15-442/15-642: Machine Learning Systems

Three Lessons Learned from ML Systems

Tianqi Chen and Zhihao Jia

Carnegie Mellon University

4/14/2025



Challenges of Building ML Systems



ML Systems

Increasingly diverse models

Large Language Models Transformers, Vision Language Models, Graph Neural Networks, Mixture of Experts, Sparse NN, Dynamic NN,

. . .

12 5GB/s 12 5 GB/s 16GB/s 16 GB/ DRAM sysment DRAM DRAM | s¥sMer CPU CPU X Bus 64GB/s HBM HBM HBM

Increasingly heterogeneous hardware

CPUs, GPUs, TPUs, Al accelerators, FPGAs, CGRAs, Programmable networks, and their combinations

. . .

CMU Automated Learning Systems Lab

Mission: Automate the design and optimization of ML systems by leveraging

- Statistical and mathematical properties of ML algorithms
- 2. Domain knowledge of modern hardware platforms





- 1. Automated approaches can offer 3-10x improvement for most tasks
- 2. Joint optimization is critical

catalyst

3. Combing systems and ML optimizations is promising but challenging

Lesson 1: Automated Approaches Offer 3-10x Improvement



Case Study 1: TASO: Up to 3.1x Case Study 2: FlexFlow: Up to 10x

Case Study: Current Rule-based Graph Optimizations

Case Study: Current Rule-based Graph Optimizations

Fuse conv + relu

TensorFlow currently
includes ~200 rulesFuse conv +
batch normalization

Fuse multi. convs

Lesson 1: Automated Approaches Offer 3-10x Improvement

(~<u>53,000</u> LOC)

namespace tensorflow { namespace graph_transforms { // Converts Conv2D or MatMul ops followed by column-wise Muls into equivalent // ops with the Mul baked into the convolution weights, to save computation // during inference. Status FoldBatchNorms(const GraphDef& input_graph_def, const TransformFuncContext& context, GraphDef* output_graph_def) { GraphDef replaced_graph_def; TF_RETURN_IF_ERROR(ReplaceMatchingOpTypes(input_graph_def, // clang-format off {"Mul", // mul node {"Conv2D|MatMul|DepthwiseConv2dNative", // conv_node {"*"}, // input_node {"Const"}, // weights_node {"Const"}, // mul_values_node // clang-format on }, // clang-format on
[](const NodeMatch& match, const std::set<string>& input_nodes, const std::set<string>& output_nodes, std::vector<NodeDef>* new_nodes) { // Find all the nodes we expect in the subgraph. const NodeDef& mul_node = match.node; const NodeDef& conv_node = match.inputs[0].node; const NodeDef& input_node = match.inputs[0].inputs[0].node; const NodeDef& weights_node = match.inputs[0].inputs[1].node; const NodeDef& mul_values_node = match.inputs[1].node; // Check that nodes that we use are not used somewhere else. for (const auto& node : {conv_node, weights_node, mul_values_node}) { if (output nodes.count(node.name())) // Return original nodes. new_nodes->insert(new_nodes->end(), {mul_node, conv_node, input_node, weights_node, mul_values_node}); return Status::OK(); Tensor weights = GetNodeTensorAttr(weights_node, "value"); Tensor mul_values = GetNodeTensorAttr(mul_values_node, "value"); // Make sure all the inputs really are vectors, with as many entries as // there are columns in the weights. int64 weights_cols; if (conv_node.op() == "Conv2D") { weights_cols = weights.shape().dim_size(3); } else if (conv_node.op() == "DepthwiseConv2dNative") { weights cols = weights.shape().dim_size(2) * weights.shape().dim_size(3); } else { 81 weights_cols = weights.shape().dim_size(1); 83 84 if ((mul_values.shape().dims() != 1) ||
 (mul_values.shape().dim_size(0) != weights_cols)) { 85 return errors::TnvalidArgument(87 "Mul constant input to batch norm has bad shape: ", mul_values.shape().DebugString()); 89 91 // Multiply the original weights by the scale vector.
auto weights_vector = weights.flat<float>(); Tensor scaled weights(DT FLOAT, weights.shape()); 93 auto scaled weights vector = scaled weights.flat<float>(); 95 for (int64 row = 0; row < weights_vector.dimension(0); ++row) {</pre> scaled_weights_vector(row) = weights_vector(row) * mul_values.flat<float>()(row % weights_cols); 99 100 101 // Construct the new nodes. NodeDef scaled_weights_node; scaled_weights_node.set_op("Const"); 104 scaled_weights_node.set_name(weights_node.name()); 105 SetNodeAttr("dtype", DT_FLOAT, &scaled_weights_node); 106 SetNodeTensorAttr<float>("value", scaled_weights, &scaled_weights_node); 107 new_nodes->push_back(scaled_weights_node); 108 109 new_nodes->push_back(input_node); 110 NodeDef new_conv_node; new_conv_node = conv_node; new_conv_node.set_name(mul_node.name()); new_nodes->push_back(new_conv_node); return Status::OK(); {{}, &replaced_graph_def)); *output_graph_def = replaced_graph_def; return Status::OK(): REGISTER_GRAPH_TRANSFORM("fold_batch_norms", FoldBatchNorms);

// namespace graph_transforms
// namespace tensorflow

8

TASO: Tensor Algebra SuperOptimizer

Key idea: replace manually-designed graph optimizations with *automated generation and verification* of graph substitutions for tensor algebra

- Less engineering effort: <u>53,000</u> LOC for manual graph optimizations in TensorFlow → <u>1,400</u> LOC in TASO
- Better performance: outperform existing optimizers by up to 3x
- Stronger correctness: formally verify all generated substitutions

End-to-end Inference Performance (Nvidia V100 GPU)

Lesson 1: Automated Approaches Offer 3-10x Improvement

Case Study 1: TASO: Up to 3.1x Case Study 2: FlexFlow: Up to 10x

FlexFlow: Automatically Optimizing DNN Parallelization

Better Performance

Up to 10x faster than manually designed strategies

Fast Deployment

Minutes of automated search to discover performant strategies

No Manual Effort

Automatically find strategies for new DNN models or hardware platforms

FlexFlow: Searching for Efficient Parallelization Strategies

A **search space** of possible parallelization strategies

A cost model and a search algorithm

Fast and Scalable Parallelization strategies

- Samples
- Operators
- Attributes
- Parameters

- Samples: partitioning training samples (Data Parallelism)
- Operators
- Attributes
- Parameters

Parallelizing a 1D convolution in Sample

- Samples: partitioning training samples (Data Parallelism)
- Operators: partitioning ML operators (Model Parallelism)
- Attributes
- Parameters

Parallelizing multiple convolutions in *Operator*

- Samples: partitioning training samples (Data Parallelism)
- Operators: partitioning ML operators (Model Parallelism)
- Attributes: partitioning attributes in a sample (e.g., pixels)
- Parameters

Parallelizing a 1D convolution in Attribute

- Samples: partitioning training samples (Data Parallelism)
- Operators: partitioning ML operators (Model Parallelism)
- Attributes: partitioning attributes in a sample (e.g., pixels)
- Parameters

- Samples: partitioning training samples (Data Parallelism)
- Operators: partitioning ML operators (Model Parallelism)
- Attributes: partitioning attributes in a sample (e.g., pixels)
- Parameters: partitioning parameters in an operator (Tensor Model Parallelism)

Parallelizing a 1D convolution in *Parameter*

Hybrid Parallelism in SOAP

Example parallelization strategies for 1D convolution

Different strategies perform the same computation.

Data parallelism

A parallelization strategy in SOAP (1.2x faster)

Data parallelism

A parallelization strategy in SOAP (1.2x faster)

Challenges of Discovering Fast Strategies in SOAP

1. SOAP contains billions or more possible strategies

MCMC search algorithm

2. Evaluating a strategy on hardware is too slow

Execution simulator

FlexFlow Overview

Deep Learning Recommendation Model (DLRM) facebook A deep learning model for ads recommendation

Lesson 1: Automated Approaches Offer 3-10x Improvement

PET: Up to 2.5x

Unity: Up to 3.6x

Lux: Up to 10x Roc: Up to 4.1x

Common Advantages of Automated Approaches

- Better runtime performance: discovering novel optimizations hard to manually designed, 3-10x speedup over manual optimizations
- Less engineering effort: code for discovering optimizations is generally much less than manual implementation of these optimizations
- Stronger correctness guarantees: using formal verification techniques

Lesson 2: Joint Optimization is Critical to Performance

1. Unity: Accelerating DNN Training Through Joint Optimization of Algebraic Transformations and Parallelization. OSDI'22.

- 2. TopoOpt: Optimizing the Network Topology for Distributed DNN Training. NSDI'23.
- 3. MetaFlow: Optimizing DNN Computation with Relaxed Graph Substitutions. MLSys'19

Revisit Tensor Model Parallelism

- y = X X W output input parameters
- Partition parameters/gradients within a layer

Tensor Model Parallelism (partition output)

1. Representation

2. Scalability

Representation -Parallel Computation Graph (PCG) Unity Scalability – Hierarchical Search Algorithm

Joint Optimization Enables Better Performance and Scalability

Lesson 2: Joint Optimization is Critical to Performance

1. Unity: Accelerating DNN Training Through Joint Optimization of Algebraic Transformations and Parallelization. OSDI'22.

2. TopoOpt: Optimizing the Network Topology for Distributed DNN Training. NSDI'23.

3. MetaFlow: Optimizing DNN Computation with Relaxed Graph Substitutions. MLSys'19

Lesson 3: Combining ML and Systems Optimizations is Promising but Challenging

Systems Optimizations

- Graph Transformations
- Auto Parallelization
- Kernel Generation
- Data Layout and Placement

- Quantization
- Low-Rank Adaptation
- Distillation
- Neural Architecture Search

Lesson 3: Combining ML and Systems Optimizations is Promising but Challenging

Systems Optimizations

Pro: preserve equivalence

- Pro: better runtime performance
 - Faster ML operators
 - Less Computation

Achieve the best of both worlds?

Lossless Optimizations

Approximate Optimizations — Predictive Performance

Lesson 3: ML and Systems Optimizations is Promising but Challenging

1. Automated approaches can offer 3-10x improvement on most tasks

2. Joint optimization is critical

3. Combing systems and ML optimizations is promising but challenging