

# Meta-learning methods for runtime reduction in Deep Equilibrium Models

Tanya Dixit

Australian National University

*tanya.dixit@anu.edu.au*

August 24, 2020

## 1 DEQ Origins

- Universal Transformer
- TrellisNet

## 2 Deep Equilibrium Models - Overview

- Forward Pass
- Backward Pass
- Convert any model to DEQ

## 3 Meta-learning

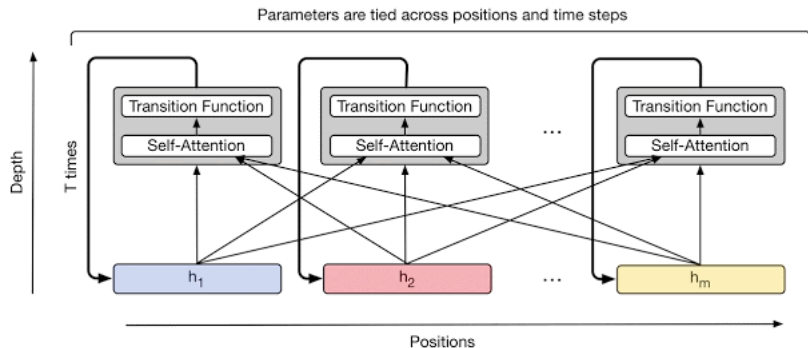
## Motivation

- Recurrence important for some tasks
- Transformer is not computationally universal (limits on seq len on num layers)
- Same computation to all inputs in transformer (complexity)

## How is Universal Transformer different and similar to Transformer?

- Parallel-in-time
- Self-attention
- Weight-tied
- Global receptive field is intact
- Each recurrent step refines representations

# DEQ Origins - Universal Transformer



# DEQ Origins - TrellisNet

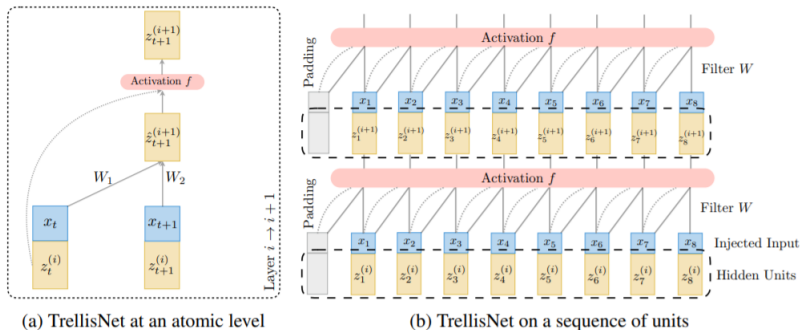
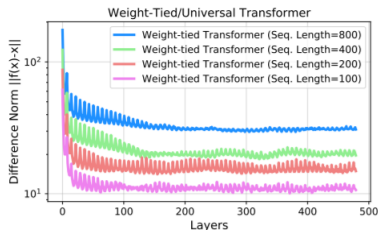
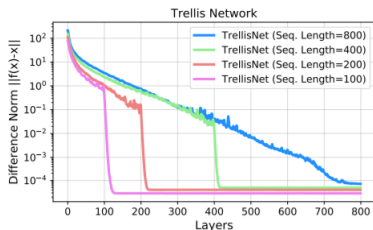


Figure 1: The interlayer transformation of TrellisNet, at an atomic level (time steps  $t$  and  $t + 1$ , layers  $i$  and  $i + 1$ ) and on a longer sequence (time steps 1 to 8, layers  $i$  and  $i + 1$ ).

The activations of different layers in TrellisNet and Universal Transformer converge, while oscillating around a fixed point/value.



# Deep Equilibrium Models - Forward Pass

As seen before, we apply a set of non-linear transformations  $T$  times and see the activations converge. We hypothesize and observe that such weight-tied models tend to converge to a fixed point as depth increases towards infinity

Take a network as a set of non-linear transformations on an input.

$$z_{1:T}^{[i+1]} = f_{\theta}(z_{1:T}^{[i]}; x_{1:T}), i = 0, \dots, L - 1, z_{1:T}^{[0]} = 0 \quad (1)$$

If considering equilibrium conditions,

$$\lim_{i \rightarrow \infty} z_{1:T}^{[i]} = \lim_{i \rightarrow \infty} f_{\theta}(z_{1:T}^{[i]}; x_{1:T}) \equiv f_{\theta}(z_{1:T}^*; x_{1:T}) = z_{1:T}^* \quad (2)$$



# Deep Equilibrium Models - Convert any model to DEQ

A traditional deep learning network is written as:

$$z^{[i+1]} = \sigma^{[i]}(W^{[i]}z^{[i]} + b^{[i]}), i = 0, \dots, L-1, z^{[0]} = x \quad (3)$$

can be written as a weight-tied, input-injected network like this. Note the repeat application of this equation will result in the same network as in equation (3)

$$\tilde{z}^{[i+1]} = \sigma(W_z \tilde{z}^{[i]} + W_x x + \tilde{b}), i = 0, \dots, L-1 \quad (4)$$

$$W_z = \begin{bmatrix} 0 & 0 & \dots & 0 & 0 \\ W^{[1]} & 0 & \dots & 0 & 0 \\ 0 & W^{[2]} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & W^{[L-1]} & 0 \end{bmatrix} \quad (5)$$

# Deep Equilibrium Models - Convert any model to DEQ

$$W_x = \begin{bmatrix} W^{[0]} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (6)$$

$$\tilde{b} = \begin{bmatrix} b^{[0]} \\ b^{[1]} \\ \vdots \\ b^{[L-1]} \end{bmatrix} \quad (7)$$

$$\sigma = \begin{bmatrix} \sigma^{[0]} \\ \sigma^{[1]} \\ \vdots \\ \sigma^{[L-1]} \end{bmatrix} \quad (8)$$

# Deep Equilibrium Models - Convert any model to DEQ

$$\tilde{z}^{[1]} = \begin{bmatrix} \sigma^{[0]} \\ \sigma^{[1]} \\ \sigma^{[2]} \\ \sigma^{[3]} \end{bmatrix} \left( \begin{bmatrix} 0 & 0 & 0 & 0 \\ W^{[1]} & 0 & 0 & 0 \\ 0 & W^{[2]} & 0 & 0 \\ 0 & 0 & 0 & W^{[3]} \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} W^{[0]} \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} x \\ 0 \\ 0 \\ 0 \end{bmatrix} + \tilde{b} \right) \quad (9)$$

$$\tilde{z}^{[2]} = \begin{bmatrix} \sigma^{[0]} \\ \sigma^{[1]} \\ \sigma^{[2]} \\ \sigma^{[3]} \end{bmatrix} \left( \begin{bmatrix} 0 & 0 & 0 & 0 \\ W^{[1]} & 0 & 0 & 0 \\ 0 & W^{[2]} & 0 & 0 \\ 0 & 0 & 0 & W^{[3]} \end{bmatrix} \begin{bmatrix} \tilde{z}^{[1]} \\ 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} W^{[0]} \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} x \\ 0 \\ 0 \\ 0 \end{bmatrix} + \tilde{b} \right) \quad (10)$$

Think of  $\tilde{z}^L$  as

$$\tilde{z}^{[4]} = \begin{bmatrix} z^{[1]} \\ z^{[2]} \\ z^{[3]} \\ z^{[4]} \end{bmatrix} \quad (11)$$

# Deep Equilibrium Models - Forward Pass

$$g_{\theta}(z_{1:T}^*; x_{1:T}) = f_{\theta}(z_{1:T}^*; x_{1:T}) - z_{1:T}^* \rightarrow 0 \quad (12)$$

$z_{1:T}^*$  is the root of  $g_{\theta}$

Equation (4) is a set of non-linear equations. Solution methods include Newton's iterative method.

## Newton iterations

$$z_{1:T}^{[i+1]} = z_{1:T}^{[i]} - \alpha B g_{\theta}(z_{1:T}^{[i]}; x_{1:T}) \quad (13)$$

But we don't want to calculate inverse jacobian ( $B$ ) at every iteration of the method - expensive operation. Solution - **Broyden Method**

# Deep Equilibrium Models - Broyden Method

Makes low-rank updates to approximate the Jacobian inverse:

$$J_{g\theta}^{-1} \big|_{z_{1:T}^{[i+1]}} \approx B_{g\theta}^{i+1} = B_{g\theta}^i + \frac{\Delta z^{[i+1]} - B_{g\theta}^{[i]} \Delta g_{\theta}^{[i+1]}}{\Delta z^{[i+1]T} B_{g\theta}^{[i]} \Delta g_{\theta}^{[i+1]}} \Delta z^{[i+1]T} B_{g\theta}^{[i]} \quad (14)$$

where

$$\Delta z^{[i+1]} = z_{1:T}^{[i+1]} - z_{1:T}^{[i]}, \Delta g_{\theta}^{[i+1]} = g_{\theta}(z_{1:T}^{[i+1]}; x_{1:T}) - g_{\theta}(z_{1:T}^{[i]}; x_{1:T})$$

# Deep Equilibrium Models - The issue with Broyden

To quote the authors - "Unlike conventional deep networks that come with a fixed number  $L$  of layers, the runtime of DEQ depends strongly on the number of Broyden steps to reach the equilibrium."

As Broyden is an iterative algorithm like Gradient Descent, we model this project on the paper "Learning to learn by gradient descent by gradient descent", where



# Meta-Learning Basic Idea

As Broyden is an iterative algorithm like Gradient Descent, we model this project on the paper "Learning to learn by gradient descent by gradient descent", where

$$\theta_{t+1} = \theta_t + g_t(\nabla f(\theta_t), \phi) \quad (15)$$

How does it learn?

We have an **optimizee** network (params  $\theta$ ) (the original neural network) and an **optimizer** network (params  $\phi$ )

We update  $\theta$  based on output from the Optimizer network - An LSTM meta-learner

Training the meta-learner by minimizing loss.

The expected loss is:

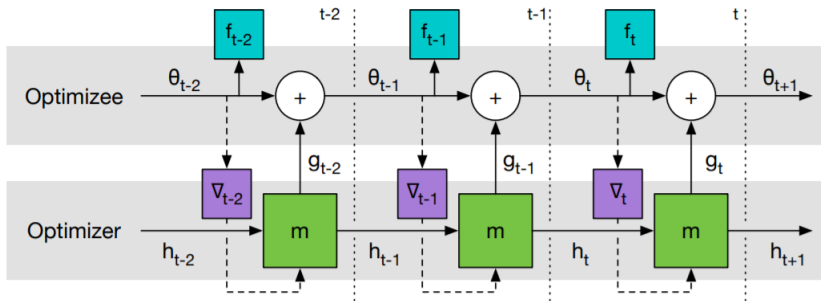
$$\mathcal{L}(\phi) = \mathbf{E}_{\mathbf{f}}[f(\theta^*(f, \phi))] = \mathbf{E}_{\mathbf{f}}[\sum_{t=1}^T w_t f(\theta_t)]$$

$$\begin{aligned} \theta_{t+1} &= \theta_t + g_t, \\ \begin{bmatrix} g_t \\ h_{t+1} \end{bmatrix} &= m(\nabla f_t, h_t, \phi) \end{aligned} \quad (16)$$

For using the above concept in this research, we formulate the LSTM meta-optimizer to take context information as inputs and try to reduce the number of iterations.

# Meta-Learning

Computational Graph unrolled at every step



# Meta-Learning Extension - Generalization

---

**Algorithm 1** Train Meta-Learner

---

**Input:** Meta-training set  $\mathcal{D}_{meta-train}$ , Learner  $M$  with parameters  $\theta$ , Meta-Learner  $R$  with parameters  $\Theta$ .

```
1:  $\Theta_0 \leftarrow$  random initialization
2:
3: for  $d = 1, n$  do
4:    $D_{train}, D_{test} \leftarrow$  random dataset from  $\mathcal{D}_{meta-train}$ 
5:    $\theta_0 \leftarrow c_0$  ▷ Initialize learner parameters
6:
7:   for  $t = 1, T$  do
8:      $\mathbf{X}_t, \mathbf{Y}_t \leftarrow$  random batch from  $D_{train}$ 
9:      $\mathcal{L}_t \leftarrow \mathcal{L}(M(\mathbf{X}_t; \theta_{t-1}), \mathbf{Y}_t)$  ▷ Get loss of learner on train batch
10:     $c_t \leftarrow R((\nabla_{\theta_{t-1}} \mathcal{L}_t, \mathcal{L}_t); \Theta_{d-1})$  ▷ Get output of meta-learner using Equation 2
11:     $\theta_t \leftarrow c_t$  ▷ Update learner parameters
12:  end for
13:
14:   $\mathbf{X}, \mathbf{Y} \leftarrow D_{test}$ 
15:   $\mathcal{L}_{test} \leftarrow \mathcal{L}(M(\mathbf{X}; \theta_T), \mathbf{Y})$  ▷ Get loss of learner on test batch
16:  Update  $\Theta_d$  using  $\nabla_{\Theta_{d-1}} \mathcal{L}_{test}$  ▷ Update meta-learner parameters
17:
18: end for
```

---



John Smith (2012)

Title of the publication

*Journal Name* 12(3), 45 – 678.

# The End