# Vision Transformer Computation and Resilience for Dynamic Inference

Kavya Sreedhar[1]
*Stanford University*
Stanford, USA
skavya@stanford.edu

Jason Clemons
*NVIDIA*
Austin, USA
jclemons@nvidia.com

Rangharajan Venkatesan
*NVIDIA*
Santa Clara, USA
rangharajanv@nvidia.com

Stephen W. Keckler
*NVIDIA*
Austin, USA
skeckler@nvidia.com

Mark Horowitz
*Stanford University*
Stanford, USA
horowitz@ee.stanford.edu

*Abstract*—State-of-the-art deep learning models for computer vision tasks are based on the transformer architecture and often deployed in real-time applications. In this scenario, the resources available for every inference can vary, so it is useful to be able to dynamically adapt execution to trade accuracy for efficiency. To create dynamic models, we leverage the resilience of vision transformers to pruning and switch between different scaled versions of a model. Surprisingly, we find that most FLOPs are generated by convolutions, not attention. These relative FLOP counts are not a good predictor of GPU performance since GPUs have special optimizations for convolutions. Some models are fairly resilient and their model execution can be adapted without retraining, while all models achieve better accuracy with retraining alternative execution paths. These insights mean that we can leverage CNN accelerators and these alternative execution paths to enable efficient and dynamic vision transformer inference. Our analysis shows that leveraging this type of dynamic execution can lead to saving 28% of energy with a 1.4% accuracy drop for SegFormer (63 GFLOPs), with no additional training, and 53% of energy for ResNet-50 (4 GFLOPs) with a 3.3% accuracy drop by switching between pretrained Once-For-All models.

*Index Terms*—vision transformers, dynamic inference, model resilience, pruning

## I. INTRODUCTION

Deep learning models in computer vision have shifted from convolutional neural networks (CNNs) [1]–[12] to transformers [13]–[33] for higher model accuracy. The transformer architecture [34] uses attention to understand global image contexts and effectively capture spatial information. It also underlies general-purpose backbones and foundation models for language [35]–[40] and vision [18]–[22] tasks. For example, in 2021, Microsoft introduced Swin Transformer [18], which has since been adopted as a backbone for various vision tasks [31]–[33], while Meta recently released the Segment Anything Model [22] as a foundation model for segmentation tasks.

These models can be computationally expensive, requiring millions of parameters and billions of floating point operations (FLOPs) [41]. It is well-known that increasing the number of parameters and FLOPs in a model can lead to better model accuracy, as seen by recent trends in larger and more accurate GPT models, as an example [19], [36]–[39]. Furthermore, these models typically have a fixed execution path and assume
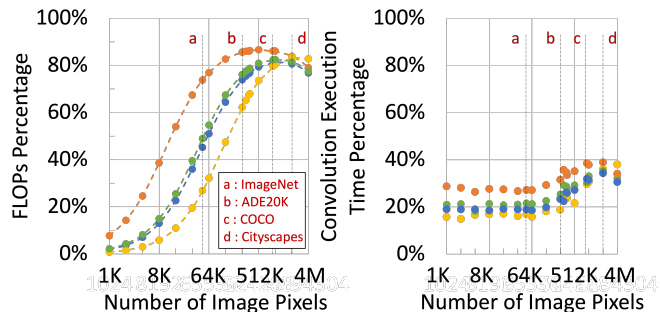
Fig. 1. FLOPs and NVIDIA RTX A5000 GPU execution time in convolutions (dots) and ResNet-50 backbone (dashed lines) for inference with DETR [13] (orange), Conditional DETR [16] (green), DAB DETR [14] (blue), and Anchor DETR [15] (yellow). For larger image sizes, convolutions dominate FLOPs but not GPU execution time.

that their needed computational resources will be available when they are run.

In contrast, real-time systems for applications such as autonomous driving [42], [43] and video conferencing [44] have limited hardware resources and dynamic system loads that change as the surrounding environment changes [45]–[48]. Choosing a static model that matches the worst-case resource utilization would leave performance on the table when more resources are available in this scenario. As a result, these systems need to leverage dynamic vision transformers.

Most prior work on dynamic models shortens the model execution based on complexity of classifying the input. These approaches exit early and remove the computation of later layers when internal predictions have already become stable [48]–[60]. While this prior work reduces average latency and energy for "easier" inputs, it does not ensure that model execution meets a given dynamic resource constraint.

Some work has addressed this challenge by trading accuracy for efficiency in order to adapt the cost of inference to not exceed an input resource constraint [47], [48], but this prior work focuses on CNNs and BERT. In this paper, we extend that work for vision transformers. We build upon work that scales static model architectures and prunes redundant computation to achieve different levels of model computation and accuracy [13], [61]–[64]. For example, Once-For-All (OFA) has developed efficient training techniques that produce many competitive subnetworks after training one model [64].
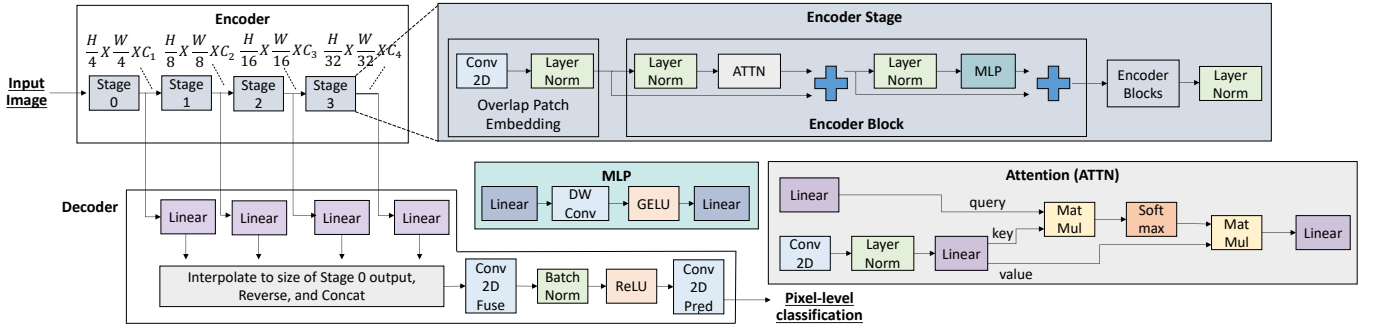
Fig. 2. Layers in SegFormer [17] model. The Swin Transformer [18] model follows the same high-level structure, with a more optimized attention module and the UPerNet decoder head [65]. The UPerNet decoder has a layer similar to SegFormer's *Conv2DFuse*, which we refer to as *fpn_bottleneck_Conv2D*.

We examine resource-dependent dynamic (RDD) inference [47], [48] for vision transformers by profiling this class of applications and identifying lower-cost execution paths. When examining the full application pipeline for state-of-the-art vision transformers [13]–[18], we find that FLOPs are dominated by convolutions, not attention. This result is due to two reasons: first, transformer models for detection often use CNN backbones to extract input features and second, many models have incorporated convolutions in the transformer encoder-decoder structure to achieve higher model accuracy and enable lower-cost attention. For complex vision applications beyond classification, it is crucial to analyze the substantial computation required for backbones and task-specific decoder heads.

Furthermore, we find that the distribution of FLOPs across model layers is not a good estimator of relative GPU execution time since GPU hardware and software have been well-optimized for exploiting the inherent locality and parallelism in convolutions [66]–[68]. This efficiency of convolutional operations is one reason for their continued use. Figure 1 shows these trends in FLOPs and GPU performance for detection models [13]–[16] using a ResNet-50 [9] backbone.

Given the importance of efficiently executing attention *and* convolutions for this class of models, we use MAGNet [69], an accelerator framework previously demonstrated for CNNs and attention-dominated transformers [70], to estimate the performance of these applications with customized hardware.

We then explore the impact of pruning in convolutional and attention layers in pretrained models on the model accuracy, execution time, and energy on an NVIDIA RTX A5000 GPU and the MAGNet-generated accelerator. Unintuitively, we find that larger models are not necessarily more resilient to accuracy loss when bypassing computation. Instead, how computation is split between the transformer encoder and decoder is a better indicator.

Finally, we augment vision transformers with lower-cost execution paths that either leverage the resilience of pretrained models to pruning or rely on switching between retrained models for RDD inference. These paths achieve a range of execution time and energy savings, and can be alternatively executed to enable meeting dynamic resource constraints.

We make the following contributions:

- We show that convolutions, not attention layers, dominate FLOPs for state-of-the-art vision transformer applications, since transformer models have integrated convolutions for accuracy and performance.
- We find that for these applications, the distribution of FLOPs across model layers is not a good estimator of relative GPU runtime.
- We identify alternative lower-cost execution paths in these models and identify indicators for the resilience of pretrained vision transformers to dynamic pruning.
- We leverage a CNN accelerator framework and our dynamic computation bypassing approach to save 28% of energy with a 1.4% accuracy loss for SegFormer B2 with no additional training and 53% of energy for RestNet-50 with a 3.3% accuracy loss by switching between pretrained OFA models.

## II. BACKGROUND

### A. RDD inference

RDD inference targets a computational platform with a finite amount of resources [47], [48]. Occasionally, there are not enough resources available for the full model execution, since there are other tasks that also need to be completed. In these scenarios, it is better to perform a degraded version of inference that requires less resources rather than to skip a frame and perform no inference. While statically selecting a model that matches the worst-case resource availability will ensure that no frames are missed, it does not allow for performing the full model execution and returning more accurate results when enough resources are available. Thus, the goal of RDD inference is to adjust the cost of inference to match the resources available for the computation, and thus maximize accuracy given input resource constraints.

### B. Vision Tasks and Models

We consider semantic segmentation [71] and object detection [72] since these tasks are widely used [42]–[44], [73] and more complicated than image classification. Semantic segmentation assigns class labels to each pixel in an input

| Model | Task | Parameters (Millions) | Dataset | Input Image Size | GFLOPs | mIoU (SS) / AP (OD) |
|---|---|---|---|---|---|---|
| SegFormer ADE B2 [17] | SS | 28 | ADE20K [74] | 512 by 512 | 63 | 0.4651 |
| SegFormer City B2 [17] | SS | 28 | Cityscapes [75] | 1024 by 1024 | 290 | 0.8098 |
| Swin Tiny [18] | SS | 60 | ADE20K | 512 by 512 | 237 | 0.4451 |
| Swin Small [18] | SS | 81 | ADE20K | 512 by 512 | 259 | 0.4764 |
| Swin Base [18] | SS | 121 | ADE20K | 512 by 512 | 297 | 0.4813 |
| DETR [13] | OD | 41 | COCO-2017 [76] | 800 by 1200 | 92 | 0.4200 |
| DAB DETR [14] | OD | 44 | COCO-2017 | 800 by 1200 | 97 | 0.328 |
| Anchor DETR [15] | OD | 37 | COCO-2017 | 800 by 1200 | 99 | 0.4188 |
| Conditional DETR [16] | OD | 43 | COCO-2017 | 800 by 1200 | 96 | 0.4161 |

TABLE I

OUR STATE-OF-THE-ART VISION TRANSFORMER CASE STUDIES FOR SEMANTIC SEGMENTATION (SS) AND OBJECT DETECTION (OD).

image, while object detection identifies objects by providing bounding boxes and class labels.

State-of-the-art models for these tasks have shifted to using transformer-based architectures, building on the success of using transformers for natural language processing [34], [35] and image classification [20], [21], [24]. Most segmentation models and some detection models use the transformer encoder as a backbone for feature extraction, pairing it with different decoders and task-specific heads to produce the final output [7], [18], [31]–[33], [77]. Other detection models models employ a transformer encoder-decoder structure, often with a CNN backbone to extract visual features [13]–[16], [25], [31]–[33]. Since we analyze the first type of models with our segmentation model case studies, we focus on the second type of models for our detection model case studies.

Following from ViT, the first vision transformer, models reshape 2D images into a 1D sequence of flattened image patches, which are linearly projected into an embedding dimension [20]. Throughout the model, layers progressively increase the embedding dimension and reduce the spatial dimensions to capture details at various image resolutions for higher accuracy. ViT is convolution-free since the original transformer layers consist of multi-head self-attention layers and multi-layer perceptrons (MLPs) [34]. However, many modern vision transformers have incorporated convolutions with the transformer architecture for better accuracy results [7], [13]–[18], [22], [25], [78].

We use NVIDIA's SegFormer [17] and Microsoft's Swin Transformer [18] as case studies for semantic segmentation since they achieve state-of-the-art accuracy and are commonly used. Figure 2 shows these model architectures. These models build from the encoder in ViT, and further optimize attention and integrate convolutions to preserve local continuity information. Both models have four encoder stages followed by a decoder. SegFormer B2 has three, four, six, and three encoder blocks in stages zero through three, respectively. Swin Tiny has six encoder blocks in stage two and two encoder blocks in all the other encoder stages. Individual layers inside these blocks are shown in Figure 2, with breakout boxes showing the individual layers in the MLP and attention components.

Prior work has focused on improving encoder backbones [18], [77] that can be used with different decoders and task-specific heads for various vision tasks. In these backbones, attention dominates FLOPs [17], [18], [20], [21], [24], so prior work has heavily focused on reducing the cost of attention [65], [79]–[81]. SegFormer contributes a simpler attention-free decoder, while Swin directly uses the UPerNet decoder head [65]. Both of these decoders use convolutions to extract and fuse features from different spatial resolutions and reduce channel dimensions, which results in higher accuracy.

For object detection, many models build from the DETR architecture [13]. DETR uses a CNN backbone to extract image features, a conventional transformer, and a feed forward network to return the detection prediction. We examine DETR and three recent models based on DETR: DAB DETR [14], Anchor DETR [15], and Conditional DETR [16]. These models improve aspects of the transformer in DETR. We consider the base variants of these models, as included in detrex, an open-source toolbox for detection models [82]. The base versions of these models use ResNet-50 [9] as the CNN backbone.

Prior work has also explored the resiliency of ResNet-50 [61], [64]. One example is OFA [64], which trains the model once to produce many sets of weights for different model subsets, resulting in models that significantly reduce FLOPs with modest decreases in accuracy. In Section V-C, we leverage switching between OFA ResNet-50 models to enable RDD inference. We denote the most accurate OFA ResNet-50 model as OFA-ResNet-50.

Finally, we use standard accuracy metrics in our analysis: mean intersection over union (mIoU) for semantic segmentation and average precision (AP) for object detection. IoU is the area of overlap between the predicted segmentation and the ground truth divided by their total area. Thus, mIoU is the average of the IoU for every class. AP measures the precision-recall curve across different IoU confidence thresholds, which specify the minimum IoU to consider a positive match. For the COCO dataset [76], AP is the average precision for IoU from 0.5 to 0.95, in increments of 0.05. Both mIoU and AP range from zero to one, and higher values indicate better accuracy.

## III. COMPUTATION IN VISION TRANSFORMERS

Table I provides an overview of our model case studies. To understand the computation required for inference with modern vision transformer applications, we profile the FLOPs
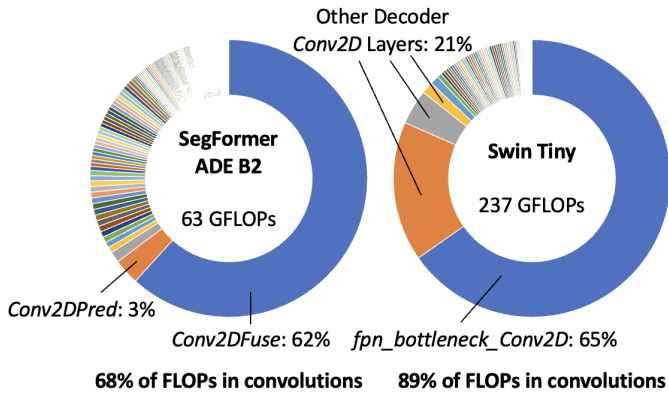
Fig. 3. FLOPs distribution across SegFormer ADE B2 model layers and Swin Tiny model layers for inference with a 512 by 512 input image size.
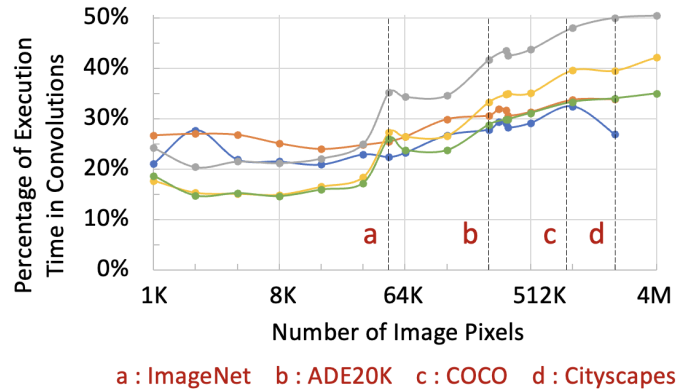


a : ImageNet  b : ADE20K  c : COCO  d : Cityscapes

Fig. 4. Image pixels versus NVIDIA RTX A5000 GPU execution time spent on convolutions for inference with the SegFormer ADE B2 (blue), SegFormer City B2 (orange), Swin Tiny (gray), Swin Small (yellow), and Swin Base models (green).

and GPU execution time distribution among model layers. We use an NVIDIA RTX A5000 GPU.

Unlike prior work that focuses on classification and a 224 by 224 image size (49K pixels) [83], our evaluation is across various image sizes (1K to 4M pixels), including those used in standard datasets for classification (ImageNet [83]), segmentation (ADE20K [74], Cityscapes [75]), and detection (COCO [76], Cityscapes [75]). These results become important for emerging applications that call for larger image resolutions [71], [72].

### A. FLOPs in Semantic Segmentation Models

We show the distribution of FLOPs across model layers for SegFormer B2 and Swin-Tiny with a 512 by 512 image size in Figure 3. We find that 68% and 89% of the total FLOPs are in convolution layers in SegFormer B2 and Swin-Tiny, respectively, in stark contrast to the zero convolutions in ViT [20] and BERT [35]. Relatively large convolutional layers are labeled in the FLOPs distributions in Figure 3. While encoders use convolutions to overlap image patches to incorporate local continuity information and reduce the size of inputs for attention, only 5% of convolution FLOPs are in the SegFormer encoder and 1% are in Swin encoder. Instead, as shown in Figure 2, decoders replace attention with large convolutions to fuse information at different image resolutions from the output of each encoder stage. We denote these convolutions as *Conv2DFuse* in SegFormer and *fpn_bottleneck_conv* in Swin.

Looking at SegFormer's decoder, where nearly 70% of the FLOPs are used, the *Conv2DFuse*, *Conv2DPred*, and *Decoder Linear connected to output of encoder stage zero* (denoted as *DecodeLinear0*) layers individually comprise 62%, 3%, and 1.3% of the total number of FLOPs for a 512 by 512 image, respectively, as shown in Figure 3. In particular, *Conv2D Fuse* alone constitutes the majority of FLOPs, considerably more so than any other layer in the model. The *Conv2D Fuse* layer has 3072 input channels, 768 output channels, and a 1 by 1 stencil.

When considering the Swin Tiny model for segmentation, we find that 89% of FLOPs are in the decoder, and 99% of

convolutions are in this decoder. The *fpn_bottleneck_Conv2D* layer alone comprises 65% of FLOPs. This layer has 2048 input channels, 512 output channels, and a 3 by 3 stencil.

From these results, it follows that these segmentation models have high operational intensity, at 130+ operations/byte, and that they are not memory bound like traditional transformers for language [35] and classification [20], [21] tasks. We also see that the task-specific decoder heads are computationally dominant for transformer models for segmentation, and critical to include in profiling analyses for this task.

### B. FLOPs in Object Detection Models

Figure 1 shows the FLOPs distribution between convolution and non-convolution layers as well as the FLOPs distribution between the ResNet-50 backbone and transformer for DETR-based models across various image sizes. Note that the percentage of FLOPs in convolutions and in ResNet-50 are very similar, indicating that the convolutions in DETR-based models are predominantly in the ResNet-50 backbone (and that the contribution of non-convolutional layers to the ResNet-50 backbone FLOPs is negligible as expected). Solely focusing on smaller image sizes below 64K pixels would lead one to conclude that improving the transformer is important for overall model efficiency. However, for object detection, we typically need to consider much larger image sizes. The importance of the ResNet-50 backbone compared to the transformer mostly increases with larger image sizes: for more than 128K image pixels, the ResNet-50 backbone comprises at least half of the total FLOPs in all of these models, and for more than 1M pixels, which corresponds to standard image sizes for detection datasets [75], [76], this backbone requires 80+% of FLOPs.

### C. Using FLOPs to Estimate GPU Runtime

We observe that the number of FLOPs in all of these models does not directly map to GPU execution time, as shown in Figure 4 for the segmentation models and Figure 1 for the detection models. For 512 by 512 image sizes, convolutions

comprise 28% and 42% of the total execution time for Seg-Former B2 and Swin Tiny, respectively, despite requiring 68% and 89% of the FLOPs. With the DETR-based models, 30% to 40% of the total execution time for larger image sizes is in convolutions, which make up 80+% of the total FLOPs. For all models, the importance of convolutions generally increases with larger image sizes, but the proportion of GPU runtime spent in convolutions is much lower than the relative FLOP counts. These results show that the GPU implementation is well-designed for reusing convolution weights and exploiting convolution parallelism. Thus, integrating convolutions results in better model accuracy *and* performance.

Figure 4 shows that convolutions comprise a smaller portion of FLOPs and GPU performance in larger Swin models compared to Swin Tiny. This result makes sense because all models have a similar decoder, which majorly comprises of convolutions, but Swin Small and Swin Base have two times the number of encoder blocks compared to Swin Tiny. We also observe that matrix multiplications comprise about an equal proportion of execution time compared to convolutions on the GPU for SegFormer and Swin for larger image sizes commonly used for segmentation.

Thus, while many prior papers use FLOPs to estimate the runtime savings for inference [49]–[54], [60], we find that FLOPs alone is not a good predictor of runtime on a GPU. GPU runtime is instead a function of FLOPs, layer parallelism, and layer operational intensity. Convolutions enable spatial parallelism since each output pixel depends on a relatively small receptive field in the input image. As a result, these layers can be efficiently executed on the highly parallel architecture in a GPU. Attention can be computed in parallel for several image patches, enabling temporal parallelism, and matrix multiplication can be easily parallelized. However, attention mechanisms also involve sequential operations where these patches attend to each other, making the overall operation less parallelizable. In addition, attention layers in these vision transformers have lower operational intensity compared to convolutions, resulting in more memory accesses. Thus, it is important for inference engines for vision transformers to efficiently execute CNNs and attention.

## IV. PROFILING ON A HARDWARE ACCELERATOR

We next analyze the performance and energy of these applications on the type of compute that one is likely to see in embedded platforms and real-time systems for RDD inference. For applications, we use SegFormer B2 and Swin Tiny for segmentation with a 512 by 512 input image size. Given the reliance of DETR-based models on the ResNet-50 backbone from Section III, we include OFA-ResNet-50 with a 224 by 224 input image size. These applications have a range of model sizes (28 to 60 million parameters), computation requirements (4 to 237 GLOPs) and model structures (68% to 95+% of FLOPs in convolutions). With these applications, we show how to tune CNN accelerators for modern vision transformers by exploring throughput, area, and energy tradeoffs with different compute organization and memory sizes. This profiling lays
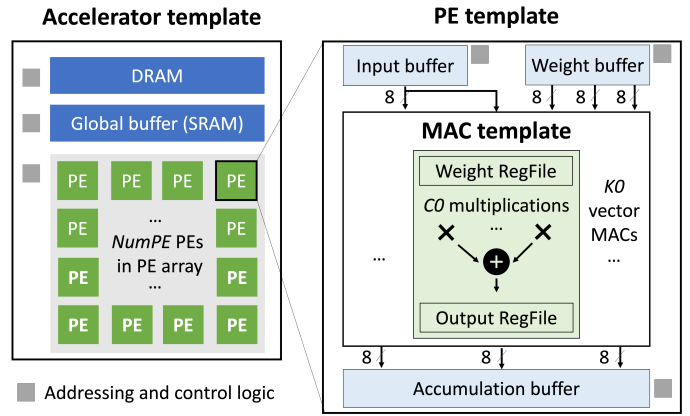


Fig. 5. Parameterizable MAGNet [69] accelerator and PE architecture templates.

the groundwork for executing dynamic vision transformers on hardware accelerators in Section V.

### A. Hardware Overview

We use MAGNet [69], a CNN accelerator framework, which has been previously extended to execute attention-dominated transformers for language and classification tasks as well [70]. As a result, its parameterizable accelerator template provides a reasonable architecture for profiling our applications on a hardware accelerator. Most other accelerators for transformers focus on attention and custom detection components [79]–[81], [84], often without support for convolutions, and are thus not well-optimized for more recent vision transformers.

Figure 5 shows a block diagram of the MAGNet accelerator and PE architecture templates. The accelerator has three levels of compute with vector multiply-accumulate (MAC) units, processing elements (PEs), and the PE array. All levels allow for parallel execution in order to exploit data parallelism present in model layers: each vector MAC unit can execute $C0$ multiplications in parallel, every PE has $K0$ parallel vector MAC units, and the PE array has $NumPE$ PEs. Each PE also contains a post-processing unit to fuse activation and pooling layers with the preceding convolutions.

The memory hierarchy has four levels: small register files within each vector MAC unit (labeled as "RegFile" in Figure 5), local buffers within each PE, a global buffer at the PE array level, and an off-chip DRAM. Inputs are shared across the vector MAC units with one input buffer per PE. Each vector MAC unit has its own weight buffer, with $K0$ total weight buffers in every PE. We can temporally tile the weights (split by output channels) and activations (split by image height and width). The paper on the MAGNet framework provides further detail on how convolutions and matrix multiplications can be mapped to this accelerator template [69].

MAGNet uses Mentor Graphics Catapult HLS to generate RTL from a synthesizable SystemC and C++ architecture description. We report results in a 5nm technology, using Synopsys Design Compiler for synthesis for execution time

| Label | NumPE | K0 = C0 | Weight buffer size (kB) | Input buffer size (kB) | PE array area (mm$^2$) |
|---|---|---|---|---|---|
| A | 32 | 32 | 1024 | 64 | 16.7 |
| B | 32 | 32 | 128 | 64 | 4.5 |
| C | 16 | 32 | 1024 | 64 | 8.3 |
| D | 16 | 32 | 128 | 64 | 2.3 |
| E | 16 | 32 | 128 | 32 | 1.9 |
| F | 16 | 32 | 64 | 64 | 2.0 |
| G | 16 | 32 | 64 | 32 | 1.7 |
| H | 64 | 16 | 128 | 32 | 6.1 |
| I | 64 | 16 | 128 | 16 | 5.4 |
| J | 64 | 16 | 64 | 32 | 4.2 |
| K | 64 | 16 | 64 | 16 | 3.5 |
| L | 64 | 16 | 32 | 32 | 3.3 |
| M | 64 | 16 | 32 | 16 | 2.6 |

TABLE II
THE MAGNET ACCELERATOR PARAMETERIZATIONS WE EXPLORE.



Fig. 7. Execution time and energy distributions across layers in the SegFormer ADE B2 model for inference with a 512 by 512 input image size on accelerator E.



Fig. 8. Normalized energy per FLOP for SegFormer ADE B2 model layers for inference with a 512 by 512 input image size on accelerator E. The layer with highest energy per FLOP is set to one, with all layers sorted in order of decreasing energy per FLOP from left to right.
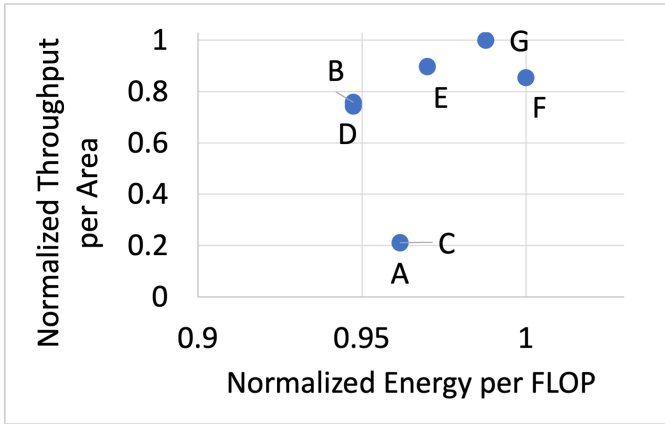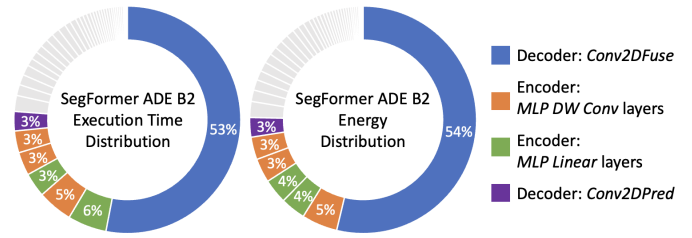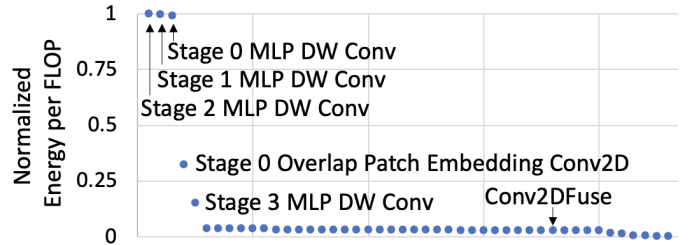


Fig. 6. Energy per FLOP versus throughput per area for SegFormer ADE B2 inference with a 512 by 512 input image size. Labels indicate the hardware accelerators in Table II.

and area results and Primetime-PX with Gaussian synthetic data for power results.

Since prior work has shown that dataflow choice minimally affects energy efficiency [85], we directly use an output-stationary local-weight-stationary dataflow that was previously shown to work well for CNNs [69] and attention layers [70]. Like these prior works, we also use 8-bit precision for data.

### B. Tradeoffs with Different Accelerator Parameters

To examine throughput, area, and energy tradeoffs for vision transformer applications, we compare different accelerator template parameters, as shown in Table II. Accelerators C through M can compute the same number of MACs in parallel, while accelerators A and B have twice the compute capability. To fairly compare the execution of this class of models on these parallel, throughput-focused compute engines, we compute the energy per operation, which to first order is independent of parallelism, and throughput/mm$^2$, to normalize throughput by the silicon cost of its hardware. As a result, accelerator B has twice the available compute compared to accelerator D and is twice as fast, but it looks identical to accelerator D with our metrics since it is also twice as big.

Figure 6 explores the tradeoff between these metrics for SegFormer ADE B2, where lower energy per operation and higher throughput per area signify more efficient inference. We also checked Swin Tiny and OFA-ResNet-50, which have similar results compared to SegFormer. All of these models can significantly reuse data across the channel dimension. As a result, accelerators with reduced vectorization ($K0 = C0 = 16$) result in lower efficiency ($1.4\times$ more energy per FLOP and $2.8\times$ more area per FLOP) compared to those with $K0 = C0 = 32$. Since the accelerators with $K0 = C0 = 16$ are not competitive with either metric, they are not shown in Figure 6.

Swin Tiny has significantly more model layers with an odd number of input and/or output channels (such as 49), due to default layer parameters or how matrix multiplications are mapped to the hardware. These channel counts are not divisible by $K0$ and $C0$ whether they are equal to 16 or 32, resulting in similar underutilization in the PEs for all the accelerator parameterizations, and similar performance across accelerators. In contrast, SegFormer and OFA-ResNet-50 have more evenly divisible channel values, resulting in slightly higher utilization and 10% faster execution with accelerators with $K0 = C0 = 16$ compared to those with $K0 = C0 = 32$.

For all of these models, accelerators D, E, and G are Pareto-optimal with our metrics. We expect this trend to continue for larger models, including CNN backbones such as ResNet-101 for DETR-based models. Note that accelerators E and G are the two smallest accelerators in Table II as well. Not surprisingly, these accelerators are similar to accelerator D, which was the design previously used for transformer models dominated by attention [70].

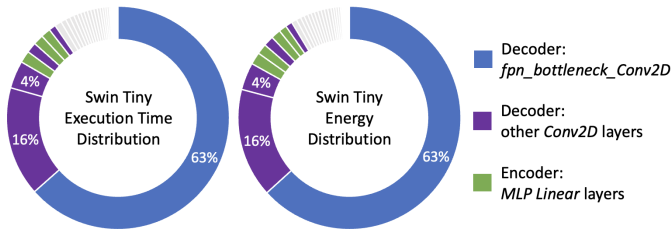In Section III, we saw that these vision transformers have

Fig. 9. Execution time and energy distributions across layers in the Swin Tiny model for inference with a 512 by 512 input image size on accelerator E.

higher operational intensity than attention-dominated transformers, due to the increased locality of the convolution operation. Thus, it makes sense that the smaller memories in accelerators E and G result in higher throughput per area compared to accelerator D. We find that 64 to 128 B/MAC and 32 to 64 B/MAC are sweet spots for the weight and input buffers, respectively, to balance throughput and area tradeoffs for these applications.

### C. Performance Distributions

To balance energy and area constraints, we profile our model case studies on accelerator E, which is between accelerators D and G on the Pareto curve in Figure 6. The synthesized clock frequency for accelerator E is 1.25 GHz.

The SegFormer ADE B2 model runs in 3.6ms on accelerator E. Figure 7 shows the execution time and energy breakdowns across model layers for inference on accelerator E. Convolutions constitute 74% of the total execution time and energy. The largest layer, *Conv2DFuse*, alone requires over half of the total execution time and energy on accelerator E.

Five convolution layers in the encoder have higher energy per FLOP compared to other layers, as indicated in Figure 8. Fortunately, these layers comprise only 14% of the total energy. These layers have a small number of input channels: the *Overlap Patch Embedding Conv2D* layer has three input channels for an input image to the model, and the *MLP DW Conv* layers each have one input channel due to how we exploit parallelism in mappings for depthwise convolutions. This results in hardware underutilization, since all 32 parallel multiplications in a vector MAC unit are not utilized as often as in other layers. Thus, convolutions comprise a slightly greater portion of accelerator time and energy (74%) compared to model FLOPs (68%). Note that *Conv2DFuse*, with 3072 input channels, more fully utilizes available parallelism in the vector MAC units and thus has relatively low energy per FLOP.

The Swin Tiny model runs in 12ms on accelerator E. The Swin Tiny execution time and energy distributions on accelerator E, shown in Figure 9, closely match the model FLOPs distribution: 87% of the execution time and energy on E are in convolutions, compared to the 89% of FLOPs in convolutions. The distribution of time and energy on accelerator E across individual layers also match the distribution of FLOPs across model layers. For example, the largest layer,

*fpn_bottleneck_conv*, comprises 63% of the total execution time and energy on accelerator E, and 65% of model FLOPs.

Swin Tiny layers with relatively higher energy per FLOP are either convolutions in the decoder with a low number of input channels, due to how we map this computation to maximize the available parallelism, or matrix multiplications in attention blocks, due to odd numbers of input channels that do not fully utilize available parallelism, leading to lower utilization.

With OFA-ResNet-50, the execution time and energy distributions for inference on accelerator E are mostly evenly split among all the convolutions in the CNN. Two layers have significantly higher energy per FLOP compared to the remaining layers. These layers are the first layer (the input layer) and the last layer (the classifier) in the model. The first layer has three input channels for an input image to the model, while the last layer has one input channel for a linear layer. Thus, as with the other models, these layers with a small number of input channels have lower utilization in the vector MAC units.

## V. Enabling RDD Inference

To enable RDD inference with these models, we explore modulating the computation in the most critical layers in order to examine the resulting performance and energy on the A5000 GPU and on accelerator E from Section IV.

For the segmentation models, we examine model resilience to selectively bypassing computation in encoder blocks and convolutional layers with pretrained model weights as well as the execution time and energy savings achievable by switching between retrained model weights. For the pretrained model experiments, we various sweep model and layer parameters and run inference with the corresponding subset of pretrained model weights from the original model. We analyze SegFormer in Section V-A and Swin Transformer in Section V-B. For the DETR-based models, we leverage switching between OFA ResNet-50 models in Section V-C to modulate the computation in the CNN backbone.

Based on these explorations, we identify general principles to guide the search for identifying competitive alternative execution paths for vision transformer models in Section V-D. Finally, Section V-E briefly overviews how we can use these results for RDD inference.

### A. SegFormer B2 Model Resilience and Switching

We examine the resilience of SegFormer B2 with different image sizes, using models trained and validated with 512 by 512 images in the ADE20K dataset [74] (denoted "SegFormer ADE B2") and with 1024 by 1024 images in the Cityscapes dataset [75] (denoted "SegFormer City B2"). Since the majority of FLOPs are in the decoder for both models, we find that the model accuracy is not resilient to reducing scaling factors and input channels in individual layers in the encoder, making these modified execution paths ill-suited for RDD inference.

Next, we examine the impact of reducing the number of input and output channels for the top three layers in the FLOPs distribution in Figure 3: *Conv2DFuse*, *Conv2DPred*,
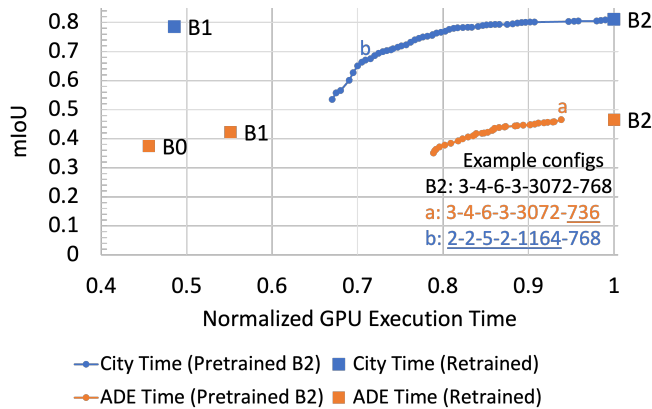
Fig. 10. Tradeoff between execution time and mIoU (model accuracy) when dynamically pruning pretrained SegFormer B2 models on the NVIDIA RTX A5000 GPU (original retrained model executions on the GPU shown as squares). Example dynamic configuration labels list the number of encoder blocks in each of the four encoder stages followed by the input channels for *Conv2DFuse* and *Conv2DPred*; values changed from the original model are underlined.
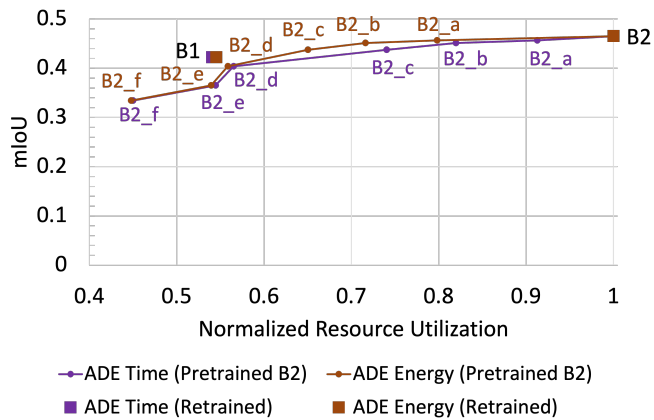


Fig. 11. Tradeoff between execution time, energy, and mIoU (model accuracy) when dynamically pruning pretrained SegFormer B2 models on accelerator E (original retrained model executions shown as squares). The model parameter values for labeled dynamic configurations B2_a to B2_f are listed in Table III.

and *DecodeLinear0*. Reducing the number of input channels to *DecodeLinear0* does not enable skipping computation in earlier layers since the full output of encoder stage zero must still be computed as input for encoder stage one. Similarly, since the interpolation of the output of the *Linear* layers after the encoder stages is reversed and then concatenated, reducing input channels from 3072 to *Conv2DFuse* only allows for skipped computation to propagate backwards if the model execution path uses less than or equal to 768 input channels (corresponding to the contribution from encoder stage three, which is not input to another encoder stage and only input to *Conv2DFuse*). Since the decoder layer outputs have only one following destination, skipped computation associated with input channels to *Conv2DPred* can be propagated backwards through preceding decoder layers (such as *ReLU*) to further reduce FLOPs.

Given a constant number of bypassed channels $x$, we found that which channels that were removed (e.g., the first $x$ versus the last $x$ versus the smallest $x$) did not substantially change the resulting accuracy, so we bypass the last $x$ in our experiments. As more channels are removed from the inputs to *Conv2DFuse*, *Conv2DPred*, and *DecodeLinear0*, the marginal accuracy drop increases. However, if we additionally bypass full encoder blocks, we can shift these trade-off curves more to the left than down (saving relatively more execution time compared to the relative accuracy drop), and produce new Pareto-optimal execution paths. Thus, Figure 10 shows the piecewise concatenation of the Pareto-optimal points on these shifted curves for SegFormer ADE B2 and SegFormer City B2 on the A5000 GPU.

We find that combinations of varying the number of encoder blocks in each stage and input channels to *Conv2DFuse* and *Conv2DPred* result in Pareto-optimal points, as shown with the "ADE Time (Pretrained B2)" and "City Time (Pretrained B2)"

lines. We generate these two lines by pruning the pretrained SegFormer B2 model weights, with no additional training. Figure 10 denotes the model parameters that two example Pareto-optimal configurations $a$ and $b$ correspond to, along with the original B2 model configuration. Surprisingly, model configuration $a$ achieves slightly better mIoU (0.4655 without any retraining and 0.4698 after training) than the original SegFormer B2 ADE model (0.4651). This point has 32 fewer input channels to *Conv2DPred* compared to the B2 model and is 6% faster on the GPU, so we can start the ADE Pareto Curve from $a$ instead of the full SegFormer ADE B2 model square.

On the GPU, we can save 11% of execution time with a 1.9% accuracy loss for SegFormer ADE B2, without any retraining. The pretrained SegFormer City B2 model is more resilient: we can save 11% of execution time with only a 0.9% accuracy loss, without any retraining. SegFormer City B2 is trained and executed on larger image sizes (1024 by 1024 vs 512 by 512), so it achieves 1.74× higher mIoU than SegFormer ADE B2 to begin with. Thus, there is likely more redundancy in these model weights, enabling a flatter tradeoff compared to SegFormer ADE B2.

We plot the original SegFormer B0, B1, and B2 models, labeled as "ADE Time (Retrained)" and "City Time (Retrained)" in Figure 10, when executed on the GPU. We can switch between SegFormer models with retrained weights to save 51% of execution time with a 4.3% accuracy loss for the ADE B2 model and 45% of execution time with a 2.5% accuracy loss for the City B2 model. As expected, switching between pruned models that have been retrained offers a better tradeoff, but the pretrained models are surprisingly resilient: pruning without retraining is competitive up until reducing 15% of the execution time for the B2 models on the GPU.

Figure 11 shows the execution time, energy, and accuracy tradeoff for the various model configurations listed in Table III, when executed on accelerator E. These points are a subset of the Pareto-optimal configurations found from the GPU

| Label | Number of Encoder Blocks in Stages 0 to 3 | *Conv2DFuse* Input Channels | mIoU |
|---|---|---|---|
| B2 | 3, 4, 6, 3 | 3072 | 0.4651 |
| B2_a | 3, 4, 6, 3 | 1920 | 0.4565 |
| B2_b | 3, 4, 6, 3 | 1664 | 0.4510 |
| B2_c | 2, 4, 6, 3 | 1408 | 0.4374 |
| B2_d | 2, 3, 6, 3 | 1024 | 0.4041 |
| B2_e | 2, 3, 5, 3 | 896 | 0.3649 |
| B2_f | 2, 3, 5, 3 | 512 | 0.3345 |

TABLE III

SEGFORMER ADE B2 MODEL EXECUTION PATH CONFIGURATIONS, AS LABELED IN FIGURE 11.



Fig. 13. Tradeoff between execution time, energy, and model accuracy on the ImageNet dataset, with a 224 by 224 input image size, when dynamically switching between OFA ResNet-50 models on accelerator E.
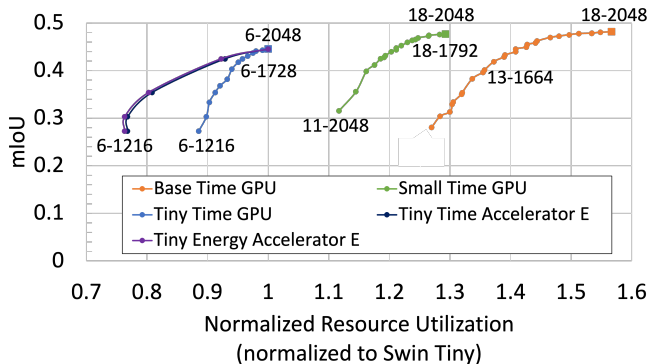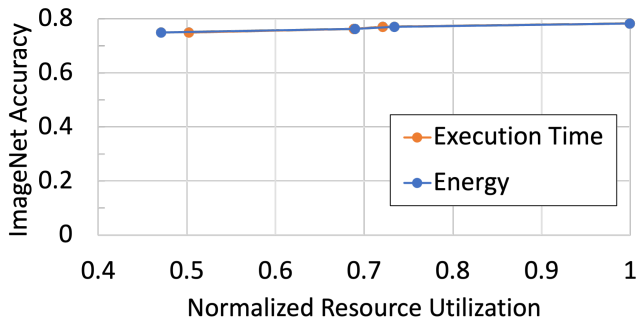


Fig. 12. Tradeoff between execution time, energy and mIoU (model accuracy) when dynamically pruning pretrained the Swin Base, Small, and Tiny models on the A5000 GPU and accelerator E (original retrained model executions on the GPU shown as squares). Dynamic configuration labels list the number of encoder blocks in encoder stage three followed by the number of channels input to *fpn_bottleneck_Conv2D*. The resource utilization numbers on the x-axis are normalized to the Swin Tiny execution results.

experiments. With specialized hardware, we can save 18% of execution time and 28% of energy with a 1.4% accuracy loss for SegFormer ADE B2 without any additional training.

We also plot the original SegFormer B1 and B2 models, labeled as "ADE Time (Retrained)" and "ADE Energy (Retrained)" in Figure 11, and see that with retraining, we can save 55% of execution time and energy for a 4.3% accuracy loss. On accelerator E, pruning without retraining reduces execution time and energy by 45% with the same accuracy loss. Thus, switching between retrained models offers an additional 10% of execution time and energy savings at the same accuracy for the SegFormer ADE model. One of these approaches can be chosen over the other depending on the range of dynamic resource constraints expected in a real-time system.

### B. Swin Model Resilience and Switching

In Swin-Tiny, almost 90% of the FLOPs are in convolutions in the decoder. As a result, we vary the input and output channels in these layers along with bypassing encoder blocks, in a similar approach to our SegFormer experiments. Skipping even a few encoder blocks in Swin Tiny leads to a higher drop in model accuracy compared to execution time on the GPU and accelerator E, as shown in Figure 12, and the marginal accuracy loss quickly outpaces the marginal savings in execution time and energy.

Swin Small and Swin Base, with $1.35\times$ and $2\times$ as many parameters as Swin Tiny, are slightly more resilient to these optimizations. While all three models have the same number of input channels to the *fpn_bottleneck_Conv2D* layer, Swin Small and Base have 18 blocks in encoder stage two compared to only six in Swin Tiny. Thus, we can bypass stage two encoder blocks to achieve a better tradeoff. However, the resilience to pruning is small even for these larger models.

While the tradeoff curve is flatter with accelerator E than with the GPU, we can only save 8% of execution time and energy with a 2% accuracy loss for Swin Tiny inference on the accelerator, without any retraining. Thus, unlike with Seg-Former, we recommend directly switching between retrained Swin model weights for RDD inference. On the GPU, this approach saves 36% of execution time for a 3.6% accuracy drop when switching from Swin Base to Swin Tiny.

While the Swin models require $3.8\times$ to $4.7\times$ more FLOPs than SegFormer ADE B2, these larger pretrained models are not more resilient to pruning. The result is because most of the extra FLOPs in Swin are from additional convolutions in the decoder and Swin Tiny has 25% fewer encoder blocks compared to SegFormer. Thus, while we could bypass encoder blocks to generate a flatter tradeoff curve with SegFormer, the Swin encoder has less redundant information, and the Swin model accuracy is more sensitive to skipping attention layers in the encoder.

### C. Once-For-All ResNet-50 Model Switching

Since the ResNet-50 backbone dominates FLOPs for the DETR-based models, we need to scale the work done in this CNN for RDD inference. Fortunately, this is a well-studied area [47], [61], [64], and we leverage switching between OFA ResNet-50 models, which scale model parameters such as the number of layers and convolution kernel sizes. The advantage of the OFA approach is that many model weight subsets are generated within a single training [64].

We execute various OFA models for ResNet-50 on accelerator E. Figure 13 shows the effectiveness of switching between these models for dynamic inference on accelerator E: we can save 58% of the execution time and 53% of energy with a 3.3%

loss in accuracy. Switching between these models produces a better tradeoff compared to switching between retrained segmentation models on accelerator E, which have comparable execution time and energy savings with a 9% accuracy loss. This result illustrates that CNNs are more resilient to pruning compared to these segmentation models with the transformer architecture.

### D. Principles for identifying alternative model execution paths

Our results in Sections V-A to V-C highlight general principles to help focus the search for competitive alternative execution paths in vision transformer models. First, we identify that these model executions greatly depend on convolutions, and that we should target convolutional paths in these models instead of pruning attention heads and attention layers as in prior work. Furthermore, our experiments in this section show that convolution weights are less critical to model accuracy compared to attention weights. Second, for modern encoder-decoder architectures, it is advantageous to prune computation in the decoder instead of the encoder. In these attention-free decoders, layers sequentially feed into each other and have one destination, while most encoder layer outputs have two destinations (to the following encoder layers and the decoder). Thus, skipped computation in the decoder can often be propagated backwards through decoder layers to further reduce FLOPs as exemplified with SegFormer in Section V-A, while this can rarely be done for skipped computation in the encoder. Third, the distribution of FLOPs between the encoder and decoder guides which alternative execution paths will be competitive. When more FLOPs are in the decoder, there is less redundancy in encoder weights, and bypassing encoder layers will lead to higher accuracy losses more quickly, as seen in Section V-B with pruning Swin Transformer, where 89% of the FLOPs are in the decoder.

### E. RDD Inference Execution

In this section, we have identified alternative model execution paths that bypass computation to meet specific resource targets for our model case studies. During runtime, the current resource availability is input to the inference engine, which can select the model execution path that matches the resource constraints for that particular inference. Note that this selection is independent of the input image. Each inference task can use a different model configuration that is looked up in real time, as opposed to a pre-selection of a static model to always execute. When enough resources are available, the original model can be executed without any modifications. This means that the average accuracy loss will usually be lower than the percentages reported for particular model configurations in this section, since RDD inference is usually using the full model accuracy.

For more resilient models like SegFormer, we only need to train the original model once and store one set of model weights. During inference, we can use the subset of the weights required for model execution. For other models, we can switch between multiple sets of model weights stored in off-chip DRAM for execution on an accelerator. The OFA approach still requires only one training for this approach with ResNet-50, while switching between the Swin models requires multiple trainings for each model configuration.

Note that analyzing the resilience of pretrained models provides a floor on the accuracy attainable given a model architecture and a model configuration, while switching between retrained models provides a ceiling. Training techniques such as OFA can be used to improve model accuracy from the floor with less training overheads for RDD inference.

Finally, since we bypass some computation in critical layers for RDD inference, the relative importance of these layers decreases for model configurations with less computation. For example, configuration $B2\_f$ in Figure 11 requires 60% fewer FLOPs than the full SegFormer ADE B2 model, and *Conv2DFuse* requires less than 25% of the total FLOPs for this configuration. While the set of critical layers change across model configurations, these layers are still convolutions. For example, 55% of the total execution time and energy on accelerator E is spent in convolutions for the smaller $B2\_f$ configuration. Thus, we find that accelerators D, E, G are still the best at balancing area, energy, and throughput tradeoffs across all these dynamic model configurations for RDD inference.

## VI. RELATED WORK

Static methods such as pruning [64], [86], [87], knowledge distillation [21], [88], [89], and quantization [90]–[92] reduce model computation but are unable to dynamically adapt to resource constraints during runtime. Input-dependent methods instead dynamically adjust computation based on input complexity and the state of internal predictions [48]–[60], but they reduce inference cost only for inputs that do not require full model execution. Both of these methods can be used with RDD inference to improve efficiency-accuracy tradeoffs.

Few RDD inference methods have been previously explored. One work adds multiple early exit paths in models [48]. Instead of only sequentially skipping layers later in the execution path, we more broadly explore skipping intermediate blocks. Another work augments ResNet and MobileNet with scaled-down versions of blocks of convolutional layers in these models that can be trained and executed instead of the original blocks [47]. We instead explore whether the existing model architecture and weights can be leveraged for RDD inference. Compared to both of these prior works, we explore a finer granularity by pruning computation associated with specific channels in critical model layers. In addition, these methods only consider CNNs, while our work is the first to profile the vision transformer applications, focus on vision tasks more complicated than classification, and explore the resilience of pretrained segmentation models compared to retrained models.

## VII. CONCLUSION

While these vision models are transformers, most of the computation comes from the backbones and task-specific decoder heads. These computations often use convolutions for

efficiency, which means that 60% to 90% of FLOPs in these applications are in convolutions. On current GPUs, computing CNN FLOPs is very cheap, so these layers take only 20% to 45% of the application running time. Thus, accelerating vision transformer applications requires hardware which can handle CNN and transformer computation well. We found that the MAGNet framework worked well for these applications.

To enable dynamic inference, we found that one needs to bypass both attention-dominated encoder blocks and computation in decoder convolutions to maintain reasonable accuracy. We find that models are more resilient to pruning without additional training when computation is more evenly split between the encoder and the decoder. Surprisingly, even in smaller model configurations, including those with less than half of the original model FLOPs, convolutions still dominate the total computation.

REFERENCES

[1] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *International Conference on Computer Vision (ICCV)*, 2017.
[2] R. Girshick, "Fast r-cnn," in *International Conference on Computer Vision (ICCV)*, 2015.
[3] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2015.
[4] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
[5] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
[6] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European Conference on Computer Vision (ECCV)*, 2016.
[7] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr, and L. Zhang, "Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers," in *Conference on Computer Vision and Pattern Recognition*, 2021.
[8] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
[9] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
[10] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *International Conference on Learning Representations (ICLR)*, 2015.
[11] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning (ICML)*, 2019.
[12] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
[13] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European Conference on Computer Vision (ECCV)*, 2020.
[14] S. Liu, F. Li, H. Zhang, X. Yang, X. Qi, H. Su, J. Zhu, and L. Zhang, "DAB-DETR: Dynamic anchor boxes are better queries for DETR," in *International Conference on Learning Representations (ICLR)*, 2022.
[15] Y. Wang, X. Zhang, T. Yang, and J. Sun, "Anchor detr: Query design for transformer-based detector," in *AAAI Conference on Artificial Intelligence (AAAI)*, 2022.
[16] D. Meng, X. Chen, Z. Fan, G. Zeng, H. Li, Y. Yuan, L. Sun, and J. Wang, "Conditional detr for fast training convergence," in *International Conference on Computer Vision (ICCV)*, 2021.

[17] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, "Segformer: Simple and efficient design for semantic segmentation with transformers," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
[18] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *International Conference on Computer Vision (ICCV)*, 2021.
[19] OpenAI, "Gpt-4v(ision) system card," 2023. [Online]. Available: https://openai.com/research/gpt-4v-system-card
[20] 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 *International Conference on Learning Representations (ICLR)*, 2021.
[21] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, "Training data-efficient image transformers & distillation through attention," in *International Conference on Machine Learning (ICML)*, 2021.
[22] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, "Segment anything," in *International Conference on Computer Vision (ICCV)*, 2023.
[23] A. Hatamizadeh, H. Yin, G. Heinrich, J. Kautz, and P. Molchanov, "Global context vision transformers," in *International Conference on Machine Learning (ICML)*, 2023.
[24] C.-F. R. Chen, Q. Fan, and R. Panda, "Crossvit: Cross-attention multi-scale vision transformer for image classification," in *International Conference on Computer Vision (ICCV)*, 2021.
[25] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable detr: Deformable transformers for end-to-end object detection," in *International Conference on Learning Representations (ICLR)*, 2021.
[26] X. Dai, Y. Chen, J. Yang, P. Zhang, L. Yuan, and L. Zhang, "Dynamic detr: End-to-end object detection with dynamic attention," in *International Conference on Computer Vision (ICCV)*, 2021.
[27] M. Zheng, P. Gao, R. Zhang, K. Li, X. Wang, H. Li, and H. Dong, "End-to-end object detection with adaptive clustering transformer," in *British Machine Vision Conference (BMVC)*, 2021.
[28] Z. Yao, J. Ai, B. Li, and C. Zhang, "Efficient detr: improving end-to-end object detector with dense prior," *arXiv preprint arXiv:2104.01318*, 2021.
[29] B. Roh, J. Shin, W. Shin, and S. Kim, "Sparse detr: Efficient end-to-end object detection with learnable sparsity," in *International Conference on Learning Representations (ICLR)*, 2022.
[30] R. Strudel, R. Garcia, I. Laptev, and C. Schmid, "Segmenter: Transformer for semantic segmentation," in *International Conference on Computer Vision (ICCV)*, 2021.
[31] H. Zhang, F. Li, S. Liu, L. Zhang, H. Su, J. Zhu, L. M. Ni, and H.-Y. Shum, "Dino: Detr with improved denoising anchor boxes for end-to-end object detection," in *International Conference on Learning Representations (ICLR)*, 2023.
[32] Z. Zong, G. Song, and Y. Liu, "Detrs with collaborative hybrid assignments training," in *International Conference on Computer Vision (ICCV)*, 2023.
[33] X. Zhou, R. Girdhar, A. Joulin, P. Krähenbühl, and I. Misra, "Detecting twenty-thousand classes using image-level supervision," in *European Conference on Computer Vision (EECV)*, 2022.
[34] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
[35] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
[36] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving language understanding by generative pre-training," 2018.
[37] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language models are unsupervised multitask learners," *OpenAI blog*, 2019.
[38] 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," *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.

[39] OpenAI, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.

[40] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *Journal of Machine Learning Research (JMLR)*, 2020.

[41] S. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. S. Khan, and M. Shah, "Transformers in vision: A survey," *ACM computing surveys (CSUR)*, 2022.

[42] D. Feng, C. Haase-Schütz, L. Rosenbaum, H. Hertlein, C. Glaeser, F. Timm, W. Wiesbeck, and K. Dietmayer, "Deep multi-modal object detection and semantic segmentation for autonomous driving: Datasets, methods, and challenges," *IEEE Transactions on Intelligent Transportation Systems*, 2020.

[43] M. Siam, M. Gamal, M. Abdel-Razek, S. Yogamani, M. Jagersand, and H. Zhang, "A comparative study of real-time semantic segmentation for autonomous driving," in *Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2018.

[44] T. Zhou, F. Porikli, D. J. Crandall, L. Van Gool, and W. Wang, "A survey on deep learning technique for video segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

[45] J. Sun, W. Yi, P. K. Varshney, and L. Kong, "Resource scheduling for multi-target tracking in multi-radar systems with imperfect detection," *IEEE Transactions on Signal Processing*, 2022.

[46] S. Kim, J. Zhao, K. Asanovic, B. Nikolic, and Y. S. Shao, "Aurora: Virtualized accelerator orchestration for multi-tenant workloads," in *International Symposium on Microarchitecture (MICRO)*, 2023.

[47] J. Clemons, I. Frosio, M. Shen, J. M. Alvarez, and S. W. Keckler, "Augmenting legacy networks for flexible inference," in *IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2023.

[48] Y. Wang, J. Shen, T.-K. Hu, P. Xu, T. Nguyen, R. Baraniuk, Z. Wang, and Y. Lin, "Dual dynamic inference: Enabling more efficient, adaptive, and controllable deep inference," *IEEE Journal of Selected Topics in Signal Processing*, 2020.

[49] J. Xin, R. Tang, J. Lee, Y. Yu, and J. Lin, "Deebert: Dynamic early exiting for accelerating bert inference," *Association for Computational Linguistics (ACL)*, 2020.

[50] Y. Kaya, S. Hong, and T. Dumitras, "Shallow-deep networks: Understanding and mitigating network overthinking," in *International Conference on Machine Learning (ICML)*, 2019.

[51] T.-K. Hu, T. Chen, H. Wang, and Z. Wang, "Triple wins: Boosting accuracy, robustness and efficiency together by enabling input-adaptive inference," *International Conference on Learning Representations (ICLR)*, 2020.

[52] A. Graves, "Adaptive computation time for recurrent neural networks," *arXiv preprint arXiv:1603.08983*, 2016.

[53] W. Liu, P. Zhou, Z. Zhao, Z. Wang, H. Deng, and Q. Ju, "Fastbert: a self-distilling bert with adaptive inference time," *Association for Computational Linguistics (ACL)*, 2020.

[54] M. Figurnov, M. D. Collins, Y. Zhu, L. Zhang, J. Huang, D. Vetrov, and R. Salakhutdinov, "Spatially adaptive computation time for residual networks," in *Computer Vision and Pattern Recognition (CVPR)*, 2017.

[55] Y. Han, G. Huang, S. Song, L. Yang, H. Wang, and Y. Wang, "Dynamic neural networks: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

[56] M. Dehghani, S. Gouws, O. Vinyals, J. Uszkoreit, and Ł. Kaiser, "Universal transformers," *International Conference on Learning Representations (ICLR)*, 2018.

[57] S. Teerapittayanon, B. McDanel, and H.-T. Kung, "Branchynet: Fast inference via early exiting from deep neural networks," in *International Conference on Pattern Recognition (ICPR)*, 2016.

[58] W. Zhou, C. Xu, T. Ge, J. McAuley, K. Xu, and F. Wei, "Bert loses patience: Fast and robust inference with early exit," *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.

[59] T. Tambe, C. Hooper, L. Pentecost, T. Jia, E.-Y. Yang, M. Donato, V. Sanh, P. Whatmough, A. M. Rush, D. Brooks, and G.-Y. Wei, "Edgebert: Sentence-level energy optimizations for latency-aware multi-task nlp inference," in *International Symposium on Microarchitecture (MICRO)*, 2021.

[60] X. Wang, F. Yu, Z.-Y. Dou, T. Darrell, and J. E. Gonzalez, "Skipnet: Learning dynamic routing in convolutional networks," in *European Conference on Computer Vision (ECCV)*, 2018.

[61] S. Jastrzkebski, D. Arpit, N. Ballas, V. Verma, T. Che, and Y. Bengio, "Residual connections encourage iterative inference," *International Conference on Learning Representations (ICLR)*, 2018.

[62] P. Michel, O. Levy, and G. Neubig, "Are sixteen heads really better than one?" *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

[63] E. Voita, D. Talbot, F. Moiseev, R. Sennrich, and I. Titov, "Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned," *Association for Computational Linguistics (ACL)*, 2019.

[64] H. Cai, C. Gan, T. Wang, Z. Zhang, and S. Han, "Once-for-all: Train one network and specialize it for efficient deployment," *International Conference on Learning Representations (ICLR)*, 2020.

[65] T. Xiao, Y. Liu, B. Zhou, Y. Jiang, and J. Sun, "Unified perceptual parsing for scene understanding," in *European Conference on Computer Vision (ECCV)*, 2018.

[66] J. Choquette, "Nvidia hopper h100 gpu: Scaling performance," *IEEE Micro*, 2023.

[67] "Nvidia tensorrt," https://developer.nvidia.com/tensorrt, 2023.

[68] S. Chetlur, C. Woolley, P. Vandermersch, J. Cohen, J. Tran, B. Catanzaro, and E. Shelhamer, "cudnn: Efficient primitives for deep learning," *arXiv preprint arXiv:1410.0759*, 2014.

[69] R. Venkatesan, Y. S. Shao, M. Wang, J. Clemons, S. Dai, M. Fojtik, B. Keller, A. Klinefelter, N. Pinckney, P. Raina, Y. Zhang, B. Zimmer, W. J. Dally, J. Emer, S. W. Keckler, and B. Khailany, "Magnet: A modular accelerator generator for neural networks," in *International Conference on Computer-Aided Design (ICCAD)*. IEEE, 2019.

[70] B. Keller, R. Venkatesan, S. Dai, S. G. Tell, B. Zimmer, W. J. Dally, C. T. Gray, and B. Khailany, "A 17–95.6 tops/w deep learning inference accelerator with per-vector scaled 4-bit quantization for transformers in 5nm," in *IEEE Symposium on VLSI Technology and Circuits (VLSI)*, 2022.

[71] H. Yu, Z. Yang, L. Tan, Y. Wang, W. Sun, M. Sun, and Y. Tang, "Methods and datasets on semantic segmentation: A review," *Neurocomputing*, 2018.

[72] Z. Zou, K. Chen, Z. Shi, Y. Guo, and J. Ye, "Object detection in 20 years: A survey," *Proceedings of the IEEE*, 2023.

[73] M. Alam, M. D. Samad, L. Vidyaratne, A. Glandon, and K. M. Iftekharuddin, "Survey on deep neural networks in speech and vision systems," *Neurocomputing*, 2020.

[74] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba, "Scene parsing through ade20k dataset," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[75] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[76] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European Conference on Computer Vision (ECCV)*, 2014.

[77] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," in *International Conference on Computer Vision (ICCV)*, 2021.

[78] E. Xie, W. Wang, W. Wang, P. Sun, H. Xu, D. Liang, and P. Luo, "Segmenting transparent object in the wild with transformer," *International Joint Conference on Artificial Intelligence (IJCAI)*, 2021.

[79] Y. Li, Y. Hu, F. Wu, and K. Li, "Divit: Algorithm and architecture co-design of differential attention in vision transformer," *Journal of Systems Architecture*, 2022.

[80] X. Wang, L. L. Zhang, Y. Wang, and M. Yang, "Towards efficient vision transformer inference: a first study of transformers on mobile devices," in *International Workshop on Mobile Computing Systems and Applications (HotMobile)*, 2022.

[81] H. You, Z. Sun, H. Shi, Z. Yu, Y. Zhao, Y. Zhang, C. Li, B. Li, and Y. Lin, "Vitcod: Vision transformer acceleration via dedicated algorithm and accelerator co-design," *International Symposium on High-Performance Computer Architecture (HPCA)*, 2023.

[82] "detrex: An research platform for transformer-based object detection algorithms," https://github.com/IDEA-Research/detrex, 2022.

[83] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.

[84] K. Zeng, Q. Ma, J. W. Wu, Z. Chen, T. Shen, and C. Yan, "Fpga-based accelerator for object detection: A comprehensive survey," *The Journal of Supercomputing*, 2022.

[85] X. Yang, M. Gao, Q. Liu, J. Setter, J. Pu, A. Nayak, S. Bell, K. Cao, H. Ha, P. Raina, C. Kozyrakis, and M. Horowitz, "Interstellar: Using halide's scheduling language to analyze dnn accelerators," in *International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2020.

[86] H. Wang, Z. Wu, Z. Liu, H. Cai, L. Zhu, C. Gan, and S. Han, "Hat: Hardware-aware transformers for efficient natural language processing," *Association for Computational Linguistics (ACL)*, 2020.

[87] T.-J. Yang, A. Howard, B. Chen, X. Zhang, A. Go, M. Sandler, V. Sze, and H. Adam, "Netadapt: Platform-aware neural network adaptation for mobile applications," in *European Conference on Computer Vision (ECCV)*, 2018.

[88] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter," *Workshop on Energy Efficient Machine Learning and Cognitive Computing co-located with Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

[89] S. Sun, Y. Cheng, Z. Gan, and J. Liu, "Patient knowledge distillation for bert model compression," *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2019.

[90] H. Bai, W. Zhang, L. Hou, L. Shang, J. Jin, X. Jiang, Q. Liu, M. Lyu, and I. King, "Binarybert: Pushing the limit of bert quantization," in *Association for Computational Linguistics (ACL)*, 2021.

[91] Z. Liu, Y. Wang, K. Han, W. Zhang, S. Ma, and W. Gao, "Post-training quantization for vision transformer," *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.

[92] O. Zafrir, G. Boudoukh, P. Izsak, and M. Wasserblat, "Q8bert: Quantized 8bit bert," in *Workshop on Energy Efficient Machine Learning and Cognitive Computing (EMC2) co-located with Advances in Neural Information Processing Systems (NeurIPS)*, 2019.