# Power and Energy Efficiency in Modern Processors

Aleksandar Milenković

November 30, 2021

Based on dissertation of Ranjan Hebbar,
VMware, Palo Alto, CA

**Department of Electrical and Computer Engineering**
**University of Alabama in Huntsville**

LaCASA

# Overview

- Introduction

- Background

- Experimental Setup

- SPEC CPU2017 Characterization

- Proposed DVFS Techniques

- Results
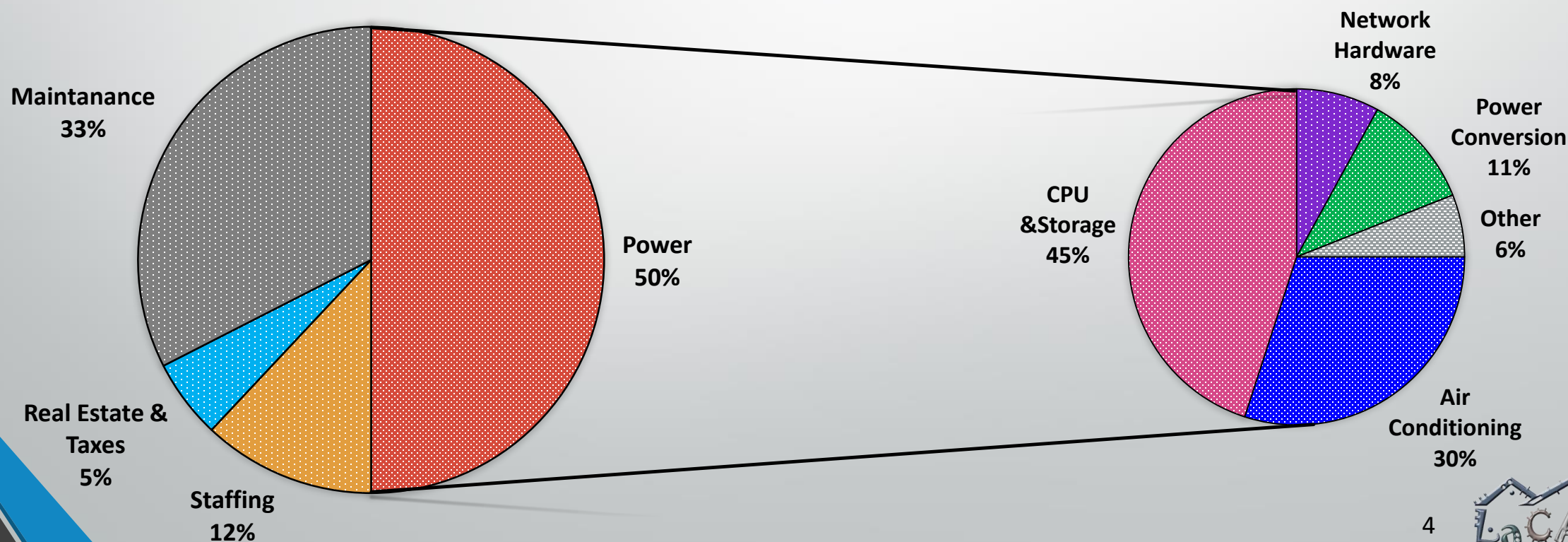
- Conclusions

# Introduction

- Technology
  - Moore's Law and Dennard's Law
  - Processor advances: Caches, Out-of-Order, Vector Processing, Prefetching
- Applications
  - Increasing sophistication and complexity
  - Machine Learning, Block Chains, Information Security, Big Data
  - Software developers are not expected to be familiar with hardware internals
- Markets
  - Shortening time-to-market
  - Demand for high performance computing and data center capacity have grown exponentially

# Typical Server Operating Cost Distribution

THE UNIVERSITY OF
ALABAMA IN HUNTSVILLE

Energy consumption challenge
- Data centers in the US consume about 2% of the country's total energy consumption (per US DoE)
- Forecast to reach ~8% by 2030 (Andrae and Edler, 2015)

## DATACENTER OPERATING COST



Maintanance 33%

Power 50%

Real Estate & Taxes 5%

Staffing 12%

Network Hardware 8%

Power Conversion 11%

CPU &Storage 45%

Other 6%

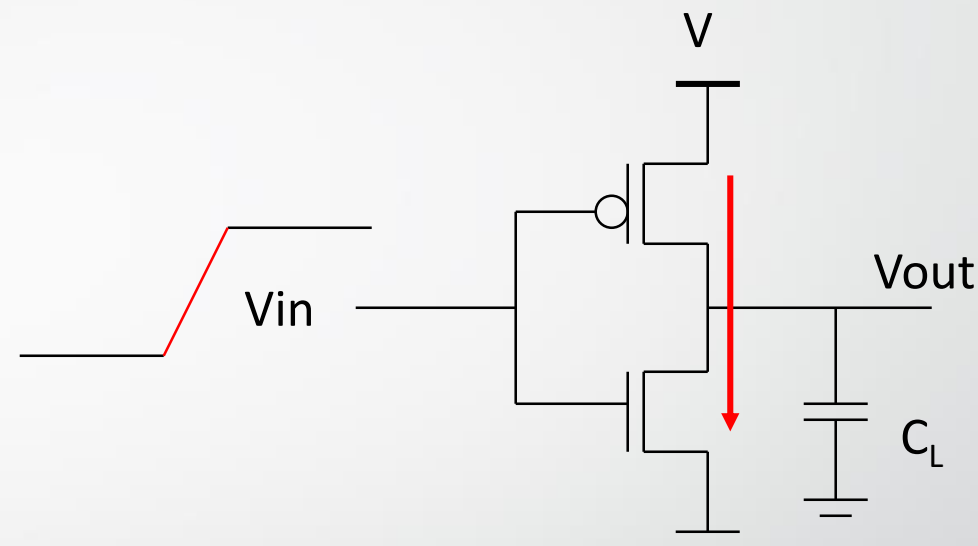Air Conditioning 30%

4

# State-of-the-art

- A majority of the servers in data centers currently utilize x86 processors

- x86 processors are complex structures integrating

  - Multiple physical cores

  - On-chip Interconnect

  - Uncore Cache

  - Memory Controllers

  - Hardware Accelerators

- Dedicated hardware resources for monitoring and management of its operating states

- DVFS: Most powerful technique in regulating CPU power consumption

# Processor Power Dissipation

- Three factors

  - Dynamic power (switching)

  - Short-circuit power

  - Leakage power

- $P_{CPU} = P_{dyn} + P_{sc} + P_{leak}$

- $P_{CPU} = ACV^2 f \ + \ A\tau I_{sc} V f \ + \ VI_{leak}$

# Dynamic Power

THE UNIVERSITY OF
ALABAMA IN HUNTSVILLE

C – Total capacitance
seen by the gate's outputs
Function of wire lengths,
transistor sizes, …

V – Supply voltage
Trend: has been dropping
with each successive fab

$$ACV^2f$$

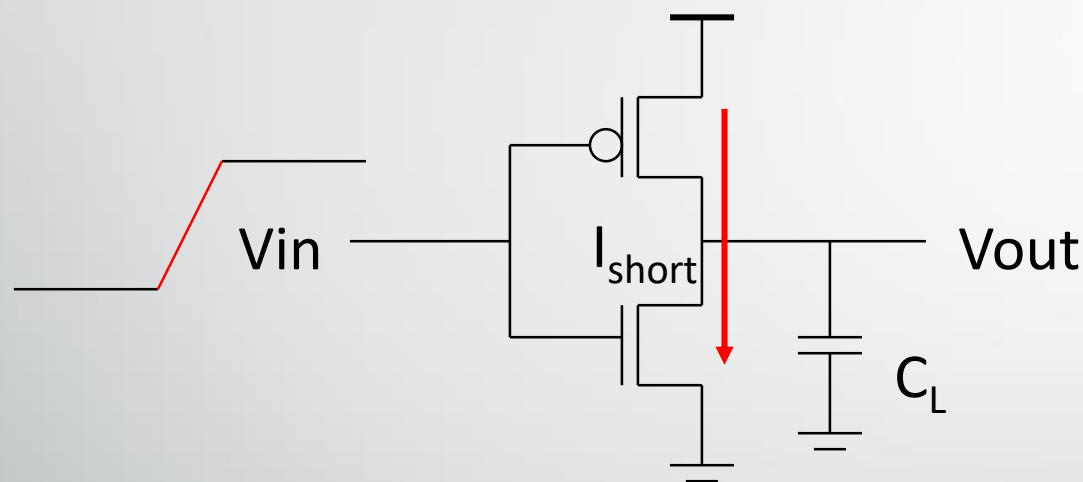A - Activity of gates
How often on average do
wires switch?

f – clock frequency
Trend: increasing? …

Reducing Dynamic Power

1) Reducing V has quadratic effect; Limits?

2) Lower C - shrink structures, shorten wires

3) Reduce switching activity - Turn off unused parts or
use design techniques to minimize number of transitions
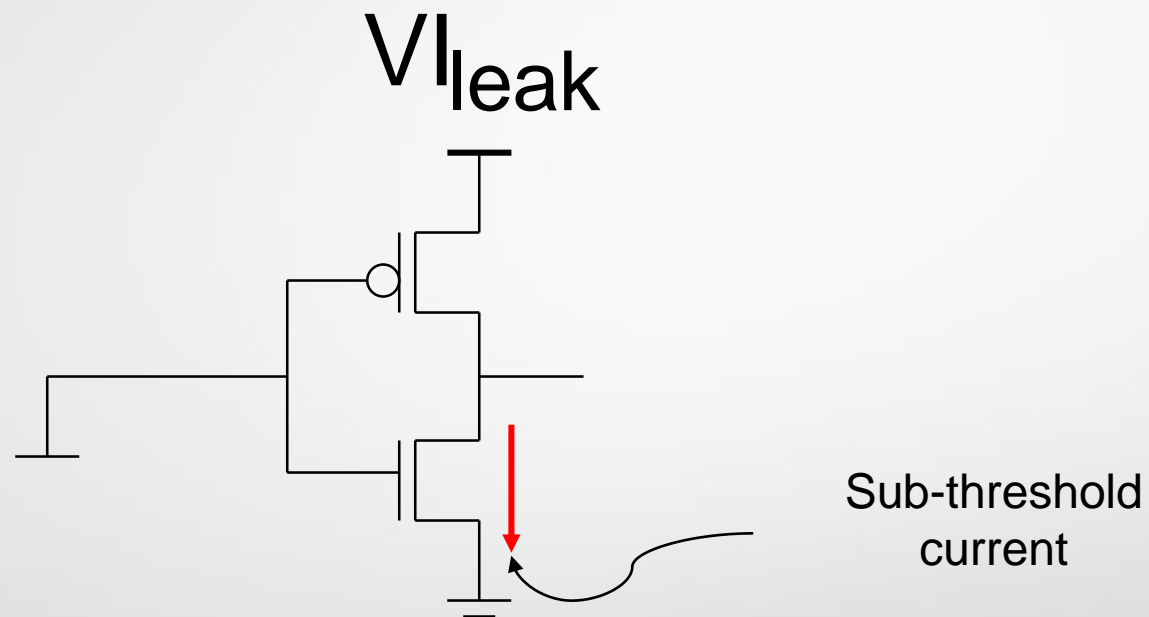
# Short-Circuit Power

$$\tau AVI_{short}f$$

Finite slope of the input signal causes a direct current path between $V_{DD}$ and GND for a short period of time during switching when both the NMOS and PMOS transistors are conducting

Vin — $I_{short}$ — Vout

$C_L$

Reducing Short-circuit

1)  Lower the supply voltage V

2)  Slope engineering – match the rise/fall time of the input and output signals

8

# Leakage Power

$$VI_{leak}$$



Sub-threshold current

Sub-threshold current grows exponentially with increases in temperature and decreases in Vt

# CMOS Power Tradeoffs

$$P = \boxed{ACV^2f} + \tau AVI_{short}f + VI_{leak}$$

Reduce the supply voltage, V

$$f_{max} \propto \frac{(V - V_t)^2}{V}$$

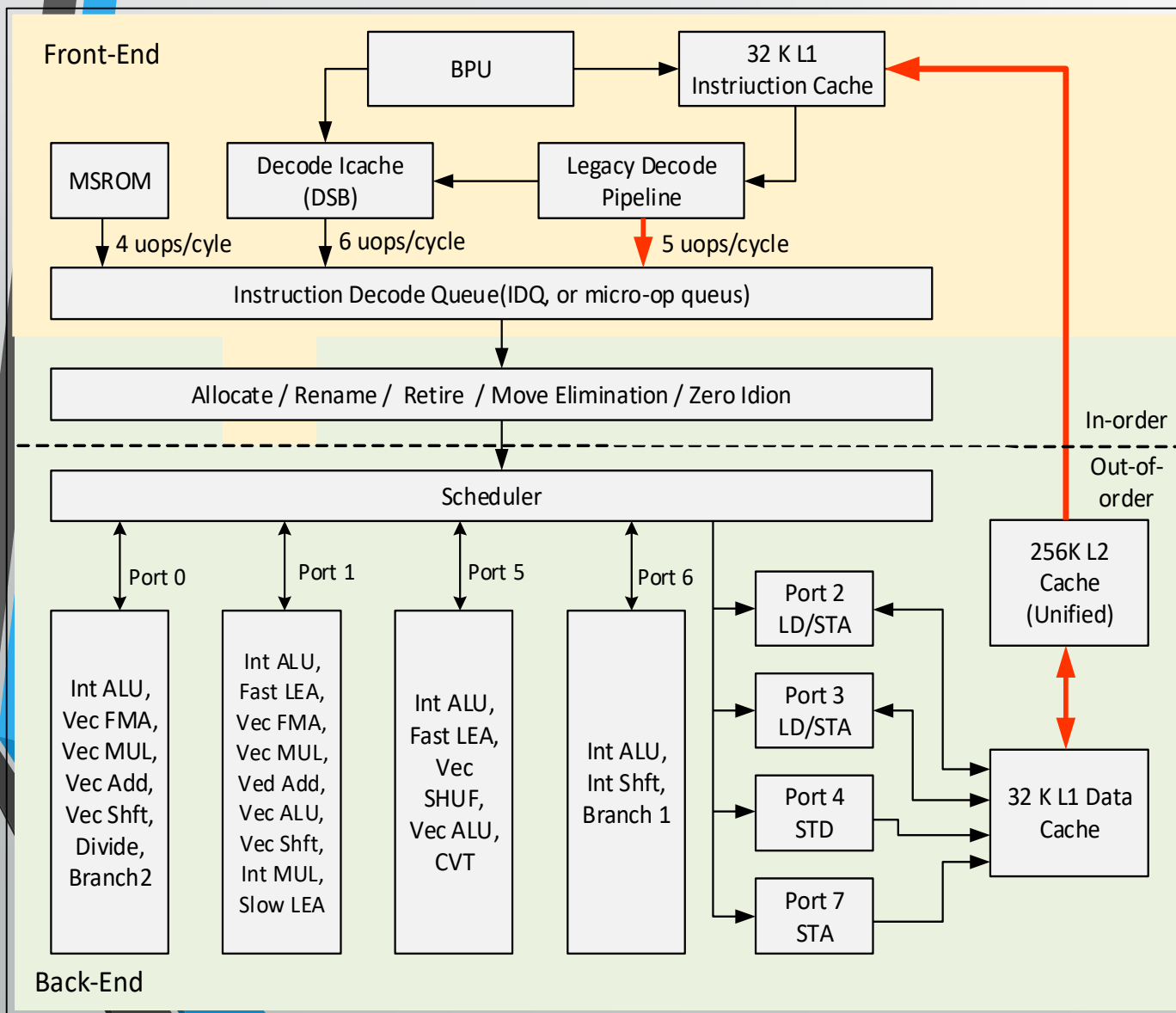$$I_{leak} \propto \exp\left(-\frac{qV_t}{kT}\right)$$

Reduce threshold $V_t$

# Overview

- Introduction

- **Background**

- Experimental Setup

- SPEC CPU2017 Characterization

- Proposed DVFS Techniques

- Results

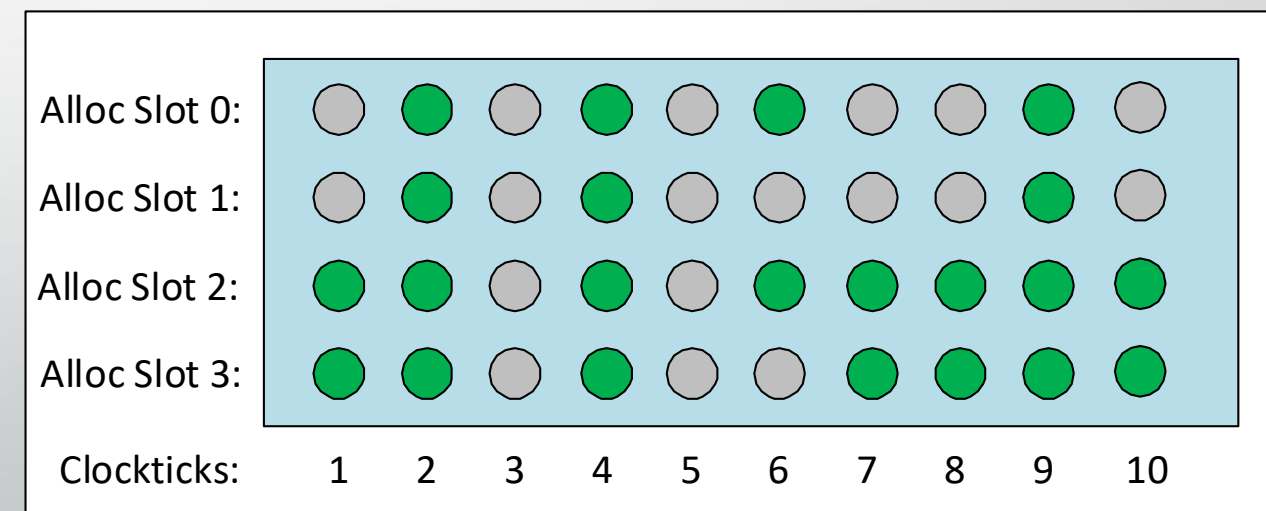- Conclusions

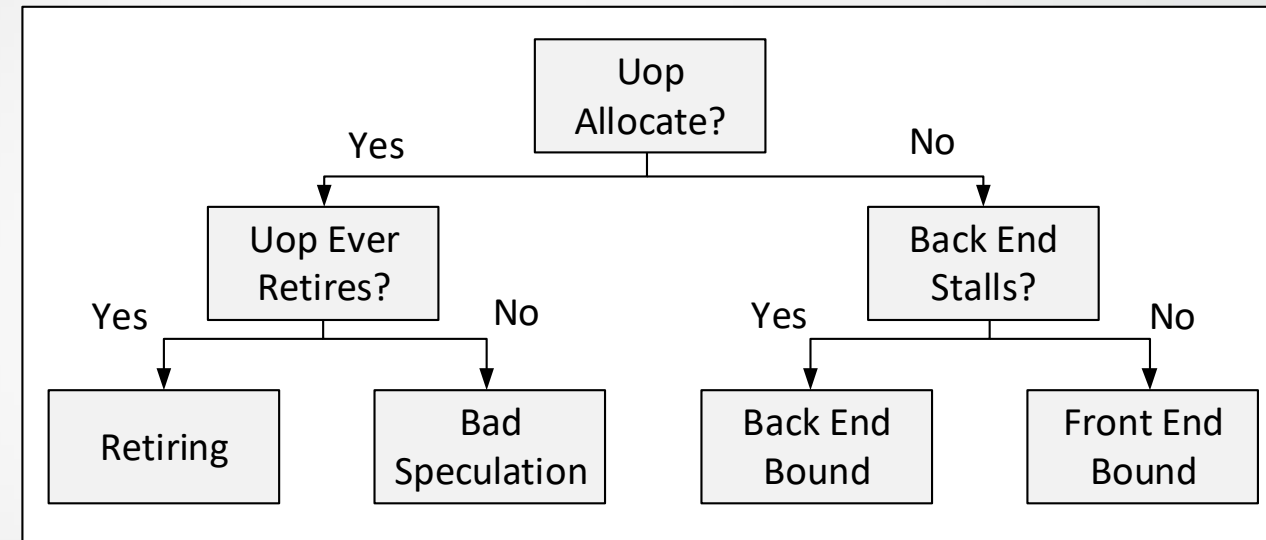# Intel Skylake Microarchitecture



- 19-stage pipeline

- Superscalar & Multi-threaded

- Front-end (in-order)
  - Fetch from memory (L1 Instruction Cache)
  - Decode: Convert to micro-operation (uops)

- Back-end (out-of-order)
  - Schedule micro-operations
  - Execute
  - Retire

12

# Top-down Microarchitectural Analysis (TMAM)

THE UNIVERSITY OF ALABAMA IN HUNTSVILLE

## Issue Stage of the Pipeline
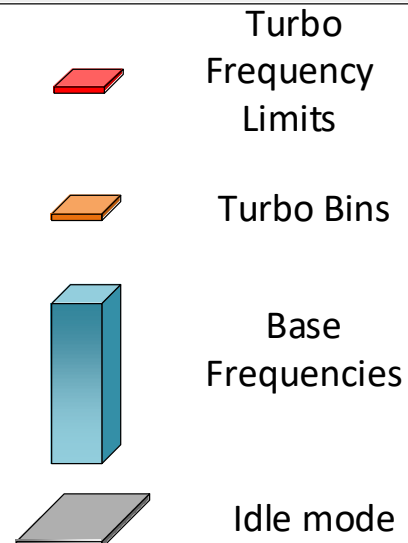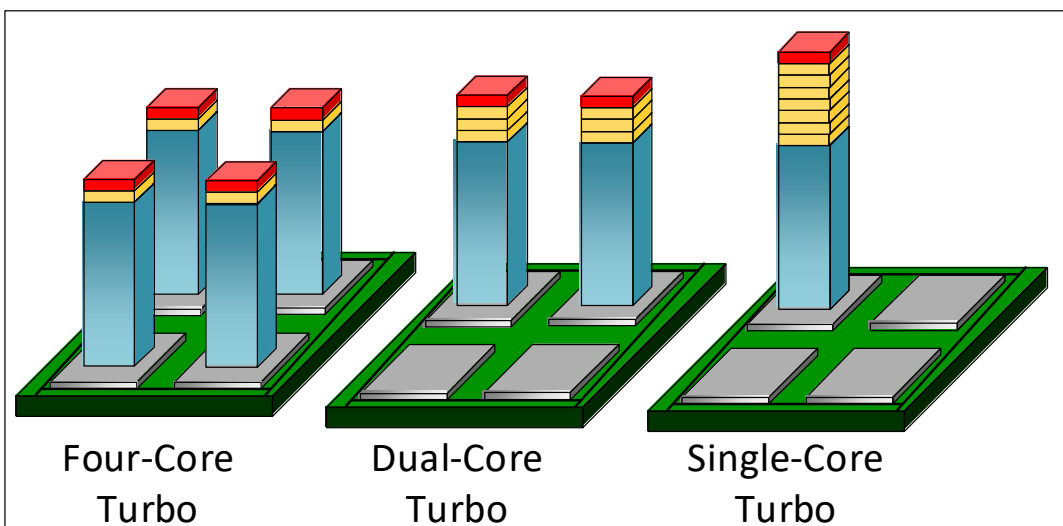
- Each stalled pipeline slot can be attributed to one of three causes
  - Front-end Bound
  - Back-end Bound
  - Bad Speculation

- Ideally high retiring is preferred

- Example: 10 Clock cycles
  - 40 available pipeline slots (4-wide pipeline)
  - Micro-operations were issues for 22 slots
  - Pipeline Slot Utilization: (22/40) : **0.55**
  - Clock Cycle Utilization: (8/10) : **0.80**



13

# Hardware Power Management Features

- External Clock (BCLK): 133 MHz
- 4 Voltage/Power Domains
  - Core Domain
  - Uncore Domain
  - Graphics Domain
  - Fixed System Agent
- Thermal Design Power (TDP)

# ACPI* Power & Performance States



- C-States: Power States or Idle States

  - C1 to Cn: No instructions are executed

  - Pros: turns off sections of the CPU $\Rightarrow$ reduce idle power

  - Cons: introduces latency to get back to operational state

- P-States: Performance/Operational States (C0 Only)

  - Each P-State directs the processor to operate at a particular clock frequency and CPU voltage

  - P-States are the basis for Dynamic Voltage & Frequency Scaling (DVFS)

  - P-States is generally assigned based on CPU load

*ACPI – *Advanced Configuration Power Interface*

15

# BIOS & OS Power Profiles

THE UNIVERSITY OF
ALABAMA IN HUNTSVILLE

| Userspace Interface |
|---|

| Governors | | | | | |
|---|---|---|---|---|---|
| performance [1] | powersave [1] | conservative [1] | ondemand [1] | schedutil [1] | userspace [1] |
| performance [2] | | | powersave [2] | | |

- - - - - - - - - - - - - - - - - - - - - - - - -

**Device Drivers**

CPUFreq Kernel Submodule

acpi-cpufreq [1]   intel-pstate [2]

**BIOS Power Management**

System Profiles

OS Control Mode

BIOS Firmware

- - - - - - - - - - - - - - - - - - - - - - - - -

**Hardware**

| Control & Status Register (CSR) | Model Specific Registers (MSR) |
|---|---|

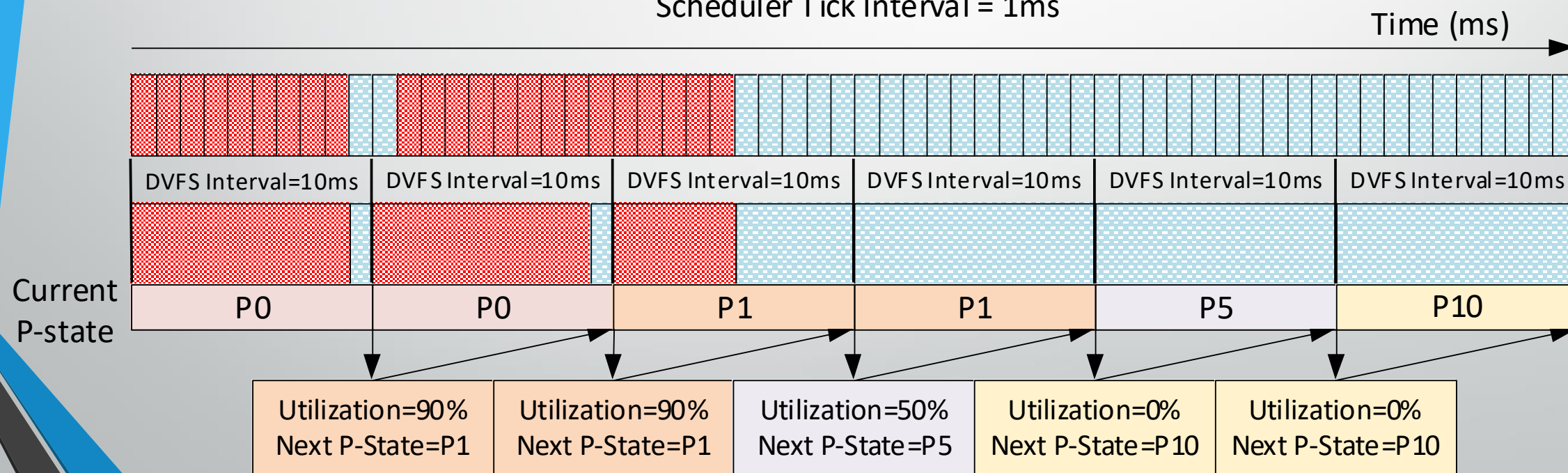| PCU/P-Unit |
|---|

16

# Utilization Based P-State Selection

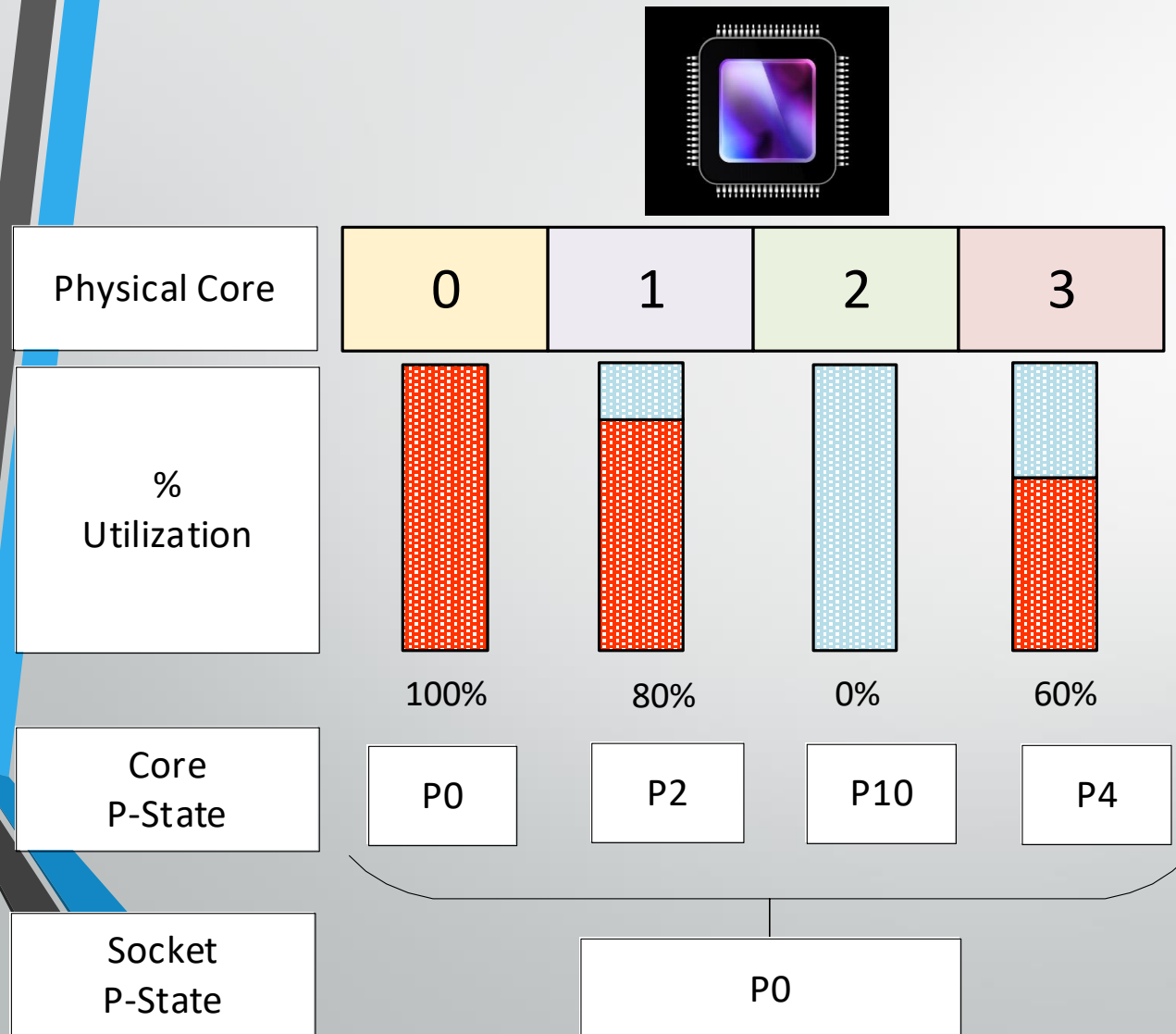*Ondemand* **governor functioning**

- Schedular provides CPU utilization metric

$$\% \text{ CPU Utilization} = \frac{time\ in\ non\ idle\ thread}{time\ duration} * 100$$

- Scheduler update interval – 1 ms to 10 ms
  - Linux default scheduler tick is 1 ms
- DVFS Interval: 10ms to 100ms
  - Linux Default is 10ms

Scheduler Tick Interval = 1ms

Time (ms)

| DVFS Interval=10ms | DVFS Interval=10ms | DVFS Interval=10ms | DVFS Interval=10ms | DVFS Interval=10ms | DVFS Interval=10ms |
|---|---|---|---|---|---|

Current P-state

| P0 | P0 | P1 | P1 | P5 | P10 |
|---|---|---|---|---|---|

| Utilization=90% Next P-State=P1 | Utilization=90% Next P-State=P1 | Utilization=50% Next P-State=P5 | Utilization=0% Next P-State=P10 | Utilization=0% Next P-State=P10 |
|---|---|---|---|---|

17

# P-State Voting Mechanism

THE UNIVERSITY OF
ALABAMA IN HUNTSVILLE

| Physical Core | 0 | 1 | 2 | 3 |
|---|---|---|---|---|

% Utilization

100%　　　80%　　　0%　　　60%

| Core P-State | P0 | P2 | P10 | P4 |
|---|---|---|---|---|

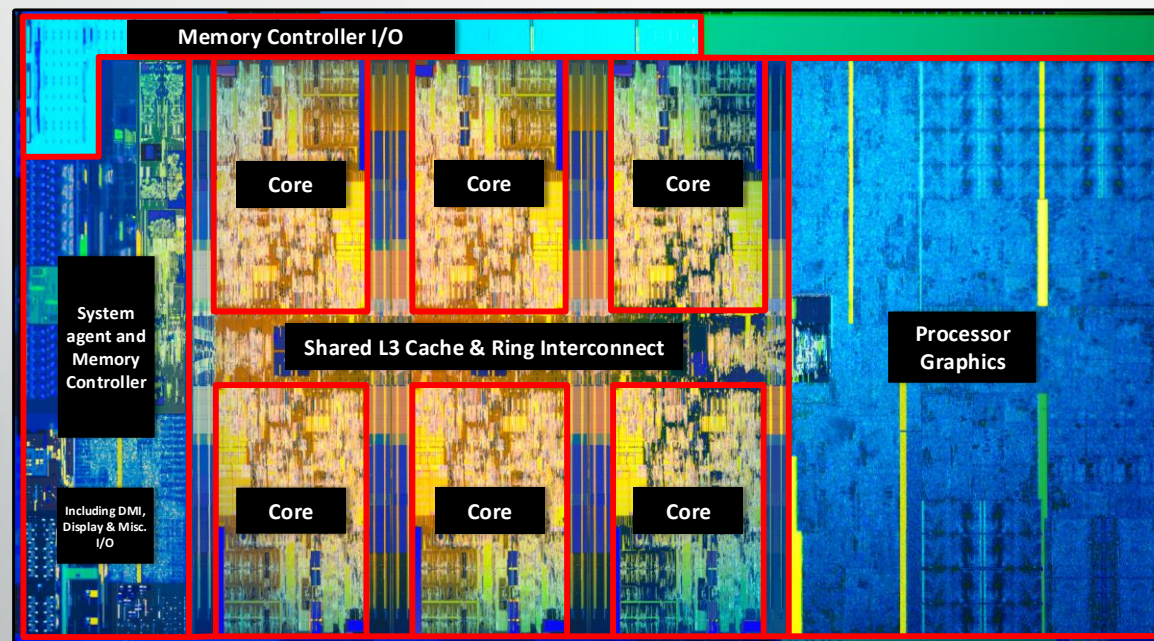| Socket P-State | P0 |
|---|---|

- P-State Management
  - Evaluate utilization for each core
  - Select P-state using linear mapping of CPU utilization to clock frequency
  - Lower P-state of logical cores is used to determine P-state of a physical core
- Processors with multiple voltage domains can have physical cores in different P-states
- Processor without multiple voltage domains select lowest P-state

# Overview

- Introduction

- Background

- Experimental Setup

- SPEC CPU2017 Characterization

- Proposed DVFS Techniques

- Results

- Conclusions

# Test Machine

- 8th Generation Intel Core i7-8700K (6 Physical Cores)

- 32 GiB of DRAM Memory: 2 Channels with a total bandwidth of 41.6 GB/s

- Coffee Lake, manufactured using Intel's 14nm++ technology node

- Nominal clock frequency of 3.70 GHz and all core turbo of 4.30 GHz

- State-of-the-art P-state management (40 P-states)



*[source-https://en.wikichip.org/wiki/intel/core_i7/i7-8700k]*

20

# Metrics

- Performance Speedup

$$P.S\ (B_i, PG_{GOV}) = \frac{T\left(B_i, OD_{GOV}\right)}{T\left(B_i,\ PG_{GOV}\right)}$$

- Energy Efficiency Improvement

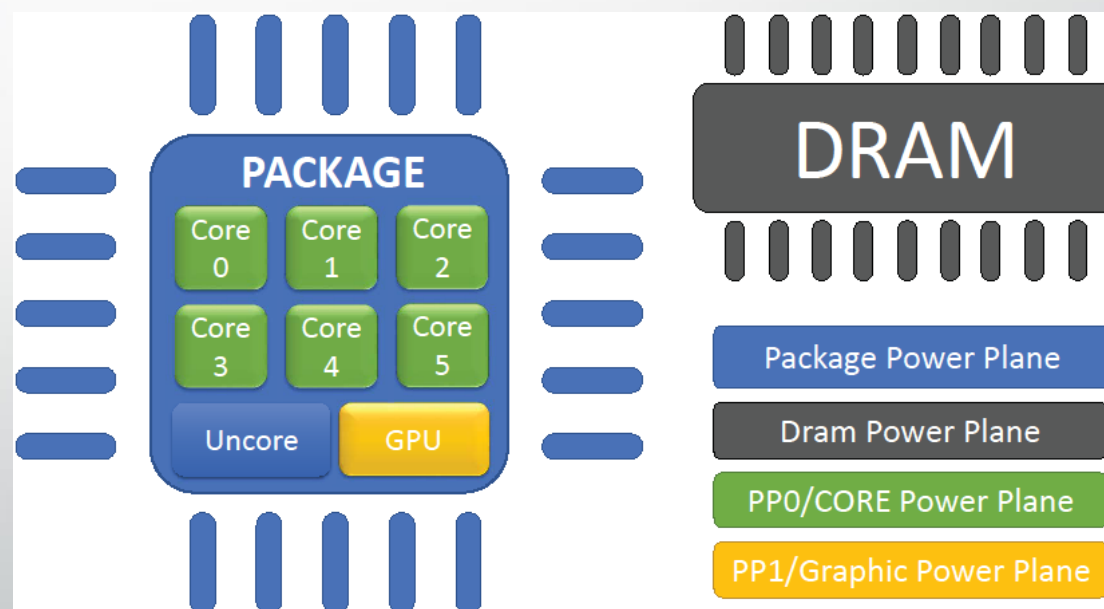$$EE.I\ (Bi, PG_{GOV}) = \frac{E(B_i, OD_{GOV})}{E(B_i, PG_{GOV})}$$

- PxEE Improvement

$$PxEE.I(B_i, PG_{GOV}) = \frac{T(B_i, OD_{GOV}) * E(B_i, OD_{GOV})}{T(B_i, PG_{GOV}) * E(B_i, PG_{GOV})}$$

# Tools: likwid

THE UNIVERSITY OF
ALABAMA IN HUNTSVILLE

- LIKWID (Like I Knew What I'm Doing): Has a set of tools with specific purposes to measure performance/energy groups

- Utilizes the RAPL counters for measuring power and energy

```
--------------------------------------------------
CPU name:       Intel(R) Core(TM) i7-8700K CPU @ 3.70GHz
CPU type:       Intel Kabylake processor
CPU clock:      3.70 GHz
--------------------------------------------------
--------------------------------------------------
Runtime: 1356.99 s
Measure for socket 0 on CPU 0
Domain PKG:
Energy consumed: 23661.9 Joules
Power consumed: 17.4371 Watt
Domain PP0:
Energy consumed: 14451 Joules
Power consumed: 10.6493 Watt
Domain PP1:
Energy consumed: 0 Joules
Power consumed: 0 Watt
Domain DRAM:
Energy consumed: 0 Joules
Power consumed: 0 Watt
--------------------------------------------------
```



22

# Workloads

- SPEC CPU2017
  - The SPEC CPU2017 benchmark suites contains standardized, CPU intensive suites for measuring and comparing performance
  - Stresses a system's processor, memory subsystem, and compiler

- SPECpower_2008ssj
  - SPECpower_ssj2008 is an industry-standard benchmark designed for experimental power and performance evaluation of server computers
  - Varying transactional load from 100% to active idle in steps of 10%

# Overview

- Introduction

- Background

- Experimental Setup

- SPEC CPU2017 Characterization
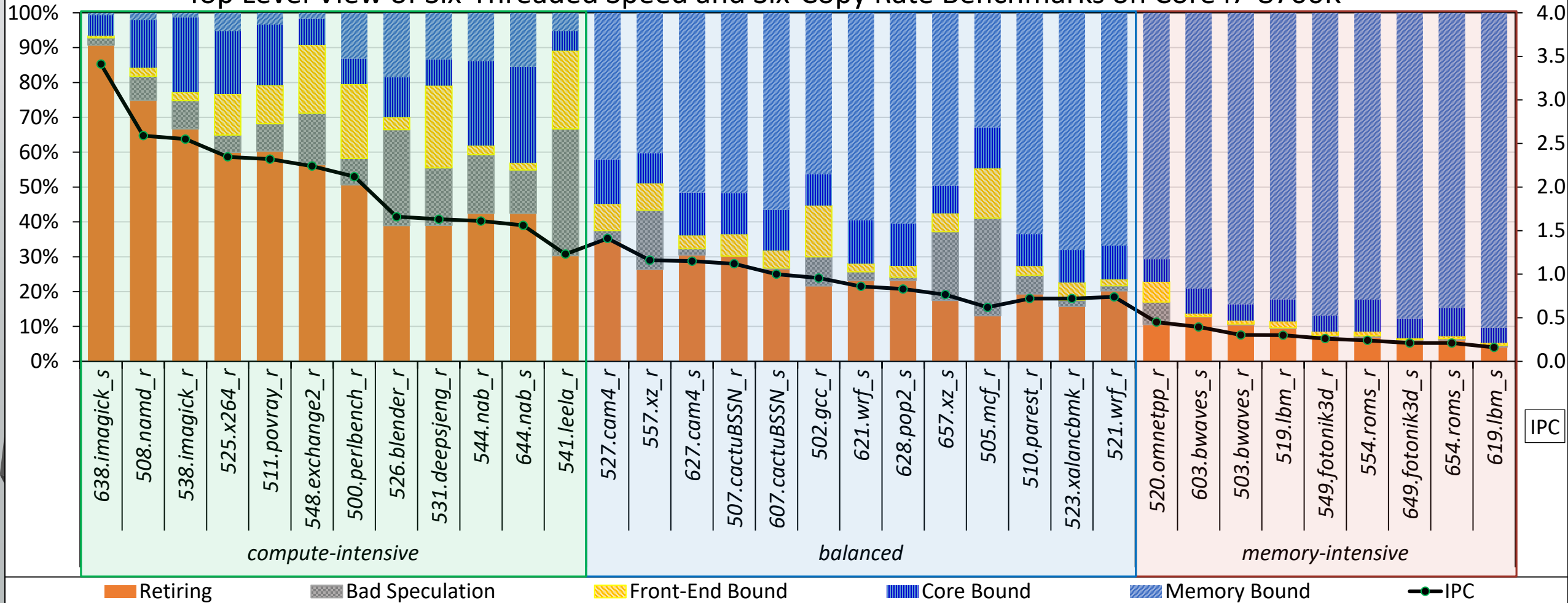
- Proposed DVFS Techniques

- Results

- Conclusions

# SPEC CPU2017 Benchmarks

| SPECrate 2017 Floating Point | SPECspeed 2017 Floating Point | Language | Application Area |
|---|---|---|---|
| 503.bwaves_r | 603.bwaves_s | Fortran | Explosion modeling |
| 507.cactuBSSN_r | 607.cactuBSSN_s | C++, C, Fortran | Physics: relativity |
| 508.namd_r | | C++ | Molecular dynamics |
| 510.parest_r | | C++ | Biomedical imaging: optical tomography |
| 511.povray_r | | C++, C | Ray tracing |
| 519.lbm_r | 619.lbm_s | C | Fluid dynamics |
| 521.wrf_r | 621.wrf_s | Fortran, C | Weather forecasting |
| 526.blender_r | | C++, C | 3D rendering and animation |
| 527.cam4_r | 627.cam4_s | Fortran, C | Atmosphere modeling |
| | 628.pop2_s | Fortran, C | Wide-scale ocean modeling (climate level) |
| 538.imagick_r | 638.imagick_s | C | Image manipulation |
| 544.nab_r | 644.nab_s | C | Molecular dynamics |
| 549.fotonik3d_r | 649.fotonik3d_s | Fortran | Computational Electromagnetics |
| 554.roms_r | 654.roms_s | Fortran | Regional ocean modeling |

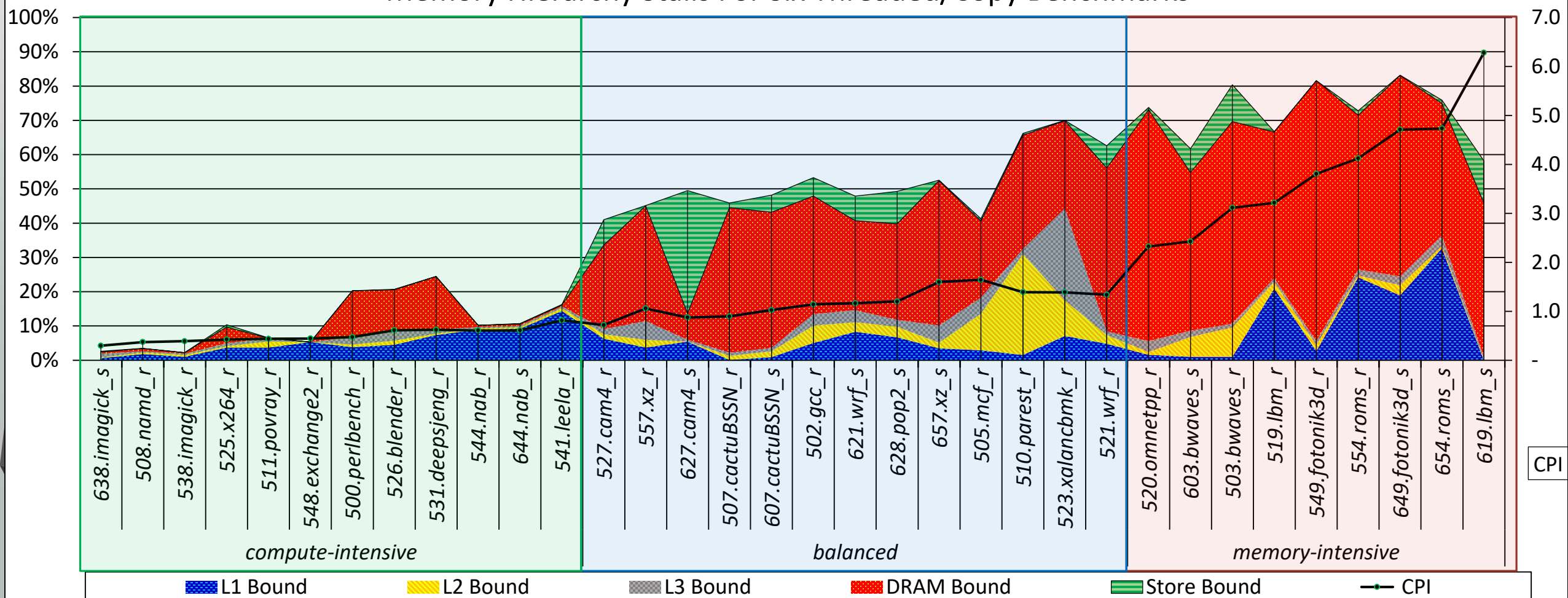| SPECrate 2017 Integer | SPECspeed 2017 Integer | Language | Application Area |
|---|---|---|---|
| 500.perlbench_r | 600.perlbench_s | C | Perl interpreter |
| 502.gcc_r | 602.gcc_s | C | GNU C compiler |
| 505.mcf_r | 605.mcf_s | C | Route planning |
| 520.omnetpp_r | 620.omnetpp_s | C++ | Discrete Event simulation - computer network |
| 523.xalancbmk_r | 623.xalancbmk_s | C++ | XML to HTML conversion via XSLT |
| 525.x264_r | 625.x264_s | C | Video compression |
| 531.deepsjeng_r | 631.deepsjeng_s | C++ | Artificial Intelligence: alpha-beta tree search (Chess) |
| 541.leela_r | 641.leela_s | C++ | Artificial Intelligence: Monte Carlo tree search (Go) |
| 548.exchange2_r | 648.exchange2_s | Fortran | Artificial Intelligence: recursive solution generator (Sudoku) |
| 557.xz_r | 657.xz_s | C | General data compression |

25

# Top-Down Analysis of SPEC CPU2017 Benchmarks

THE UNIVERSITY OF ALABAMA IN HUNTSVILLE



Top Level View of Six-Threaded Speed and Six-Copy Rate Benchmarks on Core i7-8700K

Legend: Retiring, Bad Speculation, Front-End Bound, Core Bound, Memory Bound, IPC

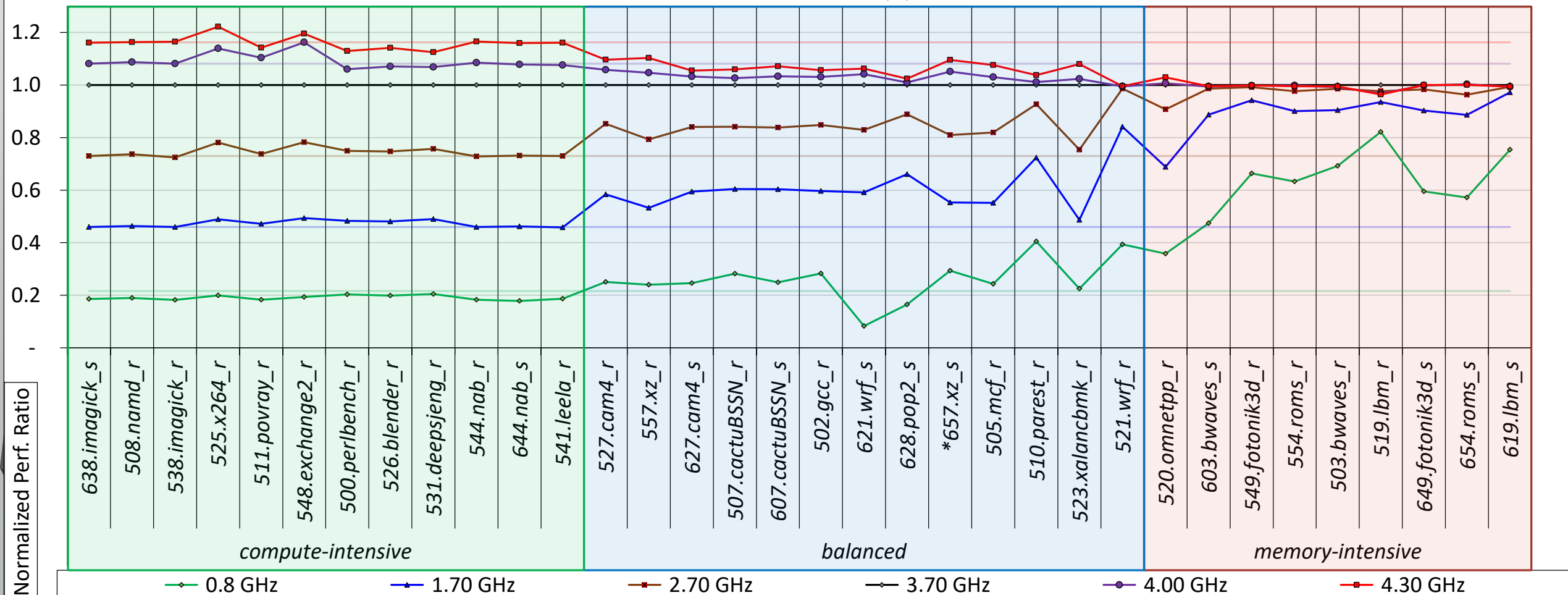Categories: compute-intensive, balanced, memory-intensive

# Memory Hierarchy Related Cycle Stall Breakdown
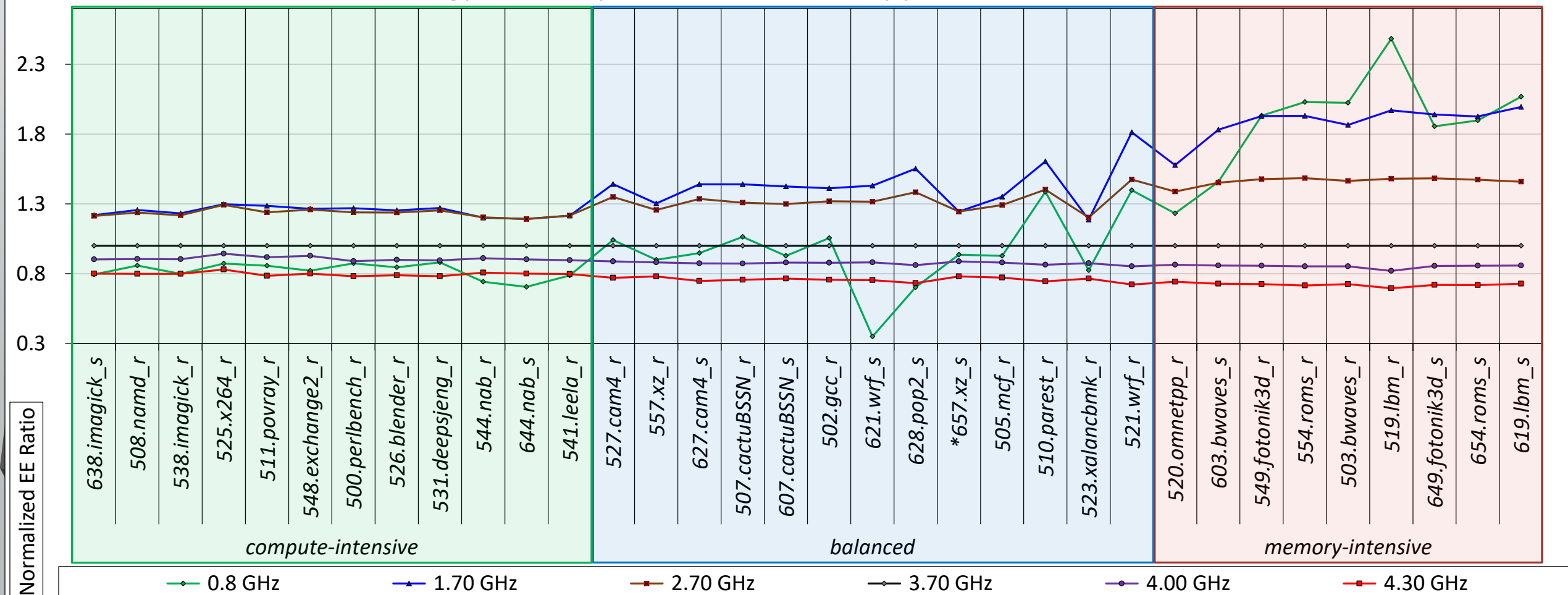


Memory Hierarchy Stalls For Six Threaded/Copy Benchmarks

# Normalized P for SPEC CPU2017 as a Function of Clock Frequency



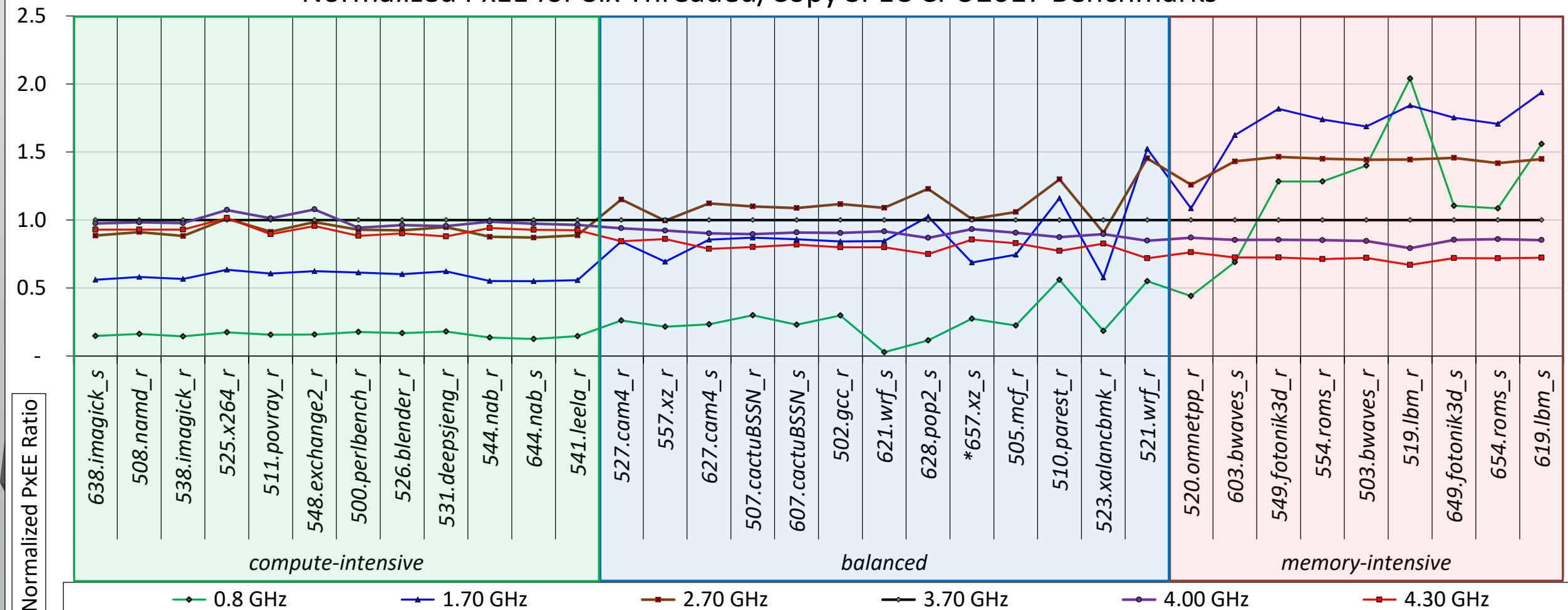Normalized Performance for Six Threaded/Copy SPEC CPU2107 Benchmarks

compute-intensive | balanced | memory-intensive

Legend: 0.8 GHz | 1.70 GHz | 2.70 GHz | 3.70 GHz | 4.00 GHz | 4.30 GHz

Normalized Perf. Ratio

# Normalized EE for SPEC CPU2017 as a Function of Clock Frequency



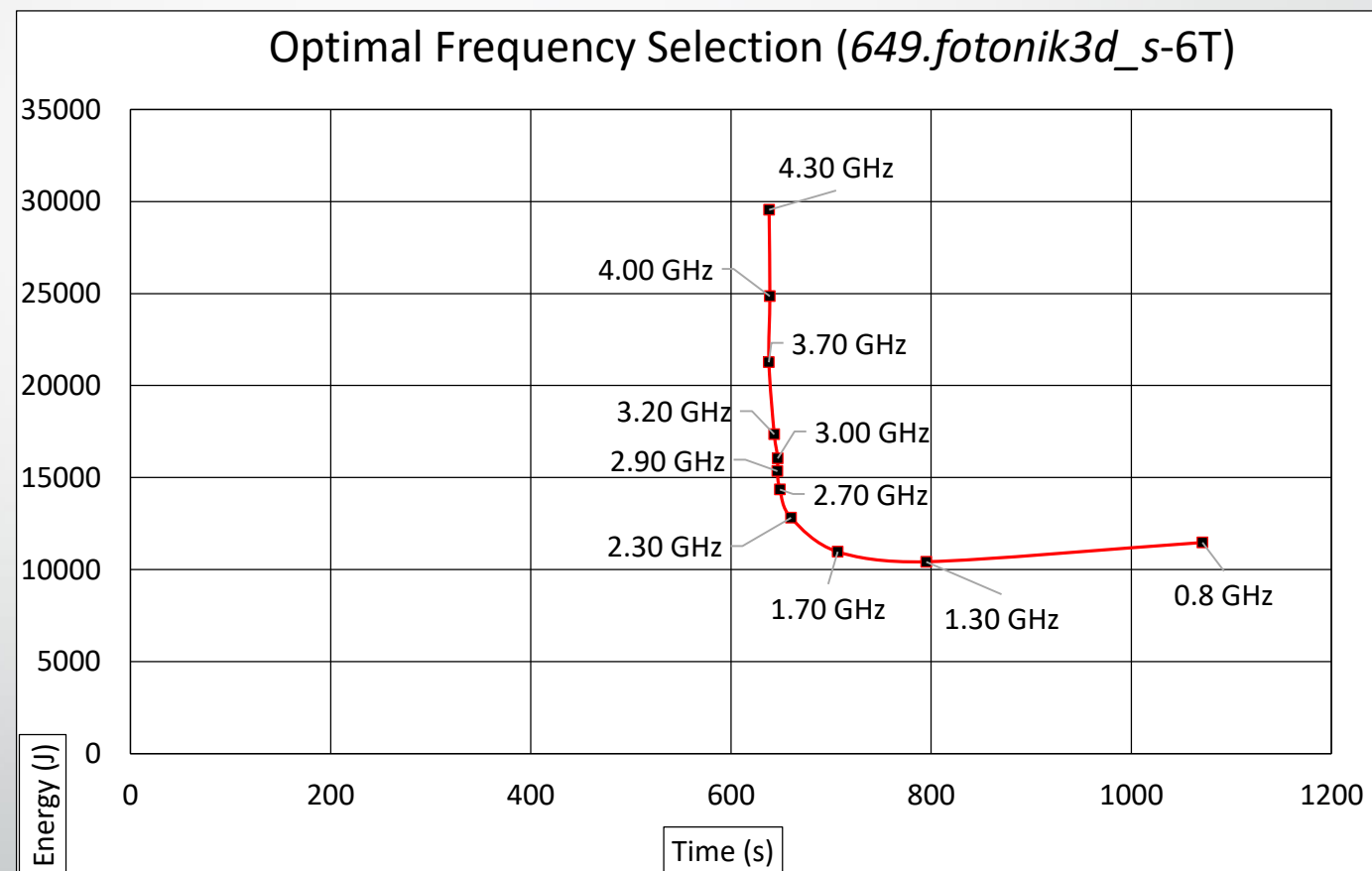Normalized Energy Efficiency for Six Threaded/Copy CPU2017 Benchmarks

# Normalized PxEE for SPEC CPU2017as a Function of Clock Frequency



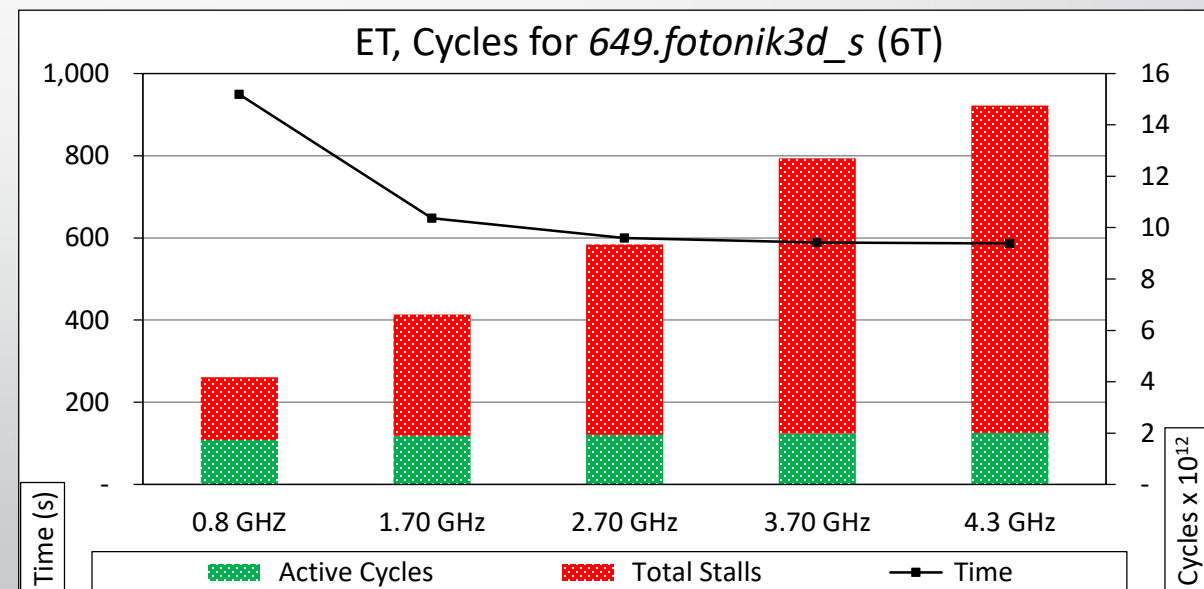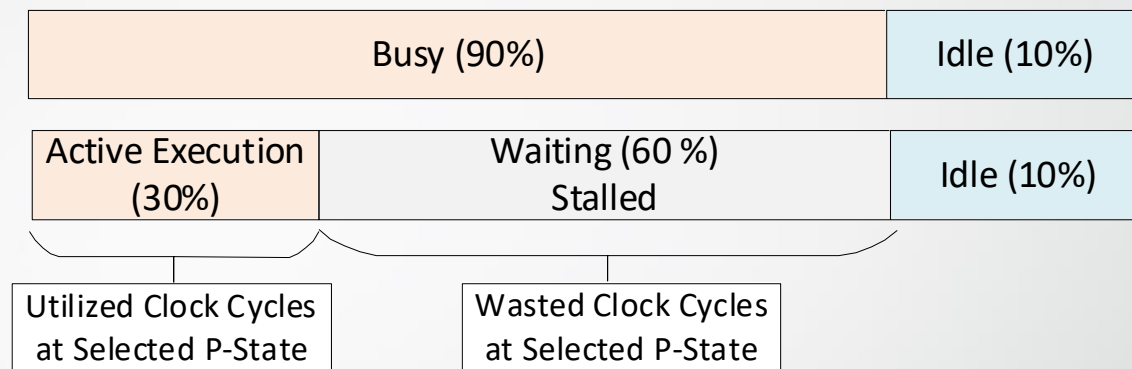Normalized PxEE for Six Threaded/Copy SPEC CPU2017 Benchmarks

# Optimal Frequency Selection
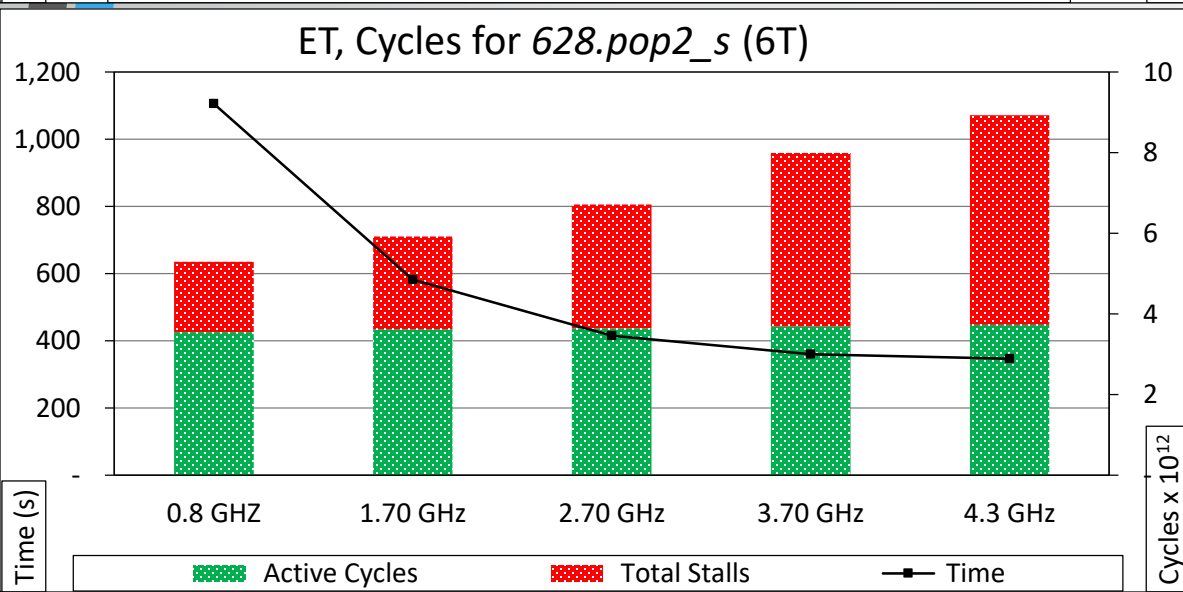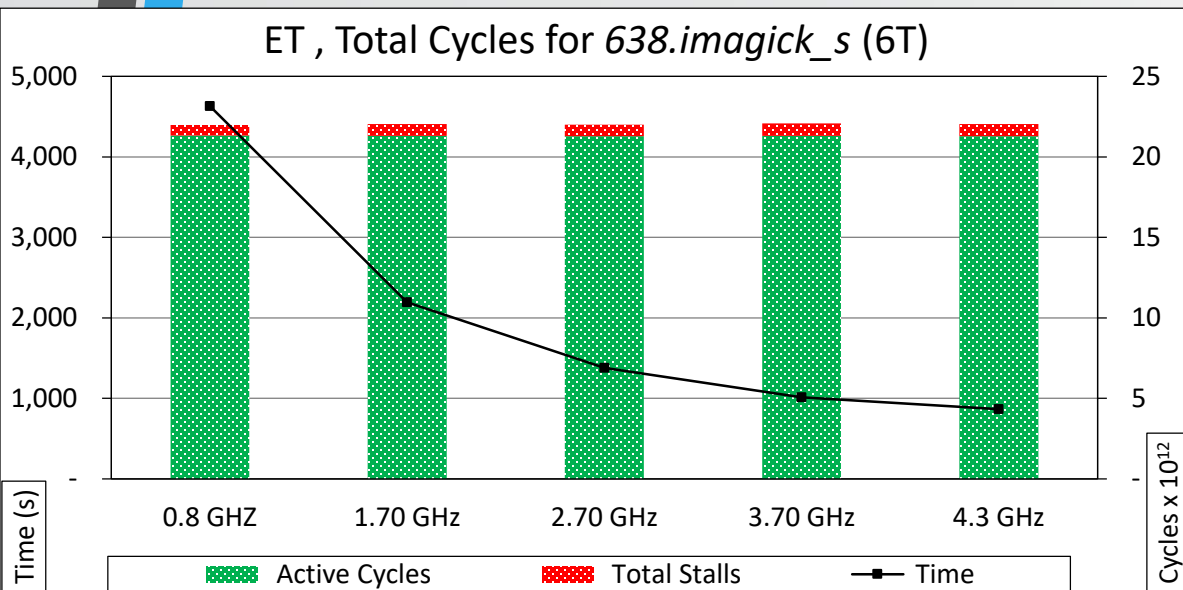
THE UNIVERSITY OF
ALABAMA IN HUNTSVILLE

- ## Performance

  - 4.30 GHz to 2.70 GHz

- ## Energy Efficiency

  - 1.70 GHz to 1.30 GHz

- ## PxEE

  - 2.70 GHz to 1.70 GHz



Optimal Frequency Selection (*649.fotonik3d_s*-6T)

31

# Motivation: Limitations of CPU Utilization Metric

**THE UNIVERSITY OF ALABAMA IN HUNTSVILLE**



ET , Total Cycles for *638.imagick_s* (6T)



| Busy (90%) | | Idle (10%) |
|---|---|---|
| Active Execution (30%) | Waiting (60 %) Stalled | Idle (10%) |
| Utilized Clock Cycles at Selected P-State | Wasted Clock Cycles at Selected P-State | |



ET, Cycles for *628.pop2_s* (6T)



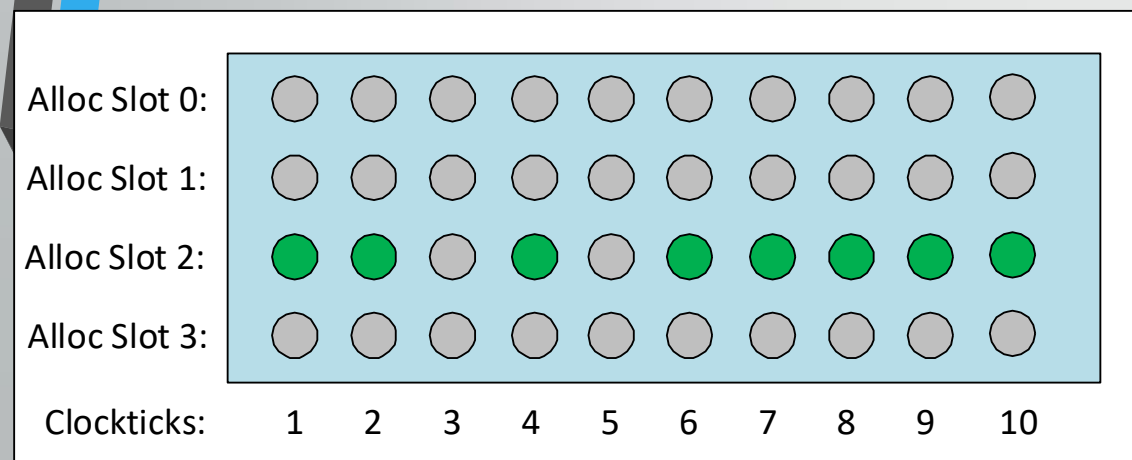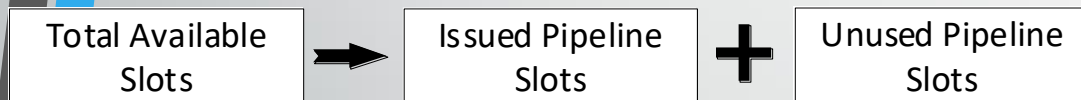ET, Cycles for *649.fotonik3d_s* (6T)

32

# Overview

- Introduction

- Background & Motivation

- Experimental Setup

- SPEC CPU2017 Characterization

- **Proposed DVFS Techniques**

- Results

- Future Work & Conclusions

# Proposed PMU-Event-Driven DVFS Techniques

1. Frequency Selection Based on Pipeline Slot Stalls: **FS-PS**

2. Frequency Selection Based on Total Cycle Stalls: **FS-TS**

3. Frequency Selection Based on Memory Related Cycle Stalls: **FS-MS**

4. Frequency Selection Based on LLC- MPKI: **FS-LLCM**

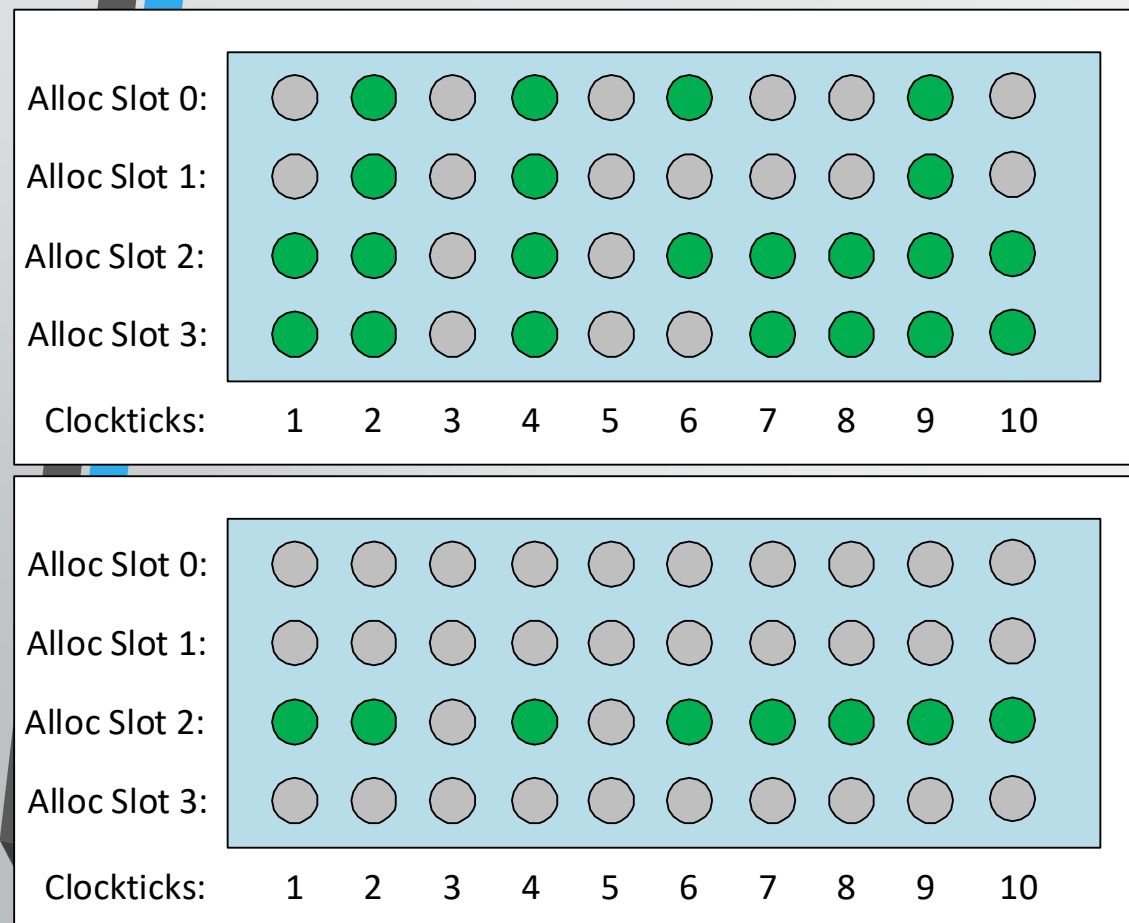# P-State Selection Based on Pipeline Slot Stalls

**Alloc Slot 0:**
**Alloc Slot 1:**
**Alloc Slot 2:**
**Alloc Slot 3:**

**Clockticks:** 1 2 3 4 5 6 7 8 9 10

| Total Available Slots | → | Issued Pipeline Slots | + | Unused Pipeline Slots |

**Alloc Slot 0:**
**Alloc Slot 1:**
**Alloc Slot 2:**
**Alloc Slot 3:**

**Clockticks:** 1 2 3 4 5 6 7 8 9 10

**Frequency Selection based on Pipeline Slot Stalls (FS-PS)**

- Example 1: 10 Clock cycles: 40 pipeline slots
  - 40 available pipeline slots (4-wide pipeline)
  - Micro-operations were issues for 22 slots
  - Pipeline Stall ratio: **0.45 $\Rightarrow$ P5-State**

- Example 2: 10 Clock cycles
  - 40 available pipeline slots (4-wide pipeline)
  - Micro-operations were issued for 8 slots
  - Pipeline Stall ratio: **0.8 (32/40) $\Rightarrow$ P8-State**

35

# P-State Selection Based on Cycle Stalls

THE UNIVERSITY OF
ALABAMA IN HUNTSVILLE



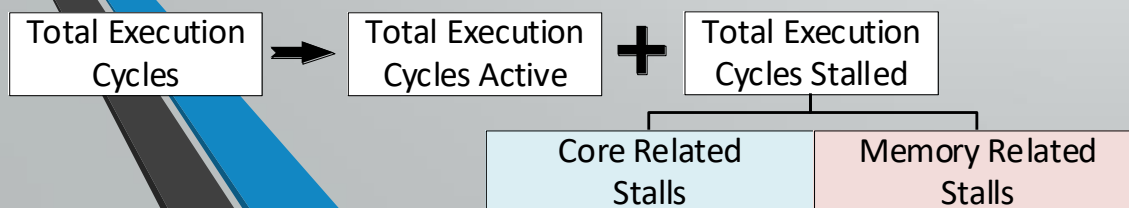- Example1 & 2: 10 Clock cycles

  - 10-clock cycles

  - Total Cycle Stalls: 2

  - Memory-related Cycle Stalls: 1

**Frequency Selection based on Total Cycle Stalls (FS-TS)**

- Total Cycle Stall ratio: 0.2 $\Rightarrow$ **P2-State**

**Frequency Selection based on Memory-Related Cycle Stalls (FS-MS)**

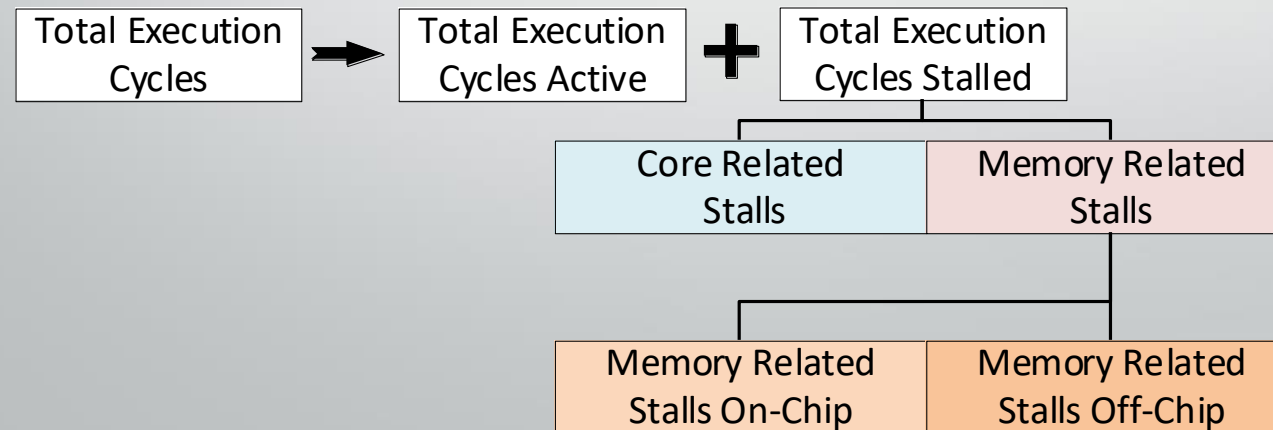- Memory-related Cycle Stall ratio: 0.1 $\Rightarrow$ **P1-State**

# P-State Selection Based on LLC Misses PKI

- Targeted DVFS for applications with significant Off-chip requests

- LLC miss denotes off-chip requests

- An LLC miss is followed by stalls associated by data fetch from main-memory

**Frequency Selection based on LLC Miss PKI (FS-LLCM)**

- Map LLC-MPKI to available P-states

- Practical Considerations

  - Range of MPKI observed on real life workload shows a range of 0-100 MPKI

```
Total Execution        Total Execution      +   Total Execution
Cycles           →     Cycles Active            Cycles Stalled
```

| Core Related Stalls | Memory Related Stalls |
|---|---|

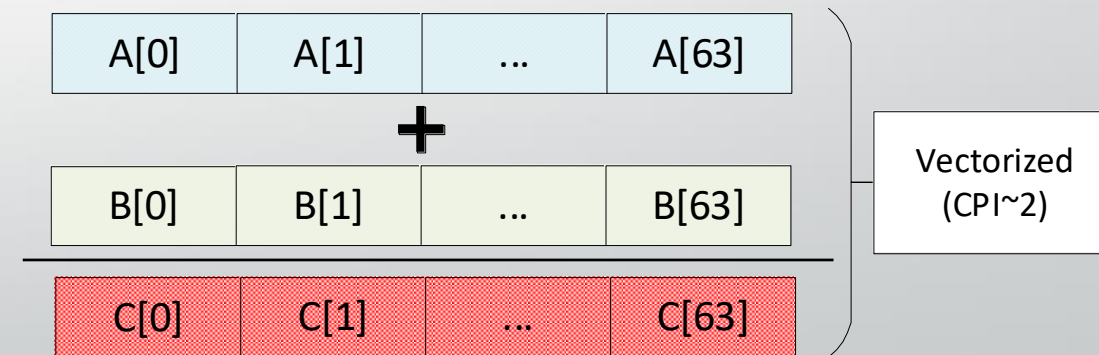| Memory Related Stalls On-Chip | Memory Related Stalls Off-Chip |
|---|---|

37

# FS-CPI: Implementation and Its Shortcomings
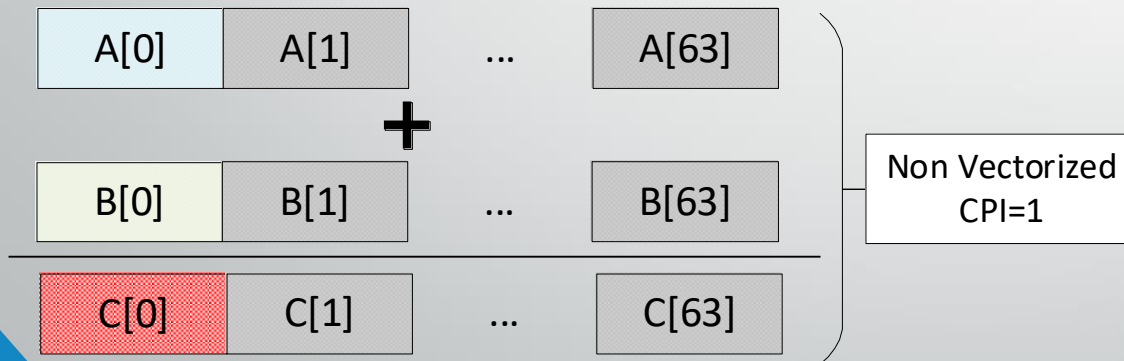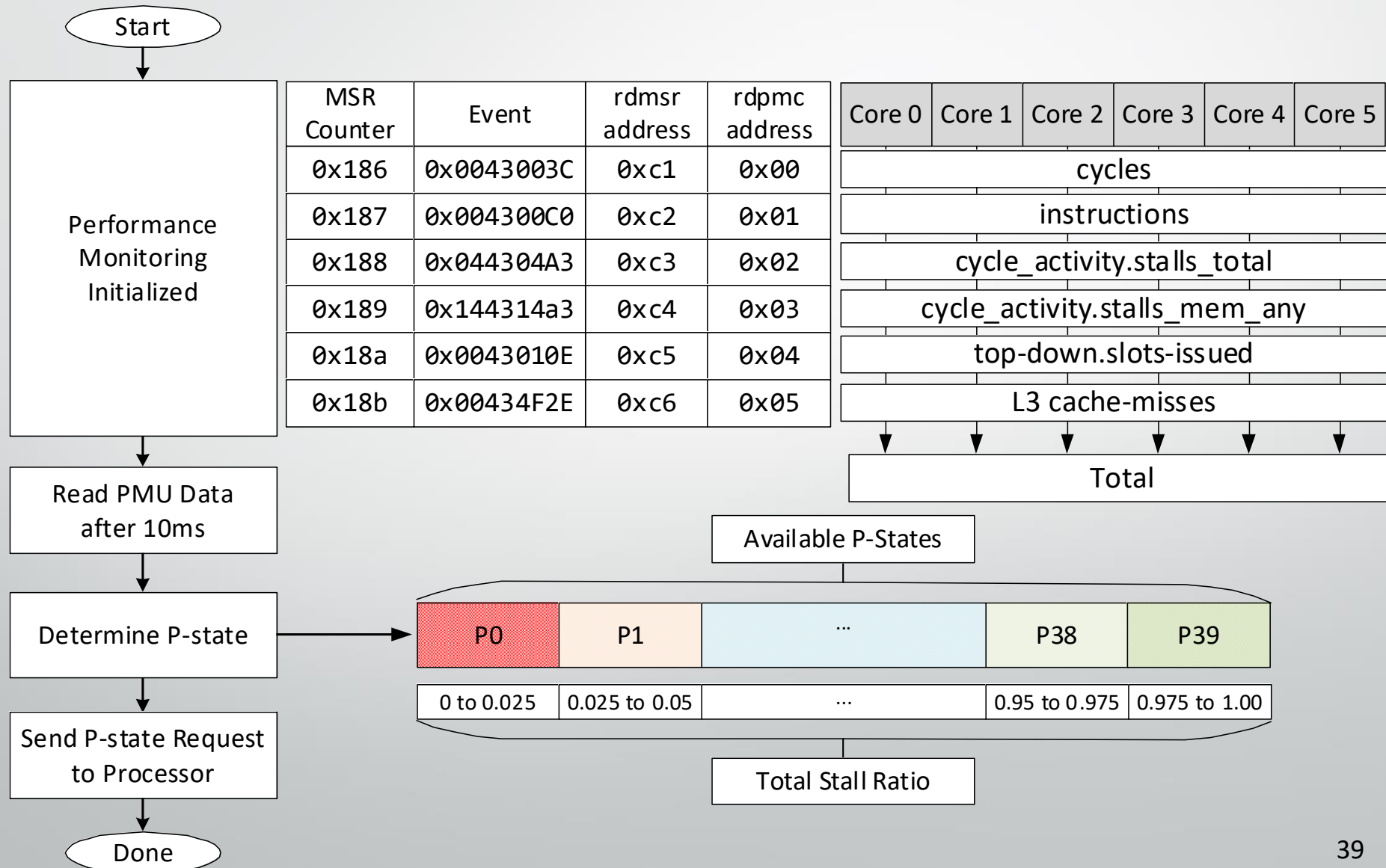
**Frequency Selection based on CPI (FS-CPI)**

- Map CPI to available P-states

- Mapping range selected based on real-life workload (0-6)

- Why CPI is not ideal?

- Metric too sensitive to operating frequency

- Vectorized Code:
  - Vectorized code may have higher CPI but can get more work done in fewer instructions.
  - Does not capture the impact of vectorization



38

# Implementation of Proposed Techniques

**Start**

Performance Monitoring Initialized

Read PMU Data after 10ms

Determine P-state

Send P-state Request to Processor

**Done**

| MSR Counter | Event | rdmsr address | rdpmc address |
|---|---|---|---|
| 0x186 | 0x0043003C | 0xc1 | 0x00 |
| 0x187 | 0x004300C0 | 0xc2 | 0x01 |
| 0x188 | 0x044304A3 | 0xc3 | 0x02 |
| 0x189 | 0x144314a3 | 0xc4 | 0x03 |
| 0x18a | 0x0043010E | 0xc5 | 0x04 |
| 0x18b | 0x00434F2E | 0xc6 | 0x05 |

| Core 0 | Core 1 | Core 2 | Core 3 | Core 4 | Core 5 |
|---|---|---|---|---|---|
| cycles | | | | | |
| instructions | | | | | |
| cycle_activity.stalls_total | | | | | |
| cycle_activity.stalls_mem_any | | | | | |
| top-down.slots-issued | | | | | |
| L3 cache-misses | | | | | |

Total

Available P-States

| P0 | P1 | ... | P38 | P39 |
|---|---|---|---|---|
| 0 to 0.025 | 0.025 to 0.05 | ... | 0.95 to 0.975 | 0.975 to 1.00 |

Total Stall Ratio

39

# Implementation Overhead & Invocation Rate

THE UNIVERSITY OF
ALABAMA IN HUNTSVILLE

- Performance monitoring period: 10 ms

- OpenMP to read PMU from all core counters

- Implementation overhead: +3 ms
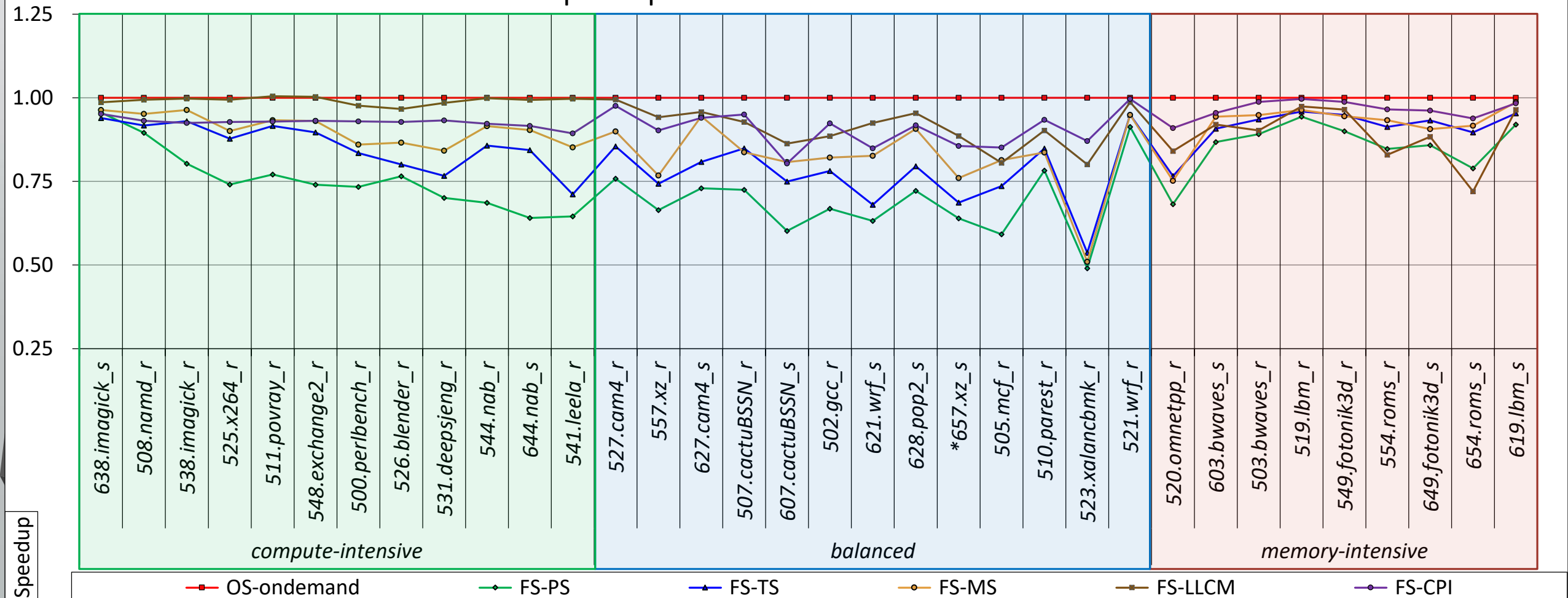
- Total duration: Idle CPU, running at 0.8 GHz: 13 ms

|  | PKG Power (W) | PP0-Power (W) |
|---|---|---|
| Ref (OS-ondemand) | 4.50 | 1.00 |
| FS-50ms (20 Hz) | 6.05 | 2.55 |
| FS-100ms (10 Hz) | 5.25 | 1.75 |

40

# Overview

- Introduction

- Background & Motivation

- Experimental Setup

- SPEC CPU2017 Characterization

- Proposed DVFS Techniques

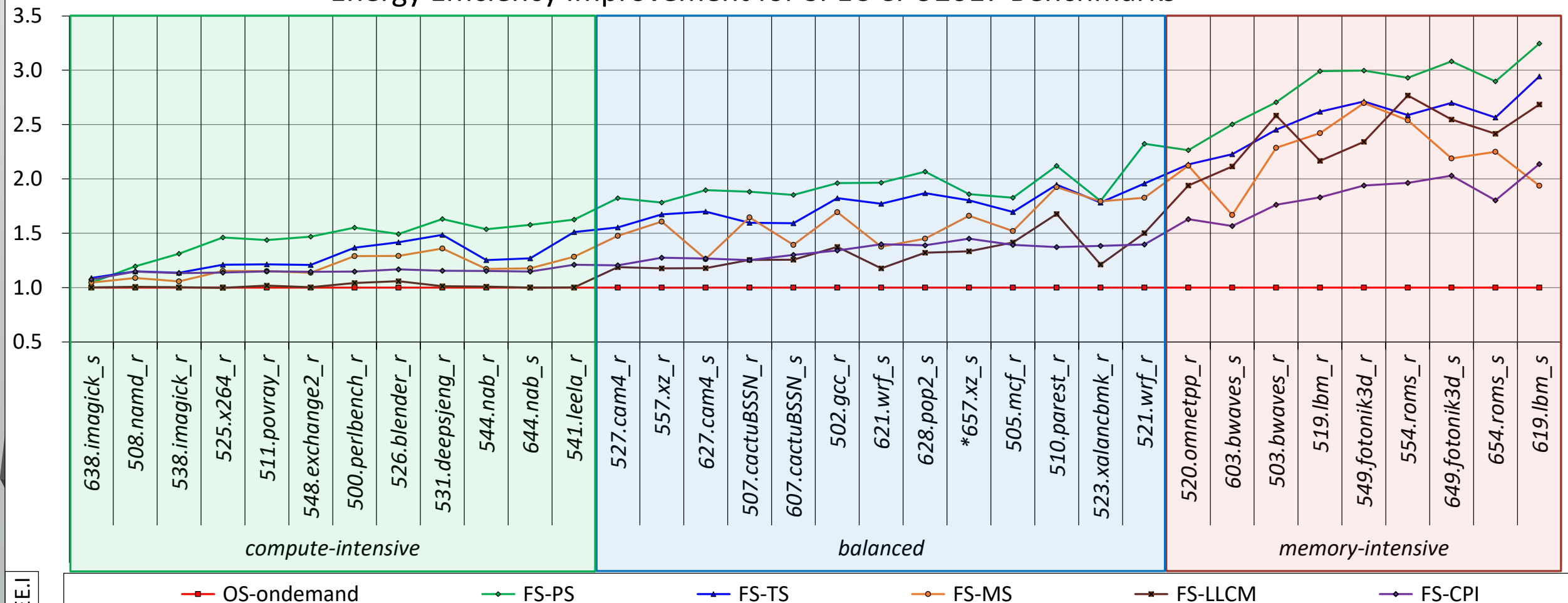- Results

- Conclusions

# SPEC CPU2017 Performance



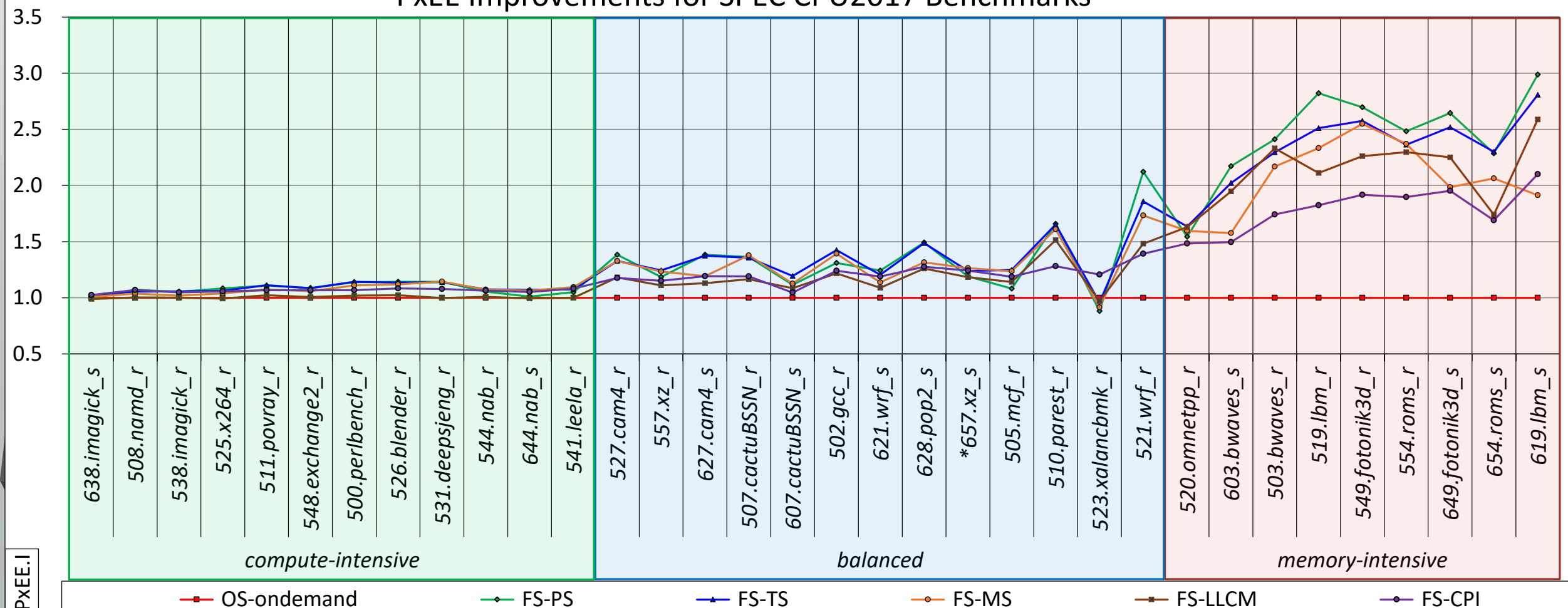Performance Speedup for SPEC CPU2017 Benchmarks

# SPEC CPU2017 Energy Efficiency



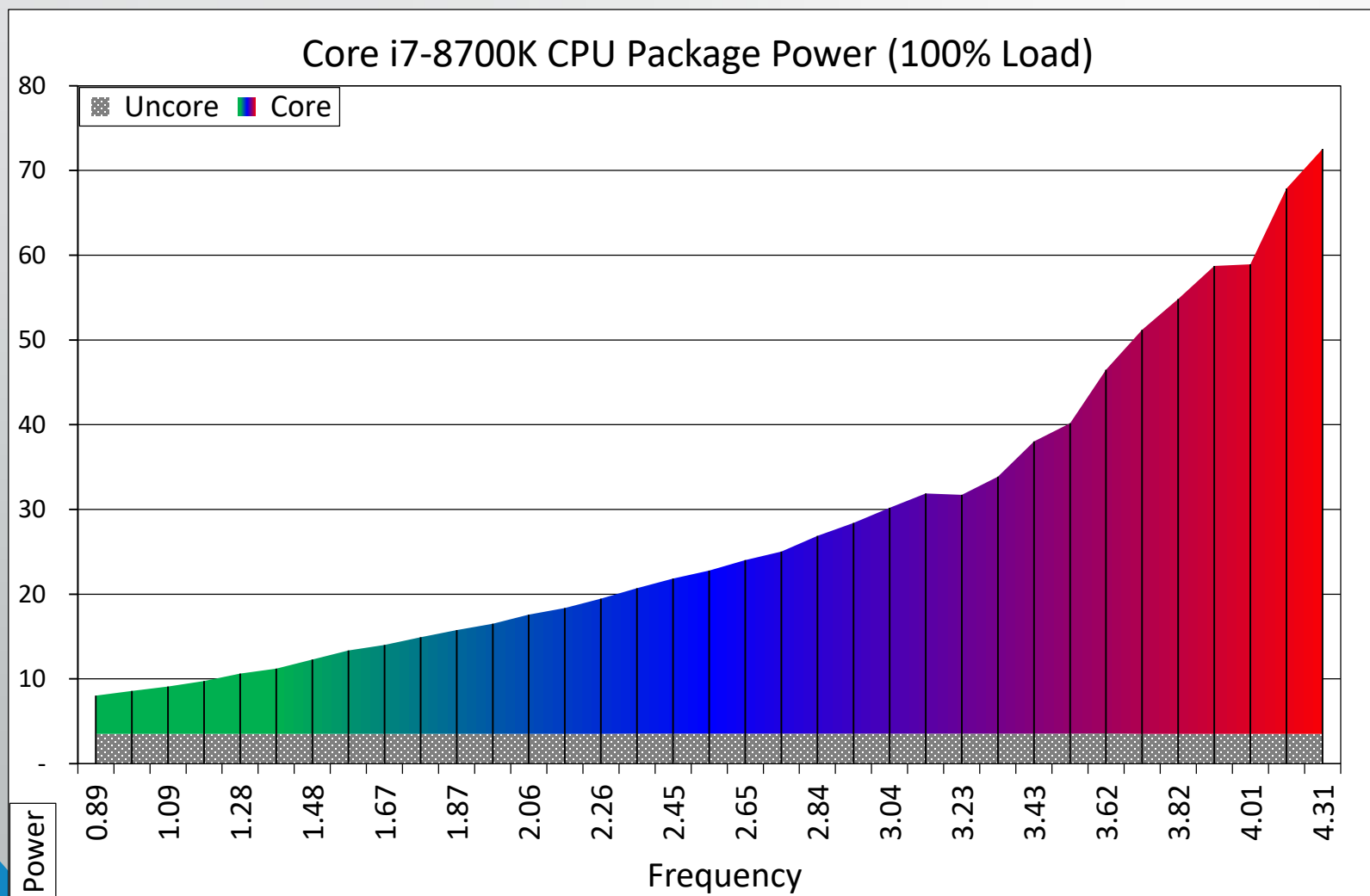Energy Efficiency Improvement for SPEC CPU2017 Benchmarks

# PxEE Summary for SPEC CPU2017 Benchmarks



PxEE Improvements for SPEC CPU2017 Benchmarks

# Future Work

Core i7-8700K CPU Package Power (100% Load)

- Mapping Techniques

- Hybrid Models

- Alternate Workloads

# Conclusions

- DVFS is one of the most important tools in regulating processor power consumption

- OS implementation relies on CPU utilization $\Rightarrow$ favors performance

- Existing FS-CPI can be unstable (vectorized code), limited benefits

- This dissertation proposes alternate techniques that significantly improve energy-efficiency of memory intensive applications

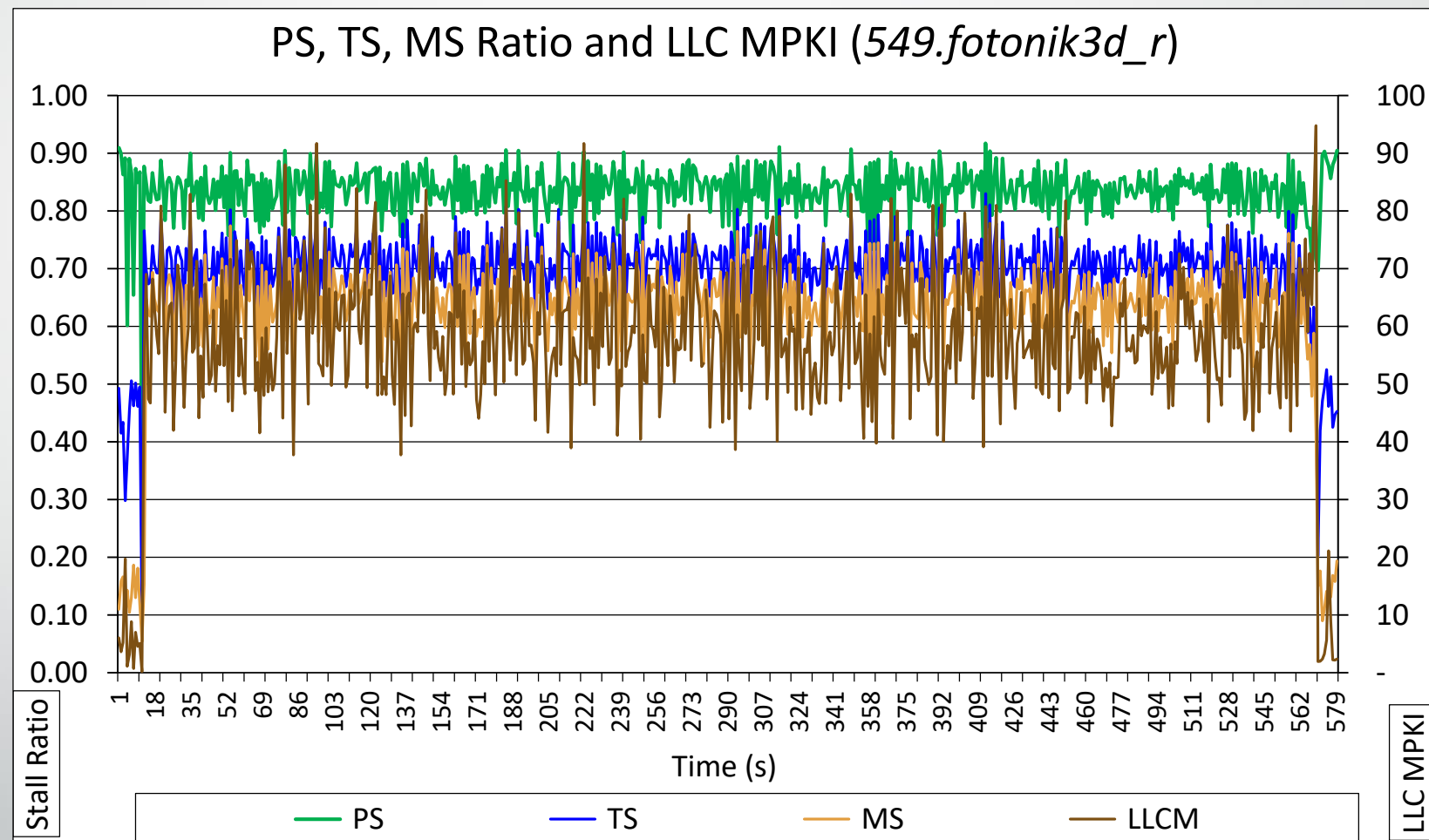| Technique | PxEE Gains w.r.t state-of-the-art |
|-----------|-----------------------------------|
| FS-PS     | 141% |
| FS-TS     | 132% |
| FS-MS     | 104% |
| FS-LLCM   | 110% |

46

# Outline

- Introduction

- Background

- Experimental Setup

- SPEC CPU2017 Characterization

- Proposed DVFS Techniques

- Results

- Conclusions

- **Backup**

# Runtime Metrics Measurements for *549.fotonik3d_r*

| | FS-PS | FS-TS | FS-MS | FS-LLCM |
|---|---|---|---|---|
| P.S | 0.90 | 0.95 | 0.94 | 0.96 |
| EE.I | 3.00 | 2.71 | 2.70 | 2.34 |
| PxEE | 2.70 | 2.58 | 2.55 | 2.26 |



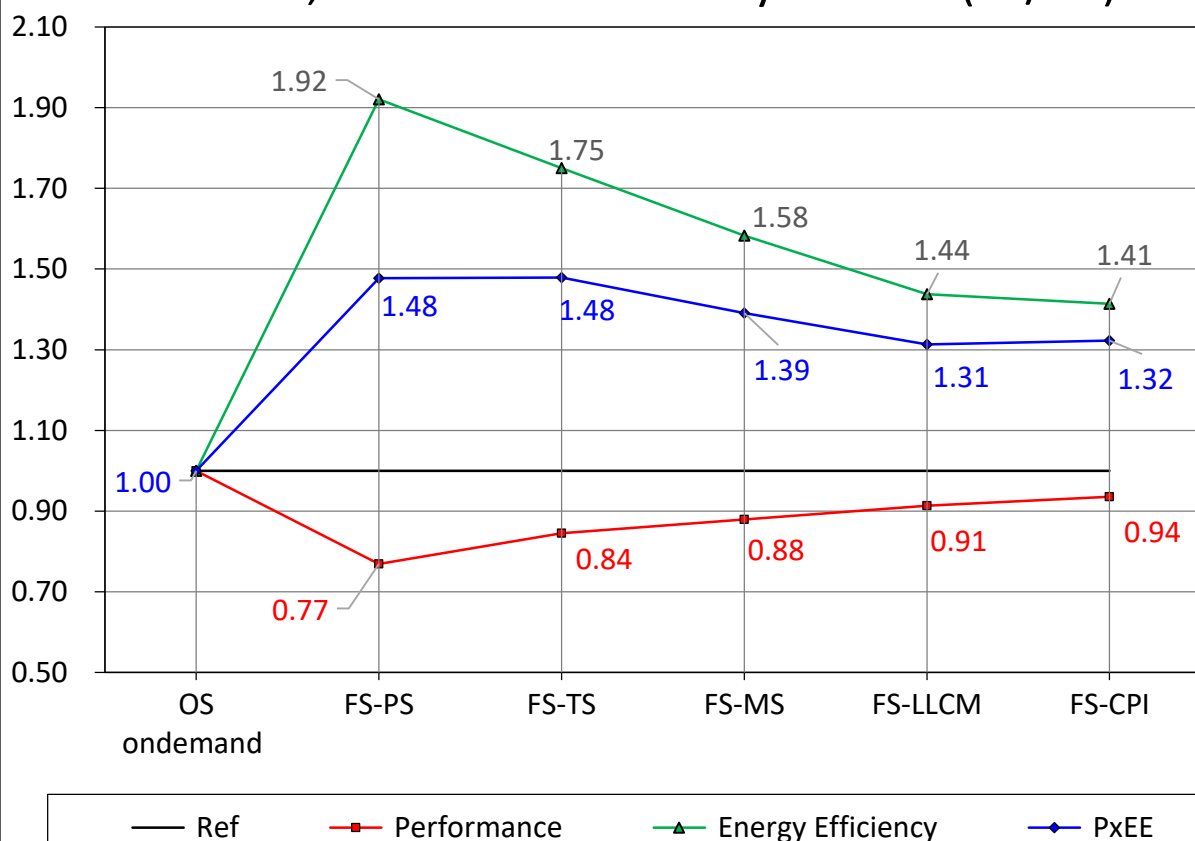PS, TS, MS Ratio and LLC MPKI (*549.fotonik3d_r*)

# Runtime Measurements of the Total Stall Ratio and the Average CPI for *654.rom_s*

- ## FS-TS
  - ### P.S: 0.90
  - ### EE.I: 2.57
  - ### PxEE.I: 2.30

- ## FS-CPI
  - ### P.S: 0.94
  - ### EE.I: 1.80
  - ### PxEE.I: 1.69



Total Stall Ratio & CPI for *654.roms_s* (FS-CPI)

# SPEC CPU2017 Summary



P.S, EE.I and PxEE.I- Fully Loaded (6T/6C) — chart with data points at OS ondemand, FS-PS, FS-TS, FS-MS, FS-LLCM, FS-CPI. Energy Efficiency: 1.92, 1.75, 1.58, 1.44, 1.41. PxEE: 1.48, 1.48, 1.39, 1.31, 1.32. Performance: 1.00, 0.77, 0.84, 0.88, 0.91, 0.94.

P.S, EE.I and PxEE.I- Partially Loaded (4T/4C) — Energy Efficiency: 1.78, 1.61, 1.52, 1.39, 1.35. PxEE: 1.36, 1.37, 1.33, 1.26, 1.27. Performance: 1.00, 0.76, 0.85, 0.87, 0.91, 0.94.

Legend: Ref, Performance, Energy Efficiency, PxEE

# SPECpower: Performance



SPECpower Performance

Legend: FS-PS, FS-TS, FS-MS, FS-LLCM, FS-CPI, OS-ondemand

Y-axis: Thousands / Performance
X-axis: CPU Load (%) — 100%, 90%, 80%, 70%, 60%, 50%, 40%, 30%, 20%, 10%, Active Idle

# SPECpower: Power



SPECpower Power Measurements

# SPECpower: Performance/Watt



SPECpower Performance Per Watt

Legend: FS-PS, FS-TS, FS-MS, FS-LLCM, FS-CPI, OS-ondemand

X-axis: CPU Load (%) — 100%, 90%, 80%, 70%, 60%, 50%, 40%, 30%, 20%, 10%, Active Idle

Y-axis: Perf/Watt