

# Node-Based Memory Management for Scalable NUMA Architectures

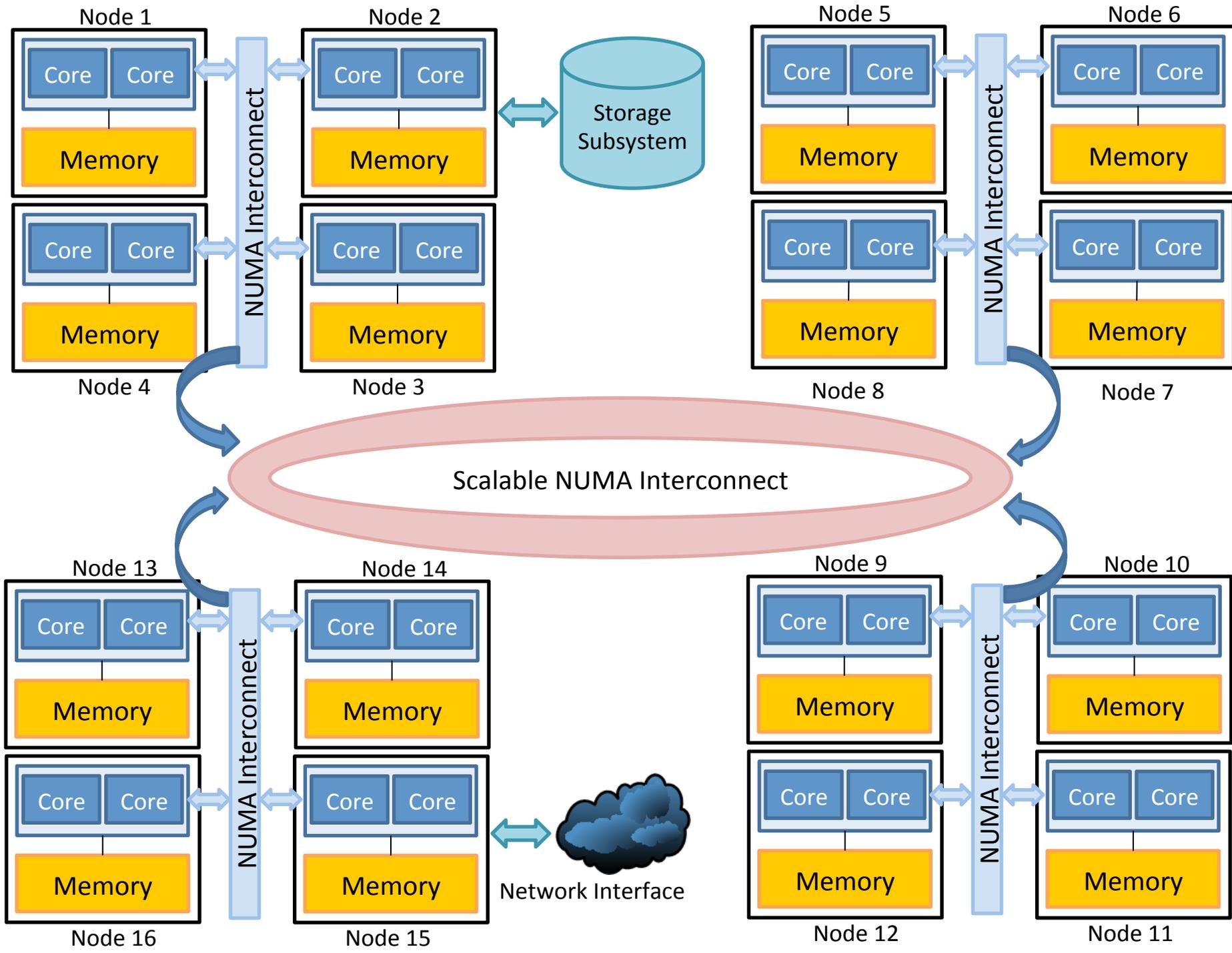
International Workshop on Runtime and Operating Systems for Supercomputers (ROSS 2012)

**Stefan Lankes<sup>1</sup>, Thomas Roehl<sup>2</sup>, Christian Terboven<sup>2</sup>, Thomas Bemmerl<sup>1</sup>**

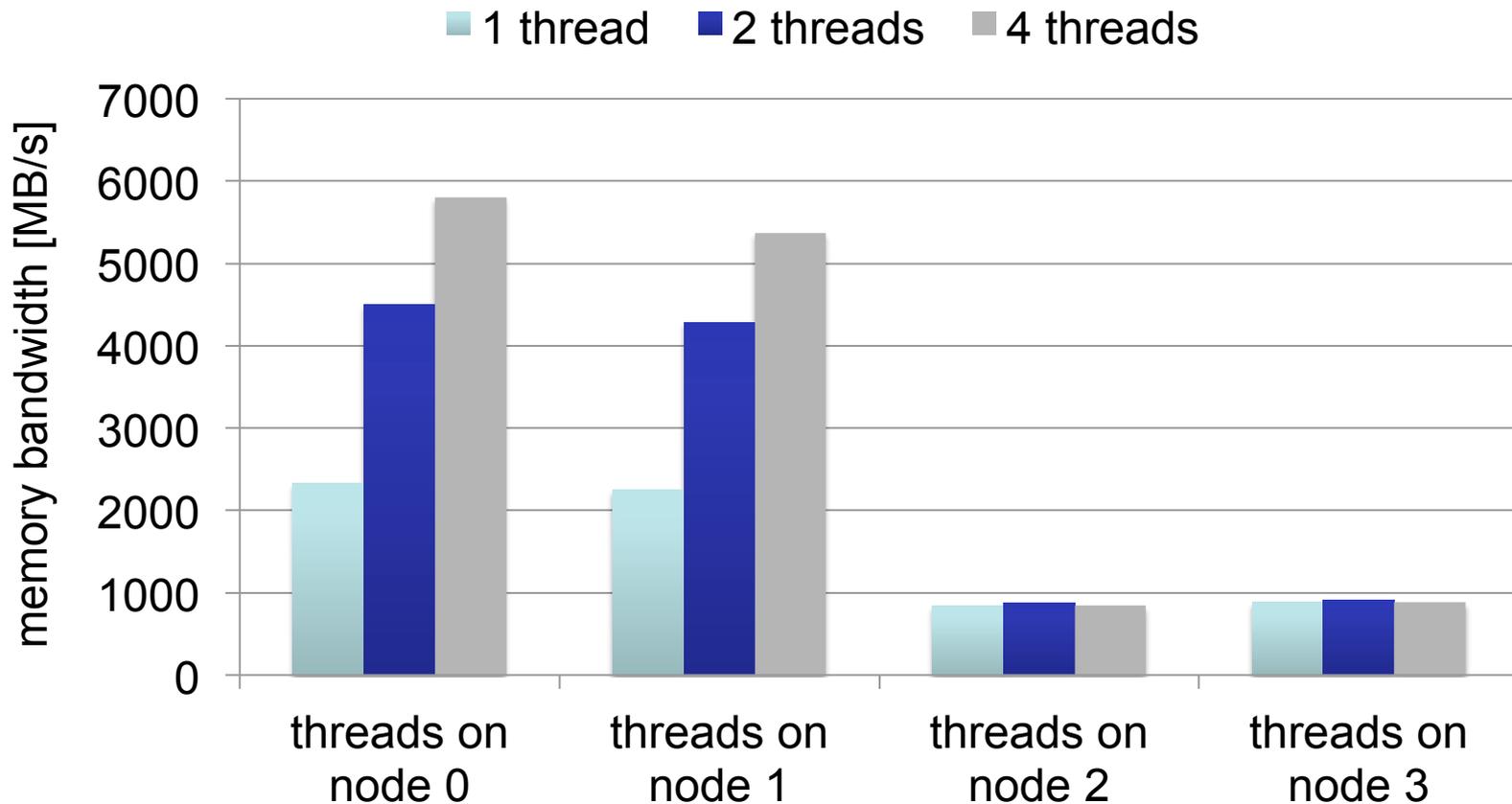
<sup>1</sup>[Chair for Operating Systems](#), RWTH Aachen University

<sup>2</sup>[Center for Computing and Communication](#), RWTH Aachen University

- Motivation
- Illustration of a common memory management
- Design of the node-based memory management
- Critical analysis
- Future prospects
- Benchmark results
- Conclusions and outlook



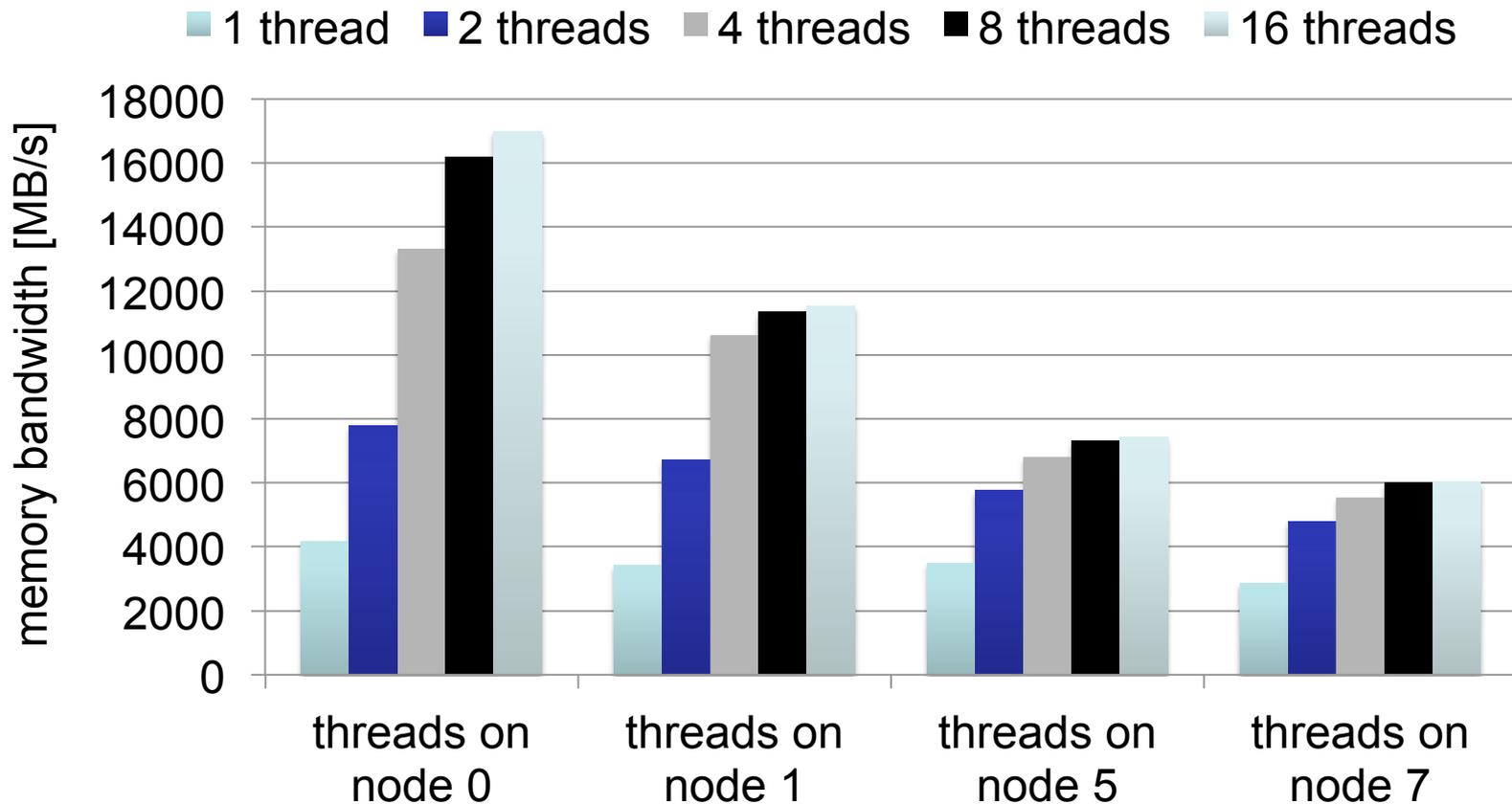
# Performance Characteristics (NumaScale-Cluster)



- 2 systems with 2 AMD QuadCores of type 8378 combined via NumaConnect
- all data on node 0



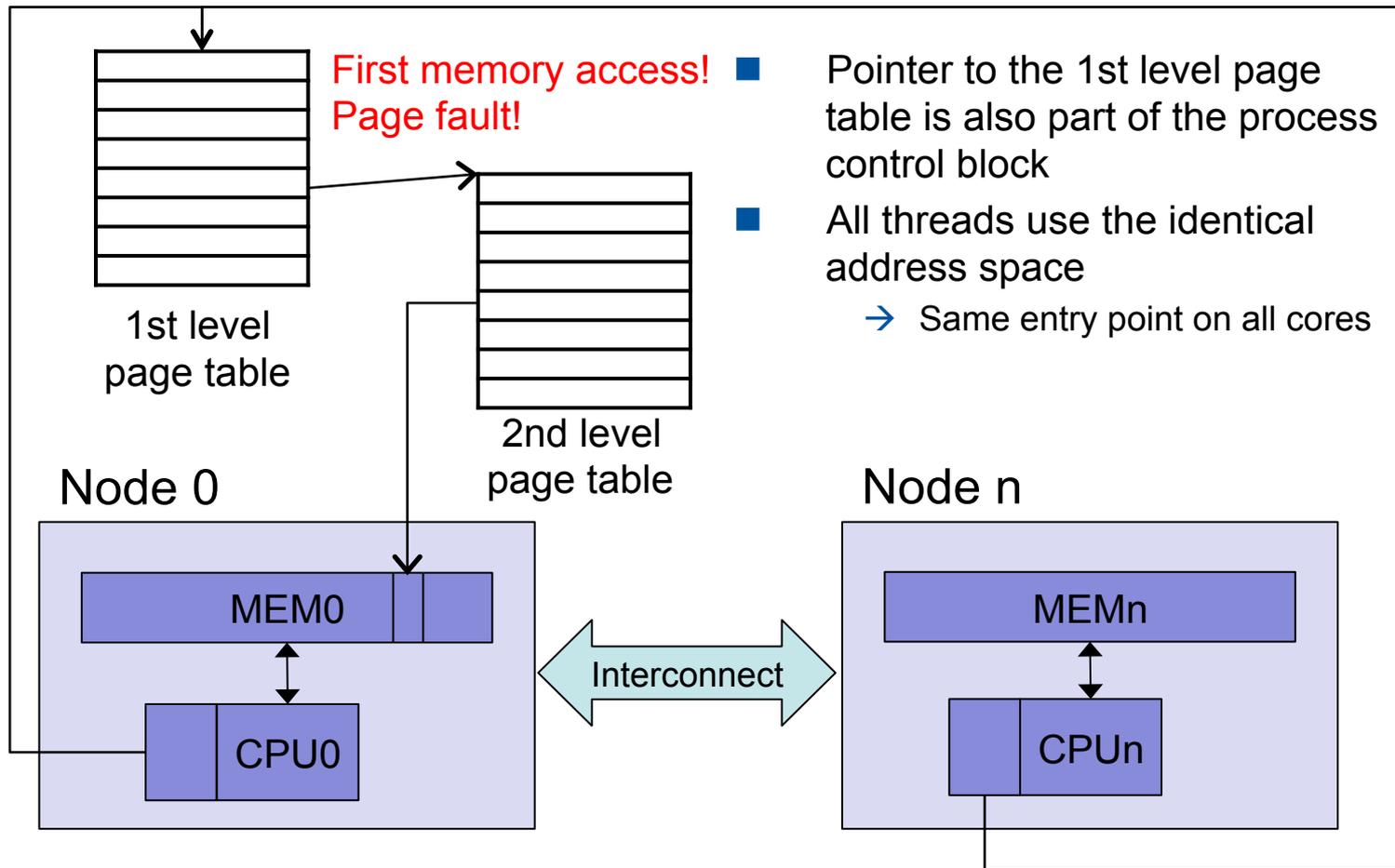
# Performance Characteristics (Westmere-EX)



- 8 Intel Xeon CPU E7-8850 (Westmere-EX)
- 8 \* 10 Cores / 8 \* 20 Cores via HyperThreading
- all data on node 0

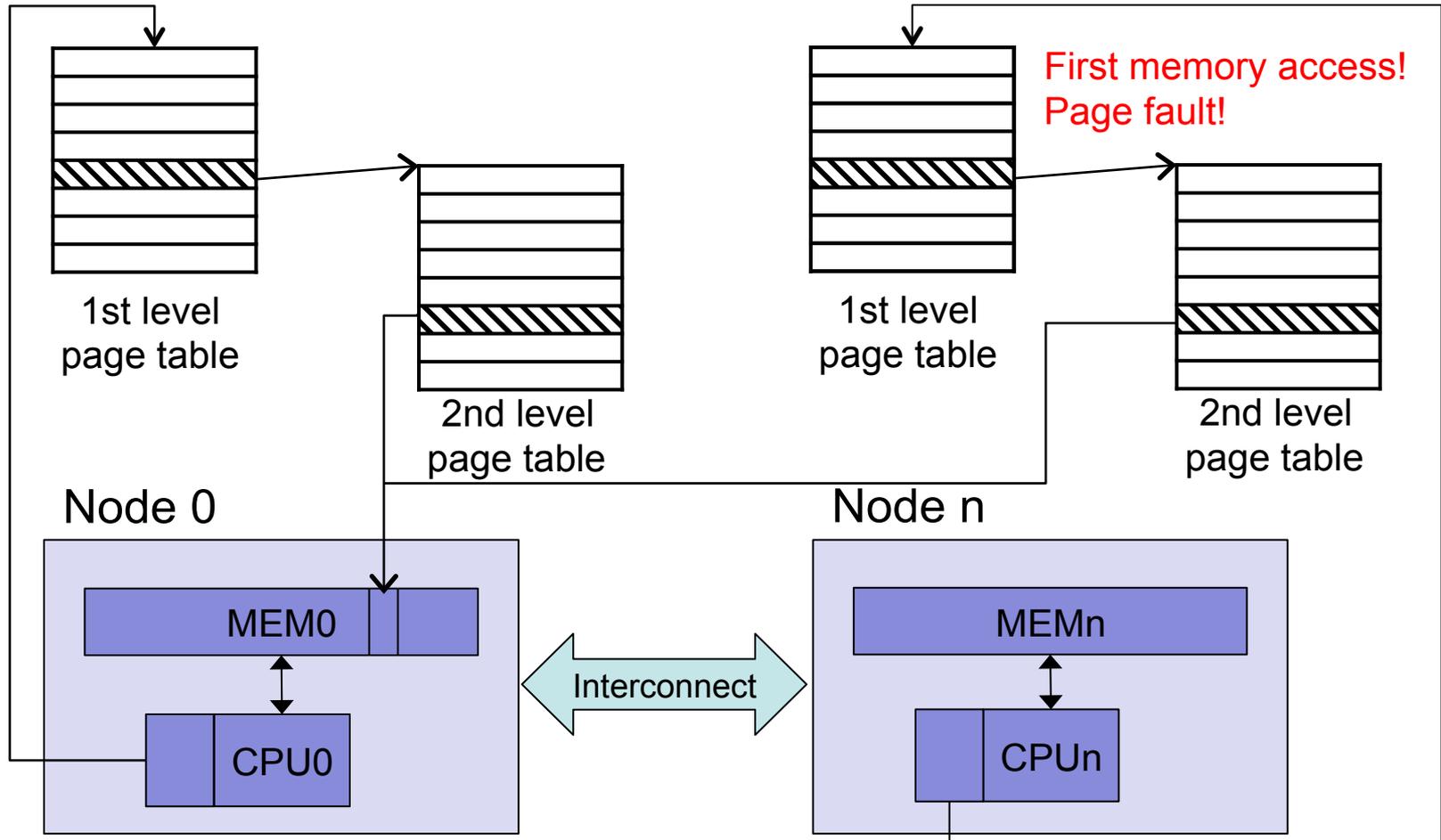
# Common Memory Management

## Process/thread creation



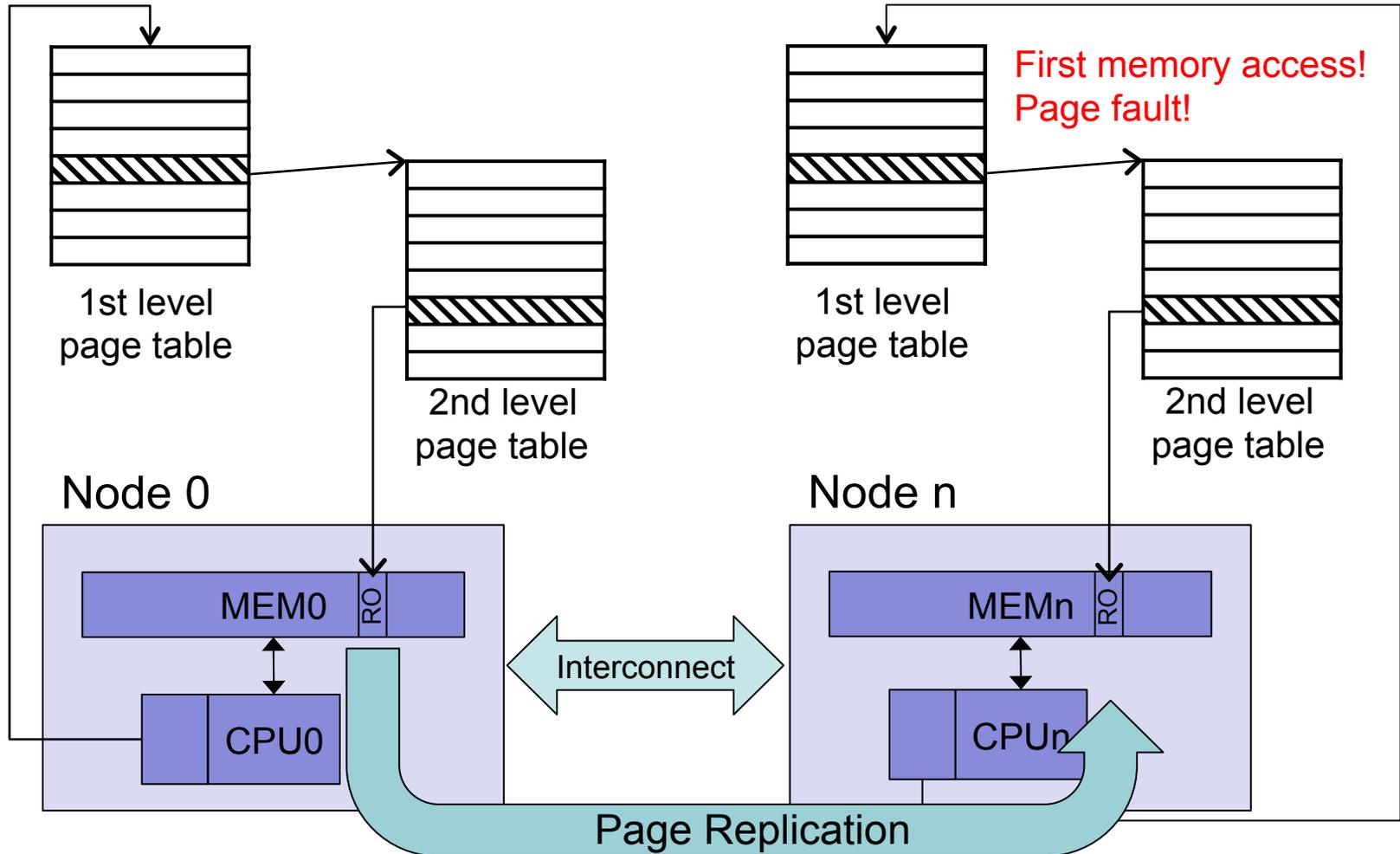
# Page Table per Node

## Basic idea



# Page Table per Node

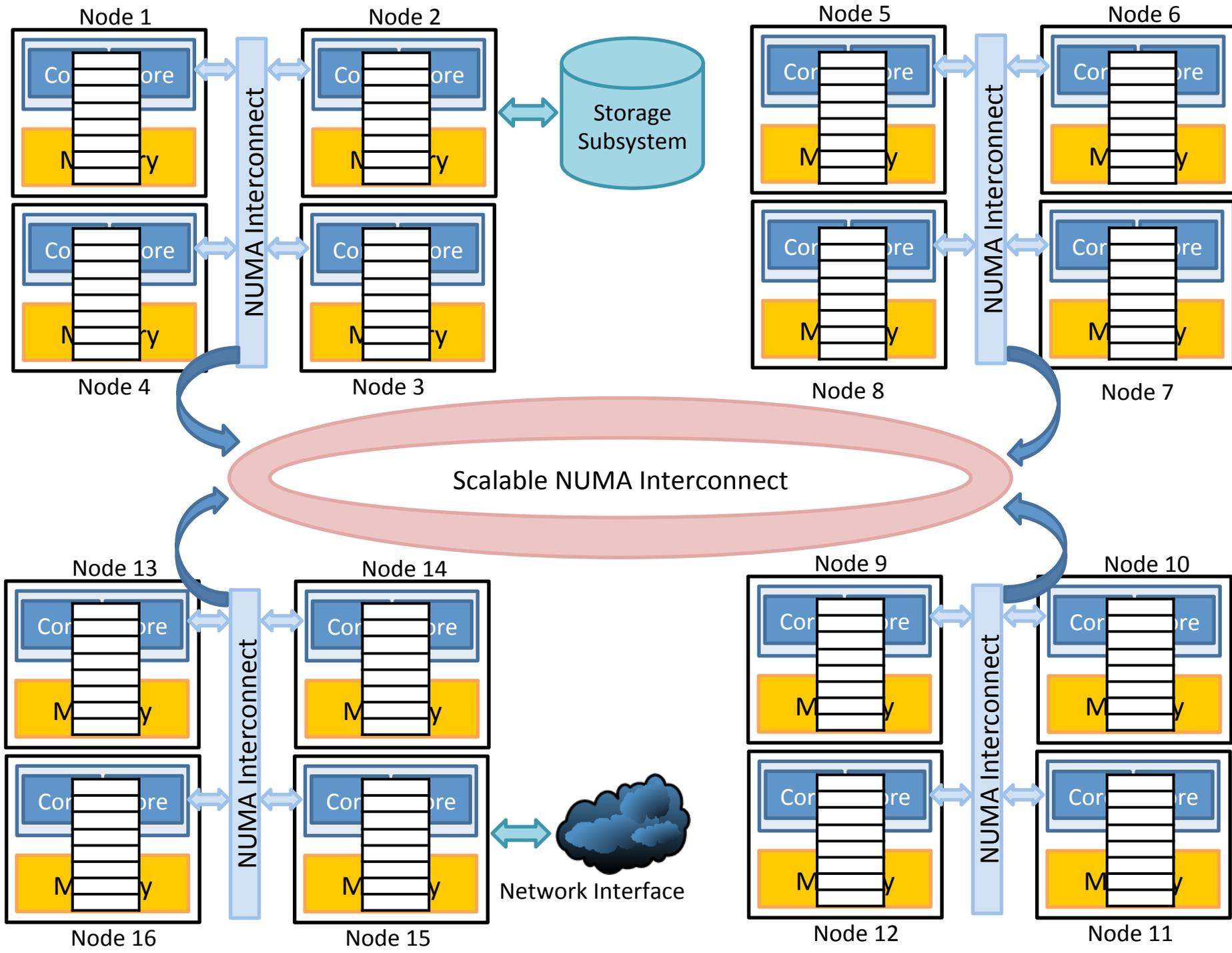
## Replication of read-only regions



- Pro:
  - Reflecting actual hardware at mapping layer
  - After duplication only accesses to local memory
  - Easy preparation of applications to use `mprotect()`
  
- Contra:
  - Memory overhead
    - » One page table per NUMA node
    - » Duplicated pages
  - Replication time
  - Searching for mappings at all NUMA nodes  
(page fault, `mprotect()`, `free()`)

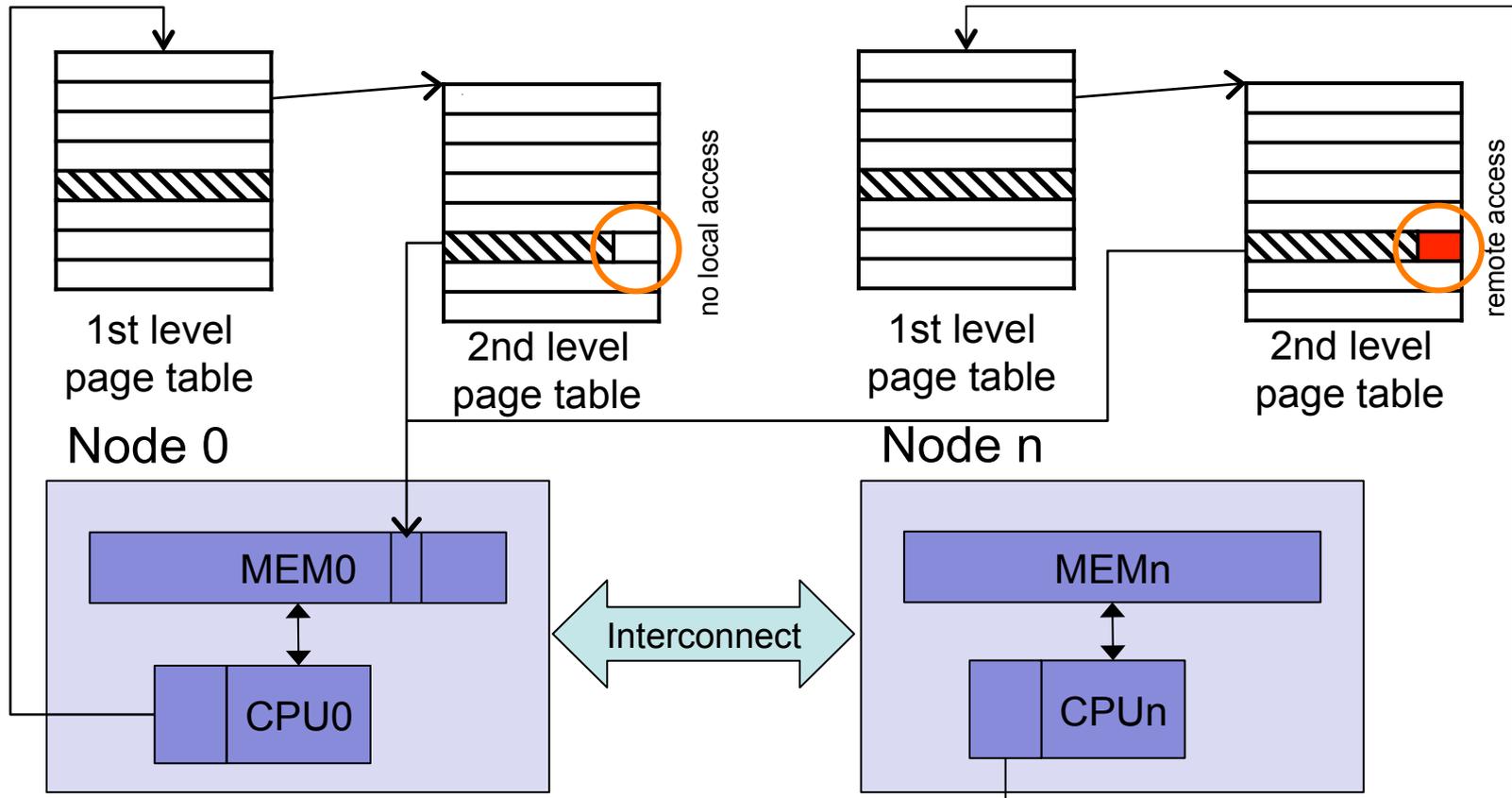
# Avoid PGT-Traversal at Mapping Search

- Current Approach
  - Searching for mappings at all NUMA nodes
  - On which node should we start?
- Under development
  - Use node-distance based search
    - » Does not guarantee less work
  - Add new management structure
    - » Derived page table stores virtual address-to-nodemask mappings
    - » Needs 2 page table traversals per search,
    - » First resolve location, then address
    - » Increases memory footprint



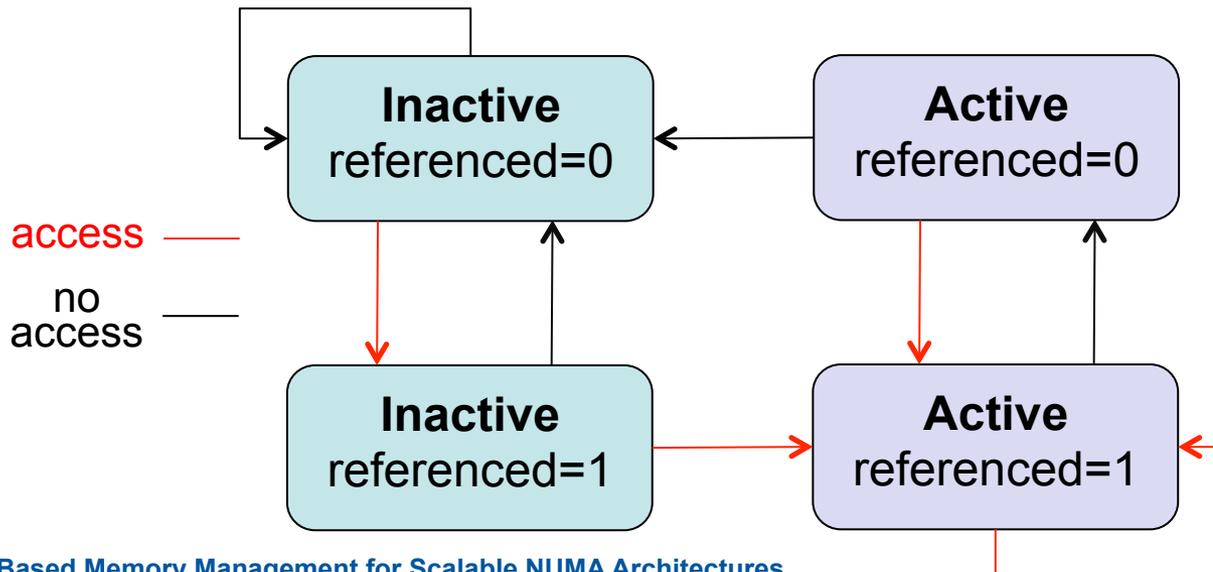
# Detection of Performance Issues

- Page tables include access/dirty bits to record memory accesses.
  - Usable to detect performance issues?



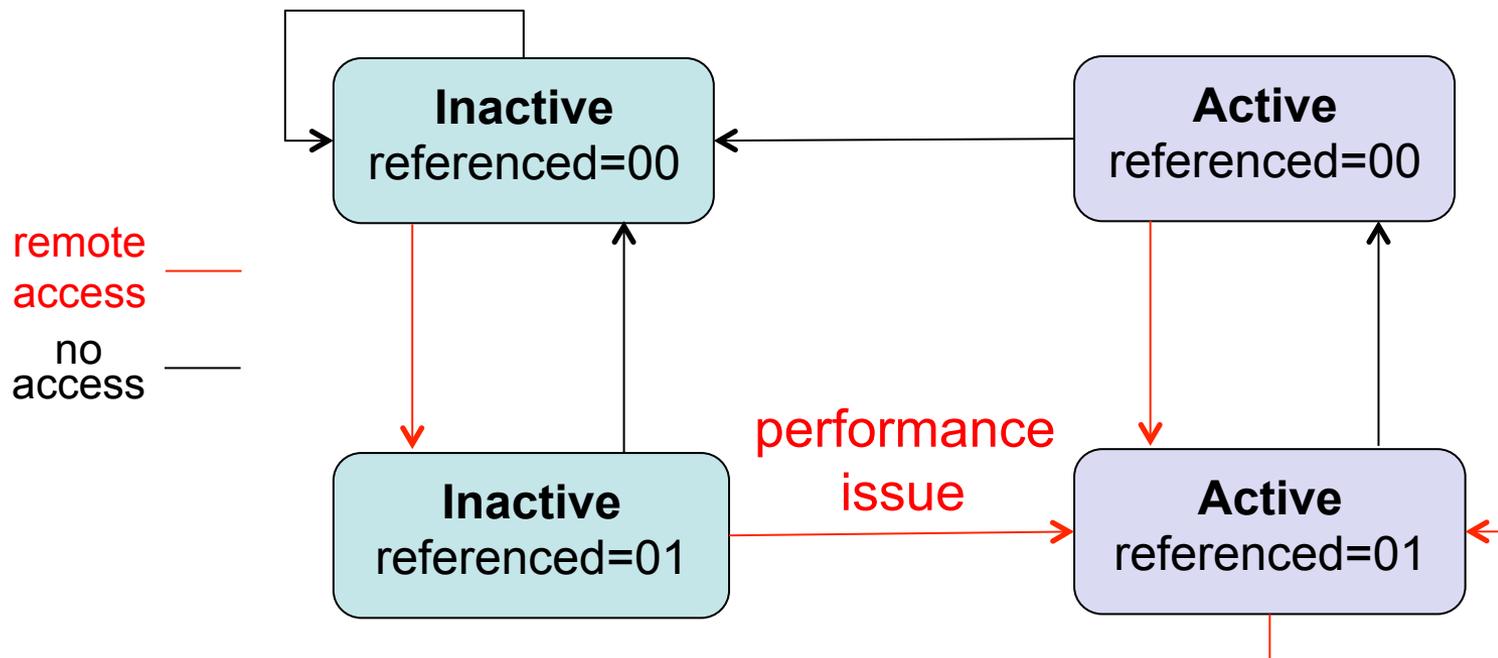
# Common usage of the access / dirty bits

- Normally used to realize demand paging.
  - Approximation of Least Recently Used (LRU)
  - Classical concept
    - » Managing of two lists of active and inactive page frames
    - » State transition realized via access bits
    - » Doubling the number of accesses via a reference bit to move pages from the inactive to active list.



# Transfer to the Node-based Memory Management

- Usage of two reference bits
  - One to signalize local and one to signalize remote memory accesses
- Abstract of the new state graph



# Jacobi solver as Application Benchmark

- Solving of  $A \cdot x = b, A \in R^{n \times n}, b \in R^n, x \in R^n$

- Iterative rule:
$$x_i^{m+1} = \frac{1}{a_{i,i}} \left( b_i - \sum_{j \neq i} a_{i,j} x_j^m \right)$$

- Abstract code for the new memory management

(sequential) initialization of A, b and  $x_0$

forbid write access to A and b

while(!found\_solution)

parallel for over the iterative rule

allow write access to A and b

- Straightforward implementation

# Jacobi solver as Application Benchmark

- Solving of  $A \cdot x = b, A \in R^{n \times n}, b \in R^n, x \in R^n$

- Iterative rule:
$$x_i^{m+1} = \frac{1}{a_{i,i}} \left( b_i - \sum_{j \neq i} a_{i,j} x_j^m \right)$$

- Abstract code

(sequential) initialization of A, b and  $x_0$

~~forbid write access to A and b~~

while(!found\_solution)

    parallel for over the iterative rule

~~allow write access to A and b~~

# Jacobi solver as Application Benchmark

- Solving of  $A \cdot x = b, A \in R^{n \times n}, b \in R^n, x \in R^n$

- Iterative rule:
$$x_i^{m+1} = \frac{1}{a_{i,i}} \left( b_i - \sum_{j \neq i} a_{i,j} x_j^m \right)$$

- Abstract code

```
(sequential) initialization of A, b and x0  
forbid write access to A and b thread binding  
while(!found_solution)  
    parallel for over the iterative rule  
allow write access to A and b
```

# Jacobi solver as Application Benchmark

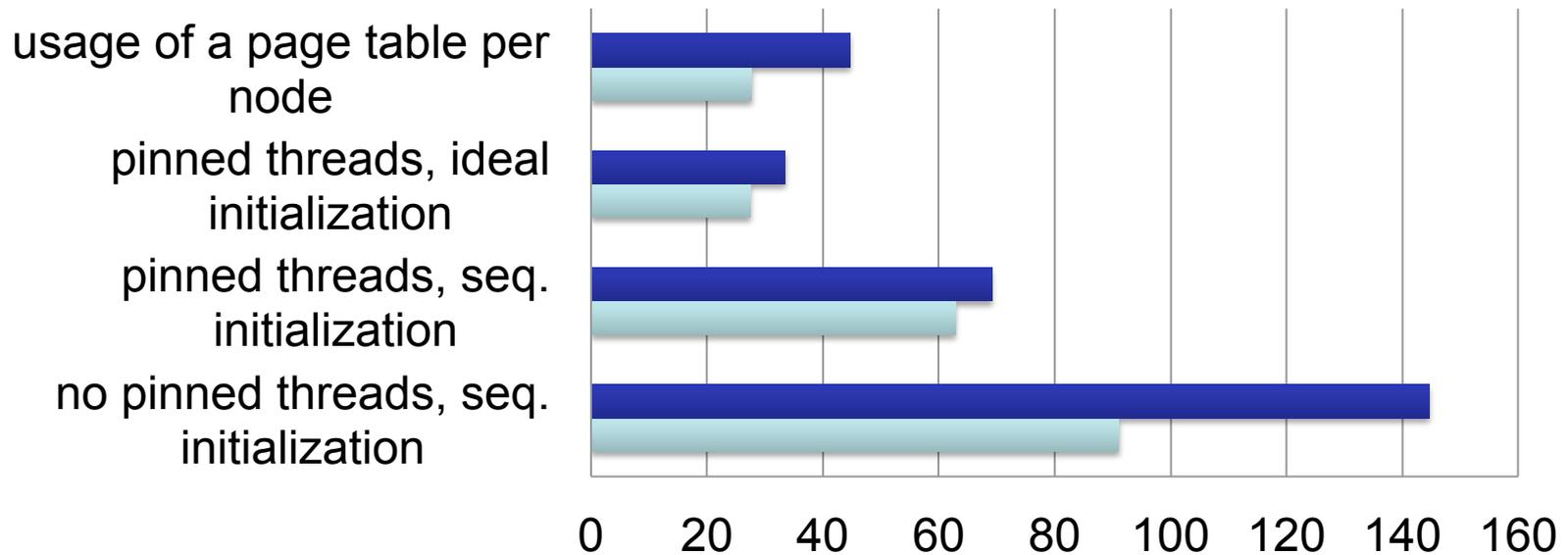
- Solving of  $A \cdot x = b, A \in R^{n \times n}, b \in R^n, x \in R^n$

- Iterative rule:
$$x_i^{m+1} = \frac{1}{a_{i,i}} \left( b_i - \sum_{j \neq i} a_{i,j} x_j^m \right)$$

- Abstract code

```
(sequentialideal) initialization of A, b and x0  
forbid write access to A and b thread binding  
while(!found_solution)  
    parallel for over the iterative rule  
allow write access to A and b
```

# Jacobi solver (Westmere-EX)



	no pinned threads, seq. initialization	pinned threads, seq. initialization	pinned threads, ideal initialization	usage of a page table per node
■ 80 threads	144,609	69,247	33,543	44,77
■ 160 threads	91,067	62,864	27,517	27,746

matrix size: 5120 x 5120  
iterations: 20000

- Memory management can reflect the actual hardware
- First performance results are promising
- Reduction of overhead by
  - usage of virtual address-to-node mapping
  - bundling of NUMA nodes
- Introduce possibilities to detect performance issues
- Simple integration into existing programming models

```
#pragma omp parallel for shared(A,B,C)  
    readonly(A,B)
```

```
for (i=0; i<0; i++)  
    C[i] = A[i] + B[i];
```



Thank you for your kind attention!

**Stefan Lankes**

Chair for Operating Systems  
RWTH Aachen University  
Kopernikusstr. 16  
52056 Aachen, Germany

[www.ifbs.rwth-aachen.de](http://www.ifbs.rwth-aachen.de)

[contact@ifbs.rwth-aachen.de](mailto:contact@ifbs.rwth-aachen.de)



## **Backup slides**

- Page placement strategies are extensively investigated
  - Page placement via hints
    - » Affinity-On-Next-Touch
      - Proposals: Nordergraaf & van der Pas
      - Variations: Shermerhorn, Goglin et al., Bircsak et al.
    - » Template library of locality management (Majo & Gross)
  - (Semi)automatic page placement
    - » profile-guided automatic page placement (Mueller et al.)
    - » dynamic page migration via counting remote memory accesses
      - Memory controller extensions: SGI Origin
      - Compiler extensions: Nikolopoulos et al.
- However, it exists room for optimizations.

# Page Table per Node

## Basic idea

- One page table per node
- Context switch: Load node-local page table
- Page fault
  - Page not mapped: allocate new page and map locally
  - Page mapped remotely:
    - » RW page: duplicate mapping
    - » RO page: duplicate page and map clone locally
- New system call to create a process, which uses our node-based memory management,
  - Per default, the processes use the traditional concept.
- Via `mprotect` the page replication could be implicitly en- or disabled for certain memory regions.

# Overhead (Westmere-EX)

	unmodified Linux kernel (3.3.8)	page table per node
time to allocate a page	1.666 $\mu$ s	6.671 $\mu$ s
time to protect a page	0.00005 $\mu$ s	0.032 $\mu$ s
time to replicate a page	—	4.479 $\mu$ s
time to unprotect a page	0.0001 $\mu$ s	0.148 $\mu$ s
time to replicate a reference	—	1.445 $\mu$ s

## Test platform

- 8 Intel Xeon CPU E7-8850 (Westmere-EX)
- 8 \* 10 Cores / 8 \* 20 Cores via HyperThreading

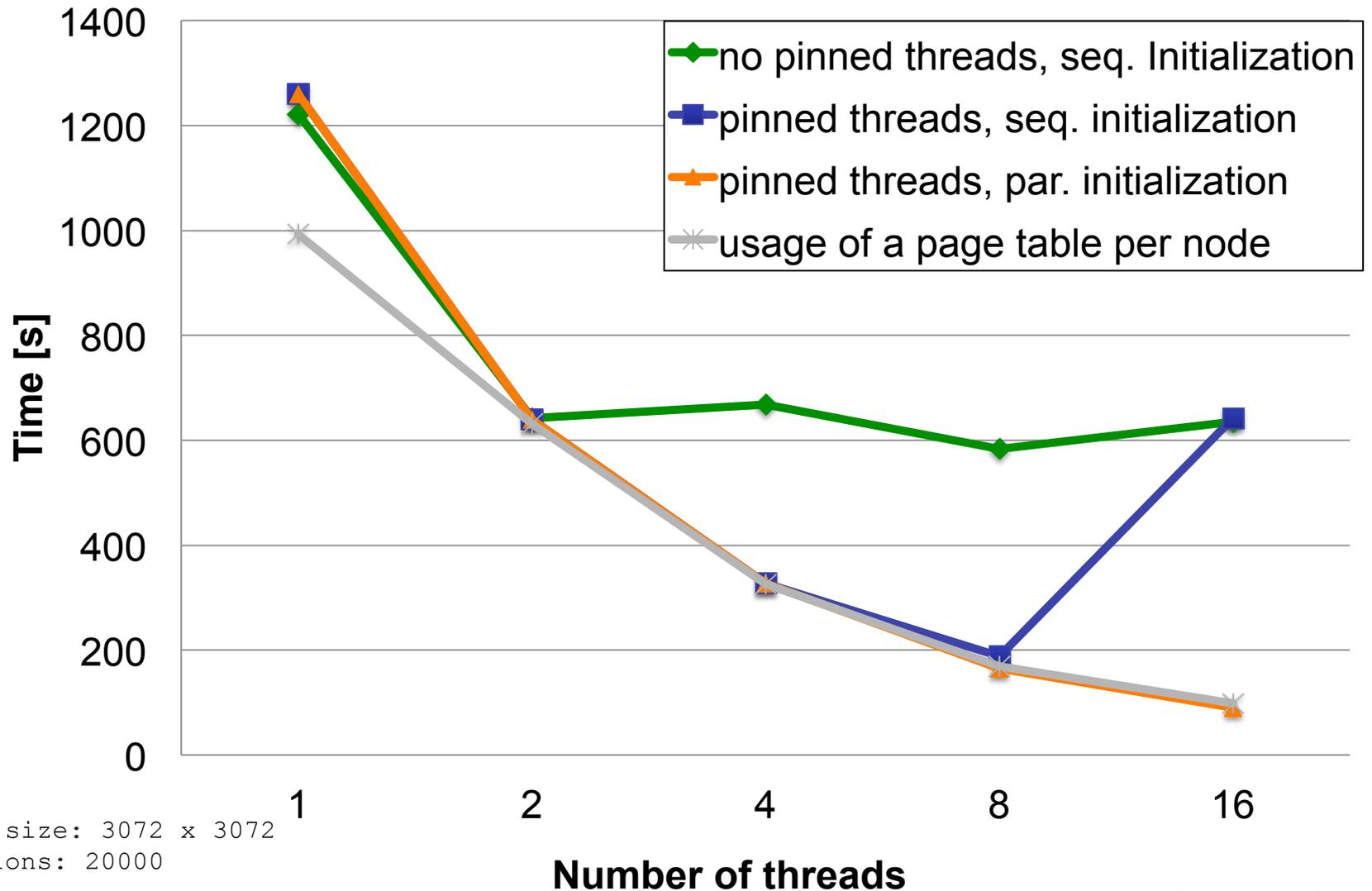
# Overhead (NumaScale-Cluster)

	unmodified Linux kernel (2.6.37)	page table per node
time to allocate a page	2.810 $\mu$ s	3.143 $\mu$ s
time to protect a page	0.034 $\mu$ s	0.110 $\mu$ s
time to replicate a page	—	26.956 $\mu$ s
time to unprotect a page	0.195 $\mu$ s	2.787 $\mu$ s
time to replicate a reference	—	6.044 $\mu$ s

## Test platform

- 2 systems with 2 AMD QuadCores of type 8378 combined via NumaConnect

# Jacobi solver (NumaScale-Cluster)



matrix size: 3072 x 3072  
iterations: 20000