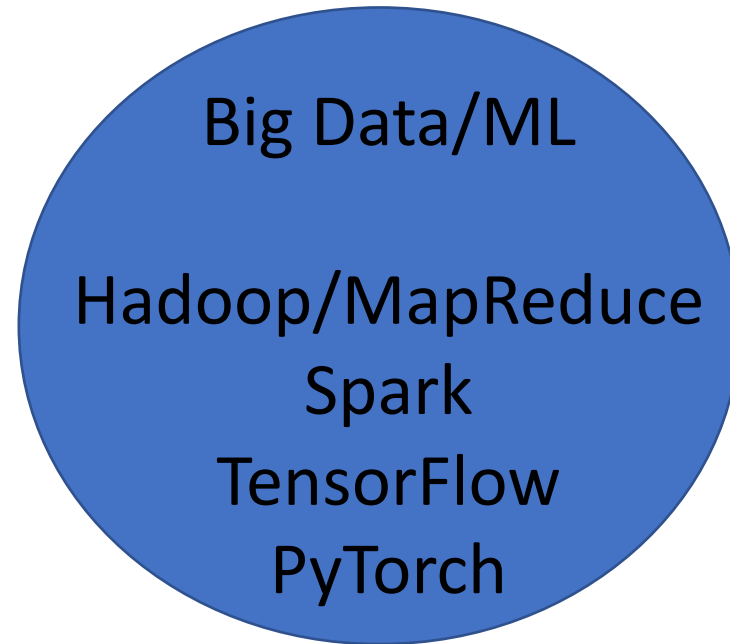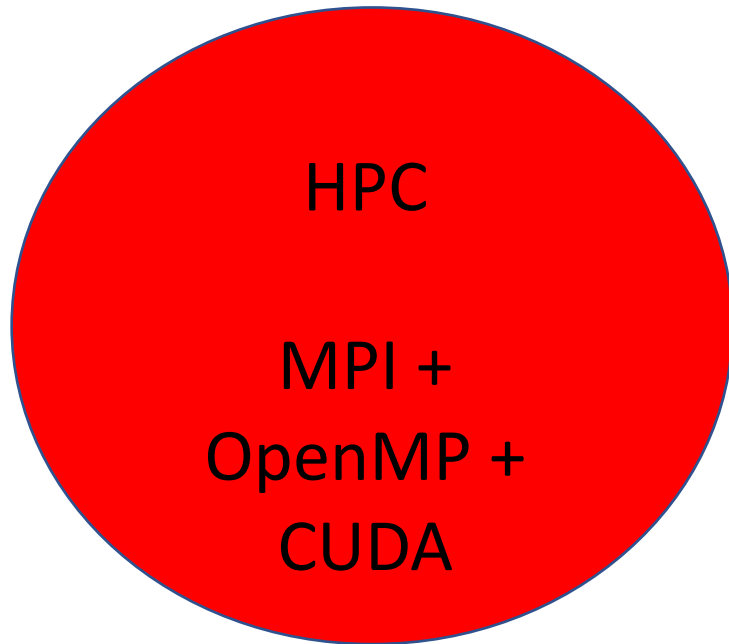# A Tale of Two Cultures

## Alex Aiken
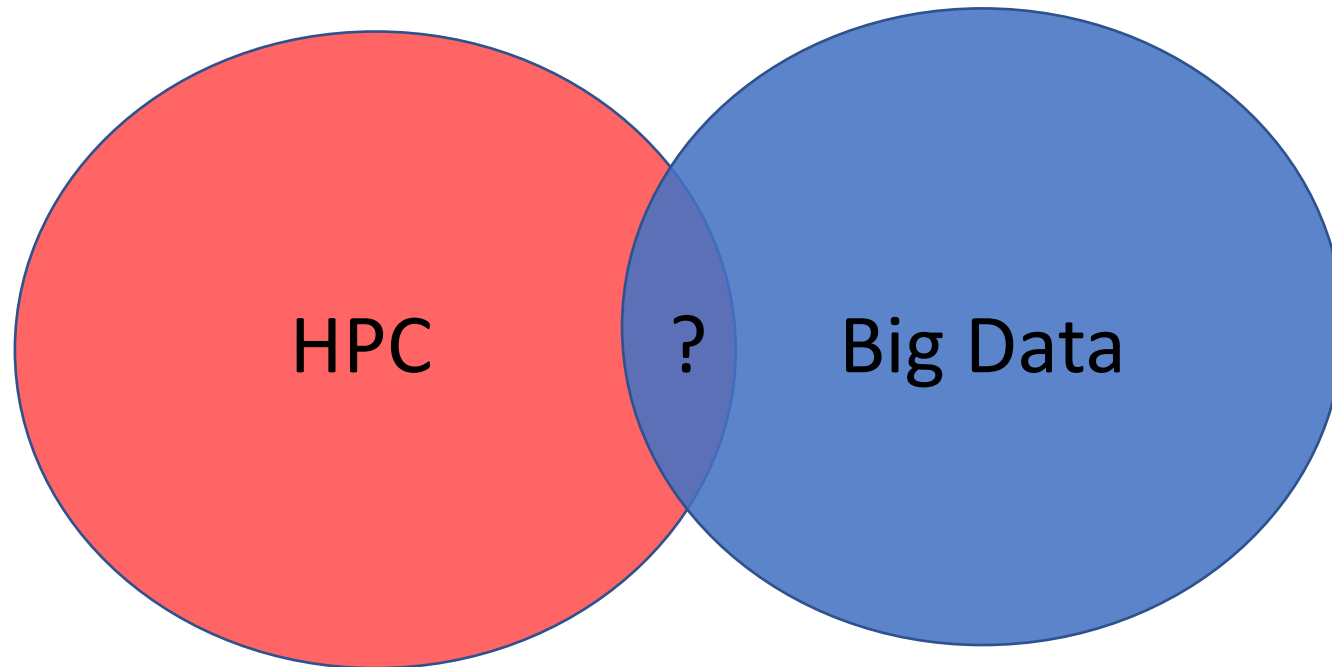## Stanford/SLAC

# A Tale of Two Software Cultures

HPC

MPI +
OpenMP +
CUDA

Big Data/ML

Hadoop/MapReduce
Spark
TensorFlow
PyTorch

# Is There Any Relationship?



HPC

Big Data
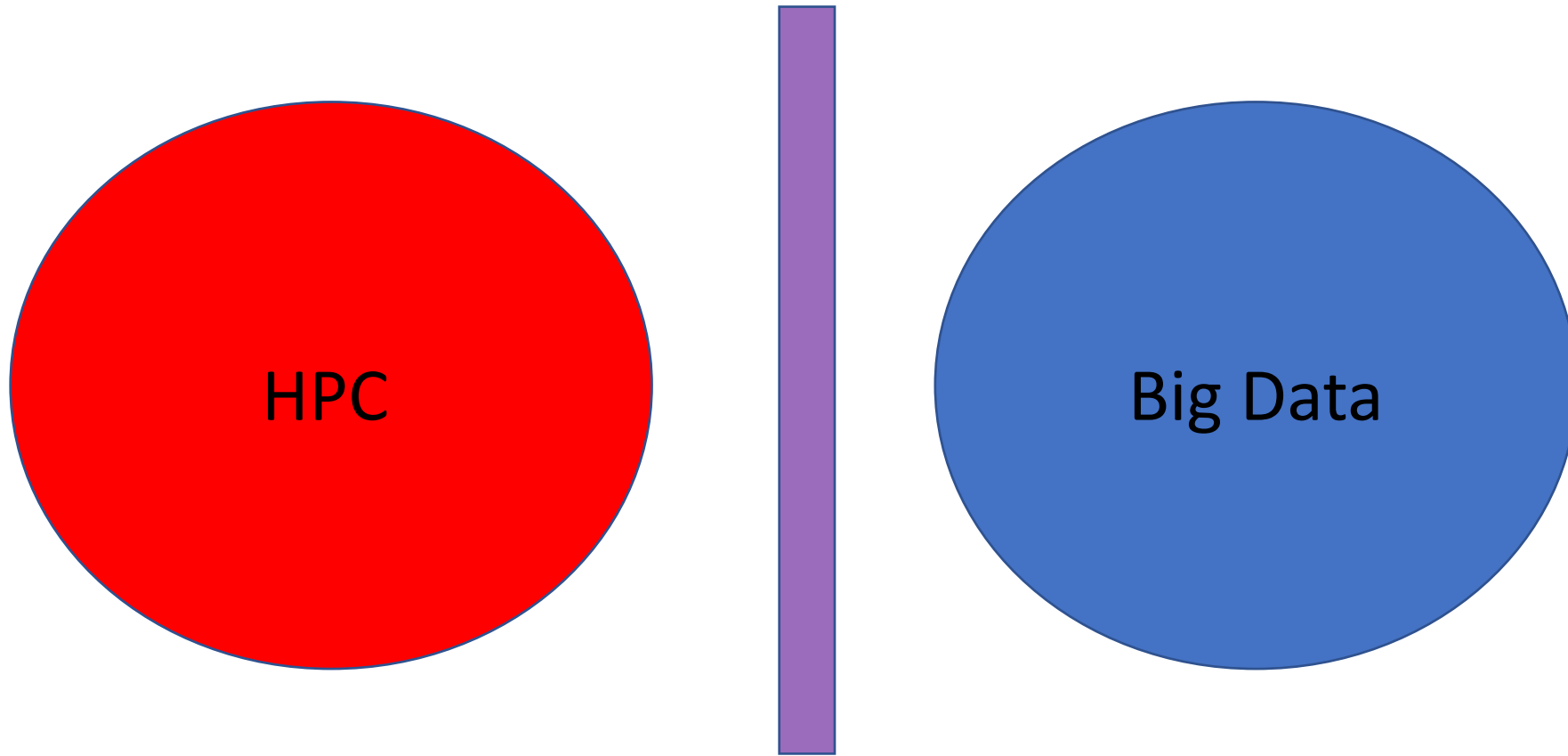
# Some Overlap …



HPC    ?    Big Data

# Some Difference in Size ...



HPC  ?  Big Data

# Will These Communities Converge?

- The stage is set:  The underlying hardware is (almost) the same

- More shortly …

# Are There Barriers to Convergence?

HPC

Big Data

# Priorities

## HPC

- Performance

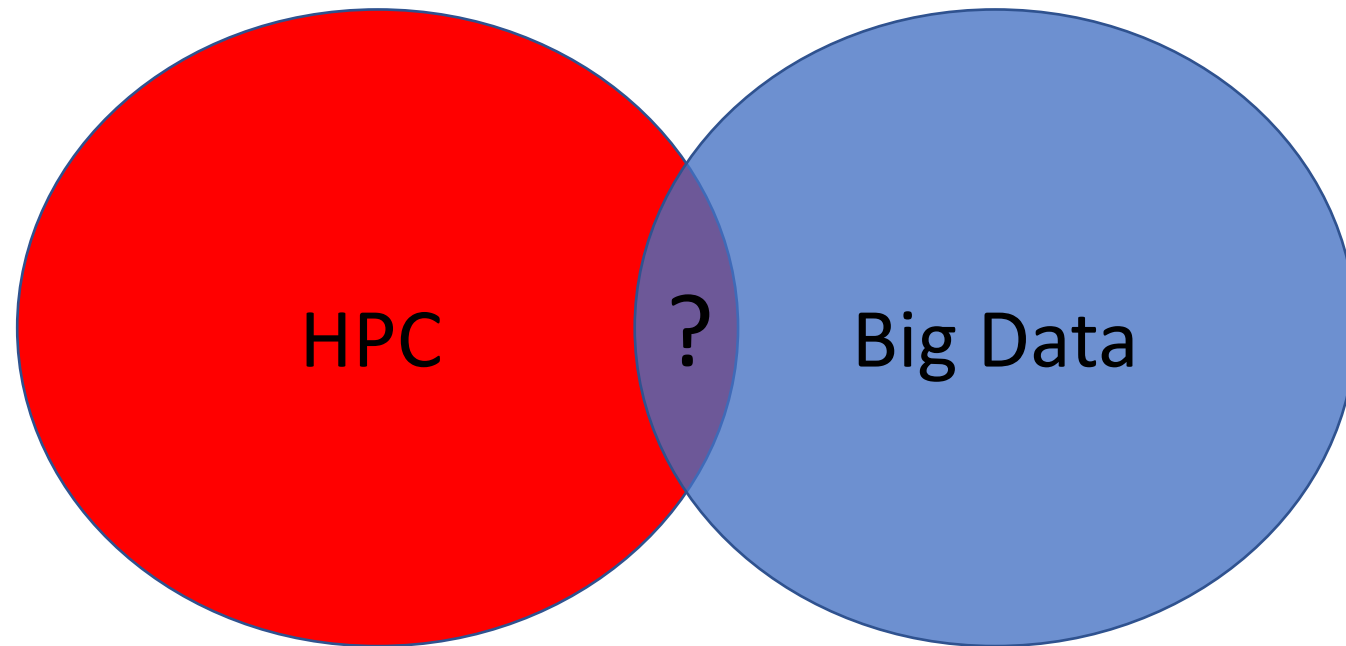- Productivity

- Correctness

## Big Data

- Productivity

- Performance

- Correctness

# Creates Significant Differences In …

- Platform performance & programmer productivity
  - Obviously!


- Scale of computations


- Economic model
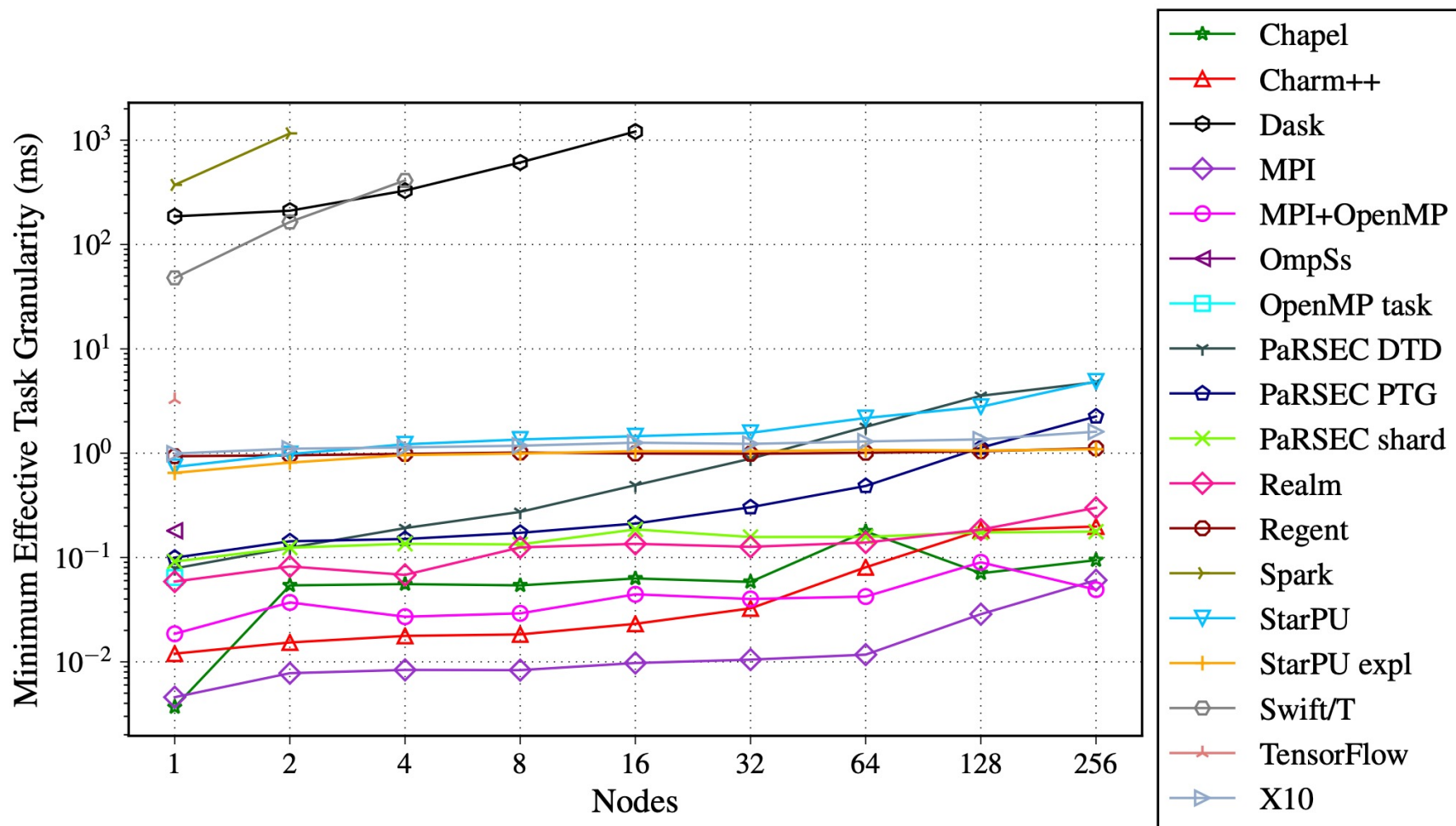
# Is There Overlap Today?

# Who Would Switch from Big Data to HPC?

# 0%

# Who Would Switch from HPC to Big Data?
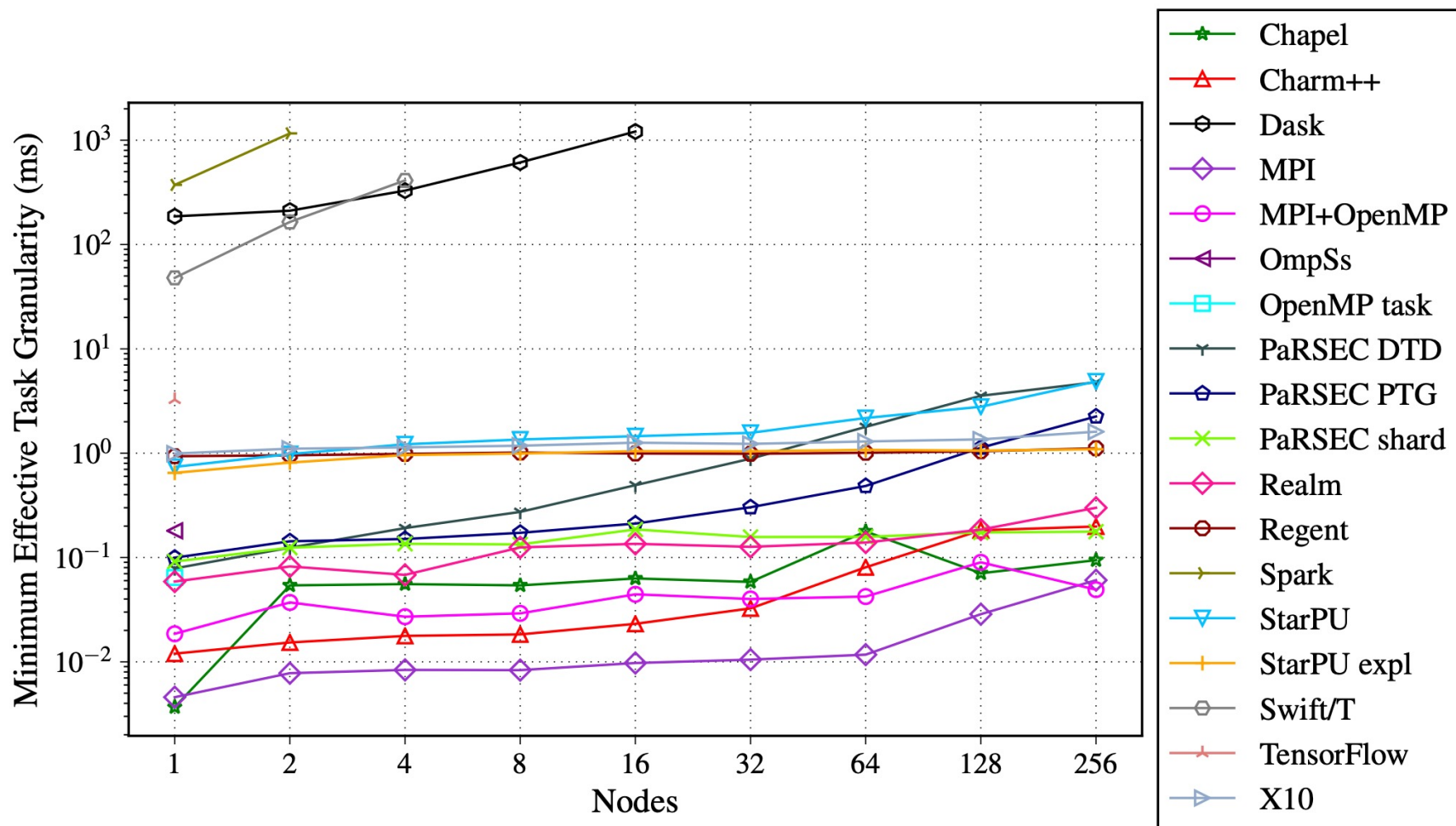
- If performance improved by switching, everyone

- If performance were comparable or not overly harmed, some

- If performance is 10X worse, none
  - And some would not switch even if performance is only 2X worse

# A Comparison: Minimum Task Granularity



(a) Stencil pattern.

Task Bench: A Parameterized Benchmark for Evaluating Parallel Runtime Performance, Slaughter et al, SC'20

# A Comparison: Minimum Task Granularity



(a) Stencil pattern.

Task Bench: A Parameterized Benchmark for Evaluating Parallel Runtime Performance, Slaughter et al, SC'20

# A Comparison: Minimum Task Granularity



(d) Nearest pattern, 5 deps/task, 4 independent graphs.

Task Bench: A Parameterized Benchmark for Evaluating Parallel Runtime Performance, Slaughter et al, SC'20

# A Brief Digression: Hardware

- The hardware platform drives the software abstractions

- The current, slow-motion revolution: accelerators
    - GPUs today
    - Other specialized hardware tomorrow

# A Key Point

- In new supercomputers, > 95% of performance is in the accelerators
  - Titan, Summit, PerlMutter, Frontier, Aurora …

- The tradeoff
  - Greatly complicates programming
  - But switching to GPUs can greatly increase performance

- This is the ground on which any convergence will happen

# An Observation

- The HPC community values performance
  - Unless it is too hard
  - Many HPC systems perform far below their potential today


- The Big Data community values productivity
  - Until the code takes forever to run
  - Organizations spend inordinate amounts of time tweaking for performance

# The Technical Issue

- The main limiter in current and future systems is data movement
  - By far the most expensive part of any computation
  - And accelerators add multiple levels of memory hierarchy

- Few programming abstractions in programming models for
  - Locality
  - Partitioning of data
  - Mapping of compute/data into a machine

# The Evidence

- S3D
  - Production chemistry combustion code
  - 7X off its potential

- Large graph analytics
  - CPU-based state of the art ~10X off potential

<span style="color:red">Switching to GPUs + good data partitioning & placement</span>

- Machine Learning
  - 10X off potential

<span style="color:red">Improved data partitioning</span>

# Where Does Productivity Come From?

- Libraries

- How many widely used parallel libraries for HPC are there?

- How many widely used libraries are there for Python?
  - Not just "big data"

# Numpy
## In One Slide

A popular Python package for (mostly) dense array computing

Common building block in other Python packages

Many drop-in replacements for one GPU

```python
import numpy as np

def cg_solve(A, b, tol=1e-10):
    x = np.zeros(A.shape[1])
    r = b - A.dot(x)
    p = r
    rsold = r.dot(r)
    for i in xrange(b.shape[0]):
        Ap = A.dot(p)
        alpha = rsold / (p.dot(Ap))
        x = x + alpha * p
        r = r - alpha * Ap
        rsnew = r.dot(r)
        if np.sqrt(rsnew) < tol:
            break
        beta = rsnew / rsold
        p = r + beta * p
        rsold = rsnew
    return x
```
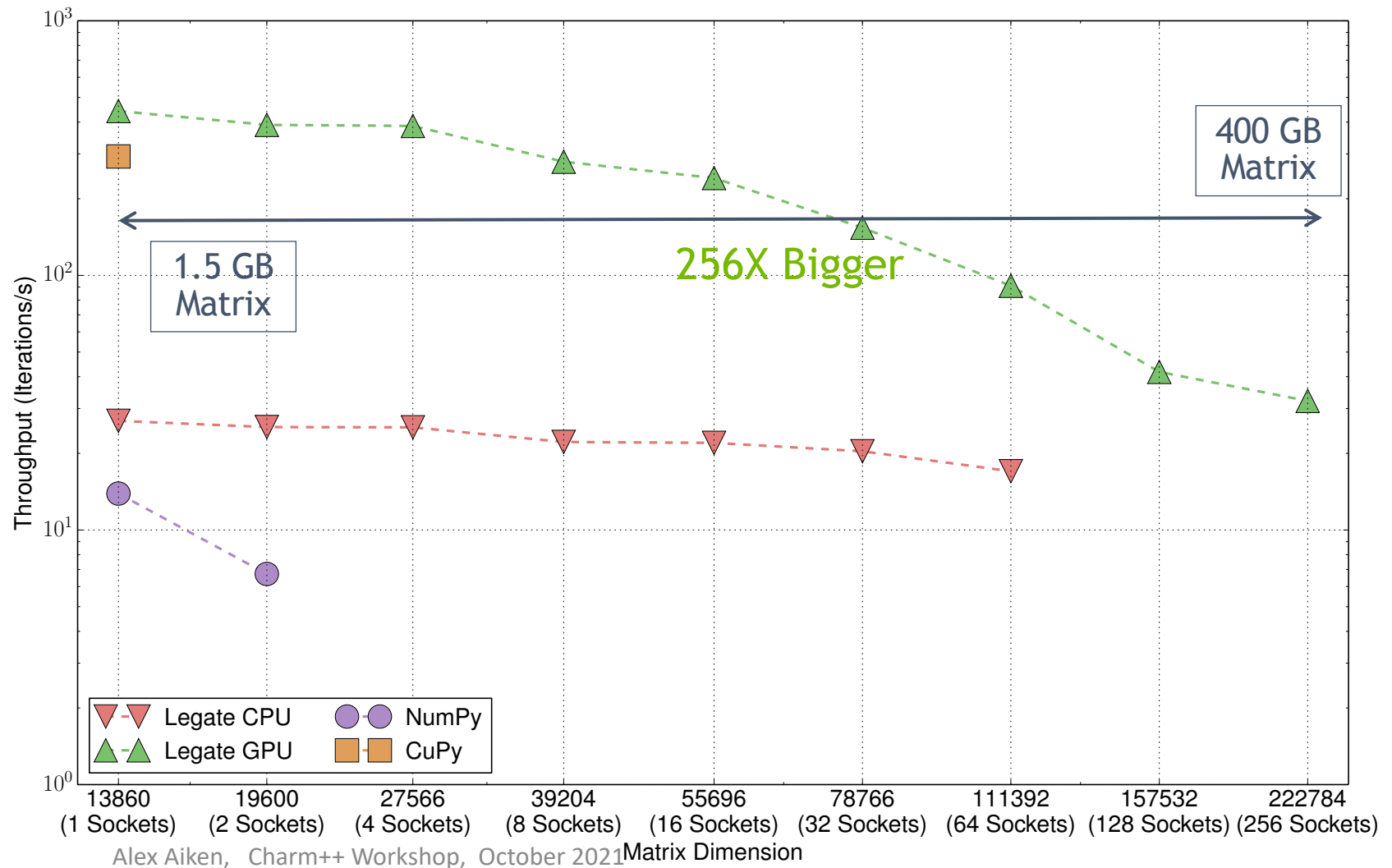
# Legate Numpy
## Accelerated and Distributed

Legate NumPy is a NumPy replacement for transparent (weak) scaling

Requires a one line code change

Same code runs on everything

Legate NumPy: Accelerated and distributed array computing,   Bauer & Garland SC'19

```python
import legate.numpy as np

def cg_solve(A, b, tol=1e-10):
    x = np.zeros(A.shape[1])
    r = b - A.dot(x)
    p = r
    rsold = r.dot(r)
    for i in xrange(b.shape[0]):
        Ap = A.dot(p)
        alpha = rsold / (p.dot(Ap))
        x = x + alpha * p
        r = r - alpha * Ap
        rsnew = r.dot(r)
        if np.sqrt(rsnew) < tol:
            break
        beta = rsnew / rsold
        p = r + beta * p
        rsold = rsnew
    return x
```

# A Simple Example: A Jacobi Solver

```python
import legate.numpy as np

A = np.random.rand(N,N)
b = np.random.rand(N)

x = np.zeros(A.shape[1])
d = np.diag(A)
R = A - np.diag(d)
for i in xrange(b.shape[0]):
    x = (b - np.dot(R,x)) / d
```



400 GB Matrix

1.5 GB Matrix

256X Bigger

Throughput (Iterations/s)

Legate CPU    NumPy
Legate GPU    CuPy

13860 (1 Sockets)   19600 (2 Sockets)   27566 (4 Sockets)   39204 (8 Sockets)   55696 (16 Sockets)   78766 (32 Sockets)   111392 (64 Sockets)   157532 (128 Sockets)   222784 (256 Sockets)

Matrix Dimension

# Legate NumPy Architecture

# Legate NumPy architecture



Map n-D arrays to Legion data model

Application

Program order
np.argmin
np.sort
np.add
np.dot
np.mul
np.norm

Legion dynamically computes dependence graph ...

A  B  C

(0, 0)

(i, j)

Legate Mapper

CPUs    GPUs

copy

copy

copy

Execution order

Alex Aiken,   Charm++ Workshop,  October 2021

# Managing Data

Each N-D array maps to a field of a Legion *logical region*

- Legion's collection data type

Different logical regions for different shapes

Dynamically allocated on demand and recycled when GC'd by Python

```python
import legate.numpy as np

A = np.random.rand(N,N)
b = np.random.rand(N)

x = np.zeros(A.shape[1])
d = np.diag(A)
R = A - np.diag(d)
for i in xrange(b.shape[0]):
    x = (b - np.dot(R,x)) / d
```
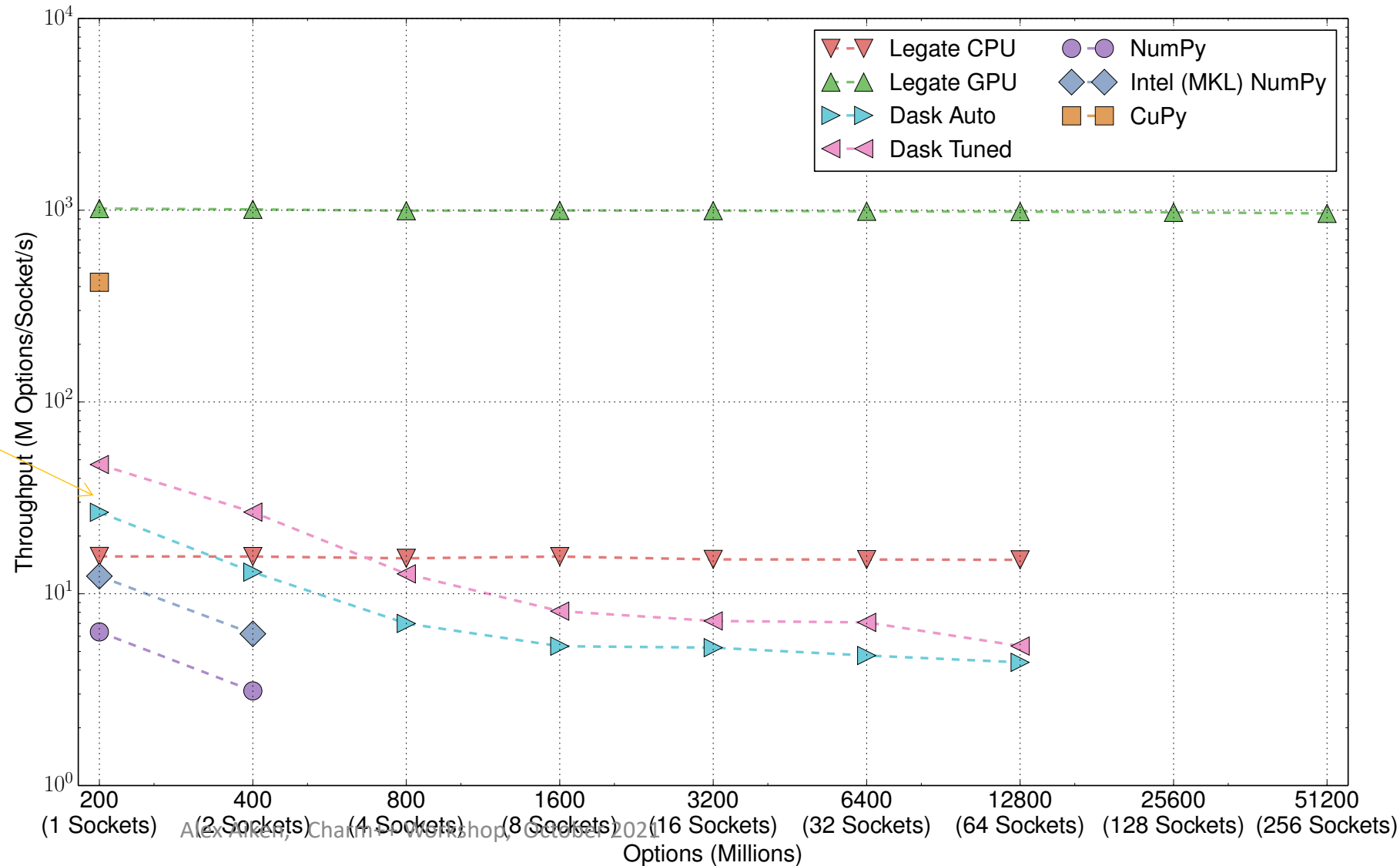
Region with Index Space (N,N)

| $A_{ij}$ | $R_{ij}$ |
|---|---|
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |

Region with Index Space (N)

| $b_i$ | $x^0_i$ | $d_i$ | $x^1_i$ | $x^2_i$ | $x^3_i$ |
|---|---|---|---|---|---|
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |

# Performance Comparison

Compare NumPy implementations:

Standard NumPy (single node)

IntelPy with MKL (single node)

Legate CPU-only

Legate CPU+GPU

Dask (CPU-only): Auto and Tuned


All plots are log-log

Experiments on a cluster of DGX-1V nodes

Weak scaling throughput on sockets

**DASK**

Popular Python library for parallel and distributed computing

dask.array similar to NumPy, except for specifying "chunk" sizes

```python
import dask.array as da

A = da.random.uniform((N,N),
        chunks=(C,C))
b = da.random.uniform(N,
        chunks="auto")
x = da.zeros(A.shape[1],
        chunks=b.chunks)
d = da.diag(A)
R = A - da.diag(d)
for i in xrange(b.shape[0]):
    x = (b - da.dot(R,x)) / d
```

# Jacobi Solver

```python
import numpy as np

A = np.random.rand(N,N)
b = np.random.rand(N)

x = np.zeros(A.shape[1])
d = np.diag(A)
R = A - np.diag(d)
for i in xrange(b.shape[0]):
    x = (b - np.dot(R,x)) / d
```

# Black Scholes

No (application) communication

Expect perfect weak scaling

Dask starts out faster... Why? Operator Fusion

... but has to trade off parallelism for task granularity to scale

# Preconditioned CG Solver

```python
def preconditioned_solve(A, M, b):
    x = np.zeros(A.shape[1])
    r = b - A.dot(x)
    z = M.dot(r)
    p = z
    rzold = r.dot(z)
    for i in xrange(b.shape[0]):
        Ap = A.dot(p)
        alpha = rzold / (p.dot(Ap))
        x = x + alpha * p
        r = r - alpha * Ap
        rznew = r.dot(r)
        if np.sqrt(rznew) < 1e-10:
            break
        z = M.dot(r)
        rznew = r.dot(z)
        beta = rznew / rzold
        p = z + beta * p
        rzold = rznew
    return x
```

# One Approach To Libraries

- Implement important Big Data libraries using HPC techniques
  - Can we get more performance for the same productivity?

- Examples
  - Legate
  - FlexFlow, replacement for TensorFlow & PyTorch
    Beyond data and model parallelism for deep neural networks, Jia et al. SysML `18

# Important Features

- Expressive data partitioning

- Ability to tune the *mapping*
  - Tasks to processors
  - Data to memories

- Runtime decision making
  - Needed to handle dynamic nature of Python

- Legion is extreme in all three dimensions
  - Sufficient, but maybe not necessary?

# Another Approach

- Demonstrate the ability to build general libraries for HPC applications
  - That compete with the best-of-class HPC implementations
  - But are more productive to write and/or use

- What are the important/novel problems in building HPC libraries?

# DISTAL: DIStributed Tensor ALgebra

Goals:

Compile tensor algebra kernels into efficient distributed implementations

Decouple computation, performance optimizations, and data distribution



Joint work with Rohan Yadav and Fred Kjolstad

# Modeling Machines

- View machines as hyper-rectangular grids of processors
  - where each processor has a local memory

- Expose any locality in the physical machine

- Structure the machine like the target computations

# Distributing Data

- State abstractly how a tensor is distributed onto a machine as part of the tensor's *format*

- Describes how dimensions of a tensor $\mathcal{T}$ map onto a machine $\mathcal{M}$

$$\mathcal{T}_{ij} \longmapsto_i \mathcal{M}$$

Name each dimension of $\mathcal{T}$ and $\mathcal{M}$

Dimensions of $\mathcal{T}$ are partitioned and mapped onto dimensions of $\mathcal{M}$ that share the same name

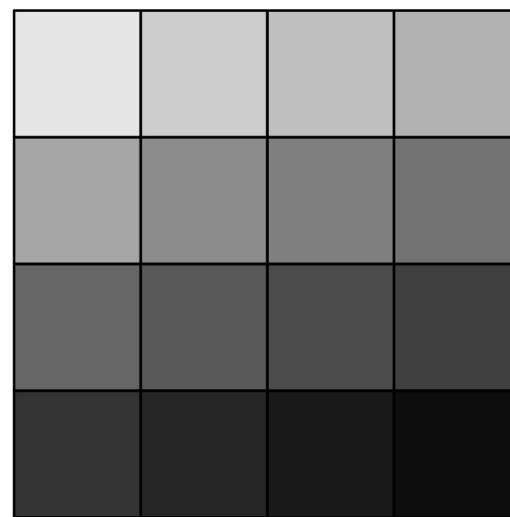$$\mathcal{T} \; {}_{\text{i}}{\longmapsto}_{\text{i}} \; \mathcal{M}$$

$\mathcal{T}$

$\mathcal{M}$

$$\mathcal{T}\ _{ij\mapsto_i}\ \mathcal{M}$$

$$\mathcal{T}_{\text{ij}} \longmapsto_{\text{ij}} \mathcal{M}$$

$\mathcal{T}$

$\mathcal{M}$

$$\mathcal{T} \underset{\mathrm{ijk} \longmapsto \mathrm{ij}}{} \mathcal{M}$$

$\mathcal{T}$

$\mathcal{M}$

$$\mathcal{T} \underset{\mathrm{ij} \mapsto \mathrm{ij}*}{} \mathcal{M}$$

$\mathcal{T}$

$\mathcal{M}$

$$\mathcal{T} \underset{\mathrm{ij} \mapsto \mathrm{ij0}}{\longrightarrow} \mathcal{M}$$
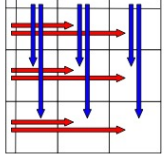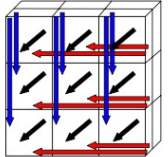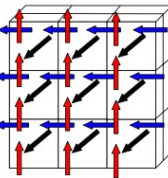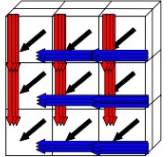
$\mathcal{T}$

$\mathcal{M}$

# Scheduling (Summary)

- Iteration spaces: hyper-rectangular grids representing points in nested loops
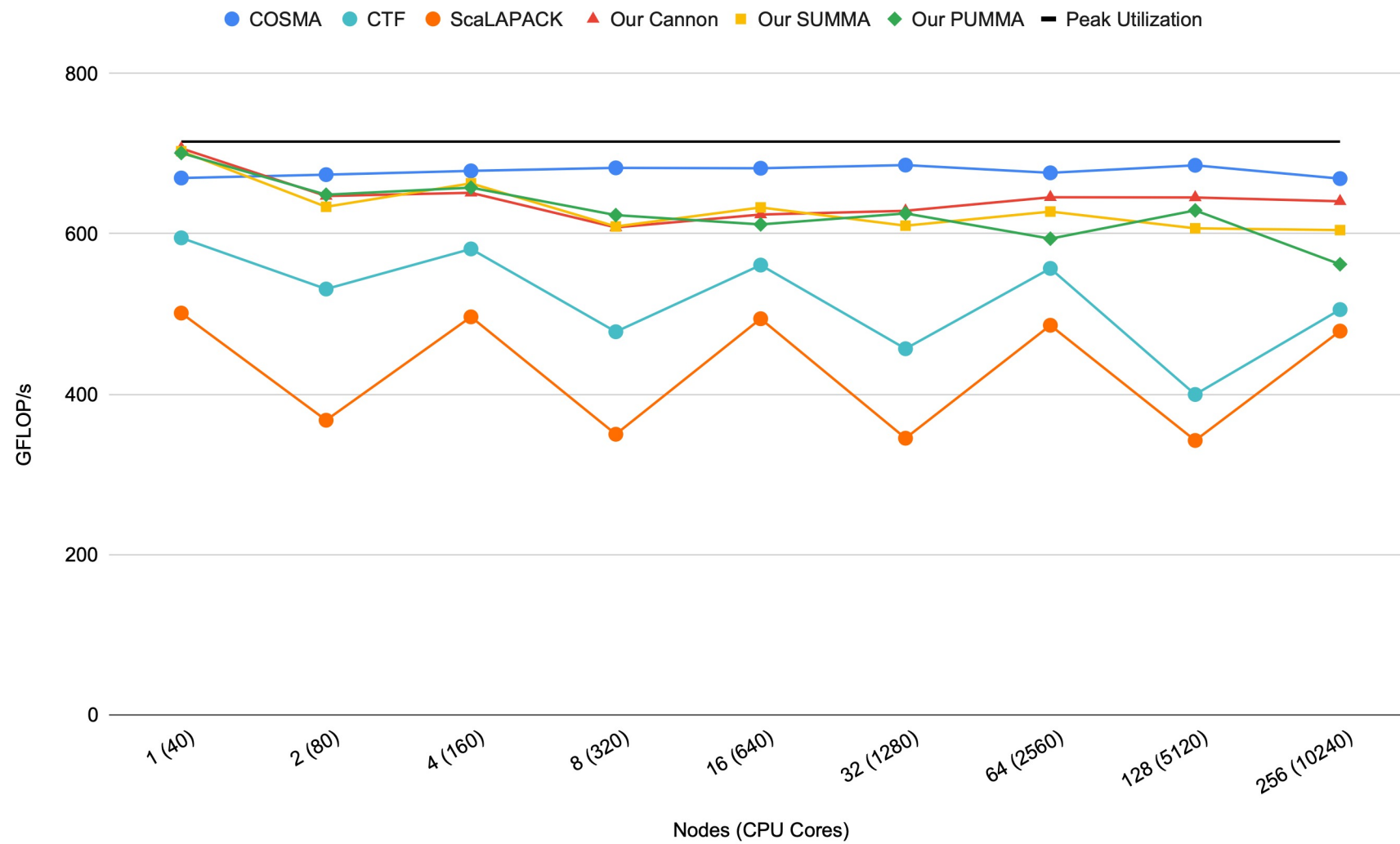
$$\forall_i \ A_i = \sum_j B_j$$

- Execution space: processors in $\mathcal{M}$ x time dimension

- Scheduling commands related to distribution change mapping of iteration space points to the execution space

- Apply scheduling commands to the computation

  - Similar to Halide schedules, with extensions for distributed computing

  - New commands: *distribute, communicate, rotate*

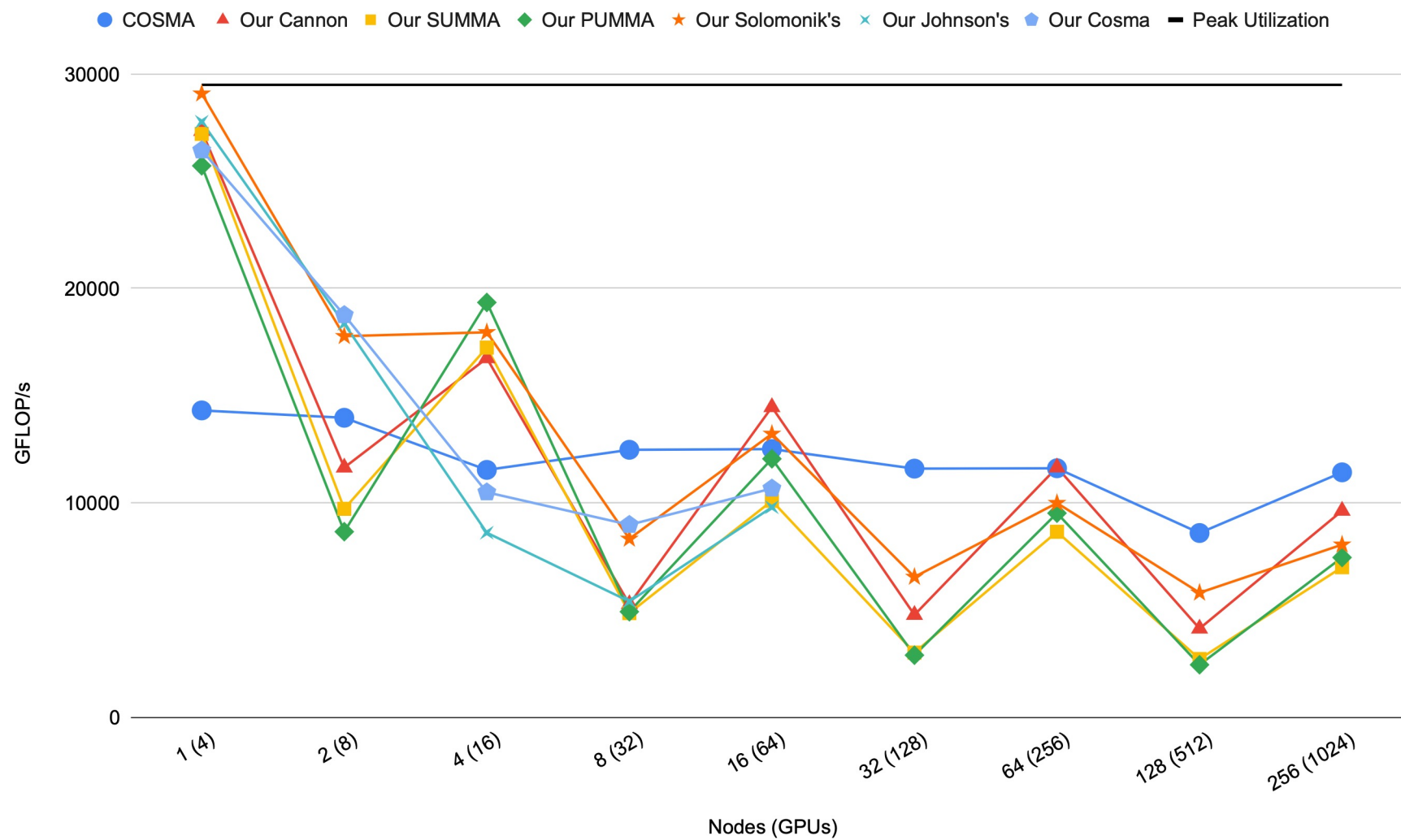| Algorithm | Comm. Pattern | Target Machine | Data Distribution | Schedule |
|-----------|---------------|----------------|-------------------|----------|
| Cannon's [7] (1969) |  | $\mathcal{M}(gx, gy)$ | $A_{ij} \mapsto_{ij} \mathcal{M}$<br>$B_{ij} \mapsto_{ij} \mathcal{M}$<br>$C_{ij} \mapsto_{ij} \mathcal{M}$ | `.distribute({i, j}, {in, jn}, {il, jl}, Grid(gx, gy))`<br>`.divide(k, ko, ki, gx)`<br>`.reorder({ko, il, jl, ki})`<br>`.rotate(ko, {in, jn}, kos)`<br>`.communicate(A, jn)`<br>`.communicate({B, C}, kos)` |
| PUMMA [10] (1994) |  | $\mathcal{M}(gx, gy)$ | $A_{ij} \mapsto_{ij} \mathcal{M}$<br>$B_{ij} \mapsto_{ij} \mathcal{M}$<br>$C_{ij} \mapsto_{ij} \mathcal{M}$ | `.distribute({i, j}, {in, jn}, {il, jl}, Grid(gx, gy))`<br>`.divide(k, ko, ki, gx)`<br>`.reorder({ko, il, jl, ki})`<br>`.rotate(ko, {in}, kos)`<br>`.communicate(A, jn)`<br>`.communicate({B, C}, kos)` |
| SUMMA [25] (1995) |  | $\mathcal{M}(gx, gy)$ | $A_{ij} \mapsto_{ij} \mathcal{M}$<br>$B_{ij} \mapsto_{ij} \mathcal{M}$<br>$C_{ij} \mapsto_{ij} \mathcal{M}$ | `.distribute({i, j}, {in, jn}, {il, jl}, Grid(gx, gy))`<br>`.split(k, ko, ki, chunkSize)`<br>`.reorder({ko, il, jl, ki})`<br>`.communicate(A, jn)`<br>`.communicate({B, C}, ko)` |
| Johnson's [1] (1995) |  | $\mathcal{M}(\sqrt[3]{p}, \sqrt[3]{p}, \sqrt[3]{p})$ | $A_{ij} \mapsto_{ij0} \mathcal{M}$<br>$B_{ik} \mapsto_{i0k} \mathcal{M}$<br>$C_{kj} \mapsto_{0jk} \mathcal{M}$ | `.distribute({i, j, k}, {in, jn, kn},`<br>`        {il, jl, kl}, Grid(`$\sqrt[3]{p}$`, `$\sqrt[3]{p}$`, `$\sqrt[3]{p}$`))`<br>`.communicate({A, B, C}, kn)` |
| Solomonik's [22] (2011) |  | $\mathcal{M}(\sqrt{\frac{p}{c}}, \sqrt{\frac{p}{c}}, c)$ | $A_{ij} \mapsto_{ij0} \mathcal{M}$<br>$B_{ij} \mapsto_{ij0} \mathcal{M}$<br>$C_{ij} \mapsto_{ij0} \mathcal{M}$ | `.distribute({i, j, k}, {in, jn, kn},`<br>`        {il, jl, kl}, Grid(`$\sqrt{\frac{p}{c}}$`, `$\sqrt{\frac{p}{c}}$`, c))`<br>`.divide(kl, k1, k2, `$\sqrt{\frac{p}{c^3}}$`)`<br>`.reorder({k1, il, jl, k2})`<br>`.rotate(k1, {in, jn}, k1s)`<br>`.communicate(A, jn)`<br>`.communicate({B, C}, k1s)` |
| COSMA [17] (2019) |  | induced by schedule | induced by schedule | `// gx, gy, gz, numSteps computed by COSMA scheduler.`<br>`.distribute({i, j, k}, {in, jn, kn}`<br>`        {il, jl, kl}, Grid(gx, gy, gz))`<br>`.divide(kl, klo, kli, numSteps)`<br>`.reorder({klo, il, jl, kli})`<br>`.communicate(A, kn)`<br>`.communicate({B, C}, klo)` |

# Experiments

- Run on Lassen
  - 4 GPUs/node, 40 CPUs/node, IB interconnect)

- All systems configured to use the same BLAS / CuBLAS

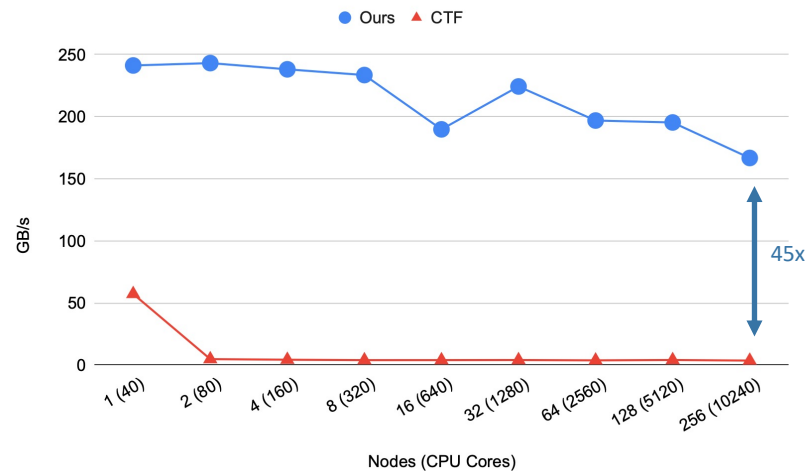- All experiments are weak-scaling (memory / node stays constant)

# GEMM (CPU)

# GEMM (GPU)

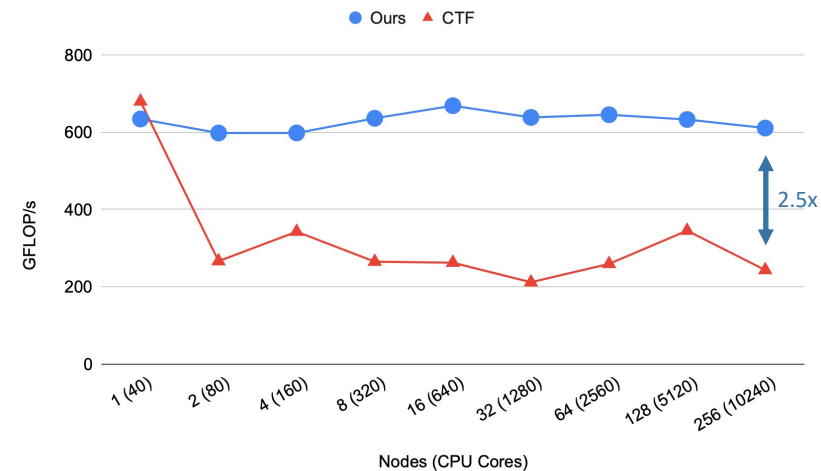# Higher Order Tensor Operations (CPU)
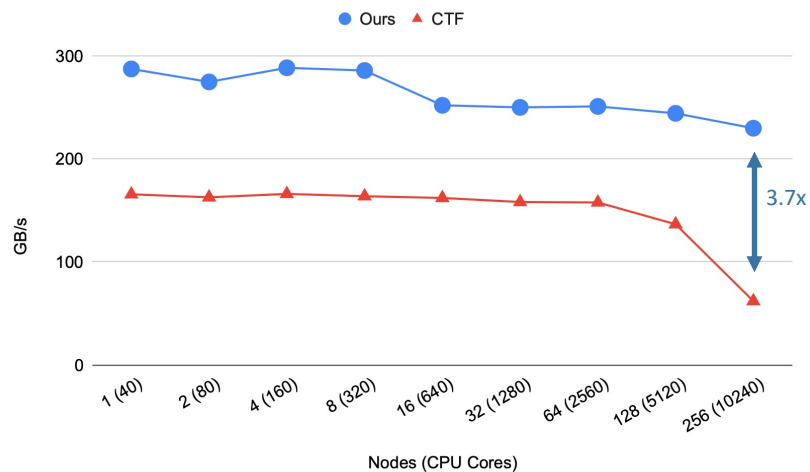


TTV

$$A_{ij} = B_{ijk} \cdot C_k$$

TTM

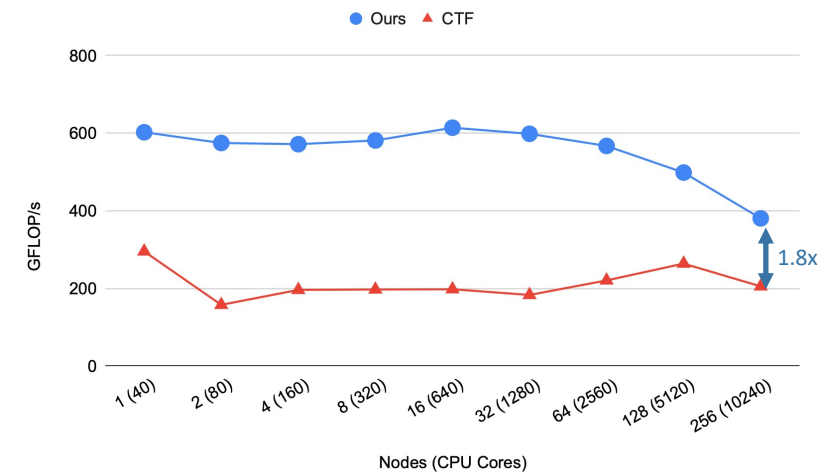$$A_{ijl} = B_{ijk} \cdot C_{kl}$$

InnerProd

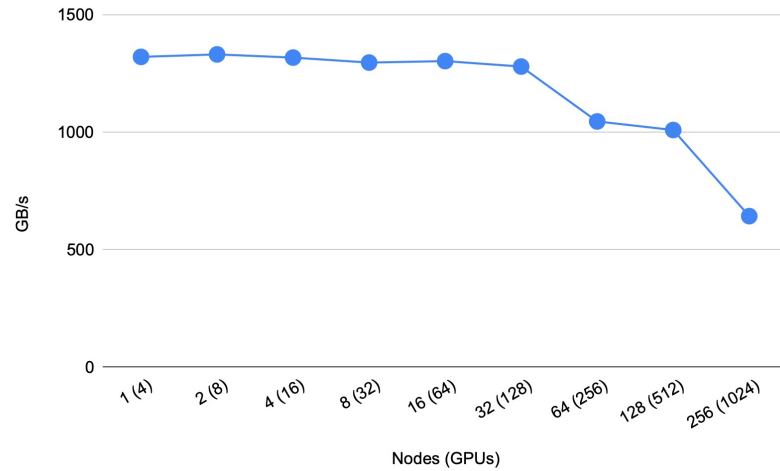$$a = B_{ijk} \cdot C_{ijk}$$

MTTKRP

$$A_{il} = B_{ijk} \cdot C_{jl} \cdot D_{kl}$$

# Higher Order Tensor Operations (GPU)
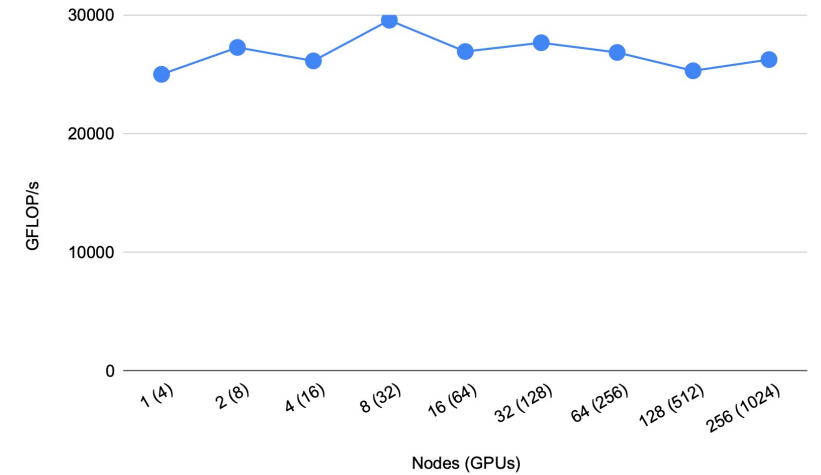
TTV

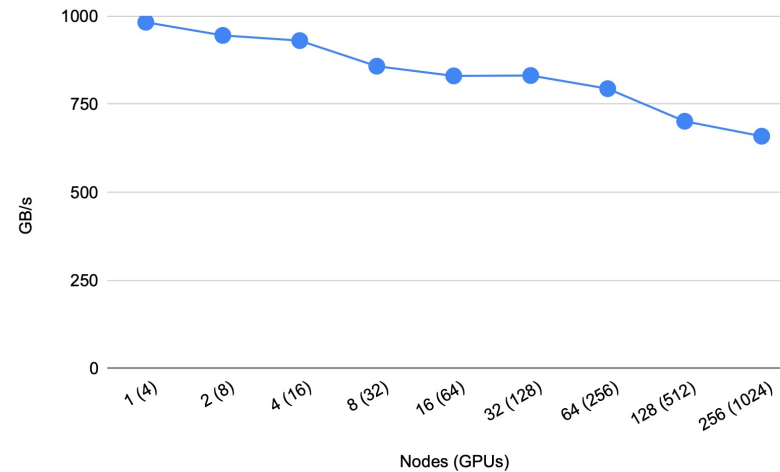$$A_{ij} = B_{ijk} \cdot C_k$$



TTM

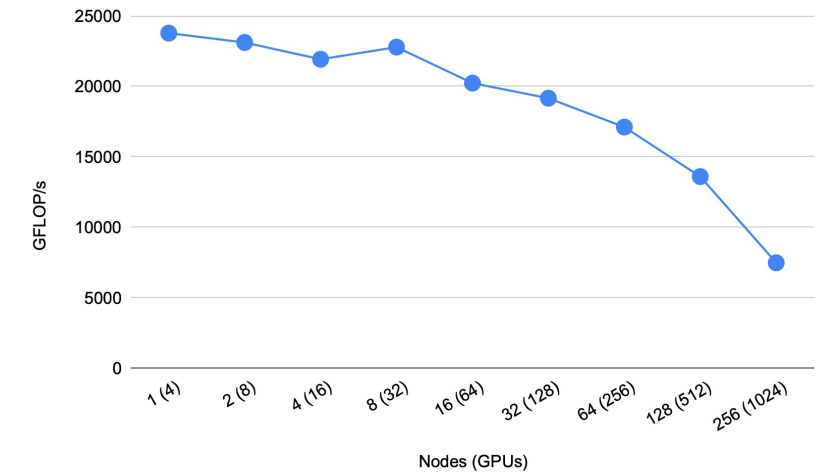$$A_{ijl} = B_{ijk} \cdot C_{kl}$$



InnerProd

$$a = B_{ijk} \cdot C_{ijk}$$



MTTKRP

$$A_{il} = B_{ijk} \cdot C_{jl} \cdot D_{kl}$$

# Lessons From DISTAL

- Expressive partitioning of data, computation and control of the mapping into the machine are all critical

- Enables writing libraries that are polymorphic in the data distribution
  - The data distribution can be different depending on the needs of the context
  - Avoids stopping-the-world and doing large copies at library boundaries
  - A form of polymorphism unique to distributed parallel programming

# Summary

- The HPC and Big Data worlds have agreed on the hardware platform
  - Parallel, accelerated, distributed (PAD) machines
  - A convergence of these two worlds is likely


- Can we have both productivity and performance?
  - There is some preliminary evidence the answer is ``yes"
  - Through libraries built on HPC programming models
  - But libraries required a degree of flexibility beyond non-library code
    - Still much to be learned about how to write reusable parallel libraries