state |

| Resume |
|:---:|
| Restore config state |
| Restore register state |
| Restore or reset DMA state |
| Re-attach/Enable device |
| Device Ready |

# Device checkpoint/restore from PM code

**Checkpoint**

**Resume**

| Save config state |
|:---:|
| Save register state |
| Save DMA state |

| Restore config state |
|:---:|
| Restore register state |
| Restore or reset DMA state |
| Re-attach/Enable device |

# Device checkpoint/restore from PM code

**Checkpoint**                                    **Resume**

| Save config state |
| :---: |
| Save register state |
| Save DMA state |

| Restore config state |
| :---: |
| Restore register state |
| Restore or reset DMA state |

# Device checkpoint/restore from PM code

**Checkpoint**

**Restore**

| Save config state |
| :---: |
| Save register state |
| Save DMA state |

| Restore config state |
| :---: |
| Restore register state |
| Restore or reset DMA state |

51

# Device checkpoint/restore from PM code

**Checkpoint**

**Restore**

| Save config state |
| :---: |
| Save register state |
| Save DMA state |

| Restore config state |
| :---: |
| Restore register state |
| Restore or reset DMA state |

**Suspend/resume code provides device checkpoint functionality**

# Fine-Grained Fault Tolerance[ASPLOS 2013]

★ **Goal: Improve driver recovery with minor changes to drivers**

★ **Solution: Run drivers as <span style="color:orange">transactions</span> using device checkpoints**

# Fine-Grained Fault Tolerance [ASPLOS 2013]

★ **Goal: Improve driver recovery with minor changes to drivers**

★ **Solution: Run drivers as transactions using device checkpoints**

**Device state**

★ **Developers export checkpoint/restore in drivers**

C    R

# Fine-Grained Fault Tolerance [ASPLOS 2013]

★ **Goal: Improve driver recovery with minor changes to drivers**

★ **Solution: Run drivers as transactions using device checkpoints**

**Device state**

**Driver state**

★ Developers export checkpoint/restore in drivers

★ **Run drivers invocations as memory transactions**

★ **Use source transformation to copy parameters and run on separate stack**

**C** **R**

network driver ⟷ SFI network driver

# Fine-Grained Fault Tolerance[ASPLOS 2013]

★ **Goal: Improve driver recovery with minor changes to drivers**

★ **Solution: Run drivers as transactions using device checkpoints**

## Device state

★ **Developers export checkpoint/restore in drivers**

**C**    **R**

## Driver state

★ **Run drivers invocations as memory transactions**

★ **Use source transformation to copy parameters and run on separate stack**

**network driver** ⟷ **SFI network driver**

## Execution model

★ **Checkpoint device**

★ **Execute driver code as memory transactions**

★ **On failure, rollback and restore device**

★ **Re-use existing device locks in the driver**

# Adding transactional support to drivers

```
If (c==0) {
.
print ("Driver
init");
}
.
.
```

Driver with
checkpoint support

**Static modifications**

53

# Adding transactional support to drivers



```
If (c==0) {
.
.
print ("Driver
init");
}
.
.
```

Source transformation
(adds driver transactions)

Driver with
checkpoint support

User supplied
annotations

**Static modifications**

# Adding transactional support to drivers

```
If (c==0) {
.
print ("Driver
init");
}
.
.
```

**Main driver module**

```
If (c==0) {
.
print ("Driver
init");
}
.
.
```

**Source transformation (adds driver transactions)**

**Driver with checkpoint support**

**User supplied annotations**

**SFI driver module**

```
If (c==0) {
.
print ("Driver
init");
}
.
.
```

**SFI = software fault isolated**

**Static modifications**

53

# Adding transactional support to drivers



If (c==0) {
.
print ("Driver init");
}
.
.

**Main driver module**

**Source transformation (adds driver transactions)**

If (c==0) {
.
print ("Driver init");
}
.
.

**Driver with checkpoint support**

**User supplied annotations**

**SFI driver module**

If (c==0) {
.
print ("Driver init");
}
.
.

**SFI = software fault isolated**

**Static modifications**

**Run-time support**

# Adding transactional support to drivers



```
If (c==0) {
.
print ("Driver
init");
}
.
.
```

1200 LOC

Main driver module

Object tracking

Marshaling/ Demarshaling

Source transformation (adds driver transactions)

```
If (c==0) {
.
print ("Driver
init");
}
.
.
```

Driver with checkpoint support

SFI driver module

Kernel undo log

User supplied annotations

```
If (c==0) {
.
print ("Driver
init");
}
.
.
```

**Communication and recovery support**

**SFI = software fault isolated**

**Static modifications**

**Run-time support**

53

# Transactional execution of drivers

# Transactional execution of drivers

# Transactional execution of drivers

# Transactional execution of drivers

**C**

**netdev->priv->rx_ring**

**netdev->priv->tx_ring**

**get ringparam**

**probe**

**xmit**

**get config**

**netdev**

**network driver**

**SFI network driver**

# Transactional execution of drivers



netdev->priv->rx_ring
netdev->priv->tx_ring

**Range Table**

| Address | Access rights |
|---------|---------------|
| 0xffffa000 | Read |
| 0xffffa008 | Write |
| 0xffffa00a | Read |

**get ringparam**
probe
xmit
get config

**netdev**

**network driver**

**SFI network driver**

# Transactional execution of drivers



**C**

**netdev->priv->rx_ring**

**netdev->priv->tx_ring**

**Range Table**

| Address | Access rights |
|---------|---------------|
| oxffffaooo | Read |
| oxffffaoo8 | Write |
| oxffffaooa | Read |

**get ringparam**

**netdev**

**network driver**

probe

xmit

get config

**SFI network driver**

# Transactional execution of drivers



C

netdev->priv->rx_ring

netdev->priv->tx_ring

**Range Table**

| Address | Access rights |
|---------|---------------|
| 0xffffa000 | Read |
| 0xffffa008 | Write |
| 0xffffa00a | Read |

get ringparam

probe

xmit

get config

netdev

**network driver**

**SFI network driver**

# Transactional execution of drivers



netdev->priv->rx_ring
netdev->priv->tx_ring

**C**

**Range Table**

| Address | Access rights |
|---------|---------------|
| oxffffaoo0 | Read |
| oxffffaoo8 | Write |
| oxffffaooa | Read |

**get ringparam**

**probe**

**xmit**

**get config**

**netdev**

**network driver**

**SFI network driver**

**Kernel Log**

**alloc**

**result**
netdev->priv->rx_ring
netdev->priv->tx_ring

# Transactional execution of drivers

**C**

**netdev->priv->rx_ring**
**netdev->priv->tx_ring**

**Range Table**

| Address | Access rights |
|---------|---------------|
| oxffffaooo | Read |
| oxffffaoo8 | Write |
| oxffffaooa | Read |

**get ringparam**

**netdev**

**network driver**

probe

xmit

get config

**SFI network driver**

**Kernel Log**

**alloc**

**result**

**netdev->priv->rx_ring**
**netdev->priv->tx_ring**

★ **Detects and recovers from:**

    ★ **Memory errors like invalid pointer accesses**

    ★ **Structural errors like malformed structures**

    ★ **Processor exceptions like divide by zero, stack corruption**

# FGFT: Failed transactions

**get ringparam**

probe

xmit

get config

**network driver**

**SFI network driver**

55

# FGFT: Failed transactions

# FGFT: Failed transactions



get ringparam

probe

xmit

get config

netdev

network driver

SFI network driver

# FGFT: Failed transactions

# FGFT: Failed transactions

**C**

netdev->priv->rx_ring

netdev->priv->tx_ring

**get ringparam**

probe

xmit

get config

**netdev**

**network driver**

**SFI network driver**

**Range Table**

| Address | Access rights |
|---------|---------------|
| 0xffffa000 | Read |
| 0xffffa008 | Write |
| 0xffffa00a | Read |

# FGFT: Failed transactions

C

netdev->priv->rx_ring

netdev->priv->tx_ring

## Range Table

| Address | Access rights |
|---------|---------------|
| 0xffffa000 | Read |
| 0xffffa008 | Write |
| 0xffffa00a | Read |

get ringparam

netdev

probe

xmit

get config

**network driver**

**SFI network driver**

# FGFT: Failed transactions

**C**

netdev->priv->rx_ring

netdev->priv->tx_ring

**get ringparam**

**netdev**

**network driver**

**probe**

**xmit**

**get config**

**SFI network driver**

### Range Table

| Address | Access rights |
|---------|---------------|
| 0xffffa000 | Read |
| 0xffffa008 | Write |
| 0xffffa00a | Read |

# FGFT: Failed transactions



**netdev->priv->rx_ring**

**netdev->priv->tx_ring**

**C**

**get ringparam**

**probe**

**xmit**

**get config**

**netdev**

**network driver**

**SFI network driver**

**Range Table**

| Address | Access rights |
|---------|---------------|
| oxffffaooo | Read |
| oxffffaoo8 | Write |
| oxffffaooa | Read |

**Kernel Log**

**alloc**

# FGFT: Failed transactions

**C**

netdev->priv->rx_ring

netdev->priv->tx_ring

**Range Table**

| Address | Access rights |
|---------|---------------|
| oxffffaooo | Read |
| oxffffaoo8 | Write |
| oxffffaooa | Read |

get ringparam

**netdev**

probe

xmit

get config

**network driver**

**SFI network driver**

**Kernel Log**

**alloc**

# FGFT: Failed transactions

**C**

**netdev->priv->rx_ring**
**netdev->priv->tx_ring**

**Range Table**

| Address | Access rights |
|---------|---------------|
| 0xffffa000 | Read |
| 0xffffa008 | Write |
| 0xffffa00a | Read |

**get ringparam**

**netdev**

probe

xmit

get config

**network driver**

**SFI network driver**

**Kernel Log**

**alloc**

# FGFT: Failed transactions



C

netdev->priv->rx_ring

netdev->priv->tx_ring

get ringparam

netdev

probe

xmit

get config

network driver

SFI network driver

## Range Table

| Address | Access rights |
|---------|---------------|
| 0xffffa000 | Read |
| 0xffffa008 | Write |
| 0xffffa00a | Read |

# FGFT: Failed transactions



netdev->priv->rx_ring
netdev->priv->tx_ring

**Range Table**

| Address | Access rights |
|---------|---------------|
| 0xffffa000 | Read |
| 0xffffa008 | Write |
| 0xffffa00a | Read |

**C**

get ringparam

**netdev**

probe

xmit

get config

**network driver**

**R**

err

**SFI network driver**

# How does this give us transactional execution?
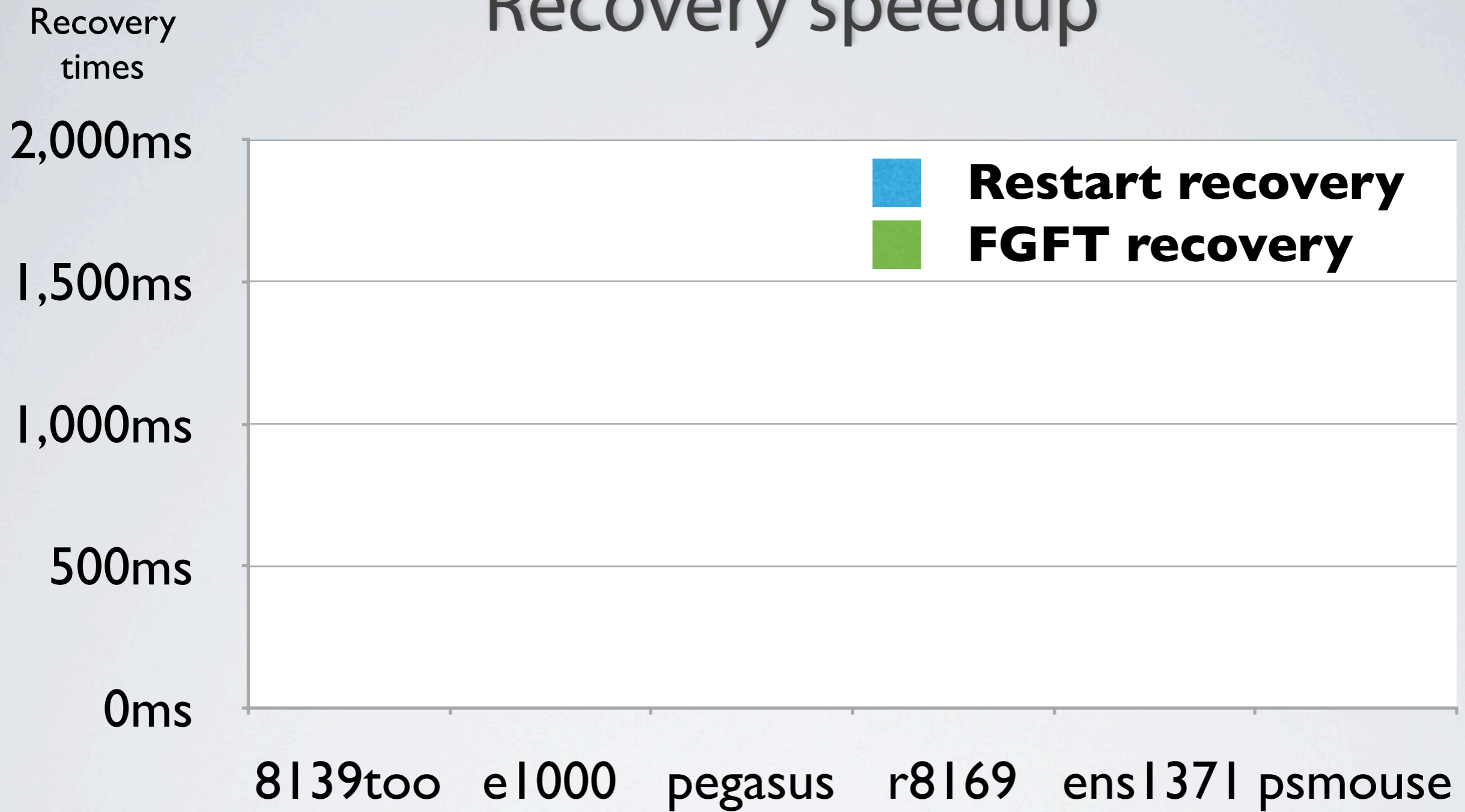
# How does this give us transactional execution?

★ **Atomicity: All or nothing execution**

    ★ **Driver state: Run code in SFI module**

    ★ **Device state: Explicitly checkpoint/restore state**

# How does this give us transactional execution?

★ **Atomicity: All or nothing execution**

   ★ **Driver state: Run code in SFI module**

   ★ **Device state: Explicitly checkpoint/restore state**

★ **Isolation: Serialization to hide incomplete transactions**

   ★ **Re-use existing device locks to lock driver**

   ★ **Two phase locking**

# How does this give us transactional execution?

★ **Atomicity: All or nothing execution**

    ★ **Driver state: Run code in SFI module**

    ★ **Device state: Explicitly checkpoint/restore state**

★ **Isolation: Serialization to hide incomplete transactions**

    ★ **Re-use existing device locks to lock driver**

    ★ **Two phase locking**

★ **Consistency: Only valid (kernel, driver and device) states**

    ★ **Higher level mechanisms to rollback external actions**

    ★ **At most once device action guarantee to applications**

# Recovery speedup

Recovery times

2,000ms

1,500ms
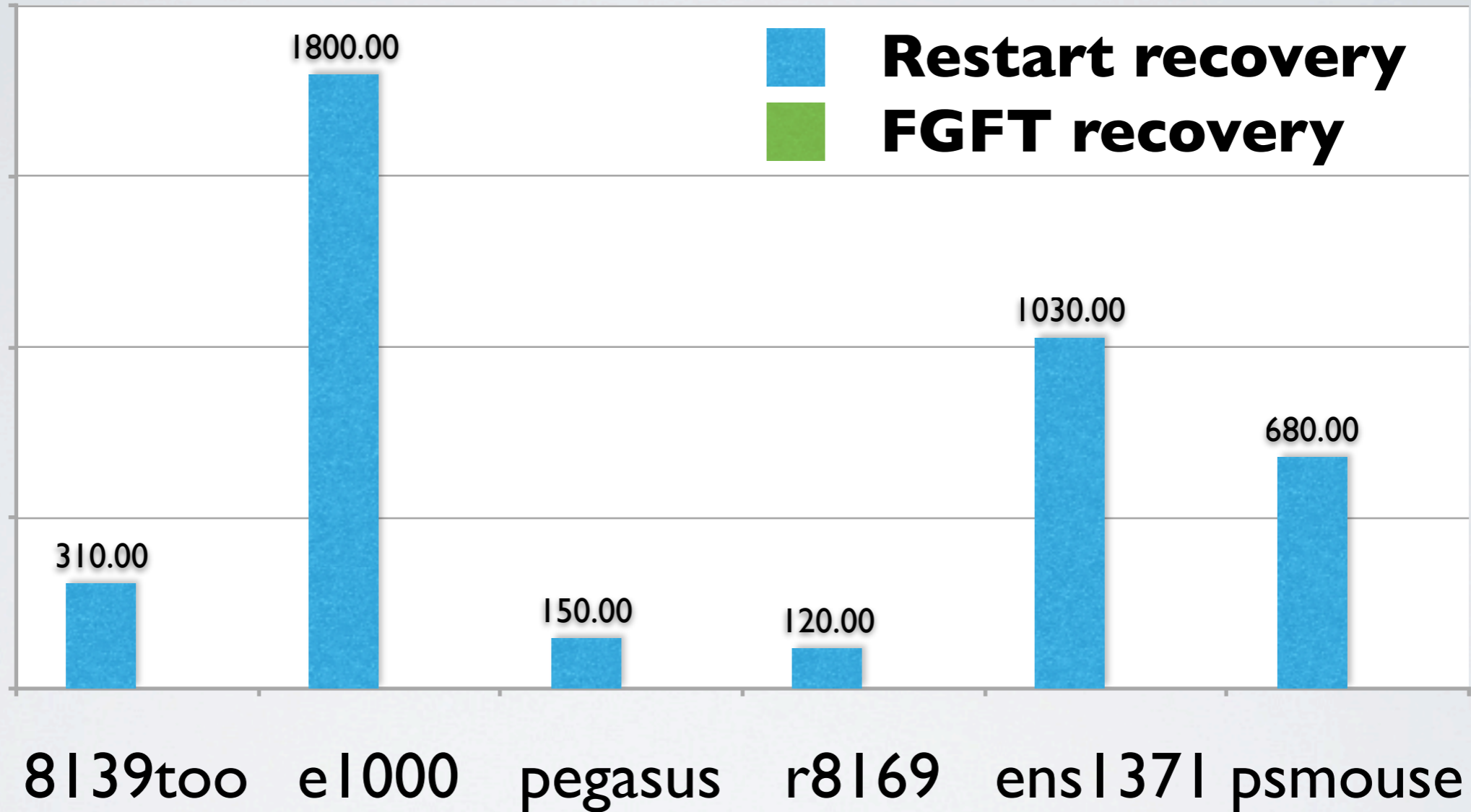
1,000ms

500ms

0ms

- Restart recovery
- FGFT recovery

8139too  e1000  pegasus  r8169  ens1371  psmouse

# Recovery speedup

**Recovery times**



Legend:
- Restart recovery (blue)
- FGFT recovery (green)

Bar chart with y-axis from 0ms to 2,000ms (gridlines at 500ms, 1,000ms, 1,500ms, 2,000ms):

- 8139too: 310.00
- e1000: 1800.00
- pegasus: 150.00
- r8169: 120.00
- ens1371: 1030.00
- psmouse: 680.00

# Programming effort

| Driver | LOC | Checkpoint/restore effort | |
| --- | --- | --- | --- |
| | | LOC Moved | LOC Added |
| 8139too | 1, 904 | 26 | 4 |
| e1000 | 13, 973 | 32 | 10 |
| r8169 | 2, 993 | 17 | 5 |
| pegasus | 1, 541 | 22 | 5 |
| ens1371 | 2, 110 | 16 | 6 |
| psmouse | 2, 448 | 19 | 6 |

**FGFT requires limited programmer effort and needs only 38 lines of new kernel code**
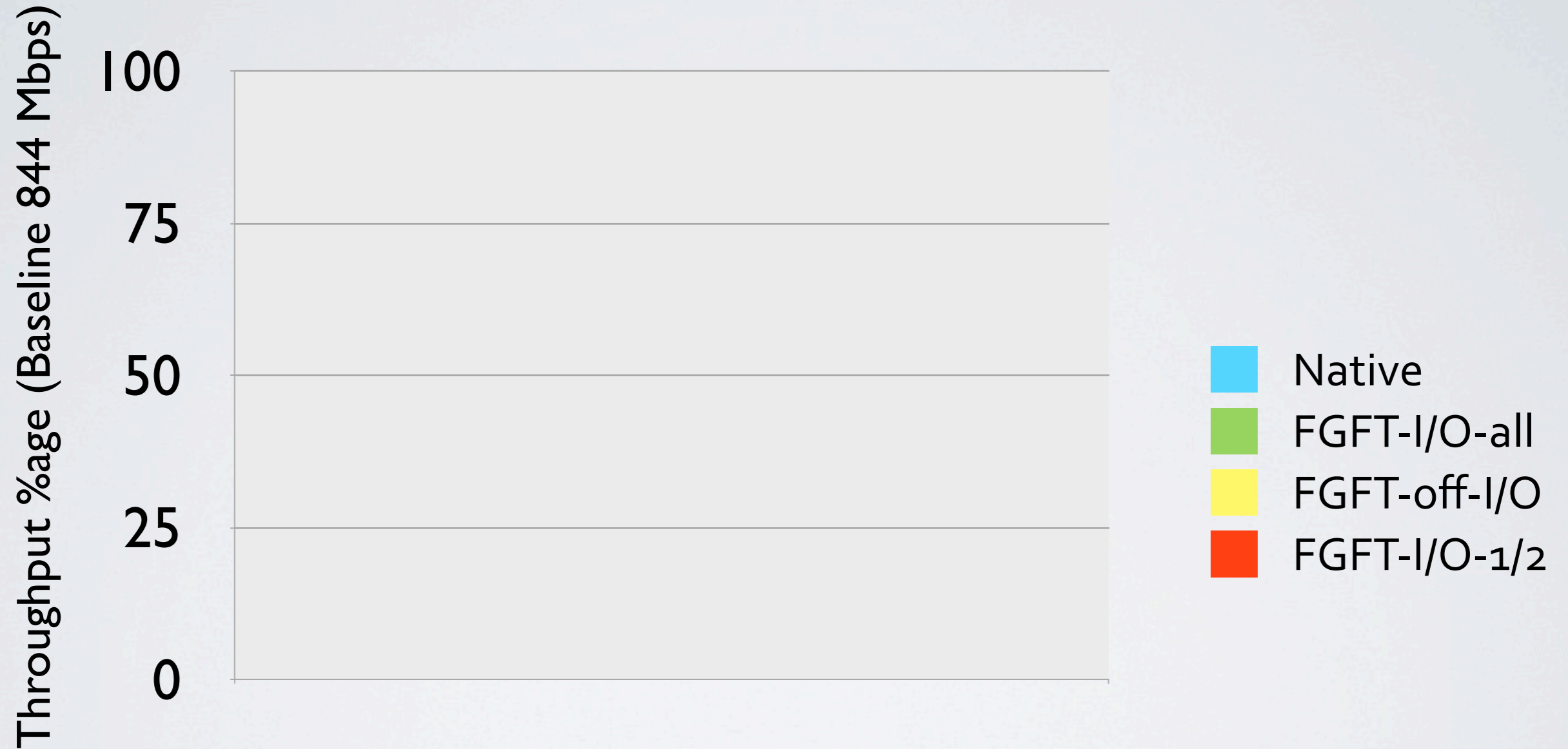
# Throughput with isolation and recovery

Native

FGFT-I/O-all

FGFT-off-I/O

FGFT-I/O-1/2

**netperf on Intel quad-core machines**

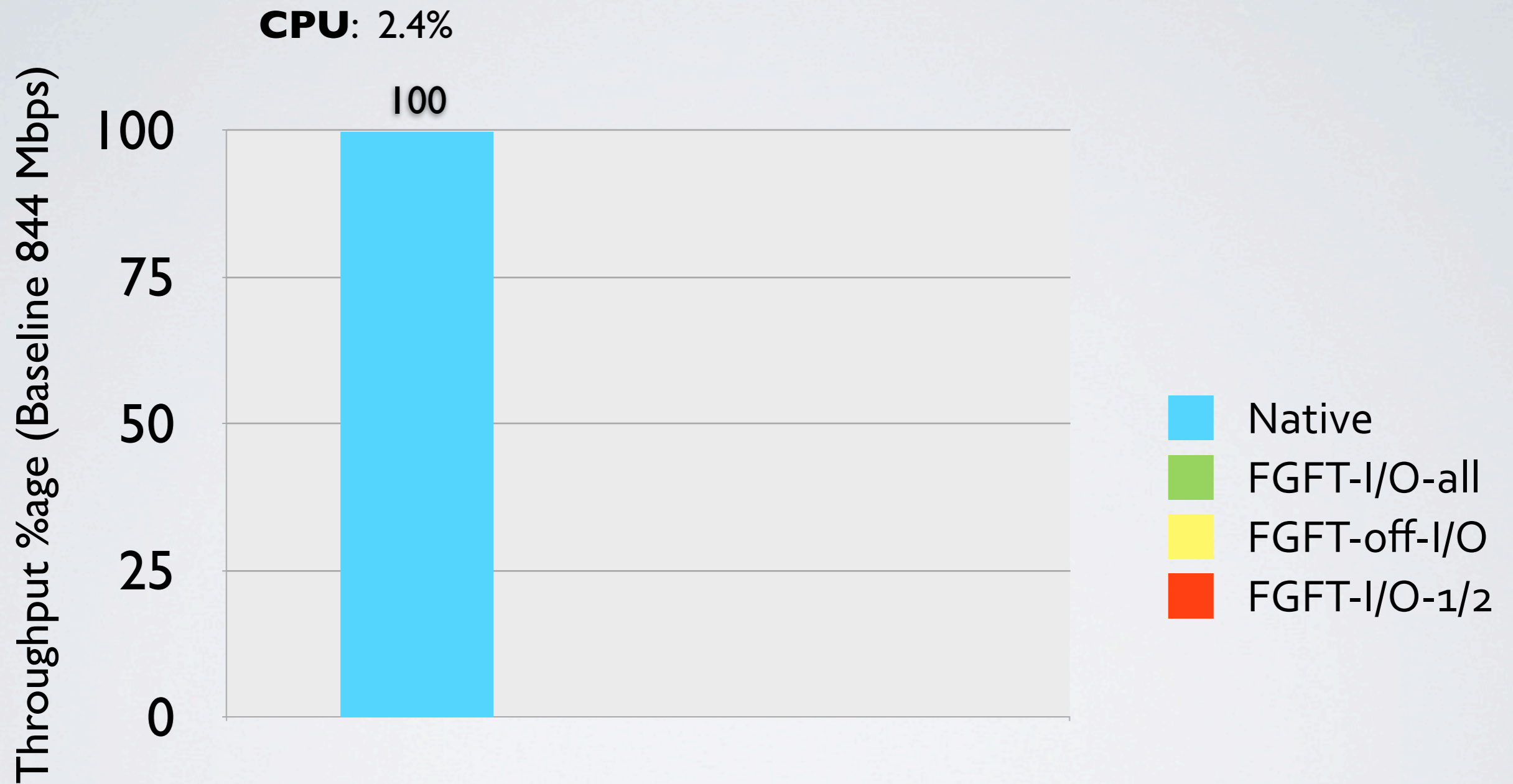# Throughput with isolation and recovery

Throughput %age (Baseline 844 Mbps)

100

75

50

25

0

Native

FGFT-I/O-all

FGFT-off-I/O

FGFT-I/O-1/2

e1000 Network Card

**netperf on Intel quad-core machines**

# Throughput with isolation and recovery

**CPU**: 2.4%

Throughput %age (Baseline 844 Mbps)

100 — 100

75

50

25

0

Legend:
- Native (blue)
- FGFT-I/O-all (green)
- FGFT-off-I/O (yellow)
- FGFT-I/O-1/2 (red)

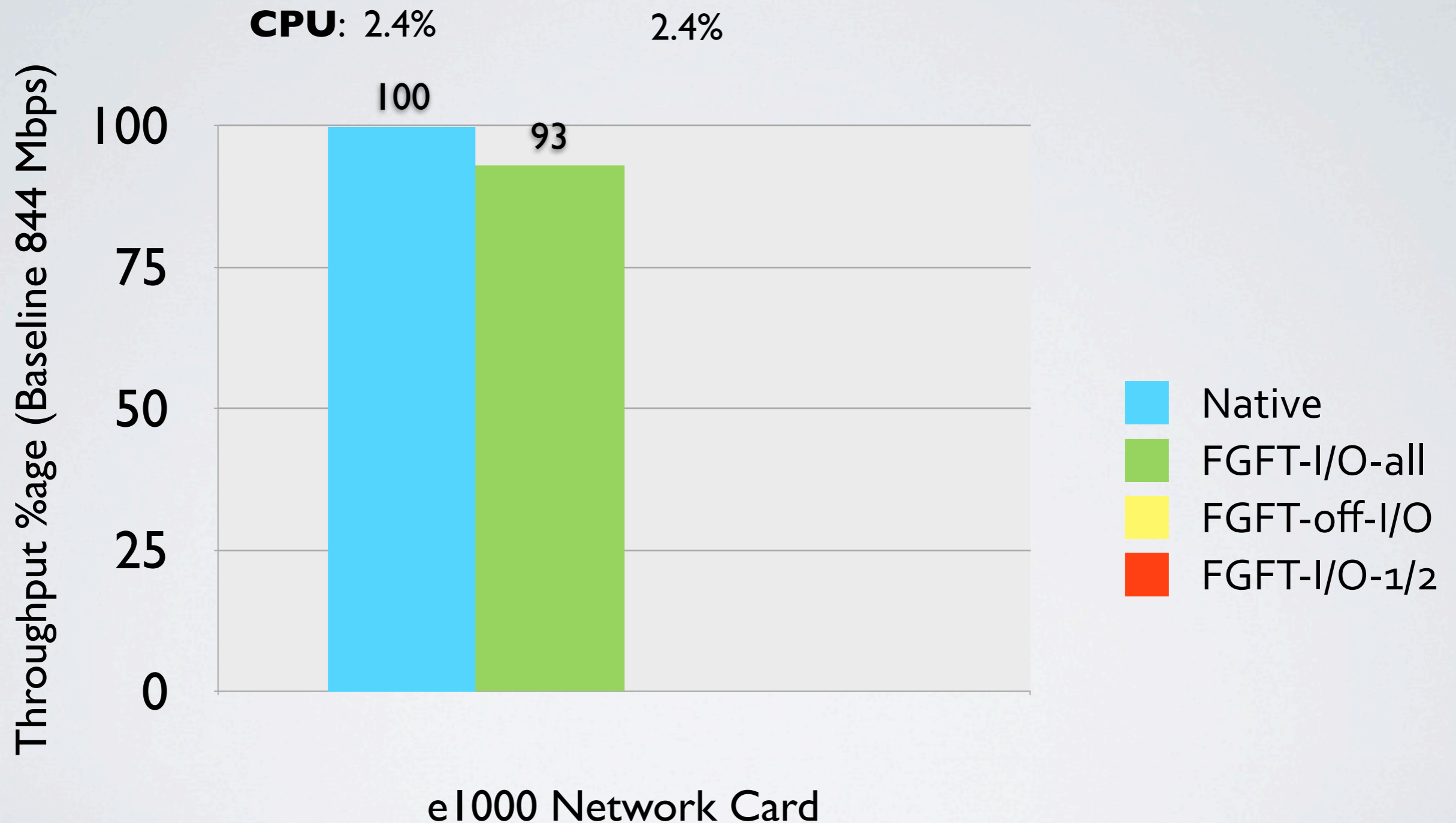e1000 Network Card

**netperf on Intel quad-core machines**

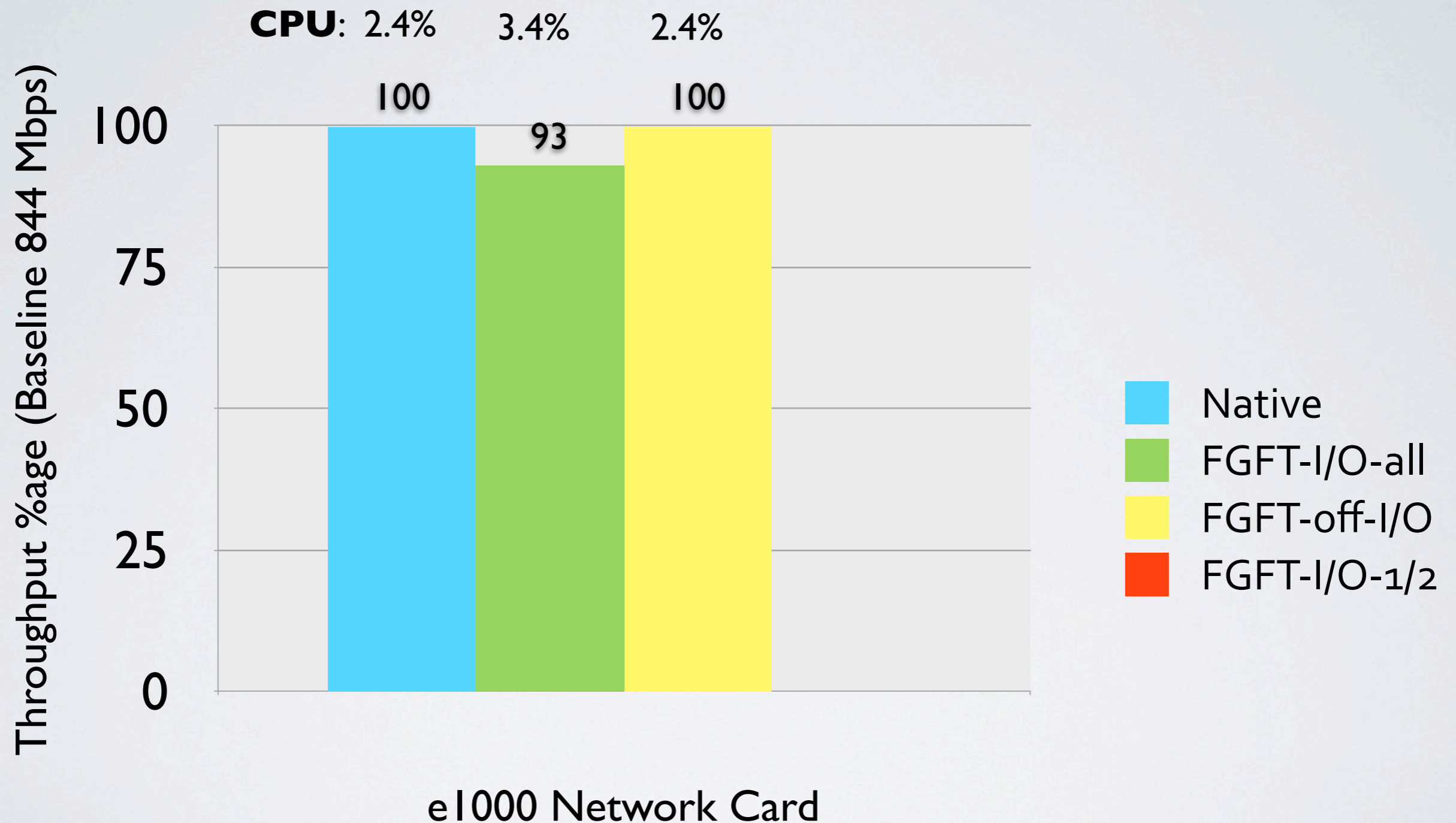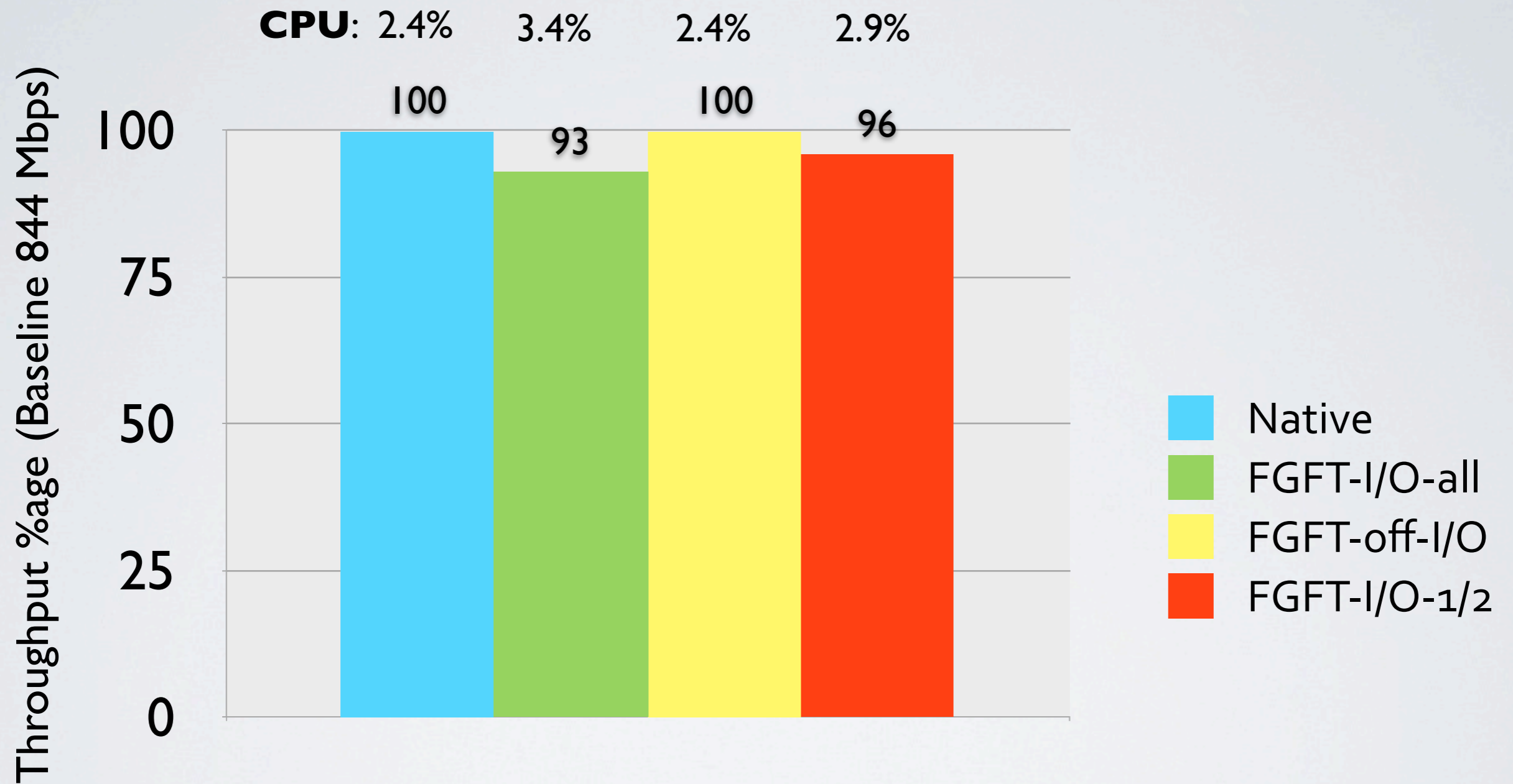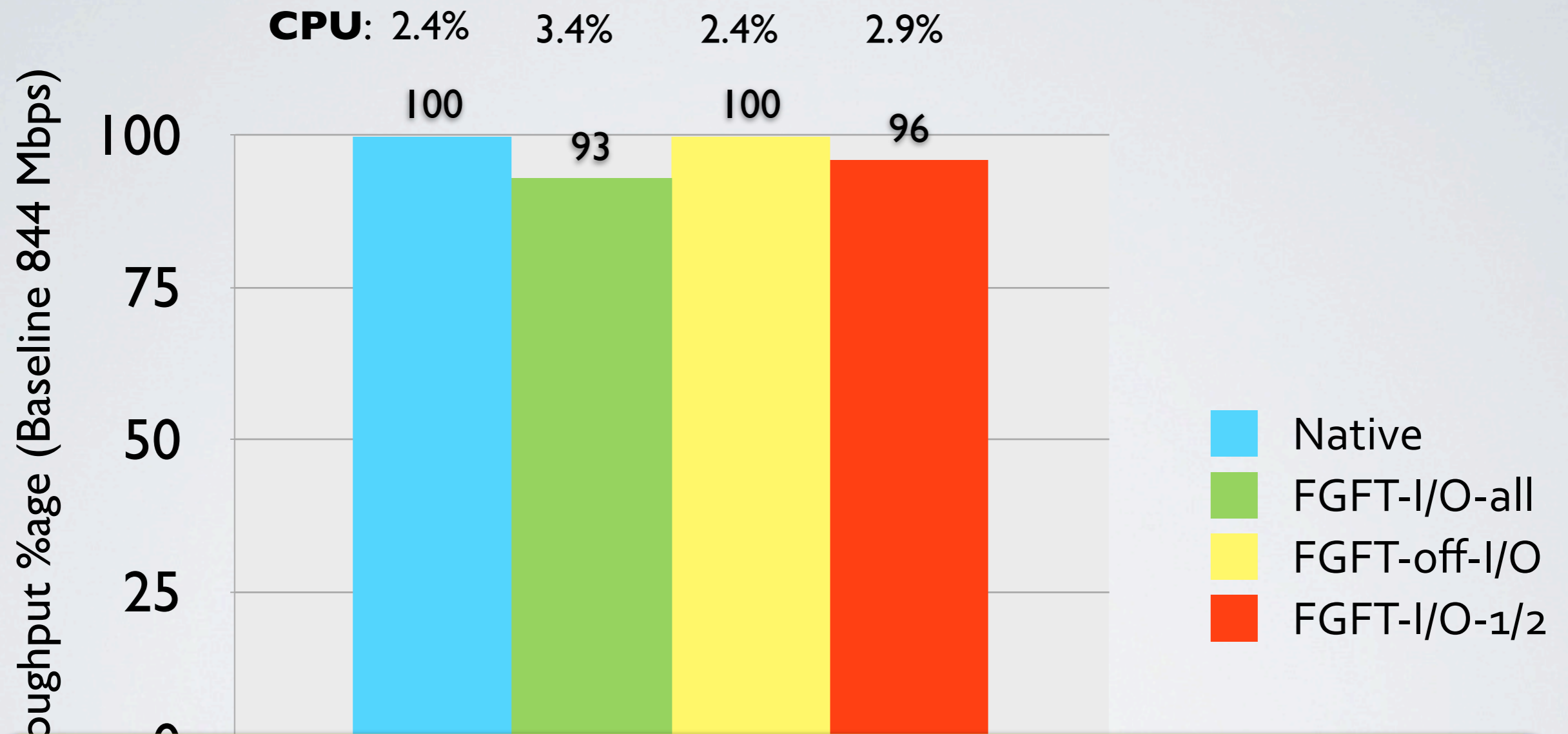# Throughput with isolation and recovery



CPU: 2.4%   3.4%   2.4%

netperf on Intel quad-core machines

Throughput with isolation and recovery

netperf on Intel quad-core machines

# Throughput with isolation and recovery

**CPU**: 2.4%  3.4%  2.4%  2.9%

Throughput %age (Baseline 844 Mbps)

- 100 — Native
- 93 — FGFT-I/O-all
- 100 — FGFT-off-I/O
- 96 — FGFT-I/O-1/2

y-axis: 100, 75, 50, 25

**FGFT can isolate and recover high bandwidth devices at low overhead without adding kernel subsystems**

**netperf on Intel quad-core machines**

# Talk summary

First research consideration of hardware failures in drivers

**SOSP '09**

Released tool, patches & informed developers

Largest study of drivers to understand their behavior and verify research assumptions

**ASPLOS '12**

Measured driver behavior & identified new directions

Introduced checkpoint/restore in drivers for low latency fault tolerance

**ASPLOS '13**

Fast & correct recovery with incremental changes to drivers

# Questions

**Asim Kadav**
★ **http://cs.wisc.edu/~kadav**
★ **kadav@cs.wisc.edu**