

## Why Multi-Threading Should No Longer Be a DIY Job



Jannes Timm, Jan S. Rellermeyer



# **Concurrency is Hard...**

|   |                        |   | -   |                        |  |
|---|------------------------|---|---|------------------------|--|
| Zen3 L2   | 32M<br>L3              | L2 Zen3   | Zen3 L2   | 32M<br>L3              | L2 Zen3  |
| Zen3 L2   |                        | L2 Zen3   | Zen3 L2   |                        | L2 Zen3  |
| Zen3 L2   |                        | L2 Zen3   | Zen3 L2   |                        | L2 Zen3  |
| Zen3 L2   |                        | L2 Zen3   | Zen3 L2   |                        | L2 Zen3  |
|   |                        |   |   |                        |  |
| Zen3 L2   |                        | L2 Zen3   | Zen3 L2   | 32M<br>L3              | L2 Zen3  |
| Zen3 L2   | 32M                    | L2 Zen3   | Zen3 L2   |                        | L2 Zen3  |
| Zen3 L2   |                        | L2 Zen3   | Zen3 L2   |                        | L2 Zen3  |
| Zen3 L2   |                        | L2 Zen3   | Zen3 L2   |                        | L2 Zen3  |
|   |                        |   |   |                        |  |
|   |                        |   |   |                        |  |
| AMD Secure<br>Processor   |                        | DDR4 Memory<br>Controllers  | Server<br>Controller Hub  |                        | PCle3/4<br>SATA3   |
| AMD Secure<br>Processor<br>Zen3 L2  |                        | DDR4 Memory<br>Controllers<br>L2 Zen3   | Server<br>Controller Hub<br>Zen3 L2   |                        | PCle3/4<br>SATA3<br>L2 Zen3  |
| AMD Secure<br>Processor<br>Zen3 L2<br>Zen3 L2   |                        | DDR4 Memory<br>Controllers  | Server<br>Controller Hub<br>Zen3 L2<br>Zen3 L2  |                        | PCle3/4<br>SATA3<br>L2 Zen3<br>L2 Zen3   |
| AMD Secure<br>Processor<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2  | 32M<br>L3              | DDR4 Memory<br>Controllers  | Server<br>Controller Hub<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2   | 32M<br>L3              | PCle3/4<br>SATA3   |
| AMD Secure<br>Processor<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2  | 32M<br>L3              | DDR4 Memory<br>Controllers  | Server<br>Controller Hub<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2  | 32M<br>L3              | PCle3/4<br>SATA3   |
| AMD Secure<br>Processor<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2   | 32M<br>L3              | DDR4 Memory<br>Controllers  | Server<br>Controller Hub<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2  | 32M<br>L3              | PCIe3/4<br>SATA3   |
| AMD Secure<br>Processor<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2   | 32M<br>L3              | DDR4 Memory<br>Controllers  | Server<br>Controller Hub<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2  | 32M<br>L3              | PCIe3/4<br>SATA3   |
| AMD Secure<br>Processor<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2   | 32M<br>L3              | DDR4 Memory<br>Controllers  | Server<br>Controller Hub<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2  | 32M<br>L3              | PCIe3/4<br>SATA3   |
| AMD Secure<br>Processor<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12                                  | 32M<br>L3<br>32M<br>L3 | DDR4 Memory   12 Zen3   12 Zen3 | Server<br>Controller Hub<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12   | 32M<br>L3<br>32M<br>L3 | PCle3/4<br>SATA3<br>12 Zen3<br>12 Zen3<br>12 Zen3<br>12 Zen3<br>12 Zen3<br>12 Zen3<br>12 Zen3<br>12 Zen3 |
| AMD Secure<br>Processor<br>2 m3 12<br>2 m3 12 | 32M<br>L3<br>32M<br>L3 | Zen3   12 Zen3                  | Zena 12   Zena 12 | 32M<br>L3<br>32M<br>L3 | PCIa1/4<br>SATA3   |

# Hardware

etc.





# **To Thread or not to Thread?**

#### (Single-Threaded) Event Driven

[V.S. Pai, P. Druschel, W. Zwaenepoel: Flash: An efficient and portable Web server. In USENIX ATC 1999]

[F. Dabek, N. Zeldovich, N., F. Kaashoek, D. Mazieres, & R. Morris: Event-driven programming for robust software. In EuroSys 2002]

#### Multithreaded

[R. Von Behren, J. Condit & E. Brewer: Why Events Are a Bad Idea (for High-Concurrency Servers). In HotOS 2003]

# Thread Pools (reusing threads, sometimes adaptive using a Watermark model)

#### JVM (Executors), MariaDB, .NET CLR, etc.

[F.W. Burton & M.R. Sleep: Executing functional programs on a virtual tree of processors. In FPCA 1981]

#### **Resource-Aware Threads**

[R. Von Behren, J. Condit, F. Zhou, G.C. Necula & E. Brewer Capriccio: scalable threads for internet services. ACM SIGOPS OSR 2003]



# **State of the Art**

### Typical setup:

| Zen3 L2   | 32M<br>L3              | L2 Zen3  | Zen3 L2  | 32M<br>L3              | L2 Zen3  |
|---|------------------------|--|--|------------------------|--|
| Zen3 L2   |                        | L2 Zen3  | Zen3 L2  |                        | L2 Zen3  |
| Zen3 L2   |                        | L2 Zen3  | Zen3 L2  |                        | L2 Zen3  |
| Zen3 L2   |                        | L2 Zen3  | Zen3 L2  |                        | L2 Zen3  |
| Zen3 L2   |                        | L2 Zen3  | Zen3 L2  | 32M<br>L3              | L2 Zen3  |
| Zen3 L2   |                        | L2 Zen3  | Zen3 L2  |                        | L2 Zen3  |
| Zen3 L2   | 132M                   | L2 Zen3  | Zen3 L2  |                        | L2 Zen3  |
| Zen3 L2   |                        | L2 Zen3  | Zen3 L2  |                        | L2 Zen3  |
|   |                        |  |  |                        |  |
| AMD Secure<br>Processor   |                        | DDR4 Memory<br>Controllers   | Server<br>Controller Hub   |                        | PCle3/4<br>SATA3   |
| AMD Secure<br>Processor<br>Zen3 L2  |                        | DDR4 Memory<br>Controllers<br>L2 Zen3  | Server<br>Controller Hub<br>Zen3 L2  |                        | PCIe3/4<br>SATA3   |
| AMD Secure<br>Processor<br>Zen3 L2<br>Zen3 L2   | 32M                    | DDR4 Memory<br>Controllers   | Server<br>Controller Hub<br>Zen3 L2<br>Zen3 L2   | 32M                    | PCle3/4<br>SATA3<br>L2 Zen3<br>L2 Zen3   |
| AMD Secure<br>Processor<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2  | 32M<br>L3              | DDR4 Memory<br>Controllers   | Server<br>Controller Hub<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2  | 32M<br>L3              | PCle3/4<br>SATA3   |
| AMD Secure<br>Processor<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2   | 32M<br>L3              | DDR4 Memory<br>Controllers   | Server<br>Controller Hub<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2   | 32M<br>L3              | PCIe3/4<br>SATA3   |
| AMD Secure<br>Processor<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2  | 32M<br>L3              | DDR4 Memory<br>Controllers   | Server<br>Controller Hub<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2  | 32M<br>L3              | PCIe3/4<br>SATA3   |
| AMD Secure<br>Processor<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2<br>Zen3 L2                                  | 32M<br>L3              | L2 Zen3<br>L2 Zen3<br>L2 Zen3<br>L2 Zen3<br>L2 Zen3<br>L2 Zen3<br>L2 Zen3<br>L2 Zen3   | Server<br>Controller Hub   | 32M<br>L3              | PCIe3/4<br>SATA3   |
| AMD Secure<br>Processor<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12<br>Zen3 12                       | 32M<br>L3<br>32M<br>L3 | DDRA Memory<br>Controllers   12 Zen3                     | Senser   Controller Hub   Zen3 12    | 32M<br>L3<br>32M<br>L3 | PCIA:J/4<br>SATA3<br>12 Zon3<br>12 Zon3<br>12 Zon3<br>12 Zon3<br>12 Zon3<br>12 Zon3<br>12 Zon3           |
| AMD Secure<br>Processor<br>2 m3 12<br>2 m3 12 | 32M<br>L3<br>32M<br>L3 | DDR4 Memory<br>Controllers   12 Zen3   12 Zen3 | Senser   Zen3 12   Zen3 12 | 32M<br>L3<br>32M<br>L3 | PCIA1/A<br>SATA3<br>12 Zan3<br>12 Zan3<br>12 Zan3<br>12 Zan3<br>12 Zan3<br>12 Zan3<br>12 Zan3<br>12 Zan3 |

# #threads = #cores

flawed!

### + manual tuning knob

### Hardware



# When Less is More



[S. Omranian Khorasani, J.S. Rellermeyer, D. Epema: Self-Adaptive Executors for Big Data Processing. In: Middleware 2019]



# **Towards Self-Adaptive Thread Pools**

Target Metric

$$rwchar-rate \quad R[t_1, t_2] = \\ \sum_{w \in W} \frac{(rchar_{w,t_2} + wchar_{w,t_2}) - (rchar_{w,t_1} + wchar_{w,t_1})}{t_2 - t_1}$$





# **Controller Logic**

Hill-climbing approach



For most experiments: 1500ms interval length and 10% adjustment



# **Experimental Results: Synthetic Workload**

- Workers read a file of 2MiB into memory and write it back to disk ("nosync")
  - Variant "sync": writeback, call fsync
  - Variant "nosync-sync" first batch nosync, second batch sync

- Compare three setups:
  - "fixed": experimental optimum, determined through parameter sweep\*
  - *"watermark"*: classic high-low watermark model
  - *"adaptive"*: our solution using hill-climbing with rwchar-rate

\*3% tolerance and favor lower thread pool size



# **Experimental Results: Synthetic Workload**





# **Performance Variability Happens...**

### In the small:



#### In the large:

[A. Uta, A. Custura, D. Duplyakin, I. Jimenez, J.S. Rellermeyer, C. Maltzahn, R. Ricci & A. Iosup: Is Big Data Performance Reproducible in Modern Cloud Networks? In NSDI 2020]



# **Experimental Results: RocksDB**

- Write-heavy workload: sequential key insertion (fillseq)
- Flush pool size has significant influence on the performance
  - Blue: default setting, Red: hand-optimized
- Interval length has an impact on the adaptive solution.
  - Green: 1500 ms, Olive: 1000 ms





# Why Thread Pool Management Should be an OS Service

- We can use the machine resources more efficiently when tuning thread pools for I/O-intensive applications
- Manual thread pool tuning is tedious...
- ... but also futile if the environment changes
  - Multi-tenant environments like containers, etc.
  - Statically tuned applications are selfish
- The OS already collects performance metrics and is supposed to be a mediator
- Providing thread pools as an OS service would allow us to implement better QOS for multi-tenant environments