- **NLP has been undergoing a revolution**

- **How do we make Transformer more training efficient?**

Devlin et al. (2019). We pretrain our model using 1024 V100 GPUs for approximately one day.

Green AI

**$4,600,000:** The full cost of training GPT-3

| Consumption | $CO_2e$ (lbs) |
|---|---|
| Air travel, 1 passenger, NY↔SF | 1984 |
| Human life, avg, 1 year | 11,023 |
| American life, avg, 1 year | 36,156 |
| Car, avg incl. fuel, 1 lifetime | 126,000 |
| | |
| **Training one model (GPU)** | |
| NLP pipeline (parsing, SRL) | 39 |
| w/ tuning & experimentation | 78,468 |
| Transformer (big) | 192 |
| w/ neural architecture search | 626,155 |

**Artificial intelligence** / Machine learning

## Training a single AI model can emit as much carbon as five cars in their lifetimes

FACEBOOK AI

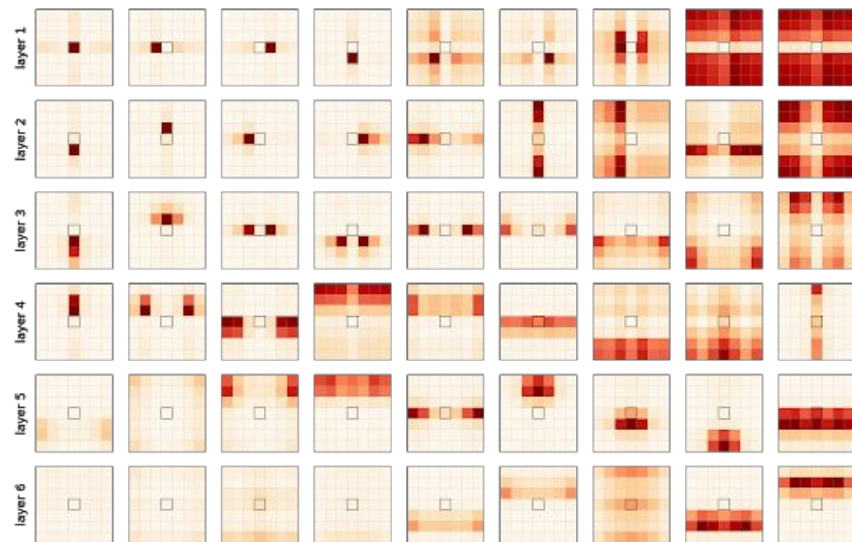- **How do we make Transformer more training efficient?**
- **What do Transformer learn?**



A Primer in BERTology: What We Know About How BERT Works

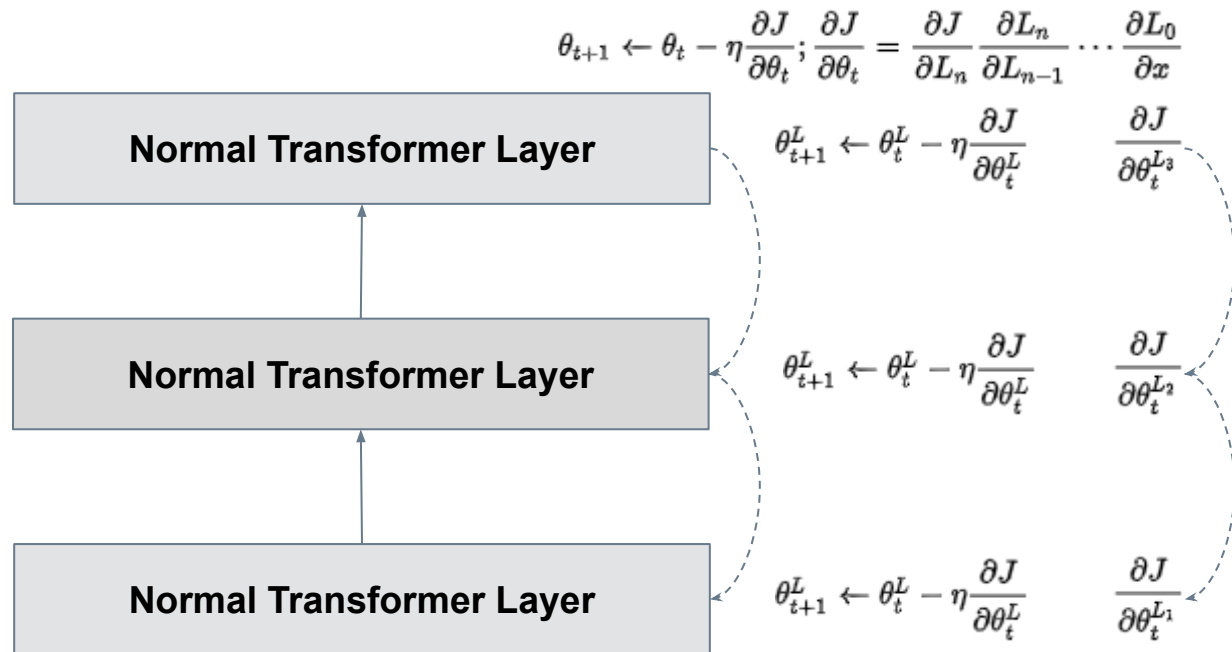Anna Rogers, Olga Kovaleva and Anna Rumshisky

Posted Online 2020
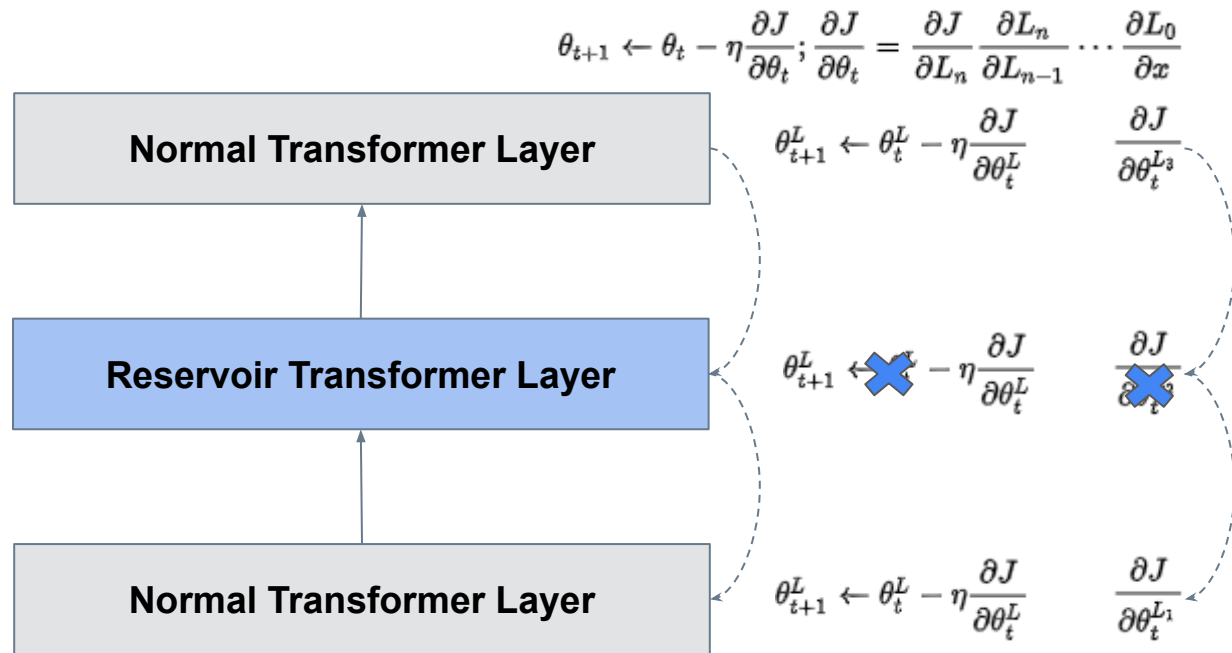https://doi.org/10.1162/tacl_a_00349

© 2020 Association for Computational Linguistics. Distributed under a CC-BY 4.0 license.

- In this work, we explore **a simple old idea** to make transformers **more training efficient**, and **understand what they learn**.

$$\theta_{t+1} \leftarrow \theta_t - \eta \frac{\partial J}{\partial \theta_t}; \frac{\partial J}{\partial \theta_t} = \frac{\partial J}{\partial L_n} \frac{\partial L_n}{\partial L_{n-1}} \cdots \frac{\partial L_0}{\partial x}$$

| Normal Transformer Layer |
| :---: |

$$\theta_{t+1}^L \leftarrow \theta_t^L - \eta \frac{\partial J}{\partial \theta_t^L} \qquad \frac{\partial J}{\partial \theta_t^{L_3}}$$

| Normal Transformer Layer |
| :---: |

$$\theta_{t+1}^L \leftarrow \theta_t^L - \eta \frac{\partial J}{\partial \theta_t^L} \qquad \frac{\partial J}{\partial \theta_t^{L_2}}$$

| Normal Transformer Layer |
| :---: |

$$\theta_{t+1}^L \leftarrow \theta_t^L - \eta \frac{\partial J}{\partial \theta_t^L} \qquad \frac{\partial J}{\partial \theta_t^{L_1}}$$

FACEBOOK AI

- In this work, we explore **a simple old idea** (**random reservoir layer**).



$$\theta_{t+1} \leftarrow \theta_t - \eta \frac{\partial J}{\partial \theta_t}; \frac{\partial J}{\partial \theta_t} = \frac{\partial J}{\partial L_n} \frac{\partial L_n}{\partial L_{n-1}} \cdots \frac{\partial L_0}{\partial x}$$

**Normal Transformer Layer**

$$\theta_{t+1}^L \leftarrow \theta_t^L - \eta \frac{\partial J}{\partial \theta_t^L} \qquad \frac{\partial J}{\partial \theta_t^{L_3}}$$

**Reservoir Transformer Layer**

$$\theta_{t+1}^L \leftarrow \qquad - \eta \frac{\partial J}{\partial \theta_t^L} \qquad \frac{\partial J}{\partial \theta_t}$$

**Normal Transformer Layer**

$$\theta_{t+1}^L \leftarrow \theta_t^L - \eta \frac{\partial J}{\partial \theta_t^L} \qquad \frac{\partial J}{\partial \theta_t^{L_1}}$$

- A random LSTM on top of pretrained word embeddings is a surprisingly good sentence encoder.

**Random**

| Model | Dim | MR | CR | MPQA | SUBJ | SST2 | TREC | SICK-R | SICK-E | MRPC | STSB |
|---|---|---|---|---|---|---|---|---|---|---|---|
| BOE | 300 | 77.3(.2) | 78.6(.3) | 87.6(.1) | 91.3(.1) | 80.0(.5) | 81.5(.8) | 80.2(.1) | 78.7(.1) | 72.9(.3) | 70.5(.1) |
| BOREP | 4096 | 77.4(.4) | 79.5(.2) | 88.3(.2) | 91.9(.2) | 81.8(.4) | **88.8(.3)** | 85.5(.1) | 82.7(.7) | 73.9(.4) | 68.5(.6) |
| RandLSTM | 4096 | 77.2(.3) | 78.7(.5) | 87.9(.1) | 91.9(.2) | 81.5(.3) | 86.5(1.1) | 85.5(.1) | 81.8(.5) | **74.1(.5)** | 72.4(.5) |
| ESN | 4096 | **78.1(.3)** | **80.0(.6)** | **88.5(.2)** | **92.6(.1)** | **83.0(.5)** | 87.9(1.0) | **86.1(.1)** | **83.1(.4)** | 73.4(.4) | **74.4(.3)** |

**Random Wide** sentence Encoder perform great on sentence classification task

Wieting, J. and Kiela, No Training Required: Exploring Random Encoders for Sentence Classification. ICLR 2019.

- Fixed random (RNN) reservoir with "readout" function trained on top.

The "echo state" approach to analysing and training recurrent neural networks – with an Erratum note[1]

Herbert Jaeger
Fraunhofer Institute for Autonomous Intelligent Systems

January 26, 2010

**Random Features for Large-Scale Kernel Machines**

**Ali Rahimi and Ben Recht**

Reservoir computing approaches to recurrent neural network training

Mantas Lukoševičius, Herbert Jaeger

FACEBOOK AI

- Goal:
  - Can we speed up transformers by freezing-and-keeping-fixed "reservoir layers"?
  - Can we do so without sacrificing (too much) task performance?

- Metrics:
  - Area under the convergence curve (f is the network, g is the eval metric, t is wall-clock):

$$\int_{t=0}^{\hat{T}} \sum_{x,y \in \mathcal{D}} g_t(f(x), y)$$

  - Train time until max validation performance
  - Test set generalization
  - Number of trainable parameters
  - Probing Task Performance

FACEBOOK AI

# Tasks

- Tasks (small, medium, big; from scratch and pretrained):
  - Machine Translation:
    - IWSLT (small)
    - WMT (medium)
  - Language Modelling:
    - Enwiki8 (big)
  - Masked Language Modelling:
    - RoBERTa pretraining (big)
    - Finetuning:
      - Sentiment: SST-2 (small)
      - Natural Language Inference: MultiNLI (medium)

FACEBOOK AI

- Comparison of different reservoir layers:
  - Regular Transformer (FFN + MultiHeadSelfAttn)

$$H = \text{MultiHeadSelfAttn}(\text{LayerNorm}(X)) + X$$
$$L = \text{FFN}(\text{LayerNorm}(H)) + H$$

  - Transformer reservoir
  - FFN reservoir, no attention
  - (CNN and BiGRU in the appendix)

- Different freezing strategies (we found **interleaving trainable and reservoir layers** to work the best)

FACEBOOK AI

- **Comparable or better final performance** of T Reservoir & FFN Reservoir



IWSLT
(test BLEU)↑

WMT
(test BLEU)↑

enwik8
(test bpc)↓

- T Reservoir & FFN Reservoir are **more training efficient** (better AUCC).



**IWSLT**
**(valid BLEU AUCC)↑**

**WMT**
**(valid BLEU AUCC)↑**

**enwik8**
**(valid bpc AUCC)↓**

FACEBOOK AI

- **21%/30% training saving** with no performance drop with FFN Reservoir.



test BLEU

Normalized Train Time Ratio

Legend: Transformer, T Reservoir, FFN Reservoir, LayerDop

- Pretraining PPL shows no big difference:



- Up to 25% time savings during pretraining for better performance:



- Downstream performance is better overall:



- Reservoir layers score better on **probing tasks**



Table 7: RoBERTa Probing Results. The line in bold text are the the frozen layers in the T Reservoir.

FACEBOOK AI

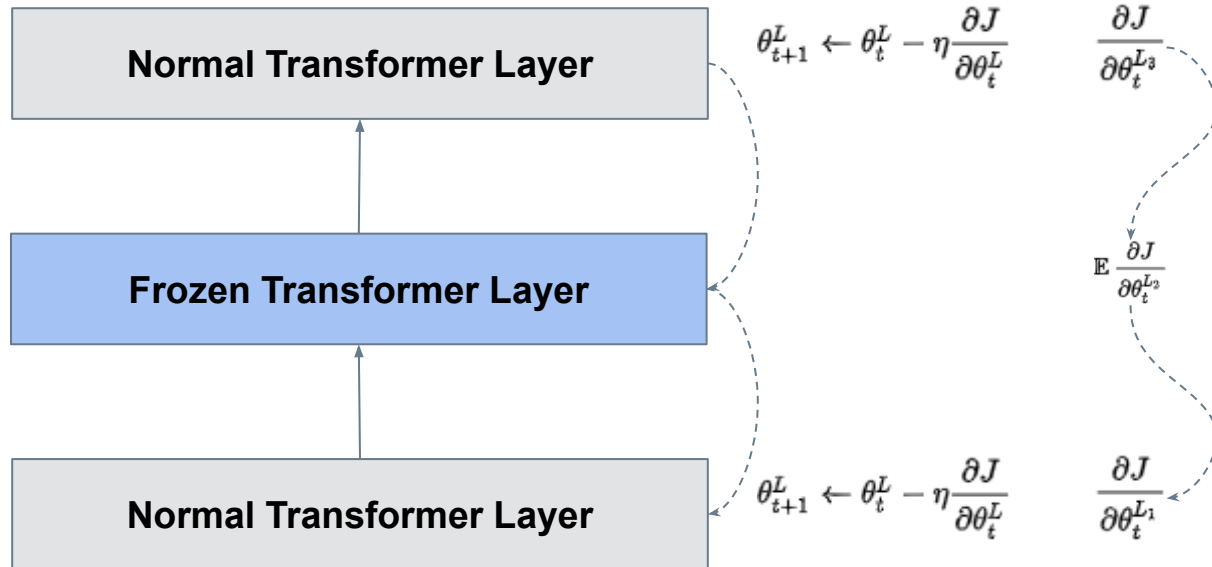# Why does it work?

- Summary:
  - **Train faster** for shorter, with **fewer parameters**, get **better results**

- Findings:
  - Reservoir layers improve efficiency
  - Appears to generalize better
    - Why?
      - Noise helps generalization
      - Detrimental noise can be removed by subsequent layers
      - "Cheap additional parameters"

- More details in the paper (ACL 2021).

FACEBOOK AI

- Estimate gradients for reservoir layers to skip **activation gradient computation**.

$$\theta_{t+1} \leftarrow \theta_t - \eta \frac{\partial J}{\partial \theta_t}; \frac{\partial J}{\partial \theta_t} = \frac{\partial J}{\partial L_n} \frac{\partial L_n}{\partial L_{n-1}} \cdots \frac{\partial L_0}{\partial x}$$

**Normal Transformer Layer**

$$\theta_{t+1}^L \leftarrow \theta_t^L - \eta \frac{\partial J}{\partial \theta_t^L} \qquad \frac{\partial J}{\partial \theta_t^{L_3}}$$

**Frozen Transformer Layer**

$$\mathbb{E} \frac{\partial J}{\partial \theta_t^{L_2}}$$

**Normal Transformer Layer**

$$\theta_{t+1}^L \leftarrow \theta_t^L - \eta \frac{\partial J}{\partial \theta_t^L} \qquad \frac{\partial J}{\partial \theta_t^{L_1}}$$

FACEBOOK AI

- Even **better training saving** with no performance drop.

| Model | Max BLEU | AUCC | Train time |
|---|---|---|---|
| Transformer | $34.59 \pm 0.11$ | $114.57 \pm 0.08$ | $142.28 \pm 1.87$ |
| T Reservoir | $34.80 \pm 0.07$ | $115.26 \pm 0.26$ | $134.49 \pm 1.70$ |
| Backskip Reservoir | $34.75 \pm 0.05$ | $115.99 \pm 0.23$ | $119.54 \pm 1.78$ |

Table 3: Validation max BLEU, AUCC at 4h and wall-clock time per epoch (averaged over multiple runs, in seconds) on IWSLT comparing backskipping with regular and reservoir transformers.

FACEBOOK AI