

Accelerated Earthquake Simulations

Alex Breuer
Technische Universität München
Germany



Acknowledgements

- Volkswagen Stiftung — Project ASCETE: Advanced Simulation of Coupled Earthquake-Tsunami Events
- Bavarian Competence Network for Technical and Scientific High Performance Computing (KONWIHR)
- Intel Corporation (IPCC, ExScaMIC)
- SuperMUC grants: pr45fi, pr63so
- NSF grant: OSI-1134872 (Stampede)

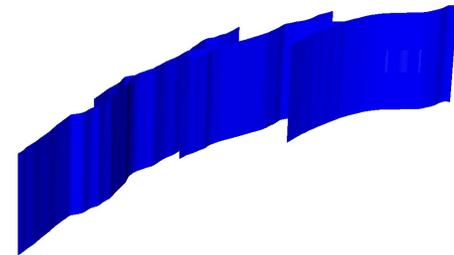
This is joint work!

- Alex Heinecke, Intel Labs
- Sebastian Rettenberger, TUM
- Michael Bader, TUM
- Alice Gabriel, LMU
- Christian Pelties, LMU
- All other colleagues and collaborators @ LRZ, TACC, NUDT, Intel Labs

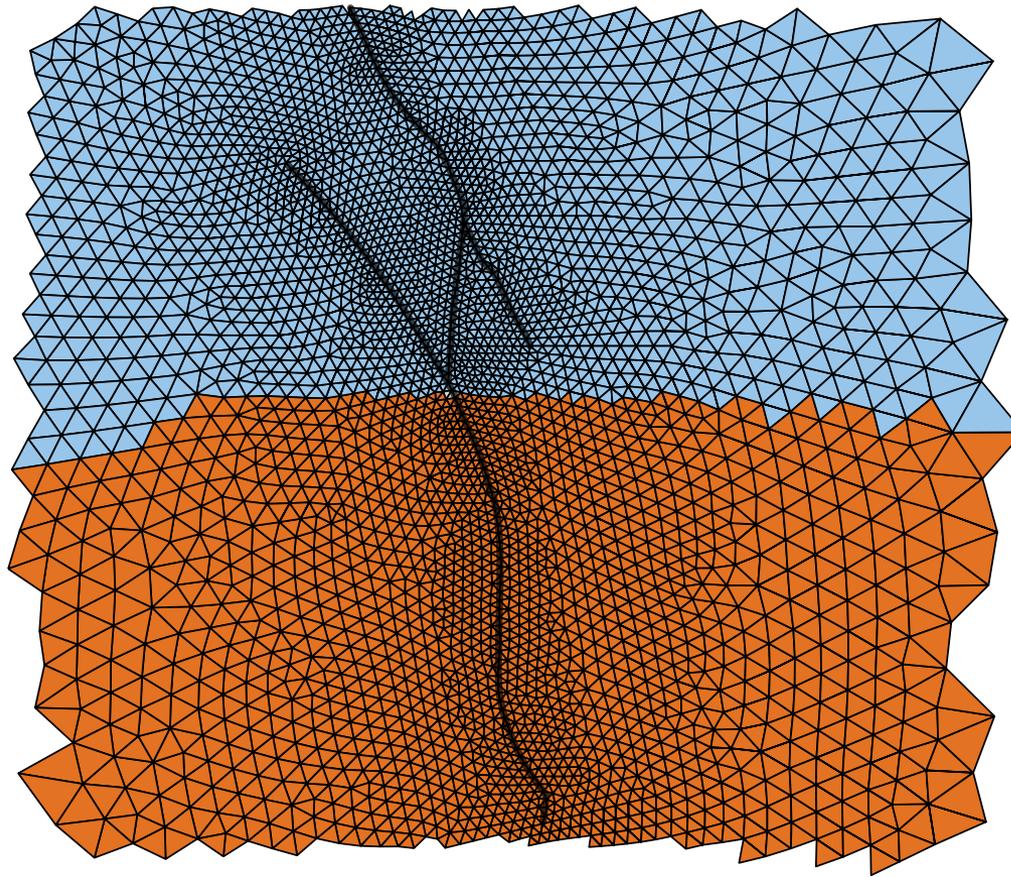
Missing graphic

Earthquake Simulations

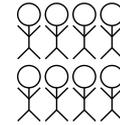
- Application: SeisSol
 - Full elastic wave equations in 3D and complex heterogeneous media
 - Dynamic Rupture w.o. artificial oscillations
 - High order: ADER(time)-DG(space)
 - Unstructured tetrahedral meshes
- Domain: Geophysics
- Execution mode: Offloaded
- Intel's LEO, Scalasca



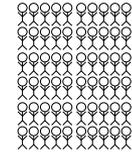
Offloading



Host

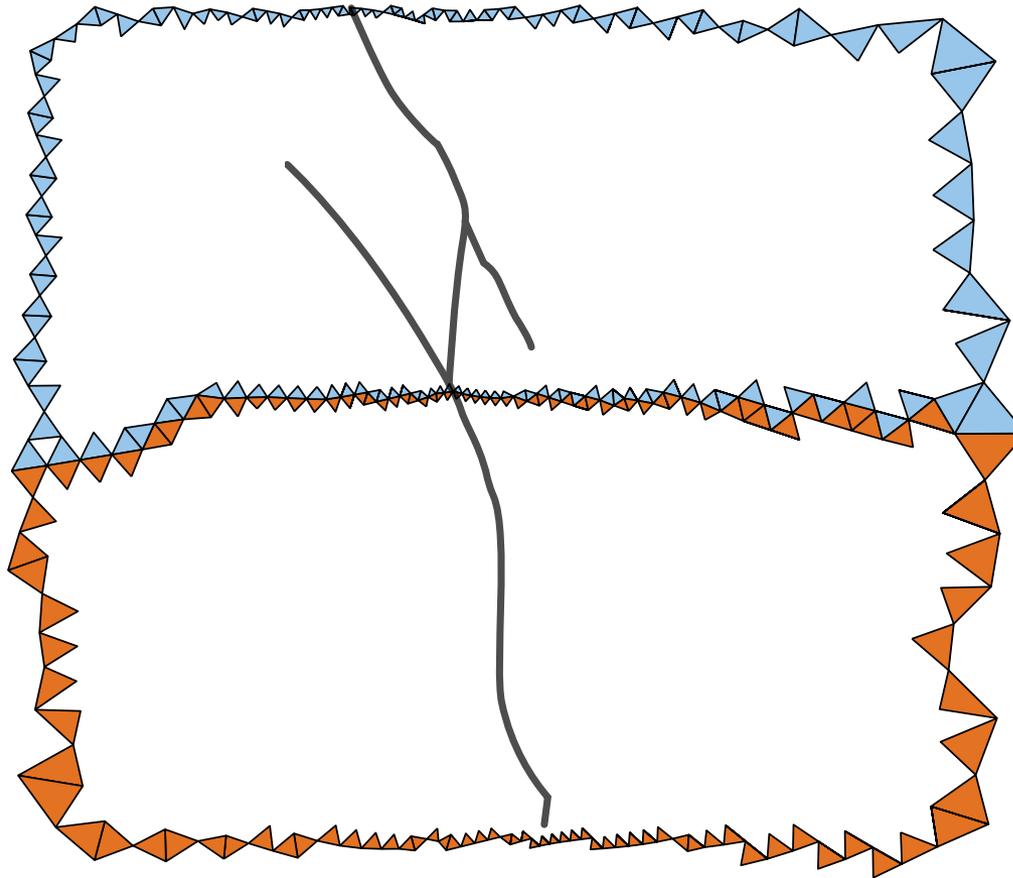


MIC

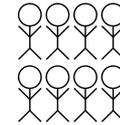


- Refactoring: Introduce control-flow aligned to hardware requirements
- Overlapping Communication / Computation
- Heterogeneous Execution: Complicated DR on fat cores
- Hybrid Parallelization
- Vectorization through custom compute kernels

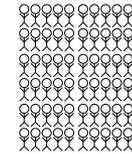
Offloading



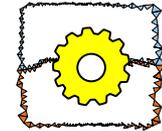
Host



MIC

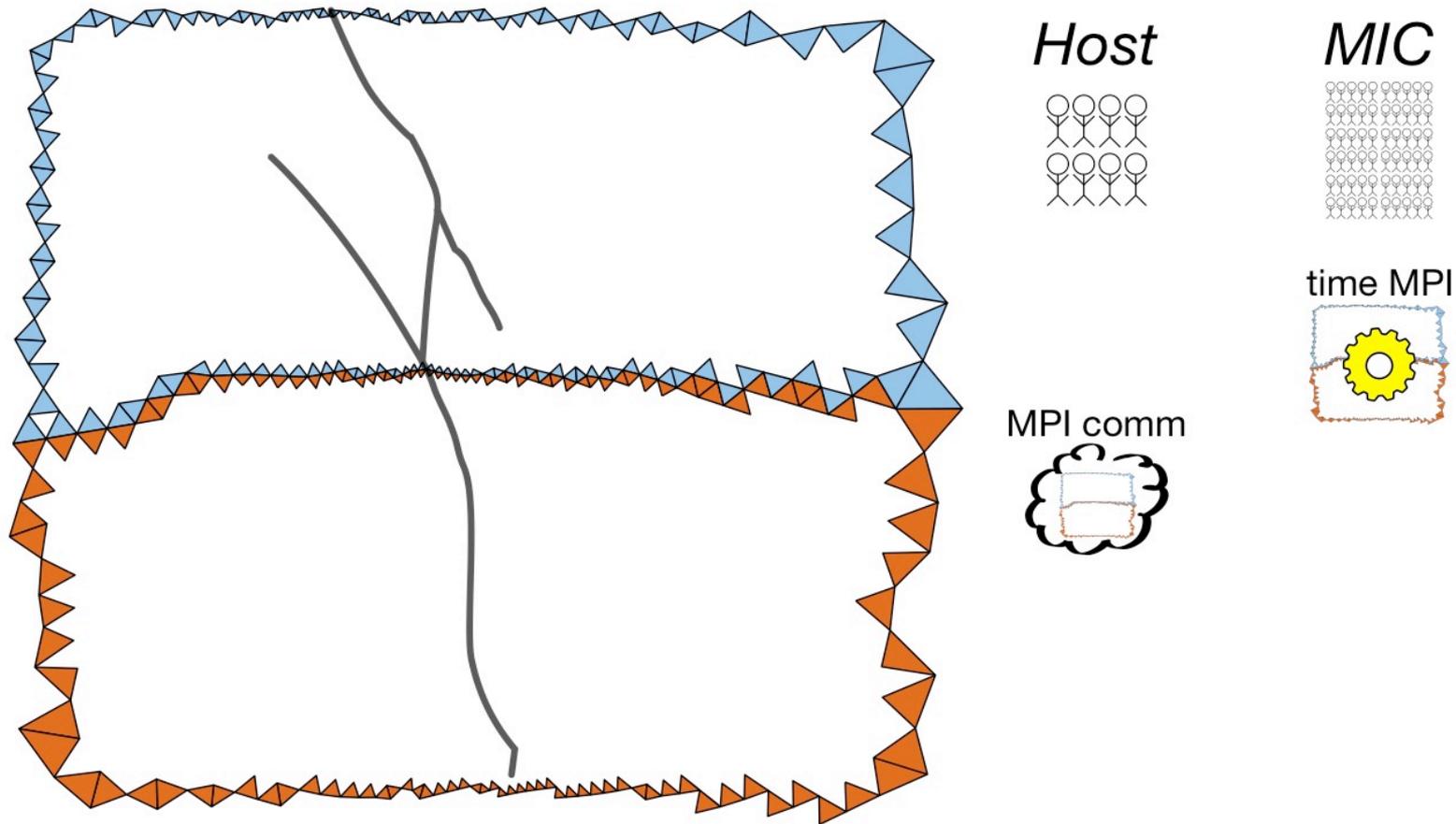


time MPI



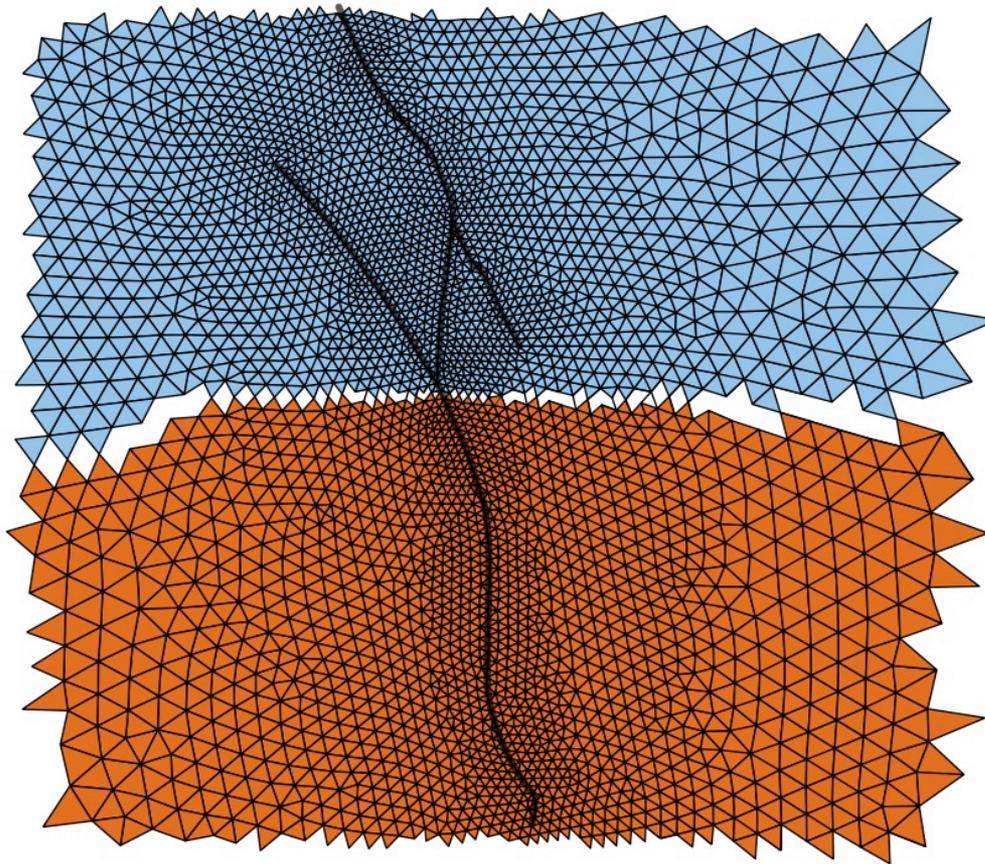
- Refactoring: Introduce control-flow aligned to hardware requirements
- Overlapping Communication / Computation
- Heterogeneous Execution: Complicated DR on fat cores
- Hybrid Parallelization
- Vectorization through custom compute kernels

Offloading



- Refactoring: Introduce control-flow aligned to hardware requirements
- Overlapping Communication / Computation
- Heterogeneous Execution: Complicated DR on fat cores
- Hybrid Parallelization
- Vectorization through custom compute kernels

Offloading



Host



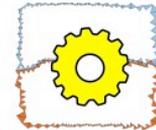
MIC



MPI comm



time MPI

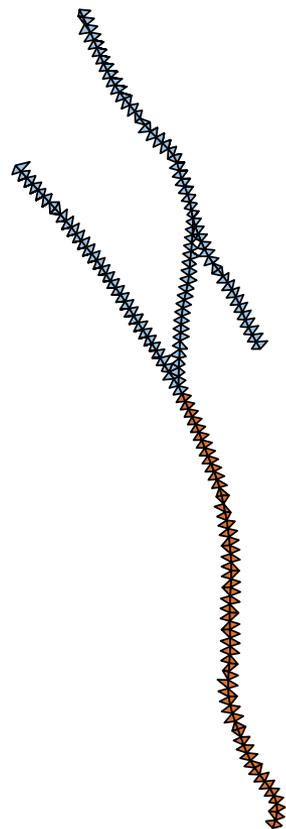


time & vol

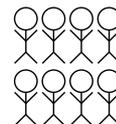


- Refactoring: Introduce control-flow aligned to hardware requirements
- Overlapping Communication / Computation
- Heterogeneous Execution: Complicated DR on fat cores
- Hybrid Parallelization
- Vectorization through custom compute kernels

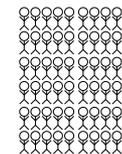
Offloading



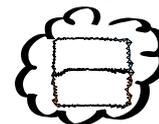
Host



MIC



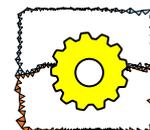
MPI comm



dyn. Rup.



time MPI

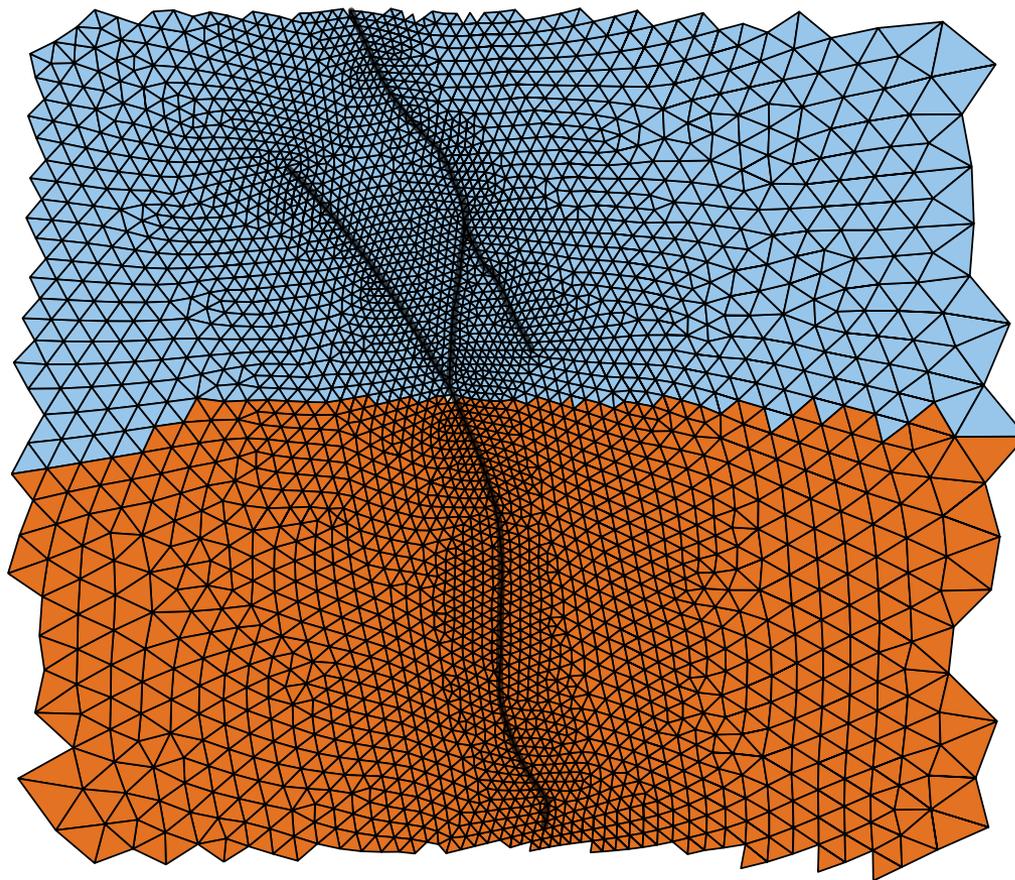


time & vol

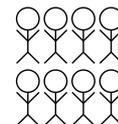


- Refactoring: Introduce control-flow aligned to hardware requirements
- Overlapping Communication / Computation
- Heterogeneous Execution: Complicated DR on fat cores
- Hybrid Parallelization
- Vectorization through custom compute kernels

Offloading



Host



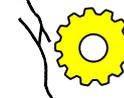
MIC



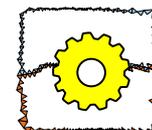
MPI comm



dyn. Rup.



time MPI



time & vol

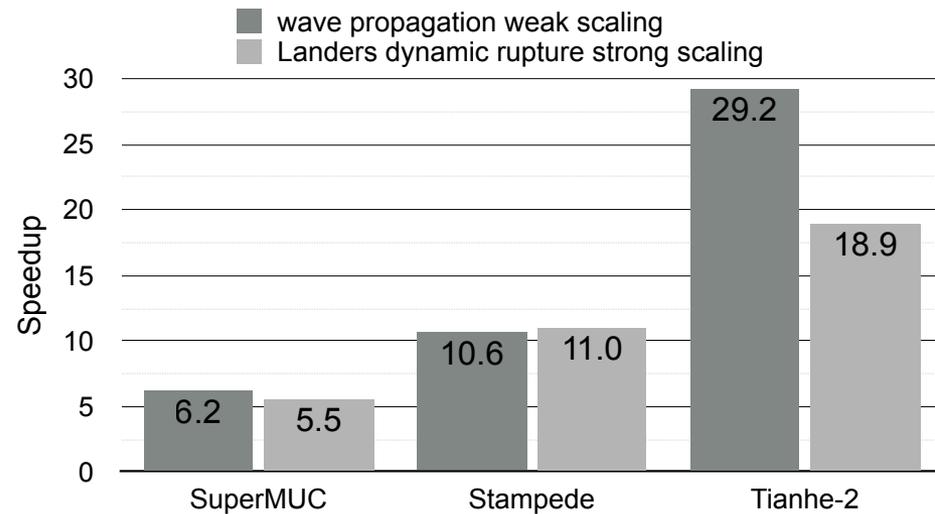
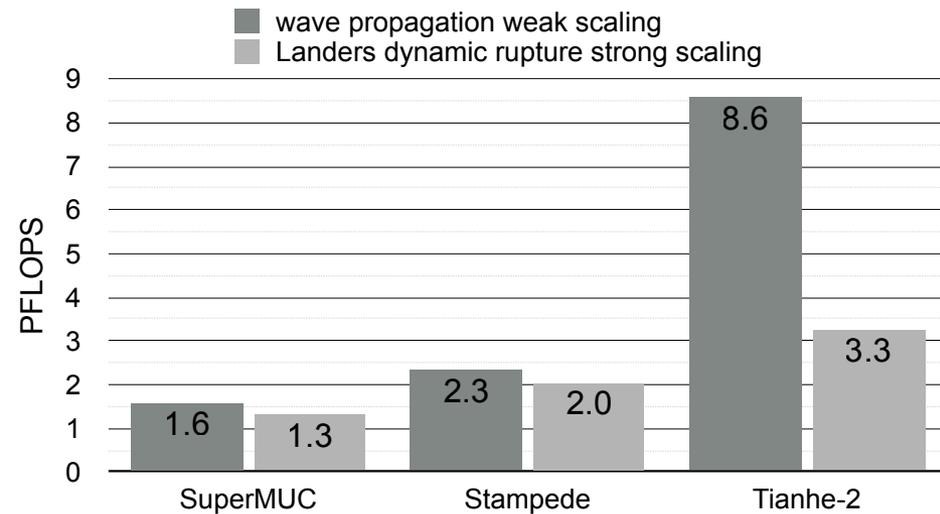


fluxes



- Refactoring: Introduce control-flow aligned to hardware requirements
- Overlapping Communication / Computation
- Heterogeneous Execution: Complicated DR on fat cores
- Hybrid Parallelization
- Vectorization through custom compute kernels

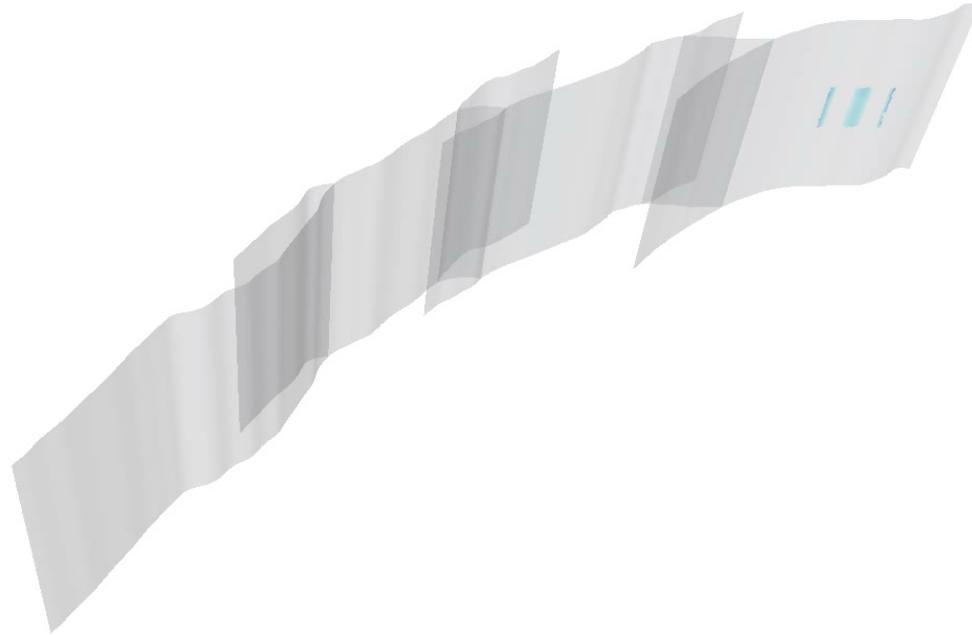
Performance



- Sustained petascale performance on SuperMUC, Stampede and Tianhe-2
- Mission accomplished: Significant reduced time solution on Xeon and Xeon Phi

Insights

- Physics!

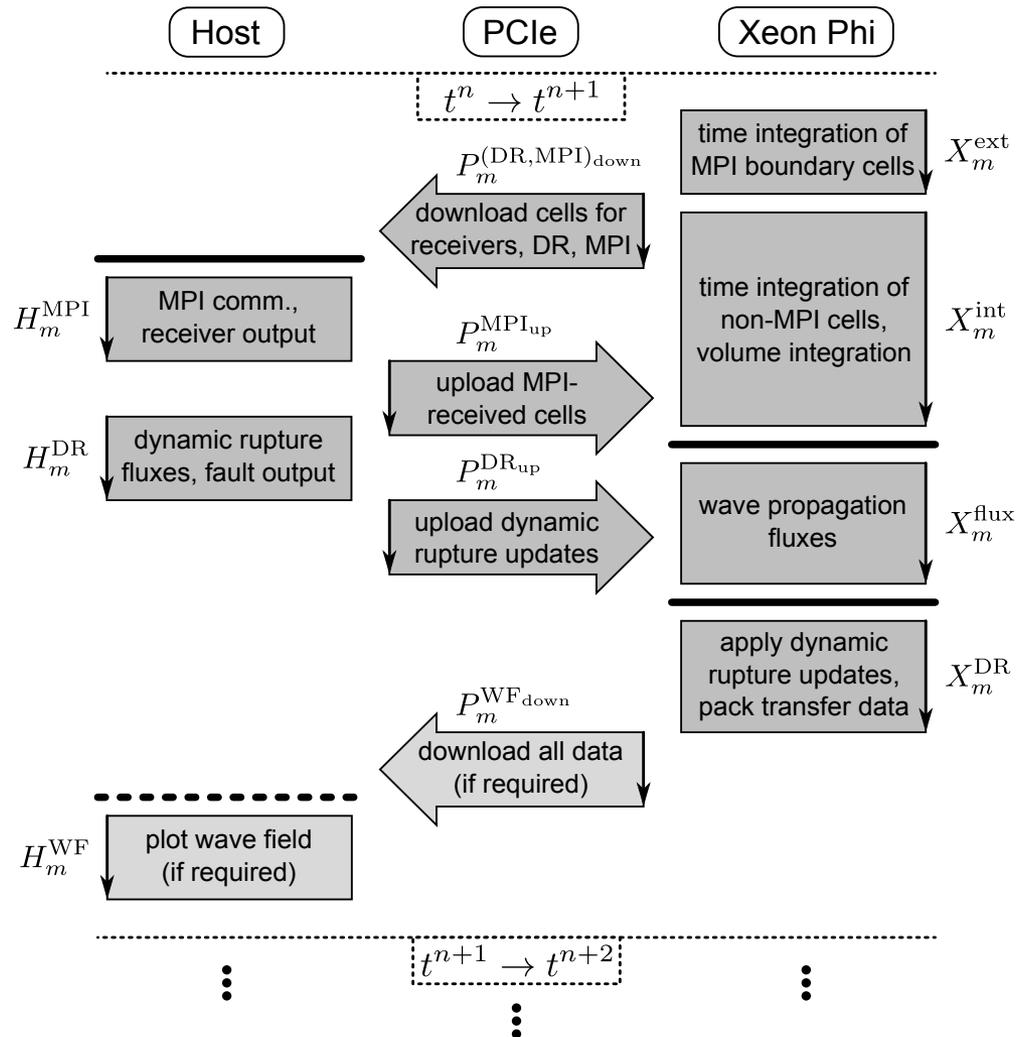


- Moving target: MIC-optimizations are general purpose optimizations
- Offload: Arbitrary interconnects at highest performance
- Performance modeling: Observation vs. Expectation
- Remaining key challenges: More physics (and all the numerics and optimizations required for that)

References

- A. Heinecke, A. Breuer, S. Rettenberger, M. Bader, A.-A. Gabriel, C. Pelties, A. Bode, W. Barth, X.-K. Liao, K. Vaidyanathan, M. Smelyanskiy and P. Dubey: Petascale High Order Dynamic Rupture Earthquake Simulations on Heterogeneous Supercomputers. In Supercomputing 2014, The International Conference for High Performance Computing, Networking, Storage and Analysis. IEEE, New Orleans, LA, USA, November 2014. Gordon Bell Finalist.
- A. Breuer, A. Heinecke, S. Rettenberger, M. Bader, A.-A. Gabriel and C. Pelties: Sustained Petascale Performance of Seismic Simulations with SeisSol on SuperMUC. In J.M. Kunkel, T. T. Ludwig and H.W. Meuer (ed.), Supercomputing — 29th International Conference, ISC 2014, Volume 8488 of Lecture Notes in Computer Science, p. 1-18. Springer, Heidelberg, June 2014. PRACE ISC Award 2014.
- A. Breuer, A. Heinecke, M. Bader and C. Pelties: Accelerating SeisSol by Generating Vectorized Code for Sparse Matrix Operators. In Parallel Computing — Accelerating Computational Science and Engineering (CSE), Volume 25 of Advances in Parallel Computing, p. 347-356. IOS Press, April 2014.

Appendix: Full-Offload Scheme



A. Heinecke, A. Breuer, S. Rettenberger, M. Bader, A.-A. Gabriel, C. Pelties, A. Bode, W. Barth, X.-K. Liao, K. Vaidyanathan, M. Smelyanskiy and P. Dubey: Petascale High Order Dynamic Rupture Earthquake Simulations on Heterogeneous Supercomputers. In Supercomputing 2014, The International Conference for High Performance Computing, Networking, Storage and Analysis. IEEE, New Orleans, LA, USA, November 2014.

Appendix: Cluster Configurations

- SuperMUC
 - Leibniz Supercomputing Centre, Munich, Germany (#12 in the Top 500 of June 2014)
 - 9,216 dual-socket Intel Xeon E5-2680 (8-core, 2.7 GHz) nodes
 - commodity InfiniBand network (fat-tree topology, FDR10 interconnect within each 512-nodes island, 4:1 bandwidth-ratio between islands)
 - Max. PFLOPS: 9216 nodes; Max. speedup: 4096 nodes vs. 4096 SuperMUC-nodes with classic version
- Stampede
 - Texas Advanced Computing Center (TACC)/Univ. of Texas, USA (#7 in the Top 500)
 - 6,400 compute nodes: Two Xeon E5-2680 processors and one Intel Xeon Phi SE10P coprocessor (61 cores, 1.1 GHz)
 - Mellanox FDR 56 Gb/s InfiniBand interconnect (2-level fat- tree topology, eight core-switches and over 320 leaf switches, 5/4 oversubscription)
 - Max. PFLOPS: 6144 nodes; Max. speedup: 4096 nodes vs. 4096 SuperMUC-nodes with classic version
- Tianhe-2
 - National Super Computer Center in Guangzhou, China (#1 in the Top 500)
 - 16,000 nodes: Two Intel Xeon E5-2692 CPUs (12-core, 2.2GHz) and three Intel Xeon Phi 31S1P coprocessors (57 cores, 1.1 GHz)
 - Interconnect, called TH Express-2, with fat-tree topology
 - Max PFLOPS: 8192 nodes (weak), 7000 nodes (strong); Max. speedup: 4096 nodes vs. 4096 SuperMUC-nodes with classic version