HardHarvest: Hardware-Supported Core Harvesting for Microservices



|SCA|25

Jovan Stojkovic, Chunao Liu*, Muhammad Shahbaz*, Josep Torrellas University of Illinois at Urbana-Champaign, *Purdue University J. Stojkovic to start at UT Austin







Microservices

Large monolithic applications decomposed into many small interdependent services Each service implements separate functionality







Scalability Design simplicity HW management

Benefits of Microservices







Microservices are Widely Used



Structure of microservices at Amazon. Looks almost like a Death Star but is way more powerful.







O To accommodate the peak load, microservices typically overprovisioned





3 large Microsoft services ~1M virtual cores



O To accommodate the peak load, microservices typically overprovisioned O Requests often stalled on synchronous RPCs to read/write to/from remote storage, or to invoke other microservices









O To accommodate the peak load, microservices typically overprovisioned O Requests often stalled on synchronous RPCs to read/write to/from remote storage, or to invoke other microservices O As a result → cores are heavily underutilized





90% of services have maximum core utilization less than 40%!



O Collocate latency-critical microservices with compute-heavy batch apps **O Primary VMs** O Run latency-critical workloads O Users request a number of cores











O Collocate latency-critical microservices with compute-heavy batch apps **O Harvest VMs O Primary VMs** O Run batch workloads in the background O Run latency-critical workloads O Steal idle Primary VM's cores & return them O Users request a number of cores









O Collocate latency-critical microservices with compute-heavy batch apps **O Harvest VMs O Primary VMs** O Run latency-critical workloads O Run batch workloads in the background O Steal idle Primary VM's cores & return them O Users request a number of cores



Core Harvesting to the Rescue









O Collocate latency-critical microservices with compute-heavy batch apps **O Harvest VMs O Primary VMs** O Run latency-critical workloads O Run batch workloads in the background O Steal idle Primary VM's cores & return them O Users request a number of cores



Core Harvesting to the Rescue

Primary VM









O Collocate latency-critical microservices with compute-heavy batch apps **O Harvest VMs O Primary VMs** O Run latency-critical workloads O Run batch workloads in the background O Steal idle Primary VM's cores & return them O Users request a number of cores









O Improved resource utilization O Higher throughput for Harvest VMs









O Improved resource utilization O Higher throughput for Harvest VMs \bigcirc Software core harvesting \rightarrow increase tail latency for Primary VMs!

O Moving core from one VM to another is expensive ONeed for two hypervisor calls (detach and attach) OCross-VM context switch

O Moving core from one VM to another is expensive ONeed for two hypervisor calls (detach and attach) OCross-VM context switch

 $\overline{\mathbb{O}}$ Q P99

No Harvesting

O Moving core from one VM to another is expensive ONeed for two hypervisor calls (detach and attach) OCross-VM context switch

ater P99

O For security \rightarrow different VMs should not observe any state left in caches/TLBs On a cross-VM context switch, need to flush and invalidate caches/TLBs

O For security \rightarrow different VMs should not observe any state left in caches/TLBs On a cross-VM context switch, need to flush and invalidate caches/TLBs

<u>O</u> 66

\odot For security \rightarrow different VMs should not observe any state left in caches/TLBs On a cross-VM context switch, need to flush and invalidate caches/TLBs

ater 66

Core Reassignment

No Harvesting

20

OHigh core utilization OHigh Harvest VM throughput OLow Primary VM tail latency

21

Controutions

O Characterize the opportunities of supporting hardware-based core harvesting in microservice environments O Propose HardHarvest OThe first architecture for hardware core harvesting • Compared to state-of-the-art software core harvesting **O Higher core utilization:** 1.5x **O Higher Harvest VM throughput:** 1.8x **O Lower Primary VM tail latency: 6.0x**

Microservices sensitive to cache misses, but tolerate Cache/TLB downsizing Tail latency remains nearly unchanged even when LLC is reduced to its ¹/₄

O Microservices sensitive to cache misses, but tolerate Cache/TLB downsizing Otail latency remains nearly unchanged even when LLC is reduced to its $\frac{1}{4}$ O Enables selective flushing and partitioning OReserve cache/TLB space for the Primary VM and preserve useful data

1. Software overheads of core reassigning high \rightarrow propose a hardware solution 2. Cache/TLB flushing expensive but needed for security \rightarrow partition caches/TLBs with smart replacement policy

Request Queue (RQ)

Primary VM2 Subqueue Primary VM1 Subqueue

27

Request Queue (RQ)

Queue

31

Harvest VM

Harvest VM

Shared across requests \rightarrow private per-invocation

- most likely to be reused in the future
- O Shared pages: program code, libraries, read-only input data
- O Private pages: allocated after a microservice request has arrived
- O Shared pages more likely to be reused in the future

• For performance, HardHarvest keeps in Non-Harvest region the data that is

If any slot empty \rightarrow Pick that slot

Addr Shared? A Prioritize Non-Harvest

Addr **Private**? A Prioritize Harvest

O Systems evaluated

- O Cycle-accurate full-system simulations: SST + QEMU O DeathStarBench microservices with Alibaba's invocation traces **OsoftHarvest (EuroSys'21):** software core harvesting
 - **O HardHarvest:** our proposal

ONOHarvest: conventional system where no VM performs core harvesting

- 10
- \mathbb{U}

NoHarvest

SoftHarvest

HardHarvest

HardHarvest Reduces Primary VMs' Tail Latency

- 1 ()

State-of-the-art Our proposal

NoHarvest	So	ftHarv	est		HardHar	/es ⁻

HardHarvest Reduces Primary VMs' Tail Latency

State-of-the-art Our proposal

P99 tail latency reduced 6x compared to state-of-the-art!

HardHarvest

	25	
	0.0	
ghput	3	
	2.5	
\hat{D}		
alized Thro	2	
	15	
	1.0	
D	1	
OLL		
Ζ	0.5	
	0	

2		
0		
25		
۷.۷		
2		
1.5		
1		
•		
0.5		
\cap		

NoHarvest

SoftHarvest

HardHarvest

HardHarvest Increases Harvest VMs' Throughput

	25			
	J.J			
+	3			
id d b	2.5			
L D C L D C	2			
Zeo Zeo	1.5			
r M G I	1			
0 Z ().5			
	0			
		NoHarvest	SoftHarvest	HardHarvest

State-of-the-art Our proposal

HardHarvest Increases Harvest VMs' Throughput

	3.5	
	3	
	2.5	
	2	
	1.5	Th
	1	СО
	0.5	

U

State-of-the-art Our proposal

roughput increased 1.8x mpared to state-of-the-art!

NoHarvest

SoftHarvest

HardHarvest

	100	
	80	
	60	
O D D D D	40	
	20	
	0	

NoHarvest

SoftHarvest

HardHarvest

	100	
	80	
	60	
O D D D D	40	
	20	
	0	

State-of-the-art Our proposal

arvest		HardHarves ⁻		

U

State-of-the-art Our proposal

Core utilization increased 1.5x compared to state-of-the-art!

NoHarvest

SoftHarvest

HardHarvest

O Microservices beneficial, but have low core utilization • How to design efficient core harvesting scheme for microservices? O HardHarvest: the first architecture for in-hardware core harvesting O 6x lower P99 tail latency of Primary VMs, 1.8x higher throughput of Harvest VMs, and 1.5x higher core utilization

HardHarvest: Hardware-Supported Core Harvesting for Microservices

|SCA|25

Jovan Stojkovic, Chunao Liu*, Muhammad Shahbaz*, Josep Torrellas University of Illinois at Urbana-Champaign, *Purdue University

