



Johannes Pekkila, PC user meeting 2020

Image credit: NASA/SDO

#### Accelerated stencil computations on GPUs

- Speaker: Johannes Pekkilä
  - Doctoral Student at Aalto University, Finland
  - Field: Computer Science (Big Data and Large-Scale Computing)
  - Worked on accelerating physical simulations on GPUs since 2014



#### Accelerated stencil computations on GPUs

- Astaroth is a multi-GPU library for high-performance stencil computations
- The API consists of an host-level interface (C, C++, Fortran) and a domain-specific language (DSL)
- Specialized in high-order stencils and computations with coupled fields → very high cache-efficiency in comparison to competitors



#### Recap

- Usable for production now
- Full MHD (hydro, magnetic, entropy) + forcing + upwinding
- All physics are implemented in the DSL
- Efficient multi-GPU MPI implementation (GPUDirect RDMA)



## This talk

- 1) Performance
- 2) Astaroth DSL
- 3) Astaroth API



### Single-node performance

Full MHD, double precision

#### Pencil Code on 2x Intel Xeon Gold 6230 CPUs, 40 cores:

23 ns / grid point / step

35x speedup

Astaroth on 4x Tesla V100 GPUs
0.65 ns / grid point / step



# Strong scaling



Table 8: Strong scaling for 6th-order stencils when varying mesh size. Times are in milliseconds per integration step.

nprocs	256 <sup>3</sup>	512 <sup>3</sup>	1024 <sup>3</sup>
1	41.74		
2	21.08	155.03	
4	10.70	72.80	
8	6.89	40.07	300.40
16	6.66	25.36	153.33
32	5.46	19.33	73.82

Figure 15: Strong scaling for 6th-order stencils. The largest submesh that could be held in the memory of a single GPU was approximately 448<sup>3</sup>.

## This talk

- 1) Performance
- 2) Astaroth DSL
- 3) Astaroth API



## Astaroth DSL

Disclaimer

Astaroth DSL (AKA Astaroth Code, AC) is a procedural stream programming language

Why?

- 1) Industry standard
- 2) Easy to translate to CUDA (less bugs)
- 3) Stream programming maps well to GPU hardware

Astaroth DSL is NOT

- 1) a functional programming language.
- 2) the simplest possible representation for physical equations.

All major design decisions are based on obtaining maximum performance with the simplest possible language, but no simpler.



## **Domain-specific languages**

```
__constant__ float alpha;

    DSL below, optimized CUDA to the right

                                                                                          static __device__ __forceinline__ float3
• CUDA code presented here is overly complex for
                                                                                          laplacian(const PreprocessedVertex3& vtx)
                                                                                             return (float3){vtx.x.laplacian, vtx.y.laplacian, vtx.z.laplacian};
   just the heat equation, but demonstrates the
   optimizations required for more complex tasks
                                                                                          static __device__ __forceinline__ float3
                                                                                          heat_equation(const PreprocessedVertex3& T)
                                                                                             return alpha * laplacian(T);
                                                                                          static __device__ __forceinline__ float3
uniform Scalar alpha;
                                                                                          stencil process(const PreprocessedVertex3& T, const float dt)
                                                                                             return T.value + heat_equation(T) * dt;
 Vector
 laplacian(in Vector T)
                                                                                          template <int step number>
                                                                                          static global void
      return (Vector){laplacian(T.x), laplacian(T.y), laplacian(T.z)};
                                                                                          __launch_bounds__(RK_THREADBLOCK_SIZE, RK_LAUNCH_BOUND_MIN_BLOCKS)
                                                                                               VertexBufferArray buffer)
 Vector
                                                                                             const int3 vertexIdx = (int3) {
 heat_equation(in Vector T)
                                                                                                threadIdx.y + blockIdx.y * blockDim.y + start.y,
                                                                                                threadIdx.z + blockIdx.z * blockDim.z + start.z
 {
     return alpha * laplacian(T);
                                                                                             if (vertexIdx.x >= end.x ||
                                                                                                vertexIdx.y >= end.y ||
                                                                                                vertexIdx.z >= end.z) {
 in VectorField T_in(X, Y, Z);
                                                                                                    return;
 out VectorField T_out(X, Y, Z);
                                                                                                                                    buffer.in[0],
 Kernel
                                                                                                                                   buffer.in[1],
                                                                                                                                    buffer.in[2]);
 solve(Scalar dt)
                                                                                             const float3 result = stencil_process(T, dt);
 {
                                                                                             buffer.out[0][IDX(vertexIdx)] = result[0];
     T out = T in + heat equation(T in) * dt;
                                                                                             buffer.out[1][IDX(vertexIdx)] = result[1];
 }
                                                                                             buffer.out[2][IDX(vertexIdx)] = result[2];
```

- Syntax of AC is C-like and similar to shading languages
- We provide a standard DSL library for computing derivatives and commonly-used math operations (curl, laplacian, etc.)
- The standard library can be extended to support nonequidistant grids, and spherical and cylindrical coordinate systems in the future



Function type qualifiers

- Kernel. The main function, analogous to \_\_\_\_\_global\_\_\_\_
- Device. Analogous to \_\_\_\_\_device\_\_\_\_
- Preprocessed. Evaluated at the preprocessing stage and cached for use in *Kernel* and *Device* functions
- <Empty>. Helper functions for *Preprocessed* functions.



- Function type qualifiers: *Preprocessed*, *Kernel*, and *Device* functions
- **Data types**: *Scalar, Vector, Complex, Matrix, ScalarField, VectorField.* Precision of real numbers is determined at compiletime. Device constants are declared with *uniform* type qualifier
- Input and output arrays are declared with in and out qualifiers

#### • Built-in variables:

- vertexIdx local index of the current vertex
- globalVertexIdx global index of the current vertex, offset based on the multi-GPU decomposition
- globalGridN dimensions of the computational domain



#include <stdderiv.h>

```
uniform ScalarField FIELD_HANDLE;
uniform Scalar dt = 1.0;
uniform int flag = 0;
in ScalarField field_in(FIELD_HANDLE);
out ScalarField field_out(FIELD_HANDLE);
Preprocessed
blur(in ScalarField field)
    return (field[vertexIdx] +
            field[vertexIdx.x + 1, vertexIdx.y,
                                                       vertexIdx.z] +
            field[vertexIdx.x - 1, vertexIdx.y,
                                                      vertexIdx.z] +
            field[vertexIdx.x, vertexIdx.y + 1, vertexIdx.z] +
field[vertexIdx.x, vertexIdx.y - 1, vertexIdx.z] +
            field[vertexIdx.x, vertexIdx.y,
                                                     vertexIdx.z + 1] +
            field[vertexIdx.x, vertexIdx.y,
                                                       vertexIdx.z - 1]) / 7.0;
Kernel
solve()
    if (flag) {
        field_out = value(field_in) + derx(field_in) * dt;
    } else {
        field_out = blur(field_in);
```

Aalto University School of Science

## This talk

- 1) Performance
- 2) Astaroth DSL
- 3) Astaroth API



- Astaroth library handles host-side functionality
  - Memory operations (initialization, transfers)
  - Queuing kernels
- Essentially all API functions are asynchronous by default
  - Required for efficient multi-GPU communication



The API is divided into two layers: Device and Node

#### Device layer

- For controlling a single GPU
- Functions start with acDevice

#### Node layer

- Abstracts the GPUs available on a node behind a single interface
- Functions start with acNode



- The library is configured by passing a config struct when initializing devices and nodes
- Devices and nodes are managed with handles

Device device; acDeviceInit(info, &device); acDeviceLoad(mesh, device); acDeviceIntegrate(device);

Node node; acNodeInit(info, &node); acNodeLoad(mesh, node); acIntegrate(node);





#### Accelerated stencil computations on GPUs

#### Code http://bitbucket.org/jpekkila/astaroth

#### Documentation http://libastaroth.bitbucket.io/

#### Thesis http://urn.fi/URN:NBN:fi:aalto-201906233993





Johannes Pekkila, PC User Meeting 2020

Image credit: NASA/SDO

#### **Backup slides**



# Single-GPU performance

#### Test case:

- Full MHD
- 6<sup>th</sup> order finite-differences
- 3<sup>rd</sup> order Runge-Kutta

#### **CPU solver:**

- Pencil Code
- 2x Intel Xeon E5-2690 v3
- Benchmarked on a total of 24
   cores
- 47 ns / grid point / step (dbl)

#### **GPU solver:**

- Astaroth
- Tesla P100 PCIe
- 4.6 ns / grid point / step (dbl)





# Domain-specific languages

Performance and productivity

#### Takeaways:

- Future is parallel
- All problems will be eventually bound by memory bandwidth (for the foreseeable future)
- The only way forward is to increase arithmetic intensity (primarily by caching)
- Domain-specific languages can offer high productivity and performance





Johannes Pekkila, PC User Meeting 2020

Image credit: NASA/SDO