#### A Preliminary Evaluation of Building Block Computing Systems

- Sayaka Terashima\*, Takuya Kojima\*, Hayate Okuhara\*,
  - Kazusa Musha\*, Hideharu Amano\*,
  - Ryuichi Sakamoto<sup>+</sup>, Masaaki Kondo<sup>+</sup>,
    - Mitaro Namiki§
  - \*Keio University, \*The University of Tokyo,
  - <sup>§</sup> Tokyo University of Agriculture and Technology
- 2019 IEEE 13th International Symposium on embedded Multicore/Manycore Systems-on-Chip (IEEE MCSoC-2019)

# Limitation of a Monolithic SoC

- Many requests for recent embedded system
  - High performance, high functionality
  - Low power consumption, low cost
- Increasing NRE cost of LSI chip
  - Due to complicated design, test, mask
- Problems
  - Hard to meet such demand with a single SoC
  - High cost to develop a LSI for each application (ASIC)

#### **Building Block Computing System**

• A technique of SiP (System in Package)

# **Building Block Computing Systems**

#### • For flexible & various systems

- Combining several basic chips depending on target apps.
- Using ThruChip Interface (TCI) for inter-chip communication



# TCI: ThruChip Interface[1]

- A wireless data transferring technique
  - Employing electromagnetic wave of coils
  - No need of special fabrication process
  - Up to 8 Gbps with 10<sup>-12</sup>
     bit error ratio
- TCI IP includes
  - Two SERDESes for Rx & Tx
  - An oscillator for trans. CLK

[1] Y. Take, *et al*, "3D NoC with Inductive-Coupling Links for Building-Block SiPs," IEEE Transactions on Computers, vol. 63, no. 3, pp. 748–763, 2014.



## Escalator Network by TCI Link

- Stacked chips form ring network
  - A packet-based network
    - The packet is composed of 1~17 of 35-bit flits



## Cube-2: A Prototype of Building Block Computing Systems

• Geyser<sup>[2]</sup>

- MIPS R3000 compatible CPU

- Accelerators
  - -CC-SOTB2<sup>[3]</sup>
    - High energy efficient CGRA
  - SNACC<sup>[4]</sup>
    - CNN accelerator
  - KVS<sup>[5]</sup>
    - Non-SQL DB accelerator



[2] L. Zhao, *et al.* "Geyser-2: The second prototype CPU with fine-grained run-time power gating", Proc of the 16th ASP-DAC 2011.
[3] T. Kojima, *et al.* "Real Chip Evaluation of a Low Power CGRA with Optimized Application Map- ping", Proc of the 9th HEART 2018.
[4] R.Sakamoto , *et al.* "The design and implementation of scalable deep neural network accelerator cores," in Proc. of IEEE MCSoC 2017
[5] Y.Tokuyoshi, , *et al.* "Key-valueStoreChipDesign for Low Power Consumption," in Proc of IEEE CoolChips 22 (2019).

## Shared Memory for Twin-Tower (SMTT)

- A bridge SRAM chip
  - Has two TCI IP
  - Shares 256KB between twin towers
  - Provides atomic operation *Fetch&Dec* for synchronization among stacked chips
  - Supports DMA transfer





# Overview of GeyerTT

- Geyser architecture
  - MIPS R3000 compatible CPU
    - General compilers are available
  - Responsible for host controller of Cube-2 system
  - Including 2-way d-cache、 2-way i-cache、 TLB
- GeyserTT
  - A real chip Implementation of Geyser for Twin-Tower
  - Three TCI IP for various stacking structure



# Overview of SNACC

- SNACC architecture
   Composed of 4 cores
- Each core consists of
  - —Custom SIMD unit
  - General-purpose ALU & Regfile
  - 5 distributed memories
    - 1. Instruction
    - 2. Input data
    - 3. Weight data
    - 4. Look-up-table
    - 5. Write buffer



9

# Memory-Mapped Chips



# Contributions of This Work

- Fabricating & evaluating Cube2-family chips
  - Focusing on GeyserTT, SNACC, SMTT
  - About power consumption & performance
  - Based on real chip measurement
- Evaluating TCI IP itself
  - About feasibility of this technology
  - About power consumption & performance
  - Based on real chip measurement
- Demonstrating possibility for practical apps.
  - With CNN application as a case study

# **Real Chip Implementation**

		GeyserTT			
Process	Renesas SOTB 65nm				
Supply voltage	0.75 V				
Design	Verilog HDL				
Synthesis	Synopsys Design Compiler 2016.03-SP4	TCI IP		SMTT	
Place & Route	Synopsys IC Compiler 2016.03-SP4		Stacked	Chip	S
Chip size	SNACC & GeyserTT <b>3mm x 6mm</b>			•	
	SMTT 6mm x 6mm				
Target Frequency	SNACC & GeyserTT <b>50MHz</b>				
	SMTT <b>100MHz</b>				
	TCI IP <b>50MHz</b>				13

#### **Evaluation: Power Consumption**



# **Evaluation: TCI performance**

- GeyserTT x SNACC case
  - Bidirectional links can work
  - Compared to design value (50MHz)
    - TCI consumes maximum 2.0x power & achieves 0.12x performance
- GeyserTT x SMTT case
  - Upward link does not work
- But the latest chip shows
  - 10~15MHz transfer
  - 1.5x power than design value







## Evaluation: TCI power consumption



#### Case study: Processing FC layers of a CNN

#### • Last two FC layers of AlexNet<sup>[6]</sup>



layer	# of input	# of output	Kernel size	Bias
FC7	4096	4096	(4096, 4096)	4096
FC8	4096	1000	(1000, 4096)	1000

[6] A. Krizhevsky, I. Sutskever and G. E. Hinton: "Imagenet classification with deep convolutional neural networks", Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1, NIPS'12, USA, Curran Associates Inc., pp. 1097–1105 (2012).

#### **Evaluation: Simulated Configurations**

- Evaluated system configurations
  - 1. <u>GeyserTT</u>
  - 2. <u>GeyserTT x2 + SMTT</u>
  - 3. GeyserTT + SNACC
  - 4. GeyserTT x2 + SNACC x2 + SMTT



#### **Evaluation: Simulated Configurations**

- Evaluated system configurations
  - 1. GeyserTT
  - 2. GeyserTT x2 + SMTT
  - 3. <u>GeyserTT + SNACC</u>
  - 4. <u>GeyserTT x2 + SNACC x2 + SMTT</u>



#### Evaluation: Execution time @50MHz



• The execution time for each configuration includes data transfer time through TCI

# Conclusion

- Evaluating some real chip fabricated with Renesas SOTB 65nm technology
  - MIPS R3000 processor ~35mW @ 50MHz
  - CNN accelerator & memory chip ~4mW @ 50MHz
- Demonstrating chip stacking with TCI
  - Communications partially work
  - Much larger power is consumed than designed one
  - A twin-tower system achieves x6.0 higher performance
- Future work
  - Optimization of TCI power using sleep mode
  - Refinement of power grid for TCI IP
    - Partially completed
  - Use of other family chip such as CC-SOTB2 & KVS