# Heterogeneous 3D Chiplet Integration for AI application

## Mitsu Koyanagi

**Global INTegration Initiative (GINTI)
Tohoku University, Japan**

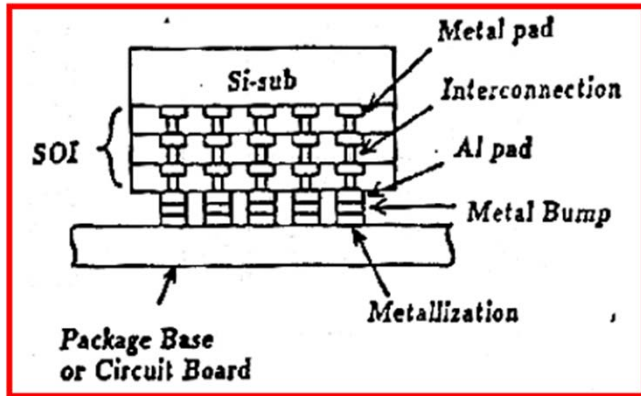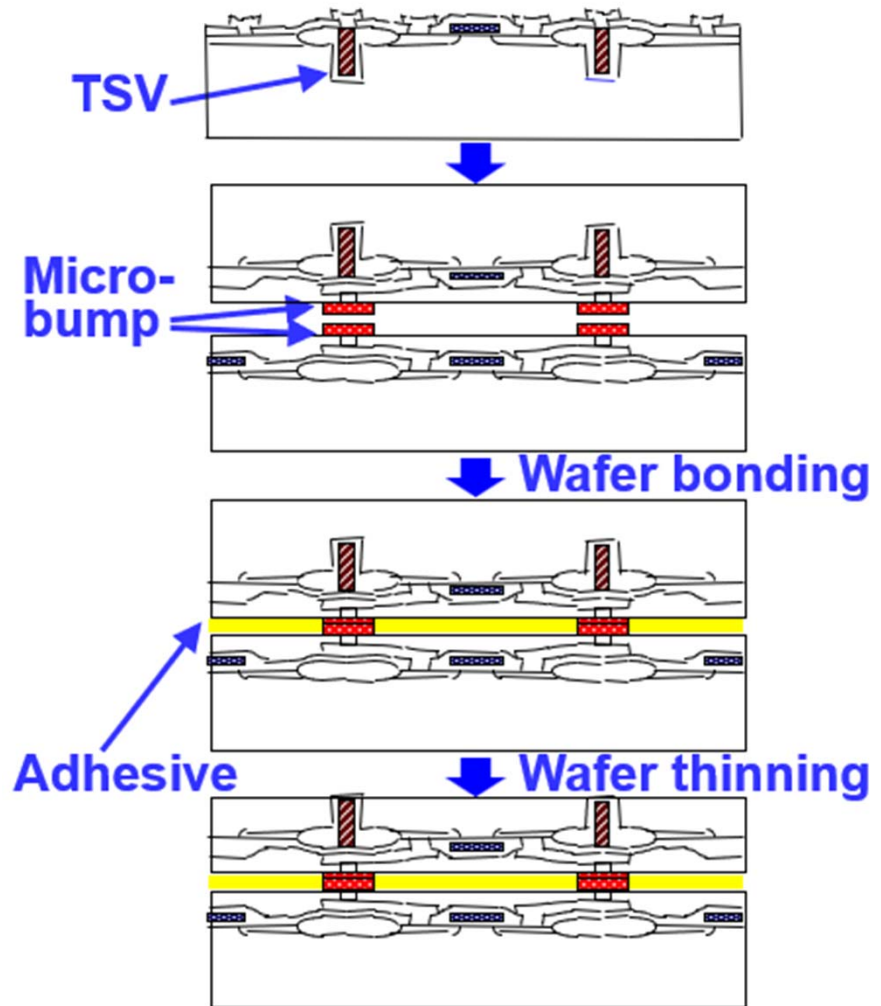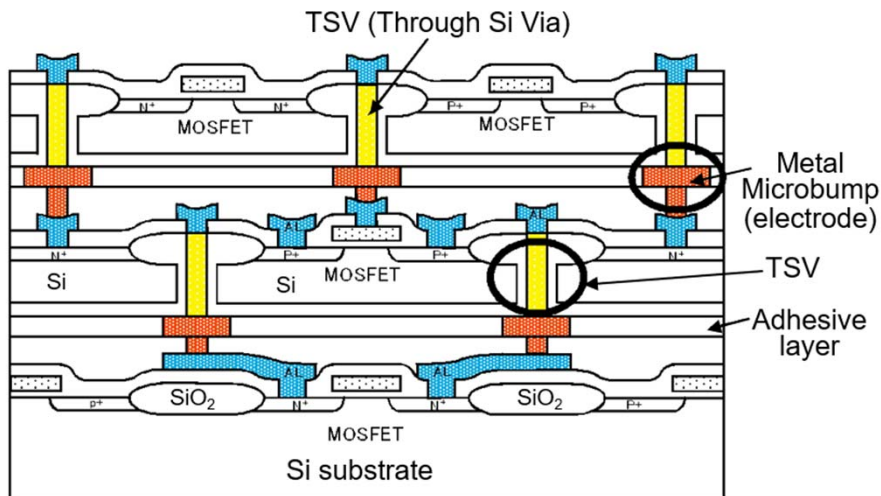**Tohoku-MicroTech. Co., Ltd, Japan**

# Outline

- ➤ Introduction
- ➤ 3D Chiplet Integration
- ➤ 3D Heterogeneous Integration
- ➤ Neuro/AI System by 3D Integration Technology
- ➤ Conclusions

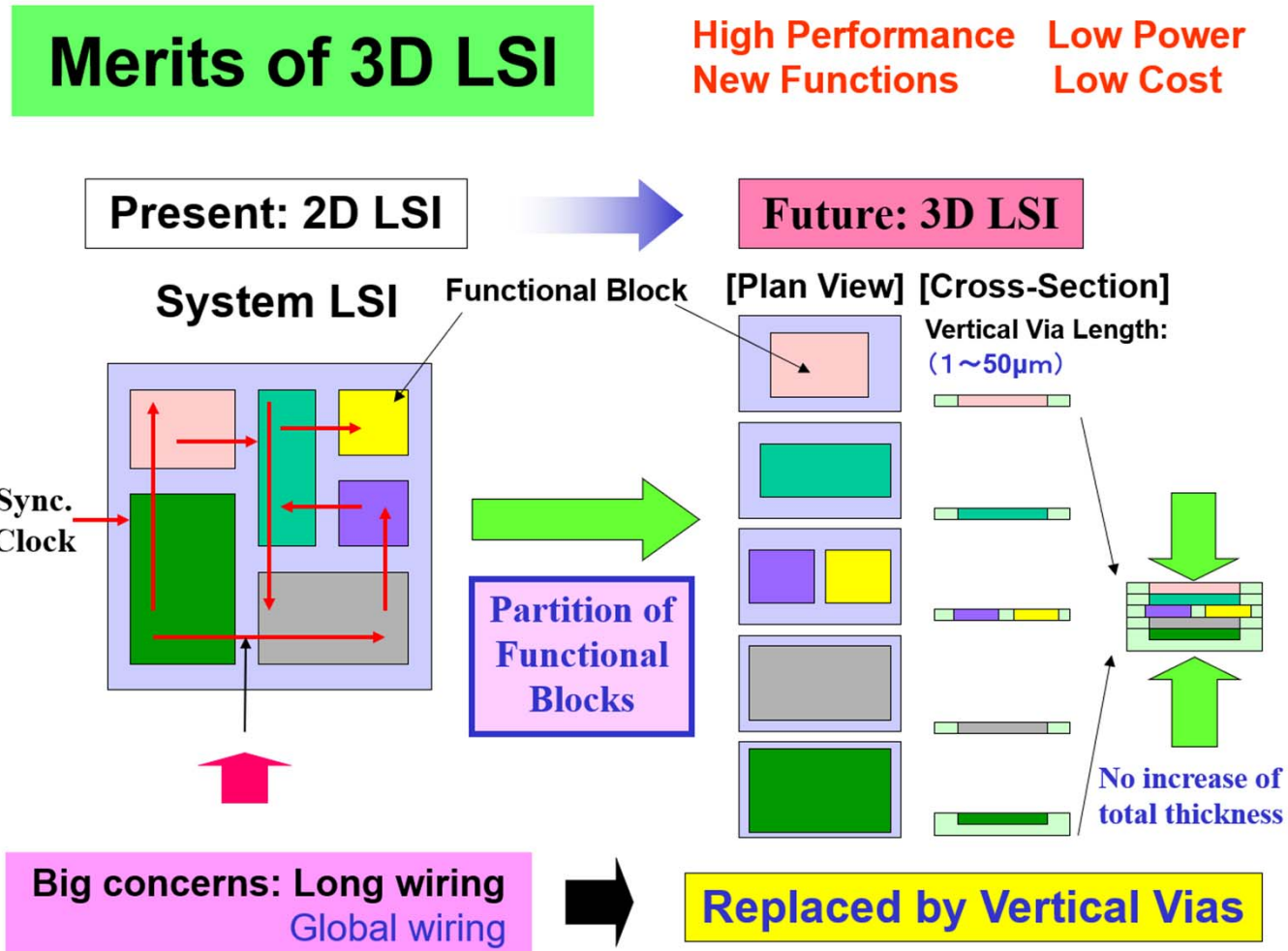# First Proposal of 3D Integration Technology in Tohoku University



The 1st proposal of 3D integration using wafer bonding in 1989

M. Koyanagi, Proc. 8th Symposium on Future Electron Devices, pp.50-60 (Oct. 1989)
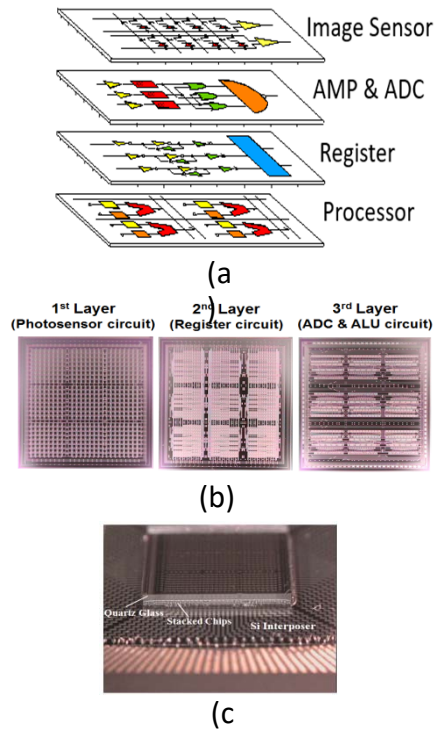
# 3D Chiplet Integration in Tohoku University



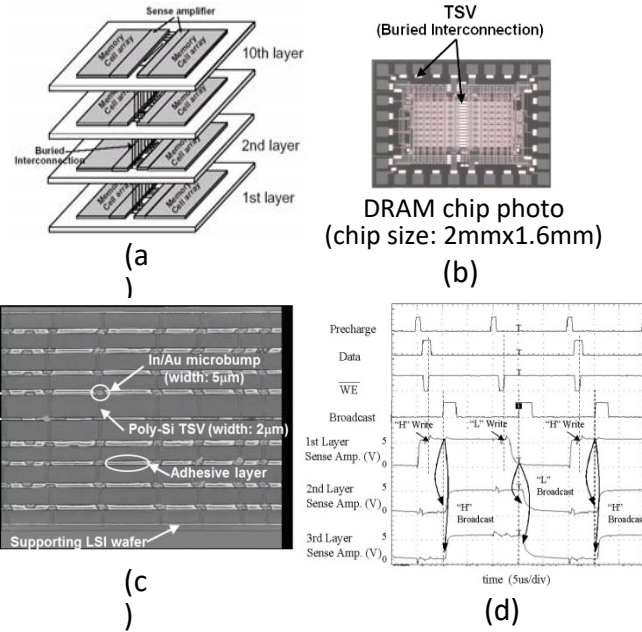M. Koyanagi, Stanford University Workshop (CIS Round Table) (2005)

# First 3D-IC Test Chips with TSVs Fabricated in Tohoku Univ.

## 3D image sensor chip



(a)



(b)



(c)

**H. Kurino, M. Koyanagi et al.**
**IEEE IEDM (1999)**

## 3D memory



(a)



(c)

DRAM chip photo
(chip size: 2mmx1.6mm)

(b)



(d)

**K. W. Lee, M. Koyanagi et al.**
**IEEE IEDM (2000)**

## 3D microprocessor chip



(a



(c
)

SRAM

Level
converter
Control logic

Processor

(b
)



(d

**T. Ono, M. Koyanagi et al., IEEE COOL Chips (2002)**

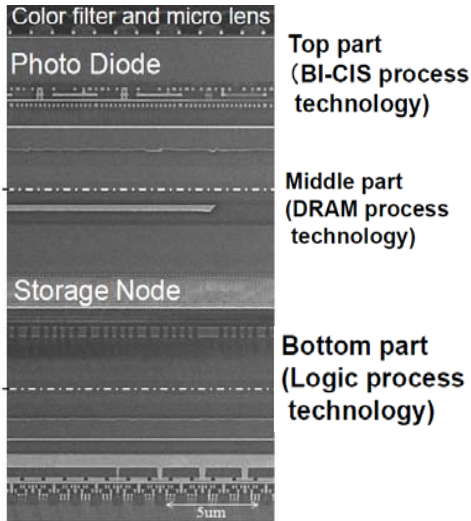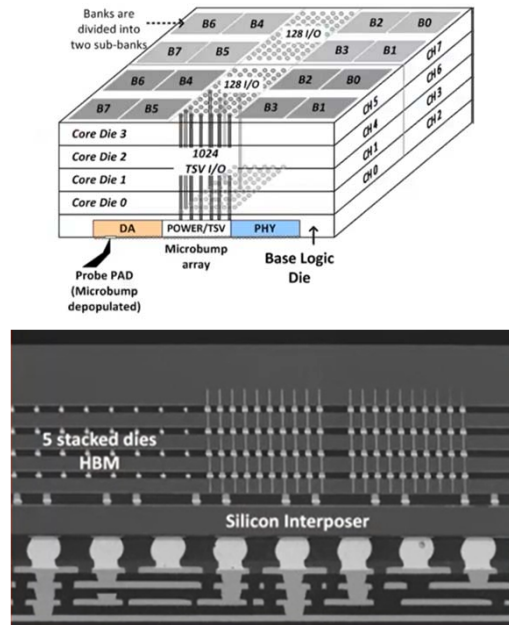# Practical Implementation of 3D-IC Chips with TSVs in Semiconductor Manufacturers
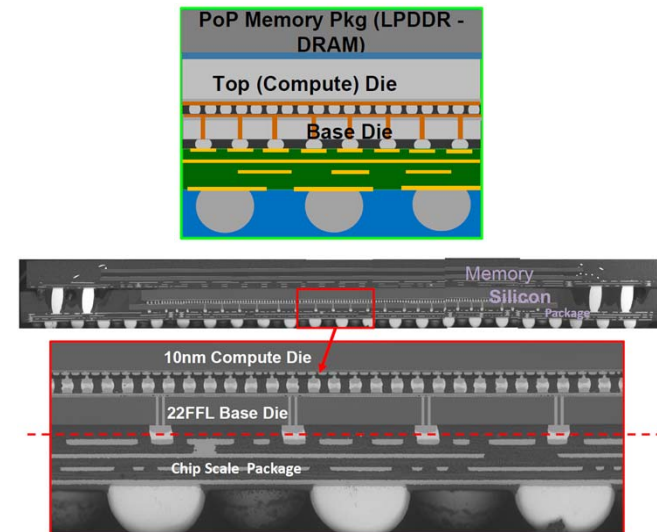
**3D image sensor chip**



**T. Haruta (Sony)**
**IEEE ISSCC (2017)**

**3D memory**



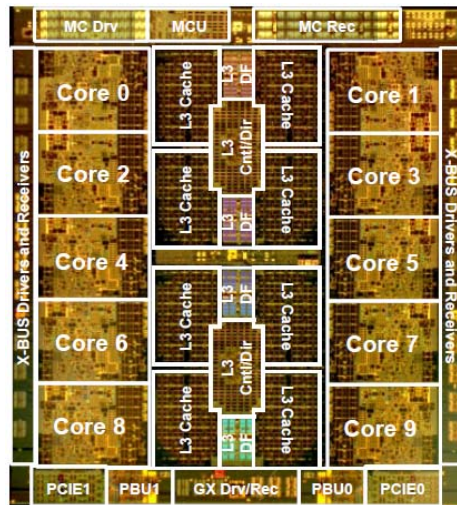**J. C. Lee et al. (SK Hynix)**
**IEEE ISSCC (2016)**

**3D microprocessor chip**



**Wilfred Gomes et al. (Intel)**
**IEEE ISSCC (2020)**

# Chiplet Integration
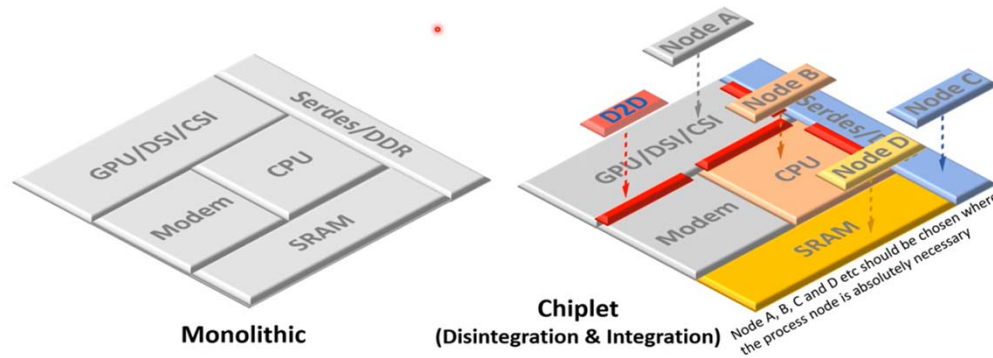
## UCIe (Universal Chiplet Interconnect Express)



Source: S-H You (IEDM2020 Short Course)



**Monolithic chip**

**IBM z14 Processor**
14nm technology node
6.1B transistors
Chip size: 696 mm2
        (ISSCC 2018)

## Chiplets & Active Interposer : Concepts

- Chiplet for:
    - Lower cost
    - Higher modularity
    - From IP-reuse to circuit-reuse

✓ Using passive interposers
    (2.5D) or organic substrate:
❖ But limitations regarding
    ➤ Chiplet connectivity (scalability),
    ➤ Less scalable function (heterogeneity)

- Active Interposer,
    the « *Smart Hub* » for :
    - Scalable System Interconnects
    - PHYs for off-chip communication
    - Power Management
    - DFT, thermal, etc

*In a mature CMOS technology (cost-performance trade off)*

D. Dutoit et al. (CEA-LETI), IEDM2020

# AMD Instinct™ MI300A Modular Chiplet Package

**I/O Die (IOD) x4**
128 Channel HBM3 Interface
256MB AMD Infinity Cache™
Infinity Fabric Network-on-Chip
2 x16 PCIe® 5 + 4th Gen Infinity Fabric™ Links
6 x16 4th Gen Infinity Fabric™ Links

**Accelerator Complex Die (XCD) x6**
228 AMD CDNA™ 3 Compute Units

**CPU Complex Die (CCD) x3**
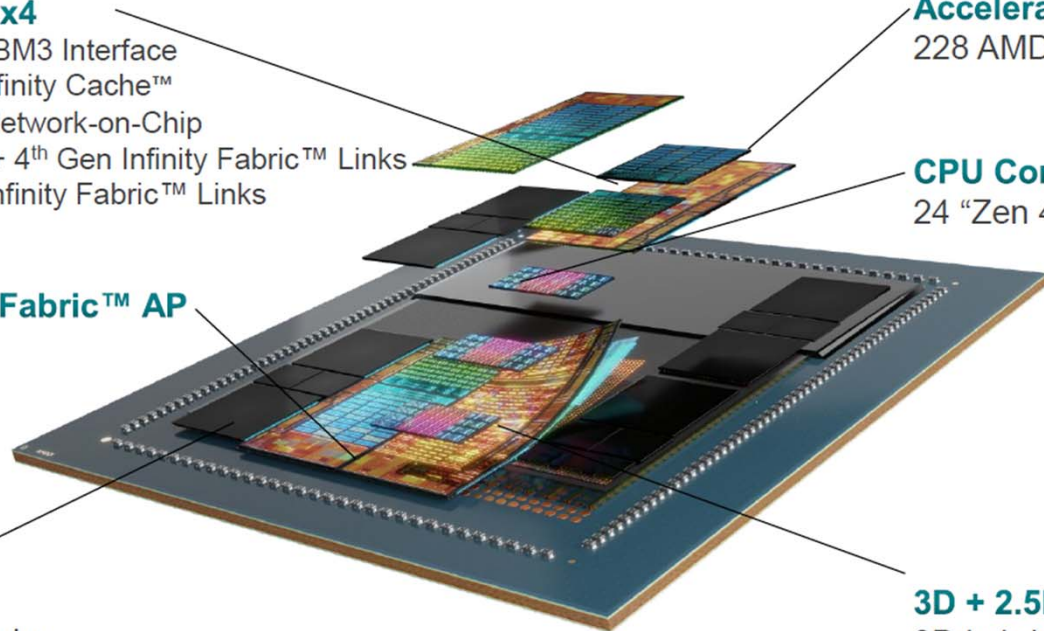24 "Zen 4" Cores [ISSCC23]

**AMD Infinity Fabric™ AP Interconnect**

**HBM3**
8 physical stacks
AMD Instinct™ MI300A: 128 GB (8-high)

**3D + 2.5D Advanced Package**
3D hybrid bonding
2.5D silicon interposer

Alan Smith et. al.. IEEE ISSCC, pp.208-209 (2024)

# Direct D2W Hybrid Bonding and Collective D2W Hybrid Bonding

## (Reconfigured W2W Hybrid Bonding)

# Process Flow of D2W Hybrid Bonding



(a) Top Wafer Process

CVD — CMP — Sawing — Plasma — Hydration

(c) D2W Bonding

D2W Bonding — Anneal

CVD — CMP — Plasma — Hydration

(b) Base wafer process

S. Lee et al., IEEE ECTC, pp.1085-1089 (2022)

# Hybrid Bonding in Tohoku University



**Alignment accuracy ~0.1μm**

Top wafer

RDL

Cu pad

Cu pad

RDL

Bottom wafer

5.0kV x7.00k                    4.29um

**Cu post size 5μm; pitch 10μm**

(a) W2W

M. Murugesan, M. Koyanagi, T. Fukushima, IEEE ECTC (2022)



Top chip

Cu pad

RDL

RDL

Bottom wafer

**Cu post size 1.5μm; pitch 2.5μm**

(b) C2W

To be published

# Cu Grain Morphology _ SEM

| | (250℃*1hr) | (300℃*1hr) | (350℃*1hr) |
|---|---|---|---|
| **Plating Chemistry A** |  |  |  |
| **Plating Chemistry B** |  |  |  |

Grain size ~2μm
(tiny and extremely random oriented)

Grain size >10μm
(very large and relatively oriented)

# Cu Grain Crystallographic Orientation _ EBSD



M. Murugesan, M. Koyanagi, T. Fukushima, IEEE ECTC (2022)

# 3D Chiplet Integration on Wafer by D2W Hybrid Bonding



(a) Low magnification image (b) Tilted SEM image, and (c) Cross-sectional SEM image

S. Lee et al., IEEE ECTC, pp.1085-1089 (2022)



HBM with 8 memory layers

Jaesik Lee, IEEE IEDM, SC2.2 (2023)

# Pick-and-Place vs Self-Assembly

● *Traditional pick-and-place assembly*

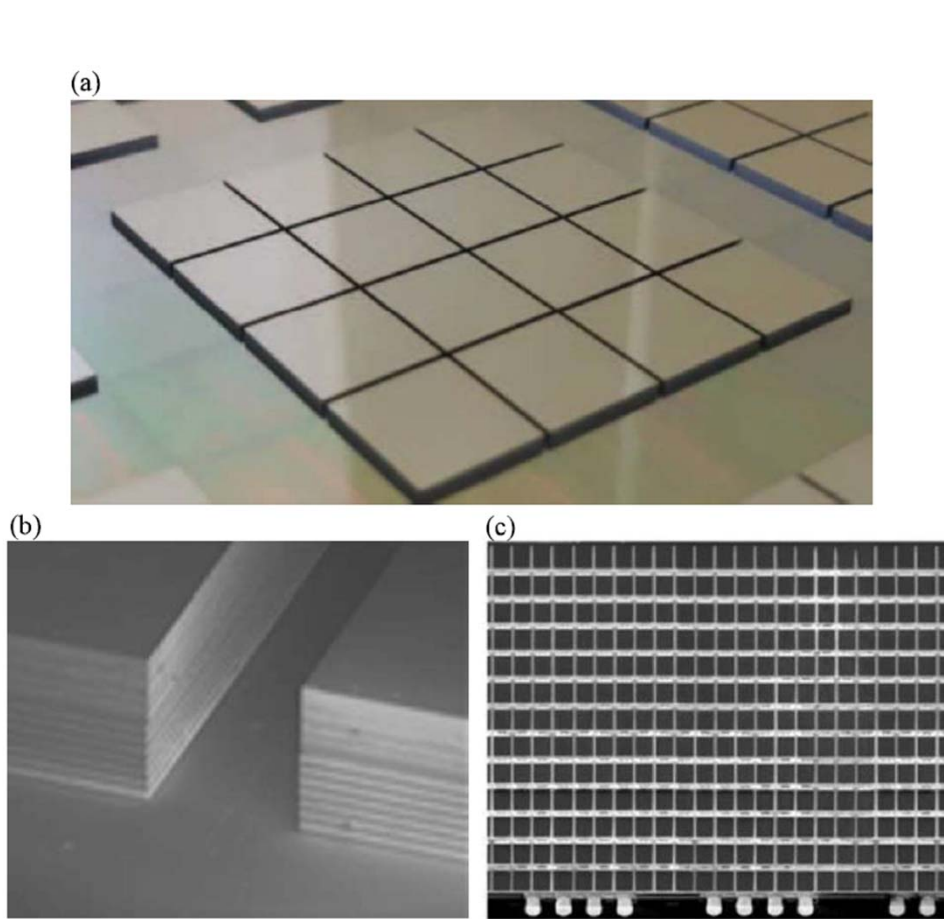Single chip pick-up tool

Tradeoff problem

Chip tray    Wafer

Alignment accuracy and time with commercial bonders



Trade-off

Target

Accuracy 3σ [um]

Time [s]

● *Self-assembly using liquid surface tension*

High accuracy & High throughput

Mult-chip pick-up tool

Liquid

Known Good Dies (KGDs)

# Simultaneous Bonding of Many Dies with Different Size by Self-Assembly



Liquid droplet on hydrophilic area

Hydrophobic area

5mm

* This movie is real-time playing.

# High-Speed Water Droplet Spray for Self-Assembly

# Water Droplets Supplied on Small Hydrophilic Area

# Photo of Various-size Self-Assembled Chips on 8-inch Wafer Prepared by Hybrid Self-Assembly



$4 \times 9$ mm$^2$ chip

$5 \times 5$ mm$^2$ chip

$10 \times 10$ mm$^2$ chip

$3 \times 3$ mm$^2$ chip

200 mm

# Photo of μLED Array Prepared by Hybrid Self-Assembly



Before self-assembly
on waffle tray

After self-assembly
on tape



Self-Assembled Micro-LED chips
(75umx125um)

# Combined Process Sequence of Self-Assembly and Hybrid Bonding
## (SA-Hybrid Bonding)

Wafer    Chips

① CMP-treated wafer and chip fabrication

Plasma

Wafer    Chips

Stage    Stage

②. Wafer and/or chip surface activation

Water supply

Water droplets

Stage

③Self-assembly with water for multichip-to-wafer bonding

Multichip-to-wafer simultaneous bonding

Stage

④Hybrid bonding

TCB conditions:
10 N/chip
200-300°C / 2h

# Results of SA-Hybrid Bonding



① ② ③

- Using this method, we have currently achieved an assembly of 6 layers. In the future, we will continue to explore the best conditions to achieve a structure of more than 12 layers.

- The current average assembly accuracy has reached a level of less than 500um, but the data is still scattered, and efforts will be made to reduce the data scatter in the future.

- The current thermal bonding can still see the bonding interface between Cu. In the future, we will continue to explore the impact of liquid on bonding and the optimal bonding conditions.

## Heterogeneous Integration

➢ Device Level

➢ Architecture/System Level

22

# Device Level Heterogeneous Integration

## 3D Heterogeneous Integration Technology in Tohoku Univ.

On-chip waveguide

TSV

Microbump

Super chip

**3D Heterogeneous Integration**

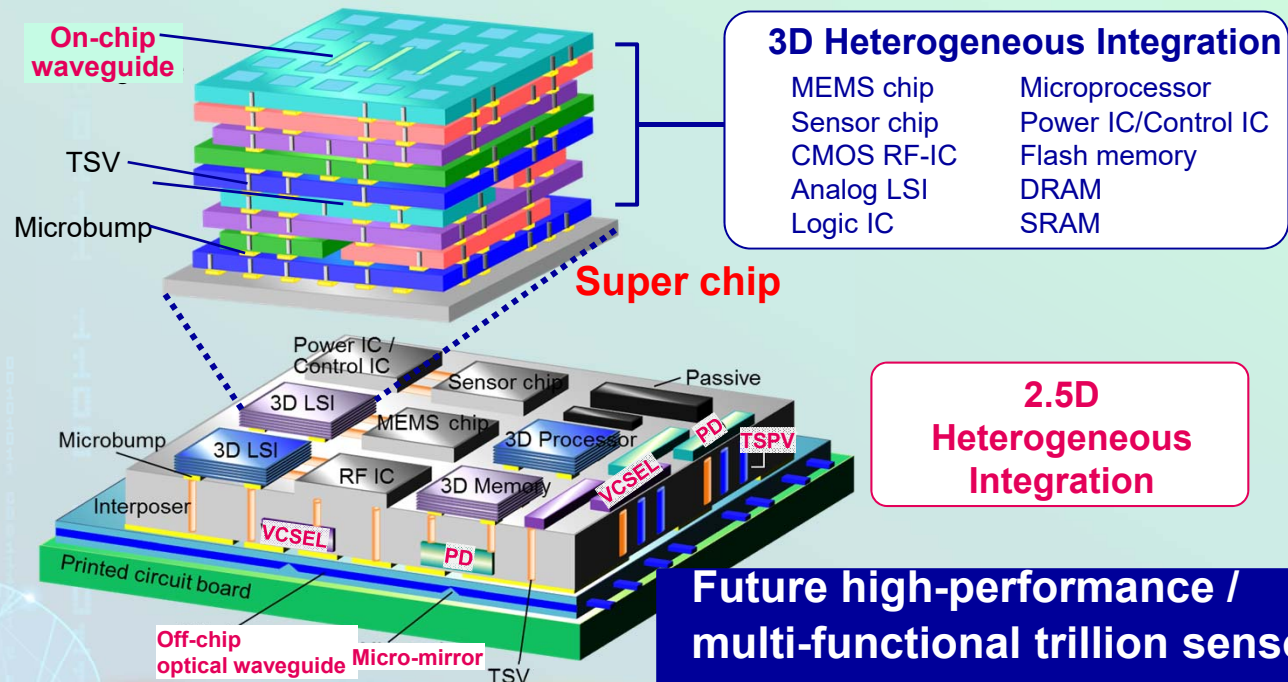| | |
|---|---|
| MEMS chip | Microprocessor |
| Sensor chip | Power IC/Control IC |
| CMOS RF-IC | Flash memory |
| Analog LSI | DRAM |
| Logic IC | SRAM |

Power IC / Control IC

3D LSI

Sensor chip

Passive

MEMS chip

3D Processor

PD

TSPV

3D LSI

RF IC

3D Memory

VCSEL

Microbump

VCSEL

Interposer

PD

Printed circuit board

Off-chip optical waveguide

Micro-mirror

TSV

**2.5D Heterogeneous Integration**

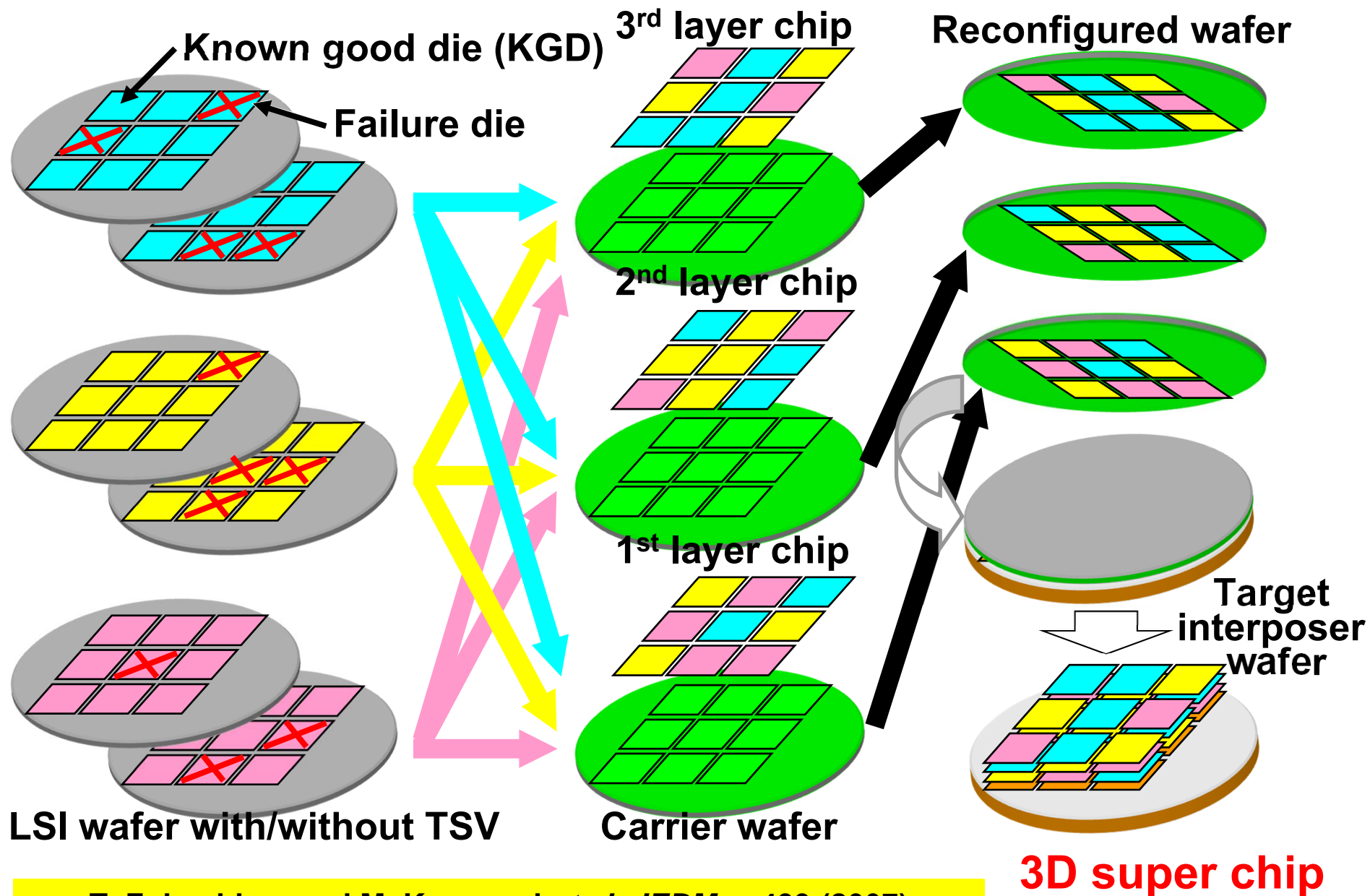**Future high-performance / multi-functional trillion sensors and IoT require high-speed data transfer & low power operation.**

Different devices

Different size chip

Different materials

T. Fukushima, M. Koyanagi et. Al., IEEE IEDM, p.359 (2005)

K-W Lee, M. Koyanagi et. al., IEEE IEDM, p.531 (2009)

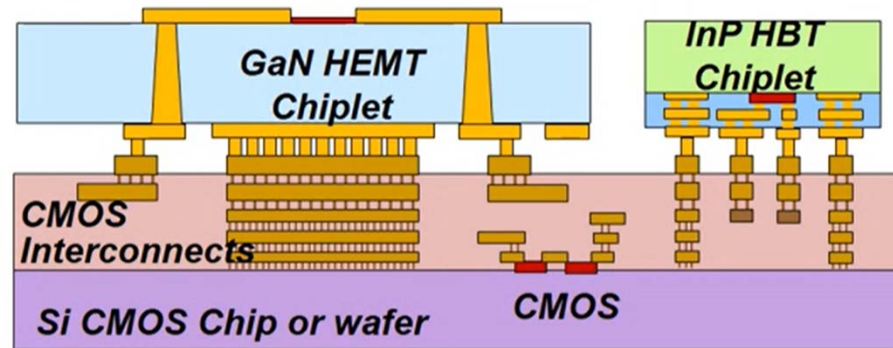# New Reconfigured Wafer-to-Wafer 3D Integration



T. Fukushima and M. Koyanagi *et al.*, *IEDM*, p.439 (2007)

# Various Kinds of Heterogeneous Integration

- **Heterogeneous Integration with Non-Si devices**
- **Heterogeneous Integration with Sensor/MEMS**
- **Heterogeneous Integration with Photonics/Optics**
- **Heterogeneous Integration with Bionics**
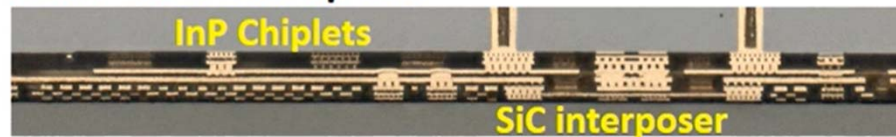- **etc.**

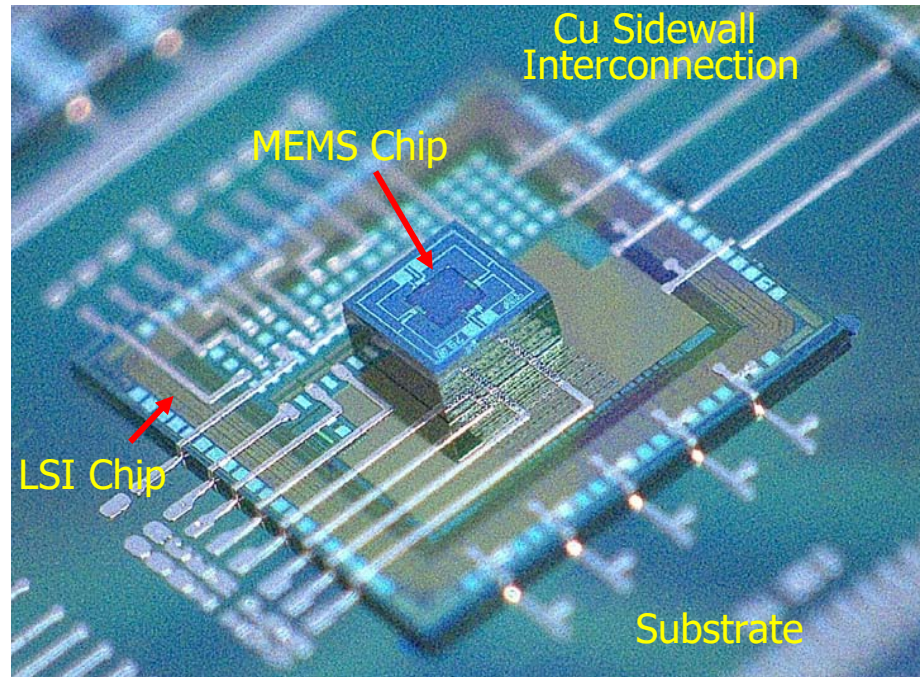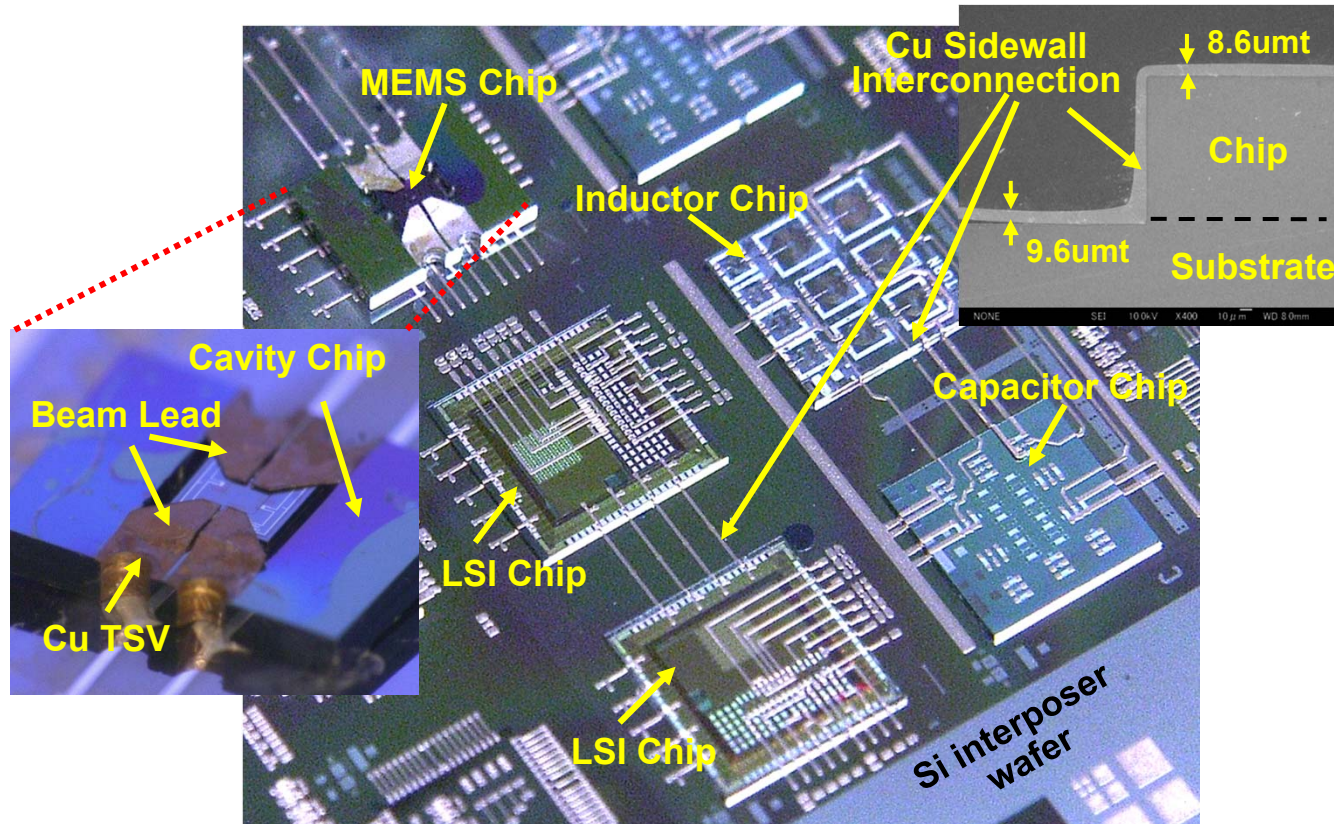# Heterogeneous Integration with Compound Semiconductor Chip (DAHI)

# 3D Heterogeneous Integration with MEMS Using Self-Assembly



K-W Lee, M. Koyanagi et al., 3D-IC, Sept. 28, 2009

# 2.5D/3D Heterogeneous Integrations of CMOS, MEMS and Passive Device Chips on Si Substrate

# Electronic-Photonic Systems-on-Chip for Compute, Communications and Sensing

## 45SPCLO process



Rakowski *et al* OFC 2020

- Same transistors as in 45nm SOI
- Number of features optimized for photonics
  - Ge photodetectors, Si dopings, Si partial etch, SiN, V-groove couplers etc.

Vladimir Stojanović, IEEE ISSCC Forum 6.8 (2024)

# Future Systems-In-Package with Optical I/O



PLATFORM ON A PACKAGE

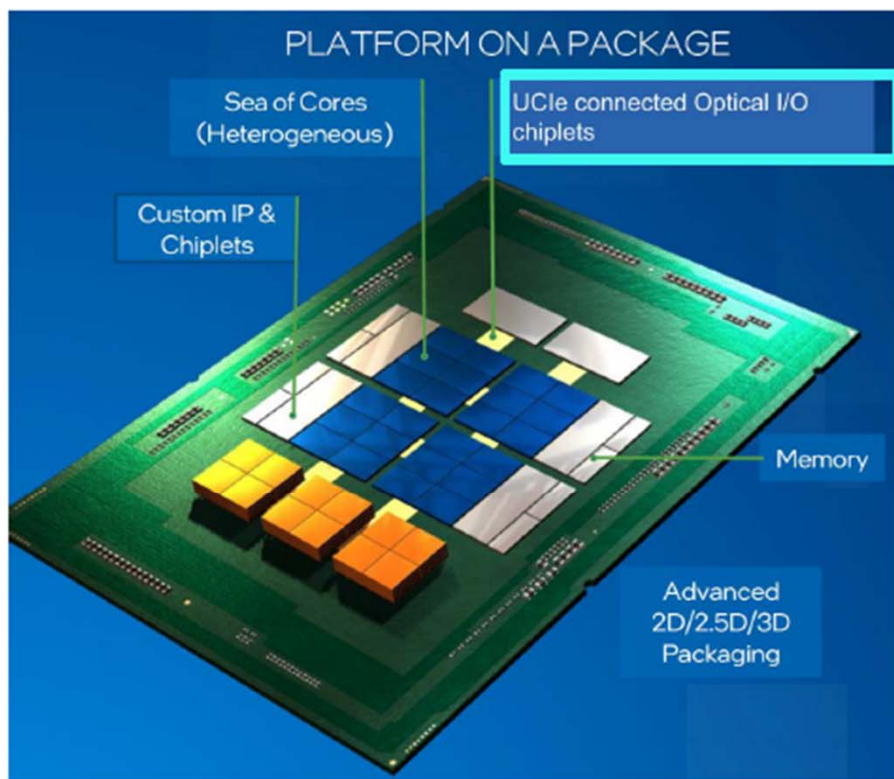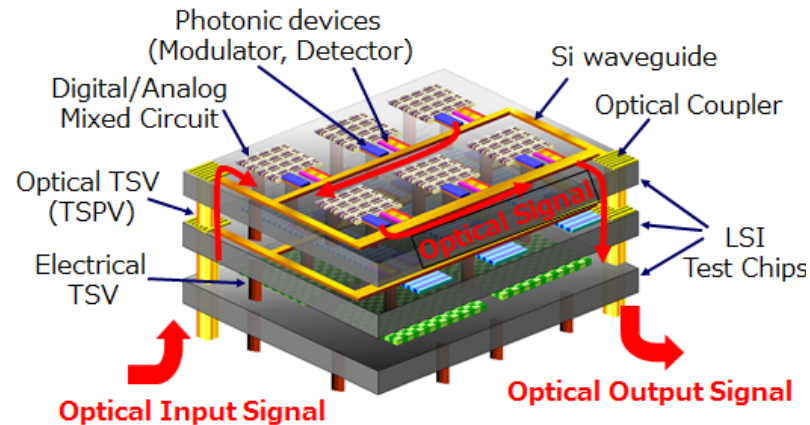- Sea of Cores (Heterogeneous)
- UCIe connected Optical I/O chiplets
- Custom IP & Chiplets
- Memory
- Advanced 2D/2.5D/3D Packaging

| Gen | Electrical I/F (Advanced Package) | | | | Optical I/F (CW-WDM) | | | Optical Chiplet BW (Tx+Rx) | Off-package IO BW (4-8 chiplets per package) |
|---|---|---|---|---|---|---|---|---|---|
| | I/F | Modules | Tx / Rx IOs | Data Rate [Gbps/IO] | Ports | λs / Port | Data Rate [Gbps/λ] | | |
| 1 | AIB | 24 | 20 / 20 | 2 | 8 | 8 | 16 | 2 Tbps | 8-16 Tbps |
| 2 | AIB | 16 | 80 / 80 | 2 | 8 | 8 | 32 | 4 Tbps | 16-32 Tbps |

- Gen 1 and Gen 2 already built and hardware validated

- 16-32 Tbps off-socket optical I/O bandwidth possible today
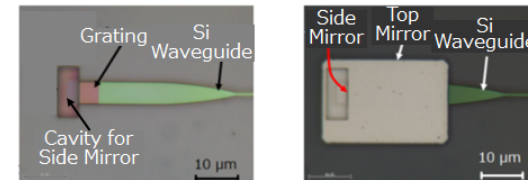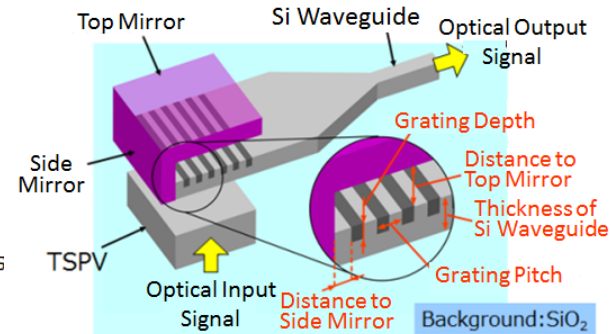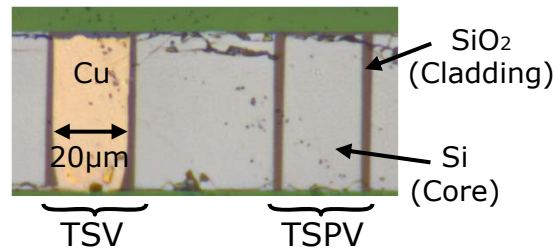
(Source: Wade *et al* HotChips 2023)

Vladimir Stojanović, IEEE ISSCC Forum 6.8 (2024)

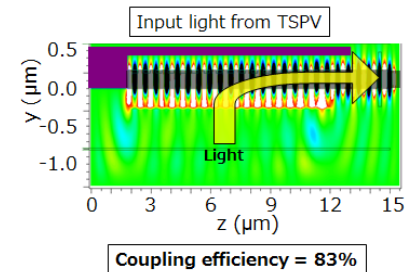# 3D Heterogeneous Integration with Photonics
## (Photonic 3D Integration)



Vertical optical interconnection (TSPV)
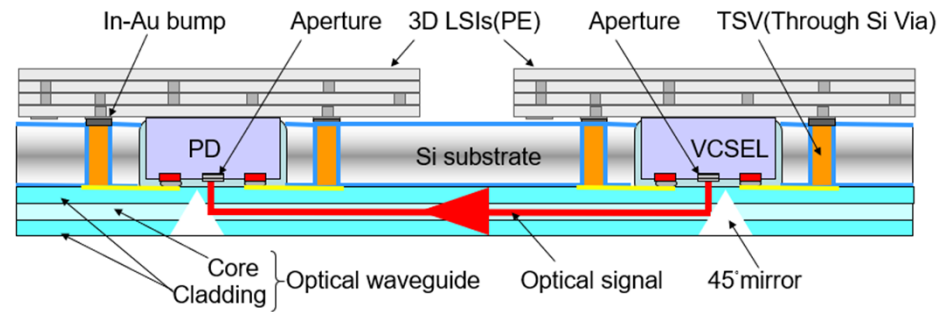
- TSPV (Through Si photonic via)

K-W Lee, M. Koyanagi et. al., IEEE Trans. on Electron Devices, vol.58, p.748 (2011)

# Photonic 2.5D/3D Heterogeneous Integration
## (Optical Interposer Embedded with VCSEL/PD Chips)



In-Au bump    Aperture    3D LSIs(PE)    Aperture    TSV(Through Si Via)

PD    Si substrate    VCSEL

Core
Cladding } Optical waveguide    Optical signal    45°mirror

(Photograph of optical interposer captured from back surface)

(Back Surface)

VCSEL    Waveguide    PD

Metal Land
VCSEL Chip
Micro-Mirror
Cavity    Beam Lead
TSV

(Back-Surface)
PD Chip    Beam Lead
Micro-Mirror
Waveguide    Cavity
TSV

A.Noriki, M. Koyanagi et. al., Jpn. Jour. Appl. Phy. Vol. 48, p.C113-1 (2009)

MEC System Module with Optical Interconnection by Heterogeneous Chiplet Integration
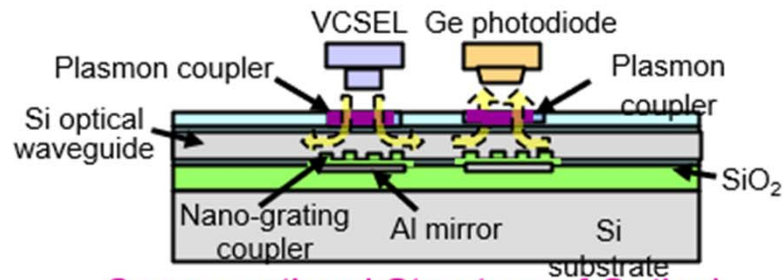
Shared memory (HBM)

Fan-out wiring

3D Processor (face-to-face))

3D DRAM

Si waveguide

VCSEL array

Ge photodiode array

Si interposer

# Fabrication Flow of Optical Interconnection with Grating Coupler and Plasmon Coupler



Cross-sectional Structure of Optical Interconnection with Grating Coupler and Plasmon Coupler
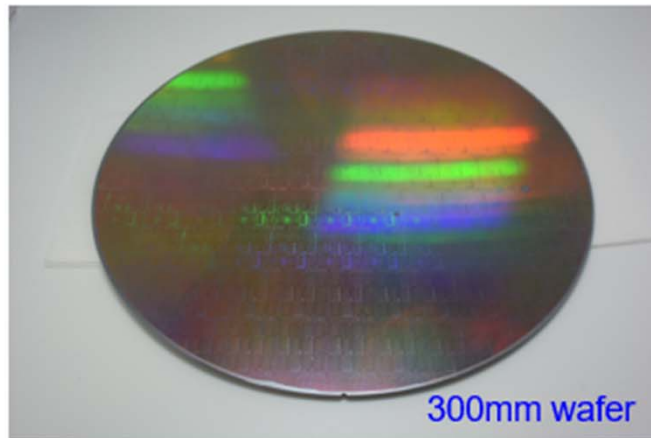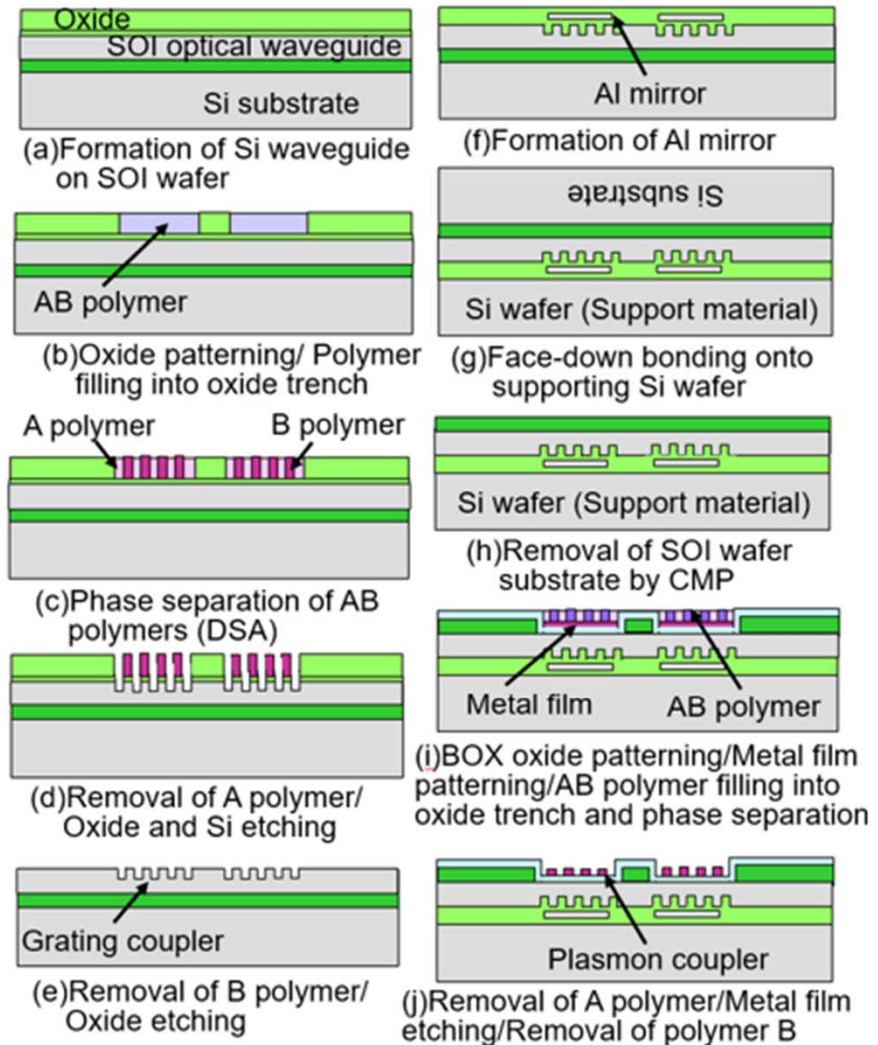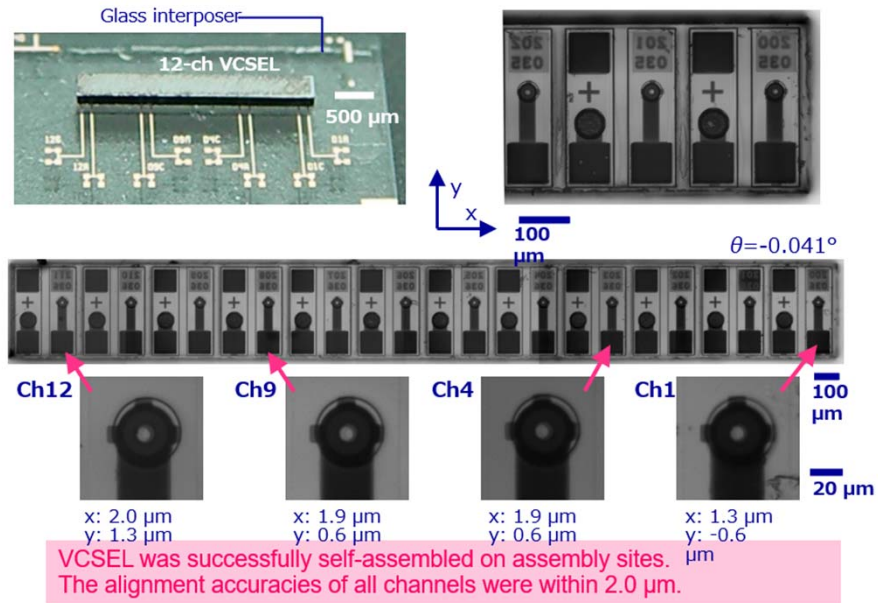
Photo after transferring Si optical waveguide patterns onto a Si interposer wafer

300mm wafer

(a)Formation of Si waveguide on SOI wafer

(b)Oxide patterning/ Polymer filling into oxide trench

(c)Phase separation of AB polymers (DSA)

(d)Removal of A polymer/ Oxide and Si etching

(e)Removal of B polymer/ Oxide etching

(f)Formation of Al mirror

(g)Face-down bonding onto supporting Si wafer

(h)Removal of SOI wafer substrate by CMP

(i)BOX oxide patterning/Metal film patterning/AB polymer filling into oxide trench and phase separation

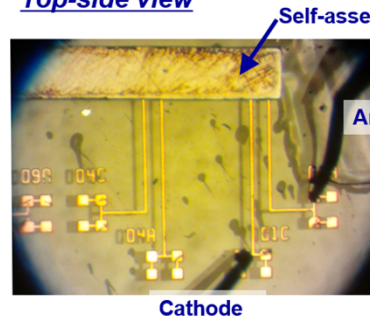(j)Removal of A polymer/Metal film etching/Removal of polymer B

# VCSEL Chiplet Integration on Glass Interposer by Self-Assembly

## Self-Assembly of 12-ch VCSEL Chip on Glass Interposer



Glass interposer

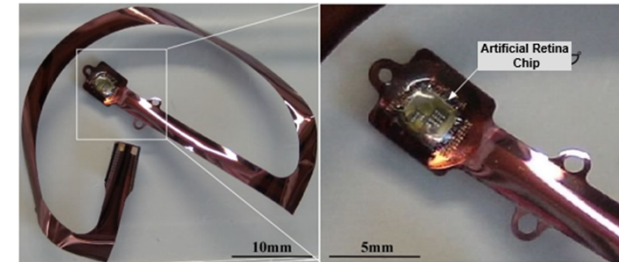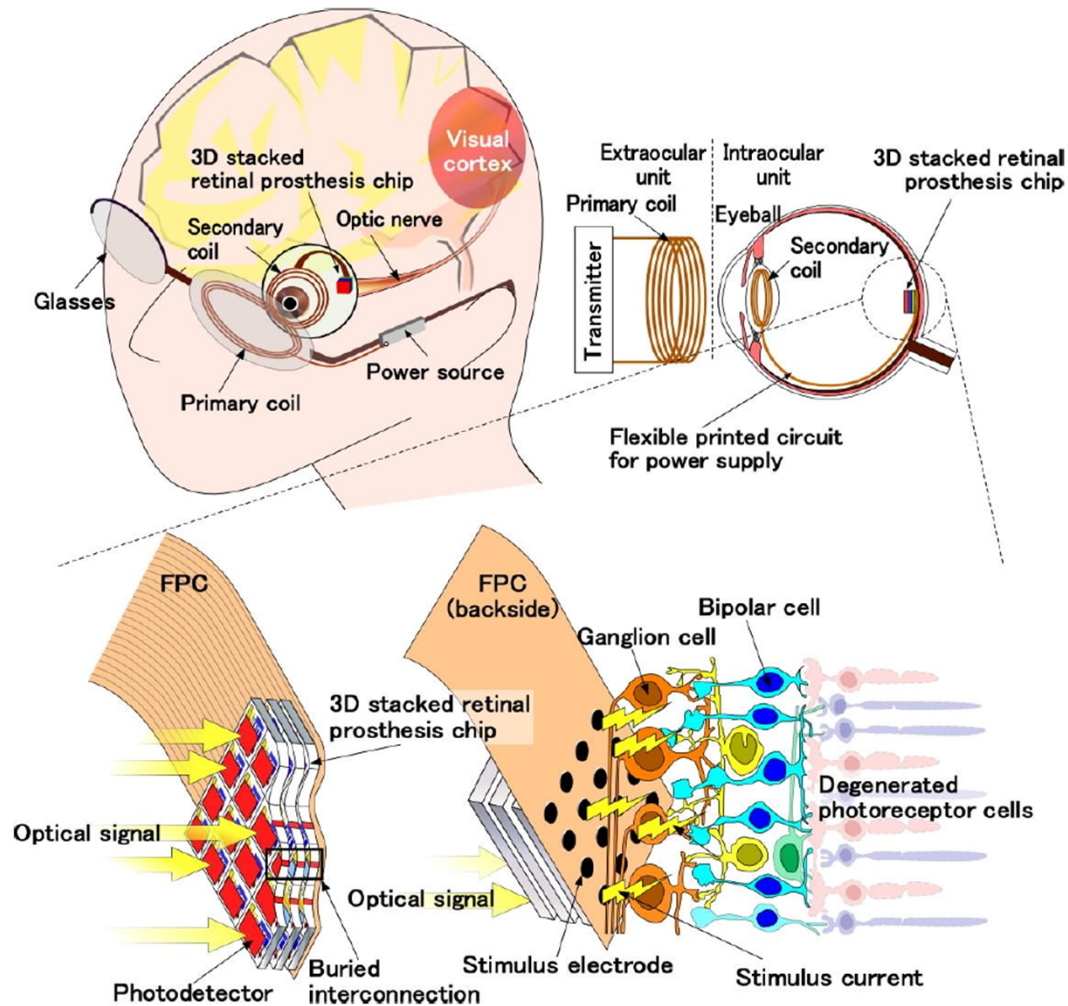12-ch VCSEL

500 μm

y
x

100 μm

$\theta = -0.041°$

100 μm

Ch12     Ch9     Ch4     Ch1

20 μm

x: 2.0 μm
y: 1.3 μm

x: 1.9 μm
y: 0.6 μm

x: 1.9 μm
y: 0.6 μm

x: 1.3 μm
y: -0.6 μm

VCSEL was successfully self-assembled on assembly sites.
The alignment accuracies of all channels were within 2.0 μm.

*Top-side view*

Self-assembled 12ch-VCSEL

Anode

Cathode

*Bottom-side view*

Cathode    Anode

Light emission
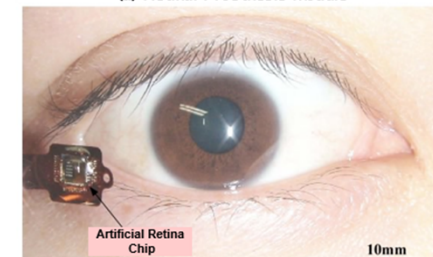
Self-assembled 12-ch VCSEL

| Applied voltage | : 2 V |
|---|---|
| Measured current value | : 5.0 mA |

# 3D Heterogeneous Integration with Bionics
## (Bionic 3D Integration)

### 3D-Retina Chip Implantation into Human Eye (Retinal Prothesis)



T. Tanaka, M. Koyanagi et. al.,
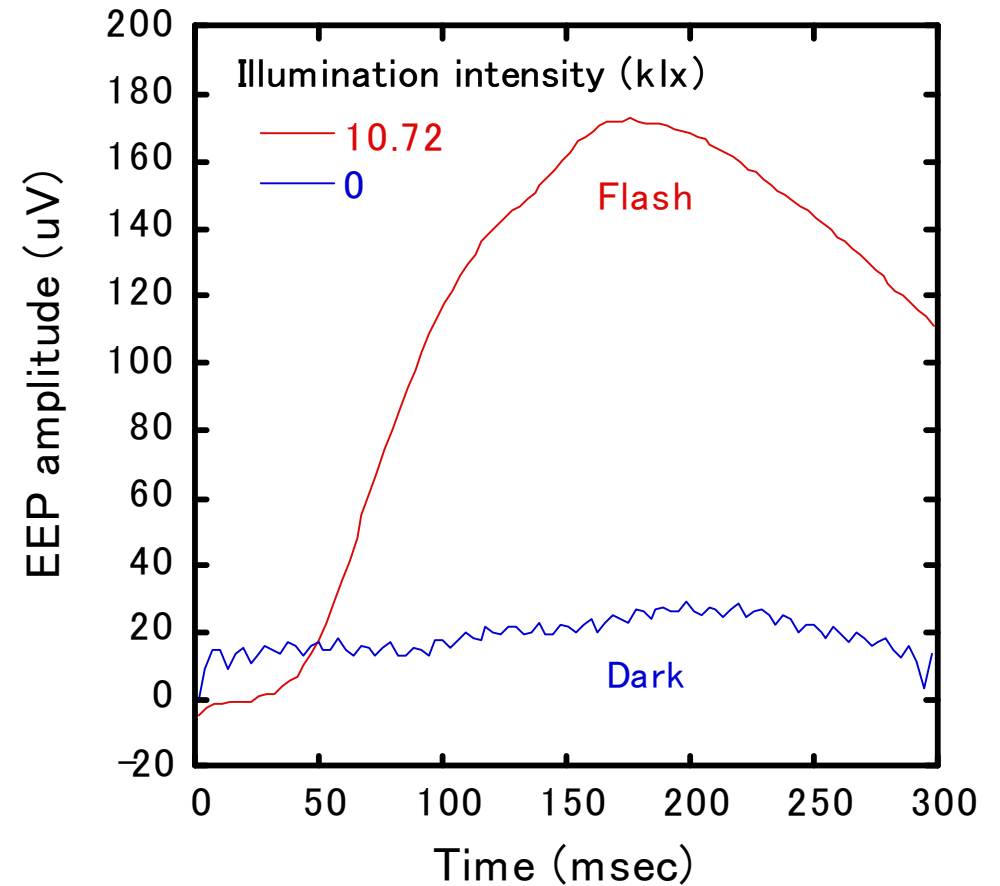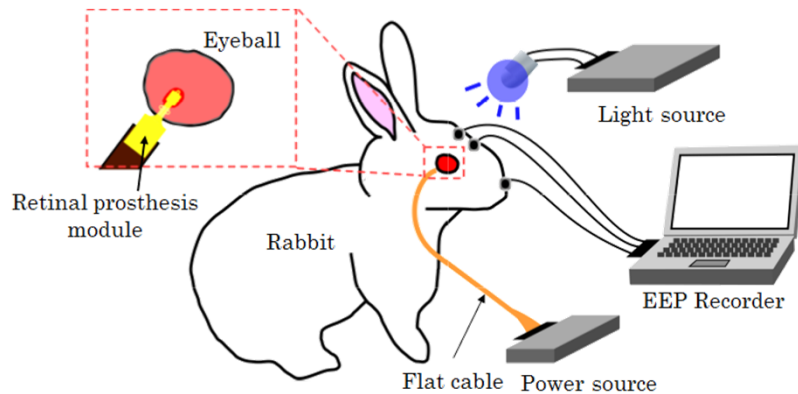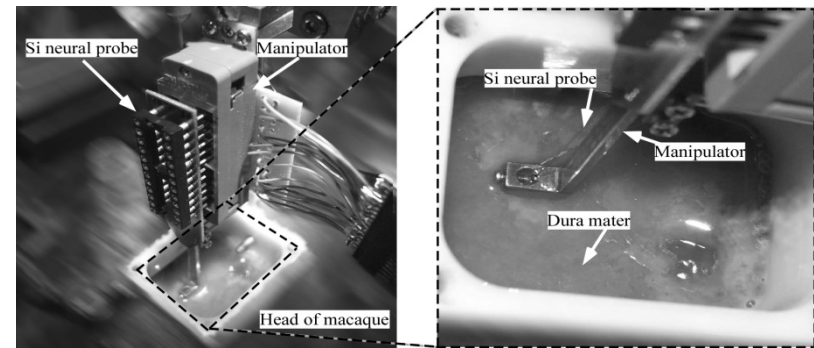IEEE IEDM, p.1015 (2007)
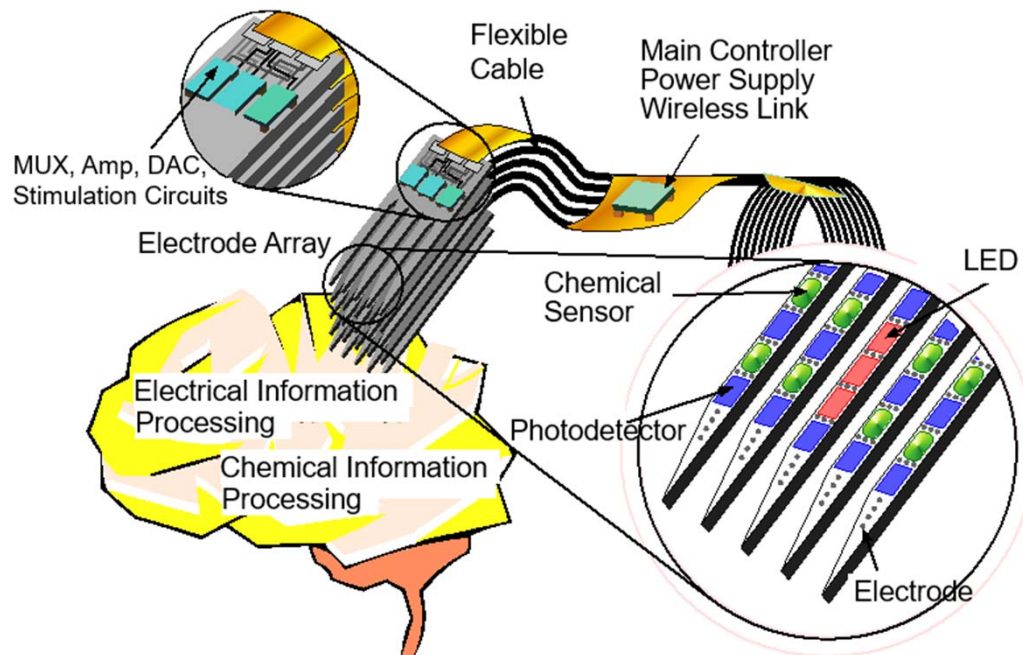
# EEP Waveforms with/without Flashlight
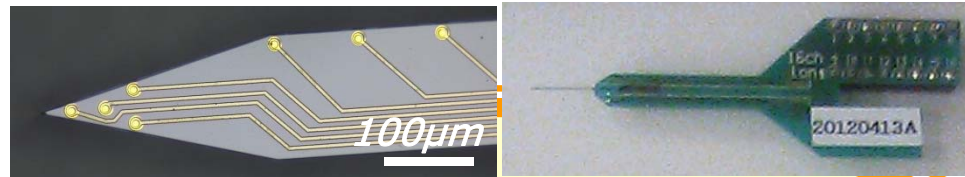
Animal experiment setup

# 3D Heterogeneous Integration with Bionics
## (Bionic 3D Integration)
### Brain-Machine Interface (BMI) and Intelligent Si Neural Probe with Multi-electrodes and Sensors
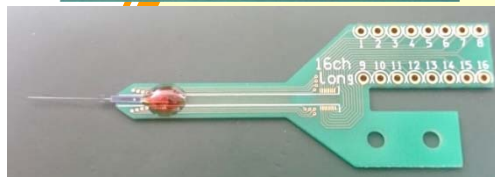


Recording of Neuron Potential in a Brain Using Si Neural Probe

S. Kanno, T. Tanaka, M. Koyanagi et. al., Jpn. Jour. Appl. Phy., Vol. 48, p. C189-1 (2009)

Si neural probe mounted on PCB

Piezoresistive force sensor

400μm

Si neural probe with piezoresistive force sensor

Tohoku Univ. Intelligent Neural Probe family

(Prof. Tetsu Tanaka)

100μm

Front

Back

600μm

4-shunk double-sided Si neural probe

Optical waveguide

1mm

Pillar-type electrode array (10x10)

Recording site

Cross-section

20μm

Outlet

100μm

Si neural probe with optical waveguides

Si neural probe with microfluidic channel

# Neural stimulation

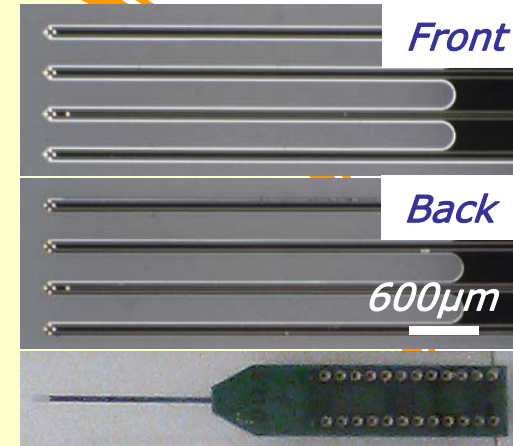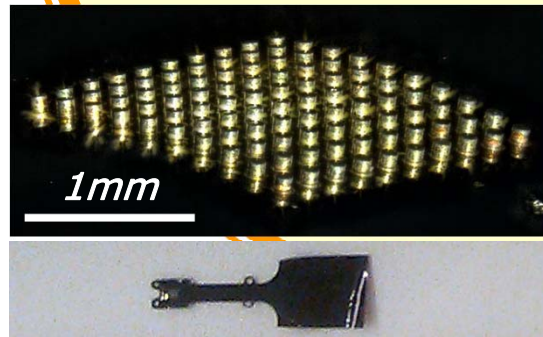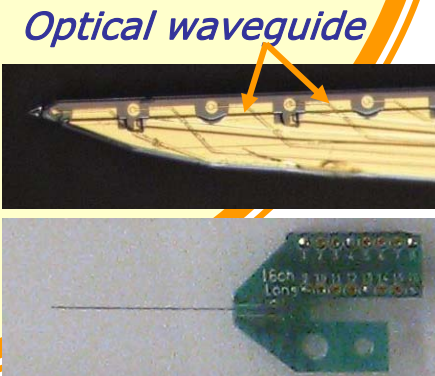| | Electrical | Chemical | Optical |
|---|---|---|---|
| Stimulus speed | Fast ☺ | Slow | Fast ☺ |
| Neural activities | Excitation | Excitation Inhibition ☺ | Excitation Inhibition ☺ |
| Cell selectivity | No | No | Yes ☺ |

■Optogenetics

Gene transfer of a protein molecule that responds to light of a specific wavelength enables the control of neuronal excitation and inhibition by light.

■Neural probe w/ light control func.

Optical fiber/Optical waveguide/µLED



Optic Fiber
100 µm
光刺激点
金属カバー
100µm
µLED

Blue
Yellow

Na⁺ Na⁺ Ca²⁺ Na⁺ Na⁺
Cl⁻ Cl⁻ Cl⁻

ChR
HR

K⁺ H⁺ Na⁺

photoactivatable proteins

O. Yizhar *et al., Neuron,* **71** (2011), 9-34

Realization of neural stimulation with high spatial and temporal resolution and cell selectivity.

By Courtesy of Prof. Tetsu Tanaka (Tohoku University)

# 3D Heterogeneous Integration with Bionics
# (Bionic 3D Integration)
## Review Bi-directional Brain Interface
## from Risk Management Perspective



Opri et al, Chronic embedded cortico-thalamic closed-loop deep brain stimulation for the treatment of essential tremor, Science Translational Medicine, 2020

Tim Denison, IEEE ISSCC, Forum 4.1 (2024)

# Architecture/System Level Heterogeneous Integration

## Examples of Heterogeneous Integration

### High-Performance Compute

7nm CPU
7nm CPU
7nm CPU
7nm CPU
14nm I/O die
7nm CPU
7nm CPU
7nm CPU
7nm CPU

Naffziger *et al.*, AMD [44]

47 tiles
5 process nodes
HBM
HBM
HBM
Xe Link

Gomes *et al.*, Intel [45]

### Automotive Microcontroller

16nm MCU
55nm EEPROM 5V Analog I/Os
NVM

Loke *et al.*, NXP [46]

Figures not drawn to scale

Alvin Loke, IEEE ISSCC, Forum 3.1 (2024)

# Architecture/System Level Heterogeneous Integration



## Homogeneous Integration

- ☐ Power, performance, area and chip cost
- ☐ Cross-IP data link and latency
- ☐ Thermal- and IR-aware die floorplanning

## Heterogeneous Integration

- ☐ System cost and time-to-market
- ☐ Cross-Die data link and latency
- ☐ Multiphysics-aware chip integration

Lawrence Loh, IEEE ISSCC, Plenary (2020)

# Evolution of IC from Device Integration to System Integration



Kevin Zhang, IEEE ISSCC, Plenary 1.1 (2024)

# Key Technologies for AI
## 3D Integration, Chiplet Integration, Heterogeneous Integration



Lip-Bu Tan, IEEE ISSCC, Plenary 1.4 (2024)

# LLM Computations (Training vs. Inference)



Joo-Young Kim, IEEE ISSCC Forum 2.5 (2024)

# Growth in application complexity



S. Bianco et al., "Benchmark Analysis of Representative Deep Neural Network Architectures," IEEE Access, 2018

Machine learning hardware:  considerations
and accelerator approaches

# Requirements for AI System and Technologies

Inference  Learning (Minimal) ⇅  Edge Devices  ← Mode  Data →  Servers  Learning ↓

Inference  Learning (Local/Personal) ⇅  Extreme Edge Devices  Tiny ML
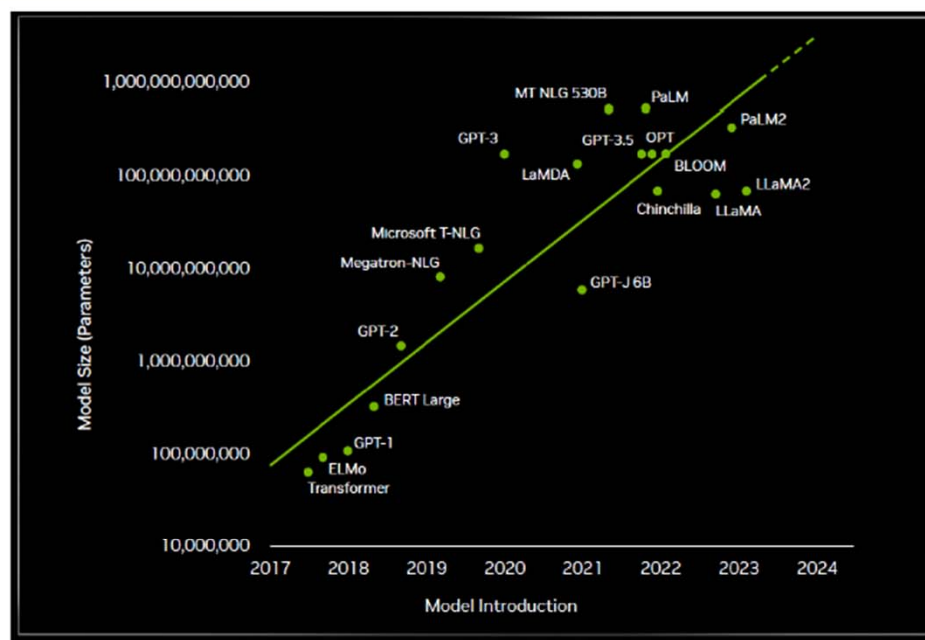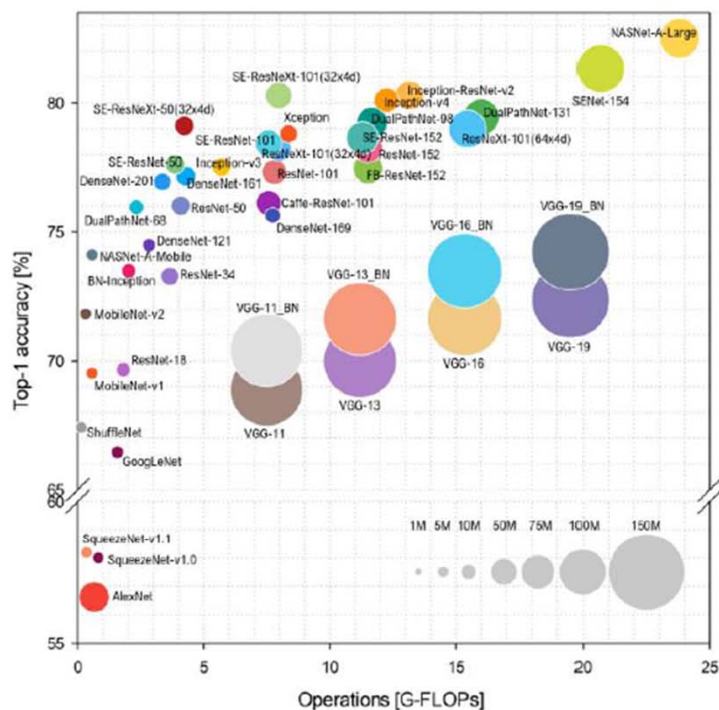
Scalability of system ⬌

> Diversity
> Flexibility
> Low power
> Mixed signal
> Compact

> More computing power: Highly efficient computing
> More memories: Large capacity memories with high bandwidth
> Connectivity with high data transfer capability

⬇

2.5D/3D heterogeneous integration
Chiplet integration
Logic-in-memory, memory-on-logic
Memory computing

⬇

Heterogeneous system integration
Wafer scale integration
Optical interconnection
Effective cooling

# Energy efficiency challenge



Rangharajan Venkatesan, IEEE ISSCC SC1 (2024)

# Optimizing Networks (for the Edge)

Every token is compared to all other tokens to compute attention map → Quadratic complexity

$$attention = softmax(\frac{QK^T}{\sqrt{d_k}})$$



Bram Verhoef, IEEE ISSCC Forum 2.7 (2024)

Bram Verhoef, IEEE ISSCC Forum 2.7 (2024)

# Memory-Base AI Processor to Achieve High Energy Efficiency and Compactness

- CIM: use memory array as a processing unit
- PIM: use embedded logic near memory array as a processing unit
- PNM: use an additional chip for processing inside a memory package or a set



Kyomin Sohn, IEEE ISSCC Forum 2.6 (2024)

# Neuro-Centric Sensor System Project in Tohoku Univ.



National Project in Japan (NEDO) of AI Chip (2019-2020)

# Cyclic Neuro Operation in 3D Stacked AI Chip
## (Forward propagation/ Backward propagation)



Memory layer

Weight
Bias
Activation
etc.

i+1 layer

i+1 layer

i+2 layer

i layer

Synapse

Neuron

Cyclic operation

Direction of neuro operation

3D Stacked AI Chip

Weight word lines

WWL$_{j-1}$    WWL$_j$    WWL$_{j+1}$    WWL$_{j+2}$

Input from neuron x$_{i-1}$

Weight bit line WBL$_{i-1}$
(Weight data input)

Input from neuron x$_i$

Weight bit line WBL$_j$
(Weight data input)

Input from neuron x$_{i+1}$

Weight bit line WBL$_{i+1}$
(Weight data input)

Input from neuron x$_{i+2}$

Weight bit line WBL$_{i+2}$
(Weight data input)

To neuron u$_{i-1}$ in next layer    To neuron u$_{i-1}$ in next layer    To neuron u$_{i-1}$ in next layer    To neuron u$_{i-1}$ in next layer

Output signals

Synapse cell array

# Cross-sectional View of 3D Stacked AL chip with Four Stacked Layers

# Neuron-Block Partition to Densify Sparse Synaptic Connections in CNN



CNN Flow

layer_2n

layer_4n

Identify or zero-padding

Identify or zero-padding

Input

Conv 1 3x3 16 channels | BN 1 16 channels | ReLU

Conv 1 3x3 16 channels | BN 1 16 channels | ReLU | Conv 1 3x3 16 channels | BN 1 16 channels | ReLU

Conv 1 3x3 32 channels | BN 1 32 channels | ReLU | Conv 1 3x3 32 channels | BN 1 32 channels | ReLU

32x32x16

16x16x32

32x32x16

3x3x16=(144)

Filtering operation in 1024 pixels.

Mapping

Weight Matrix

144

32

=

144
Input

32
Output

Weight Channel (w1-w32)

Weight Channel (w1-w32)

Weight Channel (w1-w32)

(w1-w32)

(w1-w32)

Weight Channel w1 (144)

X11 input (144)

X12 input (144)

Xij input (144)

X31,32 input (144)

X32,32 input (144)

Total Input (144 × 1024)

Synapse Circuit (Activated)

Synapse Circuit (Not activated)

X11 output (32)

X12 output (32)

Xij output (32)

X31,32 output (32)

X32,32 output (32)

Total Output (32 × 1024)

Multi-Neurochip Module (AI System Module)

Input ①

Analog/ Digital Output ①

Analog/ Digital Output ②

Input ②

3D Neurochip (Input: 96, Output: 32)

Silicon interposer

# Image Recognition Using 3D Stacked AI chip

## Paradigm Shift from CNN to ViT (Vision Transformer)



CNN

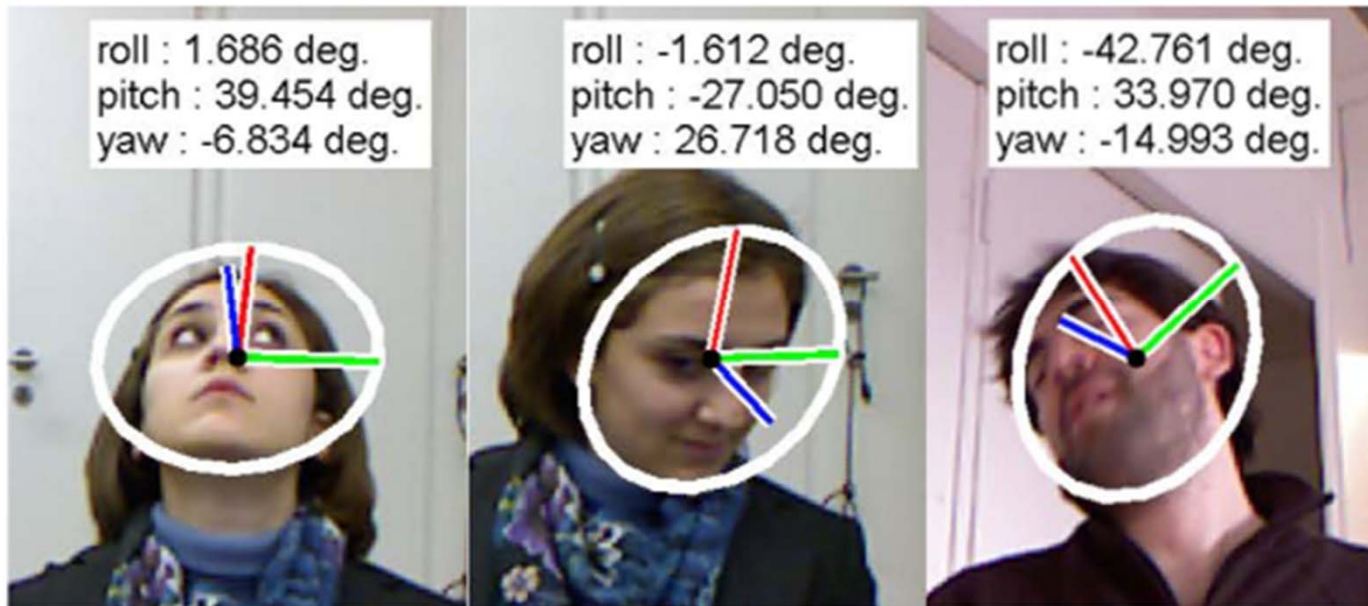CIFAR-10

Vision Transformer (ViT)

➢ We can significantly reduce the number of matrix product operation (MPO) in ViT.

➢ ViT is suitable for Edge application.

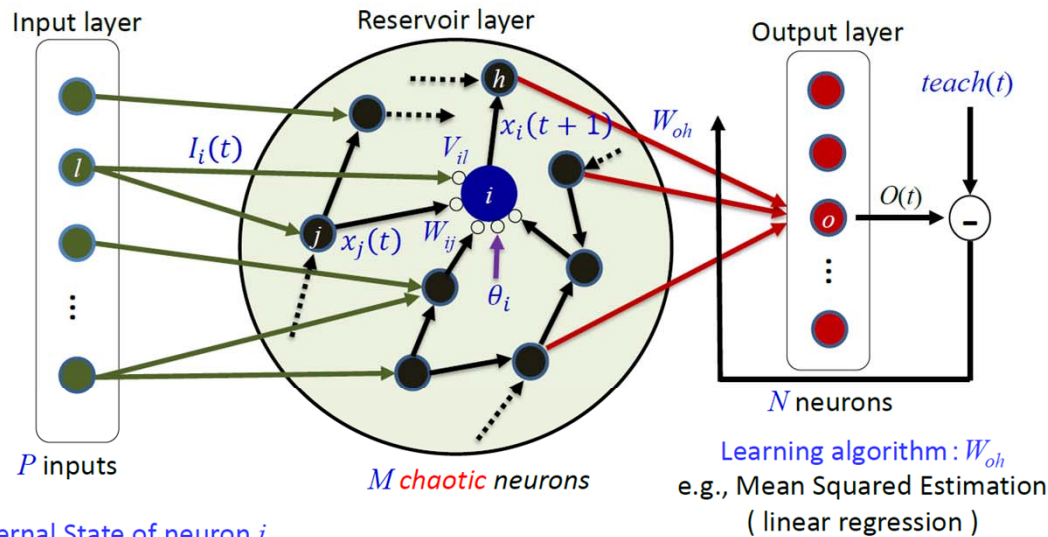| Network | Number of Matrix Product | Accuracy |
|---|---|---|
| CNN（Optimized） | 6,212 | 71.4% |
| Tiny ViT-V1 | 1,169 | 71.3% |
| Tiny ViT-V2 | 150 | 76.6% |

By Courtesy of Prof. T. Okatani, Tohoku University

# Face Recognition Using 3D Stacked AI chip

Average Error＝Yaw angle: 8.0 deg., Pitch angle: 8.7 deg., Roll angle: 7.6 deg.



By Courtesy of Prof. T. Okatani, Tohoku University

# Implementation of Reservoir Neural Network in 3D Stacked AI Chip with Cyclic Neuro Operation

Input layer

Reservoir layer

Output layer

$teach(t)$

$I_i(t)$

$V_{il}$

$x_i(t+1)$

$W_{oh}$

$O(t)$

$x_j(t)$  $W_{ij}$

$\theta_i$

$P$ inputs

$M$ *chaotic* neurons

$N$ neurons

Learning algorithm : $W_{oh}$
e.g., Mean Squared Estimation
( linear regression )

**Internal State of neuron $i$**

$$y_i(t+1) = ky_i(t) + \sum_{j=1}^{M} W_{ij}f(y_j(t)) + \sum_{l=1}^{P} V_{il}I_i(t) - \alpha x_i(t) - \theta_i(1-k)$$
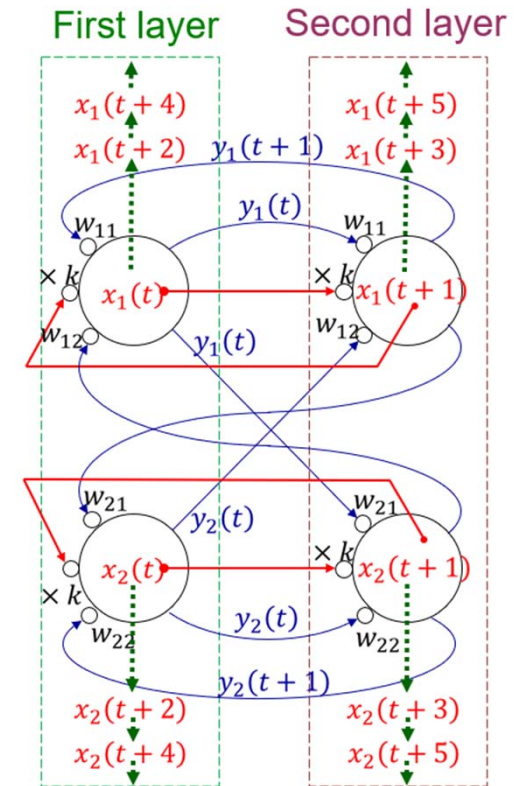
**Output of neuron $i$**

$$x_i(t+1) = f(y_i(t+1))$$

K. Aihara et al., Phys. Lett. A, vol. 144, 1990.

$k$: dumping factor of the refractriness
$\alpha$: scaling factor   $\theta_i$: threshold
$f(\cdot) = 1/(1 + \exp(-x/\epsilon))$

Configuration of Reservoir Neural Network
with Simple Learning

**First layer**   **Second layer**

$x_1(t+4)$   $x_1(t+5)$
$x_1(t+2)$   $y_1(t+1)$   $x_1(t+3)$

$w_{11}$   $y_1(t)$   $w_{11}$

$\times k$   $x_1(t)$   $\times k$   $x_1(t+1)$

$w_{12}$   $y_1(t)$   $w_{12}$

$w_{21}$   $y_2(t)$   $w_{21}$

$\times k$   $x_2(t)$   $\times k$   $x_2(t+1)$

$w_{22}$   $y_2(t)$   $w_{22}$

$y_2(t+1)$

$x_2(t+2)$   $x_2(t+3)$
$x_2(t+4)$   $x_2(t+5)$

Mapping of Reservoir Neural Network to 3D Stacked AI Chip with Cyclic Neuro Operation

K. Fukuda, Yoshihiko Horio et al. , NOLTA, IEICE (2021)

# Voice Recognition Using 3D Stacked AI chip

| Number of Reservoir Neuron | Connectivity within Reservoir |
|---|---|
| 64 | 22 % |

| Input Connectivity | Output Connectivity |
|---|---|
| 6.25 % | 100 % |

Learning (ten times) of "zero" and "one" by Linear Regression

Example of network response after learning
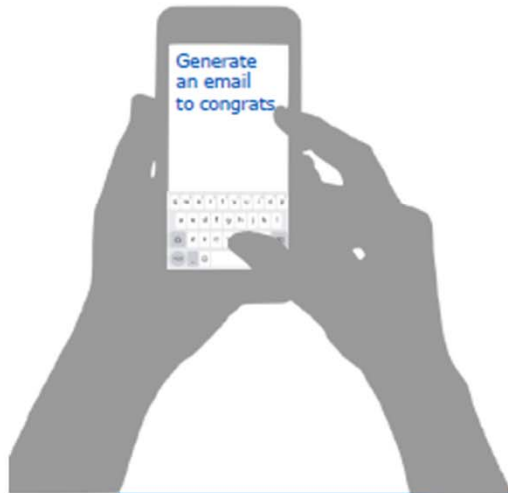
100 % recognition for ten different voices

# Next Generation AI



Bor-Sung Liang, IEEE ISSCC Forum 2.8 (2024)

# LMM AI Chip Project in Tohoku University (AI-Sensor Fusion)



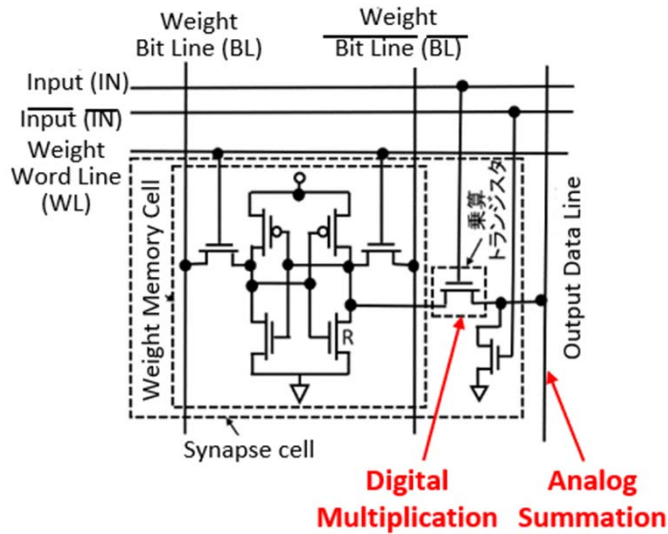Sensor Integrated AI System Module
(Final Target)



Sensor Integrated AI System Module
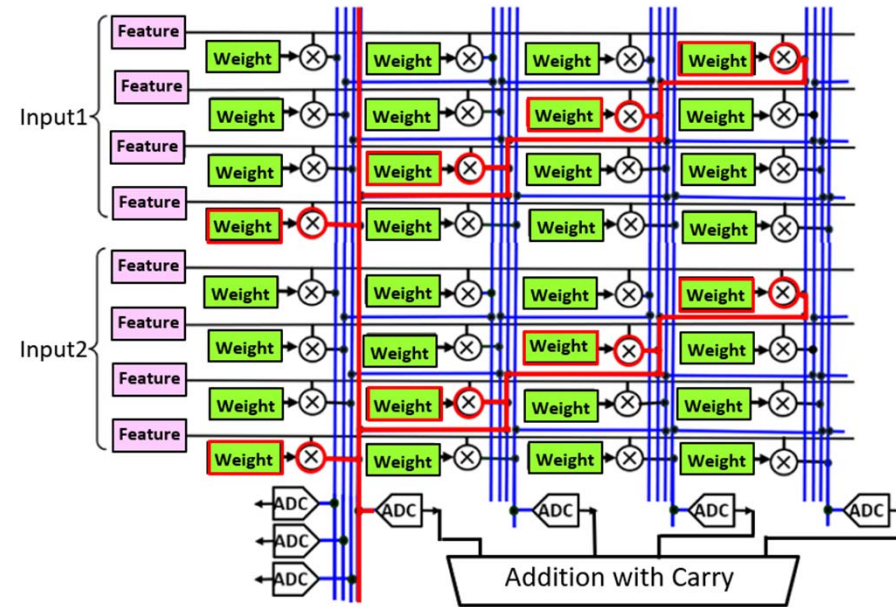by Heterogeneous 3D Chiplet Integration



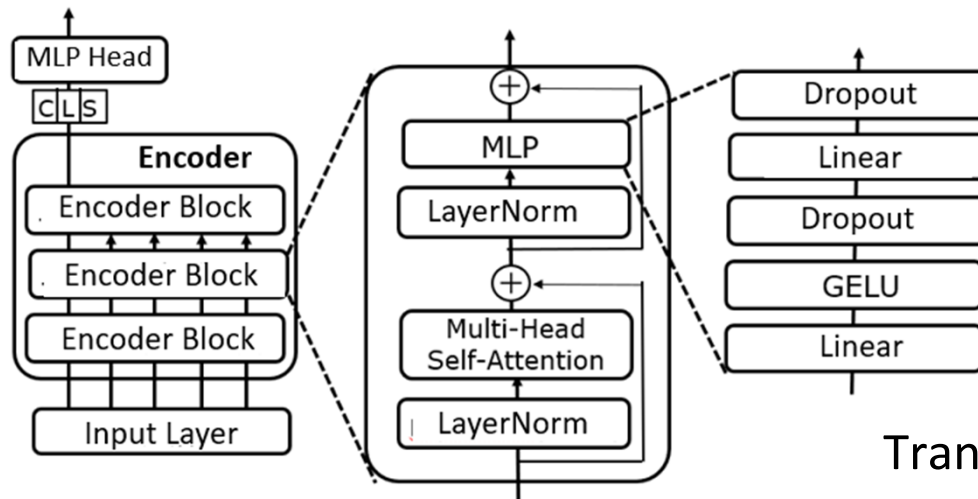Cross-sectional Structure of Sensor Integrated AI System Module

# AI Chip Based on Memory-in-Computing (CIM) with SRAM and Transformer Algorithm



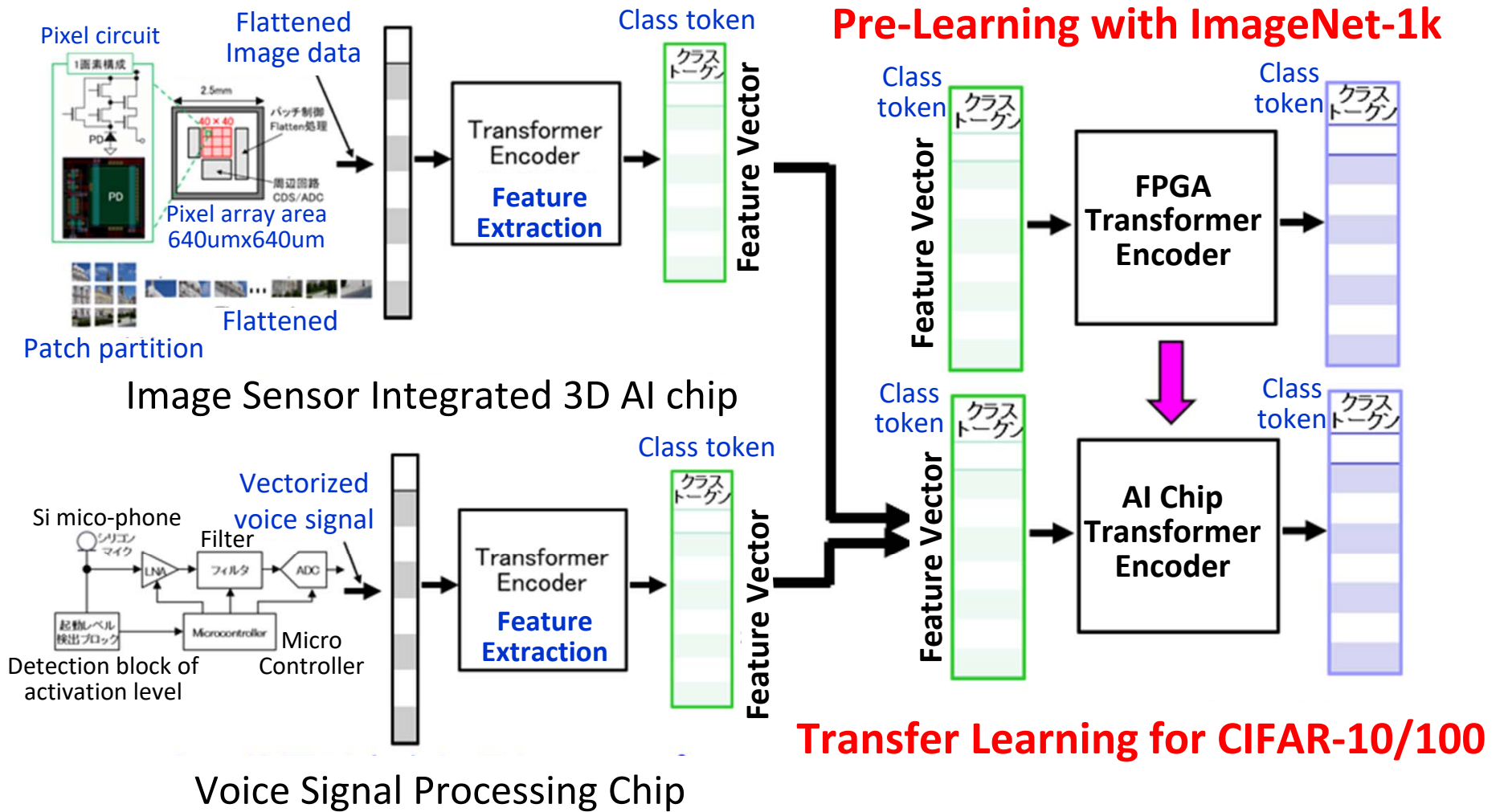CIM (Memory-in-Computing)
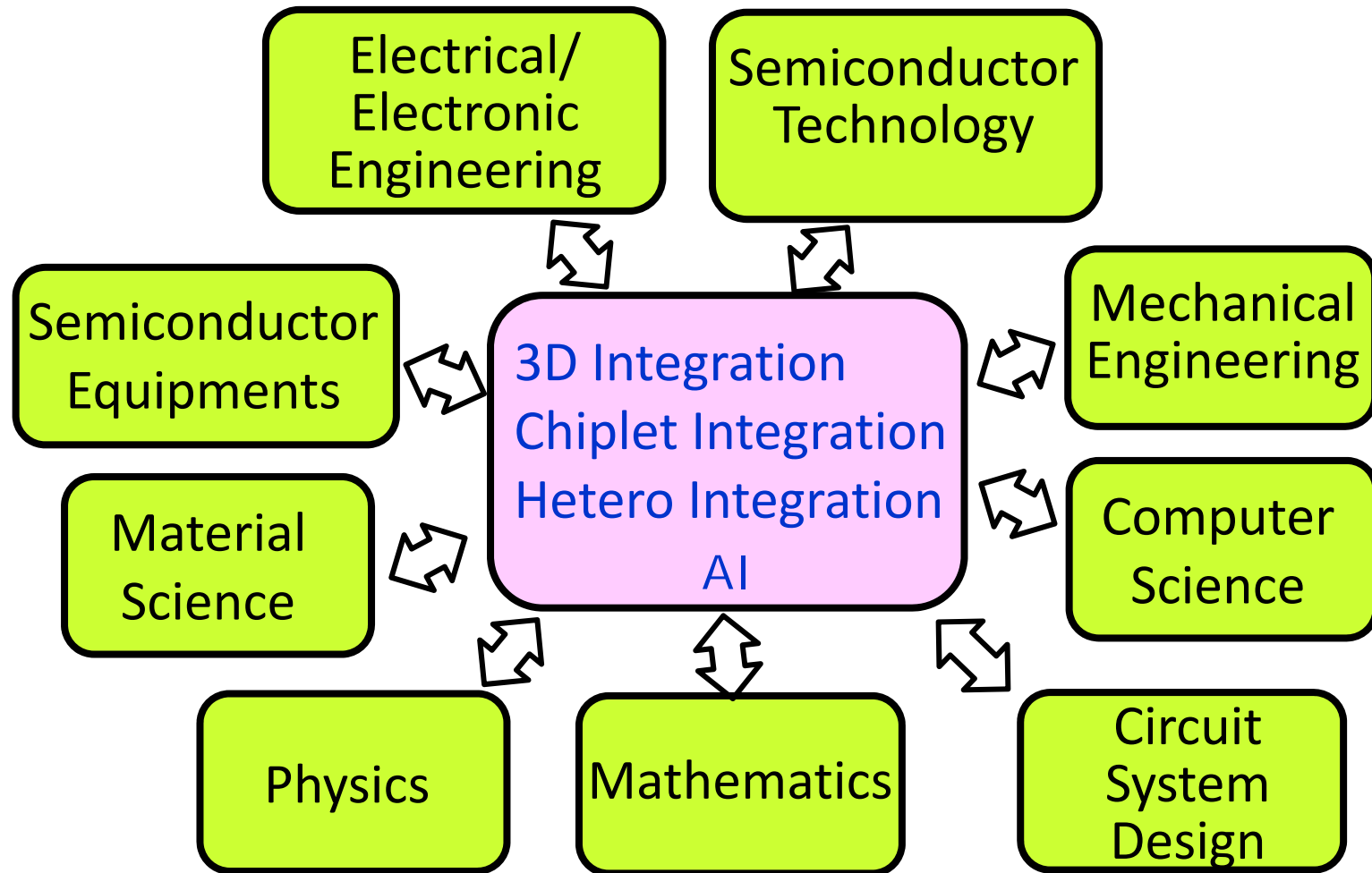by SRAM Memory Cell



CIM Basic Synapse Circuit Array



Transformer Algorithm

# LMM AI Chip Project in Tohoku University（AI-Sensor Fusion）

## Energy Efficient Neuro Processing and Reduced Data Transfer by Integrating Sensor Devices with Feature Extraction Function



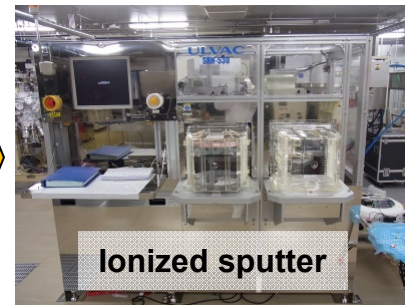**Pre-Learning with ImageNet-1k**

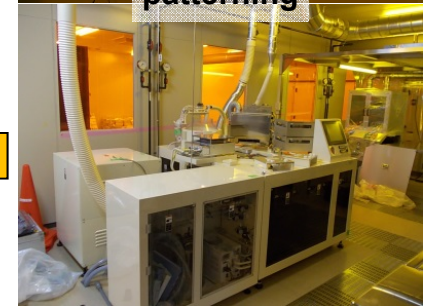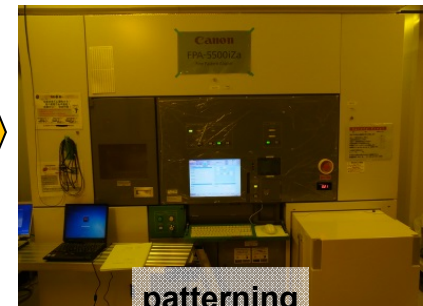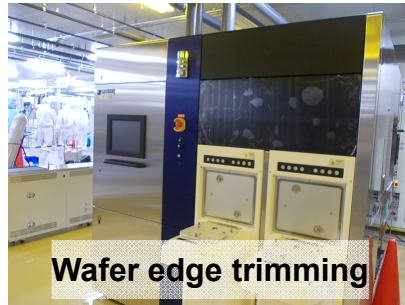**Transfer Learning for CIFAR-10/100**

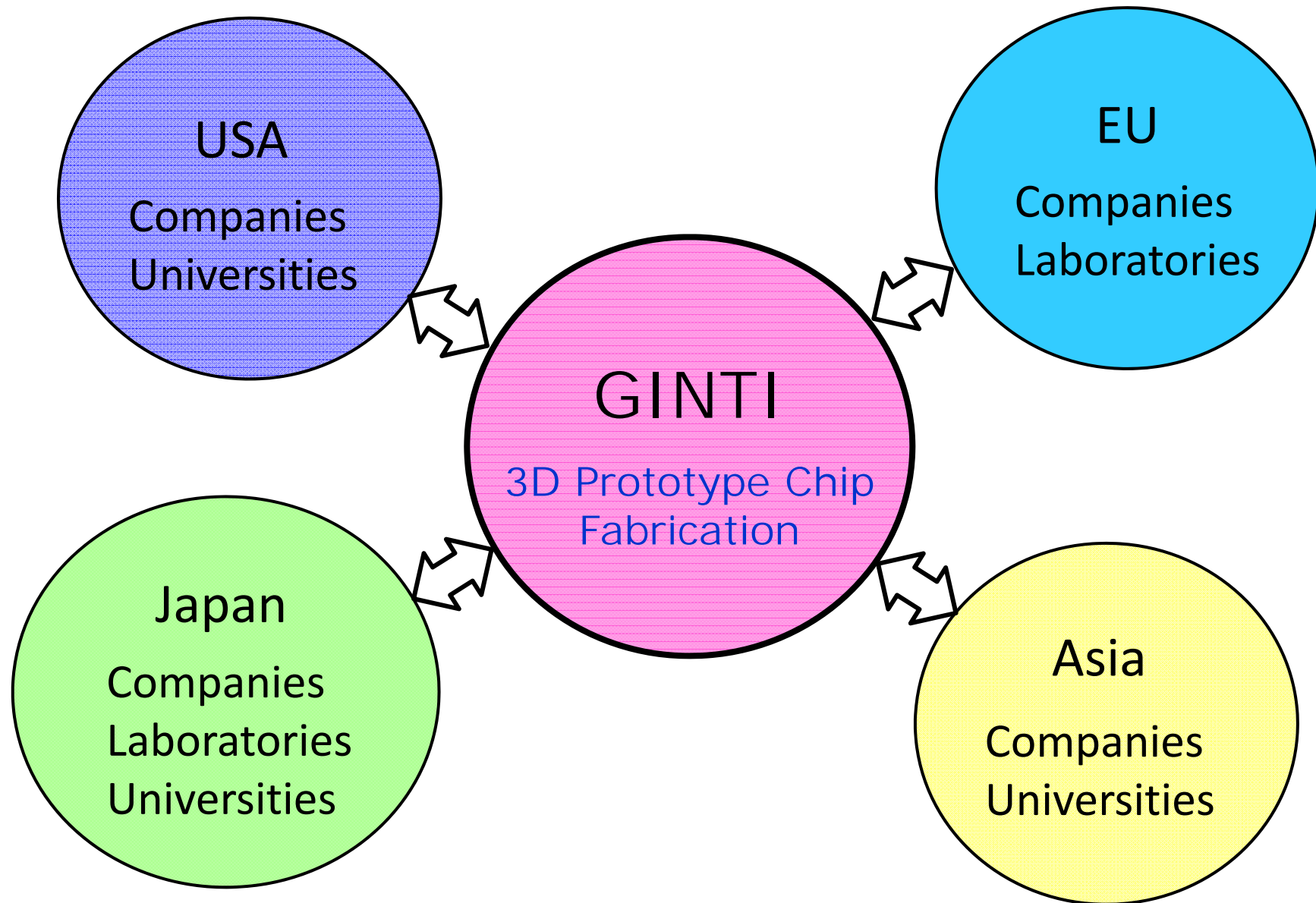# Heterogeneous Chiplet Integration/AI Need Various Kind of Knowledges and Skills

# 12-inch 3D Production Line in Tohoku Univ.
# GINTI（Global INTegration Initiative）

# International Cooperation IN GINTI

**USA**
Companies
Universities

**EU**
Companies
Laboratories

**GINTI**
3D Prototype Chip
Fabrication

**Japan**
Companies
Laboratories
Universities

**Asia**
Companies
Universities

# Conclusions

➤ 3D heterogeneous integration and chiplet integration are the key for future intelligent systems such as HPC, AI/ML, post-5G/6G systems and quantum computer systems.

➤ 3D heterogeneous integration and chiplet integration need various kinds of knowledges and skills. Therefore International collaborations are indispensable.