# Reducing Complexity in Scalable Numerical Computing with cuNumeric.jl

**David Krasowska**   Second Year Fellow

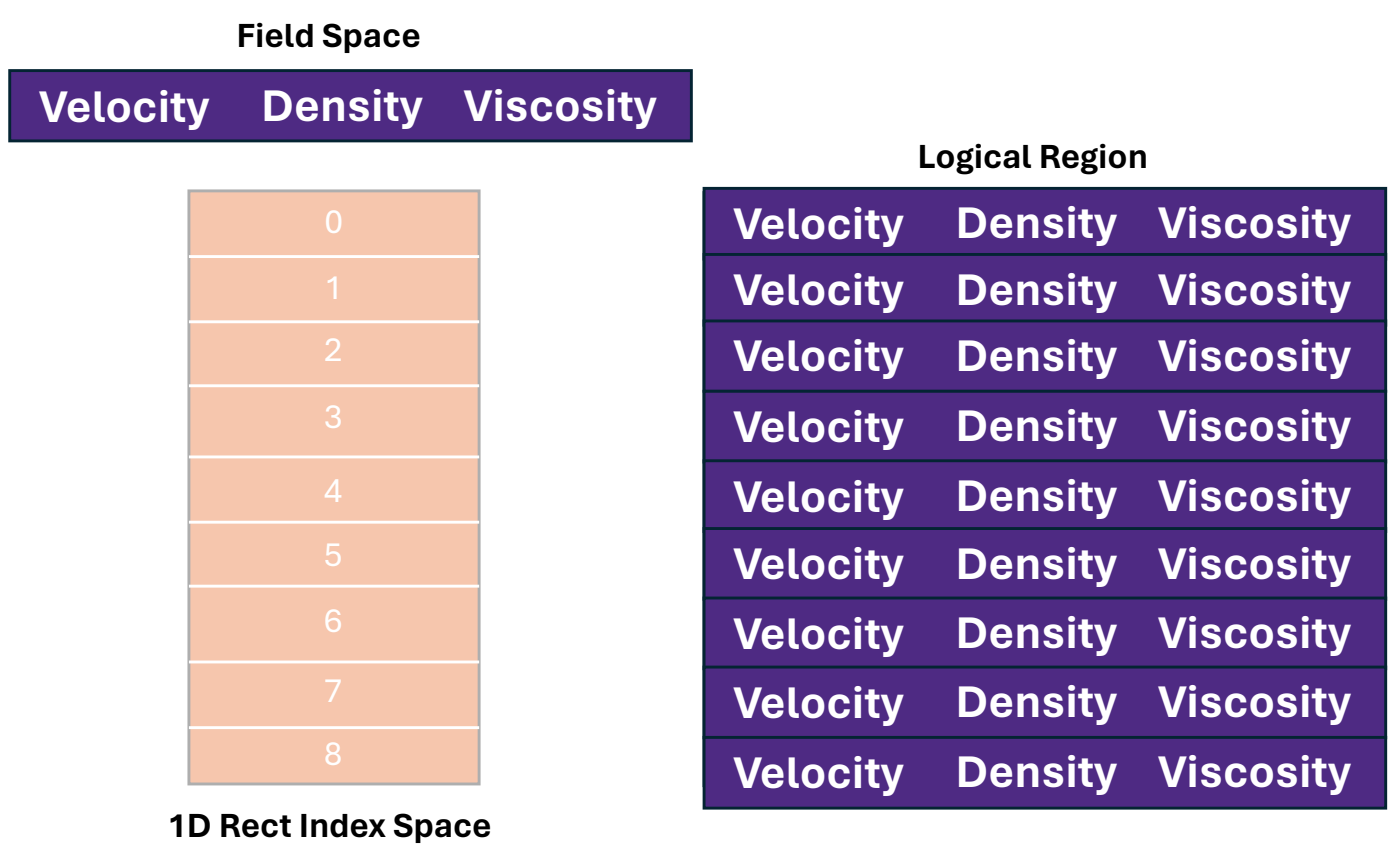Northwestern | Plab | CONSTELLATION PROJECT

## Introduction

With the end of Moore's Law, demand for specialized accelerators has surged to meet modern performance needs. Scientific computing often relies on large high-performance computing (HPC) clusters, which increasingly include many-core and heterogeneous systems. However, programming such systems introduces complexity in **memory management**, **parallelism**, and **performance portability**.

**Claim:** cuNumeric.jl will allow scientific developers to harness the power of GPU-accelerated HPC systems with minimal effort and maximum portability.
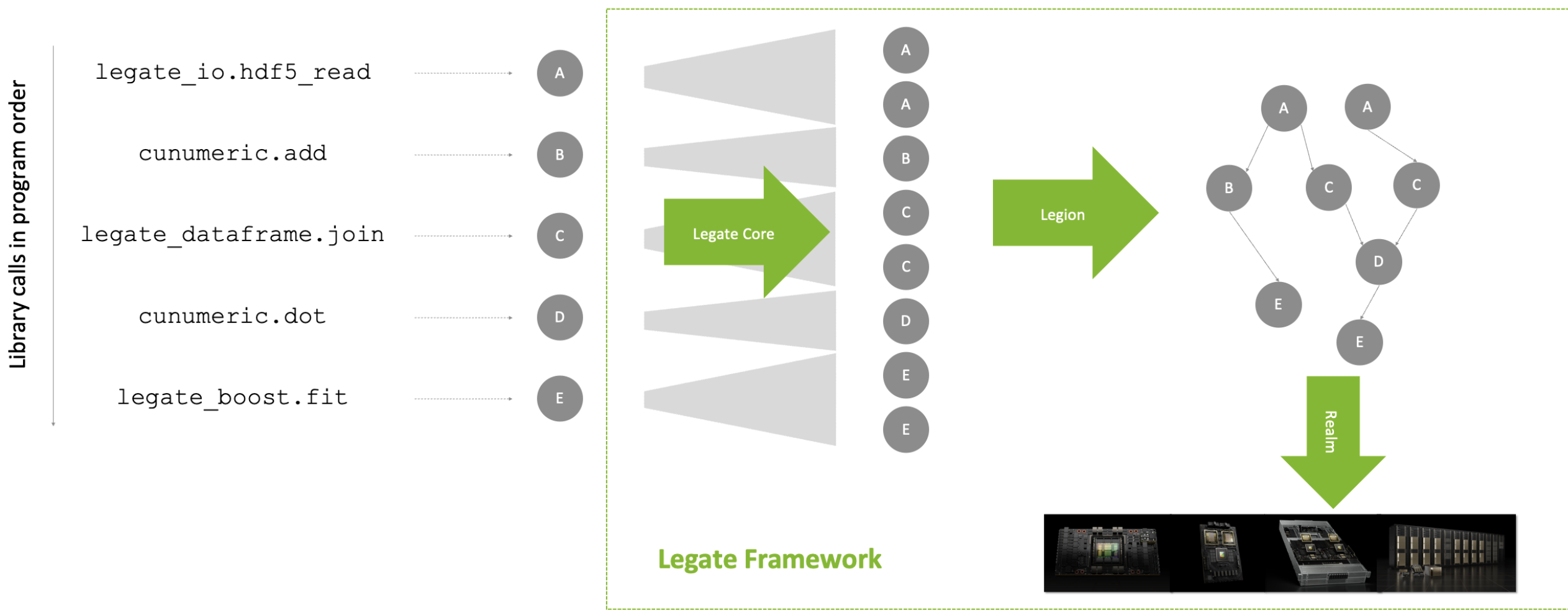
## Legion

Achieving high performance requires the programming system understand the structure of program data to facilitate efficient placement and movement of data.

**Legion** [2]: runtime system and program model for distributed heterogeneous machines. Application data is structured by:



- Implicitly extract parallelism from the logical regions
- Data movement operations in accordance with the application data properties
- All tasks specify privileges and coherence that logical regions will access

Legion utilizes large task graphs to minimize kernel invocation and communication latencies.
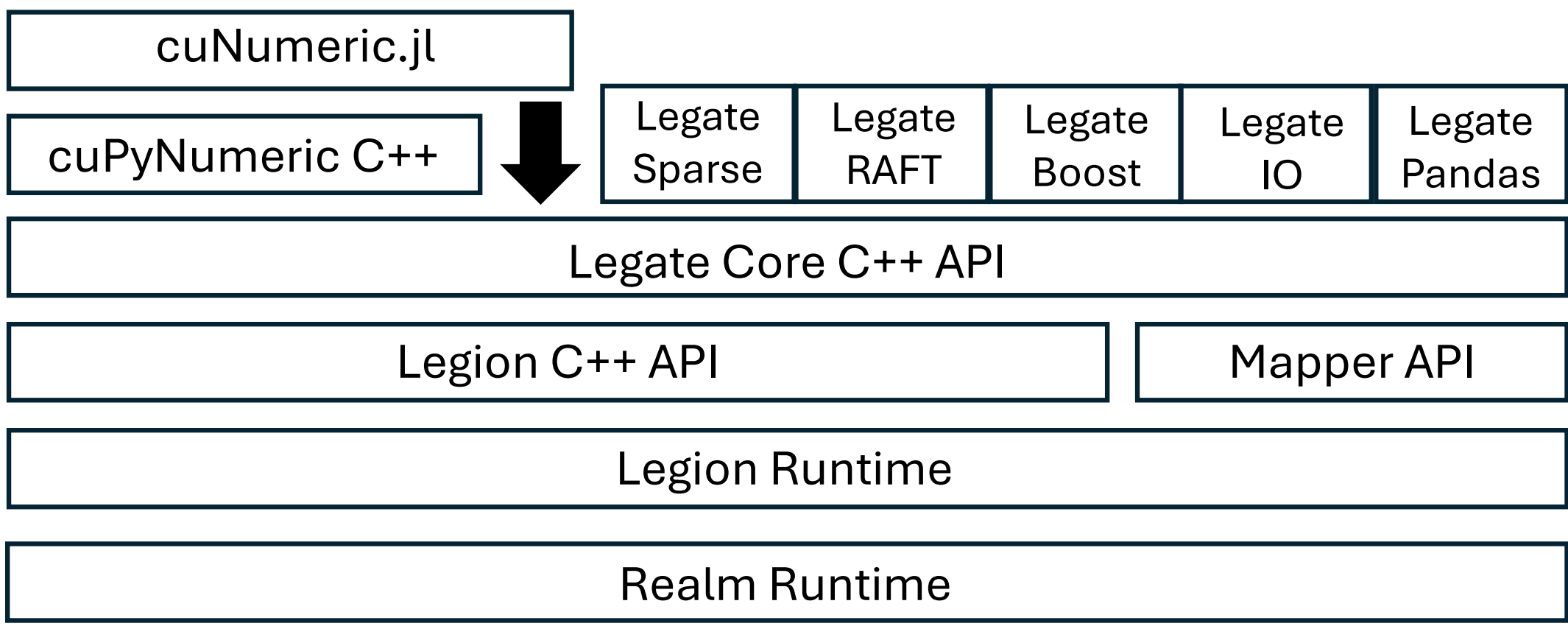
## Legate and cuPyNumeric

**Legate** [1]: Provides an abstraction over the Legion runtime through a unified library interface, enabling efficient implmenetation of complex APIs that scale. cuPyNumeric is a Legate library that provides a Python Numpy API replacement for scalable execution.

**Motivation:** "Code that runs on a single CPU should also be able to run on multiple cores, multiple GPUs and across multiple nodes"



## cuNumeric.jl

Leverages existing infrastructure for cuPyNumeric to bring seamless distributed computing to Julia.

- Legion is used to schedule Julia tasks, enabling implicit parallelism
- Matches high level array abstractions within Julia
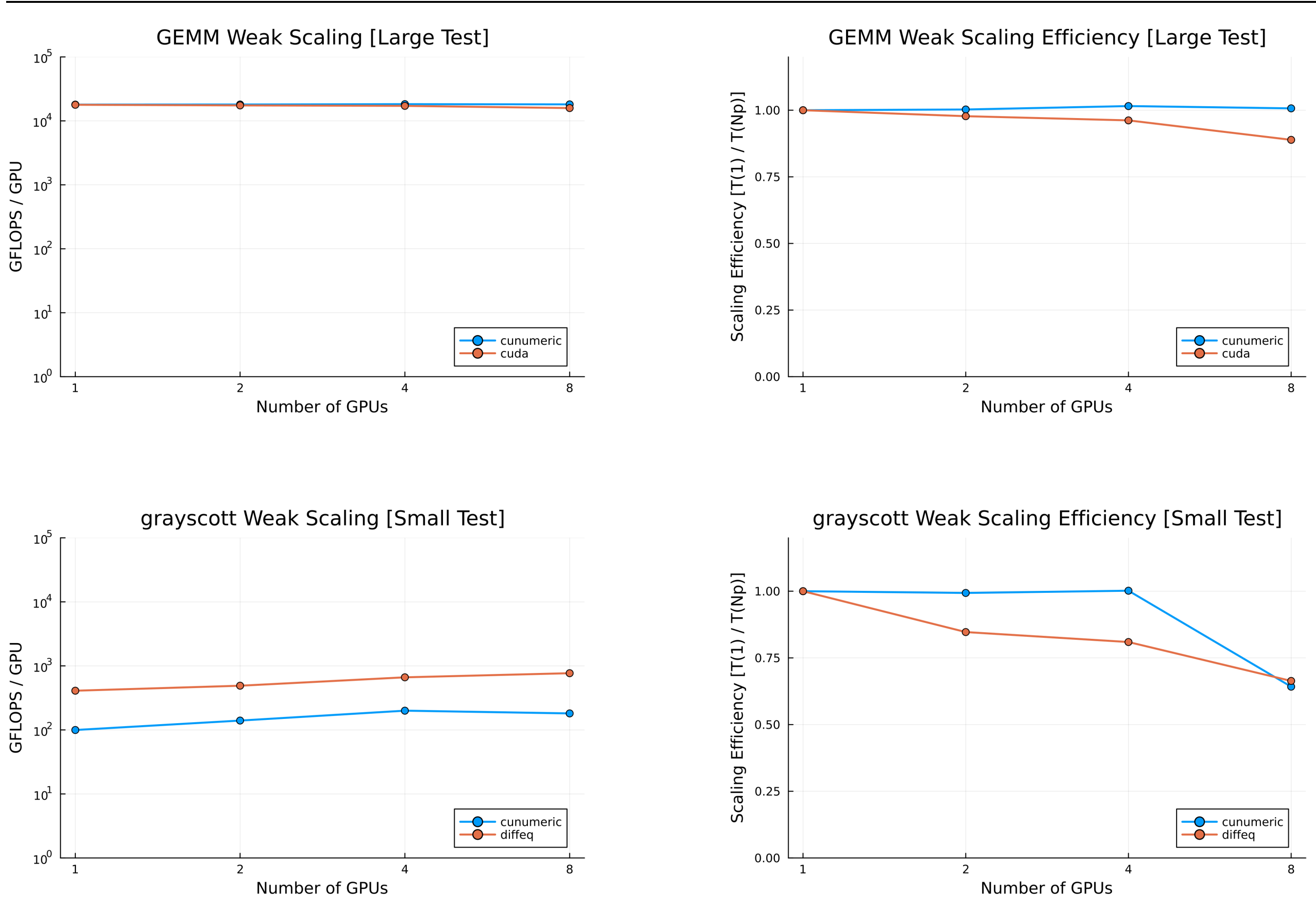- Behavior a Julia programmer would "expect"



## API examples

```julia
1  # found in examples/integrate.jl
2  using cuNumeric
3
4  integrand = (x) -> exp(-square(x))
5
6  N = 1_000_000
7
8  x_max = 5.0
9  domain = [-x_max, x_max]
10 Ω = domain[2] - domain[1]
11
12 samples = Ω*cuNumeric.rand(NDArray, N) - x_max
13 estimate = (Ω/N) * sum(integrand(samples))
```

```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
                    inputs=(a, b) outputs=c scalars=UInt32(N)
```

More details are found on our repo:
github.com/JuliaLegate/cuNumeric.jl

## Results



These results are provided using [1:8] NVIDIA A100 GPUs.

1. Good weak scaling on GEMM. Outperforms CUDA.jl multiGPU execution
2. ImplicitGlobalGrid.jl performs better on Gray-Scott 2D heat diffusion, particularly in GFLOPS

## Conclusions and Future Work

- Minimal code changes that allow users to scale code across large heterogeneous distributed systems
- Good weak scaling efficiency on a variety of applications
- Ability to register custom CUDA kernels

As this is a work in progress, we have next steps

- Support a wider range of custom CUDA kernels
- Robust and accessible package installation
- More integration with the Julia Abstraction interface + Legate

## Acknowledgments

[1] Michael Bauer and Michael Garland. Legate numpy: accelerated and distributed array computing. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, SC '19, New York, NY, USA, 2019. Association for Computing Machinery.

[2] Michael Bauer, Sean Treichler, Elliott Slaughter, and Alex Aiken. Legion: Expressing locality and independence with logical regions. In *SC '12: Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, pages 1–11, 2012.

DOE CSGF