

# Topology-aware load balancing for parallel applications on multi-core systems and beyond




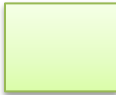
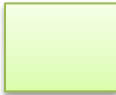
**Laércio Lima Pilla**

and Christiane P. Ribeiro, François Broquedis,  
Pierre Coucheney, Bruno Gaujal, Jean-François Méhaut,  
Philippe O. A. Navaux, Laxmikant Kale

# Overview

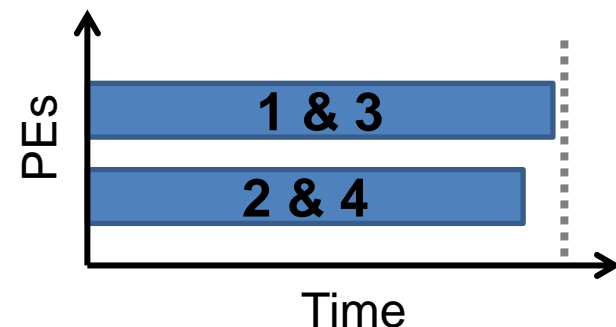
Asymptotically optimal load  
balancing algorithm  
for multi-core machines.

# Motivation

- Example of environment: **CHARM++**
  - Parallel task overdecomposition
    - Chares   
  - Platform independent
    - Processing elements  
  - Dynamic load balancing
    - Chare migration

# Motivation: example

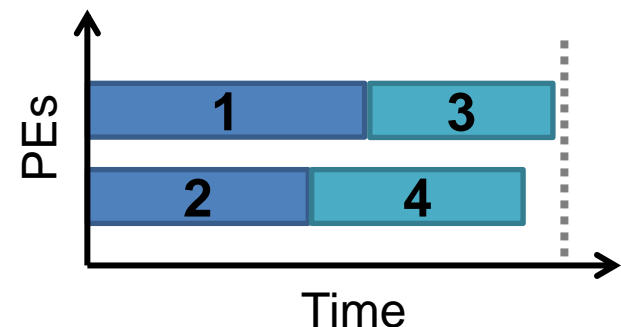
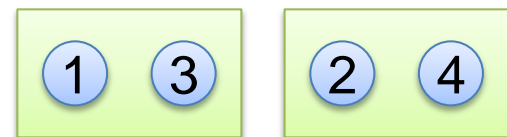
- Initial chore mapping
- Apply a load balancing algorithm
- Based on **data provided by CHARM++**



# Motivation: example

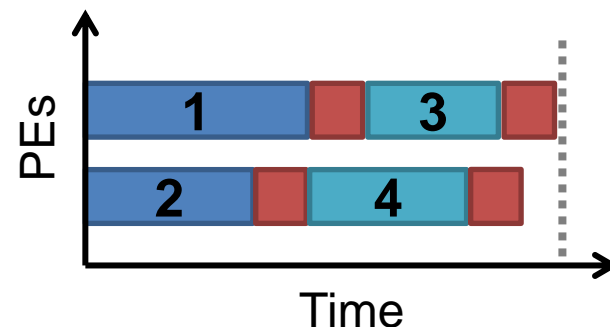
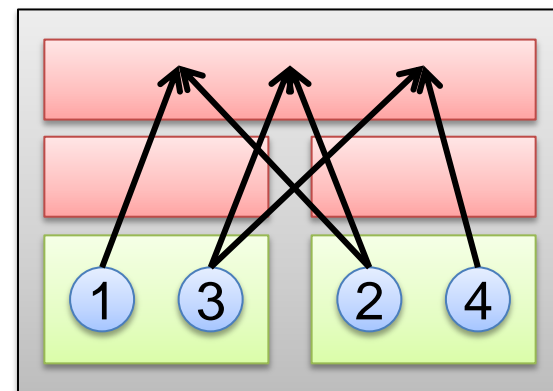


- CHARM++ load balancing data
  - Communication graph
  - Current chare mapping
  - Chares' load



# Motivation: example

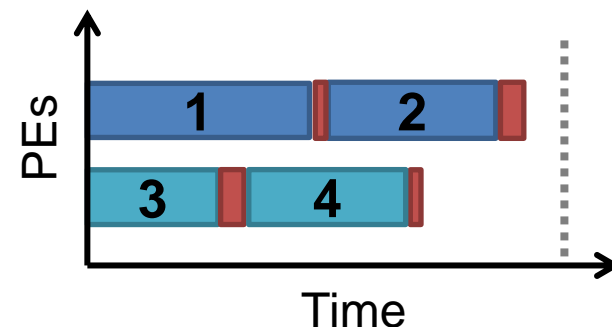
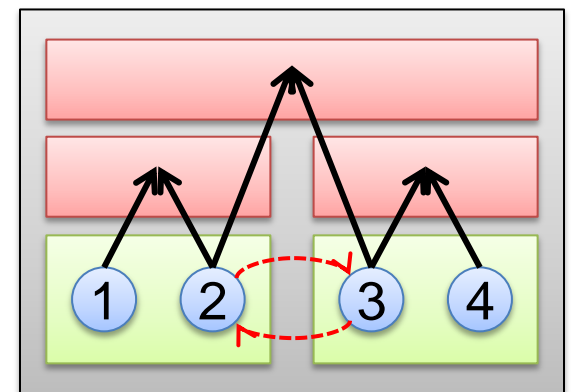
- Missing information
  - Communication costs
  - **Architectural information**
    - **Memory hierarchy**





# Motivation: example

- Knowledge about the **memory hierarchy**
- Memory access costs
  - Reduce communication costs
- **Highly hierarchical systems**
  - Even more on many-core



# Motivation

- Our approach: **TOPOLB**
  - Load balancing algorithm implemented on CHARM++
  - Combines *application information* and the *machine topology*
  - Works on UMA and NUMA machines
  - Is asymptotically optimal



# Agenda

~~Motivation~~

**TOPOLB**

Experiments

And beyond

# TOPOLB: idea

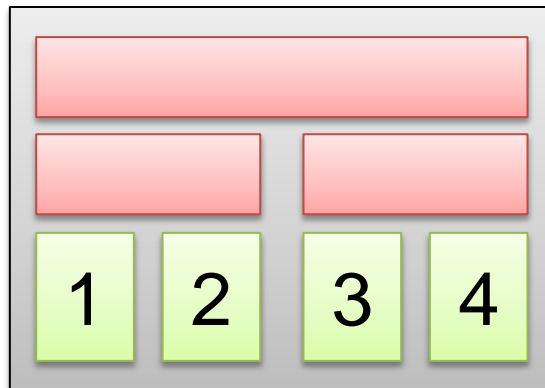
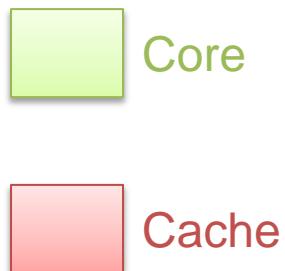
- **Heuristic**
  - Load balancing is NP-Hard
  - No initial assumption about the application
- Improve performance (reduce makespan)
  - By reducing unbalance
  - By reducing communication costs
  - While avoiding migrations (data movement costs)

# TOPOLB: data

- ***Application information*** – provided by CHARM++
  - Chares' load
    - Reduce the makespan
  - Communication graph
    - Reduce the communication overhead
    - Bring together communicating chares
  - Current chare mapping
    - Avoid migration overheads

# TOPOLB: data

- ***Machine topology*** – our library
- **Memory hierarchy**
  - Cache and memory sharing among cores
- **Memory access latencies**
  - Estimate the communication time of each message through the memory hierarchy

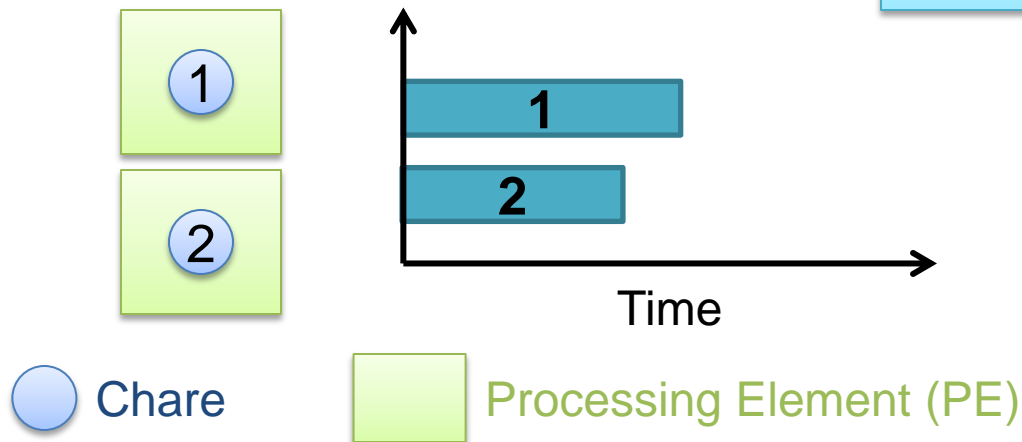


$$\begin{bmatrix} 1.63 & 1.63 & 6.67 & 6.67 \\ 1.63 & 1.63 & 6.67 & 6.67 \\ 6.67 & 6.67 & 1.63 & 1.63 \\ 6.67 & 6.67 & 1.63 & 1.63 \end{bmatrix}$$

# TOPOLB: algorithm

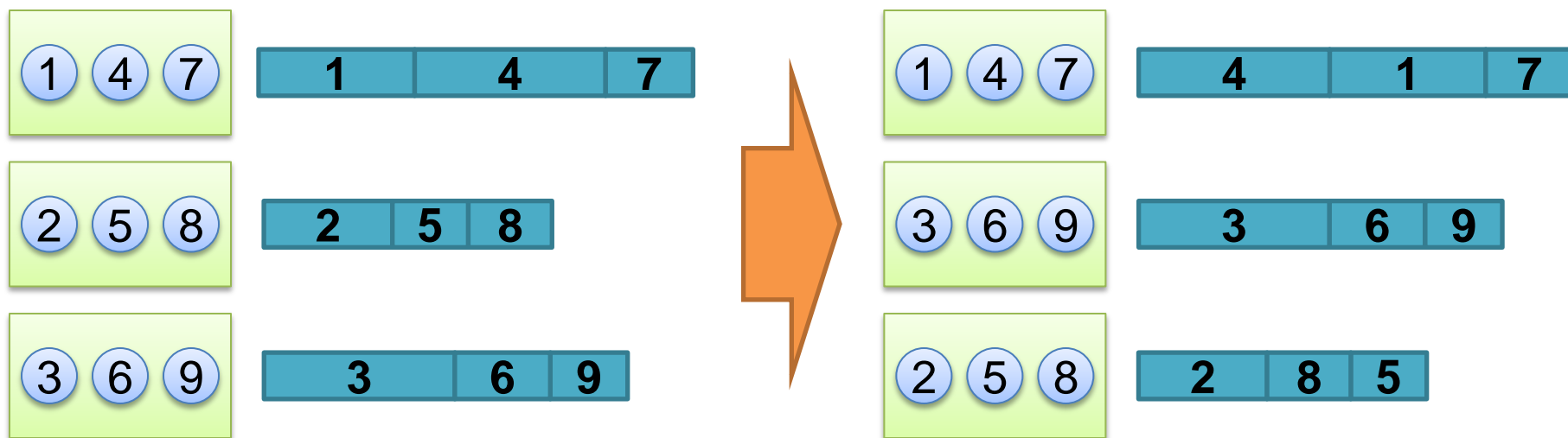
- Each chare has a **load** based on
  - Computation time
  - Communication time
    - **Number of messages** exchanged
    - Memory **latencies**

Communicating  
chares have a  
higher memory  
affinity



# TOPOLB: algorithm

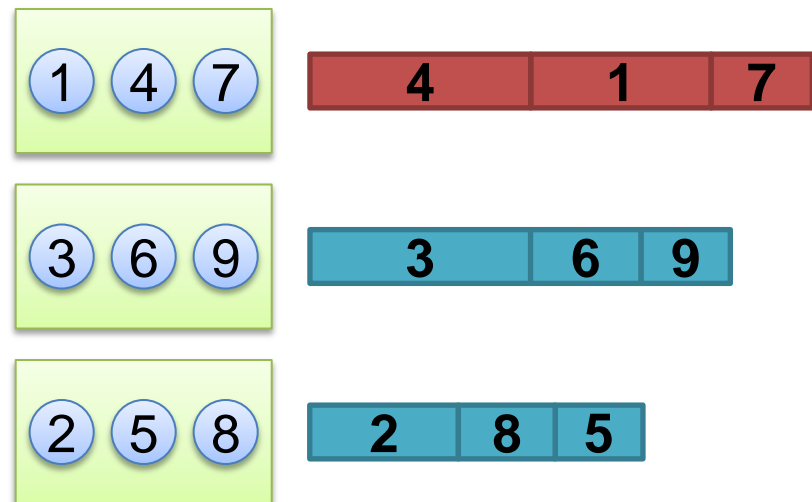
- Initial mapping
- **Sort** PEs (processors) **based on** their **load**





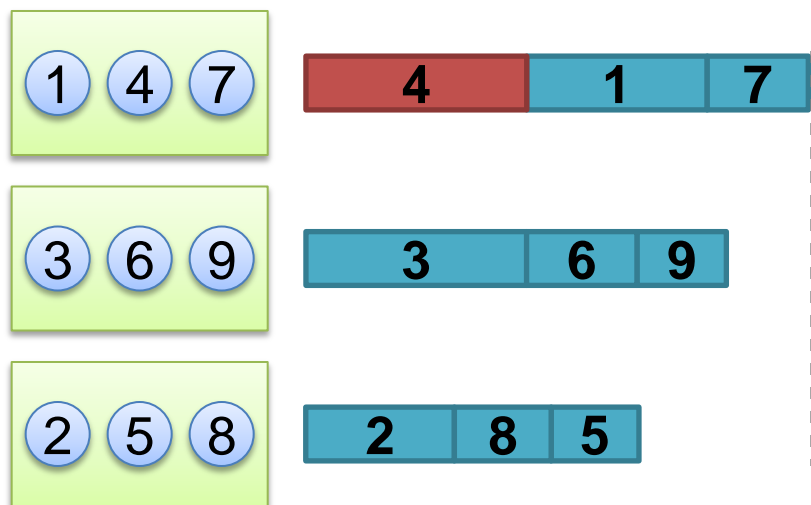
# TOPOLB: algorithm

- Tries to **migrate chares from the most loaded processor** with a **probability  $\alpha$** 
  - Or chooses another processor uniformly



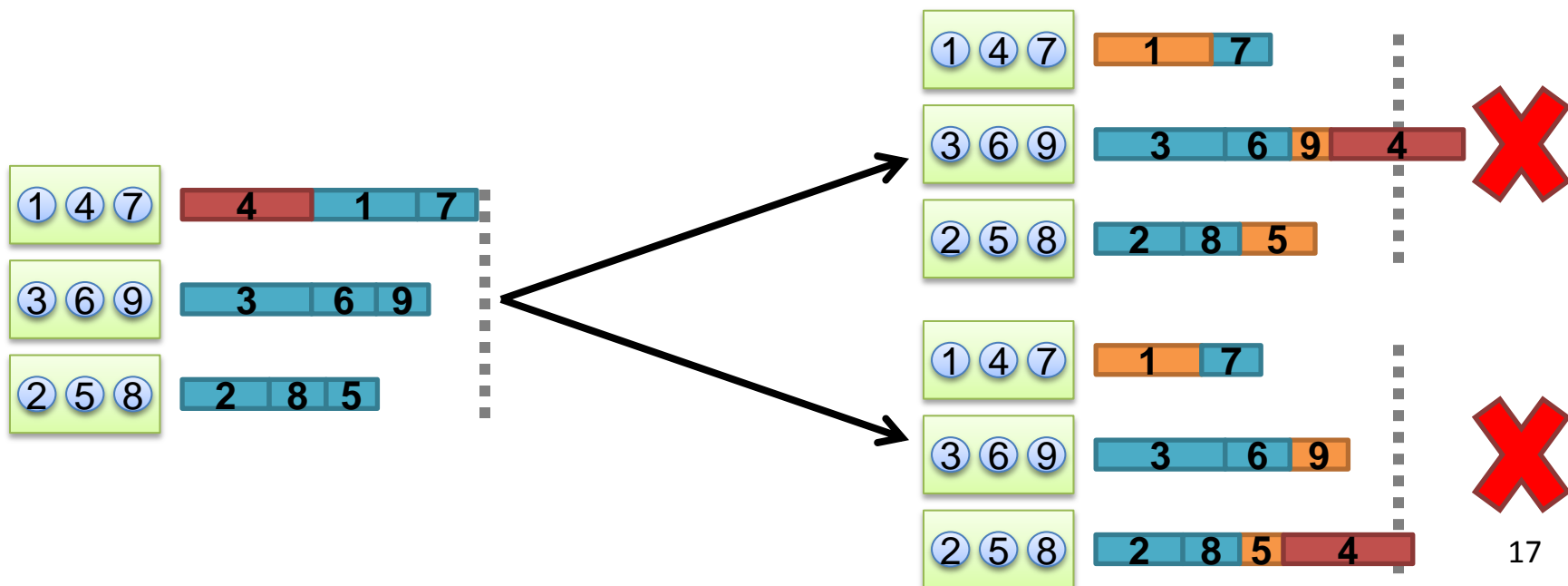
# TOPOLB: algorithm

- Starts from the **heaviest chare** with a **probability  $\beta$** 
  - Or chooses another chare uniformly
  - From heaviest chare to lightest chare



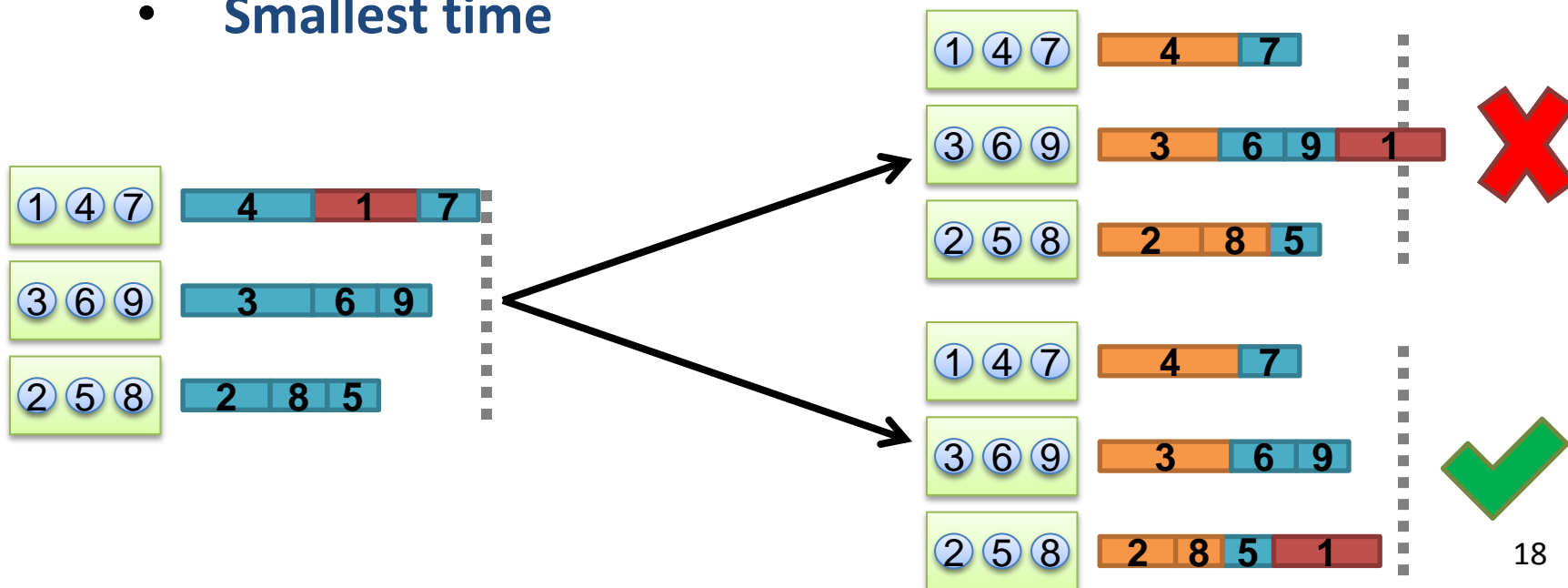
# TOPOLB: algorithm

- A migration might decrease the makespan
- A migration may affect other chares
  - Communication time may change



# TOPOLB: algorithm

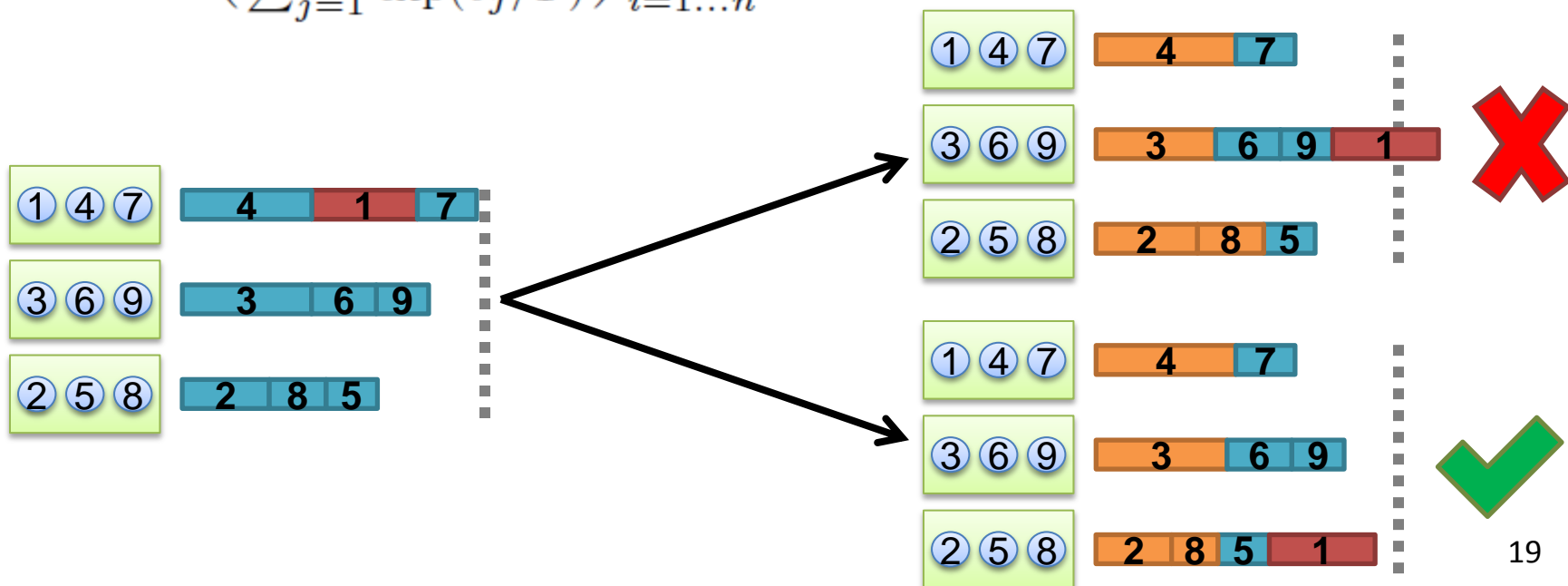
- Evaluates all possible processors
  - Migrates the chore to the **processor that minimizes the makespan** with high probability
    - **Smallest time**



# TOPOLB: algorithm

- **High probability**
  - Using a **Gibbs distribution** with a **temperature  $T$**

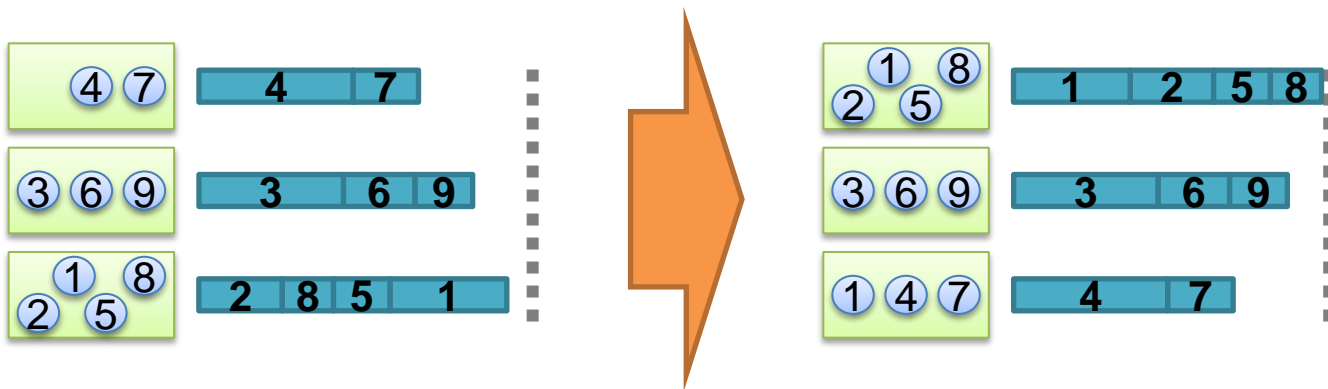
$$\left( \frac{\exp(v_i/T)}{\sum_{j=1}^n \exp(v_j/T)} \right)_{i=1 \dots n}$$





# TOPOLB: algorithm

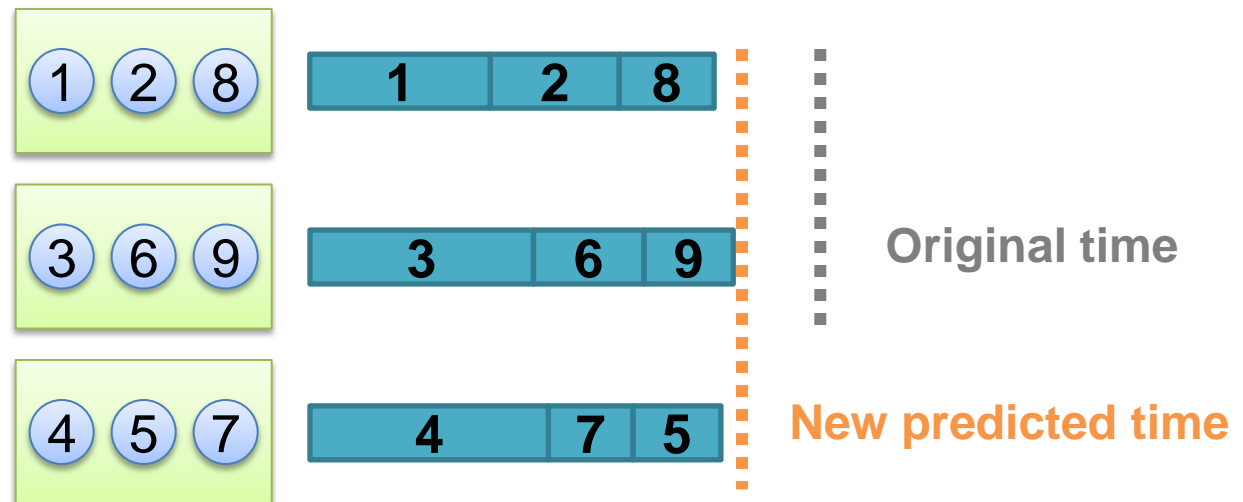
- Continues **until no chore migrates from the heaviest processor**
  - Cannot reduce the makespan
    - Bound to the heaviest processor





# TOPOLB: algorithm

- Continues **until no chare migrates from the heaviest processor**
  - Cannot reduce the makespan
    - Bound to the heaviest processor



# TOPOLB: implementation details

- **Exponential backoff**
  - Reduce the load balancing overhead by not computing the algorithm on all LB calls
  - Launched when **no migrations happen**

1

2

3

4

5

6

7

8

# Agenda

~~Motivation~~

~~TOPOLB~~

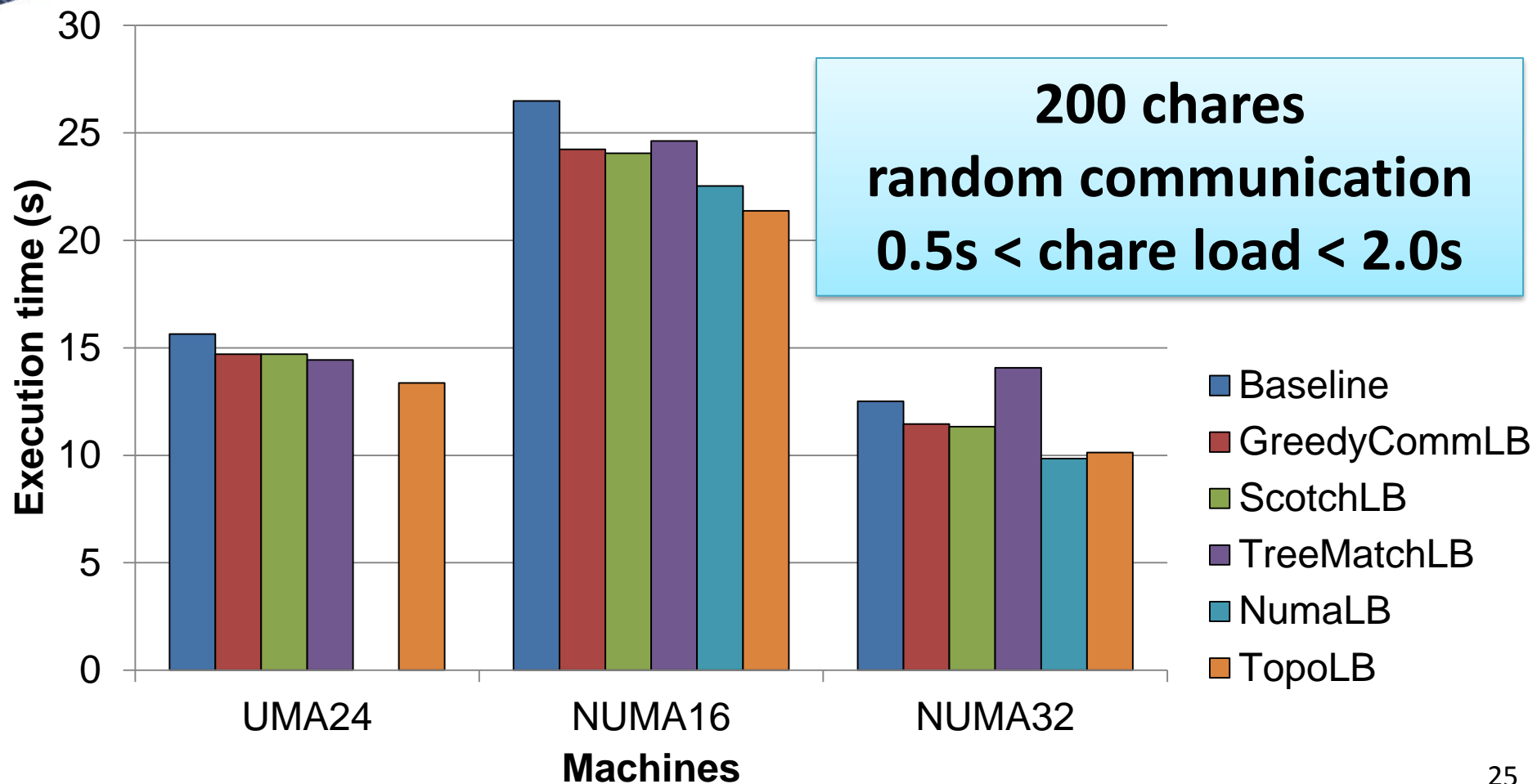
**Experiments**

And beyond

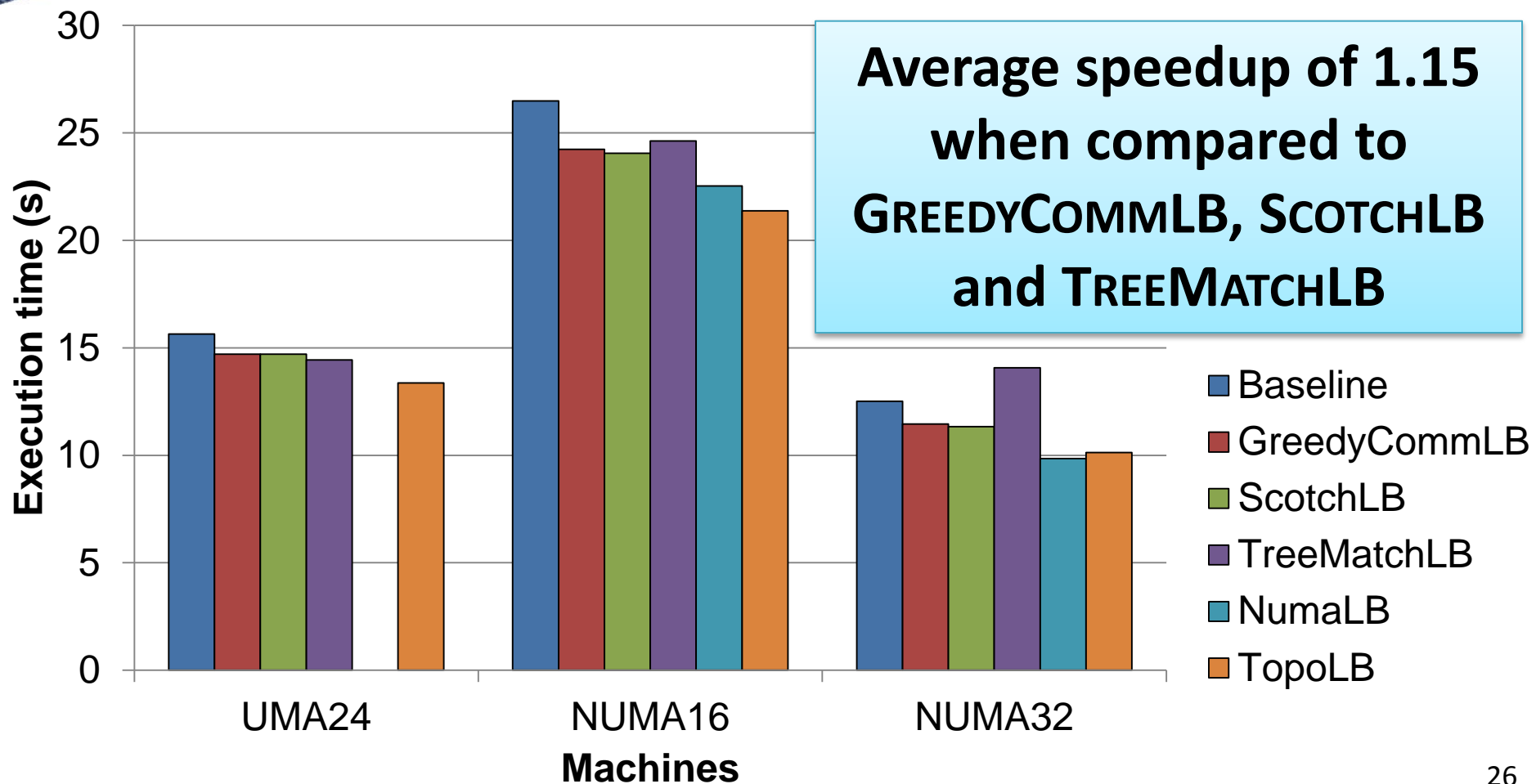
# Experiments

- 3 different machines: UMA24, NUMA16 and NUMA32
- 4 other load balancers: GREEDYCOMMLB, SCOTCHLB, TREEMATCHLB and NUMALB
- 2 benchmarks: lb\_test and mol3D
- $\alpha, \beta$  close to 1,  $T$  close to 0

# lb\_test

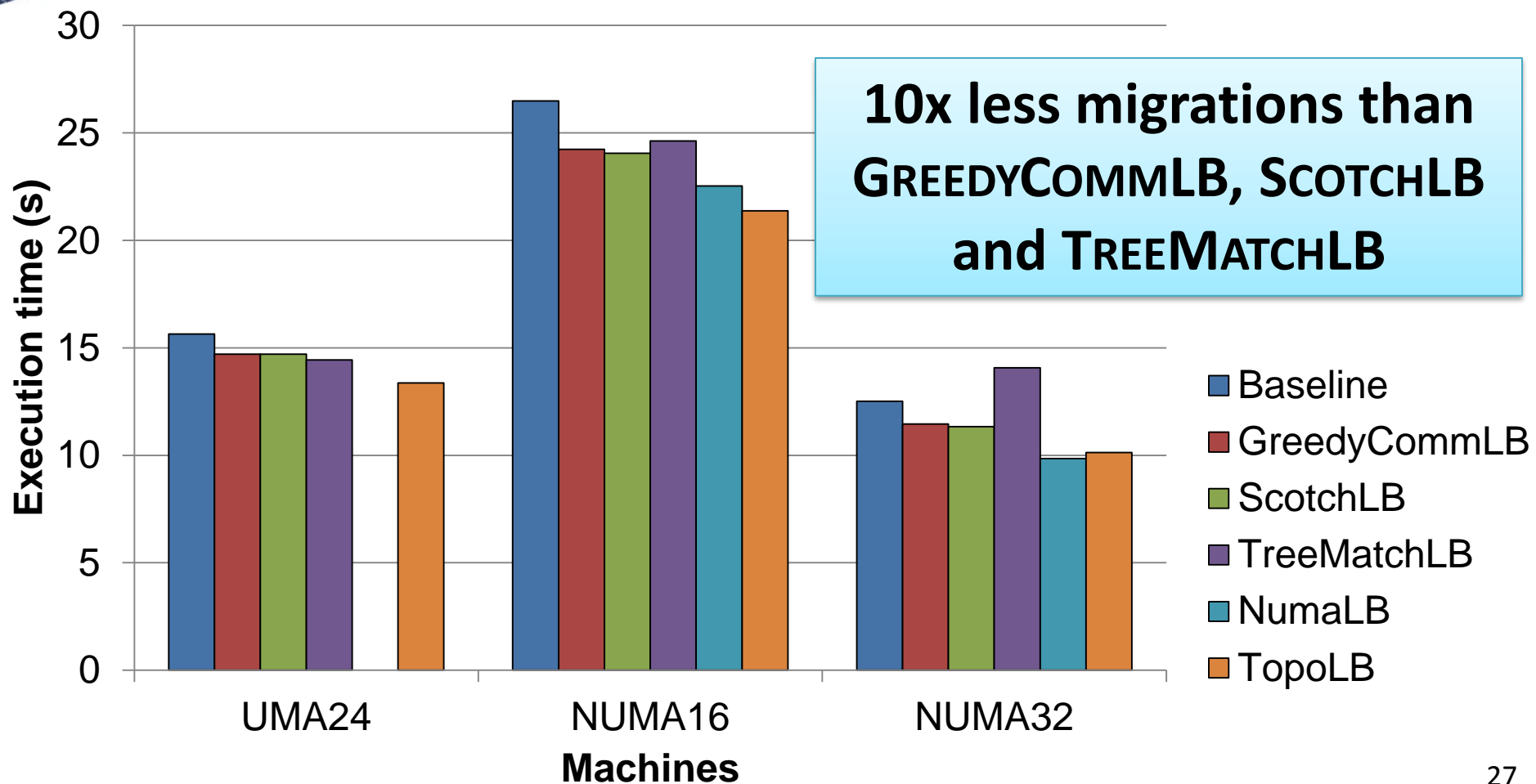


# lb\_test





# lb\_test

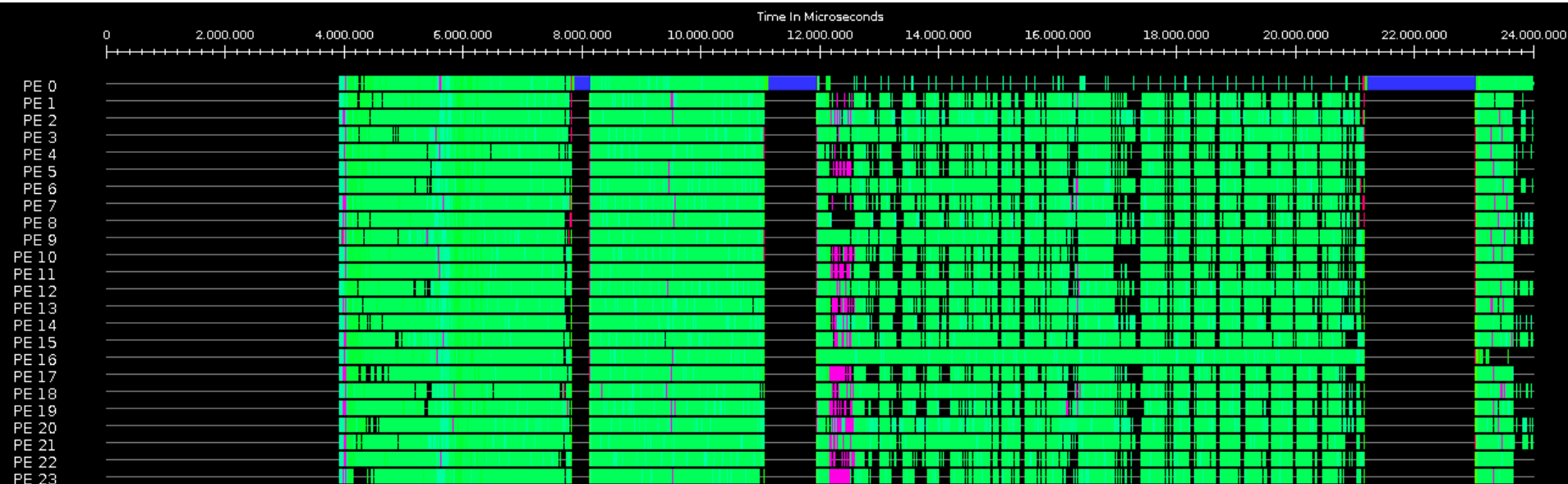


# mol3D

- Apolipoprotein-A1
- **No load balancer** improved performance
- **TOPOLB** performed **30x less migrations** than other load balancers
  - But took **2 to 3x more time on its load balancing decisions**

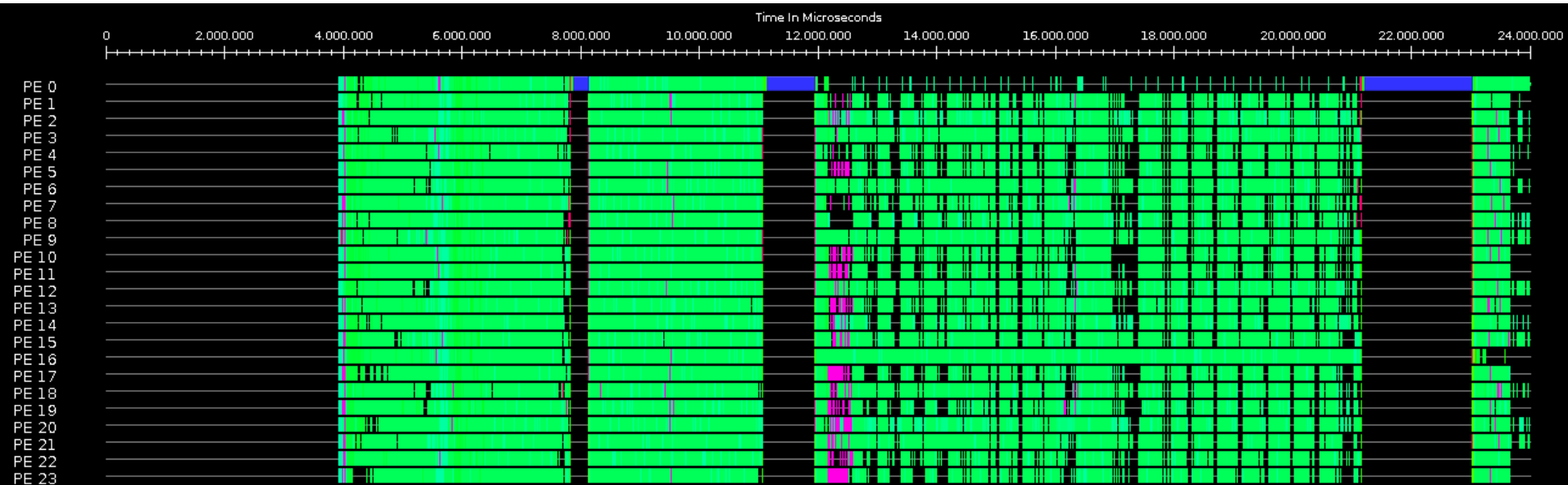
# mol3D

- View using Projections
  - 3x20 iterations, 3 calls to TopoLB
  - ~7000 chares



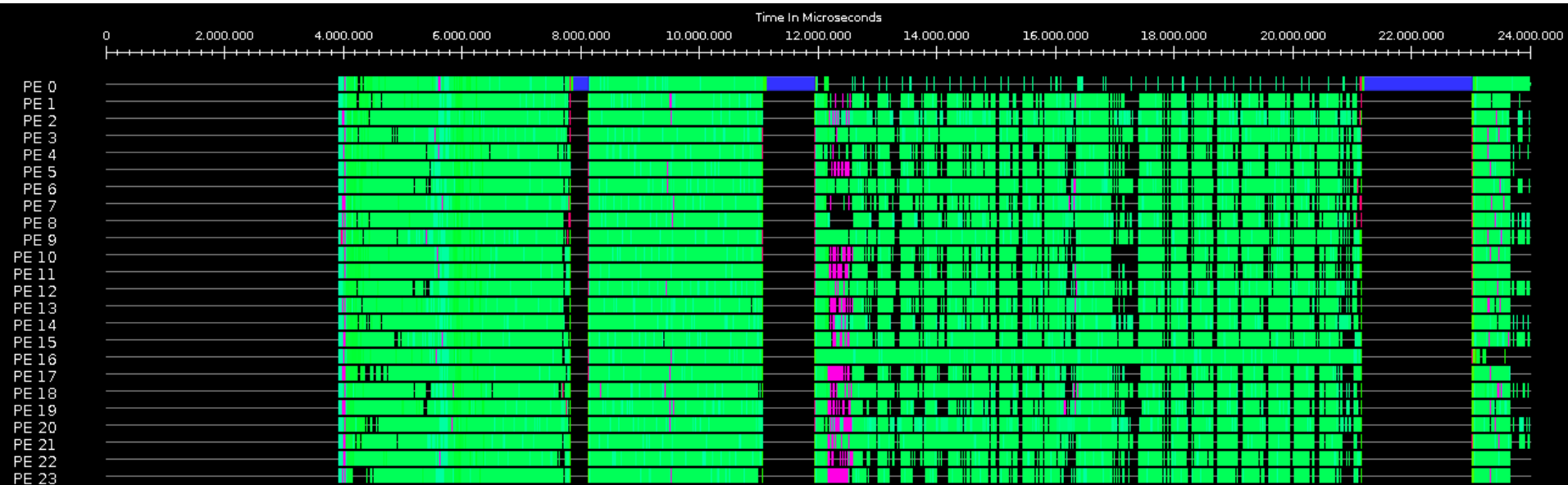
# mol3D

Application starts balanced



# mol3D

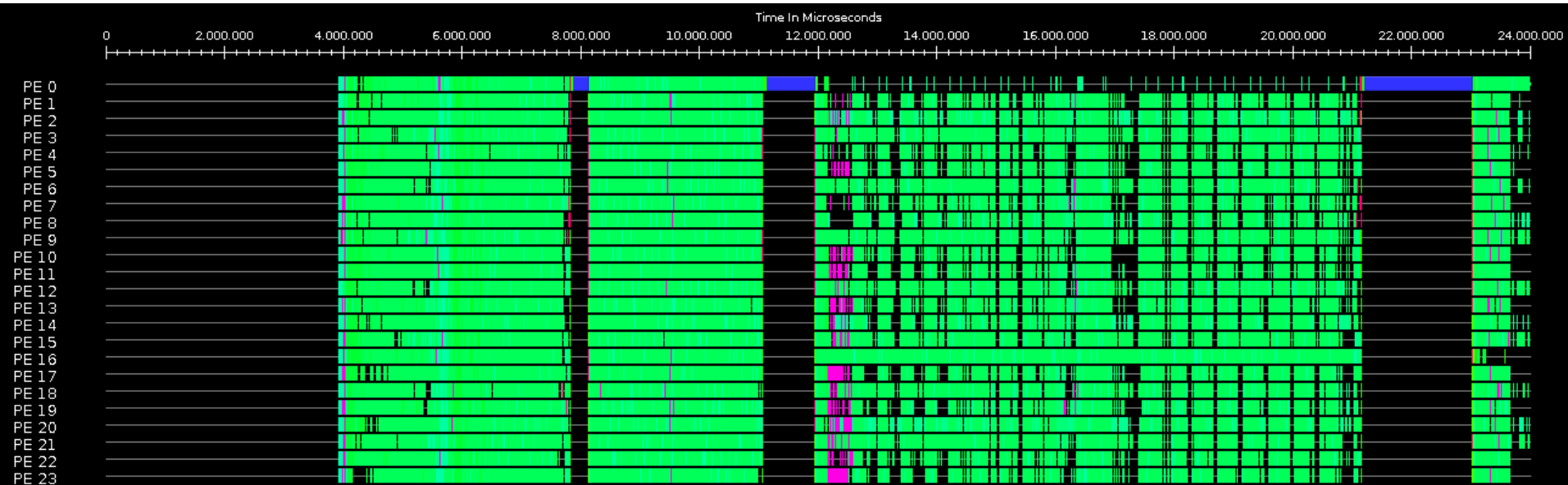
TOPOLB improves the iteration time





# mol3D

Iteration time diverges strongly





# Agenda

~~Motivation~~

~~TOPOLB~~

~~Experiments~~

**And beyond**

# And beyond

## (concluding remarks)

- TOPOLB presented a
  - Small number of migrations
  - Fast convergence
  - **High computational cost** (overhead)
    - Reduced by the exponential backoff

# And beyond

- Working to extend TOPOLB to **clusters of multi-core machines**
  - Two levels: Another LB for the cluster, TOPOLB for the compute nodes
    - TOPOLB can be **too costly** for large machines and applications

## And beyond

- Working to **better understand and evaluate** the behavior of CHARM++ **load balancers**
  - Debugging library
  - Use real applications
- Added part of our machine model to HWLOC

# Topology-aware load balancing for parallel applications on multi-core systems and beyond

Thank you.

Laércio Lima Pilla

Contact: [llpilla@inf.ufrgs.br](mailto:llpilla@inf.ufrgs.br)