

cuNumeric.jl : Automating Distributed Numerical Computing

Ethan Meitz¹, **David Krasowska**², Wonchan Lee³, Pat McCormick⁴

¹Carnegie Mellon University

²Northwestern University

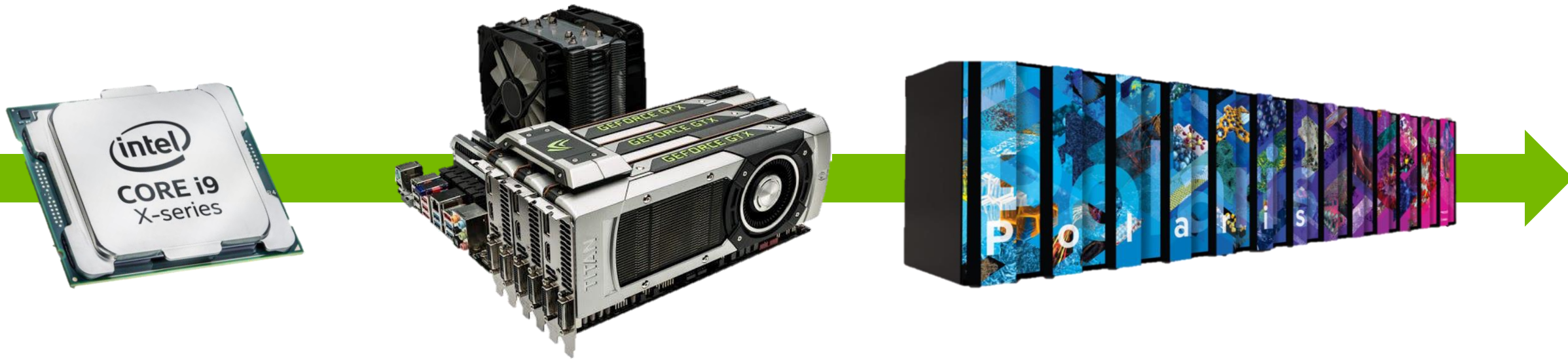
³NVIDIA

⁴Los Alamos National Laboratory



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 OpenMPI, 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

cuNumeric.jl: Distributed Code with Minimal Effort

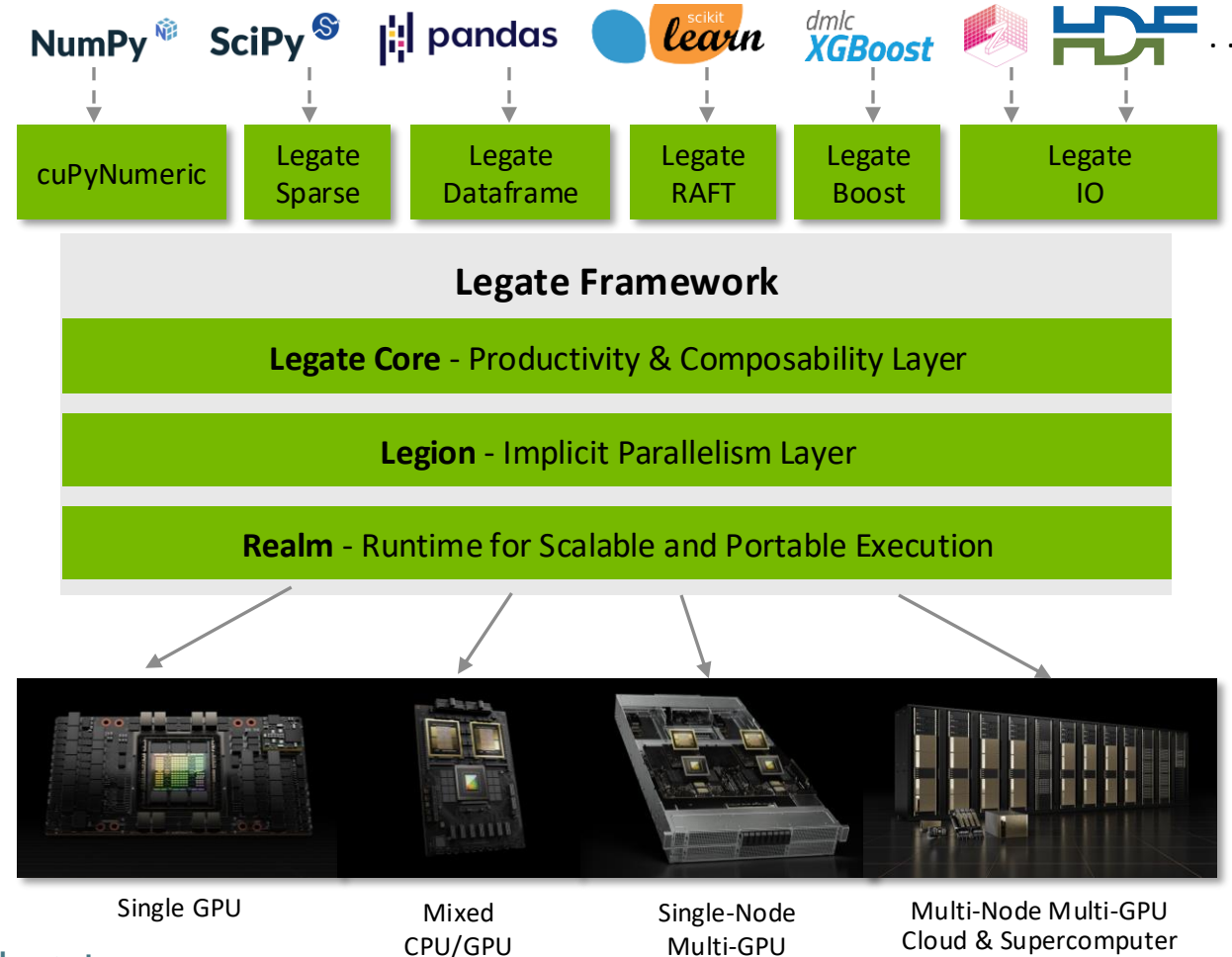
- Core Type: **NDArray**
 - Drop-in replacement for Base.Array but can represent data across multiple devices
 - Key Differences: Operations are broadcast by default, slices are always views, avoid scalar indexing

```
(3.8.0) (base) david@dubliner:~/julia-con/slides-examples/diff$ cat diff_grayscott.txt
1,2d0
< using cuNumeric
<
81,84c79,82
<     u = cuNumeric.ones(dims)
<     v = cuNumeric.zeros(dims)
<     u_new = cuNumeric.zeros(dims)
<     v_new = cuNumeric.zeros(dims)
---
>     u = ones(dims)
>     v = zeros(dims)
>     u_new = zeros(dims)
>     v_new = zeros(dims)
86,87c84,85
<     u[1:150, 1:150] = cuNumeric.random(FT, (150, 150))
<     v[1:150, 1:150] = cuNumeric.random(FT, (150, 150))
---
>     u[1:150, 1:150] = random(FT, (150, 150))
>     v[1:150, 1:150] = random(FT, (150, 150))
(3.8.0) (base) david@dubliner:~/julia-con/slides-examples/diff$
```

**7 LOC changed for 2D Gray Scott
Reaction-Diffusion simulation**

Legate Enables Composable and Distributed Libraries

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



Wouldn't this be great in Julia?

How Legate Works

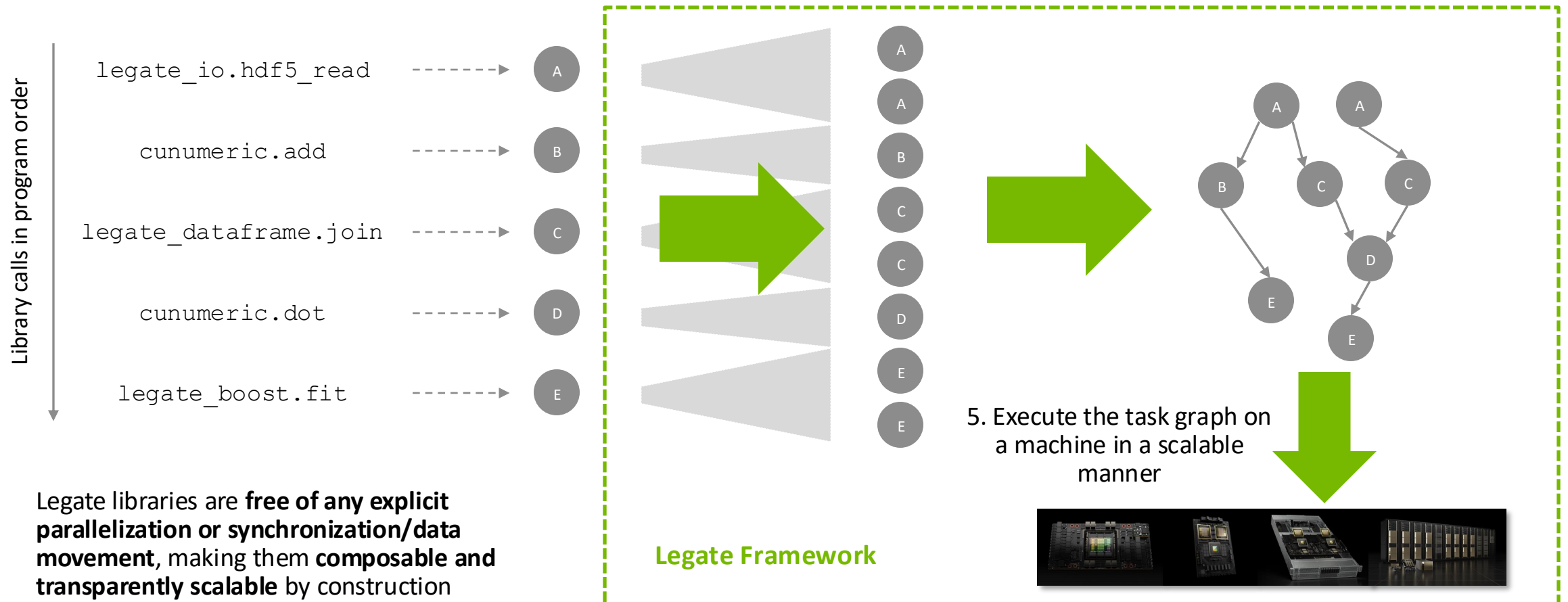
Implicit parallelism via “scale-free” tasking

1. Legate program makes API calls

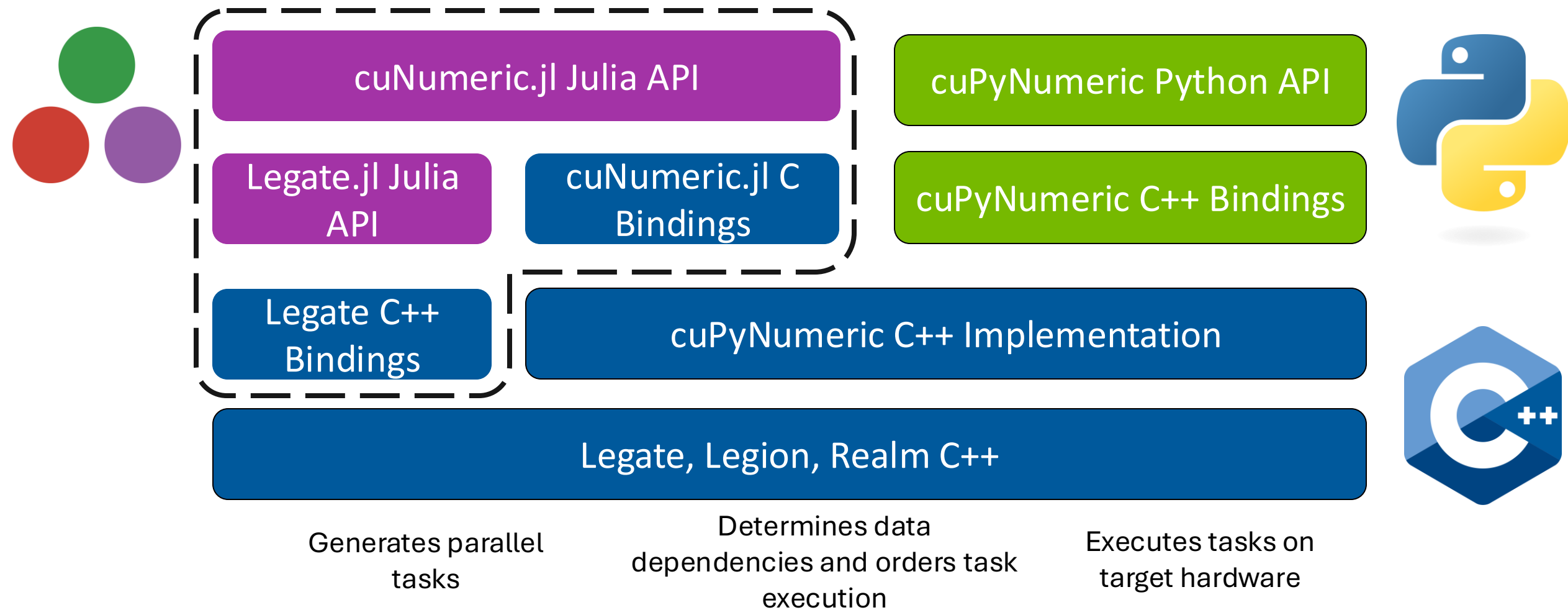
2. Legate libraries issue “scale-free” tasks

3. Convert each scale-free task to parallel tasks

4. Analyze data dependencies and constructs a task graph



Software Stack

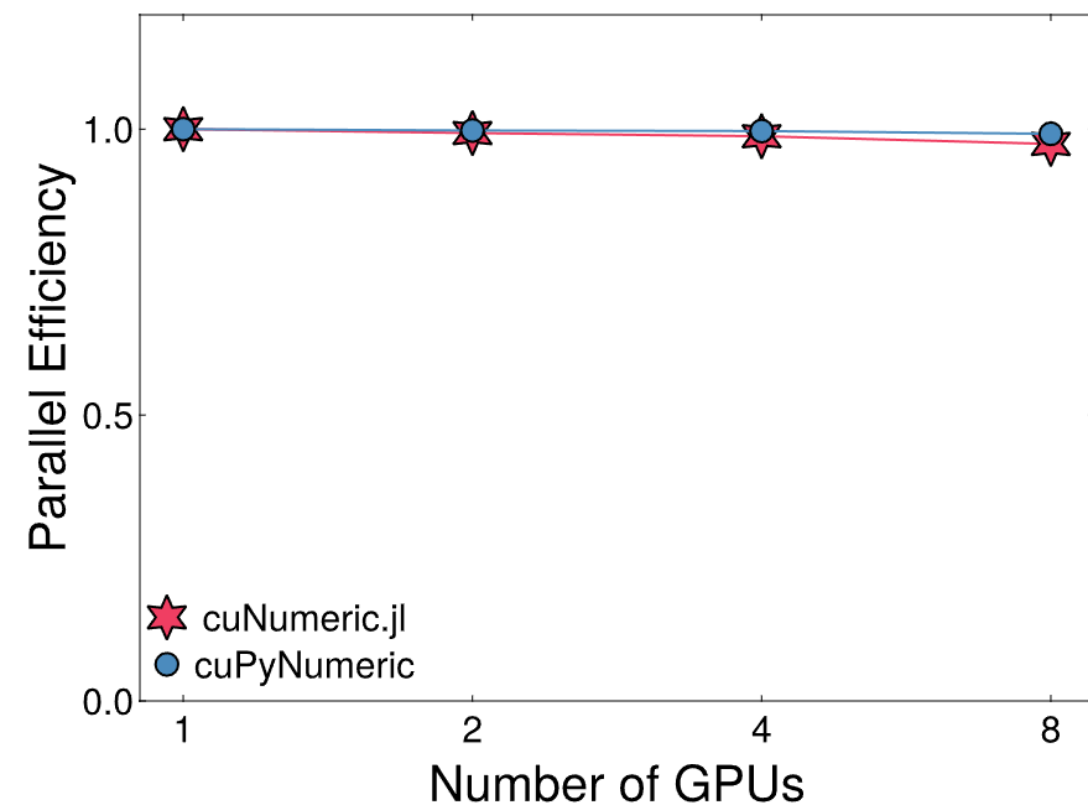


New Contributions

Monte Carlo Integration

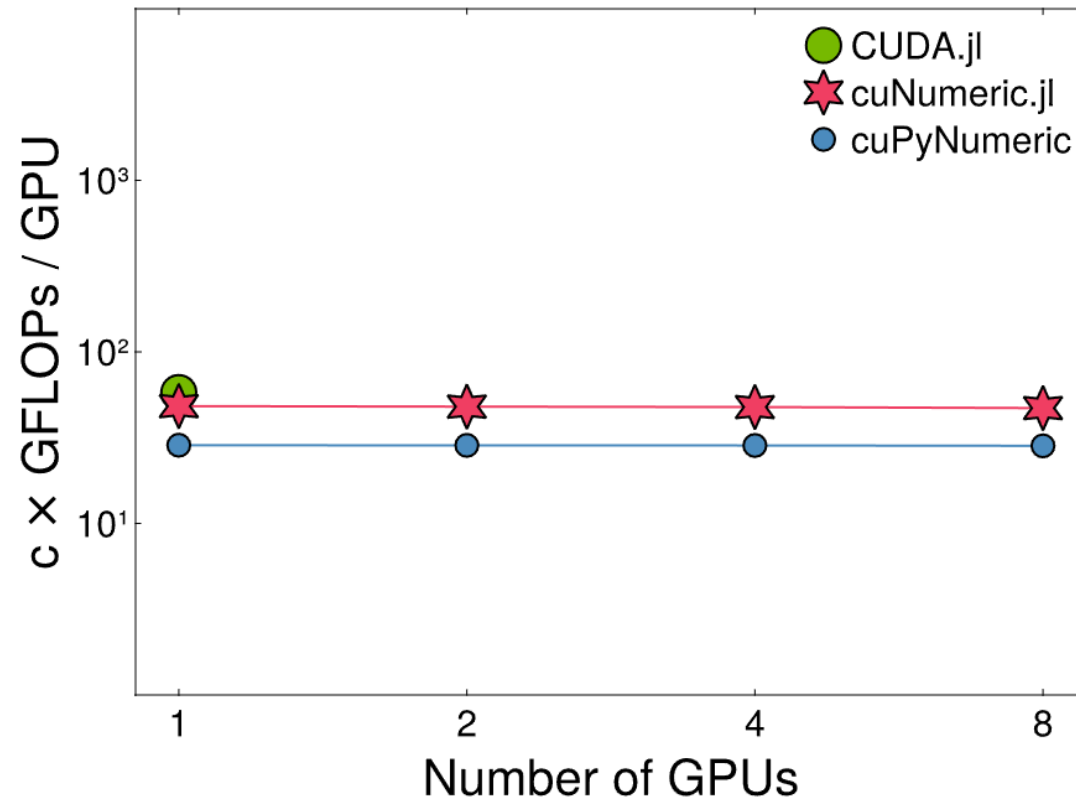
- Embarrassingly Parallel → Should scale perfectly
- cuNumeric.jl directly calls C++, cuPyNumeric passes through several layers of Python before C++

Benchmark (8x A100):



Syntax:

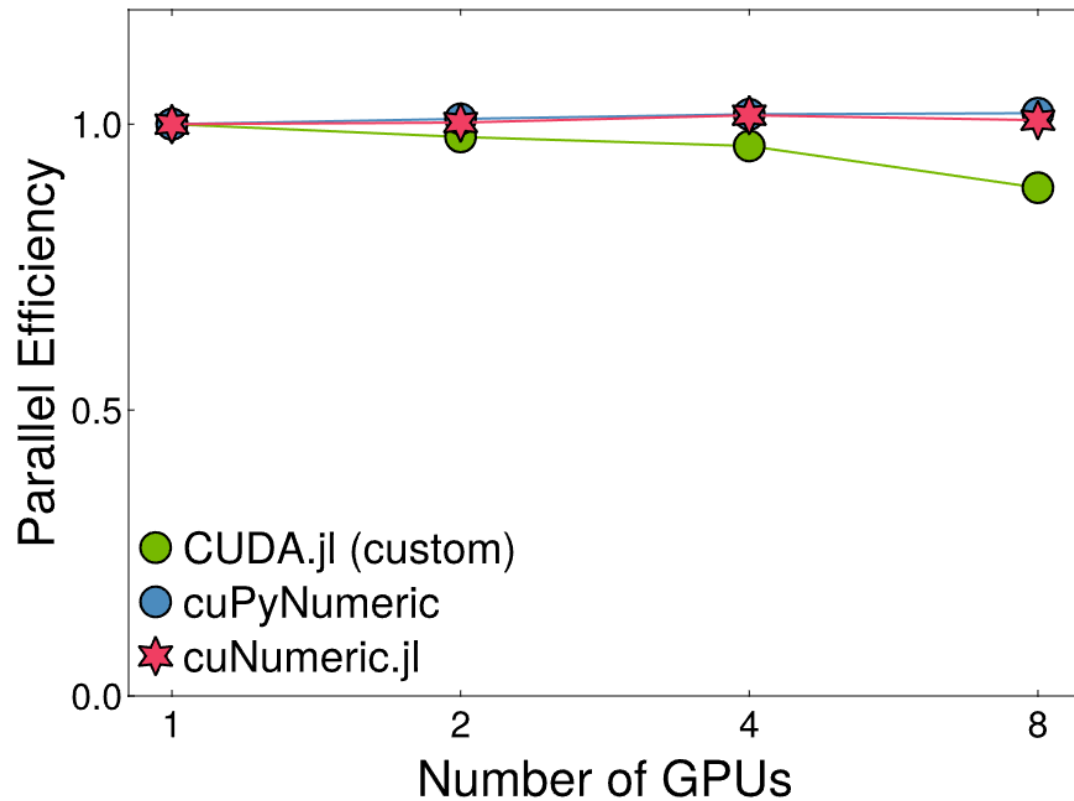
```
1 integrand = (x) -> exp(-square(x))
2 N = 1_000_000
3 x_max = 5.0
4 domain = [-x_max, x_max]
5 Ω = domain[2] - domain[1]
6 samples = Ω*cuNumeric.rand(NDArray, N) - x_max
7 estimate = (Ω/N) * sum(integrand(samples))
```



*Actual FLOPs unknown, off by constant factor ($c < 1.0$)

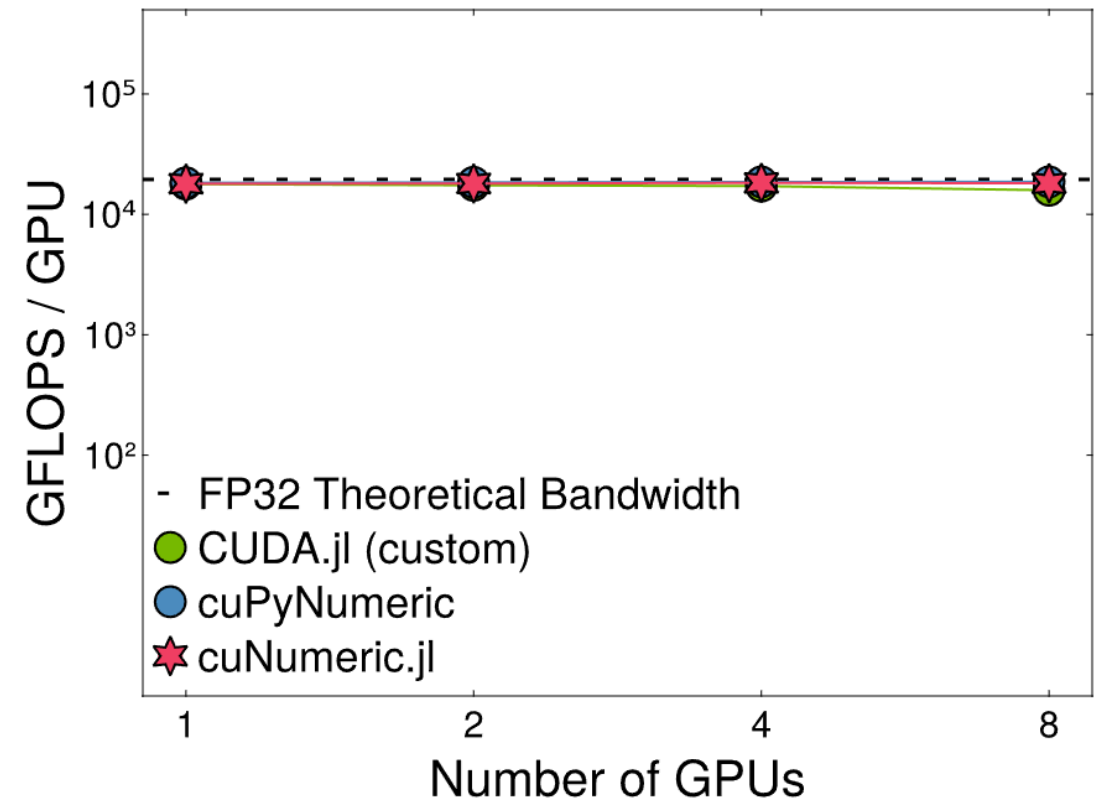
Matrix Multiplication

Benchmark (8x A100):



Syntax*:

```
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)
```



*Float32 random number generation not supported yet, requires type cast

Gray Scott Reaction Diffusion (2D)

Syntax:

```
1  N = 100
2  u = cuNumeric.ones((N, N))
3  v = cuNumeric.ones((N, N))
4  for i in iters
5      F_u = ((u[2:(end - 1), 2:(end - 1)] .* (v[2:(end - 1), 2:(end - 1)]
6              .* v[2:(end - 1), 2:(end - 1)])) - (f+k)*v[2:(end - 1), 2:(end - 1)]
7      #.... more physics
8      #.... update u and v
9      GC.gc()
10 end
```

Manual GC allows Gray Scott to run

Each object, including slices, are treated as a temporary and remain uncollected.
Why are they not collected before OOM?

Benchmark (8x A100):



Foreign Memory is not automatically garbage collected

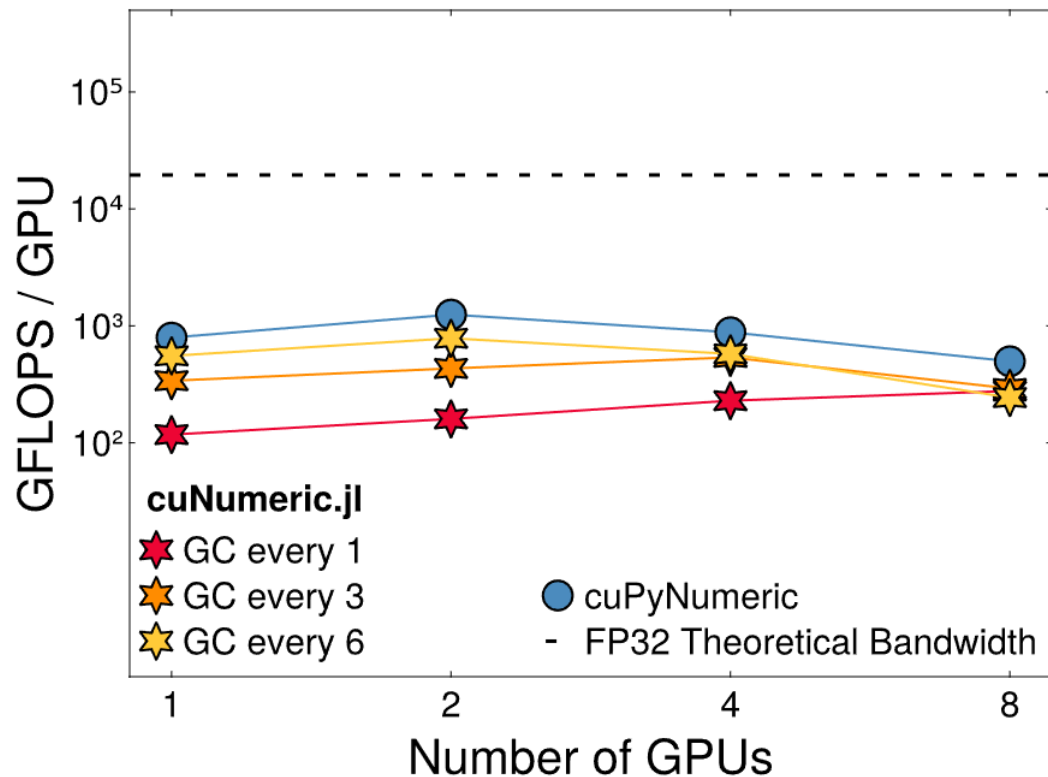
- GC sees foreign memory as a pointer → 8 bytes
- GC relies on the heap size (of Julia objects) to decide when GC is invoked

```
1  using CUDA, Profile
2
3  arr = CUDA.ones(Float32, 10_000_000) ← ~40MB
4  Profile.take_heap_snapshot()
5  Base.summarysize(arr) ← Julia sees GPU data as 184 bytes
```

Constructor	Retained Size
▼ CUDA.CuArray{Float32, 1, CUDA.Devic	0.2 kB
▼ data :: GPUArrays.DataRef{CUDA.Mi	0.1 kB
▼ rc :: GPUArrays.RefCounted{CUD	0.1 kB
▼ obj :: CUDA.Managed{CUDA.Dev	0.1 kB
stream :: CUDA.CuStream @L	0.0 kB
count :: Base.Threads.Atomic	0.0 kB
finalizer :: typeof(CUDA.poc	0.0 kB

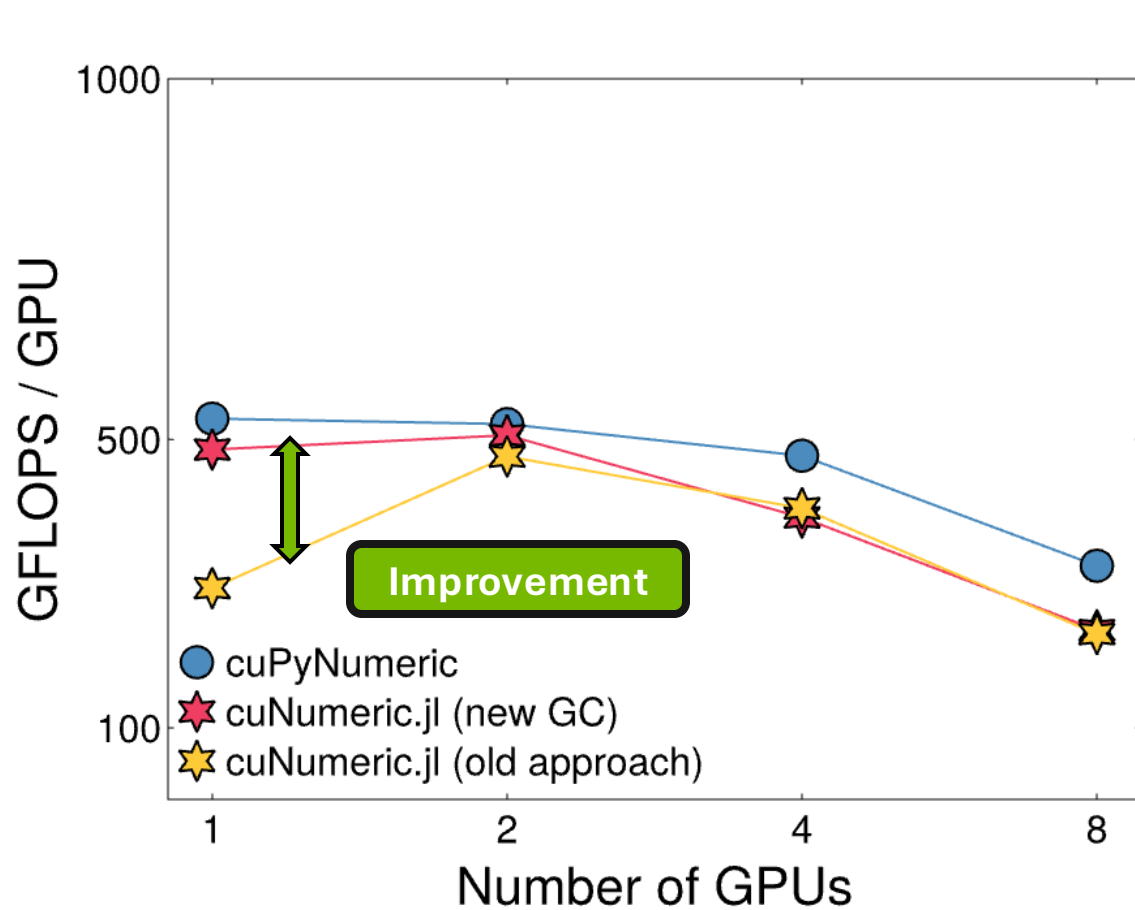
← 184 bytes in heap snapshot too

Garbage Collection heavily impacts performance



- All existing GPU backends in Julia:
 - Calculate array memory footprint in the constructor
 - Track GPU memory allocations
 - Manually calls GC based on a memory pressure heuristic
 - try-catch on allocation to avoid OOM
- NDAarray Properties:
 - Deferred execution model conceals Julia object's physical size at creation
 - Runtime cannot recover from failed mapping allocation

Custom heuristic GC improves performance



```
1 N = 100
2 u = cuNumeric.ones((N, N))
3 v = cuNumeric.ones((N, N))
4 for i in 1:iters
5     F_u = ((u[2:(end - 1), 2:(end - 1)] .* (v[2:(end - 1), 2:(end - 1)]
6           .* v[2:(end - 1), 2:(end - 1)])) - (f+k)*v[2:(end - 1), 2:(end - 1)])
7     #.... more physics
8     #.... update u and v
9 end
```

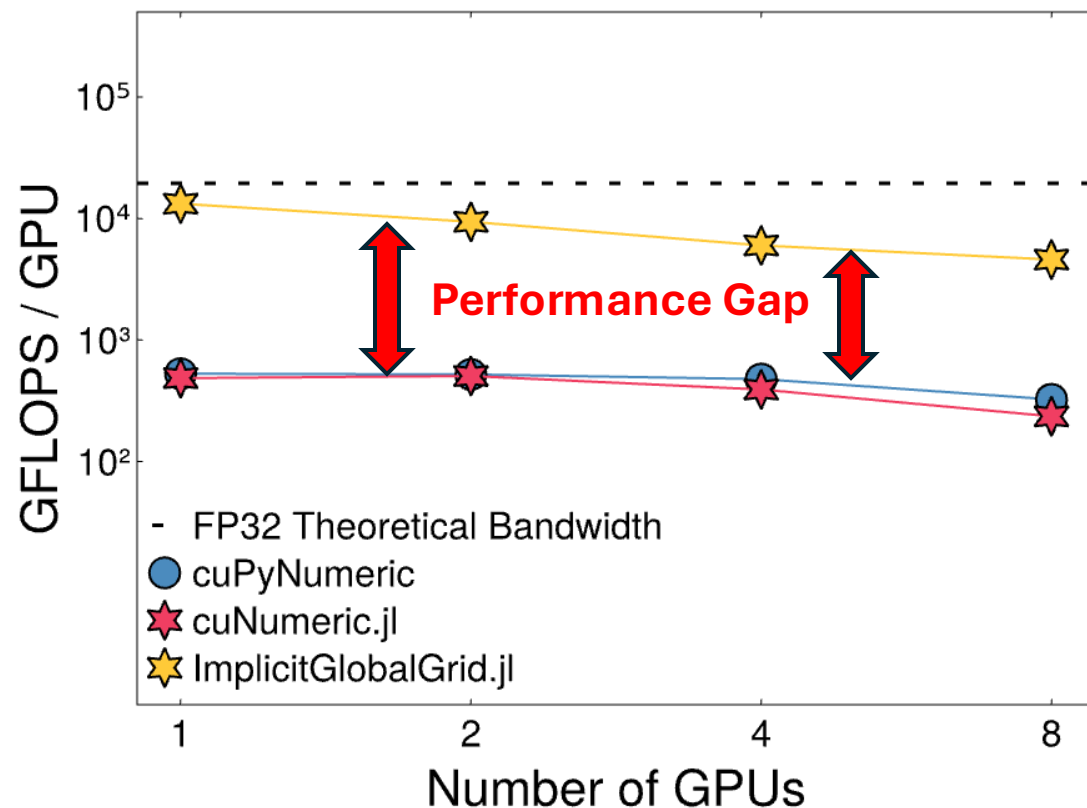
Automatic GC allows Gray Scott to run

- We are still exploring solutions to achieve better scalability

Unfused Operations Bottleneck Performance

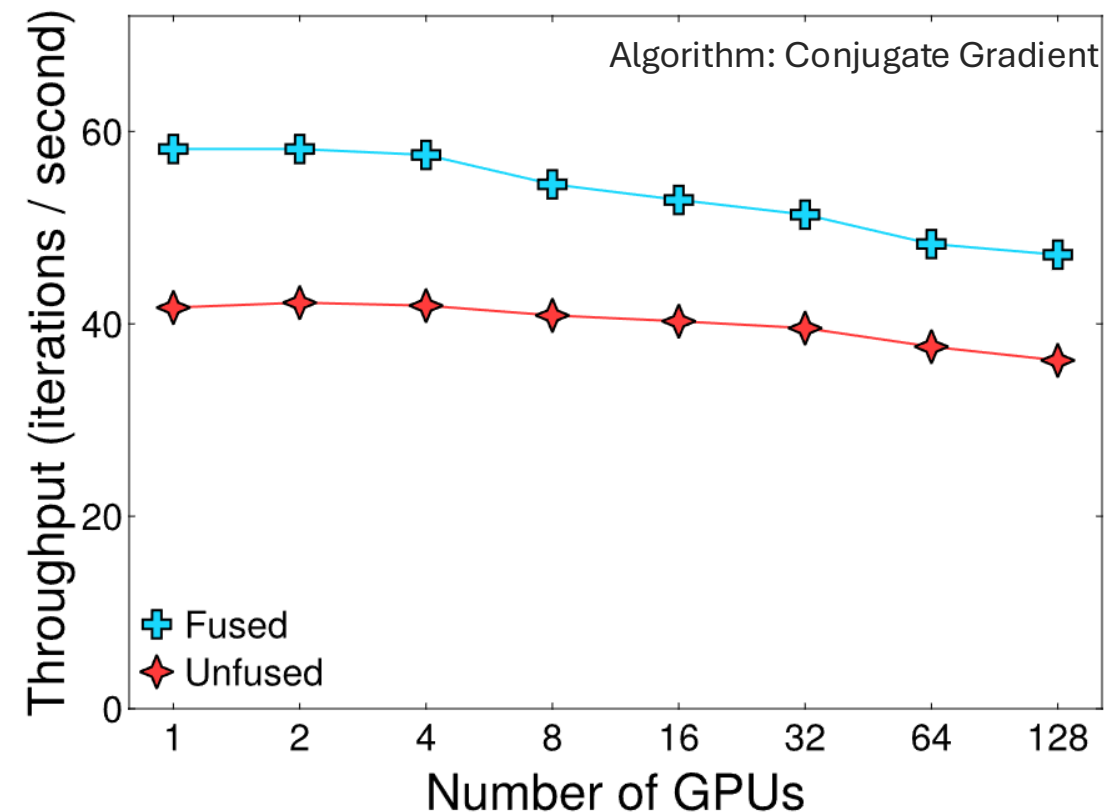
ImplicitGlobalGrid.jl:

- Stencil-based PDE solver
- Multi-node CPU and GPU (NVIDIA & AMD)
- Minimal code changes



Kernel Fusion:

- Under development in cuPyNumeric, but shown to provide 2x speed-up on average and up to 10x
- With CUDA.jl and Legate.jl we propose to generate fused, distributed kernels



Multi-GPU CUDA Kernels with cuNumeric.jl

- To enable kernel fusion, we need to be able to run CUDA.jl kernels through Legate.jl

```
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

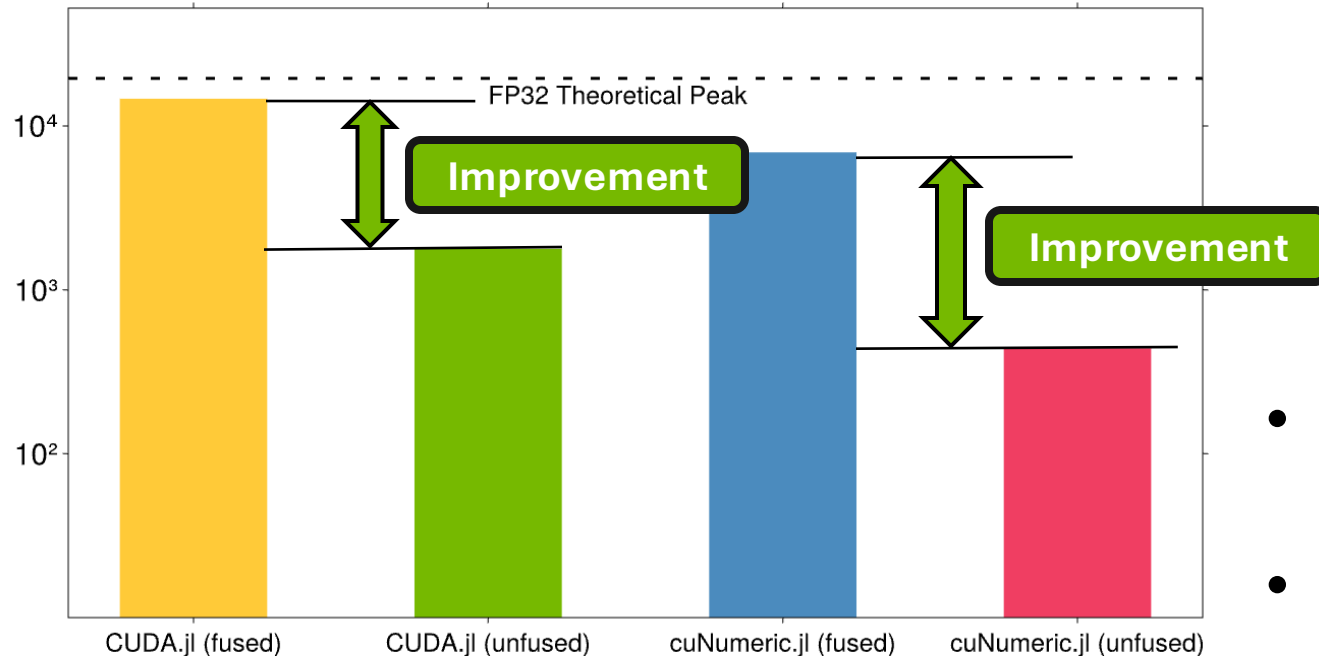
Initialize
NDArrays

Compile

Launch

Gray Scott mini 1D with custom CUDA.jl kernels

Benchmark (1x A30x):



Syntax:

```
1 function fused_kernel(u, v, F_u, F_v, N::UInt32, f::Float32, k::Float32)
2     i = (blockIdx().x - 1i32) * blockDim().x + threadIdx().x
3     if i <= (N*N-2)
4         @inbounds begin
5             u_ij = u[i + 1]
6             v_ij = v[i + 1]
7             v_sq = v_ij * v_ij
8             F_u[i] = (-u_ij * v_sq) + f*(1.0f0 - u_ij)
9             F_v[i] = (u_ij * v_sq) - (f + k)*v_ij
10        end
11    end
12
13    return nothing
14 end
```

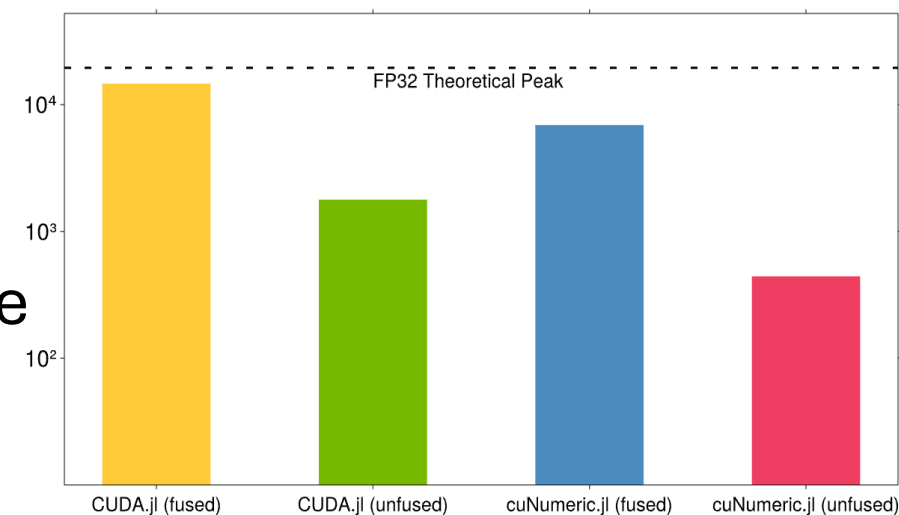
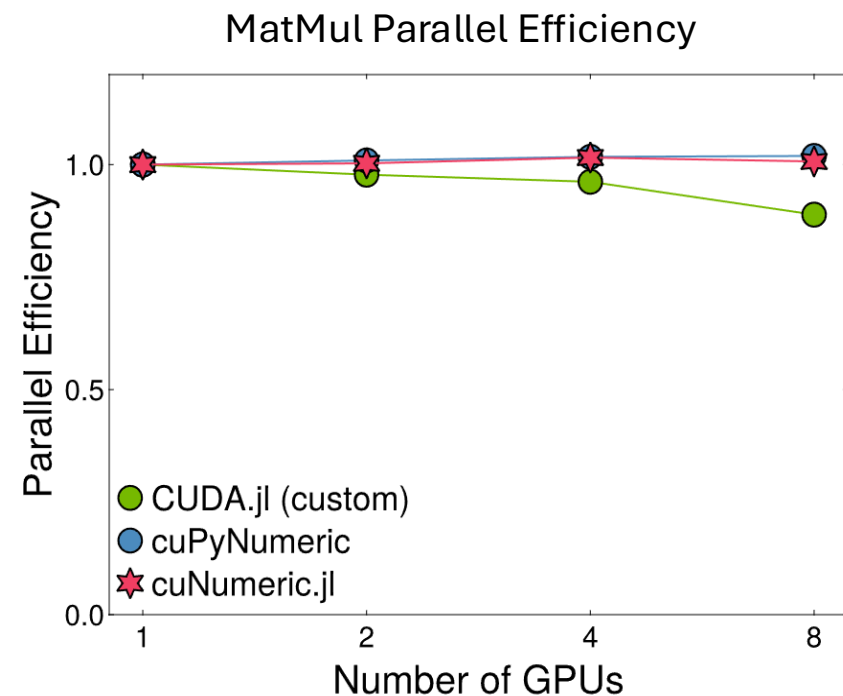
- Operator invocation overhead limits performance scalability
- Ongoing backend enhancements for multi-dimensional workloads and multi-GPU execution
- Backend overhead is amortized as GPU count increases

Conclusions and Future Work

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

Next Steps

1. Support a wider range of custom CUDA kernels
2. Improve robustness and accessibility of package installation
3. Enhance integration with Julia Abstraction interface + Legate
4. Benchmark on multi-node systems
5. Better GC heuristics



Aiming for September* beta release

Ethan Meitz emeitz@andrew.cmu.edu
David Krasowska krasow@u.northwestern.edu
Wonchan Lee wonchanl@nvidia.com
Pat McCormick pat@lanl.gov



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

*this is an estimate; we will register the package upon beta launch



NVIDIA

