RipTide: A Programmable, Energy-Minimal Dataflow Compiler and Architecture

Graham Gobieski, **Souradip Ghosh**, Marjin Heule, Todd C. Mowry, Tony Nowatzki*, Nathan Beckmann, Brandon Lucia

Carnegie Mellon University, *UCLA

SRC Task (3019.001)

Outline

- 1) Motivation
- 2) Background
 - a) Coarse-grained reconfigurable arrays (CGRAs) b) Prior work

3) **RipTide**

a) Overview → Compiler and architecture
b) Memory ordering
c) Mapping
d) Results

Motivation: Sensors at the extreme edge

Smart sensor devices at the *extreme edge* are quickly emerging. Devices need to run **complex apps** → inference, DSP, and more.



Trillions of devices¹ coming

Motivation: Computing at the extreme edge?

Compute must be *efficient* and *sustainable* at the extreme edge.



²Gobieski et al., "Intelligence Beyond the Edge: Inference on Intermittent Embedded Systems." (ASPLOS '19).

Background: Comparing existing architectures

Devices must have architectures that are *flexible* and *efficient*.



CGRAs: What is a CGRA?



Coarse-grained configurable arrays: Grid of processing elements (**PE**, to execute ops) connected by a **NoC** (to send values).

Eliminates fetch/decode & reg. file usage!

CGRAs: An end-to-end pipeline



1. Extract a *dataflow graph* (DFG) from code, optimize

2. *Map* ops to a PE mix and links on the NoC

3. Execute ops w/ "dataflow firing" or a static schedule

Prior Work: Prior ULP CGRAs are limited

Runs only affine inner loops.

No complex control-flow, irregular memory accesses, or operation ordering.

⁴Gobieski et al., "SNAFU: An Ultra-Low-Power, Energy-Minimal CGRA-Generation Framework and Architecture." (ISCA '21).

Insight: To improve efficiency, CGRAs need to run entire apps & support common PL idioms

RipTide: Overview

C code (lightly annotated)



Handles arbitrary code via
1) Complex control-flow
2) Irregular mem. access
3) Operation ordering



Generated CGRA



run control-flow. Frees PEs for ops.

RipTide: Memory ordering









a) Build an **ordering graph** for mem. deps b) **Prune** alreadyenforced arcs via existing data and control deps c) Perform a **pathsensitive** transitive reduction

d) Remaining ordering arcs are **enforced**.

RipTide: Mapping



RipTide: Results



10 apps from linalg, graph processing, and signal processing

+ Full compiler built with LLVM

+ Full RTL design and synthesized

+ Ran an entire DNN on RipTide!

Thank you!

Correspondence:

- Souradip Ghosh (<u>souradip@cmu.edu</u>)
- Graham Gobieski (gobieski@cmu.edu)

For more details, visit:

- Abstract Group (<u>abstract.ece.cmu.edu</u>)
- CORGi Group (<u>cmu-corgi.github.io</u>)