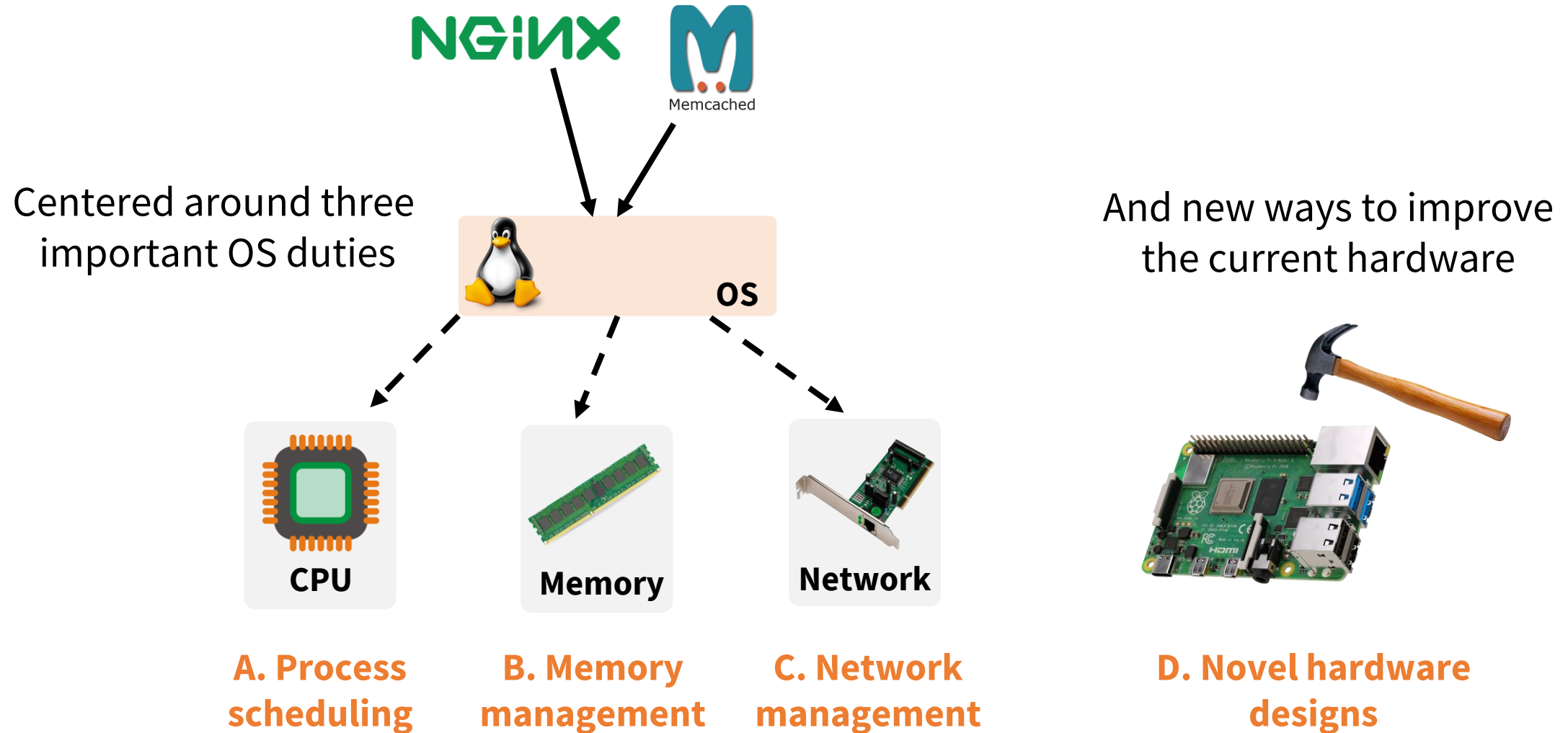


# OS and Hardware

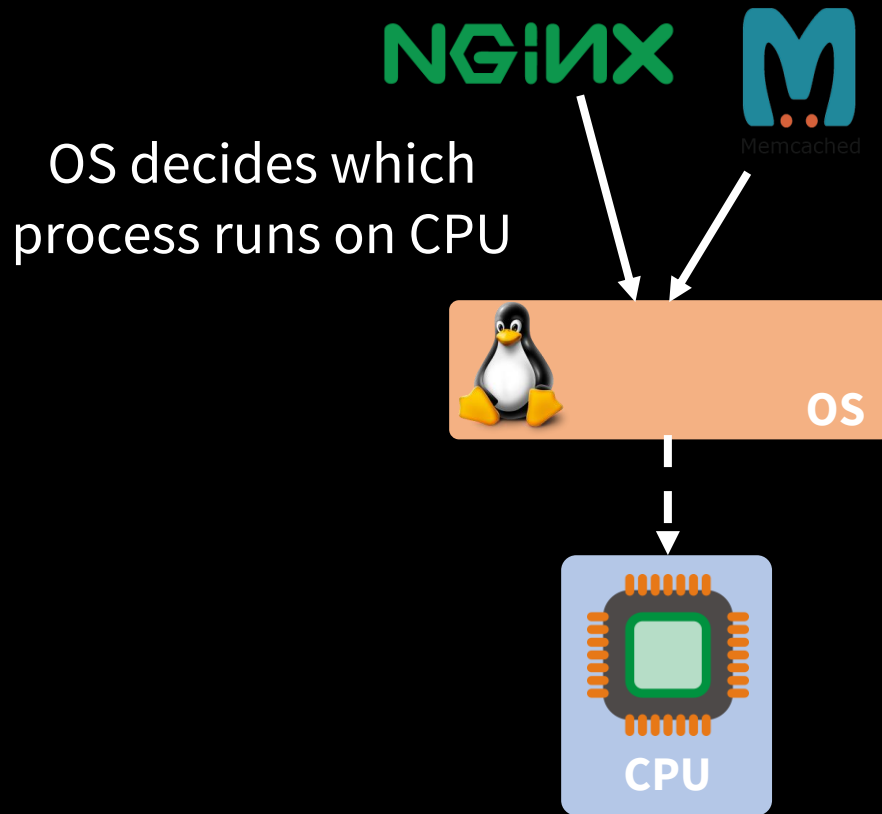
Adil Ahmad



# Categorizing the papers at OSDI and ATC



# A. Process scheduling



Two relevant concepts:

**Fairness:** each process should be allowed to use the CPU equally

**Latency:** each process should complete as soon as possible

## A: Process scheduling

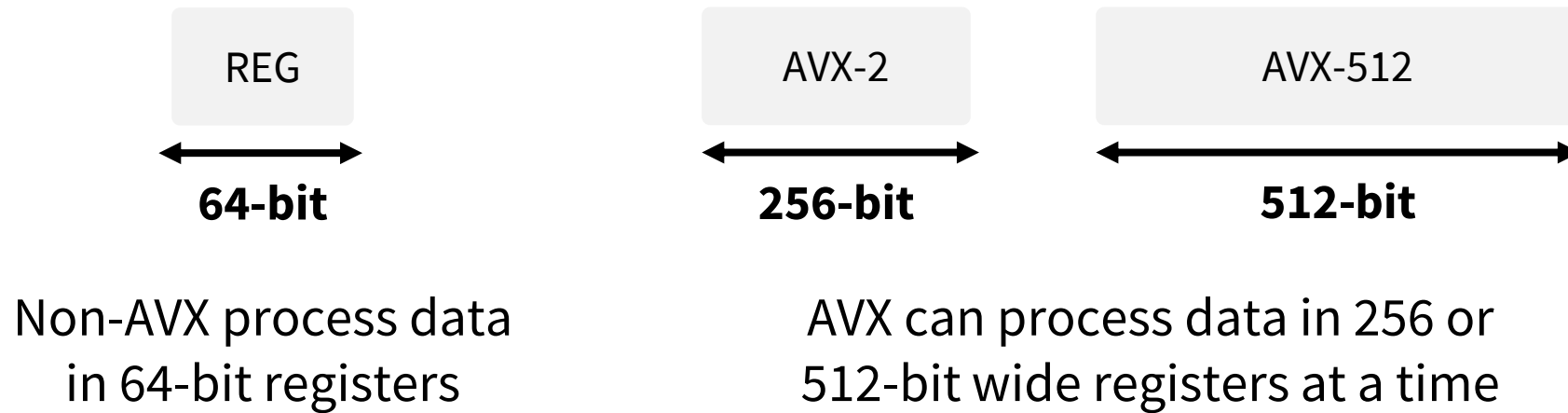
### Paper #1

**“Fair Scheduling for AVX2 and AVX-512 Workloads”**  
USENIX ATC 2021



Identify that AVX causes scheduling unfairness  
and compensate the affected processes

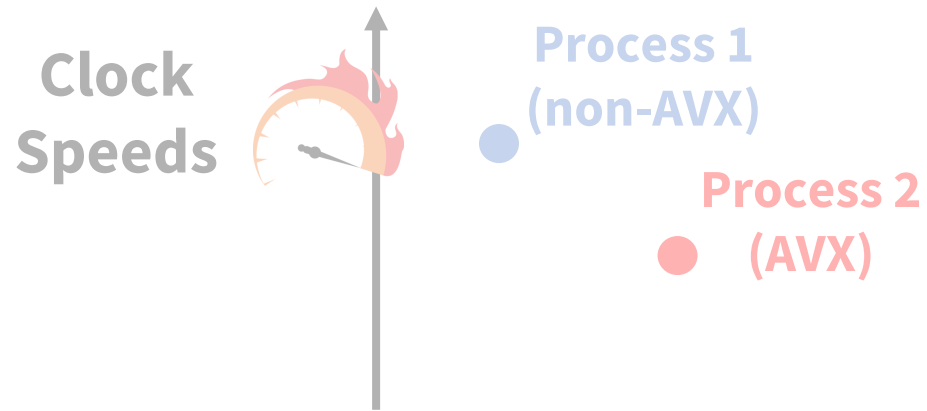
# Background: AVX instructions



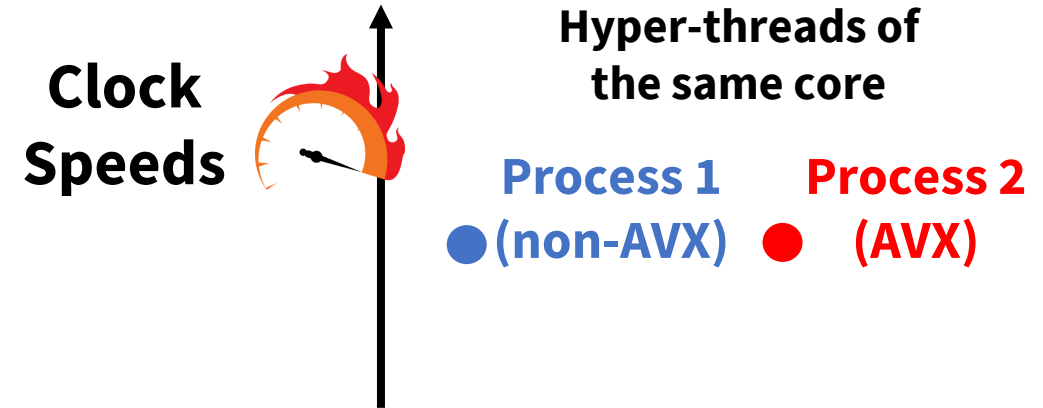
**Better performance  
for certain applications!**

## A: Process scheduling

# Problem: unfairness from AVX instructions



AVX instructions force a CPU core into a lower clock speed



**Other processes executing on the same processor core are affected!**

## A: Process scheduling

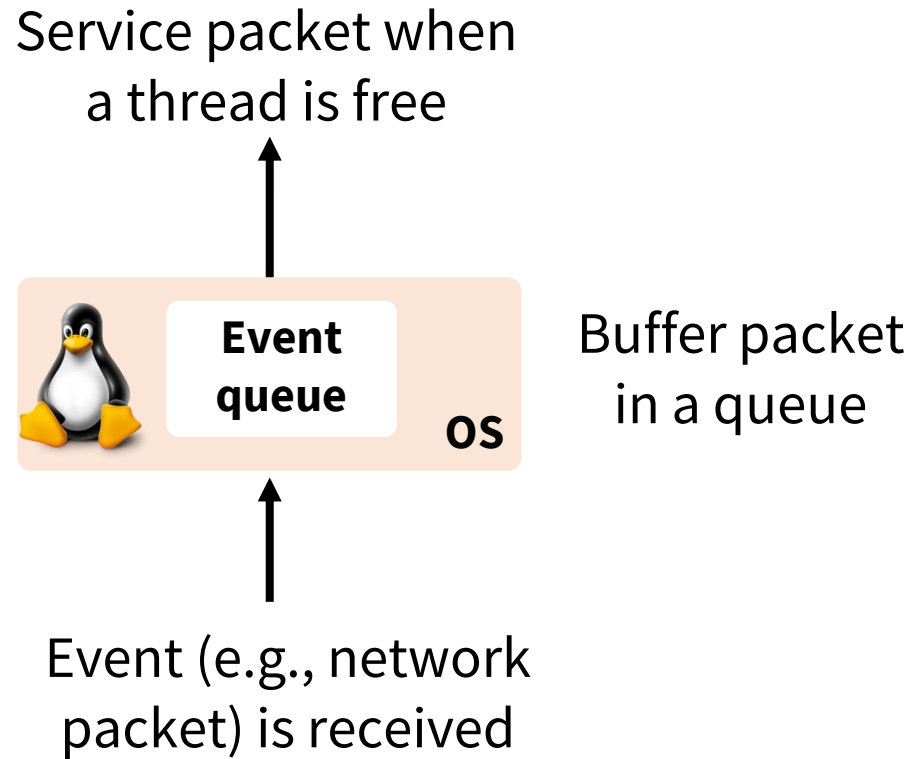
### Paper #2

**“SKQ: Event Scheduling for Optimizing Tail Latency in a Traditional OS Kernel”**  
USENIX ATC 2021



Design a scalable hybrid event queue  
to reduce application latency

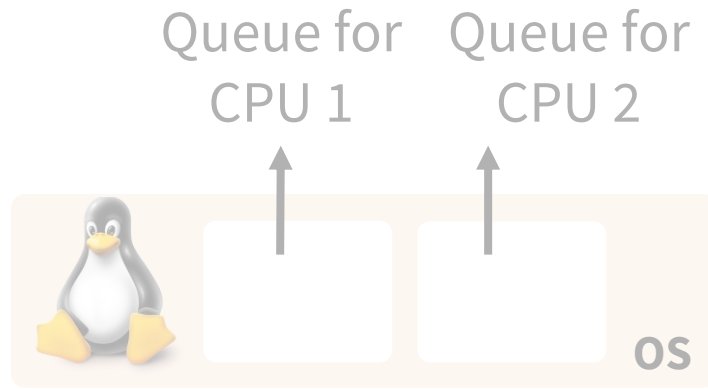
# Background: event queues





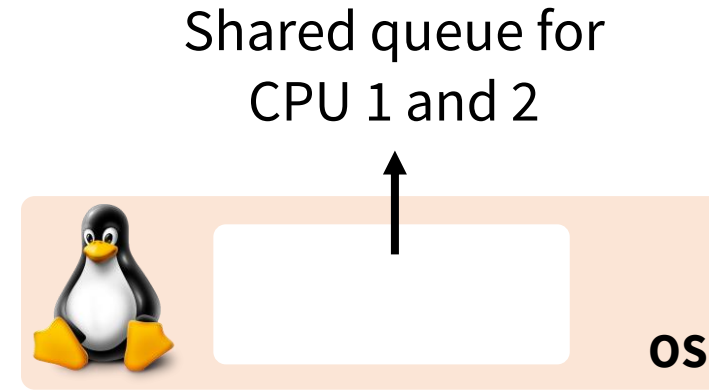
## A: Process scheduling

### Problem: latency in current event queues



**1:1 model**

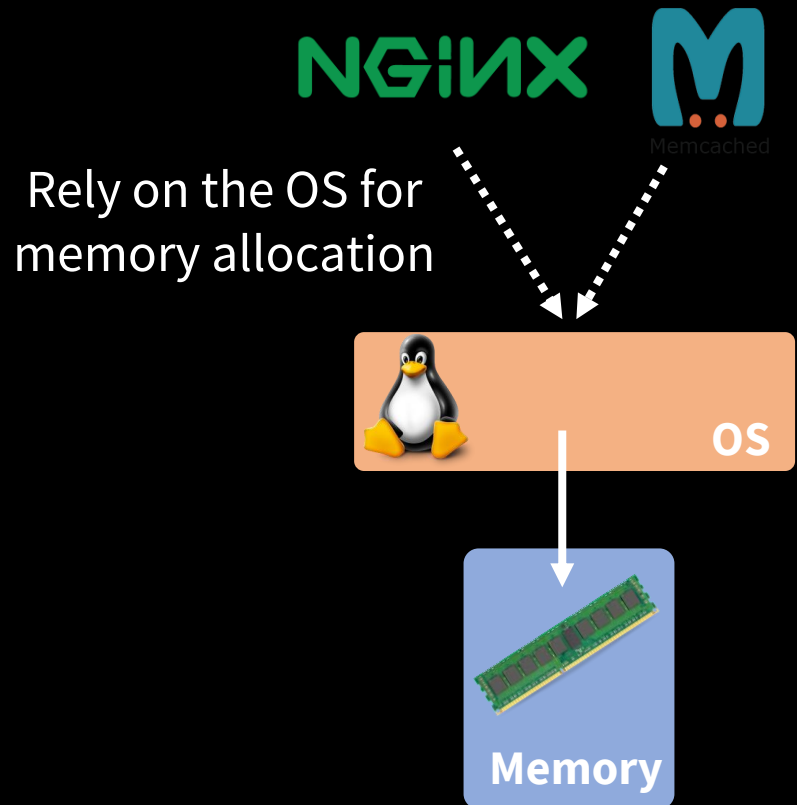
**Under-utilized CPUs because events cannot be migrated easily!**



**1:N model**

**Expensive synchronization between CPUs on each event!**

## B. Memory management



Two relevant concepts:

**Memory conservation:** support many programs with limited memory

**Maximum performance:** reduce the performance impact of memory accesses

## B: Memory management

### Paper #1

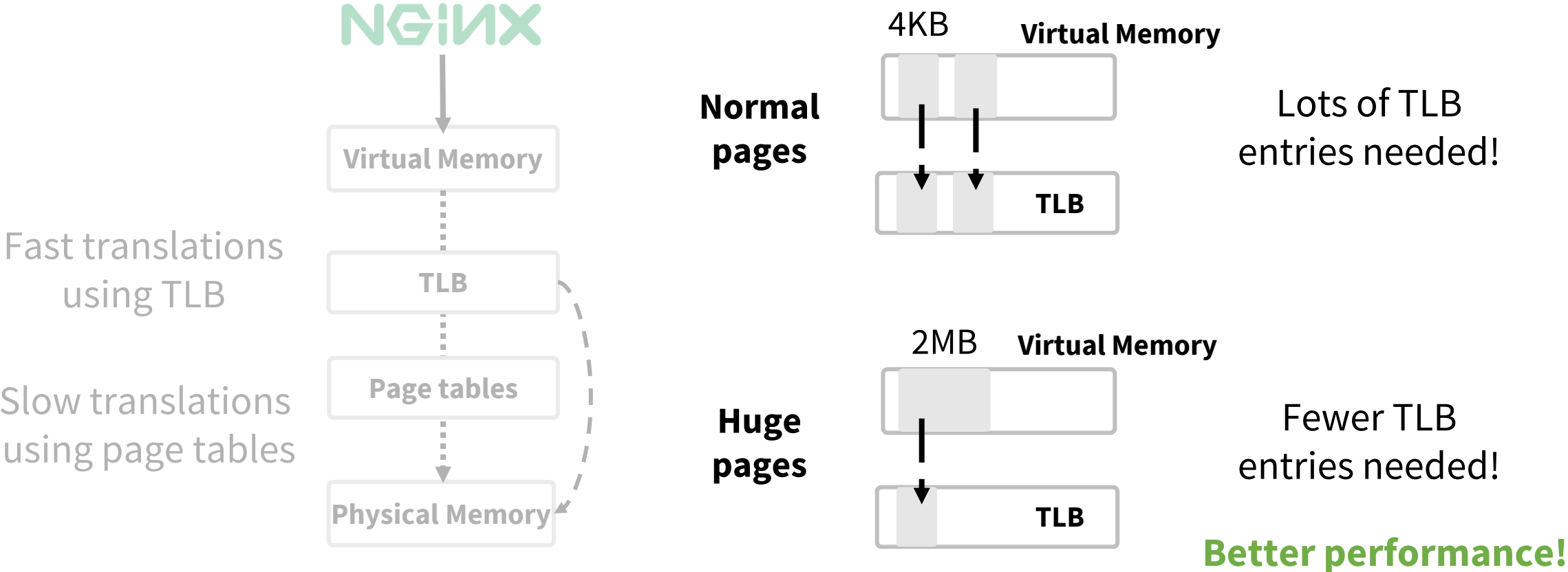
**“Beyond Malloc Efficiency to Fleet Efficiency: A Hugepage-aware Memory Allocator”** USENIX OSDI 2021



Create a memory allocator which helps conserve memory for hugepages

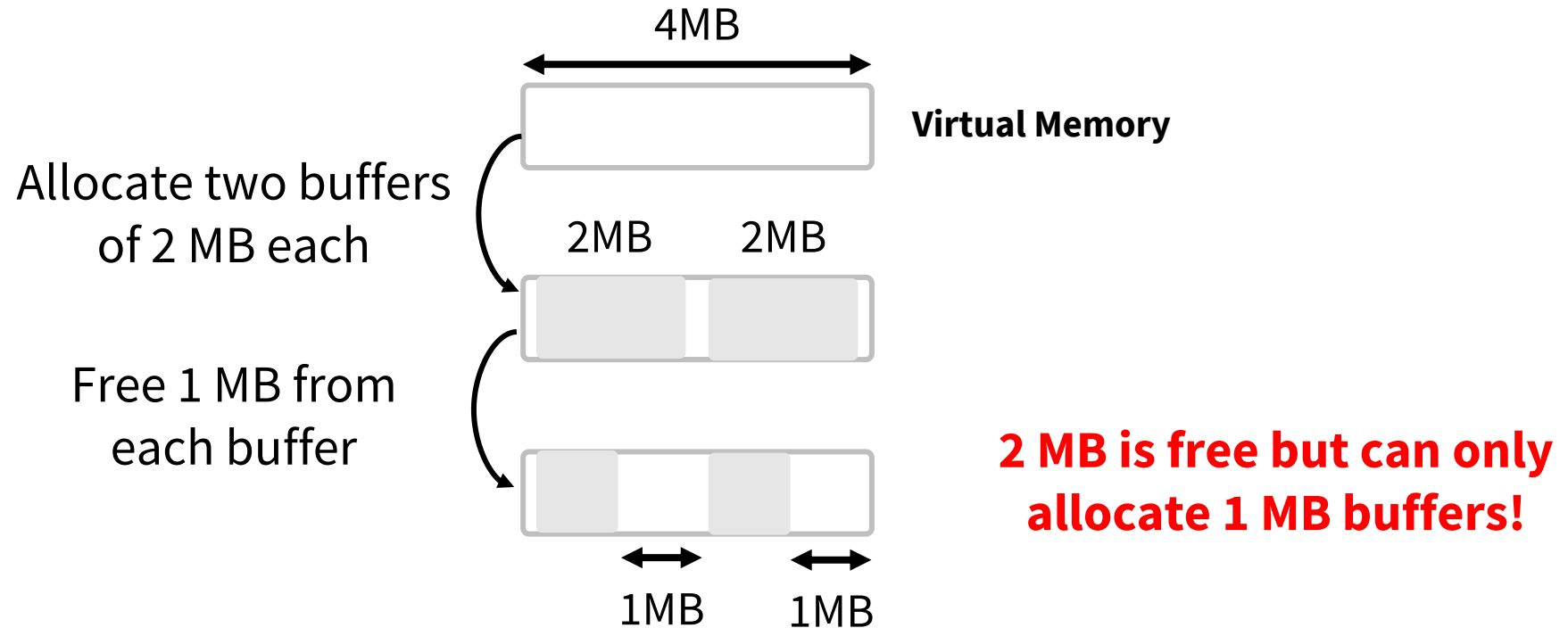
# B: Memory management

## Background: translation and hugepages



## B: Memory management

# Problem: fragmentation from hugepages



## B: Memory management

### Paper #2

**“NrOS: Effective Replication and Sharing in an Operating System”**  
USENIX OSDI 2021



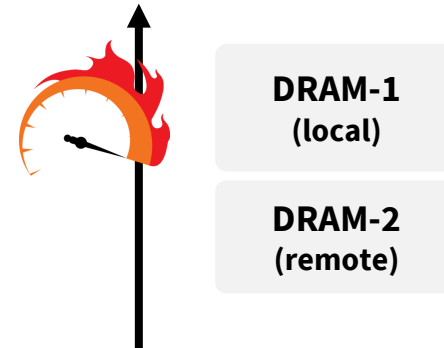
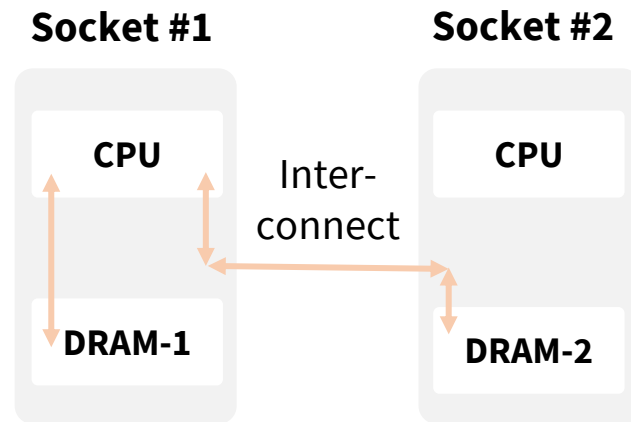
Design an OS kernel which is optimized for server architectures (i.e., NUMA)

## B: Memory management

# Background: non-uniform memory access

Can access both local  
and remote DRAM

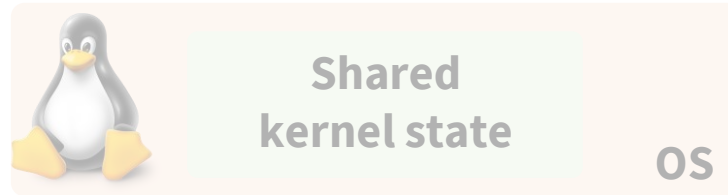
Memory access speed  
depends on locality



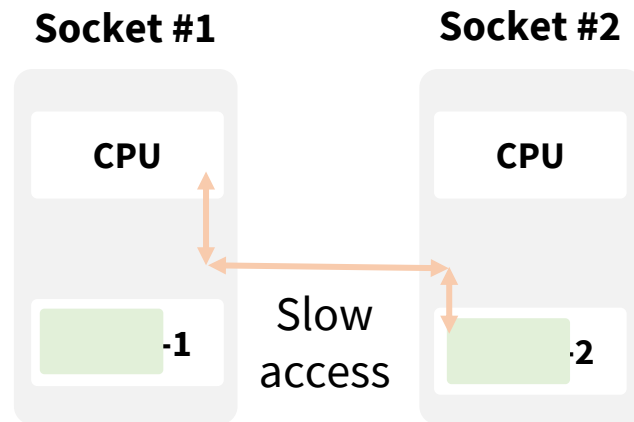
**Non-uniform memory access (NUMA)**

## B: Memory management

# Problem: NUMA slows monolithic kernels



A shared kernel state is accessed by all CPUs



**Frequent accesses to remote DRAM harm performance!**



## B: Memory management

### Paper #3

**“Exploring the Design Space of Page Management for Multi-Tiered Memory Systems”**

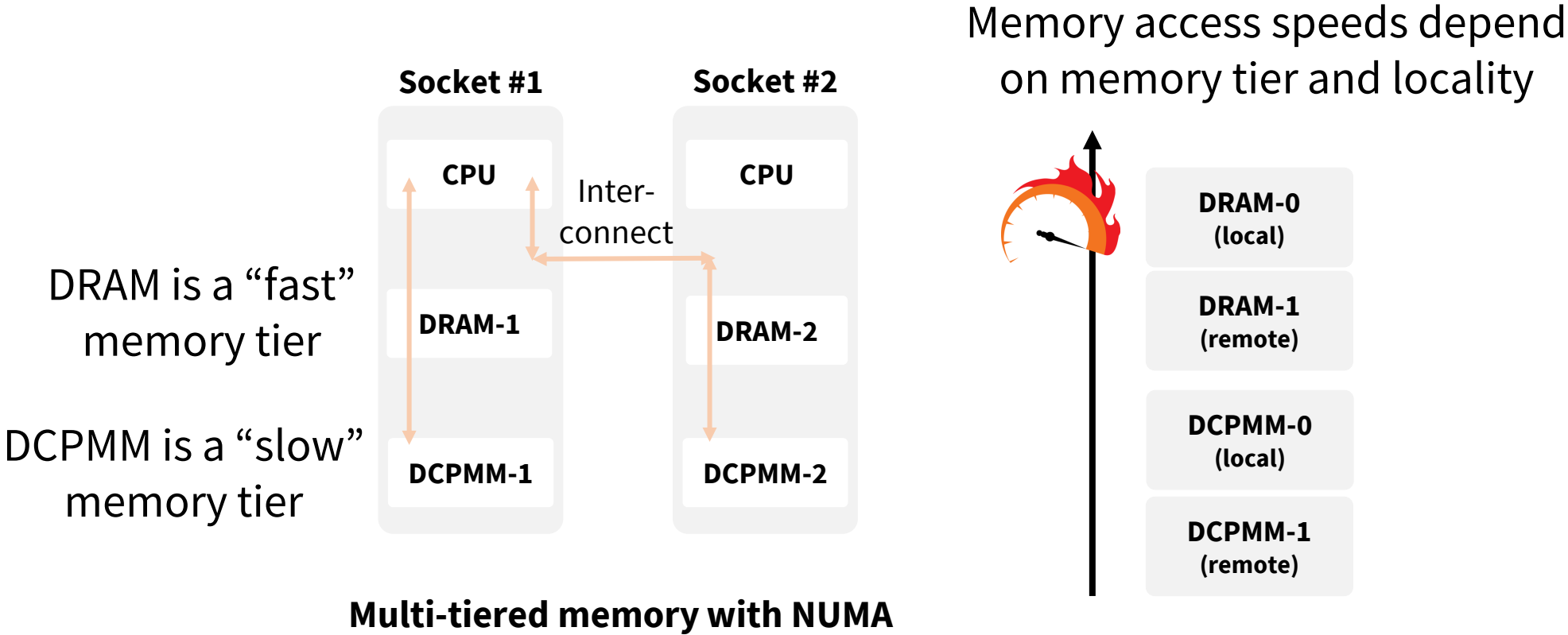
USENIX ATC 2021



Design a memory management scheme  
optimized for server memory architectures

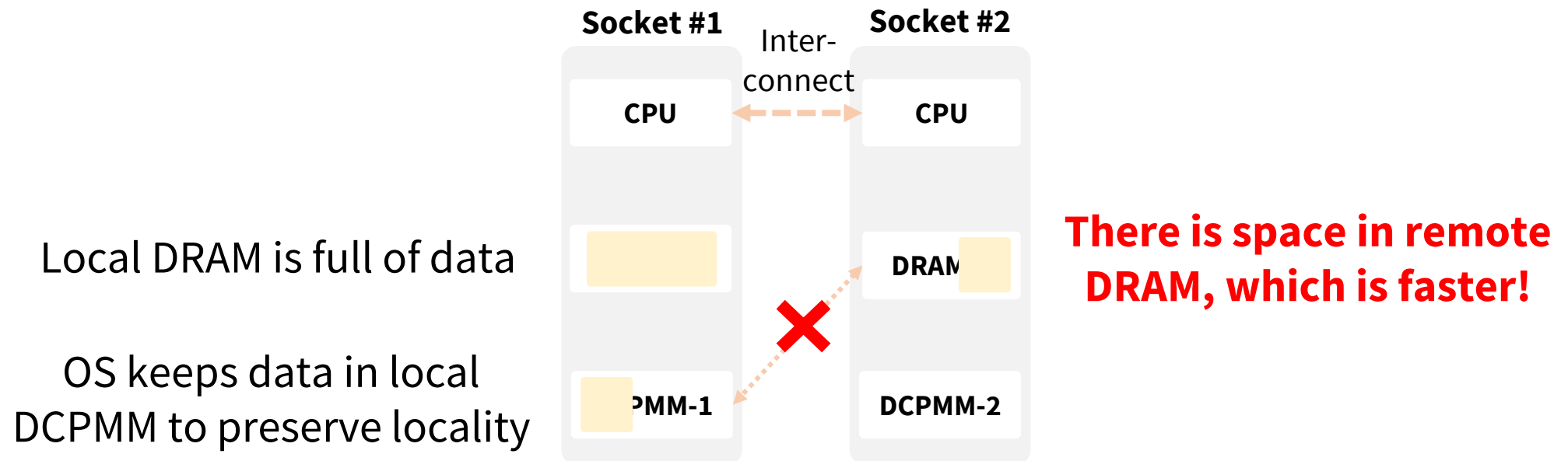
# B: Memory management

## Background: multi-tiered memory (MTM)

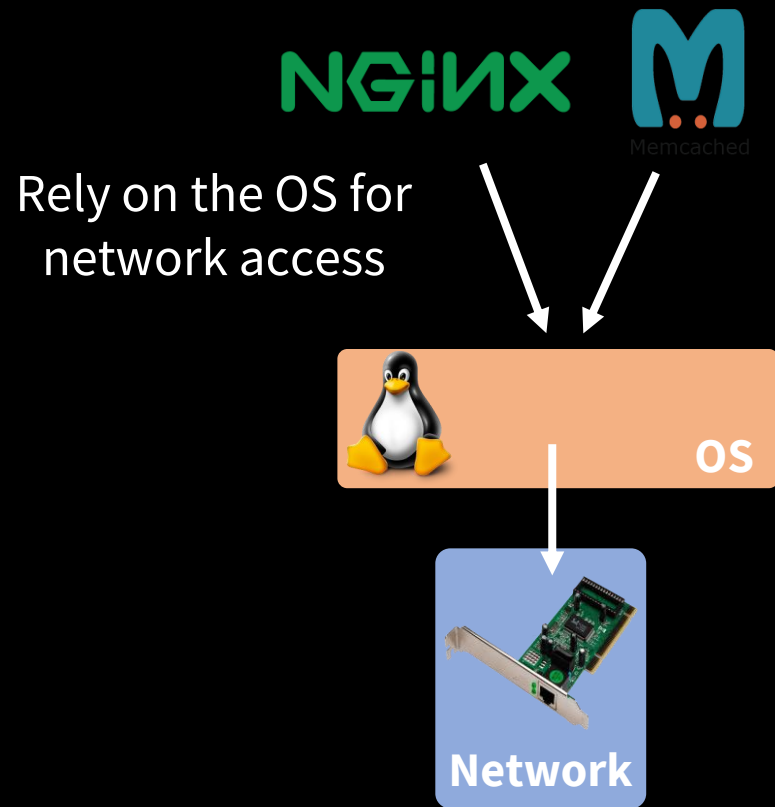


## B: Memory management

### Problem: inefficient MTM usage



# C. Network management



Two concepts are relevant:

1. Network flows are **short** (e.g., web requests) or **long** (e.g., database snapshots)

2. **Short** network flows are more **latency-sensitive** than long ones

## C: Network management

### Paper #1

**“A Linux Kernel Implementation of the Homa Transport Protocol”**

USENIX ATC 2021



Implements Homa, a transport protocol that prioritizes short flows to reduce latencies

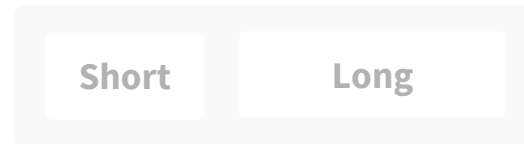
## C: Novel hardware designs

# Background: Homa transport protocol

### TCP:

Sender does not control how packets are handled by network

### Normal router queue

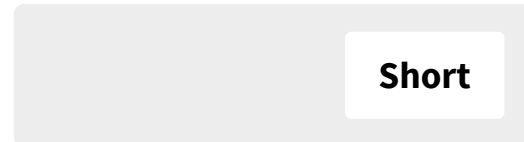


Latency-sensitive short flows could be delayed!

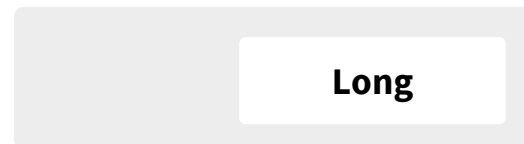
### Homa:

Sender configures packets to use priority queues in routers

### High priority router queue



Short flows are prioritized, and latency is reduced!



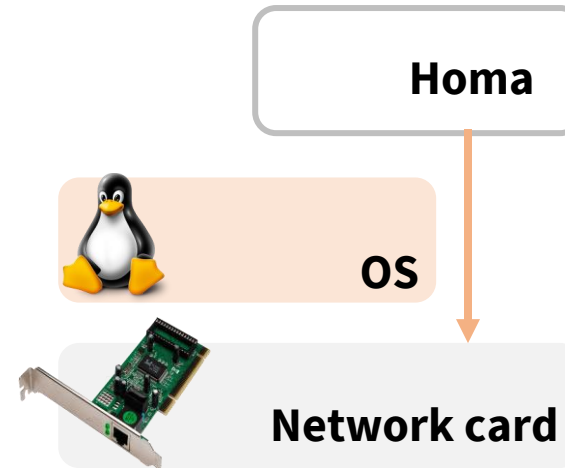
### Low priority router queue

## C: Novel hardware designs

# Problem: Homa is not tested practically



**Simulations through software  
network simulators**



**User-level implementation  
with kernel-bypass**

## D: Novel hardware designs



Improve **security** of  
remote computation



Improve **performance** in  
emerging real-world scenarios



## D: Novel hardware designs

### Paper #1

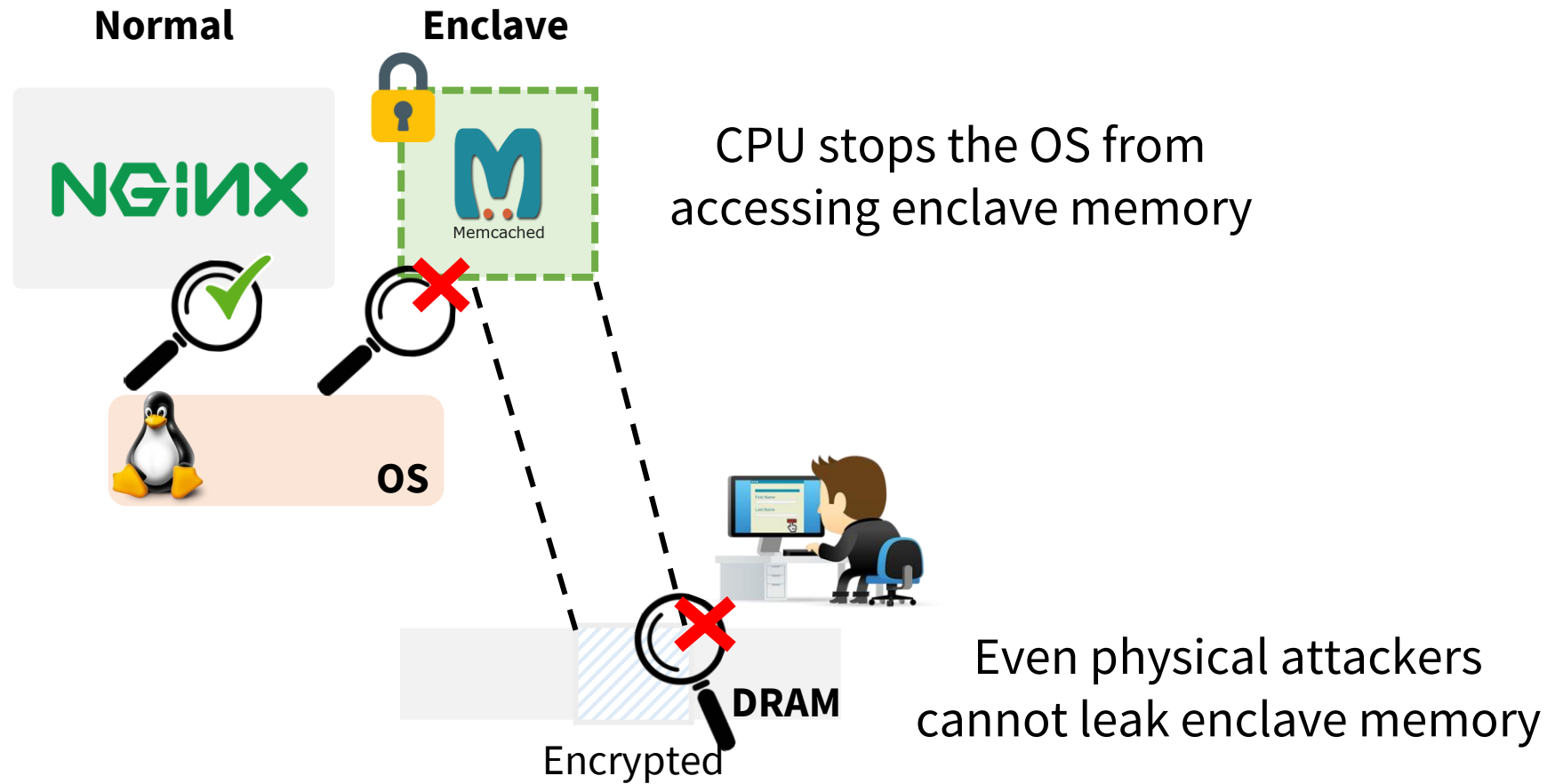
**“Scalable Memory Protection in the Penglai Enclave”**  
USENIX OSDI 2021



A new trusted hardware that is more scalable and robust than previous solutions

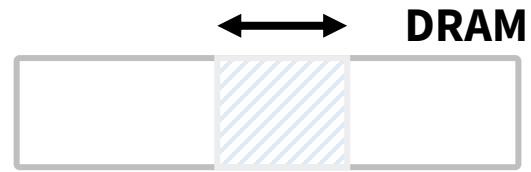
## D: Novel hardware designs

# Background: enclaves



## D: Novel hardware designs

# Problem: existing enclaves do not scale



**Limited enclave memory  
(e.g., 128-256 MB)**



**Slow startup speeds  
(e.g., many seconds)**

## D: Novel hardware designs

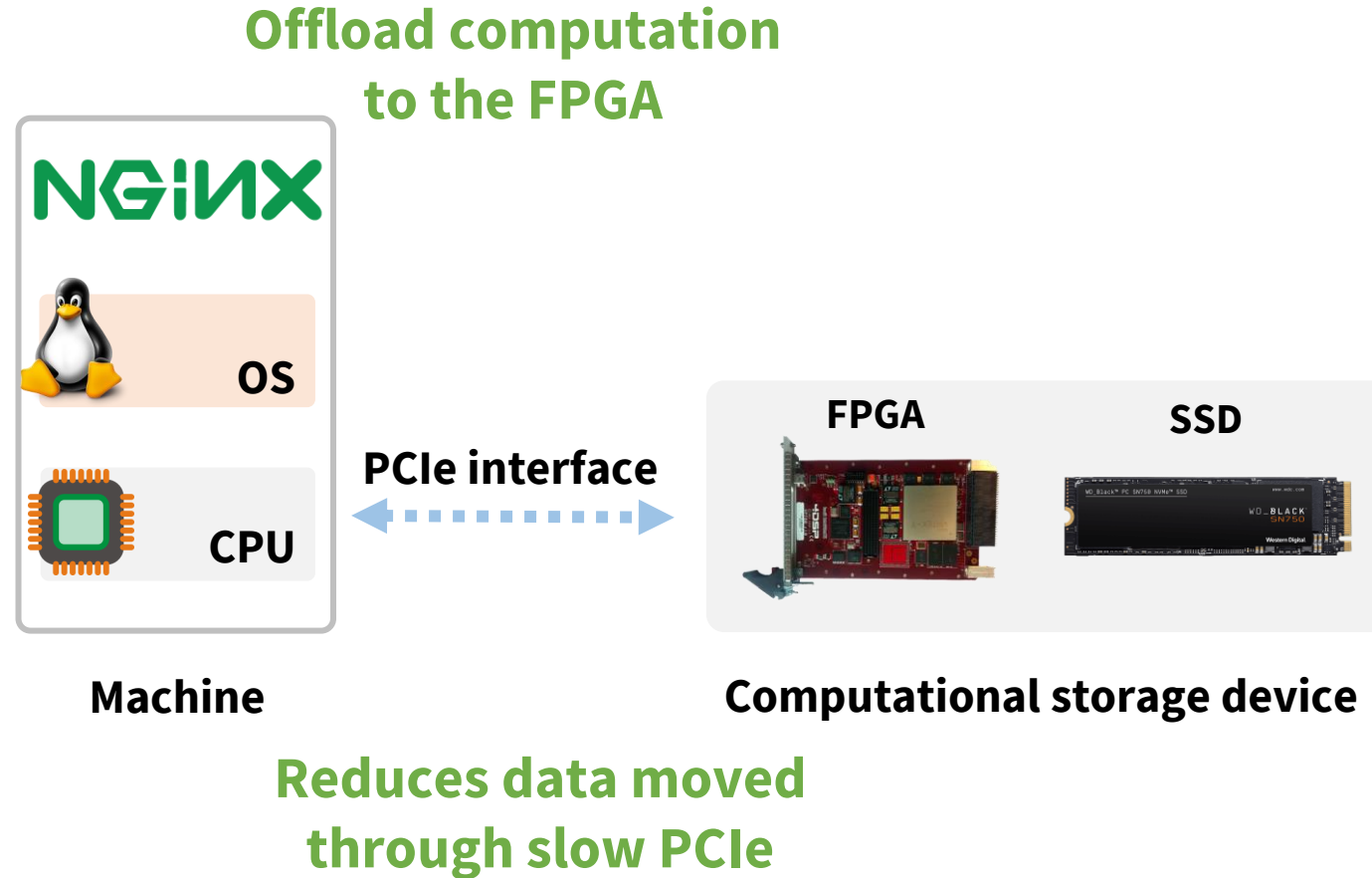
### Paper #2

**“A Fast and Flexible Hardware-based Virtualization Mechanism for Computational Storage Devices”** USENIX ATC 2021



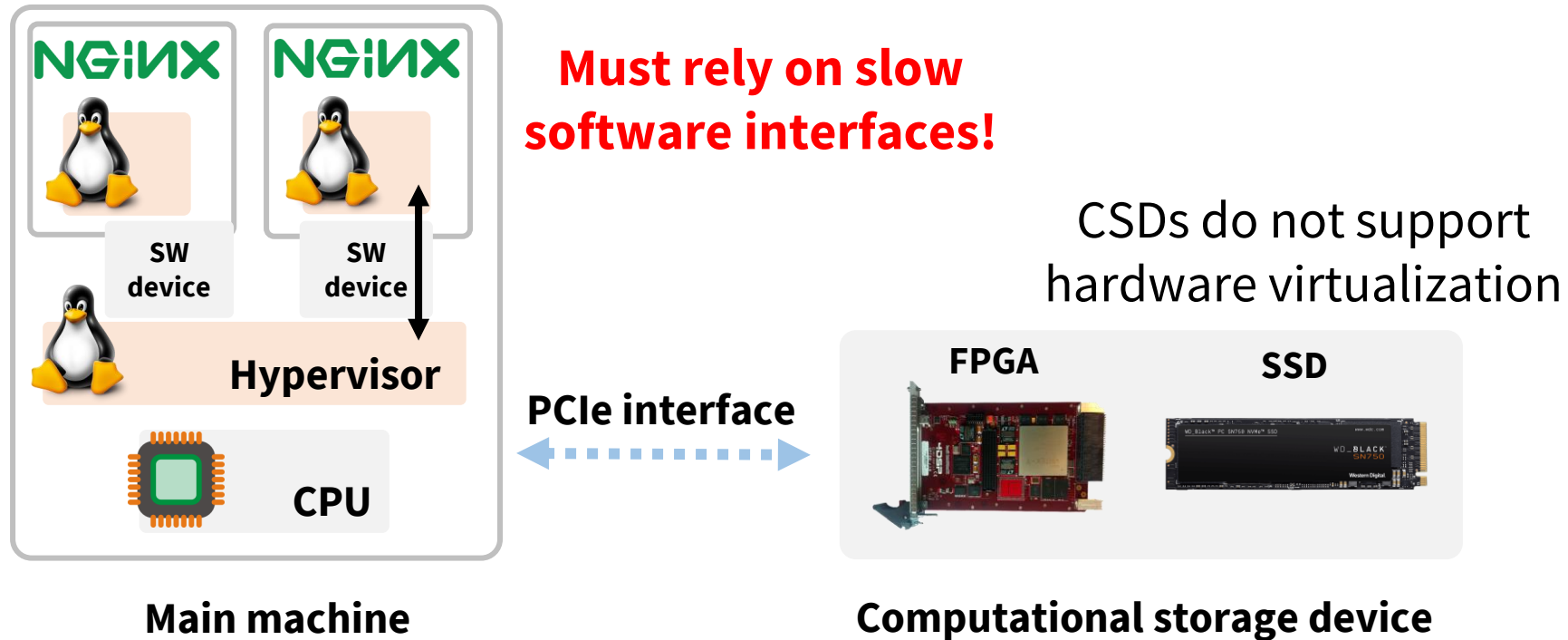
A new hardware to virtualize computational storage devices without relying on slow software

# Background: computational storage



## D: Novel hardware designs

# Problem: CSD virtualization is inefficient



## D: Novel hardware designs

### Paper #3

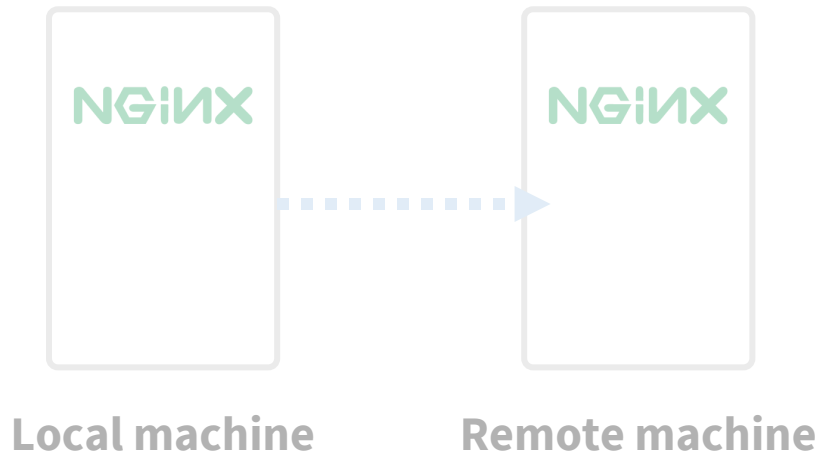
**“The nanoPU: A Nanosecond Network Stack for Datacenters”**  
USENIX OSDI 2021



A new hardware to realize nanosecond response times for very short-lived datacenter requests

## D: Novel hardware designs

# Background: remote procedure calls



**Few nano-second (ns) –  
1 micro-second (us)**

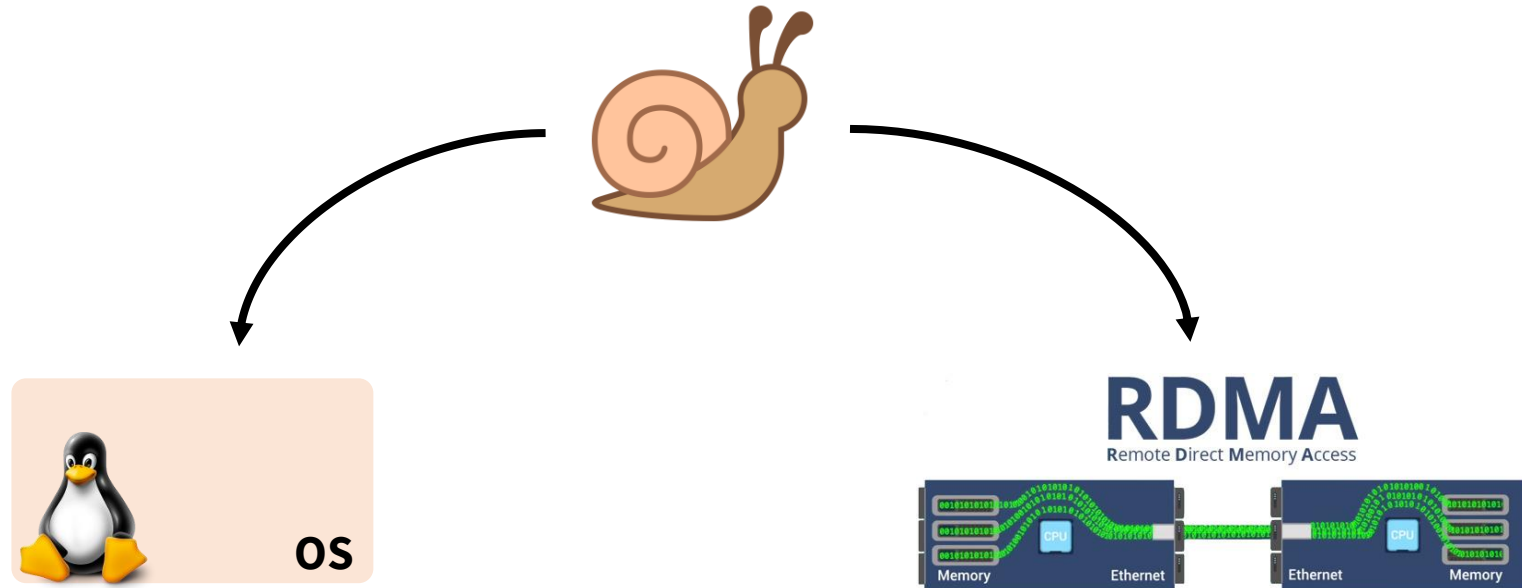
Can call a function (or procedure) on a remote machine as if it was a local function

RPCs in datacenters are becoming more and more short-lived



## D: Novel hardware designs

# Problem: existing network stacks are slow



**Software stacks (e.g., Linux)  
take ~100s of microseconds**

**Hardware stacks (e.g., RDMA)  
take 1-2 microseconds**

# Conclusion

## OS and hardware

Thursday, July 15 @ 7 AM PDT



15th USENIX Symposium on  
Operating Systems Design and Implementation

JULY 14-16, 2021  
VIRTUAL EVENT

Sponsored by USENIX in cooperation with ACM SIGOPS

[www.usenix.org/osdi21](http://www.usenix.org/osdi21)

## My tail never has any latency

Friday, July 16 @ 8.30 AM PDT

# USENIX ATC '21

2021 USENIX  
Annual Technical Conference

JULY 14-16, 2021  
VIRTUAL EVENT

[www.usenix.org/atc21](http://www.usenix.org/atc21)