

Enabling the Quantum Collisions ‘PFARM’ code on Xeon Phi

Refactoring MPI communications for Symmetric Mode

Michael Lysaght¹, Andrew Sunderland², Martin Plummer²

¹Irish Centre for High End Computing (ICHEC), Ireland

²STFC, Daresbury, UK

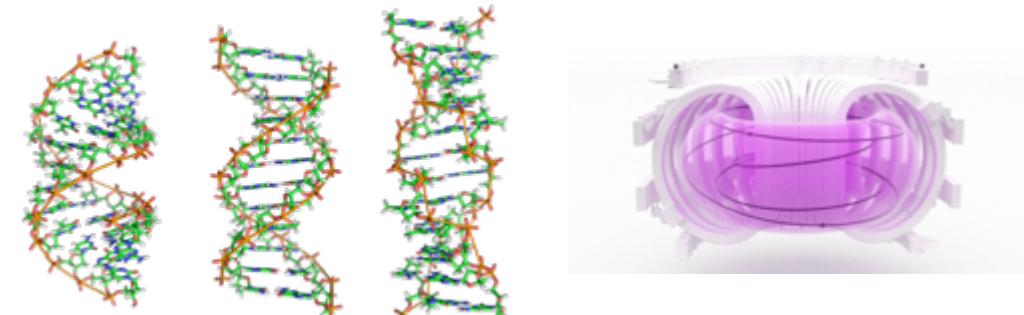
The CCPQ electron-atom and electron-molecule code suites

PRMAT and UKRMol

- Used to simulate physics of Quantum Collisions

- Scientific applications include data for:

- Studies on DNA strand-breaking
- Modeling of fusion plasmas
- Modeling of astrophysical and industrial plasmas



- Physics described by R-matrix Theory

- developed by Wigner & Eisenbud, developed for atomic physics by P. G Burke and colleagues

- Codes developed variously at Queen's University Belfast, University College London, Royal Holloway London, The Open University, STFC and internationally

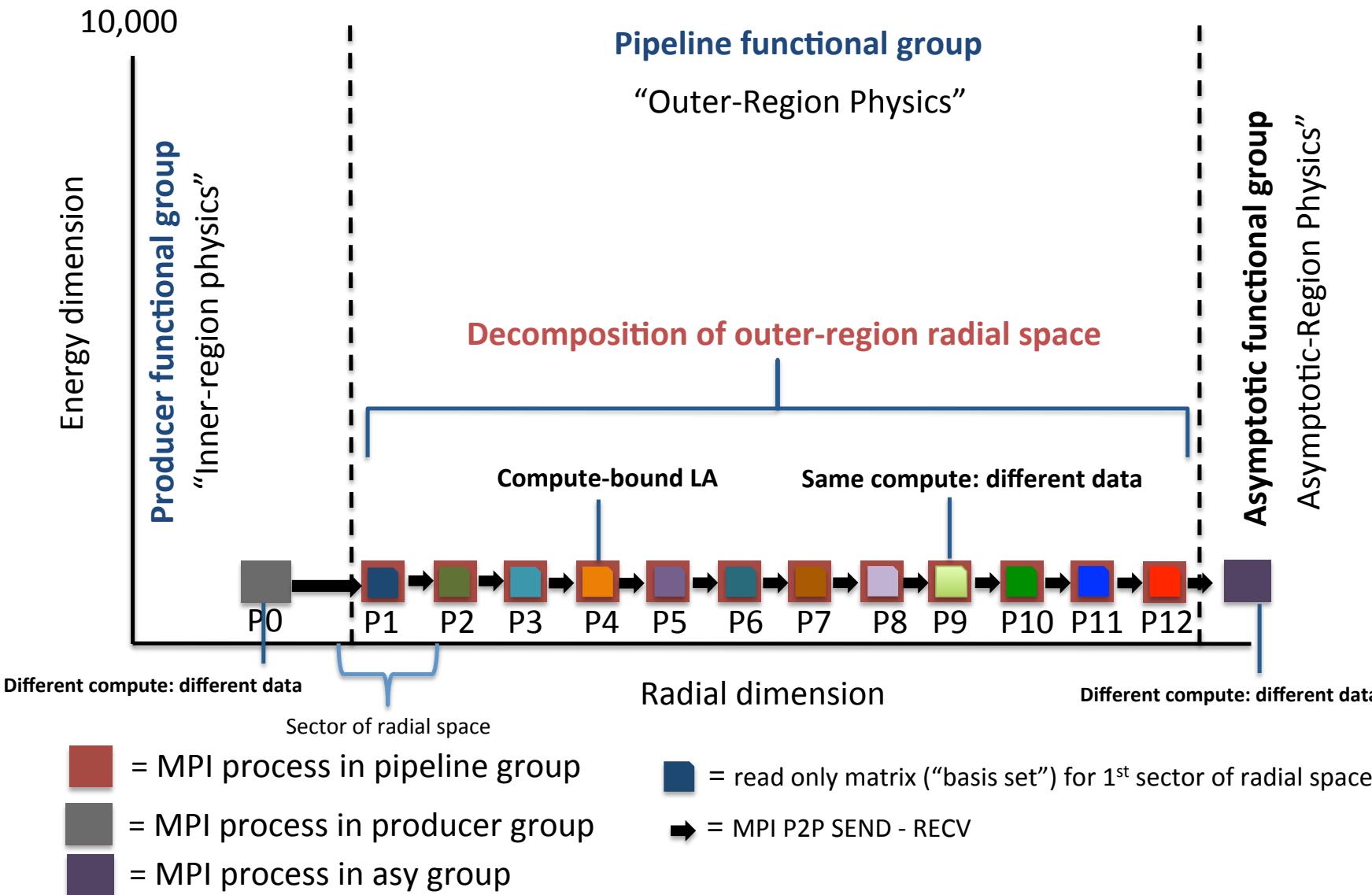
- PFARM, common to PRMAT and UKRMol, is a PRACE Benchmark Code

- Shown to scale well on up to ~30k cores (IBM BGQ)
- Enabling work on EXDIG stage has resulted in significant speedup on Xeon Phi in offload mode

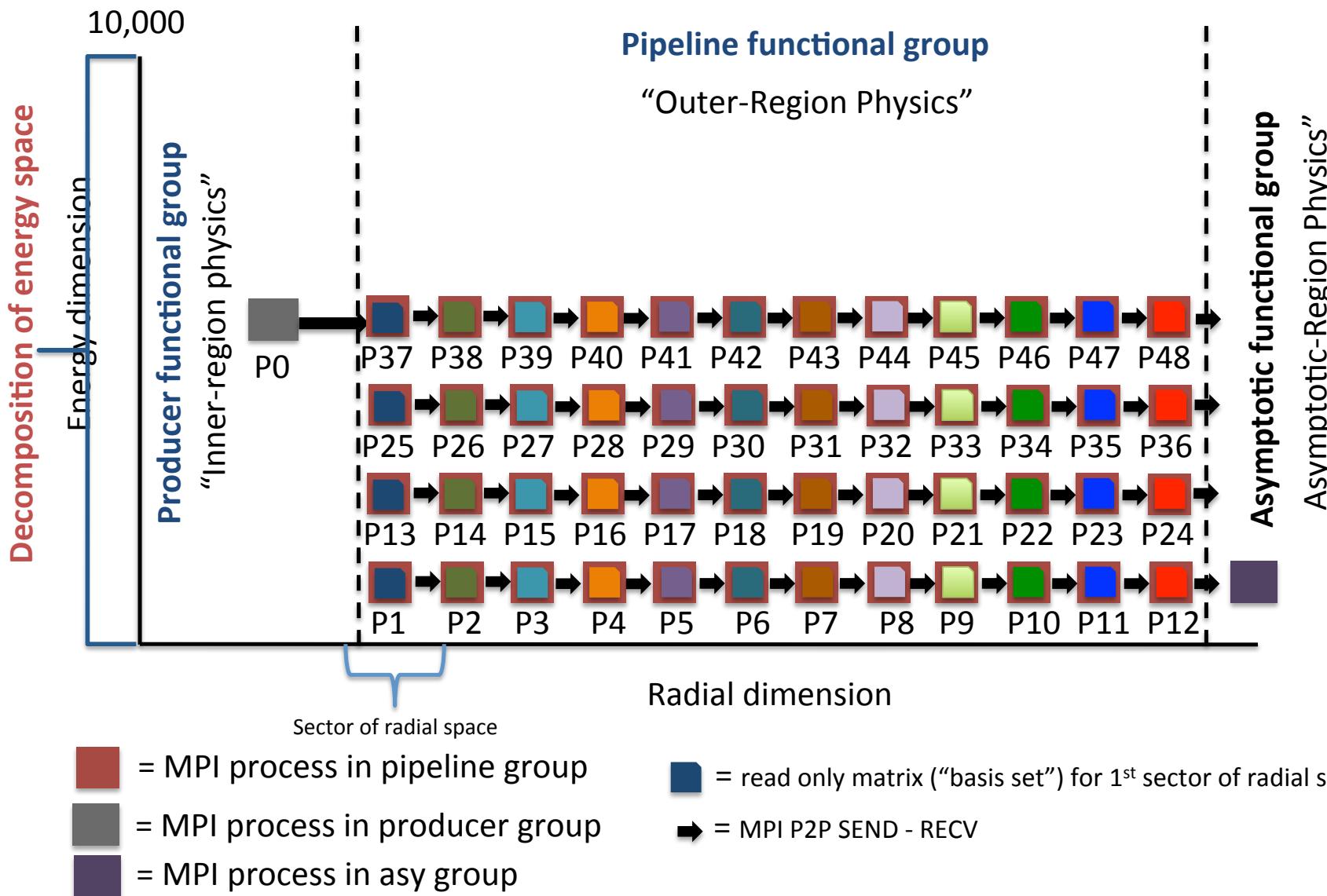
- Developed using Fortran & MPI (since mid-90s), plus OpenMP, CUDA, MAGMA-MIC

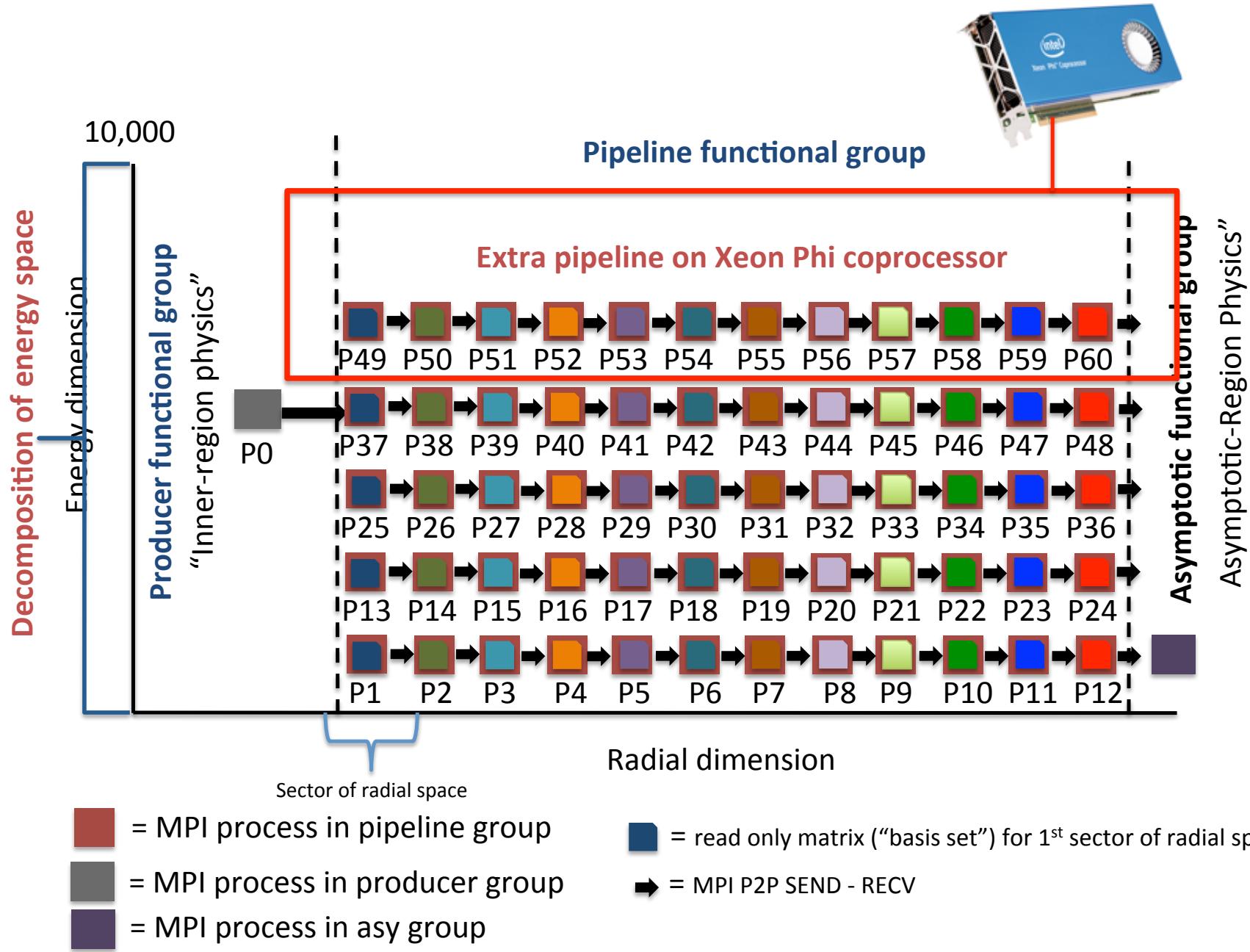
- Available for download from CCPForge

Solving coupled PDEs in multi-dimensional configuration space



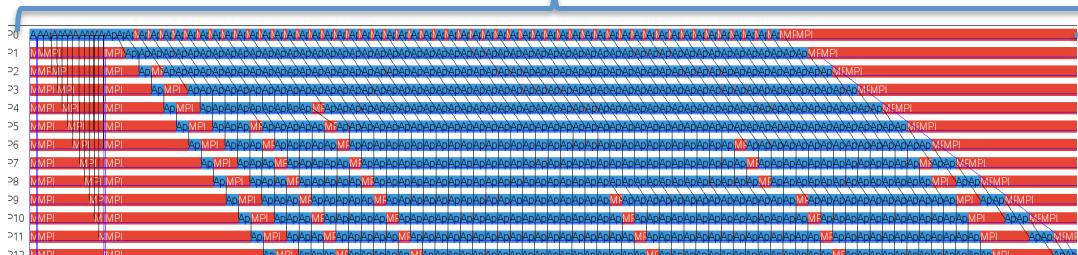
**How can we exploit Xeon Phi Coprocessor with this design?
High core count on Xeon Phi suggests more pipelines possible!**





Exploiting ITAC for Deeper Insight

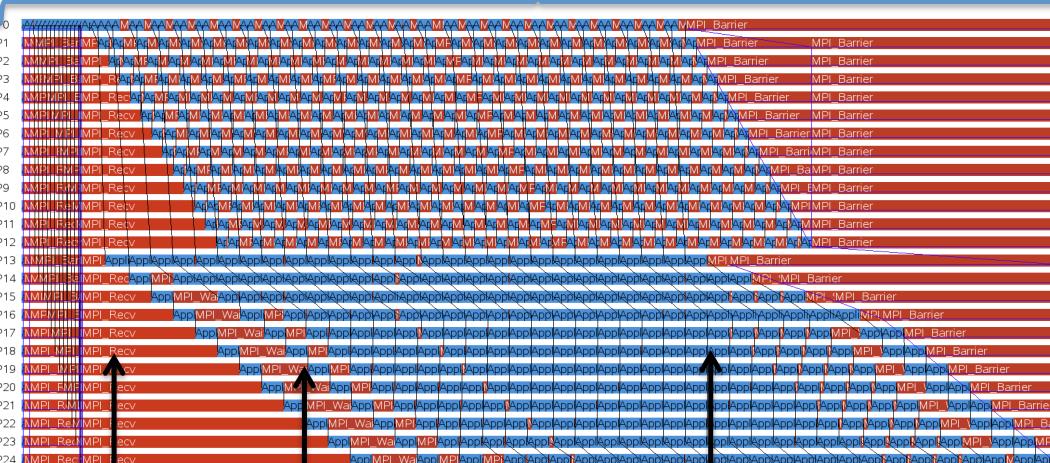
Time-to-solution: 546 secs



Original PFARM Code on Host only

Host Processes (0-12)
Single pipeline

Time-to-solution: 727 secs



Original PFARM Code Host + 1 x Xeon Phi Coprocessor

Host Processes (0-12)
1st pipeline

Xeon Phi Processes (13-24)
2nd pipeline

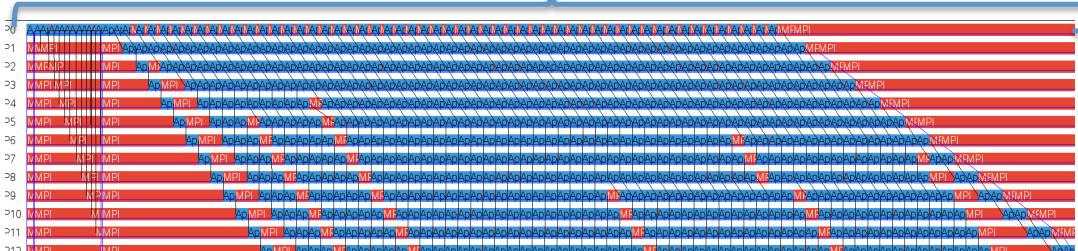
Compute phase takes longer on Xeon Phi

Higher overhead of MPI comms and MPI_WAIT on Xeon Phi

Poor load-balancing between Host and Xeon Phi

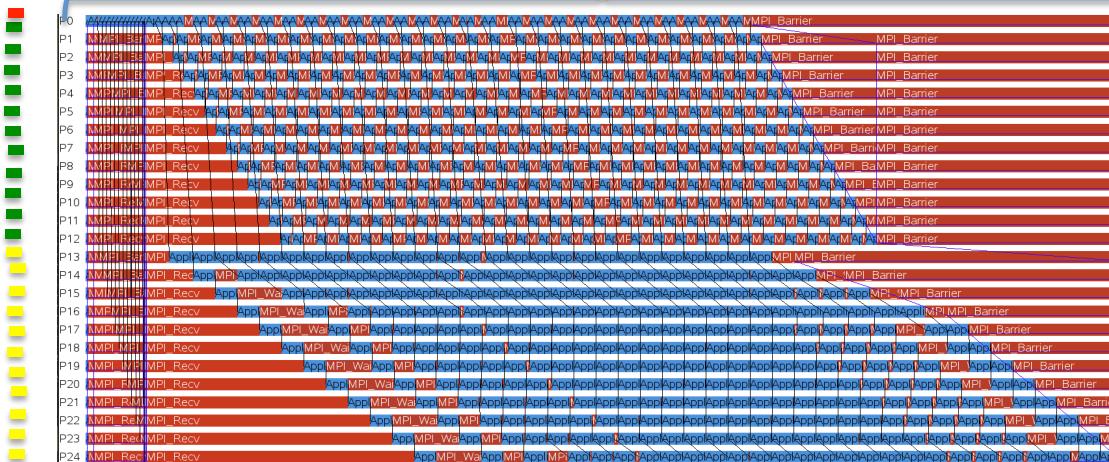
Exploiting ITAC for Deeper Insight

Time-to-solution: 546 secs



Original PFARM Code on Host only

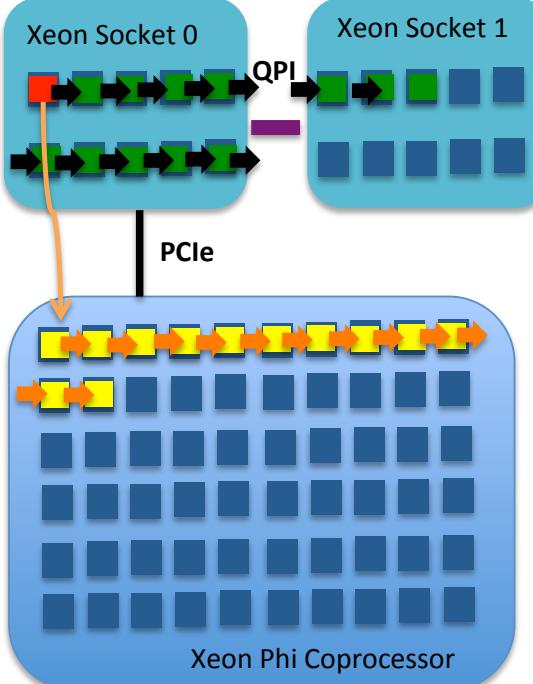
Time-to-solution: 727 secs



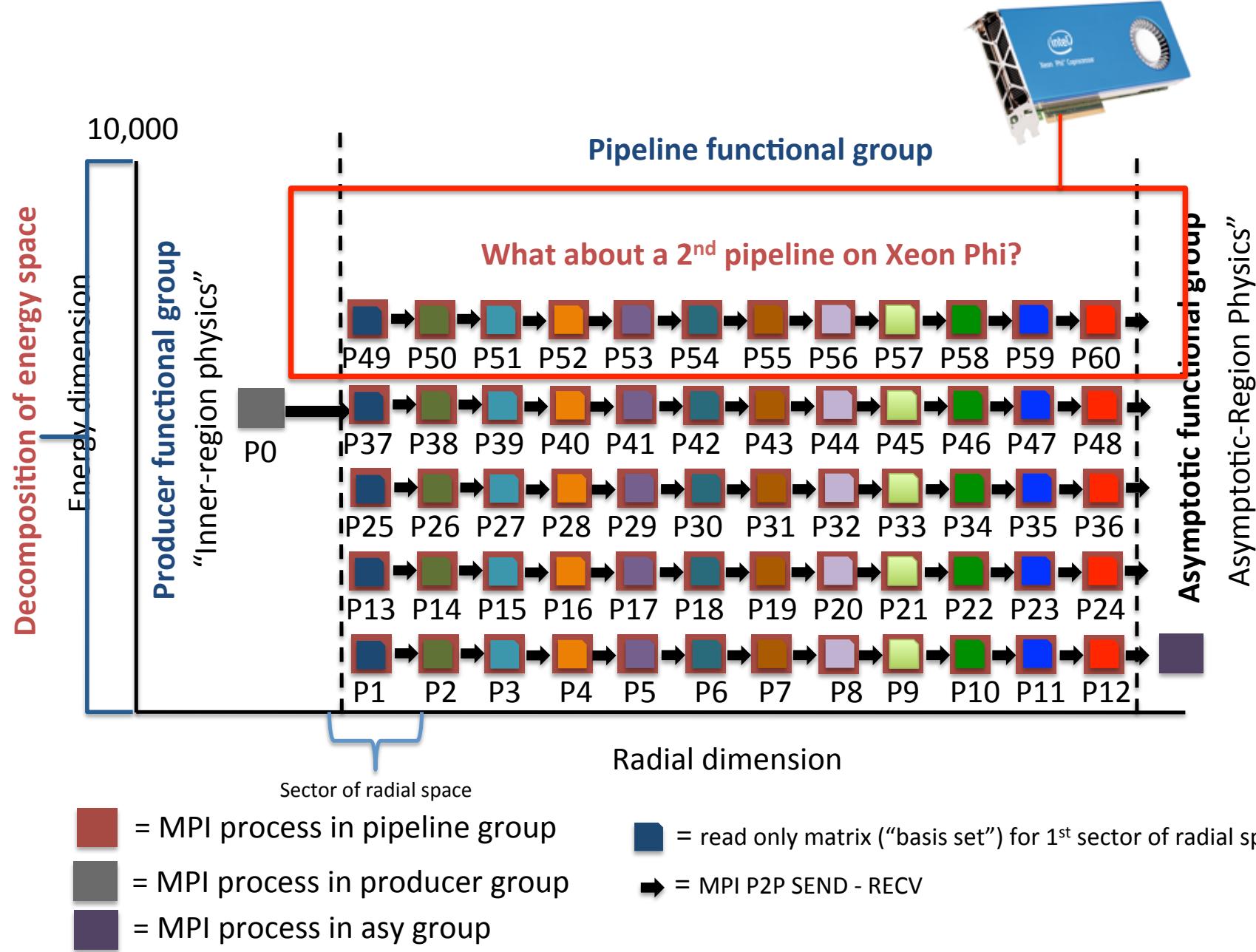
Original PFARM Code Host + 1 x Xeon Phi Coprocessor

- = producer process (fast compute)
- = 1st pipeline stage process (fast compute)
- = 2nd pipeline stage process (slow compute)
- = fast MPI comms (on Xeon)
- = slow MPI comms (on Xeon Phi)

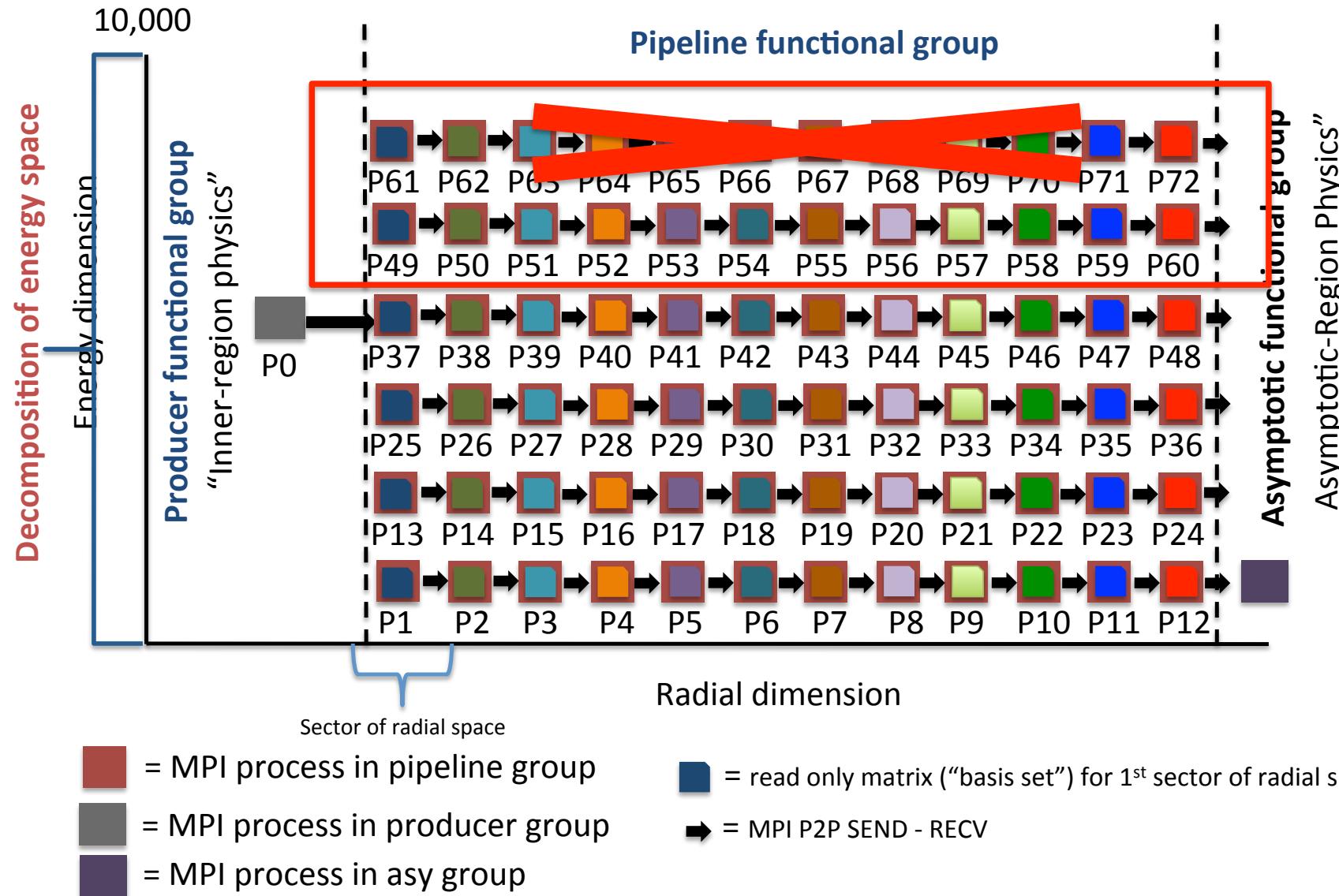
Sees resources as a homogeneous system



ITAC results strongly suggest that the only way to benefit from Xeon Phi is to increase the number of pipelines on the Xeon Phi – 1 pipeline is not enough!



Problem: Memory footprint grows too rapidly with increase in pipeline count!

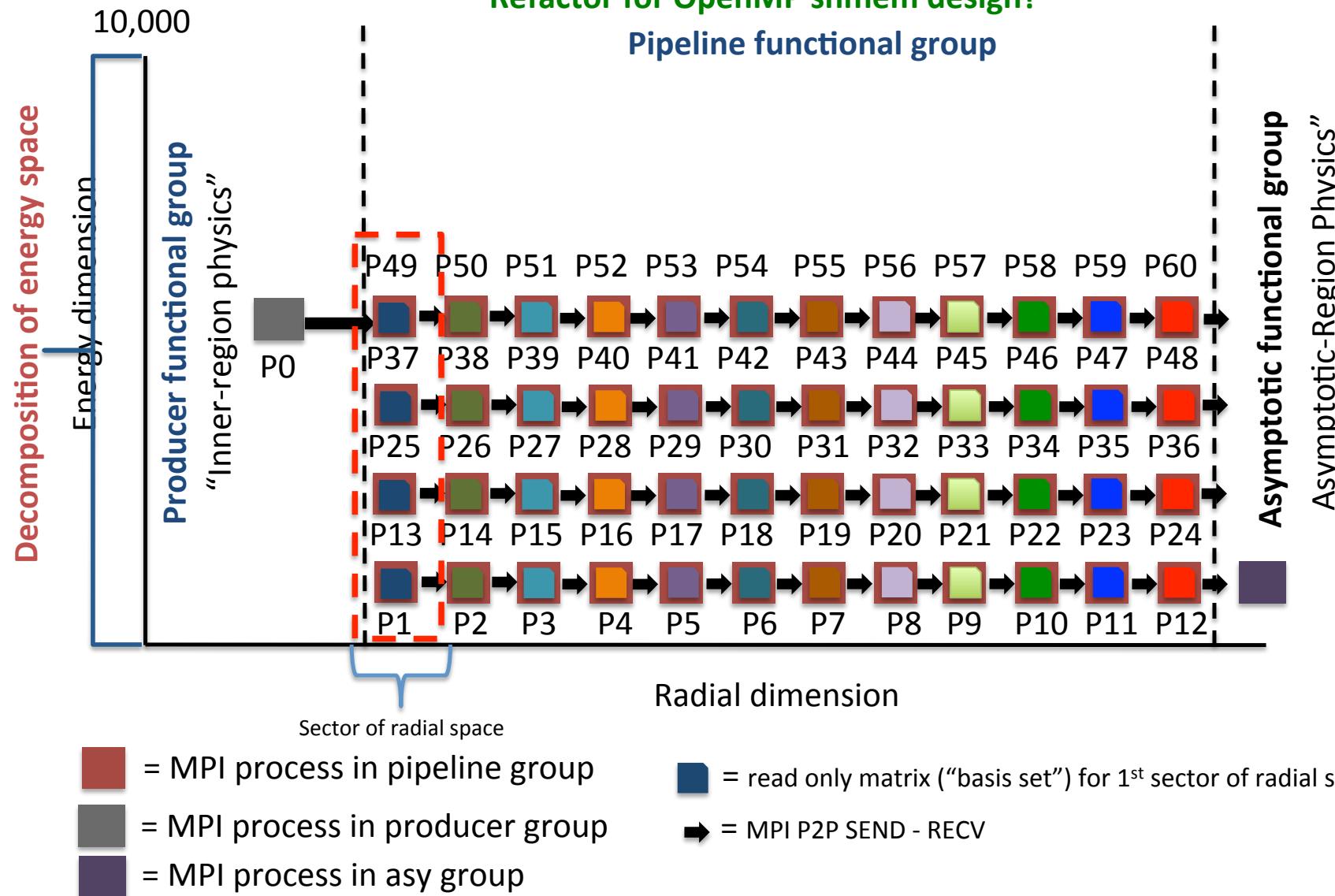


Can we share the replicated data across parallel tasks?

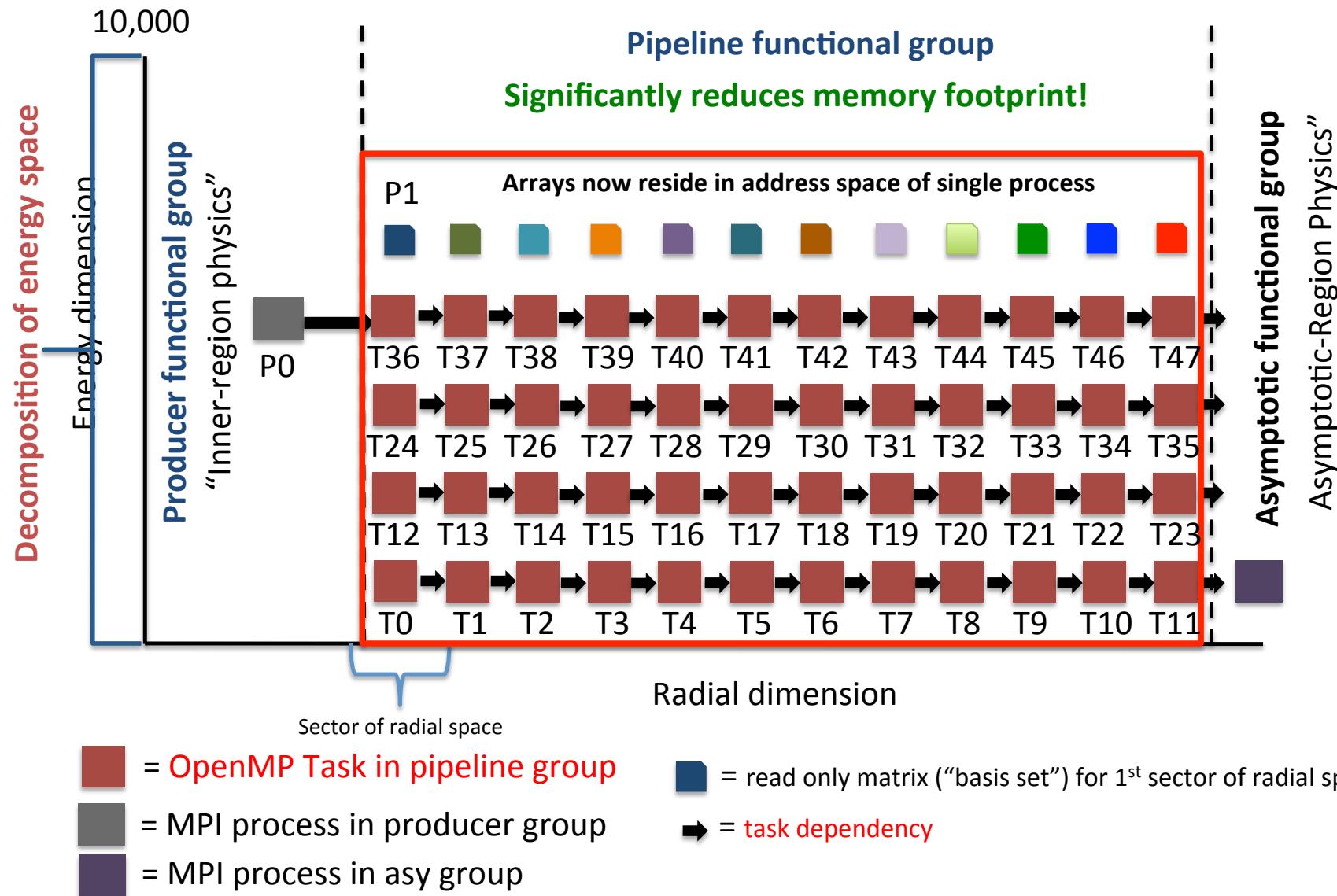
MPI 3.0 approach = tricky

Refactor for OpenMP shmem design?

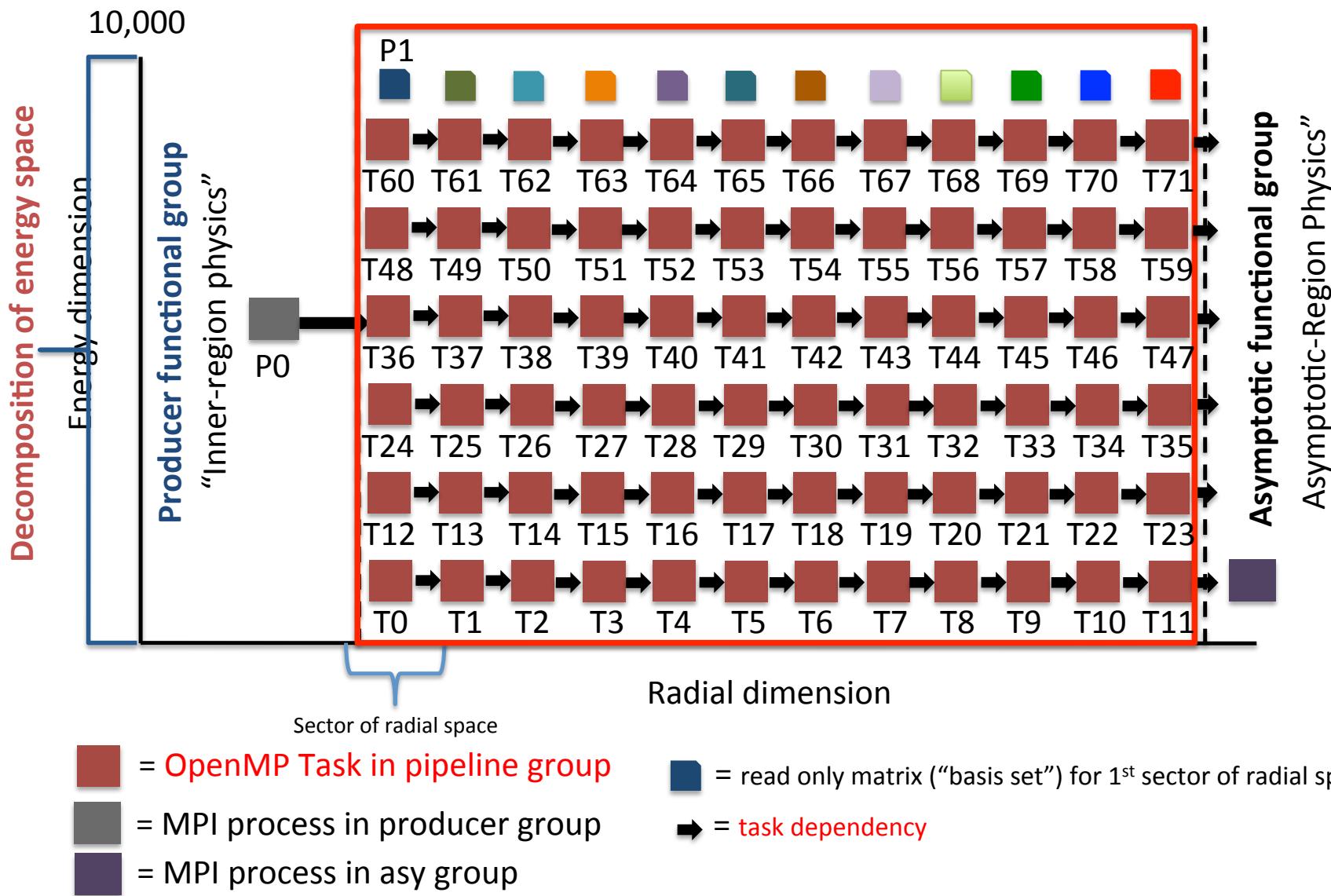
Pipeline functional group



Can we share the replicated data across parallel tasks? Refactoring for OpenMP shmem design

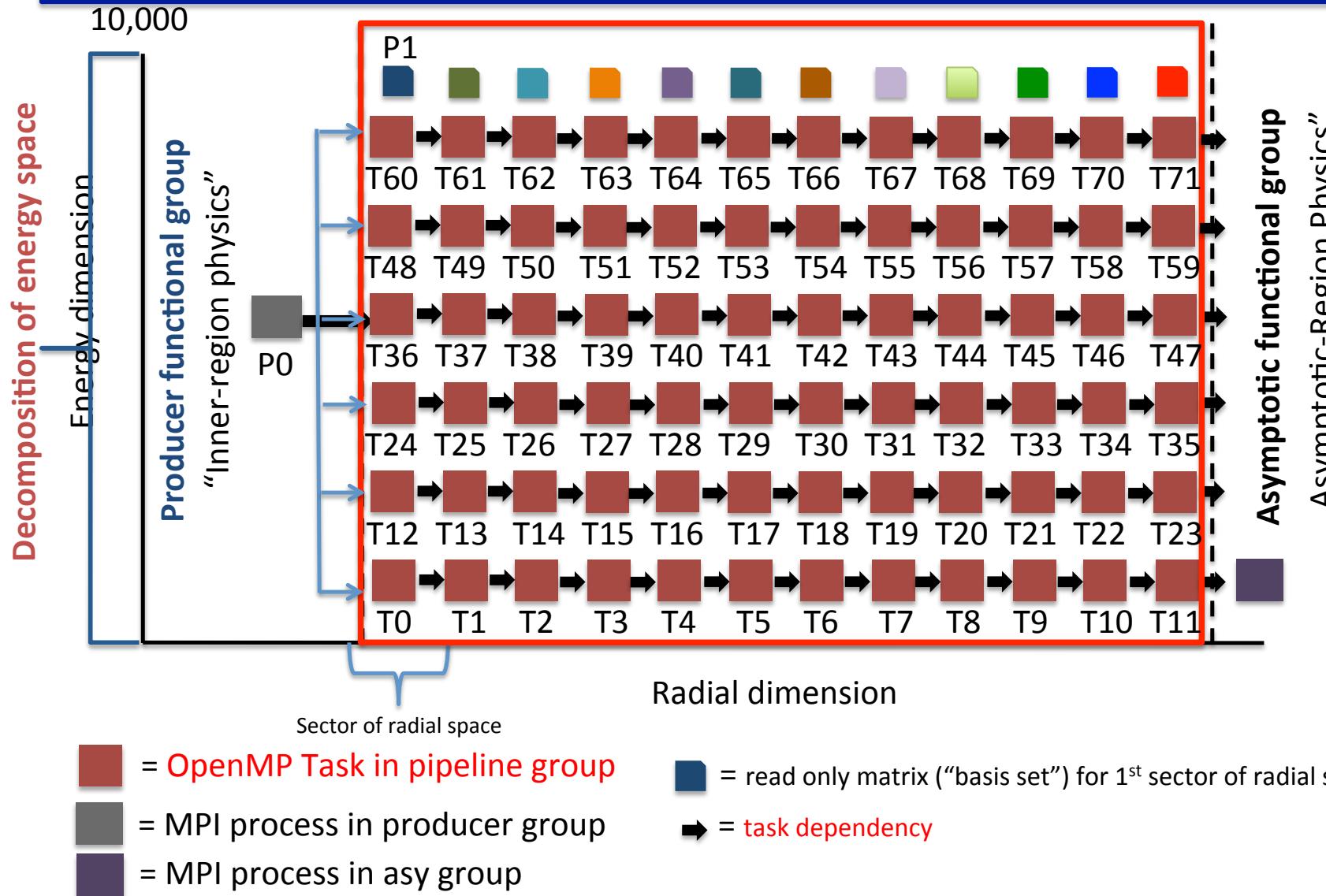


Can we share the replicated data across parallel tasks?
Refactoring for OpenMP shmem design



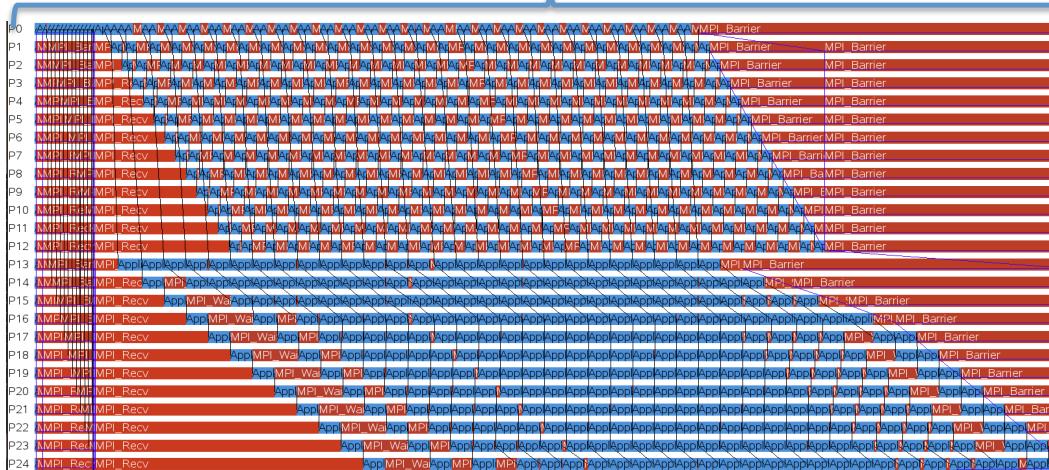
Solution: Enable MPI_THREAD_MULTIPLE support

OpenMP tasks/threads can each call MPI_RECVs from within OMP Parallel Region (on P1)
Can maintain the MPI-based Producer/Consumer model at interface of functional groups



Results and Deeper Insight from ITAC

Time-to-solution: 727 secs

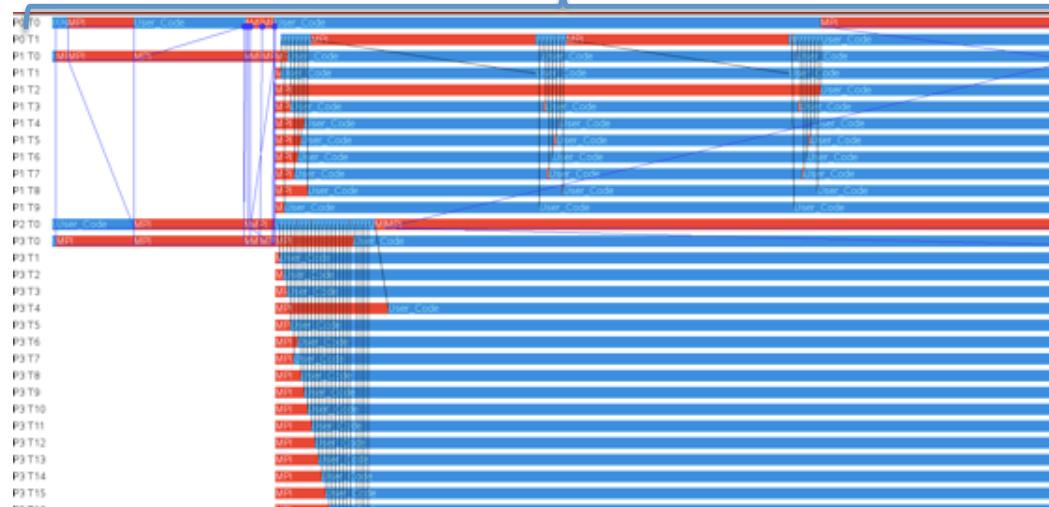


Original PFARM Code Host + 1 x Xeon Phi Coprocessor

Host Processes (0-12)
1st pipeline

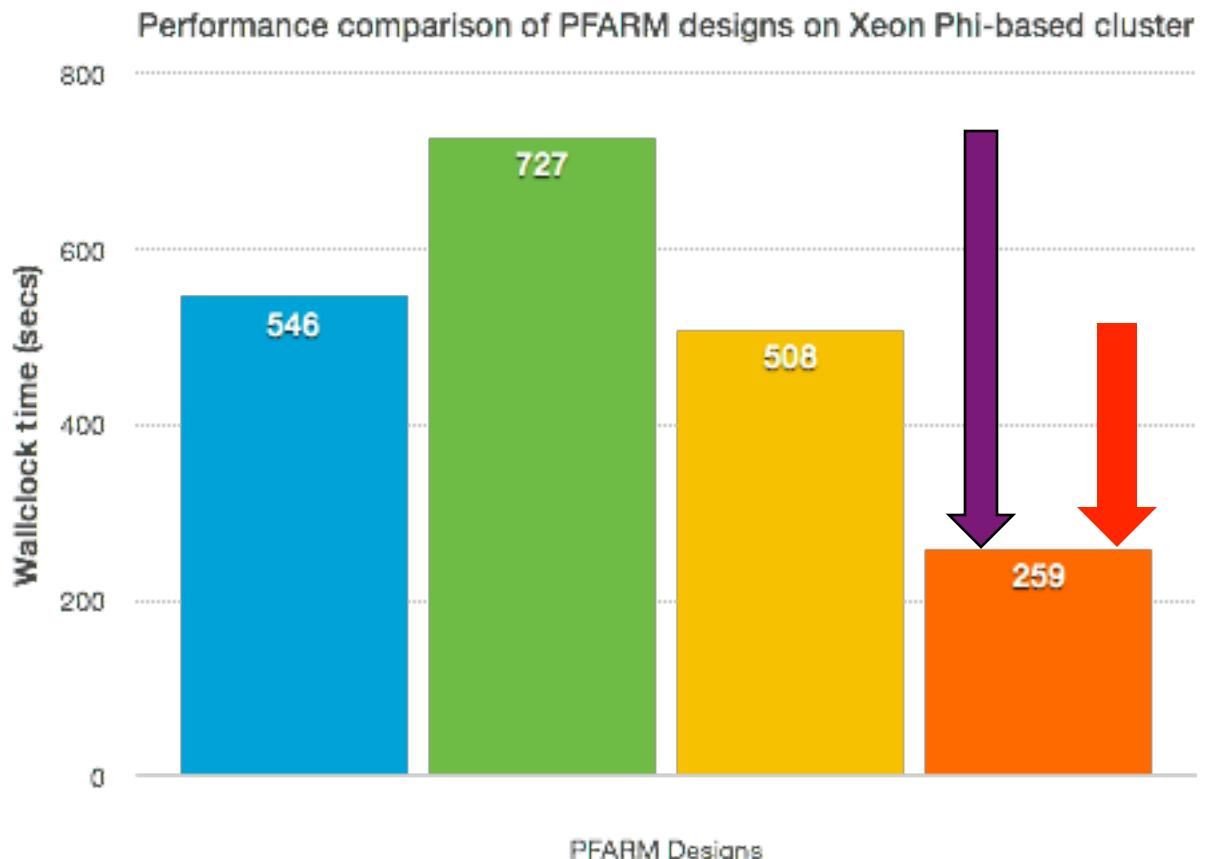
Xeon Phi Processes (13-24)
2nd pipeline

Time-to-solution: 259 secs

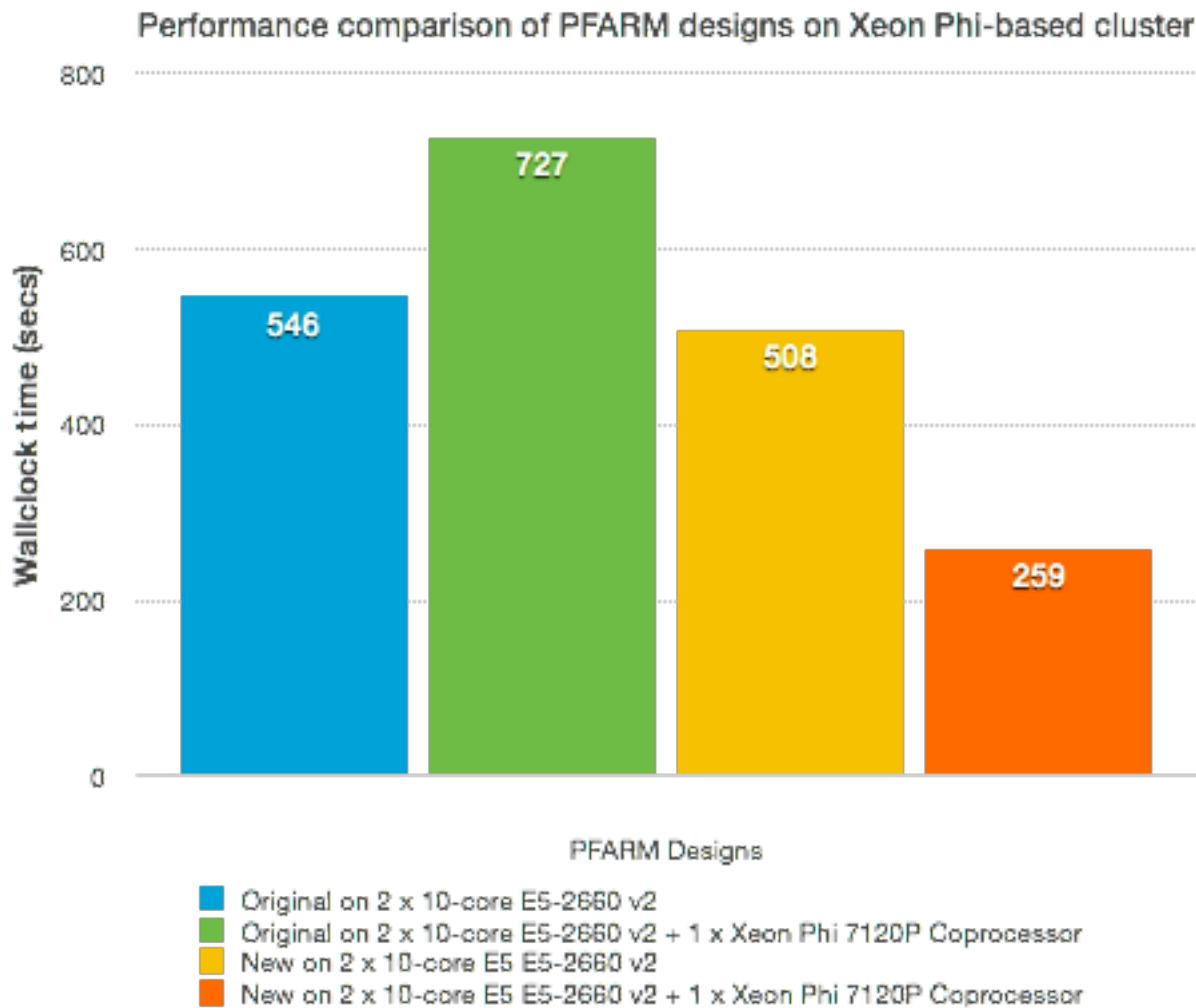


Host Threads (0-10)
1st 10 pipelines

Xeon Phi Threads (0-29)
Next 30 pipelines



- New design leads to 2.8x speedup vs original design on Xeon + 1 x Xeon Phi coprocessor
- New design on 2S Xeon IVB + 1 x Xeon Phi coprocessor leads to ~2x improvement in time-to-solution vs new design on 2S Xeon IVB only
- *New design exposes more parallelism to further improve performance on Xeon Phi*



- ITAC can help to inform and guide the non-trivial redesign of MPI codes on heterogeneous systems
- Redesigning to exploit OpenMP shmem can lead to increase in parallelism
- MPI_THREAD_MULTIPLE support can be useful when redesigning for OpenMP shmem parallelism
- Currently exploring load-balancing issues and the exposure of nested parallelism
- Also preparing for KNL – code correctness with SDE and opportunities for HBM
- One design of many – feedback from IXPUG community welcome!

BACKUP

Why did we design the trunk version of PFARM like this ?

- 1. Limited memory per node on past architectures (still is a limitation for some problems). Pipelines can take full advantage of overlapping comms/comps.**
- 2. Very flexible - parallel groups adjusted/multiplied according to target datasets.**
- 3. It generally gets ~30% peak flops performance per node.:
http://www.hpcx.ac.uk/research/hpc/technical_reports/HPCxTR0703.pdf**