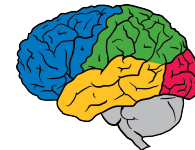# The Evolved Transformer

David R. So     Chen Liang     Quoc V. Le

# Motivation

Can we apply **Neural Architecture Search** to **feedforward sequence models**?

# Methods



- Evolution
  - simple
  - works well in vision domain

# Methods

- Evolution
  - simple
  - works well in vision domain
- **Obstacles**
  - large search space
  - high compute task

# Methods

- Evolution
  - simple
  - works well in vision domain
- **Obstacles**
  - large search space
  - high compute task
- **Solutions:**
  - First Warm Start NAS
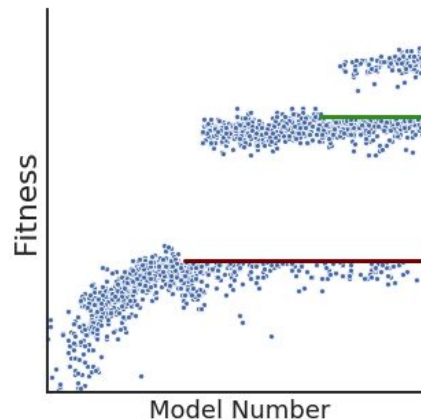
# Methods

- Evolution
    - simple
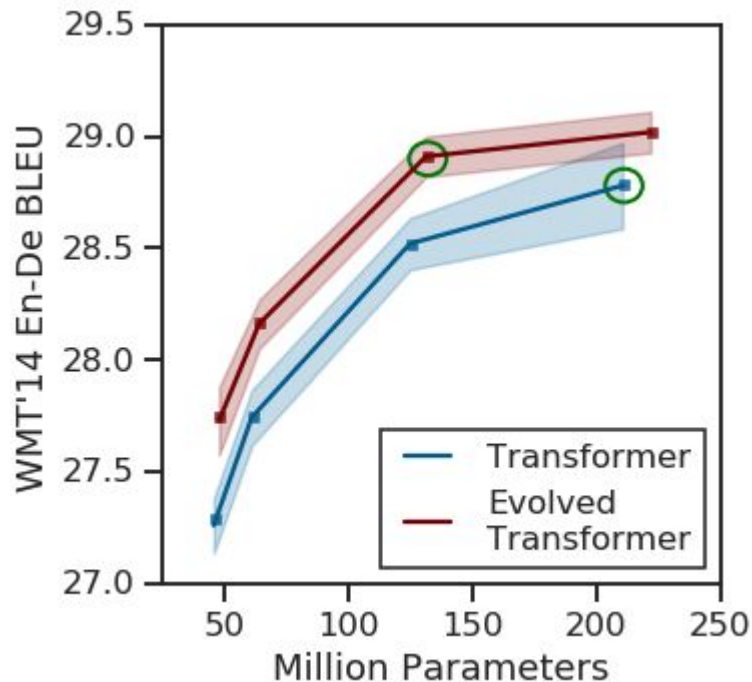    - works well in vision domain
- **Obstacles**
    - large search space
    - high compute task
- **Solutions:**
    - First Warm Start NAS
    - Progressive Dynamic Hurdles (PDH): discard bad models for cheap

# Evolved Transformer Performance



| Model | Embedding Size | BLEU | Δ BLEU |
|-------|----------------|------|--------|
| Transformer | 128 | $21.3 \pm 0.1$ | - |
| ET | 128 | $\mathbf{22.0} \pm 0.1$ | **+ 0.7** |
| Transformer | 432 | $27.3 \pm 0.1$ | - |
| ET | 432 | $\mathbf{27.7} \pm 0.1$ | + 0.4 |
| Transformer | 512 | $27.7 \pm 0.1$ | - |
| ET | 512 | $\mathbf{28.2} \pm 0.1$ | + 0.5 |
| Transformer | 768 | $28.5 \pm 0.1$ | - |
| ET | 768 | $\mathbf{28.9} \pm 0.1$ | + 0.4 |
| Transformer | 1024 | $28.8 \pm 0.2$ | - |
| ET | 1024 | $\mathbf{29.0} \pm 0.1$ | + 0.2 |

# Evolved Transformer Performance

- **State of the Art on WMT En-De**

| Work | Model | Params | BLEU | SacreBLEU (Post, 2018) |
|---|---|---|---|---|
| Gehring et al. (2017) | Convolutional Seq2Seq | 216M | 25.2 | - |
| Vaswani et al. (2017) | Transformer | 213M | 28.4 | - |
| Ahmed et al. (2017) | Weighted Transformer | 213M | 28.9 | - |
| Chen et al. (2018) | RNMT+ | 379M | 28.5 | - |
| Shaw et al. (2018) | Relative Attention Transformer | 213M | 29.2 | - |
| Ott et al. (2018) | Transformer | 210M | 29.3 | 28.6 |
| Wu et al. (2019) | Dynamic Lightweight Convolution | 213M | 29.7 | - |
| - | Evolved Transformer | 218M | **29.8** | **29.2** |

- **Generalizes to Other Tasks**

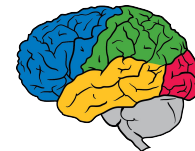| TASK | SIZE | TRAN PERP | ET PERP | TRAN BLEU | ET BLEU |
|---|---|---|---|---|---|
| WMT'14 EN-FR | BASE | $3.61 \pm 0.01$ | $\mathbf{3.42} \pm 0.01$ | $40.0 \pm 0.1$ | $\mathbf{40.6} \pm 0.1$ |
| WMT'14 EN-FR | BIG | $3.26 \pm 0.01$ | $\mathbf{3.13} \pm 0.01$ | $41.2 \pm 0.1$ | $\mathbf{41.3} \pm 0.1$ |
| WMT'14 EN-CS | BASE | $4.98 \pm 0.04$ | $\mathbf{4.42} \pm 0.01$ | $27.0 \pm 0.1$ | $\mathbf{27.6} \pm 0.2$ |
| WMT'14 EN-CS | BIG | $4.43 \pm 0.01$ | $\mathbf{4.38} \pm 0.03$ | $28.1 \pm 0.1$ | $\mathbf{28.2} \pm 0.1$ |
| LM1B | BIG | $30.44 \pm 0.04$ | $\mathbf{28.60} \pm 0.03$ | - | - |

# Architecture Comparison
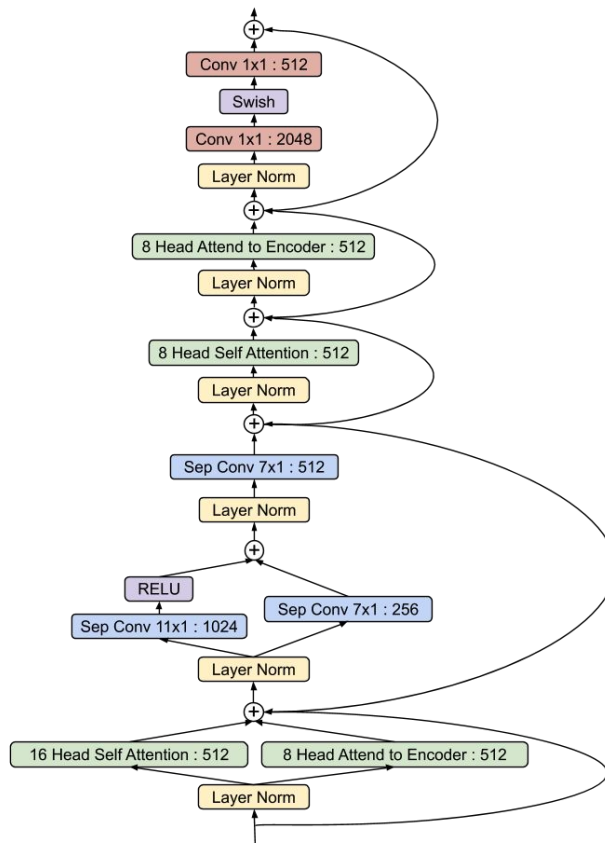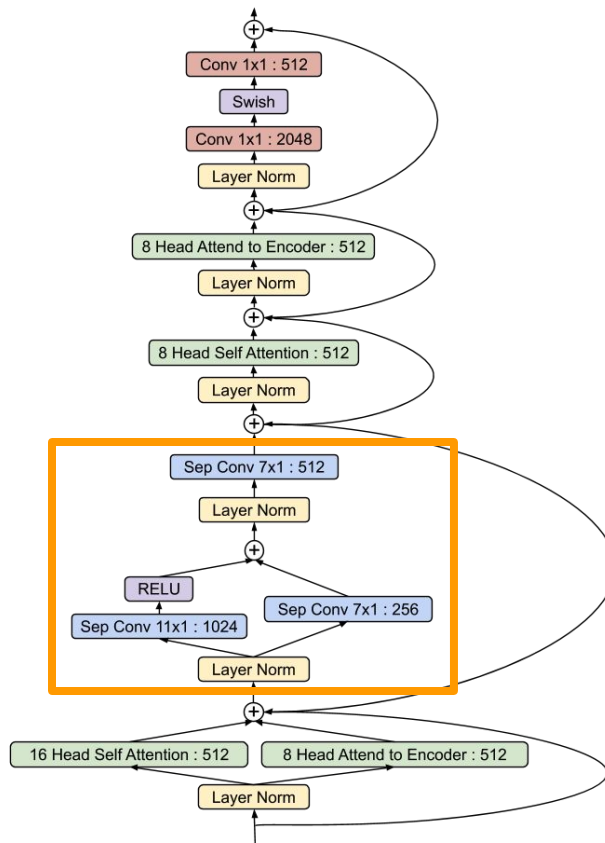
**Transformer**

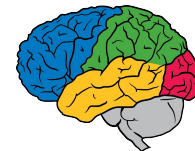# Architecture Comparison

**Evolved Transformer**

# Architecture Comparison

**Evolved Transformer**

# Summary

- First work applying NAS on feedforward sequence model.
- Discovered the Evolved Transformer, which shows better efficiency.
- Open sourced in Tensor2Tensor.

**Scan to see the paper and code.**

**Poster: Pacific Ballroom 6:30 pm to 9:00 pm**