# Improving CHARM++ Performance with a NUMA-aware Load Balancer

Laércio Lima Pilla[1,2], Christiane Pousa[2], Daniel Cordeiro[2,3], Abhinav Bhatele[4], Philippe O. A. Navaux[1], Jean-François Méhaut[2], Laxmikant V. Kale[4]

[1]Federal University of Rio Grande do Sul – Porto Alegre, Brazil
[2]Grenoble University – Grenoble, France
[3]University of São Paulo – São Paulo, Brazil
[4]University of Illinois at Urbana-Champaign – Urbana, IL, USA

# Summary

How we used NUMA architectural information to build a CHARM++ load balancer and obtained improvements on overall performance.

# Agenda
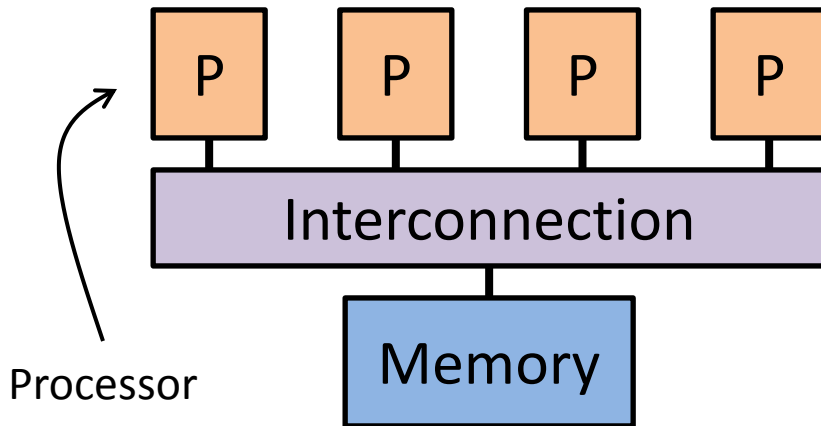
## NUMA

Our Load Balancer: NUMALB

Experimental Setup
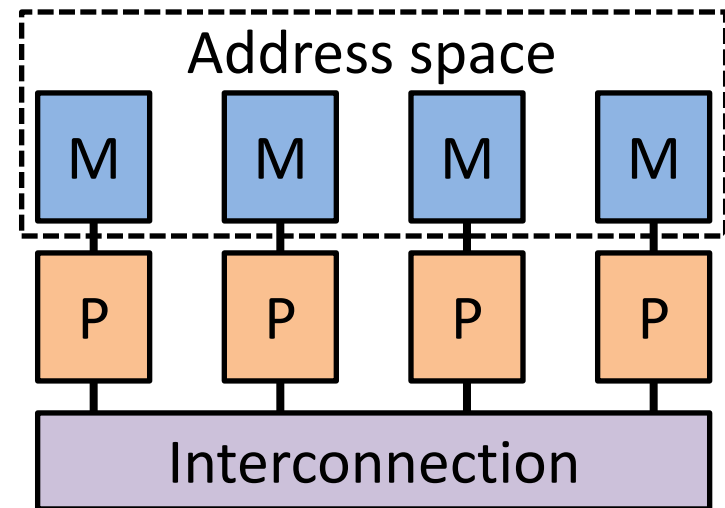
Results

Concluding Remarks

# UMA x NUMA

## Uniform Memory Access

- Centralized shared memory
  - Uniform latencies
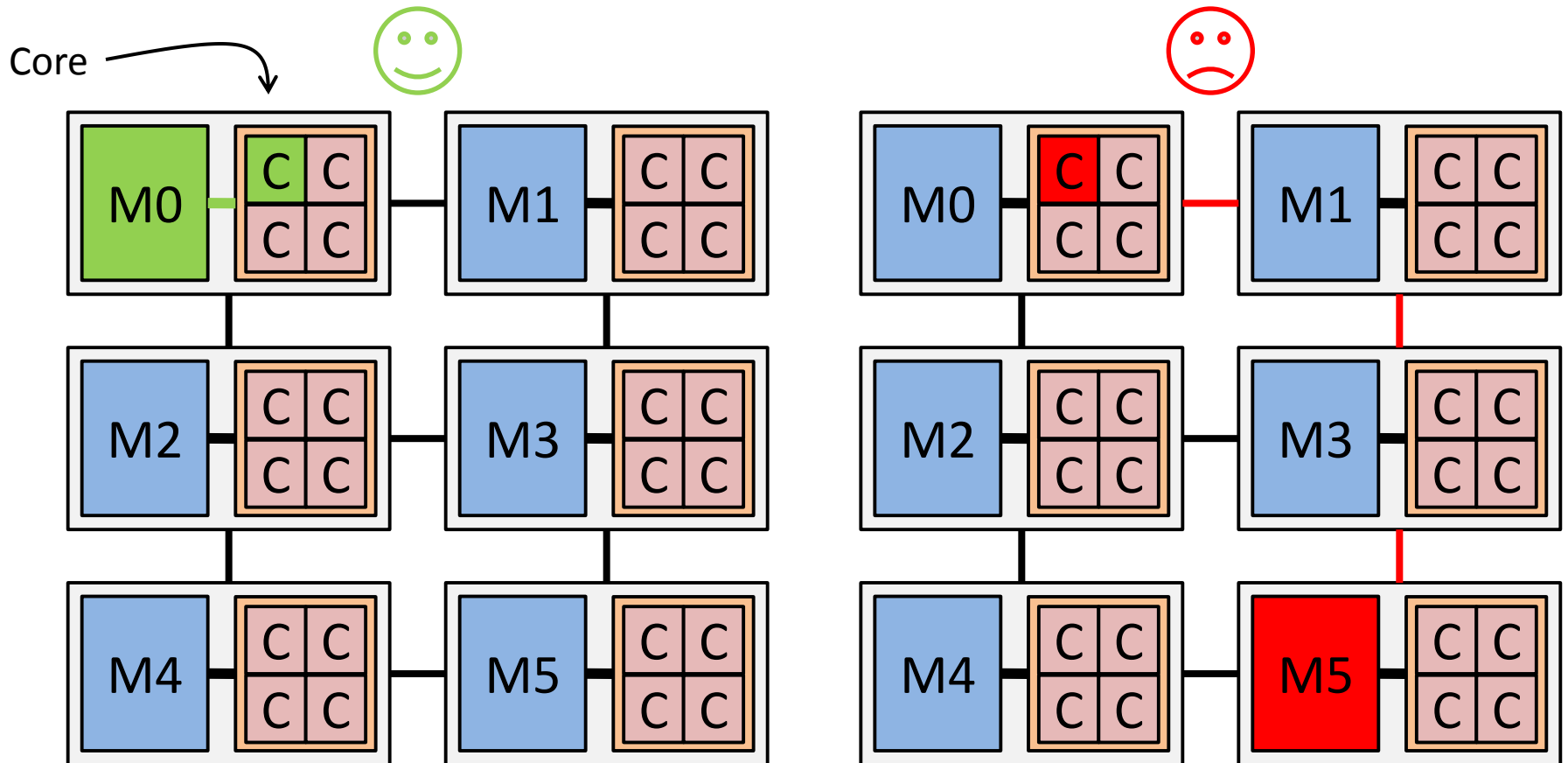- Data placement does not matter



Processor

## Non-Uniform Memory Access

- Distributed shared memory
  - Non-uniform latencies
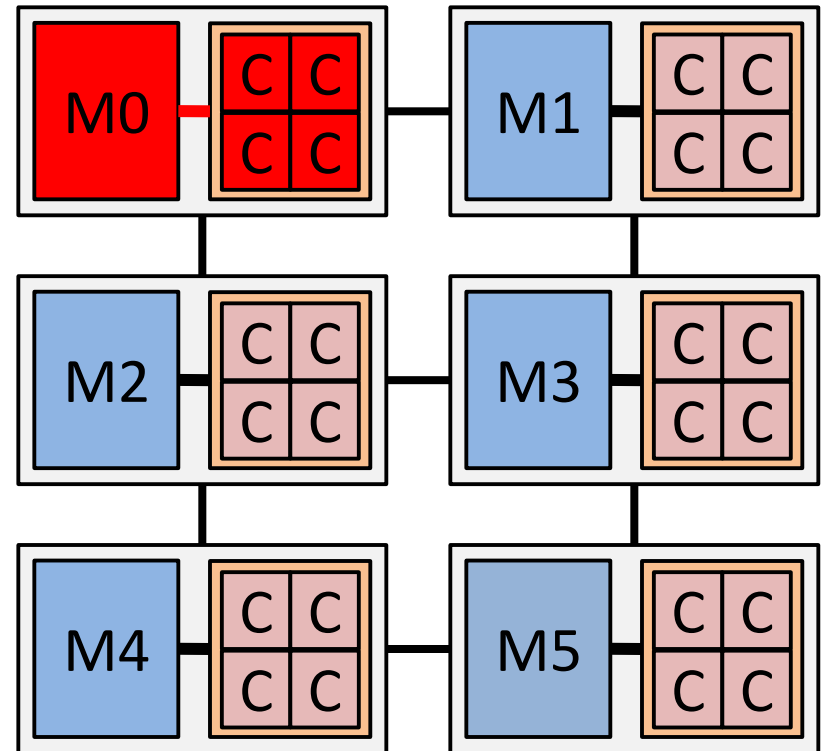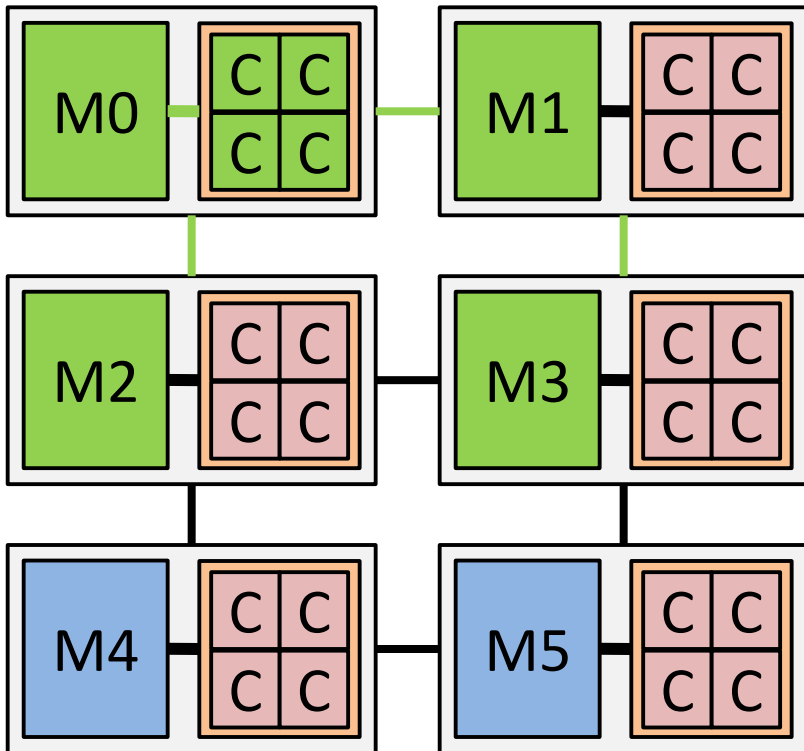- **Data placement matters**

# NUMA

## Reduce latencies
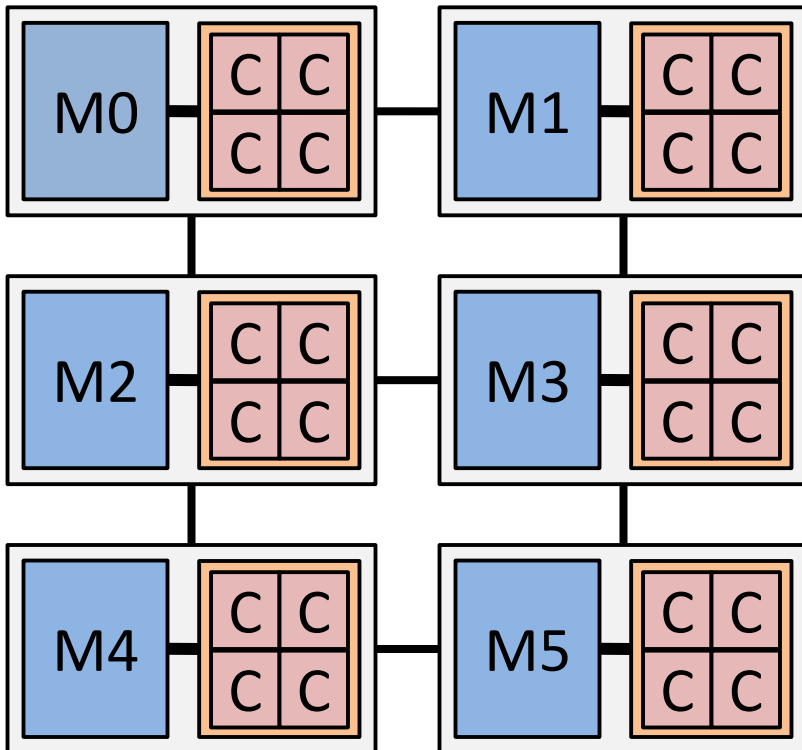
# NUMA

## Reduce contention/improve bandwidth

# NUMA

## CHARM++ does not consider these characteristics

### Physical organization



### CHARM++'s vision (UMA)

No memory hierarchy

No locality

# Agenda

~~NUMA~~

## Our Load Balancer: NumaLB

Experimental Setup

Results

Concluding Remarks

# Load Balancer

- Application data – CHARM++ LB framework
  - Processor load: execution time
  - Chare load: execution time
  - Communication graph: size and number of messages
- NUMA topology – archTopology (our library)
  - Core to NUMA node (socket) hierarchy mapping
  - NUMA factor

$$NUMA\ factor\ (i, j) = \frac{Read\ latency\ from\ i\ to\ j}{Read\ latency\ on\ i}$$

# Load Balancer

- Heuristic
  - Task mapping is NP-Hard
  - No initial assumptions about the application
- List scheduling
  - Put tasks on a priority list by load
  - Assign tasks to the processor with the smallest cost on a greedy fashion
- Improve performance
  - by reducing unbalance
  - by reducing remote communication costs
  - while avoiding migrations (data movement costs)

# Load Balancer

- Cost function

$$cost(c,p) = load(p) +$$
$$a \times ( r_{comm}(c,p) \times NUMA\ factor$$
$$- l_{comm}(c,p) )$$

Where
  $c$: chare
  $p$: core
  $load(p)$: load (execution time) on core $p$
  $r_{comm}(c,p)$: number of messages sent by chare $c$ to chares on other NUMA node
  $l_{comm}(c,p)$: number of messages sent by chare $c$ to chares on the same NUMA node
  $a$: communication weight

# Load Balancer

**NumaLB's Algorithm**

**Input**: $C$ set of chares, $P$ set of cores, $M$ mapping

**Output**: $M'$ mapping of chares to cores

1.    $M' \leftarrow M$

2.    **while** $c \neq \emptyset$ **do**       for the number of chares

3.    $c \leftarrow v \mid v \in \arg \max_{u \in C} load(u)$    take heaviest chare

4.    $C \leftarrow C \setminus \{c\}$

5.    $p \leftarrow q, q \in P \wedge \{(c,q)\} \in M$    get its core

6.    $load(p) \leftarrow load(p) - load(c)$    remove its load from its core

7.    $M' \leftarrow M' \setminus \{(c,p)\}$    remove from mapping

8.    $p' \leftarrow q \mid q \in \arg \min_{r \in P} cost(c,r)$    find core with smallest cost

9.    $load(p') \leftarrow load(p') + load(c)$    add chare load to new core

10.   $M' \leftarrow M' \cup \{(c,p')\}$    map to new core

# Agenda

~~NUMA~~

~~Our Load Balancer: NumaLB~~

[Experimental Setup](#)

Results

Concluding Remarks

# Experimental Setup

- 2 NUMA machines

- 3 CHARM++ benchmarks

- 4 other CHARM++ load balancers

- Statistical confidence of 95%
  - 5% relative error
  - Student's t-distribution
  - Minimum of 25 executions

- Performance
  - Gains: Average iteration time (baseline = no LB)
  - Costs: Load balancing overhead

# Experimental Setup: Machines

- ## NUMA16
  - AMD Opteron
  - 8×2 cores @ 2.2 GHz
  - 1 MB private L2 cache
  - 32 GB main memory
  - Low latency for memory access
  - Crossbar
  - NUMA factor: 1.1–1.5

# Experimental Setup: Machines

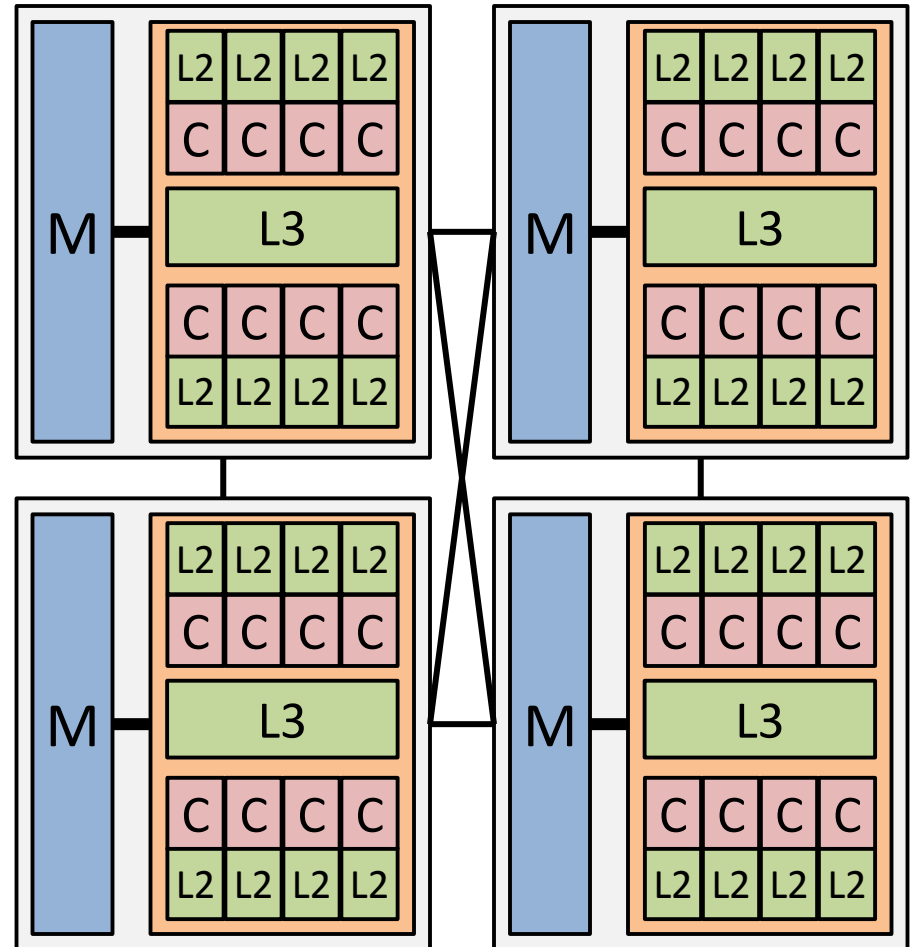- NUMA32
  - Intel Xeon X7560
  - 4×8 cores @ 2.27 GHz
  - 256 KB private L2
  - 24 MB shared L3
  - 64 GB main memory
  - QuickPath
  - NUMA factor: 1.36–3.6

# Experimental Setup: Benchmarks

- kNeighbor
  - Synthetic iterative benchmark where a chare communicates with other $k$ chares at each step
  - Completely I/O bound
  - 200 chares, 16 KB messages, $k = 8$
- lb_test
  - Synthetic unbalanced benchmark with different possible communication patterns
  - 200 chares, random communication graph, load between 50 and 200 ms
- jacobi2D
  - Unbalanced two-dimensional five-point stencil
  - 100 chares, $32^2$ data array

# Experimental Setup: LBs

- GREEDYLB
  - Iteratively maps the most loaded chares to the least loaded cores

- RECBIPARTLB
  - Recursive bipartition of the communication graph
  - Breadth-first traversal until groups the required load

- METISLB
  - Graph partitioning algorithms from METIS

- SCOTCHLB
  - Graph partitioning algorithms from SCOTCH

- Neither consider the current chare mapping

# Agenda

~~NUMA~~

~~Our Load Balancer: NumaLB~~

~~Experimental Setup~~

## Results

Concluding Remarks

# Results: kNeighbor

# Results: kNeighbor



Homogeneous distribution

Shared cache, faster communication

Group chares and migrate them together to the same core

Average iteration time (in ms)

**NUMA16:** Baseline 32.0, 30%, NumaLB 22.4, GreedyLB 22.1, MetisLB 21.5, RecBipartLB 22.8, ScotchLB 22.1

**NUMA32:** Baseline 26.9, 45%, NumaLB 14.9, GreedyLB 13.5, MetisLB 17.9, RecBipartLB 16.9, ScotchLB 16.1

Legend: Baseline, NumaLB, GreedyLB, MetisLB, RecBipartLB, ScotchLB

# Results: lb_test



Best performance by communication-aware LBs

Best average performance

Average iteration time (in s)

**NUMA16**
- Baseline: 1.01 (17%)
- NumaLB: 0.83
- GreedyLB: 0.93
- MetisLB: 0.84
- RecBipartLB: 0.83
- ScotchLB: 0.88

**NUMA32**
- Baseline: 0.6 (28%)
- NumaLB: 0.43
- GreedyLB: 0.51
- MetisLB: 0.46
- RecBipartLB: 0.47
- ScotchLB: 0.43

Baseline ■ NumaLB ■ GreedyLB ■ MetisLB ■ RecBipartLB ■ ScotchLB

# Results: jacobi2D

# Results: jacobi2D - Projections

**METISLB: 75% efficiency**
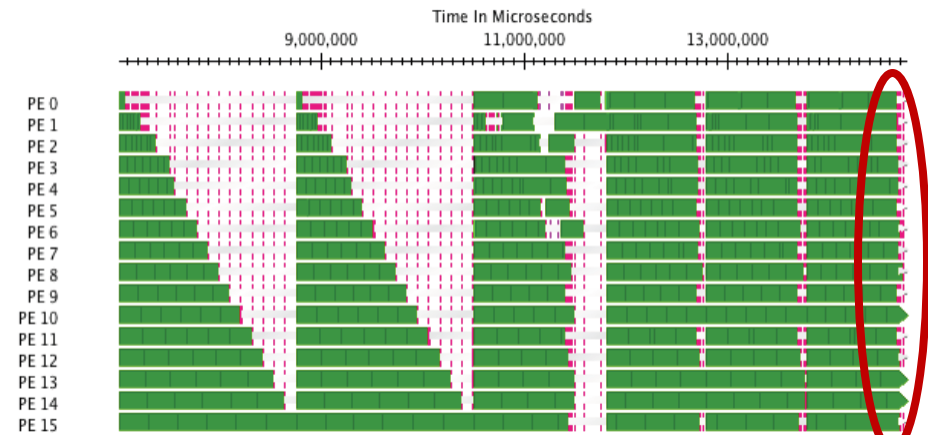
- jacobi2D on NUMA16
  - 2 steps before LB
  - 4 steps after LB

- The smaller the idle parts, the higher the efficency

**NUMALB: 93.5% efficiency**

Improving Charm++ Performance with a NUMA-aware Load Balancer

# Results: overheads

## Average number of chares migrated

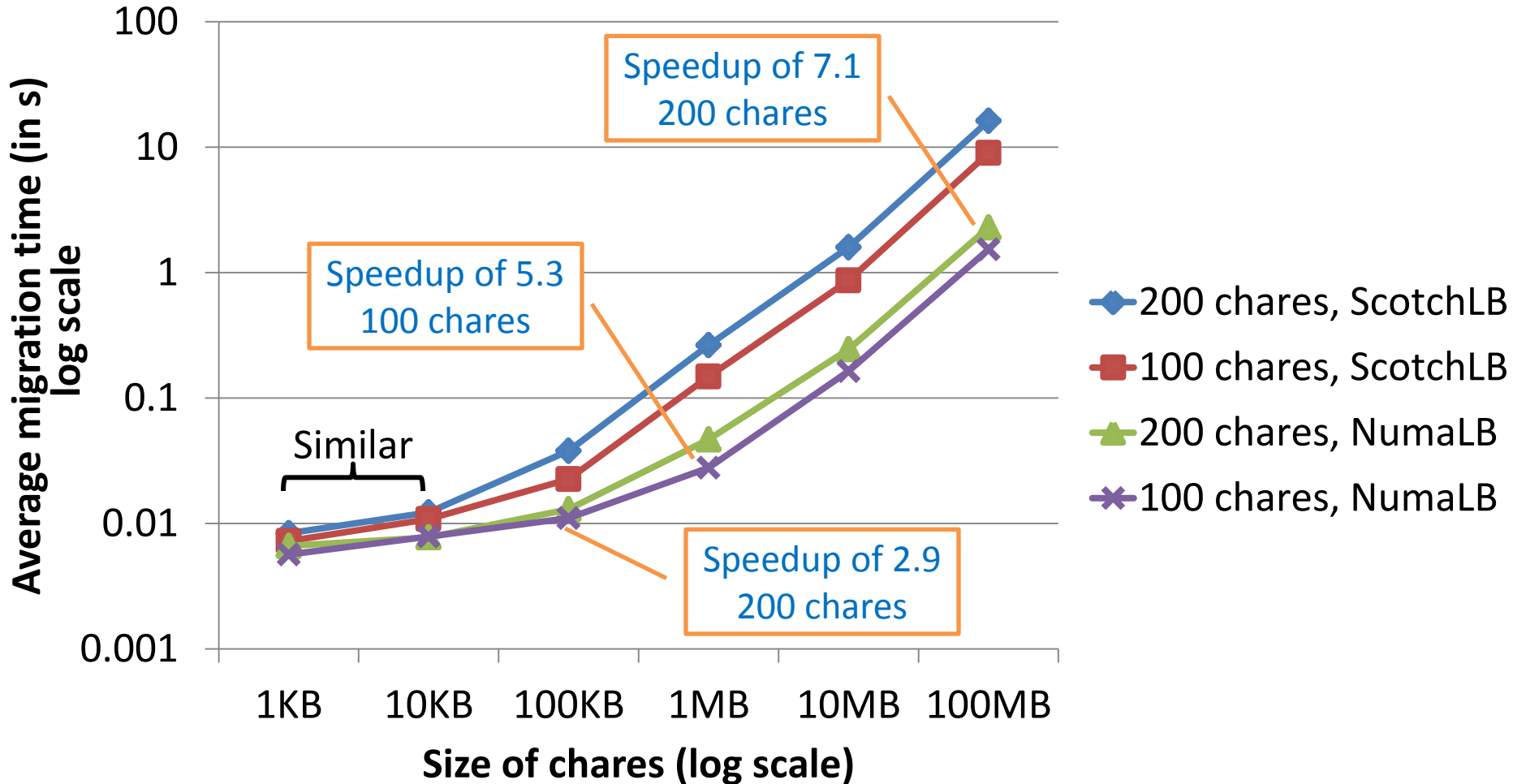| Benchmark | Machine | Load Balancer | | | | |
|---|---|---|---|---|---|---|
| | | NUMALB | GREEDYLB | METISLB | RECBIPARTLB | SCOTCHLB |
| kNeighbor | NUMA16 | 25 | 189 | 188 | 176 | 185 |
| | NUMA32 | 57 | 194 | 195 | 185 | 194 |
| lb_test | NUMA16 | 40 | 188 | 187 | 184 | 184 |
| | NUMA32 | 48 | 194 | 194 | 192 | 192 |
| jacobi2D | NUMA16 | 26 | 94 | 94 | 91 | 93 |
| | NUMA32 | 33 | 97 | 96 | 93 | 98 |

**Maximum migrations = 33%**  😊

**Minimum migrations = 88%**  ☹

All load balancers took less than 7 ms for their algorithms.

# Results: migration times for NUMA16



Speedup of 7.1
200 chares

Speedup of 5.3
100 chares

Speedup of 2.9
200 chares

Similar

Average migration time (in s) log scale

100
10
1
0.1
0.01
0.001

Size of chares (log scale)

1KB    10KB    100KB    1MB    10MB    100MB

200 chares, ScotchLB
100 chares, ScotchLB
200 chares, NumaLB
100 chares, NumaLB

# Agenda

NUMA

Our Load Balancer: NumaLB

Experimental Setup

Results

Concluding Remarks

# Conclusions

- Multi-core machines with NUMA design introduce new challenges for their efficient use

- CHARM++ does not consider NUMA asymmetries

- With our NUMA-aware LB we obtained
  - An average speedup of 1.51 over the baseline
    - Transparent to the user, no previous knowledge
  - 10% improvement over most LBs
  - Migration overheads up to 7 times smaller
    - Migrating at most 33% of all chares

# Future Work

- Multi-core load balancer
  - UMA and NUMA machines
  - Communication latencies among cores
  - Use HWLOC representation of cache hierarchy
- Distributed multi-core load balancer
  - For clusters of multi-core machines
- Gather and organize communication information
  - Latencies, bandwidth
  - Provide this data to other libraries (like SCOTCH)

# Improving CHARM++ Performance with a NUMA-aware Load Balancer

# Thank you.

Laércio Lima Pilla

Contact: llpilla@inf.ufrgs.br