

# Deep Equilibrium Models



authors:

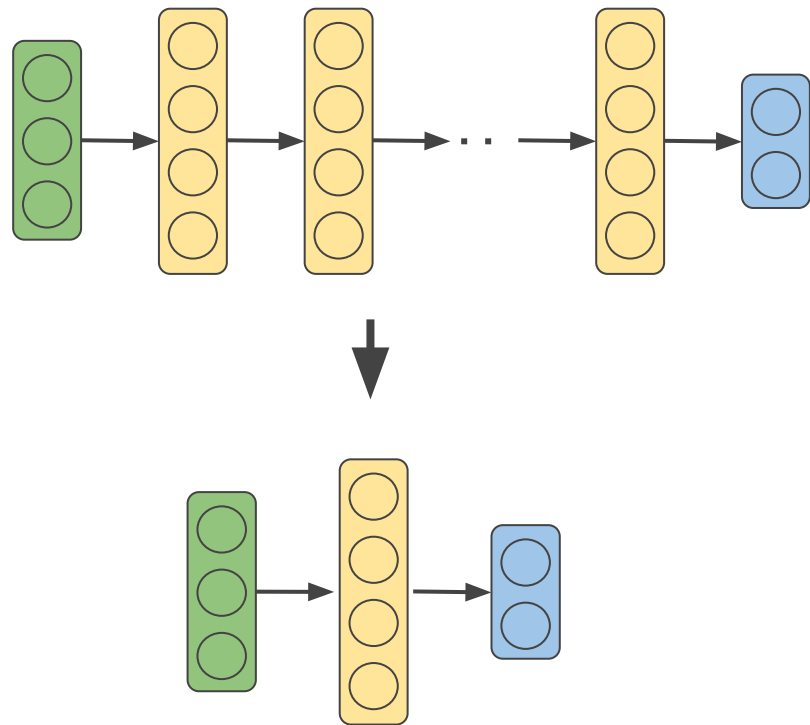
Shaojie Bai, J. Zico Kolter, Vladlen Koltun

presented by:

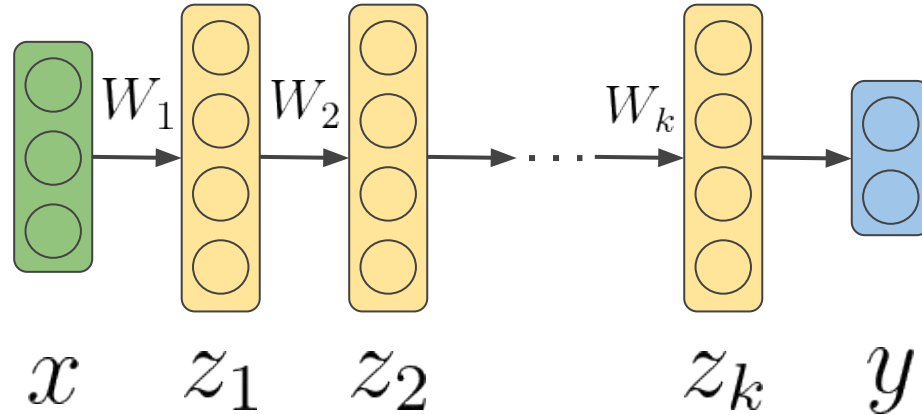
Matthias Otth

# Overview

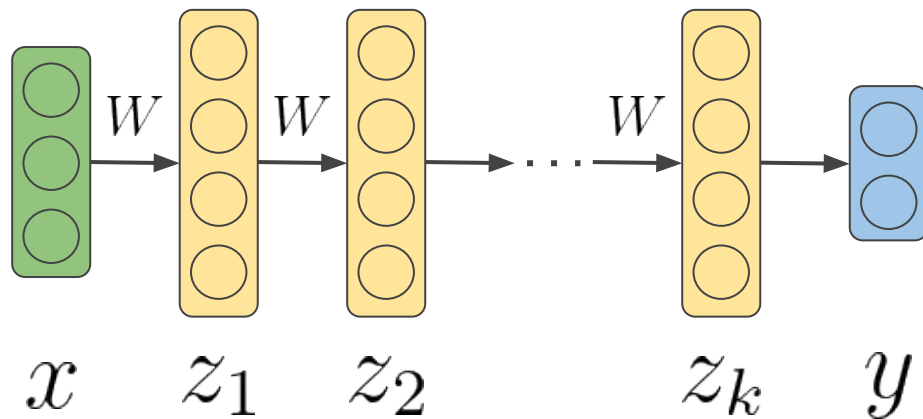
- Reinterpret deep NN
- Same performance with less memory consumption



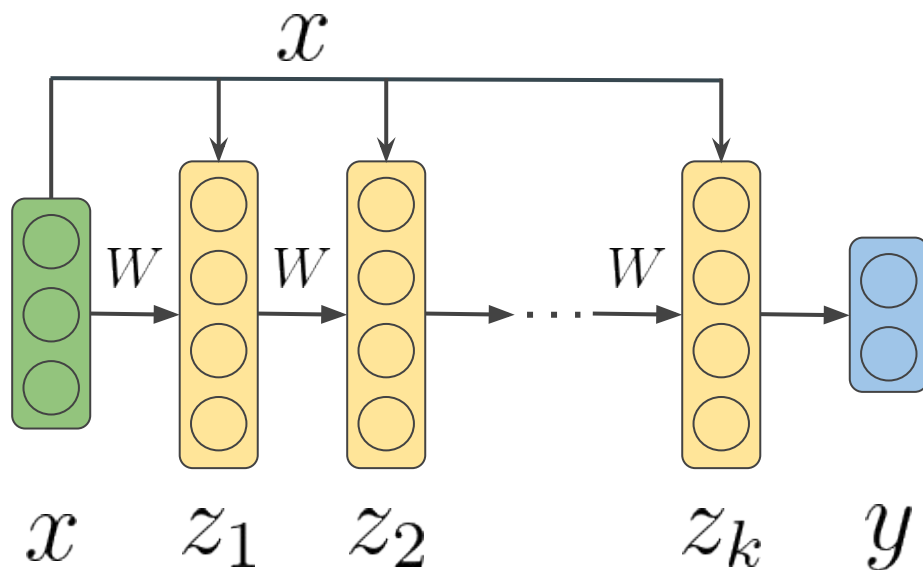
# Classical deep feedforward NN



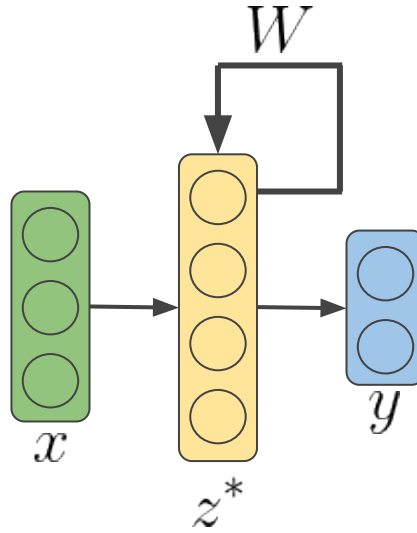
# Weight-tied Network



# Weight-tied, input-injected Network



# Infinite depth network



# Equilibrium formulation

- Almost any non-linear function:

$$z^{[i+1]} = f_{\theta}(z^{[i]}; x)$$

# Equilibrium formulation

- Almost any non-linear function:

$$z^{[i+1]} = f_{\theta}(z^{[i]}; x)$$

- Assume equilibrium point exists:

$$\lim_{i \rightarrow \infty} z^{[i]} = z^* = f(z^*; x)$$



# Equilibrium formulation

- Almost any non-linear function:

$$z^{[i+1]} = f_{\theta}(z^{[i]}; x)$$

- Assume equilibrium point exists:

$$\lim_{i \rightarrow \infty} z^{[i]} = z^* = f(z^*; x)$$

- Reformulate as root finding problem:

$$g_{\theta}(z^*; x) = f_{\theta}(z^*; x) - z^* = 0$$

# Equilibrium formulation

- Almost any non-linear function:

$$z^{[i+1]} = f_{\theta}(z^{[i]}; x)$$

- Assume equilibrium point exists:

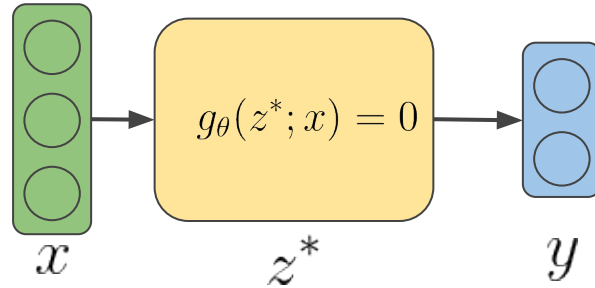
$$\lim_{i \rightarrow \infty} z^{[i]} = z^* = f(z^*; x)$$

- Reformulate as root finding problem:

$$g_{\theta}(z^*; x) = f_{\theta}(z^*; x) - z^* = 0$$

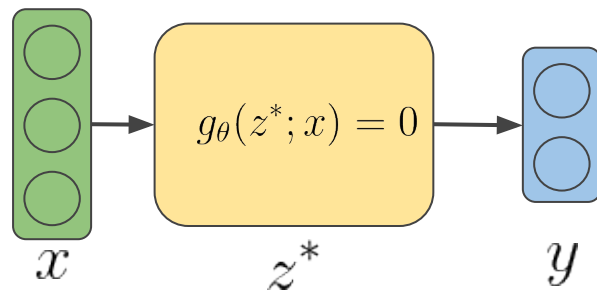
- Gives rise to a single (implicit) layer model

# Implicit layer formulation



# DEQ

- Equivalent to infinite-depth network!
- Different interpretation of deep networks
- We can backpropagate through equilibrium point:  $O(1)$  memory

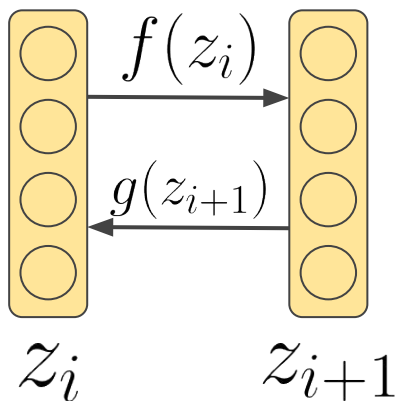


# Previous Work

# Previous work: Implicit Layers

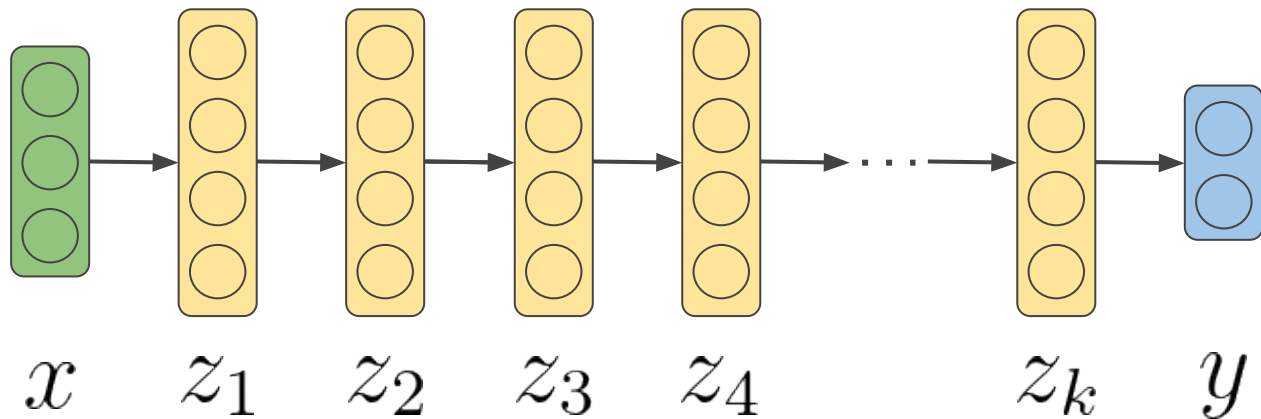
- Applied to small scales
- Very specific models and tasks

# Previous work: Reversible Networks



- $O(1)$  memory consumption
- Strong restriction in model architecture

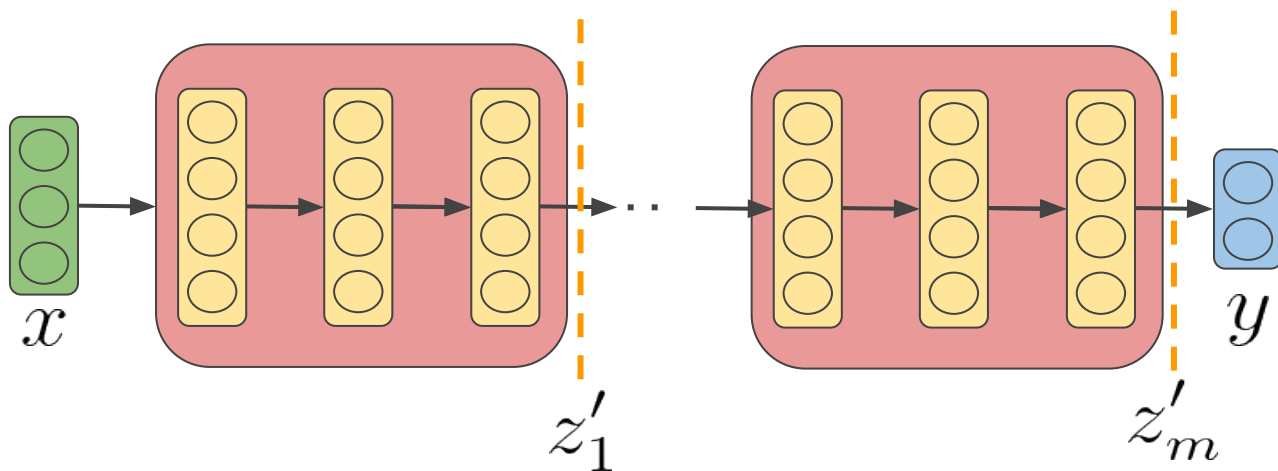
# Previous work: Gradient Checkpointing





# Previous work: Gradient Checkpointing

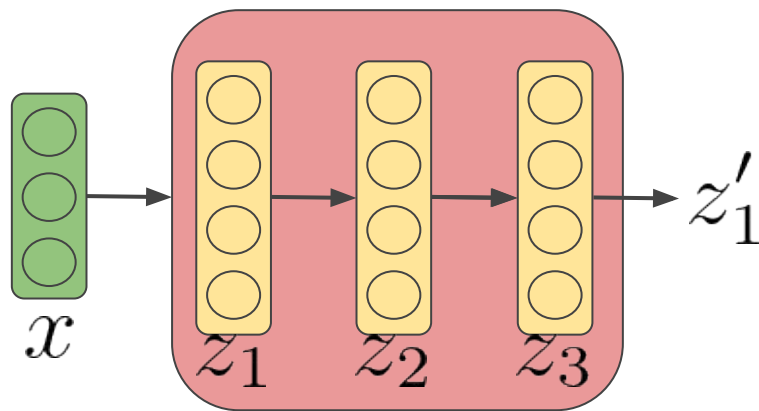
Step 1: High-level backpropagation



→ Can calculate  $\frac{\partial \mathcal{L}(y, y')}{\partial z'_i}$  in  $O(m)$  memory

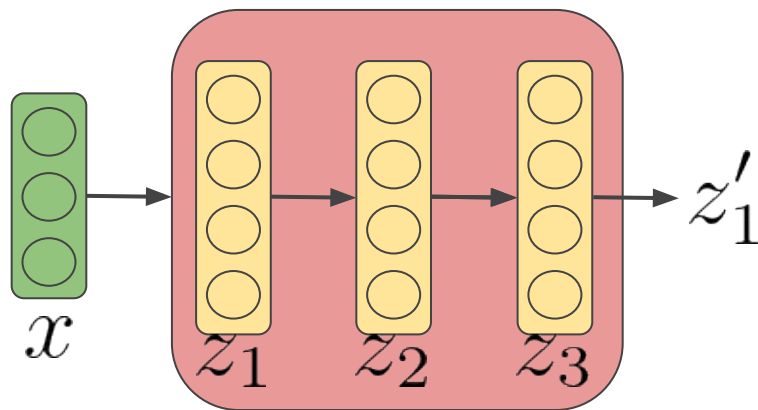
# Previous work: Gradient Checkpointing

Step 2: Low-level backpropagation



# Previous work: Gradient Checkpointing

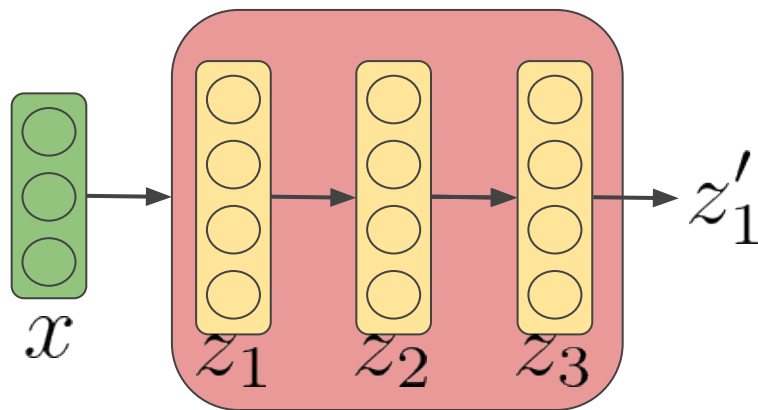
Step 2: Low-level backpropagation



$$\frac{\partial \mathcal{L}(y, y')}{\partial z_i} = \overbrace{\frac{\partial \mathcal{L}(y, y')}{\partial z'_1}}^{\checkmark} \frac{\partial z'_1}{\partial z_i}$$

# Previous work: Gradient Checkpointing

Step 2: Low-level backpropagation



$$\frac{\partial \mathcal{L}(y, y')}{\partial z_i} = \frac{\partial \mathcal{L}(y, y')}{\partial z'_1} \frac{\partial z'_1}{\partial z_i}$$

→ Can calculate  $\frac{\partial \mathcal{L}(y, y')}{\partial z_i}$  in  $O(S)$  memory

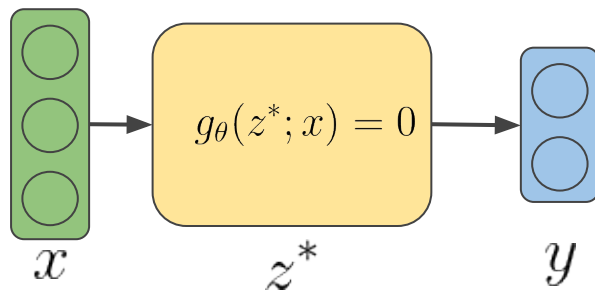
# Previous work: Gradient Checkpointing

## Summary:

- Cost:  $O(S + m)$
- $L = S \times m$
- Can achieve  $O(\sqrt{L})$  memory usage for 2x training time
- Can theoretically achieve  $O(\log L)$  memory usage, if applied recursively

# DEQ

- Equivalent to infinite-depth network!
- Different interpretation of deep networks
- We can backpropagate through equilibrium point:  $O(1)$  memory



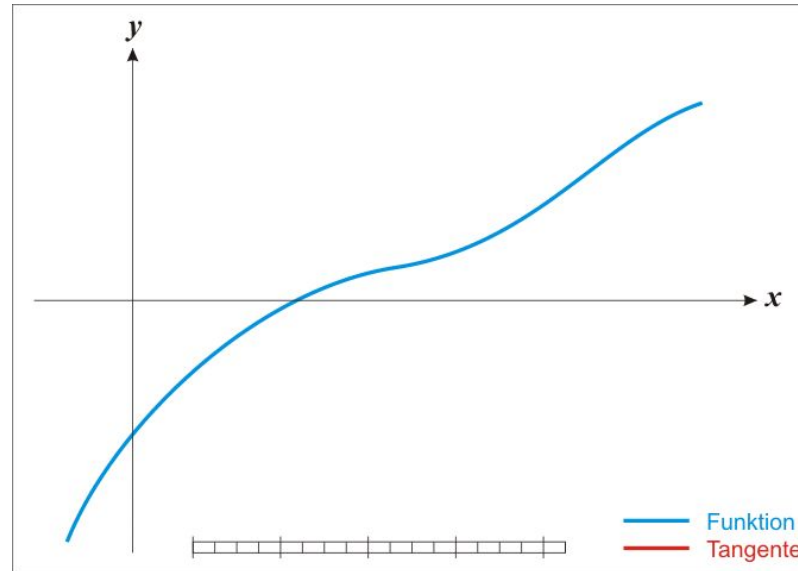
# Forward pass

Find fixpoint:  $g_{\theta}(z^*; x) = f_{\theta}(z^*; x) - z^* = 0$

# Forward pass

For example with Newton's method:

$$z^{[i+1]} = z^{[i]} - \alpha \left( J_{g_\theta}^{-1} \Big|_{z^{[i]}} \right) g_\theta(z^{[i]}, x)$$





# Forward pass

Find fixpoint:  $g_\theta(z^*; x) = f_\theta(z^*; x) - z^* = 0$

For example with Newton's method:

$$z^{[i+1]} = z^{[i]} - \alpha(J_{g_\theta}^{-1} |_{z^{[i]}})g_\theta(z^{[i]}, x)$$

Can use any black-box root-finding algorithm  $z^* = \text{RootFind}(g_\theta; x)$

# Backward pass: 1st Approach

## Procedure:

1. Fix a RootFind algorithm (e.g. Newton's method)
2. Unroll the Newton iterations
3. Do backpropagation through all iterations

# Backward pass: 1st Approach

## Procedure:

1. Fix a RootFind algorithm (e.g. Newton's method)
2. Unroll the Newton iterations
3. Do backpropagation through all iterations

## Problems:

- Need knowledge of RootFind algorithm (Not a blackbox)
- Need to store intermediate results (Not  $O(1)$ )

# Backward pass: 2nd Approach

Procedure:

- Find root:  $z^* = \text{RootFind}(g_\theta; x)$
- Calculate loss:  $\mathcal{L}(z^*, y)$
- Theorem 1: 
$$\frac{\partial \mathcal{L}}{\partial \theta} = -\frac{\partial \mathcal{L}}{\partial z^*} (J_{g_\theta}^{-1} |_{z^*}) \frac{\partial f_\theta(z^*; x)}{\partial \theta}$$

# Backward pass: 2nd Approach

## Procedure:

- Find root:  $z^* = \text{RootFind}(g_\theta; x)$
- Calculate loss:  $\mathcal{L}(z^*, y)$
- Theorem 1:  $\frac{\partial \mathcal{L}}{\partial \theta} = -\frac{\partial \mathcal{L}}{\partial z^*} (J_{g_\theta}^{-1} |_{z^*}) \frac{\partial f_\theta(z^*; x)}{\partial \theta}$

## Advantages

- Independent of RootFind!
- Single step to backpropagate through ‘infinite depth’ network.

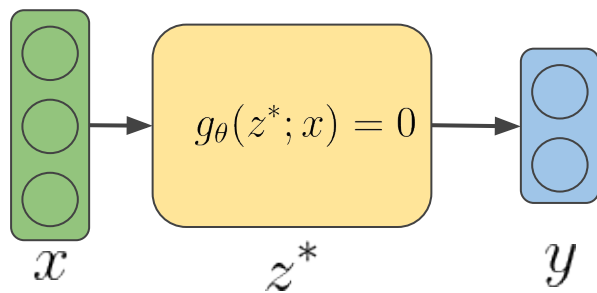
# Broyden's Method

**Problem:** Calculating Jacobian Inverse is expensive

**Solution:** Use quasi-Newton methods.

$$J_{g\theta}^{-1} \Big|_{\mathbf{z}_{1:T}^{[i+1]}} \approx B_{g\theta}^{[i+1]} = B_{g\theta}^{[i]} + \frac{\Delta \mathbf{z}^{[i+1]} - B_{g\theta}^{[i]} \Delta g_{\theta}^{[i+1]}}{\Delta \mathbf{z}^{[i+1] \top} B_{g\theta}^{[i]} \Delta g_{\theta}^{[i+1]}} \Delta \mathbf{z}^{[i+1] \top} B_{g\theta}^{[i]}$$

# DEQ



# Guarantees:

## Memory consumption independent of depth

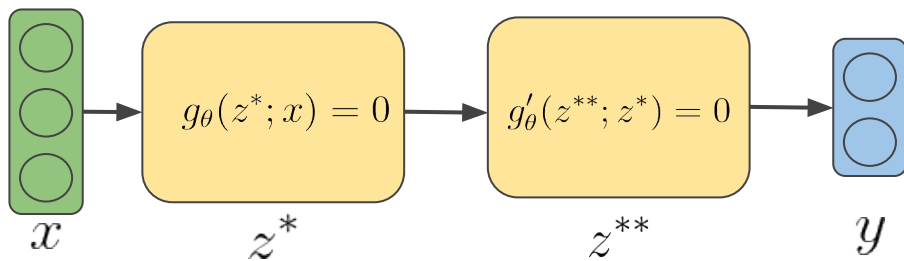
- $O(1)$  memory consumption for backpropagation



# Guarantees:

## Sufficiency of a Single DEQ “Layer”

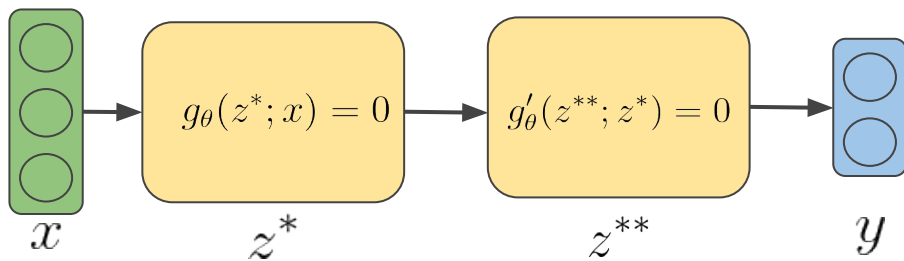
**Idea:** Stack multiple DEQs together, to get more representational power.



# Guarantees:

## Sufficiency of a Single DEQ “Layer”

**Idea:** Stack multiple DEQs together, to get more representational power.



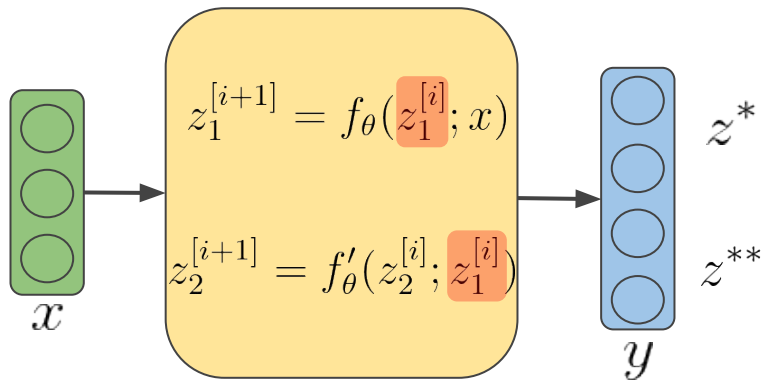
**Theorem:** A single DEQ “layer” is enough.

# Guarantees:

## Sufficiency of a Single DEQ “Layer”

Proof sketch:

- Stack the two layers
- Use output of first layer as input to second layer

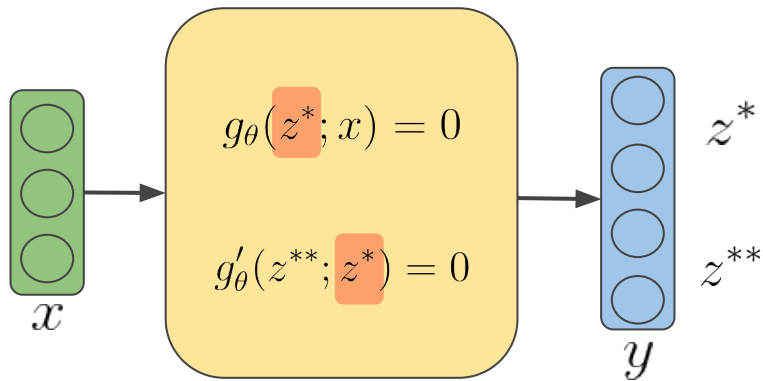


# Guarantees:

## Sufficiency of a Single DEQ “Layer”

Proof sketch:

- Stack the two layers
- Use output of first layer as input to second layer



# Guarantees:

## Universality of Weight-tied, Input-injected Networks

**Theorem:** Any traditional  $L$ -layer deep network can be represented by an  $L$ -layer deep weight tied, input-injected network with linear increase in width.

# Guarantees:

## Universality of Weight-tied, Input-injected Networks

**Theorem:** Any traditional  $L$ -layer deep network can be represented by an  $L$ -layer deep weight tied, input-injected network with linear increase in width.

### Setting:

We have:  $\mathbf{z}^{[i+1]} = \sigma^{[i]}(W^{[i]}\mathbf{z}^{[i]} + \mathbf{b}^{[i]})$ ,  $i = 0, \dots, L - 1$ ,  $\mathbf{z}^{[0]} = \mathbf{x}$

# Guarantees:

## Universality of Weight-tied, Input-injected Networks

**Theorem:** Any traditional  $L$ -layer deep network can be represented by an  $L$ -layer deep weight tied, input-injected network with linear increase in width.

### Setting:

We have:  $\mathbf{z}^{[i+1]} = \sigma^{[i]}(W^{[i]}\mathbf{z}^{[i]} + \mathbf{b}^{[i]})$ ,  $i = 0, \dots, L - 1$ ,  $\mathbf{z}^{[0]} = \mathbf{x}$

We want:  $\tilde{\mathbf{z}}^{[i+1]} = \sigma(W_z \tilde{\mathbf{z}}^{[i]} + W_x \mathbf{x} + \tilde{\mathbf{b}})$ ,  $i = 0, \dots, L - 1$ .

# Guarantees:

## Universality of Weight-tied, Input-injected Networks

**Proof:**

We have:  $\mathbf{z}^{[i+1]} = \sigma^{[i]}(W^{[i]}\mathbf{z}^{[i]} + \mathbf{b}^{[i]})$ ,  $i = 0, \dots, L-1$ ,  $\mathbf{z}^{[0]} = \mathbf{x}$

We want:  $\tilde{\mathbf{z}}^{[i+1]} = \sigma(W_z \tilde{\mathbf{z}}^{[i]} + W_x \mathbf{x} + \tilde{\mathbf{b}})$ ,  $i = 0, \dots, L-1$ .

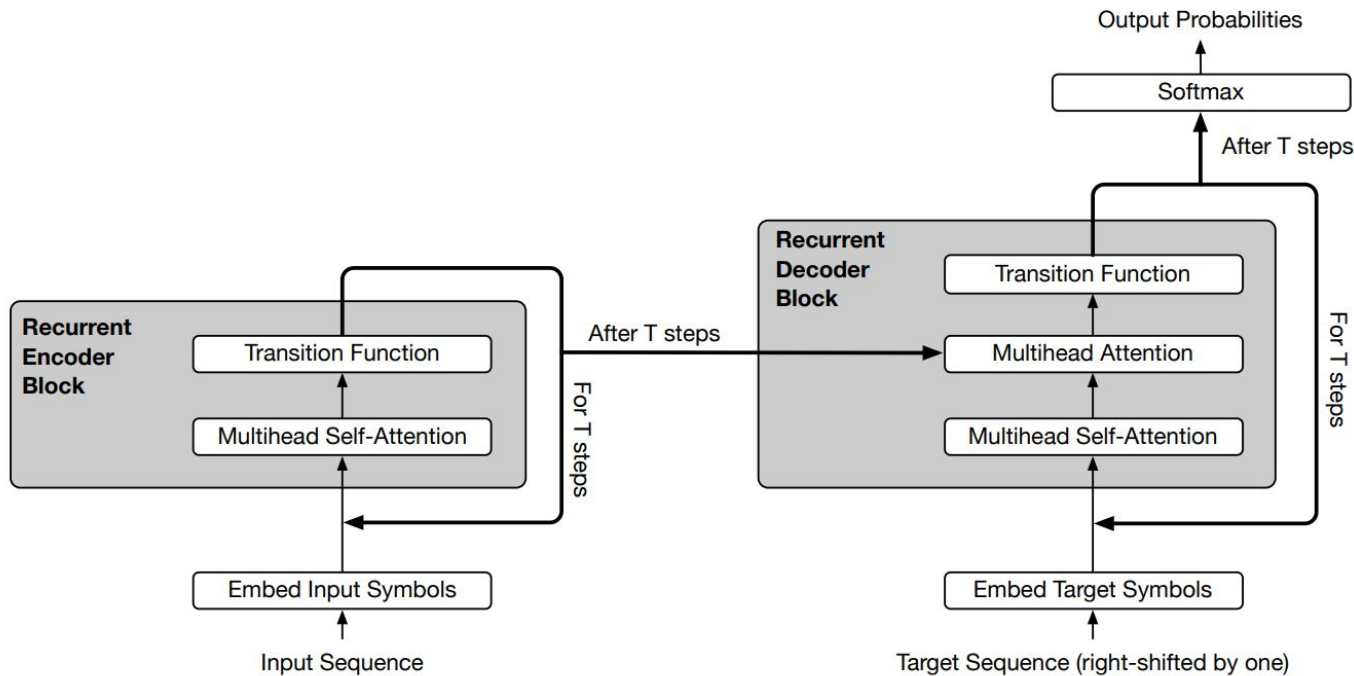
$$W_z = \begin{bmatrix} 0 & 0 & \dots & 0 & 0 \\ W^{[1]} & 0 & \dots & 0 & 0 \\ 0 & W^{[2]} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & W^{[L-1]} & 0 \end{bmatrix}, W_x = \begin{bmatrix} W^{[0]} \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \tilde{\mathbf{b}} = \begin{bmatrix} \mathbf{b}^{[0]} \\ \mathbf{b}^{[1]} \\ \vdots \\ \mathbf{b}^{[L-1]} \end{bmatrix}, \sigma = \begin{bmatrix} \sigma^{[0]} \\ \sigma^{[1]} \\ \vdots \\ \sigma^{[L-1]} \end{bmatrix}$$

**This is not done in practice!**

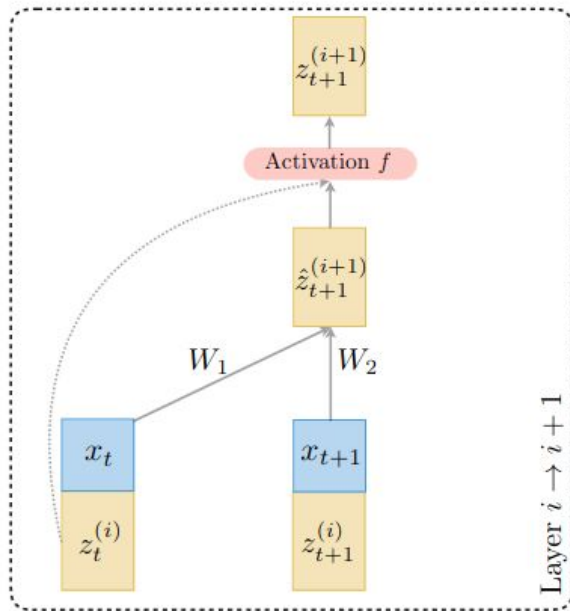


# Evaluation

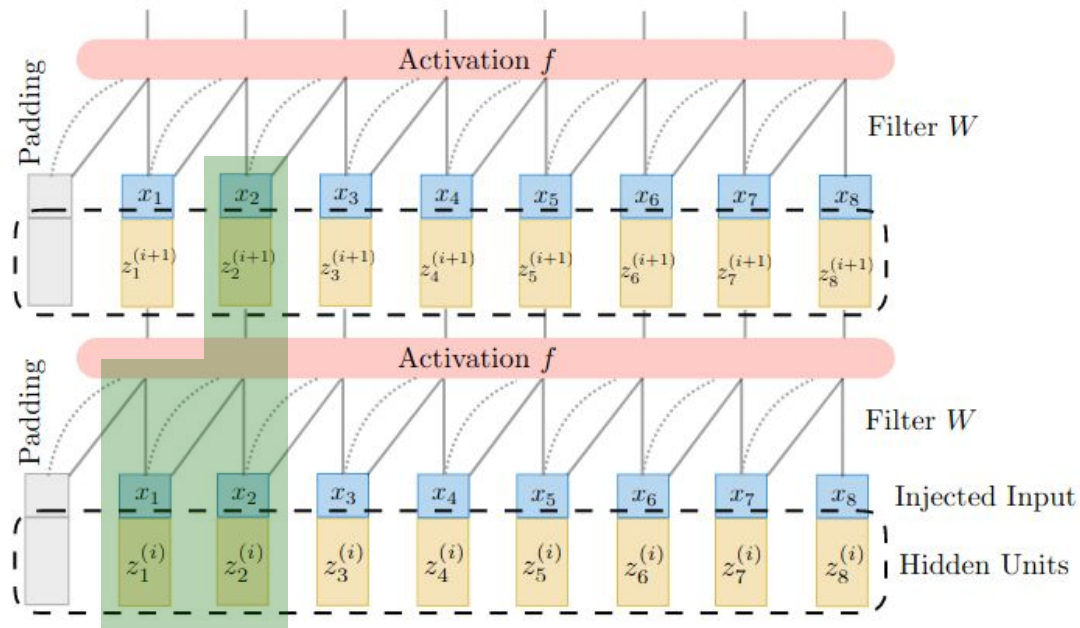
# Universal Transformer



# TrellisNet: Atomic Level



# TrellisNet



# Results: Penn Treebank

Word-level Language Modeling w/ Penn Treebank (PTB)				
Model	# Params	Non-embedding model size	Test perplexity	Memory <sup>†</sup>
Variational LSTM [22]	66M	-	73.4	-
NAS Cell [55]	54M	-	62.4	-
NAS (w/ black-box hyperparameter tuner) [32]	24M	20M	59.7	-
AWD-LSTM [34]	24M	20M	58.8	-
DARTS architecture search (second order) [29]	23M	20M	<b>55.7</b>	-
60-layer TrellisNet (w/ auxiliary loss, w/o MoS) [8]	24M	20M	57.0	8.5GB
<b>DEQ-TrellisNet (ours)</b>	24M	20M	57.1	<b>1.2GB</b>

# Results: WikiText-103

Word-level Language Modeling w/ WikiText-103 (WT103)				
Model	# Params	Non-Embedding Model Size	Test perplexity	Memory <sup>†</sup>
Generic TCN [7]	150M	34M	45.2	-
Gated Linear ConvNet [17]	230M	-	37.2	-
AWD-QRNN [33]	159M	51M	33.0	7.1GB
Relational Memory Core [40]	195M	60M	31.6	-
Transformer-XL (X-large, adaptive embed., on TPU) [16]	257M	224M	<b>18.7</b>	12.0GB
70-layer TrellisNet (+ auxiliary loss, etc.) [8]	180M	45M	29.2	24.7GB
70-layer TrellisNet with <i>gradient checkpointing</i>	180M	45M	29.2	5.2GB
<b>DEQ-TrellisNet (ours)</b>	180M	45M	<b>29.0</b>	<b>3.3GB</b>
Transformer-XL (medium, 16 layers)	165M	44M	24.3	8.5GB
<b>DEQ-Transformer (medium, ours).</b>	172M	43M	24.2	<b>2.7GB</b>
Transformer-XL (medium, 18 layers, adaptive embed.)	110M	72M	23.6	9.0GB
<b>DEQ-Transformer (medium, adaptive embed., ours)</b>	110M	70M	<b>23.2</b>	3.7GB
Transformer-XL (small, 4 layers)	139M	4.9M	35.8	4.8GB
Transformer-XL (small, weight-tied 16 layers)	138M	4.5M	34.9	6.8GB
<b>DEQ-Transformer (small, ours).</b>	138M	4.5M	<b>32.4</b>	<b>1.1GB</b>

# Results: WikiText-103

Word-level Language Modeling w/ WikiText-103 (WT103)				
Model	# Params	Non-Embedding Model Size	Test perplexity	Memory <sup>†</sup>
Generic TCN [7]	150M	34M	45.2	-
Gated Linear ConvNet [17]	230M	-	37.2	-
AWD-QRNN [33]	159M	51M	33.0	7.1GB
Relational Memory Core [40]	195M	60M	31.6	-
Transformer-XL (X-large, adaptive embed., on TPU) [16]	257M	224M	<b>18.7</b>	12.0GB
70-layer TrellisNet (+ auxiliary loss, etc.) [8]	180M	45M	29.2	24.7GB
70-layer TrellisNet with <i>gradient checkpointing</i>	180M	45M	29.2	5.2GB
<b>DEQ-TrellisNet (ours)</b>	180M	45M	<b>29.0</b>	<b>3.3GB</b>
Transformer-XL (medium, 16 layers)	165M	44M	24.3	8.5GB
<b>DEQ-Transformer (medium, ours).</b>	172M	43M	24.2	<b>2.7GB</b>
Transformer-XL (medium, 18 layers, adaptive embed.)	110M	72M	23.6	9.0GB
<b>DEQ-Transformer (medium, adaptive embed., ours)</b>	110M	70M	<b>23.2</b>	3.7GB
Transformer-XL (small, 4 layers)	139M	4.9M	35.8	4.8GB
Transformer-XL (small, weight-tied 16 layers)	138M	4.5M	34.9	6.8GB
<b>DEQ-Transformer (small, ours).</b>	138M	4.5M	<b>32.4</b>	<b>1.1GB</b>

# Results: WikiText-103

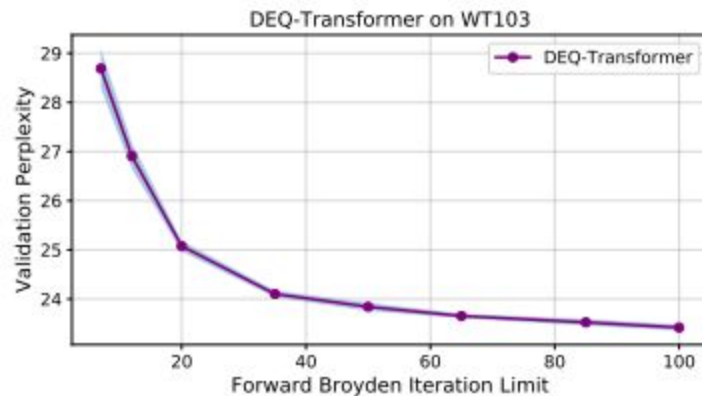
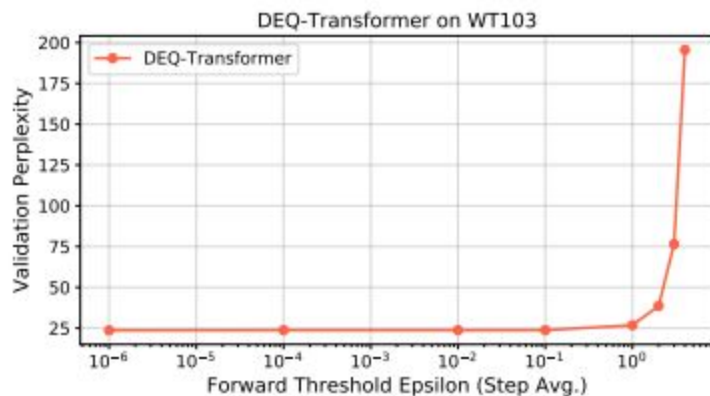
Word-level Language Modeling w/ WikiText-103 (WT103)				
Model	# Params	Non-Embedding Model Size	Test perplexity	Memory <sup>†</sup>
Generic TCN [7]	150M	34M	45.2	-
Gated Linear ConvNet [17]	230M	-	37.2	-
AWD-QRNN [33]	159M	51M	33.0	7.1GB
Relational Memory Core [40]	195M	60M	31.6	-
Transformer-XL (X-large, adaptive embed., on TPU) [16]	257M	224M	<b>18.7</b>	12.0GB
70-layer TrellisNet (+ auxiliary loss, etc.) [8]	180M	45M	29.2	24.7GB
70-layer TrellisNet with <i>gradient checkpointing</i>	180M	45M	29.2	5.2GB
<b>DEQ-TrellisNet (ours)</b>	180M	45M	<b>29.0</b>	<b>3.3GB</b>
Transformer-XL (medium, 16 layers)	165M	44M	24.3	8.5GB
<b>DEQ-Transformer (medium, ours).</b>	172M	43M	24.2	<b>2.7GB</b>
Transformer-XL (medium, 18 layers, adaptive embed.)	110M	72M	23.6	9.0GB
<b>DEQ-Transformer (medium, adaptive embed., ours)</b>	110M	70M	<b>23.2</b>	3.7GB
Transformer-XL (small, 4 layers)	139M	4.9M	35.8	4.8GB
Transformer-XL (small, weight-tied 16 layers)	138M	4.5M	34.9	6.8GB
<b>DEQ-Transformer (small, ours).</b>	138M	4.5M	<b>32.4</b>	<b>1.1GB</b>



# Results: WikiText-103

Word-level Language Modeling w/ WikiText-103 (WT103)				
Model	# Params	Non-Embedding Model Size	Test perplexity	Memory <sup>†</sup>
Generic TCN [7]	150M	34M	45.2	-
Gated Linear ConvNet [17]	230M	-	37.2	-
AWD-QRNN [33]	159M	51M	33.0	7.1GB
Relational Memory Core [40]	195M	60M	31.6	-
<b>Transformer-XL (X-large, adaptive embed., on TPU) [16]</b>	257M	224M	<b>18.7</b>	<b>12.0GB</b>
70-layer TrellisNet (+ auxiliary loss, etc.) [8]	180M	45M	29.2	24.7GB
70-layer TrellisNet with <i>gradient checkpointing</i>	180M	45M	29.2	5.2GB
<b>DEQ-TrellisNet (ours)</b>	180M	45M	<b>29.0</b>	<b>3.3GB</b>
Transformer-XL (medium, 16 layers)	165M	44M	24.3	8.5GB
<b>DEQ-Transformer (medium, ours).</b>	172M	43M	24.2	<b>2.7GB</b>
Transformer-XL (medium, 18 layers, adaptive embed.)	110M	72M	23.6	9.0GB
<b>DEQ-Transformer (medium, adaptive embed., ours)</b>	110M	70M	<b>23.2</b>	3.7GB
Transformer-XL (small, 4 layers)	139M	4.9M	35.8	4.8GB
Transformer-XL (small, weight-tied 16 layers)	138M	4.5M	34.9	6.8GB
<b>DEQ-Transformer (small, ours).</b>	138M	4.5M	<b>32.4</b>	<b>1.1GB</b>

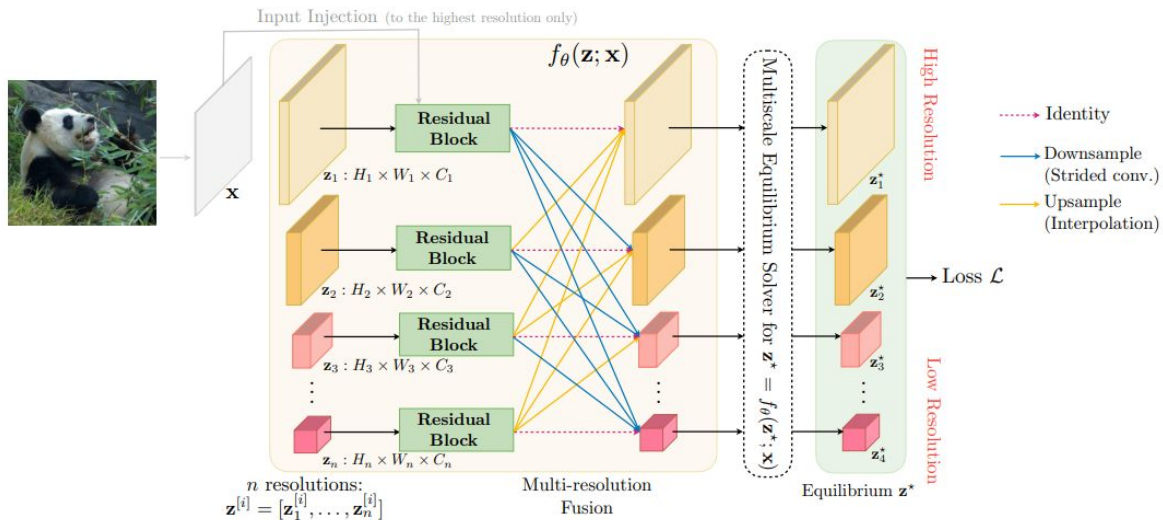
# Results: Broyden's Method



# Results: Runtime

DEQ / 18-layer Transformer		DEQ / 70-layer TrellisNet	
Training	Inference	Training	Inference
2.82×	1.76×	2.40×	1.64×

# DEQs Today



- Close to state of the art
- Very versatile (segmentation and classification)

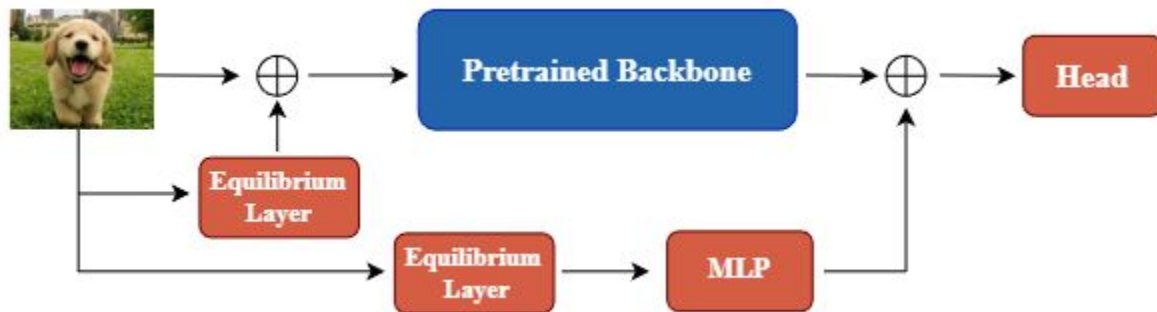
# DEQs Today

## LION: Implicit Vision Prompt Tuning

Haixin Wang<sup>1</sup> Jianlong Chang<sup>2</sup> Xiao Luo<sup>1</sup> Jinan Sun<sup>1</sup> Zhouchen Lin<sup>1</sup> Qi Tian<sup>2\*</sup>

<sup>1</sup>Peking University, Beijing, China    <sup>2</sup>Huawei Cloud & AI, Beijing, China

wang.hx@stu.pku.edu.cn, {xiaoluo, sjn, zlin}@pku.edu.cn, {jianlong.chang, tian.qi1}@huawei.com



# Conclusion

- Constant memory consumption
- New perspective on deep feed-forward NNs
- Slower to train
- Convergence to fix-point not guaranteed
- Theoretically equivalent to general network with linear width increase
- Every layer must have the same structure
- More restrictive than gradient checkpointing

# Sources

[1]: Shaojie Bai, J. Zico Kolter, Vladlen Koltun. Deep Equilibrium Models

[2]: [http://implicit-layers-tutorial.org/deep\\_equilibrium\\_models/](http://implicit-layers-tutorial.org/deep_equilibrium_models/)

[3]: Aidan N Gomez, Mengye Ren, Raquel Urtasun, and Roger B Grosse. The reversible residual network: Backpropagation without storing activations.

[4]: Matthew MacKay, Paul Vicol, Jimmy Ba, and Roger B. Grosse. Reversible recurrent neural networks.

[5]: Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost.

[6]: [https://en.wikipedia.org/wiki/Newton%27s\\_method](https://en.wikipedia.org/wiki/Newton%27s_method)

[7]: Shaojie Bai, J. Zico Kolter, Vladlen Koltun. Multiscale Deep Equilibrium Models

[8]: Wang, Haixin, Jianlong Chang, Xiao Luo, Jinan Sun, Zhouchen Lin and Qi Tian. LION: Implicit Vision Prompt Tuning.

[9]: Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, Łukasz Kaiser. Universal Transformers.

[10]: Shaojie Bai, J. Zico Kolter, Vladlen Koltun. Trellis Networks For Sequence Modeling

# Results: Fixpoint

