# Hints to improve automatic load balancing with LeWI for hybrid applications

Marta Garcia, Jesus Labarta, Julita Corbalan

- Loss of efficiency
- Hybrid programming models ($MPI + X$)
- Manual tuning of parallel codes (load-balancing, data redistribution)

# The X (in this paper)

## OpenMP

- Directives to annotate parallel code
- Fork/join model with shared memory
- Number of threads may change *between* parallel regions
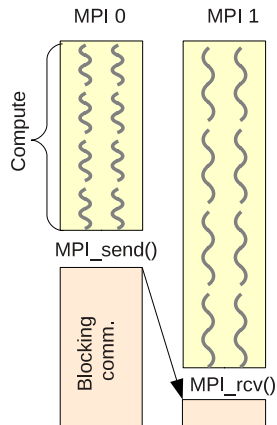
## SMPSs (SMPSuperscalar)

- Task as basic element
- Annotate *taskifiable* functions and their parameters (in/out/inout)
- Task graph to track dependencies
- Number of threads may change *any time*

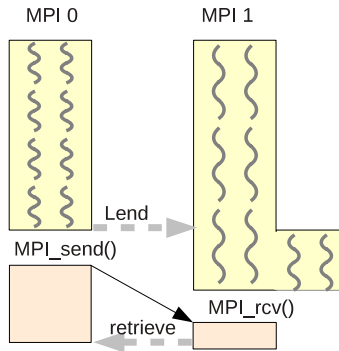# DLB and LeWI

## DLB (dynamic load balancing)

- "*Runtime interposition to [...] intercept MPI calls*"
- Balance load on the inner level (OpenMP/SMPSs)
- Several load balancing algorithms
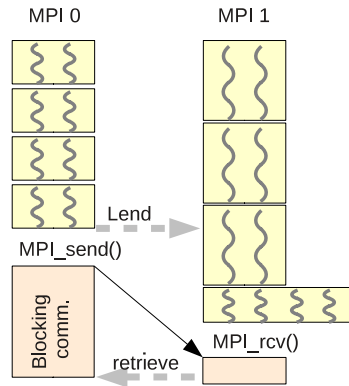
## LeWI (Lend CPU when Idle)

- CPUs of rank in blocking MPI call are idle
- Lend CPUs to other ranks and recover them after MPI call completes

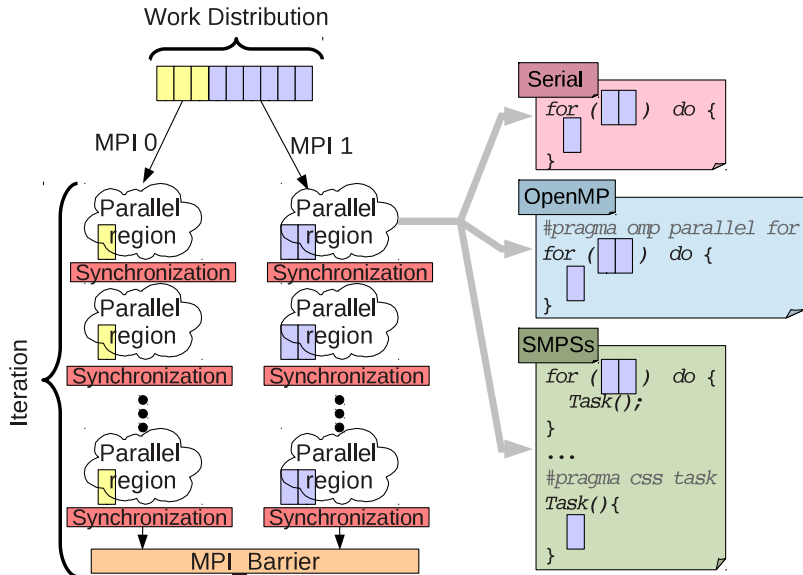(a) No load balancing.  (b) LeWI algorithm with SMPSs.  (c) LeWI algorithm with OpenMP.

# Approach

- "*Extensive performance evaluation*"
- "*Modeling parallelization characteristics that limit the automatic load balancing potential*"
- "*Improving automatic load balancing*"

# Performance evaluation

- - Marenostrum 2: $2 \times$ IBM PowerPC 970MP (2 cores); 8 GiB RAM
  - Linux 2.6.5-7.244-pseries64; MPICH; IBM XL C/C++ compiler w/o optimizations
- Metrics
  - $Speedup = \frac{parallel\_execution\_time}{serial\_execution\_time}$
  - $Efficiency = \frac{useful\_cpu\_time}{elapsed\_time * cpus}$ where
    $useful\_cpu\_time = cpu\_time - (mpi\_time + openmp/smpss\_time + dlb\_time)$
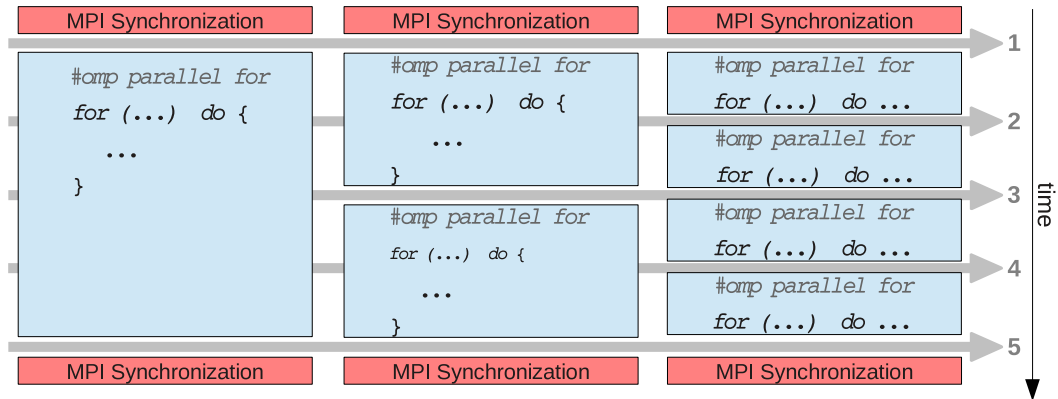  - $CPUs\_used$ to simultaneously run application code
- 3 benchmarks $+$ 2 real applications

# PILS (Parallel ImbaLance Simulation)

- Synthetic benchmark
- Core: "*floating point operations without data involved*"
- Tunable parameters
    - Programming model (MPI, MPI + OpenMP, MPI + SMPSs)
    - Load distribution
    - Parallelism grain ($= \frac{1}{\#parallel\ regions}$)
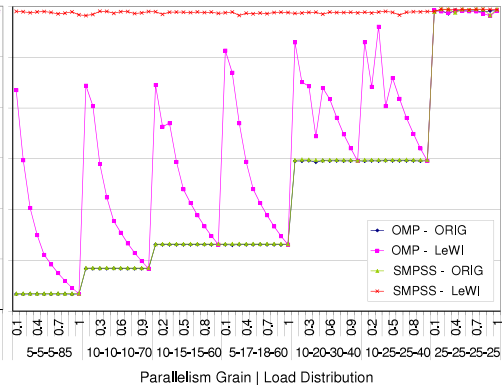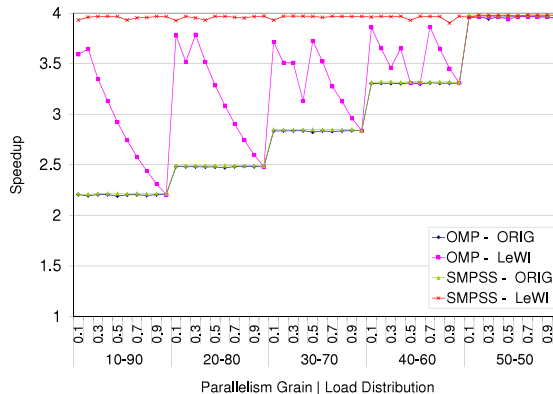    - Iterations

a) Parallelism Grain: 1. b) Parallelism Grain: 0.5. c) Parallelism Grain: 0.25.

## Other Codes

- Benchmarks
  - BT-MZ: block tri-diagonal solver
  - LUB: LU matrix factorization
- Applications
  - Gromacs: molecular dynamics, MPI-only
  - Gadget: cosmological N-body/SPH (smoothed-particle hydrodynamics) simulation

# Other Codes

- Benchmarks
  - BT-MZ: block tri-diagonal solver
  - LUB: LU matrix factorization
- Applications
  - Gromacs: molecular dynamics, MPI-only
  - Gadget: cosmological N-body/SPH (smoothed-particle hydrodynamics) simulation

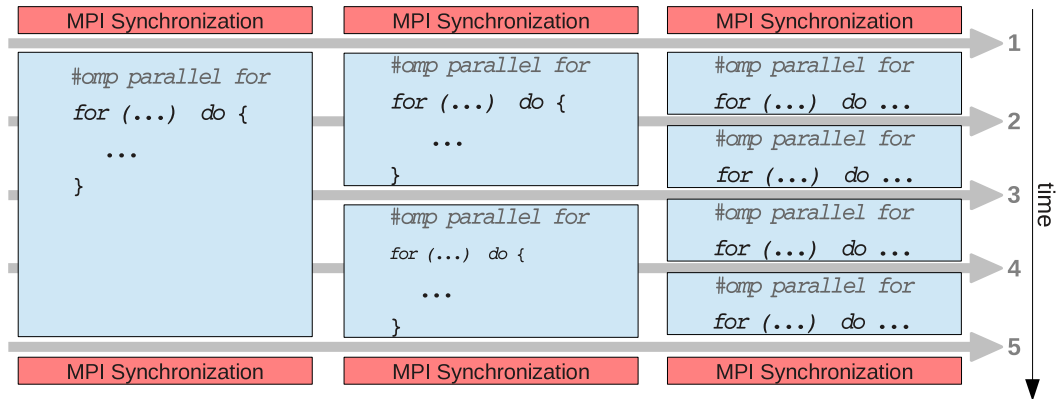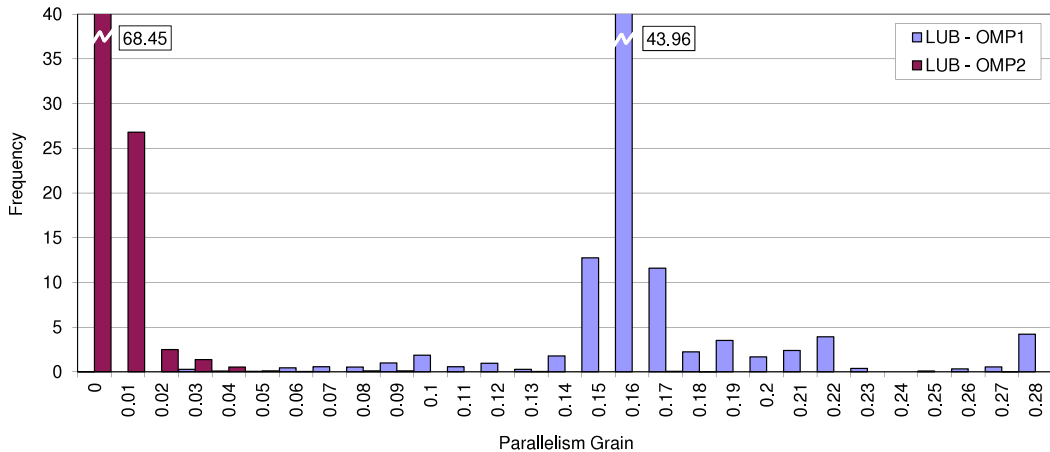| Application | Original version | MPI + OpenMP | MPI + SMPSs | Executed in nodes (cpus) |
|---|---|---|---|---|
| PILS | MPI + OpenMP<br>MPI + SMPSs | X | X | 1 (4) |
| BT-MZ | MPI + OpenMP | X | X | 1, 2, 4 (4, 8, 16) |
| LUB | MPI + OpenMP<br>MPI + SMPSs | X | X | 1, 2, 4 (4, 8, 16) |
| Gromacs | MPI | | X | 1.64 (4.256) |
| Gadget | MPI | X | | 200 (800) |

- "*Parallelism Grain in OpenMP applications*"
- "*Task duration in SMPSs applications*"
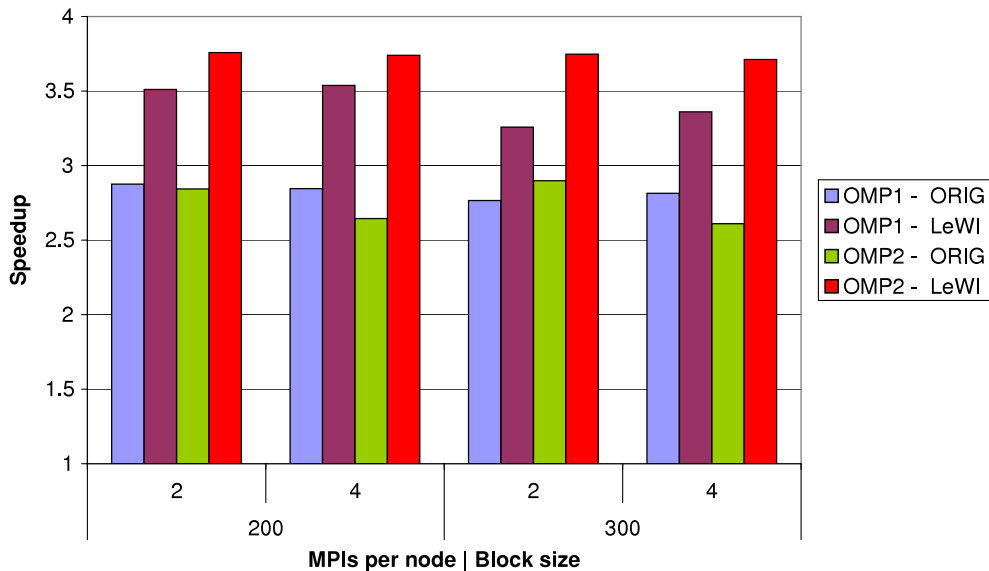- "*Distribution of MPI processes among computation nodes*"

a) Parallelism Grain: 1. b) Parallelism Grain: 0.5. c) Parallelism Grain: 0.25.

# Conclusion

## Summary

- DLB/LeWI can improve performance transparently
- Inter-node load imbalances not handled
- Granularity of parallelism and placement as important factors
- Optimal configuration with vs. without DLB/LeWI

# Conclusion

## Summary

- DLB/LeWI can improve performance transparently
- Inter-node load imbalances not handled
- Granularity of parallelism and placement as important factors
- Optimal configuration with vs. without DLB/LeWI

## Discussion

- Interaction with MPI
- Benchmarks (1.5 of 3 NPB-MZ, arbitrary load distribution)
- How to find "the right" granularity