

Agile Autotuning of a Transprecision Tensor Accelerator Overlay for TVM Compiler Stack

Dionysios Diamantopoulos¹, Burkhard Ringlein¹, Mitra Purandare¹, Gagandeep Singh², and Christoph Hagleitner¹

¹IBM Research Europe, Säumerstrasse 4, 8803 Rüschlikon, Switzerland

²Eindhoven University of Technology, Netherlands

Abstract—Specialized accelerators for tensor-operations, such as blocked-matrix operations and multi-dimensional convolutions, have been emerged as powerful architecture choices for high-performance Deep-Learning computing. The rapid development of frameworks, models, and precision options challenges the adaptability of such tensor-accelerators since the adaptation to new requirements incurs significant engineering costs. Programmable tensor accelerators offer a promising alternative by allowing reconfiguration of a virtual architecture that overlays on top of the physical FPGA configurable fabric. We propose an overlay (τ -VTA) and an optimization method guided by agile-inspired auto-tuning techniques. We achieve higher performance and faster convergence than state-of-art.

Index Terms—Neural Networks, Machine Learning, Autotuning, FPGA, Transprecision Computing, Tensor Accelerator

I. INTRODUCTION

Deep Learning (DL), a powerful set of techniques for learning in neural networks, has achieved unprecedented accuracy in numerous aspects of the digital transformation of our society. This fascinating biologically-inspired programming paradigm, which enables a computer to learn from observational data, attributes its success to a large volume of trained parameters, which, however, can contain a lot of redundant information [13]. Prior art has maintained remarkable levels of accuracy, by applying pruning techniques [11][13][31], or by using sparsification [4], or both [14][15]. A highly effective technique exploiting redundant information is moving from floating-point arithmetic to low-precision integer arithmetic [5]. Recently, transprecision computing was proposed as a paradigm shift in precision selection, and suggests the adaptability of precision according to an application’s requirements [20], as opposed to the conservative static selection of low precision processing.

DL systems rely on hardware (HW) accelerators and manually optimized, high-performance libraries to increase computational efficiency, i.e., achieving the maximum throughput while consuming the smallest possible amount of resources and energy [18], [23]. While GPUs are dominating the training and inference computations at scale, FPGAs can achieve more than $10\times$ better speed and energy efficiency than state-of-the-art GPUs [12], [21]. In addition, the inherent programmability of FPGAs at bit-level makes them ideal accelerators for ultra low-precision inference and training [29], [19].

To optimize a neural network, programmers must choose from many implementations that are logically equivalent but differ dramatically in performance due to differences in threading, memory reuse, pipelining and other HW factors.

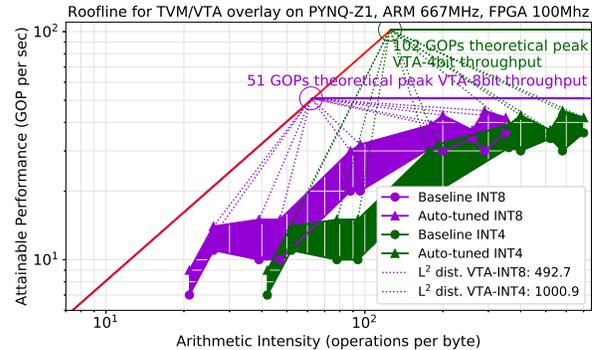


Fig. 1: Auto-tuning a model for an overlay FPGA back-end does not necessarily lead to higher performance, when the precision is decreased (e.g. INT8 to INT4). Investigating the reason and proposing path-forward motivates this work.

Supporting diverse HW back-ends therefore incurs significant engineering cost. Depending on whether it is an FPGA, a GPU or an ASIC, the accelerator’s adaptability to different models heavily relies on the software-supported libraries that bridge the gap between the programming language semantics and the HW supported intrinsics. A recent study proposed the use of a statistical cost model that predicts program run time by using a given low-level program [8]. The cost model guides the exploration of the space of possible programs.

While auto-tuning has been proved to be effective in automatic optimization, it typically assumes a fixed HW that offers tunable knobs in the software, like tiling, loop reordering, etc. However, in the case of FPGAs, the HW design space also provides both tunable micro-architecture choices, e.g. pipelining, and tunable implementation choices, e.g. device frequency. Such HW tunable choices can lead to sub-optimal auto-tuning when only a subset of them is considered for bitstream generation, particularly in low-precision DL for FPGAs [5].

To illustrate this problem, we examine the auto-tuning of an FPGA overlay developed by the community-driven TVM DL compiler stack [1], specifically VTA [22]. This overlay employs a GEMM accelerator as a computation engine for convolution operations. We configured a GEMM of 8 bits (W8A8) and 4 bits (W4A4), for weights and activations, respectively. Figure 1 depicts the roofline model [30] of VTA, overlaid on a PYNQ-Z1 device [2]. The plot shows the throughput achieved on different convolution layers of the ResNet-18 inference benchmark. Each layer has a different

arithmetic intensity, i.e. compute to data movement ratio. In the left half of the plot, convolution layers are bandwidth-limited, whereas on the right half, they are compute-limited. The operation points on the roofline, depicted as circles in Figure 1, are “optimal” since in those points neither performance nor communication is under-utilized. The goal behind designing HW architectures and compiler stacks is to bring each workload as close as possible to the roofline of the target HW and ideally to these “optimal” points.

As an experiment, we auto-tuned the VTA, using TVM’s auto-tuning flow, configured with a GEMM of W8A8 and W4A4. When VTA is configured with a 4-bit GEMM intrinsic, the theoretical performance is doubled to 102GOPS. In addition, the arithmetic intensity is doubled since half the number of bytes have to be fetched from the main memory. Please note that this is a strong advantage of overlay architectures on FPGAs [5]. However, the case of W4A4 delivers a measured performance identical to that of VTA W8A8, which wastes 100% of the theoretical possible speedup. In addition the Euclidean distance (L^2) of all auto-tuned convolutions is higher for W8A8 than W4A4 from the perspective of the “optimal” operation points for VTA of 8 bit and 4 bits. This shortcoming is attributed to the current TVM compiler stack that relies on the user for the selection of the HW parameters of the VTA overlay for different precision. In addition, the VTA auto-tuning, using TVM’s current auto-tuning flow, focuses only on software-related optimization options, assuming a fixed VTA design on the FPGA. While this assumption enables the interoperability of TVM’s auto-tuning flow to silicon-proven devices, such as GPUs, DSPs and custom ASICs, it neglects some important features of reconfigurable HW.

In this paper, we explore the following question: **Can we automatically guide the auto-tuning of a tensor accelerator overlay, for different precision settings, by leveraging knowledge from hardware design experience?** Our affirmative answer is based on a framework that makes the following contributions:

- ◊ *Engineering aspect:* We adopt the concept of agile development to the community-driven TVM github repository, by bringing-up a pipeline of engineering tasks that extends the autotuning process. When this step is finished, the best explored models can be uploaded to benefit the community.

- ◊ *Scientific aspect:* Instead of eliminating the overlay hardware design space with pruning techniques (e.g. successive halving, used in [22]), we propose a technique that builds a prediction model that quantifies the impact of a hardware design choice (feature) towards an optimization goal, e.g. increasing performance. By employing a classifier with an arbitrary differentiable loss function, we show that the features with the highest impact differ for different precisions. We further propose the use of the most important features in order to generate an overlay and then continue with auto-tuning.

In Section II we present the design-automation engineering contribution, i.e. the agile auto-tuning methodology. To demonstrate the impact of this new integral development concept, we propose the τ -VTA optimization in Section III. Experimental results are presented in Section IV and Sec-

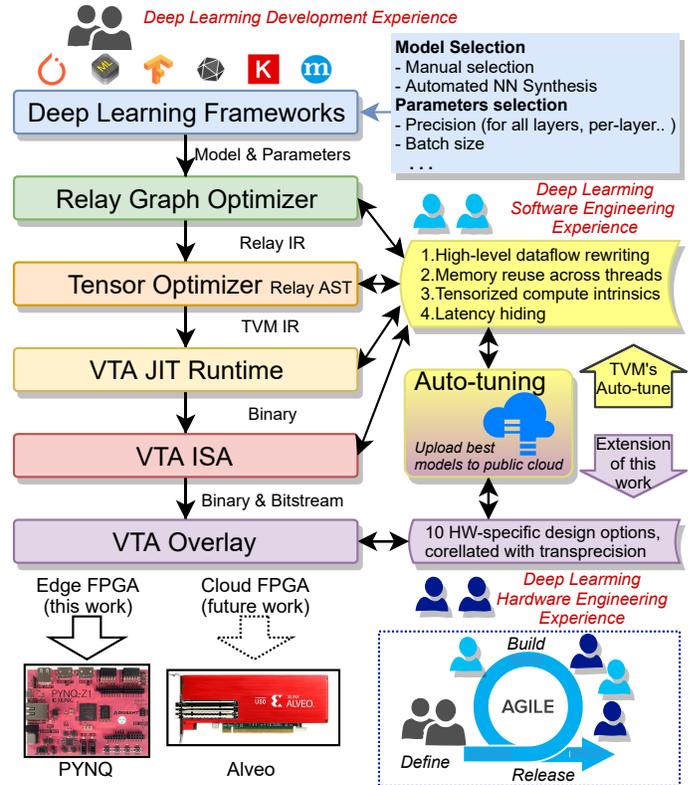


Fig. 2: Overview of the proposed approach: the auto-tuning step of the TVM toolflow for FPGAs considers also options from the implementation phase of the VTA tensor accelerator overlay.

tion VI concludes the paper.

II. AGILE AUTOTUNING METHODOLOGY

A. Integration to TVM

To handle the expanding ecosystem of DL Frameworks on the one side and specialized DL HW on the other, a group from the University of Washington proposed *TVM*, a full open source compiler-stack [7], [22], [24] that aims to “close the gap between the productivity-focused deep learning frameworks, and the performance- or efficiency-oriented hardware backends” [1]. *TVM* is built using multiple Intermediate Representation (IR) languages and therefore offers multiple layers for optimizations, as seen on the left-hand side of Figure 2. It first has multiple modules to import state-of-the-art DL Frameworks (like pytorch, Tensorflow, Keras) to the RelayIR [24]. Afterwards, it is able to perform multiple optimizations on the Abstract Syntax Tree (AST) generated out of RelayIR and lower the program to the TVM IR. The TVM IR can then be interpreted by a runtime or another compiler to finally execute the DL task on the target HW [22].

One special contribution of *TVM* is, besides the large number of supported Frameworks and HW, to perform the so-called “auto-tuning” of tensor operations in order to maximize performance [8]. During auto-tuning, a number of known optimizations are performed on the Relay AST to improve the scheduling of the arithmetic operations. That way, the auto-tuning doesn’t change the actual mathematical instructions – like ALU or GEMM operations would – of the program, but the execution order and memory accesses [8]. *TVM*’s current auto-tuning supports numerous optimizations across

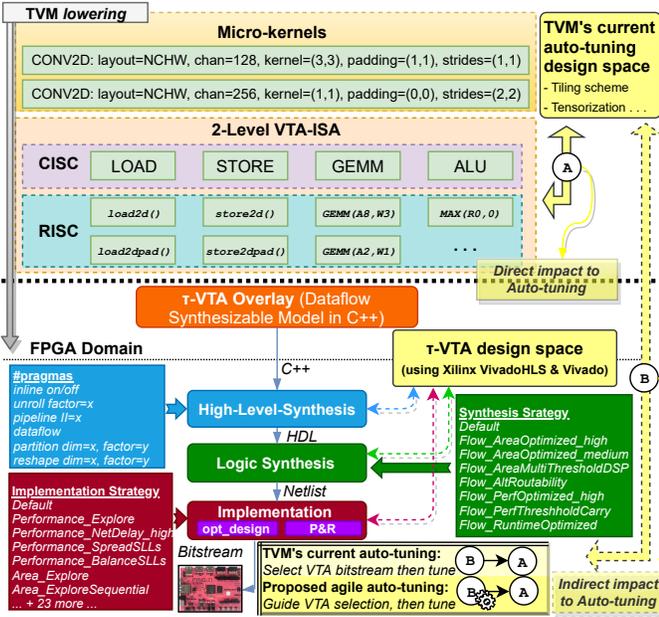


Fig. 3: The design space of τ -VTA.

the lowering from DL code to an executable, categorized in high-level dataflow rewriting, memory reuse across threads, tensorized compute intrinsics and latency hiding. These are analytically presented in [17].

We propose that HW-SW interoperability be adopted to support the VTA auto-tuning in an agile way, inspired by the positive disruption of the Hardware Agile Manifesto [16]. As depicted in Fig. 2, the processes of customizing the VTA overlay and optimizing the software in multiple IRs, are combined in a united auto-tuning task. With this approach we aim to extend TVM’s current auto-tuning flow with HW design options for VTA in an automatic way. Most importantly, the intention of our approach is to select the most important HW options during the auto-tuning, particularly the ones correlated with different precision of VTA. In a nutshell, the “agile” concept is to provide a development environment where knowledge flows easily so the best solutions are reached as quickly as possible. Fig. 2 conceptualize the steps, where experience from DL development, DL software engineering and DL hardware engineering can be combined by spawning pipelines of tasks (i.e. auto-tuning sprints) on multiple systems (i.e. auto-tuning fleet) that share libraries, tools and devices in order to decrease auto-tuning time. This step will be further explained in subsection II-C.

B. Introducing the design space of τ -VTA

Figure 3 gives a high-level overview of the VTA hardware organization, as firstly presented in [22]. In addition, the figure highlights the main contribution of this work, i.e. the extension of TVM’s auto-tuning design space with design options from HW expertise. VTA is composed of four modules: *FETCH*, *LOAD*, *COMPUTE*, and *STORE*, which communicate over command queues and on-chip shared memories, implemented using FPGA block-RAMs (BRAMs). The *FETCH* module dispatches task instructions after loading them from the DRAM. The *LOAD* and *STORE* units trans-

fer tensor-tiles from DRAM into on-chip FPGA memories (BRAMs) and vice versa, respectively. The *COMPUTE* unit performs operations on the register file. Specifically, the *tensor ALU* performs element-wise tensor operations such as activation, normalization, and pooling tasks, while the *GEMM core* performs matrix multiplication over input and weight tensors. Common deep learning *Micro-kernels*, such as 2D convolutions, are executed in *GEMM core* [22].

The VTA architecture is parameterizable, so that different shapes can be configured for tensor intrinsics, depending on the available resources. For example, the shape of input, weight, and accumulator tensors that feed the GEMM tensor intrinsic unit directly affects the utilization of multipliers and the width of BRAMs ports. The data-types of the tensors are also parameterizable e.g. 8 bits or fewer. Programmability of the VTA is based on a two-level Instruction Set Architecture (ISA): i) a CISC-like task-ISA that explicitly orchestrates concurrent compute and memory tasks and ii) a RISC-like microcode-ISA that implements a wide variety of operators with single-cycle tensor-tensor operations.

VTA is an overlay architecture programmed in synthesizable C++. Using the Xilinx Vivado HLS tool it is synthesized to a hardware description language (HDL) (either VHDL or Verilog) as a register-transfer-level (RTL) model. The downstream implementation stage, with the Xilinx Vivado tool, includes logic synthesis, place & route, optimization (e.g. timing, area, energy) and the generation of a VTA bitstream for the FPGA device.

All of these steps include many design choices that affect the trade-off between performance and resources utilization. In addition to the design options, the customization knobs of VTA define an additional hardware design space with 1000s of individual designs. VTA’s developers explore which candidate to use in a sequence of steps. First, they use a simple FPGA resource model to prune unfeasible VTA parameterizations. After pruning, each candidate hardware design is compiled, placed, and routed. They select three tunable parameters, specifically, FPGA device, precision and batch size. Typically their exploration returns a handful of promising candidates - “the rest of the designs either yield low peak performance or fail placement, routing, or timing closure” [22]. For this final set of designs, they generate optimized software, using operator auto-tuning[8], and use this software to obtain the workload’s performance profile.

While VTA’s optimization uses pruning instead of exhaustively exploring the design space to find the best candidate, the design space is limited to only three parameters. However, an overlay FPGA design can benefit from the knobs of design tools to deliver high performance. Such knobs may have an indirect impact on auto-tuning. For example, an optimal selection of the design options solely in the HLS step can lead to performance improvements of up to $29.030\times$ [9]. Figure 3 shows some important design options that are later discussed in Section III. We propose to find the impact of these parameters on VTA’s performance by using a prediction technique and then guide the exploration based on the most important parameters. Hence, instead of pruning the design space with the hardware experience of designs that fail to meet

Listing 1: The naive implementation.

```

for i in 0..1024:
  for j in 0..1024:
    C[i][j] := 0
    for k in 0..1024:
      C[i][j] += A[k][i] *
                B[k][j]

```

Listing 2: One possible optimization using loop tiling to optimize memory access leveraging the cache, where t_i and t_j are the tiling factors for each dimension.

```

for io in 0..(1024/ti):
  for jo in 0..(1024/tj):
    C[io*ti:io*ti+ti][jo*tj:jo*tj+tj] := 0
    for k in 0..1024:
      for ii in 0..ti:
        for jj in 0..tj:
          C[io*ti+ii][jo*tj+jj] +=
            A[k][io*ti+ii] * B[k][jo*tj+jj]

```

Listing 3: Another possible optimization using the intrinsic operations of VTA accelerator. This example sets the tiling factors t_i and t_j to 8 in order to match the dimensions of the available intrinsic instructions.

```

for io in 0..128:
  for jo in 0..128:
    vta.intr.fill_zero(C[io*8:io*8+8][jo*8:jo*8+8])
for ko in 0..128:
  vta.intr.gemm8x8_add(C[io*8:io*8+8][jo*8:jo*8+8], A[ko*8:ko*8+8]
                    [io*8:io*8+8], B[ko*8:ko*8+8][jo*8:jo*8+8])

```

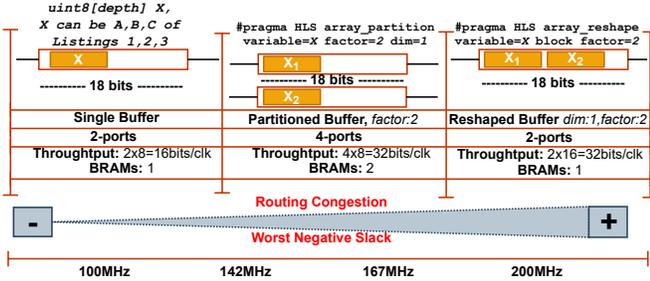


Fig. 4: Exemplary impact of features “Partition” and “Reshape” on memory throughput and BRAMs utilization for a *uint8* array.

the specs, we proactively offload the design space pruning to an algorithmic optimization problem (Section III). To establish a differentiation of the proposed VTA auto-tuning, from the current one in TVM stack, we use the term “ τ -VTA auto-tuning”, accounting for transprecision.

To highlight the difference of the current VTA auto-tuning by the “ τ -VTA auto-tuning”, we present in code Listing 1 an example of a GEMM computation. The Listing 2 is expected to be faster on a target HW that uses caches, while the Listing 3 relies on intrinsic HW specific instructions. During TVM’s auto-tuning phase, one of the optimizations is the selection of different tiling factors that accommodate for specific cache hierarchy or HW intrinsic dimensionality on the HW device (VTA). This selection can be either completely random or algorithmic-driven as shown in [8], where XGBoost [6] and TreeGRU [27] algorithms have shown to find better code-refactoring candidates in shorter time. However this analysis neglects some important features of the reconfigurable HW. For example, if the data-type of the computation is *uint8*, then during the design of VTA we can use the Vivado HLS directives `#pragma HLS ARRAY_PARTITION` and `#pragma HLS ARRAY_RESHAPE` to increase the memory throughput and BRAMs utilization, as shown in Fig. 4, assuming 18Kbits true dual-port Xilinx BRAMs. The case of *uint4* will double the throughput and so on. However, the aggressive use of those directives is known to increase the routing congestion and the difficulty to meet the target frequency. Such trade-offs are known to the HW engineering community, but not in the SW counterpart. The “ τ -VTA auto-tuning” bridges this gap by algorithmically selecting the most important VTA HW parameters, for different precision settings, before initiating the TVM’s current VTA auto-tuning.

C. Collaborative Exploration Pipeline

An end-to-end DNN framework, such as TVM, combines expertise from many levels of the computing stack, i.e. DNN front-end languages, compilers, IRs, scheduling, HW generation, etc. With it being a community-driven framework, it is expected that updates are introduced in a dynamic, non-deterministic time-plan. Automating the integration of code changes from different contributors is a key element for maintaining project stability. In this work we explore the optimal VTA designs for a given precision, based on features of the hardware design space. We propose the use of tools that allow us to enforce the concept of Continues Integration (CI) during this exploration.

Figure 5 shows the overview of this CI exploration pipeline. The proposed CI infrastructure consists of a cluster of machines, named auto-tuning fleet, on which we deploy our exploration jobs. The key technologies we use are Jenkins, Docker and Docker Swarm. **Jenkins** is the tool we use to define, run and manage the CI jobs. Jenkins supports the “Master+Agent” mode, where the Master is in charge of orchestration and acts as the user end-point, and the Agents perform the actual work, i.e. exploration. These agents can be distributed across many servers. Thanks to the Jenkins plug-in that provides integration with Docker, the agents can be even spawned on-demand as containers. **Docker** is the supporting technology for the whole infrastructure. Every component runs as a Docker container, even the Jenkins Master itself. This establishes portability and scalability for the exploration phase. **Docker Swarm** allows a seamless deployment of containers on any of the cluster’s machines in a transparent way. It is responsible for load-balancing and tracking of the status of all running containers, across all of the machines that have joined the Swarm.

The advantage of the proposed collaborative exploration pipeline is that, whenever a community member commits updates of the VTA overlay to the TVM repository (e.g. new GEMM unit, new precision in accumulator etc.), the exploration of the optimal configurations of the VTA design space is triggered automatically. As soon as an optimal VTA bitstream is explored, auto-tuning is triggered as a subsequent step, named auto-tuning sprint (inspired by the agile concept). The best VTA bitstreams and auto-tuning configurations are stored in the local repository and they can be uploaded back to the community’s repository with a merge request.

III. τ -VTA AUTO-TUNING

A. τ -VTA Features

1) *Feature Selection:* To enable a guided VTA design space exploration, by annotating knowledge from HW experience, we build a dataset of HLS and implementation results, consisting of samples across individual designs. We run each design through the complete C-to-bitstream flow for various design options. These options are used in addition to the ones of current TVM’s auto-tuning [22], i.e. fpga device, precision and batch size. The design options are the features of our formulation for predicting an optimization target. In this analysis we explore two main targets, performance, in terms of Giga Operations per second (GOPs) and resources utilization, in

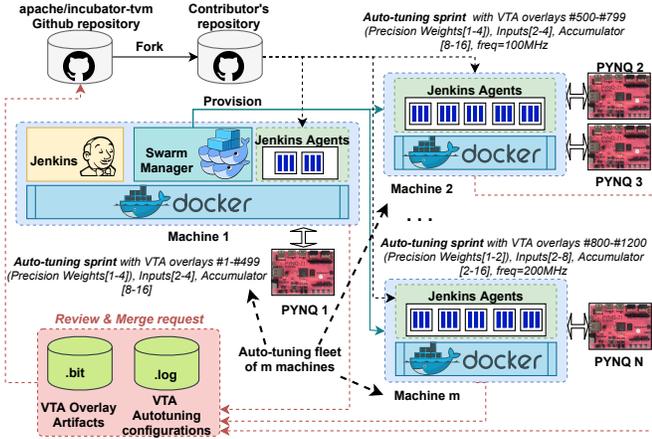


Fig. 5: The exploration pipeline using Docker Swarm and Jenkins.

terms of absolute number of utilized FPGA resources. For the latter target we explore every resource individually, i.e. Block RAMs (BRAMs), Flip Flops (FFs), Digital Signal Processors (DSPs) and Look-Up-Tables (LUTs). Feature selection is a very important step for our flow, so we decide not only to include the design options commonly used by the FPGA community, based on our experience, but also the ones used in research around design space exploration with EDA tools [9], [10]. Table I lists the features we’ve selected in this research. Similarly, we extract implementation results, known as the targets in our optimization problem, from the implementation reports. After extraction, our dataset contains features and targets for each design sample and can be used to develop prediction models that map from features to targets.

TABLE I: Features used in the architecture search

Feature	Description	Search range
HLS_freq	The frequency Vivado HLS is aiming for.	[100:20:500]MHz
Inline	Flatten RTL hierarchy by function inlining.	Enabled/Disabled
Unroll	Unroll a loop with a given factor.	[1,2,4,8,complete]
Pipeline	Pipeline(II=1) design with registers.	Enabled/Disabled
Dataflow	Use task-level pipelining.	Enabled/Disabled
Partition	Distribute the memory contents across multiple BRAMs.	cyclic/block dim=[1,2] factor=[1:32]
Reshape	First Partition, then rejoin BRAMs to decrease utilization.	cyclic/block dim=[1,2] factor=[1:32]
Impl_freq	The frequency during place and route.	[100:20:500]MHz
Syn_strategy	The Vivado synthesis strategy (heuristic).	All available (8)
Imp_strategy	The Vivado implementation strategy (heur.).	All available (32)

2) *Irrelevant Features Elimination*: Some of the design options are correlated. As such, combinations of features that are known a priori to exert little influence on the targets should be eliminated to reduce the dimensionality of the data. Having fewer features leads to simpler models, which require shorter training time, and reduces the chance of over-fitting. Feature elimination can benefit from HW knowledge. However, even if some feature choices are valid from a HW knowledge perspective, they may lead to unfeasible designs. For example, revisiting the example of Figure 4, the large factors of partitioning/reshaping can lead to more congested routing, which in combination with the frequency requirements, may result in unsuccessful timing closure. We eliminate irrelevant

features by formulating them as constraints in the τ -VTA optimization problem.

B. Problem Formulation

1) *Prediction Model for Feature Importance*: We train a classification model to predict the impact of features on the optimization goal of increasing performance while respecting constraints in resources utilization. For our study, we have a set of n training samples $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$, where $\mathbf{x}_i = [x_i^1, x_i^2, \dots, x_i^p]^\top \in \mathbb{R}^p$ is the input vector of feature values for the i th sample, and $\mathbf{y}_i = [y_i^1, y_i^2, \dots, y_i^q]^\top \in \mathbb{R}^q$ is the corresponding vector of target values. p denotes the number of input features (e.g., HLS_freq, Unroll, precision etc.), and q denotes the number of output targets (i.e., actual GOP, LUT, FF, DSP, and BRAM counts after the Vivado-implementation). In addition we define $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top$ to denote feature values for all samples and $\mathbf{y}^k = [y_1^k, y_2^k, \dots, y_n^k]^\top$ to denote values of the target k for all samples.

Each learning task corresponds to one target prediction. We train a separate model f_k for each target k , resulting in a set of mapping functions $\{f_k : \mathbb{R}^p \rightarrow \mathbb{R}\}_{k=1}^q$. We select the gradient tree boosting algorithm to build our prediction model. Specifically, we model the target as the sum of regression trees, each of which maps the features to a score for the target. Target estimation is determined by accumulating scores across all trees. By implementing gradient descent, gradient tree boosting optimizes the loss over the space of regression trees by repeatedly selecting the tree that points in the negative gradient direction. Since we eliminate irrelevant features (Section III-A2), many design options result in unfeasible implementations (e.g. timing violation, resource limitations etc.), which result in very sparse datasets. We propose the use of XGBoost [6], a gradient tree boosting algorithm for sparse datasets that employs a sparsity-aware approximate split finding technique.

The benefit of using XGBoost is that, after the boosted trees are constructed, it is relatively straightforward to retrieve importance scores for each feature and use them for our mapping functions $\{f_k : \mathbb{R}^p \rightarrow \mathbb{R}\}_{k=1}^q$. Importance provides a score that indicates how valuable each feature was in the construction of the boosted decision trees within the model. The more a feature is used to make key decisions in decision trees, the higher its relative importance. This importance is calculated explicitly for each feature in the dataset. For a single decision tree, importance is calculated through the amount that each feature-split-point improves the performance measure by, weighted by the number of observations the node of the tree is responsible for. The feature importances are then averaged across all of the decision trees within the model to obtain the feature importance space \mathcal{F}_i .

2) *τ -VTA Auto-tuning Algorithm*: For a given tensor operator specification, e.g. Listing 1, there are multiple possible low-level program implementations, each with different choices of loop order, tiling size, and other options, e.g. Listing 2 and 3. The problem of τ -VTA auto-tuning is to find which logically equivalent programs exhibit fastest execution for a VTA overlay of selected precision. Towards this goal, we extended TVM’s auto-tuning algorithm with an extra outer

Algorithm 1: Learning to Optimize TP Tensor Programs

Input: VTA precision p
Input: Feature importance space \mathcal{F}_i (discussed in III-B1)
Input: Software transformation space \mathcal{S}_e
Output: Selected VTA overlay o_p^* for precision p
Output: Selected schedule configuration s_p^* for precision p
Data: Database \mathcal{D}_o of run time statistics for VTA overlay o

```
1 while  $k\_trials\_overlays < max\_k\_trials\_overlays$  do
  /* Auto-tuning sprints on a fleet of  $m$  */
2   $M_p \leftarrow$  run  $m$ -parallel jobs of VTA bitstream generation to
   collect candidates in  $o_p^*$ , with higher probability for the
   most important features of  $\mathcal{F}_i$  for precision  $p$ 
3   $\mathcal{D}_o \leftarrow \emptyset$ 
4  while  $n\_trials < max\_n\_trials$  do
   /* Pick the next promising batch */
5    $Q \leftarrow$  run parallel simulated annealing to collect
    candidates in  $\mathcal{S}_e$  for the VTA overlay  $M_p$ 
6    $S \leftarrow$  select  $(1 - \epsilon)b$ -subset from  $Q$  using Random
    Search or TreeGRU [27] optimizer (as in current
    VTA auto-tuning)
7    $S \leftarrow S \cup \{ \text{Randomly sample } \epsilon \text{ candidates.} \}$ 
8   Run  $b$  measurements on VTA configured with  $M_p$ 
9    $\mathcal{D}_o \leftarrow$  exec. time // Update best trials
10   $n\_trials \leftarrow n\_trials + b$ 
11   $s_p^* \leftarrow$  history best schedule configuration for precision  $p$ 
12   $k\_trials\_overlays \leftarrow k\_trials\_overlays + m$ 
13  $o_p^* \leftarrow$  history best VTA overlay for precision  $p$ 
```

loop that iterates over a number of possible optimal VTA overlays. This pool is populated with VTA bitstreams, by giving priority to the configurations of the design features with higher importance per precision target. Algorithm 1 describes this algorithmic behavior. The lines 4-to-9 correspond to TVM’s auto-tuning algorithm [22] and describe the steps needed to obtain an optimal configuration from a software transformation space \mathcal{S}_e , on a target hardware device. Specifically, lines 5-to-8 account for combining quality and diversity when selecting b candidates for hardware evaluation and are analytically explained in [22]. The iteration of these steps over an outer loop (line 1) establishes the search for an optimal VTA overlay for a given precision p . Thus, the difference of Algorithm 1 with the auto-tuning algorithm in [22], is that instead of performing auto-tuning only with one VTA overlay, the Algorithm 1 proposes the auto-tuning in multiple overlays, in parallel. The output is not only an optimal configuration from a software transformation space \mathcal{S}_e , but also a VTA overlay o_p^* for a precision requirement p . The overlays with the highest probability of being optimal are prioritized in this exploration using the feature importance metric, which is described in previous subsection III-B1.

IV. EVALUATION RESULTS

A. Feature Importance

We evaluate the feature importance prediction model with a design space defined by the values listed in Table I. In addition, we use a transprecision vector of $[1, 2, 4, 8]$ bits. In this analysis, transprecision is defined as the flexibility to change precision on-demand by the application. This is a typical requirement in the area of NN architecture exploration [26]. The initial design space includes 797,442,048

evaluations (*design options listed in Table I* $\times [1, 2, 4, 8]$ bits). After the elimination of irrelevant features (Section III-A2), the design space contains 2,440 designs. After training the model described in Section III-B1 with a subset of the design space, we build the prediction model to find the importance of design features for performance and resources. Then, we use this model to select the VTA configurations that will be implemented for a given input vector (precision, freq, etc.).

We implement and train the prediction model in Python leveraging the XGBoost library [6]. All designs in the dataset are synthesized and implemented with Xilinx Vivado 2018.3 (VTA-supported) targeting Pynq-Z1. Experiments are performed on an Intel Xeon E5-2630 processor running at 2.6GHz. The model is able to complete the prediction tasks within seconds, compared to the minutes or hours that each implementation typically incurs. We measured 35 minutes for a single end-to-end implementation with default settings for VTA, Vivado HLS and Vivado. Some design options can lead to implementations that last up to eight hours. The average prediction time for all mapping functions is 2.7 seconds.

For importance classification, we compute the percentage of incorrectly classified samples out of the total number of samples. We randomly select 20% of our data as the testing set and perform a cross-validation by random permutation over 10 iterations on the remaining training/validation set. In each iteration, we randomly select 75% of the training/validation set for training and 25% for validation. The validation set is used for parameter tuning to locate better VTA models for implementation. Figure 6 depicts the feature importance on the HW design space of τ -VTA, for the objective targets a) Performance(GOPs), b) DSPs, c) BRAMs and d) Logic resources (FFs, LUTs). For every objective we plot the importance score F for every feature. The precision feature is plotted as a grouped option. Please note that the plotting is orthogonal to all features and any combination can be used. Since we focus on transprecision overlays in this study, we select an analysis that expresses the precision as a grouping option over the others.

Figure 6 also shows the prediction accuracy. DSPs and BRAMs are generally easier to estimate than logic resources because operations with LUTs and FFs experience more complicated transformations than coarse-grain elements like DSP and BRAMs. In general, XGBoost has an accuracy of 87% for predicting the importance of performance and 72%, 77% and 68% for DSPs, BRAMs, and logic, respectively.

B. τ -VTA Auto-tuning

After the feature importance analysis, we guide the selection of the final set of designs and generate optimized software using operator auto-tuning [8]. Component evaluations were based on convolution workloads in ResNet-18 for ImageNet classification. The accuracy evaluation of ResNet-18 is orthogonal to this analysis since we don’t aim to optimize the accuracy, but to optimize the performance of the convolutions given a predefined precision. Hence, the precision choice is used as input for the optimization of the τ -VTA overlay and not for its accuracy, which is a separate step. Figure 7 compares the performance of convolutions when the VTA designs are selected by using the current state-of-the-art

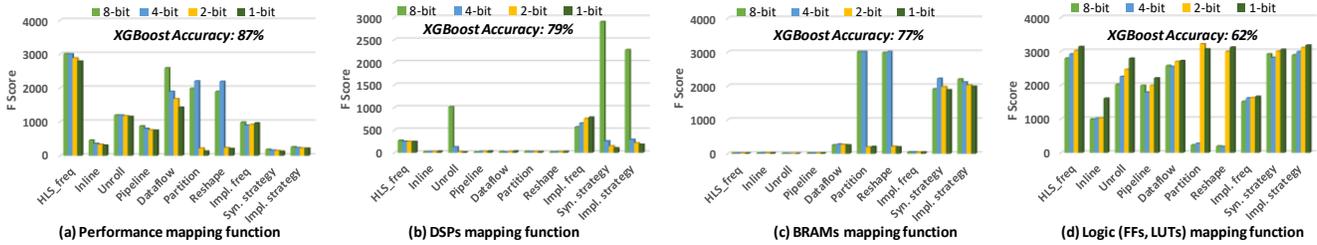


Fig. 6: Features’ impact on the HW design space of τ -VTA, using an XGBoost model. Different features express the diverse impact on performance and resources with respect to precision specification.

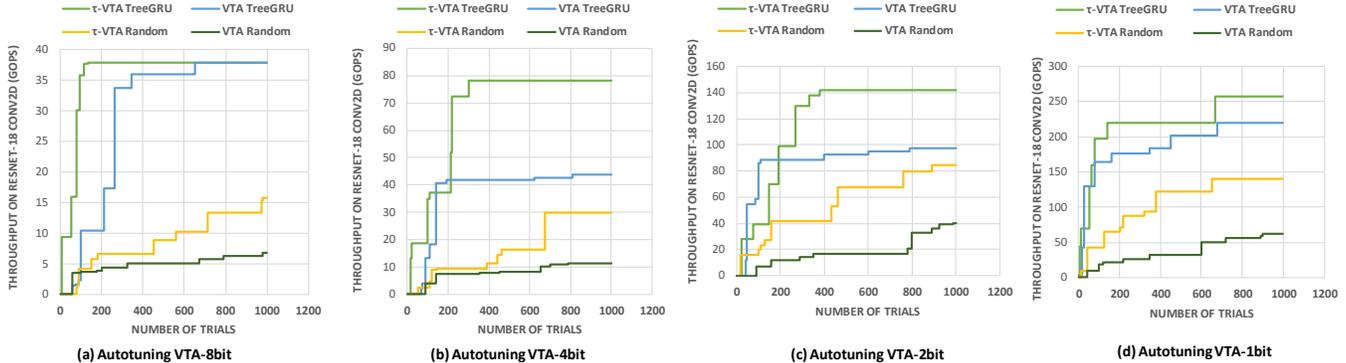


Fig. 7: Schedule exploration with a naïve random search and TreeGRU algorithm for a single convolution layer on PYNQ-Z1. Original VTA and τ -VTA are compared. Design candidates with $(2,16)\times(16,16)$ and $(8,8)\times(8,8)$ GEMM intrinsic at a)W8A8 a)W4A4 a)W2A2 d)W1A1 are considered. The layer is conv2d: IC=256, OC=256, H=W=14, KW=KH=3, stride=(1,1),padding=(0,0).

approach (VTA) or by using the feature importance metric (τ -VTA). We select two auto-tuning methods from [8] which are implemented in TVM, random search and TreeGRU [27]. For the case of τ -VTA we select the trained prediction model of Section IV-A and we get the predicted VTA designs that are proposed by the model to operate better, for an input feature vector of a precision of [1, 2, 4, 8] bits. We set the maximum number of iterations $n = 1000$. Our proposed τ -VTA designs are able to guide the auto-tuning to higher performance for the transprecision vector of [8,4,2,1] bits by $1\times$, $1.7\times$, $1.4\times$ and $1.2\times$ for TreeGRU and $2.3\times$, $2.5\times$, $2.1\times$ and $2.3\times$ for random search. In addition, the τ -VTA auto-tuning converges faster by $1.5\times$, $5\times$, $5.2\times$ and $7.1\times$ for TreeGRU and $6.6\times$, $2.2\times$, $6.2\times$ and $8.1\times$ for random search, respectively.

V. RELATED WORK

To accelerate DL models on diverse DL HW, it is important to map the computation to DL HW efficiently. On general-purpose HW, efficient computation of DL models relies on the highly optimized linear algebra libraries such as Basic Linear Algebra Subprograms (BLAS) libraries (e.g., MKL and cuBLAS). In addition, the HW vendors have released specially optimized libraries tailored for DL computations (e.g., cuDNN and MKL-DNN). More advanced tools, like TensorRT [3], support graph optimization and low-bit quantization with large collection of highly optimized GPU kernels. Due to enormous search space for parameter tuning in hardware-specific optimizations, an emerging set of techniques, namely auto-tuning, are necessary to determine the optimal parameter settings. Apart from TVM’s auto-tuning, which we used in our analysis, and was used with FPGAs [22], as well as GPUs and embedded devices [8], numerous DL works employ optimizers for parameter tuning. Tensor

Comprehensions [28] firstly uses black-box optimization to select parameters of thread blocks and secondly polyhedral optimization to generate internal loops. This task, falls within the category of hyper-parameter optimization for optimal code generation for different HW back-ends. An extensive review of such optimizers, alongside the DNN compilation frameworks to apply these optimizers, is presented in [17]. Similar to our approach, the formulation of a prediction model that quantifies the impact of a hardware design choice, has been approached several times in literature already [25]. The general category of such research direction falls within the design space exploration of HLS tools for an FPGA model. However, to the best of our knowledge, the only background work for auto-tuning with TVM compiler stack, is [8]. As a future work we aim to study how our prediction model depends on a particular DNN topology.

VI. CONCLUSION

This research explores the possibility of automatically guiding the auto-tuning of an overlay architecture, for different transprecision settings, by leveraging knowledge from hardware experience. By adopting the concept of agile development we built a pipeline of engineering tasks that support the auto-tuning process. Instead of eliminating the overlay hardware design space with pruning techniques, we propose a technique that builds a prediction model to quantify the impact of a hardware design choice towards an optimization goal. We show that the features with the highest impact differ for different precisions. By using the most important features in order to generate an overlay we manage to perform auto-tuning that succeeds in higher performance of up to $2.5\times$ and faster convergence of up to $8.1\times$.

REFERENCES

- [1] Apache (incubating) tvn: An end to end deep learning compiler stack for cpus, gpus and specialized accelerators. <https://tvm.apache.org/>.
- [2] Pynq: Python productivity for zynq. <http://www.pynq.io/>.
- [3] Tensorrt github repository. accessed april 4, 2020. <https://github.com/NVIDIA/TensorRT>.
- [4] Baoyuan Liu, Min Wang, H. Foroosh, M. Tappen, and M. Pensky. Sparse convolutional neural networks. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 806–814, June 2015.
- [5] Michaela Blott, Thomas Preusser, Nicholas Fraser, Giulio Gambardella, Kenneth O’Brien, and Yaman Umuroglu. Finn-r: An end-to-end deep-learning framework for fast exploration of quantized neural networks.
- [6] Tianqi Chen and Carlos Guestrin. XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’16*, pages 785–794, New York, NY, USA, 2016. ACM.
- [7] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Meghan Cowan, Haichen Shen, Leyuan Wang, Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. Tvm: An automated end-to-end optimizing compiler for deep learning. In *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation, OSDI18*, page 579594, USA, 2018. USENIX Association.
- [8] Tianqi Chen, Lianmin Zheng, Eddie Yan, Ziheng Jiang, Thierry Moreau, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. Learning to optimize tensor programs. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS18*, page 33933404, Red Hook, NY, USA, 2018. Curran Associates Inc.
- [9] Jason Cong, Zhenman Fang, Yuchen Hao, Peng Wei, Cody Hao Yu, Chen Zhang, and Peipei Zhou. Best-effort FPGA programming: A few steps can go a long way. *CoRR*, abs/1807.01340, 2018.
- [10] S. Dai, Y. Zhou, H. Zhang, E. Ustun, E. F. Y. Young, and Z. Zhang. Fast and accurate estimation of quality of results in high-level synthesis with machine learning. In *2018 IEEE 26th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, pages 129–132, 2018.
- [11] Julian Faraone, Nicholas Fraser, Giulio Gambardella, Michaela Blott, and Philip H. W. Leong. Compressing low precision deep neural networks using sparsity-induced regularization in ternary networks. In Derong Liu, Shengli Xie, Yuanqing Li, Dongbin Zhao, and El-Sayed M. El-Alfy, editors, *Neural Information Processing*, pages 393–404, Cham, 2017. Springer International Publishing.
- [12] Kaiyuan Guo, Shulin Zeng, Jincheng Yu, Yu Wang, and Huazhong Yang. [dl] a survey of fpga-based neural network inference accelerators. *ACM Trans. Reconfigurable Technol. Syst.*, 12(1), March 2019.
- [13] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and Huffman coding. In Yoshua Bengio and Yann LeCun, editors, *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.
- [14] Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections for efficient neural networks. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1, NIPS’15*, pages 1135–1143, Cambridge, MA, USA, 2015. MIT Press.
- [15] Forrest N. Iandola, Song Han, Matthew W. Moskewicz, Khalid Ashraf, William J. Dally, and Kurt Keutzer. SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5MB model size. pages 1–13, 2016.
- [16] Y. Lee, A. Waterman, H. Cook, B. Zimmer, B. Keller, A. Puggelli, J. Kwak, R. Jevtic, S. Bailey, M. Blagojevic, P. Chiu, R. Avizienis, B. Richards, J. Bachrach, D. Patterson, E. Alon, B. Nikolic, and K. Asanovic. An agile approach to building risc-v microprocessors. *IEEE Micro*, 36(2):8–20, 2016.
- [17] Mingzhen Li, Yi Liu, Xiaoyan Liu, Qingxiao Sun, Xin You, Hailong Yang, Zhongzhi Luan, and Depei Qian. The deep learning compiler: A comprehensive survey, 2020.
- [18] Zhen Li, Yuqing Wang, Tian Zhi, and Tianshi Chen. A survey of neural network accelerators. *Frontiers of Computer Science*, 11(5):746–761, Oct 2017.
- [19] Shuang Liang, Shouyi Yin, Leibo Liu, Wayne Luk, and Shaojun Wei. Fp-bnn: Binarized neural network on fpga. *Neurocomputing*, 275:1072–1086, 2018.
- [20] A. C. I. Malossi, M. Schaffner, A. Molnos, L. Gammaitoni, G. Tagliavini, A. Emerson, A. Toms, D. S. Nikolopoulos, E. Flamand, and N. Wehn. The transprecision computing paradigm: Concept, design, and applications. In *2018 Design, Automation Test in Europe Conference Exhibition (DATE)*, pages 1105–1110, March 2018.
- [21] Sparsh Mittal. A survey of fpga-based accelerators for convolutional neural networks. *Neural Computing and Applications*, Oct 2018.
- [22] T. Moreau, T. Chen, L. Vega, J. Roesch, E. Yan, L. Zheng, J. Fromm, Z. Jiang, L. Ceze, C. Guestrin, and A. Krishnamurthy. A hardware-software blueprint for flexible deep learning specialization. *IEEE Micro*, 39(5):8–16, 2019.
- [23] Samira Pouyanfar, Saad Sadiq, Yilin Yan, Haiman Tian, Yudong Tao, Maria Presa Reyes, Mei-Ling Shyu, Shu-Ching Chen, and S. S. Iyengar. A survey on deep learning: Algorithms, techniques, and applications. *ACM Comput. Surv.*, 51(5):92:1–92:36, September 2018.
- [24] Jared Roesch, Steven Lyubomirsky, Logan Weber, Josh Pollock, Marisa Kirisame, Tianqi Chen, and Zachary Tatlock. Relay: A new IR for machine learning frameworks. In *Proceedings of the 2nd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages, MAPL 2018*, pages 58–68, New York, NY, USA, 2018. ACM.
- [25] B. C. Schafer and Z. Wang. High-level synthesis design space exploration: Past, present and future. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, pages 1–1, 2019.
- [26] Atin Sood, Benjamin Elder, Benjamin Herta, Chao Xue, Costas Bekas, A. Cristiano I. Malossi, Debashish Saha, Florian Scheidegger, Ganesh Venkataraman, Gegi Thomas, Giovanni Mariani, Hendrik Strobelt, Horst Samulowitz, Martin Wistuba, Matteo Manica, Mihir Choudhury, Rong Yan, Roxana Istrate, Ruchir Puri, and Tejaswini Pedapati. NeuNetS: An Automated Synthesis Engine for Neural Network Design. *arXiv e-prints*, page arXiv:1901.06261, Jan 2019.
- [27] Kai Sheng Tai, Richard Socher, and Christopher D. Manning. Improved semantic representations from tree-structured long short-term memory networks. In *ACL*, 2015.
- [28] Nicolas Vasilache, Oleksandr Zinenko, Theodoros Theodoridis, Priya Goyal, Zachary DeVito, William S. Moses, Sven Verdoolaege, Andrew Adams, and Albert Cohen. Tensor comprehensions: Framework-agnostic high-performance machine learning abstractions. *CoRR*, abs/1802.04730, 2018.
- [29] E. Wang, J. J. Davis, P. Y. K. Cheung, and G. Constantinides. Lutnet: Learning fpga configurations for highly efficient neural network inference. *IEEE Transactions on Computers*, pages 1–1, 2020.
- [30] Samuel Williams, Andrew Waterman, and David Patterson. Roofline: An insightful visual performance model for multicore architectures. *Commun. ACM*, 52(4):65–76, April 2009.
- [31] J. Yu, A. Lukefahr, D. Palframan, G. Dasika, R. Das, and S. Mahlke. Scalpel: Customizing dnn pruning to the underlying hardware parallelism. In *2017 ACM/IEEE 44th Annual International Symposium on Computer Architecture (ISCA)*, pages 548–560, June 2017.