



**Hewlett Packard
Enterprise**



DHBW
Duale Hochschule
Baden-Württemberg

A load-aware scheduler for large-scale neural network autotuning

PROJECT THESIS II / T2000

for the study program

Computer Science

at the

Baden-Wuerttemberg Cooperative State University Stuttgart

by

Dominik Stiller

Submission Date

September 12, 2019

Project Period

18 Weeks

Company

Hewlett Packard Enterprise

Corporate Supervisor

Jungkuk Cho

University Supervisor

Prof. Dr. Bernd Schwinn

Matriculation Number, Course

4369179, TINF17A

Declaration of Authorship

I hereby declare that the thesis submitted with the title *A load-aware scheduler for large-scale neural network autotuning* is my own unaided work. All direct or indirect sources used are acknowledged as references.

Neither this nor a similar work has been presented to an examination committee or published.

Sindelfingen September 12, 2019

Place Date Dominik Stiller

Confidentiality Clause

This thesis contains confidential data of *Hewlett Packard Enterprise Development LP*. This work may only be made available to the first and second reviewers and authorized members of the board of examiners. Any publication and duplication of this thesis—even in part—is prohibited.

An inspection of this work by third parties requires the expressed permission of the author, the project supervisor, and Hewlett Packard Enterprise Development LP.

Abstract

Real-time computer vision applications with deep learning-based inference require hardware-specific optimization to meet stringent performance requirements. Frameworks have been developed to generate the optimal low-level implementation for a certain target device based on a high-level input model using machine learning in a process called autotuning. However, current implementations suffer from inherent resource utilization inefficiency and bad scalability which prohibits large-scale use.

In this paper, we develop a load-aware scheduler which enables large-scale autotuning. The scheduler controls multiple, parallel autotuning jobs on shared resources such as CPUs and GPUs by interleaving computations, which minimizes resource idle time and job interference. The scheduler is a key component in our proposed Autotuning as a Service reference architecture to democratize autotuning. Our evaluation shows good results for the resulting inference performance and resource efficiency.

Contents

Acronyms	V
List of Figures	VI
List of Tables	VII
List of Source Codes	VIII
1 Introduction	1
1.1 Problem	1
1.2 Scope	2
2 Background	3
2.1 Artificial Neural Networks	3
2.2 Inference Optimization	4
2.3 Manual Optimization	6
2.4 Automated Optimization	6
3 Using TVM	15
3.1 SimpleTVM	15
3.2 Parameters	18
3.3 Capabilities	19
3.4 Limitations	20
4 Autotuning Scheduler	25
4.1 Design	25
4.2 Implementation	32
4.3 Autotuning as a Service	37
5 Evaluation	40
5.1 Results	41
5.2 Discussion	43
6 Conclusion	46
Bibliography	47
Glossary	49
A Experiment Results	50

Acronyms

AaaS Autotuning as a Service

AI artificial intelligence

ANN artificial neural network

CNN convolutional neural network

DL deep learning

GPU graphics processing unit

ML machine learning

RPC remote procedure call

TC TensorComprehensions

List of Figures

1	Traditional vs. optimized machine learning workflow	5
2	Expressions and low-level code for transposed matrix multiplication .	7
3	Levels of abstractions in TVM stack	9
4	Iterative autotuning process in TVM	11
5	TVM's RPC architecture	14
6	Interface and flow of SimpleTVM	15
7	Inference performance with TensorFlow and TVM	20
8	Resource utilization during autotuning	21
9	Setups for scaling autotuning	22
10	Interleaving of multiple autotuning jobs	26
11	Client interface before and after decomposition	27
12	Autotuning process with scheduler	31
13	Layers and components of scheduler implementation	33
14	Autotuning as a Service reference architecture	38
15	Impact of interference	42
16	Results of greedy versus fair interleaving	42
17	Results of greedy versus fair interleaving with bundling	43

List of Tables

1	Comparison of scaling setups	22
2	Evaluation setups	41

List of Source Codes

1	Typical SimpleTVM flow for CPU including autotuning	17
2	Greedy interleaved scheduling pseudocode	28
3	Fair interleaved scheduling pseudocode	29
4	Sequential scheduling pseudocode	30
5	Synchronous scheduling pseudocode	30
6	Pseudocode of JobManager's stage decision logic	35
7	Comparison of default and scheduled autotuning	36

1 Introduction

In recent years, artificial intelligence (AI) has garnered tremendous success, revolutionizing the way we work and accelerating economic growth of industries and whole nation [1, p. 15 ff.]. Especially deep learning (DL), a subfield of AI, has made vast improvements and is the prime method of modern AI. In the future, AI will be applied to even more areas, where non-expert users want to benefit from AI without the technical complexity introduced by development and deployment of AI applications. This is facilitated by platforms such as BlueData and Qubole, which automate infrastructure setup and provide user-friendly interfaces to make AI and DL more accessible.

1.1 Problem

More and more applications like industrial monitoring or autonomous driving require real-time performance, most of them powered by DL. Specialized accelerator hardware such as graphics processing units (GPUs) or FPGAs are employed to speed up the computation-intensive inference. However, the model itself needs device-specific, low-level optimizations to harness the accelerator’s full potential. Currently, these optimizations are manually developed by the device vendors who have deep knowledge of their hardware. DL researchers who want to experiment with new model types and high-level optimizations are forced to wait until low-level implementations are supported by vendor libraries.

Automated performance optimization, called autotuning, creates optimized low-level implementations without the need for human experts in a vendor-agnostic way. This fosters innovation and helps manage the increasing performance demands for a growing variety of models and accelerator devices. While autotuning is already employed, it has yet to reach widespread apply, partially because it is still inconvenient to use. Offering autotuning as a service can make it accessible to a larger audience to facilitate real-time DL applications, but requires support for large-scale autotuning. However, inefficiencies in the autotuning process prohibit efficient scaling and, in turn, the

implementation of an Autotuning as a Service platform. To the best of our knowledge, there is no existing solution for scaling up autotuning.

1.2 Scope

In this thesis, we design and create a prototypical implementation of a load-aware scheduler to enable large-scale autotuning. This scheduler controls multiple autotuning jobs that share computation resources to overcome the inefficiencies of current autotuning. We show that controlling the execution of multiple jobs by a load-aware scheduler makes large-scale autotuning more efficient in terms of

- autotuning completion time,
- resulting inference performance and
- hardware requirements.

First, we discuss manual and automated performance optimization before comparing two frameworks for autotuning (Chapter 2). Next, we will develop a framework to examine the capabilities and limitations of autotuning in different scenarios. This will allow us to find the concept of interleaving which we can leverage to scale autotuning (Chapter 3). We will design and implement our scheduler which is used in our proposed reference architecture for Autotuning as a Service (Chapter 4). Finally, we evaluate our scheduler design and point to future improvements (Chapter 5). Our experiments show good results for resulting inference performance and hardware requirements.

No improvements are made to the autotuning process itself, but we base our work on the TVM [2] autotuning framework and enhance it with a further component. Also, we do not implement Autotuning as a Service. This thesis describes only a reference architecture, a prototype implementation is described in [3].

This project was conducted by the *Networking, IoT and Mobility Laboratory* of the *Hewlett Packard Labs*.

2 Background

Machine learning (ML) has become an important sub-field of computer science. It emulates human-like learning using mathematical models, so predictions can be made about new data in the future. Rather than explicitly programming how to make those predictions, the developer provides sample data during *training*. Once the accuracy of the trained model is sufficient, it can be used for *inference*. The model can be thought of as the approximation of a function mapping from the input data to some output, e.g., a label for classification, or a numerical value for regression [4, p. 164].

2.1 Artificial Neural Networks

While there are a variety of ML models in use today, artificial neural networks (ANNs) are among the most powerful and flexible, due to their ability to represent complex functions [4, p. 163]. They find application in fields as diverse as image and speech recognition, movie recommendations and medical diagnosis.

ANNs are composed of multiple layers, with the output of one layer being the input of another layer. The first layer receives the input data, and the last layer produces the final output. With an increasing number of layers, or *depth* of the network, more complex functions can be approximated. These deep networks are subject of the ML sub-field of deep learning (DL). All layers perform some computation given a set of trained or specified parameters and the input. Both parameters and inputs are tensors, a higher-dimensional generalization of vectors and matrices. The computation is implemented by a *tensor operator*. Traditional ANNs feature only fully-connected layers with some activation function.

Grid-like data such as time series (1D) or images (2D) benefit from additional layers found in convolutional neural networks (CNNs) [4, p. 326]. This makes CNNs an important tool in state-of-the-art computer vision applications. CNNs apply convolution and pooling to a region of the input tensor in a sliding fashion, so values only interact with other values that are located in their neighborhood. Convolution applies one or more kernel matrices to the input, which are element-wise multiplied with the

current region and then summed up into a single output value. Pooling averages or finds the maximum of the region as output value. Both operations support a variable stride (step size) and padding.

While neural network models logically consist of a series of layers, machine learning frameworks usually represent them in a computation graph. The computation graph's first vertex is the input node, followed by a number of tensor operators and their parameters that perform the layer's computations, and finally an output node. The edges describes how data flows between the vertices.

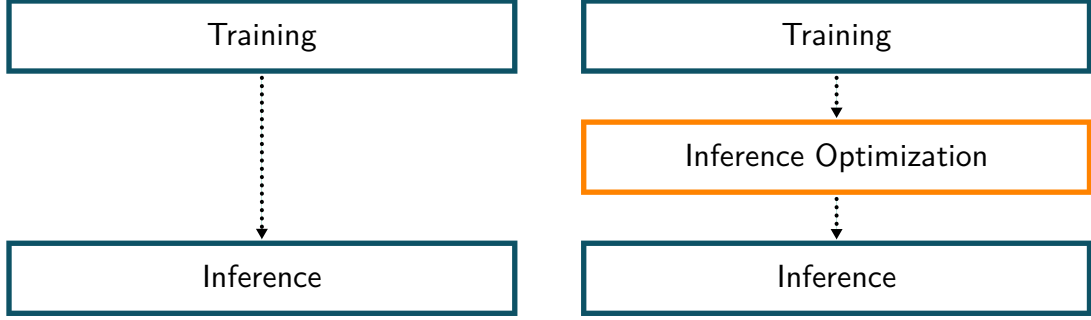
2.2 Inference Optimization

Training a deep CNN can take very long. However, the amount of inferences typically outweighs the amount of trainings heavily, since training only needs to be done once (albeit model re-training is usually done periodically when new training data is available). For this reason, speeding up inference has a larger impact than training acceleration and is a worthwhile endeavor. Reduction of the inference time has a number of advantages:

- less hardware is required to achieve the same inference rate
- a higher inference rate can be achieved with same hardware
- real-time applications are facilitated

In real-time applications with a high inference rate, even small improvements in inference performance (in the order of milliseconds) can be critical to guarantee the required throughput. For example, a major hard drive manufacturer detects defects in their products early using a CNN-based smart manufacturing solution [5, p. 11]. They perform inference on 3 million images every day, so if they were able to save only 5 ms per image due to some performance optimization, that would amount to over 4 h less total inference time every day [6]. Alternatively, they could save cost by needing less servers that are equipped with expensive accelerator devices.

Accelerator devices such as GPUs, FPGAs, or ASICs like tensor processing units are used to speed up both training and inference. However, generic ML models cannot make full use of accelerator capabilities and fall short of leveraging the full potential. Every device has different features such as specialized instructions, memory size and layout, caching and parallelization support. This means that models need to be attuned to the *target device* to achieve the best inference performance. But even if no special accelerator devices are used but only a conventional CPU, adapting to



(a) Traditional without inference optimization (b) Improved with inference optimization

Figure 1: Machine learning workflow

the specific architecture can yield great performance benefits [7, p. 1]. In a traditional machine workflow, the trained model is deployed as-is (Figure 1a). Inference optimization adds an additional step, turning the trained model into a functionally equivalent but optimized version before inference (Figure 1b). In this step, high-level transformations that rewrite the computation graph by, for example, fusing tensor operators, pre-computing constant parts or transforming the data layout in memory are applied first [2, p. 1–3]. More importantly, however, the low-level implementation of tensor operators can be changed.

The model determines what tensor operations are calculated, but it does not specify how they are calculated. Deliberately choosing the actual implementation offers great optimization potential. There is always a generic naïve implementation, which is the straightforward way of performing the computation. However, it does not consider, e.g., memory sharing between threads or cache access patterns, which can have a significant adverse effect on performance [8]. Implementations that use techniques such as loop unrolling, reordering and tiling as well as multi-dimensional threading and tensor compute instructions can help harness the accelerator’s full potential, but there is an abundance of combinations of settings for these techniques, the best of which is very much specific to the target device [2, p. 2]. Finding the optimal such combination is the key aspect of tensor operator optimization.

Convolution operators are very computationally intensive and make up the majority of modern CNNs, such as Inception [9] and ResNet [10]. Therefore, tensor operator optimization should focus on convolution over other types like pooling and fully-connected. It is not possible to optimize convolutions in general, but optimization needs to be done for every distinct parameter set that is present in the computation graph, i.e. combination of input shape, kernel shape, padding, and stride. This means that the effort increases with a higher variety of layer configurations.

2.3 Manual Optimization

Optimized implementations for tensor operators with a specific parameter set are provided by accelerator vendors in libraries like cuDNN for NVIDIA GPUs and Intel Math Kernel Library for Intel CPUs. The vendors possess the hardware-specific knowledge to write good implementations by hand, but human expertise is required for this approach. While state-of-the-art, manual optimization has a number of inherent shortcomings:

- slow support for new devices
- slow support for new graph-level optimizations
- no support for unconventional layers
- vendor lock-in

These limitations hinder innovation, which is undesirable in a field so fast-evolving and young as DL. Researchers need to choose between avoiding devices, high-level optimizations and new network architectures that are not supported by those predefined operator libraries, and falling back to unoptimized implementations [2, p. 1].

2.4 Automated Optimization

Automated tensor operator optimization, or *autotuning*, overcomes these shortcomings by eliminating the need for human experts. Vendor-agnostic frameworks can discover good implementations regardless of hardware, model or graph optimizations. This enables innovation by fostering experimentation with novel or unconventional layers and high-level transformations that are not yet supported by manual libraries. Autotuning can achieve the same, in some cases even better inference performance than state-of-the-art vendor-provided operator libraries. Compared to these libraries, autotuning delivers speedups of $0.98\times$ to $3.5\times$ on CPU [7, p. 9] and $1.6\times$ to $3.8\times$ on server-class GPUs [2, p. 10] for commonly used CNNs. Even a slightly worse performance is impressive considering that no domain-specific expert knowledge has been applied but only a few hours of autotuning.

Autotuning works by exploring the space of possible implementations in an organized fashion. Functionally equivalent implementations can be generated by a *schedule* which defines a series of parametrized transformations called *schedule primitives* that are applied to the naïve implementation. The *search space* is defined by the set of permutations of parameter settings. These settings control, for example, loop

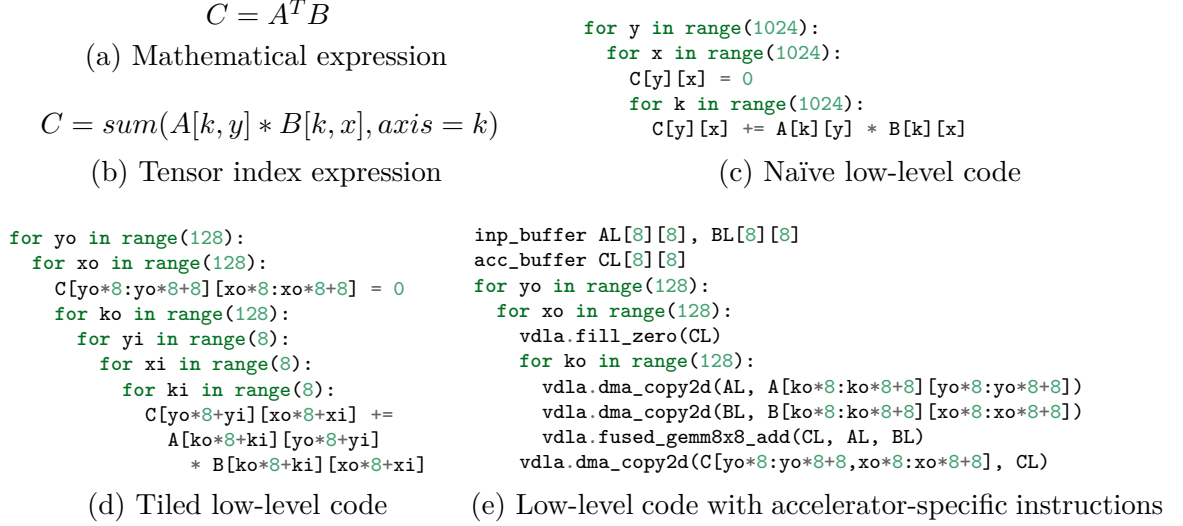


Figure 2: Expressions and low-level code for transposed matrix multiplication [2, p. 4]

unrolling factors, loop order, loop tiling sizes and thread numbers, and can usually be adjusted in steps of powers of 2 [2, p. 5] [11, p. 16]. One specific combination of settings, i.e. one element of the search space, is called *configuration*. Defining the values a setting can take is done manually for each class of target devices, but the search is guided by an algorithm that proposes candidate configurations. This is necessary since the size of the search space makes the brute-force approach of trying all configurations infeasible. As an example, the search space size for a ResNet-18 on an NVIDIA GPU exceeds 172 million possible configurations, any one of which could be the best. ML-based or genetic algorithms help with rapid convergence to a decent, or ideally the best configuration without need of exhausting the whole search space.

Figure 2 provides an example of how different configurations affect the generated *low-level code*. The operator functionality is some mathematical calculation, in our example a transposed matrix multiplication (2a). Before autotuning, that functionality is specified in a tensor expression language, which describes how to compute each element of the output tensor from the input tensors using a concise notation (2b). Note that this notation is implicit, meaning that it does not prescribe implementation details. The autotuning framework then makes the computation explicit by applying a schedule with specific parameters from the configuration to the operator’s default code. This simple but naïve reference implementation can be used to check the correctness of outputs after a complex transformation (2c). The low-level code allows transformation through schedule primitives, e.g. tiling for memory locality (2d) or accelerator-specific instructions for buffers and specialized tensor operators (2e). The specific tiling factors and buffer sizes can be varied and are determined by the applied configuration [2, p. 4 ff.] [11, p. 9 ff.].

Low-level code is only an intermediate representation from which target-specific code, e.g., LLVM assembly for CPU or a CUDA kernel for NVIDIA GPUs, needs to be generated. The appropriate compiler then builds that code, possibly in parallel for multiple configurations in a batch, after which the implementation can be executed. For autotuning, the execution time is then profiled on the target device to evaluate the performance. The profiling result is stored alongside the implementation in the *autotuning database* and fed back to the algorithm that selects candidate configurations. This allows the algorithm to improve its proposals for the next batch [11, p. 15 f.]. The iterative autotuning process can be stopped when a sufficiently fast implementation has been found or no better one has been discovered in a long time. Then, the full computation graph can be used for inference, utilizing the best implementations that have been found in the autotuning process for all tensor operators.

There are two frameworks that implement autotuning, which will be described now.

2.4.1 TensorComprehensions

TensorComprehensions¹ (TC) has been developed by Facebook’s AI Research team and comprises three main components: a language to express tensor computations (similar to Figure 2b), an optimizing compiler to generate efficient GPU code from expressions in that language, and an autotuner that finds good implementations and stores them in a compilation cache. It uses a polyhedral compiler to reason about and manipulate the loop structures of an implementation [11, p. 3]. However, only tensor-operators are considered, the framework is designed to be independent of computation graphs [11, p. 4].

Autotuning in TC starts from configurations that worked well for similar expressions, and some predefined strategies. The autotuner determines the configuration parameters and admissible value ranges. Then, a genetic algorithm generates a batch of candidate configurations. The value for each configuration parameter is randomly selected from one of three parent configurations that are selected probabilistically based on their fitness. Furthermore, there is a low probability of mutation, which means that a random value is assigned to some parameters. Configurations are then compiled in parallel and profiled on an available GPU. A fitness value inversely proportional to the execution time is assigned to the configuration and stored in the autotuning database. Then, the process starts anew by selecting the next candidates using the updated database. This is repeated for a set amount of time [11, p. 15 f.].

¹<https://github.com/facebookresearch/TensorComprehensions>

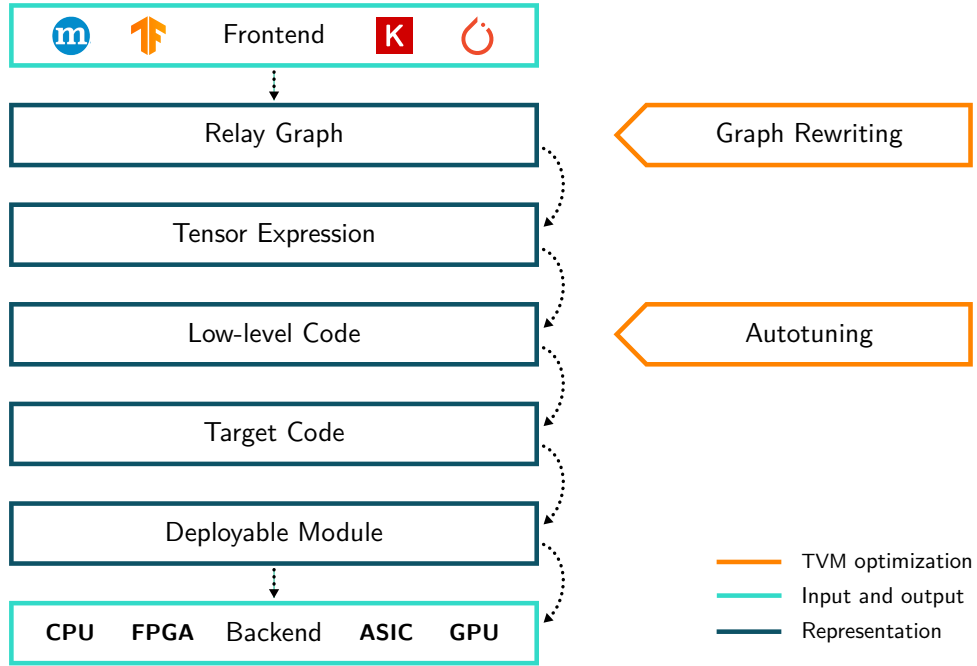


Figure 3: Levels of abstractions in TVM stack

2.4.2 TVM

TVM² started as a research project at the University of Washington but is now supported and used by a large open-source community and companies such as Amazon and Facebook. Unlike TC, which only represents and optimizes tensor operators, TVM is an end-to-end DL compiler stack. It can import whole models from a *frontend* framework and build minimal, optimized modules that can be deployed to *backends* like CPUs, GPUs or FPGAs. Figure 3 shows how the layers of the stack provide different levels of abstraction.

The top layer in the TVM stack is Relay. Relay is an intermediate model representation that enhances traditional computation graphs with concepts of functional programming to form a more powerful language. Relay supports shape-dependent tensor types and automatic differentiation, which is essential for DL training [12, p. 61]. Additionally, a runtime to execute Relay programs in various programming languages is provided and needs to be present whenever executing TVM-based models. Relay programs can be created programmatically or from a textual source code. More convenient for users, however, is the import from diverse frontends, including TensorFlow, Keras, PyTorch and MXNet, which enables the use and optimization of existing models. Graph-level optimization in TVM is pass-based, with each pass inspecting or rewriting the syntax tree of the Relay program in some way. Standard passes are provided and perform, for example, automatic differentiation, type inference, operator

²<https://github.com/dmlc/tvm/>

fusion or tensor layout transformations [2, p. 3]. Beyond that, writing custom passes is facilitated by an extensible design.

Next in the stack is a tensor expression language, which has similar features as TC’s. It allows user to describe computation rules that generate a tensor without specifying loop structures and other details. The rules are composed of primitive mathematical operations like addition and multiplication and are expressive enough to describe tensor, matrix and vector operations. TVM comes with tensor expressions for common computations used in DL such as various activation functions, convolution, pooling, and matrix multiplication [2, p. 4 f.]. The tensor expression language is used to describe the functionality of tensor operators from the model. In the usual TVM workflow, the required operators are extracted from the Relay graph and matched with existing tensor expressions, so there is no need to write them manually.

Implicit tensor expressions need to be mapped to explicit, backend-independent low-level code. TVM, again, uses a pass-based design, which is different from TC’s polyhedral approach. Each pass applies a schedule primitive to the naïve implementation. This design is based on the Halide language for image processing, which works with similar multi-dimensional data as DL, but enhances it with more primitives to optimize accelerator performance. TVM leverages nested and cooperative parallelism to make effective use of GPU memory structure by enabling data reuse across threads through shared memory regions. This is done in a special memory scoping pass. TVM also equips the low-level code with hardware-specific instructions through a tensorization pass which matches computation patterns with a corresponding intrinsic from the target (such as general matrix multiply), making it extensible for new hardware architectures. A latency hiding pass introduces explicit management of fine-grained synchronization for memory and computation instructions on specialized DL accelerators [2, p. 5 f.]. Default schedule templates are provided for every hardware type, but users can create their own templates to incorporate their knowledge of the backend.

Low-level code cannot be executed, but it can directly be converted to target-specific code and then compiled for the target device. Backend-specific code generators create the source files, which are then built by the respective compiler and packed into a module which contains the implementation of all tensor operators in the model. This module can be deployed along with a JSON description of the Relay graph and a parameter file containing the weights for all operations. The TVM runtime (300 kB to 600 kB) needs to be installed on the target system to execute the model. However, a full DL framework is not required, making TVM modules very lightweight to integrate into applications.

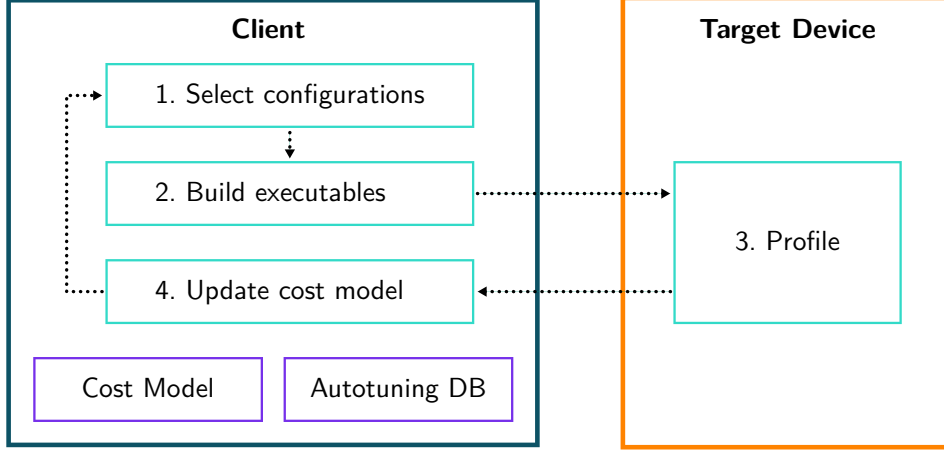


Figure 4: Iterative autotuning process in TVM

While TCs’s autotuner is guided by a genetic algorithm, TVM uses a ML-based cost model to predict the performance of an implementation. Specifically, gradient boosted trees are employed because of their advantage in training and prediction speed over neural network-based models. Since the model is queried frequently, the inference overhead must be smaller than the profiling it seeks to replace. While profiling can be in the order of seconds, the gradient boosted trees model performs prediction in 0.67 ms on average. Model training time needs to be considered as well; the cost model is updated periodically as more configurations have been explored, which improves the accuracy for further predictions with more experimental trials. This learning-based approach is preferable to static, predefined cost models for every new hardware target, which is infeasible due to the increasing complexity and diversity of modern accelerators [2, p. 8 f.]. The input for the cost model is not the configuration but the low-level code, which needs to be encoded into vector space first. This encoding is a representation which is invariant between programs to enable transfer learning. Encoding works by extracting context features from each loop level such as memory access count, but also relation features for generalization across different loop nest patterns [13, p. 4].

We call the execution of the autotuning process for one model a *job*. The component that executes the autotuning logic for one job is the *autotuning client*. Profiling the implementation requires execution on the actual target device. This can make autotuning a distributed process if the client is another device. Autotuning is not performed for a whole model at once, but rather for a set of *tasks* which correspond to autotunable tensor operators with a specific configuration (shapes, padding, stride). These tasks need to be extracted from the model before starting the process for each of them. Autotuning in TVM is an iterative process consisting of four stages that depend on each other, making it necessary to execute them in sequential order. This

is illustrated in Figure 4. Understanding the stages and their dependencies is key for enabling large-scale autotuning.

Initialization At the start of each task, profiling results from previous jobs are loaded from the global autotuning database, a file that contains data from all previous jobs along with information about the target and configuration. The loaded results are passed to the cost model for transfer learning. This yields good cost model from the beginning and improvement in model quality over time. Then the autotuning loop can be launched.

1. Select candidate configurations At the start of each iteration, a batch of candidate configurations that have a promising performance is selected using the cost model.³ A simple strategy such as enumerating and running every configuration through the model, then selecting the top performers is impracticable with large search spaces. Rather, candidates are selected using parallel simulated annealing, which is a heuristic optimization algorithm that trades off finding an exact optimum for a much improved speed. Additionally, exploration is ensured by random selection of some configurations. If no training data exist yet, random candidates are picked.

2. Build executables The client combines the batch of configurations from the previous stage with the schedule template, then applies the schedule to the tensor expression for the operator of the current task. The resulting low-level code is then translated into backend-specific code and compiled. In case the target hardware is different from the client hardware, cross-compilation is necessary. The result is a tar file that contains everything that is necessary to run the executable, namely the compiled tensor operator itself and backend-specific code such as the CUDA driver library for NVIDIA GPUs.

3. Profile on target device Since the cost model’s prediction of the implementation’s performance is not completely reliable, the real performance needs to be measured on the target device. Before, the tar files from the build stage are uploaded to the target device. Then the implementation is profiled by running the executable a number of times with random data. The measured execution times are averaged and returned to the client, which stores the results in the autotuning database for this job.

³In TVM’s implementation, the selection of the next candidates is actually performed at the end of the iteration after updating the model, and the first stage just picks a batch from the already selected configurations. However, it is logically simpler to think of the candidate selection using simulated annealing as part of the first stage. Furthermore, it is not wrong per se due to the iterative nature of autotuning.

4. Update cost model The cost model is updated with the measurements from the profiling stage to improve the proposed configurations in the next iteration. This is only done after a sufficient number of new profiling data has been collected, so this stage might be skipped in some iterations.

Finalization After a certain number of trials, the loop is stopped. The best configuration that was discovered can now be used to build a faster implementation of the tensor operator that was optimized. Usually, the best configuration is also written into a separate database that contains only the best known configurations. The autotuning database for this job is merged with the global one. This concludes the autotuning process in TVM.

While the target device is usually specialized for DL workloads, it is desirable to run the client on a machine that features a strong CPU to accelerate the compute-intensive build and model update stages. This distribution across multiple machines requires an remote procedure call (RPC) infrastructure that makes it possible to profile on a different server. TVM’s RPC architecture (Figure 5) comprises three components:

Client A client runs an autotuning job and is responsible for selecting the candidate configurations, building the executables and updating the model with the profiling results. This means that the client contains both the cost model and the autotuning database. It also controls the profiling, but the actual execution is happening on servers.

Server A server can receive and execute TVM modules, which basically makes it an RPC-enabled TVM runtime that runs on the target device. The interface on the client side for running remote TVM operations does not change because the framework transparently handles remote execution like local execution. A server has a *device key*, which is an arbitrary identifier for a certain device type, but usually is it based on the accelerator’s name. Multiple servers can have the same device key if they run on identical target devices.

Tracker A tracker keeps a list of servers to help clients discover unused servers for profiling. The tracker matches incoming requests from clients with free servers using a FIFO-based scheduling algorithm.

TVM’s RPC is enabled by two distinct protocols. The control plane protocol is used for communication involving the tracker, namely server registration and requests from clients. The data plane protocol facilitates remote execution on a server, with connections being initiated by clients.

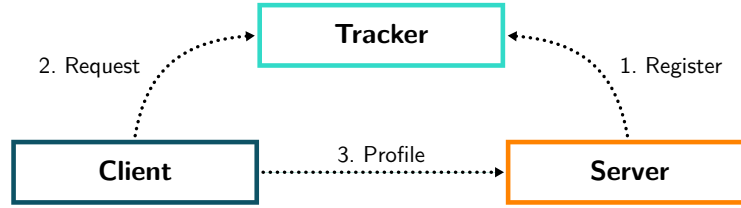


Figure 5: TVM’s RPC architecture

First, the tracker is started. Then, one or multiple servers are started and register with the tracker by transmitting the device key, and the address and port which clients can use to connect. The tracker puts them in a queue, with separate queues for every device key. At this point, clients can request a server with a specific device key, and the request put in a queue as well. The tracker matches the first-registered free server having that device key with the request that was received first, i.e. simple FIFO scheduling. Once a client has acquired a server, it is marked as busy in the tracker and the client initiates a connection to the server to use its TVM runtime for profiling.

Since autotuning works in batches, usually not a single but multiple servers are requested to run profiling in parallel. This can speed up profiling if multiple target devices are available. For example, if a machine is equipped with 4 GPUs of the target device type, 4 RPC servers can be launched on that machine, with each one being assigned to a different one of the GPUs.

In this project, we use TVM instead of TC because of the novel, machine learning-based approach, which promises better results than a genetic algorithm due to intelligent guidance by the cost model. We are using the TVM version from June 11, 2019 (commit 8f219b9) for comparable results throughout the project. We made some modifications:

- Add decomposed version of autotuner with separate methods for stages
- Add time measurement for autotuning stages
- Add loading of autotuning records from multiple files
- Fix Tensorflow import for models including PlaceholderWithDefault

3 Using TVM

For our end goal of enabling large-scale autotuning, we need to explore the current capabilities and limitations of TVM first, especially with regard to the execution of multiple autotuning jobs simultaneously. The modern DL landscape is very diverse in terms of models and hardware, so to evaluate TVM in a diverse range of scenarios is crucial for gaining a proper understanding. To this end, we developed a framework that enables us to perform a large number of experiments rapidly.

3.1 SimpleTVM

Since using TVM follows a similar flow every time, we created SimpleTVM which exposes the individual steps through a convenient, chainable interface. This makes it easy for researches who are new to TVM to get started. Since a lot of the experiments include benchmarking, time measurements are taken for most steps and automatically saved. SimpleTVM enforces flow dependencies; for example, a model needs to be imported before building. The interface including possible flows is depicted in Figure 6. The methods that are exposed are now regarded closer.

from_model Loading the Relay representation for the model is the beginning of a TVM flow. To that end, TVM supports the import from various frontends. Before the import of the model, however, it needs to be loaded and prepared for import. How exactly this is done differs even inside the same framework.

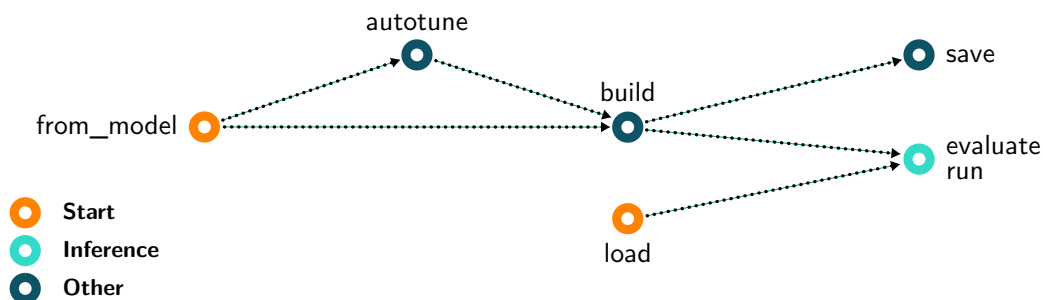


Figure 6: Interface and flow of SimpleTVM

SimpleTVM provides a unified import interface for TVM testing models, TensorFlow saved models (.pb files) and TensorFlow hub models. `from_model` can easily be extended to support more frontends.

autotune Once a model has been imported, autotuning can be run for each of the tensor operators. This step is optional, since TVM falls back to a default implementation if no records for that tensor operator exist in the autotuning database. Since this step is rather complex, there is a plethora of configuration options, the most important of which are exposed in `autotune`'s interface. SimpleTVM is designed to get new users started quickly so there are appropriate default values, but more experienced users can adjust the values to their circumstances.

build Before a TVM model can be executed, the target-specific executable needs to be built from the Relay model and tensor operators. When building, SimpleTVM can either automatically use records from the global autotuning database, or the results from a specific autotuning job can be used.

save After building, the library containing the operators can be saved along with the graph description and weights.

load Beyond starting from an imported model, SimpleTVM can also load a previously saved TVM modules, which makes it possible to use a model that has been autotuned earlier. Saved modules can only be loaded if they have been built for the same device.

run Inference can be run using this method, which accepts and returns a NumPy array. The input can, for example, be an image. However, loading the image and preparing it for inference, e.g., scaling and normalizing, needs to be done by the user.

evaluate To profile the performance, `evaluate` runs inference on random data multiple times, then averages the measured times. However, in contrast to the profiling stage of autotuning, not the performance of individual tensor operators but the whole model is measured.

An example of how SimpleTVM is typically used is presented in Listing 1. First, the `BenchmarkingContext` is created (Line 1), which stores information about the current run such as the run id (a 32-character alphanumeric identifier for this execution of SimpleTVM), target device, measured times, the loaded model and the target device key to send to the tracker. When using a CPU as target device, the CPU architecture should be specified so TVM can select the proper hardware-specific tensor instructions. The benchmarking context is passed to the `SimpleTVM` object (Line 2). Here, the


```
1 ctx = BenchmarkingContext('cpu', device_key='i7', cpu_arch='skylake')
2 tvn = SimpleTVM(ctx, rpc_tracker=('tracker', 9190))
3 tvn.from_model('mobilenet.pb', output_name='out', output_size=10)
4 tvn.autotune().build().save().evaluate()
5 ctx.save()
6
7 # Saved model can be loaded later to run inference
8 tvn.load('run_id')
9 prediction = tvn.run(data)
```

Listing 1: Typical SimpleTVM flow for CPU including autotuning

address of the RPC tracker can be specified for distributed autotuning. If the address is not specified, autotuning will create a local tracker and server to perform autotuning on the client device. Next, a model is imported (Line 3). Since the name of the output layer and the size of the output vector can differ, they need to be specified explicitly. SimpleTVM’s concise, chained syntax is used to autotune, build, save and evaluate the model (Line 4). For the sake of brevity, default parameters are used, but the user can customize the actual calls to TVM functions if desired by providing more parameters. Finally, the benchmarking context is saved (Line 5). This enables analysis at a later point, e.g., to examine the autotuning process or the inference performance measured by the evaluation. Note that this step is distinct from the saving of the TVM module. At a later point and usually by another application, the saved module which is identified by the run id can be loaded back (Line 8). Then it can be used to run inference on any data (Line 9).

Additionally to SimpleTVM, we developed an automated hyper-benchmarking script called *superb*. *superb* allows testing of different configurations without human intervention, so it performs benchmarking on a higher level than SimpleTVM’s `evaluate`. The user can specify the values for all parameters that should be tested. *superb* enumerates all possible combinations, effectively determining the product set of all value lists, then executes SimpleTVM with each configuration. Additionally, it sets up the required servers and the tracker. The results from all configurations are collected and can then be processed by another script. This script evaluates the resulting inference performances, aggregates some information and writes them into a file, enabling further analysis with other tools such as Jupyter notebooks.

Since we want to test TVM on a variety of machines, we created Docker images to be able to easily deploy TVM with all dependencies on any server. The GPU version also includes the CUDA libraries, and a helper script for using the images mounts some folders into the container and sets up the environment. The Docker images in conjunction with SimpleTVM and *superb* form the foundation for our experiments.

3.2 Parameters

Autotuning with TVM offers a plethora of configuration options that affect both the autotuning process itself and the result. Setting these parameters to adequate values for the given job and hardware requires knowledge of how TVM works, but in some cases it is a matter of trial and error. However, guidelines and descriptions of the most important parameters can help. All of the following parameters can be specified when using SimpleTVM.

Number of trials This determines the number of configurations to try for each autotuning task. A higher number will generally result in a better inference performance since the search space can be explored more extensively, but this results in an increased autotuning completion time. The resulting inference speed starts to converge to the optimum after about 500 iterations, so there is a limit to the performance that can be achieved. Especially with CPUs, that have a small search space compared to GPUs, there might not even be more options to try. Practically, the optimal result can be expected with the number of trials set to 1500–2000.

Profiling timeout This determines the time after which the profiling for one configuration is killed if it runs too long. Since every tensor operator has a different computational intensity and performance varies across types of hardware, this timeout needs to be adjusted accordingly. A high profiling timeout will allow longer execution, which drives up total autotuning completion time and might not yield better results since long-running implementations are not good and can safely be killed. A low profiling timeout might also kill off good implementations. It should be noted that the optimal timeout does not depend on the actual execution time, since profiling runs the implementation multiple times and might even dynamically adjust the number of executions. In practice, a low timeout should be set first. If the log shows too many timeout errors, the timeout can be increased. 5 seconds seems to be a good value for GPU target devices, while 20 seconds or more are appropriate for CPU autotuning.

Batch size This determines how many configurations are selected and built in parallel for every autotuning iteration, which can speed up autotuning considerably, especially if a large number of CPU cores are available on the client to run many compiler processes in parallel. The number of cores is the default value for this parameter. For larger batches, the model is updated and queried less frequently, but in general there seems to be no detrimental effect of having a high batch

size. It should be noted that this is not the same as the batch size of the model, which would change the shape of the tensor operators.

Transfer learning This determines whether or not transfer learning is used between jobs, i.e. if the data from the global autotuning database are used to train the cost model at the start of each task. Between tasks, there is always transfer learning. Usually, transfer learning should be enabled for the most optimal inference performance results. However, we disable transfer learning for experiments to guarantee a fair comparison between earlier and later ones.

3.3 Capabilities

Using SimpleTVM and our knowledge about proper parameter settings as foundation, we evaluated how TVM’s inference performance compares to state-of-the-art manual tensor operator libraries. We use TensorFlow 1.14 as baseline since it is a popular framework for DL applications. The XLA optimizer as well as CPU-specific instructions and cuDNN for GPU are enabled. Autotuning with TVM was executed with 2000 trials, so the numbers should represent the optimal implementation. For evaluation, we test a Mobilenet with a batch size of one on two mobile-grade CPUs (Intel Core i5-5300U and i5-7300U), a server-grade CPU (Intel Xeon E5-2650 v3) and a high-end GPU (NVIDIA Tesla K80). The same two images were used as model input every time.

Figure 7 shows how, in all cases, TVM with autotuning has a performance advantage over TensorFlow. Especially on CPU, inference takes 44% (18.4 ms) less time for the i5-5300U and 30% (5.9 ms) for the i5-7300U. For the server-grade CPU, inference takes 87% (19.9 ms) less time because TensorFlow is very slow, presumably due to the lack of modern SSE4 instructions which TVM can much better work around. But also for the GPU, the autotuned version takes 17% (0.9 ms) less time, albeit measurement inaccuracies are possible on this small scale. Nonetheless, even a similar performance is impressive considering that no human expert knowledge was required and autotuning took less than 6 h for CPU and less than 12 h for GPU (due to the much larger search space). Further results for a wider variety of devices and models, including recurrent neural networks, are provided in [13] and show similar improvements.

TVM performs better than TensorFlow on CPU even without autotuning. Graph-level optimizations alone are enough to result in faster inference, taking 24% (9.9 ms) less time for the i5-5300U and 7% (1.4 ms) for the i5-7300U. Since only a series of pre-defined transformation passes are applied to the Relay program, graph-level

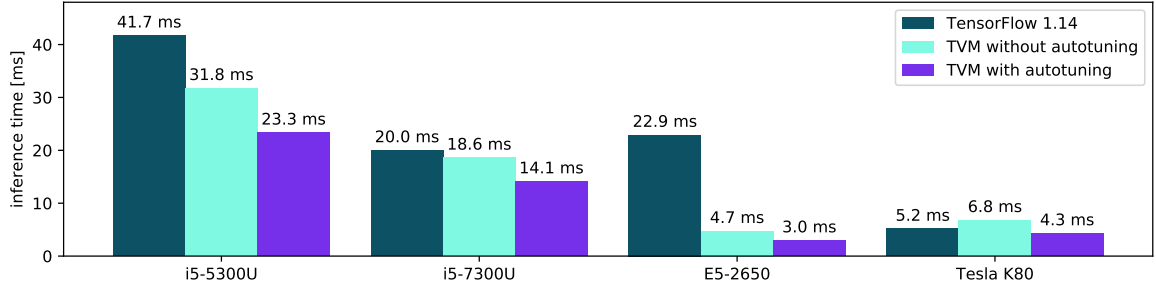


Figure 7: Inference performance with TensorFlow and TVM

optimization is performed in a matter of seconds. However, non-autotuned TVM cannot keep up with TensorFlow on GPU, it takes 31% (1.6 ms) longer.

These results show that TVM is on par with manual optimization, at least for our limited set of evaluation scenarios. Furthermore, TVM is under active development and can be expected to show further performance improvements in the future.

3.4 Limitations

While the autotuning results are promising, we found that the autotuning process suffers from some fundamental restrictions inherent to the current design which limit its efficiency. For all real measurements in this section, we evaluated autotuning of a ResNet-18 (12 convolutional layers, 1 fully-connected layer) with 2000 iterations per task on two machines with two Intel Xeon E5-2650 v3 CPUs and four Tesla K80 GPUs each, one machine as client and the other as target device.

3.4.1 Resource Utilization

Since stages in autotuning depend on results of the previous stages (configurations are required for building, executables are necessary for profiling, time measurements are used to update the model), they need to be executed in sequence. Because stages are distributed over two different machines, the result is a lot of idle time on both the client and the target device. This sub-optimal autotuning resource¹ utilization is exemplified in Figure 8.

Measurements with our test setup showed that, for a total autotuning completion time of 14.5 h, the client is idle for 6.6 h (45%) and the target device for 7.9 h (55%). Since computation resources, especially DL accelerators, are rather costly, we want

¹We define *resource* as a machine that executes some stage of the autotuning process, e.g., the target device or the hardware that the client runs on.



Figure 8: Resource utilization during autotuning

to minimize resource idle time. If existing resources are utilized better, it is not necessary to acquire new hardware. Given that the individual occurrences of idle time are long enough, it might be possible to use the slack for other computations in between. Indeed, the mean execution time for each stage is as follows: 8.3s for building, 39.5s for profiling, and 44.0s for updating the model and selecting the next batch. This is long enough to reasonably assume that resource utilization can be improved by sharing the device, but interference must be prevented.

3.4.2 Scalability

In preparation for enabling autotuning on a larger scale, we need to examine the scalability of the current design. Scalability in this section refers to the ability to run an arbitrary number of autotuning jobs at the same time without sacrificing efficiency and result quality. We define objectives, that a scalable solution should satisfy:

1. High inference performance, since it is the ultimate goal of autotuning
2. Low amount of required hardware, since additional devices are costly
3. Low autotuning time, since autotuning takes long

These objectives are listed in order of priority. Good inference performance is the primary objective. It also obviates the need to buy new hardware, and because there usually is a large amount of inferences, a longer execution time for autotuning is negligible in the long-term. Rapid autotuning is desirable nonetheless.

We compare two setups for scaling that are possible using only the components that TVM comes with by default, schematically depicted in Figure 9. For the sake of simplicity, we only show two jobs, but this generalizes to any higher number. The evaluation of both setups with regard to the previously defined objectives is summarized in Table 1.

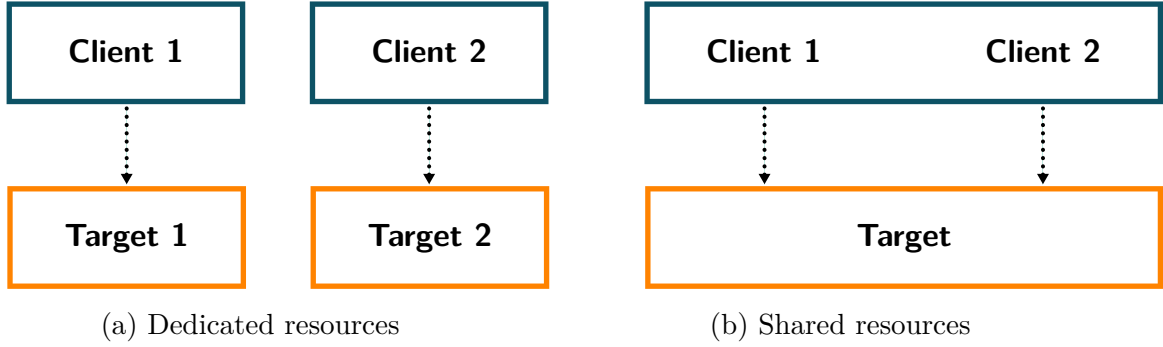


Figure 9: Setups for scaling autotuning

Setup	Hardware req.	Inference speed	Autotuning time
Dedicated resources	2x	High	Low
Shared resources	1x	Low	High
Optimum	1x	High	Low

Table 1: Comparison of scaling setups

Dedicated resources In the first setup (Figure 9a), each job is run on its own set of resources. This means they are completely independent and do not affect each other. Both the resulting inference performance and the autotuning time are as good as possible, since they are equal to single-job autotuning. However, additional resources need to be acquired for every new job, which is not an economically feasible approach on a larger scale. Adding resources on demand from a cloud computing platform such as Amazon EC2 or Azure VM might work for the client machine, but since the actual target device needs to be used for profiling, which is likely to not be available on those cloud platforms, this is not a satisfactory solution. Alternatively to having separate sets of machines, the same machines could be used to run the jobs in sequence. This would trade off the amount of hardware that is required for autotuning time. However, this will result in excessive autotuning time with an increasing number of jobs.

Shared resources In the second setup (Figure 9b), all jobs run in parallel on the same resource. To share the target device, we launch one RPC server for each client per GPU in our test setup. Due to resource sharing, only one set of machines is required, which is good in terms of cost. However, this is very probable to lead to interference, when multiple jobs execute a stage on the same resource simultaneously. Interference has detrimental effects on both inference performance and autotuning time. Interference on the client will slow down compute-intensive stages like building or model updates (50%–70% CPU usage). Modern CPUs are able to parallelize using multiple cores, but because building and model training uses all cores, time-sharing between multiple processes needs

to be employed by the kernel. Process idle times and context switching overhead increase the autotuning time considerably. Assigning dedicated CPU cores per job will also make the process slower due to decreased resources per job and will not work when scaling up.

Even more critical is interference on the target device. This distorts the profiling results, so good implementations might seem to perform bad if another job profiles on the same target device in parallel. Deceiving measurements will then also lead to an inaccurate cost model. Since good inference performance is the prime objective of large-scale autotuning, this setup is also not satisfactory.

Both setups do not meet all objectives we set for large-scale autotuning. Especially when scaling up to more parallel jobs, efficiency deteriorates significantly. We can conclude that the current implementation and architecture of autotuning in TVM does not scale well.

Having regarded possible scaling approaches and the reasons for their benefits and shortcomings, we can formulate two aspects an optimal solution needs to consider in order to satisfy the objectives:

- Resources must be shared and utilized fully before adding new servers to minimize the required hardware for cost saving
- Interference must be prevented to guarantee a high inference performance and low autotuning time

3.4.3 Similar Problems

Our literature review was not successful in finding a solution to scale up autotuning. However, we can generalize the problem statement; we are looking for a solution to share available resources optimally between multiple tasks that are partially idle due to some dependency. From this point of view, we find two papers solving a similar problem.

[14] increases parallelism of a hierarchy of tasks and subtasks on multiprocessor platforms. Subtasks have control and data dependencies, but tasks are independent of each other. They employ a mix of exact and heuristic scheduling algorithms at design-time to interleave sub-tasks of different tasks while respecting the dependencies between the subtasks of a single task. The result is a 37% shorter execution time with increased resource utilization. The hierarchical structure of tasks and subtasks is similar to jobs and stages in autotuning.

[15] enables sharing of GPU cores by multiple kernels. In current GPU architectures, concurrently launched kernels use separate cores. However, interleaving of code from multiple kernels on the same core allows them to minimize core idle time introduced by memory latency and increase the throughput of benchmarking applications by 7%.

Both [14] and [15] improve parallelism for processing units with multiple cores while we want to improve parallelism for jobs running on distributed machines. While the scale is different, the problem is the same, making the interleaving approach relevant for our solution.

4 Autotuning Scheduler

Interleaving the stages of multiple jobs is our key concept for enabling large-scale autotuning. [14] uses a design-time scheduler to create a program with good concurrency. We need to dynamically schedule incoming jobs, so the schedule cannot be predetermined. This motivates the need for an additional component in the autotuning architecture that orchestrates the execution of multiple jobs. In this chapter, we describe the design and implementation of our central scheduler that controls stage execution.

4.1 Design

Our scheduler possesses the two features that have been determined to be imperative for optimal large-scale autotuning:

- Computation resources are shared between jobs by *interleaving*. This facilitates good resource utilization since the idle time of one job is leveraged to execute another job, which saves hardware and costs as a result. However, stage dependencies of a single job must be maintained.
- Interference between jobs is prevented. This guarantees that inference performance and autotuning time are as good as possible. The scheduler needs to check if the resource that will be used by the next stage is free before execution. This might necessitate the postponing of stage executions if the stage is ready before the resource becomes free.

These two features not only make it match the optimal solution, but also do they solve the problem of bad resource utilization of single-job autotuning by leveraging that shortcoming.

Figure 10 illustrates round-robin-based interleaving with an example of two jobs. Significant events are marked with numbers. Job A and Job B are started at the same time, but assume that the scheduler knows about A earlier. The first stage of A is executed, then the first stage of B. Once B finishes, the scheduler decides it is A's

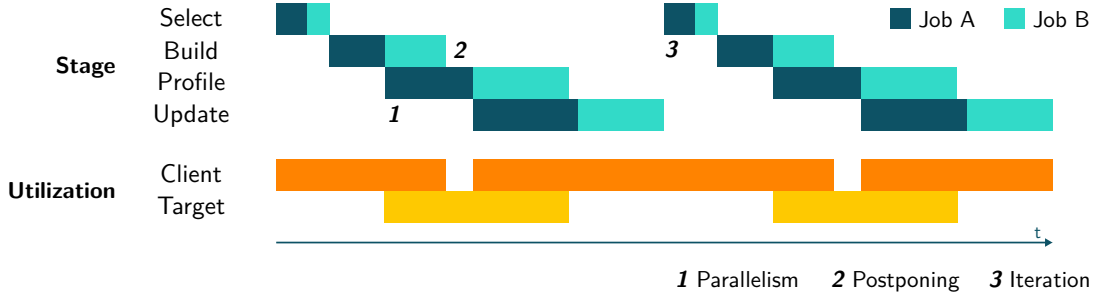


Figure 10: Interleaving of multiple autotuning jobs

turn again and executes its second stage. Once A finishes the second stage, the client machine is free and B can execute the second stage. At the same time, A is ready to execute the third stage which will run on the target device. Since the target device is not in use, A can execute the profiling there in parallel to B’s building since they use separate resources (**1**). Building does not take as long as profiling, so A is ready to profile before B finishes its profiling stage. Therefore, A’s third stage is postponed until the target device is free (**2**). A’s profiling and B’s update model can, once again, execute simultaneously since they use distinct resources. After one iteration of all four stages, the process starts anew with the first stage (**3**). This continues until both jobs are done. In a real scenario, new jobs might appear while other jobs are already running. The scheduler simply adds them to its list of jobs and includes them in the interleaving. Note how the resource utilization in Figure 10 is much improved over the single-job autotuning in Figure 8 due to overlapping and sharing. Especially on the client device, utilization has almost been maximized since three of the four stages use the client.

4.1.1 Autotuning Decomposition

The default autotuning process is monolithic and can be regarded as a blackbox from the outside (Figure 11a). This means, the autotuning loop can be started, and it does not finish until the whole job is completed. Once a stage finishes, the next one is executed immediately. However, the scheduler needs to be able to control the execution of the individual stages because it needs to prevent interference by means of delayed execution. This necessitates the decomposition of the autotuning process into schedulable units, corresponding to the stages (Figure 11b). The client does not execute any of the stages on its own. Rather, it provides an interface to execute schedulable units and waits for an external trigger to do so. The controlling part of the autotuning loop can now run in another component, such as the scheduler.

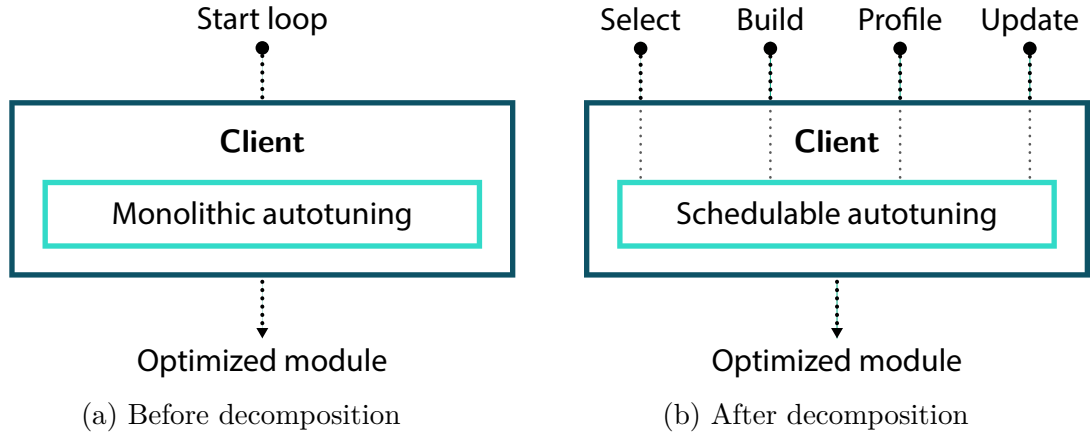


Figure 11: Client interface before and after decomposition

4.1.2 Scheduling Algorithm

To fulfill the task of interleaving while preventing interference, the scheduler needs to know four pieces of information for every job:

- Whether the current stage is done; lets scheduler know if the job is ready for next stage
- Whether the whole job is done; lets scheduler know if the job needs to be considered in the future
- Which resource the current stage runs on; lets scheduler know which resources are free and which are in use
- Which resource the next stage will run on; lets scheduler postpone stage execution to prevent interference

We call this *load-awareness*; being aware of the state of jobs and resources as well as their relationships. Theoretically, this allows the scheduler to work not only with TVM's autotuning but any process that can provide this information. Designing our scheduler to be agnostic of the underlying process also simplifies the algorithm since no state of that process, i.e. the progress or position in the loop, needs to be considered.

In case multiple jobs are ready to execute a stage, the scheduler needs to decide which one to run. The simplest approach is a round-robin algorithm which iterates over the jobs in the order they were started, picking the first one that is ready. More sophisticated approaches might apply some logic to decide on a job which would maximize resource utilization but keep the average autotuning time low. However, we choose the round-robin algorithm for our first version. It is easy to implement and

works reasonably well for an arbitrary number of jobs, which allows us to proof our concept.

We present two algorithms that perform interleaving. The greedy algorithm lets a job execute the next stage directly after the previous one finishes, provided the resource is free. This might preempt the resource from another job who is already waiting to execute a stage on it. On the other hand, the fair scheduler features a queue for every resource, so the first job to be ready to run a stage on a resource will be the first one to actually use it. This coincides with the interleaving previously shown in Figure 10. Each job is controlled using an interface that provides the four pieces of information necessary for scheduling (`is_stage_done()`, `is_complete()`, `previous_resource`, `next_resource`). The scheduler keeps a list of jobs (`jobs`) in the order in which they were registered. With each call of `next()`, the subsequent job is returned, wrapping around after the end is reached. Furthermore, the scheduler keeps a list of resources which can be marked as free or busy (`mark_as_free()`, `mark_as_busy()`).

Greedy interleaving The algorithm for greedy interleaved scheduling is presented in Listing 2. The algorithm is an infinite loop, with each iteration operating on one job (Line 1). First, the next job in the job list is retrieved (Line 2). If there is no job, the algorithm tries again until one is registered (Lines 3–4). If the job is still busy executing a stage, it is not considered further in this iteration (Line 6). If the stage has finished, the resource it used is marked as free (Line 7). If the stage was the last stage in the job, the job can be removed from the list of jobs (Line 8–9). Otherwise, the job is ready to continue. If the next stage’s resource is free, that resource is marked as busy and the stage is executed immediately (Lines 11–13).

```
1 while True:
2     job = jobs.next()
3     if not job:
4         continue
5
6     if job.is_stage_done():
7         mark_as_free(job.previous_resource)
8         if job.is_complete():
9             jobs.remove(job)
10            continue
11        if is_free(job.next_resource):
12            mark_as_busy(job.next_resource)
13            job.next_stage.execute()
```

Listing 2: Greedy interleaved scheduling pseudocode

Fair interleaving The algorithm for fair interleaved scheduling is presented in Listing 3. Each resource has a queue which contains stages that are ready and will use that resource (Line 1). The algorithm is an infinite loop (Line 2), with each iteration consisting of two phases: scheduling and execution. In the first phase, the scheduler checks for each job if the current stage is done (Lines 4–5). Busy jobs are not regarded further. If the stage is done, the resource that was used is marked as free (Line 6). If the job is complete, it can be removed from the list of jobs (Line 7–8). Otherwise, the job is ready to continue and the next stage is added to the queue of the resource that it will run on, if it not already there (Line 9–10). In the second phase, the scheduler iterates over each resource and the corresponding queue (Line 12). If the resource is free and there are pending stages for that resource, the first stage in the queue is dequeued and executed (Lines 13–15). The respective resource needs to be marked as busy (Line 16).

```

1 queues = {r: [] for r in resources}
2 while True:
3     # Phase 1: Round-robin scheduling
4     for job in jobs:
5         if job.is_stage_done():
6             mark_as_free(job.previous_resource)
7             if job.is_complete():
8                 jobs.remove(job)
9             elif not job.next_stage in queues[job.next_resource]:
10                 queues[job.next_resource].enqueue(job.next_stage)
11     # Phase 2: Execution
12     for resource, queue in queues:
13         if is_free(resource) and len(queue) > 0:
14             stage = queue.dequeue()
15             stage.execute()
16             mark_as_busy(resource)

```

Listing 3: Fair interleaved scheduling pseudocode

Additionally to interleaving, our scheduler supports two other strategies for executing multiple jobs, which will be used in the evaluation for comparison.

Sequential scheduling (Listing 4) works identical to single-job autotuning. Multiple jobs do not run in parallel but after another, so the next job is only started when the previous one finishes. This renders consideration of resource free/busy state unnecessary and stages do not need to be postponed. However, since it is controlled by the scheduler, the sequential strategy has the same overhead introduced by adding the scheduler component as interleaving, e.g., communication between scheduler and client or scheduling itself. This allows for a fair comparison.

```
1 while True:
2     job = jobs.next()
3     if not job:
4         continue
5
6     while not job.is_complete():
7         if job.is_stage_done():
8             job.next_stage.execute()
9     jobs.remove(job)
```

Listing 4: Sequential scheduling pseudocode

```
1 while True:
2     current_jobs = jobs
3     if len(current_jobs) < 2:
4         continue
5
6     while not any([j.is_complete() for j in current_jobs]):
7         if all([j.is_stage_done() for j in current_jobs]):
8             [j.next_stage.execute() for j in current_jobs]
9     jobs.remove_all(current_jobs)
```

Listing 5: Synchronous scheduling pseudocode

Synchronous scheduling (Listing 5) forces parallel execution of the same stage of multiple jobs on the same resource, making it the exact opposite of interleaved scheduling. Postponed stage execution is applied here to guarantee full interference. We use this strategy to evaluate the worst case effect of interference. However, this only works for equal jobs, since there needs to be symmetry between stages of all jobs.

4.1.3 Autotuning Process

Figure 12 shows the autotuning process with our scheduler (compare with scheduler-less, single-job autotuning in Figure 4). There are two discrete control flows at work. The per-job autotuning control flow is the same as before, spanning client and target machines. The scheduling control flow is on a higher level, incorporating multiple jobs and spanning scheduler and clients. A new job needs to be made known to the scheduler first so it can consider the job in scheduling. One client is responsible for running exactly one job, registering that job with the scheduler when launching. The client is then ready to receive control commands to execute individual stages. The client exposes the interface which is required by the scheduling algorithm while calling the respective TVM methods internally. Clients can share target devices since

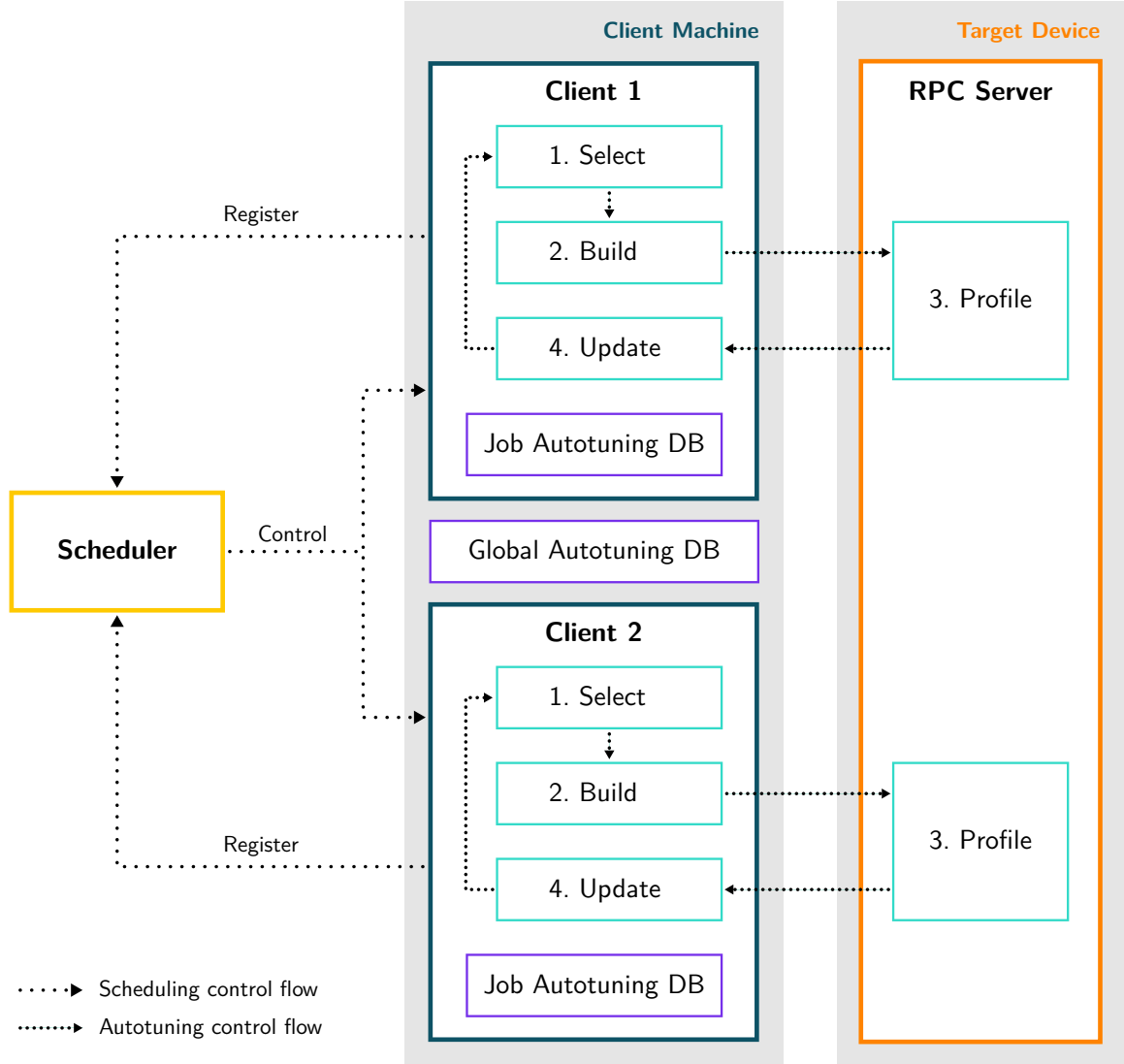


Figure 12: Autotuning process with scheduler

the scheduler prevents interference of the profiling stages of multiple jobs on the same resource. The normal autotuning process is then executed, however not in one monolithic step, but stage for stage, enabled by the decomposition. Possibly there are some waiting times between stages introduced by postponed execution.

Scheduler and clients live in different processes, usually even separate containers or physical servers. While the autotuning control flow already exists in form of in-process method calls and RPC for remote profiling, the scheduling control flow requires its own RPC infrastructure to enable communication between the scheduler and the clients.

In multi-client scenarios, the cost model of each client is initialized by transfer learning from the global autotuning database on the client machine, which contains data from all jobs that have previously been executed on that machine. During autotuning,

measurement results are written into a job-specific database which is merged back into the global database when the job is complete.

4.2 Implementation

Since TVM provides the API for autotuning in Python only, we use Python 3.6 for our implementation. It is intended as a proof of concept which we want to develop rapidly to see if the interleaving scheduler delivers the expected results. Therefore, we create an implementation that does not offer much flexibility or fault tolerance. However, it is sufficient to perform experiments in our test environment. The scheduler implementation is built on top of SimpleTVM for interfacing with TVM's autotuning.

For the beginning, only a single client machine is supported, but an arbitrary number of clients can run on it. Multiple target devices can be utilized for profiling, but the scheduler regards all target devices as a single resource. This coarse granularity is another decision to facilitate simple implementation.

4.2.1 Components

The implementation of our scheduler is distributed over multiple components. Same-machine components interact via in-process method calls, but RPC is required for cross-machine calls. We created a simple HTTP-based RPC protocol to support communication between scheduler and clients, with both scheduler and client acting as HTTP server and client. However, the protocol is very specific to the required interface and not general-purpose. If requests fail, they are retried three times with exponential backoff.

The components act as layers of abstraction sitting on top of TVM to eventually hook into the actual autotuning process. Figure 13 shows the whole stack including method interfaces. We look closer at each component, from top to bottom.

Scheduler The **Scheduler** implements the interleaved scheduling algorithms from Listings 2 and 3 as well as the sequential and synchronous strategies. The strategy can be specified when starting the scheduler. Instead of a list of jobs as in the algorithm, the scheduler actually keeps a list of the clients which run those jobs. The scheduler provides an RPC interface for the client to register (not depicted in the figure) and will keep a reference to that client to enable controlling of its autotuning job. Additionally, some error handling is added to,

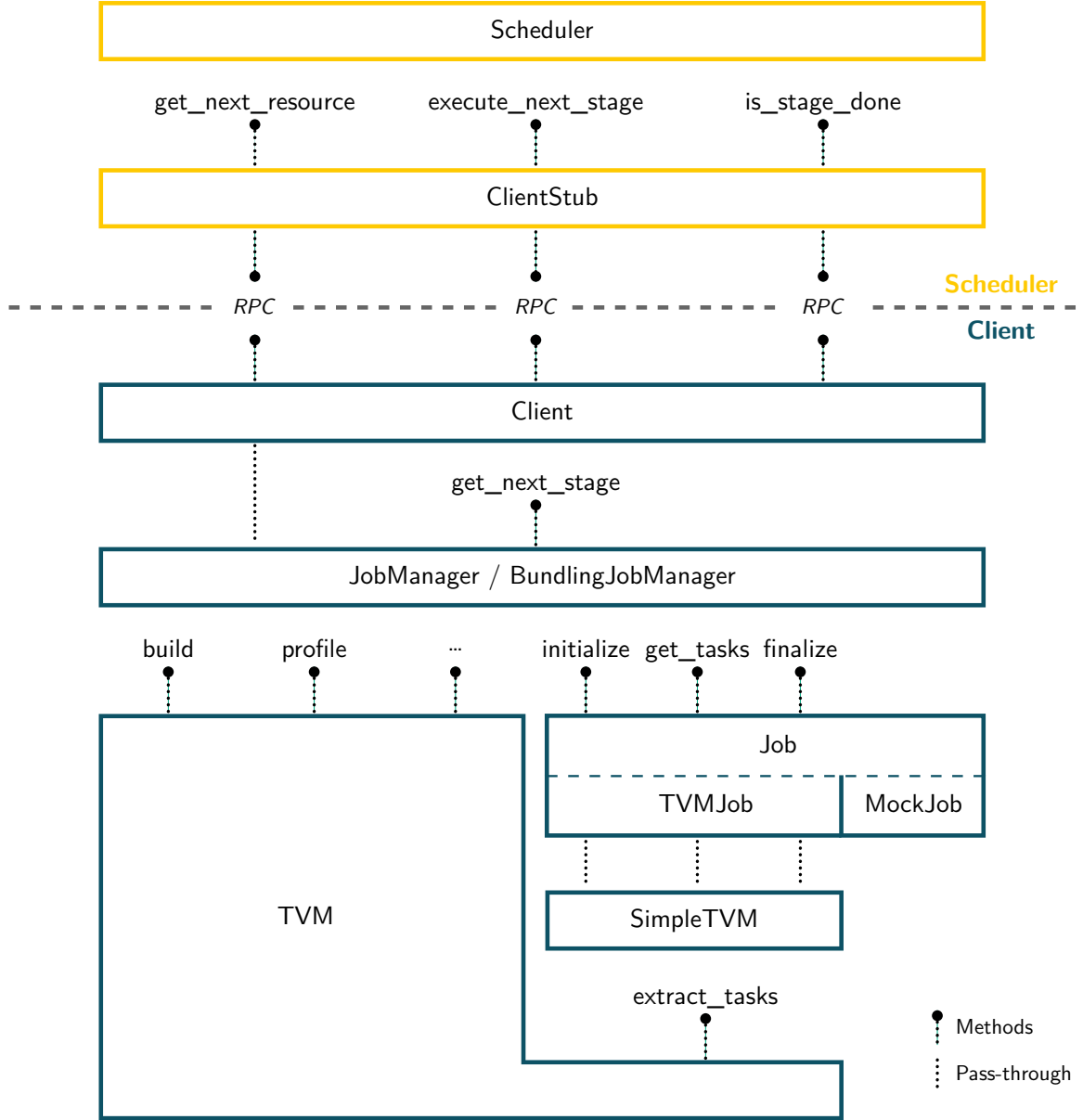


Figure 13: Layers and components of scheduler implementation

e.g., remove clients if they become unreachable. Resource state is implemented in a resource-to-boolean dictionary, indicating if a given resource is free.

ClientStub The `ClientStub` does not perform any functionality on its own, rather it acts as an abstraction of the `Client`'s RPC interface so the scheduler can call client methods as if they were an in-process object, allowing for a clean implementation of the scheduling algorithms. The `ClientStub` passes through calls to the actual `Client` but handles RPC including serialization/deserialization of variables and error handling.

Client The **Client** provides the interface that is required by the scheduler to control jobs. It registers with the scheduler upon launch and exposes three methods via RPC:

- `get_next_resource` returns the resource that the next stage will run on, or `None` if the job is complete, effectively combining two functionalities into one method. Furthermore, if the job is complete, the client shuts down after sending the RPC response.
- `execute_next_stage` calls the function for the next stage asynchronously and keeps a reference to the stage's thread in form of a future. If another stage is currently running, this methods fails.
- `is_stage_done` returns a boolean, denoting if the stage that was previously executed using `execute_next_stage` has finished. If no stage has been executed yet or the future indicates that the stage's thread has terminated, `True` is returned, `False` otherwise.

JobManager The **JobManager** is responsible for negotiating between the simple resource- and stage-based interface required by the scheduler and the more complex interface of TVM's autotuning. This conversion from a stateless to a stateful interface requires that the **JobManager** keep track of the current progress of the autotuning process, i.e. the position in the autotuning loop identified by task and stage. Effectively, it decides the order of stages in the loop, when the loop starts again, and when the loop terminates. This allows it to implement two methods: `get_next_resource` is called directly by the client to get the resource that the next stage will run on, `get_next_stage` returns a method containing the stage functionality which is executed by the client's `execute_next_stage`.

The logic by which the next stage is decided in `get_next_stage` is presented in Listing 6. `get_next_resource` follows a similar logic but returns the resource instead of methods. The job's state is determined by flags indicating if the job has been initialized and finalized as well as the current task and the stage in that task (Lines 1–2). If the job has not been initialized, the initialization method is returned and the `initialized` flag is set (Lines 4–6). Then, the autotuning loop for the first task is started. With each successive call of `get_next_stage`, `current_stage` is incremented and the respective stage's method is returned (Lines 8–22). The stage methods come from either the **Job** or the decomposed TVM autotuning API. At the end of one loop iteration, which corresponds to one batch, the `current_stage` is set to the first stage of the loop (Line 21). If the current batch was the last batch of the task (because either the

search space has been exhausted or the specified number of trials has been reached), `current_task` is incremented and `current_stage` is reset to the task initialization so autotuning can commence for the next task (Lines 15–18). If all tasks have completed, the autotuning job is done (Line 7). Next comes the job finalization method (Lines 23–25), after which `None` is always returned (Line 26).

```
1 initialized = finalized = False
2 current_task = current_stage = 0
3
4 if not initialized:
5     initialized = True
6     return initialize_job_fn
7 elif current_task < number_of_tasks:
8     current_stage += 1
9     if current_stage == 1: return initialize_task_fn
10    elif current_stage == 2: return select_batch_fn
11    elif current_stage == 3: return build_fn
12    elif current_stage == 4: return profile_fn
13    elif current_stage == 5: return update_model_fn
14    elif current_stage == 6:
15        if last_batch:
16            # Go to next task
17            current_task += 1
18            current_stage = 0
19        else:
20            # Go to select batch stage
21            current_stage = 1
22        return finish_batch_fn
23 elif not finalized:
24     finalized = True
25     return finalize_job_fn
26 else: return None
```

Listing 6: Pseudocode of `JobManager`’s stage decision logic

There is another version called `BundlingJobManager` which exposes the same interface but instead of returning each stage individually, it bundles the stages for one resource into a single, larger stage. Thus, there are effectively only two stages, one for the client (select configurations, build, update model) and one for the target device (profiling), as opposed to `JobManager` where each stage is individually schedulable.

Job The `Job` is an abstract class which acts as interface specification for `TVMJob` and `MockJob`. Jobs contain a collection of tasks as well as the initialization and finalization methods used by the `JobManager`.

TVMJob The `TVMJob` represents one autotuning job. It only passes calls through to `SimpleTVM` which contains the implementations. The initialization method only stores a timestamp of the autotuning begin for measurement. The finalization method collects error and time statistics from the autotuning process and inserts them into the benchmarking context. Additionally, it merges the job-specific autotuning database with the global database. The collection of tasks is extracted using a method from `TVM`, but modified slightly by `SimpleTVM`.

MockJob A `MockJob` is a drop-in replacement for `TVMJob`, exposing the same interface but not performing actual autotuning. Rather, it simulates work by sleeping and prints to the console when stages are started and finished. The number of mock batches and tasks can be specified, as well as a time stretch factor to slow down or speed up the simulated work. This allows rapid debugging of the scheduler algorithms because no infrastructure like servers and the tracker need to be set up and speed can be controlled. The mock stages' proportions are about the same as in real autotuning, e.g., profiling takes longer than building. It does not depend on any other components such as `TVM`.

4.2.2 Usage

Switching from scheduler-less default autotuning to scheduled autotuning is trivial if the former is already being used, as shown in Listing 7. Creation of the `SimpleTVM` object and import of the model is identical for both variants (Lines 1–2). Default autotuning is started by calling the appropriate method directly with the autotuning options such as number of trials or profiling timeout (Line 5). For scheduled autotuning, a `TVMJob` needs to be created first (Line 8). A client for the job is then started with the host names of the client and scheduler machines as parameters (Line 9). After autotuning, the optimized module can be built, saved and evaluated as in Listing 1.

```
1 tvml = SimpleTVM(BenchmarkingContext('gpu'), rpc_tracker=('host', port))
2 tvml.from_model('resnet')
3
4 # Scheduler-less
5 tvml.autotune(options)
6
7 # Scheduled
8 job = TVMJob(tvml, options)
9 Client(job, client_host='client', scheduler_host='scheduler').start()
```

Listing 7: Comparison of default and scheduled autotuning

4.2.3 Challenges

For our very first prototype, we wanted all components to run in a single multi-threaded process, so we could evaluate our approach quickly without having to implement the RPC protocol. However, Python does not support true multi-threading due to the global interpreter lock. This lock simplifies memory management for the interpreter but only one thread can execute code at a time because of it. Python’s recommended replacement is launching other interpreter processes from within the code, however this requires serialization of objects for inter-process communication. Our nested class structures and passing around of methods that are created dynamically caused problems with this serialization. That is why we had to implement RPC before the actual scheduler to enable separate processes from the beginning.

Since multi-job autotuning with a scheduler requires setup of the infrastructure and is rather slow, we created the `MockJob` class for improved debugging and testing of the scheduler functionality. This allows us to evaluate design choices in the scheduling algorithm much more rapidly. Moreover, the isolation from TVM narrows the room for errors which facilitates focusing on scheduler issues without being disrupted by problems caused by autotuning.

4.3 Autotuning as a Service

Setting up the multi-job autotuning infrastructure is cumbersome because a lot of components need to be configured and deployed, most of them imperatively in code. Additionally, provisioning of new client and target machines is not done automatically by the scheduler, which limits the scale at which autotuning can be performed. This motivates the need for an Autotuning as a Service (AaaS) platform, which hides the complexity of building and maintaining the infrastructure and thus simplifies the application of autotuning, opening up the opportunities of autotuning for any developer of DL-based software.

We imagine a solution that allows users to submit their trained model and a declarative specification, including information such as the target device type and the target inference time to meet service-level agreements. The platform then automatically sets up the required machines and components, runs the autotuning process until the user’s target inference time is achieved (or no further optimization is possible), after which the optimized version is returned to the user. To this end, we propose a reference architecture for AaaS platforms. Our architecture builds on top of a container management and resource provisioning system, however it does not prescribe any

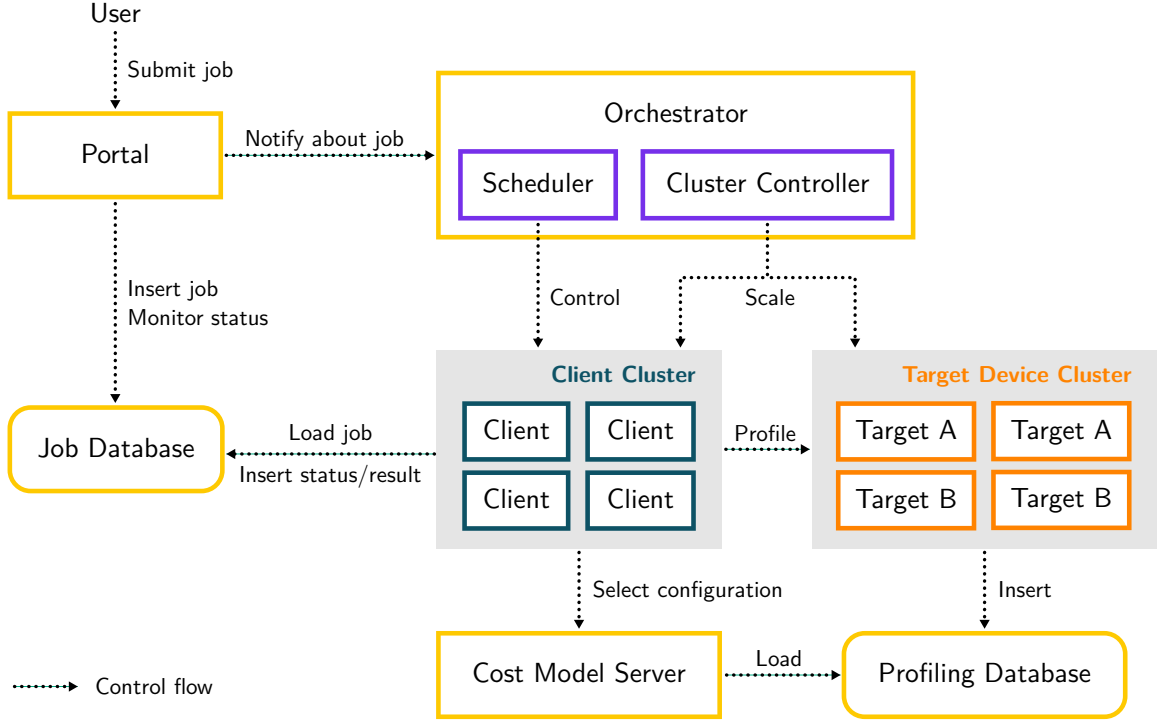


Figure 14: Autotuning as a Service reference architecture

specific product. [3] demonstrates how, for example, Kubernetes could be harnessed for these tasks, with a prototype delivering promising results.

Figure 14 shows our proposed architecture. Users can submit their jobs including the model and specification through the portal. First the portal inserts that job information into the job database, after which it notifies the orchestrator about the new job. The orchestrator, which contains both the scheduler and cluster controller, decides if the existing resources are sufficient, or if the client and target device cluster need to be scaled up, e.g., by booting new server instances on a cloud computing platform. The cluster controller, e.g., a Kubernetes Master, is responsible for executing the scaling. Then, the client is created on a client machine and the job is registered with the scheduler.

The scheduler now controls the autotuning jobs as described in the previous sections. However, the client’s responsibilities change in the AaaS scenario. The cost model is extracted to a central server that all clients share, so neither a local cost model nor a job-specific autotuning database is necessary. The cost model server keeps a separate model for every target device. Furthermore, the client does not request profiling servers from a tracker, but the target devices are explicitly assigned per profiling stage by the scheduler. Profiling results are not returned to the client, but instead inserted into a global profiling database, which the cost model server uses to update the models periodically. Clients store their state in the job database as opposed to the JobManager object, which makes the clients stateless so they can resume the job

in case of failure. This also allows the user to monitor autotuning progress through the portal. Once a job has completed, the resulting optimized TVM module is stored in the job database and the user can download it for deployment in his application.

Additionally, we make proposals to lower job completion time. Firstly, before starting the autotuning, inference performance of a newly submitted model is checked with the existing best configurations. If the user’s specifications are already met, no autotuning needs to be launched. Secondly, jobs can be parallelized by splitting them into multiple jobs. A job consist of a set of tasks that are independent of each other and can be autotuned separately. Moreover, the search space of each task can be divided to split one task into multiple. Leveraging task-level and search space-level splitting, autotuning can be accelerated if a user request a higher number of resource be allocated to them. Furthermore, unused resources can be utilized better if the total platform load is low.

One limitation of AaaS over on-premise autotuning infrastructure is the set of supported target devices types. The platform can provide support for common accelerators, but if more exotic or novel hardware is desired, AaaS will fall short. One solution would be to support profiling on devices outside of the target device cluster. However, appropriate security measures would need to be taken.

In summary, the AaaS platform makes four important changes to the existing multi-job, scheduled autotuning architecture:

- User-friendly portal instead of job specification in code
- Automatic resource provisioning/scaling instead of manual infrastructure setup
- Shared cost model between clients instead of job-local model
- External client state to recover from job failure instead of stateful clients

Instantiating such a platform is an important step towards making autotuning more accessible and enabling more real-time DL applications with little effort for developers.

5 Evaluation

To evaluate the impact of our scheduler, we compare the two approaches from Section 3.4.2 to our interleaving scheduler which aims to satisfy the “Optimum” solution. Additionally, we use the sequential scheduling algorithm as baseline while synchronous scheduling represents the worst case by forcing maximum interference. All scenarios are considered in terms of total autotuning completion time, time for the individual stages and the resulting inference performance.

Our evaluation environment consists of two identical machines with the following specifications:

- 2x Intel Xeon E5-2650 v3, 10 cores, 2.30 GHz
 - Hyper-threading enabled
 - AVX2 instruction set
- 128 GB main memory
- 4x Tesla K80 GPUs, 4992 CUDA cores, 24 GB memory
- Ubuntu 16.04.6 with Linux Kernel 4.4.0
- Python 3.6.8

Clients always run on the first machine, while profiling servers (TVM RPC servers) always run on the second machine whose GPUs are used as target devices. Scheduler and tracker run on the first machine to decrease network latency. They do not interfere with autotuning since they are not computationally intensive.

The number of clients and profiling servers changes for every experiment, as shown in Table 2. On “dedicated” servers, there is only one client. Two clients run on “shared” servers. Each profiling server of one client is assigned to a different GPU. Without interleaving, each client has its own set of four profiling servers, which might result in two servers if different clients being assigned to the same GPU. However, with interleaving, they can share the servers since they will never profile at the same time, eliminating competition.

Setup	Server	Scheduler	Bundling	Profiling servers
A	Dedicated	sequential	No	4
B	Shared	sequential	No	8
C	Shared	synchronous	No	8
D	Shared	interleaved-greedy	No	4
E	Shared	interleaved-fair	No	4
F	Shared	interleaved-greedy	Yes	4
G	Shared	interleaved-fair	Yes	4

Table 2: Evaluation setups

All setups are controlled by the scheduler so each experiment is affected by the—albeit minimal—overhead of scheduling and RPC. Two separate schedulers, one for each client, are used in B to achieve natural interference. “Bundling” indicates if stages of one resource are bundled into a single stage using the `BundlingJobManager`, or if each stage is individually scheduled.

Our test model is a ResNet-18, which consists of 12 convolution layers and one fully-connected layers, resulting in a job of 13 tasks. We autotune with 2000 trials per task and a profiling timeout of 5 s. Transfer learning from the global autotuning database is disabled so each experiment starts from the same, untrained cost model for a fair comparison. However, transfer learning is enabled between tasks. For time reasons, each experiment is only conducted once so the sample size is small.

A complete chart of all results can be found in Appendix A.

5.1 Results

First, we examine the impact of interference that was qualitatively described in Section 3.4.2. We compare A, the optimum in terms of autotuning time and inference performance, with B and C (Figure 15). B lets jobs naturally interfere, while C forces maximum interference as worst case scenario. Interference is particularly noticeable for model updates since they are rather computationally intensive, which is why we show model update time in every figure.

The baseline completion time from A is 14.5 h, 6.1 h (42%) of which are spent updating the model. Natural interference results in an increase in autotuning time by 4.1 h (+28%) while forced interference takes even 5.0 h (+34%) longer¹. This significant

¹For C, the time spent waiting for the other job to finish a stage introduced by the synchronous scheduling algorithm to force interference was subtracted from the total completion time since it would not occur naturally.

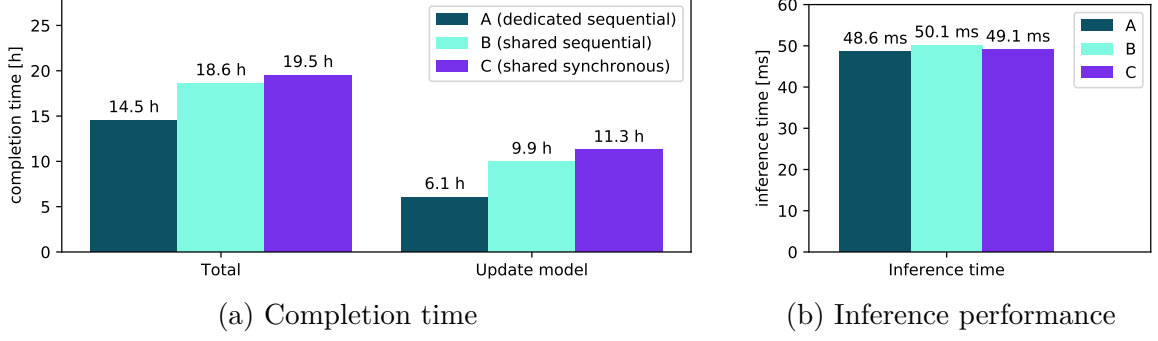


Figure 15: Impact of interference

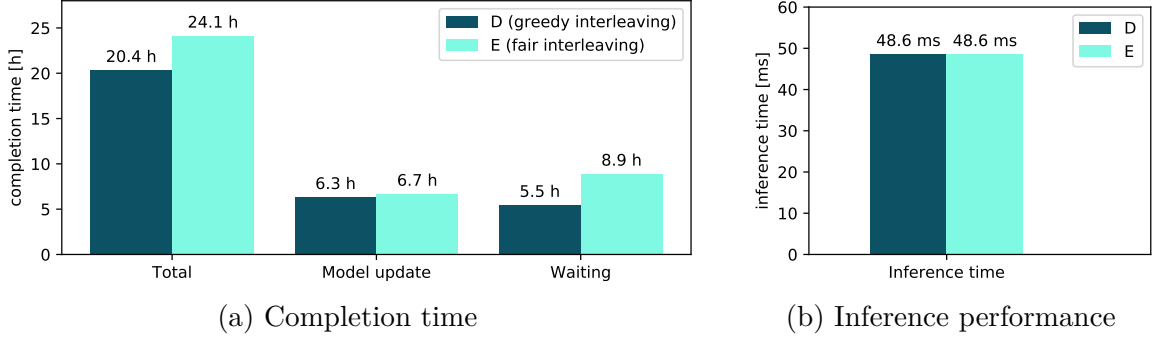


Figure 16: Greedy versus fair interleaving

increase can be attributed to slower model updates, which witness a strong decline in performance. Compared with the baseline, the total time spent updating the model increases by 3.8 h (+62%) for B and by 5.2 h (+85%) for C. Build time only increases marginally, with profiling being slightly faster, possibly due to more timeouts. The baseline inference time is 48.6 ms. B is 1.5 ms slower, C is 0.5 ms slower.

We examine experiments D and E to evaluate the greedy and fair scheduler algorithm design, comparing them with A as baseline (Figure 16). Using greedy interleaving, autotuning completes in 20.4 h, 5.9 h slower than A. While model update time stays about the same, 5.5 h (27%) of the total time is now spent waiting² for resources to become free to prevent interference. Fair interleaving is impacted even more by waiting, with total autotuning time increasing to 24.1 h, 9.6 h (+62%) slower than A. Model update time remains about the same again, but wait time accumulates to 8.9 h (37% of total autotuning). Building and profiling time do not change for both. While the total completion time deteriorates, inference performance matches the baseline of 48.6 ms for both D and E.

Finally, we examine the effect of employing stage bundling as opposed to individually scheduling them (Figure 17). The greedy algorithm is used in F, the fair algorithm is

²Only stage execution times can be measured. Wait time and non-stage times (transfer learning, file operations) need to be derived. Non-stage times are about 0.49 h, derived from A. Wait time for interleaving is calculated as follows: $t_{wait} = t_{total} - \sum_i t_{stage,i} - 0.49 \text{ h}$

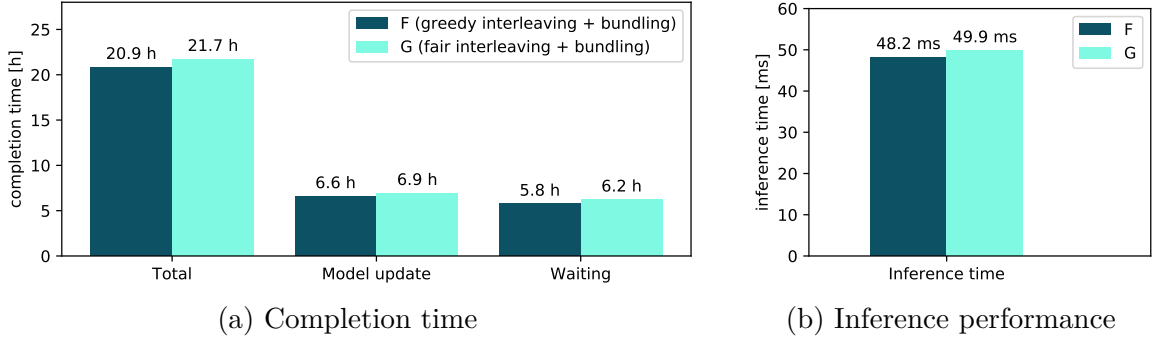


Figure 17: Greedy versus fair interleaving with bundling

used in G. With bundling, greedy interleaving performs about the same as without bundling. On the other hand, fair interleaving completes 2.4 h (-10%) earlier than the individually scheduled version due to a total wait time that is 2.7 h (-30%) shorter. As before, model update time remains relatively similar. Inference time in F is 0.4 ms faster than in D, while G’s is slower by 1.3 ms.

5.2 Discussion

The objective of large-scale autotuning is to achieve the autotuning time and inference performance of running only a single job on each server, but with shared resources to keep the amount of required hardware at a minimum. Single-job autotuning is represented by Setup A, which is why we use it as baseline for comparisons with the other experiments. Because our evaluation is limited to a single machine type and model, more experiments are required for a general analysis.

In B and C we show the detrimental effects of interference on performance which motivated the creation of a scheduler. If it were not for this, resources could be shared and jobs could be run in parallel without control by a central entity. B and C behave relatively similar because even without explicit scheduling, both jobs run approximately synchronous since they autotune the same model. In real-world scenarios with diverse workloads, there might be a more drastic difference between B and C. The fact that C has faster inference than B is presumably caused by the probabilistic nature of configuration selection. The difference in inference performance of both to the baseline is relatively small due to the large number of trials which even out the impact of interference. However, even slightly worse inference speed has a large impact with an increasing number of inferences.

Note how interference is avoided with the scheduler in D–G, which manifests in stage times (particularly model updating) that are comparable to single-job autotuning as

well as in a good inference performance. However, a significant amount of wait time is introduced which drives up total completion time. Greedy interleaving without bundling shows the shortest wait time in our experiments. Bundling does not show a big effect when comparing D and F because greedy interleaving naturally behaves like bundling. However, fair interleaving benefits much from bundling in our experiments, remedying the vast wait time. Inference performance is not optimal in G, more experiments might yield a better result.

Greedy interleaving works well in our experiments because both jobs are identical and the portion of time spent on the client machine versus time spent on the target device is roughly similar. As a result, wait time is reduced because, after an initial wait period, both jobs alternate nearly perfectly but mutually opposing between client and target. Despite these results, we believe that fair interleaving is the superior algorithm for scheduling heterogeneous jobs with varying complexity on a large scale, as would be the case with AaaS. Greedy interleaving would allow jobs to monopolize resources for an extended period of time, while fair interleaving allows for a finer control, especially if more resources are available. Bundling in conjunction with the fair algorithm is still an option to interleave homogeneous jobs on few resources, which would show the same effect as the greedy algorithm in our experiment.

Our scheduler is very rudimentary at this point, leaving much room for improvement. While the concept of interleaving seems promising, a more intelligent approach than fair round-robin could be employed to fill idle time optimally while keeping average wait time low. For example, a predictive scheduler could use knowledge from previous tasks to make a more educated decision that will maximize resource utilization. This would augment the notion of load-awareness with knowledge about the estimated stage execution time. Such a scheduler could then leverage the exact or heuristic interleaving algorithm from [14]. However, a more sophisticated algorithm might require finer control over the stages, which cannot be provided by the client’s current resource/stage interface. Trading off a complex client but simple scheduler for a thin, stateless client and autotuning-aware scheduler by shifting responsibilities might become a necessity.

Furthermore, the scheduler needs to be enhanced with support for multiple resources of the same type. At this point, all target devices are regarded as a single, atomic resource. More granularity is essential for production-grade implementations in an AaaS platform. This would go hand in hand with eliminating the tracker and letting the scheduler assign target devices to clients. A trivial refinement is replacing the busy waiting for a stage to be done in the scheduler algorithm, resulting in an excess

of redundant RPC calls, by a push-based approach, where the client notifies the scheduler once a stage has completed.

We used our scheduler only for autotuning with TVM. However, other autotuning frameworks like TC that have the same scaling issues might also profit from our solution. The AaaS platform might even provide support for a variety of autotuning frameworks.

6 Conclusion

Enabling large-scale autotuning is the foundation of any AaaS platform that uses our reference architecture. We demonstrated that our interleaving scheduler can leverage the fundamental weakness of resource idle time to make sharing of computation resources possible, which facilitates fast and economical autotuning while delivering good inference performance. This is a key step towards democratizing real-time DL applications to power exciting innovations in a wide variety of fields.

Future efforts based on our work should focus on creating a more intelligent scheduler algorithm to improve the efficiency of multi-job autotuning. A resource provisioning algorithm will then complete the orchestrator component, paving the way for the development of a mature AaaS product.

Bibliography

- [1] Statista, “In-depth: Artificial Intelligence 2019: Statista Digital Market Outlook,” 2019. [Online]. Available: <https://www.statista.com/study/50485/artificial-intelligence/>.
- [2] T. Chen, T. Moreau, Z. Jiang, L. Zheng, E. Yan, M. Cowan, H. Shen, L. Wang, Y. Hu, L. Ceze, C. Guestrin, and A. Krishnamurthy, *TVM: An Automated End-to-End Optimizing Compiler for Deep Learning*, 2018. [Online]. Available: <https://arxiv.org/pdf/1802.04799.pdf>.
- [3] J. Cho, F. Ahmed, L. Cao, P. Sharma, and D. Stiller, “Resonator: ML Autotuning-as-a-Service for Edge-to-Cloud Infrastructure,” 2019.
- [4] Ian Goodfellow, Yoshua Bengio, and Aaron Courville, *Deep Learning*. MIT Press, 2016. [Online]. Available: <https://www.deeplearningbook.org/>.
- [5] Lyve Data Labs, “Seagate Edge RX: A Smart Manufacturing Reference Architecture Solution,” 2019. [Online]. Available: https://labs.seagate.com/wp-content/uploads/sites/7/2019/06/TP711-2-1905US_Smart-MFG-Ref-Architecture.pdf.
- [6] Seagate, “Smart manufacturing moves from autonomous to intelligent: Inside Project Athena: Seagate’s internal AI edge platform,” 2019. [Online]. Available: <https://www.seagate.com/www-content/enterprise-storage/it-4-0/images/cs595-1-1901-seagate-athena.pdf>.
- [7] Y. Liu, Y. Wang, R. Yu, M. Li, V. Sharma, and Y. Wang, “Optimizing CNN Model Inference on CPUs,” 2019. [Online]. Available: <http://arxiv.org/pdf/1809.02697v3>.
- [8] Y. Hu, *Optimize Deep Learning GPU Operators with TVM: A Depthwise Convolution Example*, 2017. [Online]. Available: <https://tvm.ai/2017/08/22/Optimize-Deep-Learning-GPU-Operators-with-TVM-A-Depthwise-Convolution-Example>.

- [9] C. Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich, “Going Deeper With Convolutions,” 2015. [Online]. Available: https://www.cv-foundation.org/openaccess/content_cvpr_2015/papers/Szegedy_Going_Deepier_With_2015_CVPR_paper.pdf.
- [10] K. He, X. Zhang, S. Ren, and J. Sun, *Deep Residual Learning for Image Recognition*, [Online]. Available: <https://arxiv.org/pdf/1512.03385.pdf>.
- [11] N. Vasilache, O. Zinenko, T. Theodoridis, P. Goyal, Z. DeVito, W. S. Moses, S. Verdoolaege, A. Adams, and A. Cohen, *Tensor Comprehensions: Framework-Agnostic High-Performance Machine Learning Abstractions*, 2018. [Online]. Available: <http://arxiv.org/pdf/1802.04730v3>.
- [12] J. Roesch, S. Lyubomirsky, L. Weber, J. Pollock, M. Kirisame, T. Chen, and Z. Tatlock, “Relay: a new IR for machine learning frameworks,” in *Proceedings of the 2nd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, ser. MAPL 2018, New York, NY, USA: Association for Computing Machinery, 2018, pp. 58–68, ISBN: 9781450358347. [Online]. Available: <https://dlnext.acm.org/doi/abs/10.1145/3211346.3211348>.
- [13] T. Chen, L. Zheng, E. Yan, Z. Jiang, T. Moreau, L. Ceze, C. Guestrin, and A. Krishnamurthy, *Learning to Optimize Tensor Programs*, 2018. [Online]. Available: <https://arxiv.org/pdf/1805.08166.pdf>.
- [14] Z. Ma, F. Catthoor, and J. Vounckx, “Hierarchical Task Scheduler for Interleaving Subtasks on Heterogeneous Multiprocessor Platforms,” [Online]. Available: <http://doi.acm.org/10.1145/1120725.1120765>.
- [15] M. Awatramani, J. Zambreno, and D. Rover, “Increasing GPU throughput using kernel interleaved thread block scheduling,” 2013.

Glossary

autotuning job

one execution of the autotuning process for a model

autotuning resource

a machine that executes some stage of the autotuning process; e.g., the client machine or the target device

inference

overlapped execution of two programs on the same resource; generally results in slower performance

target device

the device that inference will be performed on; usually an accelerator located at the edge

tensor

a multi-dimensional generalization of vectors and matrices; commonly used data structure in deep learning

tensor operator

a function operating on one or multiple tensors such as convolution or matrix multiplication; commonly found in deep learning models

A Experiment Results

