

---

# Redesigning the Transformer Architecture with Insights from Multi-particle Dynamical Systems

---

**Subhabrata Dutta**  
Jadavpur University  
India  
subha0009@gmail.com

**Tanya Gautam**  
IIIT-Delhi  
India  
tanya18048@iiitd.ac.in

**Soumen Chakrabarti**  
IIT-Bombay  
India  
soumen@cse.iitb.ac.in

**Tanmoy Chakraborty**  
IIIT-Delhi  
India  
tanmoy@iiitd.ac.in

## Abstract

The Transformer and its variants have been proven to be efficient sequence learners in many different domains. Despite their staggering success, a critical issue has been the enormous number of parameters that must be trained (ranging from  $10^7$  to  $10^{11}$ ) along with the quadratic complexity of dot-product attention. In this work, we investigate the problem of approximating the two central components of the Transformer — multi-head self-attention and point-wise feed-forward transformation, with reduced parameter space and computational complexity. We build upon recent developments in analyzing deep neural networks as numerical solvers of ordinary differential equations. Taking advantage of an analogy between Transformer stages and the evolution of a dynamical system of multiple interacting particles, we formulate a temporal evolution scheme, TransEvoLve, to bypass costly dot-product attention over multiple stacked layers. We perform exhaustive experiments with TransEvoLve on well-known encoder-decoder as well as encoder-only tasks. We observe that the degree of approximation (or inversely, the degree of parameter reduction) has different effects on the performance, depending on the task. While in the encoder-decoder regime, TransEvoLve delivers performances comparable to the original Transformer, in encoder-only tasks it consistently outperforms Transformer along with several subsequent variants. Code is available in: <https://github.com/LCS2-IIITD/TransEvoLve>.

## 1 Introduction

Neural networks have evolved from early feed-forward and convolutional networks, to recurrent networks, to very deep and wide ‘Transformer’ networks based on attention [Vaswani et al., 2017]. Transformers and their enhancements, such as BERT [Devlin et al., 2019], T5 [Raffel et al., 2020] and GPT [Brown et al., 2020] are, by now, the default choice in many language applications. Both their training data and model sizes are massive. BERT-base has 110 million parameters. BERT-large, which often leads to better task performance, has 345 million parameters. GPT-3 has 175 billion trained parameters. Larger BERT models already approach the limits of smaller GPUs. GPT-3 is outside the resource capabilities of most research groups. Training these gargantuan models is even more challenging, with significant energy requirements and carbon emissions [Strubell et al., 2020].

In response, a growing community of researchers is focusing on post-facto reduction of model sizes, which can help with the deployment of pre-trained models in low-resource environments. However, training complexity is also critically important. A promising recent approach to faster training uses a

way of viewing layers of attention as solving ordinary differential equations (ODEs) defined over a dynamical system of interacting particles [Lu et al., 2019, Vuckovic et al., 2020]. We pursue that line of work.

Simulating particle interactions over time has a correspondence to ‘executing’ successive layers of the Transformer network. In the forward pass at successive layers, the self-attention and position-wise feed-forward operations of Transformer correspond to computing the new particle states from the previous ones. However, the numeric function learned by the  $i$ -th attention layer has zero knowledge regarding the one learned by the  $(i - 1)$ -th layer. This is counter-intuitive due to the fact that the whole evolution is temporal in nature, and this independence leads to growing numbers of trainable parameters and computing steps. We seek to develop time-evolution functionals from the initial condition alone. Such maps can then approximate the underlying ODE from parametric functions of time (the analog of network depth) and do not require computing self-attention over and over.

We propose such a scheme, leading to a network/method we call `TransEvo1ve`. It can be used for both encoder-decoder and encoder-only applications. We experiment on several tasks: neural machine translation, whole-sequence classification, and long sequence analysis with different degrees of time-evolution. `TransEvo1ve` outperforms Transformer base model on WMT 2014 English-to-French translation by 1.4 BLEU score while using 10% fewer trainable parameters. On all the encoder-only tasks, `TransEvo1ve` outperforms Transformer, as well as several strong baselines, with 50% fewer trainable parameters and more than  $3\times$  training speedup.

## 2 Related Work

Our work focuses on two primary areas of machine learning — understanding neural networks as dynamical systems and bringing down the overhead of Transformer-based models in terms of training computation and parameters.

**Neural networks and dynamical systems.** Weinan [2017] first proposed that machine learning systems can be viewed as modeling ordinary differential equations describing dynamical systems. Chang et al. [2018] explored this perspective to analyze deep residual networks. Ruthotto and Haber [2019] later extended this idea with partial differential equations. Lu et al. [2018] showed that any parametric ODE solver can be conceptualized as a deep learning framework with infinite depth. Chen et al. [2018] achieved ResNet-comparable results with a drastically lower number of parameters and memory complexity by parameterizing hidden layer derivatives and using ODE solvers. Many previous approaches applied sophisticated numerical methods for ODE approximation to build better neural networks [Haber and Ruthotto, 2017, Zhu and Fu, 2018]. Vuckovic et al. [2020] developed a mathematical formulation of self-attention as multiple interacting particles using a measure-theoretic perspective. The very first attempt to draw analogies between Transformers and dynamical systems was made by Lu et al. [2019]. They conceptualized Transformer as a numerical approximation of dynamical systems of interacting particles. However, they focused on a better approximation of the ODE with a more robust splitting scheme (with the same model size as Transformer and the dot-product attention kept intact). We seek to parameterize the temporal dynamics to bypass attention up to a certain degree.

**Efficient variations of Transformer.** Multiple approaches have been put forward to overcome the quadratic complexity of Transformers [Wang et al., 2020a, Choromanski et al., 2020, Peng et al., 2021, Xiong et al., 2021, Liu et al., 2018]. Kitaev et al. [2020] sought to use locality-sensitive hashing and reversible residual connections to deal with long range inputs. Some studies explored the sparsified attention operations to decrease computation cost [Liu et al., 2018, Ho et al., 2019, Roy et al., 2021]. These sparsification tricks can be done based on the data relevant to the task [Roy et al., 2021, Sukhbaatar et al., 2019] or in a generic manner [Liu et al., 2018, Ho et al., 2019]. Wang et al. [2020a] observed self-attention to be low-rank, and approximated an SVD decomposition of attention to linearize it. Peng et al. [2021] used random feature kernels to approximate the softmax operation on the attention matrix. Choromanski et al. [2020] sought to linearize Transformers without any prior assumption of low-rank or sparse distribution, using positive orthogonal random features. Lee-Thorp et al. [2021] achieved remarkable speedup over Transformers by substituting the attention operation with an unparameterized Fourier transform. A bulk of these works are suitable for encoder-only tasks and often incurs slower training (e.g, Peng et al. [2021] observed 15% slower training compared to Transformer). Our method overcomes both of these drawbacks.

Table 1: List of important notations and their denotations used.

Notation	Denotation
$d, d'$	Hidden and temporal dimension of the model
$\mathbf{X}^l$	Sequence of $d$ -dimensional vectors input to the $l$ -th encoder block
$T^l$	$d'$ -dimensional map of depth (time) $l$
$\mathbf{H}^l$	Output of softmax attention at $l$ -th encoder block
$W_q, W_k$	Query and key projection matrices
$\bar{W}_q, \bar{W}_k$	Temporal query and key projection matrices
$W_o$	Attention output projection matrix
$\mathbf{A}_0$	Query-key dot-product from initial values
$\mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3$	Time-evolution operators for attention
$U^l, V^l$	Random rotation matrices at depth $l$

**Transformer pruning and compression.** Knowledge distillation has been used to build light-weight *student* model from large, trained Transformer models [Behnke and Heafield, 2020, Sanh et al., 2019, Wang et al., 2020b]. Michel et al. [2019] experimented with pruning different attention heads of BERT to observe redundancy in computation. However, these methods still require a trained, parameter-heavy model to start with. Tensorization approach has shown efficient compression of Transformer-based language models [Ma et al., 2019, Khrulkov et al., 2019].

### 3 Transformers as Dynamical Systems

A single block of the Transformer encoder [Vaswani et al., 2017] is defined as a multi-head self-attention layer followed by two feed-forward layers, along with residual connections. The  $j$ -th head of self-attention operation in the  $l$ -th encoder block, with  $j \in \{1, \dots, m\}$ , on a given length- $n$  sequence of  $d$  dimensional input vectors  $\mathbf{X}^l := \{X_i^l | X_i^l \in \mathbb{R}^d\}_{i=1}^n$  can be defined as:

$$H_j^l = \text{Softmax}_i((\mathbf{X}^l W_q^l)(\mathbf{X}^l W_k^l)^\top / \sqrt{d_k})(\mathbf{X}^l W_v^l) \quad (1)$$

where  $W_q^l, W_k^l \in \mathbb{R}^{d \times d_k}$ , and  $W_v^l \in \mathbb{R}^{d \times d_v}$  are linear projection layers and  $d_k = d/m$ . Conventionally,  $d_v = d_k$ . Each  $H_j^l$  is aggregated to produce the output of the multi-head self-attention as follows:

$$\mathbf{H}^l = \text{Concat}(\{H_j^l\}_{j=1}^m)W_o^l + \mathbf{X}^l \quad (2)$$

where  $W_o \in \mathbb{R}^{d \times d}$  is a linear projection. The subsequent feed-forward transformation can be defined as:

$$\mathbf{X}^{l+1} = (\sigma(\mathbf{H}^l W_{ff1}^l + B_{ff1}^l))W_{ff2}^l + B_{ff2}^l + \mathbf{H}^l \quad (3)$$

where  $W_{ff1}^l \in \mathbb{R}^{d \times d_{ff}}$ ,  $W_{ff2}^l \in \mathbb{R}^{d_{ff} \times d}$ ,  $B_{ff1}^l \in \mathbb{R}^{d_{ff}}$ ,  $B_{ff2}^l \in \mathbb{R}^d$ , and  $\sigma(\cdot)$  is a non-linearity (Relu in [Vaswani et al., 2017] or Gelu in [Devlin et al., 2019]).

As Lu et al. [2019] argued, Equations 1-3 bear a striking resemblance with systems of interacting particles. Given the positions of a set of interacting particles as  $\mathbf{x}(t) = \{x_i(t)\}_{i=1}^n$ , the temporal evolution of such a system is denoted by the following ODE:

$$\frac{d}{dt}x_i(t) = F(x_i(t), \mathbf{x}(t), t) + G(x_i(t), t) \quad (4)$$

with the initial condition  $x_i(t_0) = s_i \in \mathbb{R}^d$ . The functions  $F$  and  $G$  are often called the *diffusion* and *convection* functions, respectively — the former models the inter-dependencies between the particles at time  $t$ , while the latter models the independent dynamics of each particle. Analytical solution of such an ODE is often impossible to find, and the most common approach is to use numerical approximation over discrete time intervals  $[t_0, t_0 + \delta t, \dots, t_0 + L\delta t]$ . Following Euler’s method of first order approximation and the Lie-Trotter splitting scheme, one can approach the numerical solution of Equation 4 from time  $t$  to  $t + \delta t$  as:

$$\begin{aligned} \tilde{x}_i(t) &= x_i(t) + \delta t F(x_i(t), \mathbf{x}(t), t) = x_i(t) + \mathcal{F}(x_i(t), \mathbf{x}(t), t) \\ x_i(t + \delta t) &= \tilde{x}_i(t) + \delta t G(\tilde{x}_i(t), t) = \tilde{x}_i(t) + \mathcal{G}(\tilde{x}_i(t), t) \end{aligned} \quad (5)$$

Equation 5 can be directly mapped to the operations of the Transformer encoder given in Equations 1-3, with  $\mathbf{X}^l \equiv \mathbf{x}(t)$ ,  $\mathbf{H}^l \equiv \{\tilde{x}_i(t)\}_{i=1}^n$  and  $\mathbf{X}^{l+1} \equiv \mathbf{x}(t + \delta t)$ .  $\mathcal{F}(\cdot, \cdot, t)$  is instantiated by the

projections  $W_q^l, W_k^l, W_v^l, W_o^l$  and  $\text{Softmax}(\cdot)$  operations in the multi-head self-attention.  $\mathcal{G}(\cdot, t)$  corresponds to the projections  $W_{ff1}^l, W_{ff2}^l, B_{ff1}^l, B_{ff1}^l$  and the  $\sigma(\cdot)$  non-linearity in Equation 3.

From the described analogies, it quickly follows that successive multi-head self-attention and the point-wise feed-forward operations follow a temporal evolution (here time is equivalent to the depth of the encoder). However, Transformer and its subsequent variants parameterize these two functions in each layer separately. This leads to a large number of parameters to be trained (ranging from  $7 \times 10^7$  in neural machine translation tasks to  $175 \times 10^9$  in language models like GPT-3). We proceed to investigate how one can leverage the temporal evolution of the diffusion and convection components to bypass this computational bottleneck.

## 4 Time-evolving Attention

As Equation 5 computes  $\mathbf{x}(t)$  iteratively from a given initial condition  $\mathbf{x}(t_0) = \mathbf{s} = \{s_i\}_{i=1}^n$ , one can reformulate the diffusion map  $\mathcal{F}(\cdot, \mathbf{x}(t), t)$  as  $\tilde{\mathcal{F}}(\cdot, f(\mathbf{s}, t))$ , i.e., as a functional of the initial condition and time only. When translated to the case of Transformers, this means one can avoid computing pairwise dot-product between  $n$  input vectors at each layer by computing a functional form at the beginning and evolving it in a temporal (depth-wise) manner. We derive this by applying dot-product self-attention on hypothetical input vectors with augmented depth information, as follows.

Let  $\mathbf{X}' = \{X'_i | X'_i \in \mathbb{R}^{d+d'}\}$  be a set of vectors such that  $X'_i = \text{Concat}(X_i^0, T^l)$ , where  $X_i^0 = \{x_1^i, \dots, x_d^i\}$  and  $T^l = \{\tau_1(l), \dots, \tau_{d'}(l)\}$ .  $W'_q, W'_k \in \mathbb{R}^{(d+d') \times (d+d')}$  are the augmented query and key projections, respectively, such that  $W'_q = [\omega_{ij}]_{i,j=1,1}^{d+d', d+d'}$ ,  $W'_k = [\theta_{ij}]_{i,j=1,1}^{d+d', d+d'}$ . The pre-softmax query-key dot product between  $X'_i$  and  $X'_j$  is given by  $a'_{ij} = (X'_i W'_q)(X'_j W'_k)^\top$ . We can decompose  $W'_q$  as concatenation of two matrices  $W_q, \tilde{W}_q$  such that  $W_q = [\omega_{ij}]_{i,j=1,1}^{d, d+d'}$  and  $\tilde{W}_q = [\omega_{ij}]_{i,j=d+1,1}^{d+d', d+d'}$ . Similarly, we decompose  $W'_k$  into  $W_k$  and  $\tilde{W}_k$ . Then  $a'_{ij}$  can be re-written as:

$$\begin{aligned} a'_{ij} &= (X_i^0 W_q)(X_j^0 W_k)^\top + (X_i^0 W_q)(T^l \tilde{W}_k)^\top + (T^l \tilde{W}_q)(X_j^0 W_k)^\top + (T^l \tilde{W}_q)(T^l \tilde{W}_k)^\top \\ &= a_{ij}^0 + A_{1i} T^{l\top} + T^l A_{2j} + A_3 (T^l \odot T^l) \end{aligned} \quad (6)$$

where  $A_{1i} = X_i^0 W_q \tilde{W}_k^\top$ ,  $A_{2j} = \tilde{W}_q W_k^\top X_j^0$ ,  $A_3 = \tilde{W}_q \tilde{W}_k^\top$  and  $\odot$  signify hadamard product. Detailed derivation is provided in Appendix A.

It is evident that  $a_{ij}^0$  is the usual dot-product pre-softmax attention between vector elements  $X_i^0, X_j^0$ . For the complete sequence of vector elements  $\mathbf{X}$ , we write  $\mathbf{A}_0 = [a_{ij}]_{ij}$ ,  $\mathbf{A}_1 = \{A_{1i}\}_i$  and  $\mathbf{A}_2 = \{A_{2j}\}_j$ . By definition,  $T^l$  is a vector function of the depth  $l$ . To construct  $T^l$  as a vector function of  $l$ , we formulate

$$T^l = \left[ w_1^l \sin\left(\frac{l}{P}\right), \dots, w_{\frac{d'}{2}}^l \sin\left(\frac{d'l}{2P}\right), w_{\frac{d'}{2}+1}^l \cos\left(\frac{l}{P}\right), \dots, w_{d'}^l \cos\left(\frac{d'l}{2P}\right) \right] \quad (7)$$

where  $W_t^l = \{w_i^l\}_{i=1}^{d'}$  are learnable parameters at depth  $l$ , and  $P = \frac{d'L}{2\pi}$ . Such a choice of  $T^l$  is intuitive in the following sense: let  $C = AT^l + B$  for some arbitrary  $A = [a_{ij}] \in \mathbb{R}^{p \times d}$ ,  $B = \{b_i\} \in \mathbb{R}^p$ , and  $C = \{c_i\} \in \mathbb{R}^p$ . Then we get

$$c_i = b_i + \sum_{j=1}^{\frac{d'}{2}} a_{ij} w_j^l \sin\left(\frac{j^l}{P}\right) + \sum_{j=1}^{\frac{d'}{2}} a_{ij} w_{j+\frac{d'}{2}}^l \cos\left(\frac{j^l}{P}\right) \quad (8)$$

In other words, any feed-forward transformation on  $T^l$  gives us a vector in which each component is a Fourier series (thereby, approximating arbitrary periodic functions of  $l$  with period  $P$ ). This enables us to encode the nonlinear transformations that  $\mathbf{X}$  undergoes in successive blocks. With this,  $\mathbf{A}_1, \mathbf{A}_2, A_3$  constitute the time-evolution operators to map depth information to the attention, computed only from the initial conditions  $X_i^0, X_j^0$ .

Figures 1(a) and 1(b) summarize the procedure. Given two input sequences  $\mathbf{X}$  and  $\mathbf{Y}$  (in case of self-attention, they are the same), we first compute  $\mathbf{A}_0, \mathbf{A}_1, \mathbf{A}_2, A_3$  using the linear projection matrices  $W_q, W_k, \tilde{W}_q, \tilde{W}_k$ , as described above. Additionally, we normalize  $\mathbf{A}_0$  by a factor  $\frac{1}{\sqrt{d/m}}$ , where  $m$  is the number of heads, similar to Transformer. Then, at any subsequent layer  $l \in \{1, \dots, L\}$ , instead

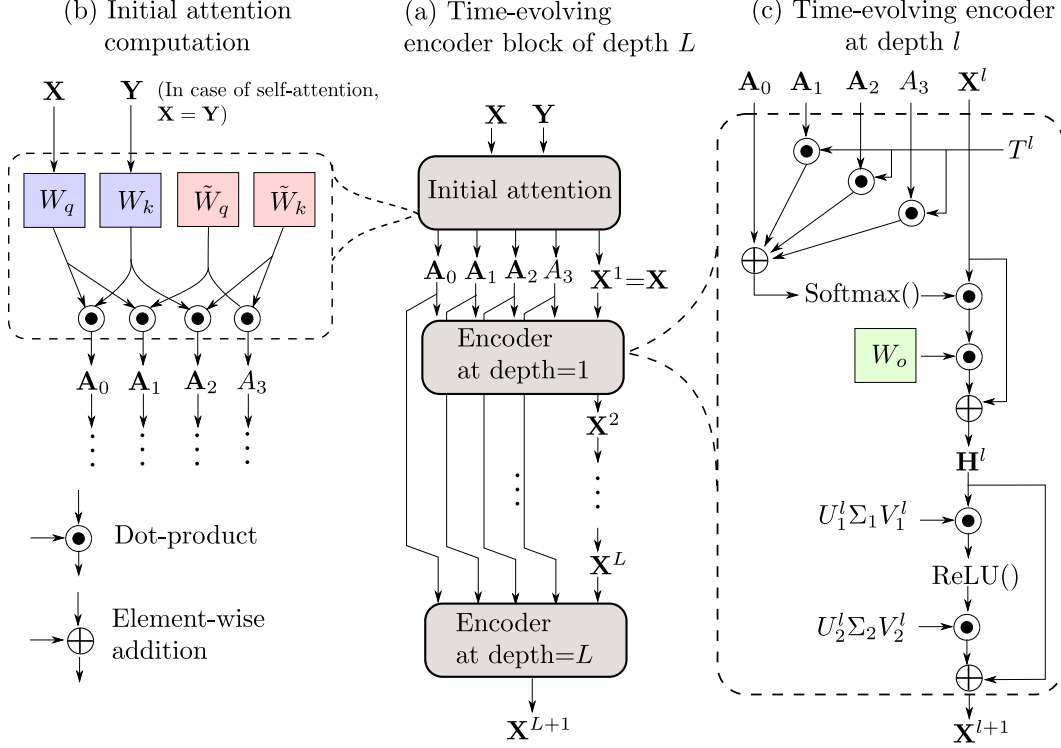


Figure 1: Dissecting the primary functional components of TransEvoLve. (a) An  $L$ -depth encoder block starts with computing (b) the initial condition matrix  $A_0$  and the evolution operator matrices  $A_1, A_2, A_3$  from the input sequence. These four are then used in (c) each encoder at depth  $l$  along with a vector function of depth,  $T^l$  to apply the attention operation on the output from the previous step. This product of attention is then passed to the feed-forward transformation actuated by the depth-dependent, random rotation matrices  $U_1^l, U_2^l, V_1^l, V_2^l$  (TransEvoLve-**randomFF**, see Section 6). In another variation, we use learnable feed-forward layers (TransEvoLve-**fullFF**).

of computing the query, key and value projections from the input  $X^l$ , we use the previously computed  $A_0, A_1, A_2, A_3$  along with  $T^l$ . Then the time-evolved attention operation becomes

$$H^l = \text{Softmax}(A_0 + A_1 T^{l\top} + T^l A_2 + A_3 (T^l \odot T^l)) X^l W_o^l + X^l \quad (9)$$

For the sake of brevity, we have not shown the concatenation of multi-headed splits as shown in Equation 2. Therefore, one should identify  $W_o^l$  in Equation 9 with Equation 2 and not as value projection as in Equation 1. We do not use value projection in our method. Also, in the multi-headed case, all the matrix-vector dot products shown in Equation 9 are computed on head-split dimension of size  $d/m$ .

**Complexity and parameter reduction.** Given a  $d$ -dimensional input sequence of length  $n$ , computing the pre-softmax attention weight matrix in a single Transformer encoder requires  $\mathcal{O}(n^2 d)$  multiplications. In our proposed method, the complexity of calculating dot-product attention (corresponding to  $a_{ij}^0$  in Equation 6) invokes computations of similar order –  $\mathcal{O}(n^2(d + d'))$ . However, this is needed only once at the beginning. The subsequent attention matrices are calculated using the components with  $T$  in Equation 6. Both  $A_1 T$  and  $T A_2$  require  $\mathcal{O}(n d')$  computation, and  $A_3 T$  requires only  $\mathcal{O}(d')$ . Therefore, if one tends to use  $L$  successive attention operations, our proposed method is computationally equivalent to a single original dot-product self-attention followed by multiple cheaper stages. In addition to this, attention weight computation using Equation 6 eliminates the need for query, key, and value projections at each self-attention block. Thereby, it ensures a parameter reduction of  $\mathcal{O}(L d^2)$  for a total stacking of  $L$ .

## 5 Time-evolving Feed-forward

The Transformer counterpart of the convection function  $\mathcal{G}(\cdot, t)$  is the point-wise feed-forward transformation in Equation 3. The complete operation constitutes two projection operations: first, the input vectors are mapped to a higher dimensional space ( $\mathbb{R}^d \rightarrow \mathbb{R}^{d_{ff}}$ ) along with a non-linear transformation, followed by projecting them back to their original dimensionality ( $\mathbb{R}^{d_{ff}} \rightarrow \mathbb{R}^d$ ). At the  $l$ -th Transformer encoder, these dimensionality transformations are achieved by the matrices  $W_{ff1}^l$  and  $W_{ff2}^l$ , respectively (see Equation 3). To construct their temporal evolution, we attempt to decompose them into time-evolving components.

Recall that any real  $m \times n$  matrix  $M$  can be decomposed as  $M = U\Sigma V^\top$  where  $U \in \mathbb{R}^{m \times m}$ ,  $V \in \mathbb{R}^{n \times n}$  are orthogonal, and  $\Sigma \in \mathbb{R}^{m \times n}$  is a rectangular diagonal matrix. However, computing exact orthogonal matrices with large dimensionality is computationally infeasible. Instead, we construct approximate rotation matrices  $U \in \mathbb{R}^{d \times d}$  as:

$$U^l = \frac{1}{\sqrt{d}} \begin{bmatrix} \sin(w_{11}^l \frac{l}{P}) & \dots & \sin(w_{1\frac{d}{2}}^l \frac{dl}{2P}) & \cos(w_{11}^l \frac{l}{P}) & \dots & \cos(w_{1\frac{d}{2}}^l \frac{dl}{2P}) \\ \vdots & & & & & \vdots \\ \sin(w_{d1}^l \frac{l}{P}) & \dots & \sin(w_{d\frac{d}{2}}^l \frac{dl}{2P}) & \cos(w_{d1}^l \frac{l}{P}) & \dots & \cos(w_{d\frac{d}{2}}^l \frac{dl}{2P}) \end{bmatrix} \quad (10)$$

where  $P = \frac{dL}{2\pi}$  and  $w_{ij}^l \in \mathcal{N}(0, d^2)$ . We discuss the properties of such a matrix  $U^l$  in the Appendix B.

We construct four matrices  $U_1^l \in \mathbb{R}^{d \times d}$ ,  $V_1^l \in \mathbb{R}^{d_{ff} \times d_{ff}}$ ,  $U_2^l \in \mathbb{R}^{d_{ff} \times d_{ff}}$ , and  $V_2^l \in \mathbb{R}^{d \times d}$  as described in Equation 10. Also, we construct two rectangular diagonal matrices  $\Sigma_1 \in \mathbb{R}^{d \times d_{ff}}$  and  $\Sigma_2 \in \mathbb{R}^{d_{ff} \times d}$  with learnable diagonal entries. With these matrices defined, one can reformulate the point-wise feed-forward operation (Equation 3) as:

$$\mathbf{X}^{l+1} = U_2^l \Sigma_2 V_2^l \sigma(U_1^l \Sigma_1 V_1^l \mathbf{H}^l + B_1) + B_2 + \mathbf{H}^l \quad (11)$$

This reformulation reduces the number of trainable parameters from  $\mathcal{O}(dd_{ff})$  to  $\mathcal{O}(d)$ .

## 6 Proposed Model: TransEvoLve

Equipped with the particle evolution based definitions of attention and point-wise feed-forward operations, we proceed to define the combined architecture of TransEvoLve. The primary building blocks of our model are the initial attention computations following Equation 6, shown in Figure 1(b) and attention operation at depth  $l$  followed by feed-forward transformations (shown together as the encoder at depth  $l$  in Figure 1(c)). An  $L$ -depth encoder block of TransEvoLve, as shown in Figure 1(a), consists of  $L$  number of attention and feed-forward operations stacked successively, preceded by an initial attention computation.

**Variations in temporal evolution of feed-forward operation.** While the re-parameterization of the point-wise feed-forward operations described in Section 5 seems to provide an astonishing reduction in the parameter size, it is natural to allow different degrees of trainable parameters for these operations. This allows us to explore the effects of time-evolving approximations of attention and point-wise feed-forward separately. We design two variations of TransEvoLve. In the original setting, the point-wise feed-forward operation is applied using random rotation matrices, following Equation 11. We call this TransEvoLve-**randomFF**. In the other variation, the time evolution process is applied for the attention operations only while the feed-forward operations are kept to be the same as Transformer (Equation 3). This variation is denoted as TransEvoLve-**fullFF**, henceforth.

**Different degrees of temporal evolution.** Recall from Section 3 that Equation 5 is a numerical approximation of Equation 4 so as to evaluate  $\mathbf{x}(t+\delta t)$  iteratively from  $\mathbf{x}(t)$ . When we attempt to approximate  $F(x_i(t), \mathbf{x}(t), t)$  as  $\tilde{F}(x_i(t), f(\mathbf{x}(t_0), t))$ , the error  $|F(x_i(t), \mathbf{x}(t), t) - \tilde{F}(x_i(t), f(\mathbf{x}(t_0), t))|$  is expected to grow as  $|t - t_0|$  grows. However, it is not feasible to compute the exact approximation error due to the complex dynamics which, in turn, varies from task to task. To compare with the original Transformer, we seek to explore this empirically. More precisely, the  $L$ -depth encoder (decoder) block with our proposed evolution strategy is expected to approximate  $L$  encoder (decoder) layers of Transformer. Following the approximation error argument, TransEvoLve is likely to bifurcate more from Transformer as  $L$  gets larger. We experiment with multiple values of  $L$  while keeping the total number of attention and feed-forward operations fixed (same as Transformer in the comparative task). To illustrate, for WMT English-to-German translation task, Transformer uses 6 encoder blocks

followed by 6 decoder blocks. Converted in our setting, one can use 1 encoder (decoder) each with depth  $L = 6$  or 2 encoders (decoders) each with depth  $L = 3$  (given that the latter requires more parameters compared to the former).

**Encoder-decoder.** To help compare with the original Transformer, we keep our design choices similar to Vaswani et al. [2017]: embedding layer tied to the output logit layer, sinusoidal position embeddings, layer normalization, etc. The temporal evolution of self-attention described in Section 4 can be straightforwardly extended to encoder-decoder attention. Given the decoder input  $\mathbf{X}_{dec}^0$  and the encoder output  $\mathbf{X}_{enc}$ , we compute the initial dot-product attention matrices  $\mathbf{A}_0^{dec}$  and  $\mathbf{A}_0^{ed}$ , corresponding to decoder self-attention and encoder-decoder attention, as follows:

$$\mathbf{A}_0^{dec} = (\mathbf{X}_{dec}^0 W_q^{dec})^\top (\mathbf{X}_{dec}^0 W_k^{dec}); \mathbf{A}_0^{ed} = (\mathbf{X}_{dec}^0 W_q^{ed})^\top (\mathbf{X}_{enc} W_k^{ed}) \quad (12)$$

If  $\mathbf{X}_{dec}^0$  and  $\mathbf{X}_{enc}$  are sequences of length  $n_{dec}$  and  $n_{enc}$ , respectively, then  $\mathbf{A}_0^{dec}$  and  $\mathbf{A}_0^{ed}$  are  $n_{dec} \times n_{dec}$  and  $n_{enc} \times n_{dec}$  matrices, respectively. We also compute the corresponding time-evolution operators  $\mathbf{A}_1^{dec}, \mathbf{A}_2^{dec}, \mathbf{A}_3^{dec}$  and  $\mathbf{A}_1^{ed}, \mathbf{A}_2^{ed}, \mathbf{A}_3^{ed}$  similar to Equation 6. Then at each depth  $l$  in the decoder block, the time-evolved attention operations are as follows:

$$\begin{aligned} \tilde{\mathbf{H}}^l &= \text{Softmax}(\mathbf{A}_0^{dec} + \mathbf{A}_1^{dec} T_{dec}^l + T_{dec}^l \mathbf{A}_2^{dec} + \mathbf{A}_3^{dec} (T_{dec}^l \odot T_{dec}^l)) \mathbf{X}_{dec}^l W_o^{l,dec} + \mathbf{X}_{dec}^l \\ \mathbf{H}^l &= \text{Softmax}(\mathbf{A}_0^{ed} + \mathbf{A}_1^{ed} T_{ed}^l + T_{ed}^l \mathbf{A}_2^{ed} + \mathbf{A}_3^{ed} (T_{ed}^l \odot T_{ed}^l)) \mathbf{X}_{enc} W_o^{l,ed} + \tilde{\mathbf{H}}^l \end{aligned} \quad (13)$$

where  $T_{dec}^l$  and  $T_{ed}^l$  are two independent vector maps of depth  $l$  representing the time-evolutions of the decoder self-attention and the encoder-decoder attention. For the sake of brevity, we have shown the full multi-head attentions in Equations 12 and 13 without showing the concatenation operation (as in Equation 2).

**Encoder-only.** For many-to-one mapping tasks (e.g., text classification), we only use the encoder part of TransEvoLve. The sequential representations returned from the encoder are applied with an average pooling along the sequence dimension and passed through a normalization layer to the final feed-forward layer to predict the classes.

## 7 Experiments

### 7.1 Model Configurations

We set  $d = d'$  across all the variations of TransEvoLve. For the encoder-decoder models, we experiment with the base version of TransEvoLve with  $d = 512$ . For encoder-only tasks, we use a small version with  $d = 256$ . In both base and small versions, the number of heads is set to 8.

As discussed in Section 6, one can vary the degree of temporal evolution (and the number of trainable parameters) in TransEvoLve by changing the value of  $L$  in the  $L$ -depth time-evolving encoder (decoder) blocks. To keep the total depth of the model constant, we need to change the number of these blocks as well. We design two such variations. Recall that the total depth of the Transformer encoder-decoder model is 12 (6 encoders and 6 decoders); in TransEvoLve, we choose (i) 1 encoder (decoder) block with depth 6 (denoted as TransEvoLve-fullFF-1, TransEvoLve-randomFF-1, etc.), and (ii) 2 encoder (decoder) blocks each with depth 3 (denoted by TransEvoLve-randomFF-2, etc.). This way, we obtain 4 variations of TransEvoLve for our experiments.

### 7.2 Tasks

We evaluate TransEvoLve over three different areas of sequence learning from texts: (i) sequence-to-sequence mapping in terms of machine translation, (ii) sequence classification, and (iii) long sequence learning. The former task requires the encoder-decoder architecture while the latter two are encoder-only.

**Machine Translation (MT).** As a sequence-to-sequence learning benchmark, we use WMT 2014 English-to-German (En-De) and English-to-French (En-Fr) translation tasks. The training data for these two tasks contain about 4.5 and 35 million sentence pairs, respectively. For both these tasks, we use the base configuration of our proposed models. We report the performance of the En-De model on WMT 2013 and 2014 En-De test sets (*newstest2013* and *newstest2014*). En-Fr model is tested only on WMT 2014 test set. Implementation details on these tasks are described in the Appendix.

**Sequence classification.** We evaluate the performance of the encoder-only version of our model on two text classification datasets: **IMDB movie-review** dataset [Maas et al., 2011] and **AGnews** topic

Model	En-De		En-Fr	#Params
	WMT 2013	WMT 2014	WMT 2014	
Transformer BASE	25.8	<b>27.3</b>	38.1	65M
TransEvoLve- randomFF-1	22.5	23.1	32.6	27M
TransEvoLve- randomFF-2	24.2	23.8	33.4	33M
TransEvoLve- fullFF-1	25.3	25.8	38.0	53M
TransEvoLve- fullFF-2	<b>26.2</b>	27.2	<b>39.5</b>	59M

Table 2: Performance of TransEvoLve variants on English-to-German (En-De) and English-to-French (En-Fr) translations in terms of BLEU scores. Transformer results are taken from the original paper [Vaswani et al., 2017].

classification dataset Zhang et al. [2015]. The IMDB dataset consists of 25000 movie reviews each for training and testing purposes. This is a binary classification task (positive/negative reviews). The AGnews dataset is a collection of news articles categorized into 4 classes (Science & Technology, Sports, Business, and World), each with 30000 and 1900 instances for training and testing, respectively. For both these tasks, we use the small version of our model. Detailed configurations and hyperparameters are given in the Appendix.

**Long sequence learning.** To test the effectiveness of our model for handling long sequences, we use two sequence classification tasks: character-level classification of the IMDB reviews and latent tree learning from sequences of arithmetic operations and operands with the **ListOps** dataset [Nangia and Bowman, 2018]. For the IMDB reviews, we set the maximum number of tokens (characters) in the training and testing examples to 4000 following [Peng et al., 2021]. In the ListOps dataset, an input sequence of arithmetic operation symbols and digits in the range 0-9 is given as input; it is a 10-way classification task of predicting the single-digit output of the input operation sequence. We consider sequences of length 500-2000 following [Peng et al., 2021]. Again, we use the small version of TransEvoLve for these two tasks. Further experimental details are provided in the Appendix.

### 7.3 Training and Testing Procedure

All experiments are done on v3-8 Cloud TPU chips. We use Adam optimizer with learning rate scheduled per gradient update step as  $lr = \frac{lr_{max}}{\sqrt{d}} \times \max(step^{-0.5}, warmup\_step^{-1.5} \times step)$ . For the MT task, we set  $lr_{max} = 1.5$  and  $warmup\_step = 16000$ . For the remaining two tasks, these values are set to 0.5 and 8000, respectively. For the translation tasks, we use a label smoothing coefficient  $\epsilon = 0.1$ . Training hyperparameters and other task-specific details are described in the Appendix. For the translation tasks, we report the BLEU scores averaged from 10 last checkpoints, each saved per 2000 update steps. For encoder-only tasks, we report the average best accuracy on five separate runs with different random seeds. Comprehensive additional results are provided in the Appendix.

## 8 Results and Discussion

**Machine Translation.** Table 2 summarizes the performance of TransEvoLve variants against Transformer (base version) on English-German and English-French translation tasks. TransEvoLve-**randomFF** versions perform poorly compared to **fullFF** versions.

Table 3: Text classification accuracy of TransEvoLve variants on AGnews and IMDB dataset. Scores of Transformer, Linformer, and Synthesizer are taken from [Tay et al., 2020a].

Model	AGnews	IMDB
Transformer	88.8	81.3
Linformer [Wang et al., 2020a]	86.5	82.8
Synthesizer [Tay et al., 2020a]	89.1	84.6
TransEvoLve-randomFF-1	90.6	87.3
TransEvoLve-randomFF-2	90.8	87.5
TransEvoLve-fullFF-1	<b>91.1</b>	86.8
TransEvoLve-fullFF-2	90.5	<b>87.6</b>

However, with less than 50% of the parameters used by Transformers, they achieve above 85% performance of that of Transformer. With all random rotation matrices replaced by standard feed-forward layers, TransEvoLve with a single 6 layers deep encoder (decoder) block performs comparably to Transformer on the En-Fr dataset. Finally, with 2 blocks of depth 3 encoders, TransEvoLve-**fullFF-2** outperforms Transformer on the WMT 2014 En-Fr dataset by 1.2

points BLEU score, despite having 10% lesser parameters. On the En-De translation task, this model performs comparable to Transformer, with 0.5 gain and 0.1 drops in BLEU scores on WMT 2013 and 2014 test datasets, respectively.



Table 4: Accuracy (%) of TransEvoLve on the long range sequence classification tasks. Speed is measured w.r.t. Transformers on the character-level IMDB dataset for input sequences of length  $1k$ ,  $2k$ ,  $3k$  and  $4k$ . All results except TransEvoLve variants are taken from [Peng et al., 2021].

Models	Tasks		Speed			
	ListOps	charIMDB	$1k$	$2k$	$3k$	$4k$
Transformer	36.4	64.3	1.0	1.0	1.0	1.0
Linformer [Wang et al., 2020a]	35.7	53.9	1.2	1.9	3.7	5.5
Reformer [Kitaev et al., 2020]	37.3	56.1	0.5	0.4	0.7	0.8
Sinkhorn [Tay et al., 2020b]	17.1	63.6	1.1	1.6	2.9	3.8
Synthesizer [Tay et al., 2020a]	37.0	61.7	1.1	1.2	2.9	1.4
Big Bird [Zaheer et al., 2020]	36.0	64.0	0.9	0.8	1.2	1.1
Linear attention [Katharopoulos et al., 2020]	16.1	65.9	1.1	<b>1.9</b>	3.7	5.6
Performers [Choromanski et al., 2020]	18.0	65.4	<b>1.2</b>	<b>1.9</b>	<b>3.8</b>	<b>5.7</b>
Random Feature Attention (RFA) [Peng et al., 2021]	36.8	66.0	1.1	1.7	3.4	5.3
TransEvoLve-randomFF-1	<b>43.2</b>	65.3	<b>1.2</b>	1.3	1.2	1.2
TransEvoLve-randomFF-2	39.1	<b>66.1</b>	1.1	1.2	1.2	1.1
TransEvoLve-fullFF-1	42.2	65.7	<b>1.2</b>	1.2	1.2	1.1
TransEvoLve-fullFF-2	37.8	65.6	1.1	1.1	1.0	1.1

**Text classification.** Table 3 summarizes the accuracy of TransEvoLve on text classification tasks. It should be noted that these results are not comparable to the state-of-the-art results on these two datasets; all four models mentioned here use no pretrained word embeddings to initialize (which is essential to achieve benchmark results for these tasks, mostly due to the smaller sizes of these datasets) or extra data for training. These results provide a comparison of purely model-specific learning capabilities. Upon that, TransEvoLve-**fullFF**-1 achieves the best performance on both datasets. For topic classification on AGnews, it scores 91.1% accuracy, outperforming Synthesizer (the best baseline) by 2%. The improvements are even more substantial on IMDB. TransEvoLve-**fullFF**-2 achieves an accuracy of 87.6% with improvements of 3% and 6.3% upon Synthesizer and Transformer, respectively.

**Long sequence tasks.** TransEvoLve shows remarkable performance in the arena of long sequence learning as well. As shown in Table 4, TransEvoLve outperforms Transformer along with previous methods on both ListOps and character-level sentiment classification on IMDB reviews. On ListOps, TransEvoLve-**randomFF**-1 establishes a new benchmark of 43.2% accuracy — beating Reformer (existing state-of-the-art) by 4.9% and Transformer by 5.8%. Moreover, all four versions of TransEvoLve show a gain in accuracy compared to the previous models. It is to be noted that we use the small version of TransEvoLve (with 256 hidden size) for these tasks, while Transformer uses the base version (512). So all these improvements are achieved while using 25% of Transformer’s parameter size. On the char-IMDB dataset, while only TransEvoLve-**randomFF**-2 outperforms the existing state-of-the-art, i.e., Random Feature Attention (RFA), other variants of TransEvoLve also turn out to be highly competitive. TransEvoLve-**randomFF**-2 achieves 66.1% accuracy, improving upon RFA by 0.1% and Transformer by 1.8%. While all the variants of TransEvoLve run faster compared to Transformer, they do not show any additional speed-up for longer sequences like RFA or Performers. This behavior is reasonable given the fact that TransEvoLve still performs softmax on  $\mathcal{O}(n^2)$  sized attention matrix at each depth. Moreover, the speedups reported in Table 4 are with the same batch size for Transformer. Practically, the lightweight TransEvoLve-**randomFF** models can handle much larger batch size (along with larger learning rates). This results in a more than  $3\times$  training speedup for all the lengths compared to Transformer. The relative gain in speed compared to Transformer or Reformer is achieved due to linear computation of pre-softmax weights (except for the initial attention calculation) and a reduced number of parameters. This is also supported by the fact that **randomFF** versions run faster than **fullFF** ones, and speed decreases with the increased number of shallower encoder blocks.

**Effects of model variations.** Random rotation matrices, instead of standard feed-forward layers and varying the number (conversely, the depth) of TransEvoLve encoder blocks, show a task-dependent effect on performance. TransEvoLve versions with a single 6-layer deep encoder perform better than 2 successive 3-layer deep encoders in the case of ListOps. In other encoder-only tasks, there are no clear winners among the two. Similar patterns can be observed among TransEvoLve-**randomFF** and TransEvoLve-**fullFF**. In the case of machine translation though, it is straightforward that having a lesser number of feed-forward parameters (**randomFF** vs. **fullFF**) and/or fewer encoder blocks

with deeper evolution deteriorate the performance. In general, the performance of TransEvoLve variations degrades with the decrease in the total parameter size on translation tasks. Encoder-only models use only self-attention, while the encoder-decoder models use self, cross, and causal attentions. From the viewpoint of multi-particle dynamical systems, the latter two are somewhat different from the former one. In self-attention, each of the ‘particles’ is interacting with each other simultaneously at a specific time-step. However, in causal attention, the interactions are ordered based on the token positions even within a specific timestep. So while decoding, there is an additional evolution of token representations at each step. Moreover, in encoder-decoder tasks as well, TransEvoLve outperforms the original Transformer in two datasets (En-De 2013 and En-Fr 2014) clearly, while providing comparable performance in another (En-De 2014). It is the random matrix versions that suffered the most in these versions. The way we designed the random matrix feedforward transformations incorporates the depth information via evolution while reducing the parameter size. In encoder-only tasks, this evolution scheme looks comparable to depth-independent, parameter-dense full versions in expressive power. Moreover, each of the tasks presents different learning requirements. Intuitively, ListOps or character-level sentiment classification tasks have a smaller input vocabulary (hence, fewer amounts of information to encode in each embedding), but longer intra-token dependency (resulting in a need for a powerful attention mechanism) compared to sentiment or topic classification at the word level. Random matrix versions provide depth-wise information sharing that may facilitate the model to better encode complex long-range dependencies, but they might remain under-parameterized to transform information-rich hidden representations. This complexity trade-off can be the possible reason behind randomFF versions, outperforming fullFF in both the long-range tasks while vice versa in text classification tasks.

These variations indicate that TransEvoLve is *not just a compressed approximation of Transformer*. Particularly in the case of encoder-only tasks, the depth-wise evolution of self-attention and feedforward projection helps TransEvoLve to learn more useful representations of the input with a much fewer number of trainable parameters. This also suggests that *the nature of the dynamical system underlying a sequence learning problem differs heavily with different tasks*. For example, the staggering performance of TransEvoLve-**randomFF-1** on the ListOps dataset implies that the diffusion component of the underlying ODE is more dominant over the convection component, and the error  $|F(x_i(t), \mathbf{x}(t), t) - \hat{F}(x_i(t), f(\mathbf{x}(t_0), t))|$  remain small with increasing  $|t - t_0|$ . One may even conjecture whether the Hölder coefficients of the functions  $F$  and  $G$  (in Equation 4) underlying the learning problem govern the performance difference. However, we leave this for future exploration.

## 9 Conclusion

Transformer stacks provide a powerful neural paradigm that gives state-of-the-art prediction quality, but they come with heavy computational requirements and large models. Drawing on an analogy between representation evolution through successive Transformer layers and dynamic particle interactions through time, we use numerical ODE solution techniques to design a computational shortcut to Transformer layers. This not only lets us save trainable parameters and training time complexity, but can also improve output quality in a variety of sequence processing tasks. Apart from building a better-performing model, this novel perspective carries the potential to uncover an in-depth understanding of sequence learning in general.

**Limitations.** It is to be noted that TransEvoLve does not get rid of quadratic operations completely, as in the case of linear attention models like Linformer or Performer. It still performs costly softmax operations on quadratic matrices. Also, at least for once (from the initial conditions) TransEvoLve computes pre-softmax  $\mathcal{O}(n^2)$  dot-product matrices. However, the significant reduction in model size compensates for this pre-existing overhead. Also, the speedup achieved by TransEvoLve is majorly relevant in the training part and many-to-one mapping. In the case of autoregressive decoding, TransEvoLve does not provide much gain over the Transformer compared to the linear versions [Peng et al., 2021]. Since the linearization approaches of the attention operation usually seeks to approximate the pair-wise attention kernel  $k(x, y)$  as some (possibly random) feature map  $\phi(x)\phi(y)$ , one may seek to design a temporal evolution of the kernel by allowing temporally evolving feature maps. We leave this as potential future work.

## Acknowledgement

The authors would like to acknowledge the support of Tensorflow Research Cloud. for generously providing the TPU access and the Google Cloud Platform for awarding GCP credits. T. Chakraborty would like to acknowledge the support of Ramanujan Fellowship, CAI, IIIT-Delhi and ihub-Anubhuti-iiitd Foundation set up under the NM-ICPS scheme of the Department of Science and Technology, India. S. Chakrabarti is partly supported by a Jagadish Bose Fellowship and a grant from IBM.

## References

- Maximiliana Behnke and Kenneth Heafield. Losing heads in the lottery: Pruning transformer attention in neural machine translation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2664–2674, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.211. URL <https://www.aclweb.org/anthology/2020.emnlp-main.211>.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>.
- Bo Chang, Lili Meng, Eldad Haber, Frederick Tung, and David Begert. Multi-level residual networks from dynamical systems view. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=SyJS-OgR->.
- Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 6572–6583, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/69386f6bb1dfed68692a24c8686939b9-Abstract.html>.
- Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamás Szepesvári, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy Colwell, and Adrian Weller. Rethinking attention with performers. *CoRR*, abs/2009.14794, 2020. URL <https://arxiv.org/abs/2009.14794>.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics, 2019. doi: 10.18653/v1/n19-1423. URL <https://doi.org/10.18653/v1/n19-1423>.
- Eldad Haber and Lars Ruthotto. Stable architectures for deep neural networks. *CoRR*, abs/1705.03341, 2017. URL <http://arxiv.org/abs/1705.03341>.
- Jonathan Ho, Nal Kalchbrenner, Dirk Weissenborn, and Tim Salimans. Axial attention in multidimensional transformers. *CoRR*, abs/1912.12180, 2019. URL <http://arxiv.org/abs/1912.12180>.
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rns: Fast autoregressive transformers with linear attention. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119

- of *Proceedings of Machine Learning Research*, pages 5156–5165. PMLR, 2020. URL <http://proceedings.mlr.press/v119/katharopoulos20a.html>.
- Valentin Khrulkov, Oleksii Hrinchuk, Leyla Mirvakhabova, and Ivan V. Oseledets. Tensorized embedding layers for efficient model compression. *CoRR*, abs/1901.10787, 2019. URL <http://arxiv.org/abs/1901.10787>.
- Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=rkgNKkHtvB>.
- James Lee-Thorp, Joshua Ainslie, Ilya Eckstein, and Santiago Ontañón. Fnet: Mixing tokens with fourier transforms. *CoRR*, abs/2105.03824, 2021. URL <https://arxiv.org/abs/2105.03824>.
- Peter J. Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. Generating wikipedia by summarizing long sequences. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=Hyg0vbWC->.
- Yiping Lu, Aoxiao Zhong, Quanzheng Li, and Bin Dong. Beyond finite layer neural networks: Bridging deep architectures and numerical differential equations. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3276–3285. PMLR, 10–15 Jul 2018. URL <http://proceedings.mlr.press/v80/lu18d.html>.
- Yiping Lu, Zhuohan Li, Di He, Zhiqing Sun, Bin Dong, Tao Qin, Liwei Wang, and Tie-Yan Liu. Understanding and improving transformer from a multi-particle dynamic system point of view. *CoRR*, abs/1906.02762, 2019. URL <http://arxiv.org/abs/1906.02762>.
- Xindian Ma, Peng Zhang, Shuai Zhang, Nan Duan, Yuexian Hou, Ming Zhou, and Dawei Song. A tensorized transformer for language modeling. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 2229–2239, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/dc960c46c38bd16e953d97cdeefdbc68-Abstract.html>.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/P11-1015>.
- Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 14014–14024, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/2c601ad9d2ff9bc8b282670cdd54f69f-Abstract.html>.
- Nikita Nangia and Samuel R. Bowman. Listops: A diagnostic dataset for latent tree learning. In Silvio Ricardo Cordeiro, Shereen Oraby, Umashanthi Pavalanathan, and Kyeongmin Rim, editors, *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 2-4, 2018, Student Research Workshop*, pages 92–99. Association for Computational Linguistics, 2018. doi: 10.18653/v1/n18-4013. URL <https://doi.org/10.18653/v1/n18-4013>.
- Hao Peng, Nikolaos Pappas, Dani Yogatama, Roy Schwartz, Noah A. Smith, and Lingpeng Kong. Random feature attention. *CoRR*, abs/2103.02143, 2021. URL <https://arxiv.org/abs/2103.02143>.

- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse attention with routing transformers. *Trans. Assoc. Comput. Linguistics*, 9:53–68, 2021. URL <https://transacl.org/ojs/index.php/tacl/article/view/2405>.
- Lars Ruthotto and Eldad Haber. Deep neural networks motivated by partial differential equations. *Journal of Mathematical Imaging and Vision*, pages 1–13, 2019.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of BERT: smaller, faster, cheaper and lighter. *CoRR*, abs/1910.01108, 2019. URL <http://arxiv.org/abs/1910.01108>.
- Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for modern deep learning research. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34 (09), pages 13693–13696, 2020.
- Sainbayar Sukhbaatar, Edouard Grave, Piotr Bojanowski, and Armand Joulin. Adaptive attention span in transformers. In Anna Korhonen, David R. Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 331–335. Association for Computational Linguistics, 2019. doi: 10.18653/v1/p19-1032. URL <https://doi.org/10.18653/v1/p19-1032>.
- Yi Tay, Dara Bahri, Donald Metzler, Da-Cheng Juan, Zhe Zhao, and Che Zheng. Synthesizer: Rethinking self-attention in transformer models. *CoRR*, abs/2005.00743, 2020a. URL <https://arxiv.org/abs/2005.00743>.
- Yi Tay, Dara Bahri, Liu Yang, Donald Metzler, and Da-Cheng Juan. Sparse Sinkhorn attention. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 9438–9447. PMLR, 13–18 Jul 2020b. URL <http://proceedings.mlr.press/v119/tay20a.html>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>.
- James Vuckovic, Aristide Baratin, and Remi Tachet des Combes. A mathematical theory of attention, 2020.
- Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *CoRR*, abs/2006.04768, 2020a. URL <https://arxiv.org/abs/2006.04768>.
- Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020b. URL <https://proceedings.neurips.cc/paper/2020/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>.
- E Weinan. A proposal on machine learning via dynamical systems. *Communications in Mathematics and Statistics*, 5(1):1–11, 2017.
- Yunyang Xiong, Zhanpeng Zeng, Rudrasis Chakraborty, Mingxing Tan, Glenn Fung, Yin Li, and Vikas Singh. Nyströmformer: A nyström-based algorithm for approximating self-attention. *CoRR*, abs/2102.03902, 2021. URL <https://arxiv.org/abs/2102.03902>.

- Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontañón, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big bird: Transformers for longer sequences. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/c8512d142a2d849725f31a9a7a361ab9-Abstract.html>.
- Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. Character-level convolutional networks for text classification. In Corinna Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 649–657, 2015. URL <https://proceedings.neurips.cc/paper/2015/hash/250cf8b51c773f3f8dc8b4be867a9a02-Abstract.html>.
- Mai Zhu and Chong Fu. Convolutional neural networks combined with runge-kutta methods. *CoRR*, abs/1802.08831, 2018. URL <http://arxiv.org/abs/1802.08831>.

## A Derivation of time-evolving attention operators

We show the full derivation of Equation 6 as follows. Let  $\mathbf{X}' = \{X'_i | X'_i \in \mathbb{R}^{d+d'}\}_{i=1}^n$  be a sequence of vectors (which is the original  $d$ -dimensional input augmented with  $d'$ -dimensional depth information). Let us further assume  $X'_i = \{x'_{ij} | x'_{ij} \in \mathbb{R}\}_{j=1}^{d+d'}$ . For two projection matrices  $W'_q, W'_k \in \mathbb{R}^{d \times (d+d')}$  where  $W'_q = [\omega_{ij}]_{i,j=1}^{d+d', d+d'}$  and  $W'_k = [\theta_{ij}]_{i,j=1}^{d+d', d+d'}$ , the query and key projections become:

$$Q_i = X'_i W'_q = \{q_{ij} | q_{ij} = \sum_{l=1}^{d+d'} x'_{il} \omega_{lj}\}_{j=1}^{d+d'}$$

$$K_i = X'_i W'_k = \{k_{ij} | k_{ij} = \sum_{l=1}^{d+d'} x'_{il} \theta_{lj}\}_{j=1}^{d+d'}$$

Then, the pre-softmax dot-product attention matrix for  $\mathbf{X}'$  becomes  $\mathbf{A}' = [a'_{ij}]_{i,j=1}^{n,n}$  where

$$\begin{aligned} a'_{ij} &= Q_i K_j = \sum_{\alpha=1}^{d+d'} (q_{i\alpha} k_{j\alpha}) \\ &= \sum_{\alpha=1}^{d+d'} \left( \sum_{\beta=1}^{d+d'} x'_{i\beta} \omega_{\beta\alpha} \sum_{\beta=1}^{d+d'} x'_{j\beta} \theta_{\beta\alpha} \right) \\ &= \sum_{\alpha=1}^{d+d'} \left( \left( \sum_{\beta=1}^d x'_{i\beta} \omega_{\beta\alpha} + \sum_{\beta=d+1}^{d+d'} x'_{i\beta} \omega_{\beta\alpha} \right) \left( \sum_{\beta=1}^d x'_{j\beta} \theta_{\beta\alpha} + \sum_{\beta=d+1}^{d+d'} x'_{j\beta} \theta_{\beta\alpha} \right) \right) \\ &= \sum_{\alpha=1}^{d+d'} \left( \sum_{\beta=1}^d x'_{i\beta} \omega_{\beta\alpha} \sum_{\beta=1}^d x'_{j\beta} \theta_{\beta\alpha} + \sum_{\beta=d+1}^{d+d'} x'_{i\beta} \omega_{\beta\alpha} \sum_{\beta=1}^d x'_{j\beta} \theta_{\beta\alpha} \right. \\ &\quad \left. + \sum_{\beta=1}^d x'_{i\beta} \omega_{\beta\alpha} \sum_{\beta=d+1}^{d+d'} x'_{j\beta} \theta_{\beta\alpha} + \sum_{\beta=d+1}^{d+d'} x'_{i\beta} \omega_{\beta\alpha} \sum_{\beta=d+1}^{d+d'} x'_{j\beta} \theta_{\beta\alpha} \right) \\ &= \sum_{\alpha=1}^{d+d'} \left( \sum_{\beta=1}^d x'_{i\beta} \omega_{\beta\alpha} \sum_{\beta=1}^d x'_{j\beta} \theta_{\beta\alpha} \right) + \sum_{\alpha=1}^{d+d'} \left( \sum_{\beta=d+1}^{d+d'} x'_{i\beta} \omega_{\beta\alpha} \sum_{\beta=1}^d x'_{j\beta} \theta_{\beta\alpha} \right) \\ &\quad + \sum_{\alpha=1}^{d+d'} \left( \sum_{\beta=1}^d x'_{i\beta} \omega_{\beta\alpha} \sum_{\beta=d+1}^{d+d'} x'_{j\beta} \theta_{\beta\alpha} \right) + \sum_{\alpha=1}^{d+d'} \left( \sum_{\beta=d+1}^{d+d'} x'_{i\beta} \omega_{\beta\alpha} \sum_{\beta=d+1}^{d+d'} x'_{j\beta} \theta_{\beta\alpha} \right) \end{aligned}$$

Recall that  $X'_i$  is the concatenation of  $X_i$  and  $T^l$ . That means, for  $1 \leq \beta \leq d$ ,  $x'_{i\beta} \in X_i = \{x_{i\gamma}\}_{\gamma=1}^d$  and for  $d+1 \leq \beta \leq d+d'$ ,  $x'_{i\beta} \in T^l = \{\tau_\gamma(l)\}_{\gamma=1}^{d'}$ . Furthermore, we decompose  $W'_q$  as concatenation of two matrices  $W_q, \tilde{W}_q$  such that  $W_q = [\omega_{ij}]_{i,j=1,1}^{d, d+d'}$  and  $\tilde{W}_q = [\omega_{ij}]_{i,j=d+1,1}^{d+d', d+d'}$ . Similarly, we decompose  $W'_k$  into  $W_k$  and  $\tilde{W}_k$ . Then the previous expression for  $a'_{ij}$  can be re-written as:

$$\begin{aligned} a'_{ij} &= \sum_{\alpha=1}^{d+d'} \left( \sum_{\gamma=1}^d x_{i\gamma} \omega_{\gamma\alpha} \sum_{\gamma=1}^d x_{j\gamma} \theta_{\gamma\alpha} \right) + \sum_{\alpha=1}^{d+d'} \left( \sum_{\gamma=1}^{d'} \tau_\gamma(l) \omega_{\gamma+d,\alpha} \sum_{\gamma=1}^d x_{j\gamma} \theta_{\gamma\alpha} \right) \\ &\quad + \sum_{\alpha=1}^{d+d'} \left( \sum_{\gamma=1}^d x_{i\gamma} \omega_{\gamma\alpha} \sum_{\gamma=1}^{d'} \tau_\gamma(l) \theta_{\gamma+d,\alpha} \right) + \sum_{\alpha=1}^{d+d'} \left( \sum_{\gamma=d+1}^{d+d'} \tau_\gamma(l) \omega_{\gamma+d,\alpha} \sum_{\gamma=d+1}^{d+d'} \tau_\gamma(l) \theta_{\gamma+d,\alpha} \right) \\ &= (X_i W_q)(X_j W_k)^\top + (X_i W_q)(T^l \tilde{W}_k)^\top + (T^l \tilde{W}_q)(X_j W_k)^\top + (\tilde{W}_q \tilde{W}_k)(T^l \odot T^l) \\ &= a_{ij} + A_{1i} T^{l\top} + T^l A_{2j} + A_3 (T^l \odot T^l) \end{aligned}$$

where  $A_{1i}$ ,  $A_{2j}$ , and  $A_3$  are  $d'$  dimensional vectors corresponding the given input vector  $X_i$ . For input vector sequence  $\mathbf{X}_i$ , these form the time-evolution operators of attention,  $\mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3$ .

## B Properties of random sine-cosine matrices

In Section 5, we redesigned a single feed-forward operation at depth  $l$  on a given input  $X_i \in \mathbb{R}^d$  to produce output  $X_{i+1} \in \mathbb{R}^{d'}$  as  $X_{i+1} = \sigma(U^l \Sigma V^l X_i + B)$  where  $U^l \in \mathbb{R}^{d \times d}$ ,  $V^l \in \mathbb{R}^{d' \times d'}$  are random sine-cosine matrices to approximate rotation,  $\Sigma \in \mathbb{R}^{d \times d'}$  is a rectangular diagonal matrix with learnable entries  $\{\lambda_j\}_{j=1}^{\min(d,d')}$ ,  $B \in \mathbb{R}^{d'}$  is a learnable bias, and  $\sigma(\cdot)$  is a non-linearity (ReLU in our case).  $U^l$  ( $V^l$ ) is defined as

$$U^l = \frac{1}{\sqrt{d}} \begin{bmatrix} \sin(w_{11}^l \frac{l}{P}) & \dots & \sin(w_{1\frac{d}{2}}^l \frac{dl}{2P}) & \cos(w_{11}^l \frac{l}{P}) & \dots & \cos(w_{1\frac{d}{2}}^l \frac{dl}{2P}) \\ \vdots & & & & & \vdots \\ \sin(w_{d1}^l \frac{l}{P}) & \dots & \sin(w_{d\frac{d}{2}}^l \frac{dl}{2P}) & \cos(w_{d1}^l \frac{l}{P}) & \dots & \cos(w_{d\frac{d}{2}}^l \frac{dl}{2P}) \end{bmatrix}$$

where  $w_{ij}^l \in \mathcal{N}(0, \sigma^2)$  and  $P = \frac{dL}{2\pi}$ .

Let  $A = U^l (U^l)^\top = [\alpha_{ij}]_{i,j=1,1}^{d,d}$ . Then for all  $1 \leq i \leq d$ ,

$$\alpha_{ii} = \sum_{j=1}^{\frac{d}{2}} \frac{1}{d} \left( \sin^2(w_{ij} \frac{jl}{P}) + \cos^2(w_{ij} \frac{jl}{P}) \right) = \frac{1}{2}$$

For all  $i \neq j$ ,

$$\begin{aligned} \alpha_{ij} &= \frac{1}{d} \sum_{k=1}^{\frac{d}{2}} \left( \sin(w_{ik} \frac{kl}{P}) \sin(w_{jk} \frac{kl}{P}) + \cos(w_{ik} \frac{kl}{P}) \cos(w_{jk} \frac{kl}{P}) \right) \\ &= \frac{1}{d} \sum_{k=1}^{\frac{d}{2}} (A_k + B_k) \end{aligned}$$

where  $A_k = \sin(w_{ik} \frac{kl}{P}) \sin(w_{jk} \frac{kl}{P})$  and  $B_k = \cos(w_{ik} \frac{kl}{P}) \cos(w_{jk} \frac{kl}{P})$ . Let  $\frac{kl}{P} = \kappa$ ; then we can rewrite  $A_k$  and  $B_k$  as:

$$\begin{aligned} A_k &= \left( \frac{\exp(\mathbf{i}w_{ik}\kappa) - \exp(-\mathbf{i}w_{ik}\kappa)}{2\mathbf{i}} \right) \left( \frac{\exp(\mathbf{i}w_{jk}\kappa) - \exp(-\mathbf{i}w_{jk}\kappa)}{2\mathbf{i}} \right) \\ &= \frac{-1}{4} (\exp(\mathbf{i}w_{ik}\kappa + \mathbf{i}w_{jk}\kappa) + \exp(-\mathbf{i}w_{ik}\kappa - \mathbf{i}w_{jk}\kappa) \\ &\quad - \exp(\mathbf{i}w_{ik}\kappa - \mathbf{i}w_{jk}\kappa) - \exp(-\mathbf{i}w_{ik}\kappa + \mathbf{i}w_{jk}\kappa)) \\ B_k &= \left( \frac{\exp(\mathbf{i}w_{ik}\kappa) + \exp(-\mathbf{i}w_{ik}\kappa)}{2} \right) \left( \frac{\exp(\mathbf{i}w_{jk}\kappa) + \exp(-\mathbf{i}w_{jk}\kappa)}{2} \right) \\ &= \frac{1}{4} (\exp(\mathbf{i}w_{ik}\kappa + \mathbf{i}w_{jk}\kappa) + \exp(-\mathbf{i}w_{ik}\kappa - \mathbf{i}w_{jk}\kappa) \\ &\quad + \exp(\mathbf{i}w_{ik}\kappa - \mathbf{i}w_{jk}\kappa) + \exp(-\mathbf{i}w_{ik}\kappa + \mathbf{i}w_{jk}\kappa)) \end{aligned}$$

Assuming  $w_{ik} \in X$  and  $w_{jk} \in Y$  where  $X$  and  $Y$  are two independent random variables with pdf defined as  $f(X) = \frac{1}{\sigma\sqrt{2\pi}} \exp(-\frac{X^2}{2\sigma^2})$  and  $f(Y) = \frac{1}{\sigma\sqrt{2\pi}} \exp(-\frac{Y^2}{2\sigma^2})$ ,

$$\begin{aligned} \mathbb{E}[\exp(\mathbf{i}w_{ik}\kappa + \mathbf{i}w_{jk}\kappa)] &= \frac{1}{2\pi\sigma^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp(\mathbf{i}X\kappa + \mathbf{i}Y\kappa) \exp(-\frac{X^2}{2\sigma^2}) \exp(-\frac{Y^2}{2\sigma^2}) dX dY \\ &= \exp(-\frac{\sigma^2}{2}\kappa) \\ &= \mathbb{E}[\exp(\mathbf{i}w_{ik}\kappa - \mathbf{i}w_{jk}\kappa)] = \mathbb{E}[\exp(-\mathbf{i}w_{ik}\kappa - \mathbf{i}w_{jk}\kappa)] \end{aligned}$$

Then

$$\mathbb{E}[A_k] = \frac{-1}{4} \left( 2 \exp(-\frac{\sigma^2}{2}\kappa) - 2 \exp(-\frac{\sigma^2}{2}\kappa) \right) = 0$$

and similarly,

$$\mathbb{E}[B_k] = \frac{1}{4} \left( 2 \exp(-\frac{\sigma^2}{2}\kappa) + 2 \exp(-\frac{\sigma^2}{2}\kappa) \right) = \exp(-\frac{\sigma^2}{2}\kappa)$$



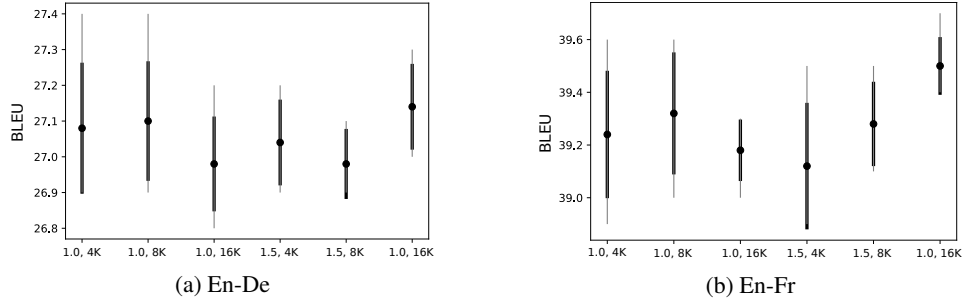


Figure 2: Variation of BLEU score for En-De (WMT 2014) and En-Fr (WMT 2014) translation with different learning rates and warmup steps. x-axis in both plots show the  $(lr_{max}, warmup\_step)$  pairs. The model variation used here in **TransEvolve-fullFF**.

Therefore,  $\mathbb{E}[\alpha_{ij}] = \frac{1}{d} \sum_{k=1}^{\frac{d}{2}} \exp(-\frac{\sigma^2}{2} \frac{kl}{P})$  which approaches 0 as  $\sigma$  gets larger. Thus, on the limiting case, we get  $\mathbb{E}[U^l(U^l)^\top] = \frac{1}{2}\mathbf{I}_d$  where  $\mathbf{I}_d$  is the  $d$ -dimensional identity matrix. This way,  $U^l$  approximates a rotation matrix as we choose  $\sigma = \mathcal{O}(d)$ .

## C Task related details

Here we describe the experimental details for encoder-decoder and encoder-only tasks. **TransEvolve** is implemented using Tensorflow version 2.4.1.

**Machine translation.** For both En-De and En-Fr tasks, we use a batch size of 512 with maximum allowed input sentence length of 256 while training and train for a total of 300,000 steps. Time needed for training varies with model configurations: **TransEvolve-randomFF-1** takes 18 hours to finish while **TransEvolve-fullFF-2** takes around 32 hours. All of these training and testings are done with 32-bit floating point precision. To find the optimal learning rate, we used the following pairs of  $(lr_{max}, warmup\_step)$  values (see Section 7.3): (1.0, 4000), (1.0, 8000), (1.0, 16000), (1.5, 4000), (1.5, 8000), and, (1.5, 16000). For all the experiments, the optimizer we use is Adam with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.98$ , and  $\epsilon = 10^{-9}$ . We used beam search with beam size 4 and length penalty 0.6. For En-De task, we used an extra decode length of 50; for En-Fr, this value is set to 35. Figure 2 summarizes the variation in performance with different  $(lr_{max}, warmup\_step)$  values; we run 5 independent training and testing with different random seeds, and choose the maximum BLEU score from each runs to plot this variation.

**Encoder-only tasks.** As mentioned in Section 7.1, we experiment with the small version of **TransEvolve** variants ( $d = 256$ ) for all the encoder-only tasks. We set the values of  $(lr_{max}, warmup\_step)$  to (0.5, 8000) and use the default parameters of Adam to optimize. All encoder-only experiments are done using a maximum input length of 512.

In the text classification regime, we use the BERT (base uncased) tokenizer from Huggingface<sup>1</sup>. The batch size is set to 80. We train each model for 15 epochs. However, the best models emerge by 7-8 epochs of training with a  $\pm 0.2\%$  error range in test accuracy over 5 randomly initialized runs.

In the long range sequence classification regime, the tokenization (character-level in IMDB and operation symbols in ListOps) and maximum input lengths are predefined. We use a batch size of 48 for the IMDB dataset, and 64 for the ListOps dataset. Again, we train all the models for 15 epochs, with best performances emerging after 9-10 epochs of training with error margins  $\pm 0.8\%$  in ListOps and  $\pm 0.3$  in IMDB datasets.

<sup>1</sup>[https://huggingface.co/transformers/model\\_doc/bert.html#berttokenizer](https://huggingface.co/transformers/model_doc/bert.html#berttokenizer)