

MASE: Abstractions, Optimizations and Implementations

Lecture 2 for Advanced Deep Learning Systems

Aaron Zhao, Imperial College London, a.zhao@imperial.ac.uk

Table of contents

1. Introduction
2. MaseGraph IR
3. Metadata and Analysis Passes
4. Transform Passes

Introduction

MASE (Machine Learning Accelerator System Exploration) is an open-source project that aims to automate the exploration of ML system software and hardware.

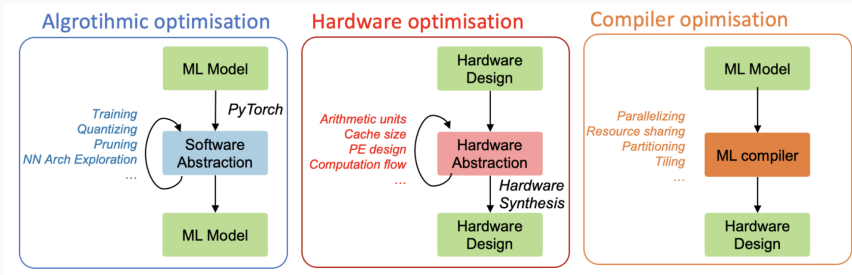
Introduction - Why re-inventing the wheel?

**Consider this image before saying "don't
reinvent the wheel"**



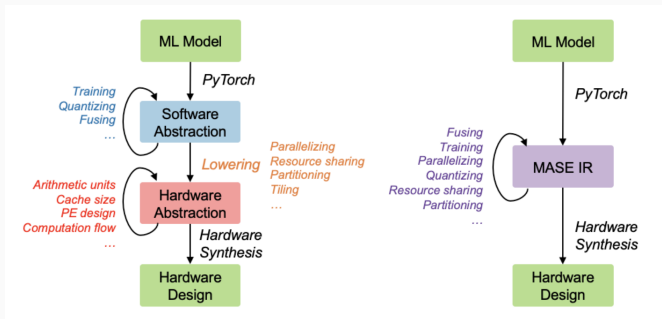
Introduction - Why re-inventing the wheel?

- Pytorch, Tensorflow (high-level python tools), algorithmic exploration, mapping mostly to CPUs and GPUs
- MLIR, TVM, compiler tools, map pre-defined network to various hardware targets
- MLIR-Circt, scheduling based HLS on top of MLIR



Introduction - Why re-inventing the wheel?

We are interested in combine these in a unified abstraction – a new graph-based **MASE Intermediate Representation (IR)**



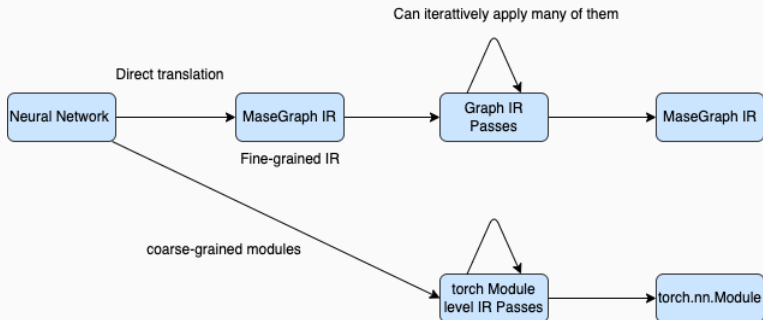
Introduction - Passes

A **pass** (can be either a transformation or an analysis), takes in the IR of the model, and returns the IR again

```
1      # pass takes a MaseGraph (the IR) and corresponding  
      ↪ pass-related arguments  
2      # graph: MaseGraph, pass_args: dict  
3      def pass_name(graph, pass_args):  
4          ...  
5      # return a MaseGraph (the IR) again, and a  
      ↪ dictionary for additional data  
6      # graph: MaseGraph, info: dict  
7      return graph, info  
8
```


Introduction - Overview

- Graph-level IR system and passes (we will cover in labs)
- Module-level passes (code is there)
- The idea of summarizing workloads into a set of IRs, and apply passes on them is the same as traditional compiler systems (eg. LLVM).

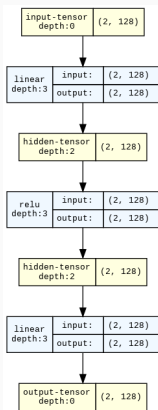


MaseGraph IR

Core idea: we represent Neural Networks as a **computation graph**, where nodes are computation blocks and **edges are data**.

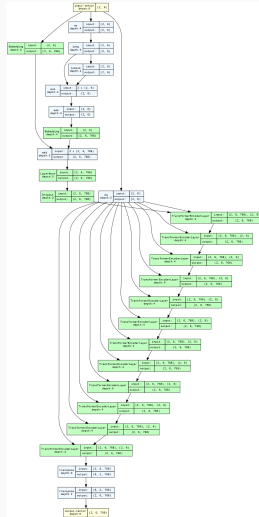
MASEGraph IR - visualization

Core idea: we represent Neural Networks as a **computation DAG** (Directed Acyclic Graph), where nodes are computation blocks and **edges** are data.



MASEGraph IR - visualization

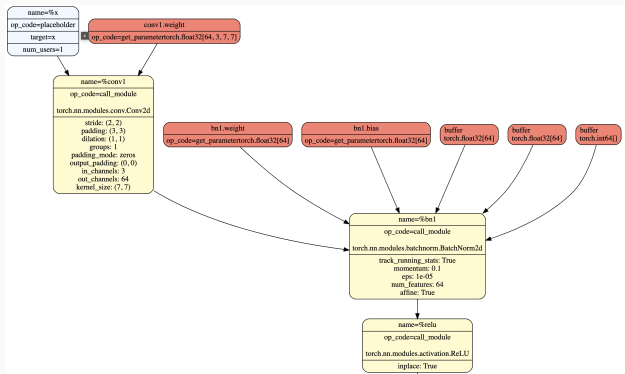
This can be very complex, notice the transformer layer is represented at a coarse granularity.



MASEGraph IR - implementation

Mase takes a **torch.fx graph representation** of a model and translates it into a customised representation (Mase graph IR).

The MaseGraph IR is a lot more complex than the previous visualization. Below is a single convolution layer. To reproduce, run *python machop/test/others/plot_graph.py*.



MASEGraph IR - implementation

The manipulation of the model requires both access to the torch fx graph and modules, you will understand this better after having the labs.

The definition can be found at

mase/machop/chop/ir/graph/mase_graph.py

```
1 class MaseGraph:
2     def __init__(self, model, cf_args) -> None:
3         ...
4
5     def fx_graph(self):
6         return self.model.graph
7
8     def modules(self):
9         return dict(self.model.named_modules())
10
```

MASEGraph IR - types

IR types are for nodes in the MaseGraph

- **placeholder**: for inputs
- **module**: for pytorch nn.Module
- **module_related_func**: some functions have the same functionality as a module, for instance, torch.nn.Conv2d (Module) and torch.nn.functional.conv2d are the same.
- **builtin_func**: what fx considers as builtin_funcs
- **implicit_func**: all other funcs that are not builtin
- **get_attr**: normally used for retrieving a parameter
- **output**: for outputs

A complete definitions of these and also supported nodes are in `mase/machop/chop/passes/graph/common.py`

Metadata and Analysis Passes

MASEMetadata - why?

IR only carries **type information** and **node relations**.

We normally need more information to perform complex operations, such information is called **metadata** and they are added to each node.

```
1 class MaseMetadata:  
2     ...  
3
```

How do we add such a class to the MaseGraph IR?

Init Metadata Pass

Implemented as a pass!

Traverse each node and append the MaseMetadata object to each node.

mase-tools/machop/chop/passes/graph/analysis/init_metadata.py

```
1 def init_metadata_analysis_pass(graph, pass_args={}):  
2     for node in graph.fx_graph.nodes:  
3         node.meta["mase"] = MaseMetadata(node=node,  
           ↪ model=graph.model)  
4     return graph, {}  
5
```

Analysis Passes

Optimizations or information gathering are implemented as **Passes** that traverse the whole or some portion of a network to either collect information or transform the network.

Generally, analysis passes are used for **collect extra information** of the network for later transformation passes.

- `add_common_metadata`
- `add_software_metadata`
- `add_hardware_metadata`

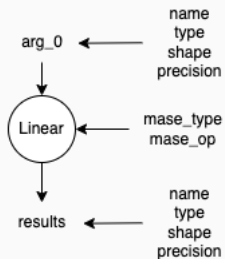
Passes are summarised at

mase-tools/machop/chop/passes/graph/__init__.py

Add a bunch of commonly used metadata to each node. This includes

- mase_type: (module_related_func, implicit_func ...)
- mase_op: (linear, relu ...)
- args: (name, type, shape and precision for all input arguments)
- results: (name, type, shape and precision for all all results)

add_common_metadata Passes



add_software_metadata and add_hardware_metadata Passes

Similarly, `add_software_metadata` and `add_hardware_metadata` passes do the same thing but add metadata for software and hardware respectively.

```
1  "common": {
2    "mase_type": "module_related_func",
3    "mase_op": "linear",
4    "args": {
5      "data_in_0": {"shape": [1, 784], "type": "float",
6        ↪ "precision": [32]},},
7      "weight": {"type": "float", "precision": [32], "shape":
8        ↪ [784, 784]},},
9      "bias": {"type": "float", "precision": [32], "shape":
10        ↪ [784]}},},
11    "results": {"data_out_0": {"type": "float", "precision":
12      ↪ [32], "shape": [1, 784], "torch_dtype":
13      ↪ torch.float32,}, },},
14    "software": {}, "hardware": {},
```

Transform Passes

Transform passes take a MaseGraph (or a network) as an input and perform certain modifications to it as an input and perform certain modifications to it.

I will use the quantization transform pass as an example.

chop/passes/graph/transforms/quantize/quantize.py

Quantize Transform Passes

```
1 def quantize_transform_pass(graph, pass_args=None):
2     by = pass_args.pop("by")
3     match by:
4         case "type":
5             graph = graph_iterator_quantize_by_type(graph,
6                 ↪ pass_args)
7         case "name":
8             graph = graph_iterator_quantize_by_name(graph,
9                 ↪ pass_args)
10        case "regex_name":
11            graph = graph_iterator_quantize_by_regex_name(graph,
12                ↪ pass_args)
13        case _:
14            raise ValueError(f'Unsupported quantize "by": {by}')
15
16 return graph, {}
```

- Transformation is also implemented as a traverse to the MaseGraph.
- You can use pass arguments to control your logic.

Quantize Transform Passes: replacing a node

```
1 def graph_iterator_quantize_by_name(graph, config):
2     ...
3     for node in graph.fx_graph.nodes:
4         ...
5         ori_module = get_node_actual_target(node)
6         # create the new quantized module
7         new_module = create_new_module(...)
8         # take the parent node based on the graph hierarchy
9         parent_name, name = get_parent_name(node.target)
10        setattr(graph.modules[parent_name], name, new_module)
11        # update meta data accordingly
12        update_quant_meta_param(node, node_config,
13                                ↪ get_mase_op(node))
13
```

What's next?

- You will go through the quantization pass and learn MASE in labs.
- Lecture 3 - 8 will cover more on recent progress in the field of ML systems.