Sandia National Laboratories

SAND2021-2166C

# Model Parallelism with Spatial Decomposition of Volumetric Data for Deep Learning

**Gary Saavedra**
Eric C. Cyr

Sandia National Laboratories

Jacob Schroder

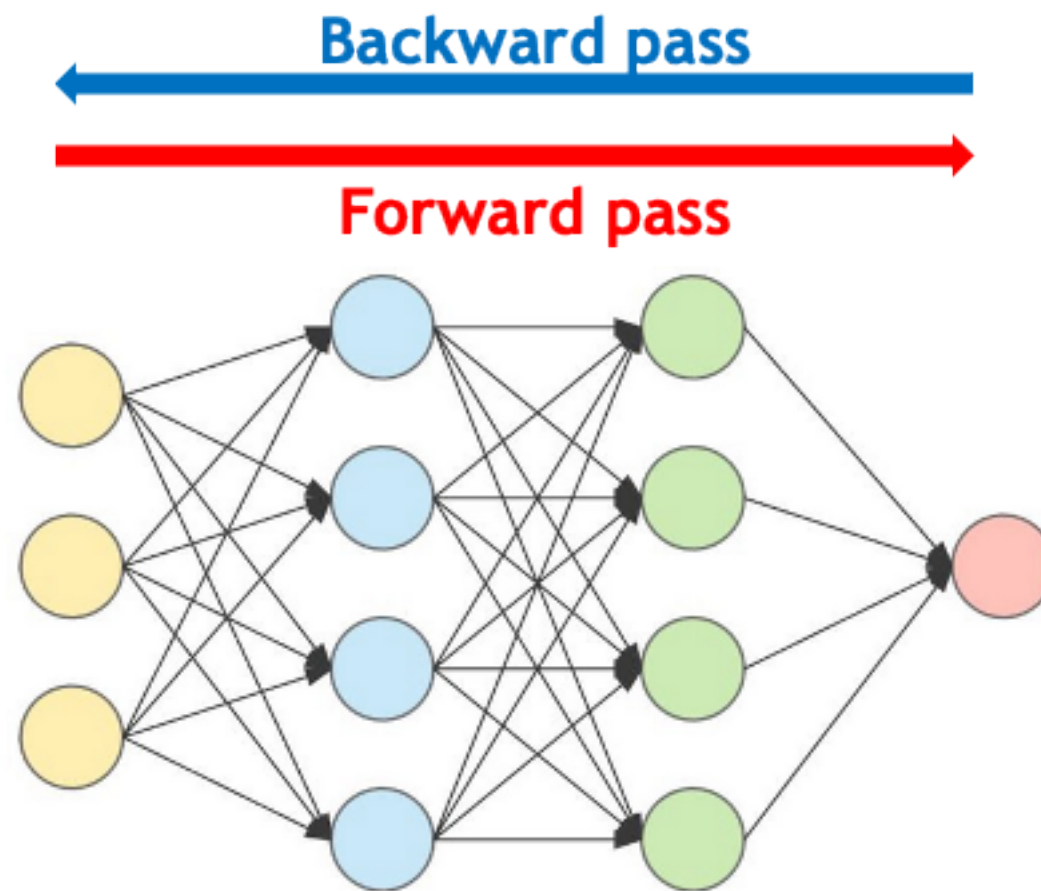University of New Mexico

Russell Hewett

Virginia Tech

ENERGY    NNSA

**SAND No:**

# Deep Learning and Time Parallelism

- More layers can improve performance

- Several groups have shown very deep networks offer improvement
  - GoogLeNet has 22 layers[1]
  - Huang et al show improvement on Cifar-10 data using up to 1200 layers[2]
  - Recurrent neural networks

- Very deep networks may be ignored due    to training limitations.

**Backward pass**

**Forward pass**

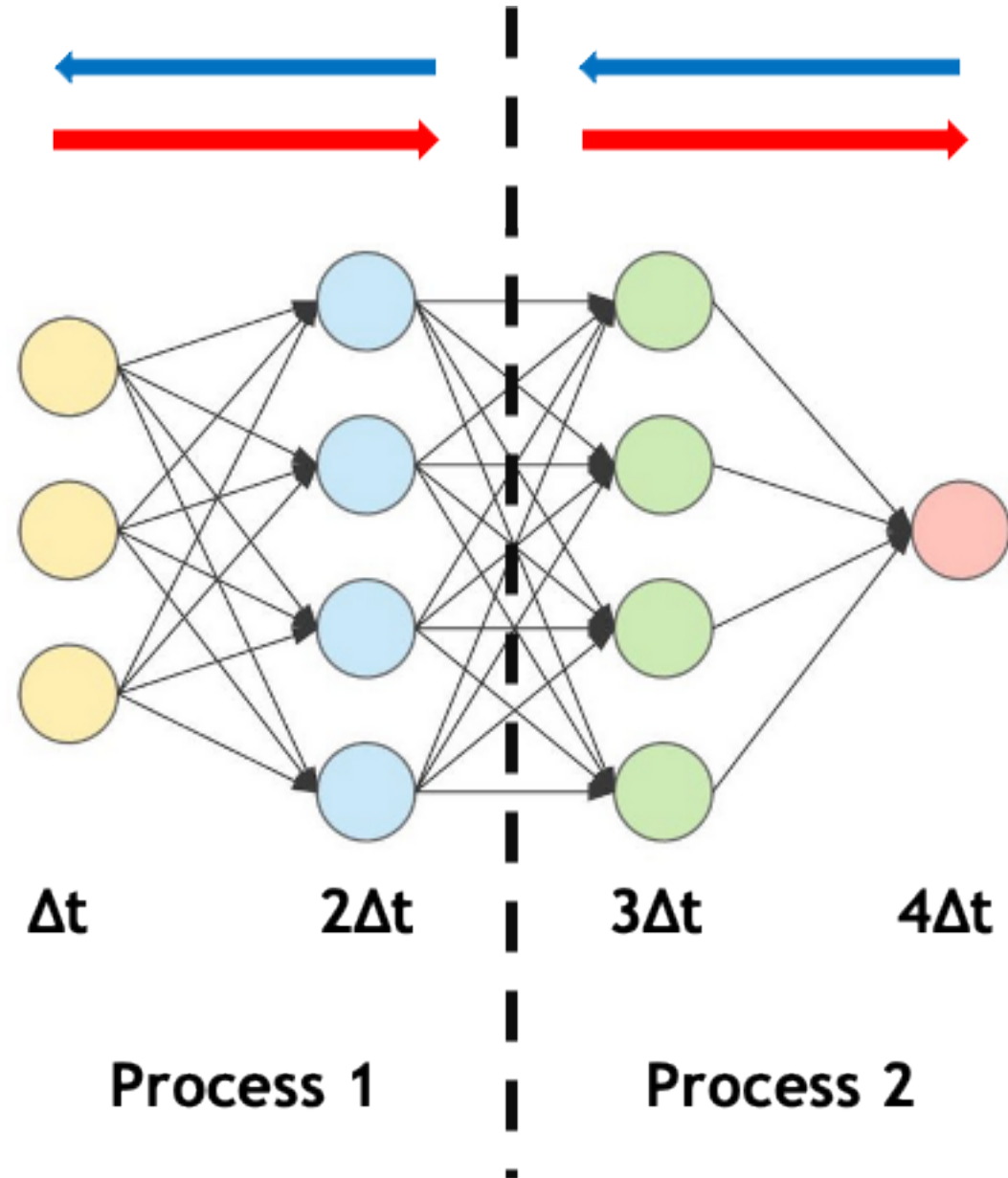**Feedforward** $\quad \mathbf{h}_{t+1} = f(\mathbf{h}_t, \theta_t)$

**ResNet** $\quad \mathbf{h}_{t+1} = \mathbf{h}_t + f(\mathbf{h}_t, \theta_t)$

1) Szegedy et al. Going Deeper with Convolutions. 2015.
2) Huang et al. Deep Networks with Stochastic Depth. 2016.
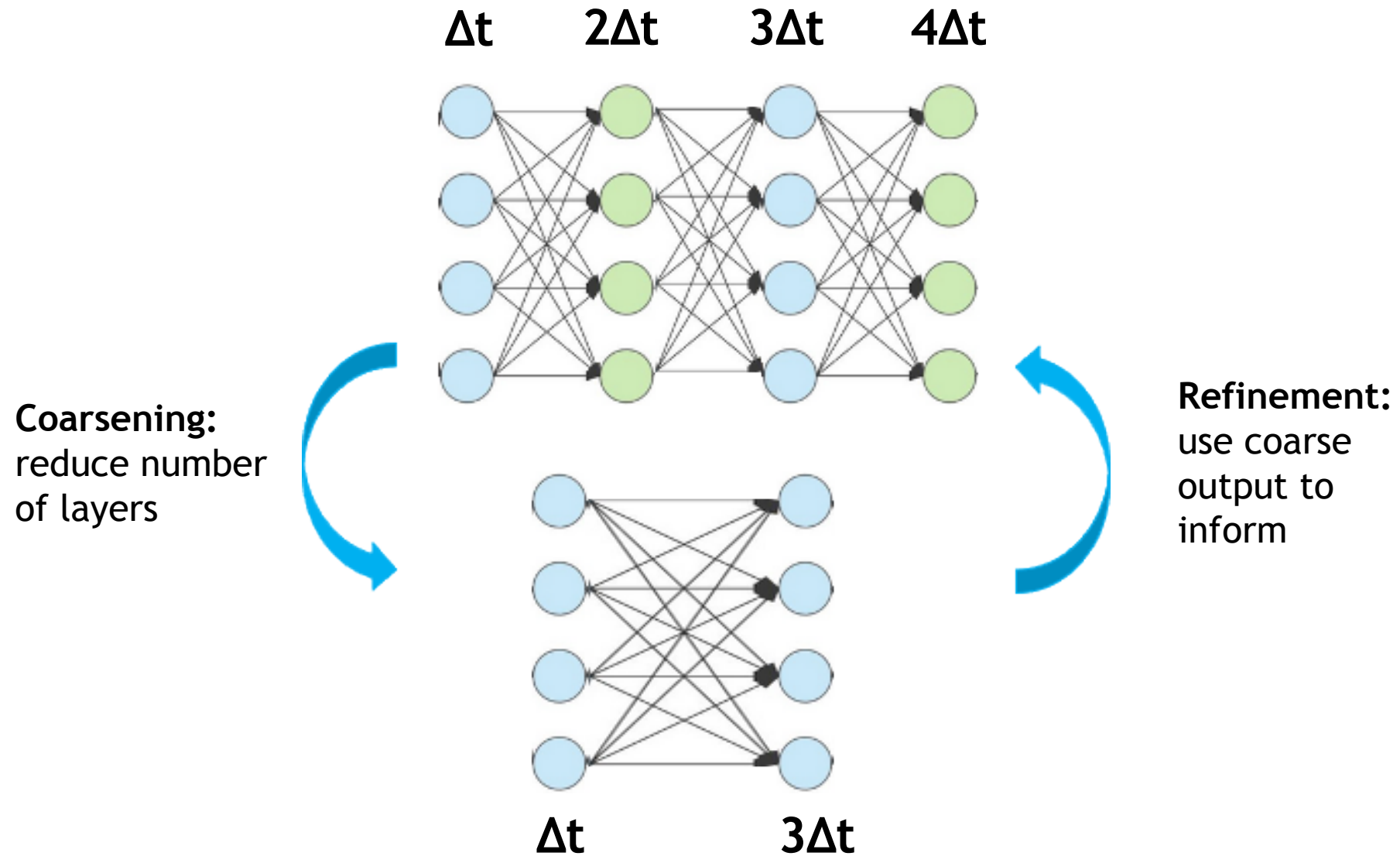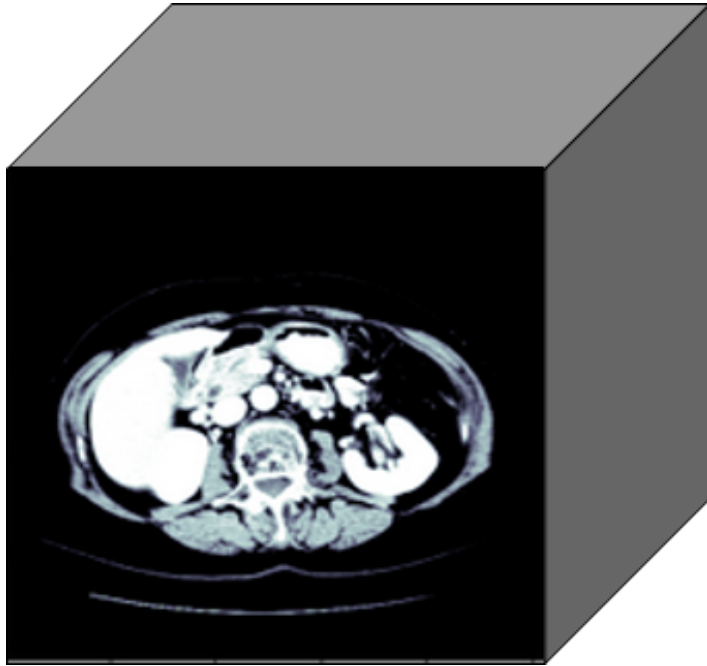3) kdnuggets.com

# Deep Learning and Time Parallelism

- Separate groups of layers between processors

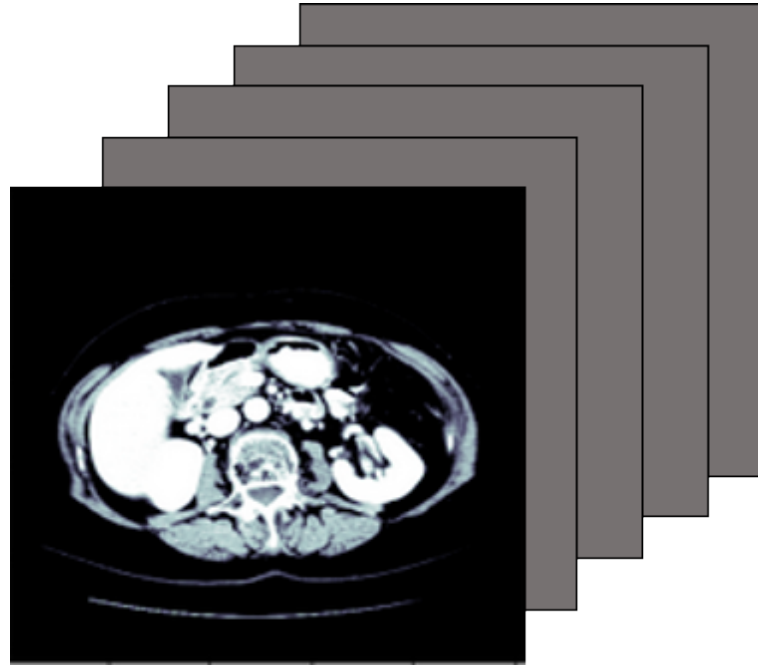- Neural networks are not naturally parallelized in time

# Time decomposition using multigrid
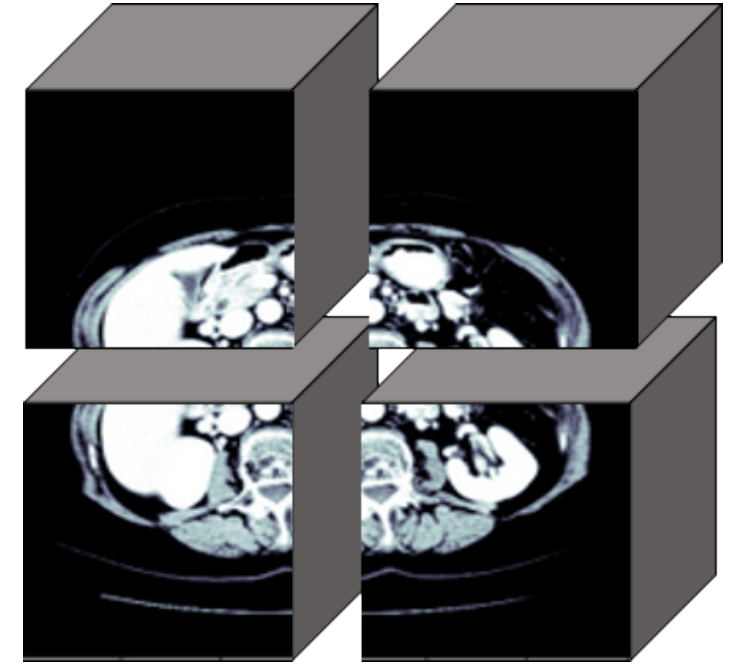


**Coarsening:** reduce number of layers

**Refinement:** use coarse output to inform

$\Delta t$     $2\Delta t$     $3\Delta t$     $4\Delta t$

$\Delta t$     $3\Delta t$

**Multigrid improves time decomposition approximations**
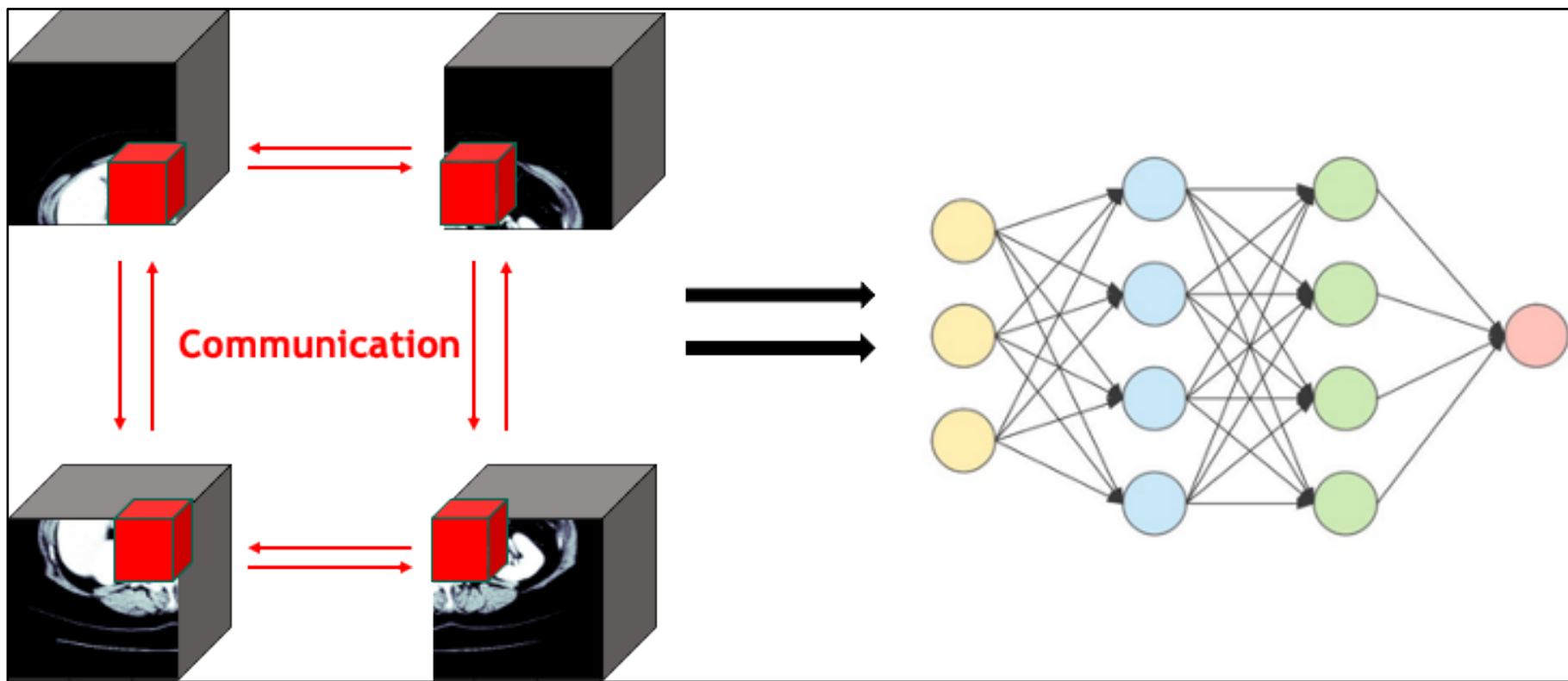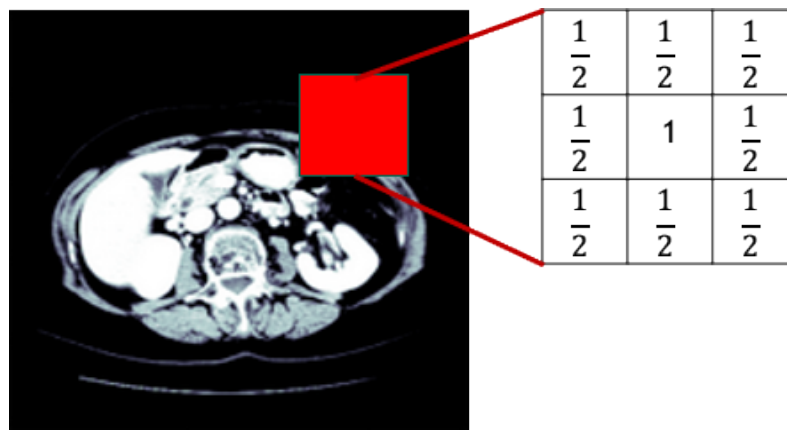
# Deep Learning and Spatial Decomposition



Original    2D slices    3D decomposition

**Medical image segmentation sees improvements when using 3D spatial decomposition.**

High Resolution Medical Image Analysis with Spatial Partitioning. Hou et al. 2019
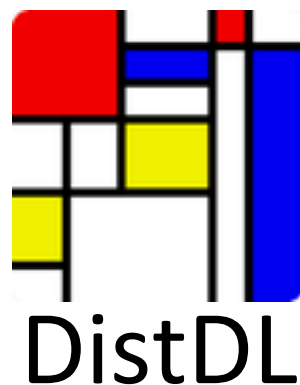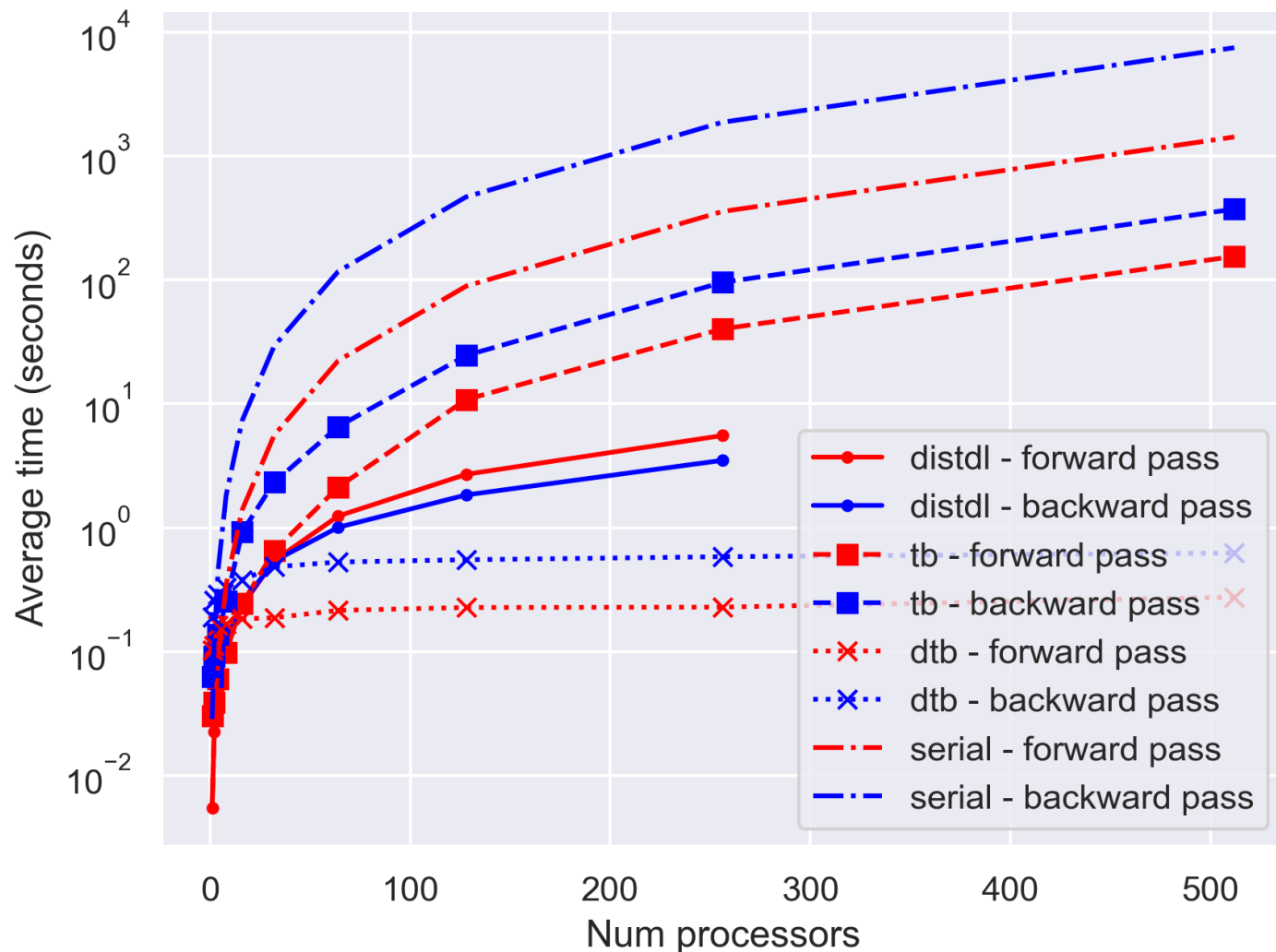
# How spatial decomposition works

# Approach – combine spatial and time parallelism

- Combine two PyTorch frameworks: **DistDL** and **TorchBRAID**.

- The combo framework is currently called **DTB.**

- Compare performance of space-time parallelism with lone space or time parallelism

- Architecture
  - Processor - 2.1 GHz Intel Broadwell E5-2695 v4 : 2 sockets : 18 cores
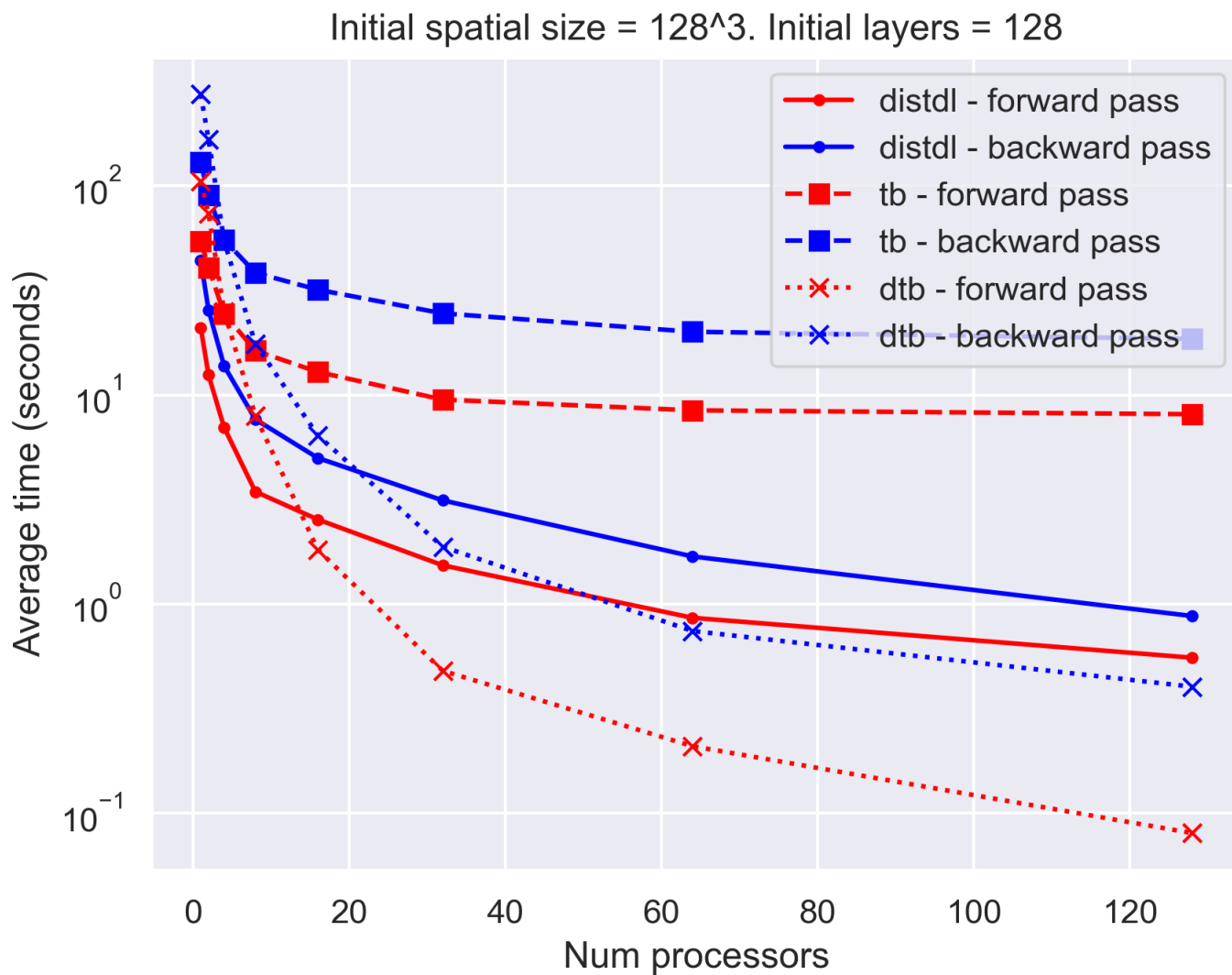  - RAM per node – 128 GB
  - 1.8 pFlops



DistDL

# Results – Convolution layers



| Problem Size per Process (Space, Layer) | | | |
|---|---|---|---|
| Procs | DistDL | TB | DTB |
| 1 | $16^3$, **8** | **$16^3$**, 8 | **$16^3$, 8** |
| 2 | $16^3$, **2x8** | **2x$16^3$**, 8 | **$16^3$, 8** |
| $2^2$ | $16^3$, **$2^2$x8** | **$2^2$x$16^3$**, 8 | **$16^3$, 8** |
| $2^3$ | $16^3$, **$2^3$x8** | **$2^3$x$16^3$**, 8 | **$16^3$, 8** |
| $2^4$ | $16^3$, **$2^4$x8** | **$2^4$x$16^3$**, 8 | **$16^3$, 8** |
| ... | ... | ... | ... |
| $2^N$ | $16^3$, **$2^N$x8** | **$2^N$x$16^3$**, 8 | **$16^3$, 8** |

## DTB exhibits weak scaling behavior while others slow down

# Results – Convolution layers



Initial spatial size = 128^3. Initial layers = 128

Legend:
- distdl - forward pass
- distdl - backward pass
- tb - forward pass
- tb - backward pass
- dtb - forward pass
- dtb - backward pass

X-axis: Num processors
Y-axis: Average time (seconds)

**DTB is up to 2 orders of magnitude faster**

# Conclusion and Future Work

- **Important takeaways**
  - Spatial + time parallelism can offer significant speedups
  - The combination enables research into deeper networks for large image segmentation
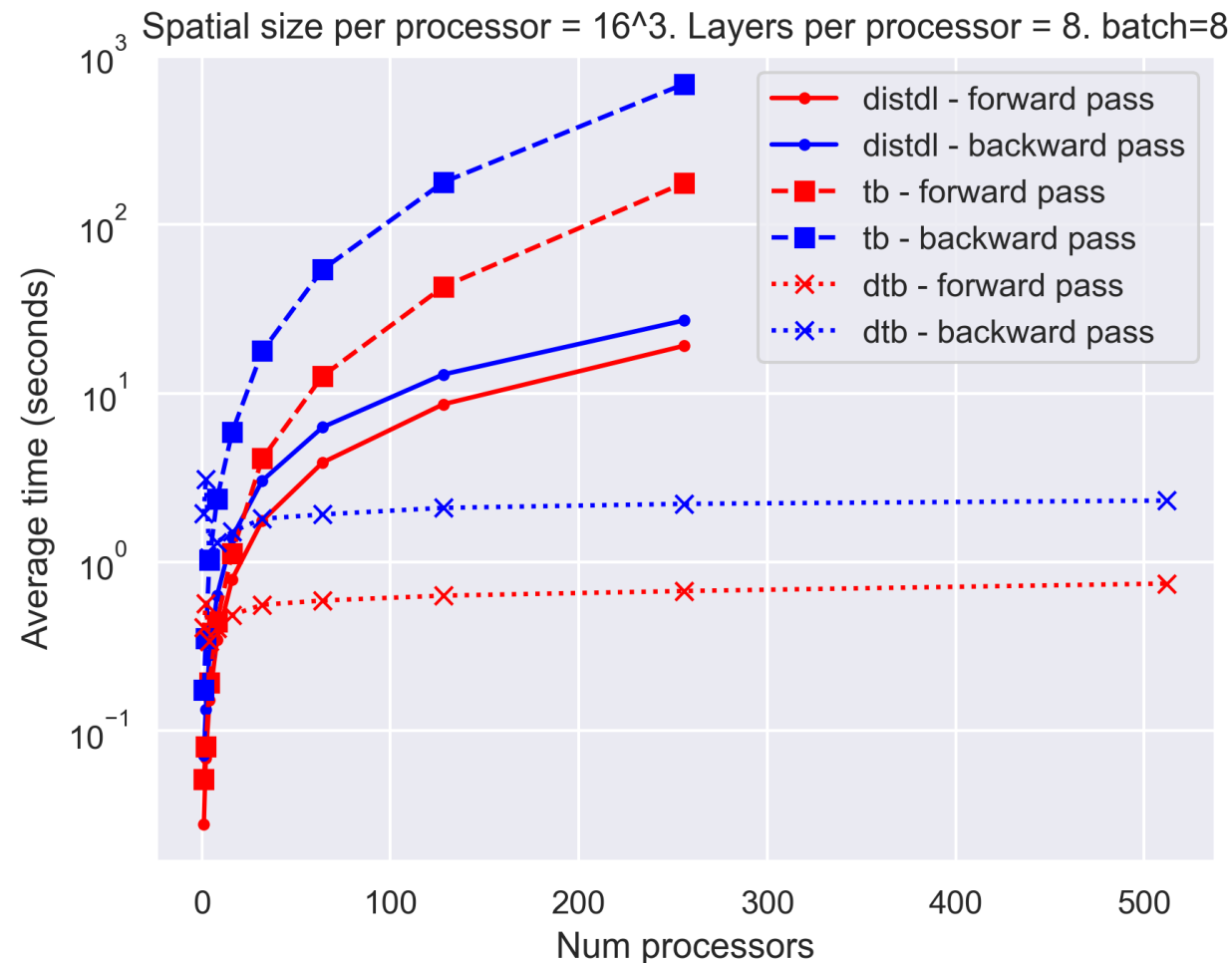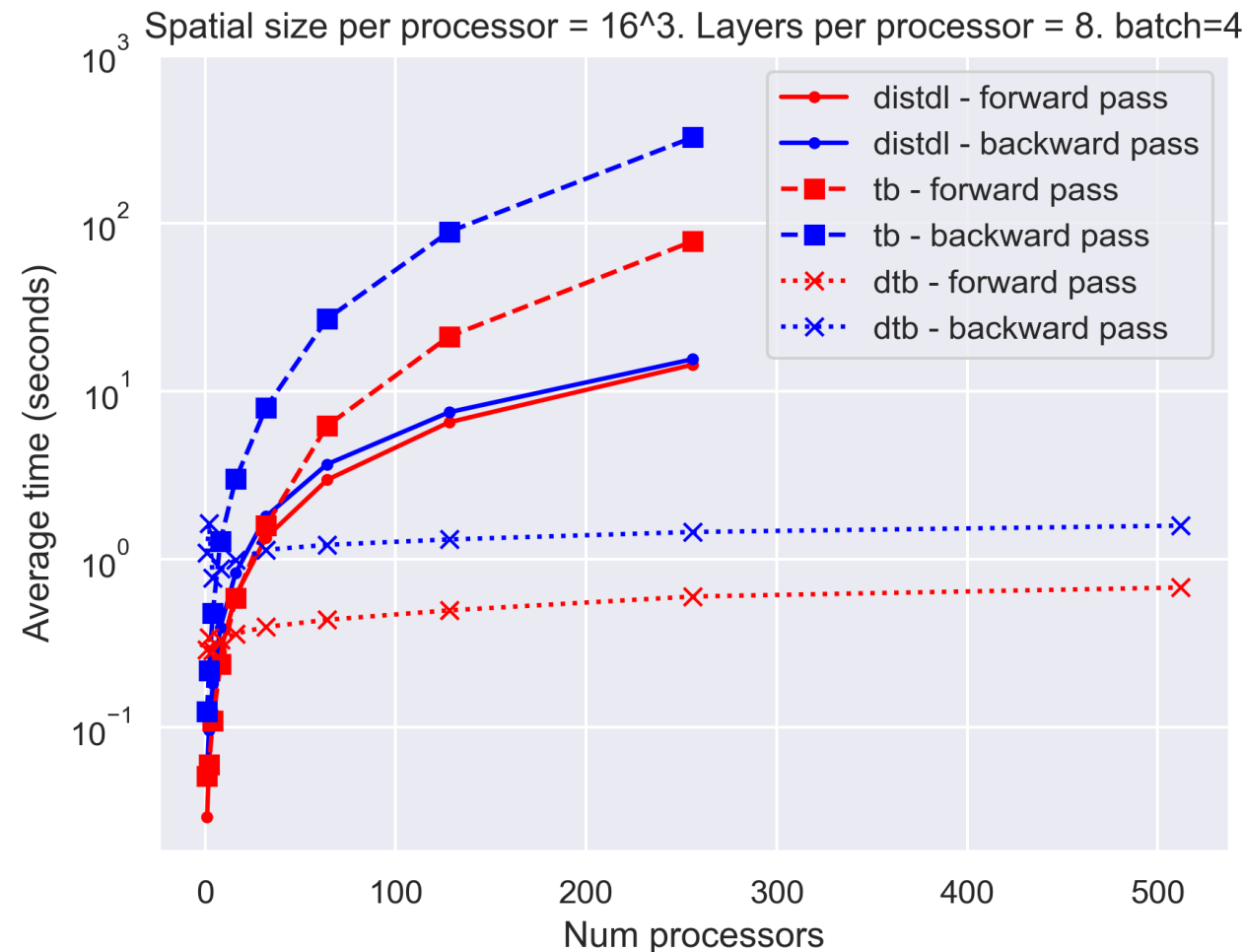
- **Future work**
  - Apply framework to real-world datasets
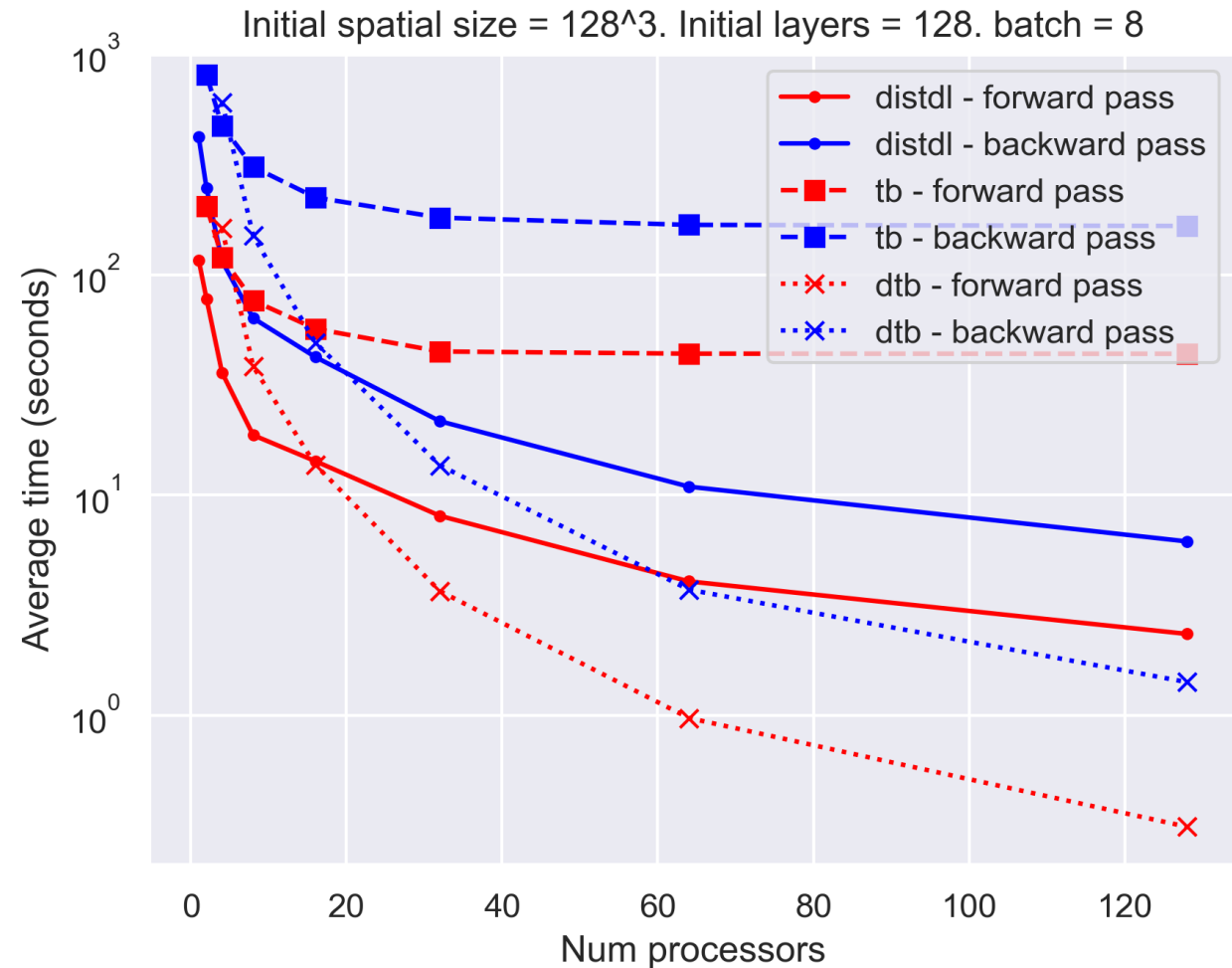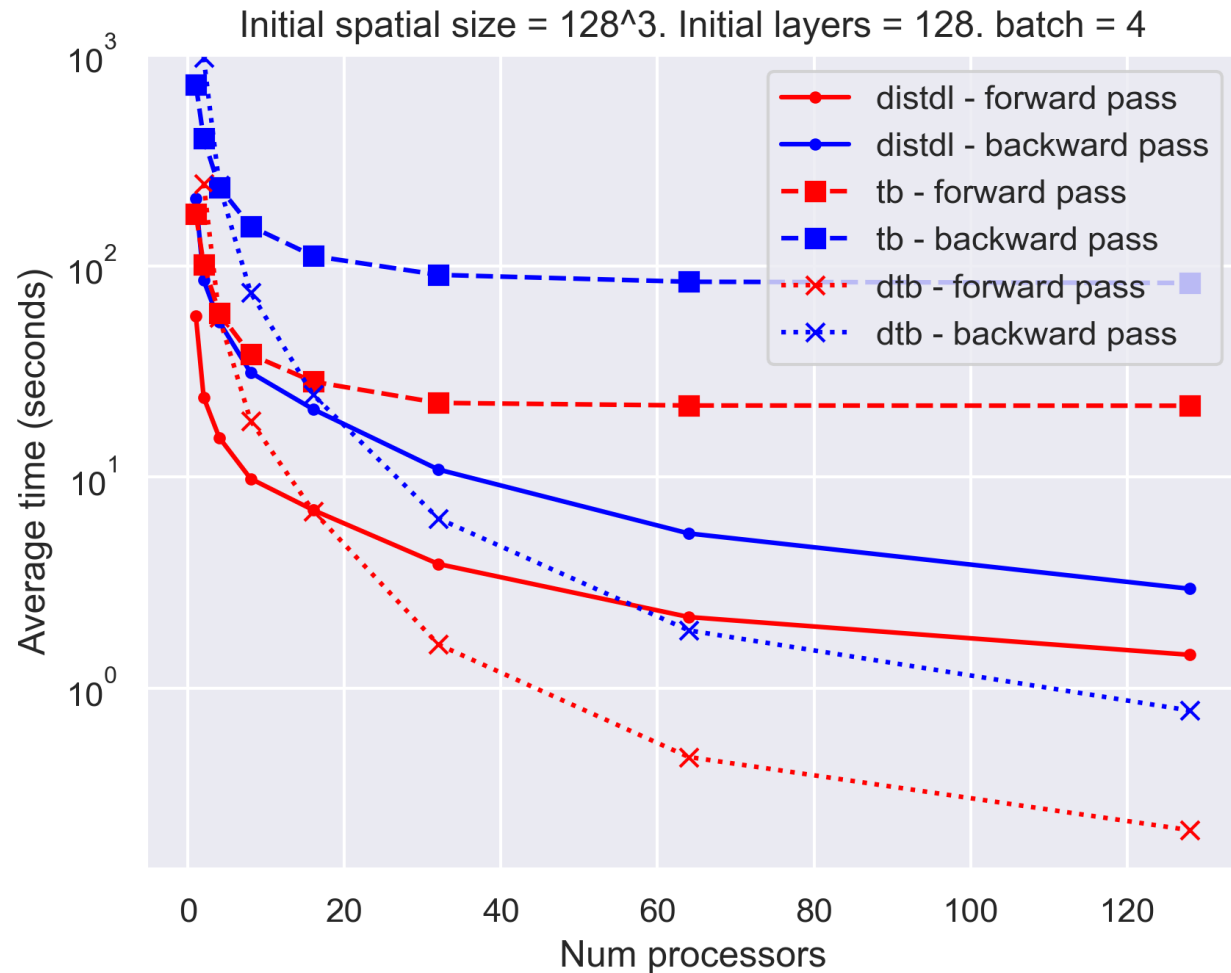  - Can we integrate spatial multigrid for further speedups?

# Backup Slides

# Results – Convolution followed by batch norm



Spatial size per processor = 16^3. Layers per processor = 8. batch=4

Spatial size per processor = 16^3. Layers per processor = 8. batch=8

**DTB exhibits weak scaling behavior while others slow down**

# Results – Convolution followed by batch norm



Initial spatial size = 128^3. Initial layers = 128. batch = 4

Initial spatial size = 128^3. Initial layers = 128. batch = 8

**DTB is up to 2 orders of magnitude faster**