# Scale-Out NUMA

ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE

PAISA PARALLEL SYSTEMS ARCHITECTURE LAB

DCSL

*Stanko Novakovic, Alexandros Daglis, Boris Grot\*, Edouard Bugnion, Babak Falsafi*
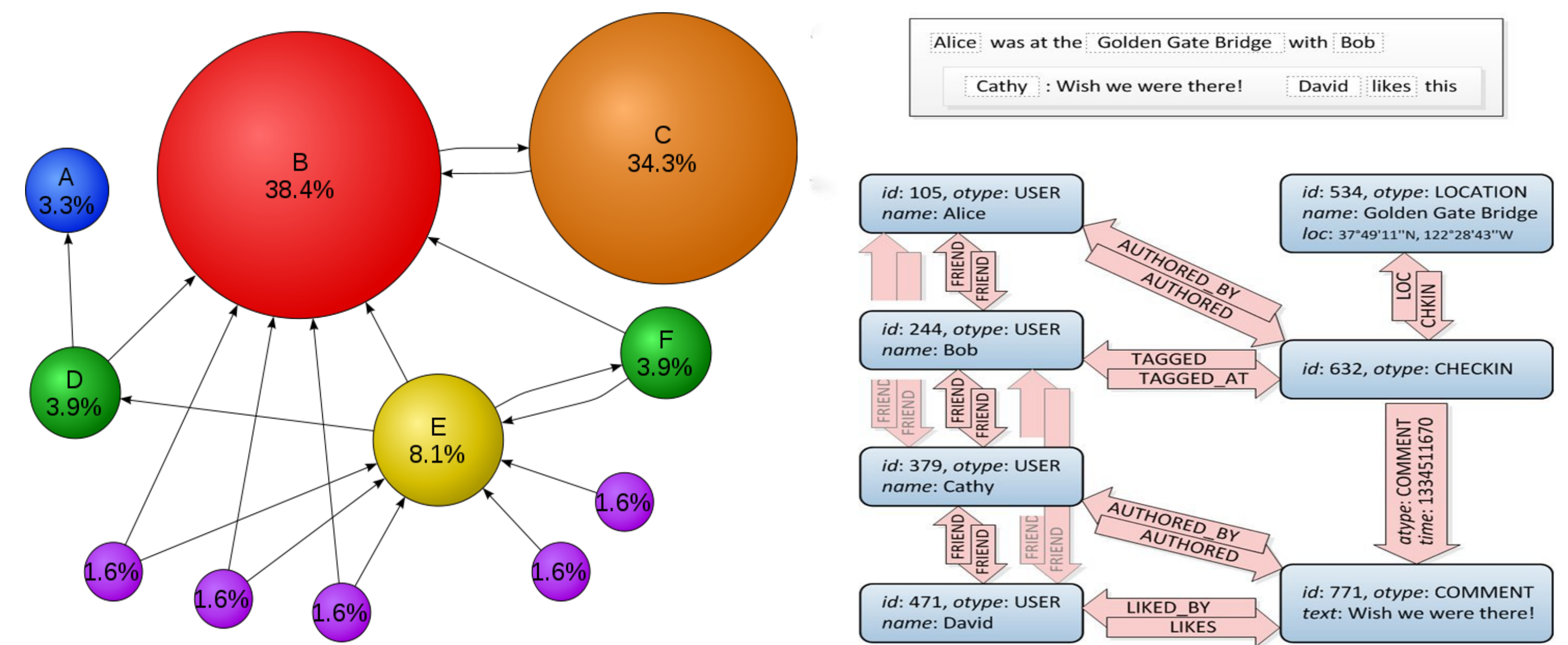
EPFL, University of Edinburgh\*

## 1. Large-scale datacenter applications

**Big-data analytics and data serving**

**Common properties:**
- Large datasets, many nodes
- Frequent accesses to non-local data
- Very little processing per query/algorithm iteration



**Most DC applications are network-bound**

## 2. Existing approaches ill-suited

**Shared memory (ccNUMA)**
  + Low latency to remote data
  - Limited scalability, high cost, single failure domain

**Distributed memory using TCP/IP over Ethernet**
  + High scalability using commodity parts
  - High remote access latencies (up to 1000x of local)

**Distributed memory using RDMA over InfiniBand**
  + High scalability, low latency
  - Remote access latency memory still high (>10x of local)

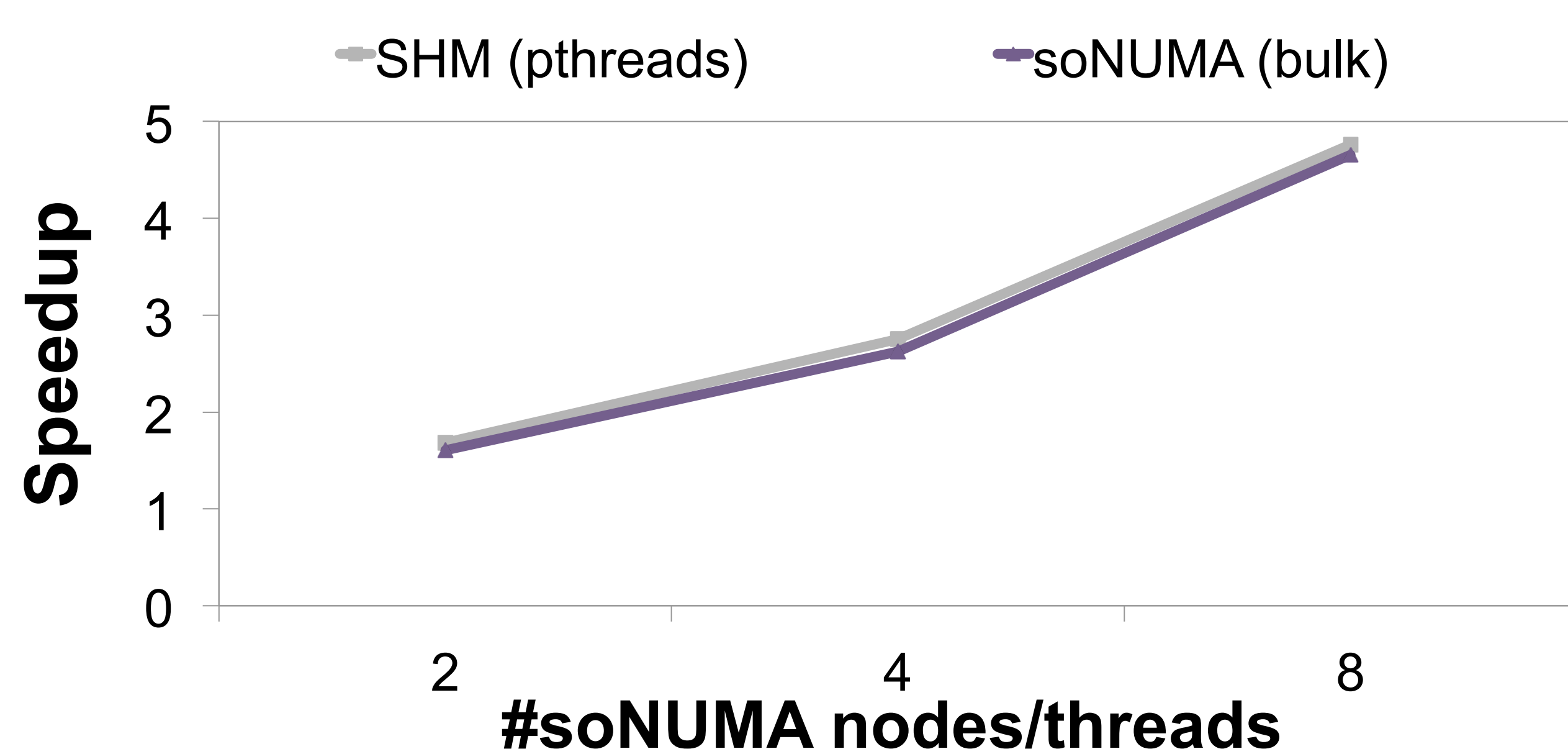**Deep network stacks, PCI/DMA limit performance**

## 3. Our proposal: Scale-Out NUMA

**Rack-scale system based on NUMA transport**
- Reliable and wide interconnect
- Integrated (locally) cache-coherent Remote MC
- Direct access via memory-mapped queue-pairs (QP)



**Remote access latency of 300ns, DDR rate, scalable**

## 4. Rack-scale graph processing

**Bulk Synchronous Parallel processing**
- Iterative computation
- Servers exchange graph updates across iterations
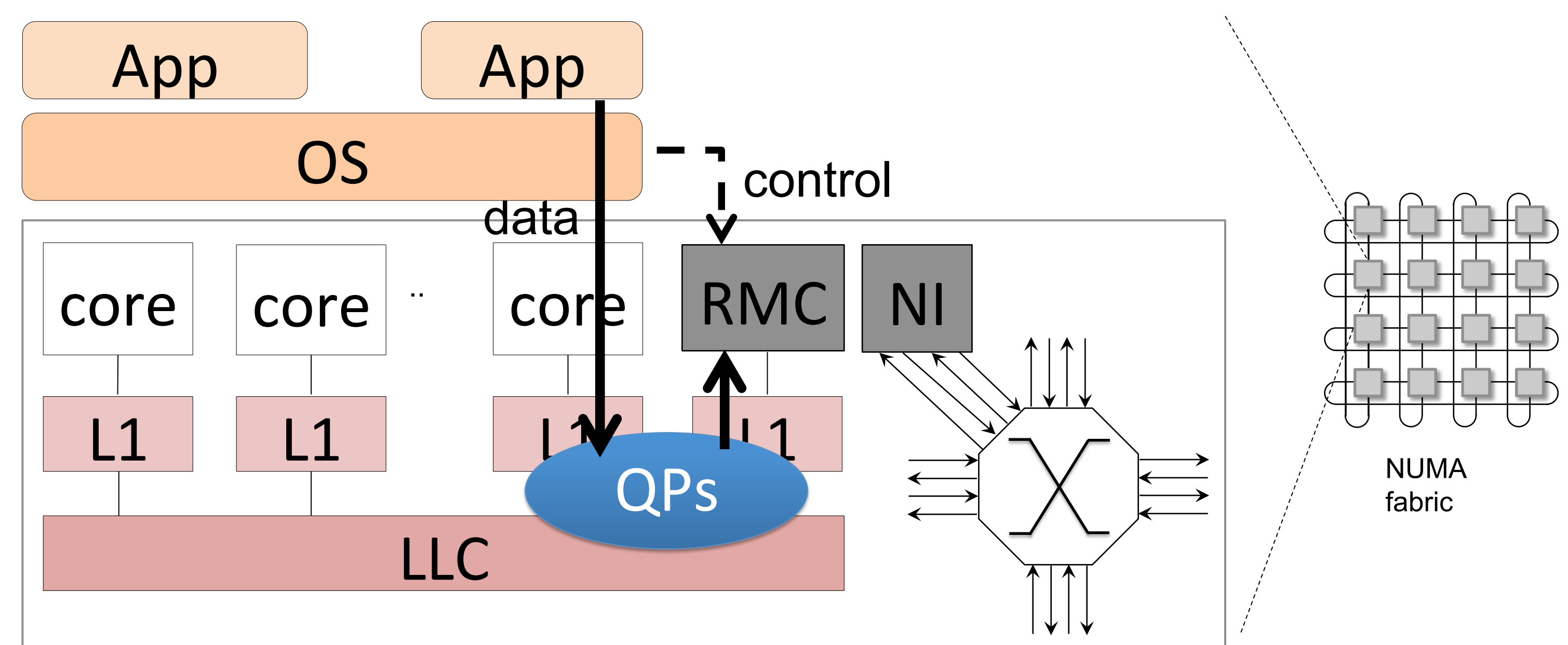
**PageRank on Twitter graph study**
- soNUMA (bulk) → BSP implementation on soNUMA
- SHM (pthreads) → shared-memory implementation



## 5. Rack-out data serving

**Rack-out: shard data at rack-scale granularity**
- Skewed access patters create hotspots in scale-out
- Group servers into racks → more compute/network
  → Deliver higher throughput w/o violating SLA



Super-shard