A Scalable Multi-Chiplet Deep Learning Accelerator with Hub-Side 2.5D Heterogeneous Integration

Zhanhong Tan¹, Yifu Wu², Yannian Zhang², Haobing Shi², Wuke Zhang², Kaisheng Ma¹

¹Tsinghua University, ²Polar Bear Tech







Abstract

With the slowdown of Moore's law, the scenario diversity of specialized computing, and the rapid development of application algorithms, an efficient chip design requires modularization, flexibility, and scalability. In this study, we propose a Chiplet-based deep learning accelerator prototype that contains one HUB Chiplet and six extended SIDE Chiplets integrated on an RDL layer for the 2.5D package. The SIDE and the HUB contain one and four AI cores, respectively.

Given that our Chiplet-system targets diverse scenarios via scalable connected SIDE Chiplets, we need to handle three challenges: a) devise a flexible architecture design supporting diverse shapes, b) search for a workload mapping with low die-to-die communication, and c) adopt a high-bandwidth die-to-die interface to maintain efficient data transfer.

This study proposes a flexible neural core (FNC) featuring dynamic bit-width computing and flexible parallelism. Next, we use a hierarchy-based mapping scheme to decouple different parallelism levels and help analyze the communication. A 12Gbps D2D interface is introduced to achieve 192Gb/s bandwidth per D2D port with 1.04pJ/bit efficiency and 55µm bump pitch.

The proposed seven-Chiplet accelerator achieves a peak performance of **10/20/40 TOPS for INT16/8/4**. When enabling 0~6 SIDE Chiplets, the system power ranges from 4.5W to 12W. The power efficiency of the FNC is **2.02TOPS/W** while that of the overall system is **1.67TOPS/W**.



Background and Challenges



Decouple a monolithic SoC into Chiplets

- Better die yield
- Scalability for diverse scenarios
- Rapid development pace to deliver new products

Challenges





Overall Architecture



Flexible Neural Core (FNC)

• Reconfigurable architecture for the shape diversity

Mapping dataflow

• Die-to-Die communicationaware workload generator

Interconnection

- High-bandwidth Die-to-Die based on 2.5D package
- Efficient chiplet routing unit (CLRU)



Flexible Neural Core

Flexible Interconnect

• Arbitrary tile-based workload assignment to 8 cores via a configurable interconnect fabric





Flexible Neural Core

■ The MAC Pair Supporting for Dynamic Bit-width



- Support three quantization modes
 - ➢ 8b-Acitvation × 4b-Weight
 - ➢ 8b-Acitvation × 8b-Weight
 - > 16b-Acitvation × 8b-Weight
- Each INT-8 MAC-Pair has eight 4×4 multipliers for mode reuse
- In three modes, the bandwidth and compute resources of one MAC-pair are fully utilized



Flexible Neural Core

Flex-Interconnect and Configurable Weight Buffer



• Support diverse eight-PE compositions

- 8-tile mode: share weights across 8 PEs for independent output in height/width
- 4-tile mode: share weights across 4 PEs and 2 4-PE groups process 2 chunks of output channels
- 2-tile mode: 4 2-PE groups for 4 chunks of output channels
- 1-tile mode: 8 PEs for 8 chunks of output channels



Pojar Beai Tech



Dynamic Workload Parallelism



Loop order in the temporal primitive

The overhead bias helps to search for a low D2D communication mapping

Reuse Region Critical Position	IS		
for h2 = [0 : H2):	for c2 = [0 : C2):	for h2 = [0 : H2):	for $c2 = [0 : C2)$:
for w2 = [0 : W2):	for h2 = [0 : H2):	for w2 = [0 : W2):	for $h2 = [0 : H2)$:
for c2 = [0 : C2):	for $w^2 = [0 : W^2)$:	for c2 = [0 : C2):	for w2 = [0 : W2):
for h1 = [0 : H1):	for $h1 = [0 : H1)$:	for h1 = [0 : H1):	for $h1 = [0 : H1)$:
for w1 = [0 : W1):	for $w1 = [0 : W1):$	for w1 = [0 : W1):	for w1 = [0 : W1):
for c1 = [0 : C1):	for c1 = [0 : C1):	for c1 = [0 : C1):	for c1 = [0 : C1):

Critical Position for the "X" buffer: the

inner-most loop related to the index of the

X-buffer data (decide the data size on-core)

Reuse Region for the "X" buffer: indicate

the reuse efficiency when caching the data

Search for an optimized loop range with

the highest memory utilization (the largest

data size that can be buffered on-core) and

reuse efficiency for each buffer

in inner loops decided by critical position

Example-1 for the W-Buf analysis Example-2 for the W-Buf analysis Example-3 for the L1-Buf analysis Example-4 for the L1-Buf analysis

Notation: HO, WO, CO: height, width, and channel of the output tensor; X_t : the tile for a Chiplet; X_c : the sub-tile for a PE



A Scalable Multi-Chiplet Deep Learning Accelerator with Hub-Side 2.5D Heterogeneous Integration

Chiplet Interconnection and Package

High-Bandwidth D2D Interface



Bandwidth per D2D	RX: 192Gb/s TX: 192Gb/s	RX(TX) Lane	2(2)
		Data width per lane	8bit
		Data Rate	12Gbps
Bump Pitch	55µm	Package	2.5D
Area	2.2×0.5mm	Power	1.04pJ/bit

Chiplet Router Unit (CLRU)



- Four FIFO queues to deal with burst transfer
- Data parser: support the data request from another Chiplet (access memory / other CLRU)
 - > The head packages indicate the transfer mode



Chiplet Interconnection and Package



- Non-conflict IO layout in the HUB Chiplet to improve the fan-out efficiency
- 2.5D integration with a **high-density 65nm RDL** layer providing 55µm bump pitch
- The RDL layer contributes to a simpler 8-layer substrate of 3-2-3









Evaluation

Evaluation on computing-bound workload **Evaluations on one Flexible Neural Core** 800 600 DeepLabv3+ 59.6ms, 16.8fps 400 SqueezeNet 1.5ms, 666.7fps

Runtime (ms) Speedup by ×9.95 200 ResNet101 15.1ms, 66.2fps 0 2 Cores 4 Cores 6 Cores 8 Cores 10 Cores 1 Core ResNet50 8.7ms, 114.9fps **Evaluation on IO-bound workload** YOLOv3 119.3ms, 8.4fps 50 YOLOv5s 13.1ms, 76.3fps Runtime (ms) **Speedup by** YOLOv5n 6.2ms, 161.3fps 25 ×6.25 **DrivableNet** 89.7ms, 11.1fps 0 utilization 0% 60% 30% 90% 6 Cores 8 Cores 10 Cores 1 Core 2 Cores 4 Cores



System Board and Demo

System PCIe-based Board



Demo for running concurrent 4 models



Model 4: DeepLab v3+





HUB Chiplet



RDL Layer for 2.5D Package

Items		Specifications	
Technology		CMOS 12nm	
Die Area	HUB Chiplet	8.5mm × 6.8 mm = 57.8 mm ²	
	SIDE Chiplet	3.5mm × 2.8 mm = 9.8 mm ²	
Supply Voltage		0.8V ~ 1.2V	
Frequency		100MHz – 1GHz	
Peak Performance	INT4	40TOPS (8b-A × 4b-W)	
	INT8	20TOPS (8b-A × 8b-W)	
	INT16	10TOPS (16b-A × 8b-W)	
NPU Core Efficiency		2.02TOPS/W	
Power		4.5W ~ 12W	
D2D Bandwidth		6×24GB/s for TX/RX	
External Memory Bandwidth		64GB/s (GDDR6)	
Bump Pitch for 2.5D Pkg		55µm	



Comparison with Prior Multi-Chiplet Accelerator Works

	Simba (NVIDIA)	CHIMERA (Stanford)	NetFlex (A*STAR)	Ours
Year	2019	2021	2022	2023
Technology	16nm	40nm	22nm	12nm
Area	6mm²	29.2mm²	11.1mm ²	HUB: 57.8mm ² Side: 9.8mm ²
Memory Size	752KB SRAM	0.5MB SRAM 2MB RRAM	2492KB SRAM	HUB: 1.7MB Side: 439KB
Voltage	0.42V ~ 1.2V	1.1V	0.6V ~ 0.89V	0.8V ~ 1.2V
Frequency	161MHz – 2001MHz	200MHz	190.3 – 492.3MHz	600MHz – 1.2GHz
Power	30 – 4160mW	126mW	57.6 – 499.8mW	Side: 0.72W Hub: 4.75W
Performance (TOPS)	0.32 – 4.01 (INT8)	2.2 (INT8, FP16)	0.41 – 1.07 (INT16)	Side Die: 1/2/4 for INT16/8/4, Hub Die: 4/8/16 for INT16/8/4, Total: 10/20/40 for INT16/8/4
Package	Organic MCM	PCB	HD-FOWLP	2.5D RDL
D2D I/O	GRS	C2C Links	AIB	12Gbps Parallel Interface
I/O Energy	0.82 – 1.75pJ/b	77pJ/b	3.07pJ/b	1.04pJ/bit



Thank You

tanzh@mails.tsinghua.edu.cn



A Scalable Multi-Chiplet Deep Learning Accelerator with Hub-Side 2.5D Heterogeneous Integration