10-414/714 – Deep Learning Systems: Algorithms and Implementation

Model Deployment

Fall 2024 J. Zico Kolter and Tianqi Chen (this time) Carnegie Mellon University

Outline

Model deployment overview

Machine learning compilation

Outline

Model deployment overview

Machine learning compilation

What we have learned so far in this class

How to build a deep learning system that trains deep learning models efficiently on a standard computing environment (with GPUs).

Automatic differentiation

Deep learning modeling techniques

Hardware accelerations and scale up

Normalization, initialization, optimization

Model deployment

Training



Bring learned models to different application environments

Model deployment considerations

Application environment may bring restrictions (model size, no-python)

Leverage local hardware acceleration (mobile GPUs, accelerated CPU instructions, NPUs)

Integration with the applications (data preprocessing, post processing)

Model exportation and deploy to inference engines



7

Inference engine internals



Computational graph

Many inference engines are structured as computational graph interpreters

Allocate memories for intermediate activations Traverse the graph and execute each of the operators

Usually only support a limited set of operators and programming models (e.g. dynamism)

Outline

Model deployment overview

Machine learning compilation

Limitation of library driven inference engine deployments

Need to build specialized libraries for each hardware backend

A lot of engineering efforts to optimization

Machine learning compilation



High-level IR Optimizations and Transformations

Tensor Operator Level Optimization

Direct code generation









Compiler representation of a model



Also called intermediate representation (IR)

IRModule: a collection of 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: code generation and execution



Runtime Execution

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 optimizations



Elements of low-level loop representations



Transforming loops: splitting



Transformation

for x in range(128):
 C[x] = A[x] + B[x]

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

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

Transforming loops: reorder

Code

Transformation

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]

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

Transforming loops: thread binding

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



Summary: elements of an automated ML compiler

Program representation

• Represent the program/optimization of interest, (e.g. dense tensor linear algebra, data structures)

Build search space through a set of transformations

- Cover common optimizations
- Find ways for domain experts to provide input

Effective search

Still an open research area!

• Cost models, transferability

Outline

Deploying models to different backends

Machine learning compilation