# Decoupled Neural Interfaces using Synthetic Gradients
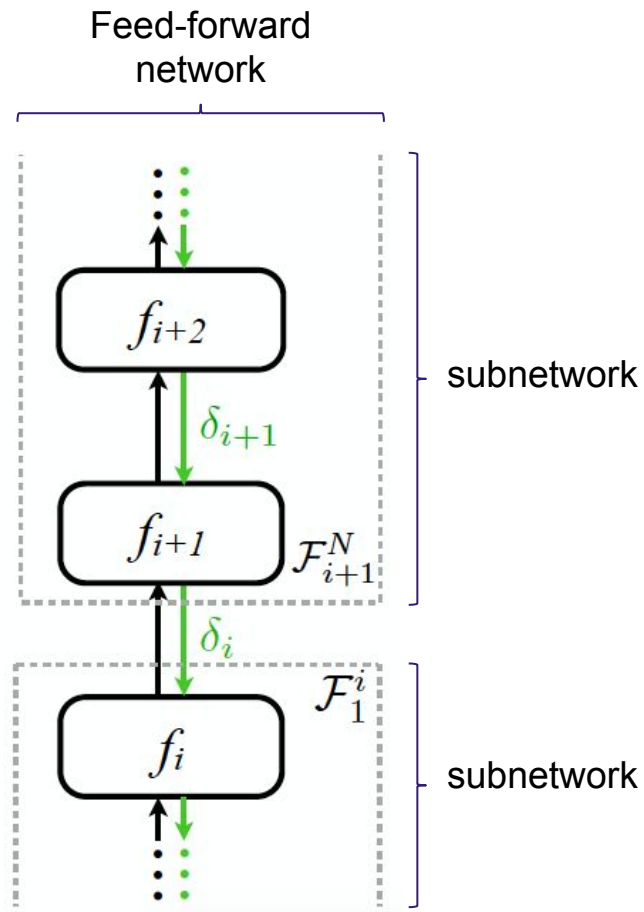
Max Jaderberg, Wojciech Marian Czarnecki, Simon Osindero,
Oriol Vinyals, Alex Graves, Koray Kavukcuoglu

**BOSTON UNIVERSITY**
**MACHINE INTELLIGENCE**
**COMMUNITY**

Justin Chen
Nov. 6, 2017

# Locking

- **Locking** - Wait for all dependent computations before parameter update at current module
- **Forward Locking** - Wait for entire forward pass
- **Update Locking** - Wait for forward pass and compute cost function
- **Backwards Locking** - Wait for forward pass and dependent backwards pass

# Backpropagation Expanded

$$\frac{\partial L}{\partial \theta_i} = f_{BProp}((h_i, x_i, y_i, \theta_i), (h_{i+1}, x_{i+1}, y_{i+1}, \theta_{i+1}), ...)\frac{\partial h_i}{\partial \theta_i} \simeq \hat{f}_{BProp}(h_i)\frac{\partial h_i}{\partial \theta_i}$$

- **Goal**: Remove **all locking**
- **Main contribution**: Backward unlocking for infinite time with RNNs
- Expanding BP exposes dependencies
- **Approximate gradient of activations** w.r.t. parameters of **adjacent layers** (modules depending on how you divide up the computation graph)
- Update **Synthetic Gradient Model** with True Gradients

M I C

# Synthetic Gradient Model

$$\hat{\delta}_A = M_B(h_A, s_B, c)$$

$\hat{\delta}_A$ — Synthetic gradient **for previous layer/module**

$M_B$ — Synthetic Gradient Model

$h_A$ — Activation of **previous layer/module**, which we refer to as module A

$s_B$ — Activation of **next layer/module**, which we refer to as module B

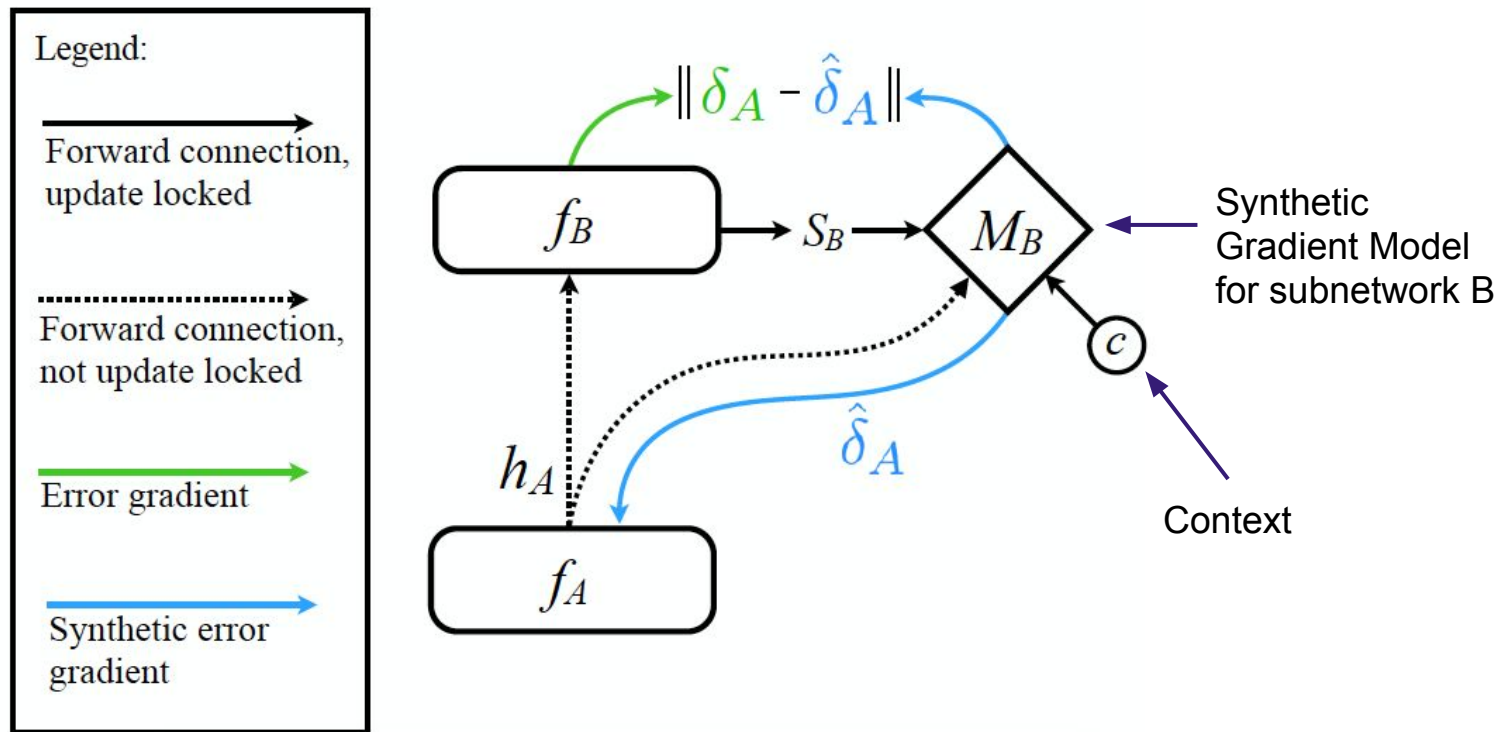$c$ — Context information - additional information such as label information if it's available

$\delta_A$ — **True gradient** from backpropagation

MIC

# Decoupled Neural Interface (DNI)



Legend:

⟶ Forward connection, update locked

⟶ (dotted) Forward connection, not update locked

⟶ (green) Error gradient

⟶ (blue) Synthetic error gradient

$\| \delta_A - \hat{\delta}_A \|$

$f_B \rightarrow S_B \rightarrow M_B$

$h_A$

$f_A$

$\hat{\delta}_A$

$c$

Synthetic Gradient Model for subnetwork B

Context

# Updating DNI

- DNI are update locked
- Trained to minimize $L_2$
- Can incorporate context information if available

$$L_{\delta_i} = d(\hat{\delta}_i, \delta_i)$$

$$\hat{\delta}_i = M_{i+1}(h_i, c)$$
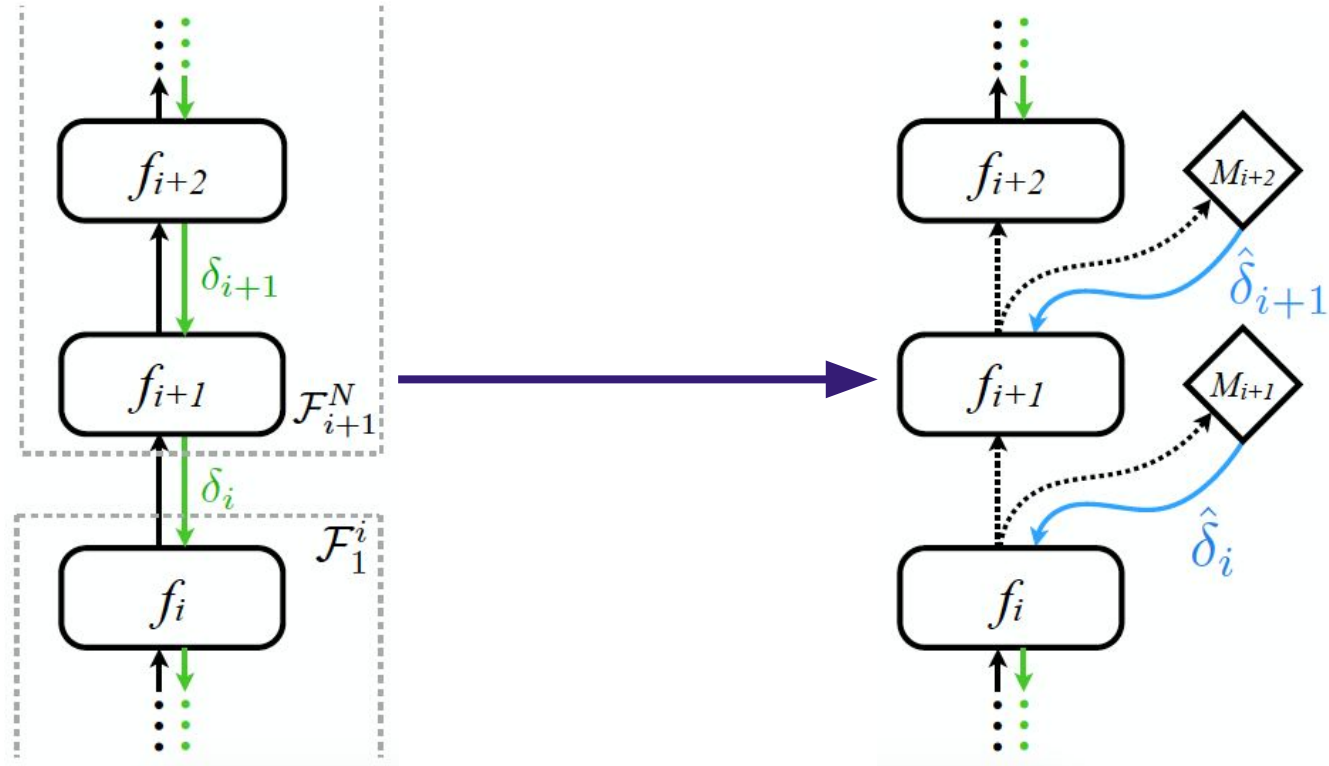
# Unlocked Subnetworks

# Expanding Backpropagation

$$\theta_i \leftarrow \theta_i - \alpha \delta_i \frac{\partial h_i}{\partial \theta_i}; \delta_i = \frac{\partial L}{\partial h_i}$$

Parameters
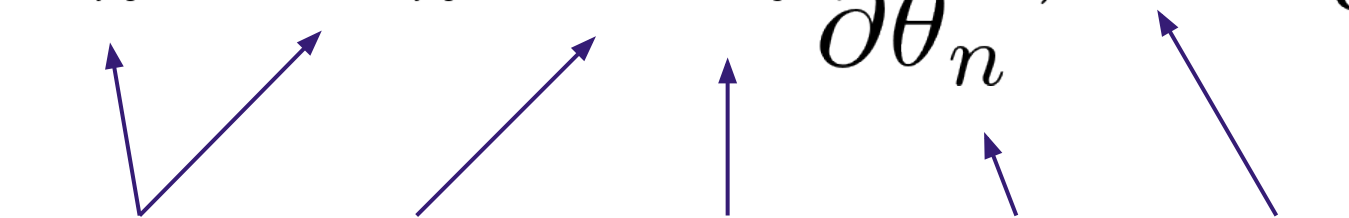
Learning Rate

Local Gradient
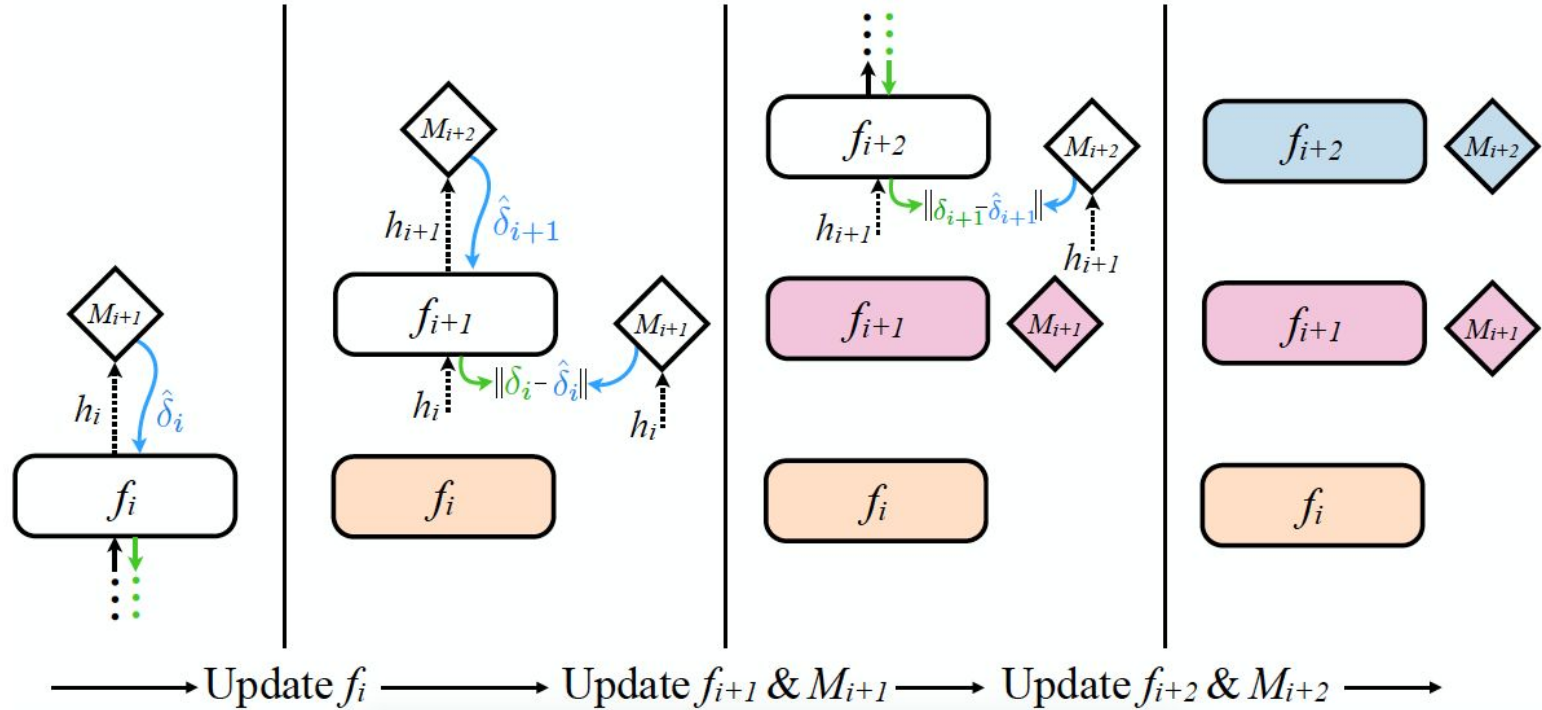
True Gradient

# Multiple DNI Feedforward Network

# Multiple DNI Update Rule

$$\theta_n \leftarrow \theta_n - \alpha \hat{\delta}_i \frac{\partial h_i}{\partial \theta_n}, n \in \{1, ..., i\}$$

Parameters

Learning
Rate

Synthetic
Gradient

Local
Gradient

Layers

# Updating Multiple DNI Feedforward Networks

# Expanding BP for RNN

$$\theta - \alpha \sum_{\tau=t}^{\infty} \frac{\partial L_\tau}{\partial \theta}$$

Parameters

Learning Rate

Sum over infinite timesteps

Gradient at timestep $\tau$

# Expanding BP for RNN

$$= \theta - \alpha \left( \underbrace{\sum_{\tau=t}^{t+T} \frac{\partial L_\tau}{\partial \theta}} + \left( \underbrace{\sum_{\tau=T+1}^{\infty} \frac{\partial L_\tau}{\partial h_T}} \right) \frac{\partial h_T}{\partial \theta} \right)$$

Unroll for only
T timesteps

Timesteps after T
(future timesteps)

# Expanding BP for RNN

$$= \theta - \alpha \left( \sum_{\tau=t}^{t+T} \frac{\partial L_\tau}{\partial \theta} + \delta_T \frac{\partial h_T}{\partial \theta} \right)$$

- Calculating infinite timesteps is intractable
- Typically ignore timesteps after T by multiplying future gradients by 0

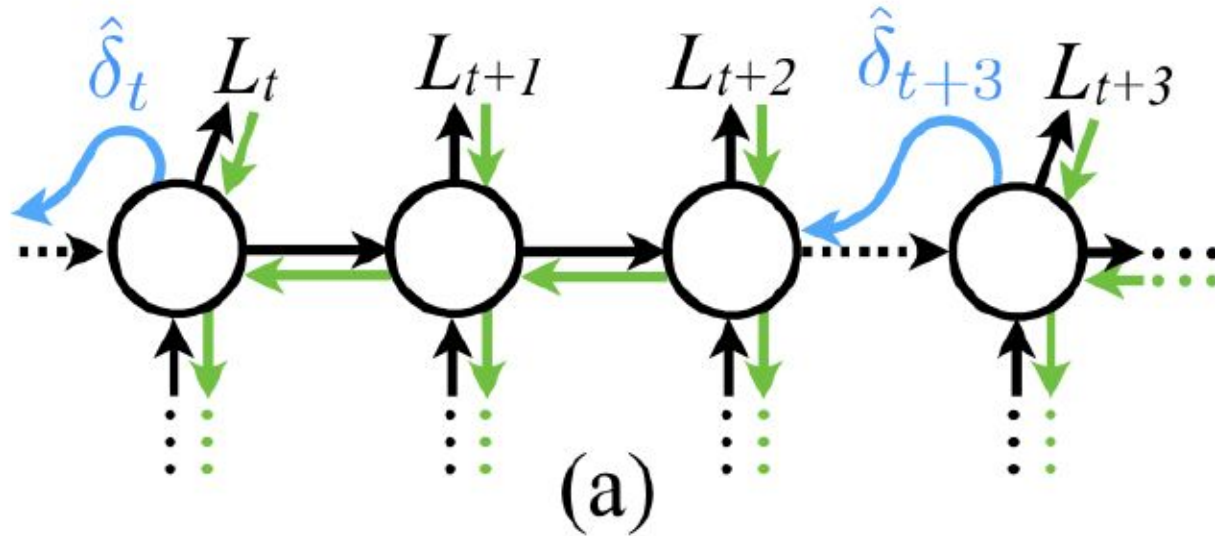Summation of true gradients over infinite timesteps after timestep T

# Breaking RNN Time Boundaries

$$\theta - \alpha \left( \sum_{\tau=t}^{t+T} \frac{\partial L_\tau}{\partial \theta} + \hat{\delta}_T \frac{\partial h_T}{\partial \theta} \right)$$
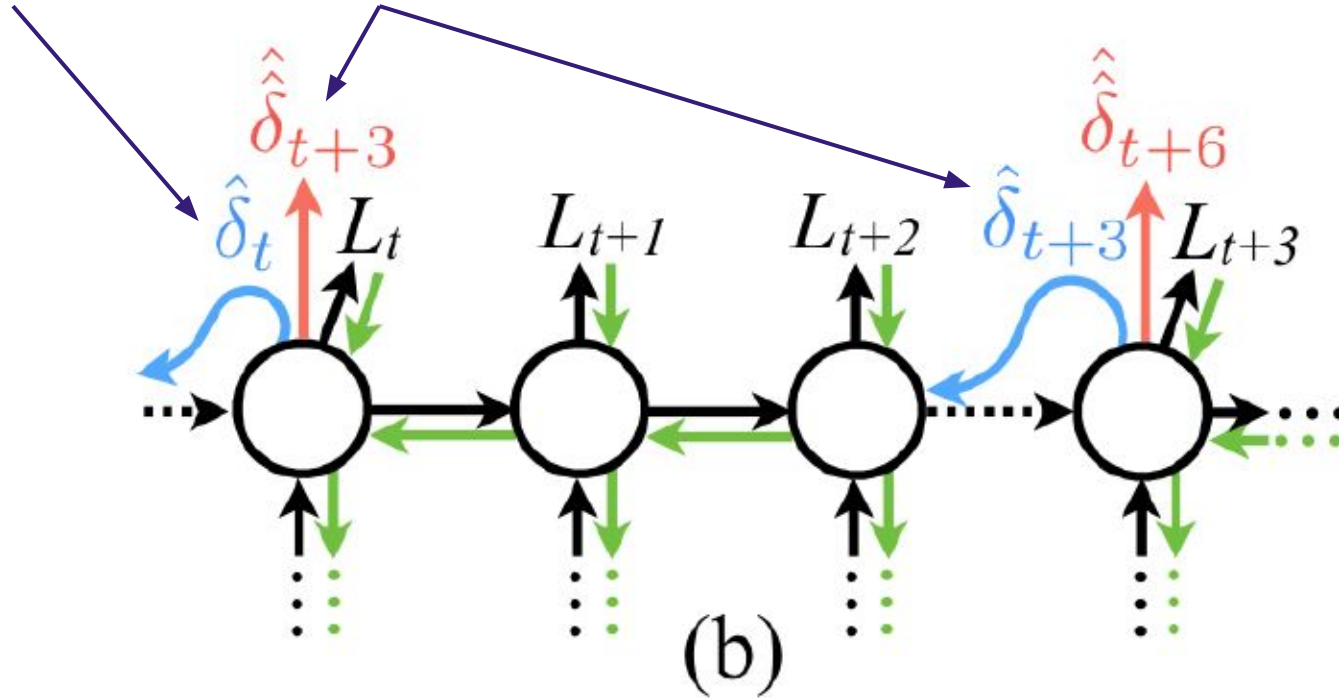
- Equivalent to unrolling for infinite timesteps with infinite subnetworks
- DNI allows RNN to asynchronously communicate with future self
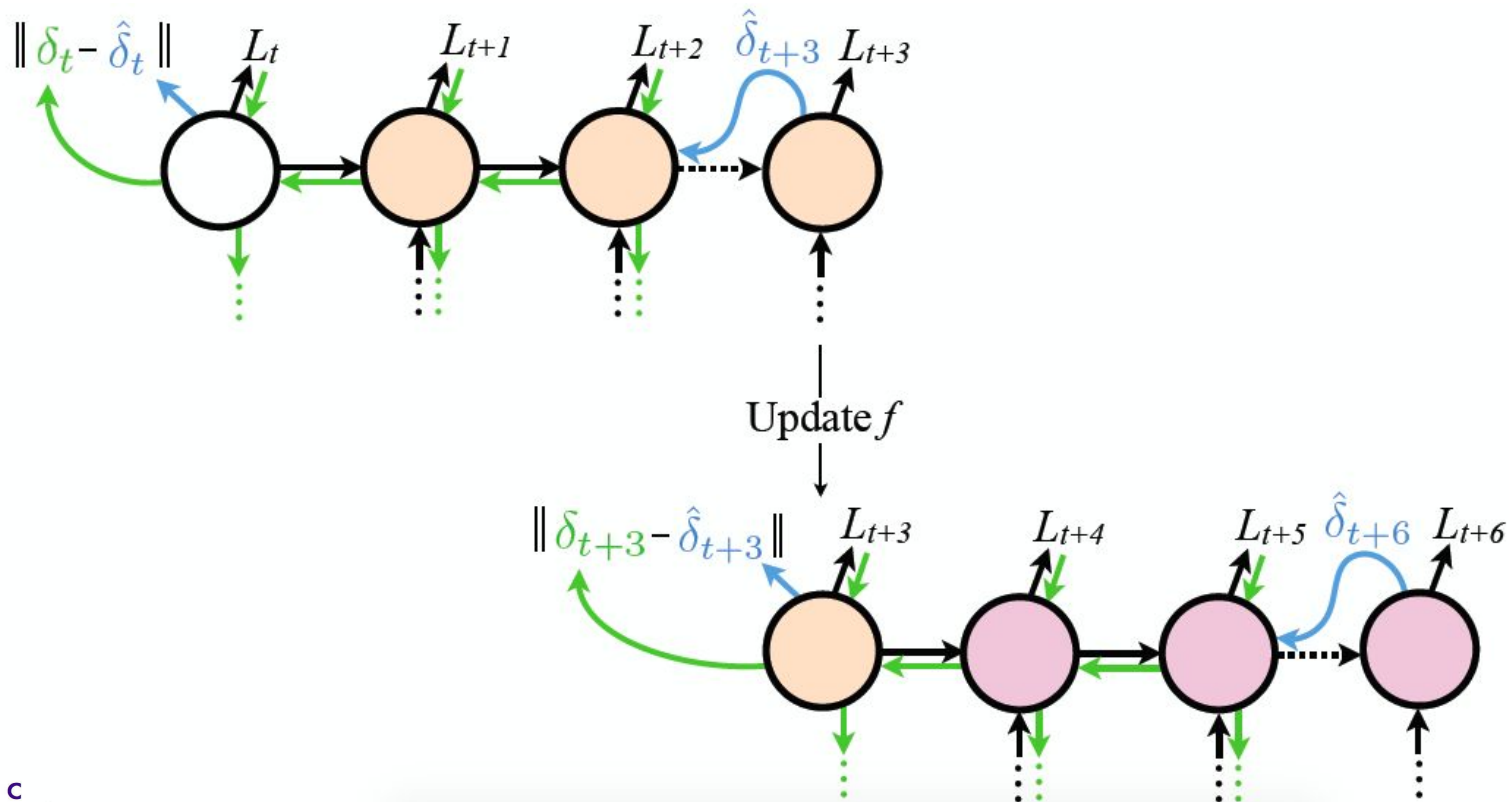
Approximate with DNI

M I C

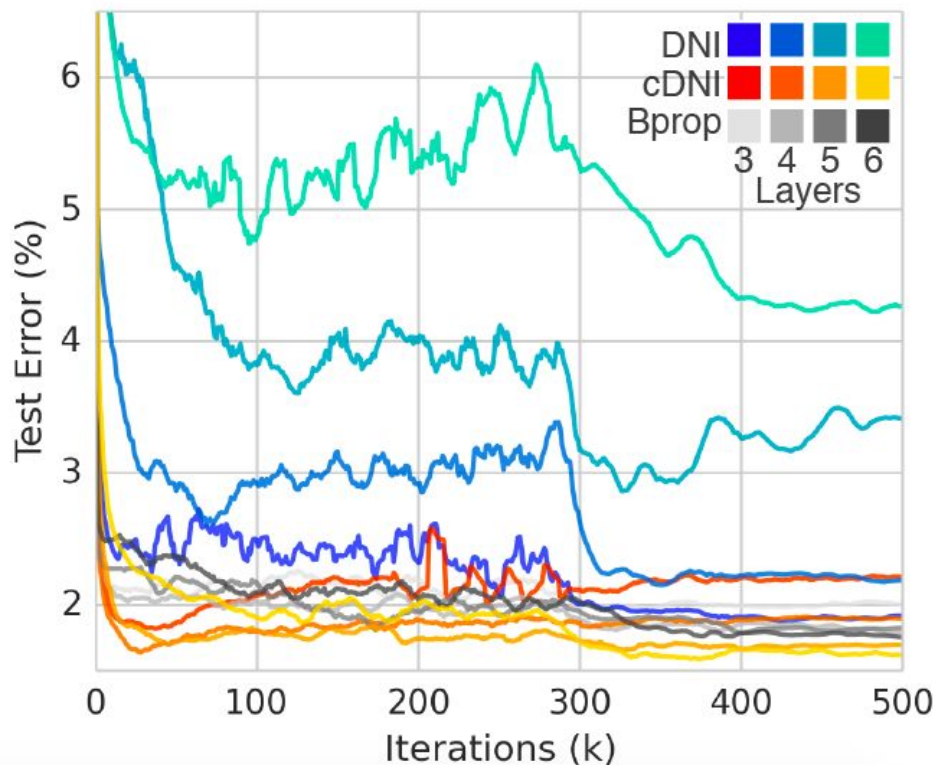15

# Breaking RNN Time Boundaries



(a)

# Short-term & Long-term Synthetic Gradients



(b)

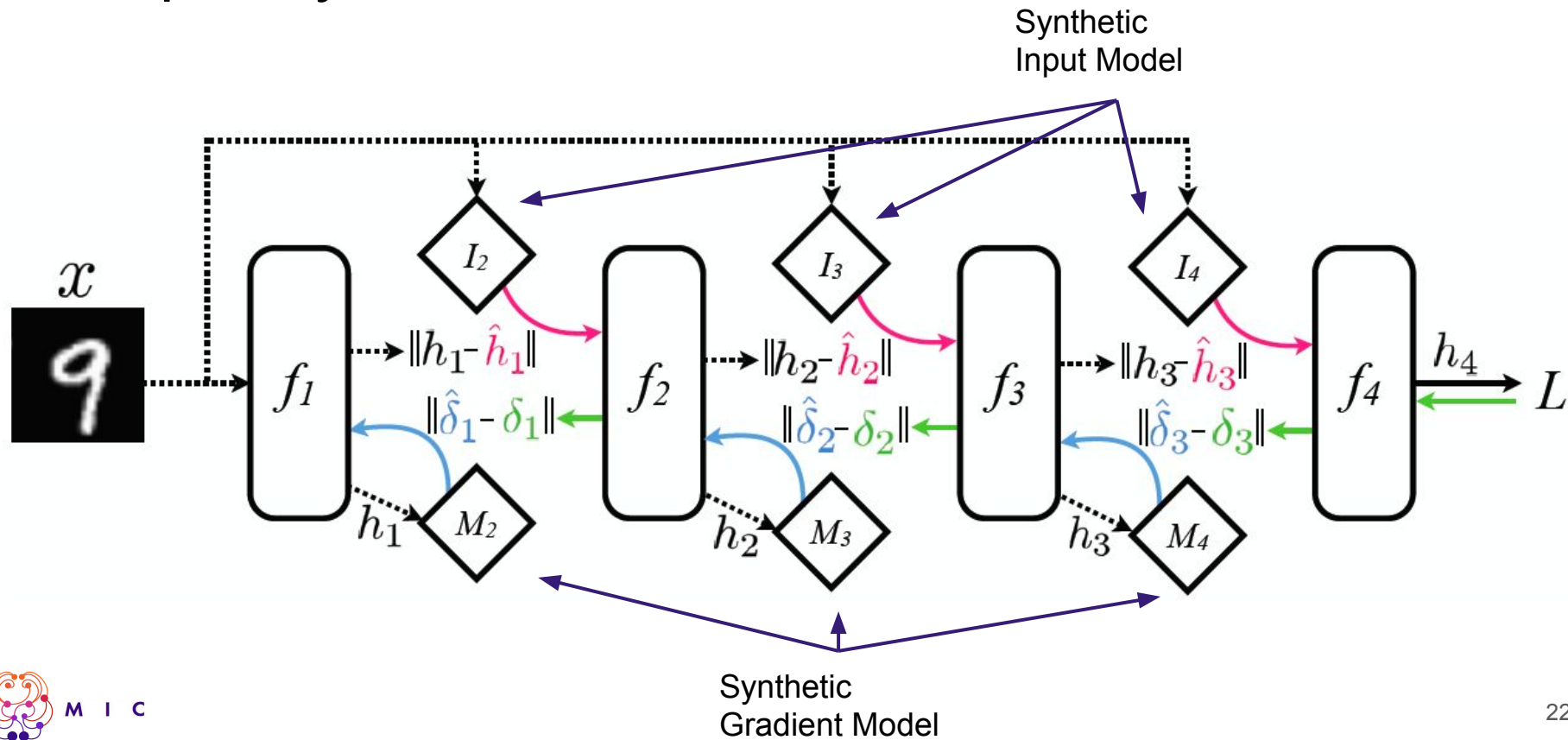# Updating Multiple DNI Recurrent Networks

# DNI between every layer in FCN

# DNI between every layer in FCN and CNN

| | Layers | MNIST (% Error) | | | | CIFAR-10 (% Error) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | No Bprop | Bprop | DNI | cDNI | No Bprop | Bprop | DNI | cDNI |
| FCN | 3 | 9.3 | 2.0 | 1.9 | 2.2 | 54.9 | 43.5 | 42.5 | 48.5 |
| | 4 | 12.6 | 1.8 | 2.2 | 1.9 | 57.2 | 43.0 | 45.0 | 45.1 |
| | 5 | 16.2 | 1.8 | 3.4 | 1.7 | 59.6 | 41.7 | 46.9 | 43.5 |
| | 6 | 21.4 | 1.8 | 4.3 | 1.6 | 61.9 | 42.0 | 49.7 | 46.8 |
| CNN | 3 | 0.9 | 0.8 | 0.9 | 1.0 | 28.7 | 17.9 | 19.5 | 19.0 |
| | 4 | 2.8 | 0.6 | 0.7 | 0.8 | 38.1 | 15.7 | 19.5 | 16.4 |

M I C

# 20% Chance Backwards Unlocking

# Completely Unlocked Network

# Forwards Unlocking



Forwards and Update Decoupled
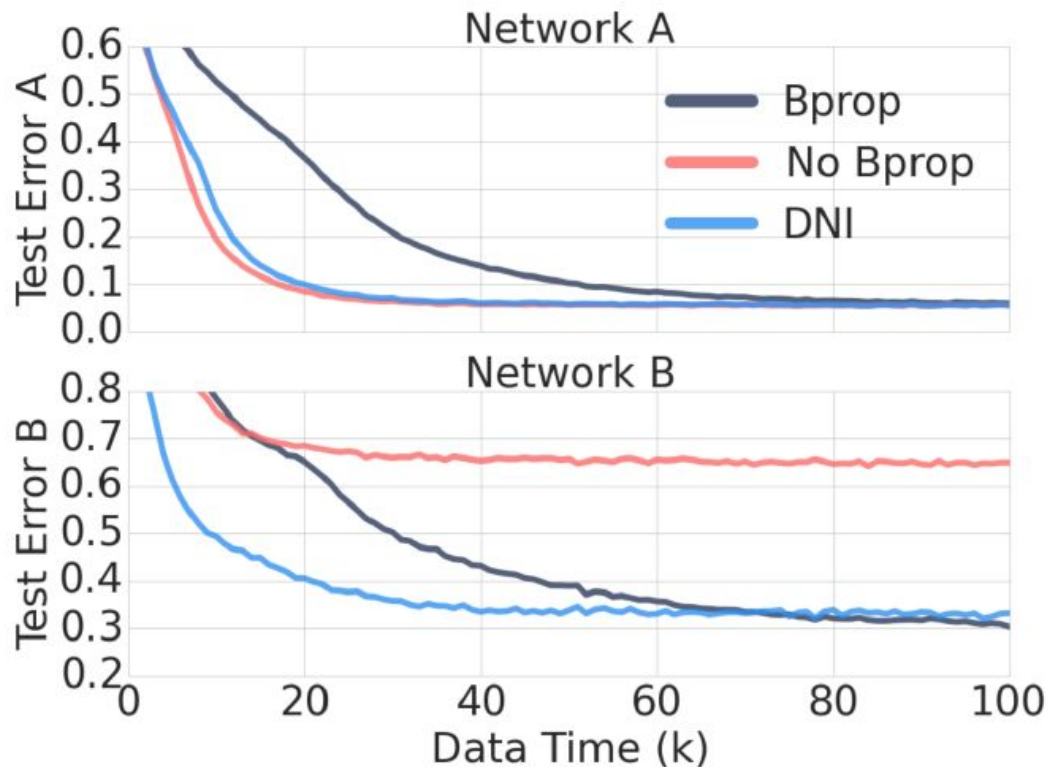
# Dynamic Computation Graphs and Ensembles

# Two RNNs communicating with DNI



- One RNN trained to count number of 3s seen

- Second RNN trained to count number odd numbers seen

# References and Further Reading

1.  Jaderberg, Max, et al. "**Decoupled neural interfaces using synthetic gradients**." arXiv preprint arXiv:1608.05343 (2016).
2.  Czarnecki, Wojciech Marian, et al. "**Understanding Synthetic Gradients and Decoupled Neural Interfaces**." arXiv preprint arXiv:1703.00522 (2017).
3.  Miyato, Takeru, et al. "**Synthetic Gradient Methods with Virtual Forward-Backward Networks**." (2017).
4.  Czarnecki, Wojciech Marian, et al. "**Sobolev Training for Neural Networks**." arXiv preprint arXiv:1706.04859 (2017).