# Task-parallel Programming for Reactive Numerical Simulation

Christian Terboven <terboven@itc.rwth-aachen.de> & Matthias S. Müller <mueller@itc.rwth-aachen.de>

November 7th 2019, Tokyo, Japan

## The future is *dynamic*

- Dynamic Variability in HPC systems continues to increase
  - Processor features (example: Intel Turbo)
  - Energy Management (example: Power Capping)
  - Detection and Correction of Errors

## Tasking to the rescue



**Bulk Synchronous Execution (now)**

**Bulk Synchronous Execution (future)**

- Tasking is well-positioned to react to dynamic system behavior
  - Less global synchronization
  - More p2p synchronization

Image Source: John Shalf

Image Source: Jack Dongarra

# Agenda

- Intra-node: Task Affinity

- Inter-node: Task Migration and Replication

- Outlook: AI- and Simulation-Based Engineering at Exascale

- Conclusions

Task-parallel Programming for Reactive Numerical Simulation | Dr. Christian Terboven

# Intra-node: Task Affinity

## Support for task affinity is part of OpenMP 5.0 released on November 8th, 2019

```
#pragma omp task [clause…] affinity(list)
```

```
int a[N];    // N is large
...
#pragma omp task affinity(a[x-y])
{
    // task that makes use of a[x], ...
}
```

- Programmer specifies data used by task
- Recommended to execute task closely to data location
  - Do not prohibit task stealing & load balancing
- Runtime identifies the location of the data and schedules task to a close thread
- Clear separation between dependencies and affinity

# Selected implementation details



Jannis Klinkenberg, Philipp Samfass, Christian Terboven, Alejandro Duran, Michael Klemm, Xavier Teruel, Sergi Mateo, Stephen L. Olivier, and Matthias S. Müller. **Assessing Task-to-Data Affinity in the LLVM OpenMP Runtime**. Proceedings of the 14th International Workshop on OpenMP, IWOMP 2018. September 26-28, 2018, Barcelona, Spain.

## How much can this improve applications?

- Little improvements on standard 2-socket systems, more improvement on larger systems



(d) Merge sort on 8-socket

(h) Health on 8-socket

- Works well working with a lot of data & single task creator scenarios & tasks created in parallel but not all close to data
- Not much room for improvement when: parallel task creator scenarios & tasks are already created where data is located

**Intel® Xeon® E7-8860v4 (codename Broadwell)**
8 sockets, 18 cores per socket = 144 cores
2.2 GHz base frequency, 1 TB memory

# Inter-node: Task Migration and Replication

# Motivation: Real-world code example

## Dynamic variability caused by application

- Showcase application: **sam(oa)²**
  - Finite-Element and Finite-Volume simulations of dynamic adaptive meshes
  - Space Filling Curves (SFC) and Adaptive Meshes for
    Oceanic And Other Applications (Tohoku Tsunami 2011)
  - Developed at TU Munich

- Depending on situation either refinement
  or coarsening of cell / section

- Refinement leads to load imbalances
  - after each iteration
  - intra <u>and</u> inter node

# Chameleon Approach: Migratable Tasks + Self Introspection

- **Migratable task**
  - Basic unit of work without side effects
  - Action + data items (input and/or output)
  - Can be executed locally or migrated to another rank



1. Based on periodically collected introspection data detect imbalance dynamically at runtime

   **Result:** Rank 0 is significantly slower or has more work

2. Migrate tasks and data to Rank 1

3. Prioritized execution of migrated tasks at Rank 1 + send back results or outputs

➢ Desired: Migrate as soon as possible to overlap communication and computation

**Tasked-based Execution Environment**

- ❑ Create, queue and execute migratable tasks
- ❑ Allows early task migration for load balancing between ranks/nodes

**Self Introspection**

- ❑ Continuous monitoring of the current rank
- ❑ Determine runtime conditions, load or performance metrics

**Consolidation and Analysis**

- ❑ Consolidates information from all ranks
- ❑ Decision making
  - o Migrate tasks?
  - o Victim selection

# Results Experiments – SW-induced Imbalances with sam(oa)$^2$



Figure 3: Load imbalances between ranks per time step in sam(oa)$^2$ for an application run with 32 nodes/ranks



Figure 4: Strong scaling experiments with Tohoku tsunami in 2011 for complete application. Relative speedup to single node base line

- Simulated 60 minutes of Tohoku tsunami in 2011

➢ Reduce degree of imbalance

# Outlook: AI- and Simulation-Based Engineering at Exascale

# Challenges at Exascale

## Tasking may be employed to provide efficient and scalable coupling of SW components

- CFD simulations cannot live without modeling approaches
  - Becomes worse in multi-physics and multi-scale phenomena, or with interactions such as combustion
  - Will be complemented with data-based models

- At Exascale, the amount of data may exceed the Exabyte range for single simulation runs
  - In-situ data reduction, extraction and interpretation will hence be unavoidable

- To utilize HPC resources efficiently, software and workflows must scale to high CPU counts
  - In compute-drive applications, analyses are frequently a posteriori, necessitating to have the data on disk
  - As the field of parallel and scalable ML and DL is progressing, those algorithms become feasible to be intertwined with simulation codes implementing full loops

- FZJ's Modular Supercomputing as a prominent heterogeneous pre-Exascale architecture

# Challenges at Exascale

## Tasking may be employed to provide efficient and scalable coupling of SW components

- Key expectation: As the field of parallel and scalable ML and DL is progressing, those algorithms become feasible to be intertwined with simulation codes implementing full loops

# Conclusions

# Reactive Task-parallel Programming

## Tasking model is becoming more attractive

- Tasking brings advantages for dynamic systems

- Affinity brings performance improvements
  - Including support for complex memory hierarchies

- Reactive MPI+OpenMP task migration for fine-granular load balancing
  - Robustness against HW- and work-induced imbalances

- Key expectation: As the field of parallel and scalable ML and DL is progressing, those algorithms become feasible to be intertwined with simulation codes implementing full loops

## Invitation to collaborate

- Future research direction: runtime work for intra-node and inter-node tasking
- Exchange with RIKEN expected to continue
- Also: see proposals from 2018 meeting

i12    High Performance Computing    RWTH AACHEN UNIVERSITY

# Vielen Dank
# für Ihre Aufmerksamkeit