



# Layer-Parallel Training with GPU Concurrency of Deep Residual Networks via Nonlinear Multigrid

September 24, 2020

Andrew C. Kirby

Siddharth Samsi, Michael Jones, Albert Reuther  
Jeremy Kepner, Vijay Gadepally

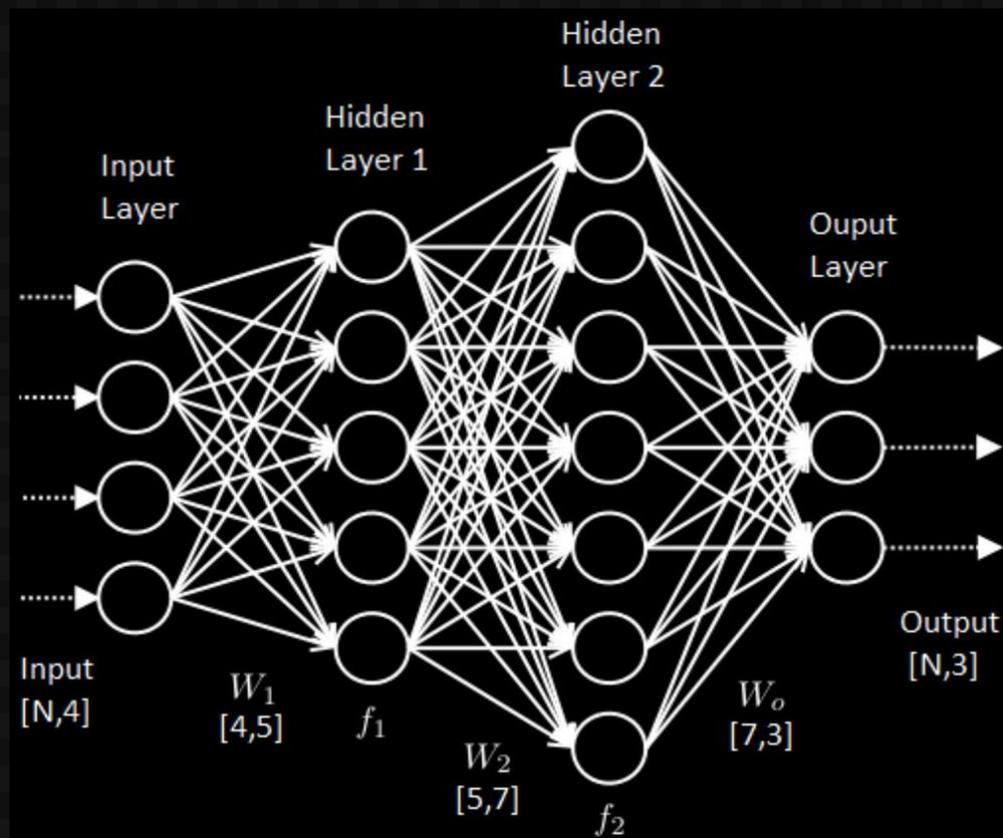
---

MIT Lincoln Laboratory Supercomputing Center

---



- Information travels sequentially through network
  - Forward propagation sets the neural states
  - Back propagation updates the weights and bias
  - Amdahl's Law
    - Theoretical speedup limited by serial execution
- Goal: minimize loss by updating the weight coefficients and bias vectors iteratively
  - Updates to unknowns are *approximate*
    - Optimization opportunities for accelerating training





# Deep Residual Networks

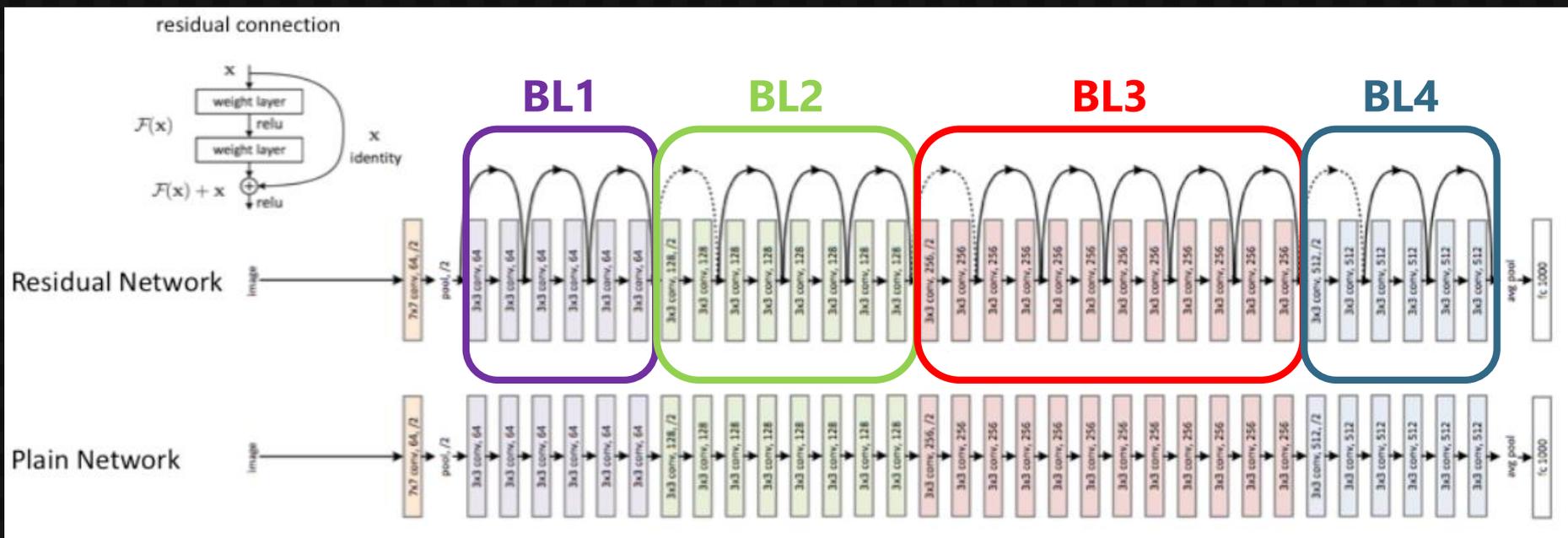
$$\mathbf{u}^{n+1} = \mathbf{u}^n + hF(\mathbf{u}^n; \boldsymbol{\theta}^n), \quad \text{for } n = 0, \dots, N - 1$$

## Parallelization Approaches

- Data-Parallel (easy – e.g. Horovod)
- Model-Parallel (hard – not *model-partitioned* serial propagation of data)

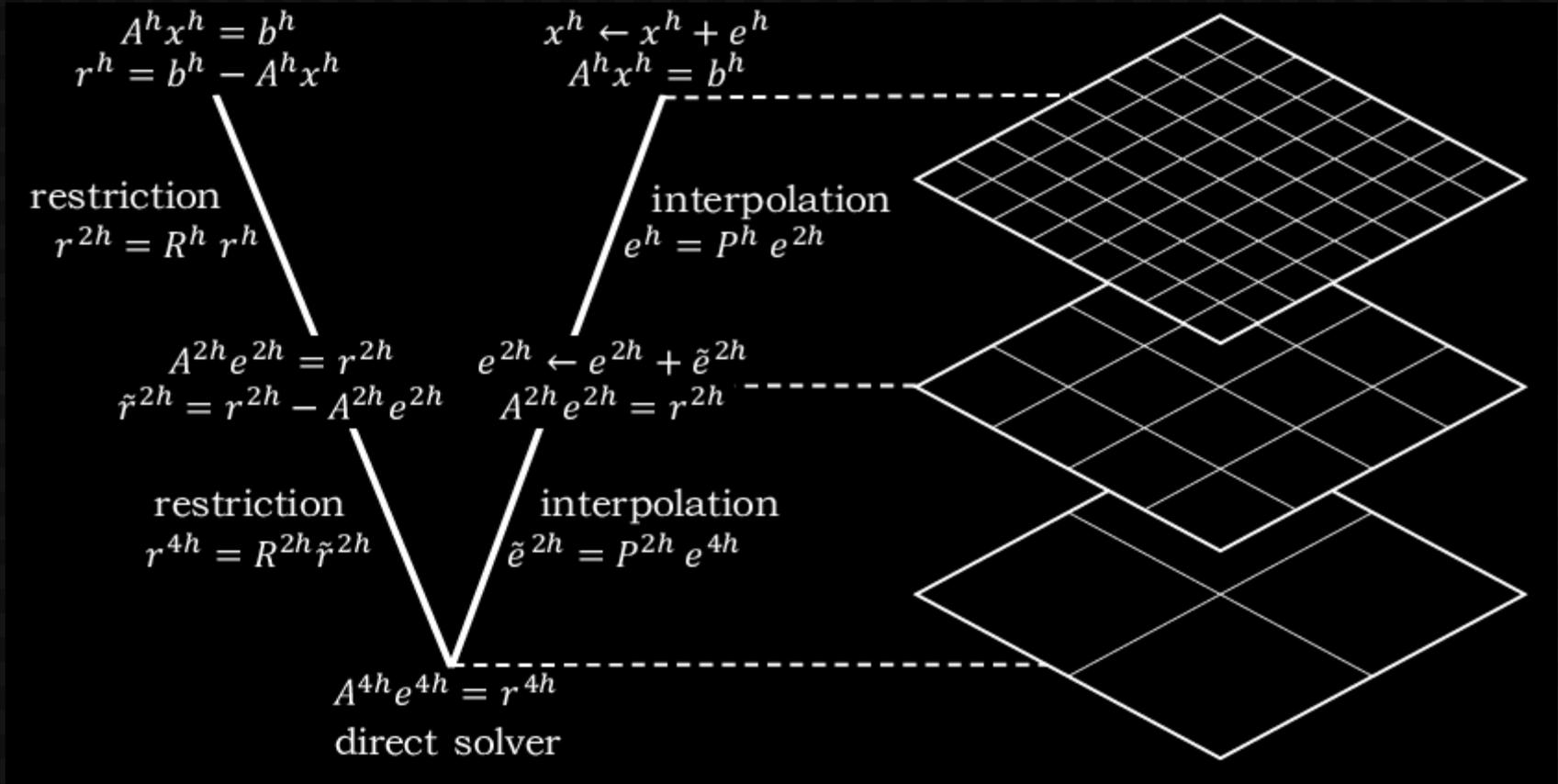
## True Model Parallelism via Iterative Approach: Solve Each Block in Parallel

- Combine with Data-Parallel Techniques for multiplicative parallelization
  - e.g. DP = 4-way parallelism, MP = 2-way parallelism, Total = 8-way parallelism





# Multigrid



[Reference] S. Gunther, et. al, "Layer-Parallel Training of Deep Residual Neural Networks", SIAM.

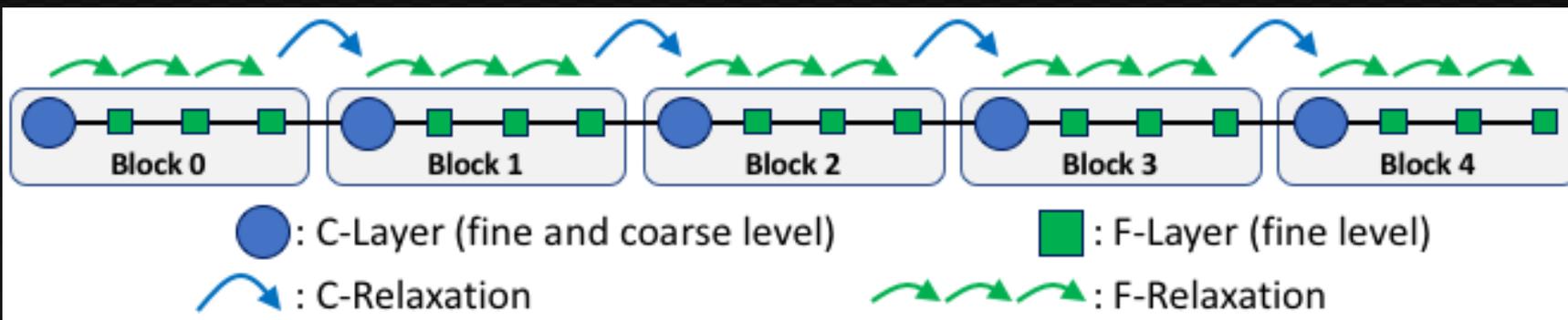
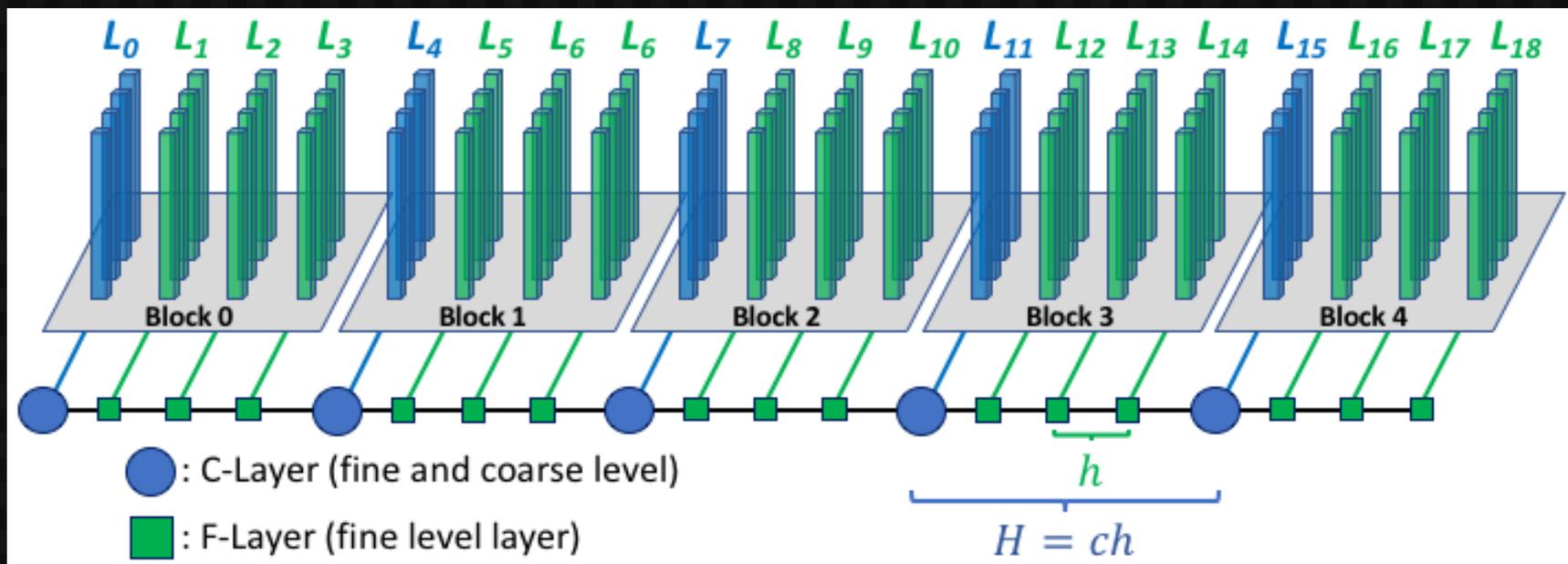
Source: [https://www.researchgate.net/figure/Illustration-of-the-multigrid-V-cycle\\_fig2\\_328599327](https://www.researchgate.net/figure/Illustration-of-the-multigrid-V-cycle_fig2_328599327)



# Network Partitioning

- Iterative Solution Procedure

- Propagate guess in block and 1<sup>st</sup> neighbor layer, restrict, solve coarse solution, update



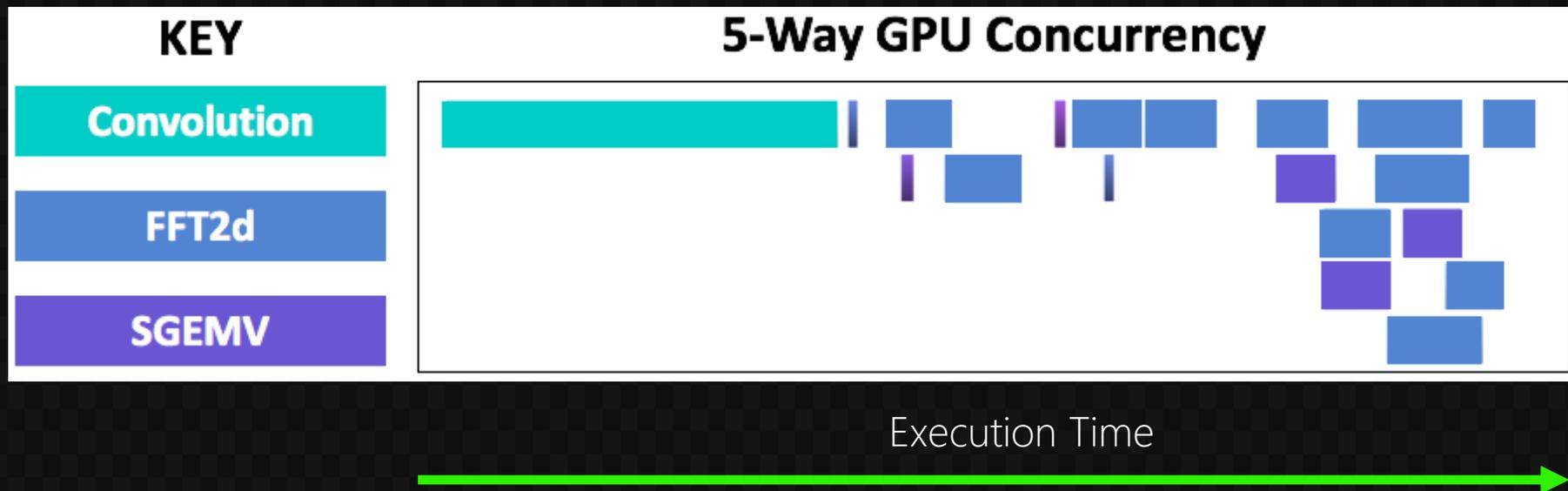




# Results: GPU Concurrency

## Implementation

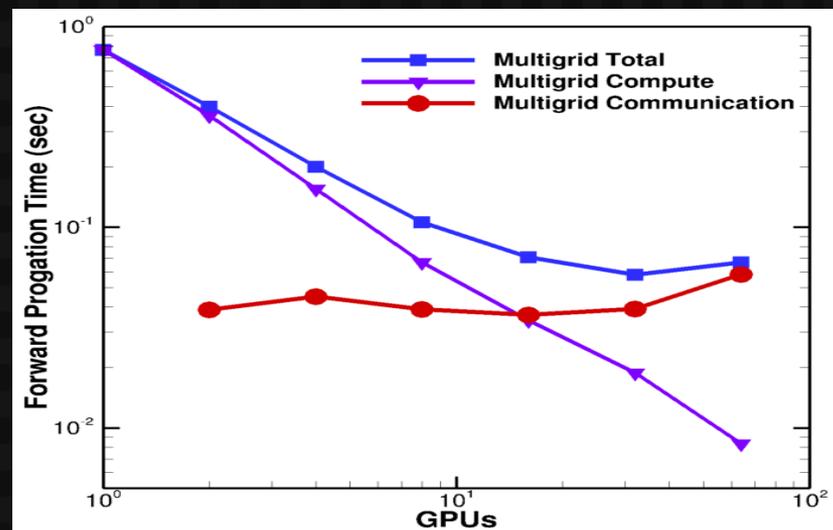
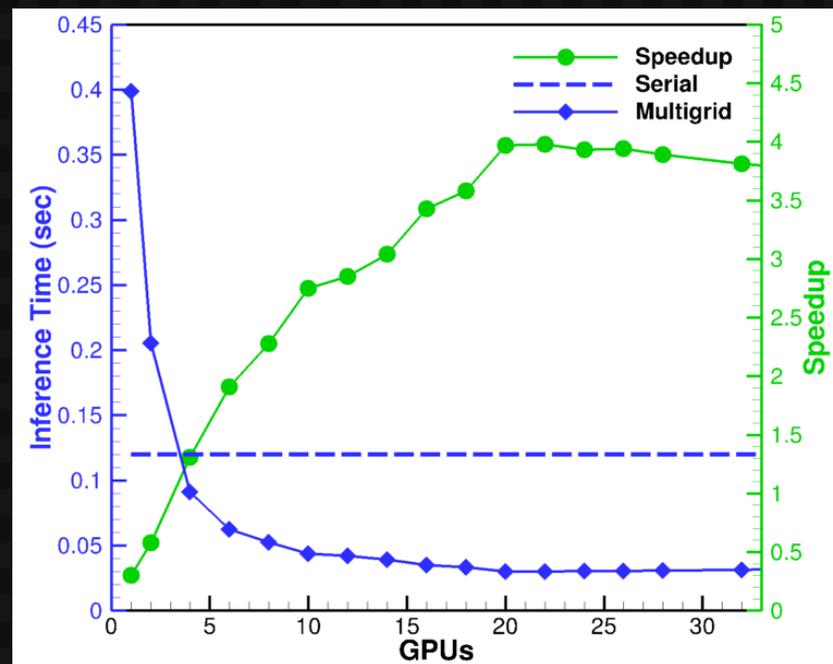
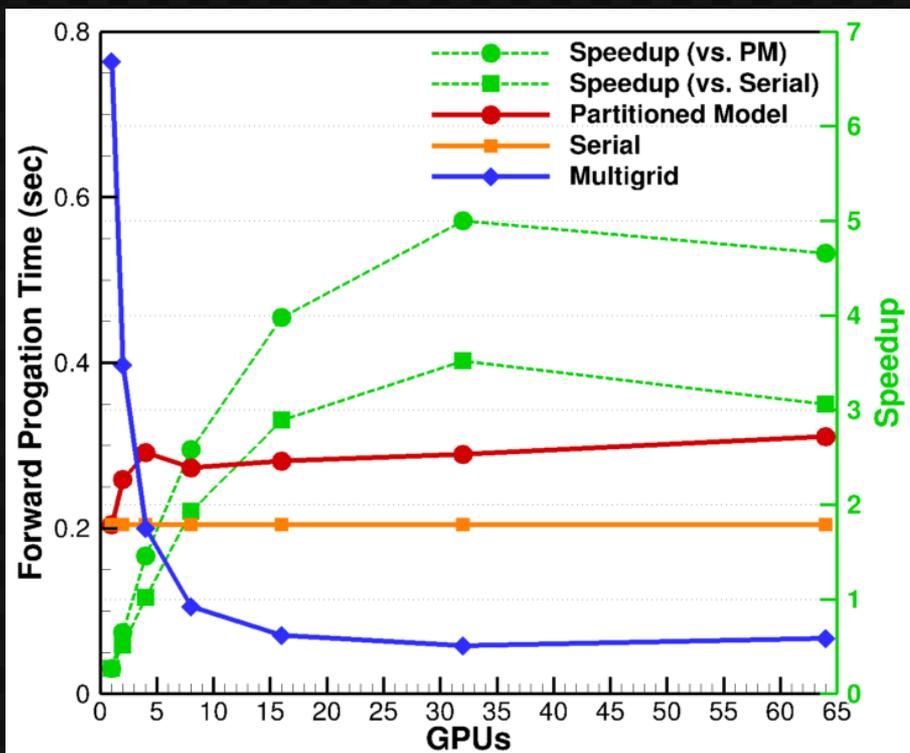
- Wrapped CuDNN kernels
  - Enabled CUDA Streams for asynchronous execution
- Placed multiple layer-blocks on same GPU – each with own CUDA Stream
  - Used OpenMP to parallel launch blocks from different CPU-threads





# Results

- MNIST Data Set
- 4,096 Layers
  - 7 x 7 convolution layers
  - 50 output channels
  - Padding size: 1
- 3,248,524 parameters



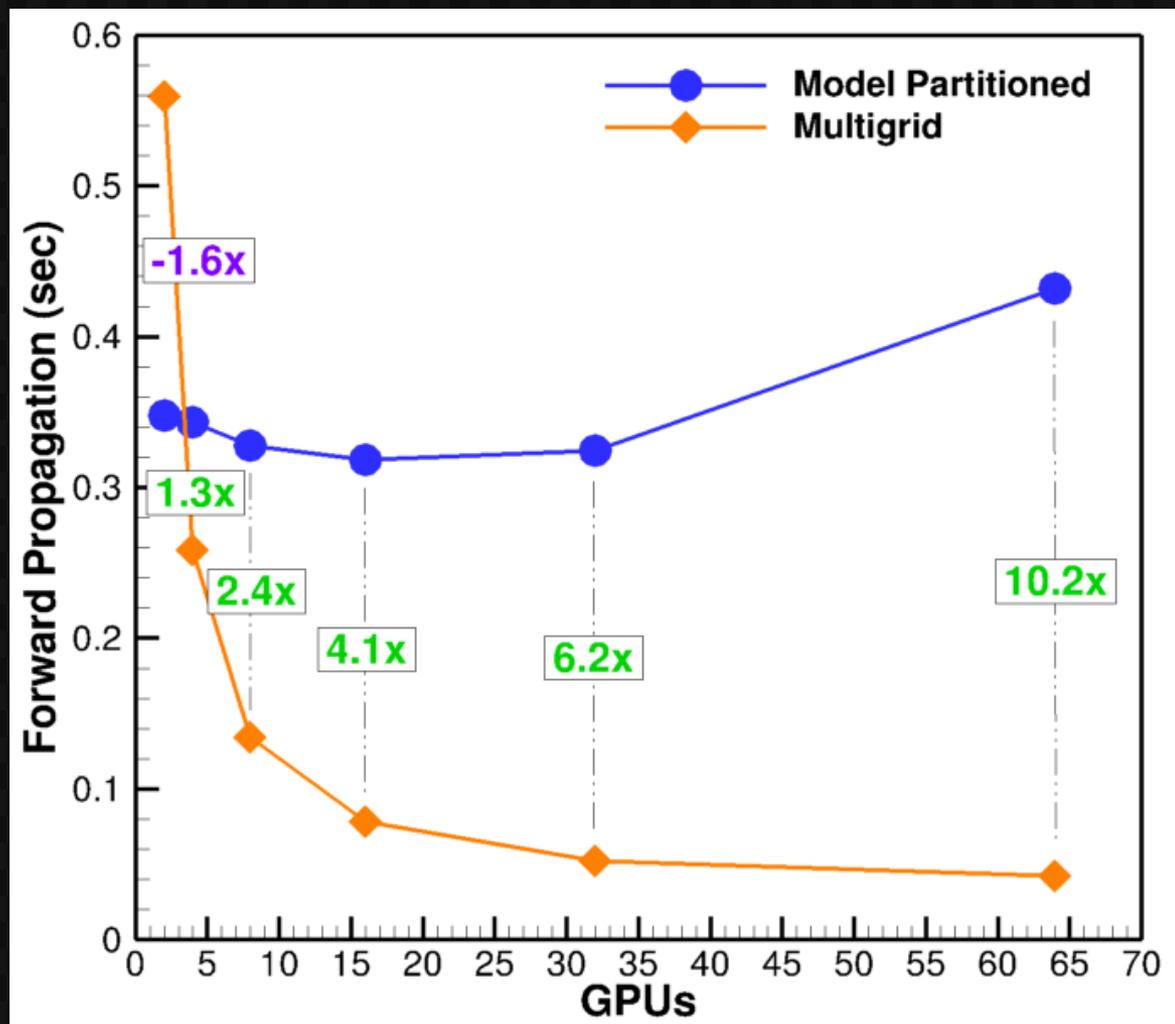


# Results

- MNIST Data Set
- 4,115 Layers
  - 7 x 7 convolution layers
  - 20 output channels
  - 16 fully connected layers
  - Padding size: 1
- 2,071,328,150 parameters

4 GPUs: Computation to communication ratio is 92.8%

64 GPUs: Computation to communication ratio is 34.5%





# Acknowledgements

## Lincoln Laboratory Leadership

Dave Martinez, Steve Rejto, Marc Zissman

## Lincoln Laboratory Supercomputing Center

William Arcand, David Bestor, William Bergeron, Chansup Byun, Matthew Hubbell, Michael Houle, Anna Klein, Peter Michaleas, Lauren Miliechin, Julie Mullen, Andrew Prout, Antonio Rosa, Charles Yee

## Compute Time

TX-Green AI Accelerator (TX-GAIA)

## Funding

MIT-Air Force AI Innovation Accelerator



**U.S. AIR FORCE**

Thank You  
Questions?