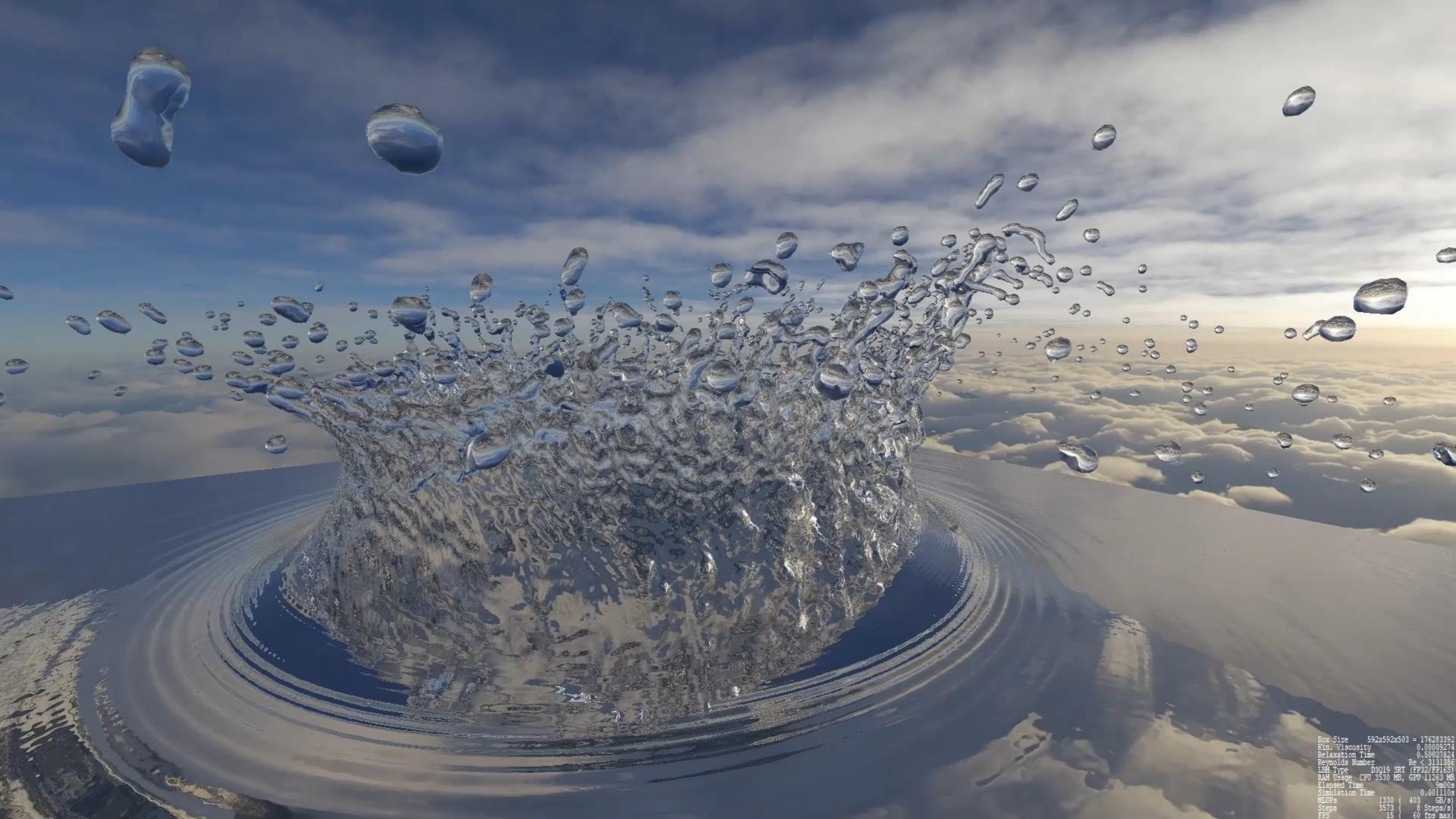


Combined scientific CFD simulation and interactive raytracing with OpenCL

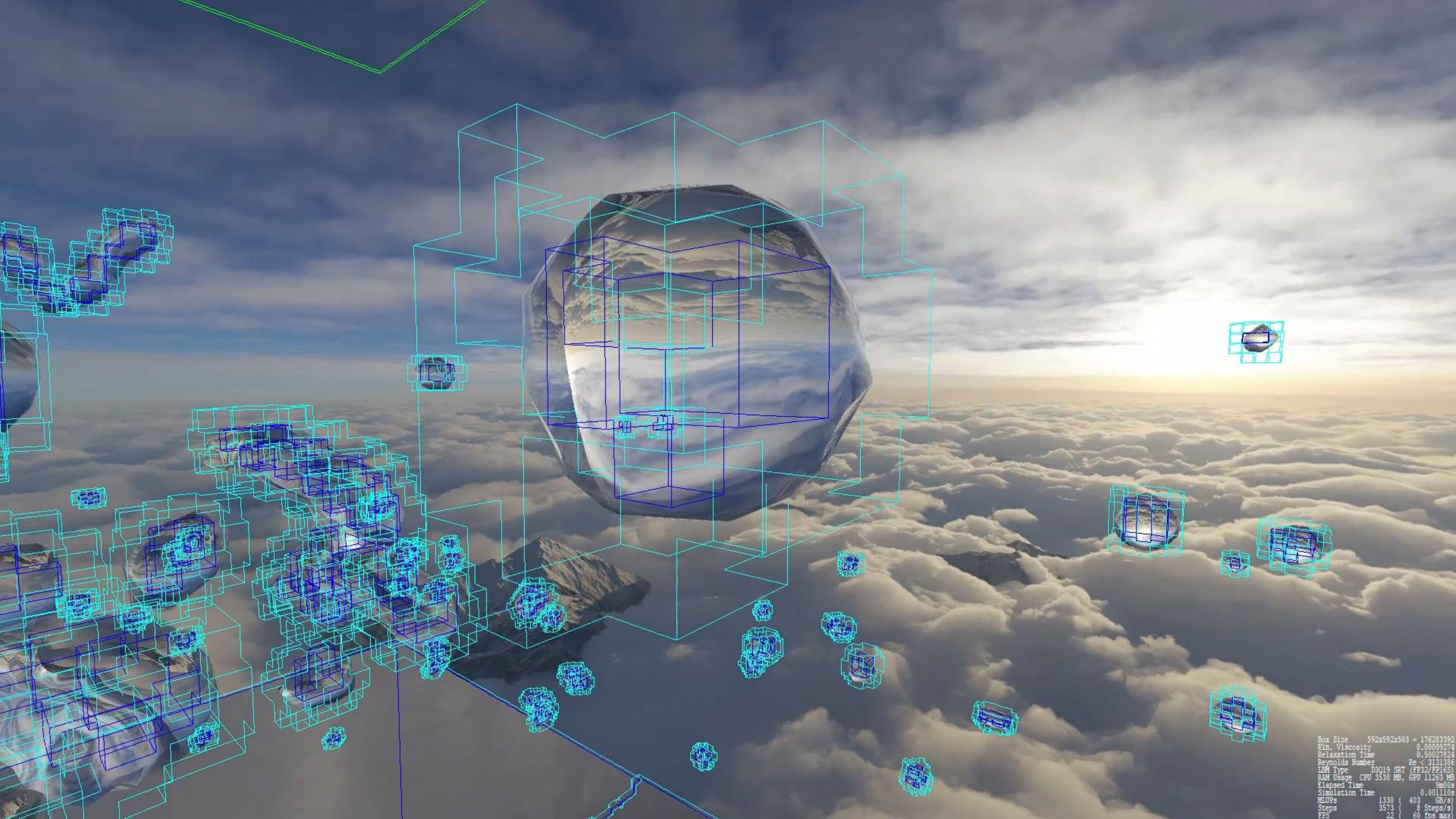
The FluidX3D Software



Box Size 96x352x96 = 3244032
Kin. Viscosity 0.00700000
Relaxation Time 0.52100004
Reynolds Number Re < 7918
LEM Type D3Q19 SRT (FP32/FP16S)
RAM Usage CPU 64 MB, GPU 207 MB
Elapsed Time 1m14s
Simulation Time 3552s
MLPs 1605 486 GB/s
Steps 35525 495 Steps/s
FPS 84 60 fps max



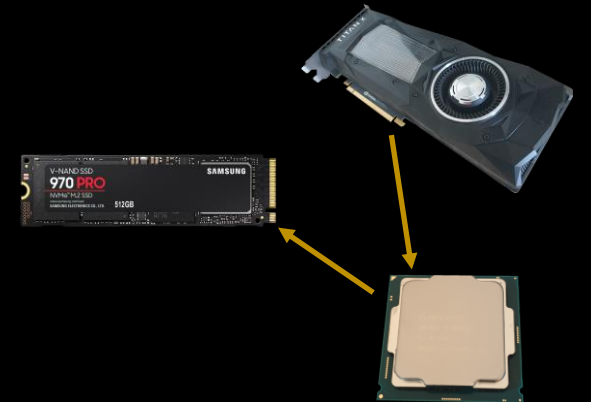
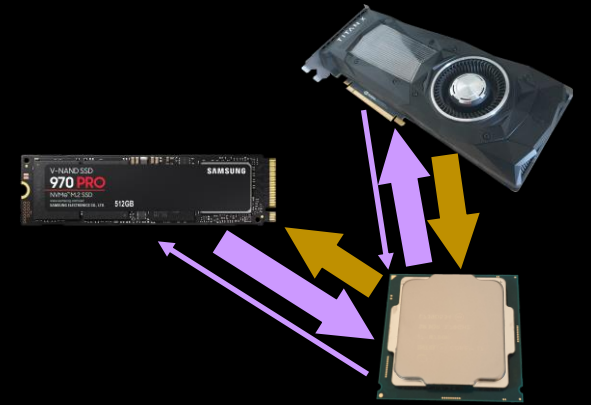
Box Size 592x592x503 = 176283392
Kin. Viscosity 0.000024
Relaxation Time 0.58027824
Reynolds Number $Re < 3131386$
LEM Type D3016 SRT (FP32/FP16S)
RAM Usage CPU 3530 MB, GPU 11263 MB
Elapsed Time 5m0s
Simulation Time 0.00110s
MLUs 1330 403 GB/s
Steps 3573 8 Steps/s
FPS 15 60 fps max



Box Size 592x592x503 = 176283392
Kin. Viscosity 0.00009274
Relaxation Time 0.50027824
Reynolds Number $Re < 3131386$
LBM Type D3Q19 SRT (FP32/FP16S)
RAM Usage CPU 3530 MB, GPU 11263 MB
Elapsed Time 8m00s
Simulation Time 0.001110s
MLLEs 1330 (403 GB/s)
Steps 3573 (8 Steps/s)
FPS 22 60 fps max

Both Simulation and Graphics with OpenCL

- traditionally
 - run simulation on GPU
 - **store volumetric frames on hard disk**
 - **render with external software**
 - > simulation would be fast, but IO takes forever
- much better alternative
 - run simulation on GPU
 - use GPU to render images from raw data in VRAM
 - **use IO only for rendered images**
 - > fast enough for interactive graphics

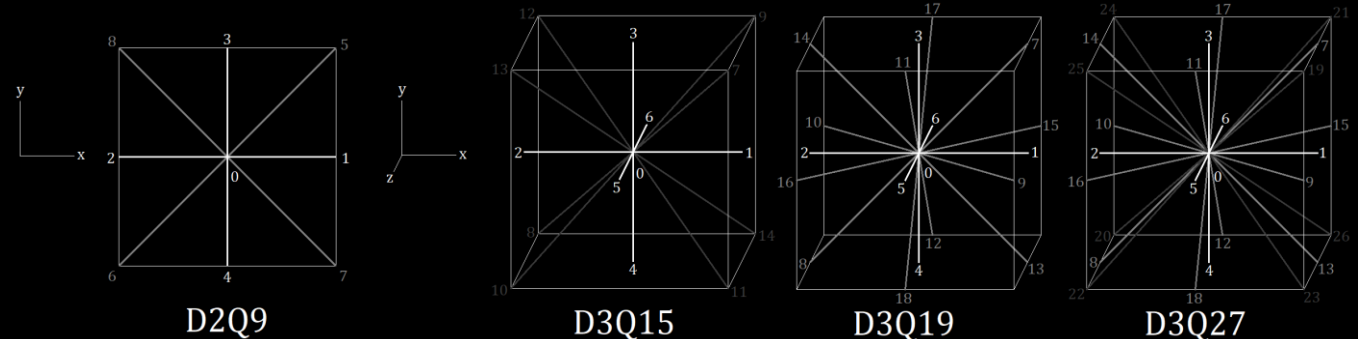


Why OpenCL™?

- OpenCL 1.2 runs on any device since ~2012
 - gaming GPUs
 - „professional“ GPUs
 - datacenter GPUs
 - Xeon Phi
 - CPUs
 - my phone
- lots of math/vector functionality already built-in
- 100% efficiency on some GPUs (A100/V100, GeForce Ampere/Turing)

Physically Accurate CFD Simulation

- lattice Boltzmann method
 - in-place streaming with implicit bounce-back
 - optional FP32/FP16 mixed precision for memory compression
--> together: 55 Bytes/node (LBM), 67 Bytes/node (FSLBM)
- Free Surface LBM
 - Volume-of-Fluid model
 - surface tension via [analytic PLIC](#)
 - improved mass conservation



On the accuracy and performance of the lattice Boltzmann method with 64-bit, 32-bit and novel 16-bit number formats

Moritz Lehmann^{1*}, Mathias J. Krause², Giorgio Amati³, Marcello Sega⁴, Jens Harting^{4,5} and Stephan Gekke¹

February 1, 2022

*Correspondence: moritz.lehmann@uni-bayreuth.de
¹Biofluid Simulation and Modeling – Theoretische Physik VI, University of Bayreuth
²Institute of Mechanical Process Engineering and Mechanics, Karlsruhe Institute of Technology
³SCAI – SuperComputing, Applications and Innovation Department, CINECA
⁴Helmholtz Institute Erlangen-Nürnberg for Renewable Energy, Forschungszentrum Jülich
⁵Department of Chemical and Biological Engineering and Department of Physics, Friedrich-Alexander-Universität Erlangen-Nürnberg

Keywords: LBM, floating-point, FP16, Posit, mixed precision, customized precision, GPU, OpenCL

1 Introduction

The lattice Boltzmann method (LBM) [1–4] is a powerful tool to simulate fluid flow. The parallel nature of the underlying algorithm has led to (multi-)GPU implementations [5–6] becoming a popular choice as speedup can be up to two orders of magnitude compared to CPUs at similar power consumption. However, most GPUs have only poor FP64 (double precision) arithmetic capabilities and thus the vast majority of GPU codes has been implemented in FP32 (single precision), while most CPU codes are written in FP64. This difference, and in particular whether FP32 is sufficient for LBM simulations compared to FP64, has been a point of persistent discussion within the LBM community [13–18, 28–33, 49–55, 57–60]. Nevertheless, only few papers [17, 32, 33, 49, 57, 61] provide some comparison on how floating-point formats affect the accuracy of the LBM and mostly find only insignificant differences between FP64 and FP32 except at very low velocity and where floating-point round-off leads to spontaneous symmetry breaking. Besides the question of accuracy, a quantitative performance comparison across different hardware microarchitectures is missing as the vast majority of LBM software is either written only for CPUs [62–74] or only for Nvidia GPUs [27–33] or CPUs and Nvidia GPUs [16–26].

A second point of concern has been the amount of video memory on GPUs, which is in general smaller than standard memory on CPU systems and can thus lead to restrictions in domain size. LBM solely works on density distribution functions (DDFs) f_α (also called fluid populations) – floating-point numbers [75–78] – that

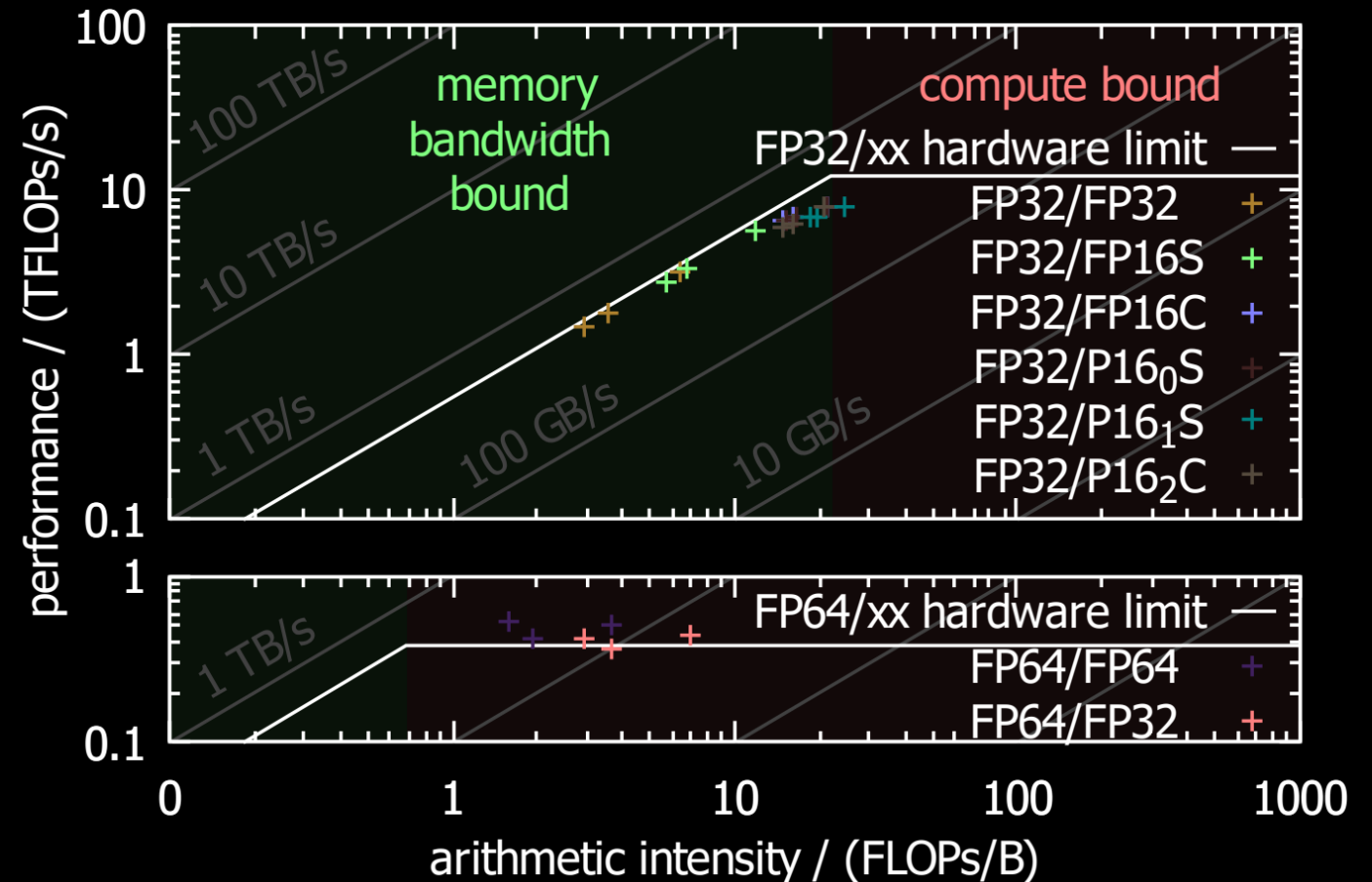
Abstract

Fluid dynamics simulations with the lattice Boltzmann method (LBM) are very memory-intensive. Alongside reduction in memory footprint, significant performance benefits can be achieved by using FP32 (single) precision compared to FP64 (double) precision, especially on GPUs. Here, we evaluate the possibility to use even FP16 and Posit16 (half) precision for storing fluid populations, while still carrying arithmetic operations in FP32. For this, we first show that the commonly occurring number range in the LBM is a lot smaller than the FP16 number range. Based on this observation, we develop novel 16-bit formats – based on a modified IEEE-754 and on a modified Posit standard – that are specifically tailored to the needs of the LBM. We then carry out an in-depth characterization of LBM accuracy for six different test systems with increasing complexity: Poiseuille flow, Taylor-Green vortices, Karman vortex streets, lid-driven cavity, a microcapsule in shear flow (utilizing the immersed-boundary method) and finally the impact of a raindrop (based on a Volume-of-Fluid approach). We find that the difference in accuracy between FP64 and FP32 is negligible in almost all cases, and that for a large number of cases even 16-bit is sufficient. Finally, we provide a detailed performance analysis of all precision levels on a large number of hardware microarchitectures and show that significant speedup is achieved with mixed FP32/16-bit.

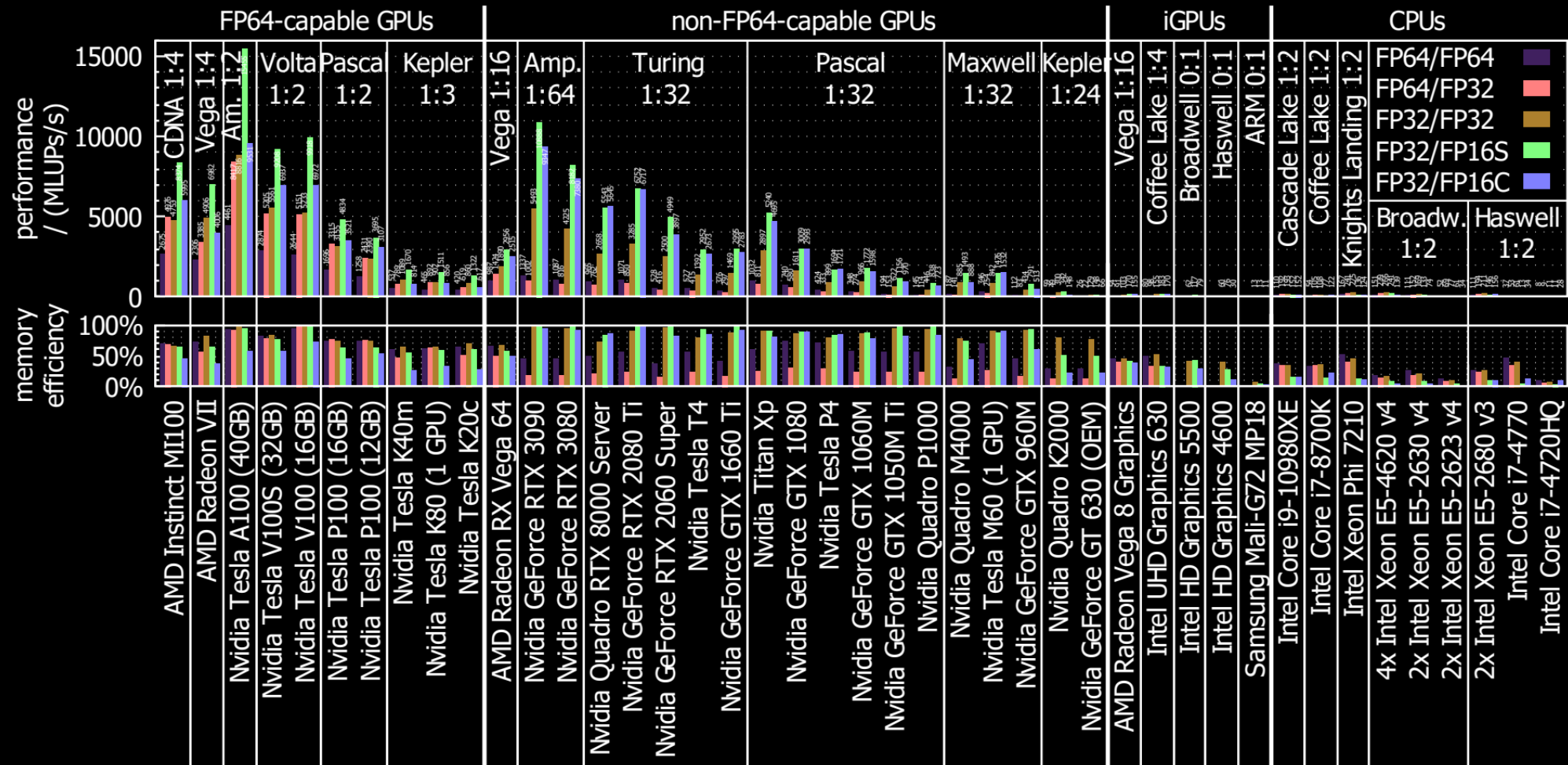
arXiv:2112.08926v2 [physics.comp-ph] 31 Jan 2022

Simulation Efficiency – Roofline Model

- Nvidia Titan Xp
- D3Q19, SRT/TRT/MRT (left to right)
- one-step-pull streaming
- FLOPs+IOPs counted combined
- mixed precision:
arithmetic/memory
- memory precision formats
 - FP16S: shifted-range IEEE FP16
 - FP16C: custom FP16 format
 - P16: 16-bit Posit formats




Simulation Efficiency – Hardware Survey



Application: Microplastic Transfer

Research Article | [Open Access](#) | [Published: 12 November 2021](#)

Ejection of marine microplastics by raindrops: a computational and experimental study

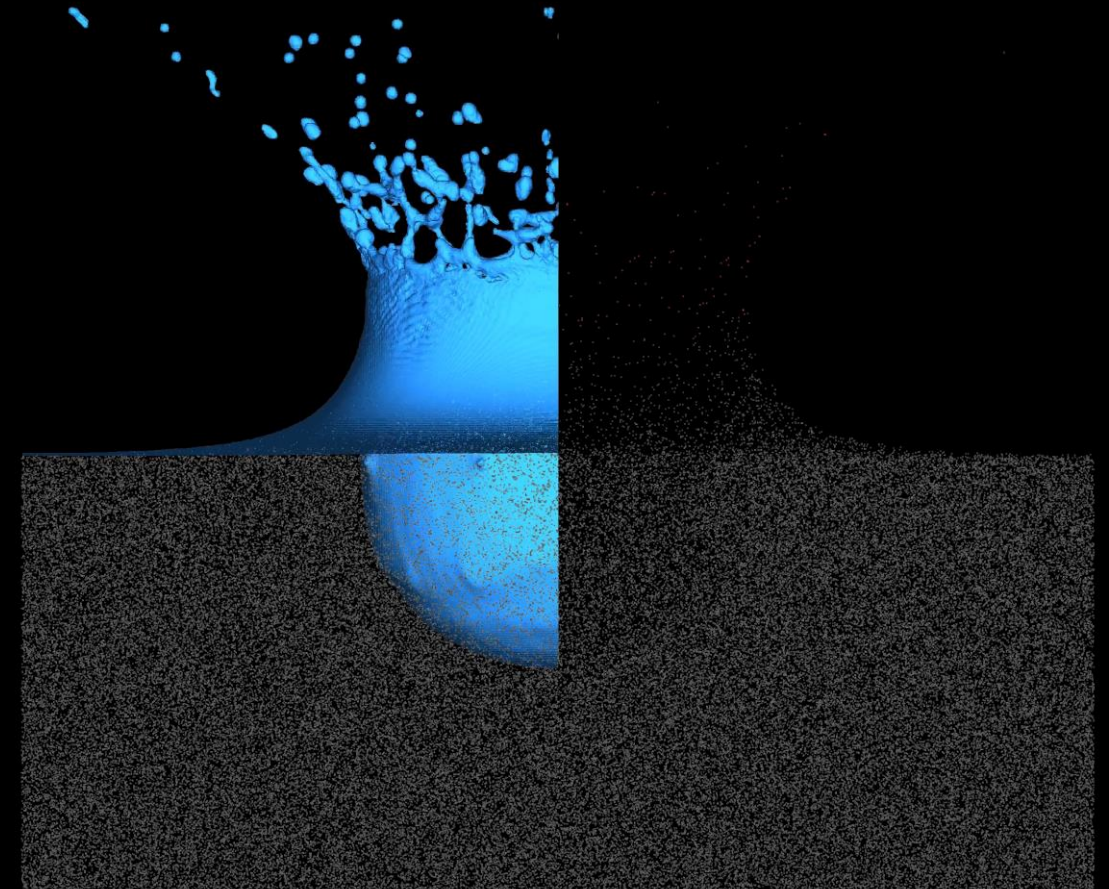
[Moritz Lehmann](#) , [Lisa Marie Oehlschlägel](#), [Fabian P. Häußl](#), [Andreas Held](#) & [Stephan Gekle](#)

[Microplastics and Nanoplastics](#) **1**, Article number: 18 (2021) | [Cite this article](#)

2153 Accesses | **2** Citations | **78** Altmetric | [Metrics](#)

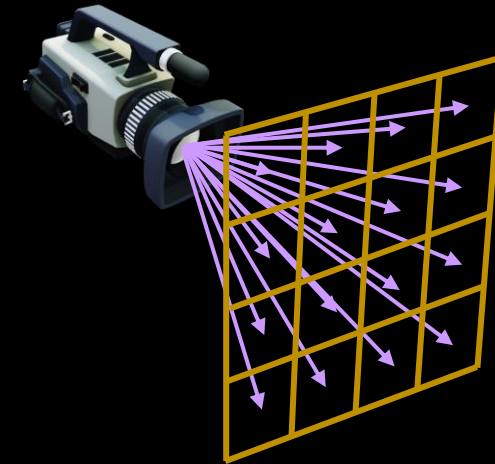
Abstract

Raindrops impacting water surfaces such as lakes or oceans produce myriads of tiny droplets which are ejected into the atmosphere at very high speeds. Here we combine computer simulations and experimental measurements to investigate whether these droplets can serve as transport vehicles for the transition of microplastic particles with diameters of a few tens of μm from ocean water to the atmosphere. Using the Volume-of-Fluid lattice Boltzmann method, extended by the immersed-boundary method, we performed more than 1600 raindrop impact simulations and provide a detailed statistical analysis on the ejected droplets. Using typical sizes and velocities of real-world raindrops – parameter ranges that are very challenging for 3D simulations – we simulate straight impacts with various raindrop diameters as well as oblique impacts. We find that a 4mm diameter raindrop impact on average ejects more than 167 droplets. We show that these droplets indeed contain microplastic concentrations similar to the ocean water within a few millimeters below the surface. To further assess the plausibility of our simulation results, we conduct a series of laboratory experiments, where we find that microplastic particles are indeed contained in the spray. Based on our results and known data – assuming an average microplastic particle concentration of 2.9 particles per liter at the ocean surface – we estimate that, during rainfall, about 4800 microplastic particles transition into the atmosphere per square kilometer per hour for a typical rain rate of $10 \frac{\text{mm}}{\text{h}}$ and vertical updraft velocity of $0.5 \frac{\text{m}}{\text{s}}$.



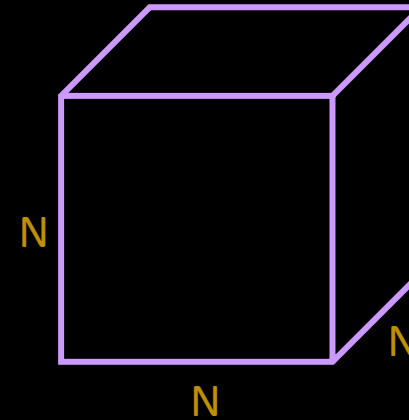
Graphics with OpenCL

- image = int array
- OpenCL can write to an array, so it can render; no OpenGL/Vulkan required
- rasterization
 - points/lines/triangles/circles... --> Bresenham
- raytracing
 - generate rays from camera
 - ray-triangle/rhombus intersection --> Möller-Trumbore
 - reflection/refraction --> Snell's law
 - UV-mapping for skybox
- ~700 lines for OpenCL rendering engine, plus ~300 lines for fluid rendering kernels
 - transfer rendered frames to CPU
 - draw on screen / write to hard drive



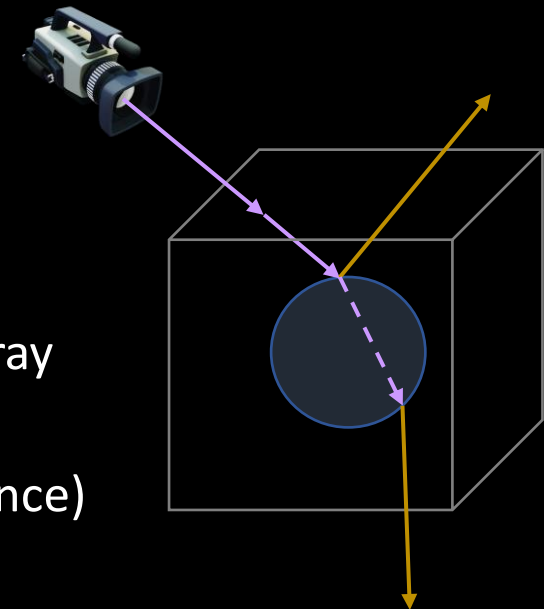
Runtime Scaling - LBM vs. Rasterization vs. Raytracing

- LBM: $\sim N^3$
- triangles of isosurface: $\sim N^2$
- surface rasterization: $\sim N^2$
- surface raytracing - BVH: $\sim \log(N)$, but generating the BVH tree is $\sim N^2$
- surface raytracing - ray-grid traversal: $\sim N$
--> benefit: runs on any hardware at peak performance, not just RTX/DXR



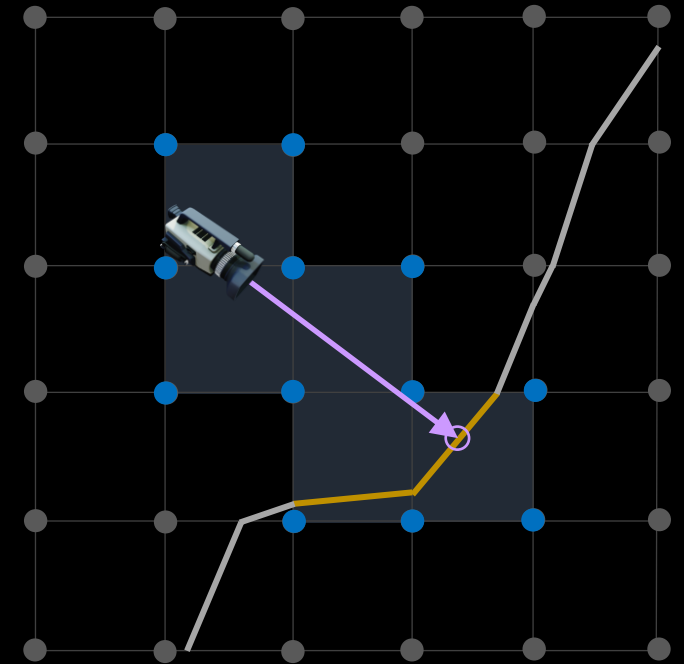
Raytracing Kernel

- for each pixel, shoot 1 ray out of the camera
- if camera is outside of grid, intersect with grid bounding box
- call ray-grid traversal on camera ray
 - no intersection --> skybox color of camera ray
 - intersection
 - reflect reflection ray & refract internal ray
 - call ray-grid traversal on internal ray & refract to refraction ray
 - get **skybox colors** of reflection ray & refraction ray
 - mix colors based on reflectivity (depends on angle of incidence)
- optional: repeat with reflection and refraction ray



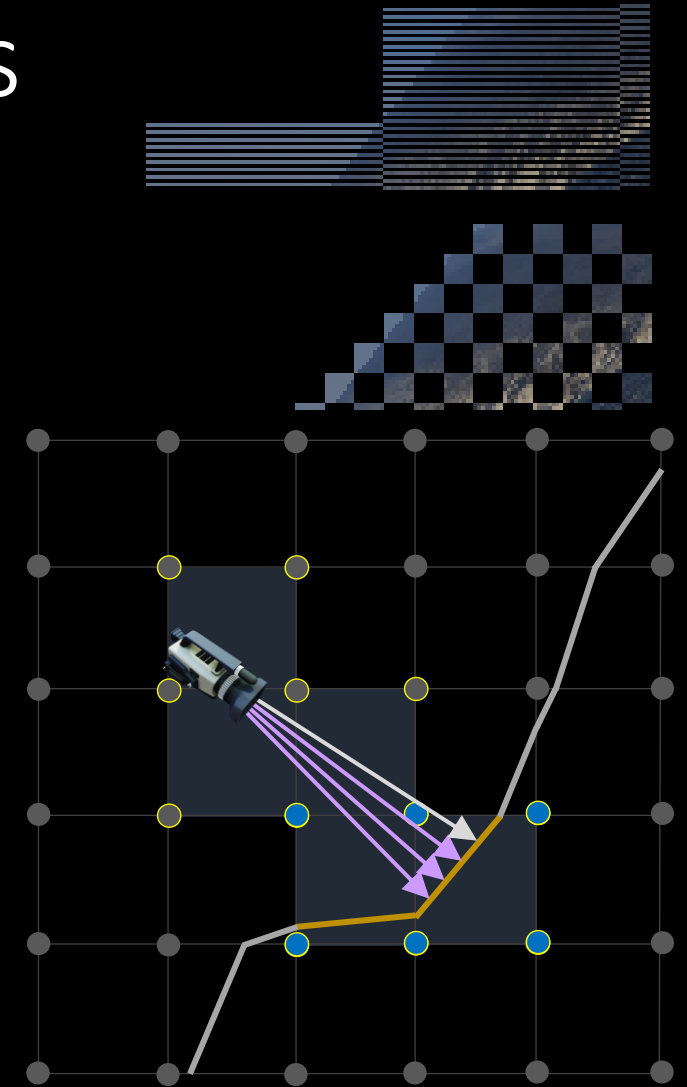
Ray-Grid Traversal

- idea: only consider **grid cells** through which a ray traverses
- for each traversed grid cell
 - check if cell contains isosurface
 - generate **triangles** with marching-cubes on-the-fly
 - check each of them for ray-triangle intersection
 - intersection
 - compute isosurface normal on all 8 grid points
 - interpolate surface normal at intersection point



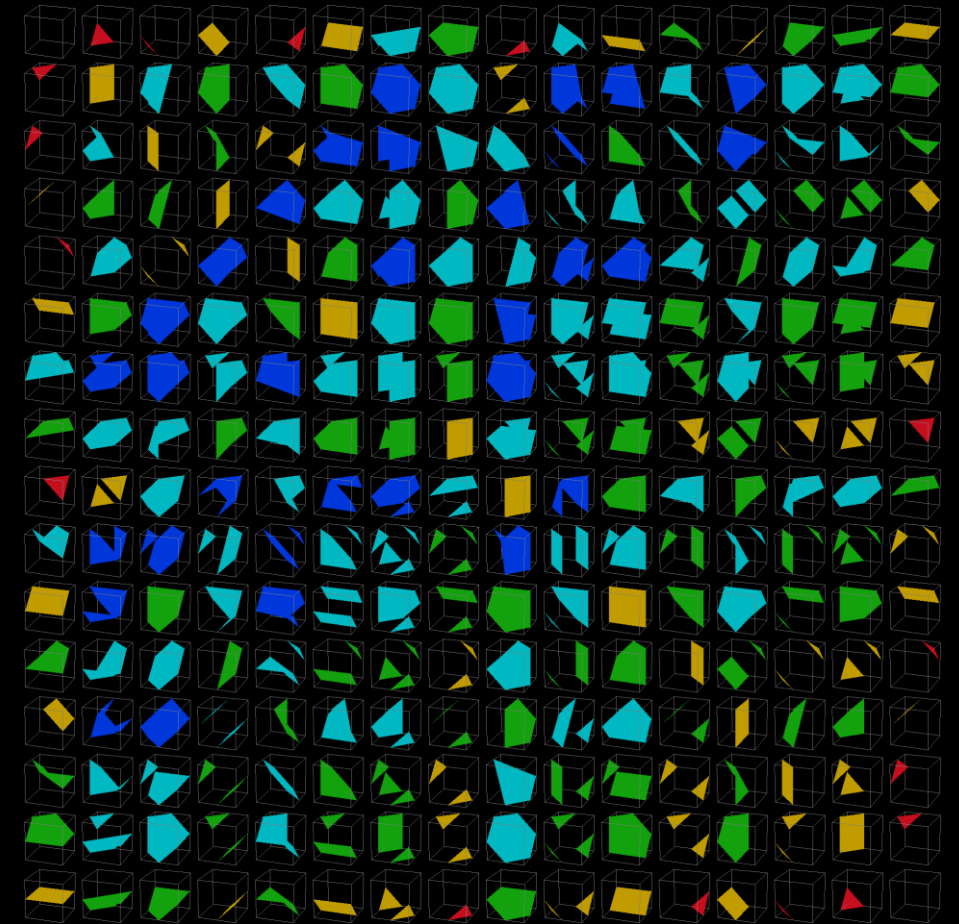
Ray-Grid Traversal - Optimizations

- workgroup alignment is critical
 - make workgroups not horizontal stripes of pixels,
 - but 8x8 pixel tiles (50% faster)
- while traversing the grid
 - instead of LBM isovalues (4 Byte), load LBM **flags** (1 Byte)
 - only if cell contains isosurface, load **isovalues** (15% faster)
- ~~• only load unknown isovalues when traversing into the next cell~~
 - ~~--> breaks memory coalescence and leads to slowdown, despite less memory access~~



Generating Triangles: Marching-Cubes

- implementation from Paul Bourke
 - 0-5 triangles/cell
- optimizations
 - reduce table size to 1/8
 - edge table: mirror symmetry + ushort, 1024B -> 256B
 - triangle table: bit-packing + uchar, 16384B -> 1920B
 - compute cubeindex branchless (bit-shifting)
- interesting:
 - edge table is not required and causes branching, but leads to ~30% speedup
 - tables must be in “constant” OpenCL memory space



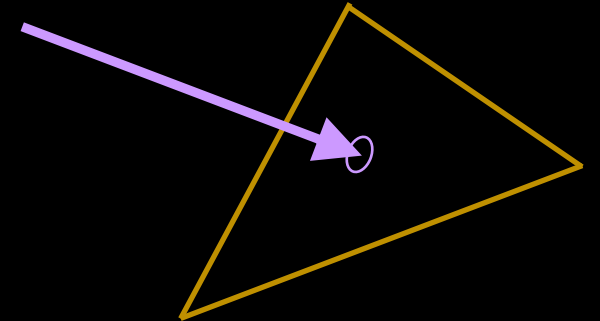
Generating Triangles: Marching-Cubes

- implementation from Paul Bourke
 - 0-5 triangles/cell
- optimizations
 - reduce table size to 1/8
 - edge table: mirror symmetry + ushort, 1024B -> 256B
 - triangle table: bit-packing + uchar, 16384B -> 1920B
 - compute cubeindex branchless (bit-shifting)
- interesting:
 - edge table is not required and causes branching, but leads to ~30% speedup
 - tables must be in “constant” OpenCL memory space

[illegible]

Ray-Triangle Intersection

- branchless version of Möller-Trumbore algorithm
- bidirectional intersection
- returns distance from ray origin to intersection point



```
float intersect_triangle_bidirectional(const ray r, const float3 p0, const float3 p1, const float3 p2) {  
    const float3 u=p1-p0, v=p2-p0, w=r.origin-p0, h=cross(r.direction, v), q=cross(w, u);  
    const float f=1.0f/dot(u, h), s=f*dot(w, h), t=f*dot(r.direction, q);  
    return (s<0.0f||s>1.0f||t<0.0f||s+t>1.0f) ? -1.0f : f*dot(v, q);  
}
```

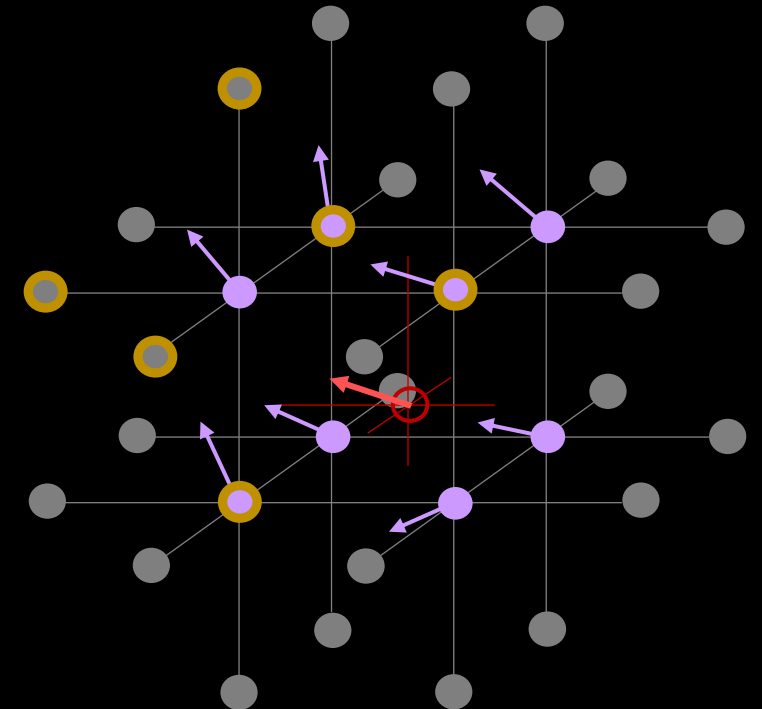
Normal Interpolation

- calculate normals on grid points
 - trilinear interpolation to hitpoint
- > comes at almost no additional cost



Normal Interpolation

- calculate normals on **grid points**
 - **central difference stencil**
 - reuse isovalues that have already been loaded
--> only 24 new isovalues instead of 56
- trilinear interpolation to **hitpoint**



Skybox

- UV-mapping (3D vector to xy-coordinates)
- bilinear interpolation for pixel color



<https://www.hdri-hub.com/hdri-skies-aviation-aerospace>

```
uint skybox_color(const ray r, const global uint* skybox) {  
    const float fu = (float)def_skybox_w*fma(atan2(r.direction.x, r.direction.y), 0.5f/3.1415927f, 0.5f);  
    const float fv = (float)def_skybox_h*fma(asin(r.direction.z), -1.0f/3.1415927f, 0.5f);  
    const int ua=clamp((int)fu, 0, (int)def_skybox_w-1), va=clamp((int)fv, 0, (int)def_skybox_h-1), ub=(ua+1)%def_skybox_w, vb=min(va+1, (int)def_skybox_h-1);  
    const uint s00=skybox[ua+va*def_skybox_w], s01=skybox[ua+vb*def_skybox_w], s10=skybox[ub+va*def_skybox_w], s11=skybox[ub+vb*def_skybox_w];  
    const float u1=fu-(float)ua, v1=fv-(float)va, u0=1.0f-u1, v0=1.0f-v1;  
    return color_mix(color_mix(s00, s01, v0), color_mix(s10, s11, v0), u0);  
}
```

Conclusions

- if the simulation is limited by PCIe/IO, you have already lost
- much faster alternative: combine compute and graphics
- rendering is just another form of compute
- OpenCL enables simulations at peak efficiency
- OpenCL can render faster than anything else
- (micro-)optimization adds up



