# µTune: Auto-Tuned Threading for OLDI Microservices
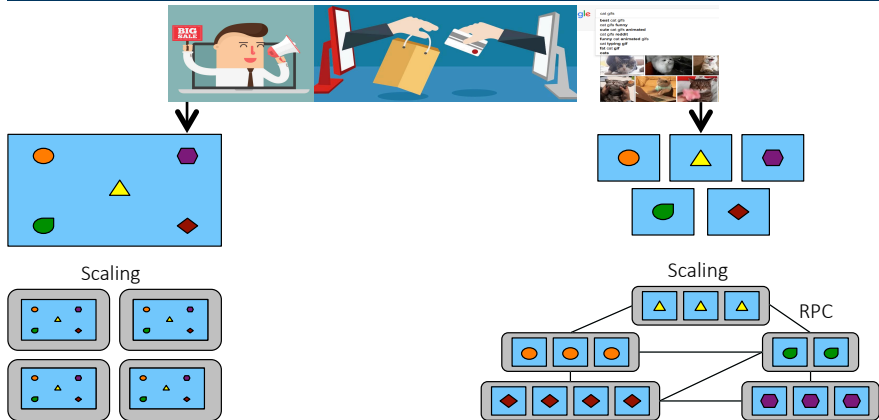
Akshitha Sriraman, Thomas F. Wenisch
University of Michigan

## On-Line Data Intensive Applications
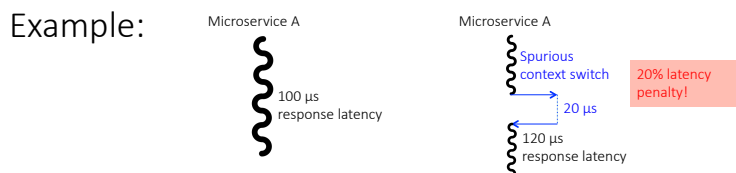
Scaling

Scaling

RPC

Monoliths (>100 ms SLO) ⟶ Microservices (sub-ms SLO)

## Impact of Threading-Induced Overhead

µ

Threading overheads

Lock contention

Thread wakeups

Spurious context switch

Impact: Minor for monoliths & major for microservices

Example:

Microservice A

100 µs response latency

Microservice A

Spurious context switch

20% latency penalty!

20 µs

120 µs response latency

## Threading Impact on Mid-Tier

Front-End Microserver

Mid-Tier Microserver

Leaf Microserver 1

Leaf Microserver 2

Mid-Tier: Heavily impacted by threading
- Server & client
- Fans queries to many leaves
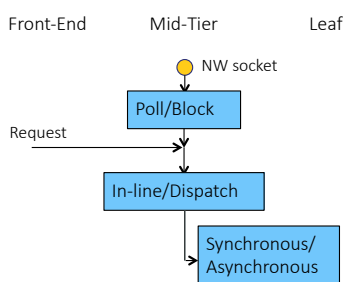- RPC layer interactions dominate compute

## Contributions

- A taxonomy of threading models
  - Structured understanding of threading implications
  - Reveals tail inflection points across load
  - Peak load-sustaining model is subpar at low load
- µTune:
  - Uses tail inflection insights to optimize tail latency
  - Tunes model & thread pool size across load
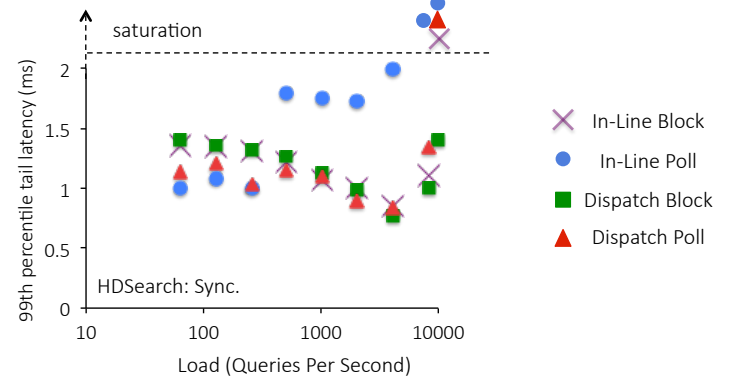  - Simple interface: Abstracts threading from RPC code

## Taxonomy of Threading Models

Front-End          Mid-Tier          Leaf

NW socket

Request

Poll/Block

In-line/Dispatch

Synchronous/Asynchronous

Threading dimensions:
- Block vs. Poll
- In-line vs. Dispatch
- Sync. vs. Async.

## Taxonomy Characterization

HDSearch: Sync.

saturation

× In-Line Block
● In-Line Poll
■ Dispatch Block
▲ Dispatch Poll

99th percentile tail latency (ms)

Load (Queries Per Second)
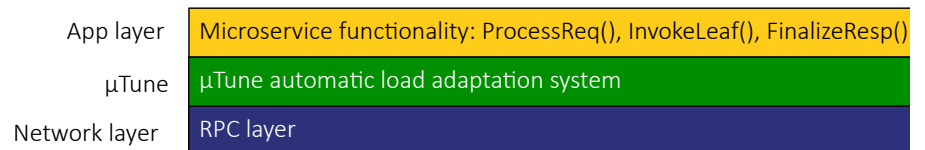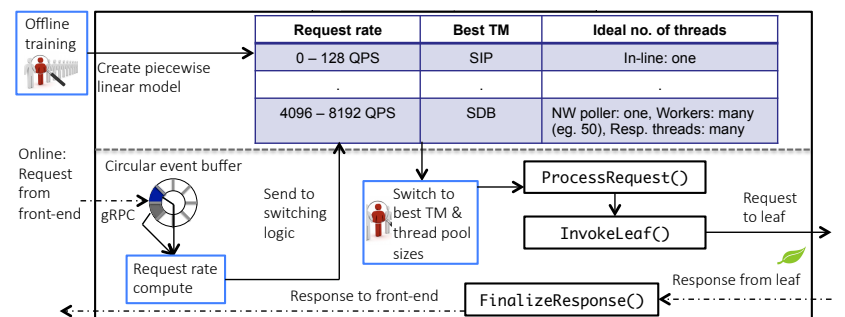
No single threading model works best at all loads

## µTune: Automatic Load Adaptation

Abstracts threading boiler-plate code from RPC code

| App layer | Microservice functionality: ProcessReq(), InvokeLeaf(), FinalizeResp() |
|---|---|
| µTune | µTune automatic load adaptation system |
| Network layer | RPC layer |

System design: offline training + run-time adaptation

Offline training

Create piecewise linear model

| Request rate | Best TM | Ideal no. of threads |
|---|---|---|
| 0 – 128 QPS | SIP | In-line: one |
| . | . | . |
| 4096 – 8192 QPS | SDB | NW poller: one, Workers: many (eg. 50), Resp. threads: many |

Online: Request from front-end

gRPC

Circular event buffer

Request rate compute

Send to switching logic

Switch to best TM & thread pool sizes

ProcessRequest()

InvokeLeaf()

FinalizeResponse()

Request to leaf

Response from leaf

Response to front-end

## Result: µTune Under Steady-State Load

■ In-Line Poll  ■ Dispatch Poll  ■ Dispatch Block  ■ Haque '15  □ Abdelzaher '99  ■ µTune

HDSearch: Async.

saturation

1.9x

∞

∞ 6.17

<5% mean overhead

99th percentile tail latency (ms)

Load (Queries Per Second)

- µTune converges to best threading model and thread pool size to improve tail latency by up to **1.9x** over static peak load-sustaining threading model with < 5% mean overhead