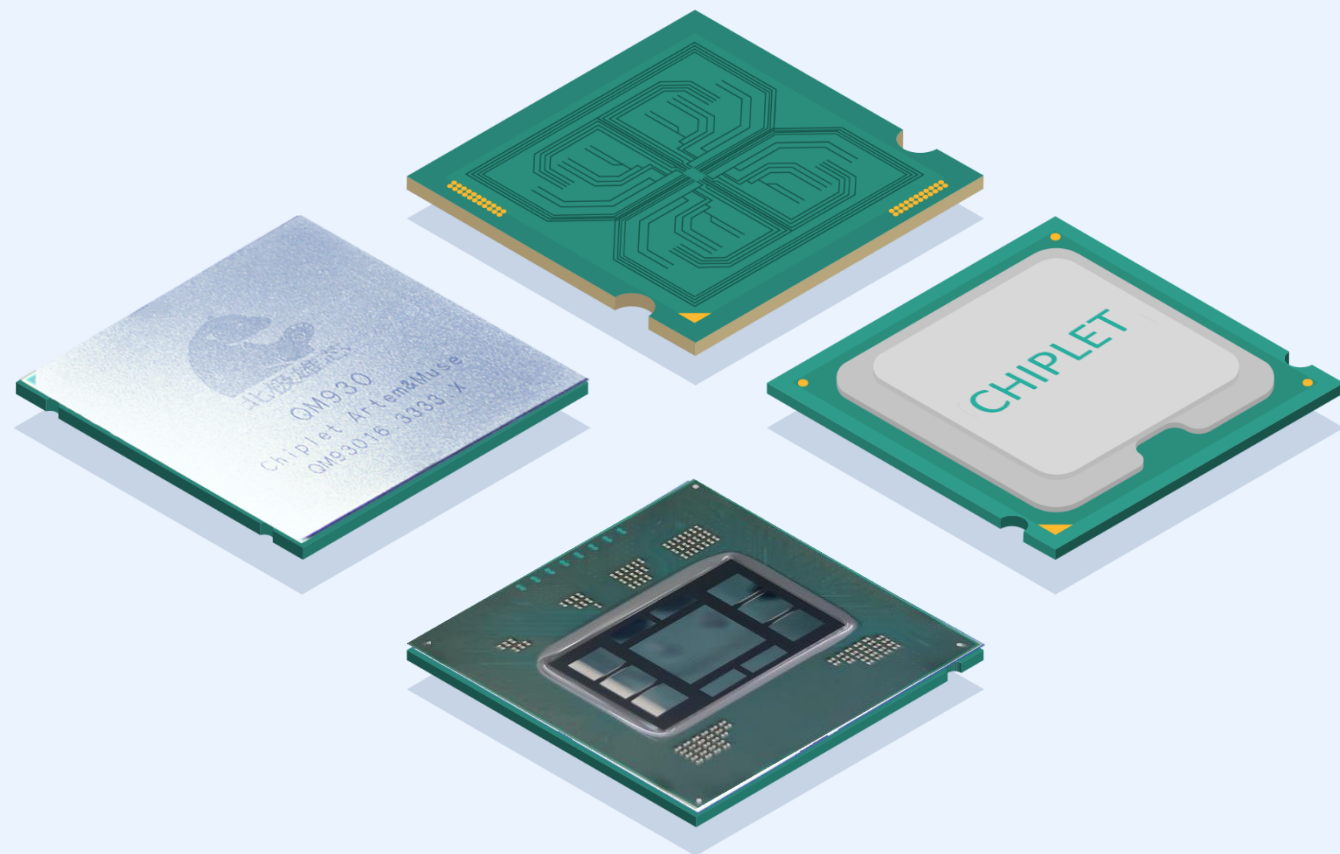


# A Scalable Multi-Chiplet Deep Learning Accelerator with Hub-Side 2.5D Heterogeneous Integration

Zhanhong Tan<sup>1</sup>, Yifu Wu<sup>2</sup>, Yannian Zhang<sup>2</sup>,  
Haobing Shi<sup>2</sup>, Wuke Zhang<sup>2</sup>, Kaisheng Ma<sup>1</sup>

<sup>1</sup>Tsinghua University,

<sup>2</sup>Polar Bear Tech



# Abstract

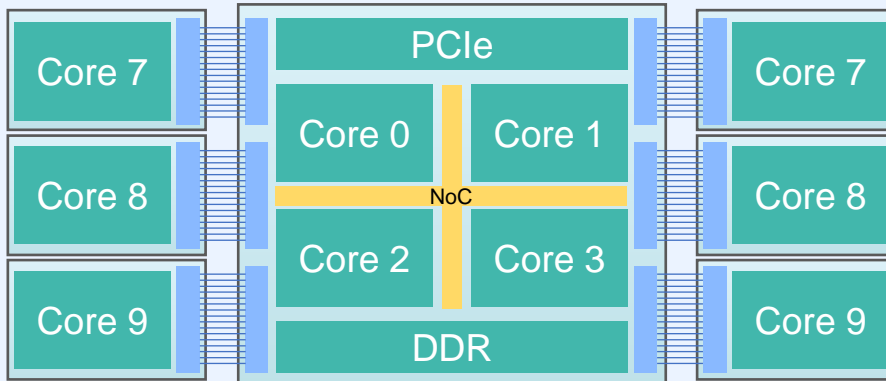
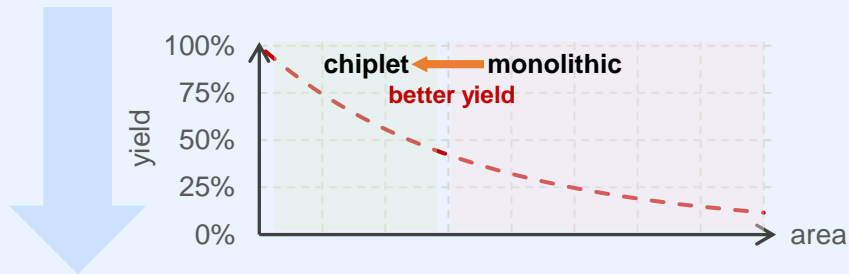
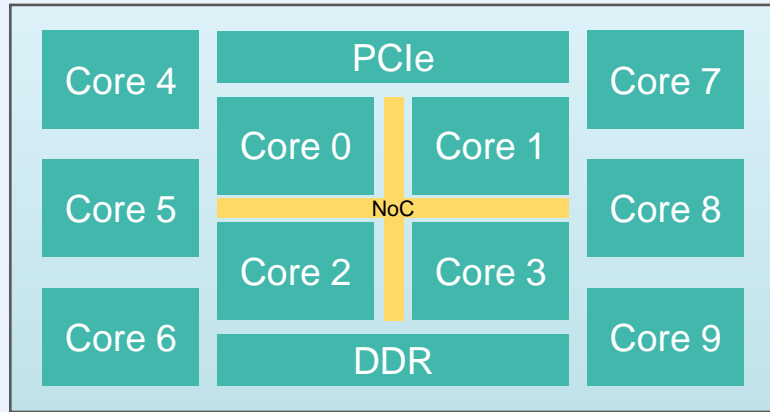
With the slowdown of Moore's law, the scenario diversity of specialized computing, and the rapid development of application algorithms, an efficient chip design requires modularization, flexibility, and scalability. In this study, we propose a Chiplet-based deep learning accelerator prototype that contains 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 55 $\mu$ m bump pitch.

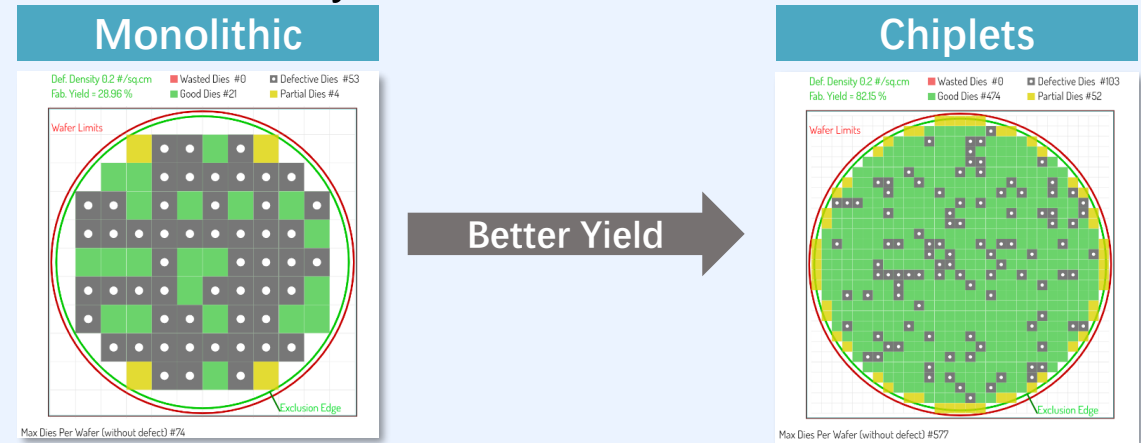
The proposed seven-Chiplet accelerator achieves a peak performance of **10/20/40 TOPS for INT16/8/4**. When enabling 0~6 SIDE Chiplets, the system power ranges from 4.5W to 12W. The power efficiency of the FNC is **2.02TOPS/W** while that of the overall system is **1.67TOPS/W**.

# Background and Challenges



## Decouple a monolithic SoC into Chiplets

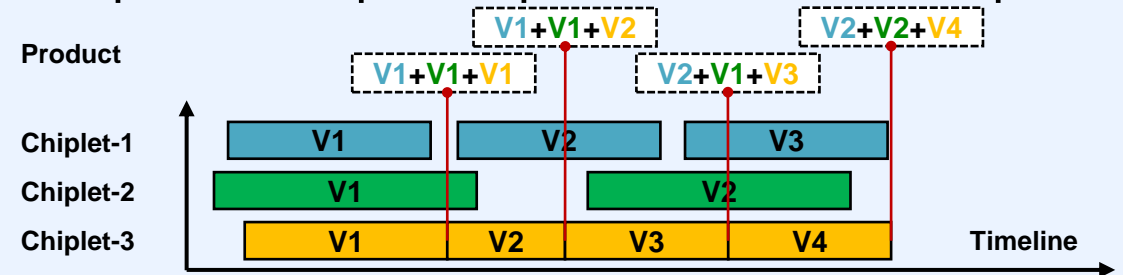
- Better die yield



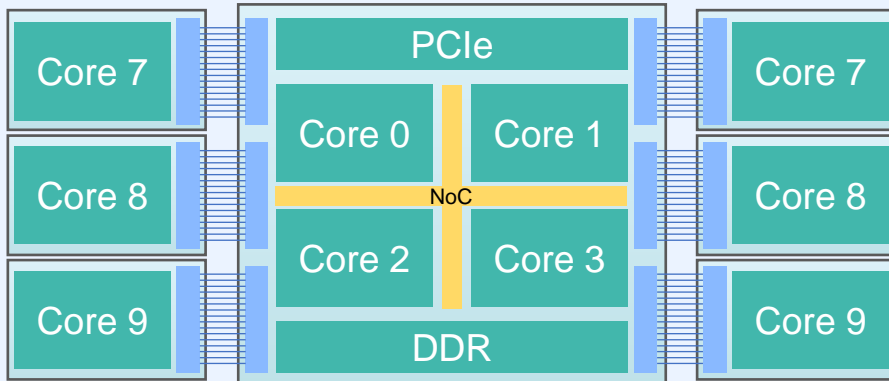
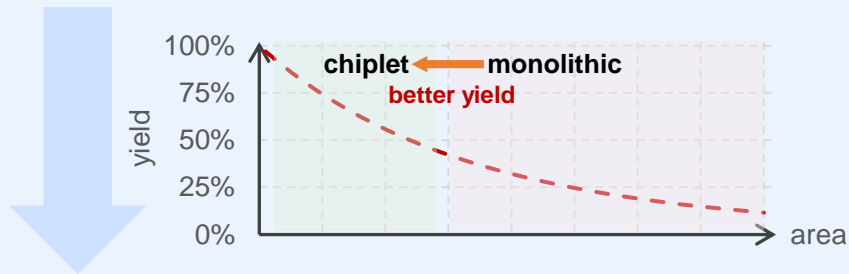
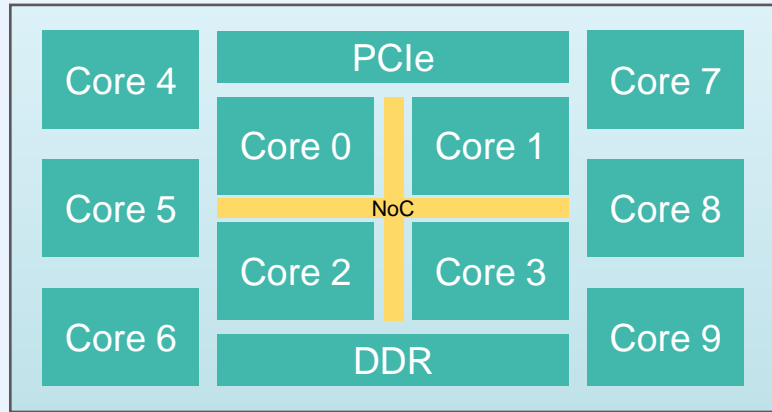
- Scalability for diverse scenarios

$$M \times \text{Chiplet1} + N \times \text{Chiplet2} + K \times \text{Chiplet3}$$

- Rapid development pace to deliver new products



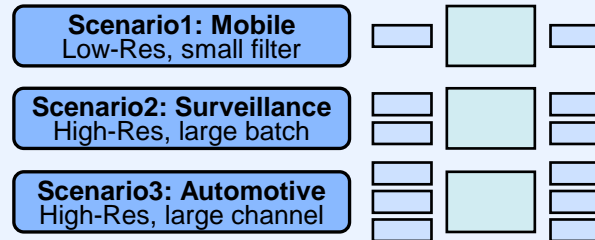
# Background and Challenges



## Decouple a monolithic SoC into Chiplets

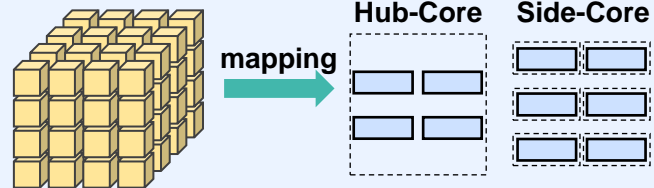
- Better die yield
- Scalability for diverse scenarios
- Rapid development pace to deliver new products

## Challenges



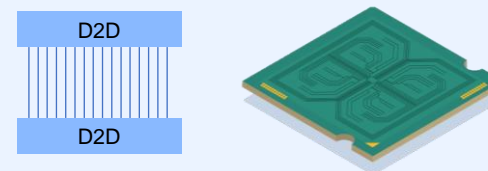
### Challenge-1

Flexible architecture design supporting diverse shapes



### Challenge-2

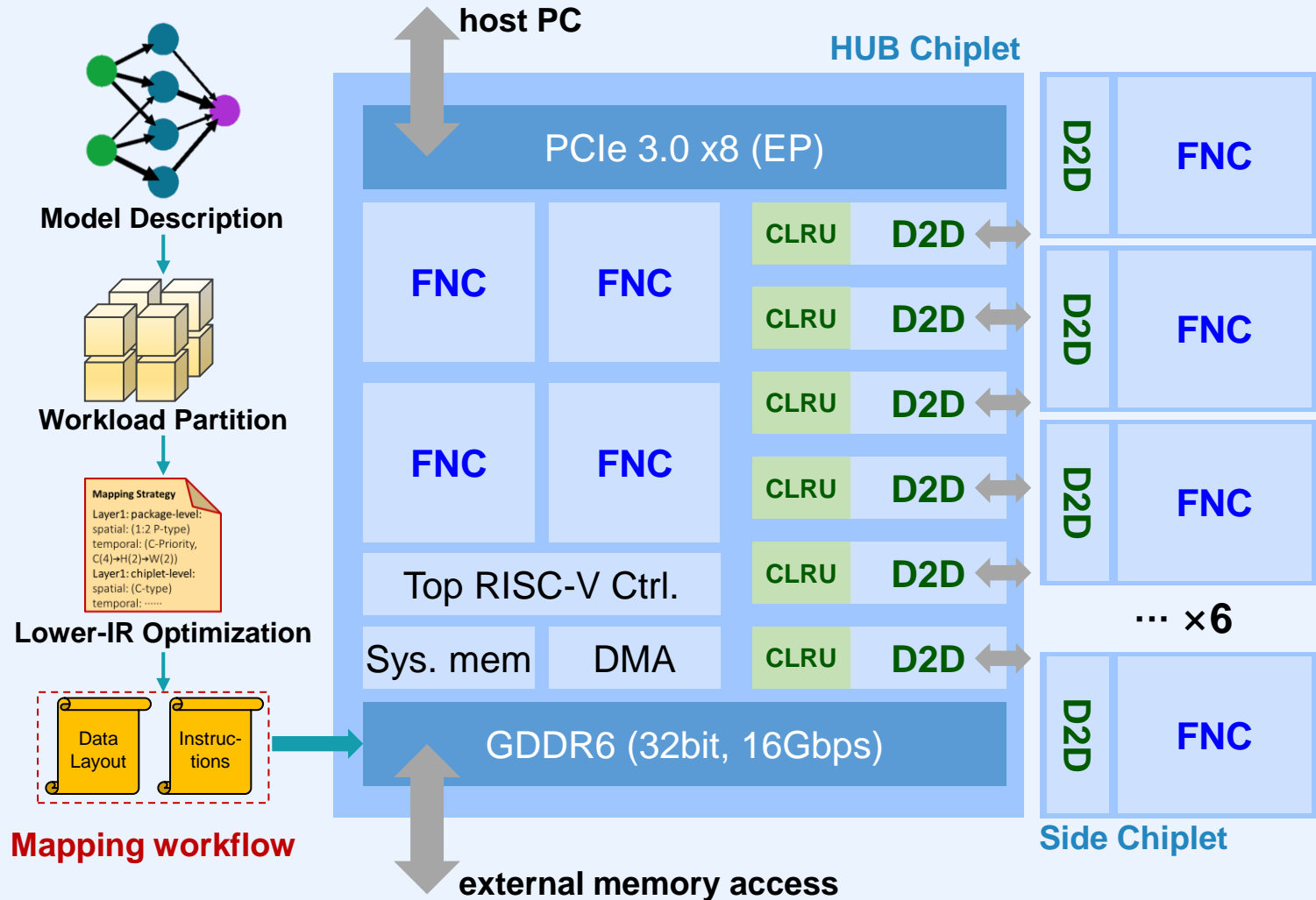
Efficient workload mapping to optimize Die-to-Die communication



### Challenge-3

High-bandwidth D2D and high-density package

# Overall Architecture



## Flexible Neural Core (FNC)

- Reconfigurable architecture for the shape diversity

## Mapping dataflow

- Die-to-Die communication-aware workload generator

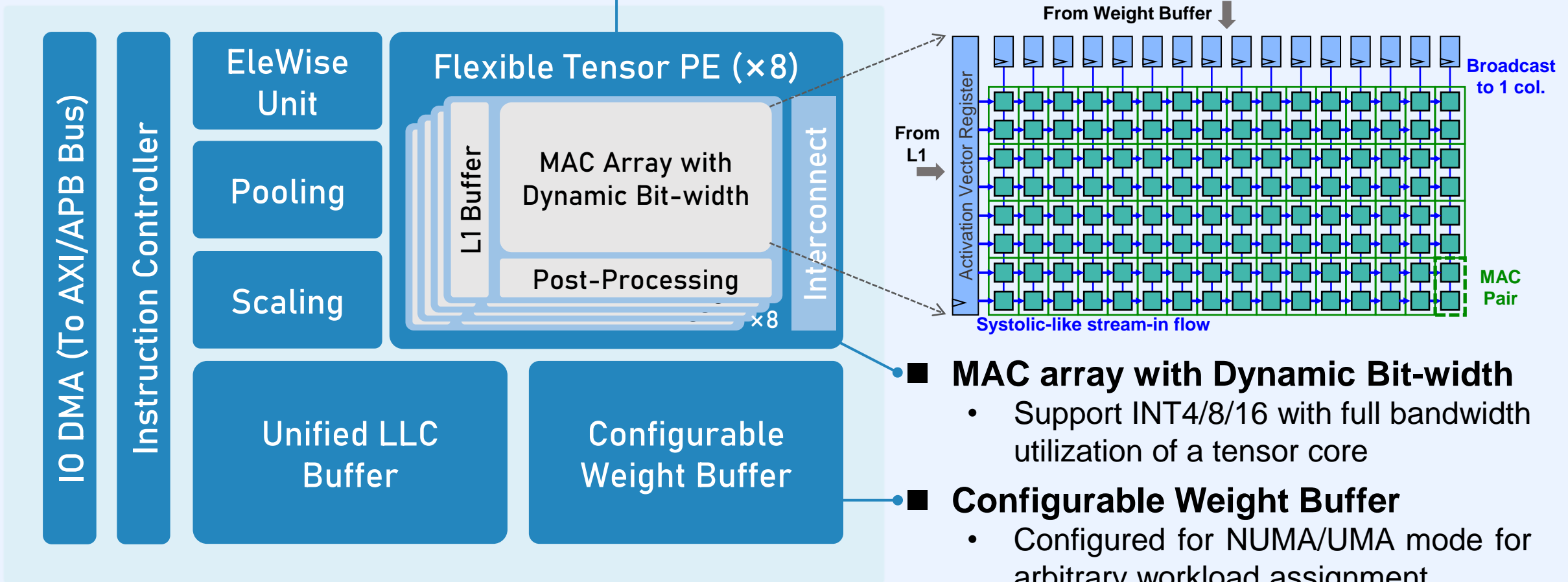
## Interconnection

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

# Flexible Neural Core

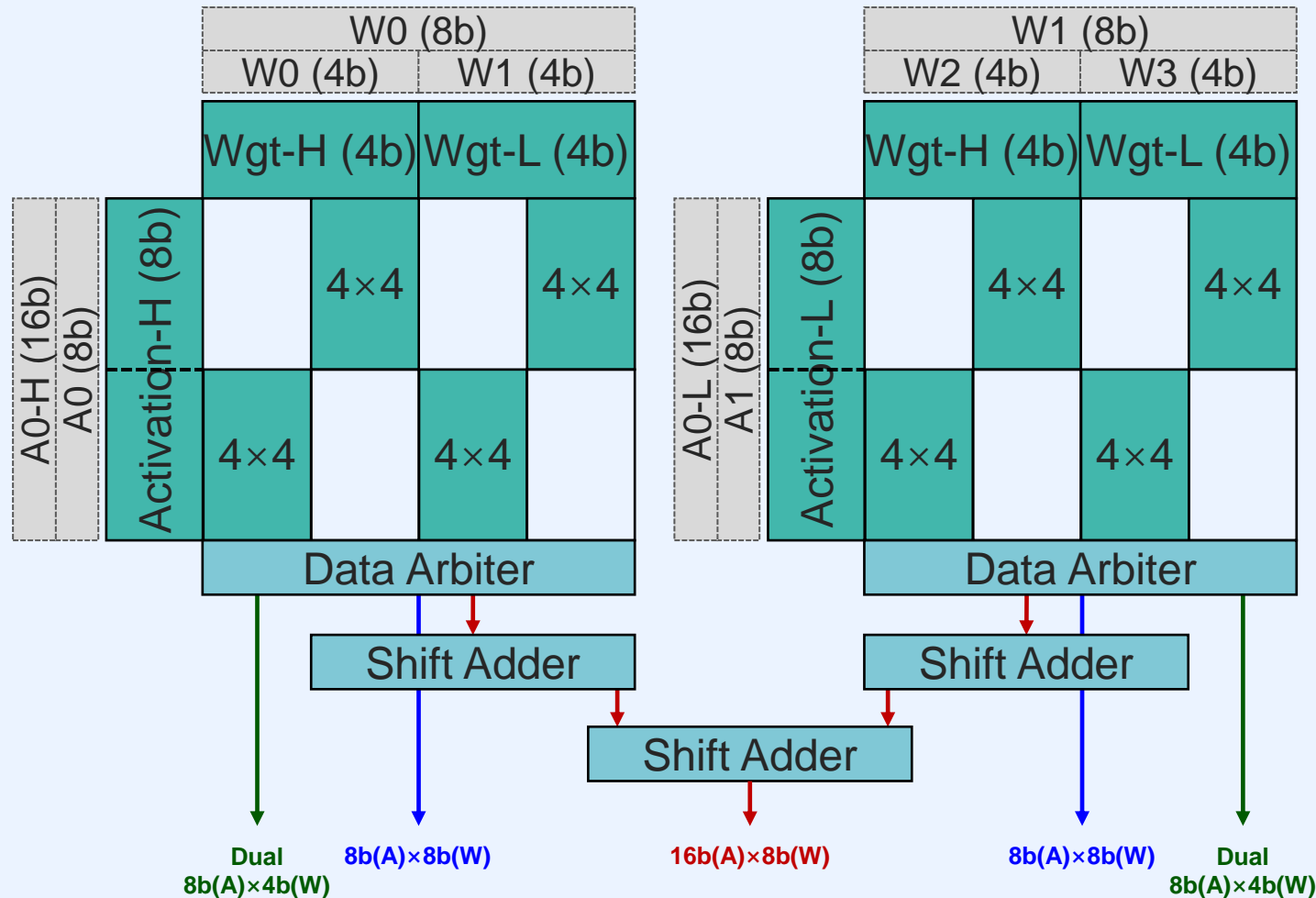
## Flexible Interconnect

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



# Flexible Neural Core

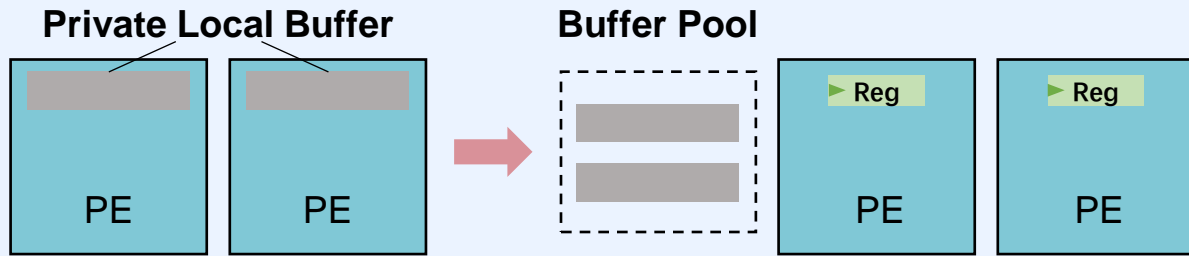
## ■ The MAC Pair Supporting for Dynamic Bit-width



- **Support three quantization modes**
  - 8b-Activation × 4b-Weight
  - 8b-Activation × 8b-Weight
  - 16b-Activation × 8b-Weight
- **Each INT-8 MAC-Pair has eight 4x4 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

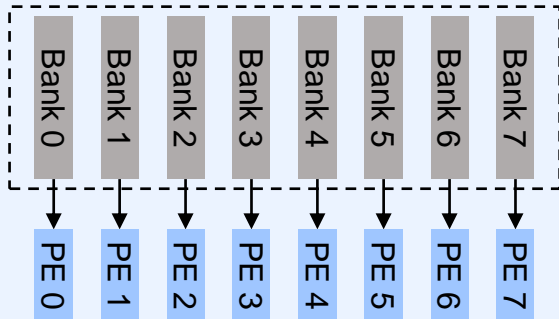


- ☹ Duplication for shared data
- ☺ Low access latency

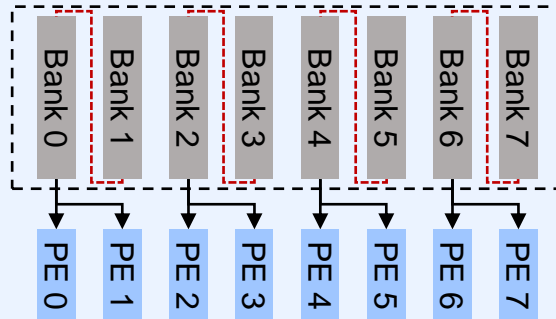
- ☺ Unified into a pool to share data
- ☺ Weight-stationary in each PE to overlap weight-loading in local registers

- **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

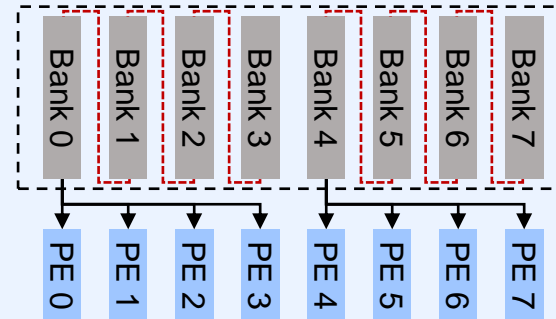
1. UMA Mode, each bank for 1 core



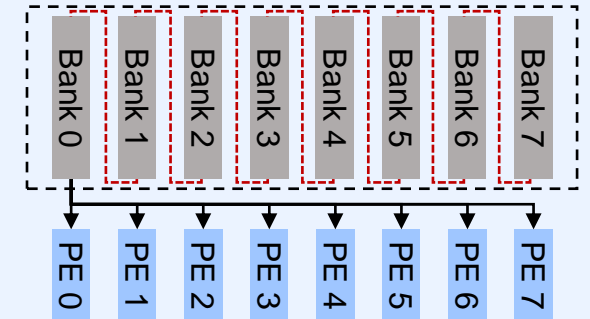
2. Dual-NUMA, each bank for 2 cores



3. Quad-NUMA, each bank for 4 cores



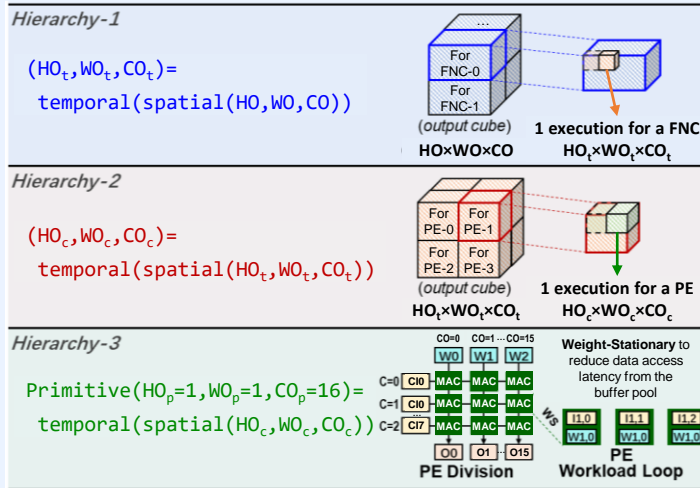
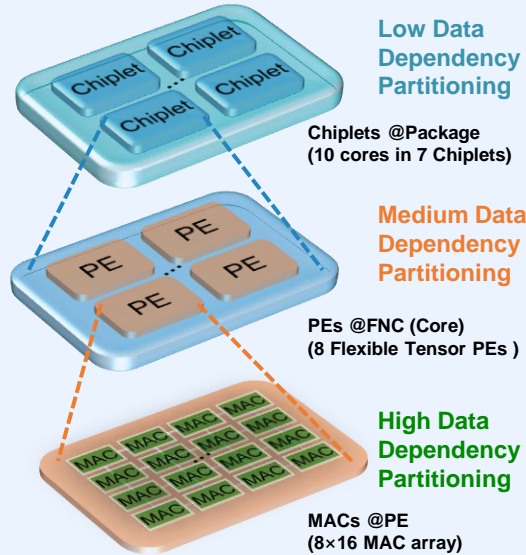
4. Full-NUMA, each bank for 8 cores



Each bank has 16 sub-banks for 16 columns of MACs



# Dynamic Workload Parallelism

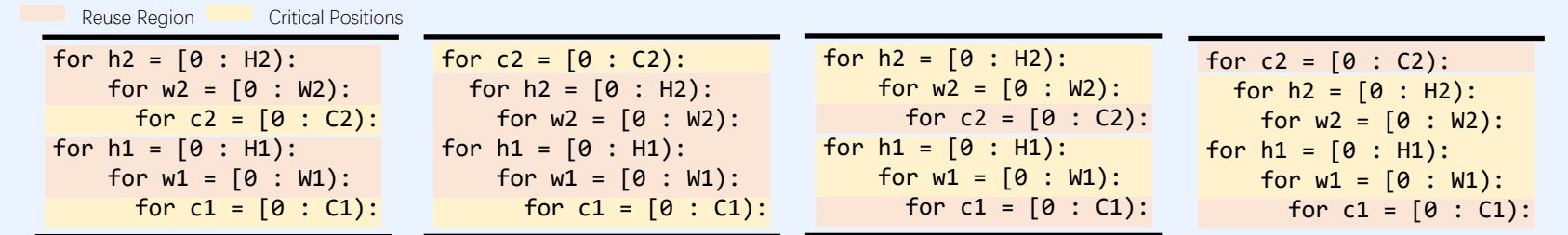


- **Critical Position for the “X” buffer:** the inner-most loop related to the index of the X-buffer data (decide the data size on-core)
- **Reuse Region for the “X” buffer:** indicate the reuse efficiency when caching the data in inner loops decided by *critical position*
- Search for an optimized loop range with the highest memory utilization (the largest data size that can be buffered on-core) and reuse efficiency for each buffer

```
// Package-Level for Chiplet Parallelism
(HOt, WOt, COt) = temporal(spatial(HO, WO, CO))
for c2 = [0 : C2): high data access overhead
  for h2 = [0 : H2): // H2 * HOt = HO
    for w2 = [0 : W2): // W2 * WOt = WO
      for c2 = [0 : C2): // C2 * COt = CO
// Chiplet-Level for PE Parallelism
(HOc, WOc, COc) = temporal(spatial(HOt, WOt, COt))
for c1 = [0 : C1): low data access overhead
  for h1 = [0 : H1): // H1 * HOc = HOt
    for w1 = [0 : W1): // W1 * WOc = WOt
      for c1 = [0 : C1): // C1 * 8 = COt
```

Loop order in the temporal primitive

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

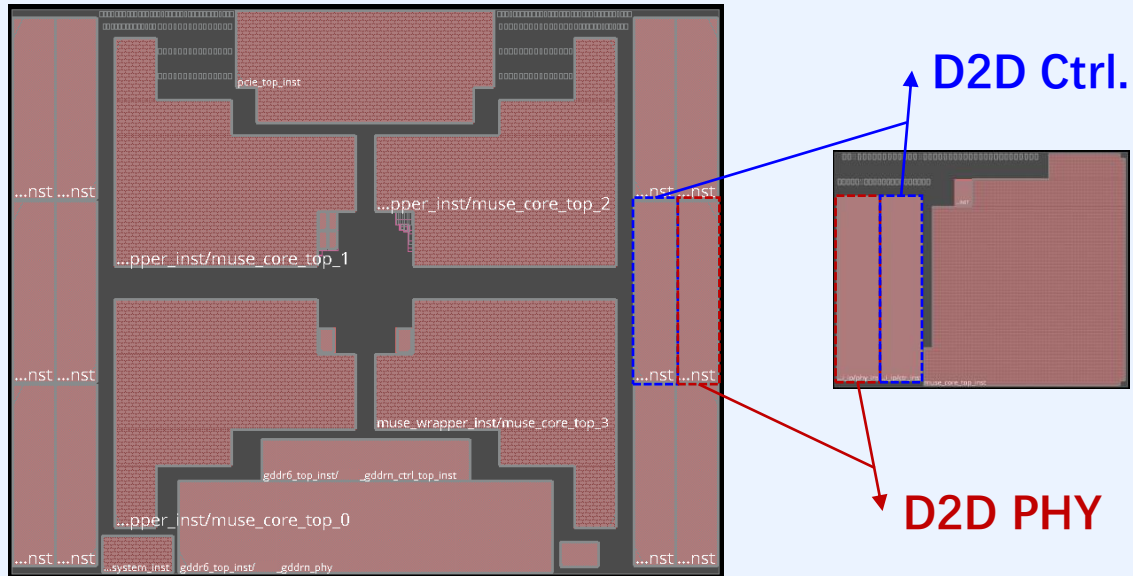


Example-1 for the W-Buf analysis Example-2 for the W-Buf analysis Example-3 for the L1-Buf analysis Example-4 for the L1-Buf analysis

**Notation:** HO, WO, CO: height, width, and channel of the output tensor; X<sub>c</sub>: the tile for a Chiplet; X<sub>c</sub>: the sub-tile for a PE

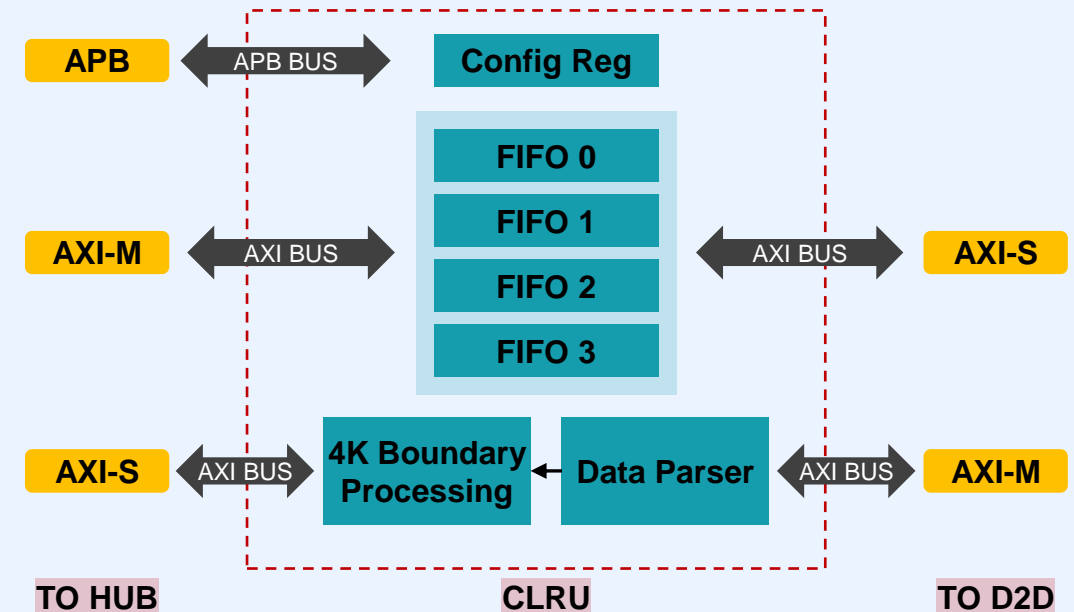
# Chiplet Interconnection and Package

## High-Bandwidth D2D Interface



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

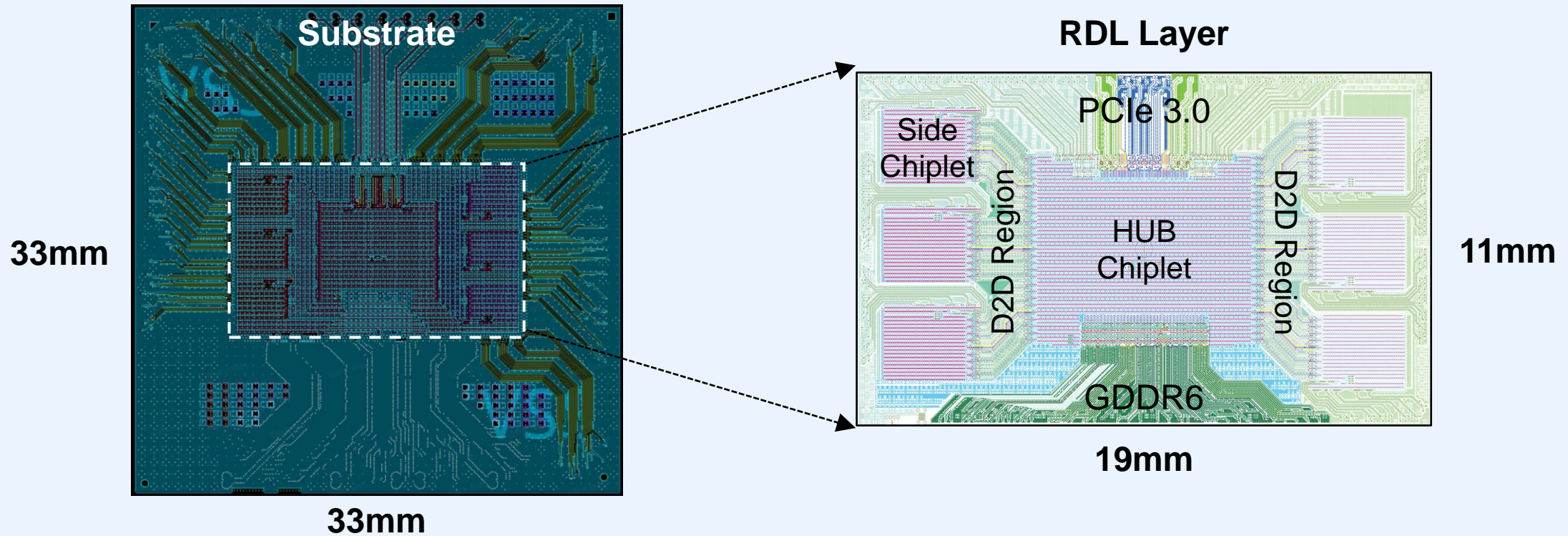
## Chiplet Router Unit (CLRU)



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

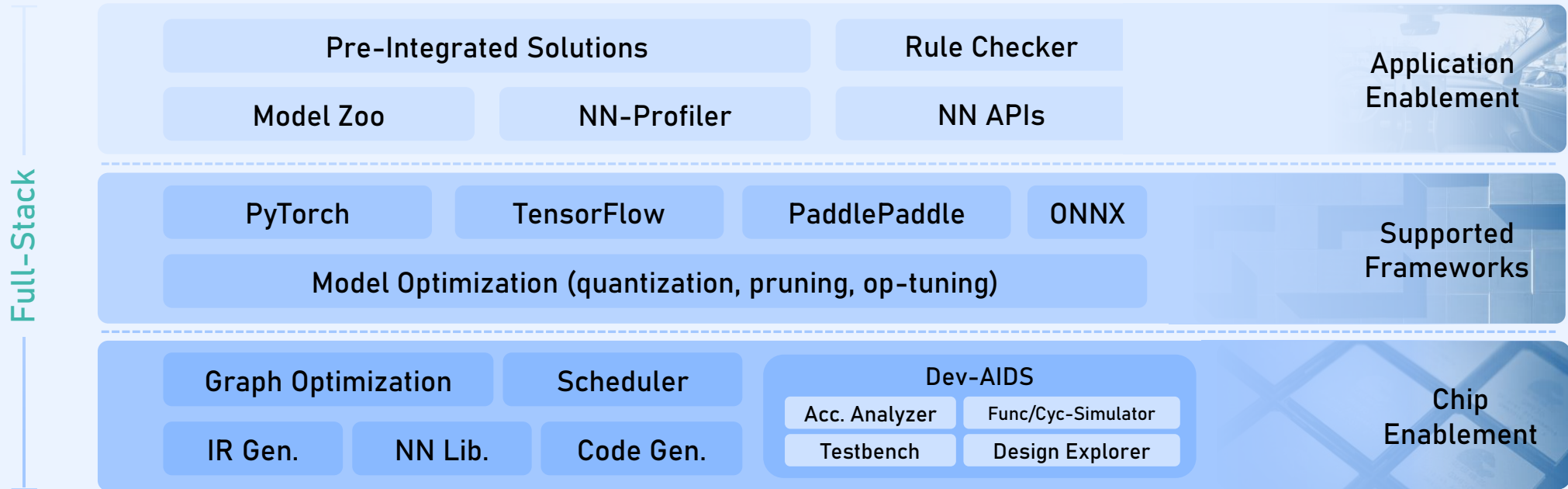
➤ *The head packages indicate the transfer mode*

# 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 55 $\mu$ m bump pitch
- The RDL layer contributes to a simpler 8-layer substrate of 3-2-3

# Software Stack



Latency-Constraint Optimization

- Core-level task scheduling



Sample-per-Sec Optimization

- Batch-level Parallelism



Query-Per-Sec Optimization

- Workload pipeline for high throughput

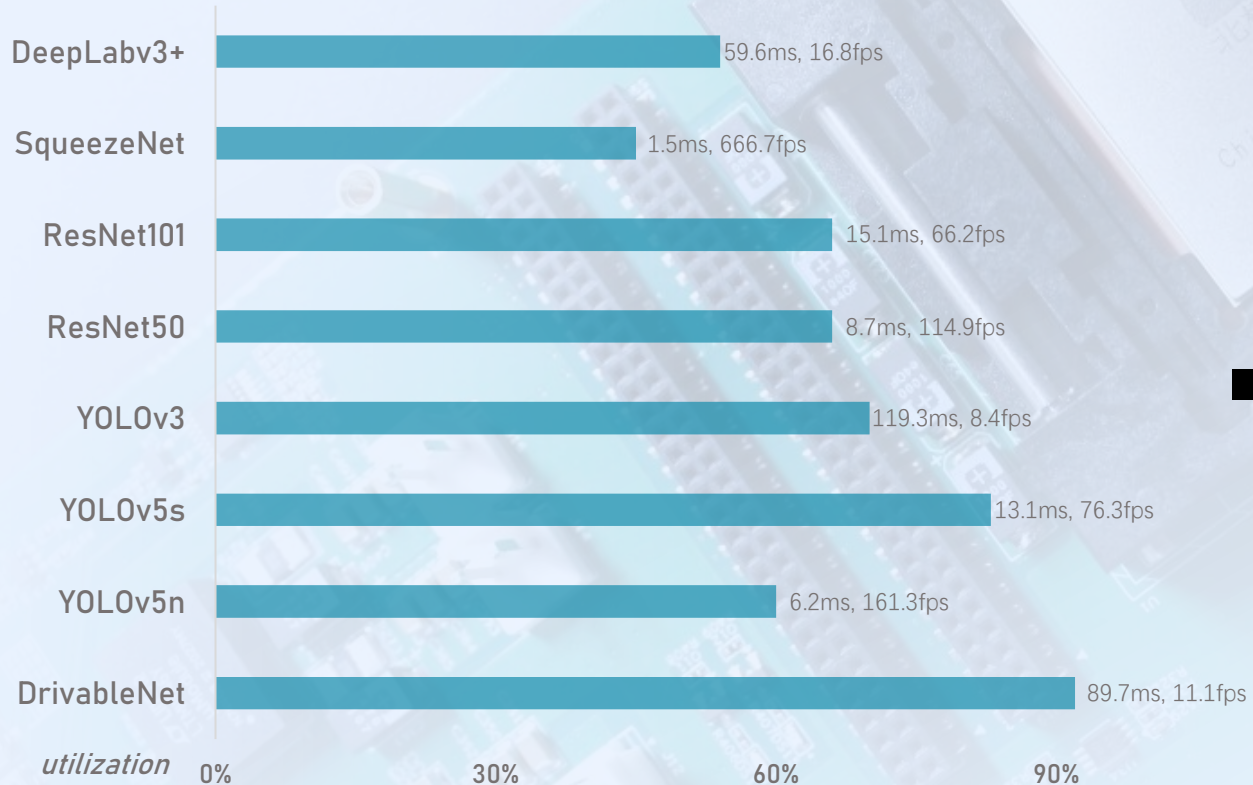


Response Time Optimization

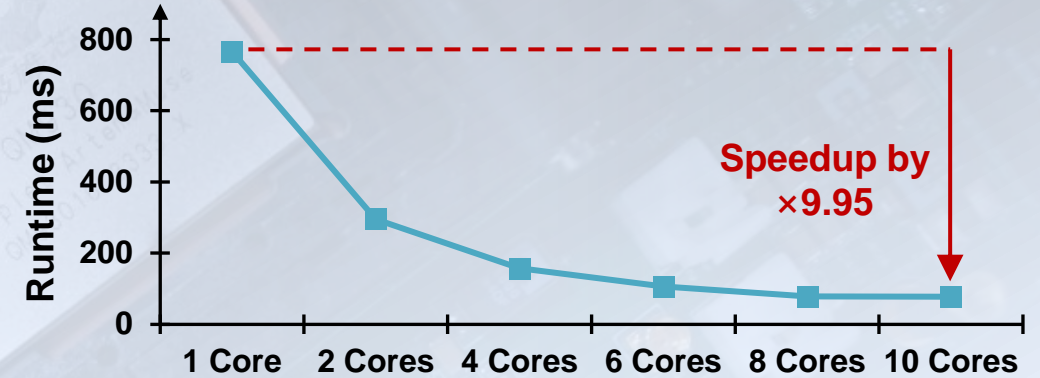
- Multi-level parallelism for high utilization

# Evaluation

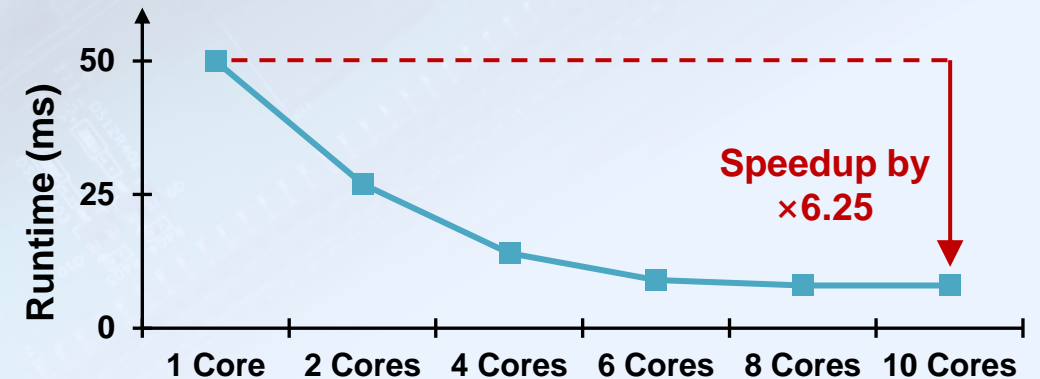
## Evaluations on one Flexible Neural Core



## Evaluation on computing-bound workload

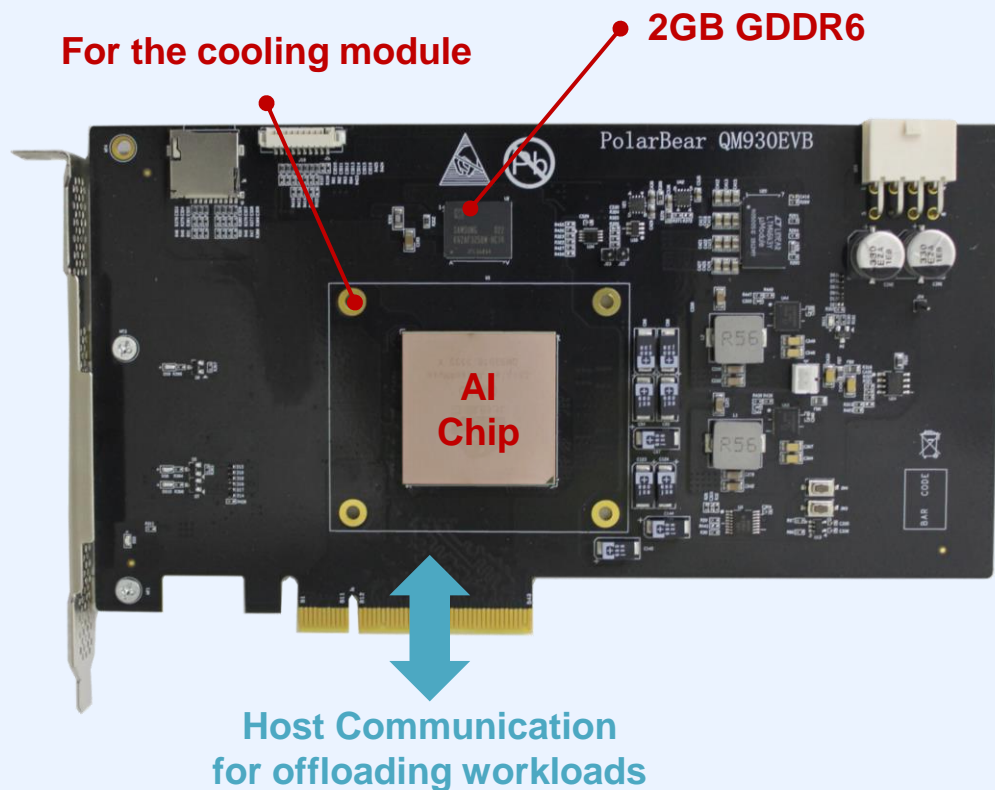


## Evaluation on IO-bound workload

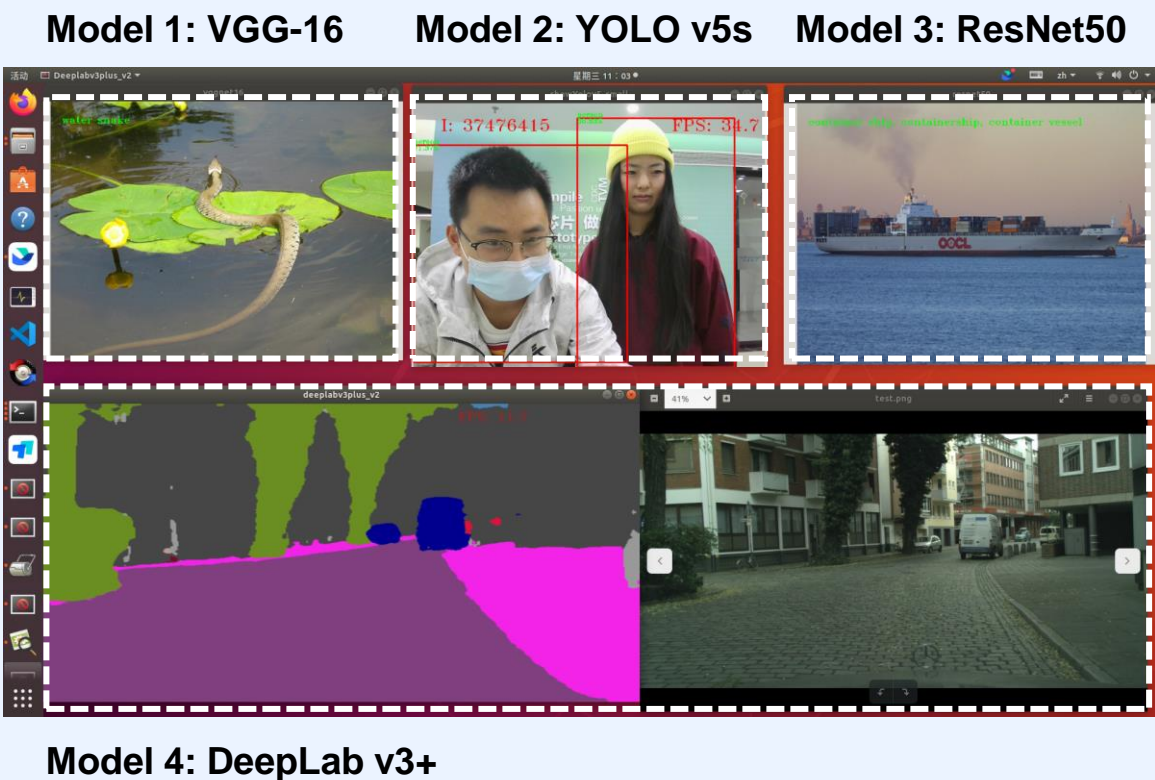


# System Board and Demo

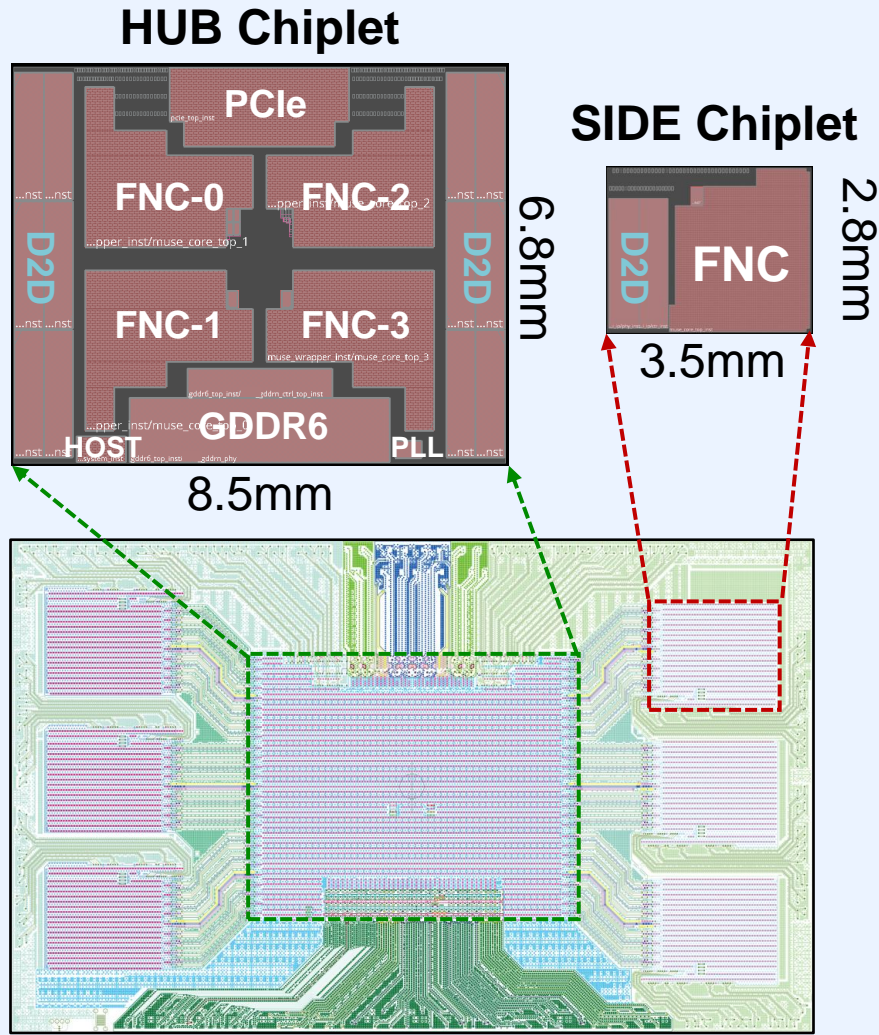
## ■ System PCIe-based Board



## ■ Demo for running concurrent 4 models



# Chip Summary



| Items                     |              | Specifications                      |
|---------------------------|--------------|-------------------------------------|
| Technology                |              | CMOS 12nm                           |
| Die Area                  | HUB Chiplet  | 8.5mm × 6.8mm = 57.8mm <sup>2</sup> |
|                           | SIDE Chiplet | 3.5mm × 2.8mm = 9.8mm <sup>2</sup>  |
| 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

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



# Thank You

[tanzh@mails.tsinghua.edu.cn](mailto:tanzh@mails.tsinghua.edu.cn)

