
Week 4 Reading Summaries: TVM + Triton Readings

Amogh Patankar

University of California, San Diego
3869 Miramar Street Box #3037, La Jolla, CA- 92092
apatankar@ucsd.edu

Abstract

1 The articles for this week’s reading summary describe two compilers custom-made
2 for deep learning and neural network computations- TVM, and Triton. TVM is a
3 compiler that allows for graph and operator level optimizations for deep learning
4 algorithms to improve performance. The result is a high performant system that is
5 competitive with hand-tuned libraries for low-power CPU, GPU, and server class
6 GPUs. Similarly, Triton is a language, and a compiler, built around the concept of
7 tiling. Triton is a C-based language for tensor expression on tiles, and tile-level
8 optimizations for execution of efficient GPU code.

9 1 TVM

10 1.1 Introduction

11 With a wider and wider range of deep learning models in the present, and with the advent of
12 generative AI, deep learning frameworks are important in order to optimize workload execution.
13 However, these frameworks cannot handle graph level optimizations due to the high-level nature
14 of those programs. Moreover, operator optimizations require manual tuning, which is not scalable.
15 TVM addresses key challenges, and solves them; namely, by 1) leveraging specific hardware features
16 and abstractions, and 2) Creating a large search space for optimization. Deep learning accelerators
17 require specialized hardware to efficiently run workloads, and systems in turn need to generate code
18 that controls pipeline dependencies. With respect to manual tuning, there *must* be a large enough
19 search space to come up with the correct combinatorial choices of memory access, threading pattern,
20 and other HW primitives.

21
22 This paper addresses these challenges with tensor expression language, an automatic program
23 optimization framework to find optimized tensor operators, and a graph rewriter to take advantage of
24 high/operator level optimizations.

25 1.2 Overview

26 TVM converts a model from an existing framework into a computational graph, then performs
27 high-level dataflow rewrites to produce an optimized graph. Fused operators in the graph are then
28 efficiently coded using a declarative tensor expression language, and still leaves execution details
29 unspecified. TVM uses a ML based cost model to identify the best operator in order to navigate the
30 vast space of possible optimizations for a given hardware target. Only then is the code packaged into
31 a deployable module for execution.

32 1.3 Optimizing Computational Graphs

33 Computational graphs in DL are representative of programs using multi-dim tensors, allowing for
34 high-level optimizations without operator implementations. TVM uses graphs for optimizations like
35 operator fusion, and more, to improve efficiency.

36 1.3.1 Operator Fusion

37 Operator fusion combines multiple operators into one kernel, hence not requiring intermediate results
38 to be stored in memory, significantly reducing execution time on GPUs and accelerators. Operators
39 can be injective (1-to-1 maps), reduction (sum), complex-out-fusable (conv2d with element-wise
40 output fusion), and opaque (non-fusable, sort). Injective operators fuse into another injective operator,
41 reduction operators fuse with input injective operators, and complex-out-fusable operators like conv2d
42 fuse with element-wise operations. So, computational graphs are transformed into an optimized,
43 fused version.

44 1.3.2 Data Layout Transformation

45 Data layout optimization is a method TVM uses to use efficient internal data layouts tailored to
46 hardware like column/row major, or tiling for accelerators by transforming a computational graph.
47 TVM specifies preferred layouts for operators and transforms data layouts between producers and
48 consumers when mismatches occur. High-level optimizations like operator fusion improve efficiency;
49 but, the effectiveness of said optimizations depend on operator library support. Customizing fused
50 kernels for various operators, data layouts + types, and hardware backends isn't sustainable, and so, a
51 code generation approach is taken. By generating diverse operator implementations and enabling
52 scalable and efficient support, the problem can be solved for workloads across multiple hardware
53 targets.

54 1.4 Generating Tensor Operations

55 TVM produces efficient code for operators by generating implementations, and choosing one. Here
56 are the TVM-specific features:

57 1.4.1 Tensor Expression and Schedule Space

58 TVM uses tensor expression, namely through a tensor expression language, to enable automatic
59 code generation. This process occurs by describing tensor operations through index formulas, with
60 output shapes and element-wise computations without defining loops or execution details. TVM
61 uses schedules to map expressions to code, and applies transformations for hardware optimization. It
62 extends existing primitives and introduces new ones for GPUs and accelerators that support CPUs,
63 GPUs, and TPU-like hardware. TVM tracks loop structures to generate efficient code, with automatic
64 schedule derivation discussed later.

65 1.4.2 Nested Parallelism with Cooperation

66 In order to optimize compute-intensive DL workloads, parallelism is very important especially on
67 GPUs with massive parallelism. Traditionally, the model is a fork-join model for nested parallelism
68 where tasks recursively get divided to exploit thread hierarchies, but threads can't share data within the
69 same stage. TVM improves this by supporting cooperative data fetching; threads can collaboratively
70 load data into shared memory, allowing reuse and optimizing GPU performance. Memory scopes are
71 introduced to mark compute stages as shared, ensuring that synchronization is done properly. The
72 approach benefits GPUs and aids in targeting accelerators by tagging memory buffers and creating
73 custom rules.

74 1.4.3 Tensorization

75 TVM introduces tensorization, which is the optimization of DL workloads by incorporating special-
76 ized tensor compute primitives, such as matmuls and 1D convolution. TVM tries to improve perfor-
77 mance and requires a flexible compilation framework to accommodate the various tensor instructions
78 across emerging accelerators. TVM creates an extensible approach by separating hardware-specific

instructions from the compilation schedule through a tensor-intrinsic declaration mechanism. This allows TVM to easily support new hardware architectures. TVM also creates a tensorized schedule primitive to map computations to hardware intrinsics, which can result in significant performance gains, such as a 1.5× speedup when using a low-precision micro-kernel for mobile CPUs.

1.4.4 Explicit Memory Latency Hiding

Latency hiding is the idea of overlapping memory ops with computation, in order to maximize memory/compute utilization. CPUs use multithreading and prefetching, while GPUs and specialized DL accelerators use DAE or threadwarps. Programming DAE accelerators requires synchronizing; so, TVM uses a virtual threading scheduling primitive. TVM starts with a high level multi-threaded program schedule, then inserts synchronization ops for correct execution. In hiding latency, the performance of algos improves for ResNet layers.

1.5 Automating Optimization

1.5.1 Schedule Space Specification

TVM utilizes schedule template specification to declare schedule options, incorporating domain-specific knowledge when needed. This unlocks automated exploration of a vast search space of configurations, and the schedule optimizer searches through billions of potential configurations, guided by the template and the tensor expression language.

1.5.2 ML-Based Cost Model

TVM utilizes a machine learning model to predict the performance of different schedule configurations as opposed to exhaustive auto-tuning. The model that is trained with real-time measurement data, helps reduce the search space. The focus on relative runtime predictions rather than absolute execution times helps to improve accuracy over time. TVM balances the need for quality predictions and fast execution, and uses various methods like gradient tree boosting (XGBoost) for efficient predictions.

1.5.3 Schedule Exploration

Using the cost model, the schedule explorer selects promising configurations, employing parallel simulated annealing algorithm to iteratively find better configurations. TVM can then search more efficiently by rejecting configurations with higher predicted costs and converging on lower-cost configurations. The exploration process is continuous and updates as new data from experiments becomes available.

1.5.4 Distributed Device Pool and RPC

TVM uses a distributed device pool and RPC framework to scale hardware trials, then managing optimization jobs across multiple devices. It automates compilation, execution, and profiling processes, especially useful for embedded devices, where tuning and cross-compilation are hard and inefficient.

1.6 Evaluation

TVM evaluates their framework, noting that it optimizes DL system SW-HW codesign as compared to other options.

2 Triton

2.1 Introduction

Deep Neural Networks (DNNs) and generative AI has now changed AI accelerators and HW, through parallel computing devices like GPUs, and vendor libraries like cuBLAS/cuDNN. These libraries are limited in supporting tensor operations, which has led to the development of DSLs, and pairing them with microkernels and high level tiling abstractions. Triton is a system comprising a C-like language for tensor programs, LLVM-based Intermediate Representation for tile-level optimization and a compiler for efficient GPU code generation.

123 2.2 Related Work

124 Triton has a few different approaches to optimizing deep learning computations. Existing frameworks
125 rely on hand-selected and optimized libraries like cuBLAS, some use domain-specific languages
126 (DSLs). Pre-existing Tensor-level IRs use predefined templates, the polyhedral model automates
127 compilation, and loop synthesizers enable manual schedule optimization. Triton integrates tile-level
128 operations, offers greater flexibility, and automatic schedule inference.

129 2.3 Triton-C

130 2.3.1 Syntax

131 Triton-C uses ANSI-C and CUDA-C, but has a few changes for semantics and programming model.
132 Primarily, special syntax for tile declarations, built-in functions (eg. dot, trans, etc.), and broadcasting
133 using the "newaxis" keyword. Some additional changes include the "" prefix for control flow.

134 2.3.2 Semantics

- 135 • Tile: Tile semantics simplifies tensor programs by hiding performance details regarding
136 memory, cache, and HW utilization.
- 137 • Broadcasting: Provides rules to perform padding and broadcasting conversions from scalar
138 to array implicitly.

139 2.3.3 Programming Model

140 Triton follows the SPMD programming model, just like many other support libraries and frameworks;
141 however, it makes each kernel single-threaded and parallelized. Triton assigns each kernel with a set
142 of global ranges that vary. This simplifies kernel execution by eliminating concurrency primitives.

143 2.4 Triton-IR

144 2.4.1 Structure

- 145 • Modules: Each unit of compilation is called a module, compiling independently and linked
146 by a linker.
- 147 • Functions: Consist of return type, name, and arguments list if required. Function and
148 parameter attributes can be specified.
- 149 • Basic Blocks: Straight line code sequences that may contain terminator instructions.

150 2.4.2 Support for Tile-Level Dataflow Analysis

- 151 • Types: Multi-dim types are declared using similar syntax to vectors in LLVM-IR.
- 152 • Instructions: These instructions are used to retile, and support broadcasting semantics.

153 2.4.3 Support for Tile-Level Control-Flow Analysis

154 Triton-IR doesn't have native support for divergent control flow in tiles; using Predicated SSA (PSSA)
155 form and ψ -functions, adding cmpp for dual predicates and psi for merging predicated instruction
156 streams.

157 2.5 Triton-JIT Compiler

158 2.5.1 Machine-Independent Passes

- 159 • Prefetching: Tile level memory ops are avoided by detecting loops and adding prefetching
160 code.
- 161 • Tile-level Peephole Optimization: Algebraic properties introduced as tiles are used as
162 operations.

163 **2.5.2 Machine-Dependent Passes**

- 164 • Hierarchical Tiling: TritonIR allows automatic enumeration and optimizatio of nested
165 tiling without polyhedral machinery.
- 166 • Memory Coalescing: TritonIR’s backend can order threads internally within microtiles.
- 167 • Shared Memory Allocation: Tile level operations store operands in fast shared memory, by
168 calculating the live range.
- 169 • Shared Memory Synchronization: R/W from and to shared memory are asynchronous; RAW
170 and WAR hazards are detected then barriers are inserted.

171 **2.5.3 Autotuner**

172 As compared to hand-written parameterized code templates, Triton-JIT extracts spaces from Triton-IR
173 by concatenating meta-parameters.

174 **2.6 Experiments**

175 In all experiments (MatMul, Convolutions, etc.), Triton is on-par with vendor librarise like cuBLAS,
176 cuDNN, etc, in some cases even more performant.

177 **2.7 Conclusions**

178 Triton’s language and compiler are effective to compiling tiled-NN computations into machine code.
179 Triton-IR also enables optimization passes, and Triton-C is a useful language to implement efficient
180 kernels.