This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

SAND2019-9518C

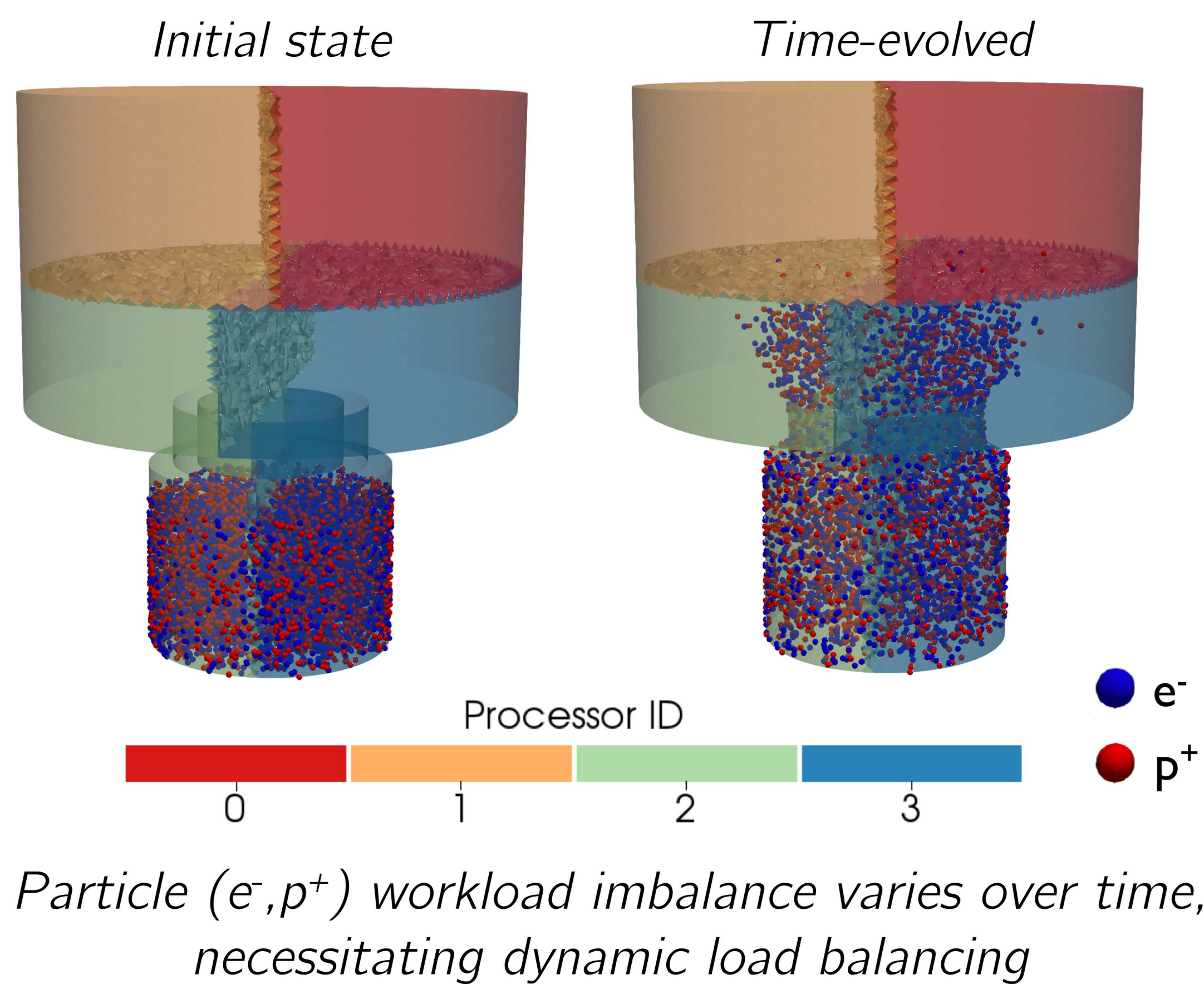# Sandia National Laboratories

# Dynamic, Task-based Load Balancing using DARMA

Presented by: *Jonathan Lifflander & Nicole Slattengren*    Contributors: *Philippe Pébaÿ & Robert Clay*

## Problem

EMPIRE is an ATDM plasma physics application that includes a Particle-In-Cell (PIC) algorithm:

- Initial particle distributions can be spatially concentrated, creating **heavy load imbalance**
- Particles may move rapidly across the domain, inducing **workload variation over time**
- Existing MPI-based EMPIRE code **does not support load balancing (LB)**
- Future **Hybrid PIC/Fluid** configurations present a difficult challenge for LB (multi-objective)
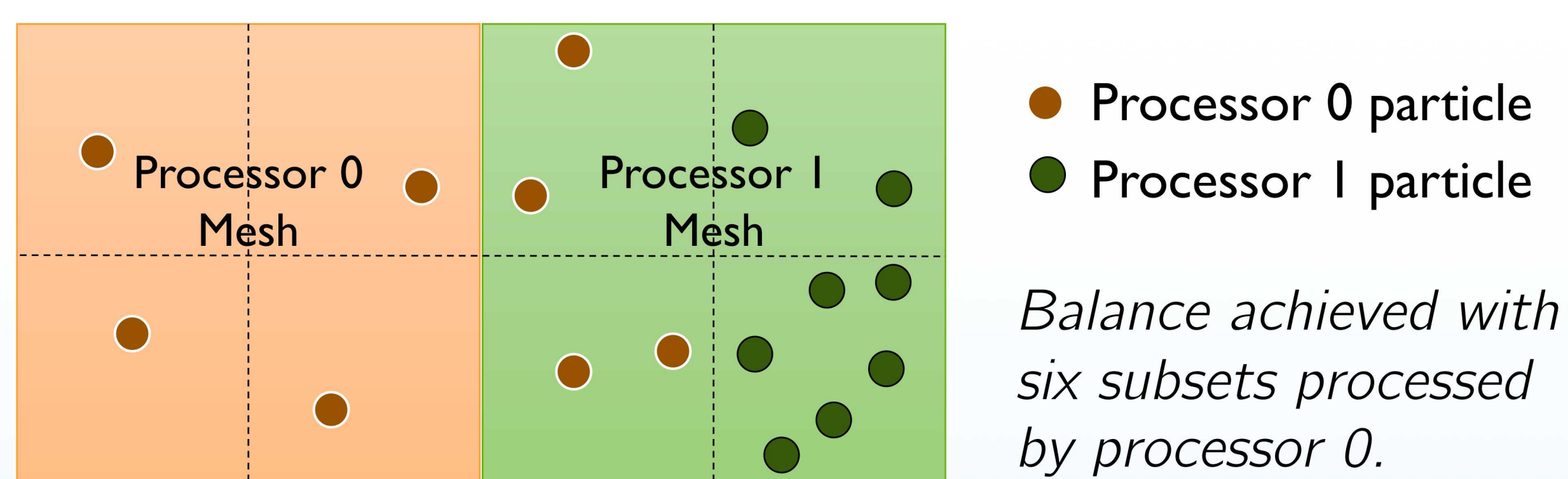


*Initial state*    *Time-evolved*

Processor ID
0  1  2  3

$e^-$
$p^+$

*Particle ($e^-$, $p^+$) workload imbalance varies over time, necessitating dynamic load balancing*

## Approach

**Conventional approach:** infrequently change the mesh decomposition to offset particle imbalance.

- − Synchronous process
- − Large volumes of data must be migrated to new processors or recomputed from the new mesh

**Our approach:** maintain the **static, balanced mesh decomposition**, but **split the particles** on *each* rank-decomposed mesh block into $k$ subsets.



- Processor 0 particle
- Processor 1 particle

*Balance achieved with six subsets processed by processor 0.*

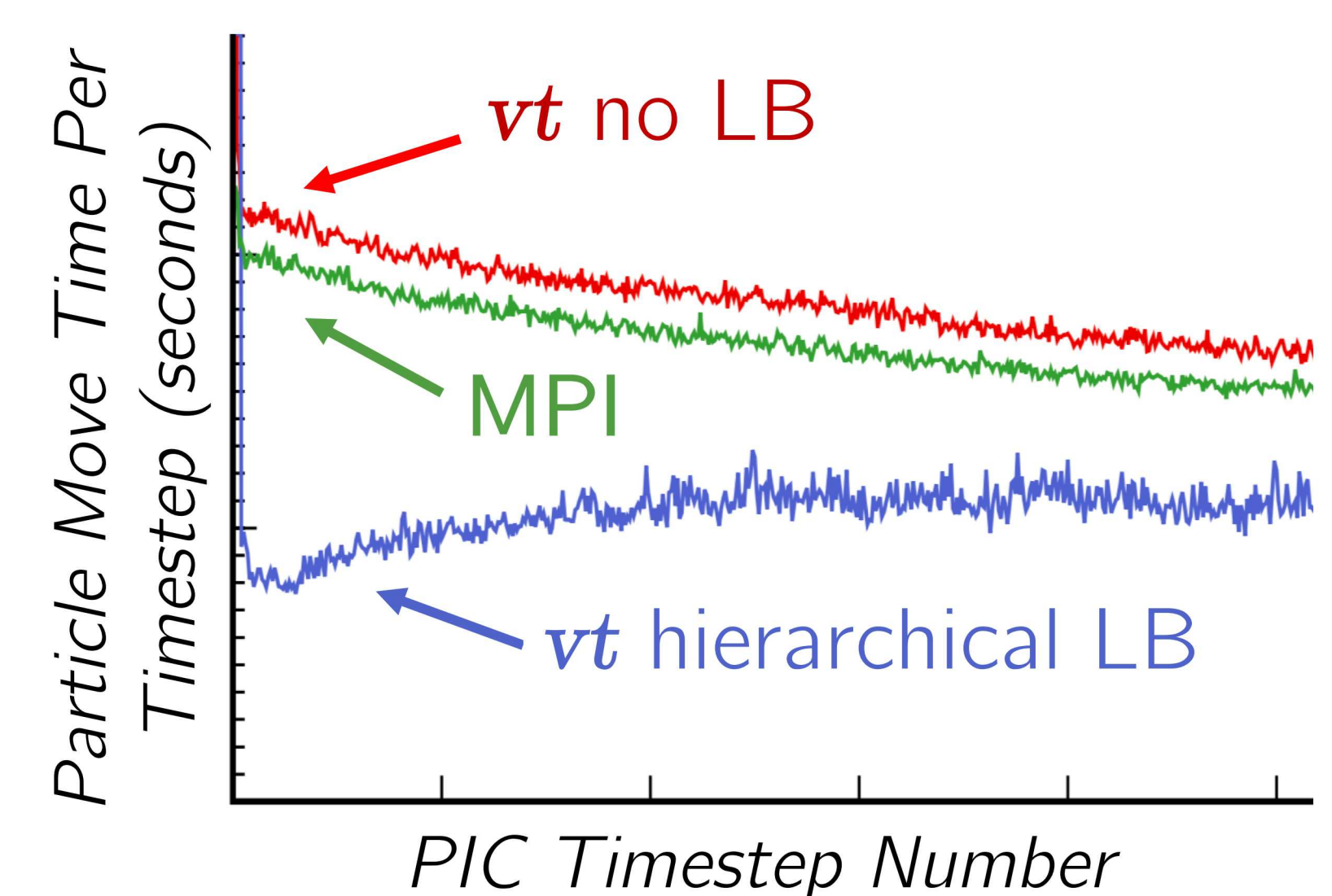Fine-grained, **dynamic** LB of particles:

- + Decreases data migration cost
- + Facilitates communication/computation overlap

## DARMA's Virtual Transport (*vt*) Tasking Library:

- Interoperable with MPI
- Incremental adoption model for C++ "taskification"
- **Dynamically** migrate data and work off-processor
- Includes scalable load balancers
- Developing a fully-distributed, measurement-driven, communication-aware LB
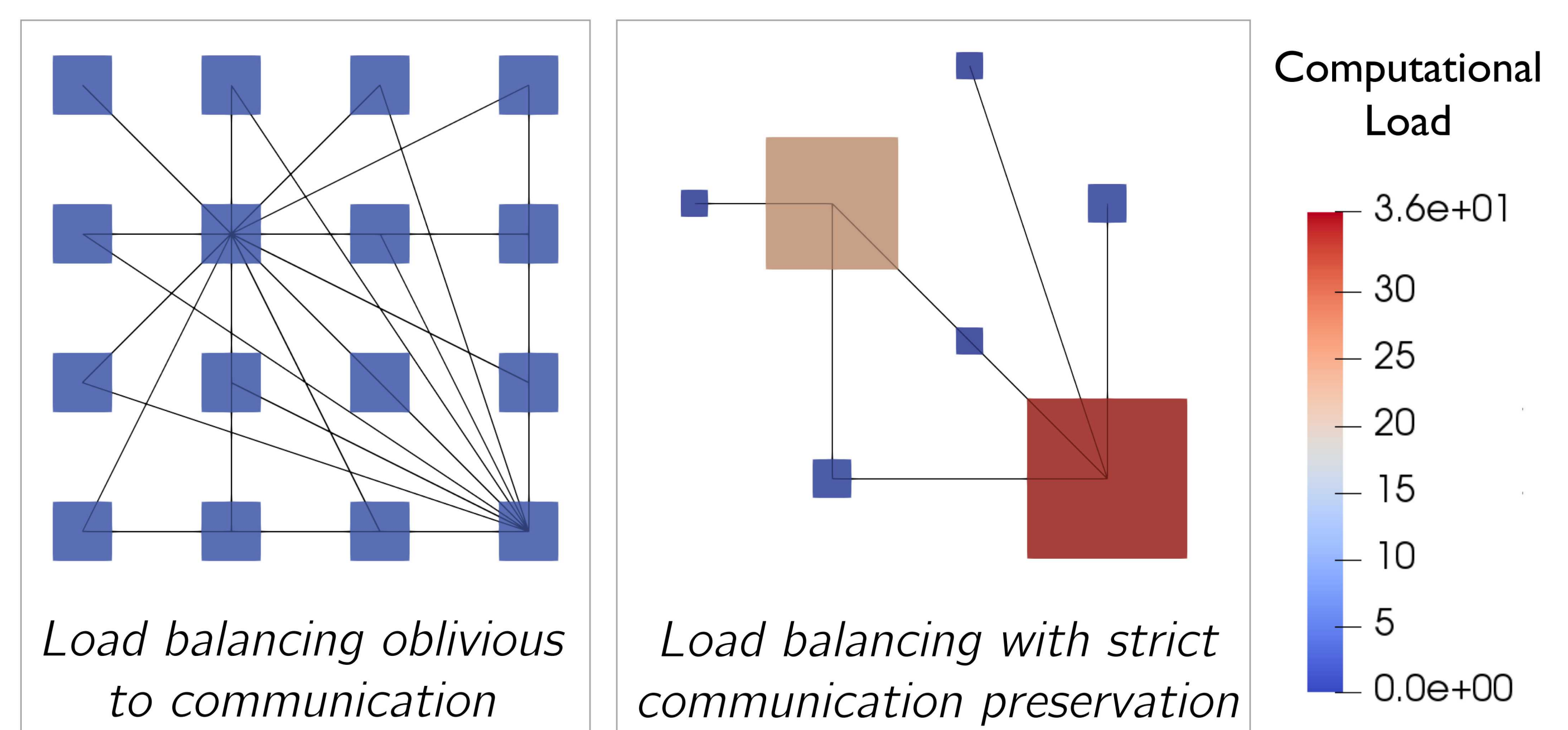- Development and tuning are driven by EMPIRE-PIC

## Results

**Proof of concept:** **demonstrated better than 2x speedup** compared to the MPI baseline by arbitrarily over-decomposing particles on an unbalanced problem.



*vt* no LB
MPI
*vt* hierarchical LB

Particle Move Time Per Timestep (seconds)

PIC Timestep Number

**Novel load balancing algorithm in development:**
Standalone LB simulation and analysis framework demonstrates the benefit of communication-aware LB:

- **Iteratively refine workloads** with incremental changes
- Preserve **localized communication graphs**
- Optimize load balance by trading off communication vs. computation imbalance



*Load balancing oblivious to communication*    *Load balancing with strict communication preservation*

Computational Load

3.6e+01
30
25
20
15
10
5
0.0e+00

## Significance

- Enables dynamic load-balancing for imbalanced, time-varying workloads in codes like EMPIRE
- Mitigates performance imbalances on heterogeneous architectures

## Funding

ASC/CSSE: ~3 FTEs, DARMA/EMPIRE integration effort started in FY19