# Deep Equilibrium models

Shaojie Bai[1], J. Zico Kolter[1,2] and Vladlen Koltun[3]

Carnegie Mellon University[1], Bosch Center for AI[2], Intel Labs[3]

**TL;DR: "One layer" is all you need! An implicit-depth (or infinite-depth) sequence model that solves for equilibriums and achieves SOTA, with only constant memory.**

## Introduction

▸ Deep networks have long been built on a core concept: *layers*.

▸ **The story we all tell**: Deep learning algorithms build *hierarchical* models of input data, where earlier layers extract "simple" features and later layers create high-level abstractions of the data.

▸ Layers are usually picked by model designers: e.g., ResNet-18/101.

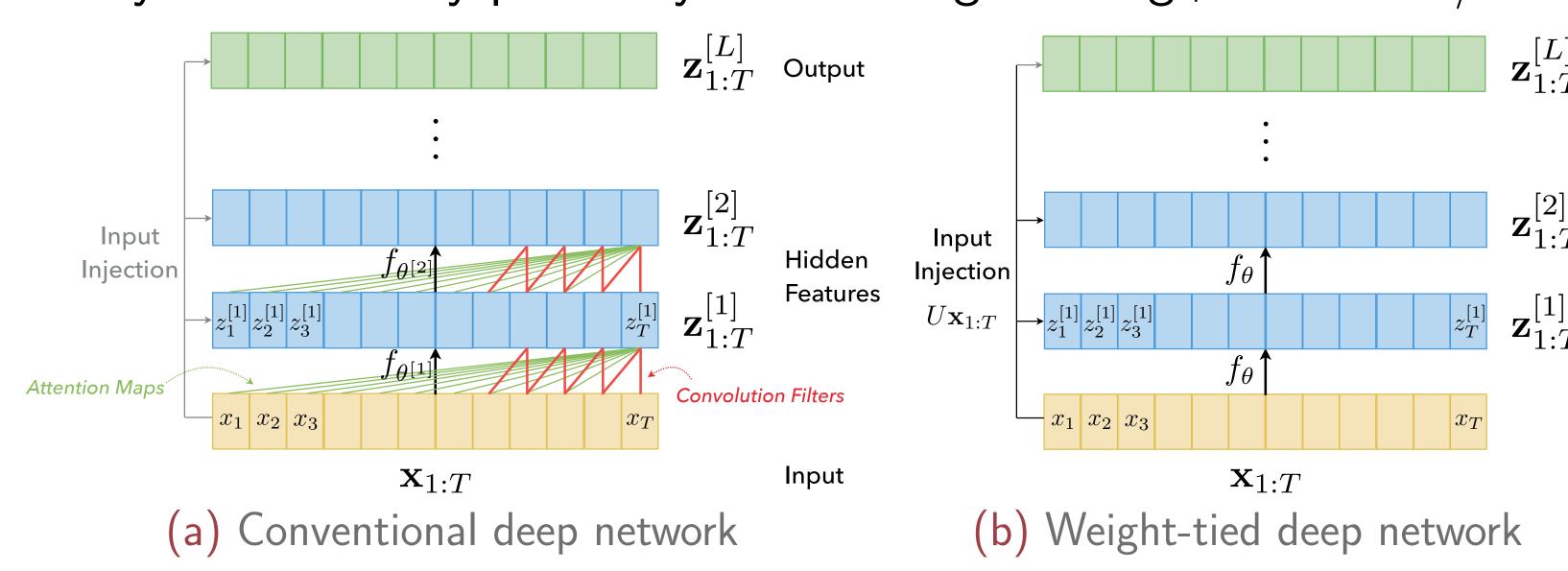(a) Conventional deep network    (b) Weight-tied deep network

Figure 1: A conventional vs. a weight-tied deep network, where the same set of parameters is shared across layers.

▸ We focuses on sequence modeling, where a typical deep feed-forward network (e.g., Transformer) can be written as:
$$\mathbf{z}_{1:T}^{[i+1]} = f_\theta^{[i]}(\mathbf{z}_{1:T}^{[i]}; \mathbf{x}_{1:T}) \quad \text{for } i = 0, 1, 2, \ldots, L-1$$

▸ Our work derives its motivation from surprising recent works that employ the *same* transformation in each layer (known as weight tying) and still achieve results competitive with the state-of-the-art:
$$\mathbf{z}_{1:T}^{[i+1]} = f_\theta(\mathbf{z}_{1:T}^{[i]}; \mathbf{x}_{1:T}), \quad \text{for } i = 0, 1, 2, \ldots, L-1$$

▸ Note that we choose to explicitly model the passthrough connections from the input $\mathbf{x}_{1:T}$ (for reasons to be clear later).

▸ What is the limit of this process? How do we model it? We propose the deep equilibrium model (DEQ) that directly computes the eventual hidden layer values of this infinite limit:
$$\mathbf{z}_{1:T}^\star = f_\theta(\mathbf{z}_{1:T}^\star; \mathbf{x}_{1:T})$$

▸ DEQ offers constant memory cost, "infinite" layers, much faster convergence to (and backpropagation from) the fixed point, and results on par with (or better than) SOTA!

▸ Code at: `https://github.com/locuslab/deq`

## Related Work

▸ Prior works have used gradient-checkpointing [4] and reversible networks [6] to reduce memory cost.

▸ Neural-ODE [3] focuses particularly on deep ResNets and models them as dynamic systems with latent, continuous trajectories.

▸ Implicit layers in deep learning can also trace back to some original work on training RNNs with recurrent backpropagation [1].

## Deep Equilibrium Sequence Models

▸ We broadly consider the class of weight-tied, input-injected deep models:
$$\mathbf{z}^{[i+1]} = f_\theta(\mathbf{z}_{1:T}^{[i]}; \mathbf{x}), \quad i = 0, \ldots, L-1, \quad \mathbf{z}^{[0]} = \mathbf{0}$$

▸ Isn't weight-tying a big restriction?
  1) Theoretically, no...
     **Theorem**[informal]: Any deep feedforward network can be represented by a weight-tied, input-injected network of equivalent depth.
  2) Empirically, no...
     **Evidence**: Recent works have well-demonstrated that weight-tied models achieve performances just as good (e.g., TrellisNet [2], Universal Transformer [5] and ALBERT [7]).

▸ Ideally, the network could have infinite depth (which is impossible in practice, due to limited memory availble in our GPUs).

▸ But *what is the limit of this process*? In practice, we observe these models tend to converge to a fixed point.
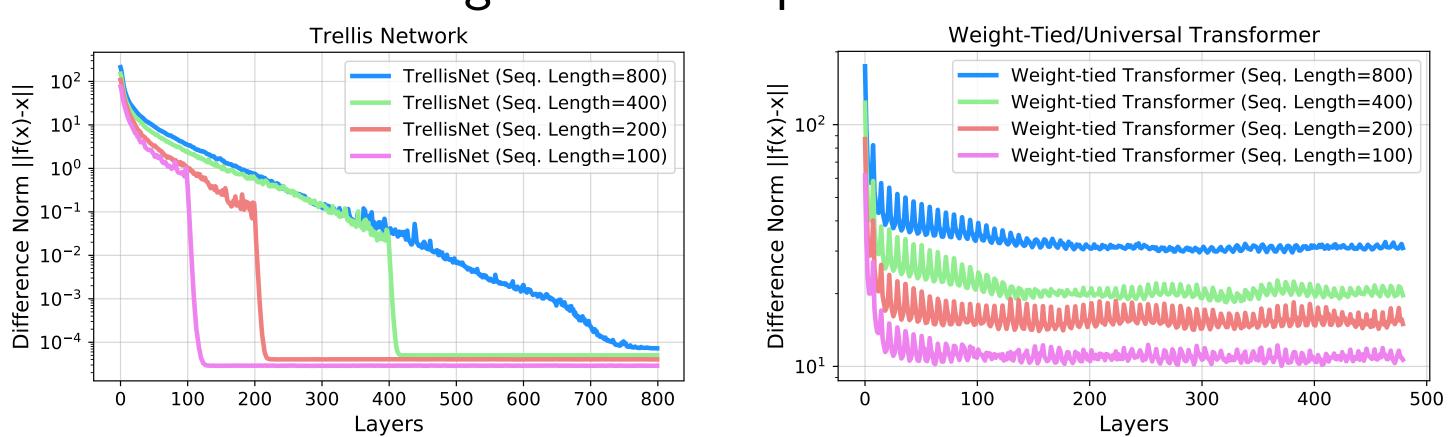
Figure 2: The convergence of intermediate activations in TrellisNet and weight-tied transformers.

▸ Intuitively, increasing depth brings "diminishing returns": each additional layer has a smaller and smaller contribution until the network reaches an equilibrium:
$$\lim_{i\to\infty} \mathbf{z}_{1:T}^{[i]} = \lim_{i\to\infty} f_\theta(\mathbf{z}_{1:T}^{[i]}; \mathbf{x}_{1:T}) \equiv f_\theta(\mathbf{z}_{1:T}^\star; \mathbf{x}_{1:T}) = \mathbf{z}_{1:T}^\star$$

▸ **Deep Equilibrium (DEQ) Model**: directly solve for this equilibrium/stable point via root-finding (e.g., Broyden's method), rather than just iterating the forward model; and directly differentiate *through* the equilibrium state.

▸ Define a single layer $f_\theta$, and $g_\theta(\mathbf{z}_{1:T}; \mathbf{x}_{1:T}) := f_\theta(\mathbf{z}_{1:T}; \mathbf{x}_{1:T}) - \mathbf{z}_{1:T}$.

Forward pass: Given input $\mathbf{x}_{1:T}$, compute the equilibrium point $\mathbf{z}_{1:T}^\star$ such that $f_\theta(\mathbf{z}_{1:T}^\star; \mathbf{x}_{1:T}) - \mathbf{z}_{1:T}^\star = \mathbf{0}$; i.e.,
$$\mathbf{z}_{1:T}^\star = \text{RootFind}(g_\theta; \mathbf{x}_{1:T})$$
(via any black-box root-finding algorithms.)

Backward pass: Implicitly differentiate through the equilibrium state, directly using the Jacobian at the equilibrium:
$$\frac{\partial \ell}{\partial (\cdot)} = -\frac{\partial \ell}{\partial \mathbf{z}_{1:T}^\star} \left(J_{g_\theta}^{-1}|_{\mathbf{z}_{1:T}^\star}\right) \frac{\partial f_\theta(\mathbf{z}_{1:T}^\star; \mathbf{x}_{1:T})}{\partial (\cdot)}$$
regardless of the path/trajectory to this equilibrium.

## Accelerating DEQ and Properties of DEQ Models

▸ One major challenge of scaling DEQ is the cost of computing/storing the exact inverse Jacobian $J_{g_\theta}^{-1}$ at every Newton step.

▸ In the forward pass, we propose to use Broyden's method (quasi-Newton) that makes low-rank updates to estimate $J_{g_\theta}^{-1}$:
$$J_{g_\theta}^{-1}|_{\mathbf{z}_{1:T}^{[i+1]}} \approx B_{g_\theta}^{[i+1]} = B_{g_\theta}^{[i]} + \mathbf{u}^{[i]}\mathbf{v}^{[i]\top}$$

▸ In the backward pass, we alternatively solve the linear system
$$\left(J_{g_\theta}^\top|_{\mathbf{z}_{1:T}^\star}\right)\mathbf{x}^\top + \left(\frac{\partial \ell}{\partial \mathbf{z}_{1:T}^\star}\right)^\top = \mathbf{0}$$
where the first term (a vector-Jacobian product) can be efficiently computed without having to explicitly form the Jacobian (which can be prohibitively large on long and high-dimensional sequences).

▸ Memory cost of DEQ: DEQ provides a constant (i.e., $O(1)$) memory cost, as it only needs to store $\mathbf{z}_{1:T}^\star$, $\mathbf{x}_{1:T}$ and $f_\theta$ (model parameters).
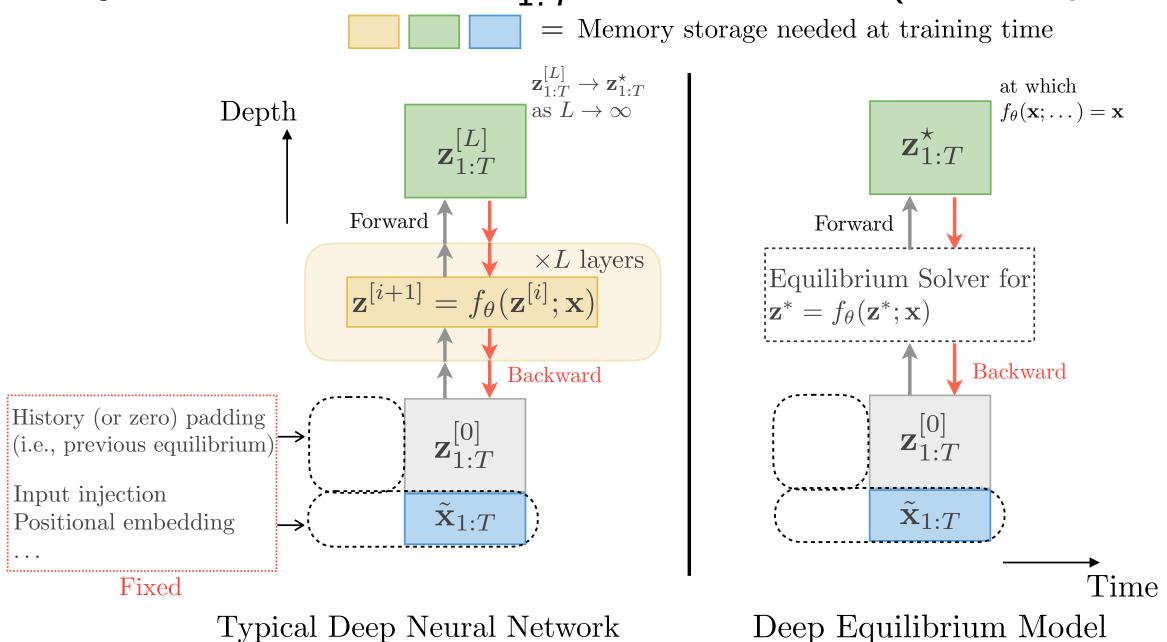
Figure 3: DEQ operates with significantly less memory than deep nets due to an analytical backward pass.

▸ Choice of $f_\theta$: The analyses above do not depend on the internal structure of $f_\theta$. Generally, $f_\theta$ should be stable and constrained so as to solve for the equilibrium in a reliable and efficient manner.

▸ If one DEQ is good, do we benefit from deeply "stacking" DEQs?
**Theorem**[informal]: Given transformations $f_{\theta[1]}$ and $\nu_{\theta[2]}$ (of potentially different function classes), $\exists$ a transformation $\Gamma_\Theta$ (where $\Theta = \theta^{[1]} \cup \theta^{[2]}$) s.t. $\text{DEQ}(\Gamma_\Theta; \mathbf{x}_{1:T}) = \text{DEQ}(\nu_{\theta[2]}; \text{DEQ}(f_{\theta[1]}; \mathbf{x}_{1:T}))$.
(In other words, stacking multiple DEQs does **not** create extra representational power over a single DEQ.)

## Instantiations of DEQ

We highlight two instantiations of $f_\theta$ based on very different SOTA sequence models.

▸ DEQ-TrellisNet [2] (weight-tied temporal convolutions that also generalize RNNs).

▸ DEQ-Transformer [5] (weight-tied multi-head self-attention).

## Experiments

▸ We test both instantiations of DEQ on both synthetic (long-range copy memory test) and realistic, large-scale and high-dimensional (word-level language modeling on PTB and WT103) sequence benchmarks.
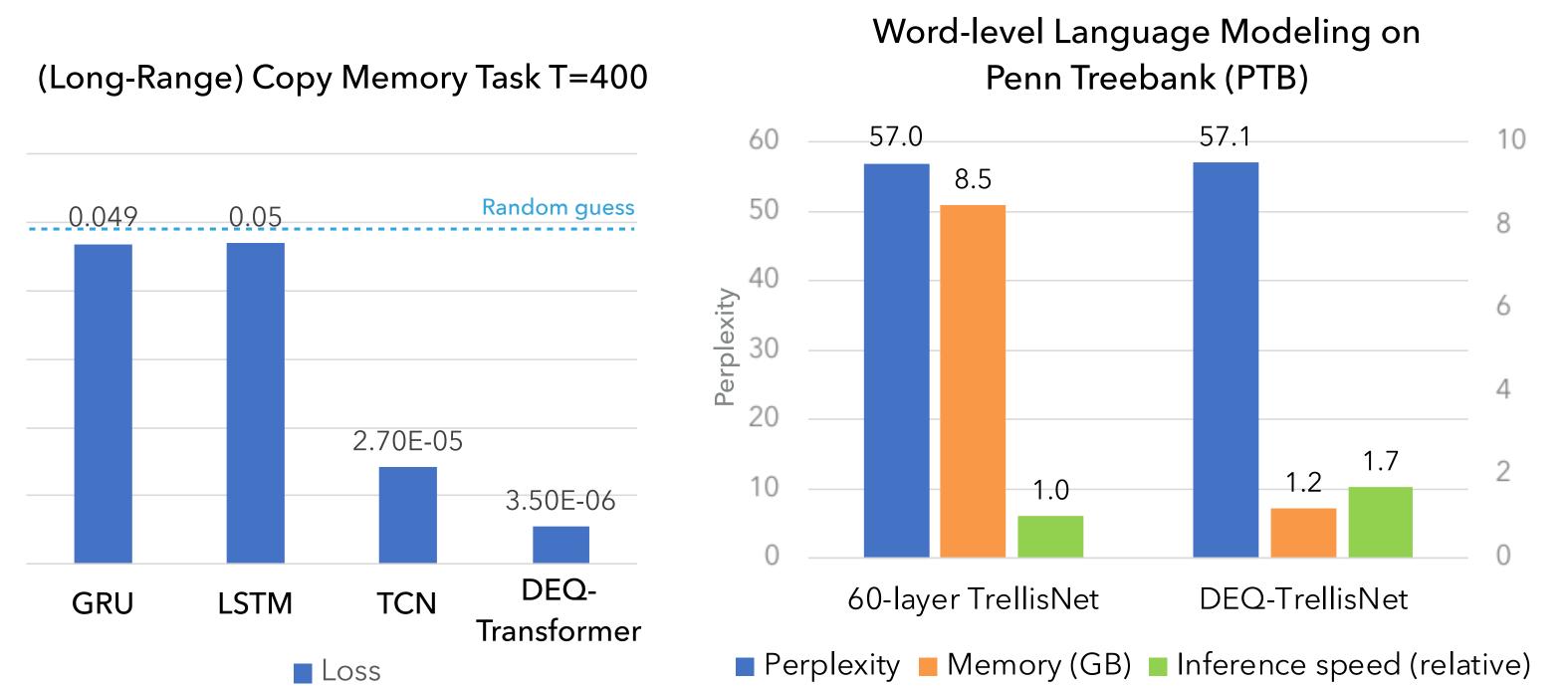
Figure 4: Left: DEQ demonstrates good memory retention over relatively long sequences. Right: DEQ achieves competitive performance on PTB corpus with > 80% reduction in memory cost.
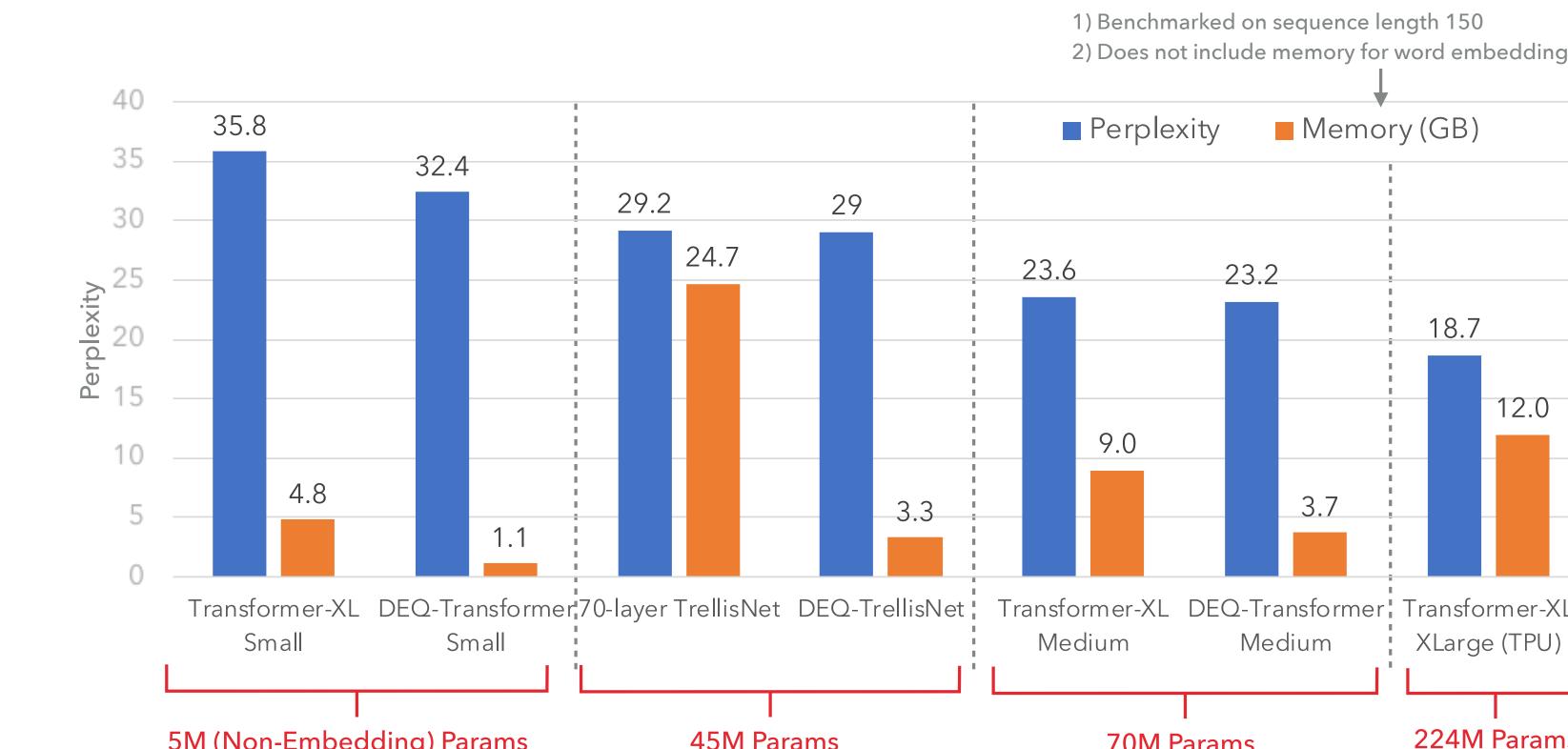
Figure 5: DEQ significantly reduces memory while improving the results on the large-scale WikiText-103 dataset.

▸ Empirically, the DEQ models achieve results on par with (or better than) the SOTA, while paying only 10-30% of the memory cost.
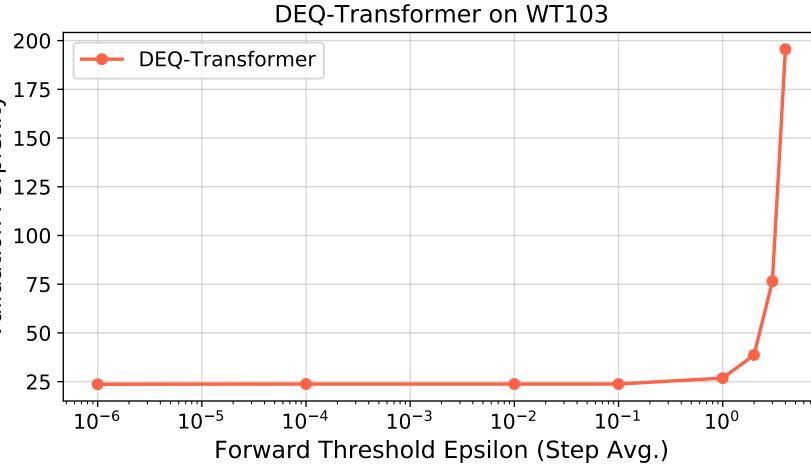
Figure 6: Left: DEQ-Transformer finds the equilibrium in a stable and efficient manner (whereas deep weight-tying could oscillate around the fixed point). Right: DEQ can be accelerated by leveraging higher tolerance $\varepsilon$ without hurting the performance.

▸ Our current models represent the largest-scale practical application of implicit layers in deep learning of which we are aware.

[1] Luis B Almeida. A learning rule for asynchronous perceptrons with feedback in a combinatorial environment. In Artificial Neural Networks, 1990.
[2] Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. Trellis networks for sequence modeling. In International Conference on Learning Representations (ICLR), 2019.
[3] Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. In Neural Information Processing Systems, 2018.
[4] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost. arXiv:1604.06174, 2016.
[5] Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. Universal transformers. International Conference on Learning Representations (ICLR), 2019.
[6] Aidan N Gomez, Mengye Ren, Raquel Urtasun, and Roger B Grosse. The reversible residual network: Backpropagation without storing activations. In Neural Information Processing Systems, 2017.
[7] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. arXiv preprint arXiv:1909.11942, 2019.