

Altocumulus: Scalable Scheduling for Nanosecond-Scale Remote Procedure Calls

Jiechen Zhao, Iris Uwizeyimana, Karthik Ganesan, Mark C. Jeffrey, Natalie Enright Jerger

55th IEEE/ACM International Symposium on Microarchitecture (MICRO), Chicago, Oct. 2022 https://www.eecg.utoronto.ca/~mcj/papers/2022.altocumulus.micro.pdf

Executive Summary

- RPC processing time has decreased down to microsecond scale
 - Therefore, bottleneck has shifted to scheduling
- Prior RPC schedulers react to imbalance ineffectively, causing
 - High tail latency at medium load
 - Low CPU utilization while satisfying strict tail latency deadline
- Altocumulus: Proactively rebalances RPC loads using HW/SW co-design
- Achieves up to 24x throughput improvement under microsecond scale tail latency deadline over state-of-the art

RPCs become ubiquitous in datacenters

- Productivity improvement
 - As a common communication API
 - Harmonizing distributed services developed by different programming languages
- Vast deployment
 - As the backbone of many latency-critical applications (KVS, SMR, RDMA, etc.)
- Massive amount
 - Each request involving 10s 1,000s RPCs







- Many datacenter applications today are interactive
 - Strict performance requirements as <u>SLO</u>: Low tail latency at high load

- Many datacenter applications today are interactive
 Strict performance requirements as <u>SLO</u>: Low tail latency at high load
- "Tail-at-scale" problem [Dean et al., CACM'13]
 - User-perceived latency determined by slowest back-end server node



- Many datacenter applications today are interactive
 Strict performance requirements as <u>SLO</u>: Low tail latency at high load
- "Tail-at-scale" problem [Dean et al., CACM'13]
 - User-perceived latency determined by slowest back-end server node



Why on-CPU RPC handling counts ?

Many datacenter applications today are interactive
 Strict performance requirements as <u>SLO</u>: Low tail latency at high load

• "Tail-at-scale" problem [Dean et al., CACM'13]



Why on-CPU RPC handling counts ?

Fig. from [Lazarev et al., ASPLOS'21]

Many datacenter applications today are interactive
 Strict performance requirements as <u>SLO</u>: Low tail latency at high load

• "Tail-at-scale" problem [Dean et al., CACM'13]



Why on-CPU RPC handling counts ?

Taking up to ~90% of end-to-end time

- Many datacenter applications today are interactive
 Strict performance requirements as <u>SLO</u>: Low tail latency at high load
- "Tail-at-scale" problem [Dean et al., CACM'13]
- "Killer microsecond" problem [Barroso et al., CACM'17]
 - Existing systems not able to handle microsecond-scale RPCs efficiently



- Many datacenter applications today are interactive
 Strict performance requirements as SLO: Low tail latency at high load
- "Tail-at-scale" problem [Dean et al., CACM'13]
- "Killer microsecond" problem [Barroso et al., CACM'17]
- Microservice queuing propagation
 - Each server receiving/sending RPCs whose delay varies



What's even worse ...

- Many datacenter applications today are interactive
 Strict performance requirements as <u>SLO</u>: Low tail latency at high load
- "Tail-at-scale" problem [Dean et al., CACM'13]
- "Killer microsecond" problem [Barroso et al., CACM'17]
- Microservice queuing propagation
 - Each server receiving/sending RPCs whose delay varies



What's even worse ...

- Many datacenter applications today are interactive
 Strict performance requirements as *SLO*: Low tail latency at high load
- "Tail-at-scale" problem [Dean et al., CACM'13]
- "Killer microsecond" problem [Barroso et al., CACM'17]
- Microservice queuing propagation
 - Each server receiving/sending RPCs whose delay varies



On-CPU RPC handling SLO ~= 10s microsecond

RPC Lifetime on a CPU

Processing

Application RPC layer

Transport

DC Network





- Processing
 - Application
 - monolithic \longrightarrow microservices (10s ns 10s µs)



• Software \longrightarrow hardware (10s ns)

Processing

Application

RPC layer







- monolithic \longrightarrow microservices (10s ns 10s μ s)
- RPC layer

Application

Processing

- Software → hardware (10s ns)
- Transport
 - TCP/IP \longrightarrow optimization (100s ns)





- Processing
 - Application
 - monolithic \longrightarrow microservices (10s ns 10s μ s)
 - RPC layer
 - Software → hardware (10s ns)
 - Transport
 - TCP/IP \longrightarrow optimization (100s ns)
- Scheduling ?



RPC Lifetime on a CPU

Processing

Nanosecond scale

- Application
 - monolithic \longrightarrow microservices (10s ns 10s μ s)
- RPC layer
 - Software → hardware (10s ns)
- Transport
 - TCP/IP \longrightarrow optimization (100s ns)
- Scheduling ?



RPC Lifetime on a CPU

Processing

- Now, Nanosecond scale
- Application
 - monolithic \longrightarrow microservices (10s ns 10s μ s)
- RPC layer
 - Software → hardware (10s ns)
- Transport
 - TCP/IP → optimization (100s ns)

Scheduling ?

Now, impact on SLO ?













• With prior work reducing the processing latency of RPCs ...

Performance bottleneck shifted from RPC processing to RPC scheduling





Scheduling (Work Stealing)



Key Questions



Minimize SLO violations

How many RPCs a server can sustain w/o violating SLO ? Maximize RPC throughput@SLO



Design Goals

Reconcile the trade-off among





High volume of requests (400+ Gbps)

Strict SLO (99th% latency@µs)



High CPU efficiency (16-256 cores)

Technical design goal: <50ns scheduling overhead --- near-ideal !

Abstracting RPC Scheduling Sub-System

Policy		When, What, Where to schedule
Runtime		Schedule scalably and adaptively
Mechanism	►	How to schedule

Decentralized First Come First Serve (D-FCFS)



Kernel-native d-FCFS runtime Shared mem.

IX [Belay et al., OSDI'14]

Decentralized First Come First Serve (D-FCFS)



Kernel-native d-FCFS runtime Shared mem.

IX [Belay et al., OSDI'14]

Decentralized First Come First Serve (D-FCFS)



Kernel-native d-FCFS runtime Shared mem.

IX [Belay et al., OSDI'14]

Decentralized First Come First Serve (D-FCFS)





[Belay et al., OSDI'14]

Decentralized First Come First Serve (D-FCFS)





[Belay et al., OSDI'14]
Decentralized First Come First Serve (D-FCFS)





[Belay et al., OSDI'14]

Decentralized First Come First Serve (D-FCFS)



Kernel-native d-FCFS runtime Shared mem.

IX [Belay et al., OSDI'14]

Decentralized First Come First Serve (D-FCFS)



Kernel-native d-FCFS runtime Shared mem.

IX [Belay et al., OSDI'14]

Load Imbalance

Decentralized First Come First Serve (D-FCFS)



Load Imbalance

Kernel-native d-FCFS runtime Shared mem.

IX [Belay et al., OSDI'14]

- SLO: Not Met
- Throughput: Low
- Core Utilization: Bad

D-FCFS + work stealing



Work stealing d-FCFS runtime Shared mem.

ZygOS [Prekas et al., SOSP'17]

D-FCFS + work stealing



Work stealing d-FCFS runtime

Shared mem.

ZygOS [Prekas et al., SOSP'17]

D-FCFS + work stealing



Work stealing d-FCFS runtime Shared mem. ZygOS

[Prekas et al., SOSP'17]

- Core Utilization: Good
- SLO: Not Met

D-FCFS + work stealing



Work stealing d-FCFS runtime Shared mem. ZygOS

ZygOS [Prekas et al., SOSP'17]

- Core Utilization: Good
- SLO: Not Met

Slow steal

• Throughput: Medium

44

Centralized First Come First Serve (C-FCFS)



Preemption **c-FCFS runtime** Shared mem.

Centralized First Come First Serve (C-FCFS)



Preemption **c-FCFS runtime** Shared mem.

- SLO: Met
- Throughput: High
- Core Utilization: Good

Centralized First Come First Serve (C-FCFS)



Preemption **c-FCFS runtime** Shared mem.

- SLO: Met
- Throughput: High
- Core Utilization: Good

Centralized First Come First Serve (C-FCFS)





Scheduling core - Bottleneck

- SLO: Not Met
- Throughput: Low
- Core Utilization: Bad

Centralized First Come First Serve (C-FCFS)





- SLO: Not Met
- Throughput: Low
- Core Utilization: Bad

Centralized First Come First Serve (C-FCFS)



JBSQ algorithm c-FCFS runtime **Cache coherence**

Nebula [Sutherland et al., ISCA'20]

JBSQ algorithm

c-FCFS runtime

Direct register

Centralized First Come First Serve (C-FCFS)



- SLO: Met
- Throughput: High
- Core Utilization: Good



Nebula [Sutherland et al., ISCA'20]

> JBSQ algorithm c-FCFS runtime **Direct register**

Centralized First Come First Serve (C-FCFS)



Coherence inscalability

Fixed scheduling algorithm

JBSQ algorithm c-FCFS runtime **Cache coherence**

Nebula [Sutherland et al., ISCA'20]

JBSQ algorithm

c-FCFS runtime

Direct register

Centralized First Come First Serve (C-FCFS)



- SLO: Not Met
- Throughput: Low
- Core Utilization: Bad



Nebula [Sutherland et al., ISCA'20]



nanoPU [Ibanez et al., OSDI'21]

Coherence inscalability

Fixed scheduling algorithm

Centralized First Come First Serve (C-FCFS)



- SLO: Not Met
- Throughput: Low
- Core Utilization: Bad



Nebula [Sutherland et al., ISCA'20]



Coherence inscalability

Fixed scheduling algorithm

Hard to scale



Abstraction



Scheduling Sub-system Abstraction



Scheduling Sub-system Abstraction



Abstraction

Altocumulus: High-Level Perspective



IX [Belay et al., OSDI'14] Work stealing

d-FCFS runtime

Shared mem.

ZygOS [Prekas et al., SOSP'17]

Preemption
c-FCFS runtime
Shared mem.

Shinjuku [Kaffes et al., NSDI'19]

JBSQ algorithm
c-FCFS runtime
Cache coherence

Nebula [Sutherland et al., ISCA'20] JBSQ algorithm c-FCFS runtime Direct register

Altocumulus: High-Level Perspective



Nebula [Sutherland et al., ISCA'20]





Altocumulus: High-Level Perspective



Nebula [Sutherland et al., ISCA'20]



[lbanez et al., OSDI'21]



<u>Altocumulus</u>

Altocumulus Architecture



- Multi-tiered global D-FCFS local C-FCFS
 - Each group: 1 queue + 1 manager (purple) + several workers (green)

Altocumulus Architecture



- Multi-tiered global D-FCFS local C-FCFS
 - Each group: 1 queue + 1 manager (purple) + several workers (green)
- Proactive & hardware-assisted work

stealing across manager cores

Altocumulus Architecture



- Multi-tiered global D-FCFS local C-FCFS
 - Each group: 1 queue + 1 manager (purple) + several workers (green)
- Proactive & hardware-assisted work

stealing across manager cores

Compatible with commodity multi-queue
NIC with RSS support







- 1. Periodically synchronize system states across manager cores
- 2. Pick appropriate threshold trained offline
- 3. Poll queues to check if threshold is met
- 4. Decide how many RPCs to migrate to which queue(s)

- 1. Periodically synchronize system states
- 2. Pick appropriate threshold trained offline
- 3. Poll queues to check if threshold is met



Heuristics: multi-queue load pattern classification

4. Decide how many RPCs to migrate to which queue(s)

- 1. Periodically synchronize system states
- 2. Pick appropriate threshold trained offline
- 3. Poll queues to check if threshold is met



Heuristics: multi-queue load pattern classification

- 4. Decide how many RPCs to migrate to which queue(s)
- 5. Pass decisions to h/w primitives and trigger migration h/w messages
- 6. Repeat Step 1 (Period as short as 50ns due to messaging offloaded to h/w)

- Periodically synchronize system states 1.
- leue Depth Pick appropriate threshold trained offline 2. How we train SLO violation prediction model?
- Poll queues to check if threshold is met 3. Heuristics: multi-queue load pattern classification
- Decide how many RPCs to migrate to which queue(s) 4.
- 5. Pass decisions to h/w primitives and trigger migration h/w messages How we offload messaging to h/w ? Repeat Step 1 (Period as short as 50ns due to messaging offloaded to h/w) 6.

Hardware Primitives: Direct Register Messaging

- Register-to-register migration messaging, bypassing cache coherence protocol
 - Inspired by [Sanchez et al., ASPLOS'10] & [Ibanez et al., OSDI'21] w/ several opt.
Hardware Primitives: Direct Register Messaging

- Register-to-register migration messaging, bypassing cache coherence protocol
 Inspired by [Sanchez et al., ASPLOS'10] & [Ibanez et al., OSDI'21] w/ several opt.
- Remain RPC message payload in LLC/memory -- Only move RPC descriptor
 - Reduce latency overhead and traffic per migration (up to 140x)

Hardware Primitives: Direct Register Messaging

- Register-to-register migration messaging, bypassing cache coherence protocol
 Inspired by [Sanchez et al., ASPLOS'10] & [Ibanez et al., OSDI'21] w/ several opt.
- Remain RPC message payload in LLC/memory -- Only move RPC descriptor
 - Reduce latency overhead and traffic per migration (up to 140x)
- Each request being migrated at most once
 - Avoid livelock and unnecessary scheduling traffic

Hardware Primitives: Direct Register Messaging

- Register-to-register migration messaging, bypassing cache coherence protocol
 Inspired by [Sanchez et al., ASPLOS'10] & [Ibanez et al., OSDI'21] w/ several opt.
- Remain RPC message payload in LLC/memory -- Only move RPC descriptor
 - Reduce latency overhead and traffic per migration (up to 140x)
- Each request being migrated at most once
 - Avoid livelock and unnecessary scheduling traffic
- Batch multiple descriptors per message
 - Improve hardware efficiency

Migration Prediction

- Periodically synchronize system states 1.
- ieue Depth 2. Pick appropriate threshold trained offline How we train SLO violation prediction model?
- Poll queues to check if threshold is met 3. Heuristics: multi-queue load pattern classification
- Decide how many RPCs to migrate to which queue(s) 4.
- 5. Pass decisions to h/w primitives and trigger migration h/w messages How we offload messaging to h/w? Repeat Step 1 (Period as short as 50ns due to messaging offloaded to h/w) 6.

- Determine a vector of SLO violation thresholds
 - Use queuing theory assisted w/ simulation Heade
- Different threshold per:
 - Load status
 - Service time distribution
 - Arrival pattern
 - Number of cores





Trade-off: Prediction accuracy V.S. migration effectiveness



Trade-off: Prediction accuracy V.S. migration effectiveness



A typical SLO value: 99th% latency <= 10 x average latency

Trade-off: Prediction accuracy V.S. migration effectiveness



A typical SLO value: 99th% latency <= 10 x average latency Naive prediction approach: Threshold = 10 x number of cores

Trade-off: Prediction accuracy V.S. migration effectiveness



Trade-off: Prediction accuracy V.S. migration effectiveness



A typical SLO value: 99th% latency <= 10 x average latency Naive prediction approach: Threshold = 10 x number of cores Aggressive prediction approach: Threshold = first SLO-violated queue length

- Trade-off: Prediction accuracy V.S. migration effectiveness
 - Prediction accuracy per migration: ~0%
 - Effectiveness (Capture ? % of SLO violation): 100%



A typical SLO value: 99th% latency <= 10 x average latency Naive prediction approach: Threshold = 10 x number of cores Aggressive prediction approach: Threshold = first SLO-violated queue length

Trade-off: Prediction accuracy V.S. migration effectiveness

We characterize thresholds for all system states offline & dynamically select threshold online



A typical SLO value: 99th% latency <= 10 x average latency Naive prediction approach: Threshold = 10 x number of cores Aggressive prediction approach: Threshold = first SLO violated queue length

Methodology

Baselines

- D-FCFS + work stealing system
 - ZygOS [Prekas et al, SOSP'17]
- C-FCFS S/W based system
 - Shinjuku [Kaffes et al, NSDI'19]
- C-FCFS H/W based system
 - Nebula [Sutherland et al, ISCA'20]
 - nanoPU [Ibanez et al, OSDI'21]

- Altocumulus configuration (Simulation)
 - Commodity RSS NIC +
 - S/W based local c-FCFS
 - comparable with S/W baselines
 - Integrated NIC +
 - H/W based local c-FCFS
 - comparable with H/W baselines











Outperform S/W solutions (ZygOS, Shinjuku) by up to 24x under 10 microsecond SLO



Outperform S/W solutions (ZygOS, Shinjuku) by up to 24x under 10 microsecond SLO



- Outperform S/W solutions (ZygOS, Shinjuku) by up to 24x under 10 microsecond SLO
- Has comparable throughput to highly-optimized H/W runtimes (Nebula & nanoPU)

















 \star <0.001% SLO violation and achieve 161 MRPS, scalable to 256 cores



 \star <0.001% SLO violation and achieve 161 MRPS, scalable to 256 cores

<5% SLO violation and achieve 216 MRPS, scalable to 256 cores



Key Evaluation: Scalability



 \star Achieve 161 MRPS while 99%th <= 5.57 us



 \star Achieve 161 MRPS while 99%th <= 5.57 us

Achieve 216 MRPS while 99%th <= 15.41 us



Conclusions



High volume of requests (400+ Gbps)







- SLO violation prediction via queueing theory, proactive migrations
 - S/W decentralized runtime with simple H/W primitives

Efficient direct register messaging w/ minimal H/W overheads

Altocumulus Scheduling Sub-system

MICRO-55, Altocumulus, <u>https://www.eecg.utoronto.ca/~mcj/papers/2022.altocumulus.micro.pdf</u>



Key Take-Aways



Altocumulus Stack

- With RPC stack components getting more optimized, other system-level components, e.g., scheduling, would become the future bottleneck
- Scheduling at <u>10s ns level is mandatory</u> for µs-scale SLO, to achieve that
 - Policy should be proactive and accurate/effective
 - Mechanism should be fast enough
- Carefully re-partitioning system stack between software or hardware can open opportunities for scalability that existing systems
 - Decentralized runtime preserves scalability
 - Software runtime offers adaptivity to various load patterns
 - Decentralized software runtime with minimal hardware overhead can reconcile design trade-offs across low tail latency, high through and high utilization (scalability)

Thank you for your attention

Jiechen Zhao, Iris Uwizeyimana, Karthik Ganesan Mark C. Jeffrey, Natalie Enright Jerger

55th IEEE/ACM International Symposium on Microarchitecture (MICRO), Chicago, Oct. 2022 https://www.eecg.utoronto.ca/~mcj/papers/2022.altocumulus.micro.pdf