#### Refinements in Data Manipulation Method for Coarse Grained Reconfigurable Architectures

#### <u>Takuya Kojima</u> and Hideharu Amano Keio University, Japan



14th International Symposium on Reconfigurable Communication-centric Systems-on-Chip (ReCoSoC 2019)

### Importance of Programmability and High Energy Efficiency

- Forthcoming
  - ■IoT devices
  - Wearable computers
  - Edge computing



- Challenges for these devices
   Programmability
   To satisfy various demands
  - High energy efficiencyTo extends long battery life



### CGRAs: Coarse-Grained Reconfigurable Architectures



#### CGRAs

- Support word-level reconfiguration ( $\leftrightarrow$  bit-level of FPGAs)
- Have many PEs (Processing Element) in 2D grid
- Change functionality for each ALU & interconnection between PEs dynamically or statically

### Power-hungy Dynamic Reconfiguration



Dynamic Reconfiguration
 Changes configuration
 <u>cycle-by-cycle</u>

Provides more flexibility

Causes large dynamic power consumption

Details of power consumption for a dynamic reconfiguration CGRA[1]

[1] Ozaki, Nobuaki, et al. "Cool mega-arrays: Ultralow-power reconfigurable accelerator chips." IEEE Micro 31.6 (2011): 6-18.

#### SF-CGRAs: Straight-Forward CGRAs



# VPCMA: Variable Pipelined Cool Mega Array [2]



[2] N.Ando, et al. "Variable pipeline structure for Coarse Grained Reconfigurable Array CMA." *Field-Programmable Technology*, 2016.

# Computation on the PE array



- Fetch registers are connected to input of the PE array
- Gather registers are connected to output of the PE array
- The micro-controller
  - Writes data to the fetch registers
  - Read result from the gather registers

Fetch Registers Gather Registers

# Computation on the PE array



- Fetch registers are connected to input of the PE array
- Gather registers are connected to output of the PE array
- The micro-controller
  - Writes data to the fetch registers
  - Read result from the gather registers

Fetch Registers Gather Registers

# Variable Pipeline Structure



- No registers in each pipeline stage
  - $\rightarrow$  Pure combinational circuit
- Clock tree only for activated pipeline registers
- Variable pipeline structure depending on application

# Multi-cycle Execution on PE Array



#### Micro-controller

- A custom tiny RISC processor controls the processing
- "Fetch" op kicks off the execution
- "Gather" op writes back the results \_
- "Delay" op specifies delay time of "Gather" execution
- "Branch" op makes a loop

Fused into

an instruction

Cycle

# Multi-cycle Execution on PE Array



- A custom tiny RISC processor controls the processing
- "Fetch" op kicks off the execution
- "Gather" op writes back the results \_
- "Delay" op specifies delay time of "Gather" execution
- "Branch" op makes a loop

**Fused** into

an instruction

# Data Manipulator of VPCMA



Data manipulator

- Placed between Dmem & PE array
- Transfers any input data to any outputs
- Loads at most consecutive 12 data from 12 mem banks
- Increments addr. automatically for next fetch

# Data Manipulator of VPCMA



Data manipulator

- Placed between Dmem & PE array
- Transfers any input data to any outputs
- Loads at most consecutive 12 data from 12 mem banks
- Increments addr. automatically for next fetch

#### Ultra Low Power Consumption of CMA

- No-Pipelined version of CMA[6]
  - Works with Lemon battery
  - Achieves 743 MOPS/mW (297MOPS/0.4mW)

#### VPCMA



- ■Keeps the same energy efficiency
- Achieves 4x higher peek performance
- Problem

#### Less flexibility because of saving too much energy

[6] M.Koichiro, *et al.* "A 297mops/0.4 mw ultra low power coarse-grained reconfigurable accelerator CMA-SOTB-2." *2015 International Conference on ReConFigurable Computing and FPGAs (ReConFig)* <sup>14</sup>

### Limitation of data handling in VPCMA



Memory allocation in bank memory

Data manipulator cannot access multiple data more than 12 step distance simultaneously → needs data rearrangement

 $\rightarrow$  often incurs extra copy of data

### Limitation of data handling in VPCMA



 ■ Data manipulator cannot access multiple data more than 12 distance simultaneously
 → needs data rearrangement
 → often incurs extra copy of data

# Other limitations of VPCMA

#### Also, VPCMA

- 1. Suffers from a lack of constant registers for the PE array
  - A PE row (12 PEs) share two const regs. or borrows from other rows via interconnection
- → Degrades mappability of complex kernels
- 2. Depends on a host processor for overall control
  - Micro-controller basically controls data transfer
     & loop counter
  - All of other controls (e.g. reconfiguration) are carried out by the host processor even if trivial change is needed

# Proposed architecture

#### A new architecture VPCMA2

- Relaxing aforementioned limitations
- 1. Improved bank access by new data manipulator
- 2. Refined connectivity of constant registers
  - PE array has 16 constant registers (same as VPCMA)
  - All PE can use any 16 registers
- 3. Introduced an extended data bus for microcontroller

### New Data Manipulator



Offset values for each bank is introduced

Relaxed the limitation of consecutive data access <sup>19</sup>

### New Data Manipulator



Offset values for each bank is introduced

Relaxed the limitation of consecutive data access <sup>20</sup>

### Extended Data Bus

Micro-controller can handle any data in other modules



# **Evaluation Setup**

- An implementation of VPCMA2
  - ■Using Renesas SOTB 65-nm technology
    - LSTP (Low STanby Power) version
  - Synthesized by Synopsys Design Compiler 2017
- A real chip of VPCMA[7]
   Fabricated same technology
   LP (Low Power) version
  - (75% slower than LSTP)
- [7] T. Kojima, *et al.* "Real chip evaluation of a low power CGRA with optimized application mapping," 9th International Symposium on Highly-Efficient Accelerators and Reconfigurable Technologies. ACM, 2018, p. 13.



Chip photo of VPCMA[7]

### Hardware overhead

	VPCMA [7]	VPCMA2	
		1-cycle f/g	2-cycle f/g
Max Frequency (MHz) (75% scaled)	87.71	95.23 ( <mark>71.42</mark> )	125.0 (93.75)
Cell Area (mm <sup>2</sup> ) without PE array	10.04	14.55	14.22

Improved data manipulator could increase critical path delay (i.e. degradation of operating freq.)

- 2 version of designs are evaluated
- 1. Fetch&Gather are performed within 1 cycle (naïve)
- 2. Fetch&Gather take 2 cycles (to divide the long critical path)

#### 2-cycle f/g

- Does not have any effects on the frequency
- Causes 40% increase of cell area

### **Comparison of Power Consumption**

Power Consumption while running gray scale processing at 30MHz				
	VPCMA[7]	VPCMA2 (sim)		
Process version	LP	LSTP		
Standard Voltage	0.55 V	0.75 V		
Static Power	0.126 mW	0.0252 mW		
Dynamic Power	3.337 mW	4.029 mW		
Total Power	3.463 mW	4.053 mW		
		17% incr		

Compared to VPCMA(real chip), VPCMA2 (simulation)

- Reduces static power consumption because of process difference (not architecture difference)
- Increases dynamic power consumption because of the improved functionality and partially due to the higher standard voltage

# Enhanced Application Mappability



Mapping result of DCT by Genetic algorithm-based mapper[6]

VPCMA2 can accommodate large & complex kernel but VPCMA cannot
<sup>25</sup>

### Performance Improvement



# Conclusion

- This work points out a problem by data handling limitations of VPCMA
- A new SF-CGRA: VPCMA2 is proposed to relax the limitations
- Evaluation results shows
  - 10% area overhead (as a whole of chip)
  - No degradation of operating frequency
  - 17% power overhead
  - 46% performance improvement
- Future work
  - Analysis of effectiveness for other architectures
  - Evaluation of real chip implementation (under fabrication)<sup>27</sup>

# End of presentation Thank you for your attention

### Any questions?