

Moses: Efficient Exploitation of Cross-device Transferable Features for Tensor Program Optimization

Zhihe Zhao[†], Xian Shuai[†], Yang Bai[†], Neiwen Ling[†], Nan Guan[‡], Zhenyu Yan[†], Guoliang Xing^{†,§}

[†]The Chinese University of Hong Kong

[‡]City University of Hong Kong

[§]Corresponding Author

Abstract

Achieving efficient execution of machine learning models has attracted significant attention recently. To generate tensor programs efficiently, a key component of DNN compilers is the cost model that can predict the performance of each configuration on specific devices. However, due to the rapid emergence of hardware platforms, it is increasingly labor-intensive to train domain-specific predictors for every new platform. Besides, current design of cost models cannot provide transferable features between different hardware accelerators efficiently and effectively. In this paper, we propose Moses, a simple and efficient design based on the lottery ticket hypothesis, which fully takes advantage of the features transferable to the target device via domain adaptation. Compared with state-of-the-art approaches, Moses achieves up to 1.53X efficiency gain in the search stage and 1.41X inference speedup on challenging DNN benchmarks.

1 Introduction

Efficient inference of deep neural networks (DNN) is of great importance for real-time AI applications such as autonomous driving and augmented reality, especially given that they usually run on embedded devices with limited compute power [Zhao *et al.*, 2021]. Existing approaches usually rely on hand-optimized acceleration libraries, e.g., NVIDIA cuDNN and Intel MKL. However, they are vendor-specific, which cannot support a wide range of diverse hardware devices and require significant engineering efforts for tuning.

Typically, the tensor computation inside deep learning operators is implemented by a set of compute-intensive nested loops. For example, Figure 1 shows an optimized tensor program of a 2D convolutional operator, which consists of multiple for-loops and involves three schedule primitives: blocking, unrolling, and vectorization. However, generating such high-performance tensor programs from a given high-level expression is extremely difficult, as the optimal organization and the parameters of the for-loops can vary significantly for different devices. Therefore, to accelerate the end-to-end model inference on various hardware platforms, existing DNN compilers first generate a large space of configura-

```
Parallel n.@oc_chunk.0@oh.0@ow.0@oc_block.0@n.1@oc_chunk.1@oh.1@ (0,240)
conv2d_NCHWc.local auto_unroll: 512
for ic.0 (0,3)
  for oc_chunk_c.2 (0,2)
    for ow_c.2 (0,2)
      for kh.1 (0,3)
        for kw.1 (0,3)
          for ow_c.3 (0,15)
            vectorize oc_block_c.3 (0,4)
              conv2d_NCHWc.local = ...
          for oc_chunk.2 (0,2)
            for ow.2 (0,30)
              vectorize oc_block.2 (0,4)
                conv2d_NCHWc = ...
```

Figure 1: An example of the program for a 2D convolution operator: $Conv2d(in_channels = 3, out_channels = 64, kernel_size = 3, stride = 1, padding = 0, bias = False)$ produced by TVM.

tions, and then search for the best-performing one based on on-device measurements [Chen *et al.*, 2018b]. This process is termed as auto-tuning.

Although DNN compilers can produce optimized programs for DNN models on specific platforms, they suffer excessively long search time. For example, although AutoTVM can outperform nearly 2 \times over default TensorFlow on ResNet-18, the auto-tuning time can take up to tens of hours on embedded devices (e.g. NVIDIA Jetson TX2). To reduce the time-consuming on-device measurements, the tensor program optimizer employs a cost model to directly predict the performance of the potential candidates in the search space. However, training a cost model offline still requires a large number of measurements. For instance, Tenset [Zheng *et al.*, 2021] provides a tensor program performance dataset collected from 6 devices, containing 52 million program performance records. Based on this dataset, it was shown that the cost model can be transferred between two Intel CPUs by fine-tuning. As a result, the online search time can be reduced without sacrificing the quality of optimized tensor programs. However, when two hardware platforms differ significantly in architecture, such a vanilla fine-tuning approach would fail to learn the runtime behaviors of a new device, and hence performs poorly at generating high-performance tensor programs, as evidenced by our results in Section 4.

Previous efforts trying to address this challenge mainly focus on either designing a new cost model [Baghdadi *et al.*, 2021][Kaufman *et al.*, 2020], or exploring effective search algorithms during auto-tuning [Ahn *et al.*, 2020][Li *et al.*, 2020]. However, these approaches still need large amounts

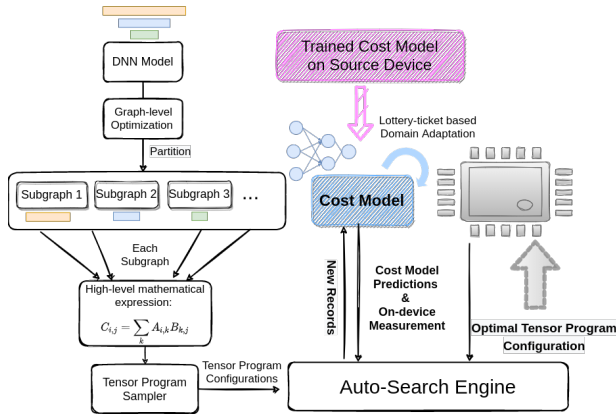


Figure 2: The complete pipeline of automatic tensor program generation for a given neural network.

of iterations during the search and the performance is usually highly program-dependent. Motivated by these challenges, we propose Moses, a novel cost model adaptation framework based on the lottery ticket hypothesis [Frankle and Carbin, 2019], which can adapt the trained cost model from a source device to a new target device with high efficiency. We summarize the contributions of this work as follows:

- 1) To the best of our knowledge, Moses is the first work that achieve highly efficient for auto-tuning between different hardware platforms with transfer learning. Moses enables the DNN compiler to generate optimized tensor programs with significantly shorter search time for a new device.
- 2) We propose a novel approach that can automatically identify the transferable hardware-independent parameters of a pre-trained cost model, and achieve cross-device cost model adaption via fine-tuning.
- 3) We conduct comprehensive experiments and show that Moses is a general and effective approach for diverse hardware platforms and DNNs.

2 Related Work and Background

2.1 DNN Compilers.

General DNN compilers optimize the computation flow of DNN tasks in two levels: graph-level and tensor-level. Some notable compilers are TVM [Chen *et al.*, 2018a], TASO [Jia *et al.*, 2019], XLA [Abadi *et al.*, 2016] and Halide [Li *et al.*, 2018]. These compilers either utilize compiler techniques such as graph substitutions to optimize the intermediate representation (IR) level graph [Bai *et al.*, 2021], or focus on tensor program optimization using heuristic and learning-based algorithms to joint-optimize the polyhedral patterns for DNN. Building on these DNN compilers, recent works treat the optimization process as a black box and propose some advanced searching or cost model training approaches based on runtime information [Chen *et al.*, 2018b][Ahn *et al.*, 2020][Li *et al.*, 2020][Haj-Ali *et al.*, 2020].

2.2 Auto-Tuning with Cost Models

Figure 2 shows the pipeline of a search-based low-level tensor code generation used by TVM [Chen *et al.*, 2018a]. The intact neural network is first partitioned by the graph-level optimizer. For instance, the ResNet-50 will be divided into 29 subgraphs and each subgraph normally includes one or two convolutional layers. Given an intact mathematical expression or a computational graph in a high-level mathematical expression, the search algorithm will search for the best low-level tensor implementation for the target hardware. Usually, the search space is in the order of millions for CPUs and billions for GPUs, as there are a variety of schedule primitives such as tiling, unrolling, vectorization, parallelization, and thread binding. Each schedule primitive can involve multiple tunable knobs. On one hand, such a large search space enables the automatic tensor compiler to find the program that is better than the hand-optimized implementation. On the other hand, the large search space can incur significant search time, especially for embedded devices with limited computation power.

To accelerate the searching process, TVM introduces the cost model to directly predict the time cost of the innermost non-loop program rather than extensively measuring program’s runtime. In this paper, we adopt the 164-*d* features in Anso [Zheng *et al.*, 2020a] to depict the program. In each iteration, based on predictions from the cost model, a batch of candidate programs are sampled by an evolutionary search engine. Then, TVM measures the actual time costs of these sampled programs. Finally, the measurements are added to the training data of the cost model. Therefore, the cost model and searched tensor program are improved iteratively.

Notably, TVM includes two search frameworks: AutoTVM [Chen *et al.*, 2018b] and Anso [Zheng *et al.*, 2020a]. AutoTVM requires hand-written scheduling templates for searching, while Anso is a fully automated framework. Due to the large number of hardware platforms and the substantial efforts to develop templates, we use Anso in this paper.

2.3 Cross-Device Cost Model Adaptation

As discussed in the previous section, a major practical drawback of the current automatic tensor program optimization approach [Chen *et al.*, 2018b] is the extremely long search time. Recent work such as [Ahn *et al.*, 2020] provides a breakdown of the search time and shows that the time for on-device measurements dominates. Therefore, a solution to shorten the online searching time is to collect an comprehensive tensor program performance dataset offline, and pre-train a cost model that can be directly utilized [Zheng *et al.*, 2021]. However, the cost model is device-specific, as its input features include configuration knobs that are closely related to the hardware architecture such as the lengths of BlockIdx and ThreadIdx. When two architectures are significantly different (e.g., server GPUs and mobile GPUs), the traditional transfer learning approach may fail.

3 Moses

In this section, we present the overview, problem formulation and design of Moses.

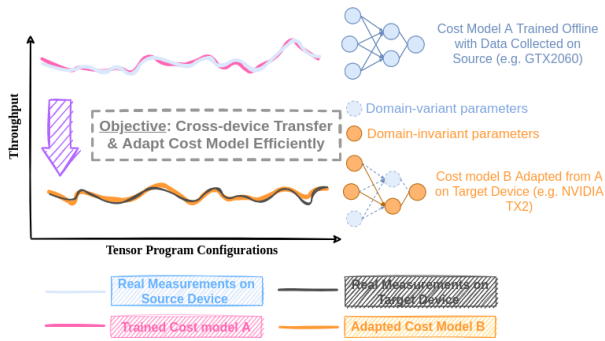


Figure 3: Given a cost model A trained on the source device, we aim to obtain a cost model B that can accurately predict the throughput of the target device under various tensor program configurations.

3.1 Overview

We propose Moses, a novel cross-device cost model transfer approach based on the lottery ticket hypothesis [Frankle and Carbin, 2019]. As shown in Figure 3, we transfer the cost model trained on source devices (e.g., Server GPUs) to target devices (e.g., mobile GPUs) by only fine-tuning portion of the model parameters while keeping the rest of parameters deactivated. The rationale of our design is two-fold. First, to accelerate the online searching instead of collecting a dataset for every new device offline, we need to take advantage of the cost model pre-trained on source devices rather than training a new model from scratch. Second, as vanilla fine-tuning approaches may fail due to substantial architecture changes, we have to utilize the prior knowledge from the pre-trained cost model wisely. Specifically, we leverage the lottery ticket hypothesis to identify the transferable parameters in the cost model space [Han *et al.*, 2021]. By distilling the transferable parameters and dropping out the untransferable ones, we can not only shrink the size of the cost model but also reduce the device-specific information, which accelerates the training of the cost model and improves the generalizability.

3.2 Problem Formulation

The Auto-tuning process in a DNN compiler aims to generate a large search space of tensor program configurations and to find the optimal one based on the on-device performance measurement records [Ryu and Sung, 2021]. We denote the transformation pass of tensor programs by t and the cost model by C . During the search process in auto-tuning, the compiler picks the top- k (where k can be one) programs with the performance predictions from the cost model for each task. The process can utilize mixed on-device measurements and cost model predictions. If the on-device measurement cost is large, the process can completely rely on the cost model; Function $Perf()$ represents the run-time hardware measurement (throughput, GFLOPs); Function $C()$ represents the cost model predictions on codes performance; i is an input task, which can be a computing subgraph in the input DNN model. For example, SqueezeNet consists of 23 tasks (a.k.a, the subgraphs) in total. Note that a subgraph is a unit with the finest granularity during the compilation process; Ψ is the set of possible tensor program configurations of a task which consists of a combination of parameters called *knobs*;

$g()$ represents the tensor programs generating functions with knobs and the transformation pass as inputs. The approximation of a cost model function to real hardware measurement can be denoted by $C() \sim Perf()$; Thus, given an input subgraph i , the objective of the auto-tuning process is finding the optimal combination of knobs ψ^* to maximize the performance defined as below:

$$\psi^* = \arg \max_{\psi \in \Psi} Perf(g(\psi, t)) \quad (1)$$

Precise estimate of the performance of tensor program candidates can effectively reduce the time-consuming on-device measurements. Thus, the objective of the cost model is to minimize the difference between the predictions and the real-world measurements, which is:

$$\min \|C(g(\Psi, t)|\psi, \theta) - Perf(g(\Psi, t))\| \quad (2)$$

Here, θ represents the parameter sets in the trained cost model. We train a high-performance cost model is to find the best parameter configuration Θ^* , which can efficiently guide the auto-tuning search process.

The on-device measurements especially for embedded/mobile devices are extremely time-consuming. For example, the time consumption for on-device measurements of a VGG16 model can be up to 10 hours on an NVIDIA TX2. In Section 3.3, we design a cross-device domain adaption method, which transfers the cost model trained on the source device to the target device. The objective of the cross-device domain adaption is to find a new parameter set Θ^\dagger which can minimize the difference between the real-world measurements and cost model predictions, and thus to guide the tuning process for each subgraph on the target device in an adaptive manner.

3.3 Feature Space Representations

In our problem, the feature space is independent of hardware architecture while the outputs of the cost model (throughput predictions for different tensor program configurations) are dependent, shown as:

$$H\{\mathcal{X}\} \equiv H\{\mathcal{X}_{DIV}\} + H\{\mathcal{X}_{DV}\} \quad (3)$$

where H maps the hidden feature space representations, \mathcal{X}_{DIV} and \mathcal{X}_{DV} present the decoupled feature space: hardware-independent and hardware-dependent information, respectively. The SOTA *Tenset* [Zheng *et al.*, 2021] shows that the model fine-tuning on the target domain (devices) is effective only when the architecture difference between two devices is small. However, Tenset does not provide any specific solutions, e.g., the cost model domain adaptation between a server level GPU and a mobile level GPU. In order to solve this problem, we propose to leverage the following lottery ticket hypothesis in this paper: only part of the parameters in the trained cost model on source devices are essential for learning the hardware-independent knowledge. In other words, Only part of the information needs to be adapted to the target device while other parameters tend to fit the domain. Motivated by this, we can transform our problem into the following question: how to learn domain invariant parameters to minimize the domain discrepancy brought by the hardware difference?

3.4 Lottery-Ticket-Based Cross-Device Adaptation

The training data for cost model updating is collected online during the search, which makes the search very time-consuming due to inevitable on-device measurements. We denote the set of labeled tensor program records (program knobs, throughput) on source device by $S = \{(x_s^i, y_s^i)\}^m$, where m is the training data points on source device and (x_s^i, y_s^i) is the i_{th} record and its corresponding label respectively in domain s . Now given a target unlabeled program performance records set $T = \{(x_T^j)\}^k$ with different configurations of program knobs x sampled from target device D_T , in which a small set of $\hat{T} = \{(x_{\hat{T}}^j, y_{\hat{T}}^j)\}^n$ can be collected by on-device measurements (which are conducted typically under a given time budgets due to the time-consuming nature). We have $n \ll k \ll m$. Formally, The expected domain adaptation error on target device can be defined as $\varepsilon(h_{Target}(\Theta^\dagger))$ where $h \in H$, a hypothesis that learned from \hat{T} , thus the following inequity holds:

$$\begin{aligned} \varepsilon(h_{Target}(\Theta^\dagger)) &\leq \varepsilon(h_{Source}(\Theta^*)) \\ &+ dist(\mathcal{D}_S(\mathcal{X}), \mathcal{D}_T(\mathcal{X})) + \varphi, \end{aligned} \quad (4)$$

In the terminology of semi-supervised domain adaptation, $dist$ represents the distribution discrepancy over the cross-device domains, and φ represents the ideal error or risks achieving cross-device domain adaptation. Since the feature representations are influenced by hardware difference in our problem, to achieve the cross-device domain adaptation, instead of minimizing the distance between feature representations and their resulting data distribution discrepancy, we show the effectiveness of bound minimization strategy on solving the cross device domain adaptation and propose to find the bound limitation by optimizing the labeling black-box functions. Such an approach is inspired by the lottery ticket hypothesis, which was originally proposed in the context of model compression, showing that only part of parameters are fit for model generalization [Frankle and Carbin, 2019].

Similar as in [Frankle and Carbin, 2019], we show experimentally in this paper the same hypothesis holds in our problem (see Section 4). That is, there exists a super-subnet, named winning ticket, with a set of essential parameters of the trained cost model on source devices, which would be the domain invariant information. In other words, training from a super-subnet on the target device would achieve the optimal transfer performance in our cross-device domain adaptation problem.

We name parameters in these super-subnets as transferable parameters, and the remaining set of parameters as untransferable parameters. A key question is how to distill these transferable parameters during each subgraph auto-tuning stage. We identify the distilling boundary criterion $\xi(ph)$ as below:

$$\xi(i) = |w(ph) * \nabla w(ph)| \quad (5)$$

where we denote the tuning phase of the subgraphs by ph . $w(ph) \in W(ph)$ represents the parameter weights and its gradient is $\nabla w(ph)$. If $\xi(ph)$ is larger than a certain threshold ϑ (e.g., 0.5), the corresponding $w(ph)$ can be viewed as transferable parameters. In contrast, $w(ph)$ would be regarded

as domain-variant parameters if $\xi(ph)$ is small (e.g., close to zero). We also provide a ranking mechanism here, specifically, we rank these parameters based on their cross-domain importance according to Eq.(5). Thus the users can set the transferable parameters ratio manually. We iteratively update the boundary of domain-invariant parameters as well as variant parameters and update these invariant ones during each online training epoch to learn invariant representations to achieve minimization of bound limitation φ . Lastly, we adopt the adversarial loss training objective function $L_{inv}(S, T)$ defined as below:

$$\begin{aligned} L_{inv}(S, T) &= L_{x \sim p(S)}(\lg(b(w(ph, inv)))) \\ &+ \beta L_{x \sim p(T, \tilde{T})}(\lg(1 - b(w(ph, inv)))) \end{aligned} \quad (6)$$

where $b()$ represents the labeling black-box functions, β as the coefficient that control the entropy effects (usually set as a small number). In such a way, the parameters with higher gradient flows, representing more beneficial to the domain-invariant information learning process, are considered. As for the parameters that represent domain variant information, We update these parameters with a weight decay mechanism as penalty, which is defined as:

$$w_v(ph + 1) \rightarrow w_v(ph) - \alpha(wd(w_v(ph))) \quad (7)$$

where α is the learning rate, $w_v(ph)$ represents the domain-variant parameters and function $wd()$ represents the weight decay process for each updating phase.

3.5 Adaptive Tuning Data Partition

To gather on-device measurements efficiently and maintain the performance of the online domain adaptation cost model, we use an adaptive controller (AC) module to early terminate the hardware tuning data collection stage. The basic idea of AC is to statistically analyze the certainty of online training cost model. For a give subgraph s which is to be tuned, we initially divide the total tuning tasks into *tuning tasks for online training* t_{train} with hardware measurement data collection and *tuning tasks for cost model predictions* t_{pred} with a ratio of p . We further divide t_{train} into $q \in 1, 2, \dots, q$ batches and collect both on-device measurement records $C(t_{train}(s))$ and $Perf(t_{train}(s))$. Then, we use the coefficient of variations (the standard deviation divided by the mean) formulated as $CV = \frac{\sigma(C(t_{train}(s))_1, C(t_{train}(s))_2, \dots, C(t_{train}(s))_q)}{\mu(C(t_{train}(s))_1, C(t_{train}(s))_2, \dots, C(t_{train}(s))_q)}$ to dynamically estimate the certainty of the existing cost model, and terminate the hardware measurement phase in advance if the value is smaller than a certain value. We empirically set these hyper-parameters based on multiple trials in our experiments.

3.6 Putting Everything Together

In previous sections, we described the design details of the lottery-ticket based cross-device cost model transfer and the adaptive online training data partition mechanism. We now put these components together and summarize the working flow of Moses.

Step 1. Pre-training a cost model on the source device: We pre-train a cost model offline using the dataset of on-device measurement records from the source device. This dataset includes randomly generated tensor programs for widely deep

learning models.

Step 2. Transferring the trained model to the target device: The learned cost model from source device is directly transferred to the target device in this step, to guide the search stage during auto-tuning.

Step 3. Adaptive training data partition with the AC module: For each tuning task, we dynamically control the hardware measurement costs by using the AC module, with which the portion of on-device measurements can be adjusted by the evaluations of cost model performance in that epoch.

Step 4. Online updating the cost model with iterative pruning: For each tuning task, we divide the parameters of the cost model into domain-invariant ones and domain-variant ones based on the calculation of $\xi(i)$, and update the domain-invariant parameters with gradient decent while letting the rest gradually decrease to zero due to weight decay. During the auto-tuning process, Moses keeps updating the cost model in an adaptive and iteratively manner based on the collected hardware measurements records, while the search algorithms keep querying the newest cost model for efficient explorations of optimal program configurations.

4 Experiments

In previous sections, we describe how we enable the cross-device cost model domain adaptation and make the auto-tuning process on a new device more adaptive. In this section, we evaluate the effectiveness of Moses, with our proposed lottery-ticket based cost model adaptation method. We implement Moses as a plug-in cross-device cost model adaptation tool in TVM auto-tuning [Chen *et al.*, 2018a]. Specifically, the cost model fine-tuning is integrated with the tensor programs random sampling and an evolutionary search algorithm [Zheng *et al.*, 2020a]. The training of the cost model is implemented in PyTorch. We set the max epoch to 30. We set the initial learning rate α to 0.001, the distilling boundary criterion threshold ϑ to 0.5.

4.1 Generated Dataset for Embedded Devices

As mentioned before, learning-based tensor compilers can greatly boost DNN model inference performance. At the core of the auto-tuning process in these compilers, is a cost model which estimates the performance of the combinations of tensor representation knobs on different devices with an input DNN model. To perform the evaluation and ease the training efforts on cost models especially on embedded devices. In this paper, we collect a comprehensive tensor program dataset for two embedded devices: NVIDIA Jetson TX2 and XAIVER. We collect tasks from over 50 DNN models including the popular mobile transformers (e.g. mobileViT [Mehta and Rastegari, 2021]). More than ten million of program records are included in this dataset.

4.2 Experimental Settings

Our experiments are conducted on NVIDIA GeForce GTX 2080 and NVIDIA Jetson TX2 with Pascal GPU architecture with 256 NVIDIA CUDA cores. The current deployments such as [Zheng *et al.*, 2020a; Zheng *et al.*, 2020b; Li *et al.*, 2020] optimize the input DNN models in a sub-graph basis, which means to optimize the fused operators one

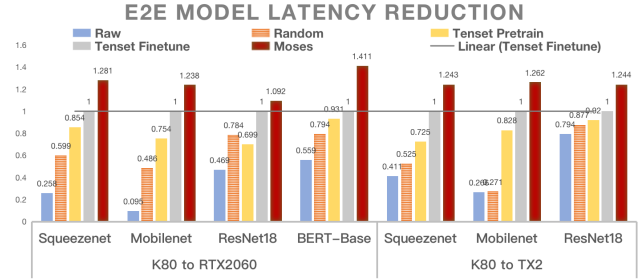


Figure 4: End-to-end DNN model inference latency reductions GAIN comparisons among MobileNet, ResNet18, BERT-Base and SqueezeNet over two domain adaptation baselines.

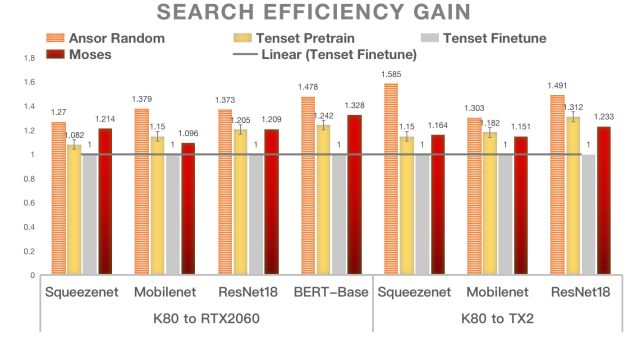


Figure 5: Auto-tuning search efficiency GAIN comparisons among MobileNet, ResNet18, BERT-Base and SqueezeNet over two domain adaptation baselines.

by one. Thus the processes are independent of the DNN model, which means we do not need to include multiple complicated DNN models to validate Moses. The current popular operators of DNN models for both academia and industry can be summarised as: convolutional layers, depthwise separable convolutional layers; multi-head attention module; residual block and other types of layers such as fully connected layers and pooling layers. We include four DNN models in our experiments: ResNet-18, MobileNet, BERT-base and SqueezeNet. We use the default settings for other hyperparameters provided by Ansor [Zheng *et al.*, 2020a]. As for the backbone of the cost model, we choose the representative one used in Ansor, which is an MLP with two hidden layers, with 512 neurons for each. We train the MLP cost model with ranking loss on NVIDIA Tesla K80, which is noted as the source device (domain), based on the dataset provided by Tenset [Zheng *et al.*, 2021]. The two main domain adaptation tasks we validate are $K80 \rightarrow 2060$ and $K80 \rightarrow TX2$.

4.3 Evaluation Metrics

We use the end-to-end latency/throughput and the end-to-end search efficiency of auto-tuning as two main metrics. Specifically, we measure the obtained speedups of tuned tensor programs and the reductions of searching time of an input DNN model over other baselines including the state-of-art cost model transfer method provided in Tenset. We also introduce a concept named Cost Model & Auto-tuning Efficiency Gain Score (CMAT) to evaluate the cost model influence on

Table 1: A summary of comparisons of CMAT under small and large trials. S, R, M and B refer to four DNNs mentioned in Fig. 5.

CMAT (%)	2060-S	2060-R	2060-M	2060-B	TX2-S	TX2-R	TX2-M
Small Trials (200)	57.2	19.6	105	66.7	28.7	66.4	64.5
Large Trials (20000/5000)	48.1	32.7	45.8	87.4	44.7	53.1	45.9

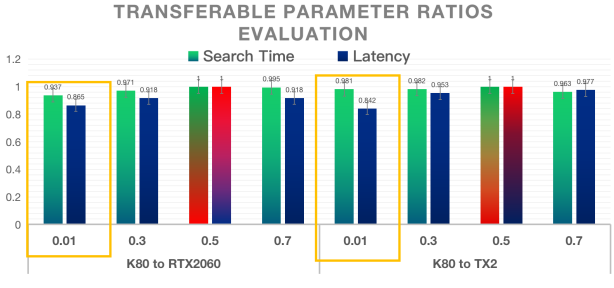


Figure 6: An illustration of Moses performance with a wider ratios of transferable parameters, the yellow box shows the results of ratio=0.01.

the end-to-end inference performance at the same time, defined as: $CMAT = (Gain\ on\ Search\ Efficiency * Reduction\ on\ Tuned\ Model\ Latency - 1) * 100\%$. As CMAT considers both search efficiency and end-to-end inference latency, it is an effective metric for evaluating the overall cross-device cost model adaptation performance.

4.4 Results

The main results we provided aim to answer the following questions: 1). Can Moses transfer the trained cost model to different hardware platforms and outperform other online fine-tuning method provided by Tenset? 2). Can Moses accelerate the auto-tuning process by adaptively scheduling the portions of on-device measuring records trials? To answer these questiones, we compare Moses with four baselines:

- 1) Raw: inference with original CUDA acceleration.
- 2) Ansor-Random [Zheng *et al.*, 2020a]: randomly initialize the cost model and train it from scratch during the auto-tuning.
- 3) Tenset-Pretrain: pre-train a cost model on TenSet dataset and directly apply it to the target device without fine-tuning.
- 4) Tenset-Finetune: utilize the cost model pre-trained on TenSet dataset and then perform the vanilla online fine-tuning.

Inference Time & Search Efficiency Comparisons.

Fig. 4 shows the comparison of the final end-to-end inference time of the input DNN model optimized on each baseline. Moses achieves up to 41.1% faster inference speed over Tenset-Finetune and up to 53% higher speed over Tenset-Pretrain on the $K80 \rightarrow 2060$ baseline. Moses also achieves up to 26.2% over Tenset-Finetune and up to 52% over Tenset-Pretrain on the $K80 \rightarrow TX2$ baseline, respectively. Overall, Moses yields the best inference performance among all other configurations and algorithms. Fig. 5 shows the auto-tuning search efficiency gains comparisons over these baselines. Moses also outperforms all other baselines for both

$K80 \rightarrow 2060$ and $K80 \rightarrow TX2$ settings. It can also be observed that, for some input DNN models such as SqueezeNet and MobileNet, Ansor-Random and Tenset-Pretrain could be more efficient than Moses. This is because these baselines provide no online learning during the auto-tuning. Therefore, the corresponding end-to-end model inference latency of these models can be greatly lower than Moses. The evaluation results show that, the search efficiency gain of the $K80 \rightarrow 2060$ setting can be up to 47.8% while up to 58.5% for the $K80 \rightarrow TX2$ setting. This is because the on-device data collection costs on TX2 is much higher than on RTX2060.

CMAT Score Comparisons.

Table. 1 shows the superior comprehensive performance of Moses over both small (200) and large (20000 for 2060, 5000 for TX2) number of trials across all input DNN models. As mentioned above, although Tenset achieves 15% auto-tuning efficiency gain on MobileNet based on the $K80 \rightarrow 2060$ setting, which is better than Moses (9.6%), the corresponding CMAT is -14.75% over Tenset fine-tuning, which is much worse than Moses (up to 45.8%). We can observe that for some cases (e.g. 2060-S), the CMAT gain under the small-trial setting can even be better than the large-trial one, due to the characteristics of the heuristic searching algorithm embedded in the auto-tuning component in TVM. The base performance of Tenset-Finetune can be extremely low with no prior knowledge during the transfer process of the cost model.

4.5 Ablation Study: Ratio of Transferable Parameters.

Here we provide more analysis on the ratio of transferable parameters. Fig. 6 shows the results of Moses on a wider setting of transferable parameters ratio: $\{0.01, 0.3, 0.5, 0.7\}$. According to the end-to-end performance results, we can observe that the optimal performance can be around 0.5. Generally speaking, the *std* value for settings of $\{0.3, 0.5, 0.7\}$ ratio is not large, which suggests that, the optimal performance produced by different ratios is not sensitive to the ratio setting when it is ranging from 0.3 to 0.7.

5 Conclusion & Future Work

We present Moses, a new framework to optimize the auto-tuning process in DNN compiler. Specifically, our approach enables cross-device domain adaptation of a trained cost model by updating the domain invariant parameters during online learning, which greatly improves the efficiency of auto-tuning process and the end to end throughput of tuned tensor programs on the target device. Besides, we generate a large-scale program performance dataset on two embedded GPUs for learning based DNN compilers. Our future work includes extending Moses to support knowledge transfer from the cross-subgraph tensor optimization perspective.

References

- [Abadi *et al.*, 2016] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek Gordon Murray, Benoit Steiner, Paul A. Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zhang. Tensorflow: A system for large-scale machine learning. *CoRR*, abs/1605.08695, 2016.
- [Ahn *et al.*, 2020] Byung Hoon Ahn, Prannoy Pilligundla, Amir Yazdanbakhsh, and Hadi Esmaeilzadeh. Chameleon: Adaptive code optimization for expedited deep neural network compilation. *CoRR*, abs/2001.08743, 2020.
- [Baghdadi *et al.*, 2021] Riyadh Baghdadi, Massinissa Merouani, Mohamed-Hicham Leghettas, Kamel Abdous, Taha Arbaoui, Karima Benatchba, and Saman P. Amarasinghe. A deep learning based cost model for automatic code optimization. *CoRR*, abs/2104.04955, 2021.
- [Bai *et al.*, 2021] Yang Bai, Xufeng Yao, Qi Sun, and Bei Yu. Autogtco: Graph and tensor co-optimize for image recognition with transformers on gpu. 2021.
- [Chen *et al.*, 2018a] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang, Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. TVM: An automated end-to-end optimizing compiler for deep learning. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, pages 578–594, Carlsbad, CA, October 2018. USENIX Association.
- [Chen *et al.*, 2018b] Tianqi Chen, Lianmin Zheng, Eddie Q. Yan, Ziheng Jiang, Thierry Moreau, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. Learning to optimize tensor programs. *CoRR*, abs/1805.08166, 2018.
- [Frankle and Carbin, 2019] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019.
- [Haj-Ali *et al.*, 2020] Ameer Haj-Ali, Hasan Genc, Qijing Huang, William S. Moses, John Wawrzynek, Krste Asanovic, and Ion Stoica. Protuner: Tuning programs with monte carlo tree search. *CoRR*, abs/2005.13685, 2020.
- [Han *et al.*, 2021] Zhongyi Han, Haoliang Sun, and Yilong Yin. Learning transferable parameters for unsupervised domain adaptation. *arXiv preprint arXiv:2108.06129*, 2021.
- [Jia *et al.*, 2019] Zhihao Jia, Oded Padon, James Thomas, Todd Warszawski, Matei Zaharia, and Alex Aiken. Taso: Optimizing deep learning computation with automatic generation of graph substitutions. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles, SOSP '19*, page 47–62, New York, NY, USA, 2019. Association for Computing Machinery.
- [Kaufman *et al.*, 2020] Samuel J Kaufman, Pithchaya Mangpo Phothilimthana, Yanqi Zhou, Charith Mendis, Sudip Roy, Amit Sabne, and Mike Burrows. A learned performance model for tensor processing units. *arXiv preprint arXiv:2008.01040*, 2020.
- [Li *et al.*, 2018] Tzu-Mao Li, Michaël Gharbi, Andrew Adams, Frédo Durand, and Jonathan Ragan-Kelley. Differentiable programming for image processing and deep learning in halide. *ACM Trans. Graph.*, 37(4), July 2018.
- [Li *et al.*, 2020] Menghao Li, Minjia Zhang, Chi Wang, and Mingqin Li. Adatune: Adaptive tensor program compilation made efficient. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 14807–14819. Curran Associates, Inc., 2020.
- [Mehta and Rastegari, 2021] Sachin Mehta and Mohammad Rastegari. Mobilevit: Light-weight, general-purpose, and mobile-friendly vision transformer. *CoRR*, abs/2110.02178, 2021.
- [Ryu and Sung, 2021] Jaehun Ryu and Hyojin Sung. Metatune: Meta-learning based cost model for fast and efficient auto-tuning frameworks. *CoRR*, abs/2102.04199, 2021.
- [Zhao *et al.*, 2021] Zhihe Zhao, Kai Wang, Neiwen Ling, and Guoliang Xing. Edgectl: An autotml framework for real-time deep learning on the edge. In *Proceedings of the International Conference on Internet-of-Things Design and Implementation, IoTDI '21*, page 133–144, New York, NY, USA, 2021. Association for Computing Machinery.
- [Zheng *et al.*, 2020a] Lianmin Zheng, Chengfan Jia, Minmin Sun, Zhao Wu, Cody Hao Yu, Ameer Haj-Ali, Yida Wang, Jun Yang, Danyang Zhuo, Koushik Sen, Joseph E. Gonzalez, and Ion Stoica. Ansor: Generating high-performance tensor programs for deep learning. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, pages 863–879. USENIX Association, November 2020.
- [Zheng *et al.*, 2020b] Size Zheng, Yun Liang, Shuo Wang, Renze Chen, and Kaiwen Sheng. Flextensor: An automatic schedule exploration and optimization framework for tensor computation on heterogeneous system. In *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS '20*, page 859–873, New York, NY, USA, 2020. Association for Computing Machinery.
- [Zheng *et al.*, 2021] Lianmin Zheng, Ruochen Liu, Junru Shao, Tianqi Chen, Joseph E. Gonzalez, Ion Stoica, and Ameer Haj Ali. Tenset: A large-scale program performance dataset for learned tensor compilers. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021.