

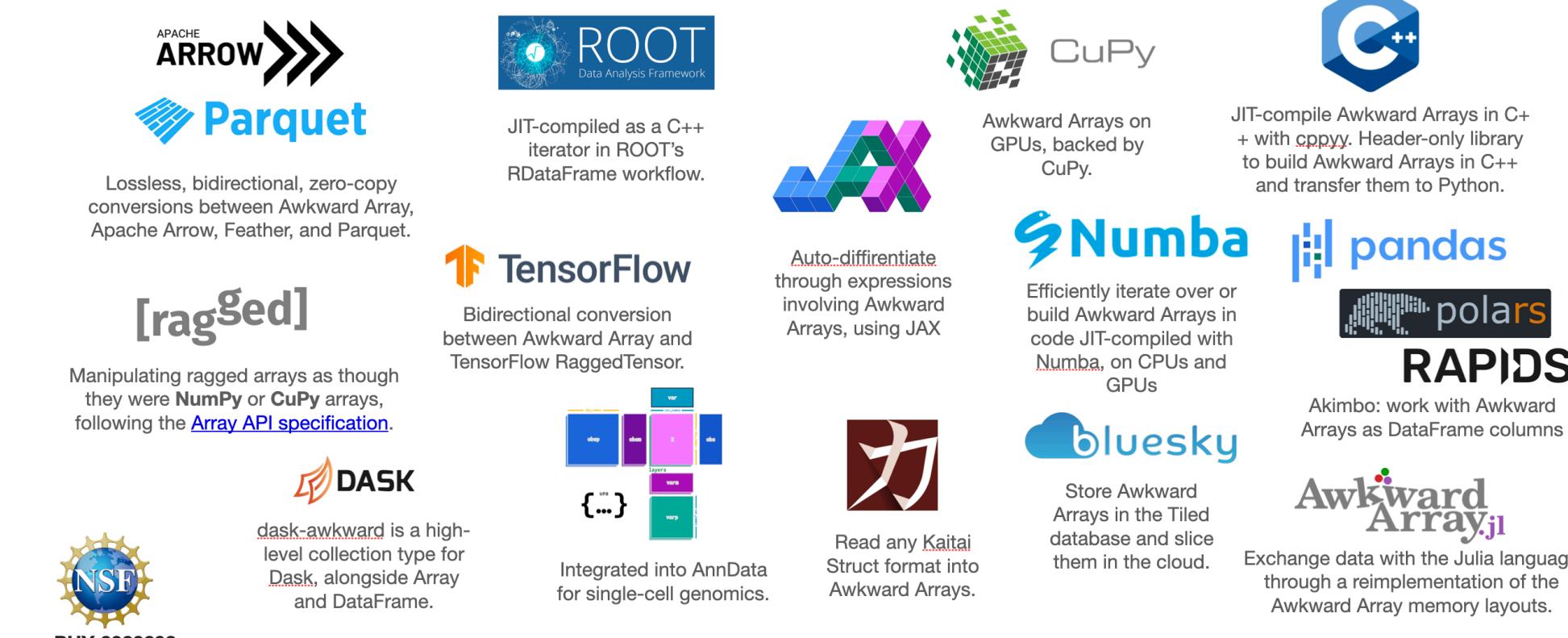


PHY-2323298

Advancing Awkward Arrays for High-Performance CPU and GPU Processing

Ianna Osborne, Princeton University, Manasvi Goyal, Harvard University

This work presents ongoing efforts to optimize Awkward Arrays for GPUs using CUDA, aiming to achieve performance parity with or surpass CPU kernel implementations. Key improvements focus on optimized memory management, leveraging CUDA-specific features, and maximizing parallelism to enhance throughput. These advancements enable faster and more scalable data processing, particularly for HL-LHC data analysis within the Python ecosystem. We will discuss the challenges, solutions, and performance benchmarks.



Awkward functions: CPU vs CUDA backend

```
array_cpu = ak.Array([[1, [6, 3]]], backend="cpu")
ak.min(array_cpu, axis=1)

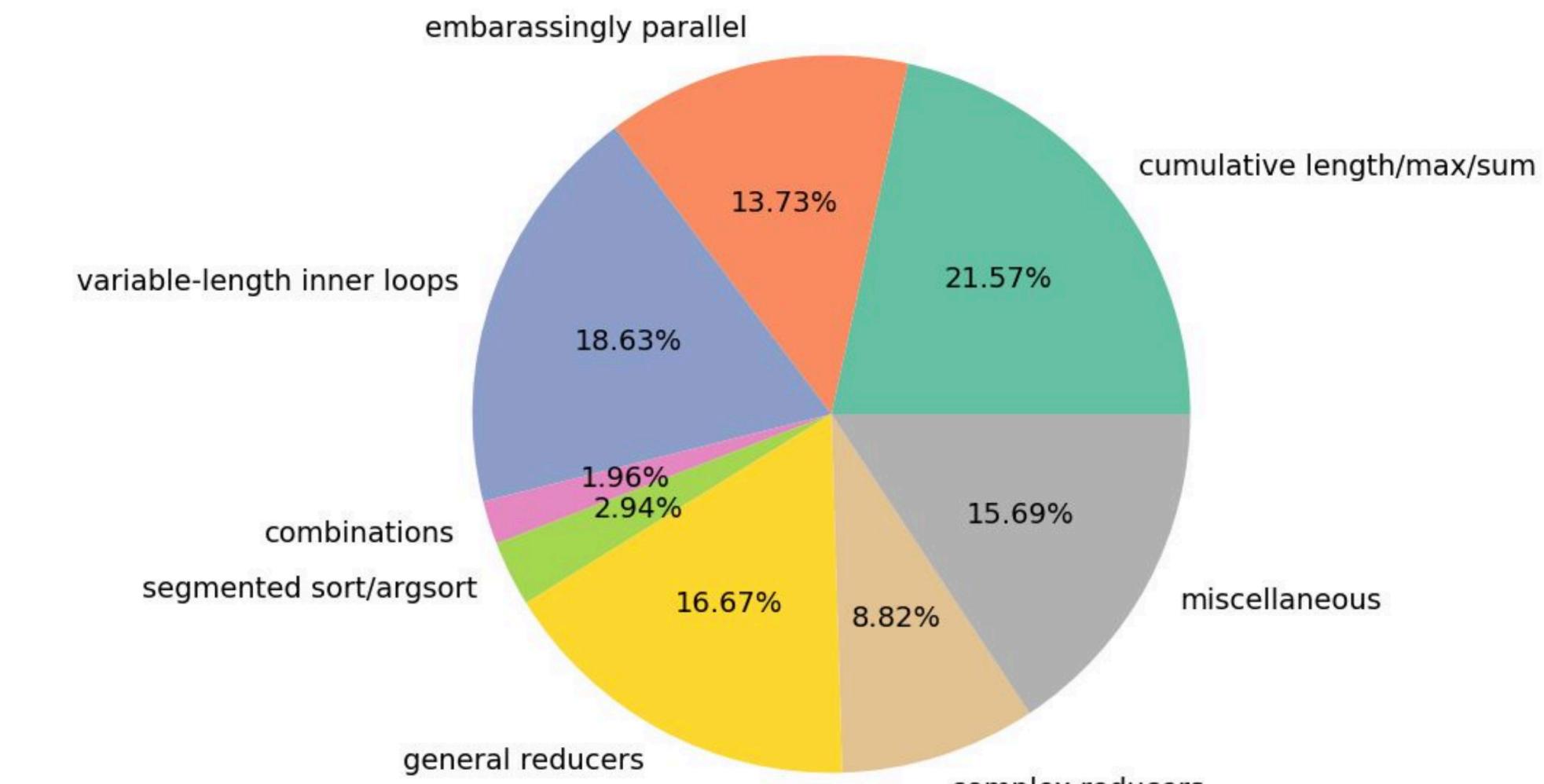
array_gpu = ak.Array([[1, [6, 3]]], backend="cuda")
ak.min(array_gpu, axis=1)

jets = ak.Array([[1, 1, 2, 1, 3, 1]], backend="cpu")
trijet = ak.combinations(jets, 3, fields=["j1", "j2", "j3"])
result = trijet.j1 + trijet.j2 + trijet.j3
ak.min(result, axis=1)

jets = ak.Array([[1, 1, 2, 1, 3, 1]], backend="cuda")
trijet = ak.combinations(jets, 3, fields=["j1", "j2", "j3"])
result = trijet.j1 + trijet.j2 + trijet.j3
ak.min(result, axis=1)
```

Awkward CPU backend

Overall 144 CPU kernels:



CuPy is used to handle the higher level functions

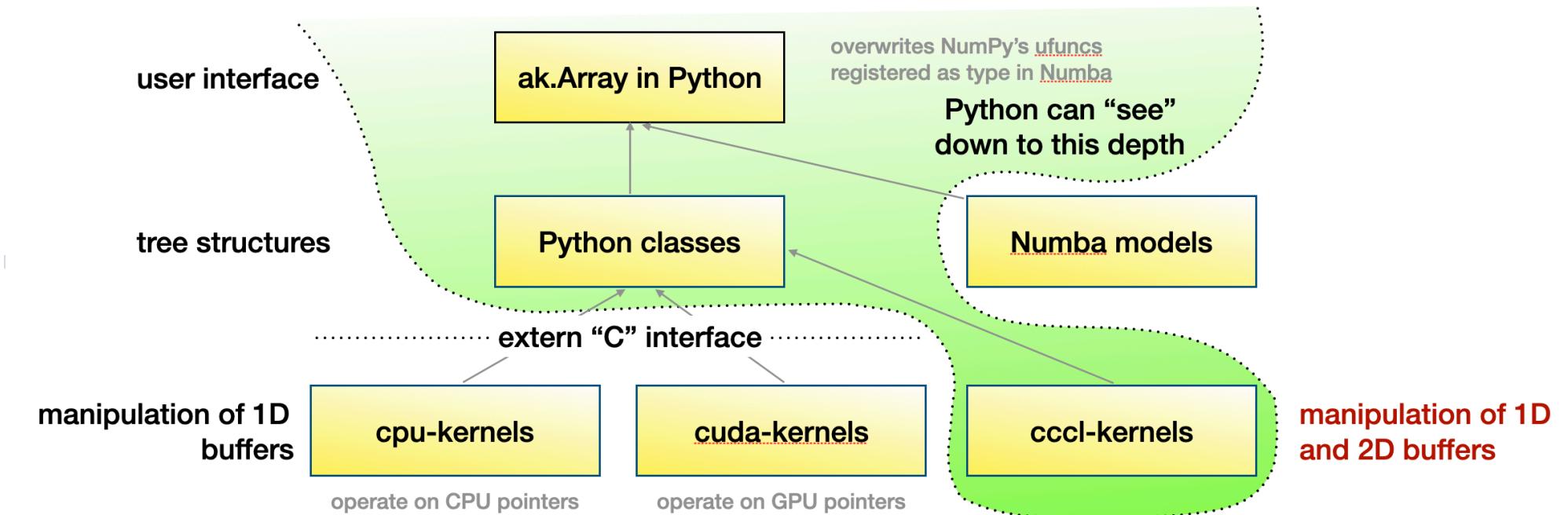
Automated infrastructure connecting CUDA kernels to Awkward operations

Testing of kernels via both generated tests:

- Python test generation scripts produce custom unit tests in Python
- they are generated from a JSON file of test cases.
- and:
- integration tests for each Awkward function on the CUDA backend

Awkward Array v2 architecture

- No direct dependency on CUDA
- Introduces an indirection as we move from the upper to the lower layers.
- Awkward functions can be used on GPU with just pip install awkward



Awkward CCCL kernels

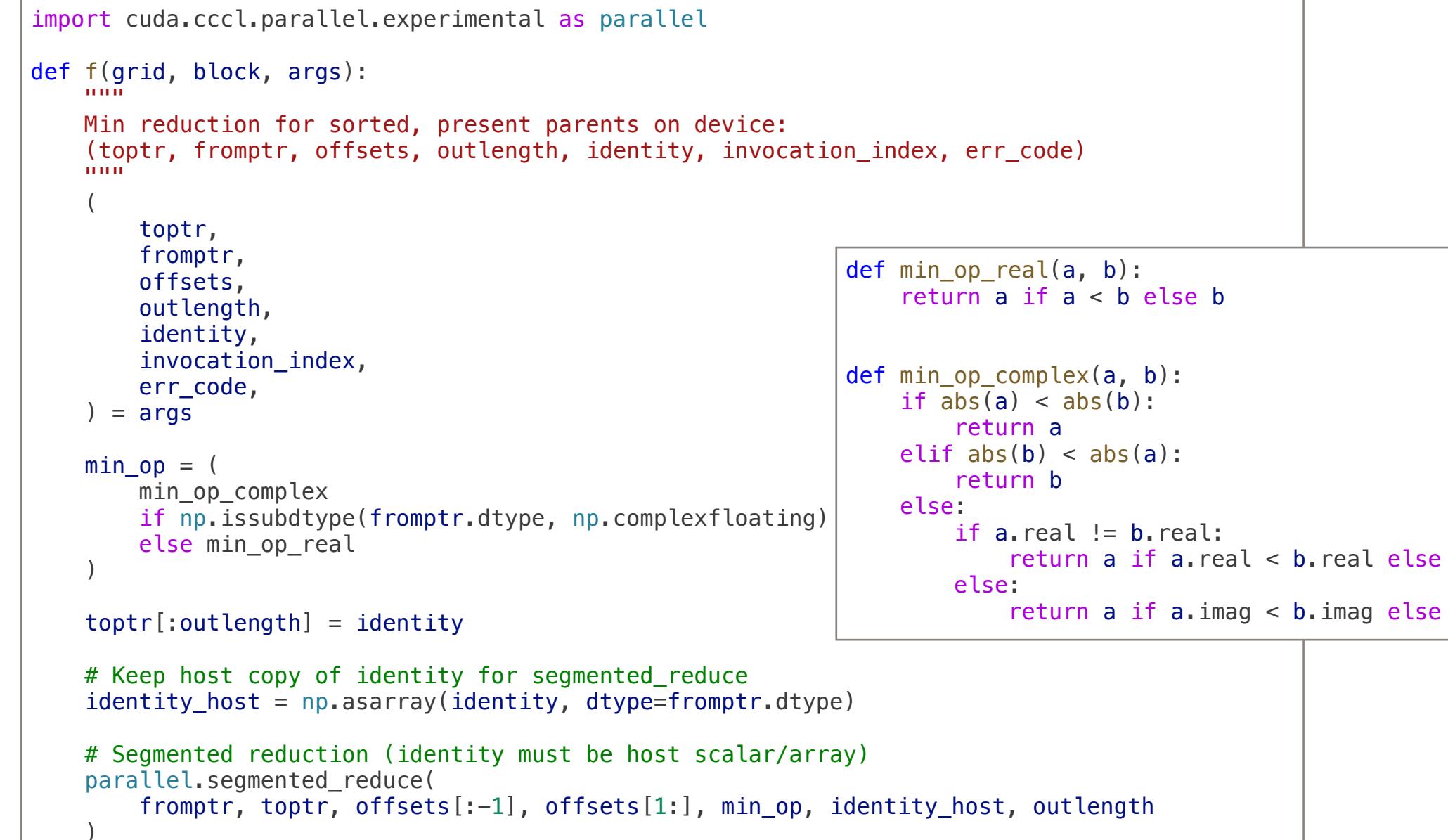


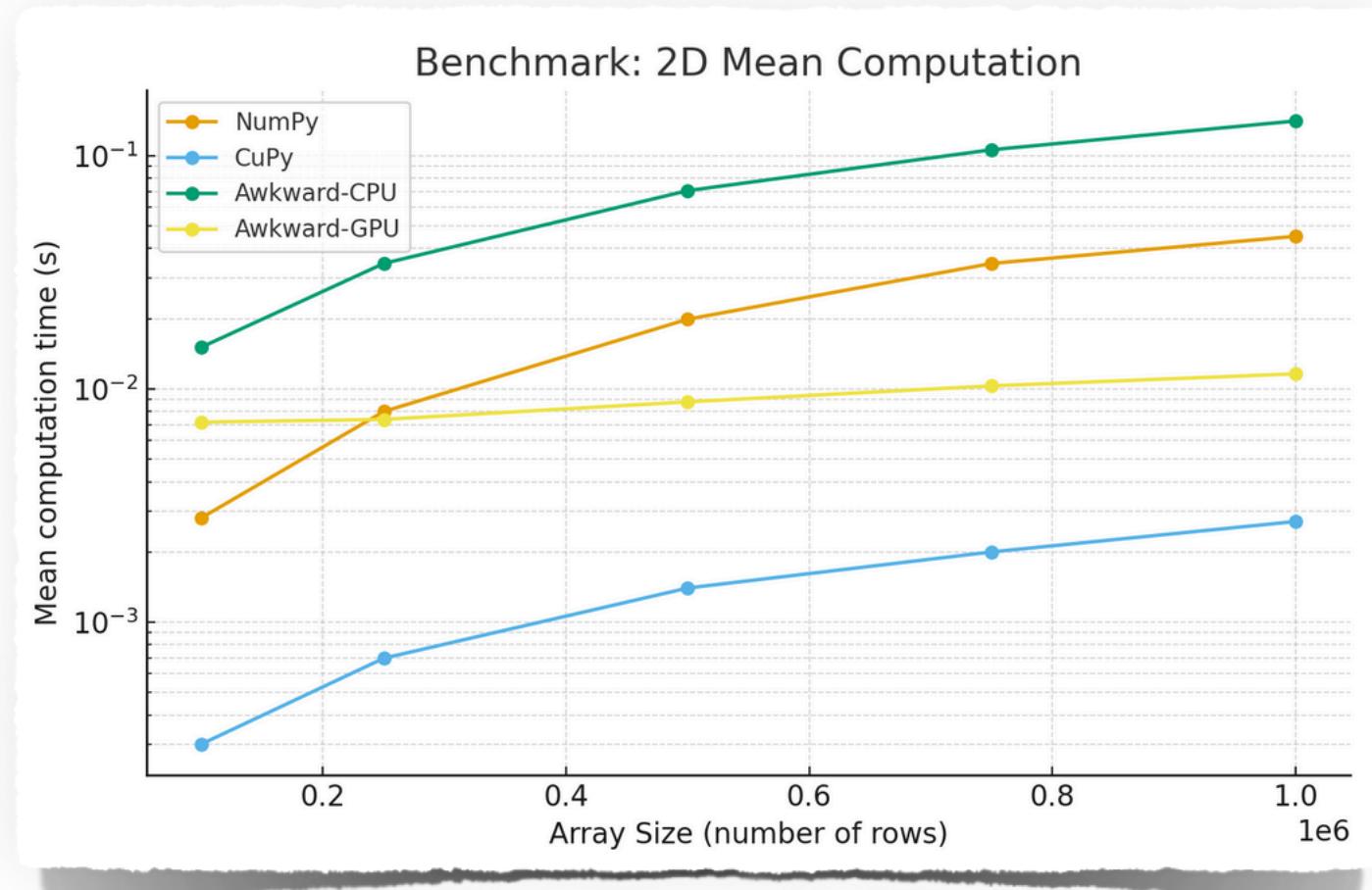
Figure 1. Architecture of CUDA-enabled Python packages like PyTorch and CuPy today (left), and the gap filled by `cuda.cccl` (right)



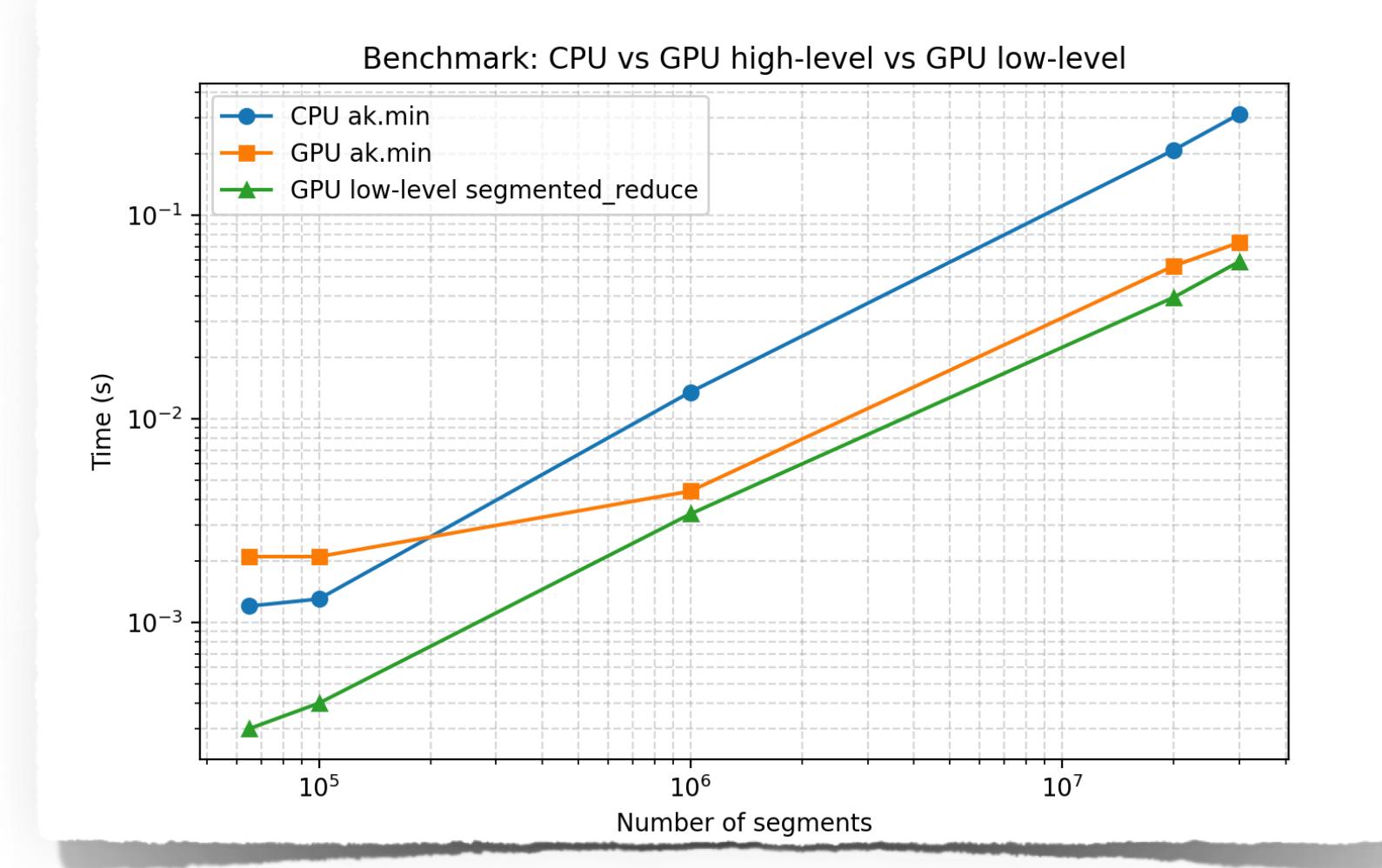
Optimizing Awkward Arrays for GPUs using the CUDA Core Compute Libraries in Python, improving memory use and parallelism to efficiently handle large, irregular data.

Performance benchmarks

NVIDIA-SMI 580.65.06		Driver Version: 580.65.06		CUDA Version: 13.0	
GPU Name	Temp	Persistence-M	Bus-Id	Disp.A	Volatile Uncorr. ECC
Fan	Pwr	Usage/Cap	Memory-Usage	GPU-Util	Compute M. MIG M.
0 NVIDIA A100-PCIE-40GB	34C	P0	On	0000:0000:21:00.0 Off	0% Default Disabled



The logarithmic y-axis makes it clear that **CuPy** and **Awkward-GPU** are significantly faster than **NumPy** and **Awkward-CPU**, especially as array size grows.



CPU times increase sharply with segment count. GPU high-level `ak.min` is slightly faster at mid-range, but slower than low-level for large arrays. GPU low-level segmented reduction is the fastest overall, especially for tens of millions of segments.

