

Introduction

Symptom: Suspend/Resume kills battery

Suspend/Resume happens frequently in modern mobile devices. For example, smart watches have to suspend/resume at least once for every one minute to update clock display.

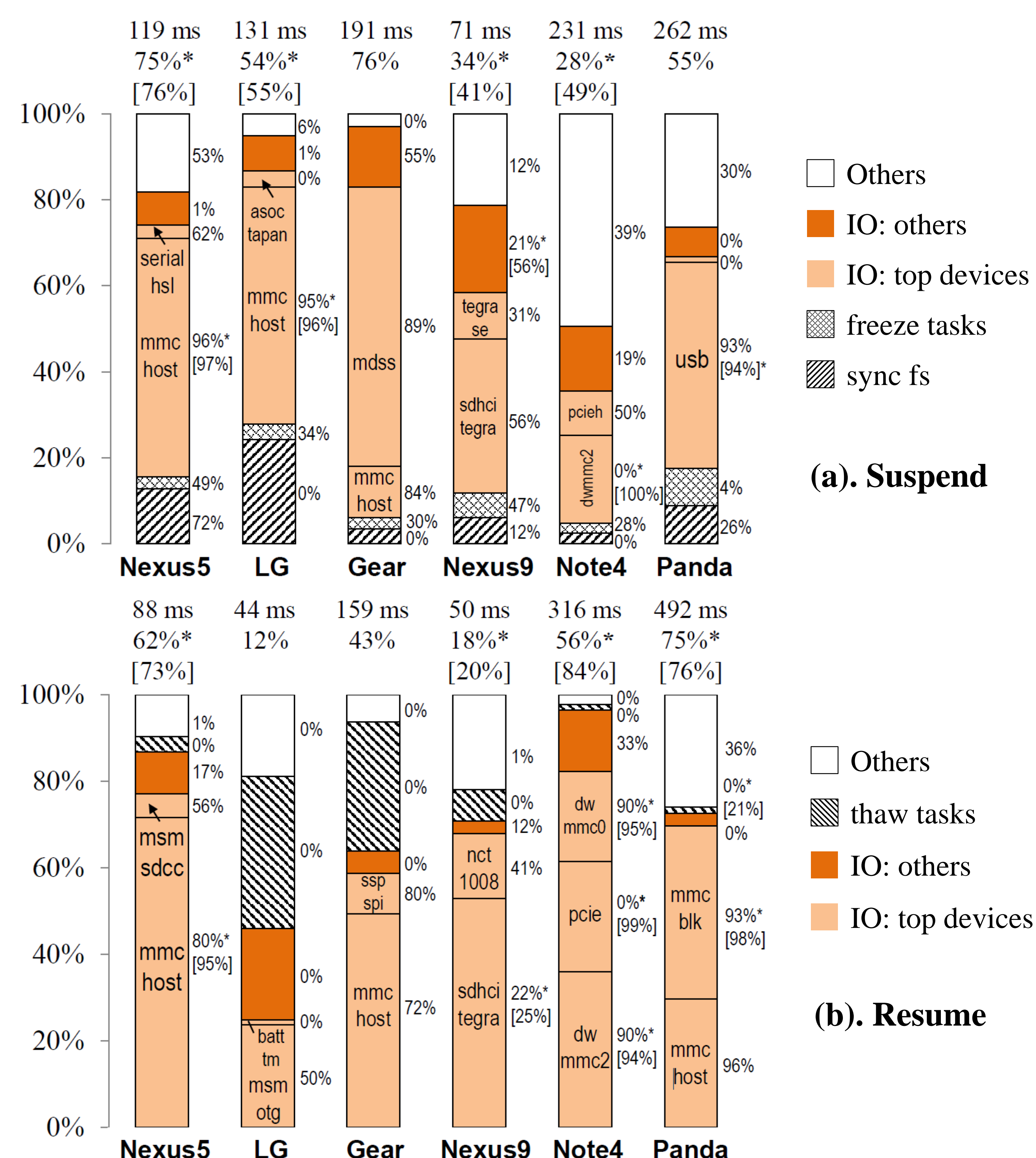
Reason: extensive IO wait

- A single execution of Suspend/Resume can take up to 500ms.
- For most of the time CPU is in idle or busy wait
 - For various IO devices

Solution: Exec. Suspend/Resume on a mini core

Motivation and Measurement

Time breakdown of Suspend/Resume on six mobile/wearable devices:



Methods

Dynamic Binary Translation(DTB)

- Retrofit an open source binary translator QEMU (>1M SLoC)
 - Reduce image size from 6.5MB to 800KB
- Build a first-of-its-kind backend for ARMv7M.
- Deep-optimize QEMU's translation with various optimizations
 - e.g. directly map status registers across translation

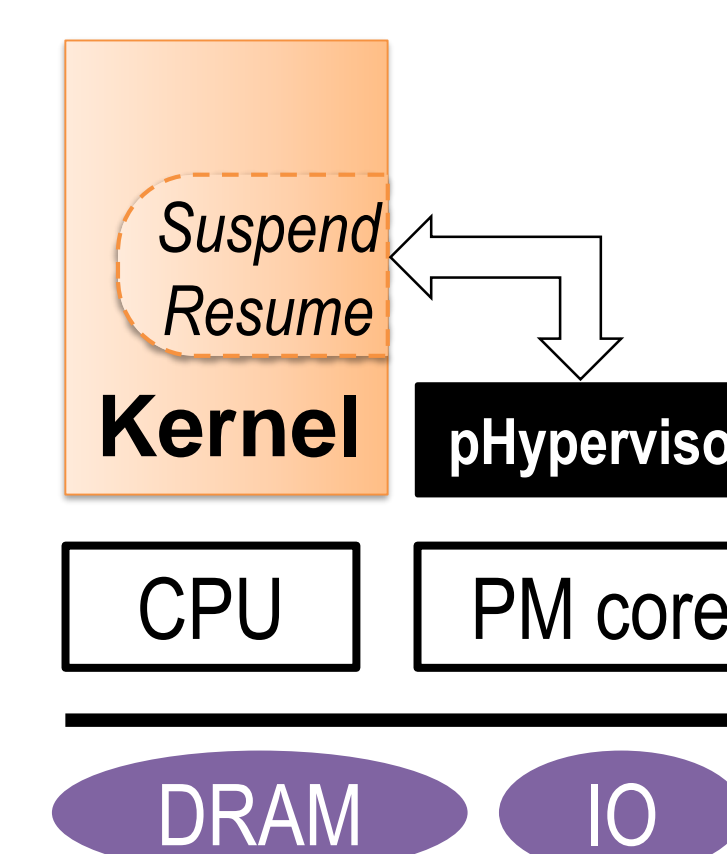
Benchmark: glob_match()

A kernel function for string pattern matching.

Test Platform

TI PandaBoard with TI OMAP 4460 Processor

- 2x ARM Cortex A9 and 2x ARM Cortex M3
- M3 cores use 10~20x less power than A9



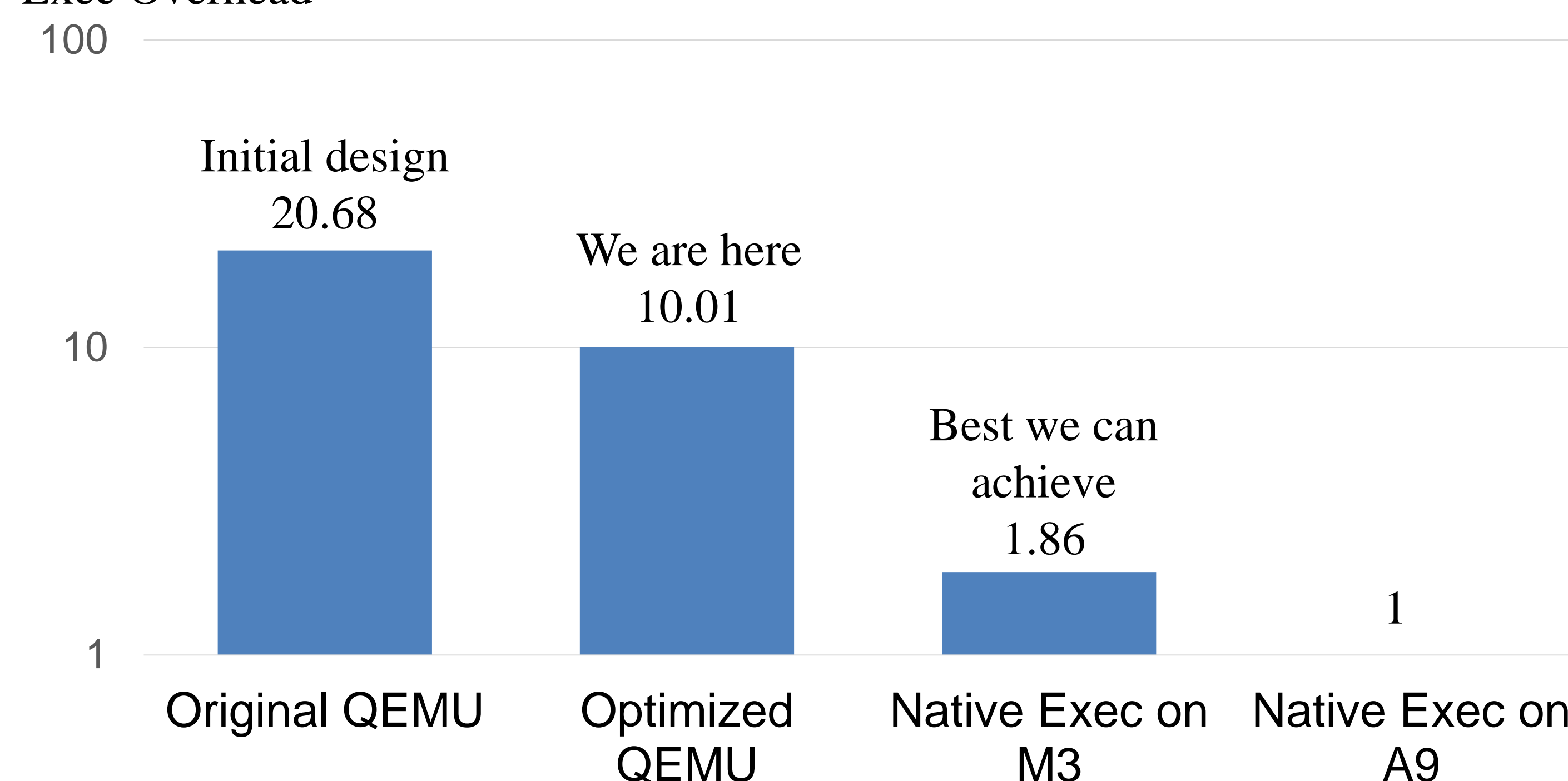
Results

Performance List

	First run	Second run	Long run
Original QEMU	22263	15292	7962
Optimized QEMU	13545	7215	3854
Native Exec on M3	1291	775	718
Native Exec on A9	1546	650	385

*All numbers listed above are in term of number of cycles

Exec Overhead



*Set Cortex A9 native execution as 1, all others are scaled accordingly

Twofold increase in performance

103% fewer instructions as compared to QEMU

Where does the overhead come from?

3.5x due to different ISA

- Not all instructions could be mapped one-to-one to other ISA
- Before translation: 94 instructions
- After translation: 326 instructions

2 cycles wasted for each conditional branch

- Cortex M3 has 3 stage pipeline and does not have branch predictor
- Whenever a conditional branch is taken, the first two pipeline states needs to be flushed

Cache difference

- Cortex M3 has smaller cache than A9
- Cortex A9: L1: 32KB instruction + 32KB data
L2: 1MB unified
- Cortex M3: L1: 32KB unified

Conclusions

- Suspend/Resume is slow across different platforms; it suffers from long wait for IO.
- Suspend/Resume should be offloaded to a mini core that does IO wait much more efficiently.
- We have built a first-of-its-kind binary translator that adopts key optimizations for 2x performance increase.

<http://xsel.rocks>

