

Performance Barriers in highly scaling earth-system models

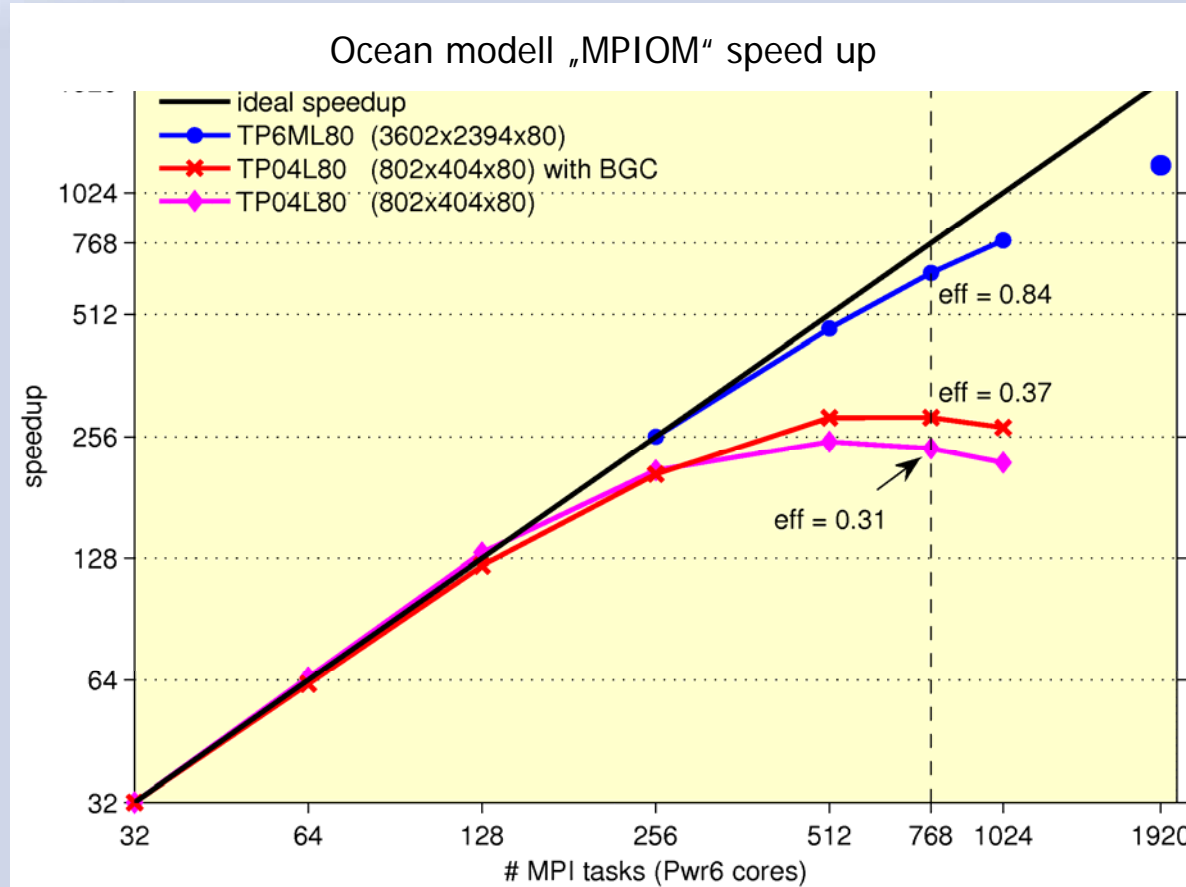
Panagiotis Adamidis
Deutsches Klimarechenzentrum GmbH

Jörg Behrens, Thomas Jahns (DKRZ),
Florian Wilhelm (KIT)

Bottlenecks

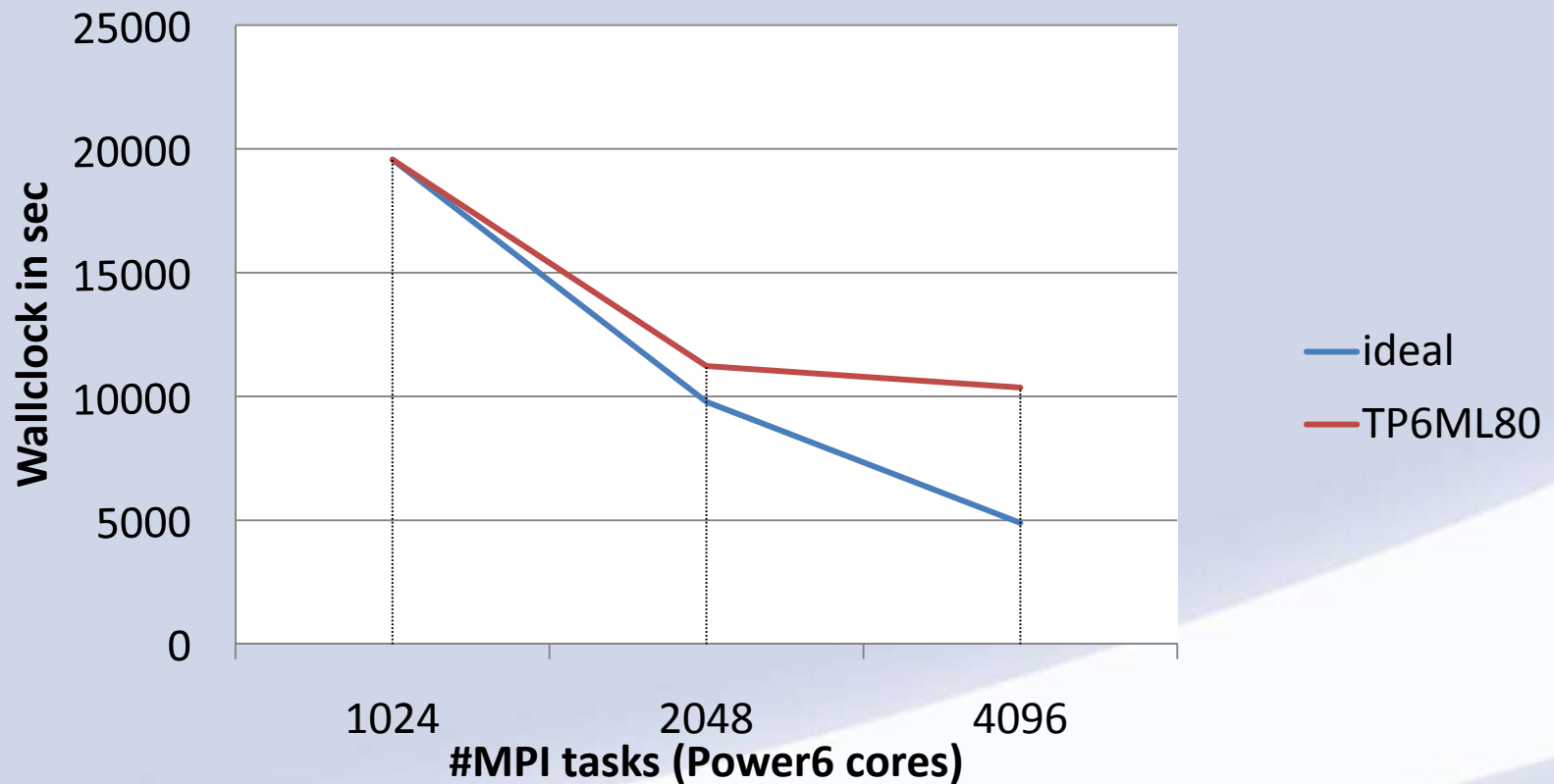
- Bottlenecks of Massively Parallel Computing Systems
 - Communication Network
 - Memory Bandwidth
 - Idle Processors

Scaling today

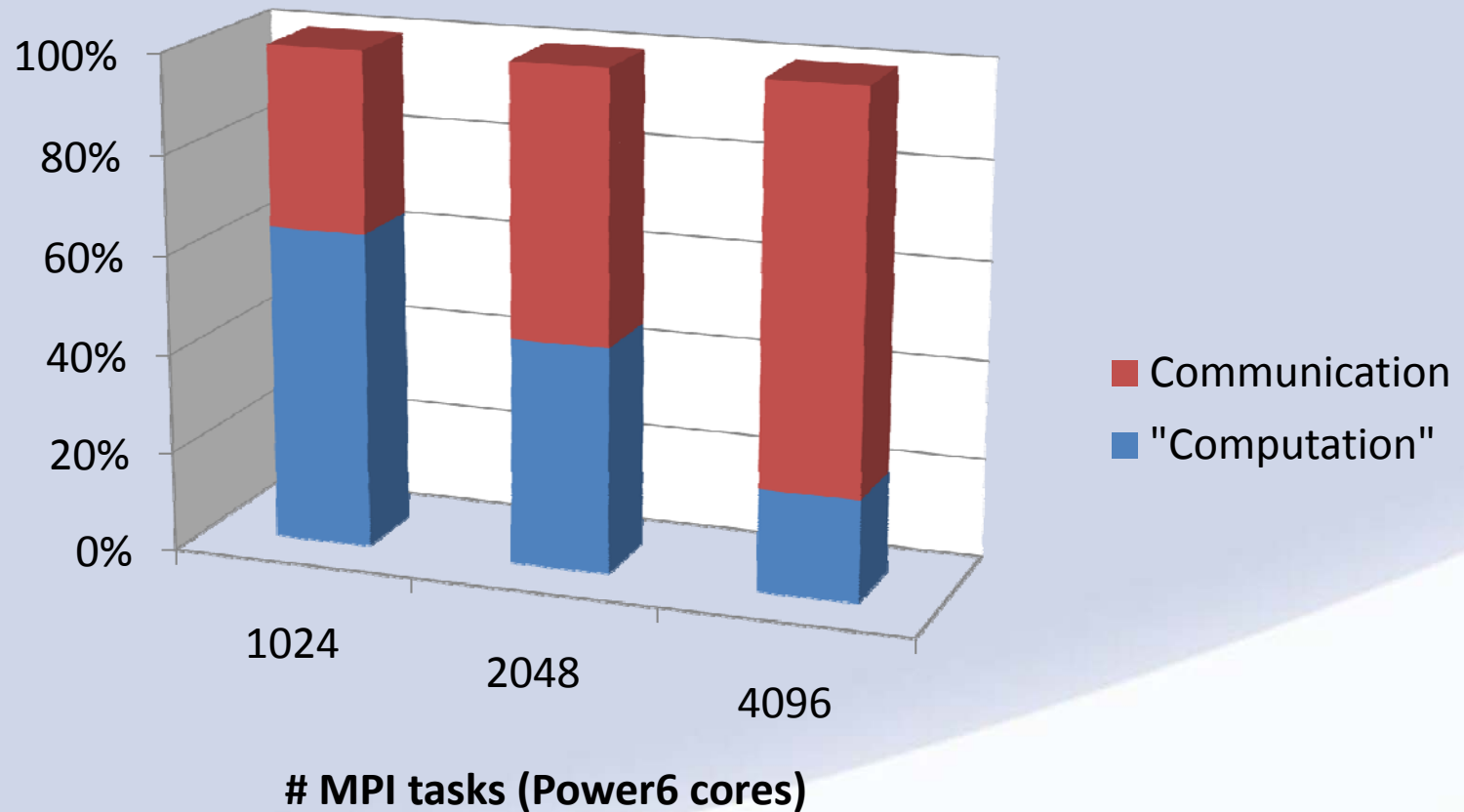


Simulating 1 month with $dt=600$

Ocean Modell MPIOM Wallclock

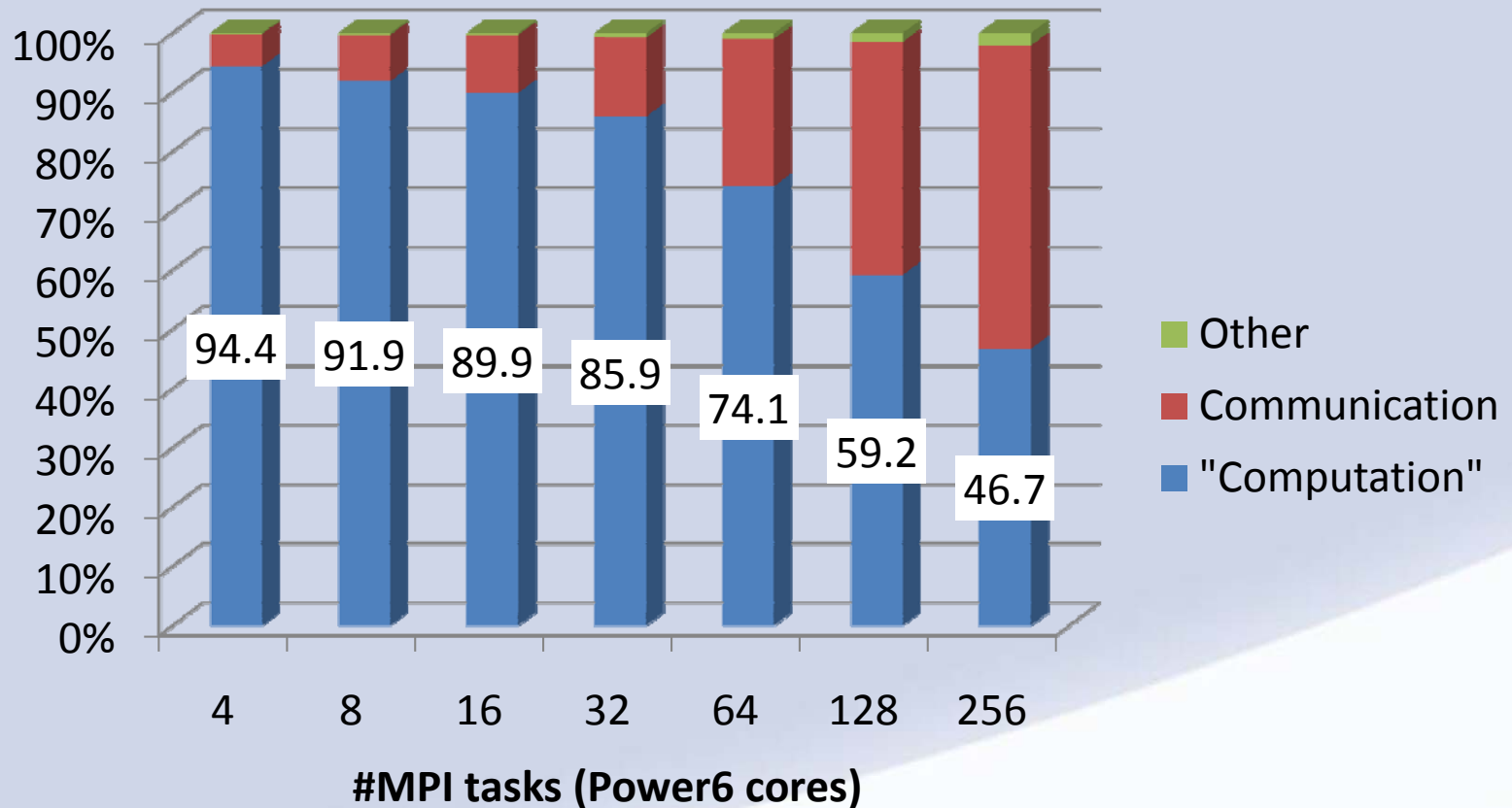


Computation vs Communication TP6ML80

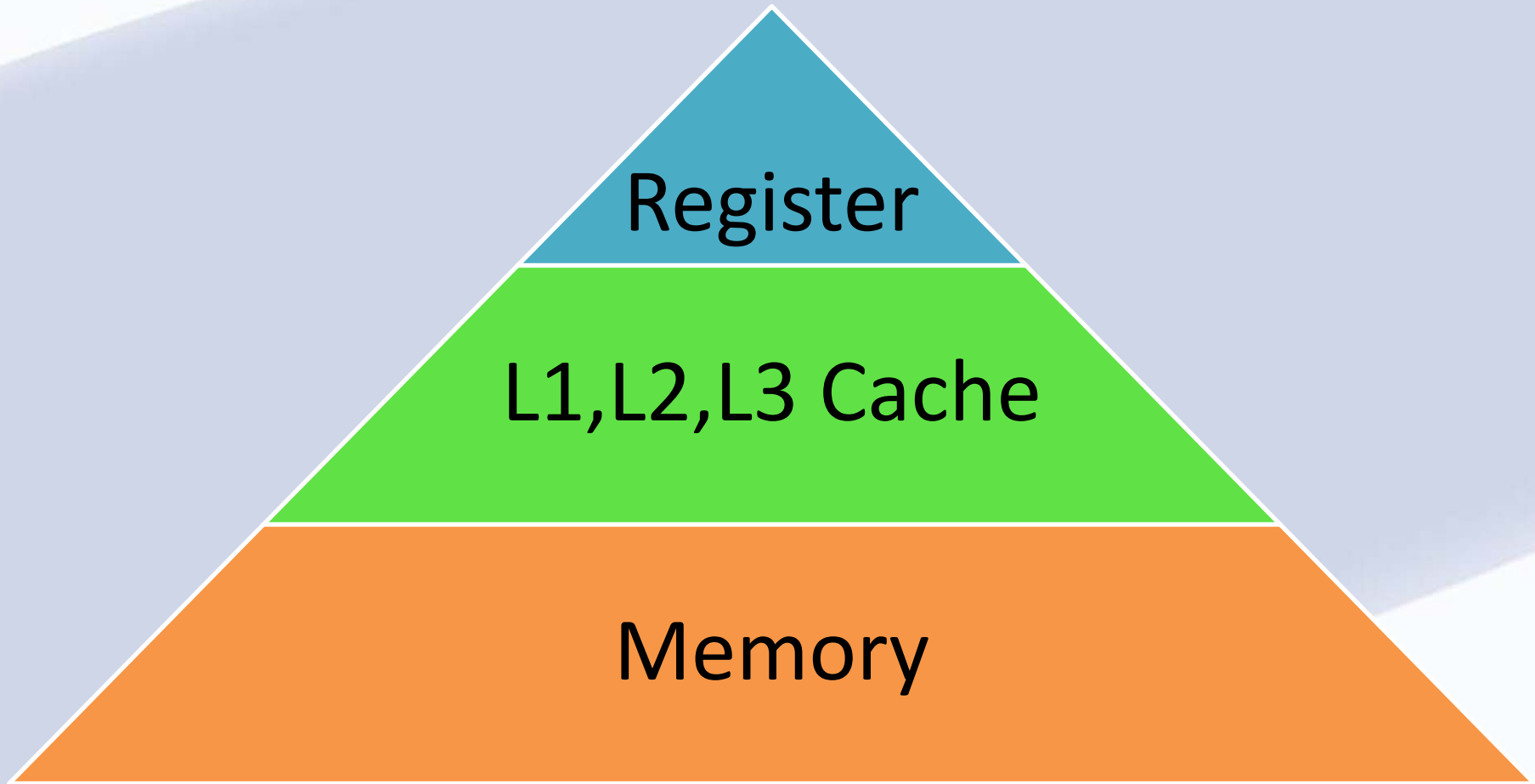




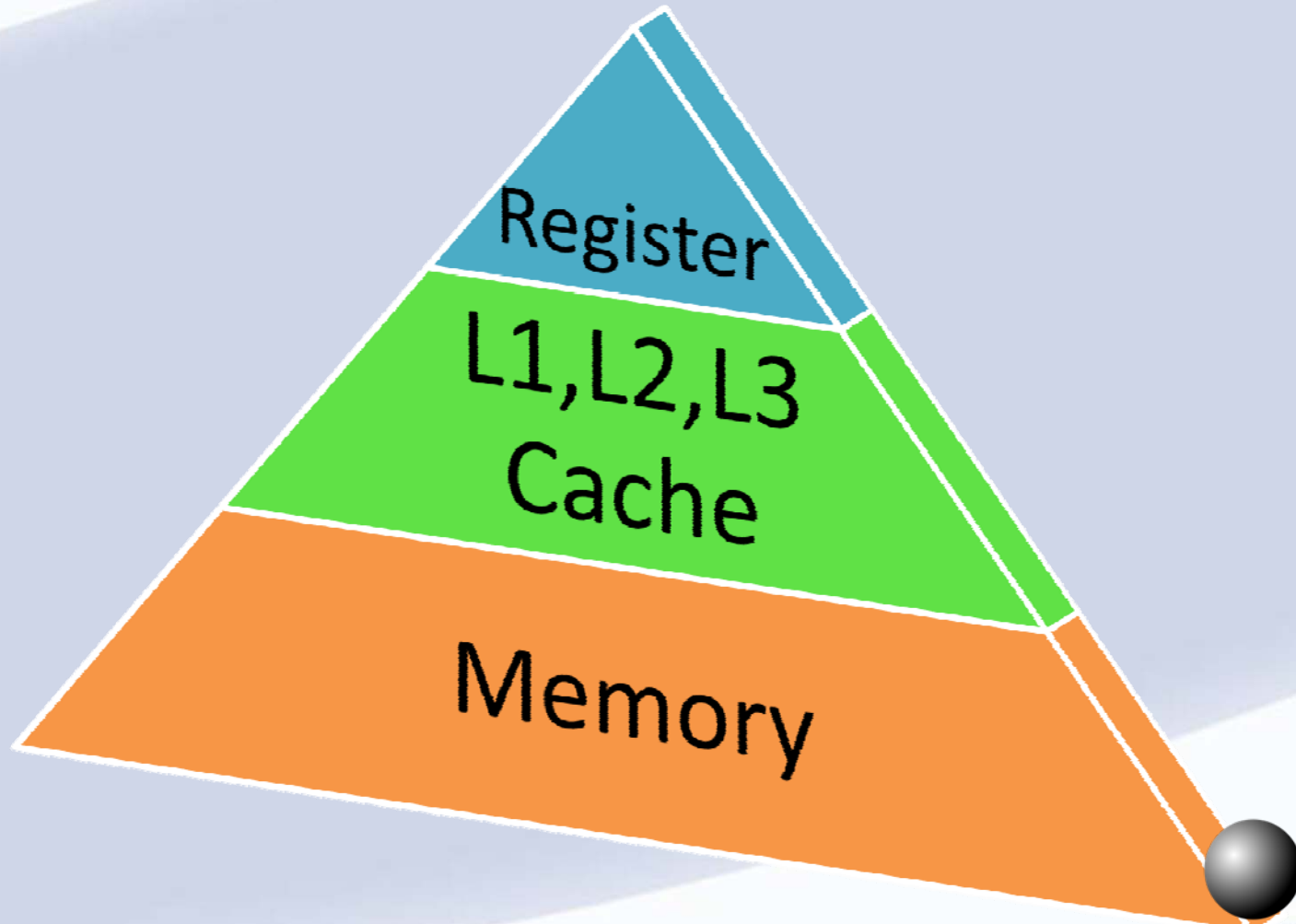
Computation vs Communication TP04L80



Memory Hierarchy

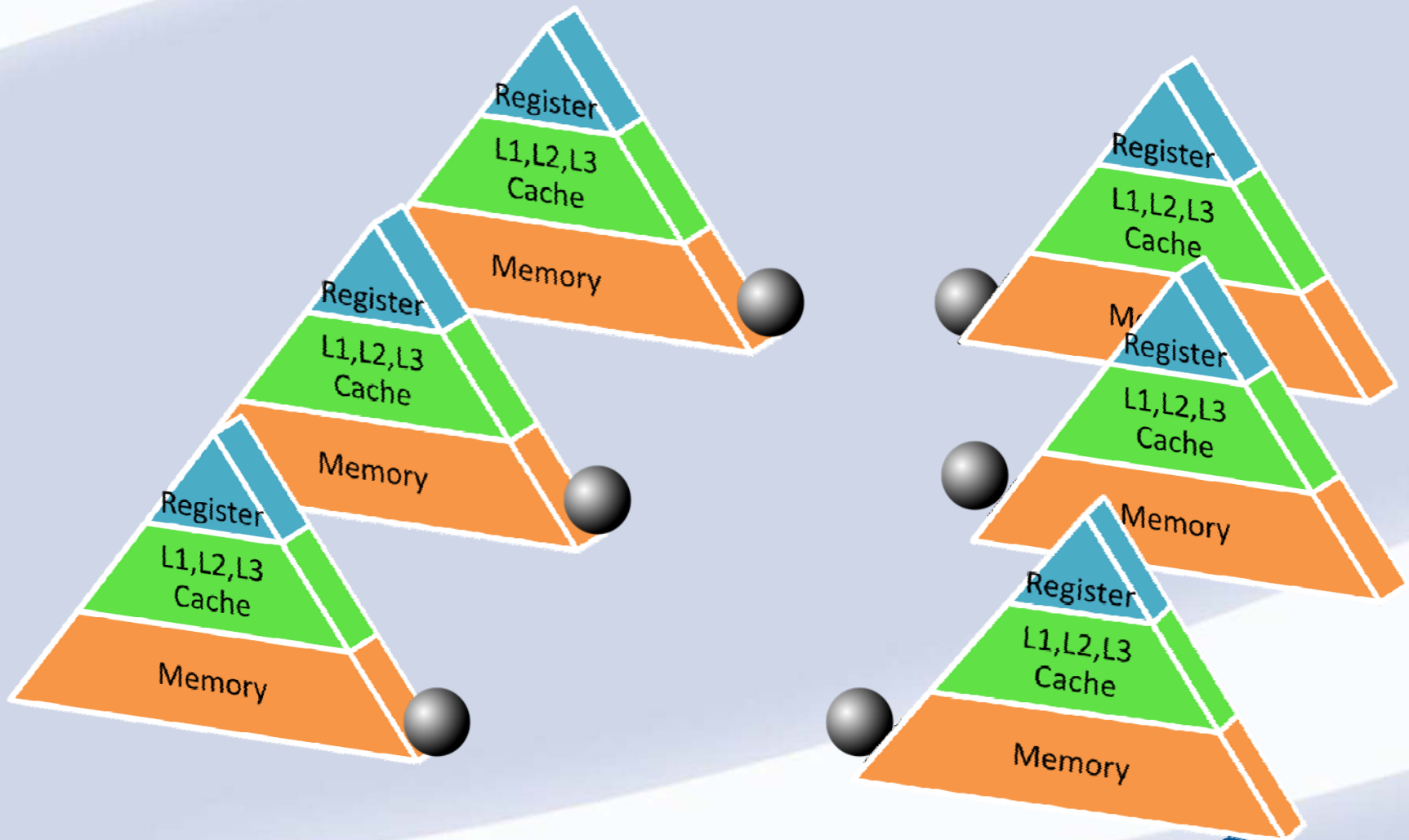


Data Movement





Data Movement in Parallel Systems



Sisyphean Challenge

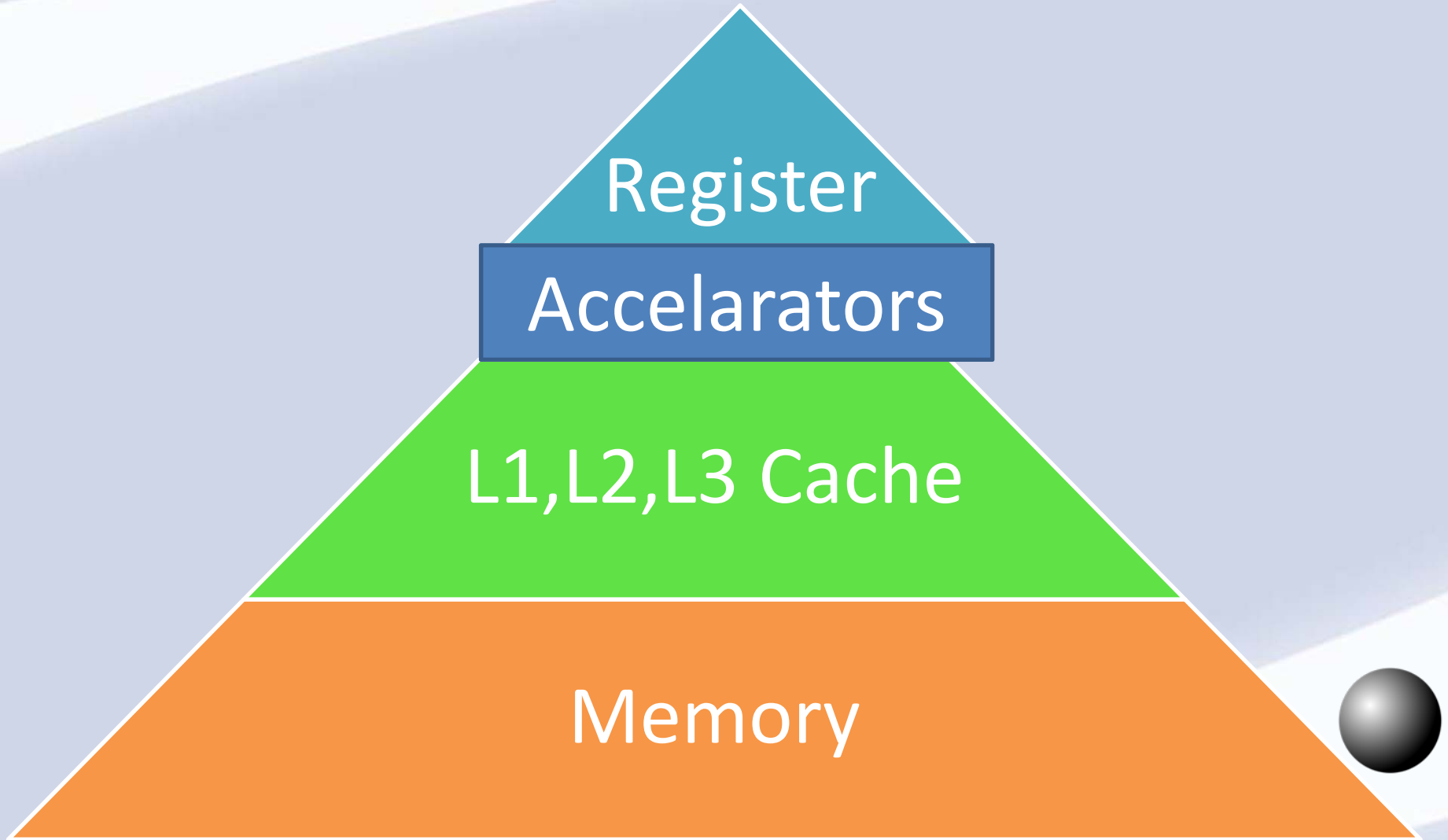
- According to greek mythology, “Sisyphus” was a greek King, famous for his cleverness.
- He played his tricks not only on humans but also on the Gods, thus being the craftiest of men
- He even managed to outwit “Thanatos” (death) himself

Sisyphean Challenge

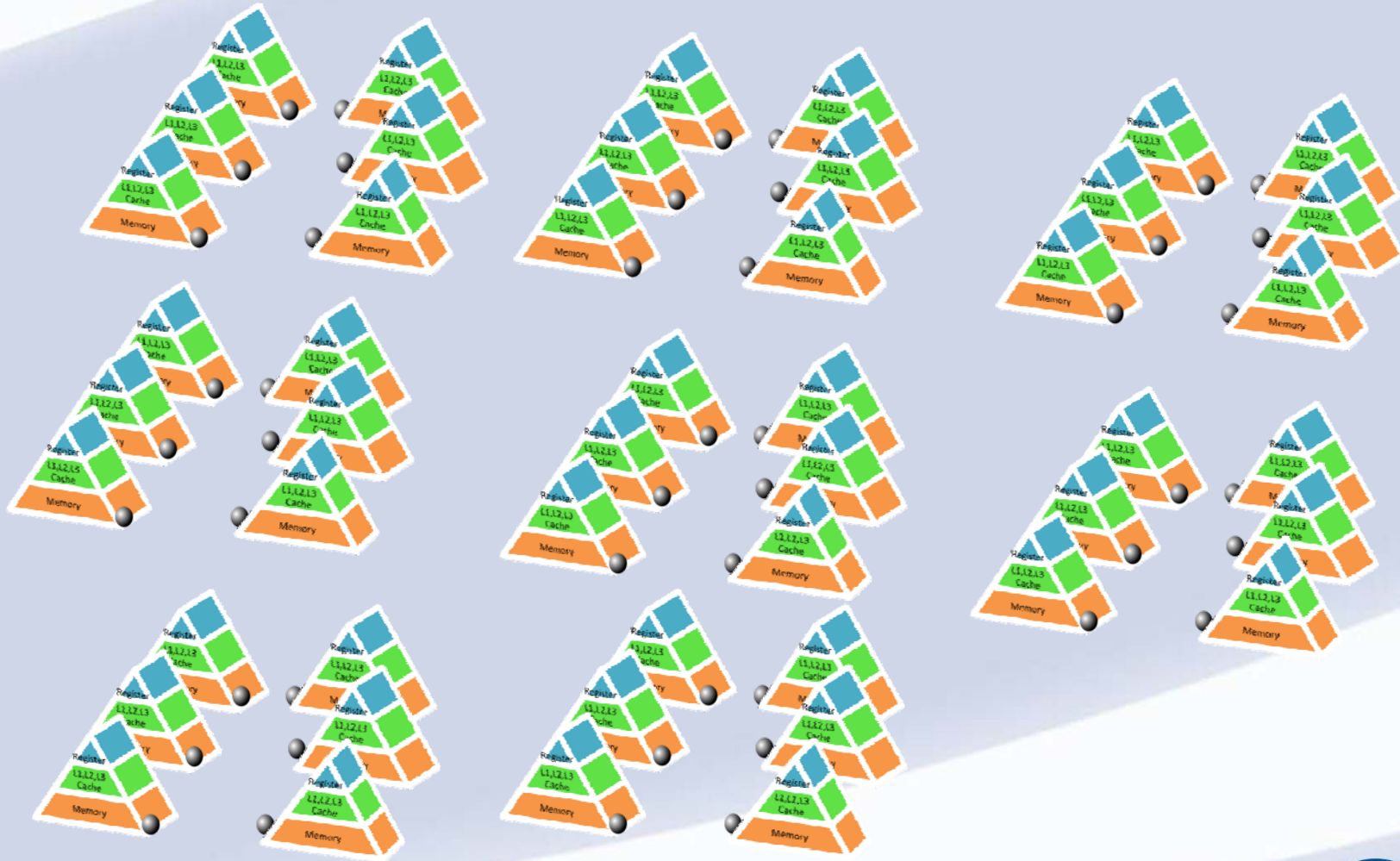
- As a punishment for his trickery, Sisyphus was compelled to roll a huge rock up a steep hill, only to watch it roll back down, and to repeat this throughout eternity



Future Processors ?



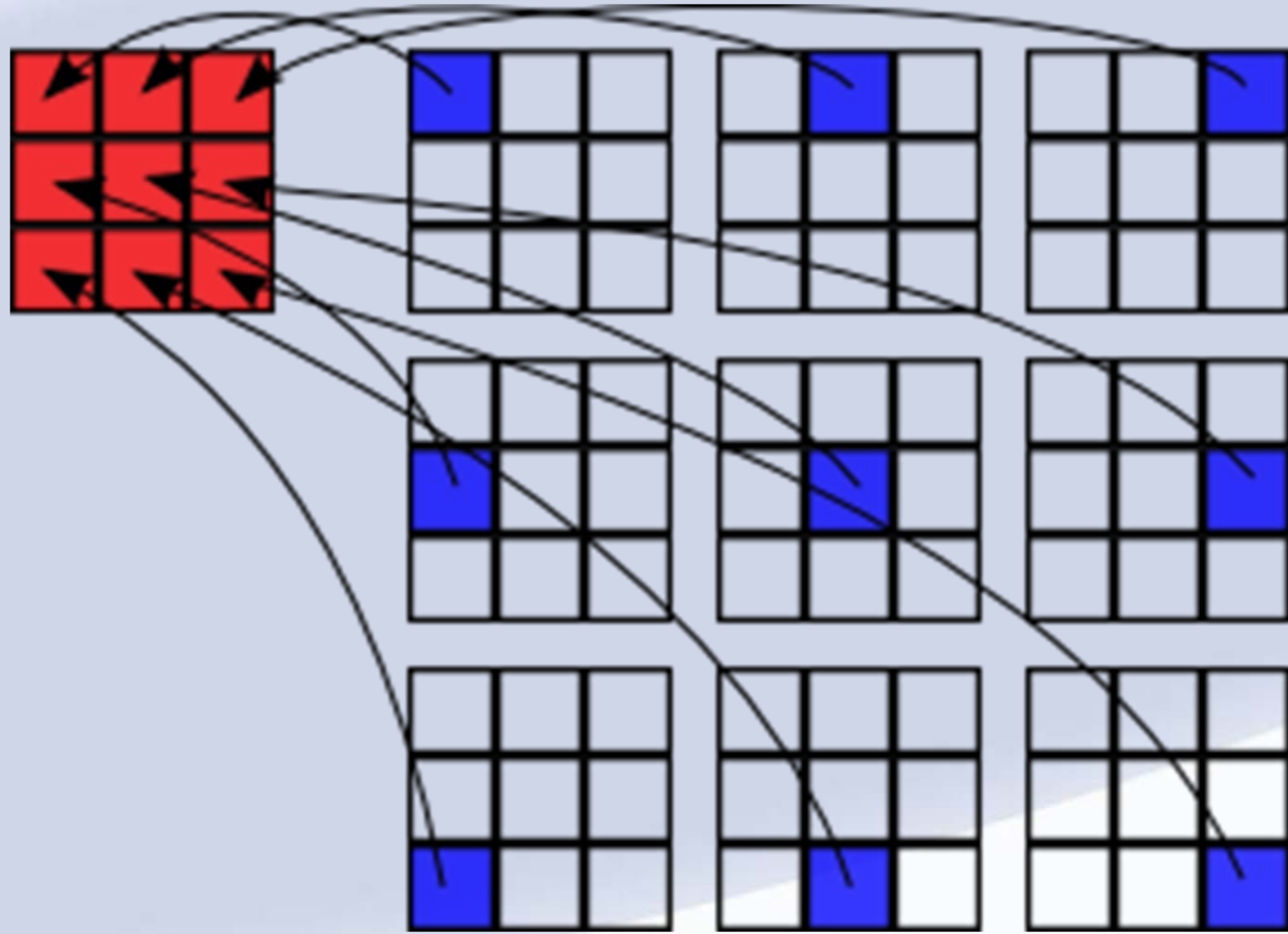
Future Parallel Systems

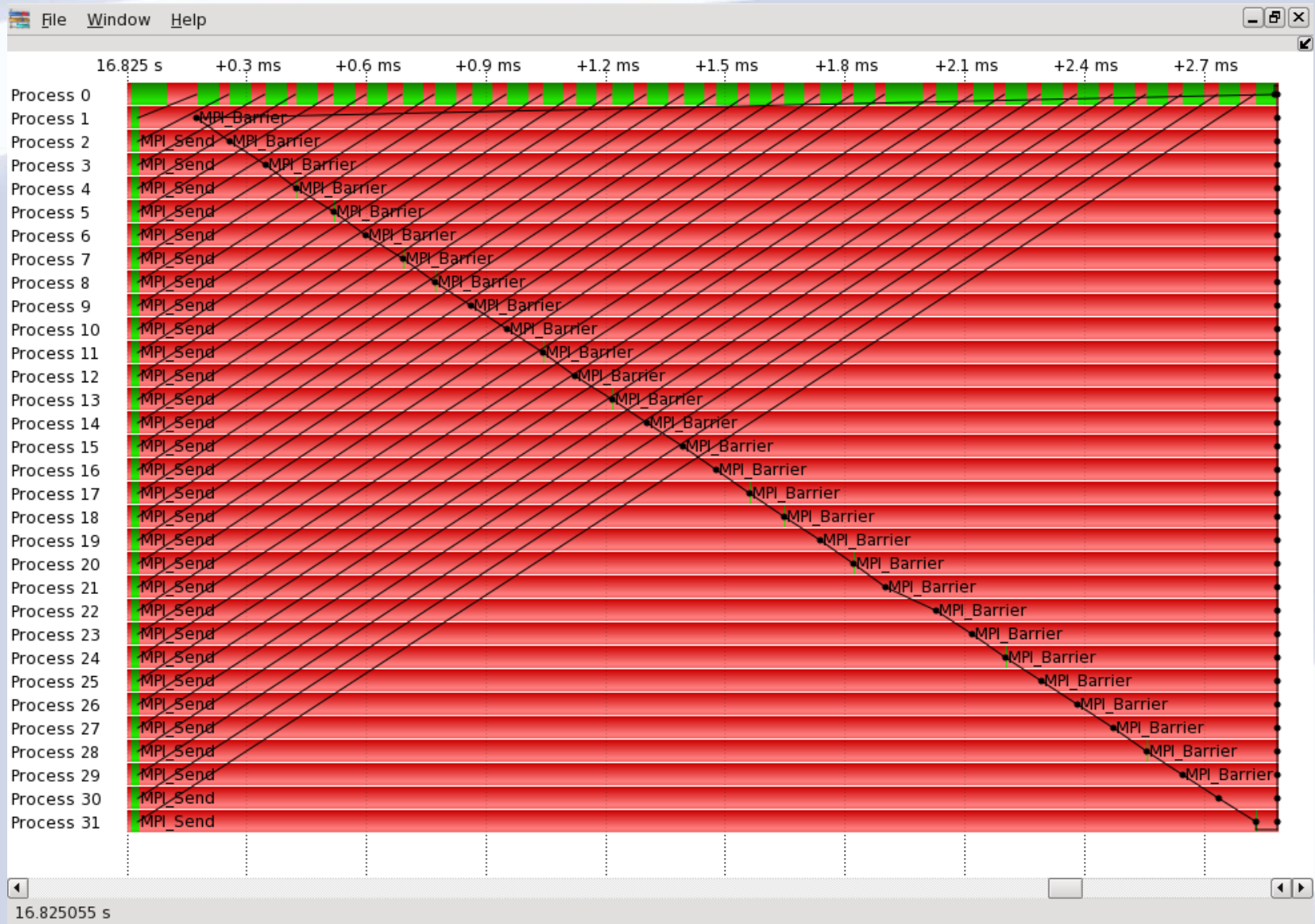


Approach

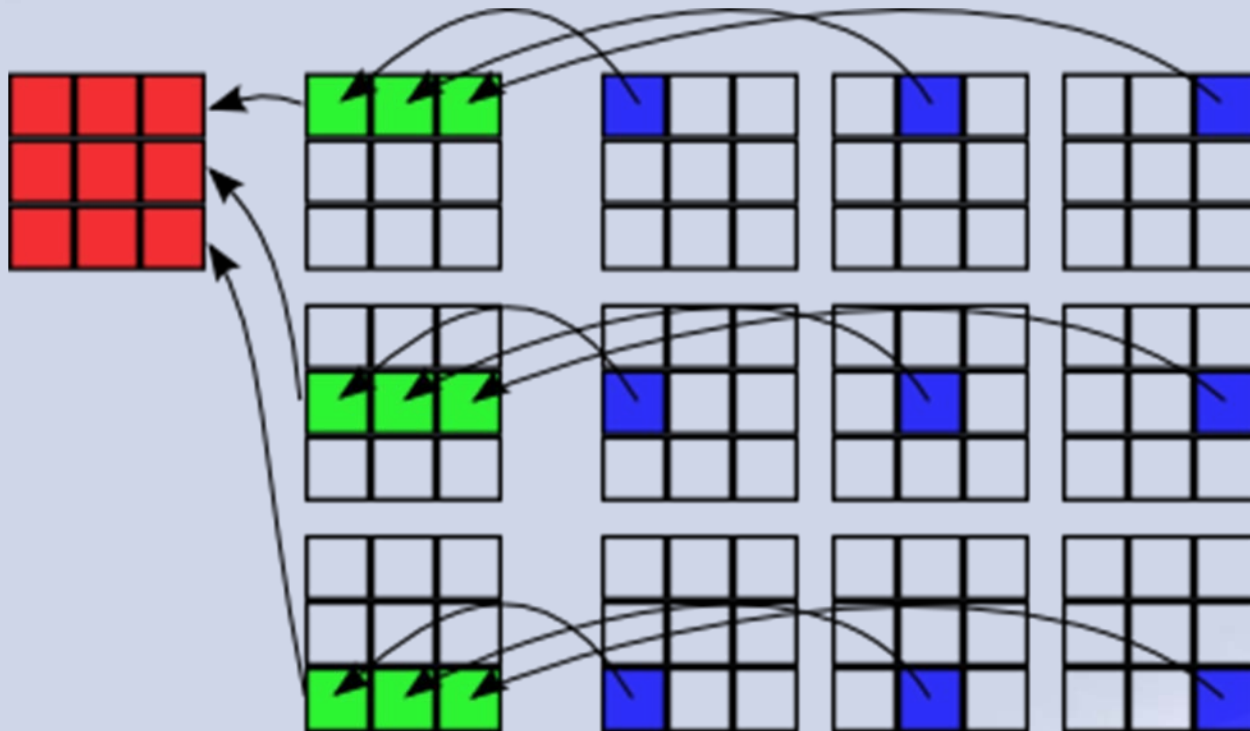
- Solutions
 - Optimizing Communication
 - Load Balancing
 - Algorithms which are better suited for the specific architecture

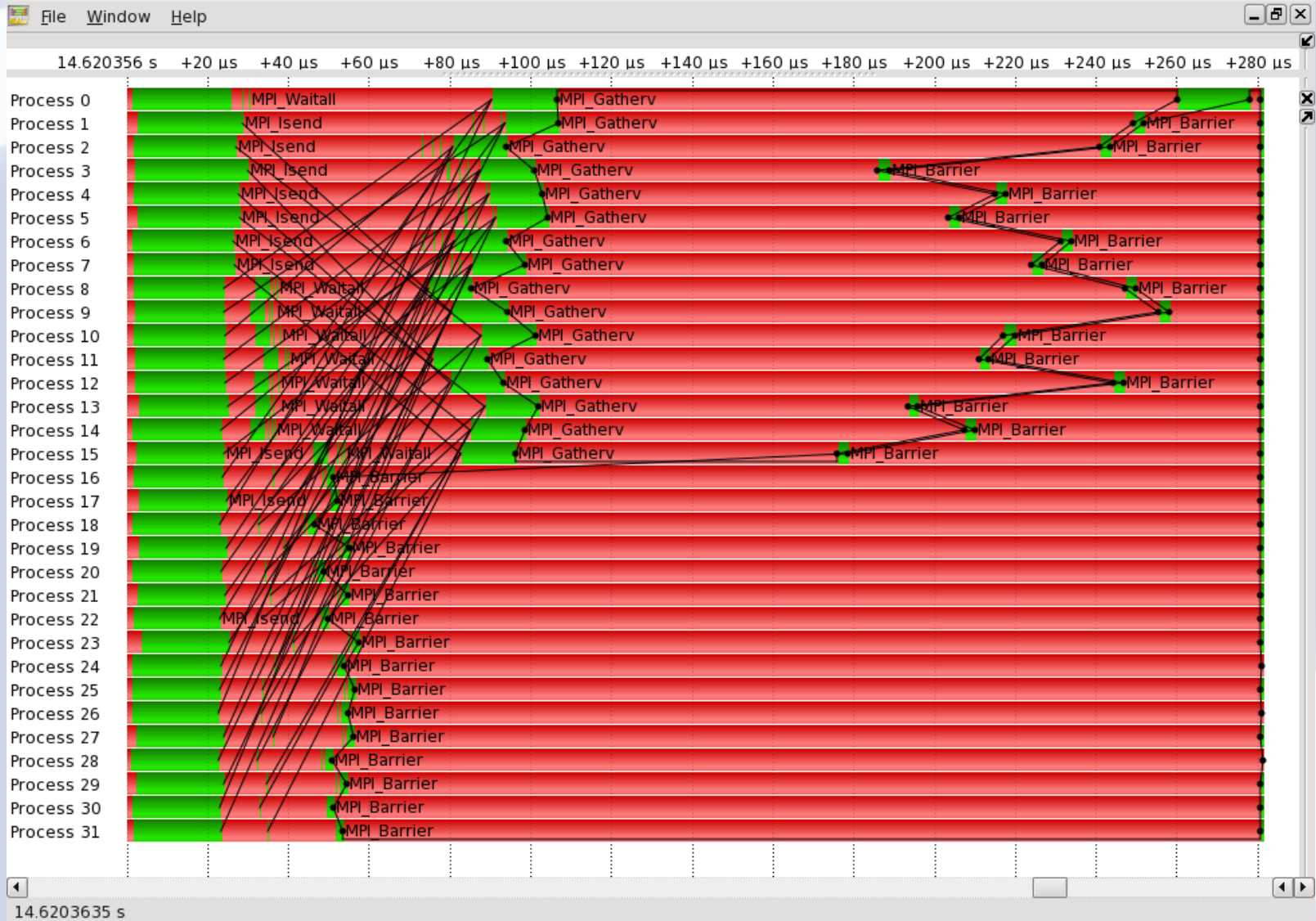
Optimizing Communication

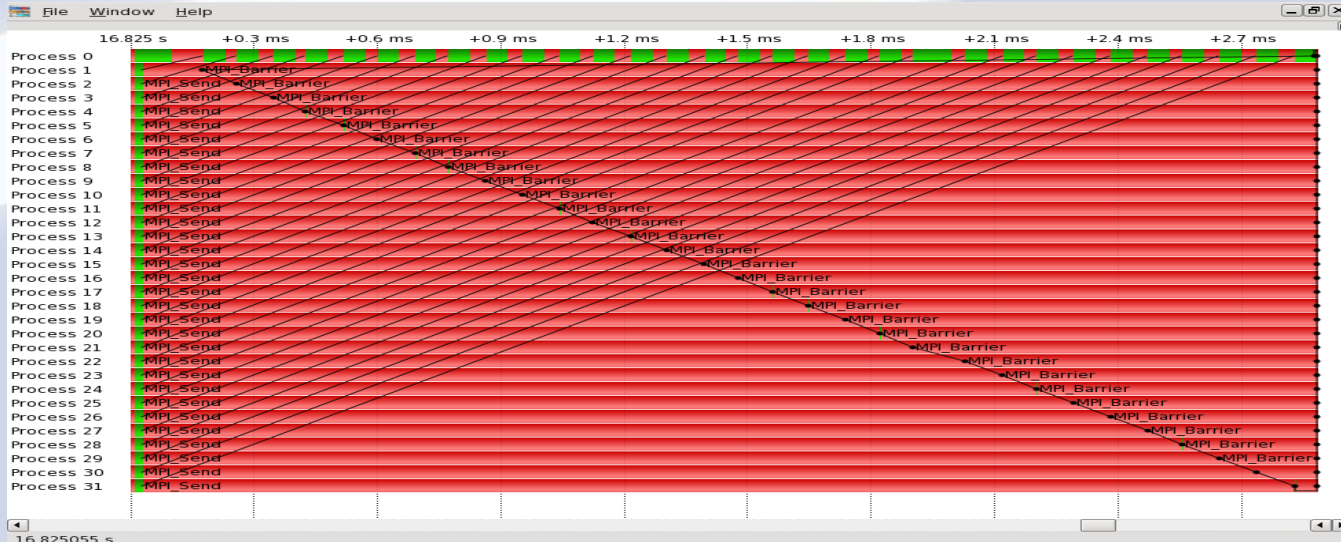




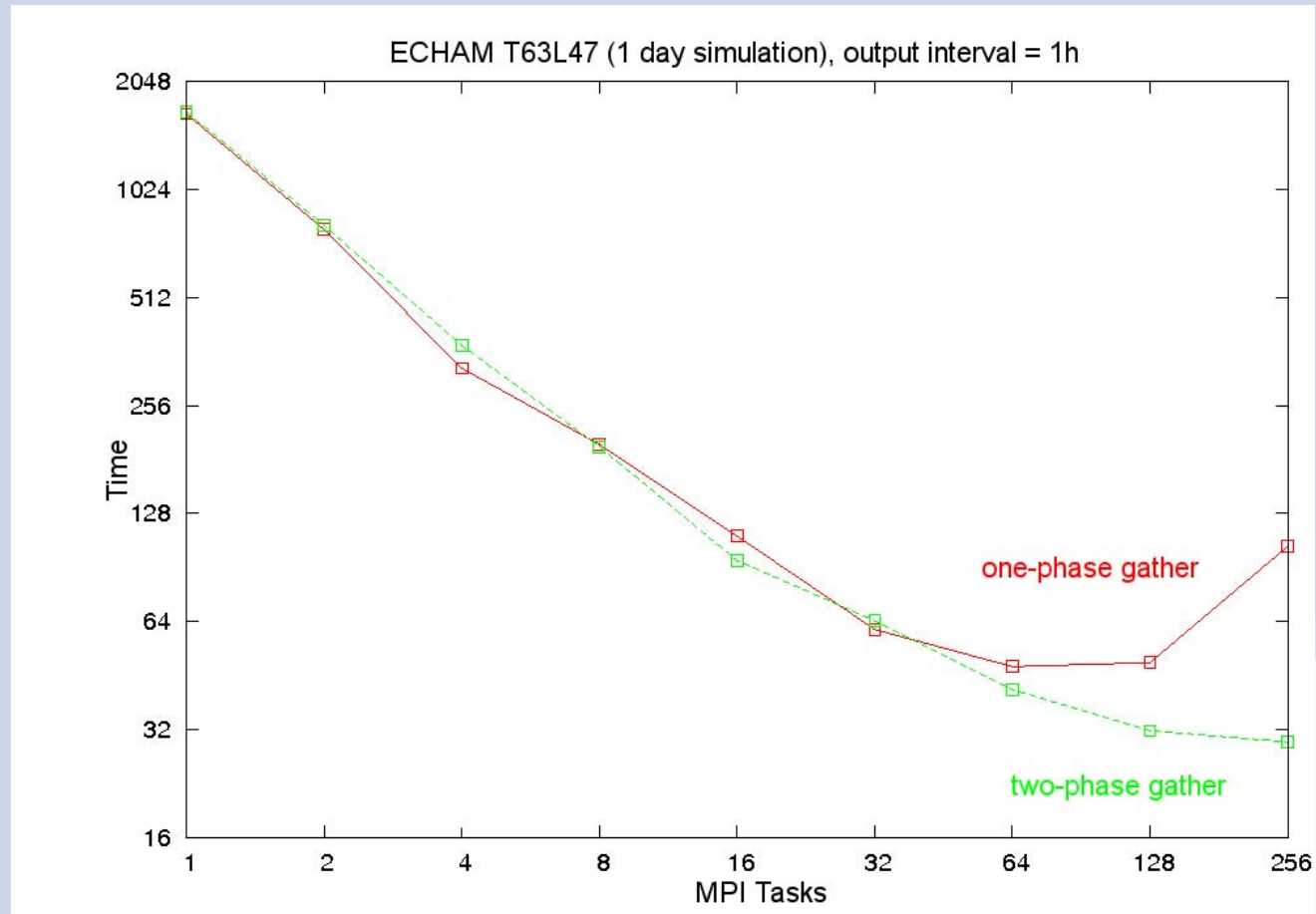
Optimizing Communication



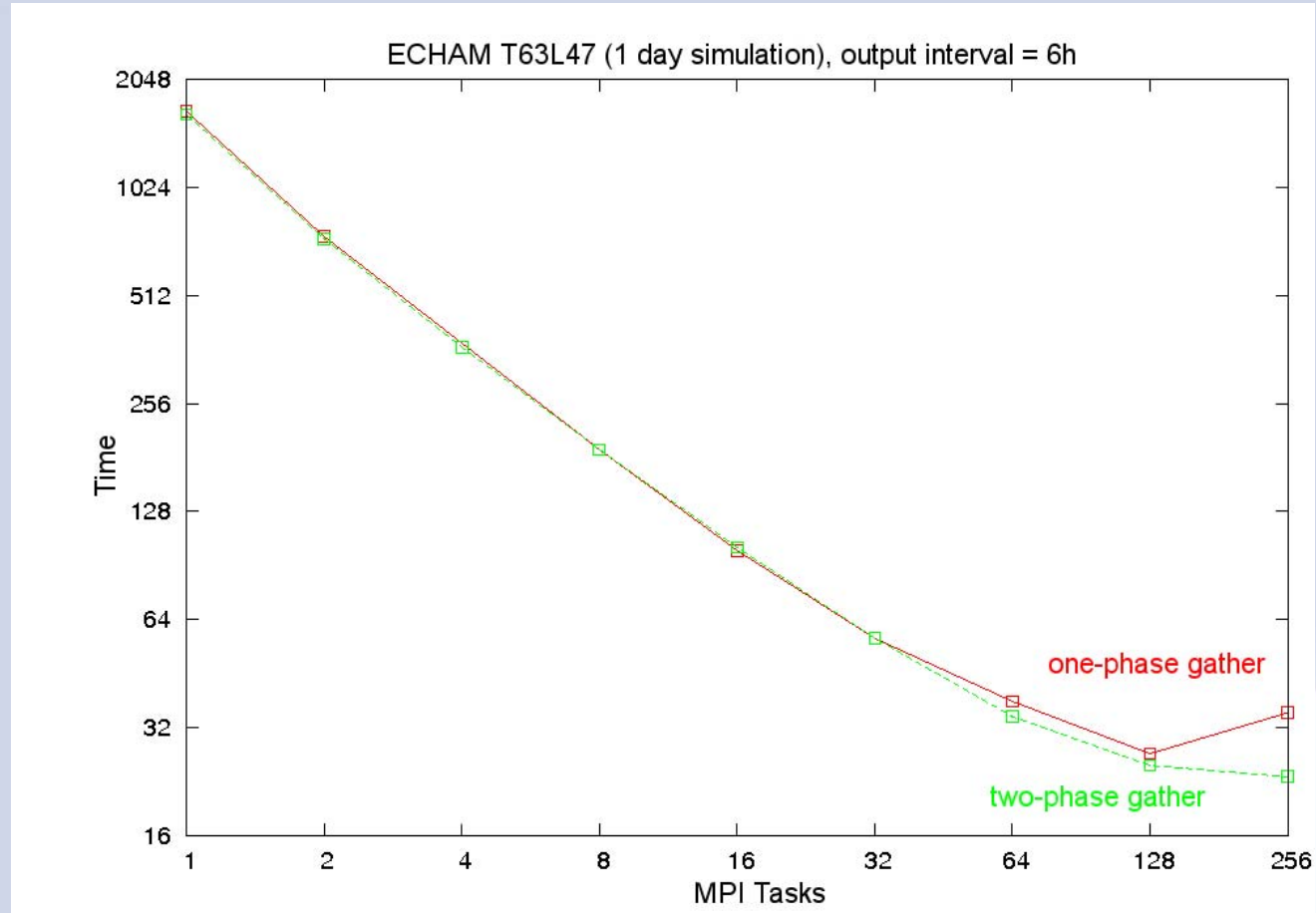




Optimizing Communication

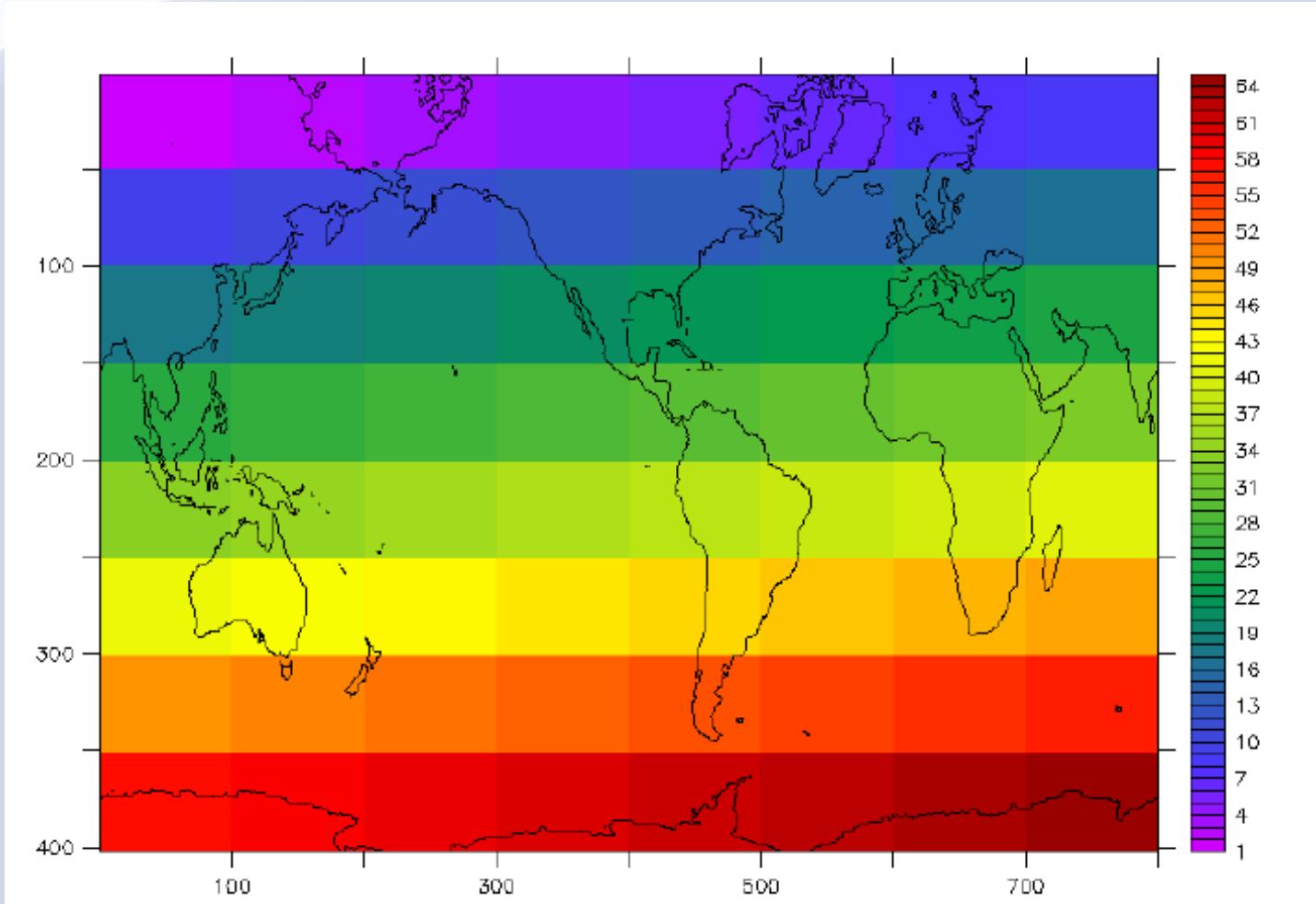


Optimizing Communication



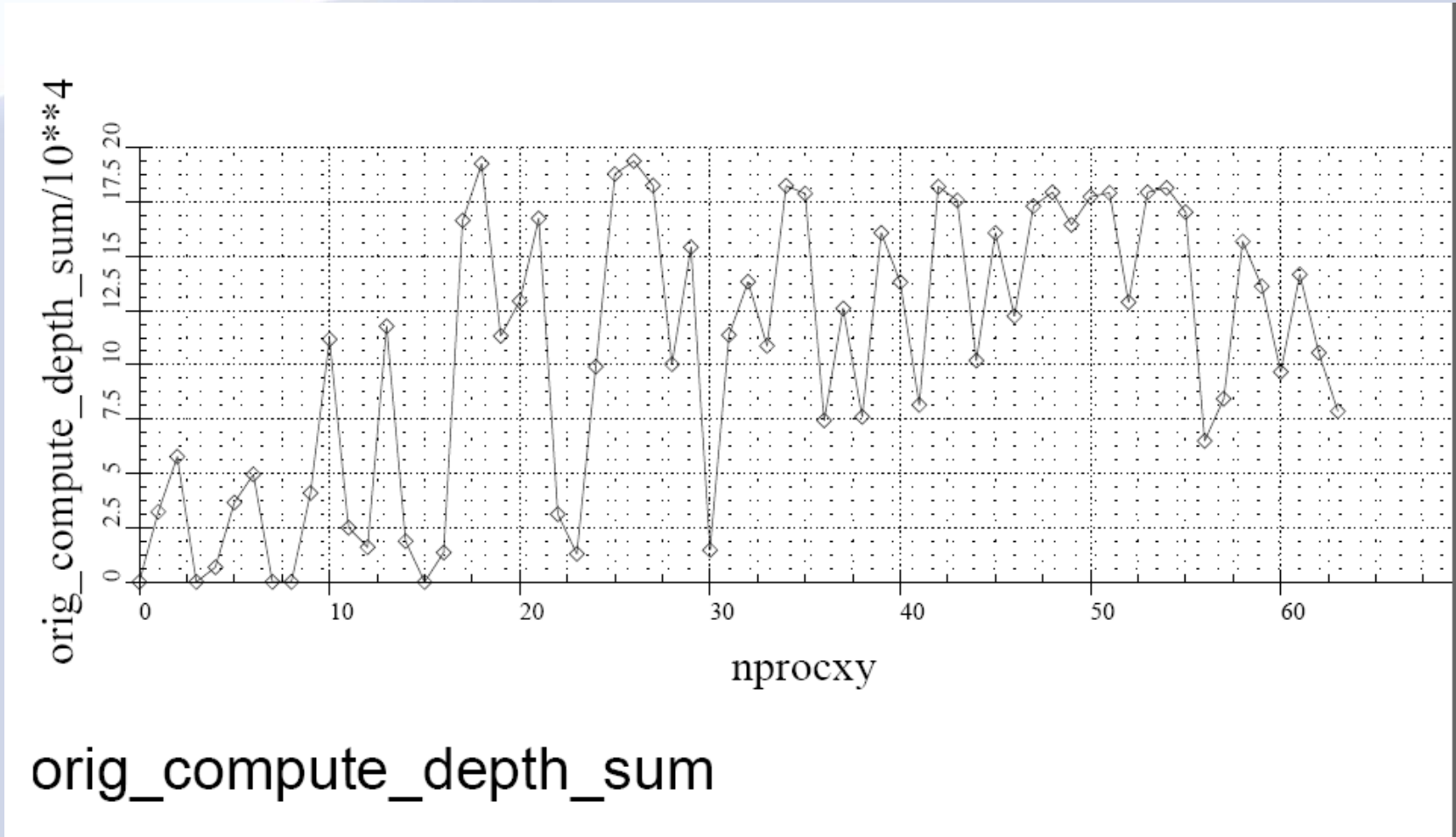


Load Balancing



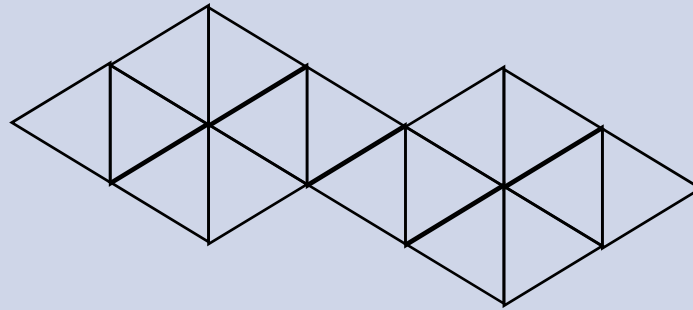


Load Imbalance

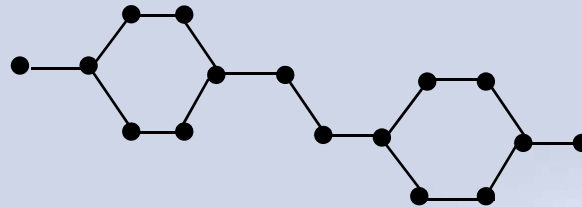


Mapping Mesh \Rightarrow Graph

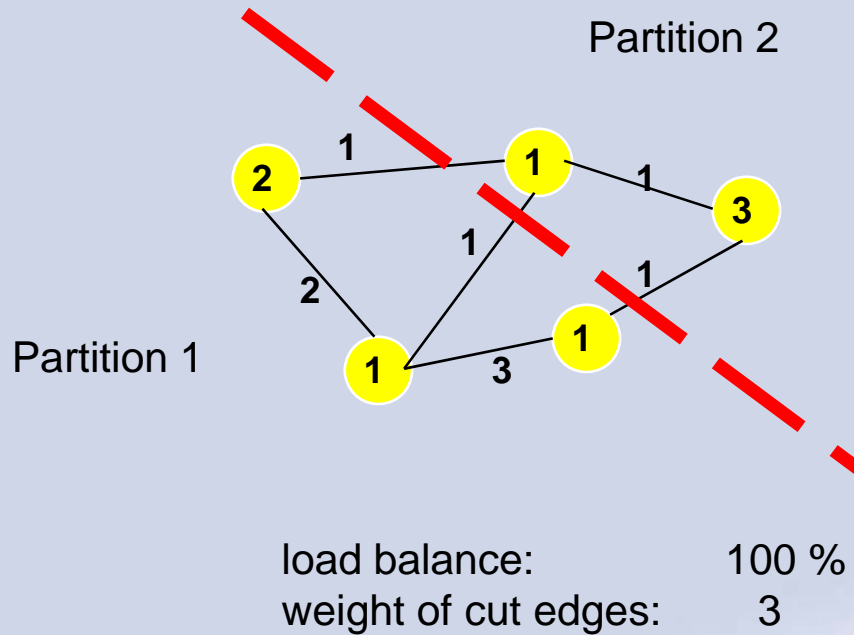
Mesh



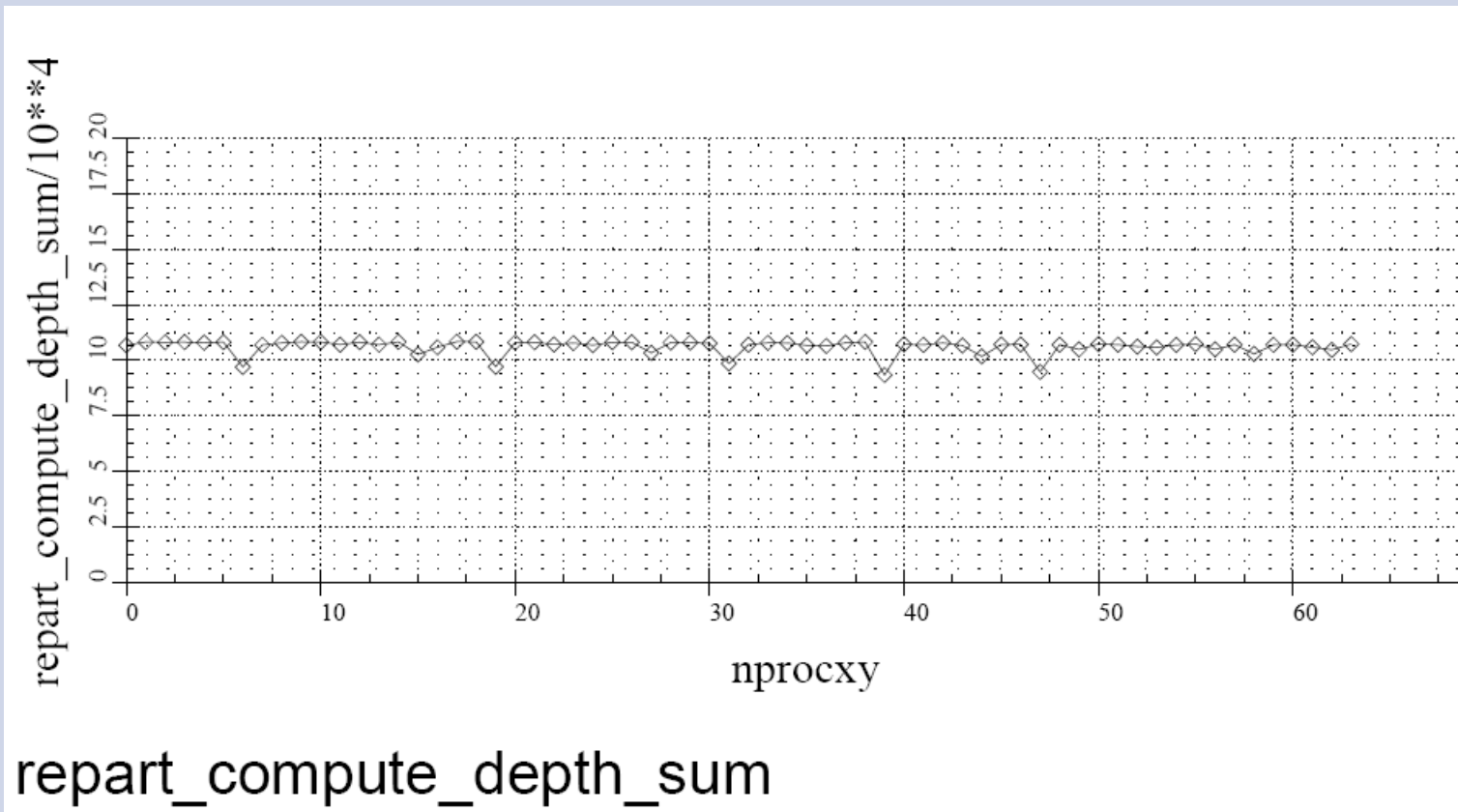
Graph



Graph Partitioning

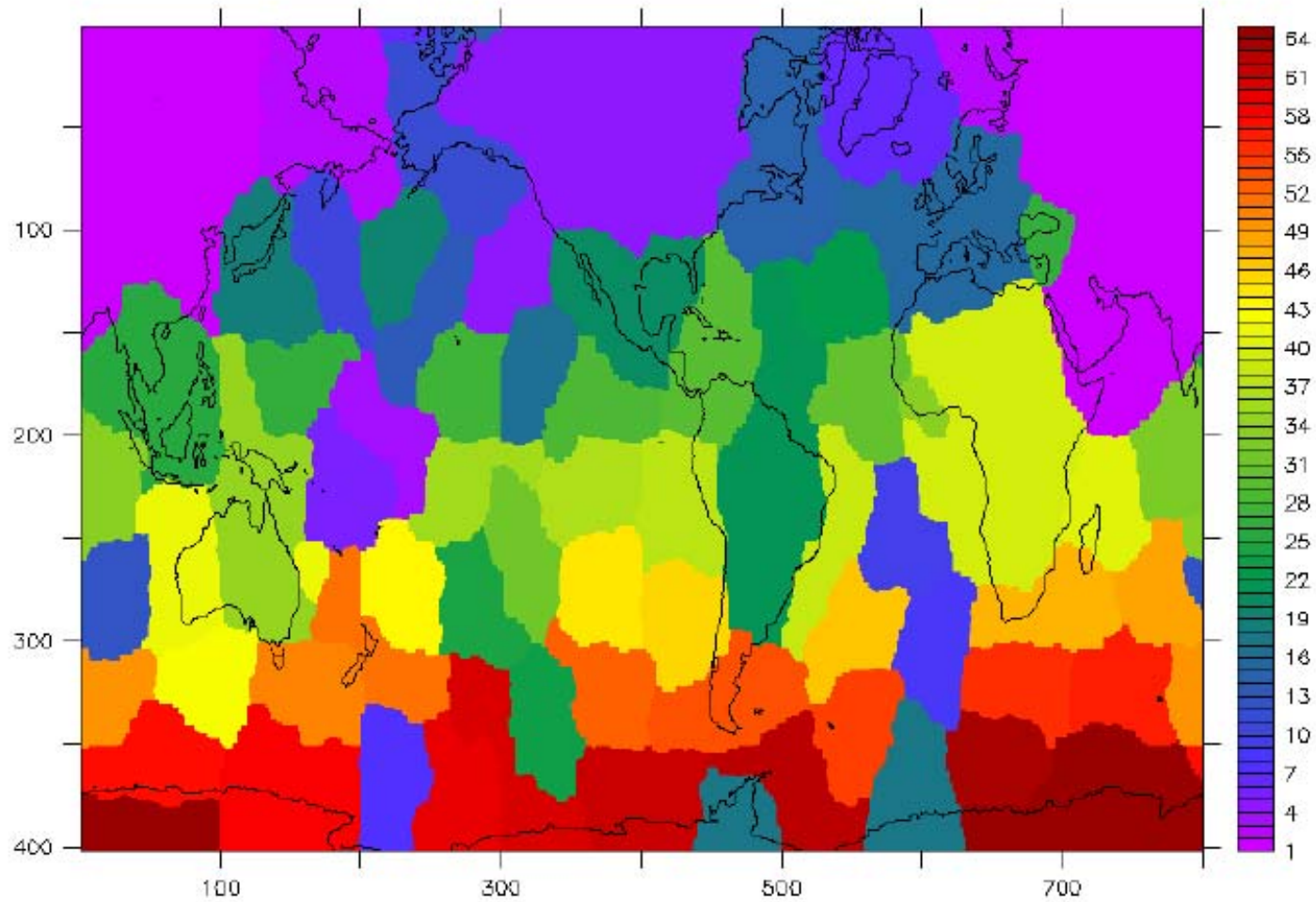


Static Loadbalancing with METIS





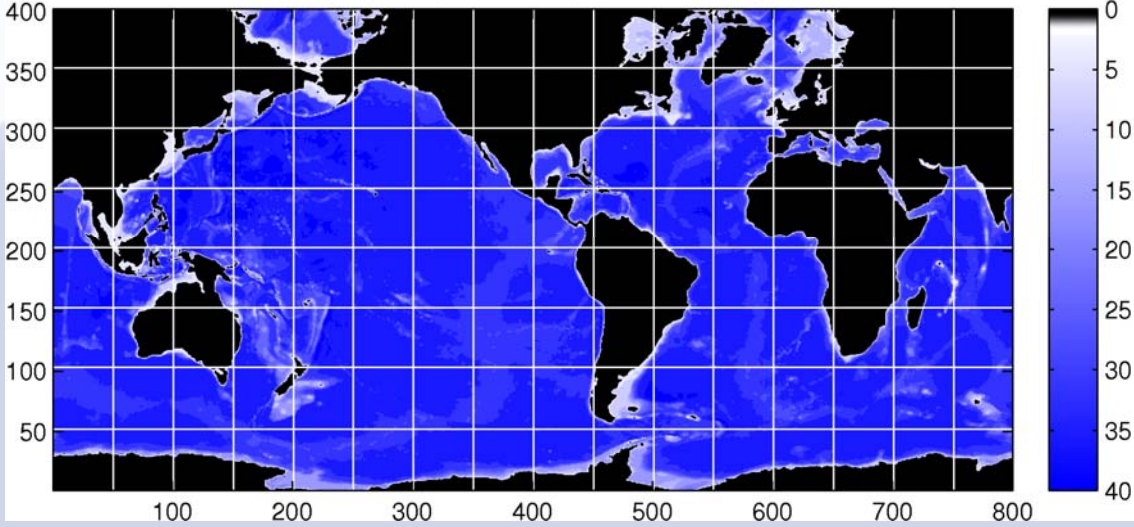
Problem of Graph Partitioning



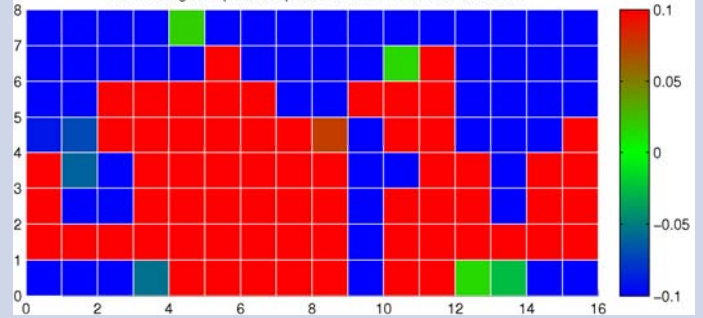


Loadbalancing

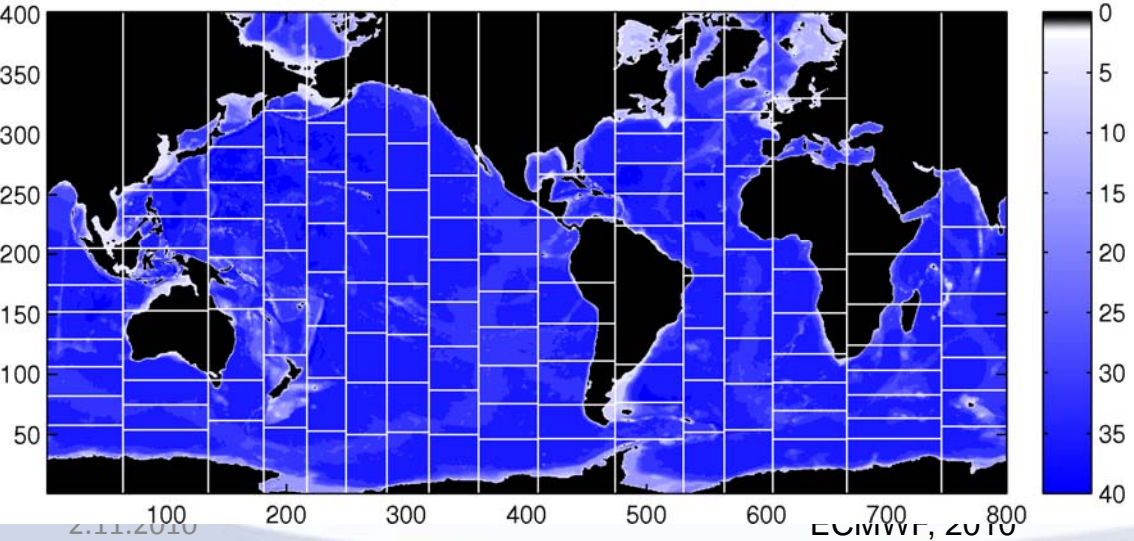
TP04L40 gridpoint space with regular 16x8 decomposition



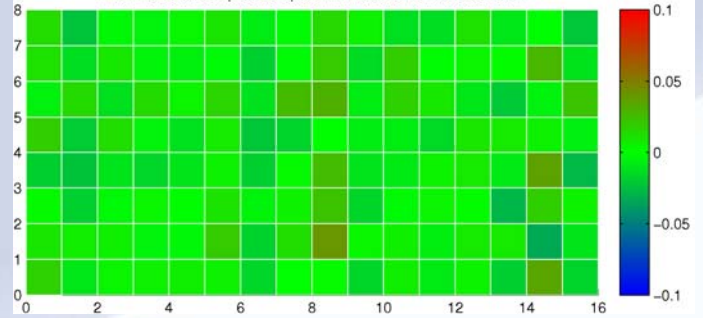
TP04L40 regular: process space with relative workload deviation



TP04L40 gridpoint space with load balanced 16x8 decomposition

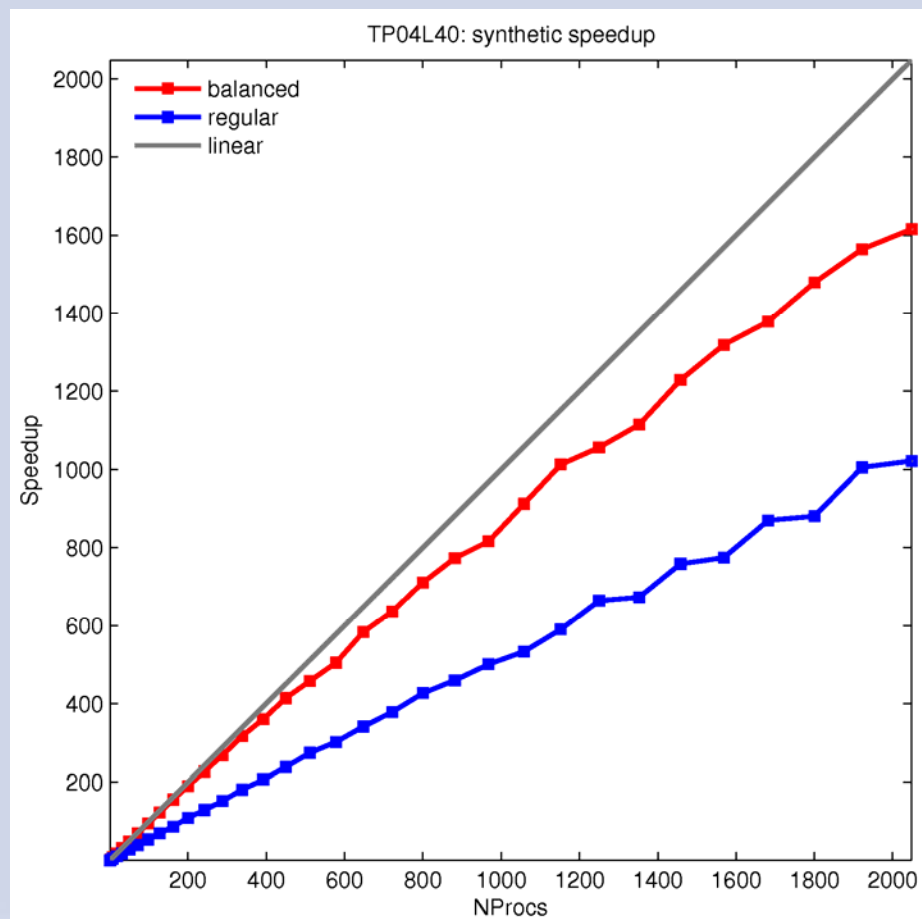


TP04L40 balanced: process space with relative workload deviation



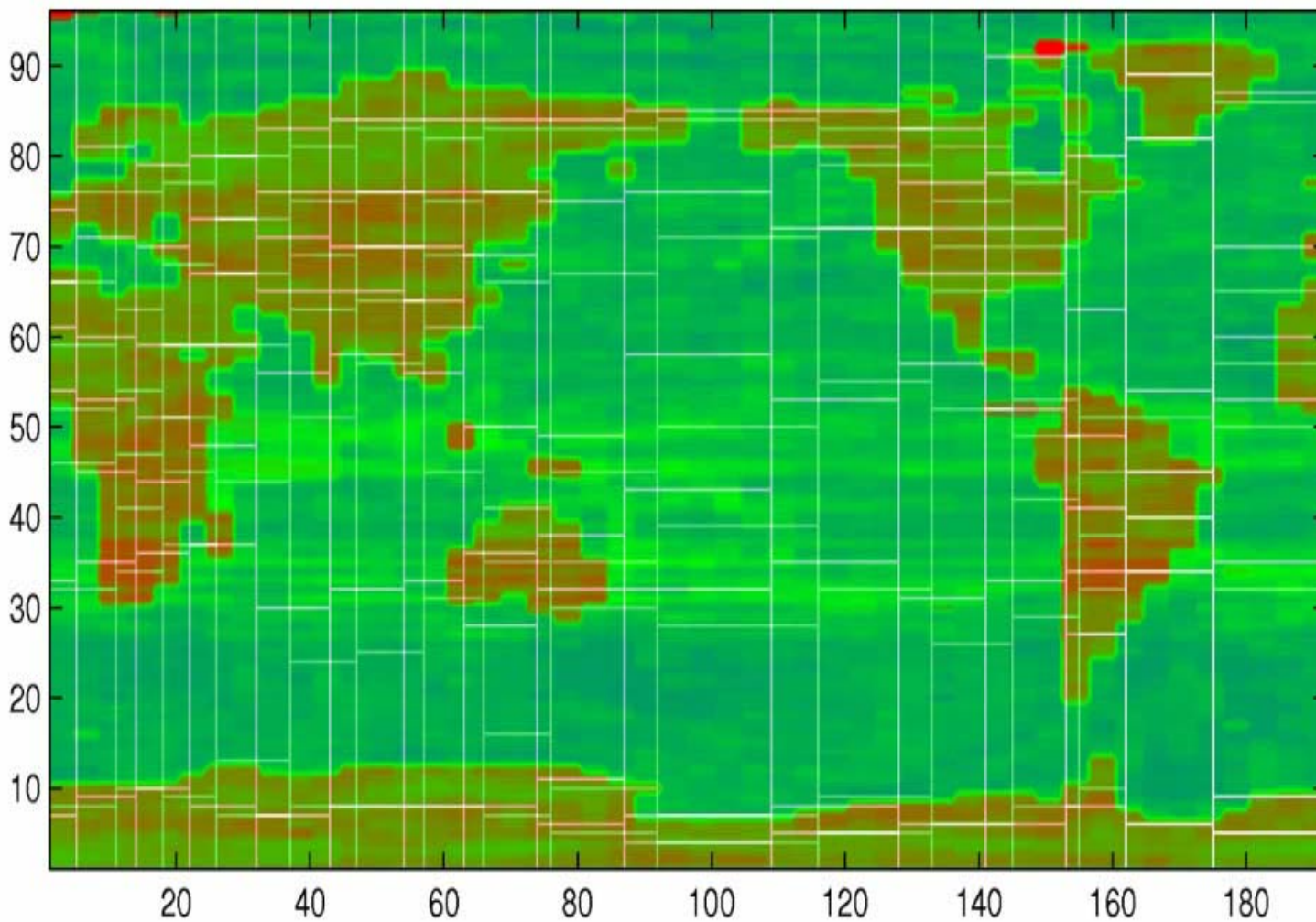


Loadbalancing



ECHAM T63L47 (192x96x47), step=001

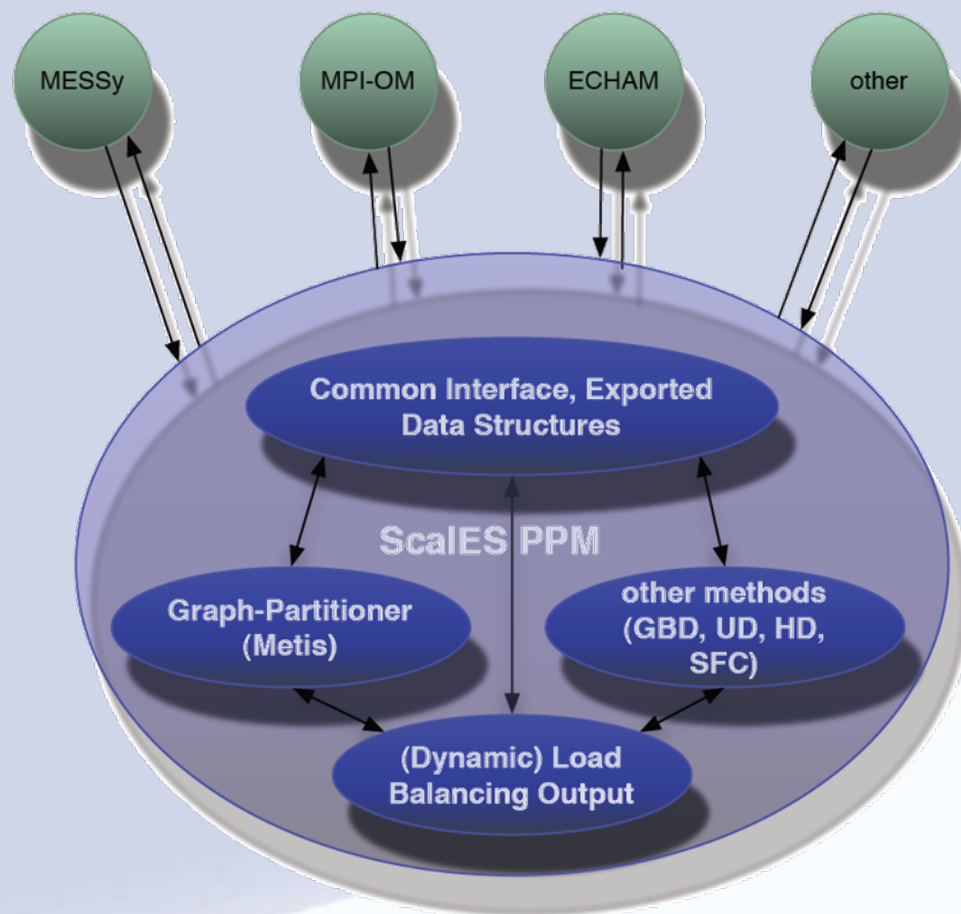
© DKRZ



ScaLES PPM (Parallel Partitioning Modul)

Features:

- Describe partitioning in model-independent data structures.
- Provide convenient API to multiple partitioning algorithms.
- Support data relocation for repartitioning
- Provide solid foundation for other partitioning dependent functionality.



ScaLES PPM schematic library design

Optimization

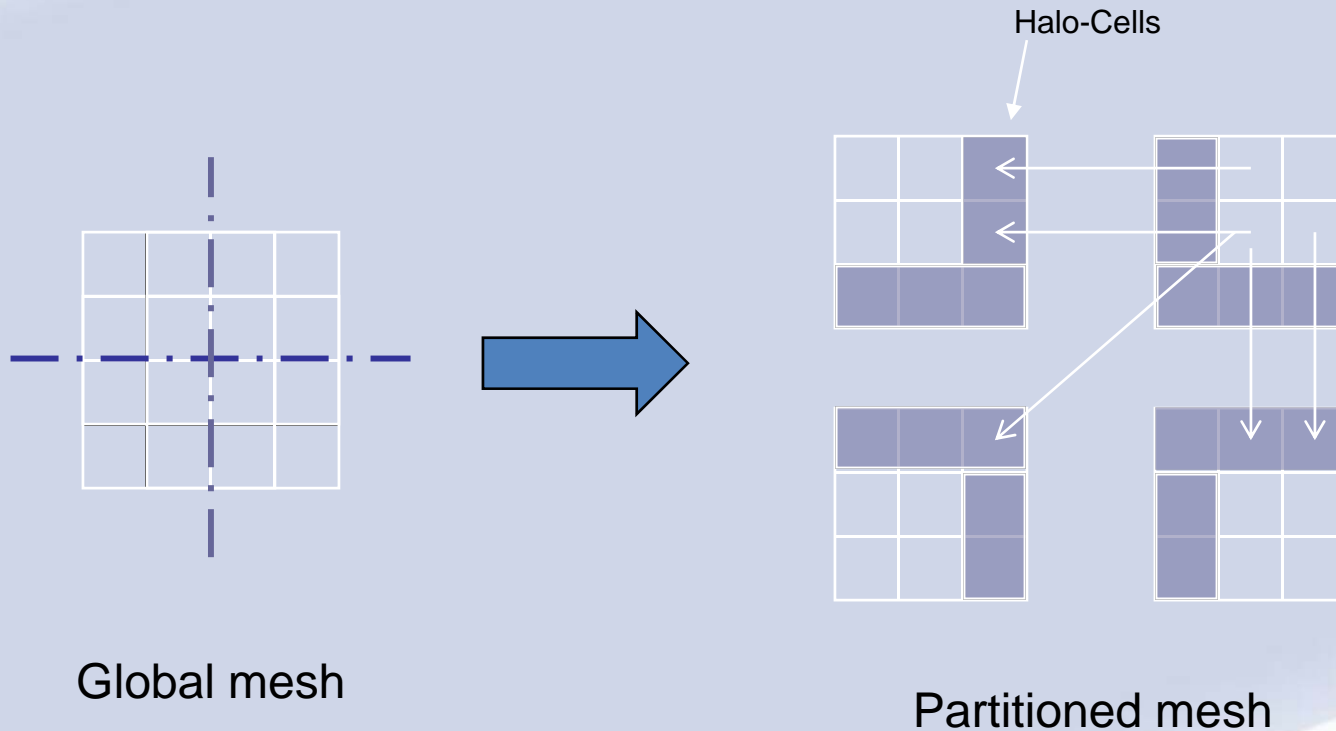
Computation vs Communication

- Expanding Halos
 - More Local Operations
 - Less Communication

CG Method with appropriate Preconditioner

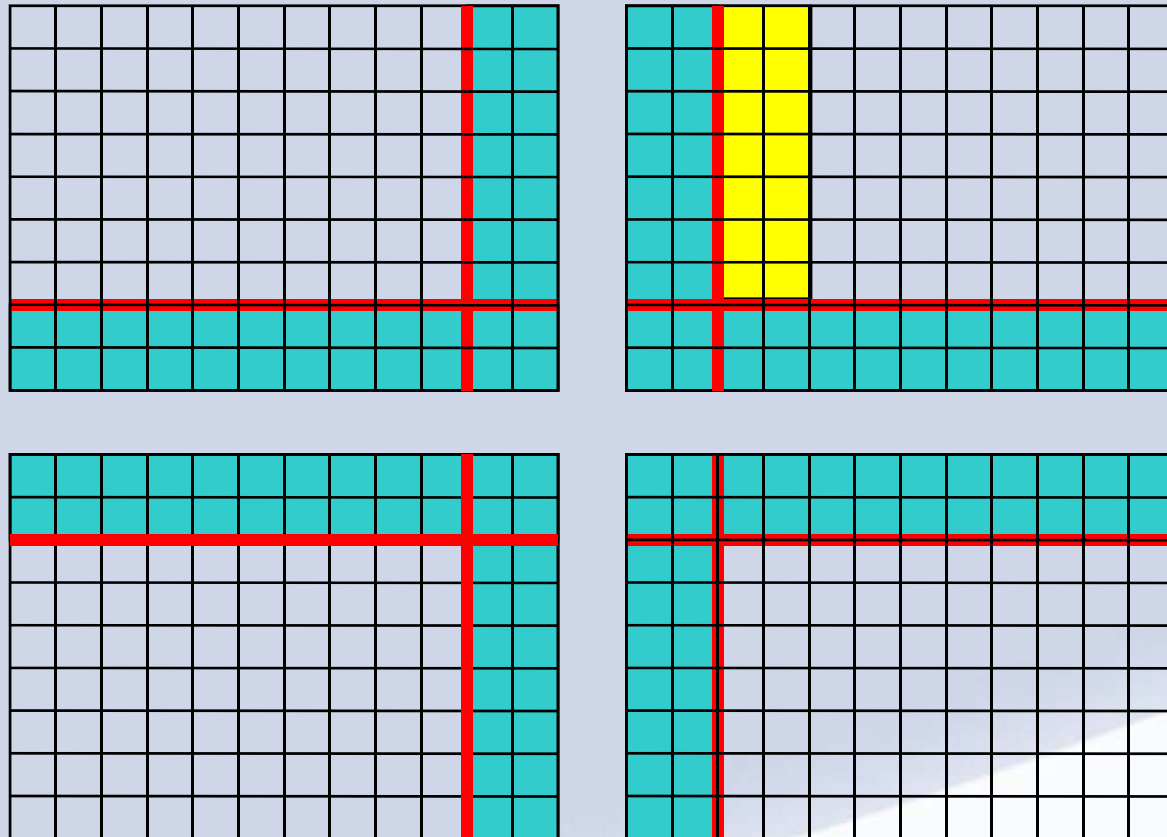
- Faster Convergence =>
 - Less Iterations
- Less Communication

Structure of a partitioned mesh



Exchange of halos at each iteration

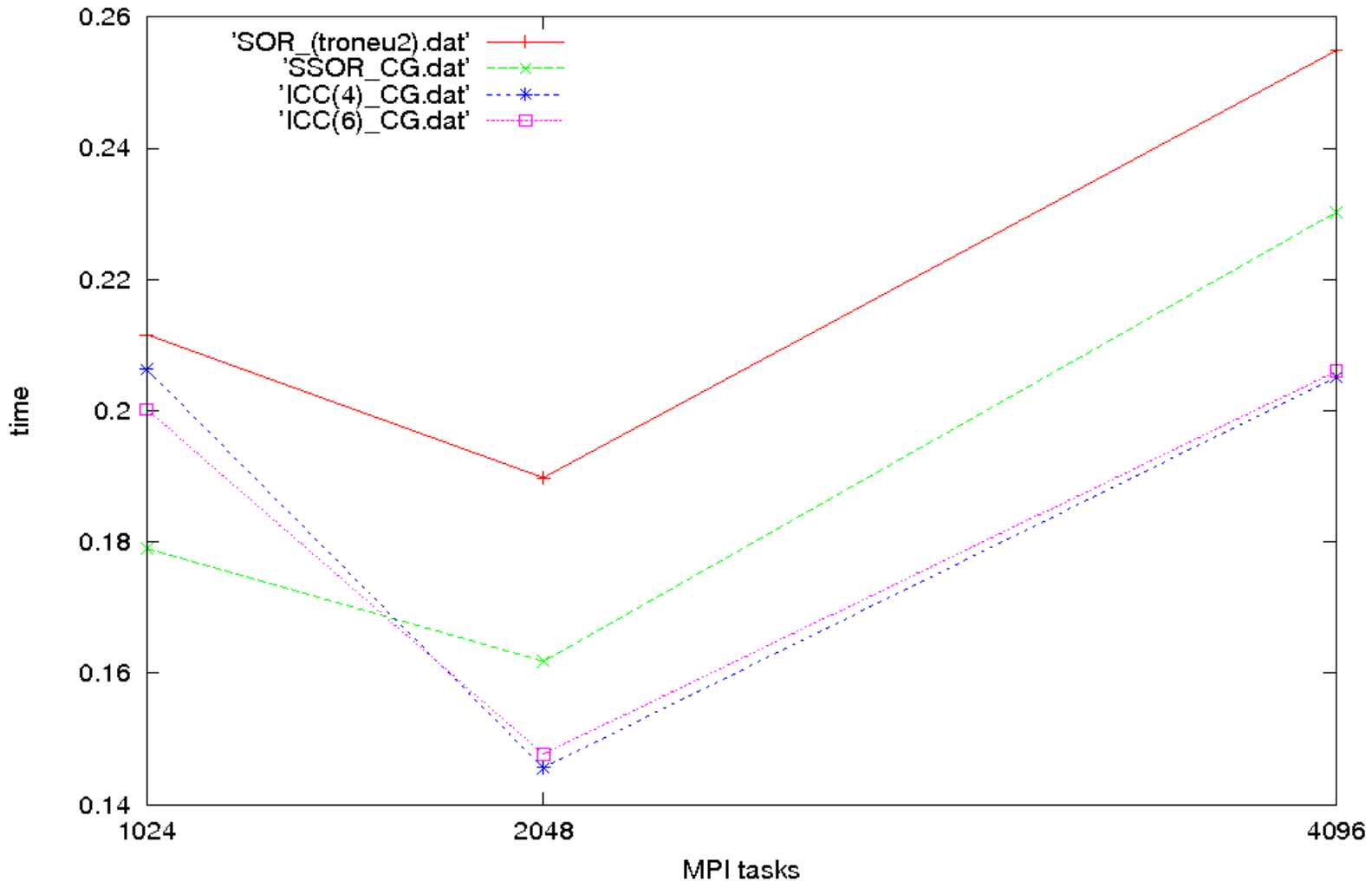
Expanding Halos



Exchange of expanded halos every 2nd iteration

Library of Solvers (KIT)

Runtime for global solvers



Optimization

- Compilers cannot optimize automatically everything
- Optimization is not just finding the right compiler flag
- Major algorithmic changes are necessary
- Solutions must be generic otherwise the whole process is a “Sisyphean Task”

Scalable Earth-System Models

- **Project:**
 - Facilitate High Productivity Climate Simulations
 - Three year program
 - Funded by Federal Ministry for Science and Research (01IH08004E)
 - Started in January 2009
- **Partners:**
 - DKRZ (German Climate Computing Centre)
 - MPIM (Max-Planck-Institute for Meteorology)
 - MPIC (Max-Planck-Institute for Chemistry)
 - AWI (Alfred-Wegener-Institute for Polar Research)
 - KIT (Karlsruhe Institute for Technology)
 - IBM (International Business Machines)

ScalES:

Use case: COSMOS

- ECHAM5 (global atmosphere)
- MESSY/MECCA (atmospheric chemistry)
- MPI-OM (global ocean)
- OASIS4 (coupler)

Work Packages

- Parallel I/O
- Load Balancing
- Architectural issues
- Coupling (e.g. of atmosphere and ocean models)

