TOSS: Tiering of Serverless Snapshots for Memory-Efficient Serverless Computing

Theodore Michailidis, Juno Kim, Linsong Guo, Steven Swanson, Jishen Zhao

University of California, San Diego



Serverless computing

► Serverless computing is a **billion-dollar** industry (exp. \$200 billion by 2035)











Serverless computing



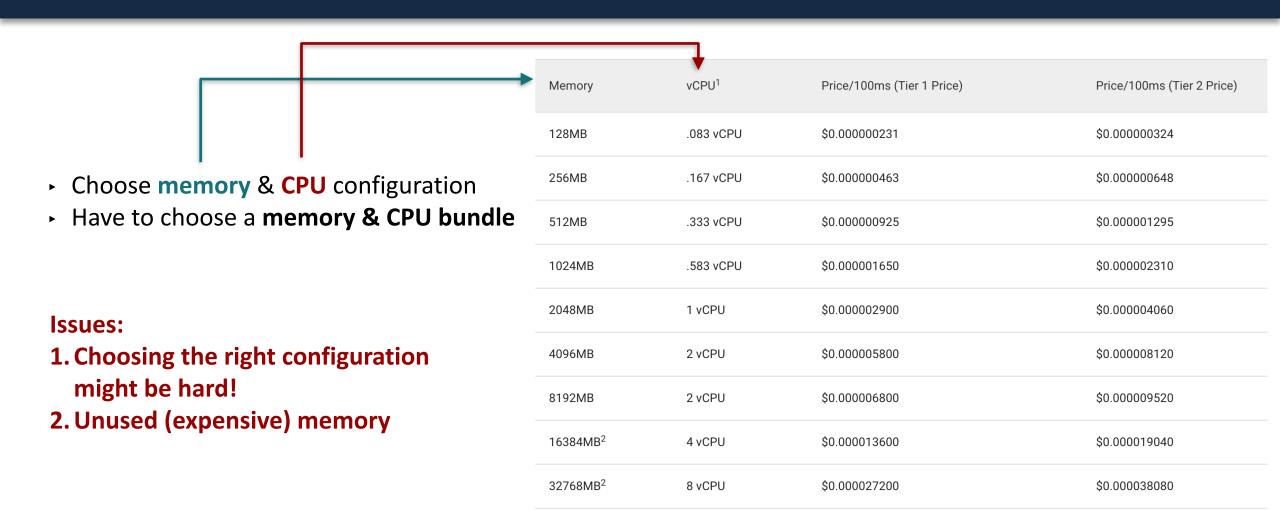
- ✓ Write code
- √ Choose Virtual Machine (VM) configuration
 - Pay what you use
- **✓** Deploy



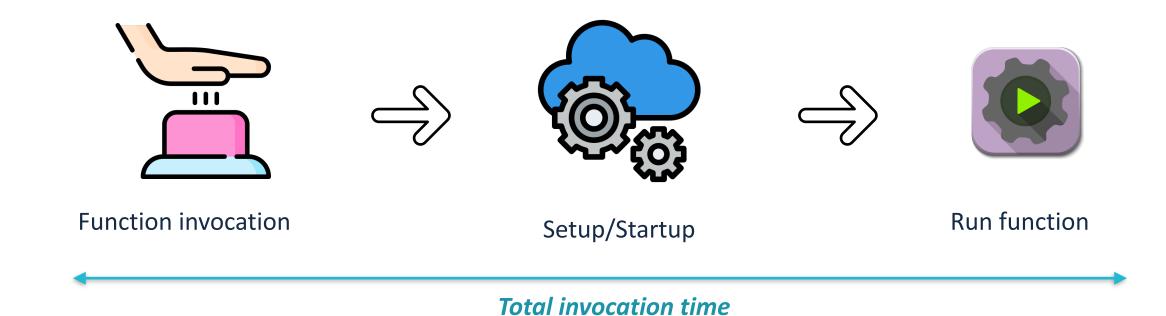
- **✓** VM Lifecycle
- ✓Infrastructure
- **✓** Security
- **✓** Storage

- **✓** Scalability
- ✓ Pricing

Serverless Pricing



Serverless invocation process



Serverless invocation process



Running time: **50%** of functions < 1s, **96%** of functions < 10s

Source: Serverless in the Wild, Shahrad et al., ATC'20

Serverless invocation process



160-3,600ms

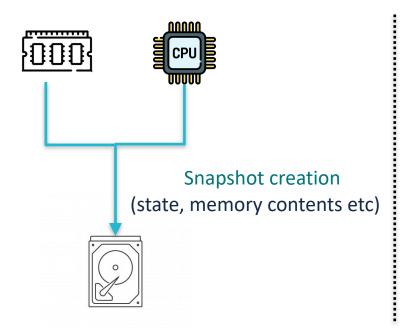
Source: Peeking Behind the Curtains of Serverless Platforms, Wang et al., ATC'18

Solutions for cold start

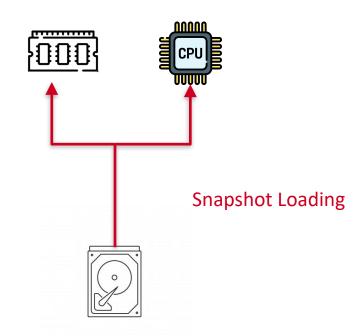
- Keep-alive
 - Functions remain in memory for predefined time
 - Standard (e.g. 30 minutes)
 - Determined by previous invocation patterns
- Snapshots
 - Fast, simple, low-cost solution
 - Lots of attraction both from industry and academia

Snapshots

Initial function invocation

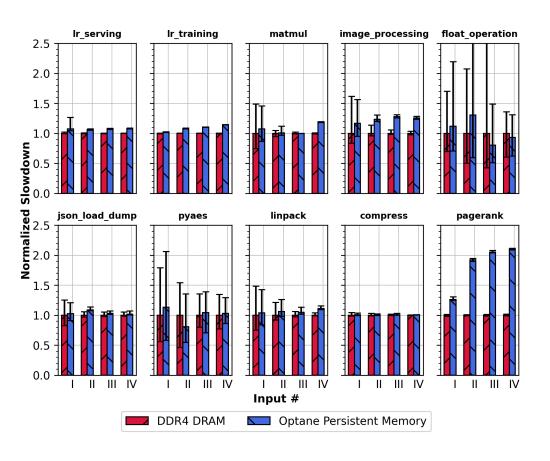


Before subsequent function invocations

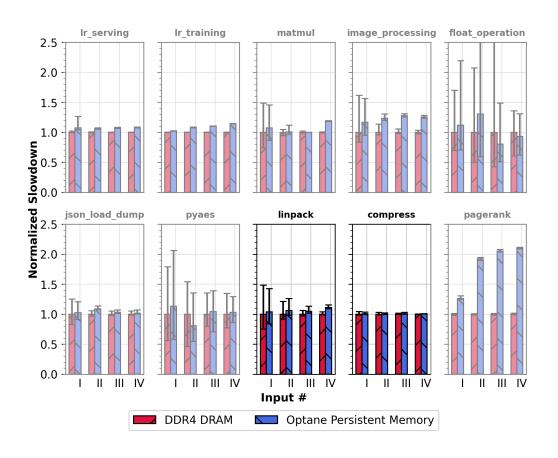


- Single pass snapshot not accounting for all inputs
- Assumes a single tier of (costly) memory

10 functions, 4 inputs per function (run entirely on fast or slow tier)

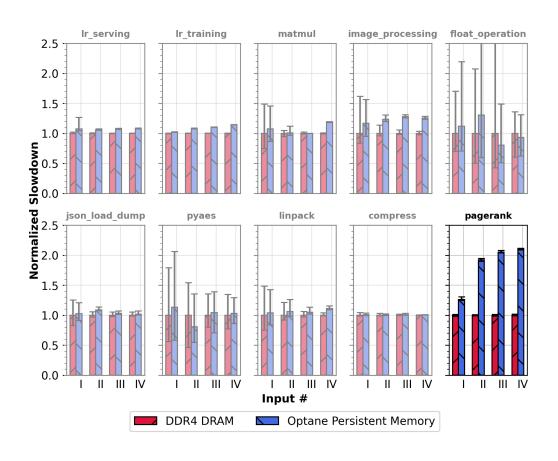


10 functions, 4 inputs per function (run entirely on fast or slow tier)



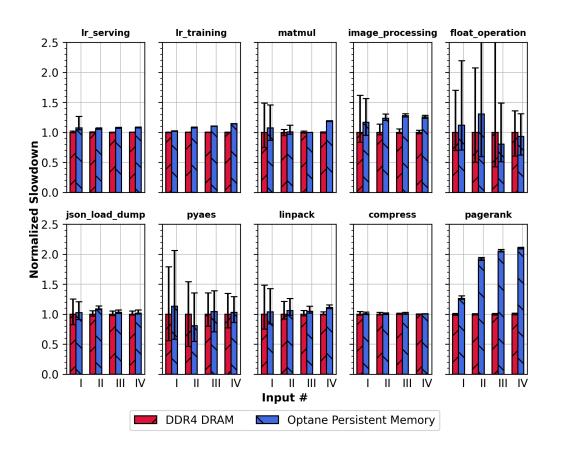
✓ Some functions display minor or no slowdown on slower memory

10 functions, 4 inputs per function (run entirely on fast or slow tier)



- ✓ Some functions display minor or no slower memory
- ✓ Different inputs display variable slowdown
 - Increased input size leads to increased memory footprint and slowdown

10 functions, 4 inputs per function (run entirely on fast or slow tier)



- ✓ Some functions display minor or no slower memory
- ✓ Different inputs display variable slowdown
 - Increased input size leads to increased memory footprint and slowdown



Need for something more sophisticated!

Pathologies on previous approaches

Single tier of memory

- Expensive: Memory accounts for 50% of total server cost
- Redundant: Usually only a subset of memory should go to the fast tier

Serveless bundle

- Hard to find optimal configuration
- Wasted memory

Single snapshot

 A single snapshot does not represent the diversity of function inputs

Tiering of Serverless Snapshots (TOSS)

TOSS

First Memory Tiering system for serverless

Goal: Reduce serverless memory cost operation

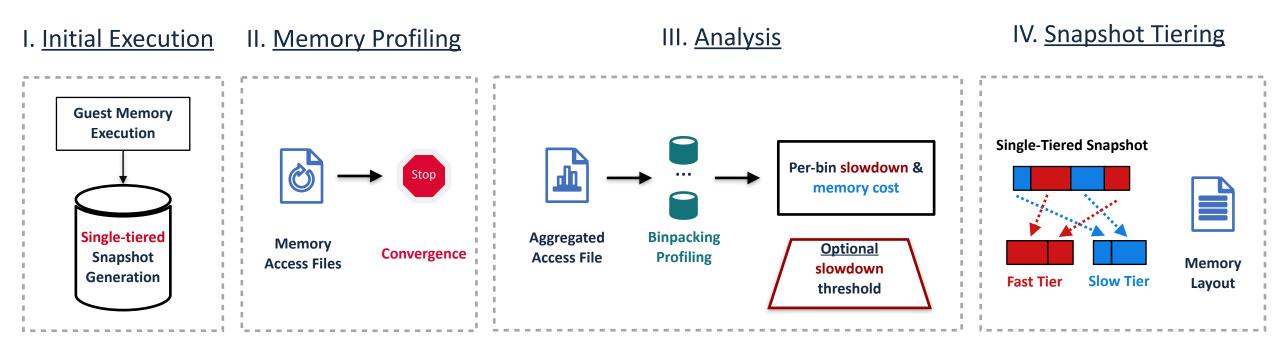
Insights about serverless memory

- Study serverless functions' memory patterns
- Correlation about slowdown and memory cost

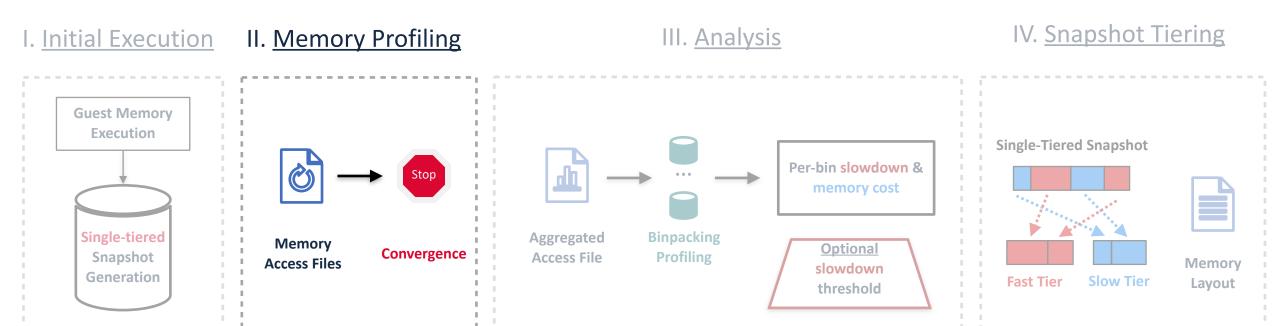
Flexible memory cost formula

- Cost per tier
- Memory per tier
- Slowdown

TOSS Overview



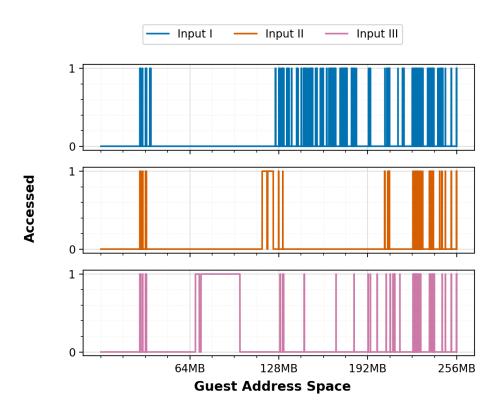
Memory profiling



- Accurate with low overhead
- Should reflect different memory patterns from different inputs

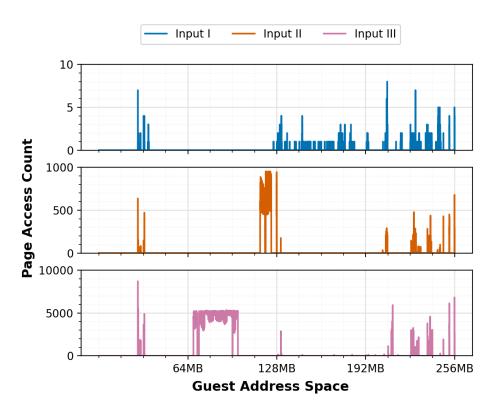
Serverless memory profiling

Previous approaches



- * Hot (accessed) vs Cold (not-accessed)
- * High overhead (15% 87%)

TOSS (DAMON)



- ✓ Lightweight, accurate and scalable
- ✓ Fine-grained access memory pattern

 $Memory\ Cost = Slowdown * (MB_{Fast} * CostPerMB_{Fast} + MB_{Slow} * CostPerMB_{Slow})$

Memory in MB that resides in the fast and slow tiers



Cost ratio per MB between tiers



Slowdown for this memory configuration

$$Memory\ Cost = Slowdown * (MB_{Fast} * CostPerMB_{Fast} + MB_{Slow} * CostPerMB_{Slow})$$

Evaluation - Setup & Systems

CPU: 2×20-core Intel® Xeon® Gold 6230

Memory: 96GB DDR4 DRAM (fast) & 192GB Intel® Optane™ Persistent Memory (slow)

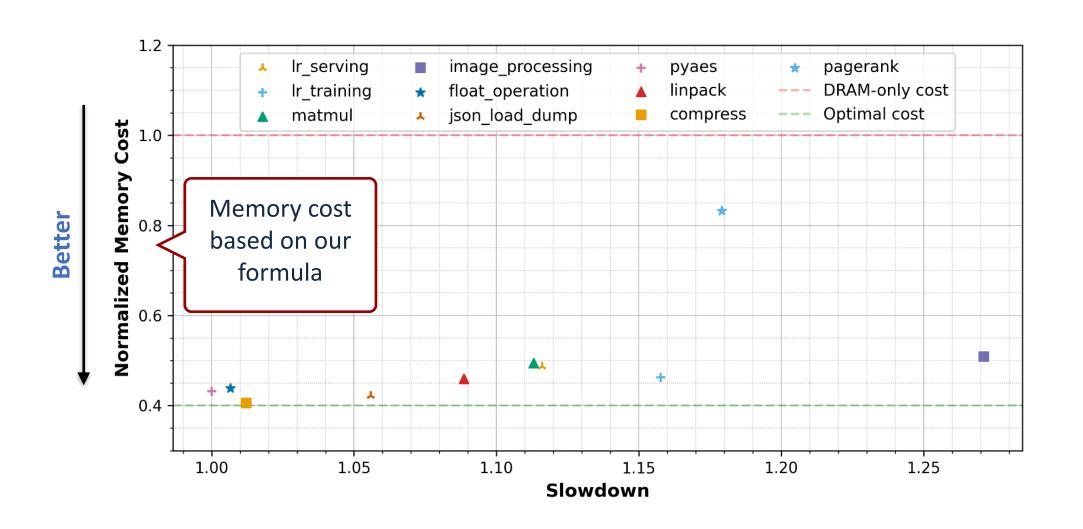
SSD: Intel® Optane™ DC (seq read and write of up to 2,500 MB/s and 2,200 MB/s)

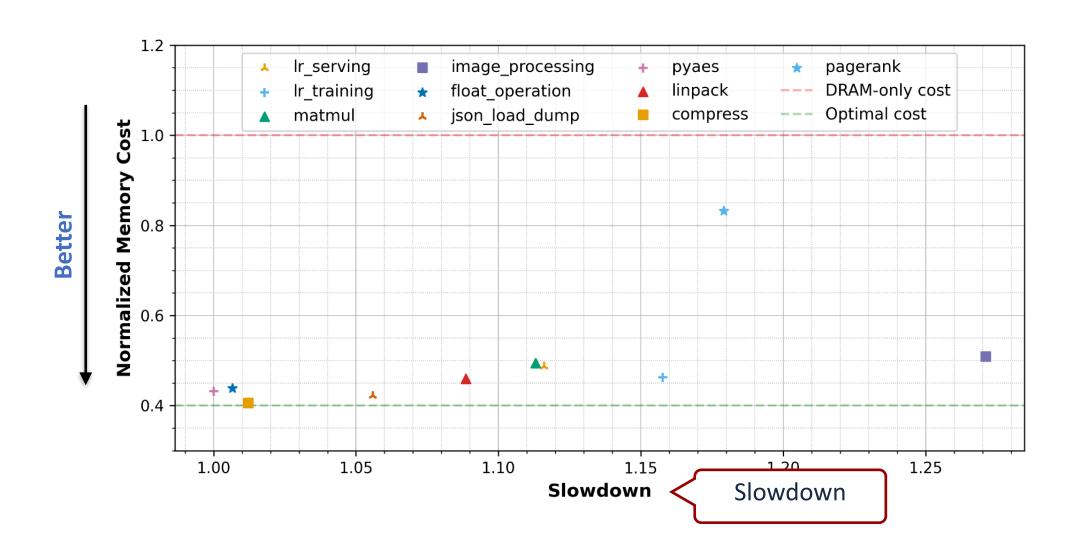
Functions from **FunctionBench** and **SeBS**

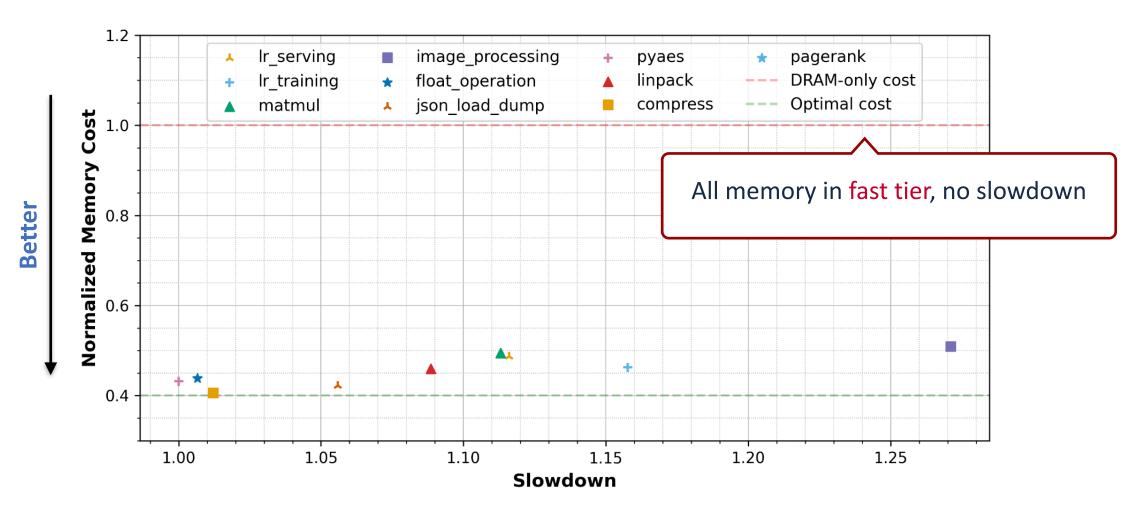
Compression, graph processing, image processing, ML training and inference,
microbenchmarks

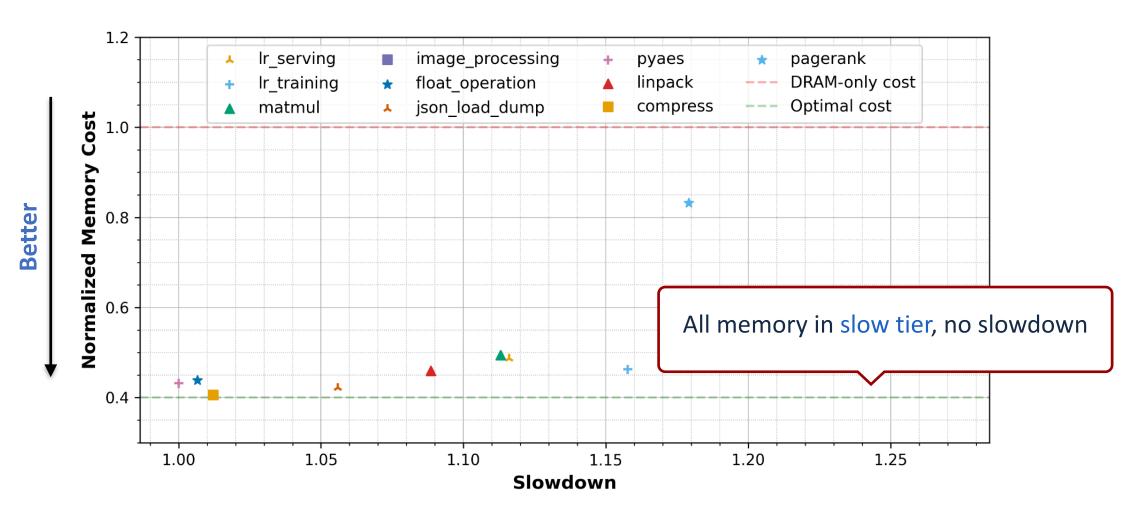
TOSS (input IV, optimal cost configuration) vs REAP

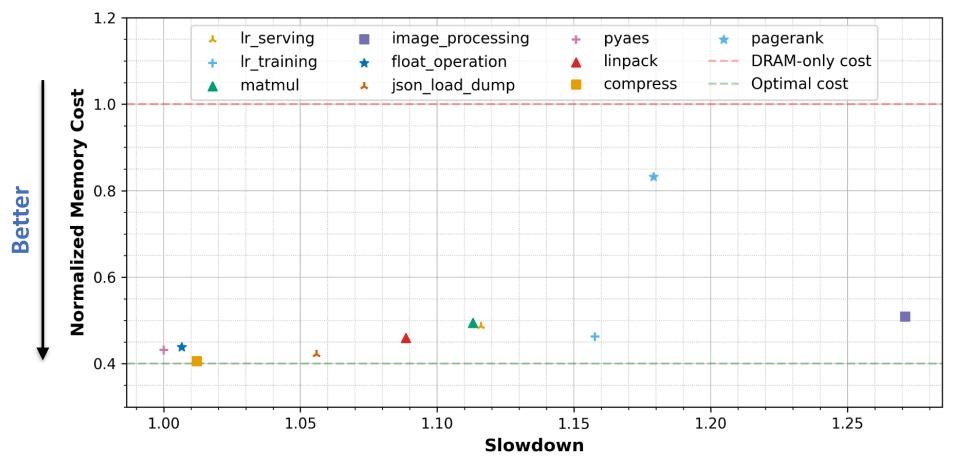
- Working set into memory, the rest loaded on-demand
- Reduced page faults & parallelized prefetching



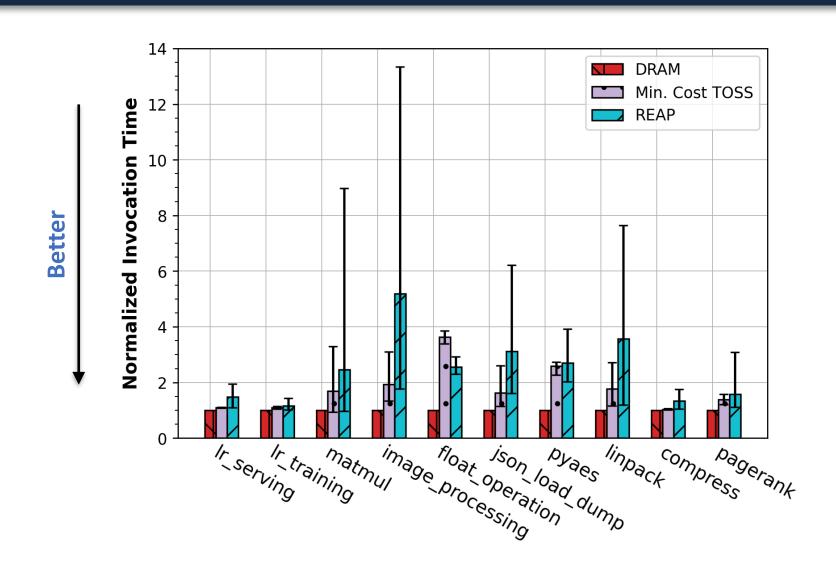


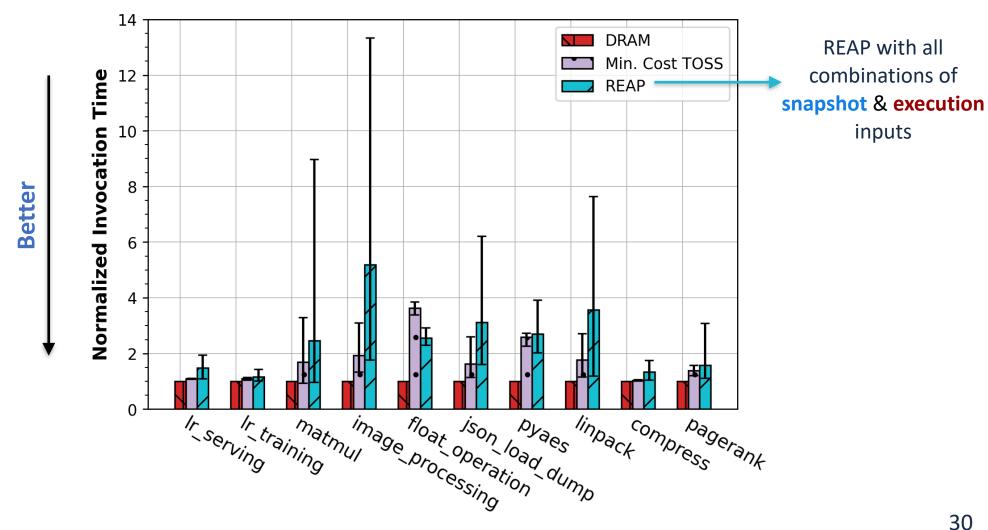


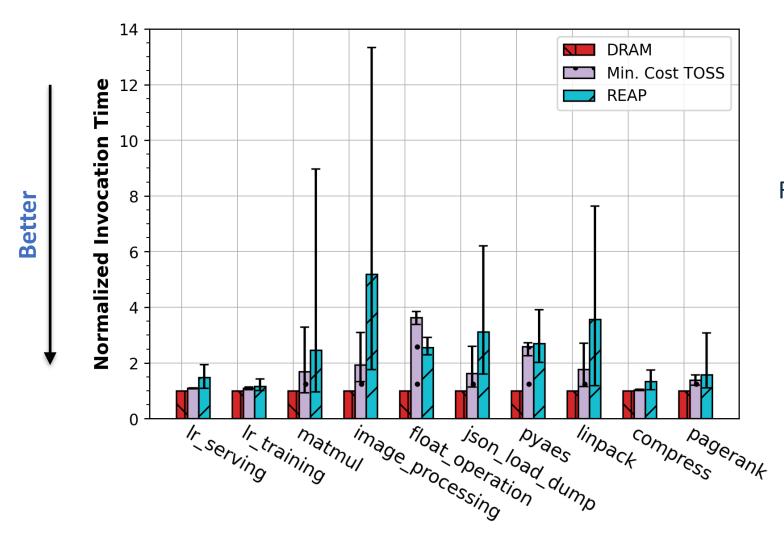




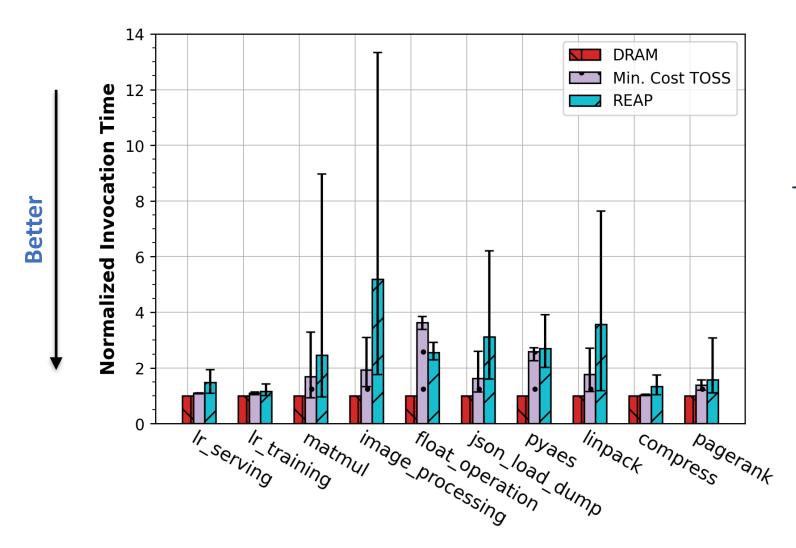
- Memory Cost Reduction
 - Average: **51%**
 - Max: **59% (~ optimal)**
- Function Slowdown
 - Average: **10%**
 - Max: 27%
- Offloaded memory
 - Average: **92%**
 - Max: **100**%



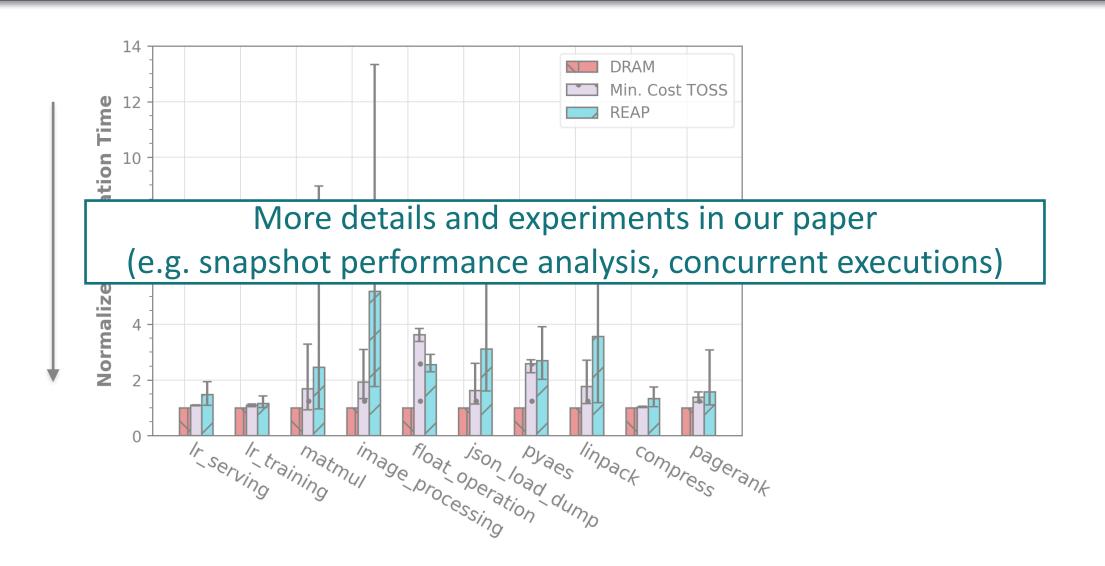




REAP's performance depends on the similarity between snapshot & execution input



TOSS achieves on average 1.7× and up to 4.2× lower total invocation time



Summary

- Problem: Inflated memory costs by assuming a single memory tier
- Key insights
 - Different inputs lead to different memory access patters
 - Functions use heavily a small memory subset → cost-effective memory tiering without major slowdowns
- TOSS: First memory tiering mechanism for serverless functions
 - Memory study & characterization
 - Cost model to guide automatic tiered snapshot generation

TOSS: Tiering of Serverless Snapshots for Memory-Efficient Serverless Computing

Questions?

