

LEARNING TO GROW PRETRAINED MODELS FOR EFFICIENT TRANSFORMER TRAINING

Peihao Wang^{1*} Rameswar Panda² Lucas Torroba Hennigen⁴ Philip Greengard³
Leonid Karlinsky² Rogerio Feris² David D. Cox² Zhangyang Wang¹ Yoon Kim⁴

¹University of Texas at Austin, ²MIT-IBM Watson AI Lab, ³Columbia University, ⁴MIT
{peihaowang, atlaswang}@utexas.edu, {rpanda, leonidka, david.d.cox}@ibm.com,
rsferis@us.ibm.com, pg2118@columbia.edu, {lucastor, yoonkim}@mit.edu

ABSTRACT

Scaling transformers has led to significant breakthroughs in many domains, leading to a paradigm in which larger versions of existing models are trained and released on a periodic basis. New instances of such models are typically trained completely from scratch, despite the fact that they are often just scaled-up versions of their smaller counterparts. How can we use the implicit knowledge in the parameters of smaller, extant models to enable faster training of newer, larger models? This paper describes an approach for accelerating transformer training by learning to grow pretrained transformers, where we learn to linearly map the parameters of the smaller model to initialize the larger model. For tractable learning, we factorize the linear transformation as a composition of (linear) width- and depth-growth operators, and further employ a Kronecker factorization of these growth operators to encode architectural knowledge. Extensive experiments across both language and vision transformers demonstrate that our learned Linear Growth Operator (LiGO) can save up to 50% computational cost of training from scratch, while also consistently outperforming strong baselines that also reuse smaller pretrained models to initialize larger models.¹

1 INTRODUCTION

The transformer architecture (Vaswani et al., 2017) has emerged as a general purpose architecture for modeling many structured domains (Devlin et al., 2019; Brown et al., 2020; Rives et al., 2021; Dosovitskiy et al., 2021; Touvron et al., 2021a). Perhaps more so than other architectures, the transformer empirically seems to have inductive biases that make it especially amenable to scaling (Rosenfeld et al., 2019; Kaplan et al., 2020), which has led to a paradigm in which larger versions of smaller, existing models are trained and released on a periodic basis (e.g., the GPT lineage of models (Radford et al., 2018; 2019; Brown et al., 2020)). New instances of such models are typically trained completely from scratch, despite the fact that they are often scaled-up versions of their smaller counterparts. Given the compute required to train even the smaller models, we argue that training each model from scratch is wasteful, and that prior knowledge implicit in the parameters of smaller pretrained models should be leveraged to enable faster training of larger models.

One approach to this problem is through the lens of *model growth*, wherein a smaller model’s pretrained parameters are used to initialize a subset of the larger model’s parameters. While earlier works generally froze the parameters initialized from the pretrained model and only trained the new (randomly initialized) parameters (Fahlman & Lebiere, 1989; Fahlman, 1990; Gutstein et al., 2008), subsequent work has shown that copying a subset of the pretrained parameters to initialize the new parameters and then finetuning the entire network significantly accelerates training and sometimes even leads to better performance (Chen et al., 2015). When applied to modern transformers, these mechanisms roughly translate to a depth-expansion operator in which pretrained models are stacked (or combined with identity layers) to initialize deeper transformers (Gong et al., 2019; Yang et al., 2020), and a width-expansion operator in which the smaller model’s matrices are copied to initialize the larger model’s matrices (e.g., in block-diagonal fashion) (Chen et al., 2021; Gu et al., 2020).

*Work done during an internship at MIT-IBM Watson AI Lab.

¹Project page: <https://vita-group.github.io/LiGO/>

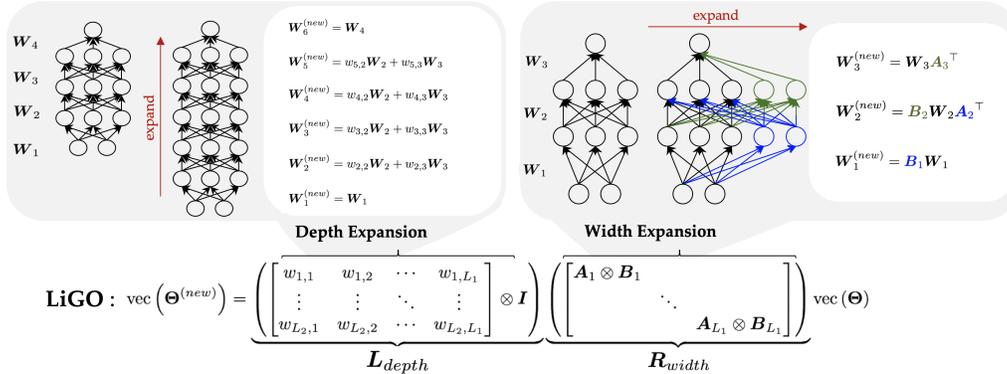


Figure 1: Our linear growth operator (LiGO) accelerates training by using the weights of a smaller model Θ to initialize the weights of the larger model $\Theta^{(new)}$. LiGO is parameterized as a sparse linear map M that can be decomposed into width- and depth-expansion operators. The width-operator R_{width} and depth-operator L_{depth} are structured matrices obtained from Kronecker products of smaller matrices which encode architectural knowledge by grouping parameters into layers and neurons. While we show the expansion operators for simple multi-layer perceptrons for illustrative purposes, in practice we apply LiGO to enable faster training of transformer networks. In our approach, we learn the growth matrix M with a 100 steps of SGD, use this to initialize the larger model, and then continue training as usual. Best viewed in color.

Noting the empirical effectiveness of such recipes, we observe that existing mechanisms generally do not have a learning component (e.g., randomly copying over neurons for width-expansion or stacking consecutive layers for depth-expansion). This paper instead proposes an efficient, data-driven approach for *learning to grow* transformers. In particular, our approach frames the problem of initializing the larger model’s parameters as learning a linear mapping from the smaller model’s parameters, i.e., $\Theta^{(large)} = M\Theta^{(small)}$ where $\Theta^{(small)}$ and $\Theta^{(large)}$ are the vectorized parameters of the small/large models. Due to the high dimensionality of the parameters, this mapping is completely intractable to learn without any restrictions on M . We thus factorize the linear mapping to be a composition of sparse width- and depth-expansion operators, $M = L_{depth}R_{width}$, where both width and depth matrices are further factorized to be a Kronecker product of smaller matrices that express architectural knowledge (e.g., through grouping parameters by layers and neurons). We show that our growth operators can represent existing approaches such as layer-stacking and neuron-copying as special cases. We find that with a small amount of learning on M (e.g., 100 gradient steps) to initialize the larger model, we can significantly accelerate training of both vision and language transformers. Figure 1 illustrates our approach.

We apply our learned linear growth operator (LiGO) to popular families of models—BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), GPT2 (Radford et al., 2019), and ViT (Dosovitskiy et al., 2021; Touvron et al., 2021a;b)—and find that LiGO can consistently improve transformer training efficiency over the traditional way of training from scratch across domains and model sizes. For instance, LiGO saves 44.7% and 22.5% FLOPs for training BERT-Base and GPT2-Medium from scratch by reusing pretrained smaller models that are half as big. Similarly, for vision transformers, when using DeiT-S (Touvron et al., 2021a) for initialization, LiGO yields 55% savings in FLOPs with no performance drop on ImageNet (Deng et al., 2009). These FLOPs savings directly translate to similar wall clock savings. We further find that models trained using LiGO achieve similar performance to the trained-from-scratch baselines when transferred to downstream tasks.

2 RELATED WORK

Efficient training. Efficient training of transformers has been studied from multiple perspectives. Some methods that are orthogonal to our work include mixed precision training (Shoeybi et al., 2019), large batch optimization (You et al., 2019), distributed training (Huang et al., 2019), and dropping layers (Zhang & He, 2020) or tokens (Hou et al., 2022). Knowledge inheritance (Qin et al., 2021) explores knowledge distillation during pretraining to efficiently learn larger transformers. Progressive training, which first trains a small transformer with few layers and then gradually expands by stacking layers, has also been applied to accelerate transformer training (Gong et al., 2019; Yang et al., 2020; Li et al., 2022; Shen et al., 2022). Net2Net Chen et al. (2015) uses function-preserving transformations to grow width by copying neurons and depth by using identity layers. Recently, bert2BERT (Chen et al., 2021) extends Net2Net to transformers. In contrast to these approaches, our approach learns to (linearly) transform the parameters of a smaller model to initialize a

larger model. While there is a line of work on learning to grow neural networks in a data-driven way, these methods are in general difficult to apply to modern-scale transformers since they (for example) involve growing a single neuron at a time or employ expensive optimization/search procedures (Wei et al., 2016; Cai et al., 2018; Wu et al., 2019; 2021; Evci et al., 2022).

Network initialization. Our work is also related to work on neural network initialization. Existing works include controlling the norm of the parameters (Mishkin & Matas, 2015; Kilcher et al., 2018; Dai et al., 2019; Wu et al., 2019; Glorot & Bengio, 2010) or replacing the normalization layers (Brock et al., 2021; Zhang et al., 2019; Huang et al., 2020). MetaInit (Dauphin & Schoenholz, 2019) proposes an automatic method that optimizes the norms of weight tensors to minimize the gradient quotient on minibatches of random Gaussian samples. GradInit (Zhu et al., 2021) learns to initialize larger networks by adjusting norm of each layer. Our work focuses on using smaller pre-trained transformers to better initialize larger transformers, which remains an understudied problem.

Structured matrices. Finally, our work is also related to structured matrices which are typically used to replace dense weight matrices for reducing training and inference computation cost. Examples include sparse and low rank matrices (Chiu et al., 2021; Han et al., 2015), Chebyshev matrices (Tang et al., 2019), Toeplitz matrices (Sindhwani et al., 2015), Kronecker-product matrices (Zhang et al., 2015), and butterfly matrices (Dao et al., 2019). A unified framework to learn a broad family of structured matrices is presented in Sindhwani et al. (2015). Dao et al. (2022) propose Monarch matrices, which inherit the expressiveness of butterfly matrices and achieve reasonable accuracy-efficiency tradeoffs in many applications. While our approach is inspired by these works, we propose to grow pretrained models by learning structured sparse linear operators with Kronecker factorization, which to our knowledge has not been explored in the literature.

3 PROPOSED APPROACH

Notation. We denote the parameters of a neural network with L layers and D dimensions as $\Theta_{L,D} = [\mathbf{W}_1 \cdots \mathbf{W}_L]^\top \in \mathbb{R}^{LD \times D}$, where $\mathbf{W}_l \in \mathbb{R}^{D \times D}$ denotes the weights for the l -th layer.² With slight abuse of notation, we denote the vectorization of $\Theta_{L,D}$ as $\text{vec}(\Theta_{L,D})^\top = [\text{vec}(\mathbf{W}_1)^\top \cdots \text{vec}(\mathbf{W}_L)^\top]^\top$.³ Our goal is to re-use the parameters $\Theta = \Theta_{L_1,D_1}$ from a pre-trained smaller model to initialize a large model $\Theta^{(new)} = \Theta_{L_2,D_2}$ through a *model growth operator* $M : \mathbb{R}^{L_1 D_1 \times D_1} \rightarrow \mathbb{R}^{L_2 D_2 \times D_2}$ that maps the weights of the smaller network to the weights of the larger one, i.e., $\Theta^{(new)} = M(\Theta)$ where $L_1 < L_2$ and $D_1 < D_2$. After model growth, we adopt $\Theta^{(new)}$ as the initialization of the large model and train it using standard recipes.

3.1 EXISTING GROWTH OPERATORS

Existing works have separately established model growth operators for depth ($L_1 < L_2, D_1 = D_2$) and width ($L_1 = L_2, D_1 < D_2$). We summarize these methods below.

Depth expansion. StackBERT (Gong et al., 2019) proposes to duplicate the smaller model to double the depth, based on the observation that upper layers share similar functionality with the lower layers. In contrast, interpolation-based depth expansion methods (Chang et al., 2017; Dong et al., 2020) interleave every layer to form a deeper model, which can be roughly interpreted as simulating a finer-grained solution to the original dynamical system from a neural ODE perspective (Chen et al., 2018). Letting $L_2 = kL_1$, the two methods’ growth operators can be formulated as:

$$\text{(StackBERT)} \mathbf{W}_i^{(new)} = \mathbf{W}_{l \bmod L_1}, \quad \text{(Interpolation)} \mathbf{W}_i^{(new)} = \mathbf{W}_{\lfloor l/k \rfloor}, \quad \forall l \in [L_2]. \quad (1)$$

Width expansion. Net2Net (Chen et al., 2015) expands the width of neural networks by randomly copying neurons while preserving output values via normalization. This can be seen as growing a matrix associated with a particular layer by duplicating the columns and rows of its weight matrix. Suppose a layer has weight matrix $\mathbf{W}_l \in \mathbb{R}^{D_1 \times D_1}$.⁴ To expand it to a matrix $\mathbf{W}_l^{(new)} \in \mathbb{R}^{D_2 \times D_2}$

²For notational brevity we assume that each hidden layer has same number of dimensions D , but LiGO can be straightforwardly generalized to layers with different dimensions (e.g., FFN layers of transformers).

³We therefore have $\text{vec}(\Theta_{L,D})^\top \in \mathbb{R}^{LD^2}$. Our approach is also agnostic with regard to vectorization order.

⁴We define a single layer as $f_l(\mathbf{x}) = \mathbf{W}_l \mathbf{x} + \mathbf{b}_l$, where the row number of \mathbf{W}_l corresponds to the output dimension, and the column number of \mathbf{W}_l corresponds to the input dimension.

($D_2 > D_1$), Net2Net copies \mathbf{W}_l to its upper-left corner of $\mathbf{W}_l^{(new)}$, fills the new columns via a random selection matrix \mathbf{S}_l , and finally duplicates and normalizes rows according to the selection matrix from the previous layer. Formally, the growth operator of Net2Net can be written as:

$$(\text{Net2Net}) \mathbf{W}_l^{(new)} = \begin{bmatrix} \mathbf{I} \\ \mathbf{S}_{l-1}^\top \end{bmatrix} \mathbf{D}_l^{-1} \mathbf{W}_l [\mathbf{I} \quad \mathbf{S}_l], \quad \mathbf{D}_l = \text{diag}(\mathbf{S}_{l-1} \mathbf{1}) + \mathbf{I}, \quad \forall l \in [L_2] \quad (2)$$

where $\mathbf{S}_l \in \{0, 1\}^{D_1 \times (D_2 - D_1)}$ is a random selection matrix. The diagonal of \mathbf{D}_l is a D_1 -dimensional histogram, whose i -th entry indicates number of times i -th column of \mathbf{W}_l was copied.

3.2 LEARNING TO GROW WITH A STRUCTURED LINEAR GROWTH OPERATOR

While existing operators have been empirically successful in accelerating transformer-based models such as BERT (Gong et al., 2019; Chen et al., 2021), we observe that generally do not have a learning component and perform the depth- and width-expansions separately. In this section we introduce a general framework for learning to grow with a linear growth operator (LiGO), which generalizes existing operators by combining the width- and depth-growth operators in a data-driven way.

We can formulate the problem of initializing the weights of the larger model $\Theta^{(new)}$ from the smaller model Θ through the following optimization problem,

$$\arg \min_M \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \mathcal{L}(\mathbf{x}; \Theta^{(new)}), \quad \text{subject to } \Theta^{(new)} = M(\Theta), \quad (3)$$

where \mathcal{D} is the data distribution and \mathcal{L} is the loss function. It is of course intractable to optimize over the entire operator space, and thus we further simplify the function M to be a linear transformation, which results in the following formulation,

$$\text{vec}(\Theta^{(new)}) = \text{vec}(M(\Theta)) = \mathbf{M} \text{vec}(\Theta), \quad \mathbf{M} \in \mathbb{R}^{L_2 D_2^2 \times L_1 D_1^2}. \quad (4)$$

This simplified objective is still completely infeasible to apply to contemporary neural networks where $L_1 D_1$ can easily be in the hundreds of millions. We therefore propose an efficient parameterization of \mathbf{M} for tractable learning.

3.2.1 DECOMPOSITION ALONG DEPTH AND WIDTH

Our first step is to decompose the LiGO operator as $\mathbf{M} = \mathbf{L}_{depth} \mathbf{R}_{width}$, where \mathbf{L}_{depth} and \mathbf{R}_{width} expand the depth and width of model separately. Concretely, we decompose \mathbf{M} as

$$\mathbf{M} = \underbrace{\begin{bmatrix} \text{diag}(\ell_{1,1}) & \cdots & \text{diag}(\ell_{1,L_1}) \\ \vdots & \ddots & \vdots \\ \text{diag}(\ell_{L_2,1}) & \cdots & \text{diag}(\ell_{L_2,L_1}) \end{bmatrix}}_{\mathbf{L}_{depth}} \underbrace{\begin{bmatrix} \mathbf{R}_1 & & \\ & \ddots & \\ & & \mathbf{R}_{L_1} \end{bmatrix}}_{\mathbf{R}_{width}}. \quad (5)$$

where $\mathbf{R}_l \in \mathbb{R}^{D_2^2 \times D_1^2}$ and $\ell_{i,j} \in \mathbb{R}^{D_2^2}$. In the above, \mathbf{L}_{depth} is an array of diagonal matrices and \mathbf{R}_{width} is a block-diagonal matrix, i.e., both matrices are highly structured and sparse. When applying \mathbf{R}_{width} to weights $\text{vec}(\Theta)$, the parameters of each layer will be transformed independently via $\text{vec}(\mathbf{W}_l^{(new)}) = \mathbf{R}_l \text{vec}(\mathbf{W}_l)$ and lifted to a higher dimension. The l -th row block of \mathbf{L}_{depth} corresponds to the growth operator of l -th layer, which amounts to linearly combining all layers of the smaller model via $\text{vec}(\mathbf{W}_l^{(new)})_k = \sum_{l'=1}^{L_1} (\ell_{l,l'})_k \text{vec}(\mathbf{W}_{l'})_k$. By this factorization, we can effectively reduce the complexity of the LiGO operator from $\mathcal{O}(D_1^2 L_1 D_2^2 L_2)$ to $\mathcal{O}(D_1^2 D_2^2 L_1)$ and encode architectural knowledge by grouping parameters by layers. Later in Section 3.4, this representation is also shown to preserve high representation power owing to its connection with Monarch matrices (Dao et al., 2022; 2019).

3.2.2 PARAMETER SHARING VIA KRONECKER FACTORIZATION

The above LiGO operator requires $\mathcal{O}(D_1^2 D_2^2 L_1)$ parameters for \mathbf{R}_{width} and $\mathcal{O}(L_1 L_2 D_2^2)$ for \mathbf{L}_{depth} . The width operator \mathbf{R}_{width} is thus still prohibitively expensive given that D_1 (and D_2) can easily be in the hundreds or thousands. In this section, we propose a Kronecker factorization to further reduce the number of learnable parameters for each growth operator.

Depth. For depth, we treat an entire layer as a single group and construct a new layer by combining existing layers, effectively tying parameters for all neurons in same layer. Formally, each block in \mathbf{L}_{depth} is simplified to be $\text{diag}(\ell_{i,j}) = w_{i,j}\mathbf{I}$. Then the entire matrix can be written as a Kronecker factorization, $\mathbf{L}_{depth} = \mathbf{w} \otimes \mathbf{I}$, where $\mathbf{w} \in \mathbb{R}^{L_2 \times L_1}$ is a matrix whose entry $w_{i,j}$ indicates blending weights of j -th layer of the small model to form i -th layer of the large model. This strategy reduces the number of parameters in \mathbf{L}_{depth} to $\mathcal{O}(L_1 L_2)$, and is shown on left-hand side of Figure 1.

Width. For width, we decompose each diagonal block of width expansion operator \mathbf{R}_{width} using the Kronecker factorization $\mathbf{R}_l = \mathbf{A}_l \otimes \mathbf{B}_l$, where $\mathbf{A}_l, \mathbf{B}_l \in \mathbb{R}^{D_2 \times D_1}$. Since $\text{vec}(\mathbf{CAB}) = (\mathbf{B}^\top \otimes \mathbf{C}) \text{vec}(\mathbf{A})$ (Schacke, 2004), we then have,

$$\mathbf{R}_{width} \text{vec}(\Theta) = \begin{bmatrix} \mathbf{A}_1 \otimes \mathbf{B}_1 & & \\ & \ddots & \\ & & \mathbf{A}_{L_1} \otimes \mathbf{B}_{L_1} \end{bmatrix} \text{vec}(\Theta) \quad (6)$$

$$= \text{vec} \left(\begin{bmatrix} \mathbf{B}_1 \mathbf{W}_1 \mathbf{A}_1^\top & \cdots & \mathbf{B}_{L_1} \mathbf{W}_{L_1} \mathbf{A}_{L_1}^\top \end{bmatrix}^\top \right). \quad (7)$$

Here we observe that $\mathbf{B}_l \mathbf{W}_l \mathbf{A}_l^\top$ performs in- and out-dimension expansion by \mathbf{A}_l and \mathbf{B}_l , respectively. Each new column/row is a linear combination of columns/rows of small model’s weight matrix. This factorization, which can be seen as grouping parameters by *neurons*, reduces the number of parameters to $\mathcal{O}(L_1 D_1 D_2)$. Figure 1 (right) illustrates LiGO’s width-expansion operator.

Altogether, we obtain the final parameterization of LiGO operator \mathbf{M} :

$$\mathbf{M} = \underbrace{\left(\begin{bmatrix} w_{1,1} & w_{1,2} & \cdots & w_{1,L_1} \\ \vdots & \vdots & \ddots & \vdots \\ w_{L_2,1} & w_{L_2,2} & \cdots & w_{L_2,L_1} \end{bmatrix} \otimes \mathbf{I} \right)}_{\text{Depth expansion}} \underbrace{\left(\begin{bmatrix} \mathbf{A}_1 \otimes \mathbf{B}_1 & & \\ & \ddots & \\ & & \mathbf{A}_{L_1} \otimes \mathbf{B}_{L_1} \end{bmatrix} \right)}_{\text{Width expansion}} \quad (8)$$

We can exploit the factorization to implement the LiGO operator (Eq. 8) efficiently.

Training. LiGO expands a model in three steps: (1) for each layer, inserting new rows by linearly combining existing rows through \mathbf{B}_l , (2) for each layer, inserting new columns by linearly combining existing columns through \mathbf{A}_l , and then finally (3) reconstructing each layer by linearly combining the weight matrices with \mathbf{w} along the depth. We then run a few steps (e.g., 100 iterations) of SGD to optimize \mathbf{M} , which has negligible compute cost relative to regular training. After obtaining \mathbf{M} , we initialize large model with $\mathbf{M} \text{vec}(\Theta)$, and train parameters $\Theta^{(new)}$ through SGD as usual. Algorithm 1 summarizes a forward pass of LiGO with transformer. Finally, as shown in Appendix A we note that StackBERT (Eq. 1), Interpolation (Eq. 1), and Net2Net (Eq. 2) are all special cases of LiGO (Eq. 8) with a particular setting of \mathbf{L}_{depth} and \mathbf{R}_{width} .

3.3 LIGO FOR TRANSFORMERS

While LiGO can be applied to any multi-layer neural network architecture, in this paper we focus on using LiGO to grow transformers which have been shown to be particularly amenable to scaling. Below we briefly describe how LiGO is applied to the main transformer embedding/attention layers and defer further details (e.g., growing bias vectors, layer norm parameters) to Appendix B.1.

Embedding layer. The embedding layer can be regarded as a linear layer whose inputs are one-hot vectors. We learn a matrix $\mathbf{B}^{(emb)}$ to extend its output dimension. This embedding layer is also used as the final output layer for our transformer language modeling experiments.

Attention and feedforward Layers. An attention layer consists of multi-head attention weights $(\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V)$ and a linear projection (\mathbf{W}^O) . Let \mathbf{A}_l^k and \mathbf{B}_l^k where $k \in \{Q, K, V, O\}$ be the l -th layer’s in- and out-dimension expansion matrices (Eq. 6) for the query, key, value, and projection matrices. To make sure new input and output channels are aligned across modules, we tie the LiGO operator as follows: for all $l \in [L_1]$, (1) $\mathbf{A}_l^k = (\mathbf{B}^{(emb)})^\top$ for $\forall k \in \{Q, K, V\}$, (2) $\mathbf{A}_l^O = (\mathbf{B}_l^V)^\top$, (3) $\mathbf{B}_l^O = \mathbf{B}^{(emb)}$. The last constraint is added to take into account the residual connections (Chen et al., 2021). We similarly tie parameters for the feed-forward networks, $\mathbf{A}_l^{(fc1)} = (\mathbf{B}^{(emb)})^\top$, $\mathbf{A}_l^{(fc2)} = (\mathbf{B}^{(fc1)})^\top$ and $\mathbf{B}_l^{(fc2)} = \mathbf{B}^{(emb)}$. Since transformers make heavy use of residual layers

with skip connections, we found that simply using the same $\mathbf{B}^{(emb)}$ to parameterize \mathbf{A}_l^k and \mathbf{B}_l^k for many layers/modules worked well in practice. This reduces the number of learnable parameters even further and enables fast learning of \mathbf{M} on a small amount of data (100 gradient steps).

3.4 CONNECTION TO MONARCH MATRICES

As shown in Section 3.2.1, our depth-width decomposition factorizes \mathbf{M} into a multiplication of two structured sparse matrices. We examine the expressiveness of this factorized representation by relating it to Monarch matrices (Dao et al., 2022), defined below.

Definition 1. Let the space of Monarch matrices be $\mathcal{M} \subseteq \mathbb{R}^{m n_1 \times m n_2}$. Then matrix $\mathbf{M} \in \mathcal{M}$ if $\mathbf{M} = \mathbf{P}_1 \mathbf{L} \mathbf{P}_2^\top \mathbf{R} = \mathbf{P}_1 \text{diag}(\mathbf{L}_1, \dots, \mathbf{L}_{n_1}) \mathbf{P}_2^\top \text{diag}(\mathbf{R}_1, \dots, \mathbf{R}_{n_2})$ where $\mathbf{L}_i \in \mathbb{R}^{b_1 \times b_2}$, $\mathbf{R}_i \in \mathbb{R}^{b_3 \times b_4}$ are dense rectangular matrices, and $n_1 b_2 = n_2 b_3$. \mathbf{P}_1 is the permutation $\pi(i) = (i - b_1 \lfloor \frac{i}{b_1} \rfloor - 1)n_1 + \lfloor \frac{i}{b_1} \rfloor + 1$ and \mathbf{P}_2 is the permutation $\pi(j) = (j - b_2 \lfloor \frac{j}{b_2} \rfloor - 1)n_1 + \lfloor \frac{j}{b_2} \rfloor + 1$.

It is clear that the block-diagonal matrix \mathbf{R} has the identical form to our width growing operator \mathbf{R}_{width} . By applying the permutation matrices \mathbf{P}_1 and \mathbf{P}_2 to \mathbf{L} , \mathbf{L} is transformed into exactly the same form with our depth-growth operator \mathbf{L}_{depth} in Eq. 5. This implies that our depth-width decomposition coincides with Monarch sparsification of dense matrices, which generalize butterfly matrices (Dao et al., 2019) and enjoy rich expressivity properties (Dao et al., 2020; 2022).

4 EXPERIMENTS

We conduct experiments to answer three key research questions. Q1: To what extent can LiGO improve the training efficiency (FLOPs and wall time) of transformers compared to training from scratch and other growth operators? Q2: Can LiGO be universally effective across transformers from different domains (e.g., language and vision) and sizes? Q3: Can models trained using LiGO achieve similar performance compared to the baselines when transferred to downstream tasks?

4.1 EXPERIMENTAL SETUP

Datasets. We follow Tan & Bansal (2020) and use the English Wikipedia corpus⁵ for training BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019). We use the public C4 (Raffel et al., 2020) dataset for training GPT2 (Radford et al., 2019). We use ImageNet (Deng et al., 2009) for training vision transformers. We use GLUE (Wang et al., 2018), SQuADv1.1 (Rajpurkar et al., 2016), and SQuADv2.0 (Rajpurkar et al., 2018) for evaluating pretrained BERT models. We test downstream performance of vision transformers (DeiT (Touvron et al., 2021a)) by performing transfer learning on 5 downstream image classification tasks, including CIFAR10 (Krizhevsky et al., 2009), CIFAR100 (Krizhevsky et al., 2009), Flowers102 (Nilsback & Zisserman, 2008), Stanford-Cars (Krause et al., 2013), and ChestXRy8 (Wang et al., 2017).

Models. We experiment with growing the following language and vision transformers: (1) BERT-Small→BERT-Base, BERT-Base→BERT-Large, BERT-Small→BERT-Large; (2) RoBERTa-Small→RoBERTa-Base for RoBERTa; (3) GPT2-Base→GPT2-Medium, (4) DeiT-S→DeiT-B, and (5) CaiT-XS→CaiT-S. BERT-Small has 6 layers with 512 hidden dimensions, while other named models are their usual sizes. See Appendix B.2 for full details.

Baselines. We compare our approach with the following baselines: (1) training from scratch baseline where we train the larger transformer without using any smaller pretrained models; (2) progressive training methods designed for growing depth in transformers (StackBERT (Gong et al., 2019) and MSLT (Yang et al., 2020)); (3) bert2BERT (Chen et al., 2021) that extends Net2Net (Chen et al., 2015) for width expansion and stacking for depth expansion; (4) KI (Qin et al., 2021) which uses distillation for transferring knowledge from the smaller model to the larger model.

Implementation details. We always use 100 gradient steps to learn the LiGO for all models, which is negligible in terms of FLOPs/wall time compared to full training after initialization. We train both BERT and RoBERTa models for 400K steps with a warmup of 10K steps. We remove the next-sentence prediction task (Liu et al., 2019) and use a fixed sequence length of 128 for pretraining

⁵While the original BERT (Devlin et al., 2019) paper also uses the Toronto Book Corpus (Zhu et al., 2015), we do not include it here since it is no longer publicly available.

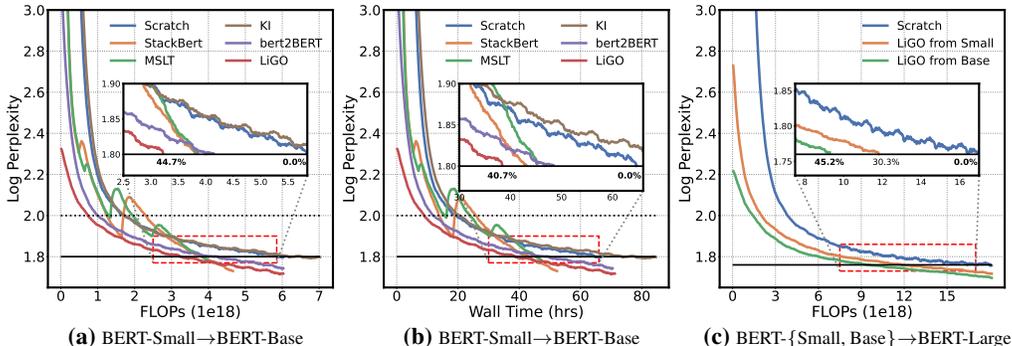


Figure 2: Results on BERT. (a-b) shows validation log perplexity vs. FLOPs and wall time respectively for training BERT-Base by reusing BERT-Small. (c) shows log perplexity vs. FLOPs in growing BERT-Small and BERT-Base to BERT-Large. The solid line indicates the final perplexity of the larger model trained from scratch, while the dotted line represents performance of the smaller model trained from scratch. LiGO offers about 45% savings in FLOPs and 40% savings in wall time over BERT-Base training from scratch. Our approach is also flexible in reusing either BERT-Small or BERT-Base for accelerating BERT-Large training.

Table 1: Downstream transfer learning performance on GLUE and SQuAD. All of the results are based on BERT-Base models trained using the different baselines. LiGO achieves similar or even better performance than the original training from scratch baseline on several downstream tasks, despite improving training efficiency.

Method	Savings (FLOPs)	Savings (Walltime)	SST-2 (Acc.)	MNLI (Acc.)	MRPC (Acc.)	CoLA (Acc.)	QNLI (Acc.)	QQP (Acc.)	STS-B (Acc.)	SQuADv1.1 (F1/EM)	SQuADv2.0 (F1/EM)	Avg. GLUE	Avg. SQuAD
Scratch	-	-	88.19	78.43	85.78	62.09	87.06	87.18	86.99	86.55 / 77.31	71.31 / 67.07	82.25	78.79 / 72.19
StackBERT	34.1%	33.3%	88.99	79.72	85.29	59.09	87.28	89.17	86.97	86.50 / 77.42	71.32 / 67.41	82.36	78.91 / 72.41
MSLT	34.9%	30.0%	88.53	78.10	82.60	64.76	83.58	88.54	85.89	86.07 / 76.73	70.68 / 67.17	81.72	78.47 / 71.95
KI	-5.7%	-13.9%	88.65	78.83	83.50	64.86	86.25	88.96	87.09	84.93 / 76.29	71.09 / 67.41	82.59	78.01 / 71.85
bert2BERT	29.0%	25.1%	88.30	80.05	85.54	61.73	88.16	86.18	87.00	86.24 / 77.09	71.52 / 66.85	82.42	78.88 / 71.97
LiGO	44.7%	40.7%	88.42	79.29	84.31	62.09	88.07	88.81	87.00	86.28 / 77.45	71.24 / 67.17	82.57	78.76 / 72.31

both models. For BERT, we use a batch size of 256 and a learning rate of $2e^{-4}$, while we use a batch size of 1024 and a learning rate of $8e^{-4}$ for training RoBERTa models.

Following Shen et al. (2022), we train GPT2 models with a batch size of 384 and sequence length of 1024. For vision transformers, we build our models based on DeiT (Touvron et al., 2021a) and CaiT (Touvron et al., 2021b), and apply their default hyper-parameters for training on ImageNet dataset. We train all our vision transformers for 300 epochs with a batch size of 1024. For transfer learning with BERT/RoBERTa, we follow Tan & Bansal (2020) and train for 3 epochs with a learning rate of $1e^{-4}$ and a batch-size of 32 for all tasks in GLUE. On SQuAD v1.1 and SQuAD 2.0, we fine-tune for 2 epochs with a learning rate of $5e^{-5}$ and a batch size of 12. We run both GLUE and SQuAD evaluations three times with different random seeds and report the mean numbers. For transfer learning experiments on DeiT, we finetune the pretrained models with 1000 epochs, batch size 768, learning rate 0.01, and use the same data augmentation in training on ImageNet. We use the same pretraining data and experimental settings for all the baselines (including our approach) for a fair comparison. Note that we include the additional compute required for training LiGO in all our tables and figures. However, since our LiGO is only trained for 100 steps, the influence on visualization and quantitative saving percentages is negligible.

4.2 RESULTS AND ANALYSIS

BERT. Figure 2 shows the comparison between the different baselines for training BERT models. As seen from Figure 2(a), LiGO saves 44.7% computational cost (FLOPs) of training BERT-Base (12 layers, 768 dimensions) from scratch by reusing BERT-Small (6 layers, 512 dimensions). LiGO offers 40.7% savings in wall time compared to training from scratch (Figure 2(b)). Among the compared methods, StackBERT is the most competitive in terms of both FLOPs and wall time, although LiGO obtains +10.6% and +7.2% improvements in FLOPs and wall time on top of StackBERT. Similarly, LiGO significantly outperforms the recent bert2BERT method which saves about 30% computational costs. We observe that KI does not provide any real savings in training as it requires additional computation for knowledge distillation. Figure 2(c) shows that our LiGO approach is flexible in growing either BERT-Small or BERT-Base for accelerating BERT-Large training. As expected, reusing BERT-Base instead of BERT-Small leads more savings in FLOPs (45.2% vs 30.3%) as BERT-Base contains more implicit knowledge in its parameters. Table 1 shows the per-task per-

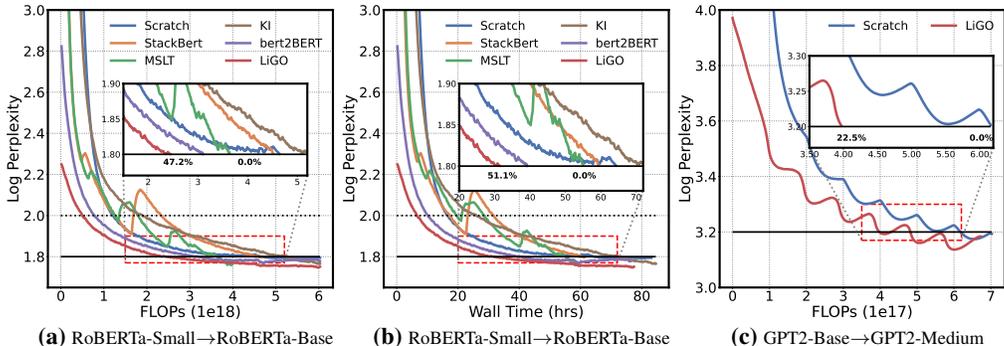


Figure 3: Results on RoBERTa and GPT2. LiGO reduces FLOPs by 47.2% and 22.5% for RoBERTa-Base and GPT2-Medium, demonstrating its effectiveness across different training strategies and architectures.

formance of different BERT-Base models on both GLUE and SQuAD benchmarks, where we find that BERT trained with LiGO achieves very similar performance compared to the baselines on both benchmarks. Finally, in Table 5 of the Appendix C.3, we show that growing BERT-Small to BERT-Base with 100 steps of LiGO and then finetuning on GLUE tasks *without* additional pretraining outperforms just directly finetuning BERT-Small.

RoBERTa and GPT2. Figure 3(a-b) shows the results on RoBERTa, whose training recipe uses larger batch size and learning rate than BERT. LiGO similarly accelerates RoBERTa training, which indicates that our method is robust to optimization hyperparameters. On GPT2, LiGO saves 22.5% computation cost of training GPT2-Medium (345M parameters) by reusing GPT2-Base (117M parameters) (Figure 3(c)). These consistent improvements show that LiGO is effective for accelerating transformer training across different model architectures and sizes.

Vision Transformers. Figure 4 shows that by growing from DeiT-S, LiGO can save 55.4% FLOPs and 52% GPU wall time to reach the same performance of 81% on ImageNet. Interestingly, the model initialized by our data-driven growth operator (w/ only 100 gradient steps of tuning) can already achieve 72% accuracy at the beginning of training and leads to the final accuracy of 81.7% at the end of the training. Compared to the next best method, bert2BERT, LiGO

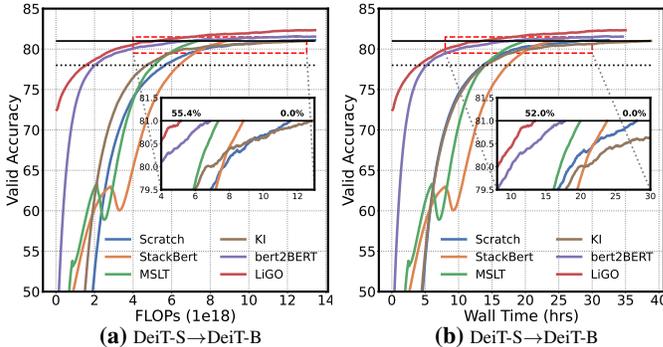


Figure 4: Results on DeiT. (a) Accuracy vs. flops and (b) accuracy vs. wall time for training DeiT-B. LiGO saves flops and wall time by more than 50% over training from scratch on ImageNet.

obtains more than 15% savings, which once again demonstrates the effectiveness of our approach in growing vision transformers as well. Table 2 shows that finetuning results on downstream tasks perform on-par with the model trained from scratch, showing that LiGO does not harm the model’s generalization capabilities when transferred to downstream datasets. We also find similar savings in CaiT-XS→CaiT-S where LiGO saves FLOPs by 52.6% and wall time by 46.1% over training CaiT-S from scratch on ImageNet (see Appendix C.2 for more details).

Combining with other training strategies. We also find that LiGO can be effectively combined with orthogonal strategies such as layer dropping (Zhang & He, 2020), token dropping (Hou et al., 2022), and staged training (Chen et al., 2021). More details are included in Appendix B.3. Figure 5 shows that LiGO can be combined with other training techniques to improve the computational savings by 4.7%, 7.4%,

Table 2: Transfer learning performance of DeiT-B. DeiT-B model trained using LiGO performs similarly to the original train-from-scratch baseline on all downstream tasks.

Method	FLOPs	Walltime	ImageNet	CIFAR10	CIFAR100	Flowers	Cars	ChestXRy8
Scratch	-	-	81.10	99.09	90.76	97.79	92.06	55.81
StackBERT	23.8%	15.1%	81.21	99.11	90.80	97.56	92.09	55.77
MSLT	36.7%	28.9%	81.27	99.07	90.21	97.71	92.11	55.79
KI	-11.2%	-36.8%	81.01	98.94	90.32	97.81	92.08	55.80
bert2BERT	40.8%	37.0%	81.59	99.14	90.69	97.67	92.15	55.82
LiGO	55.4%	52.0%	81.71	99.12	90.74	97.77	92.09	55.82

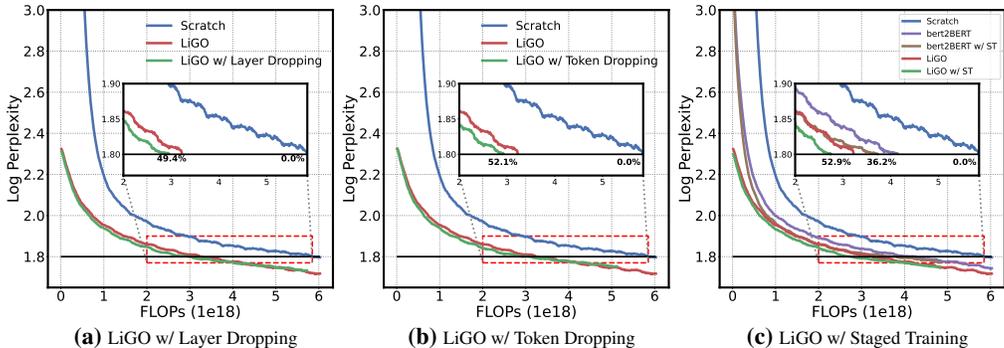


Figure 5: LiGO with other efficient training strategies. Our approach can be combined with (a) layer dropping, (b) token dropping, and (c) staged training (ST), for further accelerate BERT training.

and 8.2% with layer dropping, token dropping and staged training, respectively. Following (Chen et al., 2021), we also apply staged training strategy to bert2BERT and observe that LiGO still outperforms bert2BERT with staged training by 16.7% (see Figure 5(c)).

4.3 ABLATION STUDIES

Depth-only expansion. We examine the effectiveness of our proposed depth expansion operator (L_{depth}) by only growing the depth of BERT from 6 layers to 12 layers, i.e., $BERT(6, 768) \rightarrow BERT(12, 768)$. We compare with stacking (StackBERT, Gong et al., 2019), Interpolation (InterBERT, Chang et al., 2017; Dong et al., 2020) (see Eq. 1), and MSLT (Yang et al., 2020). For LiGO, we only apply its L_{depth} component to the pre-trained model weights. Results in Figure 6(a) show that a data-driven approach works well even when just growing across the depth dimension.

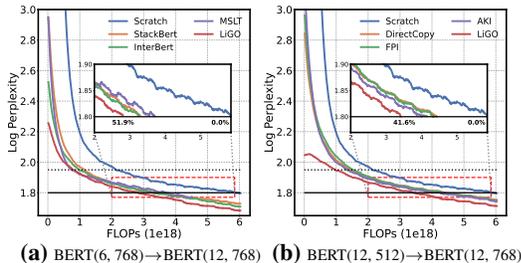


Figure 6: Results on Depth-only and Width-only growth. LiGO saves 51.7% FLOPs when expanding depth-only, and 41.6% FLOPs when expanding width-only.

Width-only expansion. We also verify the effectiveness of R_{width} by only extending BERT width from 512 to 768, i.e., $BERT(12, 512) \rightarrow BERT(12, 768)$. We compare LiGO based initialization with direct copy (Wei et al., 2016), function preserving initialization (FPI, Chen et al., 2015), and advanced knowledge initialization (AKI, Chen et al., 2021). LiGO’s width expansion component outperforms all other methods, as shown in Figure 6(b).

Table 3: Effect of number of gradient steps. “+FLOPs” stands for additional flops (in 10^{15}).

# of Steps	+FLOPs	Savings
100	3.61	44.7%
500	18.06	44.5%
1000	36.13	44.2%
10000	361.30	38.9%

Number of growing steps. Our main experiments just use 100 gradient steps to grow. We tune our LiGO on the pretraining set for 100, 500, 1000, and 10000 steps and compute the additional FLOPs for BERT-Small \rightarrow BERT-Base training. Table 3 shows that training LiGO within 1000 steps results in the identical model convergence (reaching 1.8 PPL at 215K steps). This suggests tuning model weights under the linear constraints of LiGO can achieve faster convergence. Training LiGO for more than 10000 steps can provide a model with slightly faster convergence (214K steps), but results in less saving overall.

5 CONCLUSION

This paper describes an approach for accelerating transformer training by learning to grow pre-trained transformers, where the larger transformer’s parameters are initialized as a linear mapping from the smaller pretrained model’s parameters. The linear map is factorized to be a composition of sparse width- and depth-expansion operators with a Kronecker factorization that groups parameters into layers and neurons. We demonstrate the effectiveness of our proposed approach on both language and vision transformers of different sizes, outperforming several competing methods. While our compute resources prevented us from applying LiGO to even larger transformers, it would be interesting to see if this can be applied on top of even larger models.

ACKNOWLEDGMENTS

PW sincerely thanks Zhen Wang for the insightful discussion and for providing reference repositories for language model pre-training. PW also appreciates Hao Tan’s assistance for reproducing fine-tuning results on GLUE datasets. YK and LTH were partially supported an MIT-IBM Watson AI grant and an Amazon award. We also acknowledge support from the IBM Research AI Hardware Center, and the Center for Computational Innovation at Rensselaer Polytechnic Institute for the computational resources on the AiMOS Supercomputer. The research of ZW is in part supported by the US Army Research Office Young Investigator Award (W911NF2010240).

REFERENCES

- Andy Brock, Soham De, Samuel L Smith, and Karen Simonyan. High-performance large-scale image recognition without normalization. In *International Conference on Machine Learning*, pp. 1059–1071, 2021.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Proceedings of NeurIPS*, 2020.
- Han Cai, Tianyao Chen, Weinan Zhang, Yong Yu, and Jun Wang. Efficient architecture search by network transformation. In *Proceedings of AAAI*, 2018.
- Bo Chang, Lili Meng, Eldad Haber, Frederick Tung, and David Begert. Multi-level residual networks from dynamical systems view. *arXiv preprint arXiv:1710.10348*, 2017.
- Cheng Chen, Yichun Yin, Lifeng Shang, Xin Jiang, Yujia Qin, Fengyu Wang, Zhi Wang, Xiao Chen, Zhiyuan Liu, and Qun Liu. bert2bert: Towards reusable pretrained language models. *arXiv preprint arXiv:2110.07143*, 2021.
- Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.
- Tianqi Chen, Ian Goodfellow, and Jonathon Shlens. Net2net: Accelerating learning via knowledge transfer. *arXiv preprint arXiv:1511.05641*, 2015.
- Justin Chiu, Yuntian Deng, and Alexander Rush. Low-rank constraints for fast inference in structured models. *Advances in Neural Information Processing Systems*, 34:2887–2898, 2021.
- Xiaoliang Dai, Hongxu Yin, and Niraj K Jha. Nest: A neural network synthesis tool based on a grow-and-prune paradigm. *IEEE Transactions on Computers*, 68(10):1487–1497, 2019.
- Tri Dao, Albert Gu, Matthew Eichhorn, Atri Rudra, and Christopher Ré. Learning fast algorithms for linear transforms using butterfly factorizations. In *International conference on machine learning*, pp. 1517–1527, 2019.
- Tri Dao, Nimit S Sohoni, Albert Gu, Matthew Eichhorn, Amit Blonder, Megan Leszczynski, Atri Rudra, and Christopher Ré. Kaleidoscope: An efficient, learnable representation for all structured linear maps. *arXiv preprint arXiv:2012.14966*, 2020.
- Tri Dao, Beidi Chen, Nimit S Sohoni, Arjun Desai, Michael Poli, Jessica Grogan, Alexander Liu, Aniruddh Rao, Atri Rudra, and Christopher Ré. Monarch: Expressive structured matrices for efficient and accurate training. In *International Conference on Machine Learning*, pp. 4690–4721, 2022.
- Yann N Dauphin and Samuel Schoenholz. Metainit: Initializing learning by learning to initialize. *Advances in Neural Information Processing Systems*, 32, 2019.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255, 2009.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL*, 2019.

- Chengyu Dong, Liyuan Liu, Zichao Li, and Jingbo Shang. Towards adaptive residual network training: A neural-ode perspective. In *International conference on machine learning*, pp. 2616–2626. PMLR, 2020.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *Proceedings of ICLR*, 2021.
- Utku Evci, Max Vladymyrov, Thomas Unterthiner, Bart van Merriënboer, and Fabian Pedregosa. Gradmax: Growing neural networks using gradient information. *arXiv preprint arXiv:2201.05125*, 2022.
- Scott Fahlman. The recurrent cascade-correlation architecture. In *Advances in Neural Information Processing Systems*, 1990.
- Scott Fahlman and Christian Lebiere. The cascade-correlation learning architecture. In *Advances in Neural Information Processing Systems*, 1989.
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 249–256, 2010.
- Linyuan Gong, Di He, Zhuohan Li, Tao Qin, Liwei Wang, and Tiejun Liu. Efficient training of bert by progressively stacking. In *International conference on machine learning*, pp. 2337–2346, 2019.
- Xiaotao Gu, Liyuan Liu, Hongkun Yu, Jing Li, Chen Chen, and Jiawei Han. On the transformer growth for progressive bert training. *arXiv preprint arXiv:2010.12562*, 2020.
- Steven Gutstein, Olac Fuentes, , and Eric Freudenthal. Knowledge transfer in deep convolutional neural nets. In *Proceedings of International Journal on Artificial Intelligence Tools*, 2008.
- Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- Le Hou, Richard Yuanzhe Pang, Tianyi Zhou, Yuexin Wu, Xinying Song, Xiaodan Song, and Denny Zhou. Token dropping for efficient bert pretraining. *arXiv preprint arXiv:2203.13240*, 2022.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, 2019.
- Xiao Shi Huang, Felipe Perez, Jimmy Ba, and Maksims Volkovs. Improving transformer optimization through better initialization. In *International Conference on Machine Learning*, pp. 4475–4483, 2020.
- Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Chen, HyukJoong Lee, Jiquan Ngiam, Quoc V Le, Yonghui Wu, et al. Gpipe: Efficient training of giant neural networks using pipeline parallelism. *Advances in neural information processing systems*, 32, 2019.
- Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. *arXiv preprint arXiv:2203.12119*, 2022.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Yannic Kilcher, Gary Bécigneul, and Thomas Hofmann. Escaping flat areas via function-preserving structural network modifications. 2018.
- Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pp. 554–561, 2013.

- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- Changlin Li, Bohan Zhuang, Guangrun Wang, Xiaodan Liang, Xiaojun Chang, and Yi Yang. Automated progressive learning for efficient training of vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12486–12496, 2022.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Dmytro Mishkin and Jiri Matas. All you need is a good init. *arXiv preprint arXiv:1511.06422*, 2015.
- Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pp. 722–729. IEEE, 2008.
- Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapter-fusion: Non-destructive task composition for transfer learning. *arXiv preprint arXiv:2005.00247*, 2020.
- Yujia Qin, Yankai Lin, Jing Yi, Jiajie Zhang, Xu Han, Zhengyan Zhang, Yusheng Su, Zhiyuan Liu, Peng Li, Maosong Sun, et al. Knowledge inheritance for pre-trained language models. *arXiv preprint arXiv:2105.13880*, 2021.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67, 2020.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable questions for squad. *arXiv preprint arXiv:1806.03822*, 2018.
- Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15), 2021. ISSN 0027-8424. doi: 10.1073/pnas.2016239118.
- Jonathan S Rosenfeld, Amir Rosenfeld, Yonatan Belinkov, and Nir Shavit. A constructive prediction of the generalization error across scales. *arXiv preprint arXiv:1909.12673*, 2019.
- Kathrin Schacke. On the kronecker product. *Master’s thesis, University of Waterloo*, 2004.
- Sheng Shen, Pete Walsh, Kurt Keutzer, Jesse Dodge, Matthew Peters, and Iz Beltagy. Staged training for transformer language models. *arXiv preprint arXiv:2203.06211*, 2022.
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*, 2019.
- Vikas Sindhvani, Tara Sainath, and Sanjiv Kumar. Structured transforms for small-footprint deep learning. *Advances in Neural Information Processing Systems*, 28, 2015.

- Hao Tan and Mohit Bansal. Vokenization: Improving language understanding with contextualized, visual-grounded supervision. *arXiv preprint arXiv:2010.06775*, 2020.
- Shanshan Tang, Bo Li, and Haijun Yu. Chebnet: Efficient and stable constructions of deep neural networks with rectified power units using chebyshev approximations. *arXiv preprint arXiv:1911.05467*, 2019.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pp. 10347–10357, 2021a.
- Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021b.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is All You Need. In *Proceedings of NeurIPS*, 2017.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 353–355, 2018.
- Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2097–2106, 2017.
- Tao Wei, Changhu Wang, Yong Rui, and Chang Wen Chen. Network morphism. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings of The 33rd International Conference on Machine Learning*, pp. 564–572, 2016.
- Lemeng Wu, Dilin Wang, and Qiang Liu. Splitting steepest descent for growing neural architectures. *Advances in neural information processing systems*, 32, 2019.
- Lemeng Wu, Dilin Wang, Peter Stone, and Qiang Liu. Firefly neural architecture descent: a general approach for growing neural networks. *Advances in neural information processing systems*, 2021.
- Cheng Yang, Shengnan Wang, Chao Yang, Yuechuan Li, Ru He, and Jingqiao Zhang. Progressively stacking 2.0: A multi-stage layerwise training method for bert training speedup. *arXiv preprint arXiv:2011.13635*, 2020.
- Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019.
- Hongyi Zhang, Yann N Dauphin, and Tengyu Ma. Fixup initialization: Residual learning without normalization. *arXiv preprint arXiv:1901.09321*, 2019.
- Minjia Zhang and Yuxiong He. Accelerating training of transformer-based language models with progressive layer dropping. *Advances in Neural Information Processing Systems*, 33:14011–14023, 2020.
- Xu Zhang, Felix X Yu, Ruiqi Guo, Sanjiv Kumar, Shengjin Wang, and Shi-Fu Chang. Fast orthogonal projection based on kronecker product. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2929–2937, 2015.
- Chen Zhu, Renkun Ni, Zheng Xu, Kezhi Kong, W Ronny Huang, and Tom Goldstein. Gradinit: Learning to initialize neural networks for stable and efficient training. *Advances in Neural Information Processing Systems*, 34:16410–16422, 2021.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pp. 19–27, 2015.

A UNIVERSALITY OF LIGO OPERATOR

Proposition 1. *StackBERT (Eq. 1), Interpolation (Eq. 1), and Net2Net (Eq. 2) are all the special cases of the LiGO operator (Eq. 8).*

Proof. We prove Proposition 1 by constructing parameters in L_{depth} and R_{width} .

Stacking. Stacking-based methods (Gong et al., 2019; Yang et al., 2020) duplicate the entire lower blocks on top of the small model to the form new layers (Eq. 1). Formally, we show this operation can be done by the following operator:

$$M = \underbrace{\begin{bmatrix} I & & & \\ & I & & \\ & & \ddots & \\ I & & & I \end{bmatrix}}_{L_{depth}} \underbrace{\begin{bmatrix} I & & \\ & \ddots & \\ & & I \end{bmatrix}}_{R_{width}} \quad (9)$$

Interpolation. Interpolation based methods (Chang et al., 2017; Dong et al., 2020) interleave each layer for twice. We can construct the following matrix to achieve layer interpolation (Eq. 1).

$$M = \underbrace{\begin{bmatrix} I & & & \\ I & & & \\ & I & & \\ & I & & \\ & & \ddots & \\ & & & I \end{bmatrix}}_{L_{depth}} \underbrace{\begin{bmatrix} I & & \\ & \ddots & \\ & & I \end{bmatrix}}_{R_{width}} \quad (10)$$

We remark that any rearrangement of layers to construct new layers (mathematically a permutation of existing layers with replacement) can be constructed in a similar way.

Net2Net. Since we show in Eq. 6, the Kronecker factorization on R_l amounts to decomposing the general growth operator into in-dimension and out-dimension expansion. We can construct Net2Net (Chen et al., 2015) based growth by simply letting:

$$L_{depth} = I \in \mathbb{R}^{L_1 D_2 \times L_2 D_2}, \quad R_{width} = \begin{bmatrix} A_1 \otimes B_1 & & \\ & \ddots & \\ & & A_{L_1} \otimes B_{L_1} \end{bmatrix} \quad (11)$$

$$A_l = \begin{bmatrix} I \\ \tilde{S}_{l-1} \end{bmatrix}, \quad B_l = \begin{bmatrix} I \\ S_l \end{bmatrix} \quad (12)$$

where $S_l \in \{0, 1\}^{(D_2 - D_1) \times D_1}$ is a selection matrix to enlarge the out dimension, and $\tilde{S}_{l-1} = S_{l-1} \text{diag}(\mathbf{1}^\top S_{l-1})^{-1}$ copies the selection from S_{l-1} with normalization to guarantee functionality preserving in expansion. \square

B IMPLEMENTATION DETAILS

B.1 GROWING TRANSFORMERS WITH LIGO

The transformer architecture consists of an embedding layer, multi-block attention layer, and an output layer. The core ingredient attention block consists of a Multi-Head Attention (MHA) module followed by a FeedForward Network (FFN), with a skip connection across the both blocks. Applying LiGO requires the following considerations:

Embedding layer. For both language and vision transformers, the embedding layer can be regarded as a linear layer, whose inputs are one-hot embeddings in language models. We draw a learnable matrix $\mathbf{B}^{(emb)}$ to extend its output dimension.

Multi-head attention blocks. An attention layer in transformer consists of multi-head attention weights ($\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V$) and a linear projection (\mathbf{W}^O). Let \mathbf{A}_l^k and \mathbf{B}_l^k with $k \in \{Q, K, V, O\}$ be the in- and out-dimension expansion matrices (Eq. 6) for query, key, value, and projection in the l -th layer, respectively. Applying \mathbf{B}_l^k to \mathbf{W}^k ($k \in Q, K, V$) constructs new heads by a weighted summation of rows of all existing heads. To make sure the new input and output channels are aligned across modules, we tie our LiGO operator with the following scheme: (1) $\mathbf{A}_l^k = (\mathbf{B}^{(emb)})^\top$ for $\forall k \in \{Q, K, V\}$, (2) $\mathbf{A}_l^O = (\mathbf{B}_l^{(V)})^\top$, (3) $\mathbf{B}_l^O = \mathbf{B}^{(emb)}$ for $\forall l \in [L_1]$. Both the bias and layer normalization inherit the associated linear transformations’ out-dimension expansion matrices to grow the width. For depth expansion, each module independently combines the same module from other layers (Eq. 8) with learnable coefficients w .

Feed-forward networks. Each attention block is followed by a two-layer FFN. Let \mathbf{A}_l^k and \mathbf{B}_l^k with $k \in \{fc1, fc2\}$ be the in- and out-dimension expansion matrices (Eq. 6) for the first and second FFN layer in the l -th layer, respectively. We tie the parameters for feed-forward networks: $\mathbf{A}_l^{(fc1)} = \mathbf{B}^{(emb)\top}$, $\mathbf{A}_l^{(fc2)} = \mathbf{B}_l^{(fc1)\top}$ and $\mathbf{B}_l^{(fc2)} = \mathbf{B}^{(emb)}$.

Output layer. For output head, we have $\mathbf{A}^{(out)} = \mathbf{B}^{(emb)\top}$, since the output dimension of attention layers are always aligned with $\mathbf{B}^{(emb)}$ by our construction. The output layer does not need out-dimension expansion. Algorithm 1 summarizes LiGO for growing transformers.

B.2 MODEL CONFIGURATIONS

We summarize the settings of different transformer models used for our experiments in Table 4. For BERT and RoBERTa, we re-use the code base provided by Tan & Bansal (2020). For GPT2, we follow the model configuration of OpenAI and use the pre-training code provided by Shen et al. (2022). For DeiT, we use their official codebase (Touvron et al., 2021a).

Table 4: Configuration of different transformers.

	BERT-Small	BERT-Base	RoBERTa-Small	RoBERTa-Base	GPT2-Base	GPT2-Medium		DeiT-S	DeiT-B
# layers	6	12	6	12	12	24	# layers	12	12
# hidden	512	768	512	768	768	1024	# hidden	384	768
# heads	8	12	8	12	12	16	# heads	6	12
# vocab	30522	30522	50265	50265	50257	50257	input res.	224	224
seq. length	128	128	128	128	1024	1024	patch size	16	16

B.3 ORTHOGONAL EFFICIENT TRAINING STRATEGIES

For layer dropping, we follow the same progressive dropping rate schedule with Zhang & He (2020), and set the maximum dropping rate to 0.1 to recover the performance. For token dropping, we randomly set 15% tokens aside in the middle layers. In the first 50k steps of staged training, only a sub-network is activated and trained, and afterwards, we perform full-model training for 350k steps.

C ADDITIONAL EXPERIMENTS

C.1 REUSING SMALLER MODELS TRAINED FOR ONLY FEW STEPS

LiGO focuses on utilizing the knowledge of smaller models that have already been pretrained and available. In this section, we investigate how LiGO can leverage smaller existing models that are only trained for few steps to accelerate training of a larger model. We perform an experiment on BERT-Base by reusing a BERT-Small trained for only 50k steps instead full training for 220k steps as used in our experiments. Figure 7 shows that LiGO can still save 35.2% savings in FLOPs and 30.2% savings in wall time over the BERT-Base training from scratch.

Algorithm 1 A forward pass of LiGO with transformer.

-
- 1: **Input:** A small transformer with hidden D_1 and number of layer L_1 . Denote the embedding layer as $\mathbf{W}^{(emb)} \in \mathbb{R}^{D_1 \times E}$, attention layers as $\mathbf{W}_l^Q, \mathbf{W}_l^K, \mathbf{W}_l^V, \mathbf{W}_l^O \in \mathbb{R}^{D_1 \times D_1}$, FFN layers as $\mathbf{W}_l^{(fc1)} \in \mathbb{R}^{4D_1 \times D_1}$, $\mathbf{W}_l^{(fc2)} \in \mathbb{R}^{D_1 \times 4D_1}$, LayerNorm layers as $\mathbf{W}_l^{(ln1)} \in \mathbb{R}^{D_1 \times 4D_1}$, $\mathbf{W}_l^{(ln2)} \in \mathbb{R}^{D_1 \times 4D_1}$, $\forall l \in [L_1]$, the output head $\mathbf{W}^{(out)} \in \mathbb{R}^{C \times D_1}$
 - 2: **Output:** A large transformer with hidden D_2 and number of layer L_2 . Denote the weight matrices as Ω with the corresponding superscripts as the small model.
 - 3: $\Omega^{(emb)} \leftarrow \mathbf{B}^{(emb)} \mathbf{W}^{(emb)}$
 - 4: **for** $l = 1, \dots, L_1$ **do** ▷ Width Expansion
 - 5: $\Omega_l^Q \leftarrow \mathbf{B}^{(Q)} \mathbf{W}_l^Q \mathbf{B}^{(emb)\top}$
 - 6: $\Omega_l^K \leftarrow \mathbf{B}^{(K)} \mathbf{W}_l^K \mathbf{B}^{(emb)\top}$
 - 7: $\Omega_l^V \leftarrow \mathbf{B}^{(V)} \mathbf{W}_l^V \mathbf{B}^{(emb)\top}$
 - 8: $\Omega_l^O \leftarrow \mathbf{B}^{(emb)} \mathbf{W}_l^O \mathbf{B}^{(V)\top}$
 - 9: $\Omega_l^{(ln1)} \leftarrow \mathbf{B}^{(emb)} \mathbf{W}_l^{(ln1)}$
 - 10: $\Omega_l^{(fc1)} \leftarrow \mathbf{B}^{(fc1)} \mathbf{W}_l^V \mathbf{B}^{(emb)\top}$
 - 11: $\Omega_l^{(fc2)} \leftarrow \mathbf{B}^{(emb)} \mathbf{W}_l^V \mathbf{B}^{(fc1)\top}$
 - 12: $\Omega_l^{(ln2)} \leftarrow \mathbf{B}^{(emb)} \mathbf{W}_l^{(ln2)}$
 - 13: **end for**
 - 14: **for** $l = 1, \dots, L_2$ **do** ▷ Depth Expansion
 - 15: $\Omega_l^Q \leftarrow \sum_{j=1}^{L_1} w_{l,j}^Q \Omega_j^Q$
 - 16: $\Omega_l^K \leftarrow \sum_{j=1}^{L_1} w_{l,j}^K \Omega_j^K$
 - 17: $\Omega_l^V \leftarrow \sum_{j=1}^{L_1} w_{l,j}^V \Omega_j^V$
 - 18: $\Omega_l^O \leftarrow \sum_{j=1}^{L_1} w_{l,j}^O \Omega_j^O$
 - 19: $\Omega_l^{(ln1)} \leftarrow \sum_{j=1}^{L_1} w_{l,j}^{(ln1)} \Omega_j^{(ln1)}$
 - 20: $\Omega_l^{(fc1)} \leftarrow \sum_{j=1}^{L_1} w_{l,j}^{(fc1)} \Omega_j^{(fc1)}$
 - 21: $\Omega_l^{(fc2)} \leftarrow \sum_{j=1}^{L_1} w_{l,j}^{(fc2)} \Omega_j^{(fc2)}$
 - 22: $\Omega_l^{(ln2)} \leftarrow \sum_{j=1}^{L_1} w_{l,j}^{(ln2)} \Omega_j^{(ln2)}$
 - 23: **end for**
 - 24: $\Omega^{(out)} \leftarrow \mathbf{W}^{(out)} \mathbf{B}^{(emb)\top}$
 - 25: Train transformer with parameters Ω .
-

C.2 RESULTS ON CAiT

In addition to DeiT (Touvron et al., 2021a), we perform additional experiments with CaiT (Touvron et al., 2021b) on ImageNet and find that while reusing CaiT-XS, LiGO offers about 52.6% savings in FLOPs and 46.1% savings in wall time over the CaiT-S training from scratch (see Figure 8).

C.3 TASK-SPECIFIC FINETUNING WITH LiGO INITIALIZATION WITHOUT FURTHER PRETRAINING

We perform additional experiments by directly finetuning BERT-Base initialized by LiGO (from BERT-Small) without any further pretraining. We observe in Table 5 that the LiGO-initialized model can benefit downstream tasks compared to BERT-Small trained from scratch (1st row vs 2nd row).

C.4 GLUE PERFORMANCE USING ADAPTERFUSION

LiGO is mainly proposed for improving efficiency of the pre-training stage and hence is compatible with various finetuning schemes like full model finetuning, adapters (Houlsby et al., 2019; Pfeiffer et al., 2020) or prompt tuning (Lester et al., 2021; Jia et al., 2022) for adaptation to downstream tasks. We test BERT-Base models trained using different baselines by using adapterfusion (Pfeiffer et al., 2020) instead of full finetuning on GLUE benchmark. Table 6 shows that LiGO also achieves

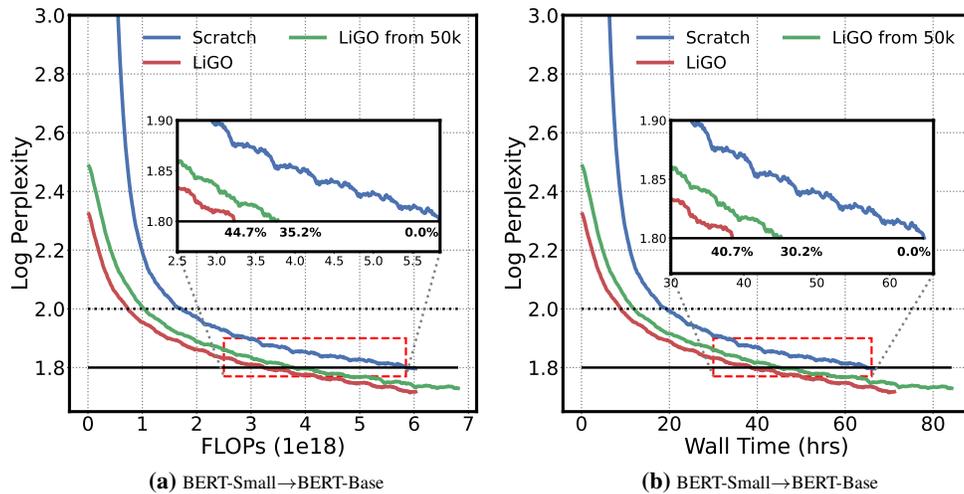


Figure 7: Results on BERT-Base by reusing BERT-Small trained for 50k steps. Instead of training BERT-Base from fully trained BERT-Small, we run LiGO on BERT-Small trained with 50k steps. LiGO offers about 35.2% savings in FLOPs and 30.2% savings in wall time over the BERT-Base training from scratch.

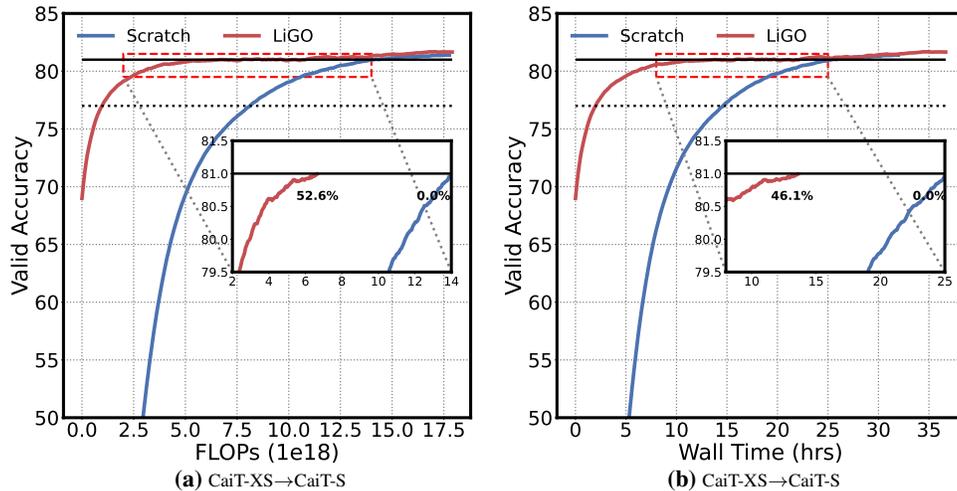


Figure 8: Results on CaiT. (a) Accuracy vs. flops and (b) accuracy vs. wall time for training CaiT-S. LiGO saves flops by 52.6% and wall time by 46.1% over training from scratch on ImageNet.

Table 5: GLUE performance of different LiGO models. All of the results are based on BERT-Base models with BERT-Small as the base model for LiGO optimization.

Method	SST-2 (Acc.)	MNLI (Acc.)	MRPC (Acc.)	CoLA (Acc.)	QNLI (Acc.)	QQP (Acc.)	STS-B (Acc.)	Average (Acc.)
BERT-Small (Scratch)	87.21	77.56	82.11	59.93	85.06	85.82	84.99	80.38
BERT-Base (LiGO Init)	88.15	77.62	82.53	60.70	85.79	86.65	85.83	81.04
BERT-Base (LiGO Init + Pretrain)	88.42	79.29	84.31	62.09	88.07	88.81	87.00	82.57
BERT-Base (Scratch)	88.19	78.43	85.78	62.09	87.06	87.18	86.99	82.25

Table 6: Downstream performance using AdapterFusion (Pfeiffer et al., 2020) on GLUE Benchmark. All of the results are based on BERT-Base models trained using different baselines.

Method	Savings (FLOPs)	Savings (Walltime)	SST-2 (Acc.)	MNLI (Acc.)	MRPC (Acc.)	CoLA (Acc.)	QNLI (Acc.)	QQP (Acc.)	STS-B (Acc.)	Average (Acc.)
Scratch	–	–	88.41	78.60	86.02	62.39	87.62	88.02	86.52	82.51
StackBERT	34.1%	33.3%	88.78	79.80	85.43	59.56	87.71	89.19	86.27	82.39
MSLT	34.9%	30.0%	88.41	78.35	83.15	63.97	86.19	88.20	86.42	82.10
KI	-5.7%	-13.9%	88.94	78.84	84.00	64.61	86.75	88.19	87.93	82.75
bert2BERT	29.0%	25.1%	88.47	80.53	85.50	62.33	88.57	86.72	87.10	82.75
LiGO	44.7%	40.5%	88.45	80.01	84.67	63.05	88.06	88.92	87.00	82.88

on-par performance with model trained from scratch under adapter-based tuning with 44.7% savings in FLOPs and 40.5% savings in wall time. This shows that LiGO does not harm the model generalization capability when adapters are used as a parameter-efficient finetuning strategy for transferring a trained model to downstream datasets.

C.5 INITIAL RESULTS ON BILLION+ PARAMETER MODELS

Our extensive experiments on BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), GPT2 (Radford et al., 2019), DeiT (Touvron et al., 2021a) and CaiT (Touvron et al., 2021b) show that LiGO can consistently improve transformer training efficiency over the traditional way of training from scratch across domains and model sizes. One interesting future direction of our work is scaling LiGO to very large models with parameters more than 100B, such as GPT3 (Brown et al., 2020). While we currently do not possess the compute resources for this extreme large-scale study, we perform a preliminary experiment on GPT2-1.5B (Radford et al., 2019) by using GPT2-Medium as the initialization. We train for 15k steps on C4 dataset (Raffel et al., 2020) and find that our proposed LiGO saves about 39% computation cost (FLOPs) of training GPT2-1.5B from scratch to reach the same log perplexity (which is 3.3). We believe that it is imperative to study the extent to which the benefits of LiGO remain at the scale on which the modern large language models are trained. We hope to cover this in our future work.