SAND2014-19851C

# Early Experiences Co-Scheduling Work and Communication Tasks for Hybrid MPI+X Applications

**Dylan Stark**, R. Barrett, R. Grant, S. Olivier, K. Pedretti, C. Vaughan
**Sandia National Laboratories**

**ExaMPI Workshop**

**SC14**

**November 17, 2014**

Approved for unclassified unlimited release
SAND2014-XXXX PE
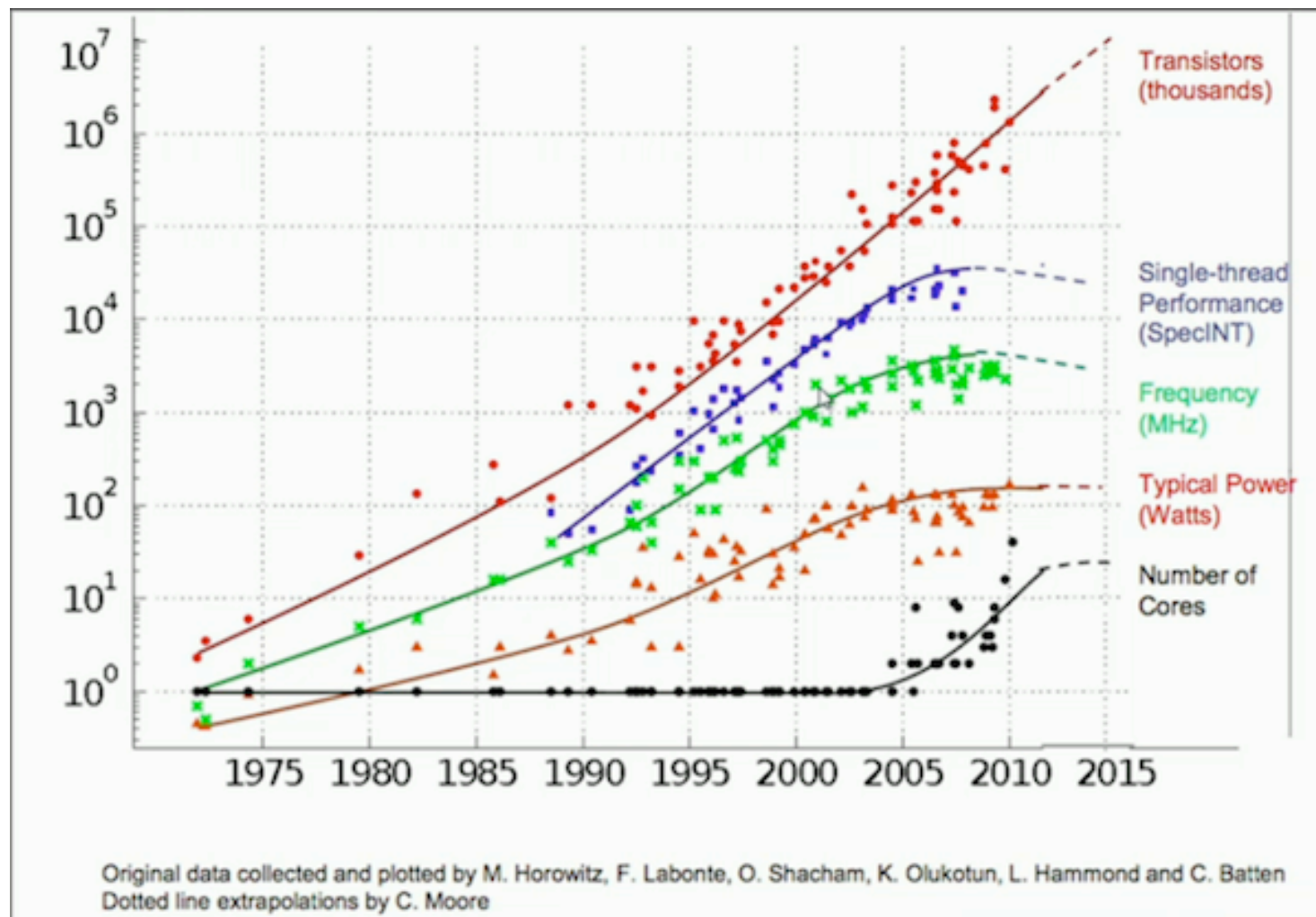
*Exceptional service in the national interest*

Sandia National Laboratories

# End of an era …

Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten
Dotted line extrapolations by C. Moore

# ... where to start?

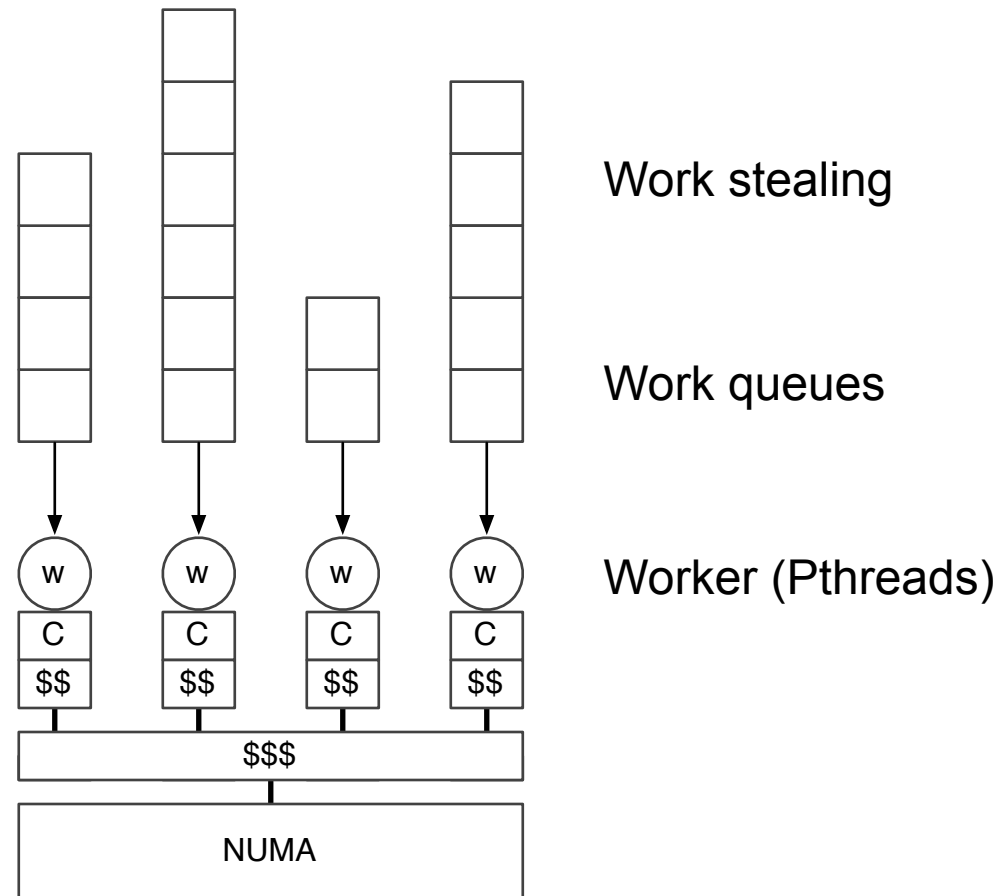| Level | Parallelism | Coordination |
|---|---|---|
| SINGLE | | |
| FUNNELED | Single call site | User's problem |
| SERIALIZED | Multiple calls sites, app-level mutex | User's problem |
| MULTIPLE | Multiple call sites, MPI-level mutex | MPI's problem |

# Leverage low-level threading runtime (LLTR)

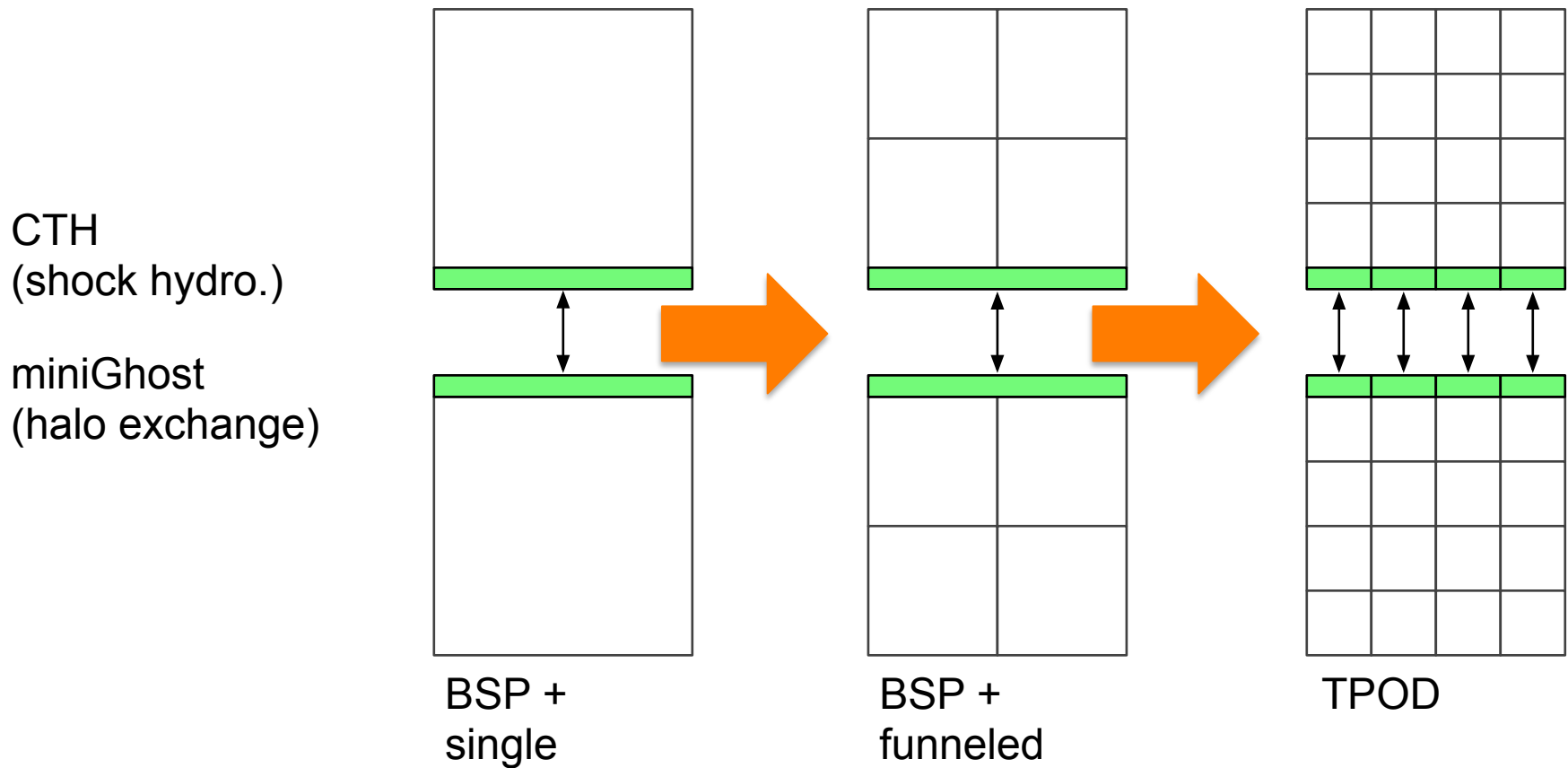| Level | Parallelism | Coordination |
|---|---|---|
| SINGLE | | |
| FUNNELED | Single call site | ~~User's problem~~ <br> **LLTR's problem** |
| SERIALIZED | Multiple calls sites, app-level mutex | ~~User's problem~~ <br> **LLTR's problem** |
| MULTIPLE | Multiple call sites, MPI-level mutex | ~~MPI's problem~~ <br> **LLTR's problem** |

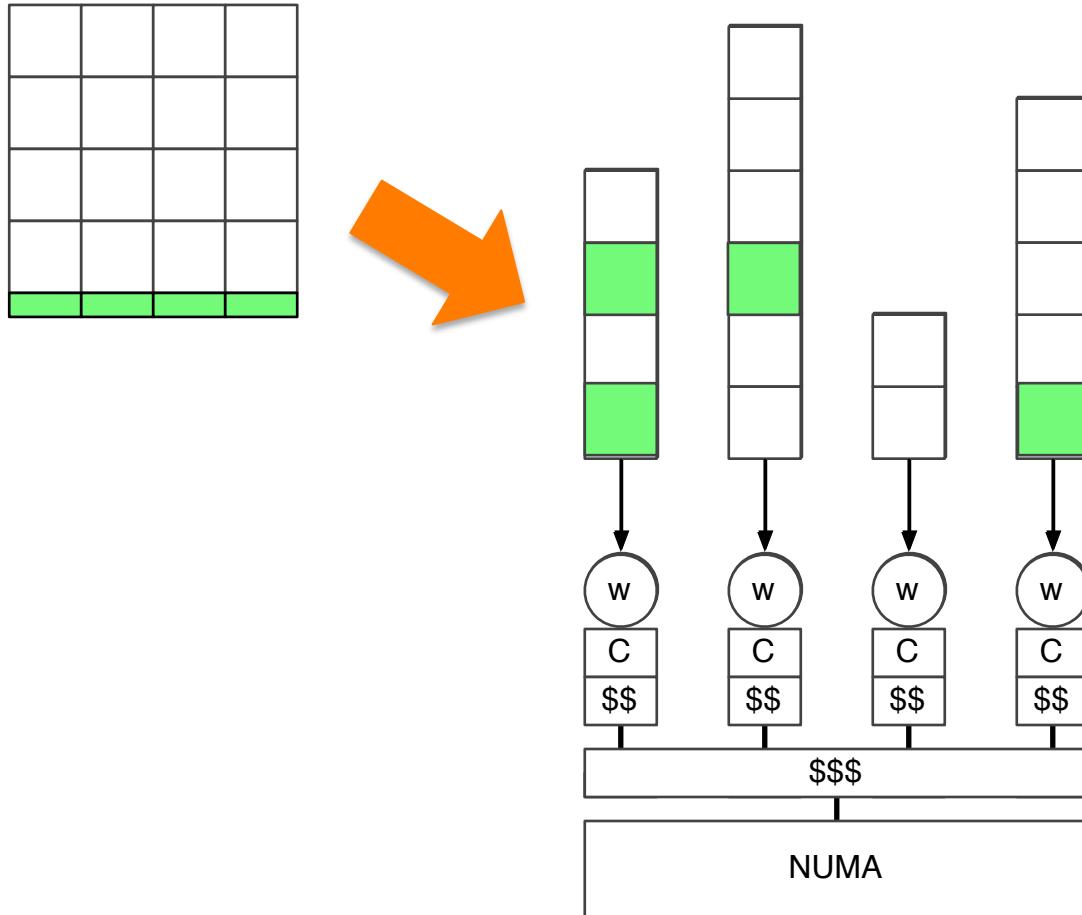# Specifically: cooperative user-level threading (tasking)

*Challenges:*

- State management
- Context swapping
- Memory allocation
- System call interception
- Scheduling
- Synchronization
- Locality
- Hardware affinity
- Concurrent data structures
- Adaptivity
- ...
- Configurability
- Portability

Work stealing

Work queues

Worker (Pthreads)

# Evaluation using a proxy application

CTH
(shock hydro.)

miniGhost
(halo exchange)

BSP +
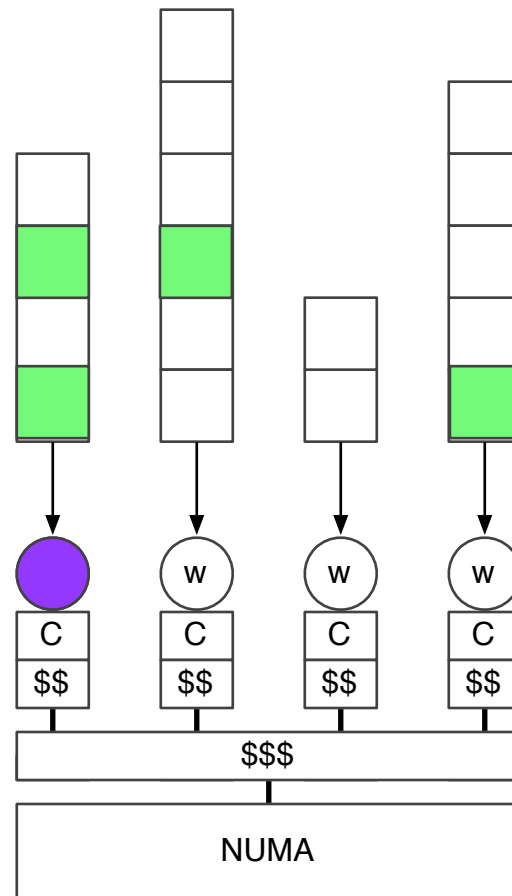single

BSP +
funneled

TPOD

# Design of LLTR

# Design of LLTR
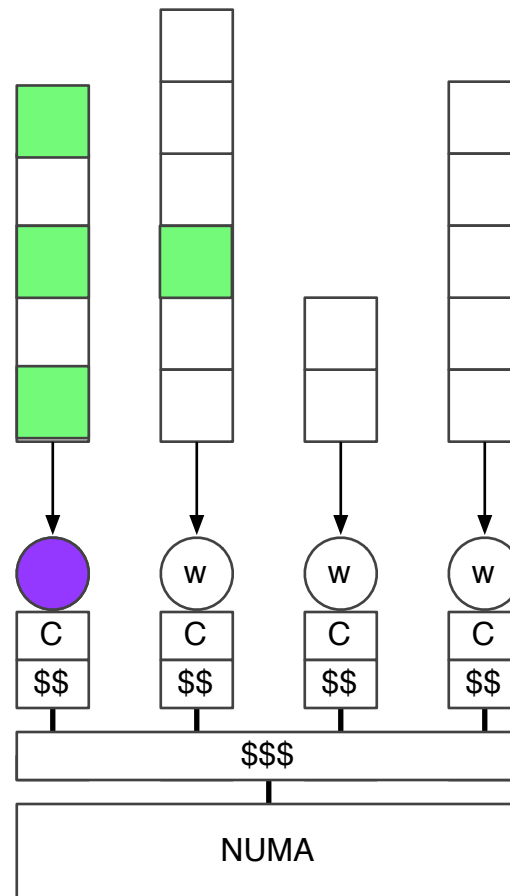
*Challenges:*

- Ensuring thread safety

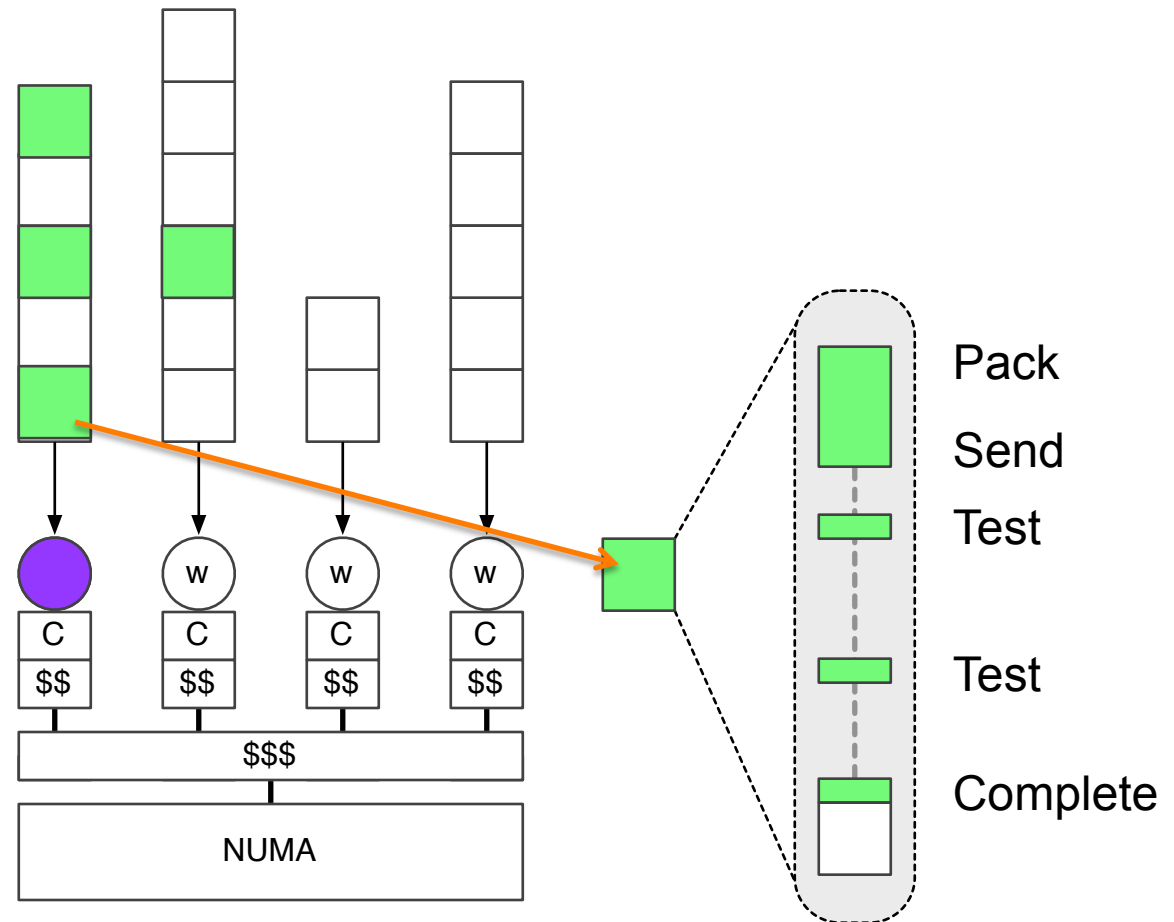# Design of LLTR

*Challenges:*

- Ensuring thread safety
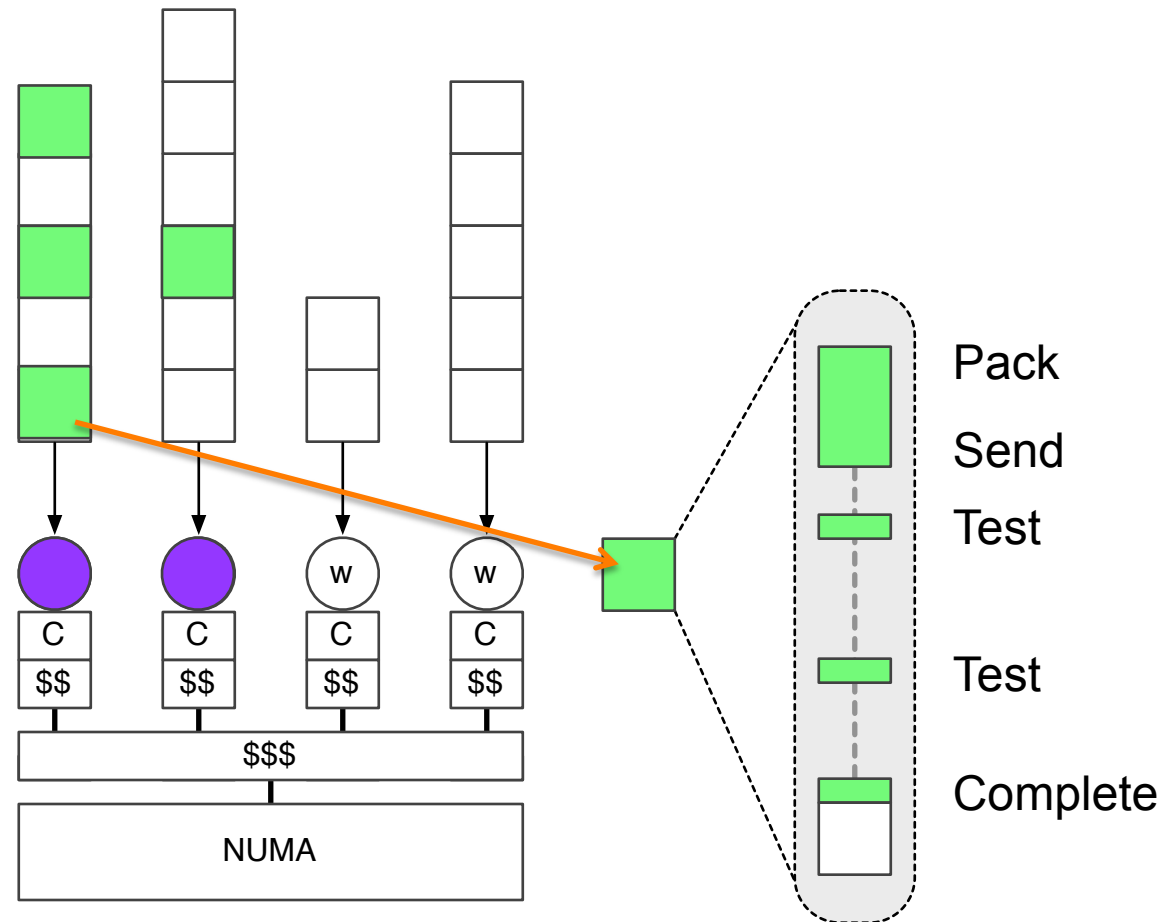
- Avoiding deadlock

# Design of LLTR

*Challenges:*

- Ensuring thread safety

- Avoiding deadlock

- Managing long-latency events
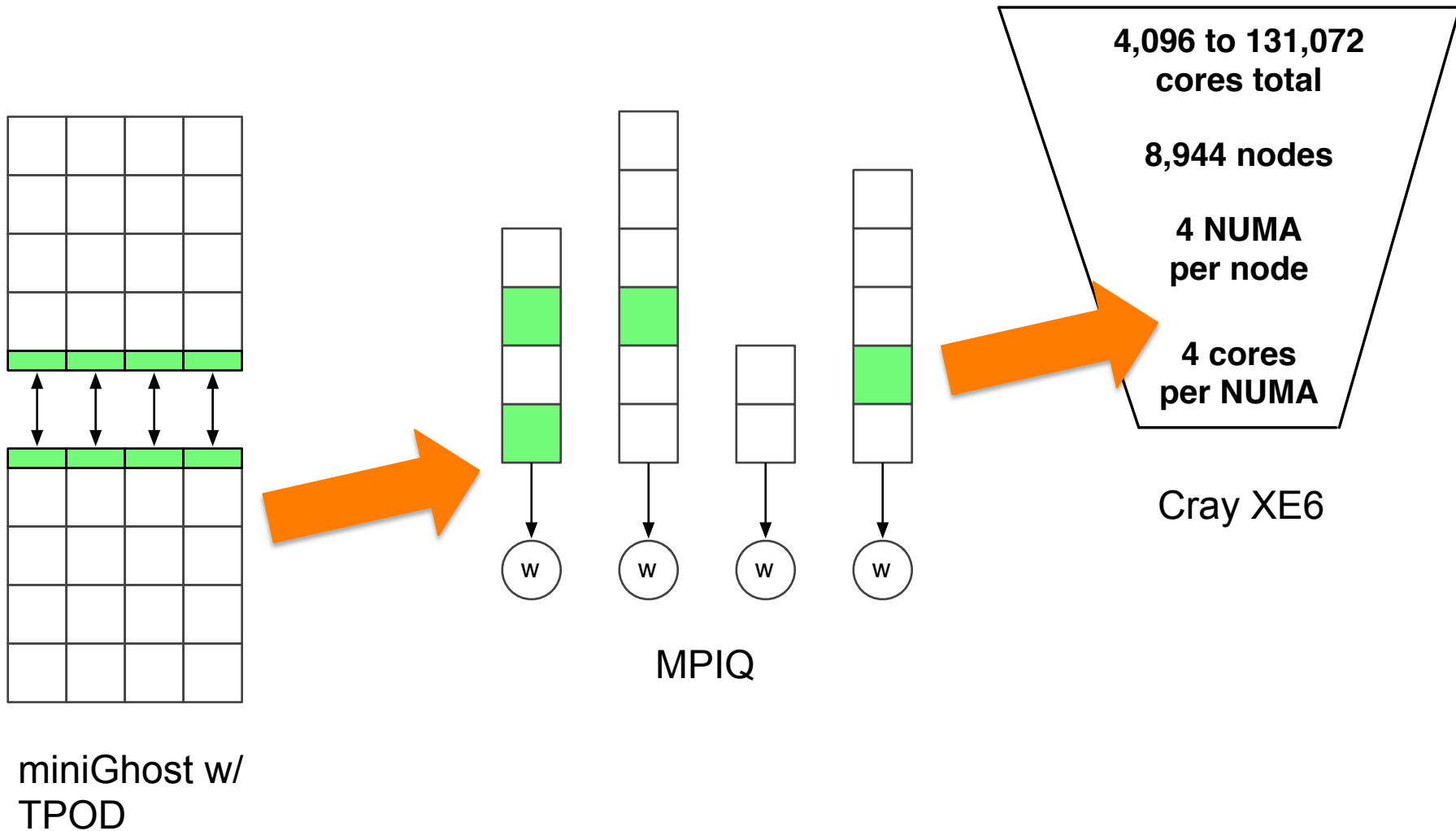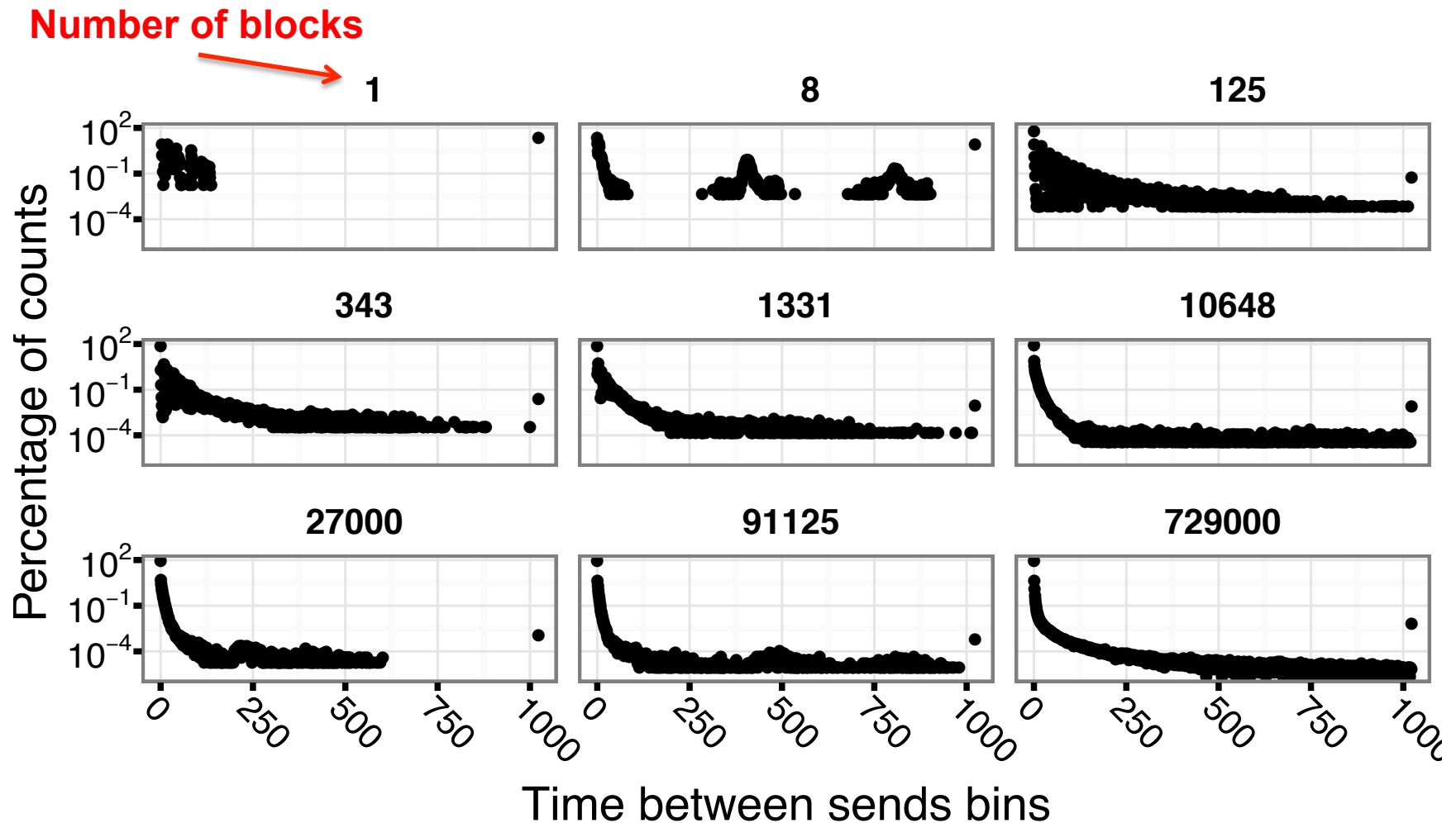
NUMA

Pack
Send
Test
Test
Complete

# Design of LLTR

*Challenges:*

- Ensuring thread safety

- Avoiding deadlock

- Managing long-latency events

- Balancing heterogeneous workload and resources



Pack

Send

Test

Test

Complete

# Bringing it all together



miniGhost w/ TPOD

MPIQ

4,096 to 131,072 cores total

8,944 nodes

4 NUMA per node

4 cores per NUMA

Cray XE6

# Spreading message injection

# Decreasing network stalls

# Increasing communication-computation overlap

# Improving performance, to a point



**Number of cores**

Better performance

More concurrent MPI calls

# Summary

- Co-schedule work and communication using a task parallel runtime

- Imposed minimal code changes on a proxy application

- Showed promising early results up to 131,072 cores

- Lots more to explore …

  - MPI thread support levels, Barriers, Endpoints, …

  - Interoperability with high-level PMs (OpenMP, Kokkos, …)

  - Runtime support for efficient event-driven wake-up, scheduled polling events, …

  - Use cases for AMR, graph analytics, …