

Appendix: Artifact Description/Artifact Evaluation

Artifact Description (AD)

I. OVERVIEW OF CONTRIBUTIONS AND ARTIFACTS

A. Paper's Main Contributions

This paper proposes M³XU, which is a multi-mode matrix processing unit that supports IEEE 754 single-precision and complex 32-bit floating-point computations. M³XU consists of (1) additions of logic to feed different parts of inputs, (2) extensions to the arithmetic units to support exact FP32 precisions, and (3) slight extensions to accumulators to accumulate numbers in double-precision formats. M3XU enables MXUs to handle standard FP32 floating-point numbers and FP32C complex numbers adeptly, achieving their theoretical throughput under current memory technologies and with relatively minor area overhead. The following is a list of the main contributions of the paper.

- C_1 HDL implementation of M³XU.
- C_2 Performance emulation framework for M³XU.
- C_3 Applications using M³XU.

B. Computational Artifacts

- A_1 <https://doi.org/YY.YYYY/zenodo.0XXXXX>

| Artifact ID | Contributions Supported | Related Paper Elements |
|-------------|-------------------------|-------------------------|
| A_1 | C_1, C_2, C_3 | Table 2-4 Figure 4-9 |

II. ARTIFACT IDENTIFICATION

A. Computational Artifact A_1

Relation To Contributions

Relationship between the artifact and contributions are as follows:

- HDL for hardware synthesise result.
- Setting up performance emulation framwork for M³XU.
- Performance emulation of case study application using M³XU.

Expected Results

The expected outcome of the artifact is as follows:

- Area energy consumption of M³XU. M³XU is expected to show a 41% area overhead and 31% less power consumption compared to the FP16 baseline.
- Raw perfomance data including throughput and latency of M³XU.
- End to end latency of application M³XU.

Expected Reproduction Time (in Minutes)

- 5 hrs for HDL synthesization.
- 120 hrs for M³XU microbenchmark on all input sizes.
- 24 hrs for case study applications.

Artifact Setup (incl. Inputs)

Hardware: Nvidia A100 DGX Station. Or equivalent system installed with Nvidia A100 PICE GPUs.

Software:

- Linux kernel version 5.4.0-81-generic
- CUDA 11.4, driver 470.57.0
- Synopsis design compiler with the 45nm FreePDK45 library
- Nvidia cutlass

Datasets: Input data are synthetic by framework or provided.

Installation and Deploymen: Configuration and build commands are provided.

Artifact Execution

- follow ReadME.md for cutlass installation and M³XU integration.
- Configure licenses for Synopsis design compiler.
- Makefiles are provided for each case study applications.