

Simultaneous Layer-Parallel Training for Deep Residual Networks

Stefanie **Günther**^{*},

Lars **Ruthotto**[†], Jacob **Schroder**[◇], Eric **Cyr**[‡], Nicolas **Gauger**^{*}

^{*}Scientific Computing Group, **TU Kaiserslautern** | Germany

[†]Dept. of Mathematics and Computer Science, **Emory University**, Atlanta | USA

[◇]Dept. of Mathematics and Statistics, **University of New Mexico**, Albuquerque | USA

[‡]Computational Mathematics Dept., **Sandia National Laboratories**, Albuquerque | USA

Feb 25, 2019

SIAM CSE, Spokane

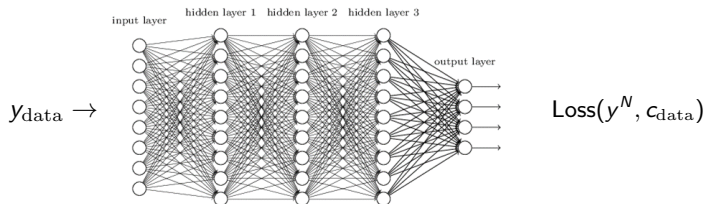
Reduce training runtime of deep residual networks!

① Enable parallelization across layers

- Iterative Multigrid-in-Time for forward and backward propagation

② Simultaneous optimization

- One-shot approach: optimization based on inexact gradients



► ResNet¹ propagation

$$\begin{aligned} y^0 &= y_{\text{data}} \\ y^{n+1} &= y^n + F(W^n y^n + b^n) \quad \forall n = 0, \dots, N-1 \quad (*) \end{aligned}$$

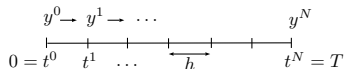
► Learning Problem

$$\min_{W^n, b^n} \text{Loss}(y^N, c_{\text{data}}) \quad \text{subject to} \quad (*)$$

¹He et. al.: *Deep Residual Learning for Image Recognition*, IEEE Conf. on CVRP, 2016

$$y^{n+1} = y^n + hF(W^n y^n + b^n) \quad \forall n$$

- Explicit Euler discretization of



$$\frac{dy(t)}{dt} = F(W(t)y(t) + b(t)) \quad \forall t \in (0, T)$$

ResNet Training \Leftrightarrow Optimal Control Problem

$$\begin{aligned} & \min_{W(t), b(t)} \text{Loss}(y(T), c_{\text{data}}) \\ \text{subject to} & \quad \frac{dy(t)}{dt} = F(W(t)y(t) + b(t)) \quad \forall t \in (0, T) \\ & \quad y(0) = y_{\text{data}} \end{aligned}$$

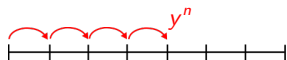
Weinan E: **A Proposal on Machine Learning via Dynamical Systems**, Comm. Math. Stats., 2017

E. Haber, L. Ruthotto: **Stable Architectures for Deep Neural Networks**, Inv. Probl., 2017

► Iterative updates for $\theta^n := (W^n, b^n)$

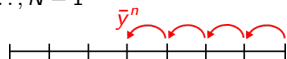
1. Solve ODE \leftrightarrow **Forward propagation**

$$y^0 = y_{\text{data}}$$
$$y^{n+1} = \Phi_h(y^n, \theta^n) \quad \text{for } n = 0, \dots, N-1$$



2. Solve adjoint ODE \leftrightarrow **Backpropagation**

$$\bar{y}^N = \partial_{y^N} \text{Loss}(y^N, c)$$
$$\bar{y}^n = \partial_{y^n} \Phi_h(y^n, \theta^n)^T \bar{y}^{n+1} \quad \text{for } n = N-1, \dots, 0$$



3. Network parameter **update**

$$\theta^n \leftarrow \theta^n - \alpha \left(\partial_{\theta^n} \Phi_h(y^n, \theta^n)^T \bar{y}^{n+1} \right) \quad \forall n$$

$$\begin{pmatrix} y^0 \\ y^1 - \Phi_h(y^0, \theta^0) \\ y^2 - \Phi_h(y^1, \theta^1) \\ \vdots \\ y^N - \Phi_h(y^{N-1}, \theta^{N-1}) \end{pmatrix} \stackrel{!}{=} \begin{pmatrix} y_{\text{data}} \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

► Layer-serial propagation



time step 1

► Layer-parallel multigrid



grid level 0



grid level 1



grid level 2

¹Pictures taken from S. Friedhoff, talk on *Multigrid reduction techniques for parallel-in-time integration*, 2016

$$\begin{pmatrix} y^0 \\ y^1 - \Phi_h(y^0, \theta^0) \\ y^2 - \Phi_h(y^1, \theta^1) \\ \vdots \\ y^N - \Phi_h(y^{N-1}, \theta^{N-1}) \end{pmatrix} \stackrel{!}{=} \begin{pmatrix} y_{\text{data}} \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

► Layer-serial propagation



time step 2

► Layer-parallel multigrid



grid level 0



grid level 1



grid level 2

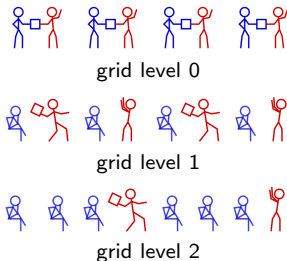
¹Pictures taken from S. Friedhoff, talk on *Multigrid reduction techniques for parallel-in-time integration*, 2016

$$\begin{pmatrix} y^0 \\ y^1 - \Phi_h(y^0, \theta^0) \\ y^2 - \Phi_h(y^1, \theta^1) \\ \vdots \\ y^N - \Phi_h(y^{N-1}, \theta^{N-1}) \end{pmatrix} \stackrel{!}{=} \begin{pmatrix} y_{\text{data}} \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

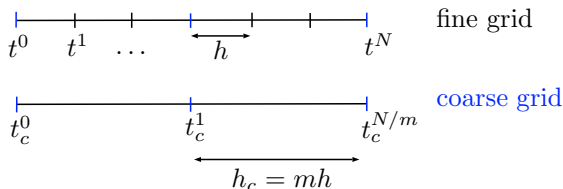
► Layer-serial propagation



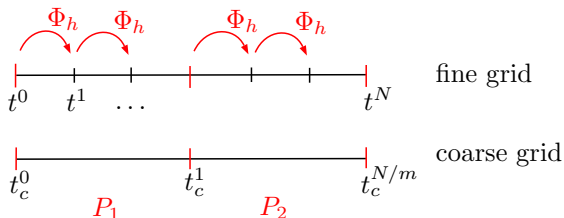
► Layer-parallel multigrid



¹Pictures taken from S. Friedhoff, talk on *Multigrid reduction techniques for parallel-in-time integration*, 2016



- ▶ Nonlinear multigrid scheme (FAS):
 1. Approximate y^1, \dots, y^N on the fine grid
 2. Solve residual equation on the coarse grid
 3. Correct y^1, \dots, y^N on the fine grid



► Nonlinear multigrid scheme (FAS):

1. Approximate y^1, \dots, y^N on the fine grid

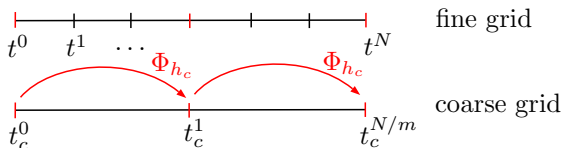
► In each interval: $y^{n+1} = \Phi_h(y^n, \theta^n)$

► Residual: $r^{mn} = y^{mn} - y^{mn-1}$

2. Solve residual equation on the coarse grid

3. Correct y^1, \dots, y^N on the fine grid

► parallel



▶ Nonlinear multigrid scheme (FAS):

1. Approximate y^1, \dots, y^N on the fine grid

▶ In each interval: $y^{n+1} = \Phi_h(y^n, \theta^n)$

▶ Residual: $r^{mn} = y^{mn} - y^{mn-1}$

▶ parallel

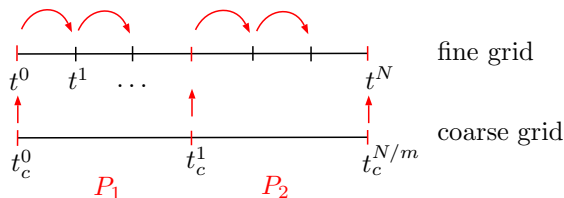
2. Solve residual equation on the coarse grid

▶ $\tilde{y}_c^{n+1} = \Phi_{h_c}(\tilde{y}_c^n, \theta_c^n) - r_c^n$

▶ Error approximation $e_c^n = y_c^n - \tilde{y}_c^n$

▶ serial on coarse grid

3. Correct y^1, \dots, y^N on the fine grid



▶ Nonlinear multigrid scheme (FAS):

1. Approximate y^1, \dots, y^N on the fine grid

- ▶ In each interval: $y^{n+1} = \Phi_h(y^n, \theta^n)$
- ▶ Residual: $r^{mn} = y^{mn} - y^{mn-1}$

▶ parallel

2. Solve residual equation on the coarse grid

- ▶ $\tilde{y}_c^{n+1} = \Phi_{h_c}(\tilde{y}_c^n, \theta_c^n) - r_c^n$
- ▶ Error approximation $e_c^n = y_c^n - \tilde{y}_c^n$

▶ serial on coarse grid

3. Correct y^1, \dots, y^N on the fine grid

▶ parallel

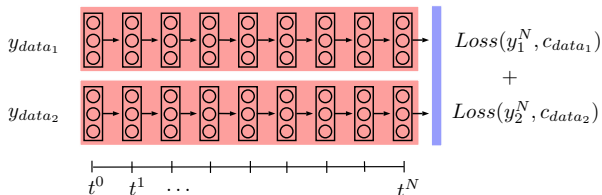
- ▶ Fixed-point iteration for $\mathbf{y} = (y^1, \dots, y^N)$:

$$\mathbf{y}_{k+1} = H(\mathbf{y}_k, \theta)$$

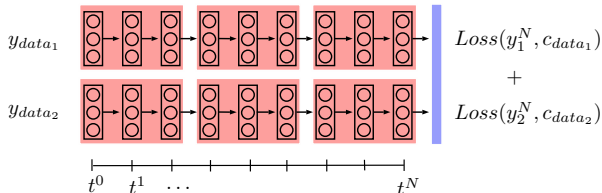
recovers at convergence same output y^N as serial propagation

- ▶ But greater concurrency \rightarrow cross over point:
Parallel speedup expected if N and n_{cores} are large

► Data-parallelization



► Additional Layer-parallelization





<https://github.com/xbraid>

- ▶ Open-source software library
- ▶ Flexible user-interface for existing time-stepping codes

my_Step: $y^{n+1} = \Phi_h(y^n, \theta^n)$

my_Objective: $\text{Loss}(y^n, \theta^n)$

my_Clone: $v^n = y^n$

my_Sum: $y^n = \alpha v^n + \beta y^n$

my_SpatialNorm: $\|y^n\|$

...

- ▶ **forward propagation** \Leftrightarrow parallel multigrid for $\mathbf{y} = (y^0, \dots, y^N)$

$$y^0 = y_{\text{data}}$$
$$y^{n+1} = \Phi(y^n, \theta^n) \quad \forall n = 0, \dots, N-1$$

$$\Leftrightarrow \mathbf{y}_{k+1} = H(\mathbf{y}_k, \boldsymbol{\theta}) \quad k = 1, 2, \dots$$

- ▶ **backpropagation** \Leftrightarrow parallel multigrid for $\bar{\mathbf{y}} = (\bar{y}^N, \dots, \bar{y}^0)$

$$\bar{y}^N = \partial_{y^N} \text{Loss}(y^N, c_{\text{data}})$$
$$\bar{y}^n = \partial_y \Phi(y^n, \theta^n)^T \bar{y}^{n+1} \quad \forall n = N-1, \dots, 0$$

$$\Leftrightarrow \bar{\mathbf{y}}_{k+1} = H(\bar{\mathbf{y}}_k, \mathbf{y}, \boldsymbol{\theta}) \quad k = 1, 2, \dots$$

► Iterative updates for θ :

1. Layer-parallel forward propagation:

$$\text{For } k = 1, 2, \dots : \quad \mathbf{y}_{k+1} = H(\mathbf{y}_k, \theta) \xrightarrow{k \rightarrow \infty} \mathbf{y}_*$$

$$\Rightarrow \text{Loss}(\mathbf{y}_*^N, \mathbf{c}_{\text{data}})$$

2. Layer-parallel backpropagation:

$$\text{For } k = 1, 2, \dots : \quad \bar{\mathbf{y}}_{k+1} = H(\bar{\mathbf{y}}_k, \mathbf{y}_*, \theta) \xrightarrow{k \rightarrow \infty} \bar{\mathbf{y}}_*$$

$$\Rightarrow \nabla_{\theta^n} \text{Loss} = \partial_{\theta^n} \Phi(\mathbf{y}_*^n, \theta^n)^T \bar{\mathbf{y}}_*^n \quad \forall n$$

3. Network parameter update:

$$\theta^n \leftarrow \theta^n - \alpha \nabla_{\theta^n} \text{Loss} \quad \forall n$$

► Iterative updates for θ :

1. Layer-parallel forward propagation:

$$\text{For } k = 1, 2, \dots: \mathbf{y}_{k+1} = H(\mathbf{y}_k, \theta) \quad \text{with } k \rightarrow \infty \text{ and } \mathbf{y}_k$$

$$\Rightarrow \text{Loss}(\mathbf{y}_k^N, \mathbf{c}_{\text{data}})$$

2. Layer-parallel backpropagation:

$$\text{For } k = 1, 2, \dots: \bar{\mathbf{y}}_{k+1} = H(\bar{\mathbf{y}}_k, \mathbf{y}_{k+1}, \theta) \quad \text{with } k \rightarrow \infty \text{ and } \bar{\mathbf{y}}_k$$

$$\Rightarrow \nabla_{\theta^n} \text{Loss}_{k+1} = \partial_{\theta^n} \Phi(\mathbf{y}_{k+1}^n, \theta^n)^T \bar{\mathbf{y}}_{k+1}^n \quad \forall n$$

3. Network parameter **update**:

$$\theta^n \leftarrow \theta^n - \alpha \nabla_{\theta^n} \text{Loss}_{k+1} \quad \forall n$$

- ▶ Convergence theory: One-shot method²

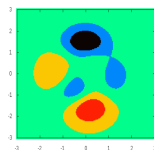
$$\theta^n \leftarrow \theta^n - \alpha H_k^{-1} \left(\frac{d\text{Loss}_k}{d\theta^n} \right)$$

$$\text{with } H \approx \partial_{\theta}^2 (L + \alpha \|\mathbf{y} - H(\mathbf{y}, \theta)\|^2)$$

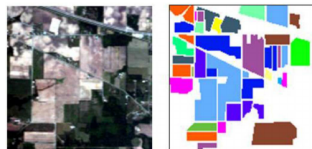
- ▶ Numerically: limited memory BFGS approximation, Identity, ...

²[Griewank, Hamdi (2011)], [Gauger, Schulz et al. (2008)], [Blommert (2016)]

1. Level set classification:
5000 grid points in $[-3, 3]^2$, 5 classes (level sets)



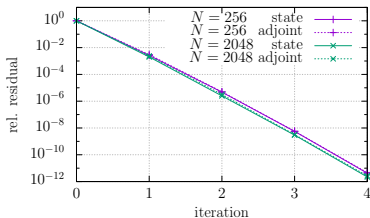
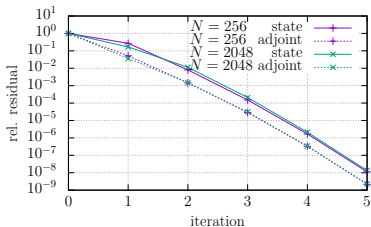
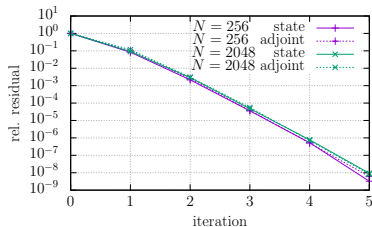
2. Hyperspectral image segmentation for Indian Pines data set:
145 × 145 pixels each with 220 spectral reflectance bands, 16 classes (land cover)



3. MNIST image classification:
28 × 28 grey scale images, 10 classes (digits)

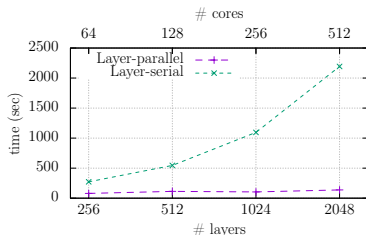
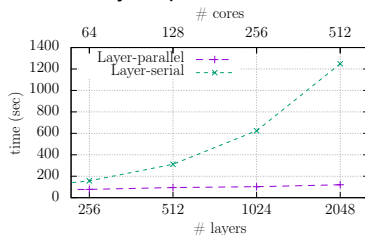
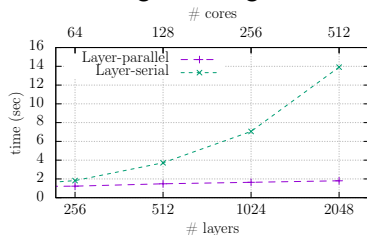


3 6 8 1 7 9 6 6 4 1
6 7 5 7 8 6 3 4 8 5
2 1 7 9 7 1 2 3 4 5
4 8 1 9 0 1 8 8 9 4
7 6 1 8 6 4 1 5 6 0
7 5 9 2 6 5 8 1 9 7
2 2 2 2 2 3 4 8 0
0 2 3 8 0 7 3 8 5 7
0 1 4 6 4 6 0 2 4 3
7 1 2 8 9 6 9 8 6 1



- Fast multigrid convergence, independent of the number of layers

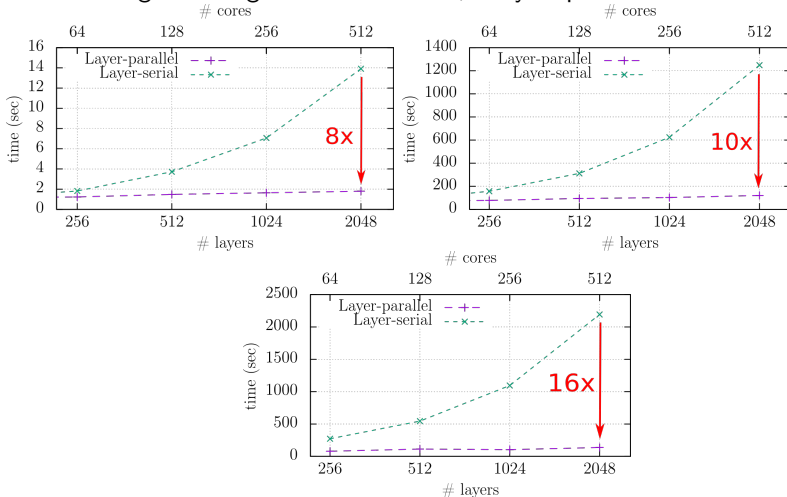
Weak scaling for one gradient evaluation, 4 layers per core



- Constant runtimes for increasing problem sizes and computational resources.

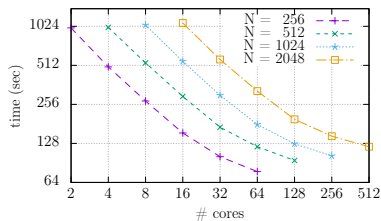
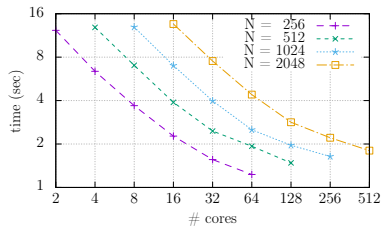
Layer-parallel Multigrid Performance

Weak scaling for one gradient evaluation, 4 layers per core



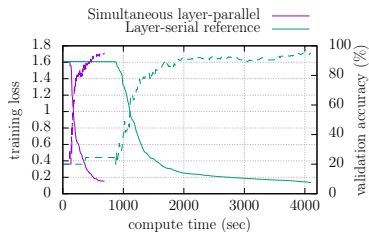
- Constant runtimes for increasing problem sizes and computational resources.

Strong scaling

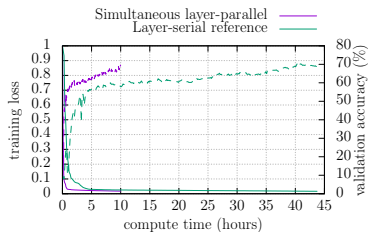


- Cross over point: ≈ 16 cores

2 multigrid cycles per optimization iteration



Level sets
 $N = 1024$



Hyperspectral images
 $N = 512$

① Deep Learning Meets Optimal Control

- Residual network training \Leftrightarrow optimal control problem

② Multigrid-in-Time for Parallelizing Across Layers

- Nonlinear multigrid iterations to solve the network propagation
- Runtime reduction from **added concurrency** across layers

③ Simultaneous Optimization

- Early stopping of multigrid iterations
- Runtime reduction from **inexact gradient** evaluation

① Deep Learning Meets Optimal Control

- Residual network training \Leftrightarrow optimal control problem

② Multigrid-in-Time for Parallelizing Across Layers

- Nonlinear multigrid iterations to solve the network propagation
- Runtime reduction from **added concurrency** accross layers

③ Simultaneous Optimization

- Early stopping of multigrid iterations
- Runtime reduction from **inexact gradient** evaluation

Thank you! Questions?