# µDPM: Dynamic Power Management for the Microsecond Era

Chih-Hsun Chou

cchou001@cs.ucr.edu

Laxmi N. Bhuyan bhuyan@cs.ucr.edu Daniel Wong danwong@ucr.edu







#### Computer systems efficiently support . . .







# The Killer Microseconds



"System designers can no longer ignore efficient support for microsecond-scale I/O ... <u>Novel</u> <u>microsecond-optimized system</u> <u>stacks are needed</u>"



[1] Luiz Barroso, Mike Marty, David Patterson, and Parthasarathy Ranganathan. Attack of the killer microseconds. Commun. ACM 60, 4 (March 2017), 48-54.







#### Computer systems cannot efficiently support Microsecond-scale service time

Traditional Monolithic Services















#### Computer systems cannot efficiently support Microsecond-scale service time



### **Microservice Example**











# Implications of Killer Microsecond service time on Dynamic Power Management?









8 UCRIVERSITY OF CALIFORNIA

# Opportunity for DPM – Latency Slack

- Slow down request processing (DVFS)
- Delay request processing (Sleep)







- > DVFS (Rubik [MICRO'15], Pegasus [ISCA'14])
  - Rubik adjusts f per request



- > DVFS + Sleep
  - > <u>SleepScale</u> [ISCA'14] finds optimal frequency & C-state depth for 60s epochs





10 UCRIVERSITY OF CALIFORNIA







10 UCREVERSITY OF CALIFORNIA







10 UCREVERSITY OF CALIFORNIA







10 UCREVERSITY OF CALIFORNIA







10 UCREVERSITY OF CALIFORNIA







10 IICRIVERSITY OF CALIFORNIA







10 IIC RIVERSITY OF CA







10 IIC RUNIVERSITY OF CALIFORNIA







10 IIC RUNIVERSITY OF CALIFORNIA







#### DPM ineffective w/ microsecond service time







## DPM ineffective w/ microsecond service time

>250µs: DVFS effective at slowing down request processing







### DPM ineffective w/ microsecond service time

- >250µs: DVFS effective at slowing down request processing
- > <250µs: DPM becomes ineffective</p>





#### Systems Optimization + Computer Architecture Lab

#### DPM ineffective w/ microsecond service time

HPCA 2019

- >250µs: DVFS effective at slowing down request processing
- > <250µs: DPM becomes ineffective</p>
- Surprisingly, sleep-based policies outperform DVFS-based policies







#### Fragmented idle periods $\rightarrow$ Lost Opportunities









12 UCREVERSITY OF CALIFORNIA

#### Fragmented idle periods $\rightarrow$ Lost Opportunities



Shorter Service Time





UNIVERSITY OF CALIFORNIA

#### Fragmented idle periods $\rightarrow$ Lost Opportunities







#### Fragmented idle periods $\rightarrow$ Lost Opportunities

 Short service times fragment idle periods



13





#### Fragmented idle periods $\rightarrow$ Lost Opportunities

- Short service times fragment idle periods
- Sleep states / request delaying can consolidate idle periods





#### Significant transition overheads and idle power



Idleness and transition overhead still account for up to ~25% of energy

















\* SPECjbb timing

15 UCRIVERSITY OF CALIFORNIA







\* SPECjbb timing

15 UCRIVERSITY OF CALIFORNIA



Systems Optimization + Computer Architecture Lab



\* SPECjbb timing

15 UCRIVERSITY OF CALIFORNIA



Systems Optimization + Computer Architecture Lab











\* SPECjbb timing

15 UCRIVERSITY OF CALIFORNIA





# Key Insight

Careful coordination of DVFS, Sleep state, and request delaying is the key to effective DPM with microsecond service times





# μDPM



- > Aggressively Deep Sleep
- Delay and slow down request processing to finish just-in-time, even under microsecond request service times
- Carefully coordinating DVFS, Sleep, and request delaying





UNIVERSITY OF CALIFORNIA

18 **|||** 

#### Can latency-critical workloads utilize deep sleep states?







#### Aggressive deep sleep and request delaying

> Wakeup after residency time



#### > Wakeup before residency time if needed to meet tail latency







# Estimating Service time and Latency

- Estimate tail service time
  - Statistical performance model<sup>[2]</sup>
  - Online periodic resampling (100ms)



S.

[2] Kasture, Harshad, Davide B. Bartolini, Nathan Beckmann, and Daniel Sanchez. "Rubik: Fast analytical power management for latency-critical systems." MICRO 2015







# Estimating Service time and Latency



[2] Kasture, Harshad, Davide B. Bartolini, Nathan Beckmann, and Daniel Sanchez. "Rubik: Fast analytical power management for latency-critical systems." MICRO 2015.







Detecting critical request arrival

If inter-arrival time between 2 consecutive requests are shorter than the tail service time









21 **UCREVERSITY OF CALIFORNIA** 



- Detecting critical request arrival
  - If inter-arrival time between 2 consecutive requests are shorter than the tail service time





## Detecting critical request arrival



#### Detecting critical request arrival

If inter-arrival time between 2 consecutive requests are shorter than the tail service time







## Detecting critical request arrival



#### Detecting critical request arrival

If inter-arrival time between 2 consecutive requests are shorter than the tail service time

























- Reset to lowest frequency on wake up
- > Only increase frequency on reconfiguration







#### Calculate criticality score

$$criticality = \frac{S^{tail}/f}{t_{R_i} - t_{R_{i-1}}}$$

#### Send to core that is least critical

$$\begin{split} E\left(W,f\right) &= \left(W - T_{sleep}\right)P_{idle} + \left(T_{sleep} + T_{wake}\right)\\ P_{max} &+ T_{dvfs}P_{dvfs} + \left(\frac{S_i}{f}\right)P_f \end{split}$$

Minimize state transitions

#### Algorithm 1: Criticality-aware scheduling

1: $non\_critical\_cores = \phi, non\_critical\_sleep\_cores = \phi$			
2: for each core do			
3: compute <i>core</i> <sub>i</sub> 's criticality			
4: <b>if</b> criticality $\leq 1$ <b>then</b>			
5: $non\_critical\_cores \leftarrow non\_critical\_cores \cup core_i$			
6: <b>if</b> <i>core</i> <sub><i>i</i></sub> is sleeping <b>then</b>			
7: $non\_critical\_sleep\_cores \leftarrow$			
$non\_critical\_sleep\_cores \cup core_i$			
8: end if			
9: end if			
10: end for			
11: if non_critical_cores $\neq \phi$ then			
12: if non_critical_sleep_cores $\neq \phi$ then			
13: <b>return</b> min (extra energy) in non_critical_sleep_cores			
14: <b>else</b>			
15: <b>return</b> min (extra energy) in non_critical_cores			
16: end if			
17: else			
18: return min (extra energy) in all cores			
19: end if			

#### See paper for details!

23 UCREVERSITY OF CALIFORNIA



## Evaluation



24

- In-House Simulator (similar to BigHouse)
- Empirical Power Model
  - > 10µs DVFS transition,
    - 89µs sleep transition time (double to account for cache flushing)
  - > Add 25µs to first request service time after idle period for cold miss penalty
- Baseline Linux menu idle governor and intel\_pstate driver
- > Workloads

lable II:	Workload	characteristics.	
Name	Avg. Service Time	Tail Service Time	Target Tail Latency
Memcached [43, 44]	30µs	33µs	150µs
SPECjbb [8, 11]	65µs	78µs	800µs
Masstree [8, 11, 46]	246µs	250µs	1100µs
Xapian [8, 11]	431µs	1200µs	2100µs





NIVERSITY OF\_CALIFORNIA

# Energy Savings





Load





25

Systems Optimization + Computer Architecture Lab



Systems Optimization + Computer Architecture Lab

25

#### State transition overhead reduction



0.6 0.68 0.76 0.84 0.92 1

Normalized Energy







# **Under Varying Load**



Systems Optimization + Computer Architecture Lab

27 UCRIVERSITY OF CALIFORNIA

## Sensitivity to target tail latency





28 UCRIVERSITY OF CALIFORNIA

#### Sensitivity to Transition time









## Conclusion



- Microsecond service times present challenges for Dynamic Power Management
- Careful coordination of DVFS, Sleep and Request delaying can achieve savings with µs service times
- >µDPM is able to save ~2x energy compared to state-of-the-art techniques



# Thank you!

# µDPM: Dynamic Power Management for the Microsecond Era

Chih-Hsun Chou cchou001@cs.ucr.edu

Laxmi N. Bhuyan bhuyan@cs.ucr.edu Daniel Wong danwong@ucr.edu



