# The Rocky Road To Tasking

March 21, 2019 | Ivo Kabadshow, Laura Morgenstern | Jülich Supercomputing Centre

JÜLICH
Forschungszentrum

# HPC ≠ HPC

# HPC ≠ HPC

# HPC ≠ HPC

# HPC ≠ HPC

# HPC ≠ HPC



High Frequency Trading

MD

CPU Cycle   Network Latency   Game Dev   Astrophysics   Deep Learning

ns   $\mu s$   ms   s   min   h

Critical walltime

## Requirements for MD

- Strong scalability
- Performance portability

JÜLICH
Forschungszentrum

# Our Motivation

**Solving Coulomb problem for Molecular Dynamics**

## Task: Compute all pairwise interactions of *N* particles

N-body problem: $\mathcal{O}(N^2) \rightarrow \mathcal{O}(N)$ with FMM

## Why is that an issue?

- MD targets $< 1ms$ runtime per time step
- MD runs millions or billions of time steps
- not compute-bound, but synchronization bound
- no libraries (like BLAS) to do the heavy lifting

We might have to look under the hood ... and get our hands dirty.

JÜLICH
Forschungszentrum

# Parallelization Potential



## Classical Approach

- Lots of independent parallelism

# Parallelization Potential

# Coarse-Grained Parallelization

JÜLICH
Forschungszentrum

# Coarse-Grained Parallelization



- Different amount of available loop-level parallelism within each phase
- Some phases contain sub-dependencies
- Synchronizations might be problematic

JÜLICH
Forschungszentrum

# FMM Algorithmic Flow

**Multipole to multipole (M2M), shifting multipoles upwards**

# FMM Algorithmic Flow

**Multipole to multipole (M2M), shifting multipoles upwards**

# FMM Algorithmic Flow

**Multipole to local (M2L), translate remote multipoles into local taylor moments**

# FMM Algorithmic Flow

**Multipole to local (M2L), translate remote multipoles into local taylor moments**

# FMM Algorithmic Flow

**Local to local (L2L), shifting Taylor moments downwards**

# FMM Algorithmic Flow

**Local to local (L2L), shifting Taylor moments downwards**

# CPU Tasking Framework



Queue

Scheduler

ThreadingWrapper
Thread

Core

JÜLICH
Forschungszentrum

# CPU Tasking Framework



Queue

Scheduler

ThreadingWrapper
Thread

Core

JÜLICH
Forschungszentrum

# CPU Tasking Framework

# CPU Tasking Framework



| | | | | | |
|---|---|---|---|---|---|
| Queue | | | | | Dispatcher |
| Scheduler | | | | | TaskFactory |
| ThreadingWrapper | | | | | LoadBalancer |
| Thread | | | | | |
| Core | | | ... | | |

JÜLICH
Forschungszentrum

# CPU Tasking Framework

**Task life-cycle per thread**

Queues

# CPU Tasking Framework

**Task life-cycle per thread**

task

⚙ Task execution

Queues

JÜLICH
Forschungszentrum

# CPU Tasking Framework

**Task life-cycle per thread**

JÜLICH
Forschungszentrum

# CPU Tasking Framework

**Task life-cycle per thread**

# CPU Tasking Framework

**Task life-cycle per thread**



- Tasks can be prioritized by task type
- Only ready-to-execute tasks are stored in queue
- Worksteating from other threads is possible

# Tasking Without Workstealing

103 680 **Particles on 2×Intel Xeon E5-2680 v3 (2×12 cores)**

JÜLICH
Forschungszentrum

# Tasking With Worksteaming

103 680 **Particles on 2×Intel Xeon E5-2680 v3 (2×12 cores)**

# The Rocky Road To Tasking

March 21, 2019 | Ivo Kabadshow, Laura Morgenstern | Jülich Supercomputing Centre

JÜLICH
Forschungszentrum

# GPU Tasking

**Goal**

- Provide same features as CPU tasking:
    - Static and dynamic load balancing
    - Priority queues
    - Ready-to-execute tasks

JÜLICH
Forschungszentrum

# GPU Tasking

**Uniform Programming Model for CPUs and GPUs**

Thread

issues scalar instruction

FPU

**Persistent** Thread

issues scalar instruction

FPU

JÜLICH
Forschungszentrum

# GPU Tasking

**Uniform Programming Model for CPUs and GPUs**

Thread

issues vector instruction

→

FPU

**Persistent** Warp

issues 32 scalar instructions

→

Processing Block

Warp Scheduler

| FPU | FPU |
| FPU | FPU |

JÜLICH
Forschungszentrum

# GPU Tasking

## Uniform Programming Model for CPUs and GPUs

SMT-Threads



run on →

Core

| L2 Cache |
| L1 Cache |
| FPU | FPU |
| FPU | FPU |

**Persistent** Thread Block



runs on →

Streaming Multiprocessor

| Shared Memory | | | |
| FPU | FPU | FPU | FPU |
| FPU | FPU | FPU | FPU |
| FPU | FPU | FPU | FPU |
| FPU | FPU | FPU | FPU |

JÜLICH
Forschungszentrum

# GPU Tasking

## Uniform Programming Model for CPUs and GPUs

# GPU Tasking

**Uniform Programming Model for CPUs and GPUs**

# Pitfalls

**Performance Portability**

Diverse GPU programming approaches:

- OpenCL
- CUDA
- SYCL

Our requirements:

- Strong subset of C++11
- Portability between GPU vendors
- Tasking features
- Maturity

## (Intermediate) Solution

Use CUDA for reasons of performance, specific tasking features and maturity. Take the loss of not being portable out of the box.

**JÜLICH**
Forschungszentrum

# Pitfalls

**Performance Portability**

For performance portability we consider diverse GPU programming approaches:

- OpenCL
- CUDA
- SYCL

## Unsatisfying (Intermediate) Solution

Use CUDA for reasons of performance and specific features. Take the loss of not being portable out of the box.

JÜLICH
Forschungszentrum

# Pitfalls

**Architectural Differences**

## Pitfalls for Load Balancing

- No thread pinning
- No cache coherency

## Pitfalls for Mutual Exclusion

- Weak memory consistency
- Missing forward progress guarantees

JÜLICH
Forschungszentrum

# Pitfalls

**Load Balancing**

- No possibility to pin threads to streaming multiprocessors
- No direct access to shared memory of other streaming multiprocessors
- Work stealing requires multi-producer multi-consumer queues → Mechanism for mutual exclusion?

JÜLICH
Forschungszentrum

# Pitfalls

**Mutual Exclusion**

- Weak memory consistency
- Warp-synchronous deadlocks due to lock step
- How to prove thread safety?

**JÜLICH**
Forschungszentrum

# Pitfalls

## Mutex Implementation

```
class Mutex
{
    __inline__ __device__ void lock()
    {
        while (atomicCAS(&mutex, 0, 1) != 0)
        __threadfence();
    };
    __inline__ __device__ void unlock()
    {
        __threadfence();
        atomicExch(&mutex, 0);
    };

    int mutex = 0;
};
```

JÜLICH
Forschungszentrum

# Very First Evaluation

**Conditions**

- Tasking with global queue only
- Measurements without work load to determine enqueue and dequeue overhead
- Measurements on P100 with 56 thread blocks with 1024 threads each
- Measurements on V100 with 80 thread blocks with 1024 threads each

JÜLICH
Forschungszentrum

# First Evaluation

**Tasking Overhead on P100 and V100**

JÜLICH
Forschungszentrum

# GPU Tasking

**Conclusion**

- Fine-grained task parallelism pays off on CPUs
- Developed mapping between CPU and GPU concepts
- (Partly) overcome pitfalls:
    - Lock-based mutual exclusion
    - Reusability of CPU tasking code
    - Architectural differences between CPU and GPU
- Successfully transferred parts of CPU tasking to GPUs

JÜLICH
Forschungszentrum

# Next Steps

- Analyze and solve performance issues in dependency resolution
- Use memory pool for dynamic allocations
- Implement hierarchical queues
- Transfer priority queue to GPU
- Exploit data-parallelism through warps
- Consider the use of lock-free data structures
- Implement FMM based on GPU tasking

JÜLICH
Forschungszentrum

# Thank You to Our Sponsor!

NVIDIA Tesla V100 and NVIDIA Tesla P100 where provided by

JÜLICH
Forschungszentrum

# The Rocky Road To Tasking

March 21, 2019 | Ivo Kabadshow, Laura Morgenstern | Jülich Supercomputing Centre

JÜLICH
Forschungszentrum