

## Presentation Abstract

# A Compiler Overview for A Deep Neural Network Accelerator with Software-controlled Memory Hierarchies

KAZUAKI ISHIZAKI<sup>1,a)</sup> ERI OGAWA<sup>1</sup> HIROSHI INOUE<sup>1</sup> SWAGATH VENKATARAMANI<sup>2</sup>  
JINTAO ZHANG<sup>2</sup> WEI WANG<sup>2</sup> VIJAYALAKSHMI SRINIVASAN<sup>2</sup>  
MORIYOSHI OHARA<sup>1</sup> KAILASH GOPALAKRISHNAN<sup>2</sup>

Presented: March 13, 2020

This presentation gives the design and implementation of a compiler for a deep neural network accelerator that provides high performance and power efficiency. The accelerator consists of multiple heterogeneous units and each unit has a limited set of instructions. Also it does not have hardware-controlled caches for power efficiency. It is very hard for programmers to directly write a program for such an accelerator. Thus, we develop a compiler that automatically generates native code for each unit from a program in deep learning frameworks such as TensorFlow. This compiler generates outer loops around highly-tuned hand-crafted inner kernel loops for each unit with a wide range of neural network parameters. It effectively supports various types of parameters such as zero padding sizes around an image. This compiler also generates data transfer code between hierarchical memories, and applies code optimizations to reduce the code size.

---

This is the abstract of an unrefereed presentation, and it should not preclude subsequent publication.

<sup>1</sup> IBM Research - Tokyo, Chuo, Tokyo 103-8501, Japan

<sup>2</sup> IBM Research - Watson Research Center, Yorktown Heights, NY 10598

<sup>a)</sup> ishizaki@jp.ibm.com