

1. Project Description

Research Question

Given massive parallelism, at multiple levels, and of diverse forms and granularities, how can it be **exposed**, **expressed**, and **exploited** by parallel and distributed applications such that execution times are reduced, performance targets are achieved, and acceptable efficiency is maintained?

Methodology

Multilevel scheduling (MLS) extends and bridges the most successful batch (jobs scheduled on a system), application (application processes scheduled on allocated nodes/processors/cores), and thread (application threads scheduled on the allocated cores) scheduling approaches beyond a single or a couple levels of parallelism (scaling up) and beyond their current scale (scaling out).

Objectives

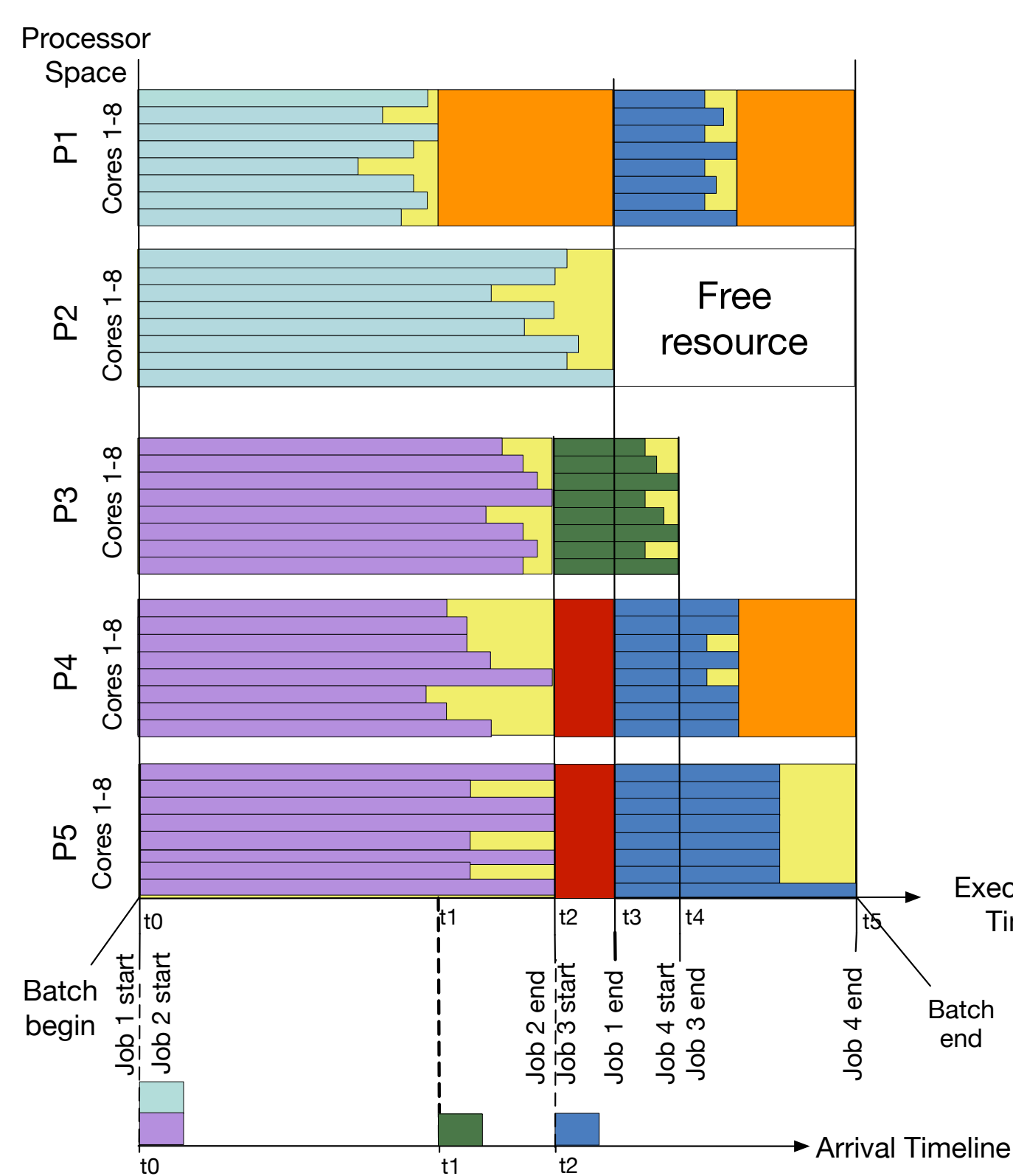
1. Leverage all available parallelism **within** each single hardware parallelism level and **across** the three hardware parallelism levels (system, node, core).
2. Achieve robustness against perturbations (including variations and failures) while minimizing execution time, maintaining acceptable efficiency, and maximizing resource utilization.

Envisioned Outcomes

A prototype multilevel scheduling solution that integrates live feedback information from three, currently disjoint, scheduling levels: job, process, and thread.

3. Exposing, Expressing, and Exploiting Multilevel Parallelism

Current State in Batch, Application, and Thread Level Scheduling

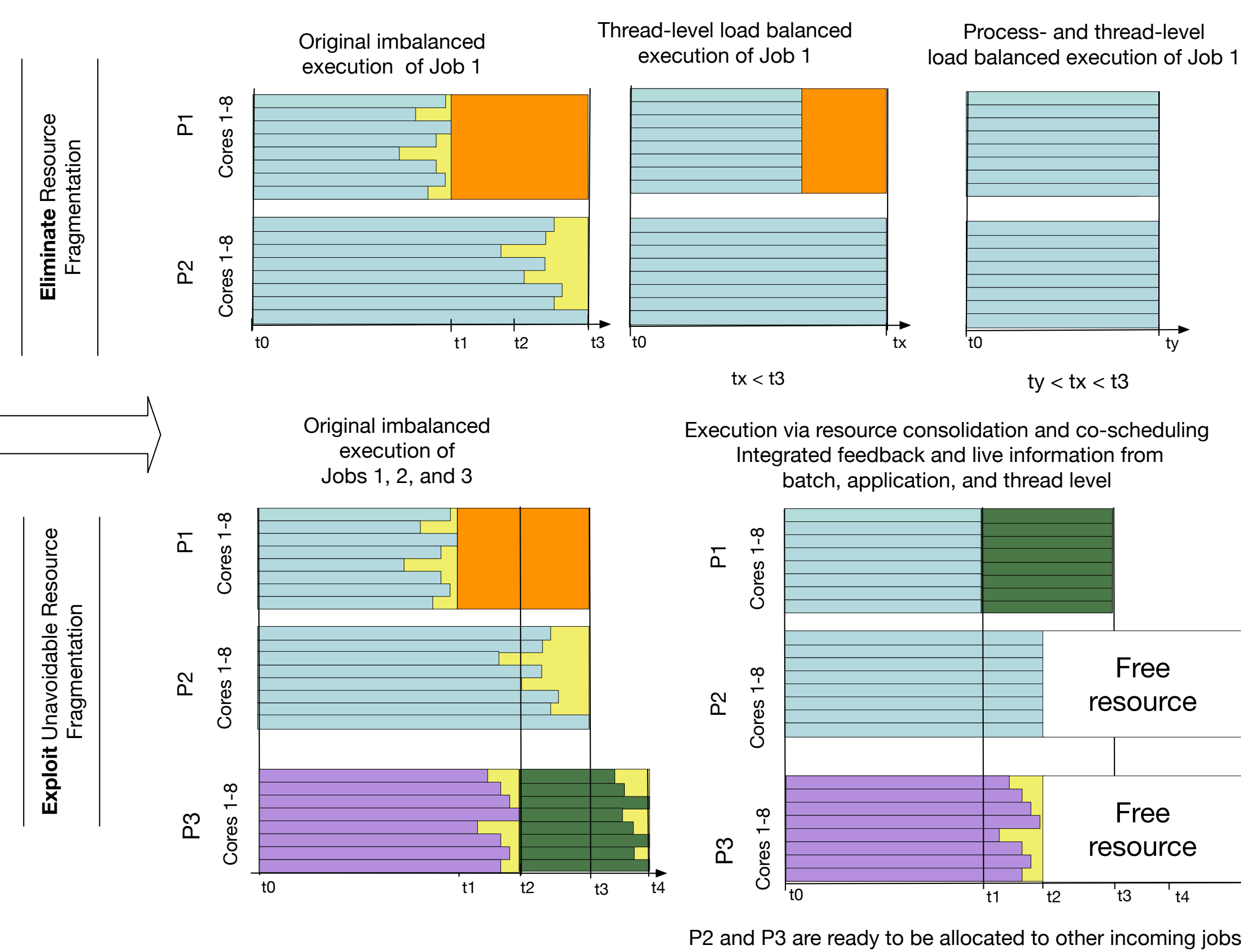


• Four jobs, five compute nodes, eight cores per compute node, one process per node and one thread per core

Legend for current state: Job 1 (blue), Job 2 (orange), Job 3 (green), Job 4 (red). Load imbalance at the thread level (yellow), Load imbalance at the batch level (grey), Load imbalance at the process level (purple), Unexploited hardware parallelism (white).

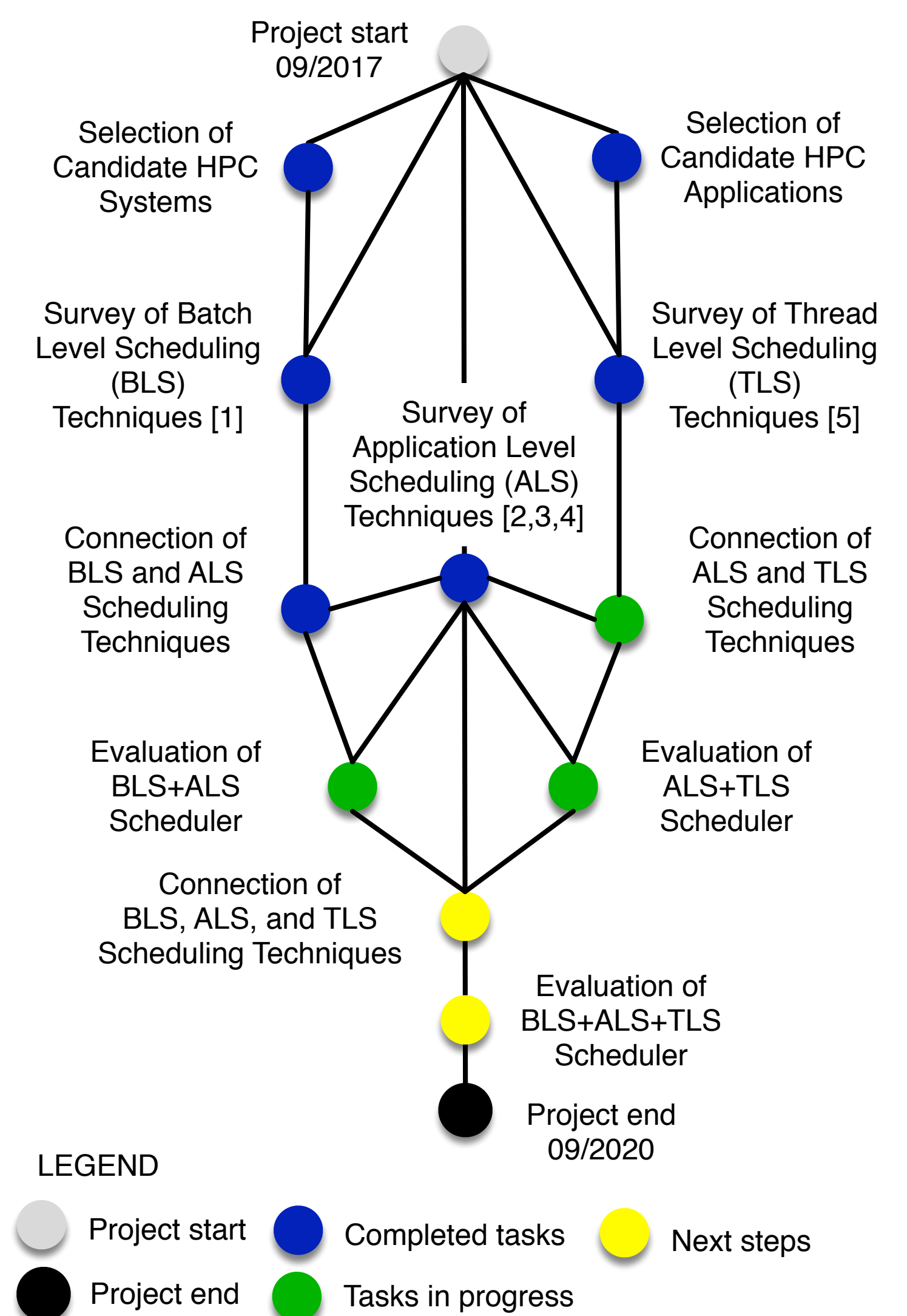
MLS

Multilevel Load Balanced Execution



Batch (job) scheduling: SLURM, OpenPBS, TORQUE, etc.
Application (MPI) scheduling: self-scheduling, work-stealing, etc.
Thread (OpenMP) scheduling: static, guided, dynamic

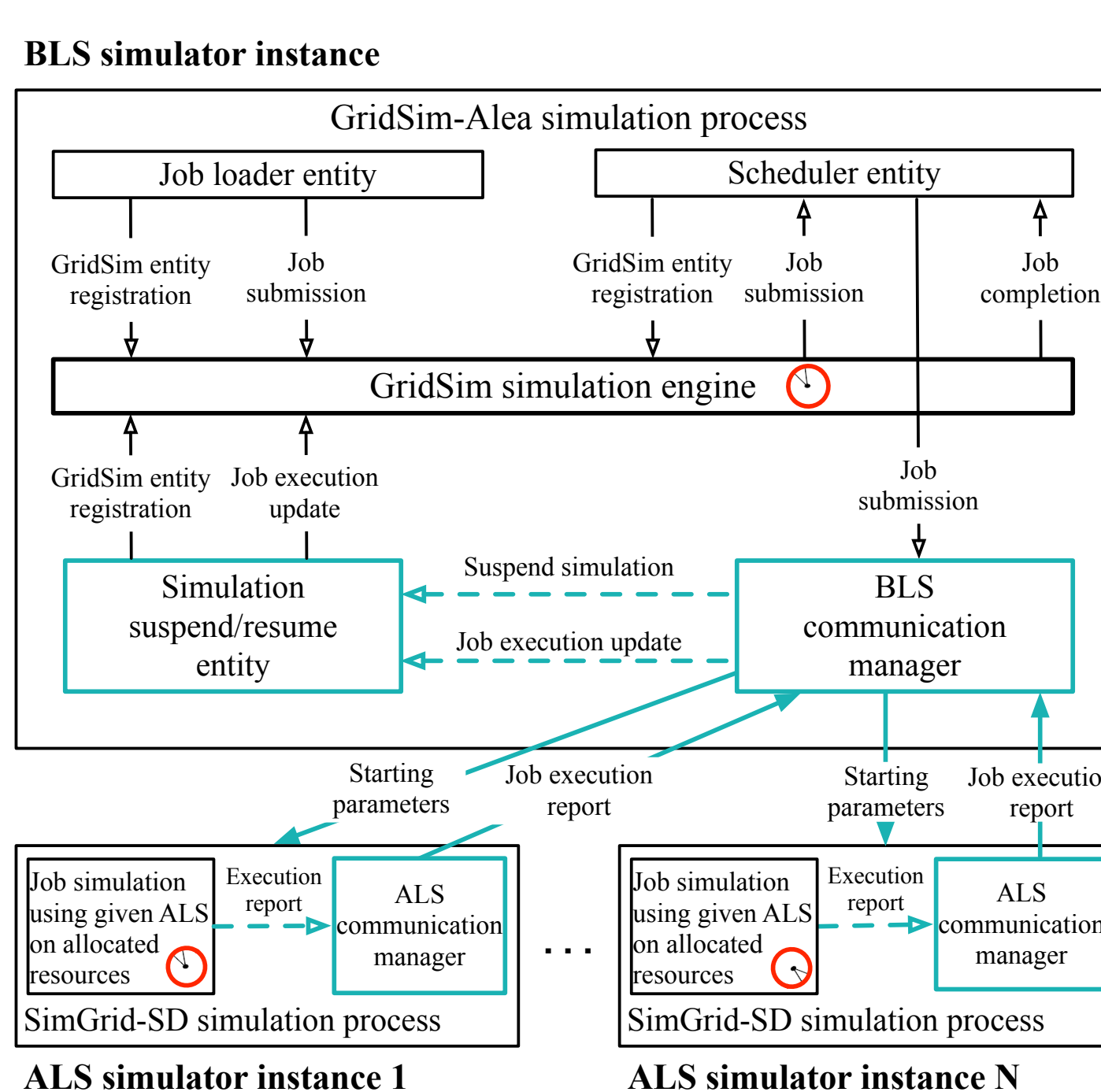
4. Project State



LEGEND: Project start (grey circle), Completed tasks (blue circle), Next steps (yellow circle), Project end (black circle), Tasks in progress (green circle).

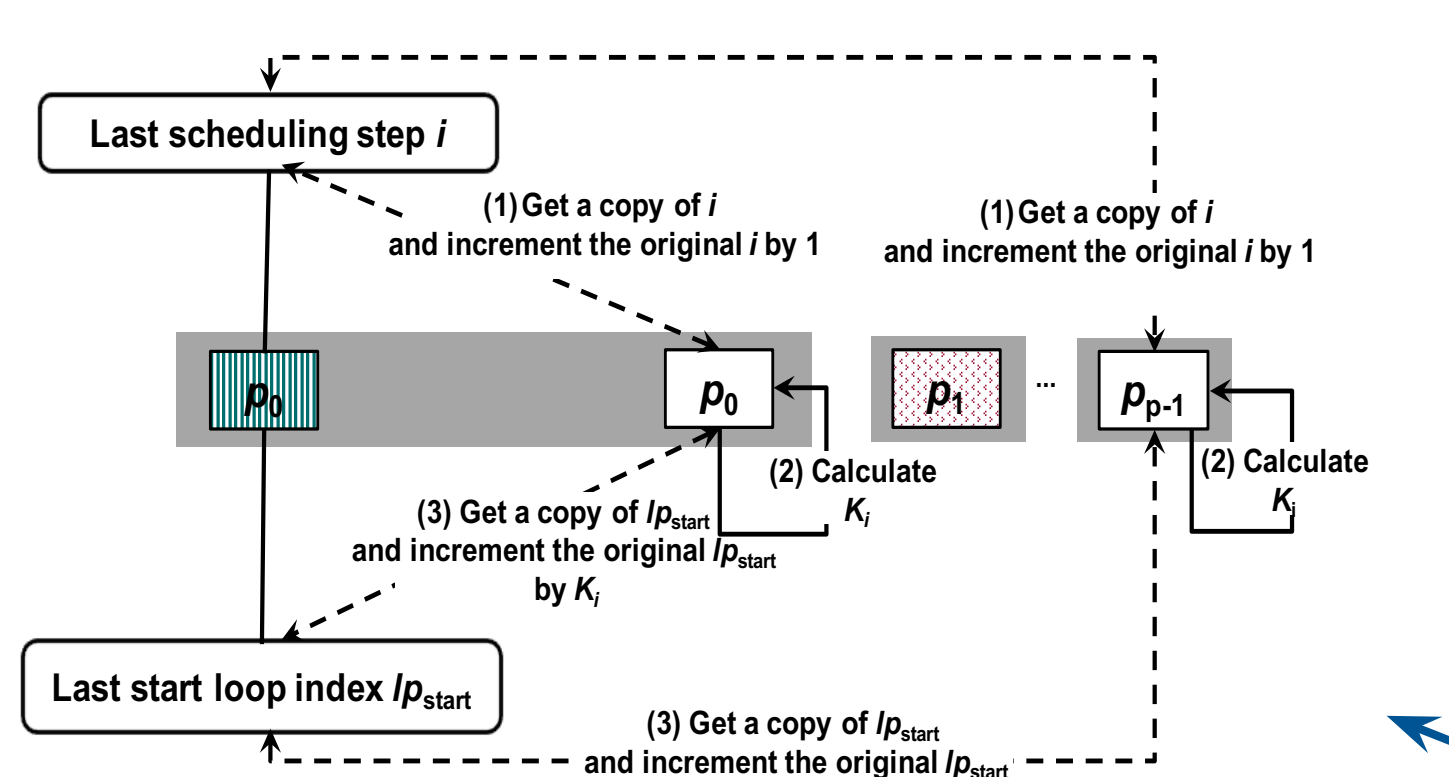
5. Selected Results

Result 1: A two-level simulator for batch and application level scheduling [ISPDC 2017]

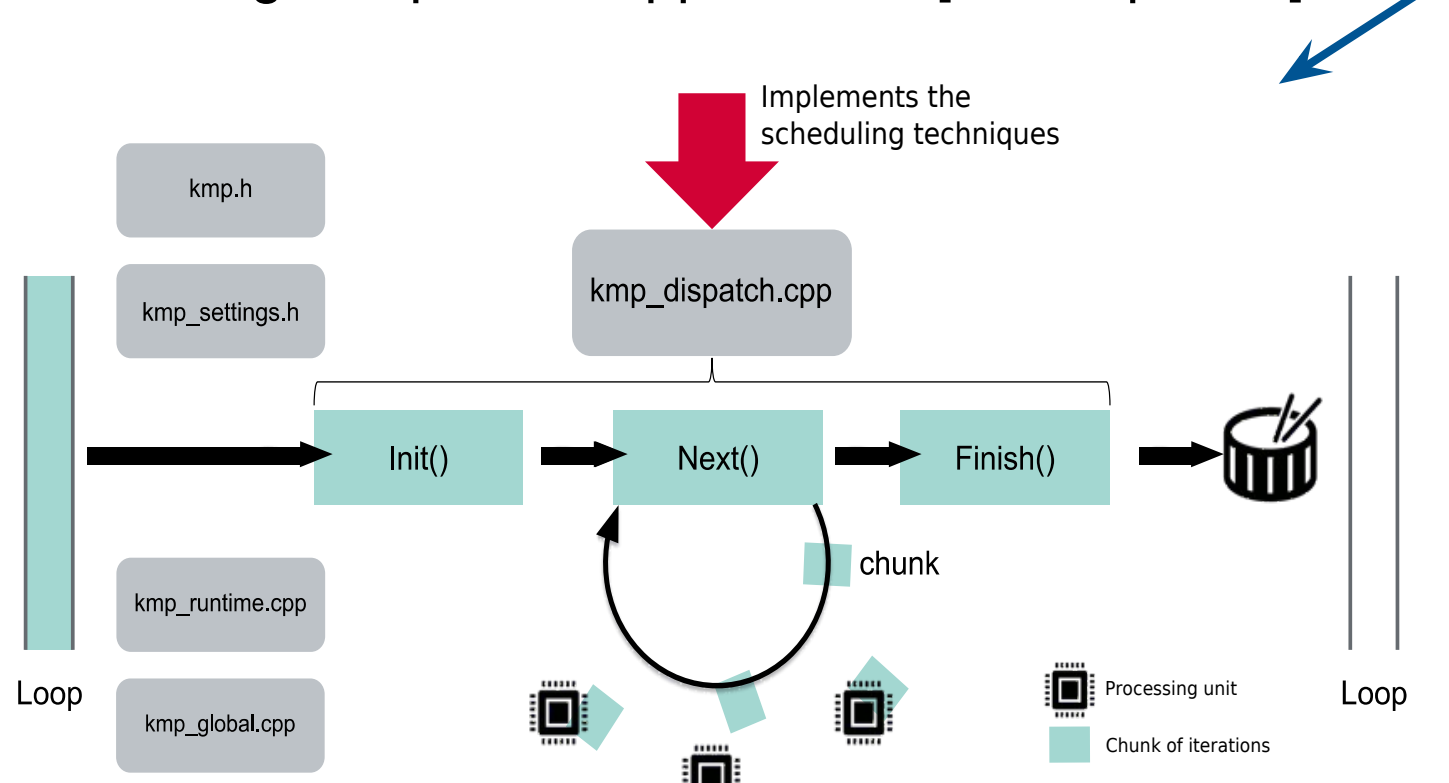


Legend: Internal GridSim events (red arrow), External messages of the connection layer (blue arrow), Internal synchronization events of the connection layer (green arrow), Connection layer entities (blue box), Simulation clock within a simulation instance (red circle).

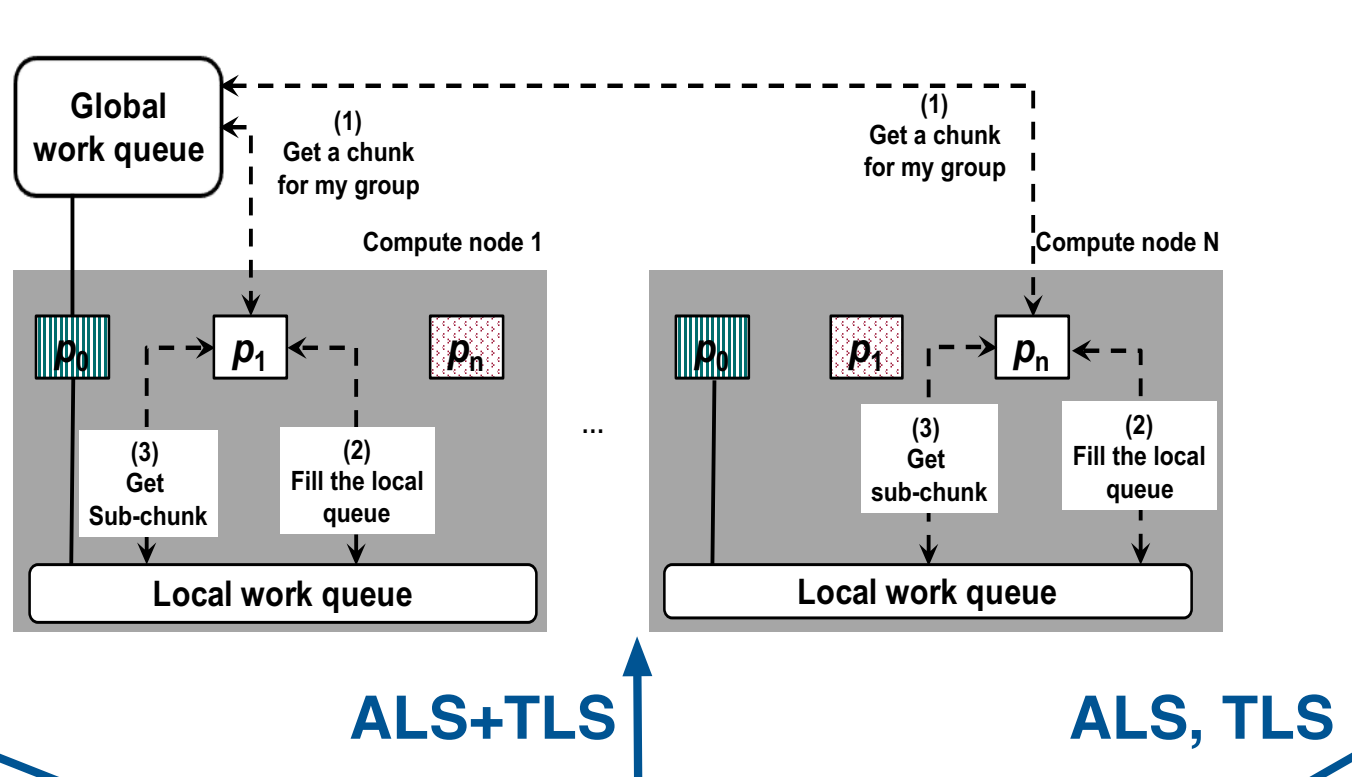
Result 2: A distributed chunk-calculation approach for dynamic loop scheduling at ALS on heterogeneous distributed-memory systems [PDP 2019]



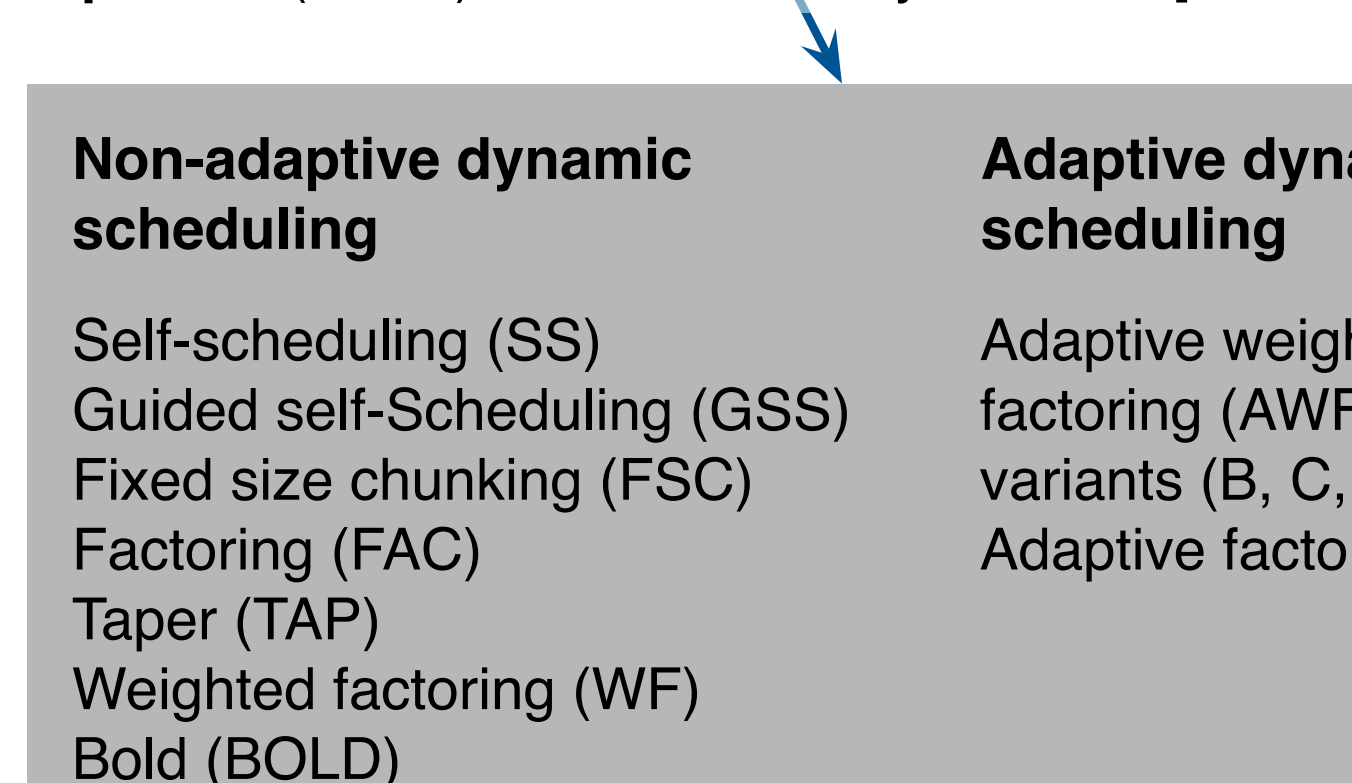
Result 5: A Runtime Approach for Dynamic Load Balancing of OpenMP Applications [SC'19 poster]



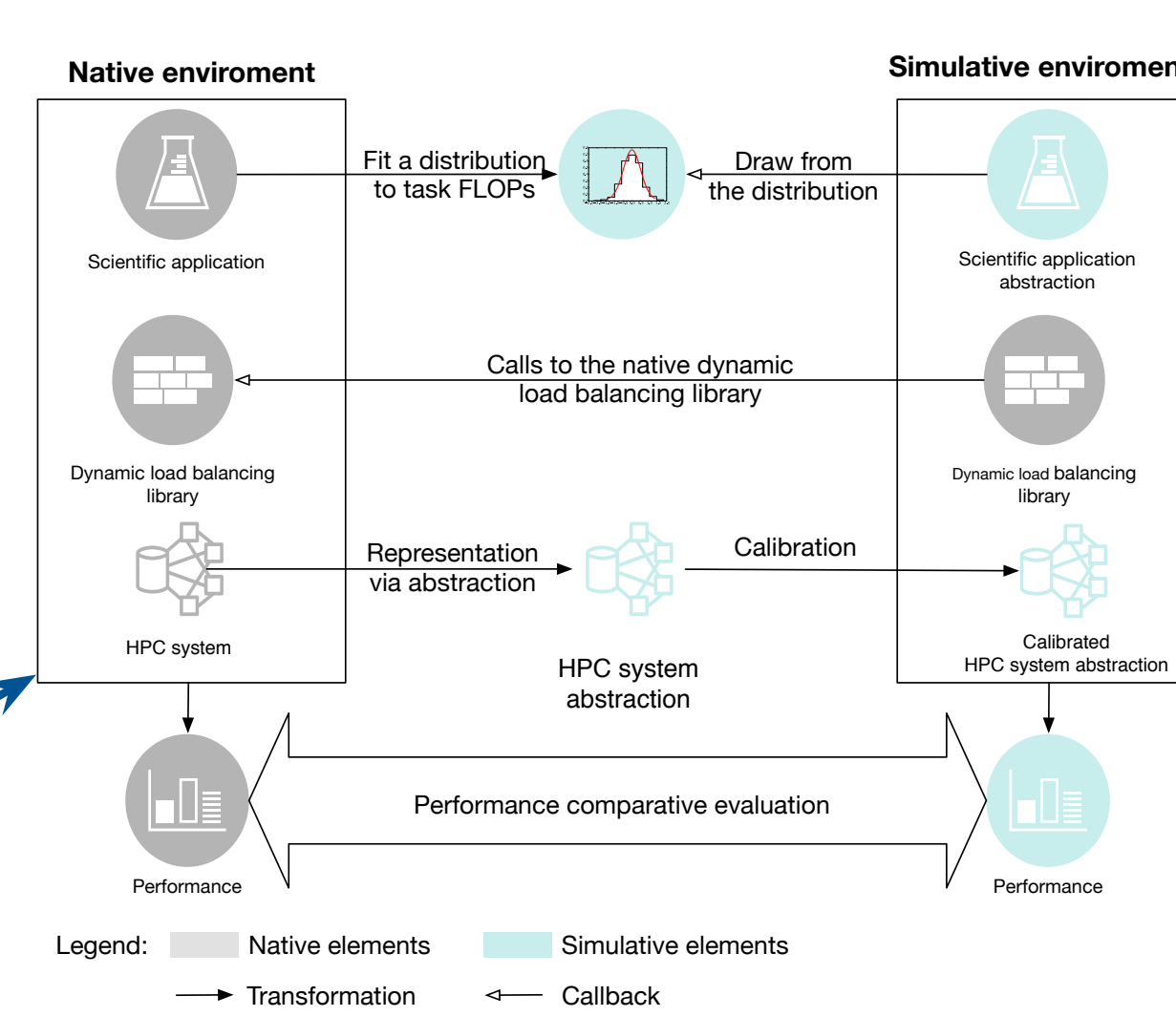
Result 3: Hierarchical distributed chunk calculation approach for intra- and inter node load balancing using an MPI+MPI approach [IPDPSW 2019]



Result 6: Implemented dynamic loop scheduling algorithms in OpenMP (LLVM) & an MPI Library DLS4LB [PASC19] [SIAM PP 20]



Result 4: Generic methodology for realistic simulations of parallel applications on HPC systems [FGCS 2019]



More MLS Details



hpc.dmi.unibas.ch/HPC/MLS.html

6. Selected Project Publications

- [1] Eleliemy, A., Mohammed, A., and Ciorba, F. M. "Exploring the Relation Between Two Levels of Scheduling Using a Novel Simulation Approach", In Proceedings of the 16th International Symposium on Parallel and Distributed Computing (ISPDC), 2017.
- [2] Eleliemy, A. and Ciorba, F. M. "Dynamic Loop Scheduling Using MPI Passive-Target Remote Memory Access", In Proceedings of the 27th Euromicro International Conference on Parallel, Distributed, and Network-Based Processing (PDP), 2019.
- [3] Eleliemy, A. and Ciorba, F. M. "Hierarchical Dynamic Loop Scheduling on Distributed-Memory Systems Using an MPI+MPI Approach", In Proceedings of the 20th IEEE International Workshop on Parallel and Distributed Scientific and Engineering Computing (PDSEC), 2019.
- [4] Mohammed, A., Eleliemy, F. M., Ciorba, F., Kasiecke, and I. Banicescu "An Approach for Realistically Simulating the Performance of Scientific Applications on High Performance Computing Systems", Future Generation Computer Systems (FGCS), 2019
- [5] J. H. Müller Korndörfer, F. M. Ciorba, A. Yilmaz, C. Iwainy, J. Doerfert, H. Finkel, and V. Kale, M. Klemm "A Runtime Approach for Dynamic Load Balancing of OpenMP Parallel Loops in LLVM", Poster at the International Conference for High Performance Computing, Networking, Storage and Analysis (SC), 2019.
- [6] Mohammed, A., Ciorba, F. M., Cavelan, A., Cabezón, R., and Banicescu, I. "Identifying Performance Challenges in Smoothed Particle Hydrodynamics Simulations", Poster at the Platform for Advanced Scientific Computing Conference (PASC), 2019.
- [7] Eleliemy, A., Mohammed, A., and Ciorba, F. M. "Efficient Generation of Parallel Spin-images Using Dynamic Loop Scheduling", In Proceedings of the 8th International Workshop on Multicore and Multithreaded Architectures and Algorithms (M2A2) of the 19th IEEE International Conference for High Performance Computing and Communications (HPCC), 2017.