

# TransCODE: Co-design of Transformers and Accelerators for Efficient Training and Inference

Shikhar Tuli <sup>✉</sup>, *Student Member, IEEE*, and Niraj K. Jha, *Fellow, IEEE*

**Abstract**—Automated co-design of machine learning models and evaluation hardware is critical for efficiently deploying such models at scale. Despite the state-of-the-art performance of transformer models, they are not yet ready for execution on resource-constrained hardware platforms. High memory requirements and low parallelizability of the transformer architecture exacerbate this problem. Recently-proposed accelerators attempt to optimize the throughput and energy consumption of transformer models. However, such works are either limited to a one-sided search of the model architecture or a restricted set of off-the-shelf devices. Furthermore, previous works only accelerate model inference and not training, which incurs substantially higher memory and compute resources, making the problem even more challenging. To address these limitations, this work proposes a dynamic training framework, called DynaProp, that speeds up the training process and reduces memory consumption. DynaProp is a low-overhead pruning method that prunes activations and gradients at runtime. To effectively execute this method on hardware for a diverse set of transformer architectures, we propose ELECTOR, a framework that simulates transformer inference and training on a design space of accelerators. We use this simulator in conjunction with the proposed co-design technique, called TransCODE, to obtain the best-performing models with high accuracy on the given task and minimize latency, energy consumption, and chip area. The obtained transformer-accelerator pair achieves 0.3% higher accuracy than the state-of-the-art pair while incurring 5.2× lower latency and 3.0× lower energy consumption.

**Index Terms**—Application-specific integrated circuits; hardware-software co-design; machine learning; neural network accelerators; transformers.

## I. INTRODUCTION

ARTIFICIAL intelligence (AI) is undergoing a paradigm shift with the rise of large language models, e.g., BERT [1], GPT-3 [2], DALL-E [3]). These models, backed by the transformer architecture [4], target many applications, including language [1], vision [5], robotic manipulation [6], reasoning [7], human interaction [8], and forecasting [9]. Training on comprehensive datasets (generally using self-supervision at scale) and finetuning on downstream tasks have enabled their widespread application. However, training and inference with such large models either involve high power consumption on graphical processing units (GPUs) or high energy and latency on off-the-shelf edge-AI devices. For instance, the lowest possible latency for transformer inference on a Raspberry Pi [10] is 2.1 seconds [11], which is too slow for real-time natural language processing (NLP) tasks. This

makes efficient training (and even inference) of such models an unsolved problem. The increasing size of state-of-the-art language models [2, 12] results in a higher memory footprint and computational complexity, exacerbating this problem.

Previous works propose many specialized hardware accelerators to address the abovementioned challenges. For instance, A<sup>3</sup> [13] is one of the first accelerators to enable efficient transformer inference by leveraging algorithmic approximation and hardware specialization. SpAtten [14] proposes a cascade token pruning mechanism to prune the weights of a transformer at runtime. Energon [15] approximates this pruning mechanism to speed up inference. AccelTran [16], a state-of-the-art transformer accelerator, executes dynamic inference by skipping *all* ineffectual multiply-and-accumulate (MAC) operations. It also *tiles* the matrices to facilitate higher parallelization and hardware utilization while executing the matrix multiplication operations. However, the above accelerators do not support transformer training, which demands a higher memory footprint for transformer execution on resource-constrained edge-AI devices.

Tuli et al. [17] showed that each NLP task has a unique optimal transformer architecture, requiring a specialized accelerator for efficient evaluation. However, designing a single accelerator that efficiently executes a diverse set of transformer architectures takes significant time and effort. Transformers involve serially-connected encoder (and sometimes decoder) layers. An accelerator with too many processing elements (PEs) would have low compute resource utilization for a deep but narrow transformer model. A PE is the basic compute module in an accelerator. A shallow and wide transformer would incur low latency and enable high parallelization [11, 18] but would require many PEs and high bandwidth memory. However, a manual search of the best accelerator architecture and transformer design decisions is computationally too expensive due to the vastness of each design space [19].

To tackle the abovementioned challenges, previous works implement automated hardware-software co-design. NAAS [20] and CODEBench [19] simultaneously search for the best design choices for the convolutional neural network (CNN) and accelerator architecture. However, a CNN workflow is different from that of a transformer, warranting substantially different accelerator design choices. Some recent works target co-design with transformer models. Qi et al. [21] use a recurrent neural network (RNN) and a reinforcement learning (RL)-based controller to guide the search using a pool of five field-programmable gate arrays (FPGAs) and adjust the pruning parameters of an input transformer model. Peng et al. [22] explore the scheduling

This work was supported by NSF Grant No. CNS-2216746. S. Tuli and N. K. Jha are with the Department of Electrical and Computer Engineering, Princeton University, Princeton, NJ, 08544, USA (e-mail: {stuli, jha}@princeton.edu).

Manuscript received —; revised —.

and sparsity decisions on a single FPGA. CODEBench [19], although a framework for CNN accelerators, shows the advantages of exploring massive CNN and accelerator search spaces, resulting in high gains in accuracy, energy consumption, evaluation throughput, and chip area. Hence, we leverage its co-design technique, BOSHCODE, in our proposed framework (details in Section II-C). On the other hand, the abovementioned works on transformer accelerator search only target one (or few) transformer models on a limited set of hardware platforms. This restricts the gains from automated co-design, leading to low resource utilization and inefficient configurations.

In order to address the above issues, we propose TransCODE, a co-design framework for transformers and application-specific integrated circuit (ASIC)-based accelerators. Our main contributions are as follows.

- For efficient on-device training, we propose DynaProp, which dynamically prunes weights, activations, and gradients to skip ineffectual MAC operations and speed up the transformer training/inference process. DynaProp leverages specialized low-overhead hardware modules to induce sparsity into transformer training and inference.
- To support vast design spaces involving *flexible* and *heterogeneous* transformer architectures [17], we propose a *flexible BERT accelerator* (ELECTOR) framework. ELECTOR supports diverse transformer architectures within the FlexiBERT 2.0 design space [11]. It efficiently implements model operations through dedicated hardware modules and a functional transformer mapper. The design space within the ELECTOR framework involves disparate accelerators that can execute the transformers in the FlexiBERT 2.0 design space. ELECTOR also effectively implements the proposed DynaProp algorithm to speed up transformer training and inference. It involves 14,850,000 accelerators, a design space much more extensive than investigated in any previous work.
- We then leverage the proposed ELECTOR and FlexiBERT 2.0 design spaces to implement co-design and obtain a transformer-accelerator pair that maximizes the performance objectives within the given user-defined constraints. We call this framework, which co-designs the transformer-accelerator pair, TransCODE. It leverages the best-performing optimization technique in the considered design spaces.

We organize the rest of the article as follows. Section II presents background on transformer and accelerator design choices along with automated hardware-software co-design. Section III illustrates the TransCODE framework that includes DynaProp, ELECTOR framework and its design space, and the co-design pipeline. Section IV describes the experimental setup and targeted baselines. Section V presents the results. Section VI discusses the limitations of the proposed work and future work directions. Finally, Section VII concludes the article.

## II. BACKGROUND AND RELATED WORK

In this section, we present background material on popular transformer and accelerator architectures and the correspond-

ing design decisions. We also describe previously proposed hardware-software co-design methods.

### A. Transformer Design Space

Previous works propose various transformer architectures. BERT is one of the most popular architectures that is widely used for language modeling [1]. Its variants leverage mechanisms other than vanilla self-attention [23] to optimize performance or reduce model size and complexity. They include RoBERTa [24] that implements robust pre-training techniques, ConvBERT [25] that uses one-dimensional convolutional operations, MobileBERT [26] that employs bottleneck structures and multiple feed-forward stacks, among many others. Further, architectures like FNet [27] and LinFormer [28] use Fourier transform and low-rank approximation, respectively, of the self-attention operation to aid efficiency and reduce the number of model parameters.

In order to search for the best-performing model for a given task, FlexiBERT [17] unifies and implements *heterogeneous* and *flexible* transformer architectures, encapsulating various self-attention operation types. Each encoder layer in its design space can have a different attention mechanism (heterogeneity) and a different hidden dimension (flexibility). Among many works that implement neural architecture search (NAS) in a design space of transformer models [29, 30, 31, 32], FlexiBERT has the largest and the most expressive design space. This results in state-of-the-art models that outperform previous architectures in accuracy. FlexiBERT 2.0 [11] extends the design space to  $1.7 \times 10^{88}$  transformer models, the largest and the most expressive transformer design space to date. We thus use the FlexiBERT 2.0 design space to implement co-design in this work. Note that no previously proposed accelerator supports heterogeneous and flexible transformer workflows. We discuss traditional transformer accelerators next.

### B. Accelerator Design Space

A transformer model’s hardware performance (characterized by latency, energy consumption, and chip area) on a given platform depends on multiple factors. These factors include memory size and bandwidth, number of MAC units (that can execute matrix multiplication operations in parallel), number of specialized hardware modules (e.g., ones for softmax and layer-norm operations), operation scheduling, dataflow, model sparsity, etc. These design decisions lead to many existing accelerators proposed in the literature.

A<sup>3</sup> [13] is one of the first ASICs to support transformer acceleration. It uses several approximation strategies to avoid computing attention scores that are close to zero. SpAtten [14] proposes the top- $k$  pruning algorithm that ranks input token and attention-head scores using a dedicated hardware module. However, it only considers part of the activations formed, not sparsity in all possible matrix multiplication operations. Further, implementing the proposed top- $k$  pruning mechanism involves high compute overhead; its time complexity is  $\mathcal{O}(N^3)$ , leading to marginal gains in energy efficiency [16]. Energon [15] approximates this pruning mechanism. However,

since it is limited to being a co-processor, it requires high off-chip memory access. Finally, OPTIMUS [33] targets sparsity in a broader scope, using a set-associative rearranged compressed sparse column format to eliminate ineffectual MAC operations, although limited to weight matrices. Here, weights correspond to the trainable transformer model parameters and activations are represented by intermediate matrices formed by the transformer model operations.

To overcome the drawbacks of the abovementioned accelerators, AccelTran [16] implements dynamic inference with a transformer while pruning *all* weights and activations. In addition, it leverages matrix tiling to improve parallelization and resource utilization. However, it only executes transformer inference and not training, uses a fixed set of design choices (e.g., a fixed tile size, number of PEs, buffer sizes), and does not support diverse models, thus leading to sub-optimal utilization. To tackle this problem, various works propose design spaces of transformer accelerators to efficiently obtain the optimal transformer architecture for the given task. However, such design spaces are limited to off-the-shelf FPGAs [21, 22] that only focus on inference. We next describe previous works on co-design of the AI model and hardware accelerator.

### C. Hardware-software Co-design

Various works target CNN-accelerator co-design [19, 20, 34]. CODEBench [19] searches over massive CNN and accelerator design spaces. However, its accelerators are geared toward CNN workflows and thus inefficient for transformer pipelines. As discussed before, Qi et al. [21] use an RNN and RL-based controller to guide search in a pool of five FPGAs and adjust the pruning parameters of an input transformer model. However, they only consider latency and accuracy constraints and do not optimize energy consumption and chip area. Peng et al. [22] explore the scheduling and sparsity decisions on an FPGA and adapt the input sequence length. SpAtten [14] implements hardware-aware NAS (HW-NAS), where it finds a sub-net of a trained super-net [18]. However, its design space only involves 72 transformers that are not flexible. Thus, there is a need for an exhaustive design space of transformer accelerators to implement co-design and obtain the best-performing transformer-accelerator pair. This pair should not only deliver high accuracy on a given task but also be energy-efficient and have a high throughput and low chip area.

In this work, we leverage Bayesian optimization using second-order gradients and a heteroscedastic surrogate model for co-design, i.e., BOSHCODE [19]. It is a scalable co-design framework that efficiently searches the hardware and software design spaces at scale. CODEBench [19] proposes and uses BOSHCODE to search over significantly large design spaces ( $4.2 \times 10^{812}$  CNNs and  $2.3 \times 10^8$  accelerators). EdgeTran [11] leverages BOSHCODE to search over the joint space of FlexiBERT 2.0 and a set of off-the-shelf edge-AI devices, including Raspberry Pi [10], Apple M1 [35] system-on-chip (SoC) with a central processing unit, CPU, and a GPU, Intel Neural Compute stick [36] (a neural processing unit), and Nvidia Jetson Nano [37] (SoC with both CPU and GPU).

BOSHCODE supports co-design with any two search spaces. It leverages second-order gradient-based optimiza-

tion [38] on an actively-trained [39] surrogate model for performance prediction (which is the optimization objective). The surrogate model combines a natural parameter network (NPN), a teacher, and a student network. The NPN predicts the mean performance of the transformer-accelerator pair along with the aleatoric uncertainty. The teacher and student networks predict the epistemic uncertainty in performance. Epistemic uncertainty is the uncertainty in performance due to an unexplored design space. In contrast, aleatoric uncertainty refers to the uncertainty due to parameter initializations and variations in model performance due to different training recipes. BOSHCODE exploits epistemic and aleatoric uncertainty estimates to obtain the best design decisions for the transformer, the accelerator, and the model training recipe that maximizes accuracy. We present more details on how we leverage BOSHCODE in our search process in Section III-C.

## III. METHODOLOGY

Fig. 1 shows an overview of the TransCODE framework. ELECTOR, in Fig. 1(a), takes the accelerator embedding and transformer computational graph as input. Using the accelerator embedding, it implements a hardware accelerator with the corresponding design decisions. Next, it converts the computational graph into a corresponding transformer model with modular operations (supported by the FlexiBERT 2.0 design space), which it then maps to specialized hardware modules. It also tiles the matrices for efficient resource allocation, operation scheduling, and data reuse. Fig. 1(b) shows how we leverage the FlexiBERT 2.0 [17] framework to convert a transformer embedding to its corresponding computational graph and employ the surrogate model to predict model accuracy. Finally, Fig. 1(c) illustrates TransCODE, which uses previous performance results to train a surrogate model and query the next transformer-accelerator pair. Finally, it feeds the output accelerator and transformer embeddings to ELECTOR and FlexiBERT 2.0, respectively.

We now discuss the dynamic inference and training technique, DynaProp, that prunes activations and gradients to skip ineffectual operations. We then present the ELECTOR simulator and the accelerator design choices it supports. Finally, we describe the TransCODE pipeline that implements co-design and obtains the best-performing transformer-accelerator pair.

### A. Dynamic Inference and Training

DynaTran [16] is a low-overhead dynamic inference method that quickly prunes ineffectual weight and activation values at runtime. However, it only targets transformer inference and not training. We propose DynaProp that induces sparsity in weights and activations at runtime (during inference) and gradients (during training). DynaProp takes an input matrix, which is either a weight matrix (loaded from memory), an activation matrix (obtained from previous MAC operations), or a gradient matrix (formed while backpropagating gradients). It then prunes values with a magnitude less than a given

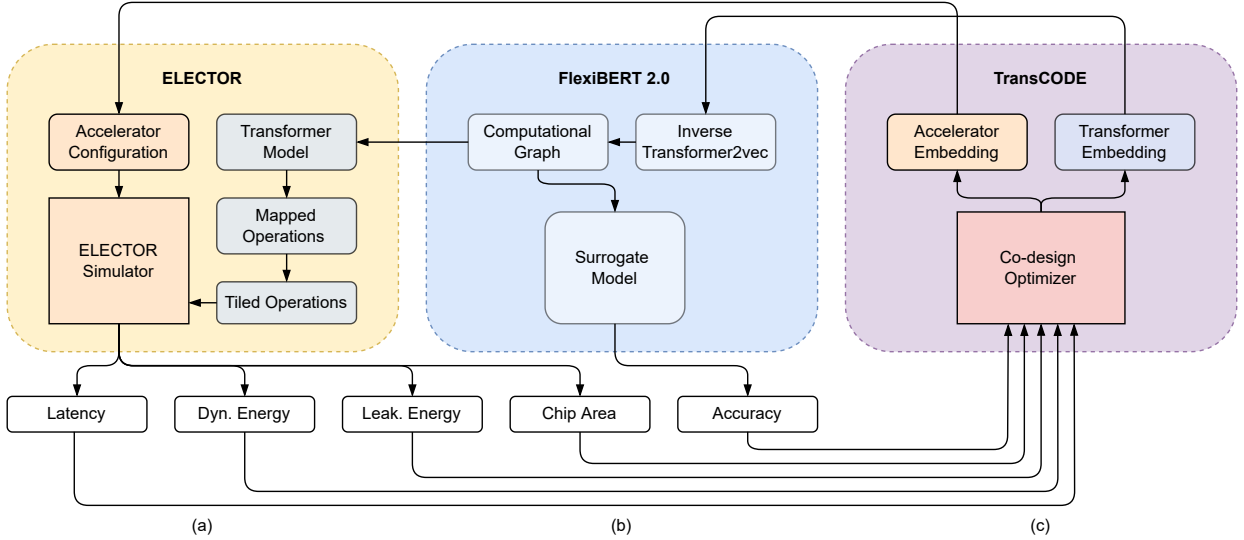


Fig. 1: Overview of the TransCODE framework. (a) ELECTOR takes an accelerator embedding and a transformer computational graph to simulate its training/inference on the given accelerator. (b) FlexiBERT 2.0 converts the input transformer embedding to a computational graph and employs a pre-trained surrogate model to predict model accuracy. (c) The TransCODE optimizer takes in the performance values of the previously evaluated transformer-accelerator pair to query another pair in the active learning loop.

TABLE I: Forward and backward pass operations for matrix multiplication and 1D convolution.

Matrix Multiplication	
Forward Pass	$\mathbf{X}_i = f_i(\mathbf{W}_i \mathbf{X}_{i-1})$
Backward Pass	$\delta_i = \mathbf{W}_{i+1}^T \delta_{i+1} \cdot f'(\mathbf{W}_i \mathbf{X}_{i-1})$
Weight Update	$\nabla_{\mathbf{W}_i} \mathcal{L} = \delta_i \mathbf{X}_{i-1}^T$ $\mathbf{W}_i = \mathbf{W}_i - \alpha \cdot \nabla_{\mathbf{W}_i} \mathcal{L}$
1D Convolution	
Forward Pass	$\mathbf{X}_i = \mathbf{w}_i * \mathbf{X}_{i-1}$
Backward Pass	$\delta_i = \nabla_{\mathbf{W}_i} \mathcal{L}$
Weight Update	$\nabla_{\mathbf{W}_i} \mathcal{L} = \delta_i * \overleftarrow{\mathbf{X}_{i-1}}$ $\mathbf{W}_i = \mathbf{W}_i - \alpha \cdot \nabla_{\mathbf{W}_i} \mathcal{L}$

threshold  $\tau$  (i.e., it forces them to zero). Mathematically, we prune an input matrix  $\mathbf{M} \in \mathbb{R}^{m \times n}$  to  $\mathbf{M}^P$  as follows:

$$\mathbf{M}_{ij}^P = \begin{cases} \mathbf{M}_{ij} & \text{if } |\mathbf{M}_{ij}| \geq \tau \\ 0 & \text{if } |\mathbf{M}_{ij}| < \tau \end{cases}$$

We implement each such comparison in parallel, thus requiring only one clock cycle for the pruning process. We define the pruning ratio (or level of sparsity) of the output matrix as:

$$\rho(\mathbf{M}^P) = \frac{\sum_{x \in \mathbf{M}^P} \delta_{x,0}}{m \times n}$$

where  $\delta$  is the Kronecker delta function. We profile the resultant sparsity in weights, activations, and gradients for different transformer models on diverse applications to obtain a desired  $\rho$ . ELECTOR stores these curves in memory. For the desired values of  $\rho$ , we determine the corresponding  $\tau$  at

runtime through a simple look-up operation. We present such curves in Section V-A.

Table I shows the operations underlying the forward and backward pass for matrix multiplication and one-dimensional (1D) convolution, respectively. The table shows that training requires the same operation types (as inference) and thus mandates identical hardware, although with a separate dataflow. We also observe that the number of backward pass and weight update operations (executed during training) is more than the number of those for the forward pass (executed during inference). This shows that training is much more computationally expensive than inference, involving more activations and gradients that the accelerator needs to account for. DynaProp prunes each such matrix before it executes the respective operation in hardware. Thus, the accelerator skips ineffectual operations, improving latency and energy efficiency.

Optimizers like Adam would require extra computation (e.g., the calculation of momentum and storage of previous weights/gradients). These computations can easily be incorporated into the accelerators supported in the proposed design space. However, second-order gradients would add much more computational overhead. We leave the application of complex optimizers to future work.

### B. The ELECTOR Framework

Accelerators in the ELECTOR design space take inspiration from previously proposed state-of-the-art accelerators, including SPRING [40] and AccelTran [16]. We divide the overall accelerator architecture into the accelerator tier and the (on-chip or off-chip) memory tier. Fig. 2 shows the organization of the accelerator tier in the proposed architecture. The control block receives the instruction stream of the

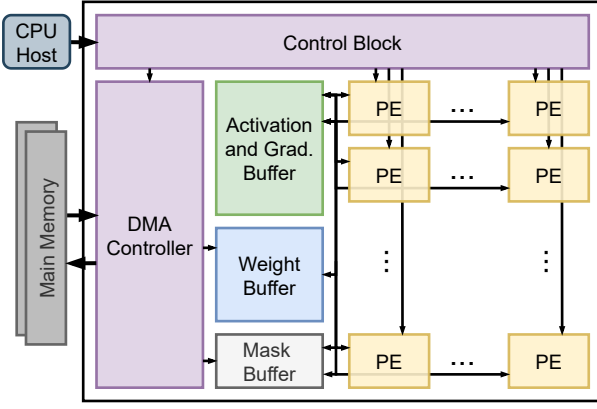


Fig. 2: Organization of a typical accelerator in the ELECTOR design space.

transformer model from the host CPU. The direct memory access (DMA) controller fetches the weights and embeddings from the main memory. Thus, the PEs communicate with the on-chip buffers while the DMA controller transfers data between the buffers and the on-chip/off-chip main memory. The activation-and-gradient buffer stores the activations and gradients formed during transformer evaluation. The weight buffer stores the transformer weights. ELECTOR stores all data in a compressed format (discussed in Section III-B1). Data compression relies on binary masks (stored in the mask buffer). The PEs employ the compressed data and associated masks to perform the main compute operations of any transformer in the FlexiBERT 2.0 design space.

1) *Hardware Modules*: We now describe various modules supported in the ELECTOR design space.

- **Main Memory**: ELECTOR supports three memory types: an off-chip dynamic random access memory (DRAM) for scalable and economical deployments, an on-chip high-bandwidth memory (HBM) for memory-intensive edge/server applications, and an on-chip monolithic-3D resistive random access memory (RRAM). Monolithic-3D integration leverages monolithic inter-tier vias, allowing much higher density than traditional through-silicon-via-based 3D integration [41]. This leaves much more logic space and permits high memory bandwidth, which are crucial for large transformer models in the FlexiBERT 2.0 design space.
- **Control Block**: The control block takes the transformer model as input. It then converts all functions in the model into hardware-mappable operations (details in Section III-B2) that it later converts to *tiled* operations. For instance, it converts the matrix multiplication operation  $\mathbf{O} = \mathbf{W} \times \mathbf{A}$  to multiple operations of the form  $\mathbf{O}[\mathbf{b}, \mathbf{i}, \mathbf{j}] = \mathbf{W}[\mathbf{b}, \mathbf{i}, \mathbf{k}] \times \mathbf{A}[\mathbf{b}, \mathbf{k}, \mathbf{j}]$ , where each tiled matrix  $\in \mathbb{R}^{b \times x \times y}$ , i.e., the tile size [16]. The control block also assigns and schedules the tiled operations to different PEs [16].
- **Processing Elements**: Fig. 3 shows the organization of a PE (the basic compute module of an accelerator) in the ELECTOR design space. The local registers of the

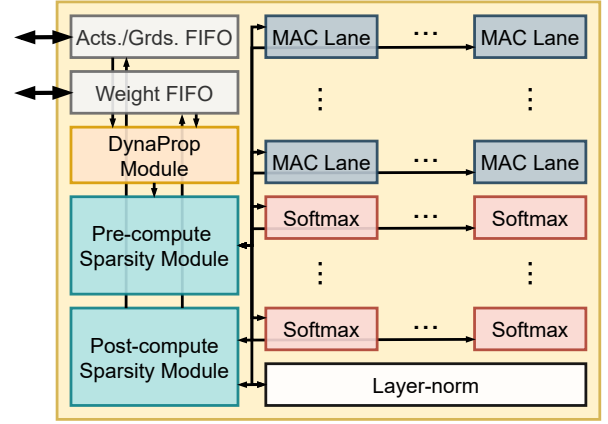


Fig. 3: Internal components of a PE.

PE store the compressed data. These are the first-in-first-out (FIFO) registers for the activations (and gradients) and weights. The data then enter the DynaProp module that induces sparsity based on the desired  $\rho$ . As explained in Section III-A, this module prunes the given activation/gradient/weight matrices based on a pre-calculated threshold  $\tau$ . The PE then feeds the sparse data to the pre-compute sparsity module with the binary masks. These binary masks have the same shape as the uncompressed data, where each binary bit in a mask depicts if the corresponding element in the original data vector is ineffectual or not. The pre-compute sparsity module converts the input data into a zero-free format based on the associated masks [16]. The PE then forwards this zero-free data to the MAC lanes (for matrix multiplication), softmax modules (for softmax operation), or the layer-norm module (for layer-norm operation). The zero-free data eliminate any ineffectual computations in these modules. Finally, the post-compute sparsity module [16] implements the inverse of this operation on the output activations before storing them in the corresponding FIFO register and, eventually, the main buffer.

- The MAC lanes execute multiplication between two tiles in a parallelized manner. We store all activation, gradient, and weight data in fixed-point format with  $(\text{IL} + \text{FL})$  bits, denoting integer length (IL) and fractional length (FL), respectively [16]. Data first reach  $M$  multipliers and then an adder tree with depth  $\log_2 M$ . The MAC lanes also include a ReLU and a GeLU [42] module for feed-forward operations.
- Fig. 4 shows the DynaProp module that executes dynamic inference and training on the transformer. It takes the input activation/gradient/weight matrix and prunes ineffectual values for efficient evaluation. As explained in Section III-A, we prune the values of the input matrix by comparing their magnitude with a pre-determined threshold  $\tau$ . The DynaProp module implements this in parallel for the entire tile. We first feed an input tile  $\mathbf{M} \in \mathbb{R}^{b \times x \times y}$  to the matrix transpose block, which carries out the transpose operation, if required. Mathematically, it

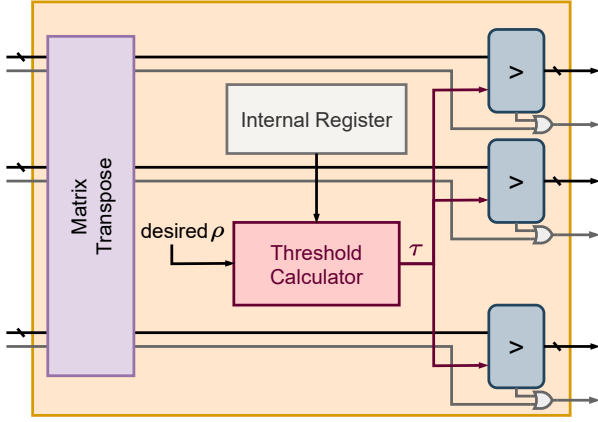


Fig. 4: Implementation of the DynaProp module. The wires for mask bits are in grey.

outputs  $\mathbf{M}^\top \in \mathbb{R}^{b \times y \times x}$ , transposing all matrices in the batch of size  $b$ . It then feeds the input tile to  $b \times x \times y$  comparators. The threshold calculator determines the required threshold using the desired  $\rho$  and the pre-profiled transfer functions for different transformer models on diverse applications (stored in the internal register; more details in Section V-A). If the output of the comparator is zero, we set the corresponding mask bit to 1. Here, we represent the lines carrying mask information in grey and those carrying activation/gradient/weight information in black.

- For all other hardware modules, we use the proposed implementation of AccelTran [16]. However, we expand the operation executability of all modules (e.g., support for different tile sizes in the softmax module), as explained in Section III-B3.

The optimal selection of the number of PEs, buffer sizes, and other design choices results in the highest possible resource utilization while minimizing the number of compute/memory stalls (when we do not execute either a compute operation or a memory fetch operation). Hence, determining the best accelerator hyperparameters is essential for energy-efficient designs with a low chip area and high throughput.

2) *The Transformer Mapper*: The FlexiBERT 2.0 [11] design space supports various operation types. We describe each operation next.

- **Self-attention**: The self-attention (SA) operation finds how much one token *attends* to another token. For an output attention head  $\mathbf{H}_i \in \mathbb{R}^{N_T \times d_{out}}$  with query  $\mathbf{Q}_i \in \mathbb{R}^{N_T \times h/n}$ , key  $\mathbf{K}_i \in \mathbb{R}^{N_T \times h/n}$ , and value  $\mathbf{V}_i \in \mathbb{R}^{N_T \times h/n}$  matrices [17]:

$$\mathbf{H}_i = \text{softmax}(\text{SA}) \mathbf{V}_i \mathbf{W}_i^o$$

where  $N_T$  is the input sequence length,  $h$  is the hidden dimension of the encoder layer, and  $n$  is the number of heads. The SA operation has two sub-types:

- The scaled dot-product (SDP) attention [4] is the de-facto standard operation in traditional transformer

architectures. Mathematically,

$$\text{SA}_{\text{SDP}} := \frac{\mathbf{Q}_i \mathbf{K}_i^\top}{\sqrt{h}}.$$

- The weighted multiplicative attention (WMA) [43] involves a trainable weight matrix  $\mathbf{W}_a \in \mathbb{R}^{h/n \times h/n}$  such that

$$\text{SA}_{\text{WMA}} := \frac{\mathbf{Q}_i \mathbf{W}_a \mathbf{K}_i^\top}{\sqrt{h}}.$$

The mapper converts the self-attention operation into various MAC and softmax operations that the corresponding hardware modules can execute in the accelerator.

- **Linear Transform**: As the name suggests, this operation implements a linear transform (LT) on the input sequence. The FlexiBERT 2.0 design space supports two sub-types:

- The discrete Fourier transform (DFT) that we implement in hardware using the corresponding Vandermonde matrix  $\mathbf{V}_{\text{DFT}} \in \mathbb{R}^{N_T \times N_T}$  for the roots of unity (also called the DFT matrix) [44] such that

$$\text{LT}_{\text{DFT}} := \mathbf{V}_{\text{DFT}} \mathbf{H}$$

where  $\mathbf{H} \in \mathbb{R}^{N_T \times d_{in}}$  represents a matrix for the input hidden states.

- The discrete cosine transform (DCT) that we again implement using an equivalent Vandermonde matrix  $\mathbf{V}_{\text{DCT}} \in \mathbb{R}^{N_T \times N_T}$  such that

$$\text{LT}_{\text{DCT}} := \mathbf{V}_{\text{DCT}} \mathbf{H}$$

We store the  $\mathbf{V}_{\text{DFT}}$  and  $\mathbf{V}_{\text{DCT}}$  matrices in the buffer for subsequent use while executing the above operations. Although these operations are slower than the fast Fourier transform (FFT) [45] and the fast cosine transform (FCT) [46], respectively, sparsification of the input matrices results in a low overall execution time. Furthermore, converting these operations to MAC operations enables the reuse of the MAC lanes, thus not requiring special hardware modules for the LT operation. Nevertheless, these methods (FFT and FCT) may lead to high gains for transformer models that support long sequences [47], due to their  $\mathcal{O}(N \log N)$  complexity. We leave their hardware implementation to future work.

- **Dynamic-span-based Convolution**: The dynamic-span-based convolution (DSC) operation implements a 1D convolution over the input. Mathematically,

$$\text{DSC}_k := \mathbf{w}_k * \mathbf{H}$$

where  $\mathbf{w}_k$  is the convolution kernel of length  $k$ . To implement this operation in hardware, we convert the convolution operation into an equivalent matrix multiplication operation. In other words, we convert the convolutional kernel to a sparse matrix that we multiply with the input. We tweak the MAC lane module to incorporate this conversion.

Now that the mapper has converted the operations in the FlexiBERT 2.0 design space to hardware-implementable formats, the control block tiles, schedules, and assigns these

TABLE II: Hyperparameters supported in the ELECTOR design space.

Hyperparameter	Permissible values
Batch tile size	1, 4
Spatial tile size	8, 16, 32
Activation function	ReLU, GeLU
#PEs	64, 128, 256, 512, 1024
#MAC lanes per PE	8, 16, 32, 64, 128
#MACs per lane	1, 16
#Softmax modules per PE	2, 4, 8, 16, 32, 64
Batch size	4, 16, 32
Act./grad. buffer size (MB)	4, 8, 16, 32, 64
Weight buffer size (MB)	8, 16, 32, 64, 128
Mask buffer size (MB)	1, 2, 4, 8, 16
	RRAM: [16, 2, 2], [8, 2, 4], [4, 2, 8], [2, 2, 16], [32, 2, 1], [1, 2, 32]
Main memory configuration [banks, ranks, channels]	DRAM: [16, 2, 2], [8, 2, 4], [32, 2, 1], [16, 4, 1]
	HBM: [32, 1, 4]

mapped operations to the accelerator for transformer evaluation.

3) *Design Space*: ELECTOR supports various accelerators in its design space. It allows adaptation of many design decisions in an ASIC-based accelerator. We describe these tunable hyperparameters next.

- **Batch Tile Size**: This is the size of a tile along the batch. Mathematically, a tile  $M \in \mathbb{R}^{b \times x \times y}$  has the batch tile size  $b$ .
- **Spatial Tile Size**: This is the size of a tile orthogonal to the batch dimension. In the above example,  $x = y$  is the spatial tile size (we assume square matrices for the tiles). A higher tile size (either  $b$  or  $x/y$ ) would imply that each hardware module (MAC lane, softmax module, or layer-norm module) could execute more operations in parallel since the module evaluates a larger tile. This enables latency reduction at the cost of higher dynamic power.
- **Activation Function**: Transformer evaluation uses a non-linear function following a feed-forward operation. We support two functions: ReLU and GeLU [42]. This is in accordance with the FlexiBERT 2.0 design space.
- **Number of PEs**: The number of PEs in the accelerator.
- **Number of MAC Lanes per PE**: The number of MAC lanes in each PE of the accelerator. We keep the number of MAC lanes constant for every PE.
- **Number of MACs per Lane**: The number of MAC units per MAC lane. Again, this is constant across all MAC units.
- **Number of Softmax Modules per PE**: The number of softmax modules in each PE. Every PE has only one layer-norm module. Therefore, the number of MAC lanes and softmax modules in each PE determines the net ratio of the number of MAC lanes, softmax modules, and layer-norm modules in an accelerator. One can tune this ratio based on the corresponding proportion of these operations in evaluating the selected transformer.
- **Batch Size**: The batch size for transformer evaluation. More compute resources and high bandwidth memory enable a larger batch, reducing evaluation latency.

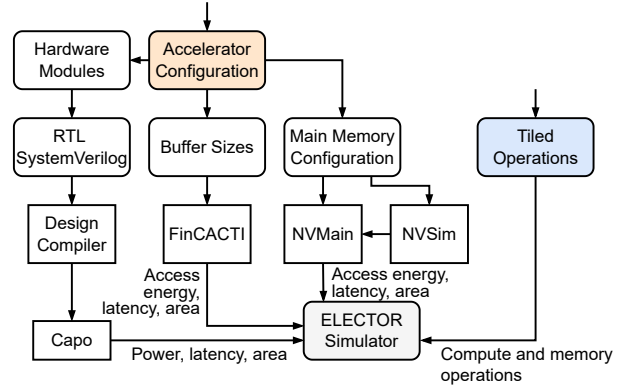


Fig. 5: Flow of simulation in ELECTOR.

- **Activation and Gradient Buffer Size**: The size of the activation/gradient buffer. Training requires more activation matrices than inference. It also has gradient matrices, requiring a larger buffer size.
- **Weight Buffer Size**: The size of the weight buffer. A larger transformer model requires a larger weight buffer.
- **Mask Buffer Size**: The size of the mask buffer that stores the binary masks for the zero-free format [16] used in the accelerators in ELECTOR.

Table II summarizes the possible design choices for accelerators in the ELECTOR design space. The possible memory configurations include the memory type (RRAM, DRAM, and HBM) along with the banks, ranks, and channels.

4) *Accelerator Embeddings*: We now describe how we convert the selected accelerator configuration (a sample from Table II) to an embedding for surrogate modeling. We generate a 12-dimensional embedding ( $e$ ) for a selected accelerator configuration as follows:

- $e_1$  denotes the batch tile size, i.e.,  $e_1 = b$ .
- $e_2$  and  $e_3$  correspond to the spatial tile sizes, i.e.,  $e_2 = x$ ,  $e_3 = y$ . For the targeted design space,  $e_2 = e_3$ .
- $e_4$  denotes the number of PEs.
- $e_5$  denotes the number of MAC lanes per PE.
- $e_6$  denotes the number of MACs per lane.
- $e_7$  denotes the number of softmax modules in each PE.
- $e_8$  denotes the selected batch size for model evaluation.
- $e_9$ ,  $e_{10}$ , and  $e_{11}$  denote the activation/gradient, weight, and mask buffer sizes, respectively, in MBs.
- $e_{12}$  denotes the index of possible memory configurations in Table II, thus ranges from 1 to 11.

We use these generated embeddings to train the TransCODE surrogate model, which also outputs the subsequent query as an accelerator embedding.

5) *Simulation Flow*: Fig. 5 shows the simulation flow for evaluating an input accelerator configuration and tiled operations (obtained after mapping and tiling the input transformer) in ELECTOR. We first select the compute modules (including the tile size for parallel operation), buffer sizes, and main memory configuration. Next, we implement different hardware modules discussed in Section III-B1 at the register-transfer level (RTL) using SystemVerilog. We use Design Compiler [48] to synthesize the RTL design based on a

TABLE III: Hyperparameter ranges in FlexiBERT 2.0 design space [11]. Super-script ( $j$ ) depicts the value for layer  $j$ .

Design Element	Allowed Values
Number of encoder layers ( $l$ )	{2, 4, 6, 8, 10, 12}
Type of attention operation used ( $o^j$ )	{SA, LT, DSC}
Number of operation heads ( $n^j$ )	{2, 4, 8, 12}
Hidden size ( $h^j$ )	{128, 256}
Feed-forward dimension ( $f^j$ )	{256, 512, 1024, 2048, 3072, 4096}
Number of feed-forward stacks	{1, 2, 3}
Operation parameters ( $p^j$ ):	
if $o^j = \text{SA}$	Self-attention type: {SDP, WMA}
else if $o^j = \text{LT}$	Linear transform type: {DFT, DCT}
else if $o^j = \text{DSC}$	Convolution kernel size: {5, 9}

14nm FinFET technology library [49]. Capo [50], an open-source floorplacer, performs floorplanning. FinCACTI [51], a cache modeling tool for deeply-scaled FinFETs, models the on-chip buffers. NVSim [52] and NVMain [53] model the main memory (either the off-chip DRAM or on-chip HBM/RRAM). ELECTOR then plugs the synthesized results into a Python-based cycle-accurate simulator. Finally, the control block segregates the tiled operations into compute and memory operations for separate execution pipelines [16].

### C. TransCODE

We use BOSHCODE to obtain the best-performing transformer-accelerator pair. BOSHCODE takes as input the accelerator and transformer embeddings and outputs the performance measure to be estimated. For the transformer embeddings, we use the embeddings used in FlexiBERT 2.0 [11] as opposed to the `Transformer2vec` encodings [17] since they are fast and efficient. This is critical for exploring the vast FlexiBERT 2.0 design space efficiently. For the accelerator embeddings, we use the embeddings from the accelerator configuration discussed in Section III-B4. We define the output performance measure as follows:

$$\begin{aligned} \text{Performance} = & \alpha \times (1 - \text{Latency}) + \beta \times (1 - \text{Area}) \\ & + \gamma \times (1 - \text{Dynamic Energy}) \\ & + \delta \times (1 - \text{Leakage Energy}) + \epsilon \times \text{Accuracy} \end{aligned}$$

where  $\alpha + \beta + \gamma + \delta + \epsilon = 1$  are hyperparameters. We normalize the values of the individual performance measures with respect to their maximum values (hence, these values reside in the  $[0, 1]$  interval). Thus, for edge applications where the power envelope of devices is highly restricted, users can set the hyperparameters  $\gamma$  and  $\delta$  high. On the other hand, for server-side deployments, where accuracy is of utmost importance, one can set  $\epsilon$  high.

TransCODE needs five performance values for the queried transformer-accelerator pair: latency, area, dynamic energy, leakage energy, and model accuracy. To obtain the first four performance values, we leverage the ELECTOR simulator. To obtain the transformer model accuracy, we employ the FlexiBERT 2.0 surrogate model, which outputs the GLUE score [54].

TABLE IV: Data statistics of datasets in the GLUE benchmark.

Task	Training Size	Metric
SST-2	67K	Accuracy
MNLI	393K	Accuracy
QQP	364K	Accuracy
QNLI	105K	Accuracy
MRPC	3.7K	Accuracy
CoLA	8.5K	Matthew's Correlation
STS-B	7K	Spearman Correlation
RTE	2.5K	Accuracy
WNLI	634	Accuracy

## IV. EXPERIMENTAL SETUP

In this section, we present the setup behind various experiments we performed, along with the baselines considered for comparison.

### A. Evaluation Models and Datasets

To test the efficacy of the DynaProp method, we evaluate transformer models in the FlexiBERT 2.0 design space. Table III shows the hyperparameter ranges supported by the FlexiBERT 2.0 design space [11]. Evidently, shallow models (e.g., with two encoder layers) incur lower latency relative to deep models (e.g., with 12 encoder layers) [11, 18]. Moreover, wide models (e.g., with 12 attention heads) require more compute resources to enable higher parallelization than narrow ones (e.g., with two attention heads). Further, different attention-head types have different latencies and energy consumption characteristics. Hence, there is a need for optimized dataflows when executing such heterogeneous architectures.

We test the models on representative natural language understanding tasks under the GLUE benchmark [54]. The included tasks are: SST-2 [55], MNLI [56], QQP, QNLI, MRPC [57], CoLA [58], STS-B [59], RTE [60], and WNLI [61]. The surrogate model trained on the FlexiBERT 2.0 design space [11] reports the overall GLUE score. We show the training sizes and used metrics in Table IV. The GLUE score represents average performance across all the tasks.

While running DynaProp, we target activation, weight, and gradient sparsity. Weight sparsity is static and depends on pruning performed during model pre-training or finetuning [62]. Activation and gradient sparsity change for every input sequence – we report their averages over the entire validation set.

### B. The ELECTOR Design Space

Table II summarizes the ELECTOR design space. Taking into account all the possible combinations presented in this table, ELECTOR supports 14,850,000 accelerators in its design space. This space includes accelerators meant for resource-constrained edge applications as well as those relevant to high-energy server settings that require high throughput. In addition, ELECTOR allows different memory configurations to support diverse user requirements, from high-bandwidth monolithic-3D RRAM to economic off-chip DRAM.



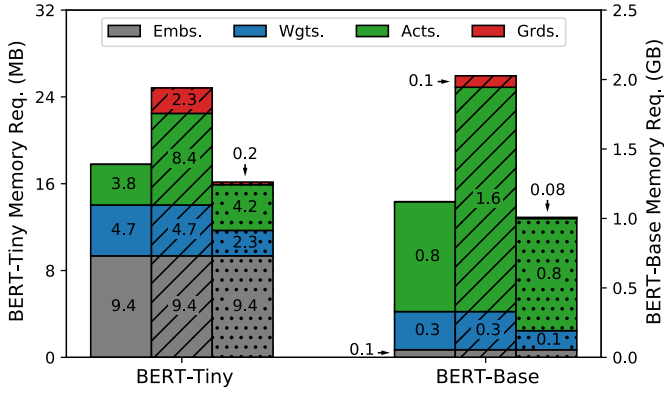


Fig. 6: Breakdown of the total memory required while running inference (solid bars), traditional training (hatched bars), and DynaProp training (dotted bars) for BERT-Tiny and BERT-Base.

### C. Co-design Pipeline

To run BOSHCODE, we use the following parameter values to obtain the net performance measure:  $\alpha = 0.1$ ,  $\beta = 0.1$ ,  $\gamma = 0.2$ ,  $\delta = 0.1$ , and  $\epsilon = 0.5$  (see Section III-C). We leverage the network and hyperparameters used in EdgeTran [11] for co-design. The BOSHCODE model takes  $x_{\text{TXF}}$  and  $x_{\text{ACC}}$  as input and outputs the predicted performance measure. Here,  $x_{\text{TXF}}$  and  $x_{\text{ACC}}$  correspond to the FlexiBERT 2.0 and ELECTOR embeddings, respectively. BOSHCODE then leverages gradient-based optimization using backpropagation to the input (GOBI) [38] while freezing the model weights.

All input embeddings obtained using GOBI from the surrogate models may not be valid. For instance,  $x_{\text{ACC}}$  should be well-defined (e.g., we allow the batch tile size,  $b$ , to only be 1 or 4). To add constraints to the optimization process, along with forcing the model to learn the performance only for valid input embeddings, we add a datapoint ( $x_{\text{TXF}}$ ,  $x_{\text{ACC}}$ ,  $P_{\text{MIN}}$ ) to the dataset if either of the input embeddings is invalid or does not adhere to user-defined constraints. Another example of an input constraint could be that transformers with only up to six layers are allowed.  $P_{\text{MIN}}$  has a low value, set to  $-1$  for our experiments (where well-defined inputs would result in  $P$  to lie in the  $[0,1]$  range).

### D. Evaluation Baselines

We compare our experimental results with previously proposed transformer-accelerator pairs. The baseline accelerators include SpAtten and AccelTran, hand-designed for a specific transformer architecture. SpAtten implements HW-NAS. For fair comparisons, we present an HW-NAS version of TransCODE in which we execute BOSHCODE while forcing gradients to the accelerator to zero, i.e., we only search for transformer models run on a given edge platform. We also include co-design baselines implemented on a set of FPGAs [22].

## V. EXPERIMENTAL RESULTS

In this section, we present the experimental results and comparisons of the TransCODE framework with relevant

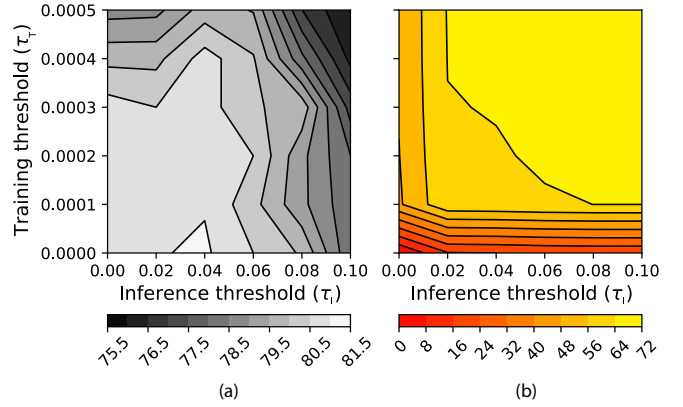


Fig. 7: Effect of training ( $\tau_T$ ) and inference ( $\tau_I$ ) thresholds on (a) accuracy (%) of the SST-2 task and (b) averaged activation and gradient sparsity (%) in BERT-Tiny.

baselines.

### A. Dynamic Pruning of Transformer Weights, Activations, and Gradients

Fig. 6 shows a breakdown of the memory required for three evaluation modes: inference, traditional training, and DynaProp training, for BERT-Tiny and BERT-Base. The evaluation mode does not affect the memory usage for the token and position embeddings. Moreover, inference and training require the same memory size for transformer model weights. However, training generates gradients and also more activation operations. Here, the  $\delta$ 's described in Table I define the gradient memory consumption. For BERT-Tiny and BERT-Base, training requires  $2.8\times$  and  $2.1\times$  more memory (for activations and gradients), respectively. However, the buffer can be smaller since it only stores the activations or gradients required by the PEs at a given time. Finally, we show the memory required for DynaProp training. We configure DynaProp to induce 50% sparsity in weights and activations (resulting in no loss in accuracy [16]) and 90% sparsity in the gradients (marginal accuracy loss, as shown below). DynaProp thus requires  $1.5\times$  and  $1.9\times$  smaller memory for BERT-Tiny and BERT-Base, respectively, while running training. This results in a smaller main memory, smaller buffers, and fewer MAC operations, thus leading to improved throughput.

To decouple and study the effects of pruning while running model inference and training, we execute DynaProp with two pruning thresholds:  $\tau_I$  and  $\tau_T$ . It prunes activation and gradient matrices using  $\tau_I$  for the forward pass and  $\tau_T$  for the backward pass. It leverages movement pruning [62] for transformer weights [16]. Fig. 7(a) presents a contour plot showing the effect of these thresholds on accuracy for the BERT-Tiny model. As previously observed [15, 16], the accuracy first increases and then decreases as we increase  $\tau_I$ . However, accuracy monotonically decreases on increasing  $\tau_T$ . Fig. 7(b) shows the average between activation and gradient sparsities (or the net sparsity) when changing  $\tau_I$  and  $\tau_T$ . The net sparsity increases as both  $\tau_I$  and  $\tau_T$  increase.

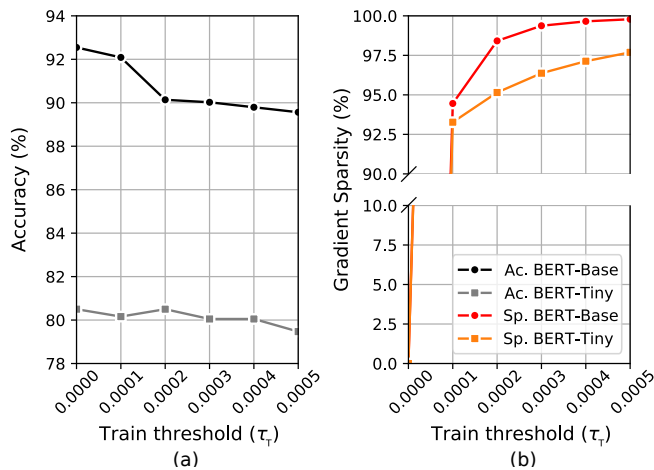


Fig. 8: Effect of changing training threshold ( $\tau_T$ ) on (a) accuracy of the SST-2 task and (b) sparsity in gradient matrices.

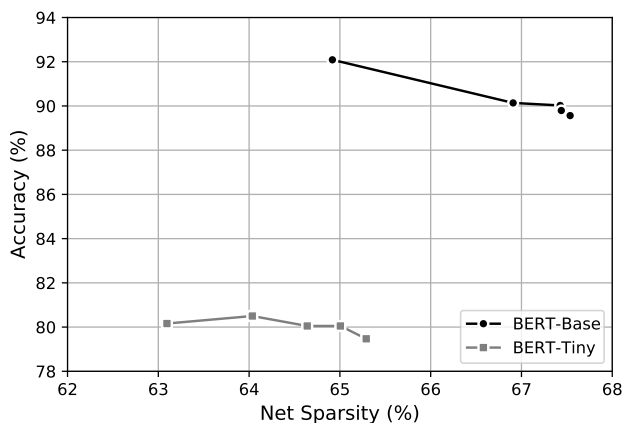


Fig. 9: Accuracy of the SST-2 task plotted against net sparsity in activation and gradient matrices.

Fig. 8(a) shows accuracy plotted against  $\tau_T$ . Unlike pruning during inference (using  $\tau_I$ ), accuracy decreases under DynaProp as we increase the pruning threshold  $\tau_T$ . However, this loss in accuracy is a result of high gradient sparsity, as shown in Fig. 8(b). This enables ELECTOR to skip many ineffectual MAC operations, reducing energy consumption and latency. We achieve 90% gradient sparsity when we set  $\tau_T$  to 0.0001 with an accuracy loss of only 0.4%. Fig. 9 shows a plot of accuracy against net sparsity. Again, we define *net sparsity* as the average of the activation and gradient sparsities (weight sparsity remains constant at 50%). The plot shows that accuracy decreases with increasing net sparsity for BERT-Base. However, for BERT-Tiny, accuracy increases and decreases as we increase net sparsity.

Fig. 10 shows a plot of the normalized time for traditional and DynaProp training. Here, we evaluate the BERT-Tiny model on an Nvidia A100 GPU and an ELECTOR-supported accelerator (AccelTran-Edge [16] with added training support). Training takes  $761.9\times$  longer than inference on a GPU. However, ELECTOR only requires  $1.6\times$  more time. This is due to optimized scheduling, tiling of operation matri-

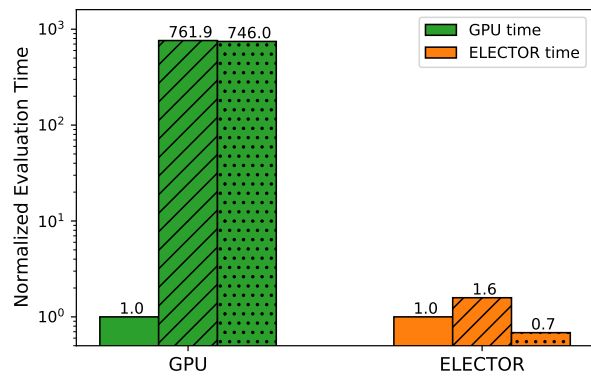


Fig. 10: Evaluation time for traditional (hatched bars) and DynaProp (dotted bars) training normalized by the inference time (solid bars) on a GPU and an ELECTOR accelerator.

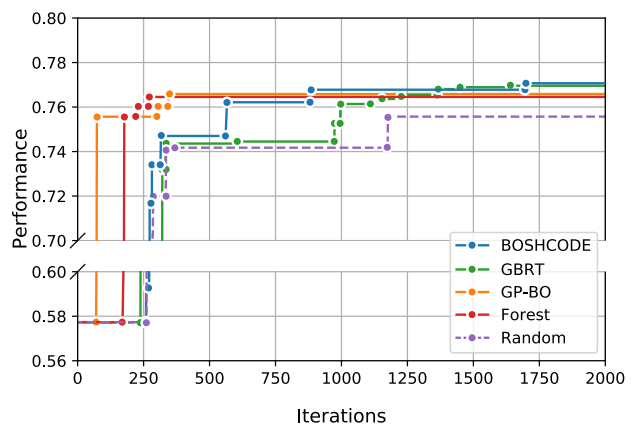


Fig. 11: Co-design convergence plots for BOSHCODE and various baselines.

ces, specialized hardware modules, and a dataflow curated for transformer workflows [16]. Since an off-the-shelf GPU does not automatically skip ineffectual computations (in other words, it is not *sparsity-aware*), DynaProp training hardly reduces evaluation time on the A100 GPU. However, due to the zero-free data format and specially designed hardware modules that skip ineffectual operations, ELECTOR reduces the training time by  $2.3\times$ . Thus, high activation, weight, and gradient sparsities enabled by DynaProp, along with ASIC-based acceleration, allow ELECTOR to substantially reduce evaluation times relative to a baseline GPU.

### B. Design Space Exploration

Fig. 11 shows convergence plots while executing co-design using BOSHCODE and various baselines. These baselines include random search, gradient-boosted regression trees (GBRT), Gaussian-process-based Bayesian optimization (GP-BO) that approximates performance through Gaussian process regression and optimizes it through the L-BFGS method [63], and random forest that fits various randomized decision trees over sub-samples of the dataset. As shown in Fig. 11, BOSHCODE achieves the highest performance. It yields the

TABLE V: Design choices of the converged TransCODE pair.

Hyperparameter		Value
<b>Transformer</b>		
Encoder Layer 1	$h^1$	256
	#SA-SDP	3
	#SA-WMA	1
	#LT-DFT	1
	#DSC-5	1
	#DSC-9	1
	#DSC-13	5
	FF	1024, 1024, 512
Encoder Layer 2	$h^2$	512
	#LT-DFT	4
	#DSC-5	5
	#DSC-9	3
	FF	256, 1024, 1024
<b>Accelerator</b>		
Batch tile size		4
Spatial tile size		32
Activation function		GeLU
#PEs		128
#MAC lanes per PE		32
#MACs per lane		16
#Softmax modules per PE		4
Batch size		4
Act./grad. buffer size (MB)		64
Weight buffer size (MB)		128
Mask buffer size (MB)		8
Main memory configuration		RRAM [8, 2, 4]

optimal transformer-accelerator pair, FB\*-ELECTOR\* (FB is an acronym for FlexiBERT 2.0). Here, performance refers to the net measure found using a convex combination of accuracy, latency, area, dynamic energy, and leakage energy (Section III-C).

Table V summarizes the design choices of the converged co-design pair, i.e., FB\*-ELECTOR\*. To optimize latency, FB\* uses only two encoder layers. However, FB\* uses 12 attention heads in each encoder layer to avoid performance loss. Thus, BOSHCODE searches for a shallow but wide model to improve throughput while not incurring a performance penalty. The converged architecture is also highly heterogeneous, with diverse attention types in each layer, leveraging the modeling capabilities of each operation type. ELECTOR\* has many PEs to parallelize the computation of 12 attention heads in each FB\* layer. It also leverages monolithic-3D RRAM, which has the highest bandwidth and lowest energy consumption. The net area of this accelerator is 359.3 mm<sup>2</sup>.

### C. Performance Improvements

We now compare the converged transformer-accelerator pairs obtained by the proposed approach with baseline pairs. Fig. 12 shows Pareto frontiers of GLUE scores with respect to hardware measures, i.e., latency, chip area, and energy consumption. We obtain GLUE scores from the surrogate model described in the EdgeTran framework [11]. We also plot state-of-the-art transformer-accelerator pairs for comparison. Our pair on the Pareto frontier with the same accuracy as

BERT-Base evaluated on AccelTran-Server incurs 44.8× lower latency. On the other hand, the pair on the Pareto frontier with the same latency as that of BERT-Tiny evaluated on AccelTran-Edge achieves a 14.5% higher GLUE score. Similarly, the pair with the same accuracy as that of BERT-Base evaluated on AccelTran-Server but on the Pareto frontier in Fig. 12(b) requires 34.5× lower chip area. The one with the same chip area as that evaluated on AccelTran-Edge finds a transformer model on the frontier that achieves a 14.8% higher GLUE score. Finally, the pair with the same accuracy as that of BERT-Base incurs 1050× lower energy consumption than that of the model evaluated on AccelTran-Server. In contrast, the *same-energy* pair with BERT-Tiny evaluated on AccelTran-Edge, but on the Pareto frontier, achieves a 13.9% higher GLUE score.

Table VI compares the proposed TransCODE approach against various baselines. These baselines include HAT [18] and AutoTinyBERT [31], which implement HW-NAS on off-the-shelf edge-AI devices. We also add a co-design method implemented on a set of FPGAs [22] and another HW-NAS approach implemented on the SpAtten ASIC-based accelerator [14]. For fair comparisons, we also include a monolithic-3D-RRAM-based transformer accelerator, i.e., AccelTran-Server [16], that evaluates BERT-Base. Finally, the table presents an ablation study in which we implement HW-NAS (by forcing the gradients to the accelerator to zero) with AccelTran-Server [16] as the base accelerator. We also include performance values for FB\*-ELECTOR\* without DynaProp training implemented. Since the baselines do not support training, we report performance values for running inference with the proposed pairs. FB\*-ELECTOR\* outperforms the state-of-the-art pair, i.e., BERT-Base/AccelTran-Server, achieving 0.3% higher accuracy, 5.2× lower latency, and 3.0× lower energy consumption.

## VI. DISCUSSIONS AND FUTURE WORK

In this section, we discuss the implications of the proposed work along with future work directions.

### A. Multi-objective Optimization

To perform co-design with the BOSHCODE framework, we model performance as a linear function of latency, energy consumption, chip area, and accuracy. This converts a multi-objective optimization problem into a single-objective optimization problem. We use this approach because BOSHCODE supports single-objective optimization only. The designer can decide the importance of each such objective when running the co-design pipeline. However, one could extend this approach to multi-objective optimization that increases/decreases a Pareto front’s hypervolume [64, 65]. In this case, the designer would obtain a set of non-dominated solutions. We leave the application of multi-objective optimization methods to the FlexiBERT 2.0 and ELECTOR design spaces to future work.

### B. In-memory and Reconfigurable Processors

The proposed framework optimizes for a specific accelerator deployed in practice for edge-based training or inference.

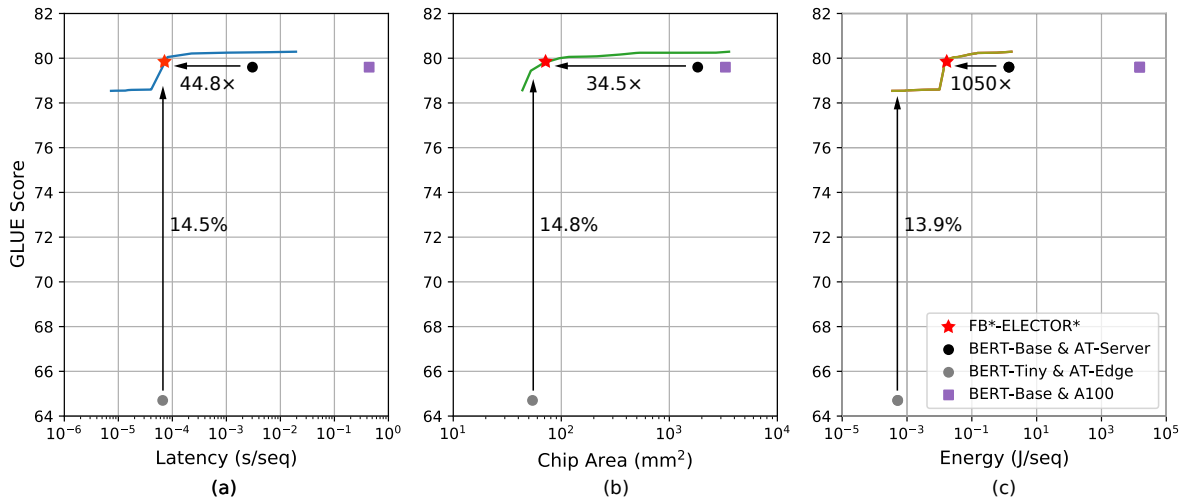


Fig. 12: Pareto frontiers of GLUE scores for models in the FlexiBERT 2.0 design space with (a) latency, (b) chip area, and (c) energy consumption on difference accelerators in the ELECTOR design space. All performance values correspond to model training. AT stands for a modified version of AccelTran that supports training.

TABLE VI: Ablation analysis and baseline comparisons for our proposed TransCODE framework.

Method	Hardware-Aware	Flex. Layers	Co-design	ASIC-based Accelerator	Platform	GLUE Score (%)	Latency (ms/seq)	Energy (J/seq)
Devlin et al. [1]	✗	✗	✗	✗	A100	79.6	10.6	0.6
<b>Baselines</b>								
Wang et al. [18]	✓	✗	✗	✗	Raspberry Pi	77.1	12,351.6	38.2
Yin et al. [31]	✗	✗	✗	✗	Raspberry Pi	78.3	10,427.7	20.7
Peng et al. [22]	✓	✗	✓	✗	FPGA	77.0	15.8	0.4
Wang et al. [14]	✓	✗	✗	✓	SpAtten	77.1	2.44	0.3
Tuli et al. [16]	✓	✗	✗	✓	AccelTran-Server	79.6	0.26	0.06
<b>Ablation Analysis</b>								
TransCODE (HW-NAS; Ours)	✓	✓	✗	✓	ELECTOR	78.4	0.23	0.08
TransCODE (without DynaProp; Ours)	✓	✓	✓	✓	ELECTOR	79.9	0.11	0.04
TransCODE (Ours)	✓	✓	✓	✓	ELECTOR	<b>79.9</b>	<b>0.05</b>	<b>0.02</b>

However, any accelerator in the proposed ELECTOR design space can execute any transformer, although it would not be the best accelerator for that transformer (in terms of hardware performance). The hardware architectures are not reconfigurable at runtime (except the pruning ratios  $\tau_I$  and  $\tau_T$ ). The architectures in the Sanger [66] design space are reconfigurable. However, Sanger is limited to only pruning a given model in the software space. Meanwhile, TransCODE leverages the FlexiBERT 2.0 design space to search for dense and small models. It also supports dynamic pruning of the model (using runtime-tunable pruning ratios) to trade off accuracy with hardware performance, while also searching for the best-performing set of accelerator design decisions. Nevertheless, adding reconfigurability to accelerators in the ELECTOR design space would benefit dynamic workloads. One could also implement co-design for a group of transformers instead of just one. We leave this to future work.

## VII. CONCLUSION

In this work, we presented TransCODE, a co-design framework for flexible and heterogeneous transformer models eval-

uated on diverse accelerator architectures. We proposed a novel, low-overhead dynamic inference-and-training scheme, DynaProp, that increases the sparsity of activations and gradients at runtime with controllable accuracy loss. DynaProp attains 90% sparsity in gradient matrices with negligible accuracy loss while improving training throughput by  $2.3\times$  relative to traditional training. We further proposed a design space of diverse ASIC-based transformer accelerators: ELECTOR. It supports accelerators targeted at various scenarios, budgets, and user-defined constraints that support flexible and heterogeneous transformer inference and training. The best transformer-accelerator pair achieves 0.3% higher accuracy than the state-of-the-art pair while enabling  $5.2\times$  lower latency and  $3.0\times$  lower energy consumption.

## ACKNOWLEDGMENTS

We performed the simulations presented in this article on computational resources managed and supported by Princeton Research Computing at Princeton University.

## REFERENCES

- [1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. Conf. North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, vol. 1, 2019, pp. 4171–4186.
- [2] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *Proc. Int. Conf. Neural Information Processing Systems*, vol. 33, 2020, pp. 1877–1901.
- [3] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, "Hierarchical text-conditional image generation with CLIP latents," *CoRR*, vol. abs/2204.06125, 2022.
- [4] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. Int. Conf. Neural Information Processing Systems*, vol. 30, 2017, pp. 5998–6008.
- [5] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *Proc. Int. Conf. Learning Representations*, 2021.
- [6] H. Kim, Y. Ohmura, and Y. Kuniyoshi, "Transformer-based deep imitation learning for dual-arm robot manipulation," in *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems*, 2021, pp. 8965–8972.
- [7] J. M. Han, J. Rute, Y. Wu, E. Ayers, and S. Polu, "Proof artifact co-training for theorem proving with language models," in *Proc. Int. Conf. Learning Representations*, 2022.
- [8] D. Zhou, Z. Liu, J. Wang, L. Wang, T. Hu, E. Ding, and J. Wang, "Human-object interaction detection via disentangled transformer," in *Proc. IEEE/CVF Conf. Computer Vision and Pattern Recognition*, 2022, pp. 19 568–19 577.
- [9] B. Lim, S. Ö. Arık, N. Loeff, and T. Pfister, "Temporal fusion transformers for interpretable multi-horizon time series forecasting," *Int. J. Forecasting*, vol. 37, no. 4, pp. 1748–1764, 2021.
- [10] Raspberry Pi 4 Model-B. [Online]. Available: <https://www.raspberrypi.com/products/raspberry-pi-4-model-b/>
- [11] S. Tuli and N. K. Jha, "EdgeTran: Co-designing transformers for efficient inference on mobile edge platforms," *CoRR*, vol. abs/2303.13745, 2023.
- [12] S. Smith, M. Patwary, B. Norick, P. LeGresley, S. Rajbhandari, J. Casper, Z. Liu, S. Prabhunoye, G. Zerveas, V. Korthikanti, E. Zheng, R. Child, R. Y. Aminabadi, J. Bernauer, X. Song, M. Shoyebi, Y. He, M. Houston, S. Tiwary, and B. Catanzaro, "Using DeepSpeed and Megatron to train Megatron-Turing NLG 530B, A large-scale generative language model," *CoRR*, vol. abs/2201.11990, 2022.
- [13] T. J. Ham, S. J. Jung, S. Kim, Y. H. Oh, Y. Park, Y. Song, J.-H. Park, S. Lee, K. Park, J. W. Lee *et al.*, "A<sup>3</sup>: Accelerating attention mechanisms in neural networks with approximation," in *Proc. Int. Symp. High-Performance Computer Architecture*, 2020, pp. 328–341.
- [14] H. Wang, Z. Zhang, and S. Han, "SpAtten: Efficient sparse attention architecture with cascade token and head pruning," in *Proc. Int. Symp. High-Performance Computer Architecture*, 2021, pp. 97–110.
- [15] Z. Zhou, J. Liu, Z. Gu, and G. Sun, "Energon: Towards efficient acceleration of transformers using dynamic sparse attention," *IEEE Trans. Computer-Aided Design of Integrated Circuits and Systems*, pp. 1–14, 2022.
- [16] S. Tuli and N. K. Jha, "AccelTran: A sparsity-aware accelerator for dynamic inference with transformers," *CoRR*, vol. abs/2302.14705, 2023.
- [17] S. Tuli, B. Dedhia, S. Tuli, and N. K. Jha, "FlexiBERT: Are current transformer architectures too homogeneous and rigid?" *CoRR*, vol. abs/2205.11656, 2022.
- [18] H. Wang, Z. Wu, Z. Liu, H. Cai, L. Zhu, C. Gan, and S. Han, "HAT: Hardware-aware transformers for efficient natural language processing," in *Proc. 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 7675–7688.
- [19] S. Tuli, C. H. Li, R. Sharma, and N. K. Jha, "CODEBench: A neural architecture and hardware accelerator co-design framework," *ACM Trans. Embedded Computing Systems*, Dec. 2022. [Online]. Available: <https://doi.org/10.1145/3575798>
- [20] Y. Lin, M. Yang, and S. Han, "NAAS: Neural accelerator architecture search," in *Proc. 58th ACM/IEEE Design Automation Conference*, 2021, pp. 1051–1056.
- [21] P. Qi, E. H.-M. Sha, Q. Zhuge, H. Peng, S. Huang, Z. Kong, Y. Song, and B. Li, "Accelerating framework of transformer by hardware design and model compression co-optimization," in *Proc. IEEE/ACM Int. Conf. Computer Aided Design*, 2021, pp. 1–9.
- [22] H. Peng, S. Huang, S. Chen, B. Li, T. Geng, A. Li, W. Jiang, W. Wen, J. Bi, H. Liu, and C. Ding, "A length adaptive algorithm-hardware co-design of transformer on FPGA through sparse attention and dynamic pipelining," in *Proc. 59th ACM/IEEE Design Automation Conference*, 2022, pp. 1135–1140.
- [23] P. Shaw, J. Uszkoreit, and A. Vaswani, "Self-attention with relative position representations," in *Proc. Conf. North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, vol. 2, 2018, pp. 464–468.
- [24] L. Zhuang, L. Wayne, S. Ya, and Z. Jun, "A robustly optimized BERT pre-training approach with post-training," in *Proc. Chinese National Conference on Computational Linguistics*, 2021, pp. 1218–1227.
- [25] Z.-H. Jiang, W. Yu, D. Zhou, Y. Chen, J. Feng, and S. Yan, "ConvBERT: Improving BERT with span-based dynamic convolution," in *Proc. Int. Conf. Neural Information Processing Systems*, vol. 33, 2020, pp. 12 837–12 848.
- [26] Z. Sun, H. Yu, X. Song, R. Liu, Y. Yang, and D. Zhou, "MobileBERT: A compact task-agnostic BERT for resource-limited devices," in *Proc. 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 2158–2170.
- [27] J. Lee-Thorp, J. Ainslie, I. Eckstein, and S. Ontanon, "FNet: Mixing tokens with Fourier transforms," in *Proc. Conf. North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2022, pp. 4296–4313.
- [28] S. Wang, B. Z. Li, M. Khabza, H. Fang, and H. Ma, "Linformer: Self-attention with linear complexity," *CoRR*, vol. abs/2006.04768, 2020.
- [29] A. Khetan and Z. Karnin, "schuBERT: Optimizing elements of BERT," in *Proc. 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 2807–2818.
- [30] L. Hou, Z. Huang, L. Shang, X. Jiang, X. Chen, and Q. Liu, "DynaBERT: Dynamic BERT with adaptive width and depth," in *Proc. Int. Conf. Neural Information Processing Systems*, vol. 33, 2020, pp. 9782–9793.
- [31] Y. Yin, C. Chen, L. Shang, X. Jiang, X. Chen, and Q. Liu, "AutoTinyBERT: Automatic hyper-parameter optimization for efficient pre-trained language models," in *Proc. 59th Annual Meeting of the Association for Computational Linguistics*, 2021, pp. 5146–5157.
- [32] J. Xu, X. Tan, R. Luo, K. Song, J. Li, T. Qin, and T.-Y. Liu, "NAS-BERT: Task-agnostic and adaptive-size BERT compression with neural architecture search," in *Proc. 27th ACM SIGKDD Conf. Knowledge Discovery & Data Mining*, 2021, pp. 1933–1943.
- [33] J. Park, H. Yoon, D. Ahn, J. Choi, and J.-J. Kim, "OPTIMUS: Optimized matrix multiplication structure for transformer neural network accelerator," in *Proc. Machine Learning and Systems*, vol. 2, 2020, pp. 363–378.
- [34] M. S. Abdelfattah, Ł. Dudziak, T. Chau, R. Lee, H. Kim, and N. D. Lane, "Best of both worlds: AutoML codesign of a CNN and its hardware accelerator," in *Proc. ACM/EDAC/IEEE Design Automation Conference*, 2020.
- [35] Apple. (2020) Apple unleashes M1. [Online]. Available: <https://www.apple.com/newsroom/2020/11/apple-unleashes-m1/>
- [36] Intel Neural Compute Stick 2. [Online]. Available: <https://www.intel.com/content/www/us/en/developer/tools/neural-compute-stick/overview.html>
- [37] NVIDIA Jetson Nano Developer Kit. [Online]. Available: <https://developer.nvidia.com/embedded/jetson-nano-developer-kit>
- [38] S. Tuli, S. R. Poojara, S. N. Srirama, G. Casale, and N. R. Jennings, "COSCO: Container orchestration using co-simulation and gradient based optimization for fog computing environments," *IEEE Trans. Parallel and Distributed Systems*, vol. 33, no. 1, pp. 101–116, 2021.
- [39] P. Ren, Y. Xiao, X. Chang, P.-Y. Huang, Z. Li, B. B. Gupta, X. Chen, and X. Wang, "A survey of deep active learning," *ACM Comput. Surv.*, vol. 54, no. 9, 2021.
- [40] Y. Yu and N. K. Jha, "SPRING: A sparsity-aware reduced-precision monolithic 3D CNN accelerator architecture for training and inference," *IEEE Trans. Emerging Topics in Computing*, vol. 10, no. 1, pp. 237–249, 2022.
- [41] P. Batude, B. Sklenard, C. Fenouillet-Beranger, B. Previtali, C. Tabone, O. Rozeau, O. Billoint, O. Turkyilmaz, H. Sarhan, S. Thuries, G. Cibrario, L. Brunet, F. Deprat, J.-E. Michallet, F. Clermidy, and M. Vinet, "3D sequential integration opportunities and technology optimization," in *Proc. Int. Interconnect Technology Conference*, 2014, pp. 373–376.

- [42] D. Hendrycks and K. Gimpel, "Bridging nonlinearities and stochastic regularizers with Gaussian error linear units," *CoRR*, vol. abs/1606.08415, 2016.
- [43] T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," in *Proc. Conf. Empirical Methods in Natural Language Processing*, 2015, pp. 1412–1421.
- [44] K. R. Rao and P. C. Yip, *The Transform and Data Compression Handbook*. CRC press, 2018.
- [45] J. W. Cooley and J. W. Tukey, "An algorithm for the machine calculation of complex Fourier series," *Mathematics of Computation*, vol. 19, no. 90, pp. 297–301, 1965.
- [46] J. Makhoul, "A fast cosine transform in one and two dimensions," *IEEE Trans. Acoustics, Speech, and Signal Processing*, vol. 28, no. 1, pp. 27–34, 1980.
- [47] I. Beltagy, M. E. Peters, and A. Cohan, "Longformer: The long-document transformer," *CoRR*, vol. abs/2004.05150, 2020.
- [48] Synopsys Design Compiler (2022). [Online]. Available: <https://www.synopsys.com/implementation-and-signoff/rtl-synthesis-test/dc-ultra.html>
- [49] A. Guler and N. K. Jha, "Hybrid monolithic 3-D IC floorplanner," *IEEE Trans. Very Large Scale Integration Systems*, vol. 26, no. 10, pp. 1868–1880, 2018.
- [50] J. A. Roy, D. A. Papa, S. N. Adya, H. H. Chan, A. N. Ng, J. F. Lu, and I. L. Markov, "Capo: Robust and scalable open-source min-cut floorplanner," in *Proc. Int. Symp. Physical Design*, 2005, pp. 224–226.
- [51] A. Shafaei, Y. Wang, X. Lin, and M. Pedram, "FinCACTI: Architectural analysis and modeling of caches with deeply-scaled FinFET devices," in *Proc. Computer Society Annual Symp. VLSI*, 2014, pp. 290–295.
- [52] X. Dong, C. Xu, Y. Xie, and N. P. Jouppi, "NVSim: A circuit-level performance, energy, and area model for emerging nonvolatile memory," *IEEE Trans. Computer-Aided Design of Integrated Circuits and Systems*, vol. 31, no. 7, pp. 994–1007, 2012.
- [53] M. Poremba, T. Zhang, and Y. Xie, "NVMain 2.0: A user-friendly memory simulator to model (non-)volatile memory systems," *IEEE Computer Architecture Letters*, vol. 14, no. 2, pp. 140–143, 2015.
- [54] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman, "GLUE: A multi-task benchmark and analysis platform for natural language understanding," in *Proc. EMNLP Workshop on BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, 2018, pp. 353–355.
- [55] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, and C. Potts, "Recursive deep models for semantic compositionality over a sentiment treebank," in *Proc. Int. Conf. Empirical Methods in Natural Language Processing*, 2013, pp. 1631–1642.
- [56] A. Williams, N. Nangia, and S. Bowman, "A broad-coverage challenge corpus for sentence understanding through inference," in *Proc. Conf. North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2018, pp. 1112–1122.
- [57] W. B. Dolan and C. Brockett, "Automatically constructing a corpus of sentential paraphrases," in *Proc. Int. Workshop on Paraphrasing*, 2005.
- [58] A. Warstadt, A. Singh, and S. R. Bowman, "Neural network acceptability judgments," *Trans. Association for Computational Linguistics*, vol. 7, pp. 625–641, 2019.
- [59] D. Cer, M. Diab, E. Agirre, I. Lopez-Gazpio, and L. Specia, "SemEval-2017 task 1: Semantic textual similarity - multilingual and cross-lingual focused evaluation," in *Proc. Int. Workshop on Semantic Evaluation*, 2017, pp. 1–14.
- [60] I. Dagan, O. Glickman, and B. Magnini, "The PASCAL recognising textual entailment challenge," in *Proc. Int. Conf. Machine Learning Challenges*, 2006, pp. 177–190.
- [61] H. Levesque, E. Davis, and L. Morgenstern, "The Winograd schema challenge," in *Proc. Int. Conf. Principles of Knowledge Representation and Reasoning*, 2012.
- [62] V. Sanh, T. Wolf, and A. Rush, "Movement pruning: Adaptive sparsity by fine-tuning," in *Proc. Int. Conf. Neural Information Processing Systems*, vol. 33, 2020, pp. 20378–20389.
- [63] D. C. Liu and J. Nocedal, "On the limited memory BFGS method for large scale optimization," *Mathematical Programming*, vol. 45, no. 1, pp. 503–528, 1989.
- [64] X. Yang, S. Belakaria, B. K. Joardar, H. Yang, J. R. Doppa, P. P. Pande, K. Chakrabarty, and H. H. Li, "Multi-objective optimization of ReRAM crossbars for robust DNN inferencing under stochastic noise," in *Proc. IEEE/ACM Int. Conf. Computer Aided Design*, 2021, pp. 1–9.
- [65] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, "A fast and elitist multiobjective genetic algorithm: NSGA-II," *IEEE Trans. Evolutionary Computation*, vol. 6, no. 2, pp. 182–197, 2002.
- [66] L. Lu, Y. Jin, H. Bi, Z. Luo, P. Li, T. Wang, and Y. Liang, "Sanger: A

co-design framework for enabling sparse attention using reconfigurable architecture," in *Proc. IEEE/ACM Int. Symp. Microarchitecture*, 2021, pp. 977–991.



**Shikhar Tuli** received the B. Tech. degree in electrical and electronics engineering from the Indian Institute of Technology (IIT) Delhi, India, with a department specialization in very large-scale integration (VLSI) and embedded systems. He is currently pursuing a Ph.D. degree at Princeton University in the department of electrical and computer engineering. His research interests include deep learning, edge artificial intelligence (AI), hardware-software co-design, brain-inspired computing, and smart healthcare.



**Niraj K. Jha** (Fellow, IEEE) received the B.Tech. degree in electronics and electrical communication engineering from IIT, Kharagpur, India, in 1981, and the Ph.D. degree in electrical engineering from the University of Illinois at Urbana-Champaign, Champaign, IL, USA, in 1985. He is a professor of electrical and computer engineering, Princeton University. He has co-authored five widely used books. He has published more than 470 papers (h-index: 83). He has received the Princeton Graduate Mentoring Award. His research has won 15 best paper awards, six award nominations, and 25 patents. He was given the Distinguished Alumnus Award by IIT, Kharagpur, in 2014. He has served as the Editor-in-Chief of TVLSI and an associate editor of several IEEE Transactions and other journals. He has given several keynote speeches in the areas of nanoelectronic design/test, smart healthcare, and cybersecurity. He is a fellow of ACM. His research interests include machine learning algorithms/architectures and smart healthcare.