# IO bound or CPU bound?

2014.10.31.

Dongjun Shin

Samsung Electronics

# Contents

- **Background**
- **Understanding CPU behavior**
- **Experiments**
- **Improvement idea**
- **Revisiting Linux I/O stack**
- **Conclusion**

## CPU bound

- A computer is CPU-bound (or compute-bound) when the time for it to complete a task is determined principally <u>by the speed of the central processor</u>: processor utilization is high, perhaps at 100% usage for many seconds or minutes (wikipedia)

## I/O bound

- I/O bound refers to a condition in which the time it takes to complete a computation is determined principally <u>by the period spent waiting for input/output operations</u> to be completed (wikipedia)
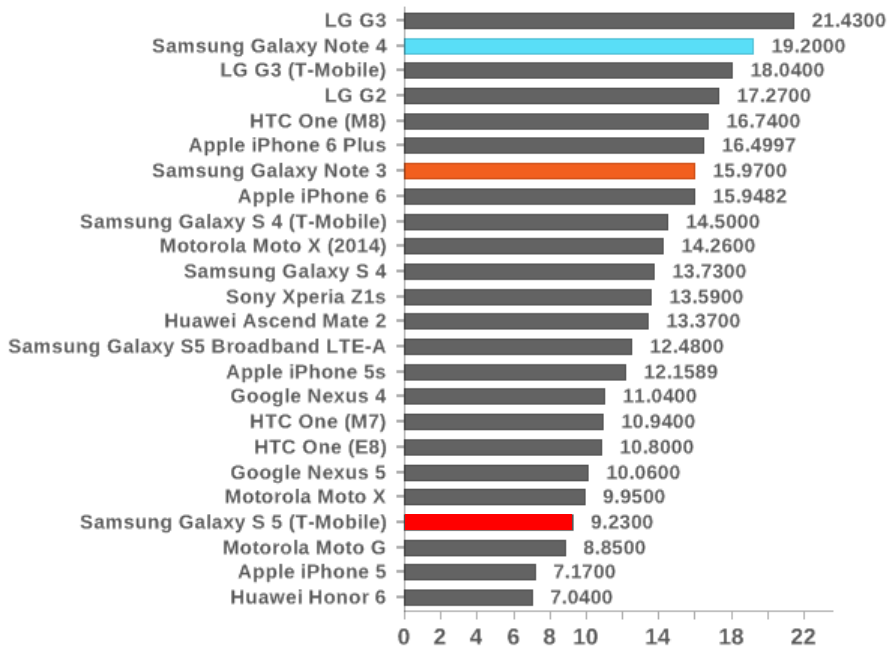
## In reality

- I have a lagging application. Who is to be blamed?
- I have a lagging application **and it seems that there are lots of I/O. It must be I/O bound.**
- I have a lagging application and it seems that there are lots of I/O. **I can't believe it because it's running on ultra fast SSD!**
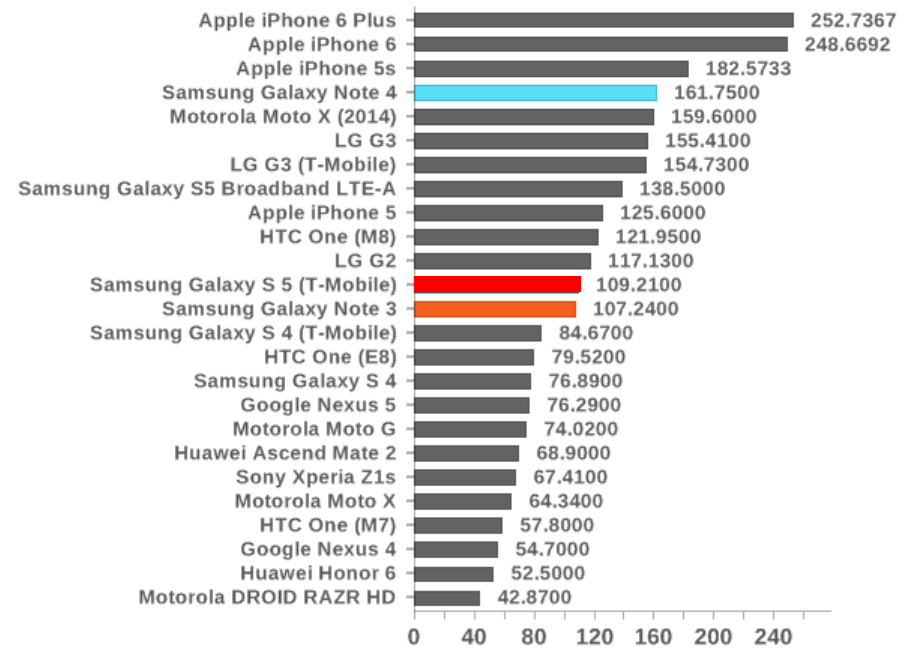
■ **AnandTech: Samsung Galaxy S5 vs. Galaxy Note 4**

– What's the cause of read I/O regression?

**Internal NAND - Random Read**
4KB Random Reads in MB/s - Higher is Better

| Device | MB/s |
|---|---|
| LG G3 | 21.4300 |
| Samsung Galaxy Note 4 | 19.2000 |
| LG G3 (T-Mobile) | 18.0400 |
| LG G2 | 17.2700 |
| HTC One (M8) | 16.7400 |
| Apple iPhone 6 Plus | 16.4997 |
| Samsung Galaxy Note 3 | 15.9700 |
| Apple iPhone 6 | 15.9482 |
| Samsung Galaxy S 4 (T-Mobile) | 14.5000 |
| Motorola Moto X (2014) | 14.2600 |
| Samsung Galaxy S 4 | 13.7300 |
| Sony Xperia Z1s | 13.5900 |
| Huawei Ascend Mate 2 | 13.3700 |
| Samsung Galaxy S5 Broadband LTE-A | 12.4800 |
| Apple iPhone 5s | 12.1589 |
| Google Nexus 4 | 11.0400 |
| HTC One (M7) | 10.9400 |
| HTC One (E8) | 10.8000 |
| Google Nexus 5 | 10.0600 |
| Motorola Moto X | 9.9500 |
| Samsung Galaxy S 5 (T-Mobile) | 9.2300 |
| Motorola Moto G | 8.8500 |
| Apple iPhone 5 | 7.1700 |
| Huawei Honor 6 | 7.0400 |

**Internal NAND - Sequential Read**
256KB Sequential Reads in MB/s - Higher is Better

| Device | MB/s |
|---|---|
| Apple iPhone 6 Plus | 252.7367 |
| Apple iPhone 6 | 248.6692 |
| Apple iPhone 5s | 182.5733 |
| Samsung Galaxy Note 4 | 161.7500 |
| Motorola Moto X (2014) | 159.6000 |
| LG G3 | 155.4100 |
| LG G3 (T-Mobile) | 154.7300 |
| Samsung Galaxy S5 Broadband LTE-A | 138.5000 |
| Apple iPhone 5 | 125.6000 |
| HTC One (M8) | 121.9500 |
| LG G2 | 117.1300 |
| Samsung Galaxy S 5 (T-Mobile) | 109.2100 |
| Samsung Galaxy Note 3 | 107.2400 |
| Samsung Galaxy S 4 (T-Mobile) | 84.6700 |
| HTC One (E8) | 79.5200 |
| Samsung Galaxy S 4 | 76.8900 |
| Google Nexus 5 | 76.2900 |
| Motorola Moto G | 74.0200 |
| Huawei Ascend Mate 2 | 68.9000 |
| Sony Xperia Z1s | 67.4100 |
| Motorola Moto X | 64.3400 |
| HTC One (M7) | 57.8000 |
| Google Nexus 4 | 54.7000 |
| Huawei Honor 6 | 52.5000 |
| Motorola DROID RAZR HD | 42.8700 |

| | AP | eMMC |
|---|---|---|
| Galaxy S5 | S801 2.5GHz x 4 | 5.0 |
| Galaxy Note4 | S805 2.7GHz x 4 | 5.0 |

http://www.anandtech.com/show/7903/samsung-galaxy-s-5-review
http://www.anandtech.com/show/8613/the-samsung-galaxy-note-4-review

# Understanding CPU Behavior (1/2)

■ **CPUFreq governor**

– Performance

▪ This locks the phone's CPU at maximum frequency

– Powersave

▪ This locks the CPU frequency at the lowest frequency

– Ondemand

▪ Boost clock speed to maximum on demand and step down if CPU load is low

– Interactive

▪ Similar to ondemand, but this governor dynamically scales CPU clock speed in response to workload

▪ Interactive is significantly more responsive than ondemand, because it's faster at scaling to maximum

■ **io_is_busy**

– Flag that determines if waiting for IO should increase CPU utilization in bump up CPU frequency (for ondemand and interactive)

– Tradeoff: performance vs. power

https://android.googlesource.com/kernel/common/+/android-3.4/Documentation/cpu-freq/governors.txt

## Characteristics of ARM big.LITTLE scheduling

- All interrupts are handled by CPU0
  - Load-balancing of interrupts across cores is not always the best solution*
- Designed for power efficiency
  - Only use big cores when it is necessary**

## What's the impact of this scheduling on I/O intensive app?

* Migrating software to multicore SMP systems (by Satyaki Mukherjee, ARM)
** Update on big.LITTLE scheduling experiments (by Morten Rasmussen, ARM)

Software R&D Center

# Experiments

- **Hardware: ODROID XU3**
  - Exynos5422 (4x A15 1.2-2GHz, 4x A7 1-1.5GHz)
    - Little(A7): CPU0-3, Big(A15): CPU4-7
  - 2GB LPDDR3 DRAM
  - eMMC 5.0 HS400 64GB
- **Software: Android 4.4.4**
  - Linux 3.10.9
- **Benchmark: fio**
  - Single thread: SW→RW→SR→RR (3 loops for each)
  - File size: 100MB (direct I/O), 1GB (buffered I/O)
  - I/O chunk: 256KB for sequential, 4KB for random
- **Parameters**
  - Governor: interactive (default), powersave (min), performance (max)
    - io_is_busy: toggle for interactive
  - Affinity: big vs. little

Software R&D Center

# Experimental Results – 100MB Direct I/O

- **I/O throughput scales with CPU clock**
  - Performance vs. powersave: **+30% for RR & RW**, +20% for SR, +15% for SW
  - Interactive & io_is_busy=0: almost same with powersave
- **Effects of big.LITTLE**
  - +15% for SW



(numbers in the box means CPU clock frequency in GHz)
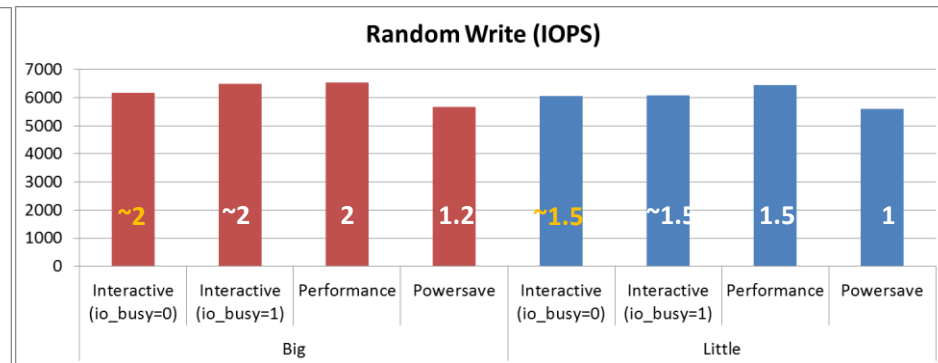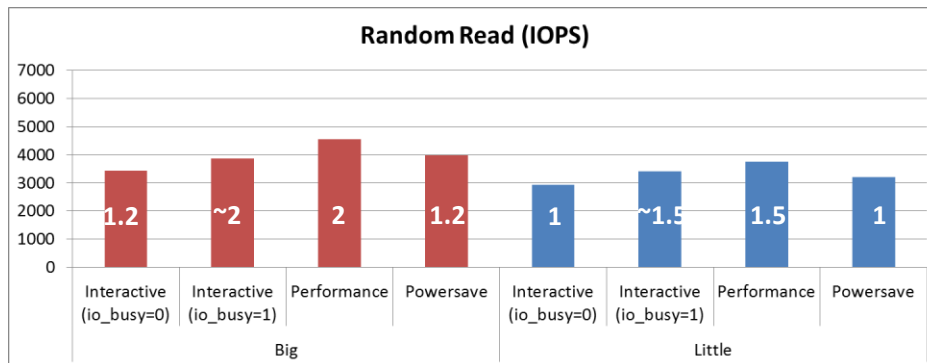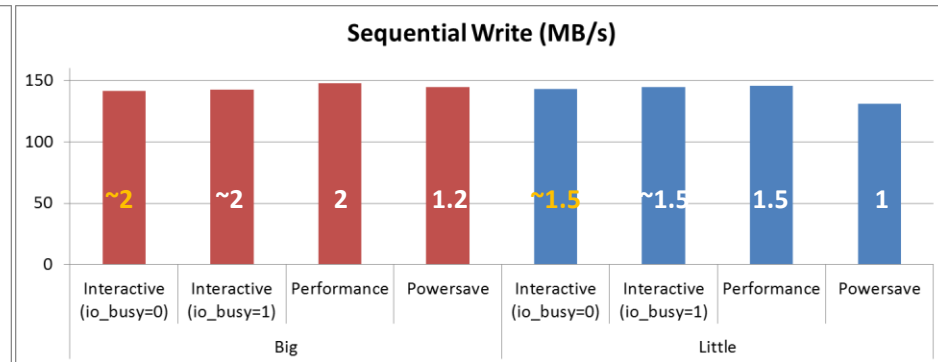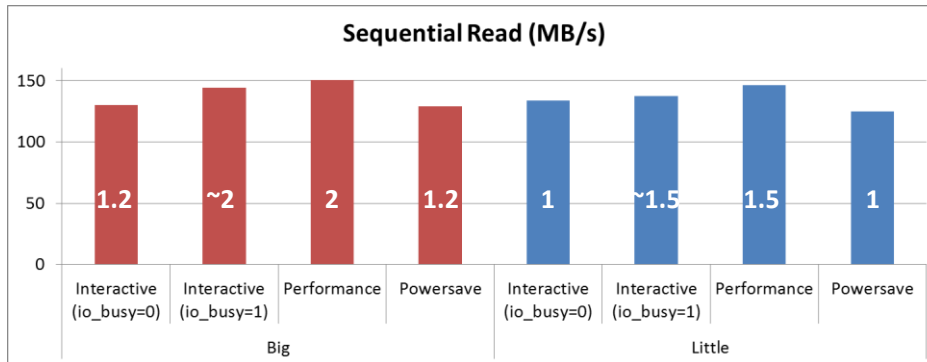
Software R&D Center

# Experimental Results – 1GB Buffered I/O

■ **Benchmark results are higher and less variable than direct I/O**

– Buffered vs. direct: **+100% for SW & RW**, +50% for SR

■ **RR is still CPU bound**

– Performance vs. powersave: +30%

– big vs. LITTLE: +20%


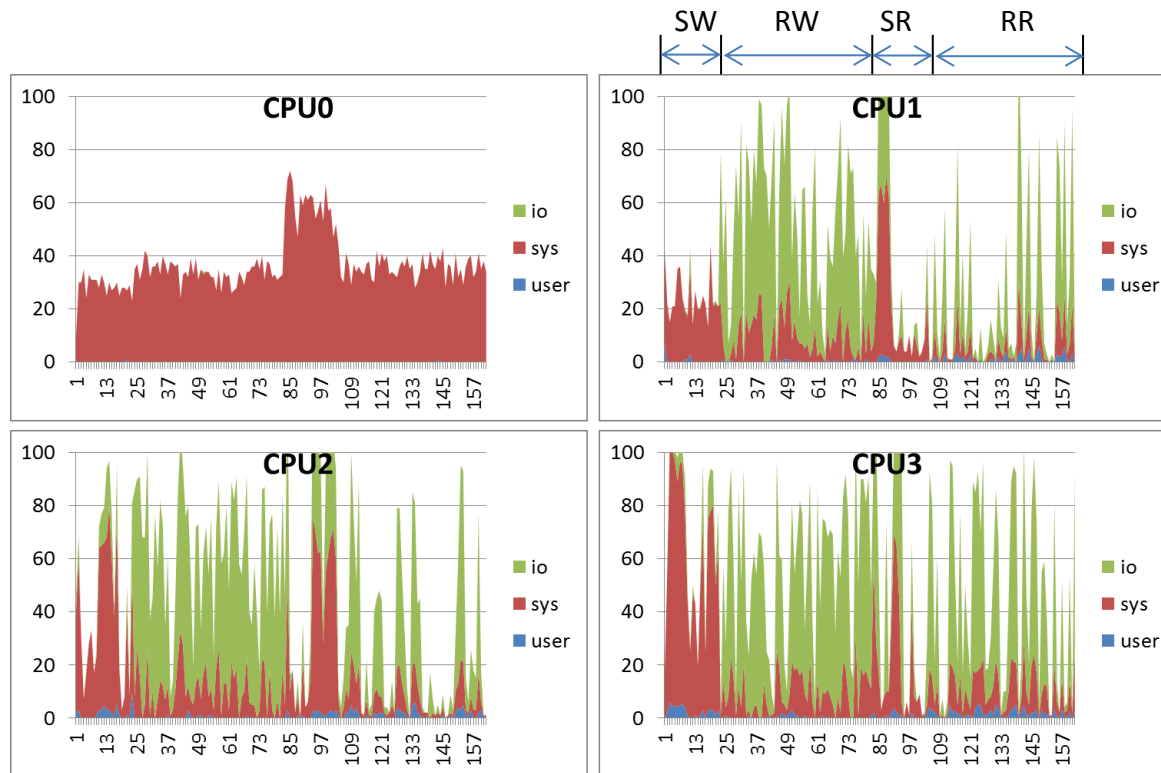
(numbers in the box means CPU clock frequency in GHz)

Software R&D Center

# CPU Load & Scheduling Analysis

■ **fio runs on 3x A7 only although all 8 cores are available**
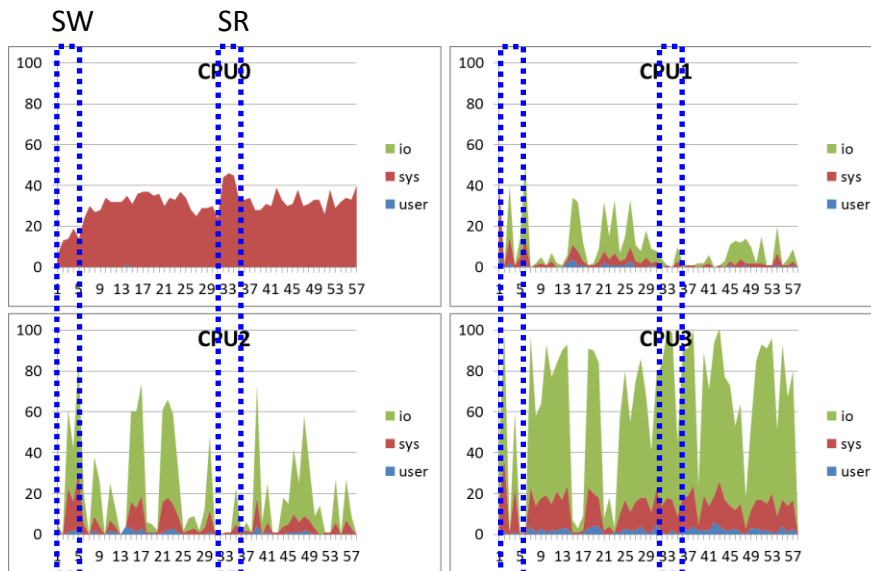
– fio process migrates among A7 cores

■ **Issues**

– CPU migration may be harmful for I/O intensive workload (D-cache efficiency)
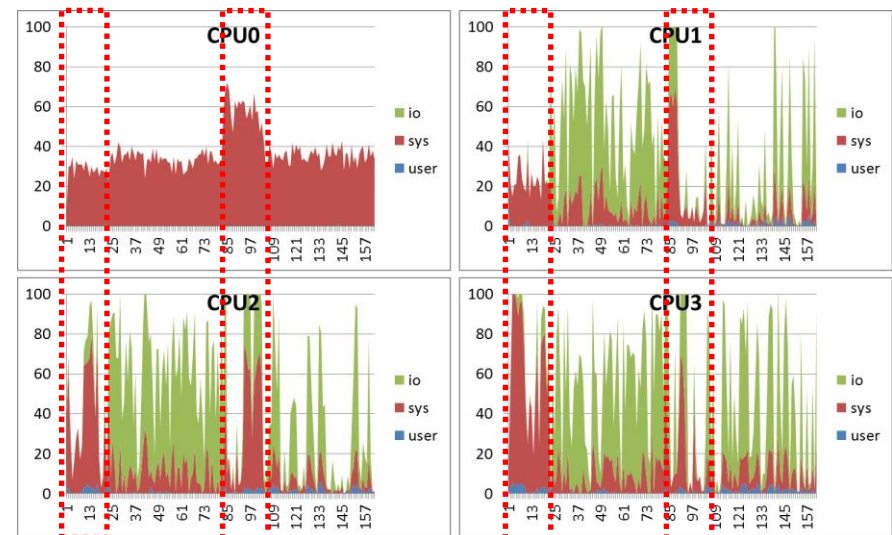
– A15 is faster at I/O handling



<Buffered I/O, interactive governor>

- **Overall CPU utilization of direct I/O is lower by imbalanced %sys vs. %io**
  - Balanced means "well-pipelined"
- **Buffered sequential I/O is much faster when %sys is higher**
  - End-to-end pipeline: readahead, delayed write
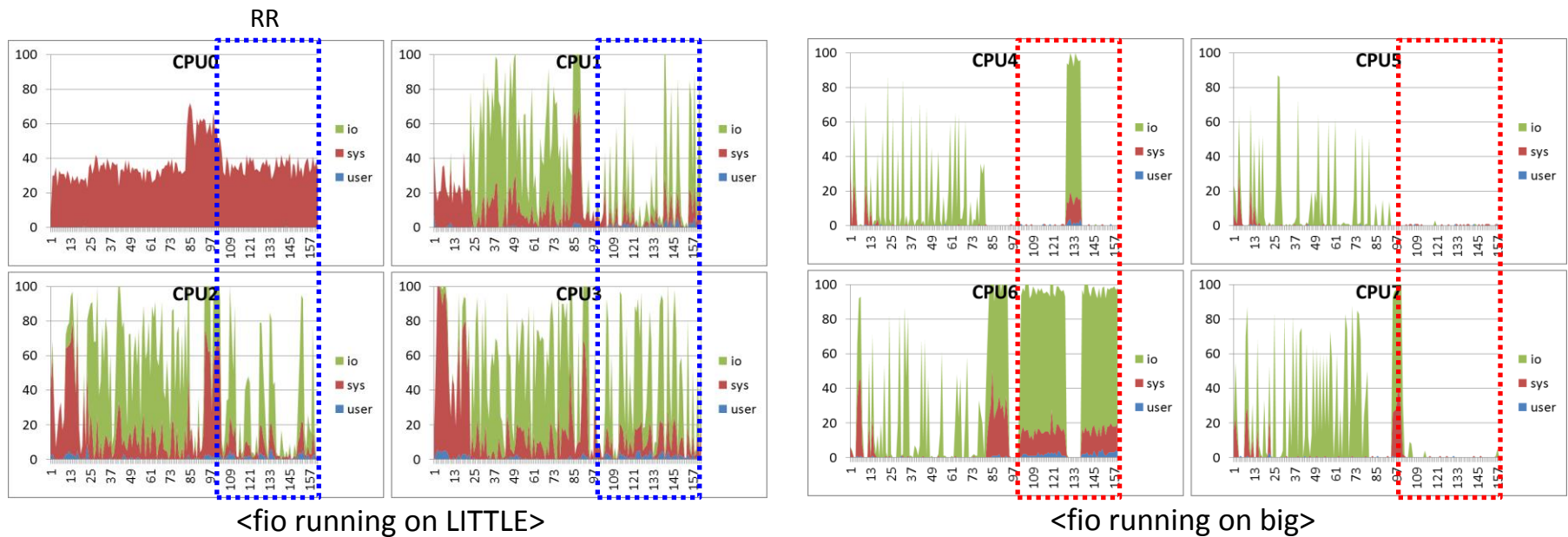- **Buffered RW is faster mainly due to eMMC cache (not CPU dependent)**



&lt;fio in direct I/O&gt;

&lt;fio in buffered I/O&gt;

# Little vs. Big

- **Buffered I/O performance is almost the same except RR**
  - CPU load is different: big has higher %io
  - **Big has potential room for improvement** if %io is balanced with %sys (more pipeline)
- **RR throughput has some relationship with CPU migration policy**
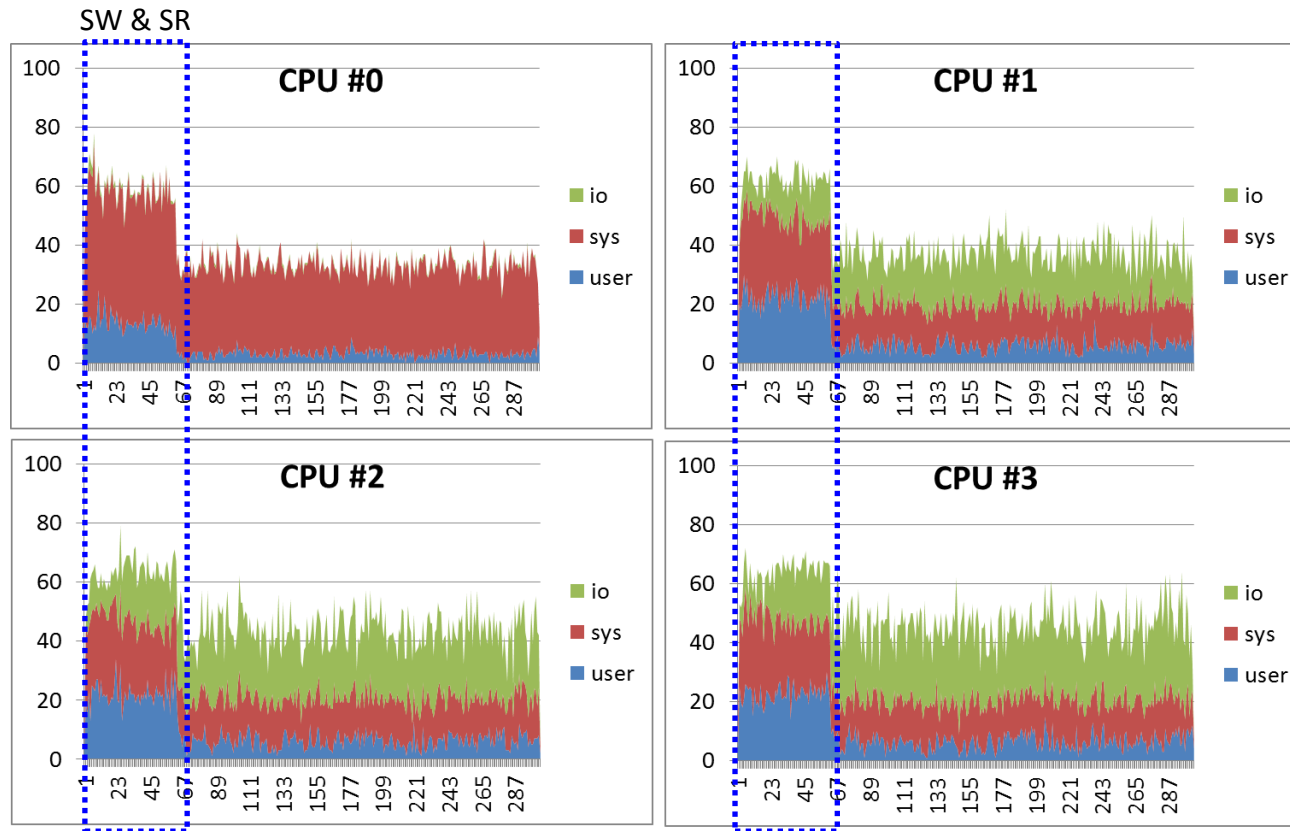  - CPU migration: big << little



<fio running on LITTLE>

<fio running on big>

■ **I/O performance is lower than fio (direct I/O)**

– App keeps migrating among little cores

– CPU utilization is balanced, but is underutilized → **app is slow**

# Improvement Idea

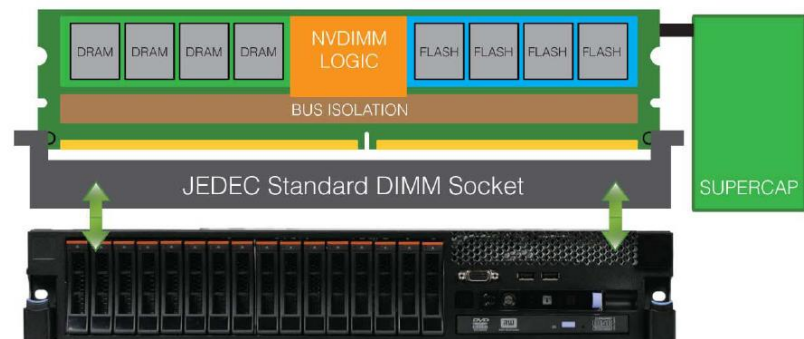- **I/O friendly CPU scheduling**
  - ARM big.LITTLE scheduling is still in work-in-progress
- **Command queueing**
  - End-to-end parallelism by multiple I/O threads or async I/O
  - Good for benchmark vs. real user benefit
- **NVDIMM**
  - Move NVM from I/O bus to memory bus (no DMA!)
    - SNIA NVDIMM SIG (http://www.snia.org/forums/sssi/NVDIMM)
  - **OS & BIOS support is necessary**
    - Linux persistent memory API (https://github.com/pmem/linux-examples)



Legacy Command

Command Queuing (UFS)

&lt;UFS vs. eMMC @ NVRAMOS 2013&gt;

&lt;How it works @ SNIA NVDIMM Tutorial&gt;

# Revisiting Linux I/O Stack

- **Design of Linux I/O**
  - Designed when CPU >> DRAM >> I/O
  - POSIX I/O results in memory operation
  - Buffered I/O, unified VM, DMA, …
- **CPU technology**
  - Clock speed race has been stopped
  - Mobile computing trend puts more emphasis on **power-efficiency**
- **Storage**
  - Flash is much faster than HDD, but still **trying to mimic HDD** (FTL, position in I/O stack)
  - SATA/SCSI → NVMHCI → NVDIMM(?)



The Linux Storage Stack Diagram

Software R&D Center

# Conclusion

■ **What is the bottleneck if flash storage is fast enough?**

– I/O bound: total I/O latency by software barrier – sync(), journaling by FS & DB

– CPU bound: when CPU utilization is not balanced

■ **Which I/O methods to use for benchmark?**

|  | Pros | Cons |
|---|---|---|
| Buffered I/O | Closer to device-level number, less CPU-bound | Need large file for benchmark to get consistent results |
| Direct I/O | Get consistent result in short time | Gap between benchmark and device number, more CPU-bound |

■ **Research trend keeps changing**

– New storage → optimizing SW stack → **new SW & HW architecture**

|  | Improvements | Issues |
|---|---|---|
| CPU bound | Multi-core (homo vs. hetero) | No more CPU clock speed scaling Trends toward power-efficiency |
| I/O bound | Flash memory, I/O stack optimization, clustering | I/O is getting faster Deciding scale-up or scale-out |