# **Evaluation of Versal AI Engines** for Simple Neural Network Inference

### Yilin Shen, Caroline Johnson, Scott Hauck (University of Washington)

## Introduction

- Al deployment requires hardware assistance (CPU/GPU/FPGA, etc.)
- FPGA is a popular solution for its low latency
- Next-gen FPGA (AMD Versal) introduced AI Engine(AIE)



The AIE looks promising:



Should we use it? Let's compare AIE to FPGA!

# **Metrics for Evaluation**

- Initiation interval (1/throughput)
- Latency
- Power
- Price
- **Resource utilization**
- Silicon Area utilization



The smaller, the better!



\* FPGA Implementation: see [1]









# Conclusion

# **Future Work**

- Map more complex models Convolution with stride
  - Multi-channel convolution
- Evaluate AIE-ML: next-gen AIE

### Reference

2023.

• If throughput and latency is the priority, go for FPGA • If cost or replication is the priority, go for FPGA • AIE saves silicon area and power for basic ML workload

Optimized interconnection and enhanced MAC unit

