# Compiler for DNN Workloads on AMD GPUs
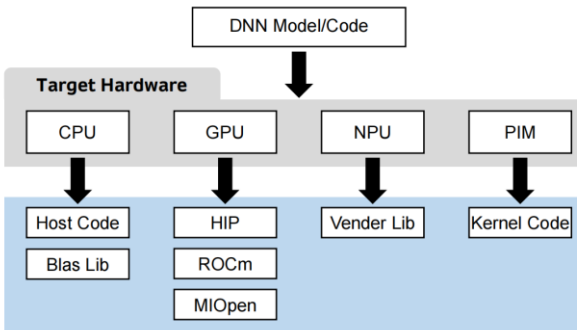
**Jumin Lee**
**Yonsei University, Korea**

OpenHW2023

AMD

*On board test by AMD GPUs*

## INTRODUCTION

The motivation for a DNN Compiler lies in the **diverse architecture environments**, where each hardware platform utilizes unique APIs. This diversity requires the creation of specific code generating paths for different hardware configurations. The challenge arises from the need to adapt and optimize code generation paths to **meet the specific requirements of each hardware environment**. Therefore, the development of a DNN Compiler addresses this challenge by providing a flexible and efficient solution that can dynamically adjust to the varied architectures encountered in different hardware platforms.
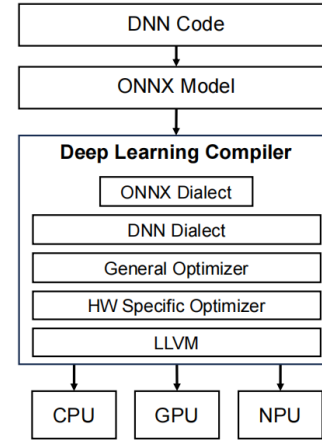


**Motivation of DNN Compiler**



**ONNX MLIR based system**

The ONNX graph-to-code converter serves as a publicly available compiler facilitating the transformation of ONNX models to the MLIR framework. **This converter** operates by providing an intermediate representation of DNN operations, including the definition of DNN-related operations and an abstract API for each hardware. Leveraging LLVM calls for **backend APIs**, the compiler ensures a 1-to-1 mapping of DNN dialect to AMD GPU code. Furthermore, it handles translation for hardware-specific libraries and code, while optimizations are conducted at the DNN dialect level. This integrated approach **enables efficient and hardware-adaptive compilation of ONNX models** into the MLIR framework.

## CREATIVE DESING

## RESULT

### *Correctness verification*

– Evaluated with **1000 iterations** each with random value input

– Convolution with maxpooling and ReLU is found to be correct

– **Extension with ROCblas** is expected in future work



### *Future Work*

• Extension on other benchmarks

– Current implementation is limited on fixed sized CNNs

– Extension on transformer or LLMs is expected

• Optimizations on Multi-tenant benchmarks

– Scheduling in multi tenant workloads are needed

– Memory optimization for scaled situations are considered

• Expansion of backends

– Not only GPU, but NPU and cluster target backends