

Auto-Tuned Threading for OLDI Microservices

Akshitha Sriraman
University of Michigan

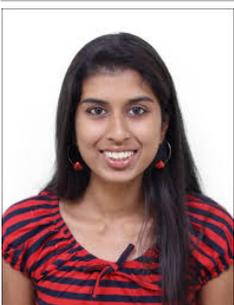
Thursday, October 11th, 2018

2:00 pm
EEB 248

Modern On-Line Data Intensive (OLDI) applications have evolved from monolithic systems to instead comprise numerous, distributed microservices interacting via Remote Procedure Calls (RPCs). Microservices face sub-ms RPC latency goals, much tighter than their monolithic ancestors that must meet ≥ 100 ms latency targets. Sub-ms-scale threading and on currency design effects as well as OS and network overheads that were once insignificant for such monoliths, can now come to dominate in the sub-ms-scale microservice regime. It is therefore vital to characterize the influence of threading design, OS, and network effects on microservices. Unfortunately, widely used academic data center benchmark suites are unsuitable to aid this characterization as they use monolithic rather than microservice architectures.

We first investigate how OS/network overheads impact microservice tail latency by developing a complete suite of microservices called μ Suite that we use to facilitate our study. Our characterization reveals that the relationship between optimal OS/network parameters and service load is complex. Our primary finding is that non-optimal OS scheduler decisions can degrade microservice tail latency by up to $\sim 87\%$.

Secondly, we investigate how threading design critically impacts microservice tail latency by developing a taxonomy of threading models – a structured understanding of the implications of how microservices manage concurrency and interact with RPC interfaces under wide-ranging loads. We develop μ Tune, a system that has two features: (1) a novel framework that abstracts threading model implementation from application code, and (2) a novel automatic load adaptation system that curtails microservice tail latency by exploiting inherent latency trade-offs revealed in our taxonomy to transition among threading models. We study μ Tune in the context of μ Suite to demonstrate up to 1.9x tail latency improvements over static threading choices and state-of-the-art adaptation techniques.



Akshitha is a fourth year Ph.D. student at the University of Michigan, where she is advised by Dr. Thomas F. Wenisch. Her primary research interests are in software systems and computer architecture. Her research focuses on developing software and hardware optimizations to improve the performance of large-scale distributed data center systems.

Hosted by Xuehai Qian, x04459, xuehai.qian@usc.edu