# cuNumeric.jl :
# Automating Distributed Numerical Computing

**David Krasowska**[1], Ethan Meitz[2], Wonchan Lee[3]

[1]Northwestern University
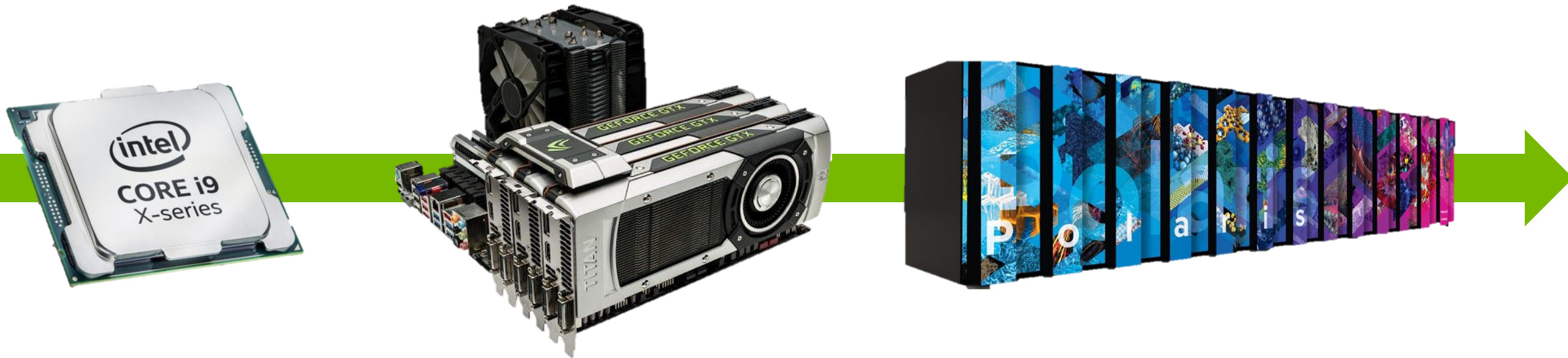[2]Carnegie Mellon University
[3]NVIDIA

SC25
St. Louis, MO | hpc ignites.

Julia for HPC BoF

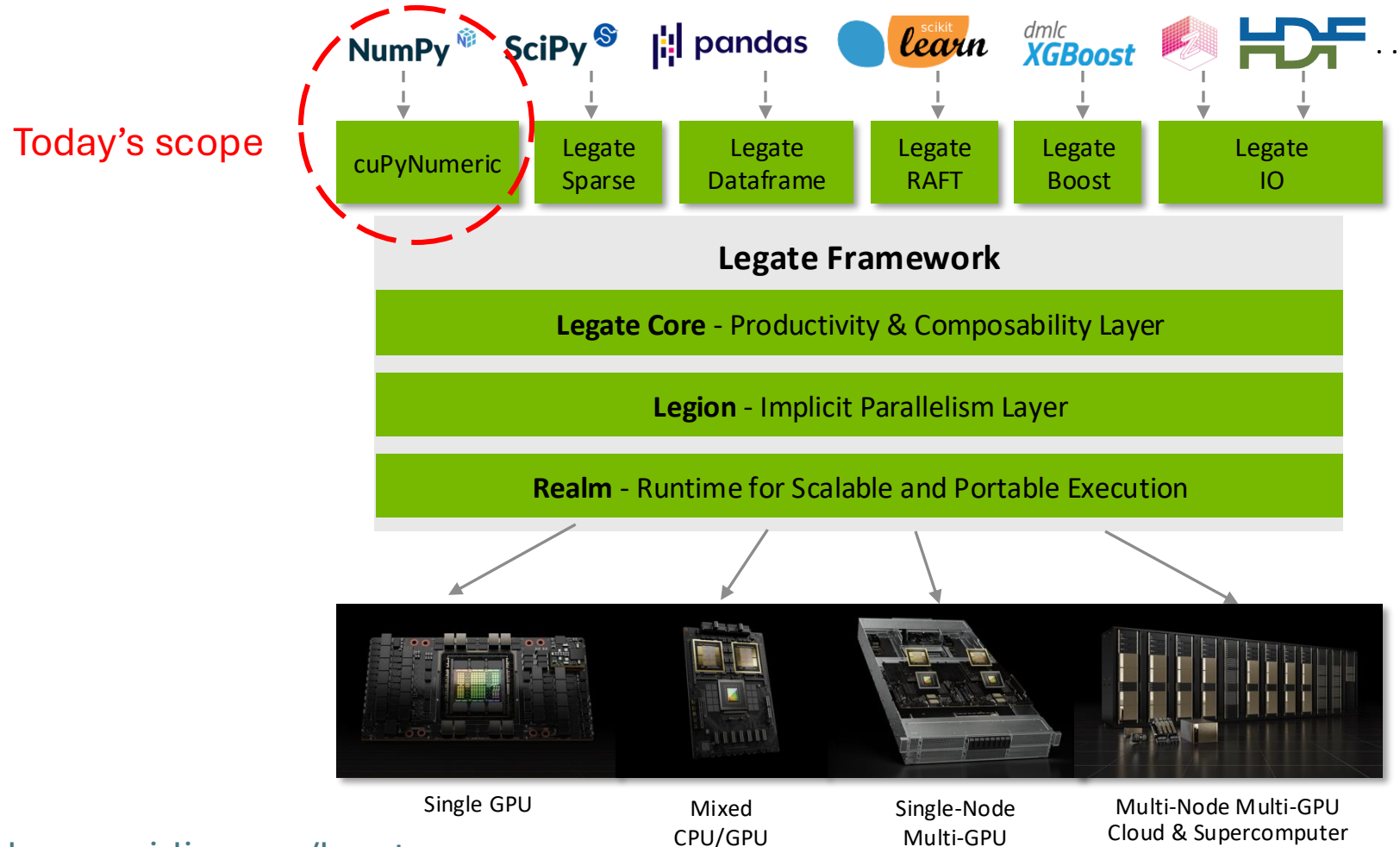# The Goal: Scale with Zero Code Changes

- "Easy" to implement the correct physics in a high-level language like Python, Julia, or MATLAB
- Time consuming to modify code to scale across multiple CPUs/GPUs
  - Need to learn and debug new technologies like MPI, CUDA etc.

**Code that runs on a single CPU core should also be able to run on multiple cores, multiple GPUs and across multiple nodes**
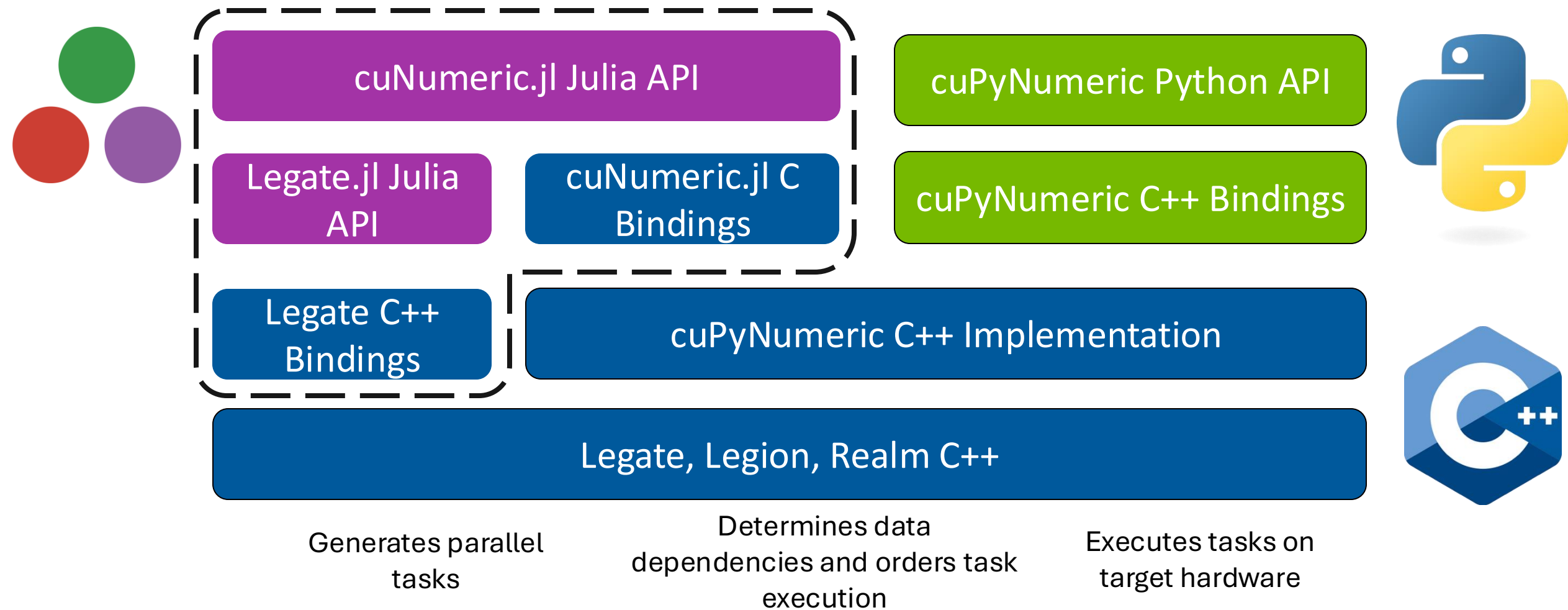
# Legate Enables Composable and Distributed Libraries

By providing data and task management abstractions, this enables the efficient implementation of complex library APIs.



Today's scope

| cuPyNumeric | Legate Sparse | Legate Dataframe | Legate RAFT | Legate Boost | Legate IO |

**Legate Framework**

**Legate Core** - Productivity & Composability Layer

**Legion** - Implicit Parallelism Layer

**Realm** - Runtime for Scalable and Portable Execution

Single GPU    Mixed CPU/GPU    Single-Node Multi-GPU    Multi-Node Multi-GPU Cloud & Supercomputer

**Wouldn't this be great in Julia?**

https://developer.nvidia.com/legate
https://legion.stanford.edu/

# Software Stack

cuNumeric.jl Julia API

cuPyNumeric Python API

Legate.jl Julia API

cuNumeric.jl C Bindings

cuPyNumeric C++ Bindings

Legate C++ Bindings

cuPyNumeric C++ Implementation

Legate, Legion, Realm C++

Generates parallel tasks

Determines data dependencies and orders task execution
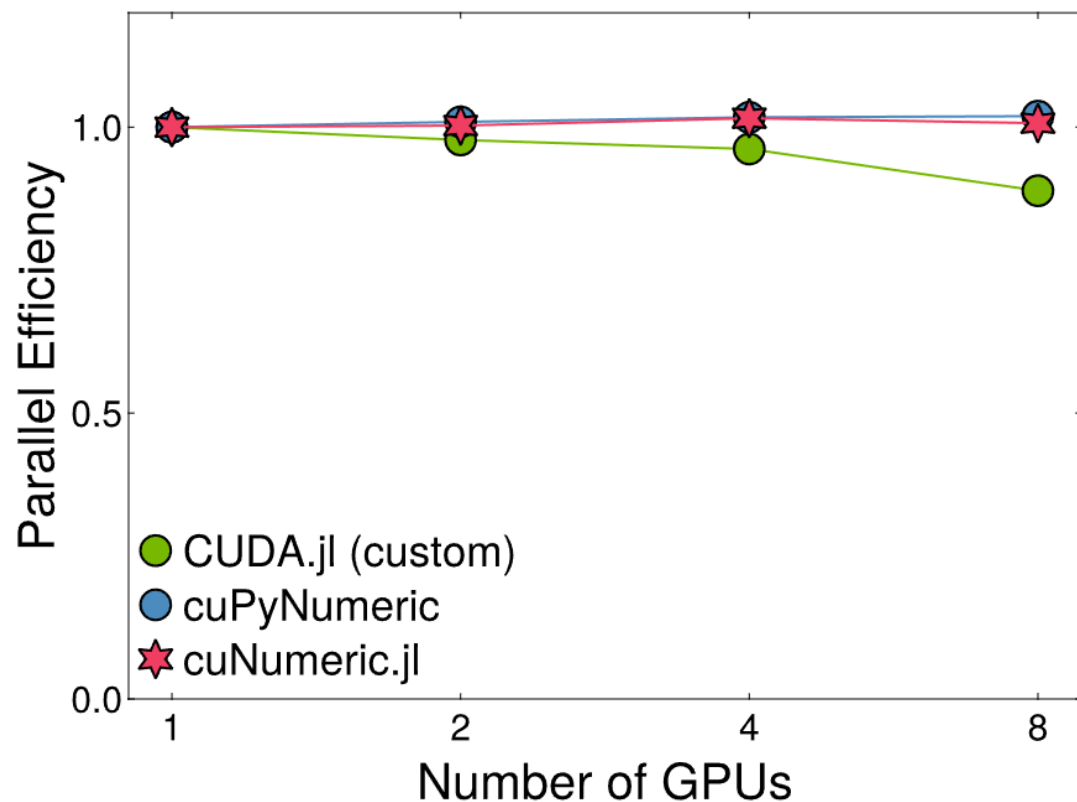
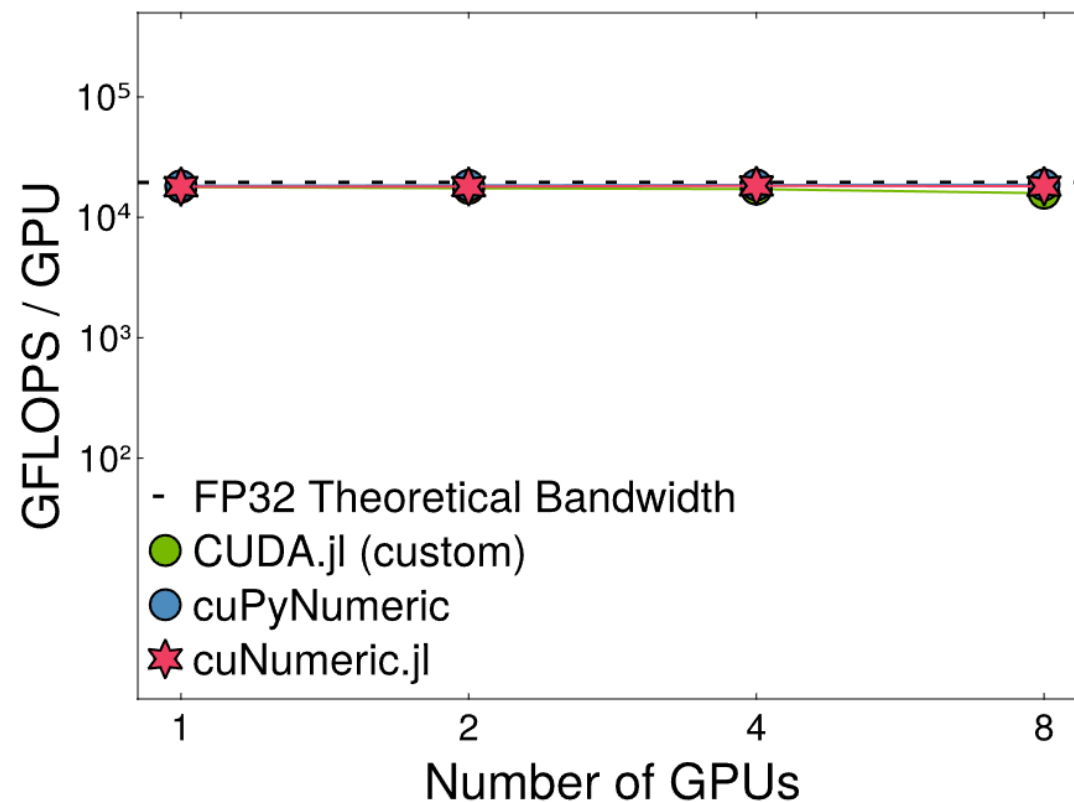Executes tasks on target hardware

New Contributions

6

# Matrix Multiplication

## Benchmark (8x A100):

```
1    N = 10
2    A = cuNumeric.rand(Float32, N, N)
3    B = cuNumeric.rand(Float32, N, N)
4    C = cuNumeric.zeros(Float32, N, N)
5    mul!(C, A, B)
```

# Running CUDA Kernels with cuNumeric.jl

```julia
1   using cuNumeric
2   using CUDA
3
4   function kernel_add(a, b, c, N)
5       i = (blockIdx().x - 1i32) * blockDim().x + threadIdx().x
6       if i <= N
7           @inbounds c[i] = a[i] + b[i]
8       end
9       return nothing
10  end
11
12
13  N = 1024
14  threads = 256
15  blocks = cld(N, threads)
16
17  a = cuNumeric.full(N, 1.0f0)
18  b = cuNumeric.full(N, 2.0f0)
19  c = cuNumeric.ones(Float32, N)
20
21  task = cuNumeric.@cuda_task kernel_add(a, b, c, UInt32(1))
22
23  cuNumeric.@launch task=task threads=threads blocks=blocks \
24                 inputs=(a, b) outputs=c scalars=UInt32(N)
25
```
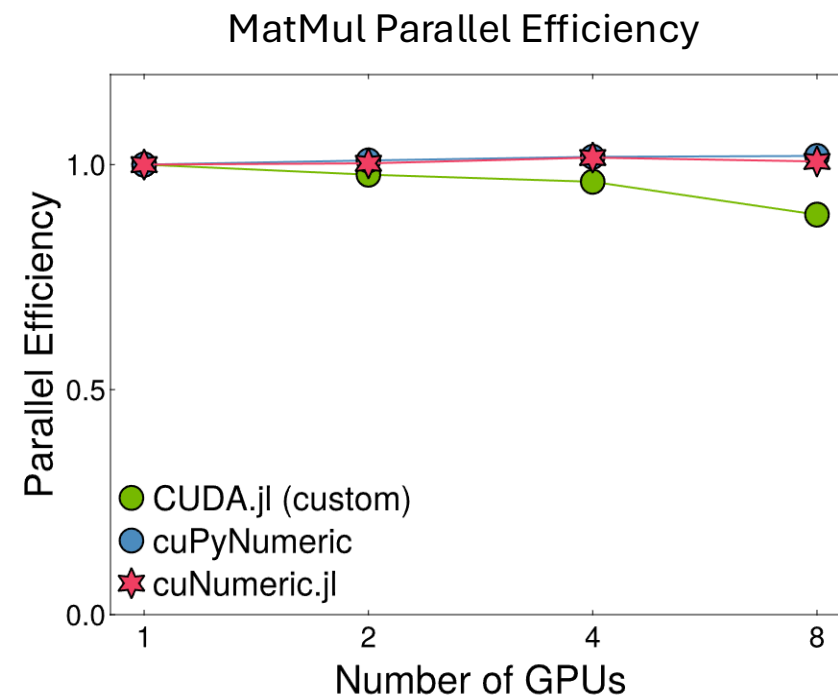
CUDA.jl
kernel

Initialize
NDArrays

Compile

Launch

# Conclusions and Future Work

- Minimal code changes for scaling across large heterogeneous distributed systems
- Good weak scaling efficiency
- Ability to register custom CUDA kernels

**Disclaimers:**

1. Still in development, we will be releasing a public beta within the coming months

2. cuPyNumeric and Legate are currently in beta


MatMul Parallel Efficiency

**David Krasowska** [krasow@u.northwestern.edu](mailto:krasow@u.northwestern.edu)

Ethan Meitz [emeitz@andrew.cmu.edu](mailto:emeitz@andrew.cmu.edu)

Wonchan Lee [wonchanl@nvidia.com](mailto:wonchanl@nvidia.com)

Check out our repo
[https://github.com/JuliaLegate/cuNumeric.jl/](https://github.com/JuliaLegate/cuNumeric.jl/)