



**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 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 3rd, 2019

Place

Date

Dominik Stiller

Abstract

Real-time computer vision applications with deep learning-based inference require hardware-specific optimization to meet stringent performance requirements. However, this approach requires vendor-specific libraries developed by experts for some particular hardware, limiting the set of supported devices and hindering innovation. The deep learning compiler stack TVM is developed to address these problems. TVM generates 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.

In this paper, we first explore the capabilities and limitations of TVM’s autotuning implementation. Then, we develop a scheduler to orchestrate multiple, parallel autotuning jobs on shared computation resources such as CPUs and GPUs, allowing us to minimize resource idle time and job interference. Finally, we reflect our design choices and compare the efficiency of our approach with the default, scheduler-less design.

Contents

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 Deep Learning	3
2.1 Artificial Neural Networks	3
2.2 Convolutional Neural Networks	3
2.3 Computation Graph Representation	3
3 Inference Optimization	4
3.1 Tensor Operator Optimization	4
3.2 Manual Optimization	5
3.3 Automated Optimization	5
4 SimpleTVM	7
4.1 Design	7
4.2 Exploration	7
4.3 Autotuning Limitations	8
5 Autotuning Scheduler	10
5.1 Design	10
5.2 Implementation	11
5.3 Autotuning as a Service	12
6 Evaluation	13
6.1 Results	13
6.2 Limitations	13
7 Conclusion	14
7.1 Future Work	14

List of Figures

List of Tables

List of Source Codes

1 Introduction

AI is increasing in popularity

AI is used in many different areas

Users aren't experts

Existing products for easier setup and deployment of training and inference infrastructure by offering AI infrastructure as a service

1.1 Problem

Common applications like industrial monitoring or autonomous driving require real-time performance

accelerator hardware with device-specific model optimizations needed

Currently manual optimization

Requires deep knowledge -> not easy for non-expert users

required: automated inference performance optimization (autotuning)

To offer it as a service so it can be used by a larger audience requires it to be scalable

Current autotuning does not scale well

To the best of our knowledge, there is no existing solution.

1.2 Scope

In this paper, we enable large-scale autotuning by sharing computation resources and controlling jobs with a central, load-aware scheduler

First step, develop framework to examine capabilities and limitations of TVM's autotuning in different configurations on multiple accelerator devices

Design and create a working proof-of-concept implementation

Evaluate our scheduler design and compare with default implementation

Propose an Autotuning as a Service architecture as base for future work

Thesis:

A load-aware scheduler makes large-scale autotuning more efficient in terms of

- autotuning speed
- resulting inference performance and
- hardware requirements

compared to the default design.

don't improve autotuning process itself, but propositions are made in future work

don't develop actual autotuning as a service product, but propose an architecture

2 Deep Learning

Machine learning

Consists of training and inference

2.1 Artificial Neural Networks

general structure with layers

layers: fully-connected, batch norm

basically a function mapping from input (image, text, data) to result (e.g. classification)

trend towards more layers to create more complex function which are more powerful

many layers: deep neural network

2.2 Convolutional Neural Networks

more layers: pooling and convolutions

important for computer vision

2.3 Computation Graph Representation

in essence just a lot of calculations

calculations can be represented in computation graph

graph describes which calculations need to be done and stays the same

tensor operators implement convolutions, pooling and need to be optimized for accelerator hardware

3 Inference Optimization

Often real-time necessary (examples)

Look closer at Seagate

requires low inference time, every saved ms makes big difference when a lot of inferences are being made

use of specific accelerator devices

Generic models perform poorly because they don't make full use of accelerator capabilities

Model needs to be attuned to accelerator

3.1 Tensor Operator Optimization

additionally to traditional training and inference deep learning workflow, we introduce inference performance optimization to meet real-time requirements

include graphic showing train-inference vs train-optimize-inference

to optimize for minimal inference time of whole network, we need to optimize every layer/tensor operator

especially conv2d, because there are many and they are very computationally intensive

many possible implementations

only few optimal ones for target device

focus on convolutions, as opposed to dense (why?)

one layer corresponds to one tensor operator with a specific shape

3.2 Manual Optimization

state of the art cuDNN and TensorRT and Intel MKL, taken as baseline

requires deep knowledge of target device, usually provided by vendor

limitations

- no support for new devices
- no support for unconventional shapes
- no support for new graph-level optimizations

elaborate limitations

3.3 Automated Optimization

vendor-agnostic

define autotuning job, task

describe autotuning process on high level

definition of search space (loop unrolling, tiling, threads), example code?

Problem: search space is very large (billions), and any one of them could be the best one

impossible to try all

look at both TVM and TC

has same or even better performance than hand-optimized libraries

show numbers

There are two frameworks that implement autotuning

3.3.1 TensorComprehensions

does not use machine learning

3.3.2 TVM

using machine learning

TVM is framework that proposed and implements autotuning

To create a scheduler, we examined TVM (commit id) with a few modifications to support measurements (check what else we changed)

written in Python and C++, interoperating

import from many frontends, compilation for many backends

has own graph-level and tensor operator-level representation

calls target-specific compiler

first extraction of tasks

schedules as abstraction with knobs

details of autotuning process

RPC allows autotuning logic to run on powerful server, but profiling to happen on target device

with figure

4 SimpleTVM

Using TVM follows the same workflow every time

To be able to quickly test different scenarios, we created a simpler interface for TVM, called SimpleTVM

4.1 Design

created wrapper for simpler usage of TVM

expose easy, chainable interface

Created automated benchmarking framework superb

enable automated testing of different configurations to be able to run multiple configurations without human intervention

Docker container to be able to easily deploy TVM with all dependencies on any server

4.2 Exploration

Using SimpleTVM and superb, it was easy to explore TVM behavior in different configurations of concurrency and resource sharing

First phase of experiments to investigate impact of interference

Evaluation in Jupyter notebooks

4.3 Autotuning Limitations

4.3.1 Resource Utilization

We noticed lots of resource idle time due to synchronous design

Show figure from poster

Want to minimize idle time because edge resources are limited (define edge)

Fundamental restriction due to dependencies of stages, cannot be changed for a single job

4.3.2 Scalability

Our goal is to enable large-scale autotuning for our AaaS, autotune multiple models at the same time

objectives:

Be able to run an arbitrary number of autotuning jobs while

1. maximizing inference performance: ultimate goal of autotuning
2. minimize hardware requirements: save cost
3. minimizing autotuning time: make autotuning worth the effort

in order of priority

State that autotuning time is not as crucial since it is rendered negligible by a large amount of inferences

With default tvm, there are two possible setups

Include figure with two setups

Include table with three experiments here

1. two completely separate autotuning jobs running independently on additional dedicated servers, one autotuning runner per server

Pros: good autotuning and inference time

Cons: Costly because we need multiple sets of the same hardware

not an economically feasible approach. We cannot simply use machines from a PaaS provider since actual target device needs to be used

Alternatively, we could use the same server and run them in sequence, trading off hardware required for autotuning time

2. two autotuning runners sharing the same server

Pros: only one set of hardware

Cons:

- interference drives up autotuning time

Autotuning takes long (in our tests anywhere between 3 and 36 hours, depending on hardware and network size)

Especially update model takes 64% longer when two jobs are running simultaneously

- results in worse inference performance because profiling is distorted (show numbers)

In both setups, we do not meet all objectives

Gets worse the more jobs we add

AaaS is not possible efficiently with current implementation and architecture of autotuning in TVM, does not scale well

Ideally:

Prevent interference, because it affects autotuning time and inference performance

Minimize hardware required by utilizing available hardware fully before adding new servers for cost reasons

However, there does not seem to be any solution yet

4.3.3 Similar Problems

In general, problem can be formulated as follows:

How can resources be shared optimally between multiple tasks that are partially idle?

Add two examples

5 Autotuning Scheduler

Enabling controlled parallel autotuning is necessary to solve those problems

necessitates central scheduler that orchestrates all jobs

5.1 Design

general idea: interleave stages while keeping dependencies and preventing interference

Make use of idle time

Want to keep both amount of hardware and inference time low, we want to run multiple jobs on the same hardware but make leverage idle time to interleave stages

include figure from poster

since only proof-of-concept, very specific to make it work quickly and non-flexible/fault-tolerant

Leverage SimpleTVM

5.1.1 Scheduling Algorithms

to keep scheduler algorithm simple, we designed it to be agnostic of stages

scheduler needs to know

- knows which job will use which resource
- knows which resource is currently available

we call this load-aware

theoretically, could work for any application that supports this interface (e.g. TC?)

allows for variable strategies to compare different designs

show scheduling pseudocode

5.1.2 Autotuning Decomposition

Necessary step before implementation

Show figure

Default TVM:

Procedure is monolithic

Start runner and it does not stop until its finished

We want to be able to control the execution of individual stages

Decompose monolith into separate units for stages

This allows us to control when which stage is being executed

Necessary for scheduler

Runner does not do anything on its own but waits for commands

5.2 Implementation

Figure with autotuning procedure with scheduler

Since TVM only provides a python interface, we are using python 3.5

5.2.1 RPC

We want clients to live in different docker containers, possibly physical servers (why?)

requires RPC infrastructure consisting of scheduler and clients

different from TVM RPC infrastructure

clients register to scheduler

describe endpoints

5.2.2 Components

Show whole stack, denote what happens in scheduler, what happens in runner

Show which communication is in-process and which is RPC

JobManager negotiates between autotuning stages interface and simple scheduler interface

show abstract scheduler and client interface

5.2.3 Challenges

evaluation of design choices takes long because autotuning is a slow process

initially wanted to run scheduler and clients in one multi-threaded process without RPC to get results quickly

not possible due to python global interpreter lock

what else?

5.3 Autotuning as a Service

imagine autotuning as a service where users can submit their trained model and receive an optimized version according to SLA

Describe as a service

More sophisticated scheduler, requires moving more autotuning logic from client to scheduler

Make client stateless

6 Evaluation

6.1 Results

Comparison of interleaved design vs synchronous and sequential in terms of autotuning time and inference time

hardware and network specifications

more and heterogenous jobs will enable better resource utilization and less wait time

6.2 Limitations

Evaluation only with limited set of hardware and models, general statement requires more experiments

Especially running jobs in parallel that vary significantly in complexity need to be examined

more limitations?

7 Conclusion

Describe results

Only used scheduler for TVM, but should work for TC as well because it also has stage dependencies

Enabled large-scale autotuning with only small sacrifices in autotuning time

7.1 Future Work

Predictive scheduler using times for task to make scheduling more intelligent

Requires more control in scheduler, not only simplified interface

running update model and build of one job directly after another will probably decrease waiting time, since that job can then already use the target device, so there is less target device idle time

Keep trained model and update it every n new entries to skip loading history time for every task

Check currently known best configurations and see if SLA is already met before actually starting autotuning

Automatically set up autotuning infrastructure

Split jobs on task and search space level to parallelize more

- make better use of unused resources
- faster autotuning, e.g. for paying customers

After best approach is found from prototype, make into mature product to enable real-time DL applications for everybody

