

GPU Optimizations of Material Point Method and Collision Detection

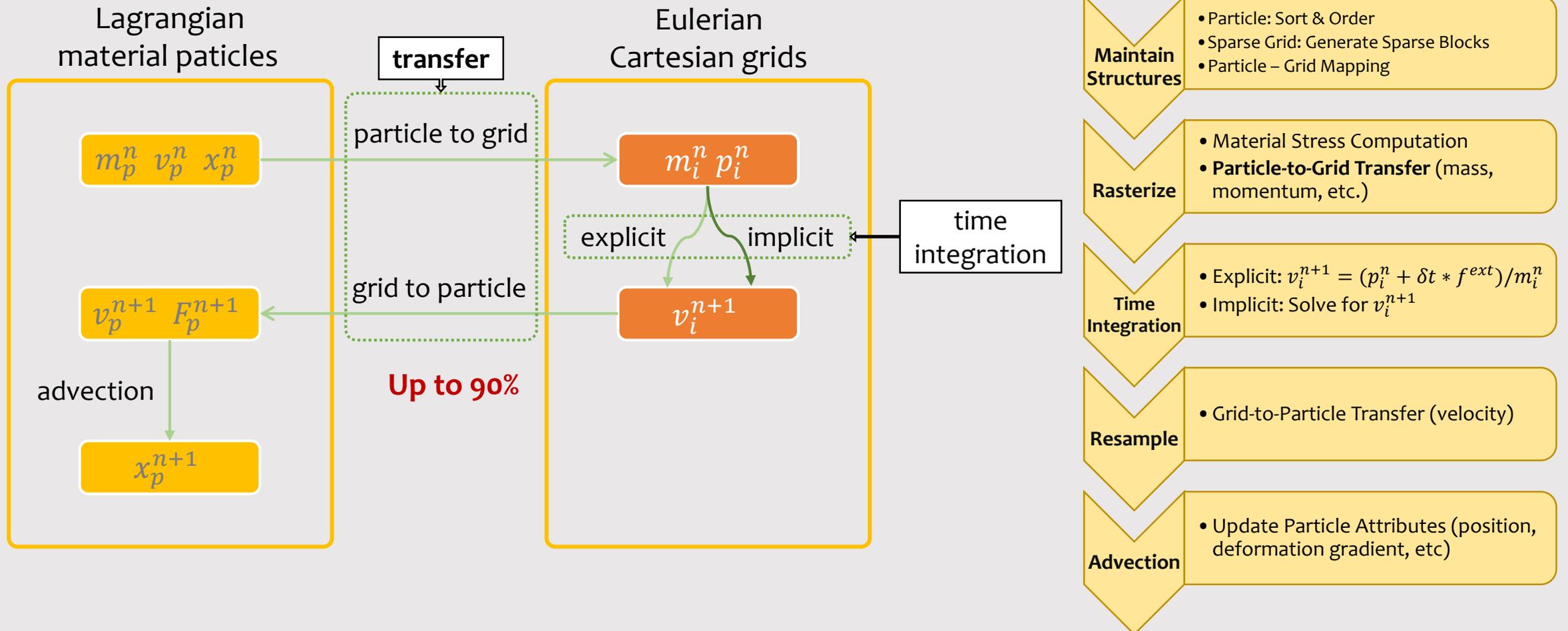
Xinlei Wang, 王鑫磊

浙江大学

Material Point Method

- Fluid
 - Smoothed-Particle Hydrodynamics
 - Grid-based Methods
- Solid
 - Finite Element Method
 - Finite Difference Method
- Material Point Method
 - large deformation, complex topology changes
 - multi-material & multiphase coupling
 - (self) collision handling

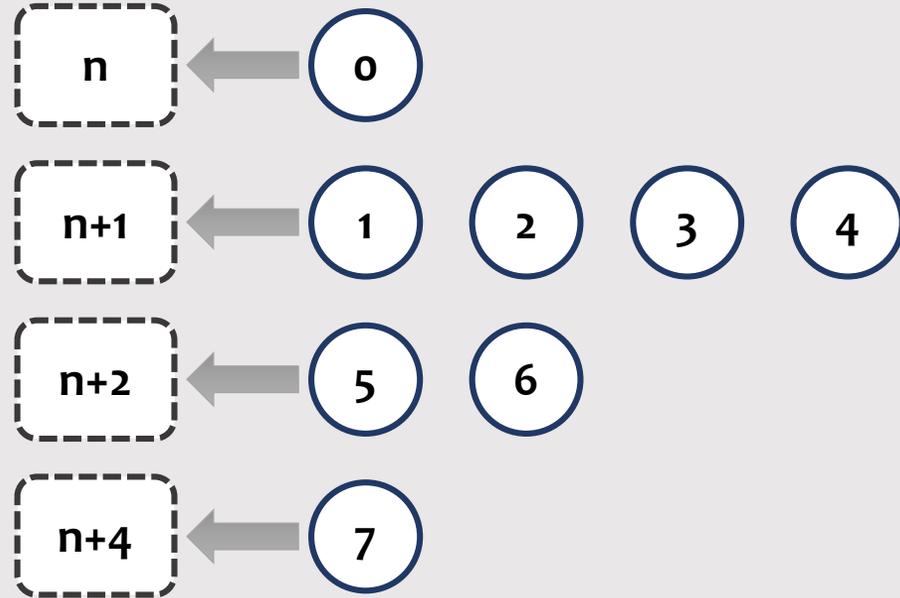
MPM Pipeline Overview



Performance is the Solution

- “dx gap”
 - a gap between adjacent models when colliding
 - increase grid resolution => more particles to achieve equal magnitude
- CFL Condition
 - for simulation stability and collision handling
 - more time steps per frame => more work to compute a frame
- Performance is the key!

Gather (node based)



transfer

notation

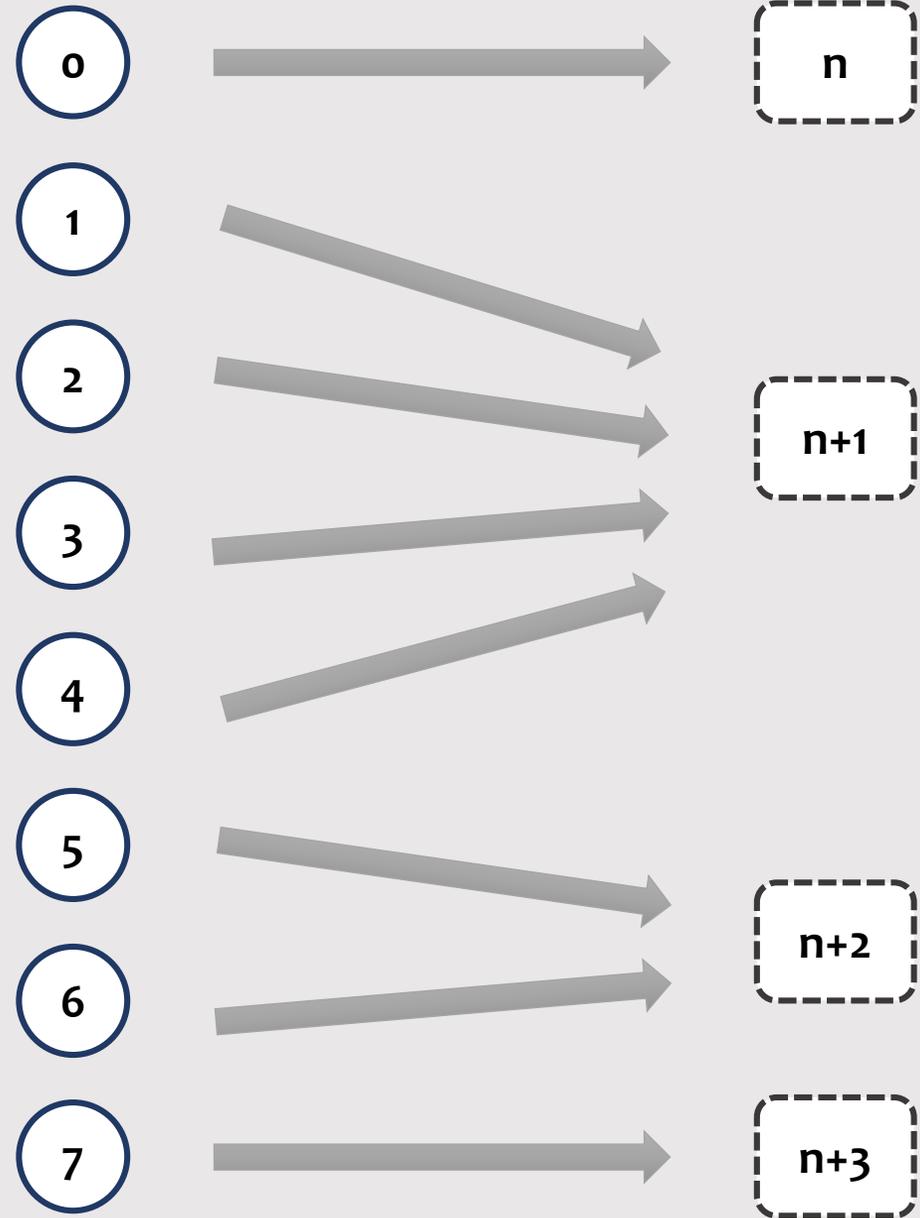


grid node



particle

Scatter (particle based)



Hardware Friendly Solutions

- **MLS MPM**

- [2018 SIGGRAPH, Hu, et al.] A Moving Least Squares Material Point Method with Displacement Discontinuity and Two-Way Rigid Body Coupling

- **Async MPM**

- [2018 SCA, Fang, et al.] A Temporally Adaptive Material Point Method with Regional Time Stepping

- **GVDB**

- [2018 EG, Wu, et al.] Fast Fluid Simulations with Sparse Volumes on the GPU

- **Warp for Cell**

- [2017 GTC, Museth, et al.] Blasting Sand with NVIDIA CUDA: MPM Sand Simulation for VFX
- <http://on-demand.gputechconf.com/gtc/2017/video/s7298-ken-museth-blasting-sand-with-nvidia-cuda-mpm-sand-simulation-for-vfx.mp4>

- **Bottleneck: Particle-to-Grid Transfer**

The Alternative of Transfer

warp intrinsics

lane id	0	1	2	3	4	5	6	7
node id	n	n+1	n+1	n+1	n+1	n+2	n+2	n+3
boundary mark	1	1	0	0	0	1	0	1
region interval	0	3	2	1	0	1	0	0

ballot
clz

region 0 region 1 region 2 region 3

attribute mass	m_0	m_1	m_2	m_3	m_4	m_5	m_6	m_7
----------------	-------	-------	-------	-------	-------	-------	-------	-------

iteration 0, stride 1

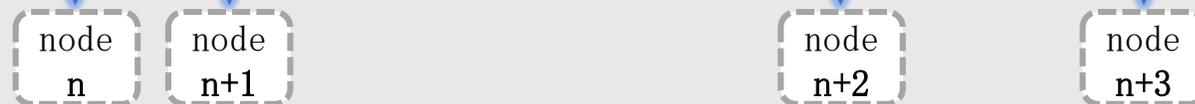
shfl

mass sum	m_0	$m_1 + m_2$	$m_2 + m_3$	$m_3 + m_4$	m_4	$m_5 + m_6$	m_6	m_7
----------	-------	-------------	-------------	-------------	-------	-------------	-------	-------

iteration 1, stride 2

mass sum	m_0	$(m_1 + m_2) + (m_3 + m_4)$	$(m_2 + m_3) + m_4$	$m_3 + m_4$	m_4	$m_5 + m_6$	m_6	m_7
----------	-------	-----------------------------	---------------------	-------------	-------	-------------	-------	-------

shared memory



Comparison

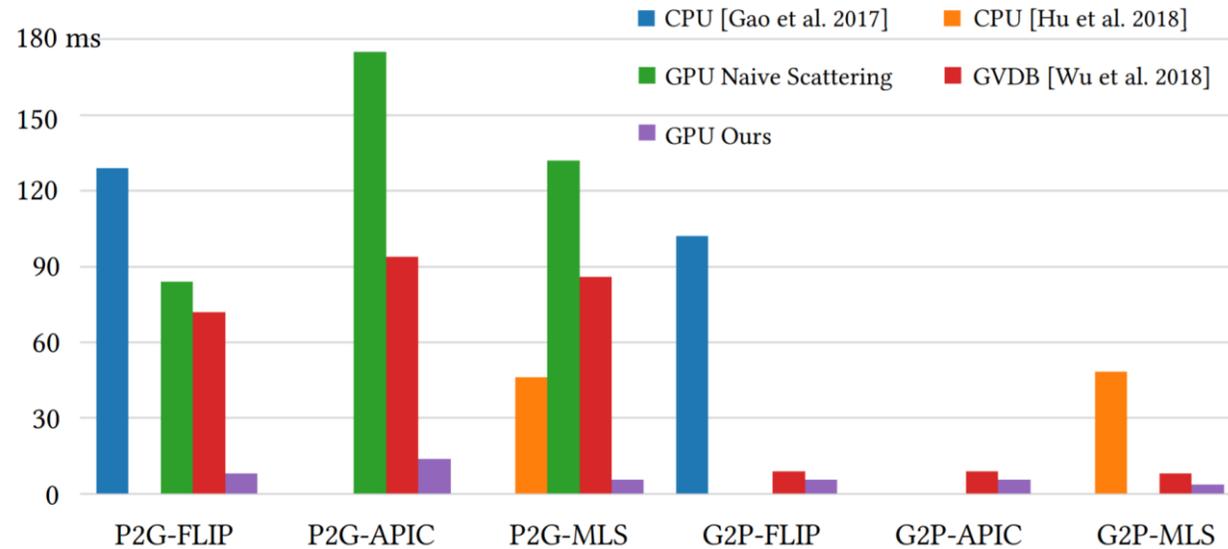
Optimized Scatter

- No auxiliary structures or memory
- Uniform workload for each thread
- Very few 'atomicAdd' write conflicts

Gather

- Additional particle list for each grid node
- Divergent workload
- No write-conflicts at all

CPU: 18-core Intel Xeon Gold 6140, ¥16000
GPU: Nvidia Titan XP, ¥8000



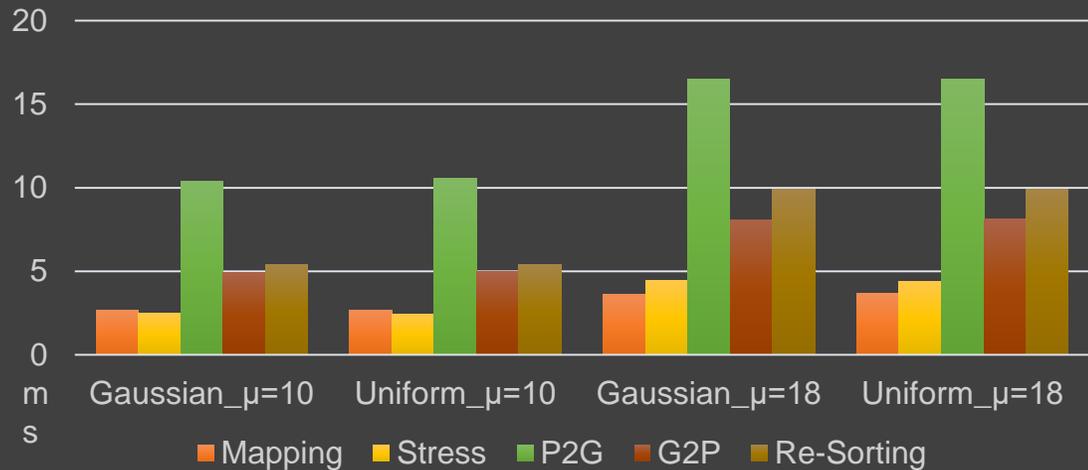
- vs. FLIP [Gao et al. 2017]
 - CPU-based, Gather-style
 - ~16X Speed-up
- vs. MLS [Hu et al. 2018]
 - CPU-based, Scatter-style
 - ~8X Speed-up
- vs. Naïve Scatter
 - GPU-based, Scatter-style
 - ~10~24X Speed-up
- vs. GVDB [Wu et al. 2018]
 - GPU-based, Gather-style
 - ~7~15X Speed-up

Performance Benchmarks

Fundamental Implementation Choices

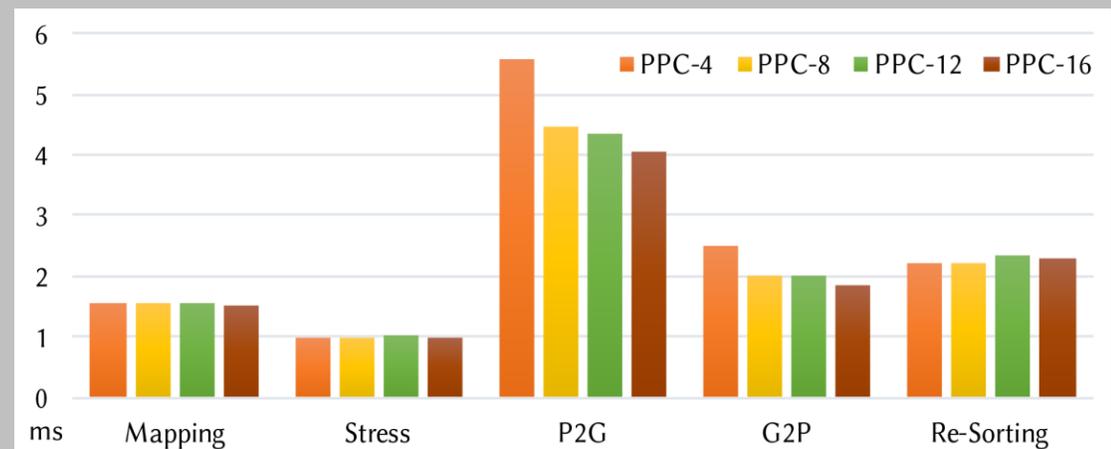
- Data Structure for Particles
 - Arrays in the SoA (Structure of Array) layout
- Data Structure for Space
 - Perceptually a sparse uniform grid
 - Support efficient interpolation operations
 - GSPGrid vs. GVDB
- Sort
 - Radix sort vs. Histogram sort

Performance Factors

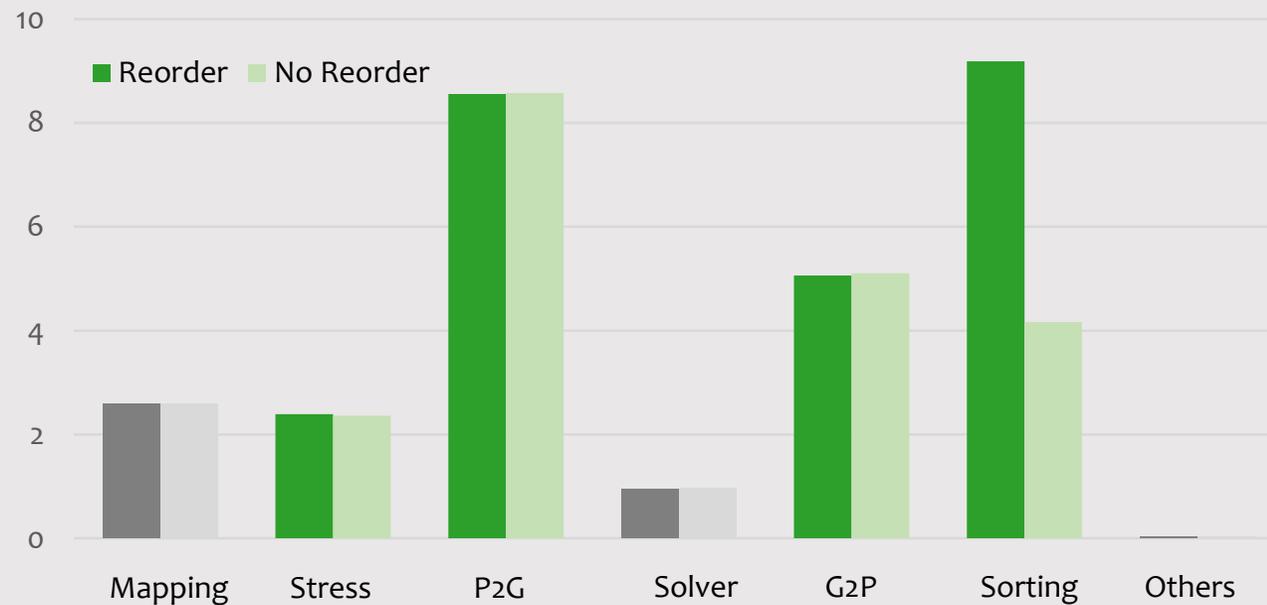


- Particle distribution doesn't matter much
- The number of particles matters

- When the number of particles is fixed,
- ppc ↑, node ↓, performance ↑



Delayed Ordering Speedup



Delayed Ordering

- Particle Attributes Classification
- By Perception
 - Intrinsic: Mass, Physical Property (Constitutive Model, etc.)
 - Extrinsic: Position, Velocity, Deformation Gradient, Affine Velocity Field (or Velocity Gradient)
- By Access (Write/ Read) Frequency
 - Mass: remains static after initialized, read once per timestep
 - Position: maintained after each timestep,
 - Everything else (Velocity, Deformation Gradient, Affine Velocity Field , etc.)

Ordering Strategy

particle index

3 4 1 2 6 0 7 5

3 1 5 4 0 2 7 6

7 1 6 4 5 2 0 3

step n-1

step n

step n+1

particle attribute

m_0^n	m_1^n	m_2^n	m_3^n	m_4^n	m_5^n	m_6^n	m_7^n
x_3^n	x_4^n	x_1^n	x_2^n	x_6^n	x_0^n	x_7^n	x_5^n

m_0^n	m_1^n	m_2^n	m_3^n	m_4^n	m_5^n	m_6^n	m_7^n
v_3^n	v_4^n	v_1^n	v_2^n	v_6^n	v_0^n	v_7^n	v_5^n
x_3^n	x_1^n	x_5^n	x_4^n	x_0^n	x_2^n	x_7^n	x_6^n

m_0^n	m_1^n	m_2^n	m_3^n	m_4^n	m_5^n	m_6^n	m_7^n
v_3^n	v_1^n	v_5^n	v_4^n	v_0^n	v_2^n	v_7^n	v_6^n
x_7^n	x_1^n	x_6^n	x_4^n	x_5^n	x_2^n	x_0^n	x_3^n

Ordering Strategy

Access times per-particle per-timestep

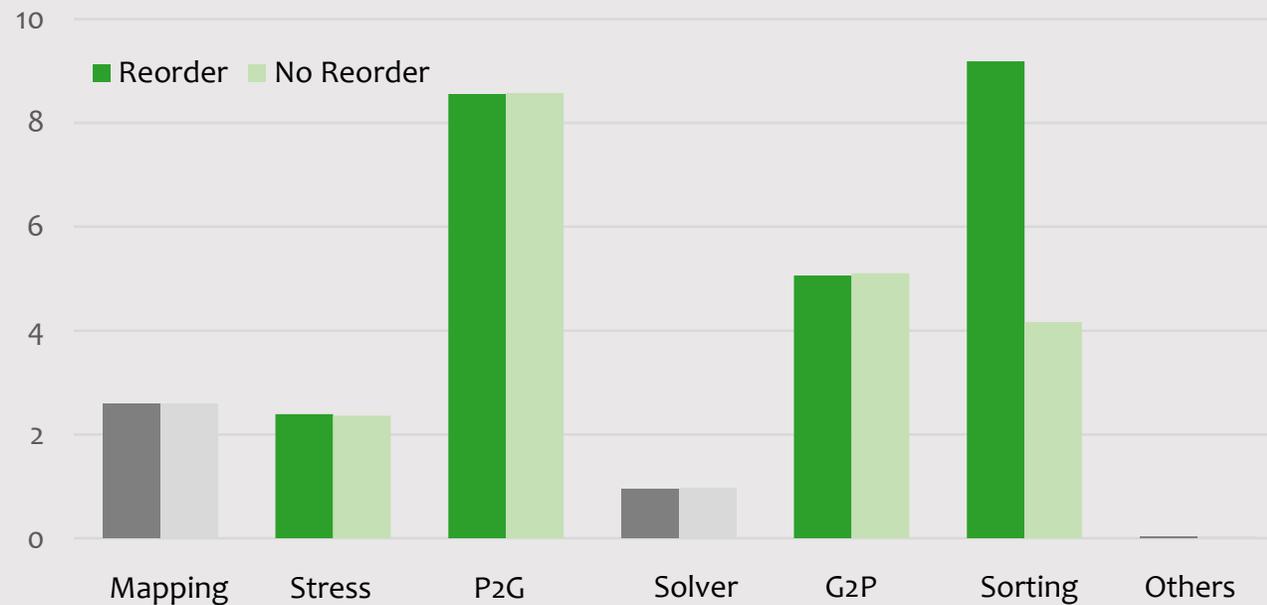
Reorder Everything

Particle Attribute (Dimension)	Read		Write	
	arbitrary	contiguous	arbitrary	contiguous
mass (1)	1	1	0	1
position (d)	1	3	0	1+1
velocity (d)	1	1	0	1+1
deformation gradient (d*d)	1	1	0	1+1
...	

Delayed Ordering

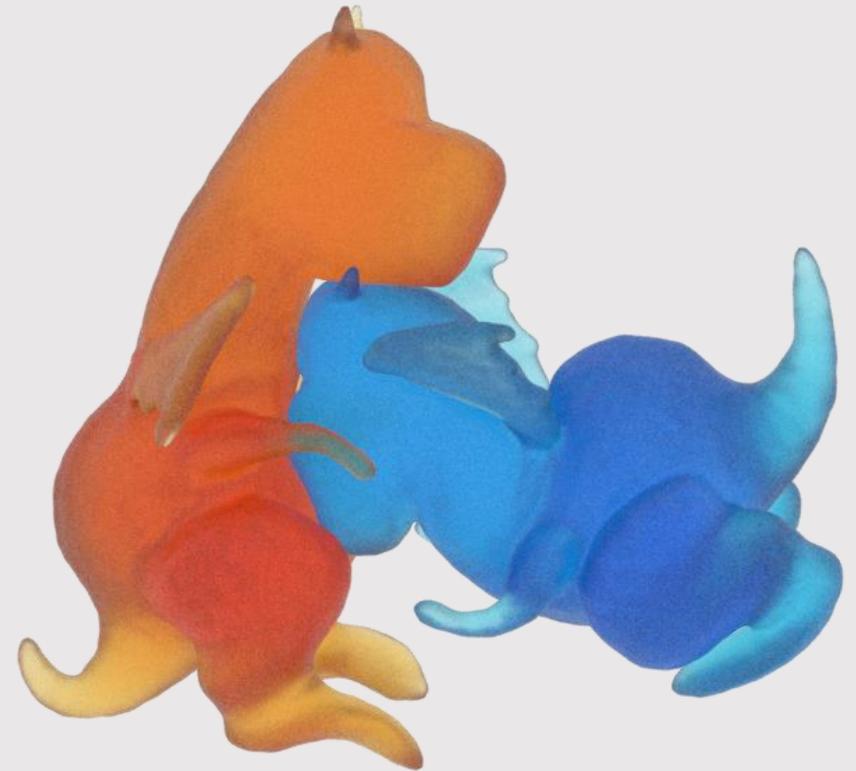
Particle Attribute (Dimension)	Read		Write	
	arbitrary	contiguous	arbitrary	contiguous
mass (1)	1	0	0	0
position (d)	1	3	0	1+1
velocity (d)	1	0	0	1
deformation gradient (d*d)	0	1	0	1
...	

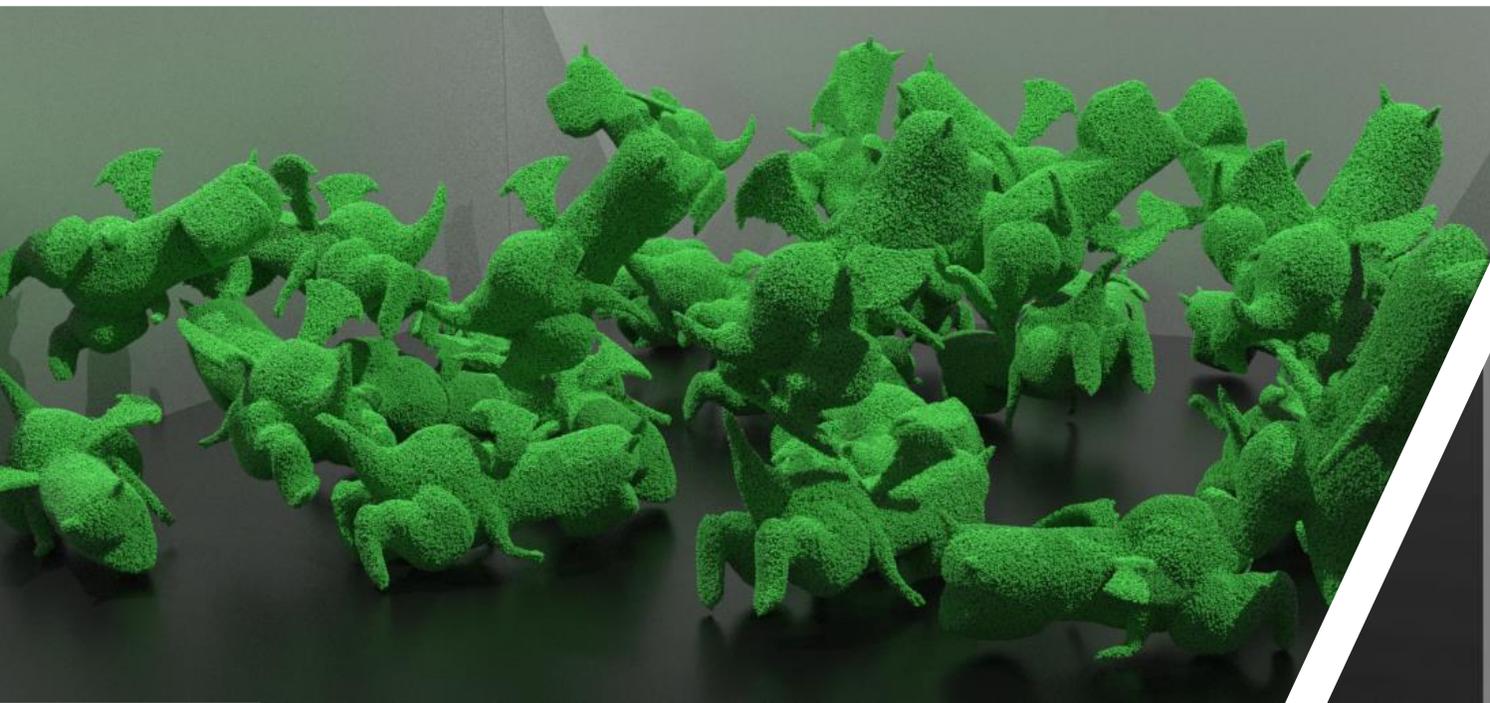
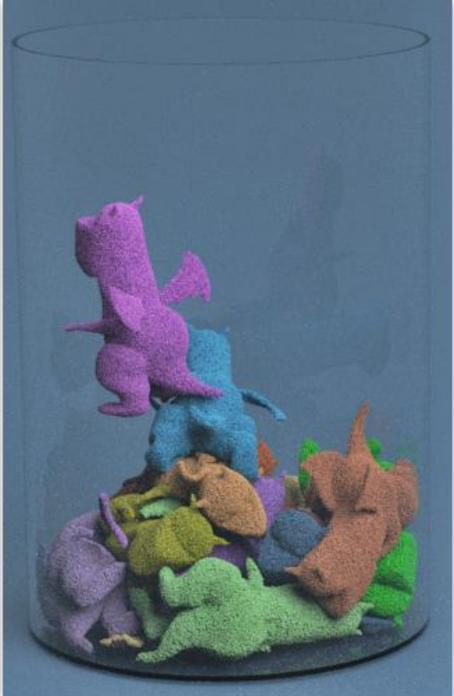
Delayed Ordering Speedup



Summary:

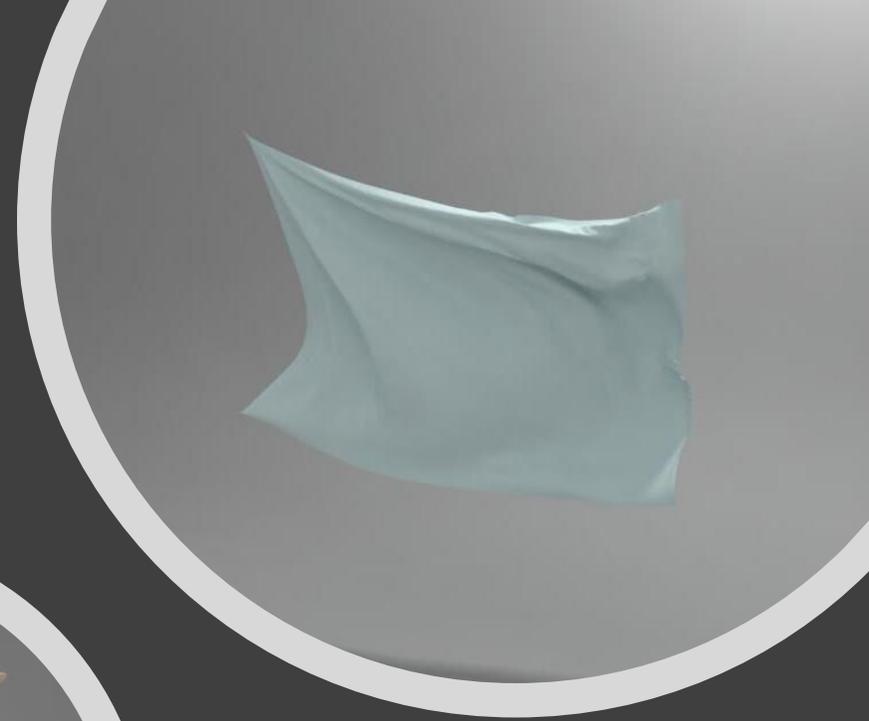
- GPU MPM pipeline
 - efficient, extensible, cross-platform
 - support multiple-materials
 - <https://github.com/kuiwuchn/GPUMPM>
- What's next?
 - Multi-GPU MPM
 - Distributed GMPM





Collision Detection

- Broad-phase Collision Detection
- Look for AABB bounding box intersections
- Typical memory-bound CUDA kernels!



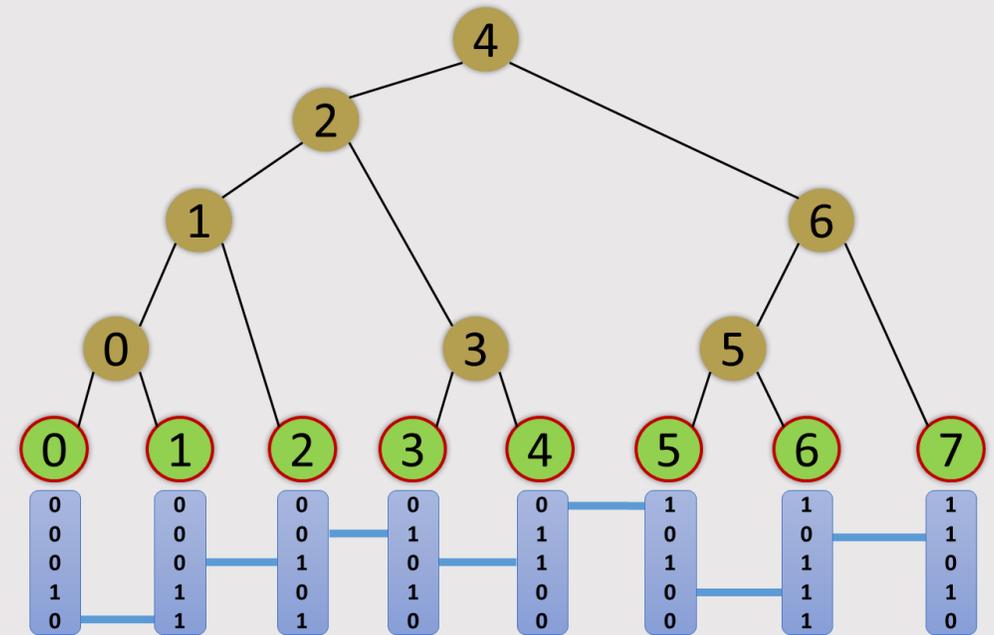
BVH (Bounding Volume Hierarchy) Construction

- **BVH Construction**

- [2012 Karras] builds all nodes in parallel
- [2014 Apetrei] builds & refits in one iteration

- BVH Stackless Traversal

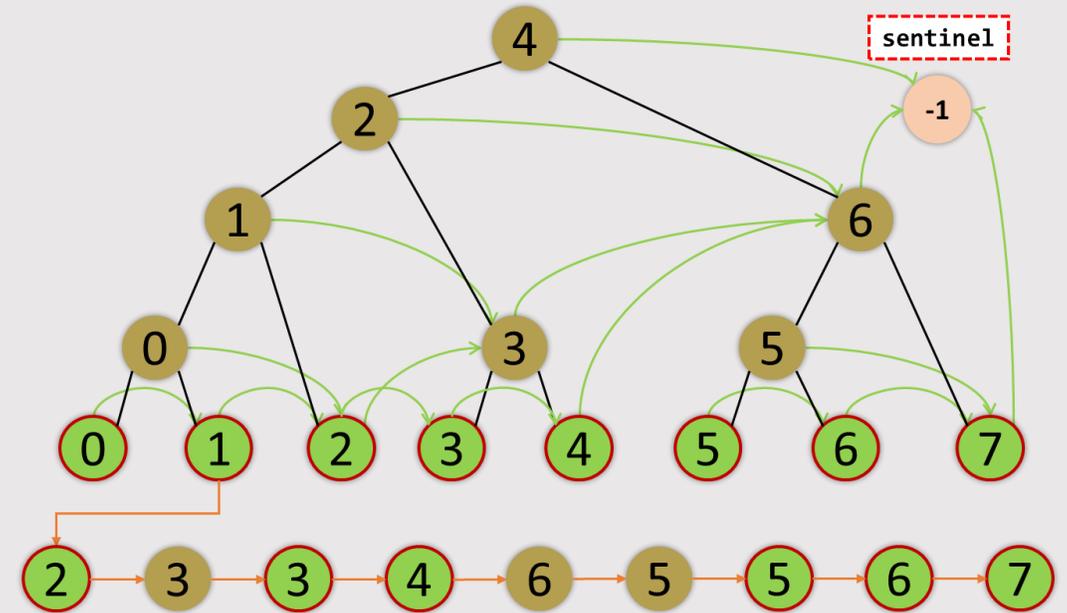
- [2007 Damkjaer] depth-first order traversal using escape index



Linear BVH built on top of primitives sorted by their Morton codes

Stackless BVH Traversal

- BVH Construction
 - [2012 Karras] builds all nodes in parallel
 - [2014 Apetrei] builds & refits in one iteration
- **BVH Stackless Traversal**
 - [2007 Damkjaer] depth-first order traversal using escape index



Depth-first order traversal track
of *Primitive-1* assuming it collides
with all the other primitives

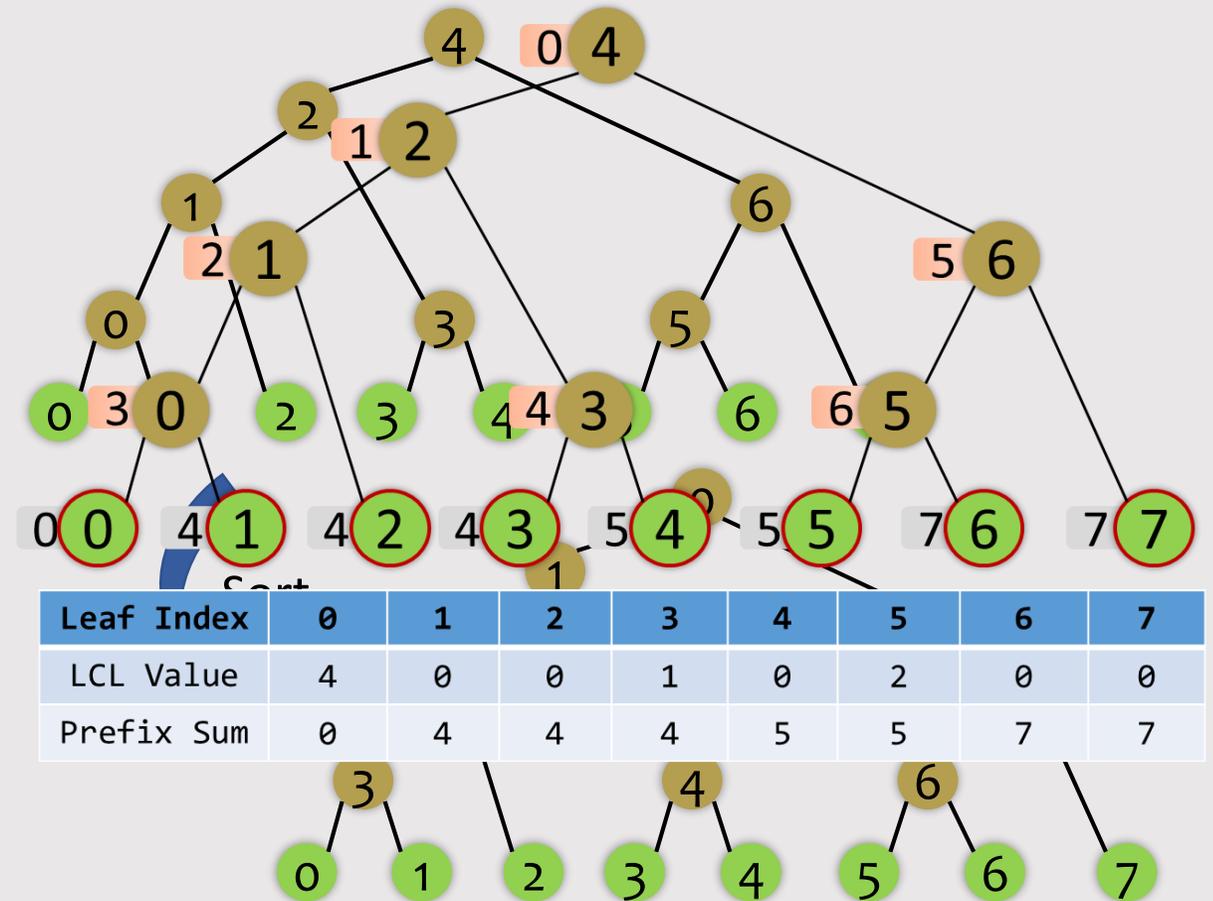
BVH-based Collision Detection

- Full traversal of the internal nodes

- Original BVH 4 2 1 0 3 6 5
- Ordered BVH 0 1 2 3 4 5 6

- How to compute BVH order

- Calculate the LCL-value of each leaf node
- Compute prefix sums of LCL-values
- Assign the indices from LCA from top to bottom



Effectiveness of ordering

- **Without ordering**

- L2 Cache Hit Rate (L1 Reads)
 - 88%
- Global Load L2 Transactions/Access
 - 31.7
- Maximum Divergence
 - 99.9%

2~3x speedup !

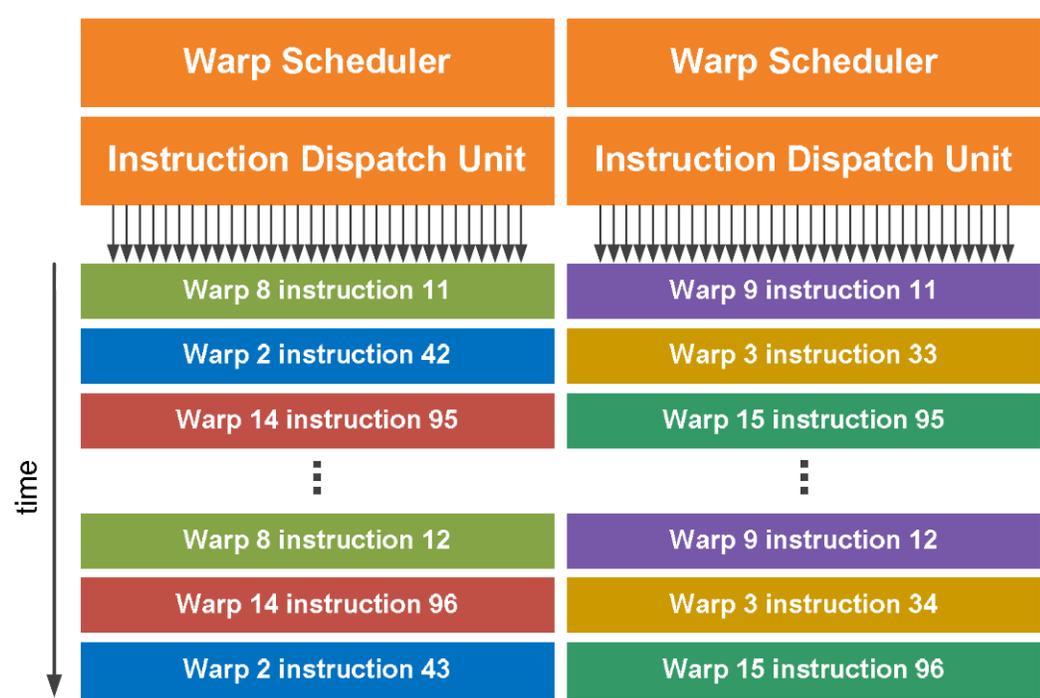
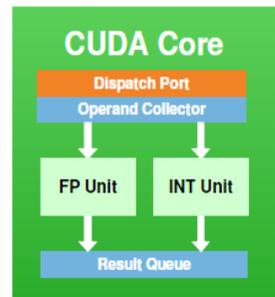
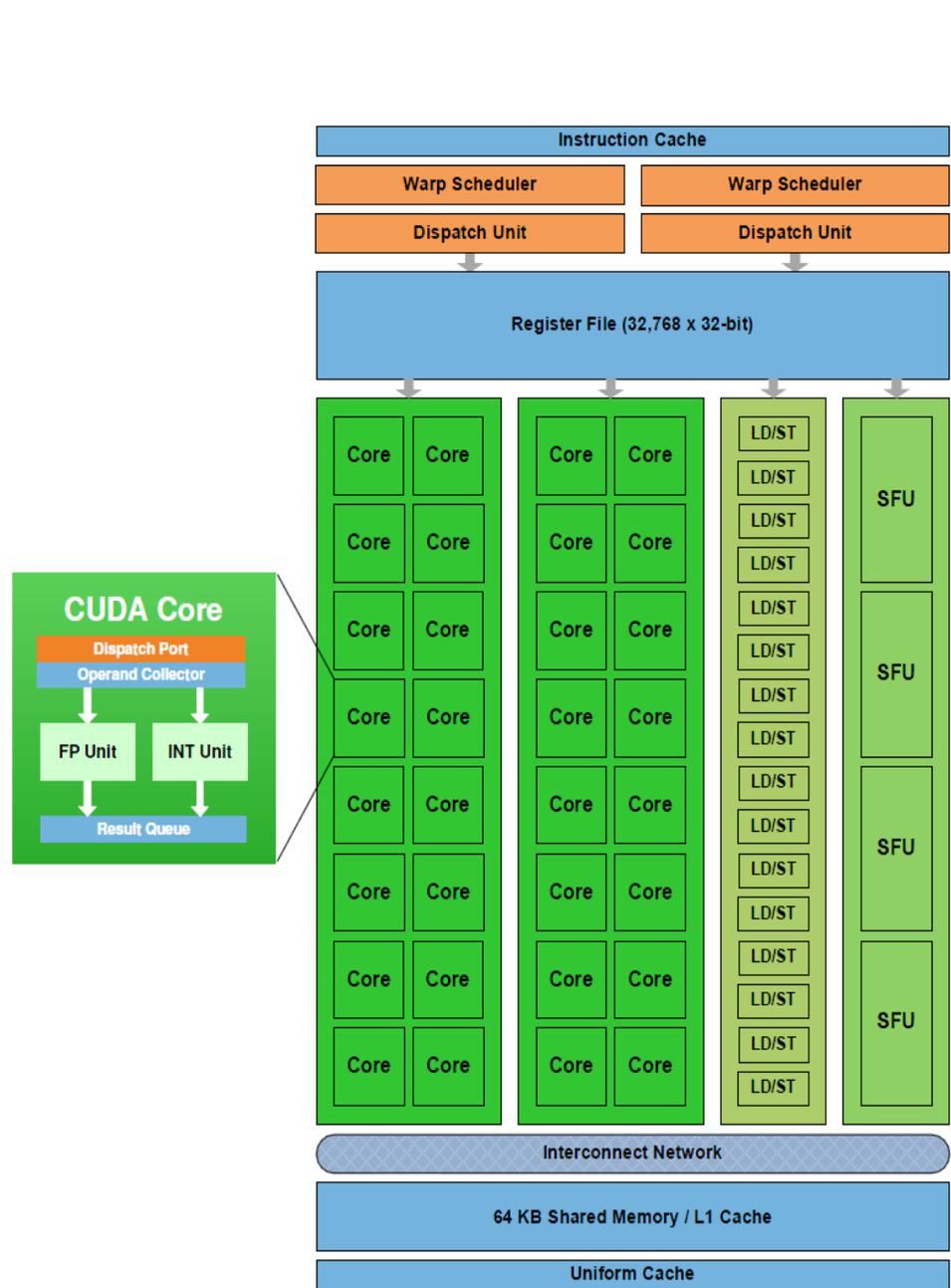
- **With ordering**

- L2 Cache Hit Rate (L1 Reads)
 - 92%
 - Global Load L2 Transactions/Access
 - 23.4
 - Maximum Divergence
 - 65.7%
- The overhead of histogram sort is low (~1ms)

Thanks!

<https://github.com/littlemine>

Xinlei Wang, 王鑫磊



GPU Execution Model

Other Useful Engineering Tips

- For Performance:
 - SoA memory layout
 - Per-material computation, separate material properties from particle attributes
- For Code Reusability:
 - Entity-Component System
 - Particle extrinsics formulation relies on certain components (MLS/non-MLS, PIC/FLIP/APIC)
 - Functional Programming
 - Implicit Time Integration involves lots of similar grid operations
 - Transfer schemes can be formulated by various submodules (kernel, transfer method)
 - Easier to make task parallel