15-442/15-642: Machine Learning Systems

Machine Learning Compilation

Spring 2025

Tianqi Chen Carnegie Mellon University

2/10/2025



Overview of Machine Learning Compilation

Example ML Compilation Flow



Overview of Machine Learning Compilation

Example ML Compilation Flow

ML System Optimization Problem



- Specialized libraries for each backend (labor intensive)
- Non-automatic optimizations



Machine Learning Compilation



High-level IR Optimizations and Transformations

Tensor Operator Level Optimization

Direct code generation









Machine Learning Compilation

Development Form

Deployment Form



An example instance of deployment form

Key Elements in Machine Learning Compilation



input: Tensor[(1, 3072)]



Tensor multi-dimensional array that stores the input, output and intermediate results of model executions.

Tensor Functions that encodes computations among the input/output. Note that a tensor function can contain multiple operations

ML Compilation Goals

There are many equivalent ways to run the same model execution. The common theme of MLC is optimization in different forms:

Minimize memory usage.

Improve execution efficiency.

Scaling to multiple heterogeneous nodes.

Example Compilation Process

Development

Deployment



In this particular example, two tensor functions are folded into one (linear-relu). With a specialized implementation (in reality, they will be implemented using low-level primitives).

Abstraction and Implementation

Abstraction refers to different ways to represent the same system interface (tensor function)



def linear_relu(x, w, out):
 for i in range(1):
 for j in range(200):
 out[i, j] = 0
 for k in range(3072):
 out[i, j] += x[i, k] * w[j, k]
 out[i, j] = max(out[i, j], 0)

Three abstraction ways to represent the same tensor function (linear_relu), each providing a different level of details. In practice, we usually say that the more specialized version is an **implementation** of higher-level abstraction.

MLC as Tensor Function Transformation (with different abstractions)

Development

Deployment



Most MLC process can be viewed as transformation among tensor functions (that can be represented with different abstractions).



Overview of Machine Learning Compilation

Example ML Compilation Flow

Compiler Representation of a ML Model



IRModule: a collection if interdependent functions

Example Compilation Flow: High-Level Transformations







High-level transformations

Example Compilation Flow: Lowering to Loop IR





Example Compilation Flow: Low Level Transformations





Low-level transformations

Example Compilation Flow: CodeGen and Execution



Runtime Execution

Discussion

- What are possible ways to represent a function in ML
- The possible set of optimizations we can perform in each type of representations.

High-level IR and Optimizations



- Computation graph(or graph-like) representation
- Each node is a tensor operator(e.g. convolution)
- Can be transformed (e.g. fusion) and annotated (e.g. device placement)
- Most ML frameworks have this layer

Low-level Code Optimization



Elements of Low-level Loop Representation



Transforming Loops: Loop Splitting



Transformation

x = get_loop("x")
xo, xi = split(x, 4)

Transforming Loops: Loop Reorder

Code for xo in range(32): for xi in range(4): C[xo * 4 + xi]= A[xo * 4 + xi] + B[xo * 4 + xi]for xi in range(4): for xo in range(32): C[xo * 4 + xi]= A[xo * 4 + xi] + B[xo * 4 + xi] Transformation

x = get_loop("x")
xo, xi = split(x, 4)
reorder(xi, xo)

Transforming Loops: Thread Binding

for xi in range(4): for xo in range(32): C[xo * 4 + xi]= A[xo * 4 + xi] + B[xo * 4 + xi]def gpu kernel(): C[threadId.x * 4 + blockIdx.x] = . . .

Code

Transformation

x = get_loop("x") xo, xi = split(x, 4) reorder(xi, xo) bind_thread(xo, "threadIdx.x") bind_thread(xi, "blockIdx.x")

Search via Learned Cost Model

