



Near Data Processing for AI & Big data (Data Hierarchy)



MEMORY SYSTEMS R&D

Eui-Cheol Lim

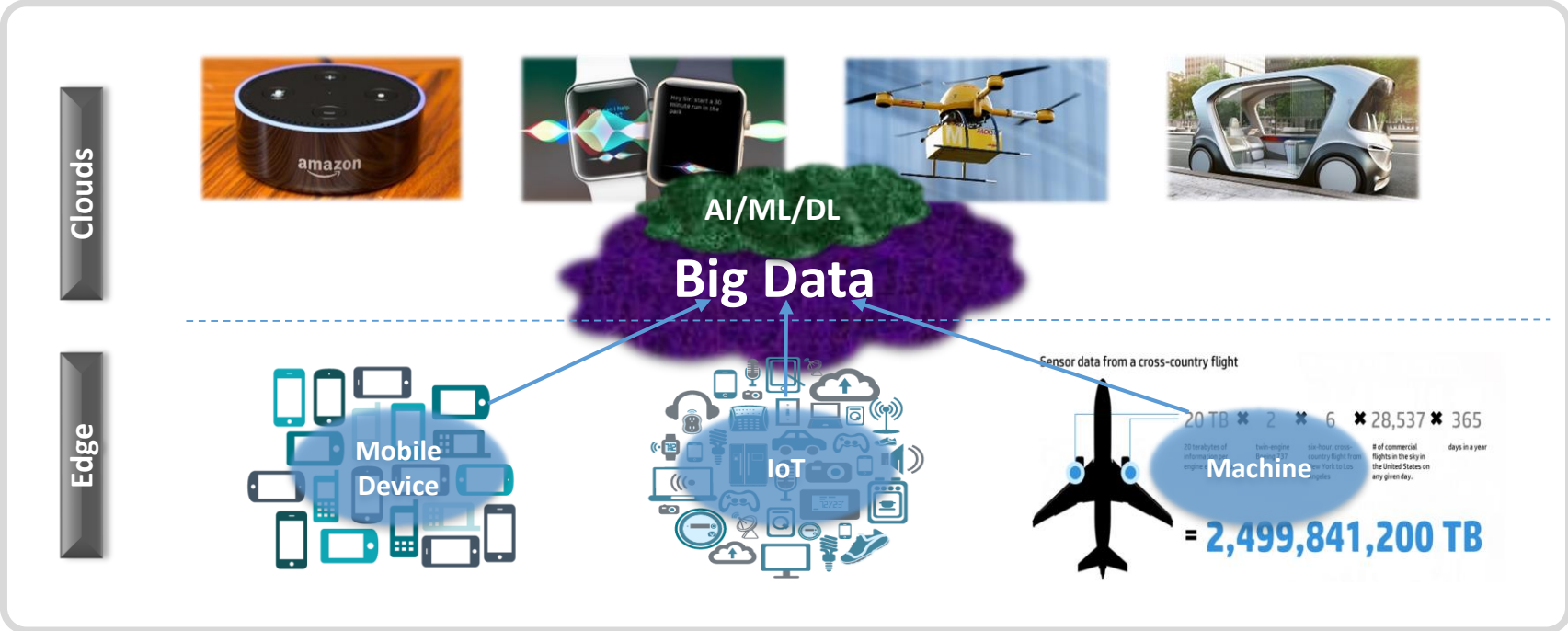
Contents



1	Computing Trend – AI, Big Data
2	Architectures for AI
3	Near Data Processing
4	Data Hierarchy – holistic approach for Near Data Processing
5	Conclusion

AI and Big Data

- AI with big data processing is one of the key technologies in modern industry
- From the edge to the clouds, AI is proven to be advantageous in various application fields



Computer vs. Human Brain

- AI has already proven that it is more capable than human in certain applications
- However in energy perspective, human brain is much more efficient than current computer system

Go match : AlphaGo vs. Lee Sedol

AlphaGo

1202 CPUs, 176 GPUs, 100+ Scientists
170kWh

Lee Sedol

Lee Sedol 1 Human Brain, 1 Coffee
20Wh



AlphaGo



Lee Sedol

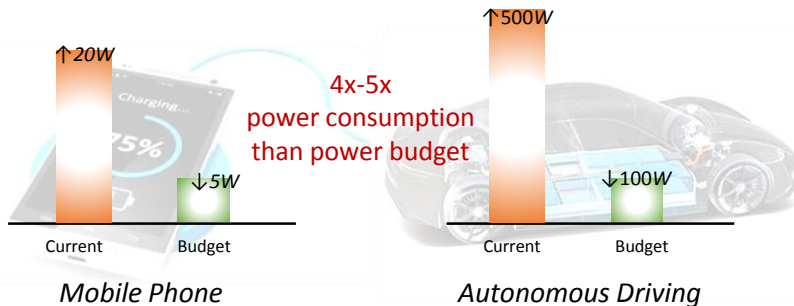


Power Efficiency – Challenges of AI

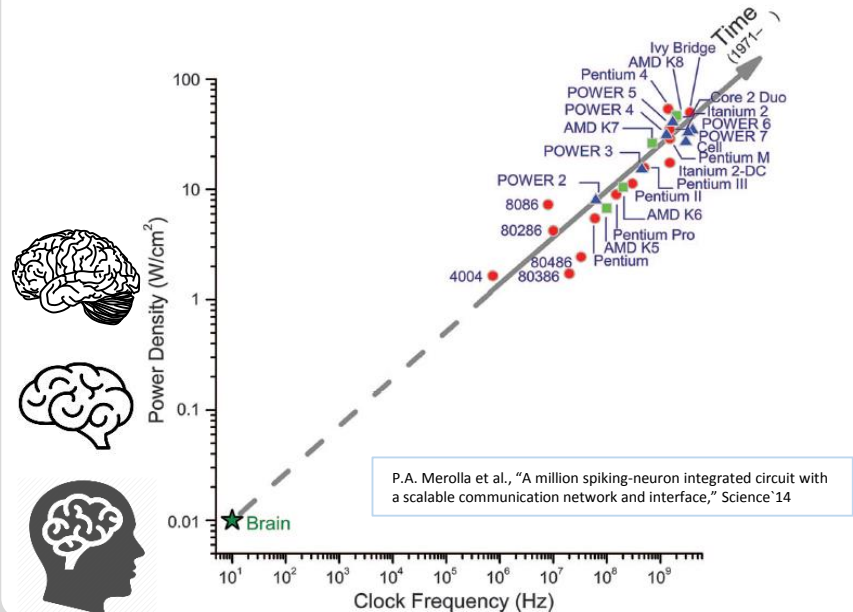
- AI performance is limited by power budget
- Power density and clock frequency of traditional computer architecture is constantly increasing while brain operates at the lowest point

Power Limitation in AI application

- Current AI applications are power hungry
 - Most of AI apps already exceed its power budget
 - AI still needs more data to enhance accuracy, and current architecture cannot satisfy that requirement



Power Efficiency: Human Brain vs. Computer




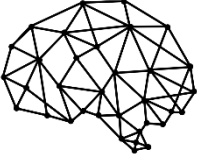
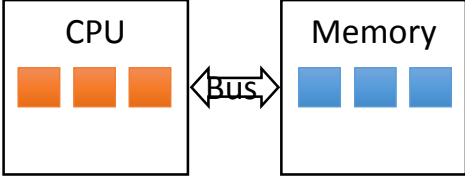
Contents



1	○	Computing Trend – AI, Big Data
2	○	Architectures for AI
3	○	Near Data Processing
4	○	Data Hierarchy – holistic approach for Near Data Processing
5	○	Conclusion

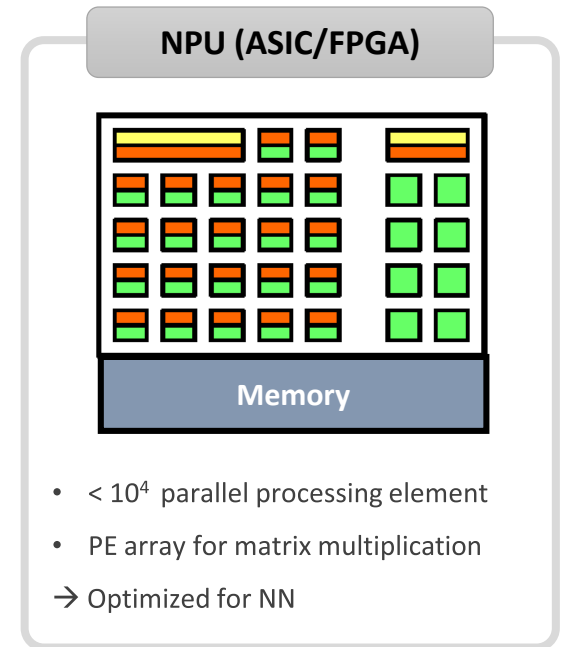
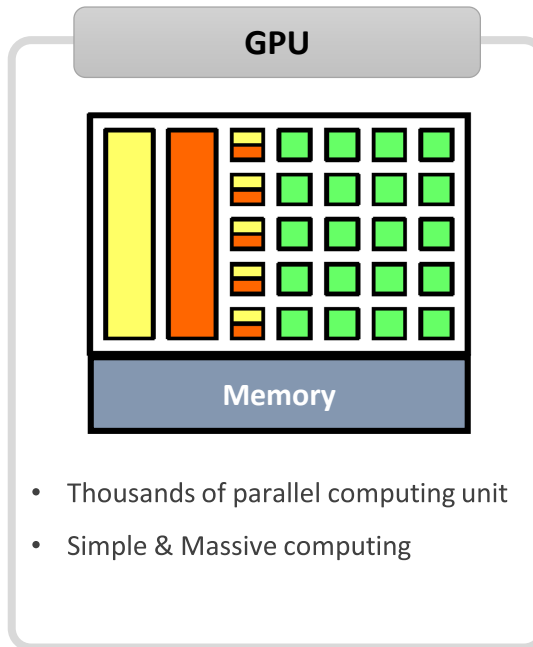
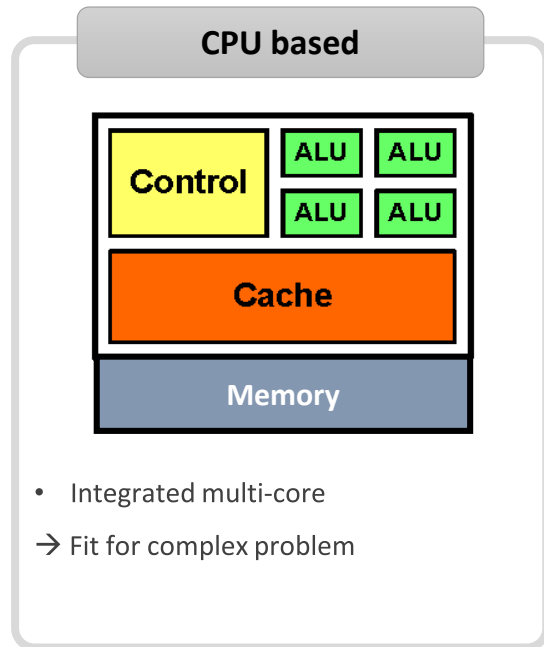
Human Brain vs. Neuromorphic vs. von Neumann

- Neuromorphic computing looks far better than von Neumann architecture, but it is in premature stage yet

	Human Brain	Neuromorphic	von Neumann
Architecture			
Device	Neuron / Synapse	Artificial Neuron / Synapse	CPUs / Memories / Bus
Features	Event-driven	Event-driven / Asynchronous	Computing-centric / Compute, memory separated
Encoding Scheme	Spiking Signals	Analog/Digital/Mixed Spiking Signals	Binary Signals
Accelerator	-	Neuromorphic Processor	GPU / FPGA / ASIC
Power / Energy Consumption	Ultra Low	Low	Relatively High

Architecture Revisited: von Neumann

- Data movement is a common issue in von Neumann architecture and causes power efficiency issue as AI requires big data calculation

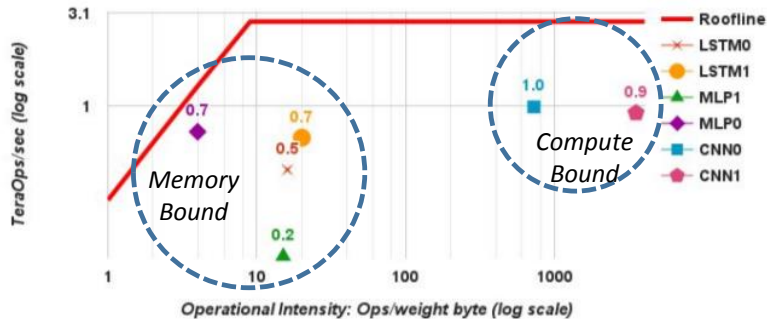


Energy Cost of Data Movement

- Most AI algorithm is memory bandwidth bounded
- Data movement consumes more than 60% of power in modern processors

ROOFLINE Graph According to Neural Network

- In most cases, NN performances are bounded by B/W
 - Recent NN Accelerators use internal memory to reduce bottleneck, but it is not sufficient.



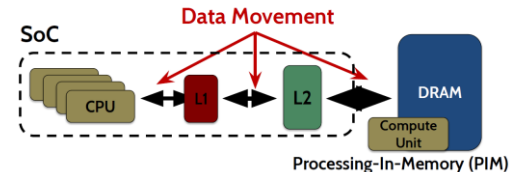
Norman P. Jouppi, "In-Datcenter Performance Analysis of a Tensor Processing Unit," ISCA'2017

Power Portion by Data Movement

- More than 60% of power is consumed by data movement
 - Analyze power consumption in popular google apps

Energy Cost of Data Movement

1st key observation: 62.7% of the total system energy is spent on data movement



Amirali Boroumand, "GoogleWorkloads for Consumer Devices: Mitigating Data Movement Bottlenecks," ASPLOS'18



Chrome



TensorFlow



Video Playback



Video Capture

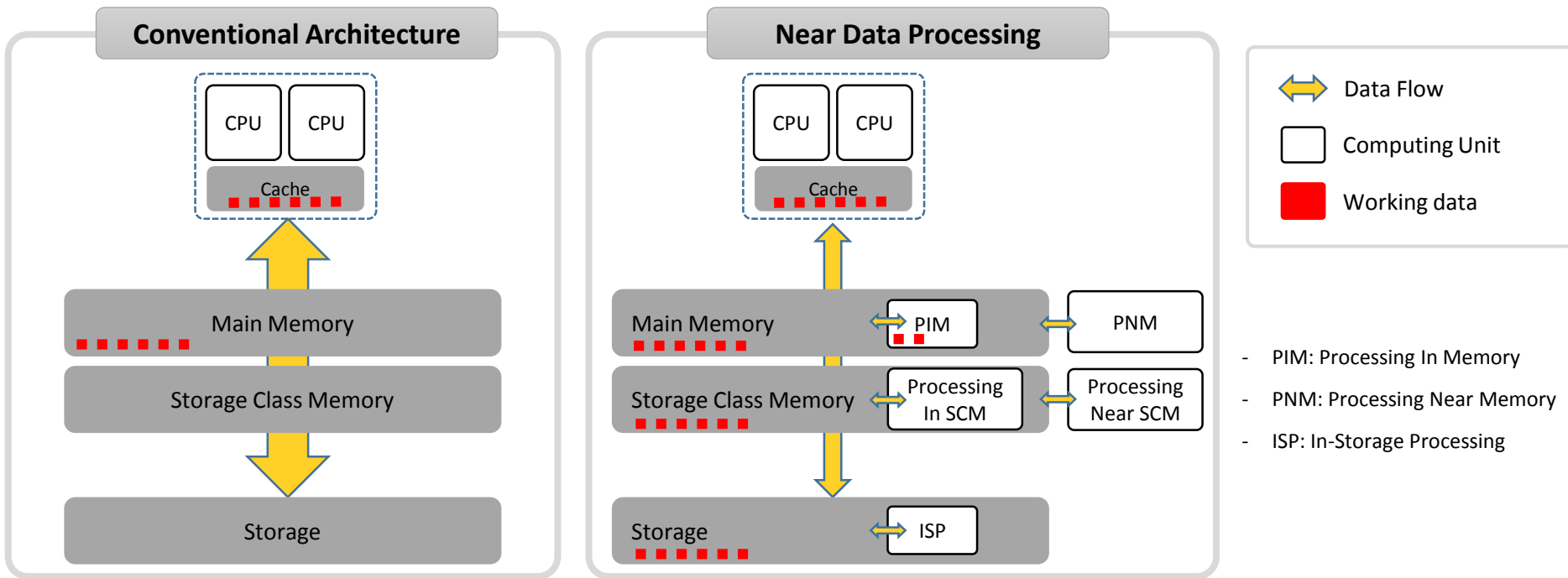
Contents



1	Computing Trend – AI, Big Data
2	Architectures for AI
3	Near Data Processing
4	Data Hierarchy – holistic approach for Near Data Processing
5	Conclusion

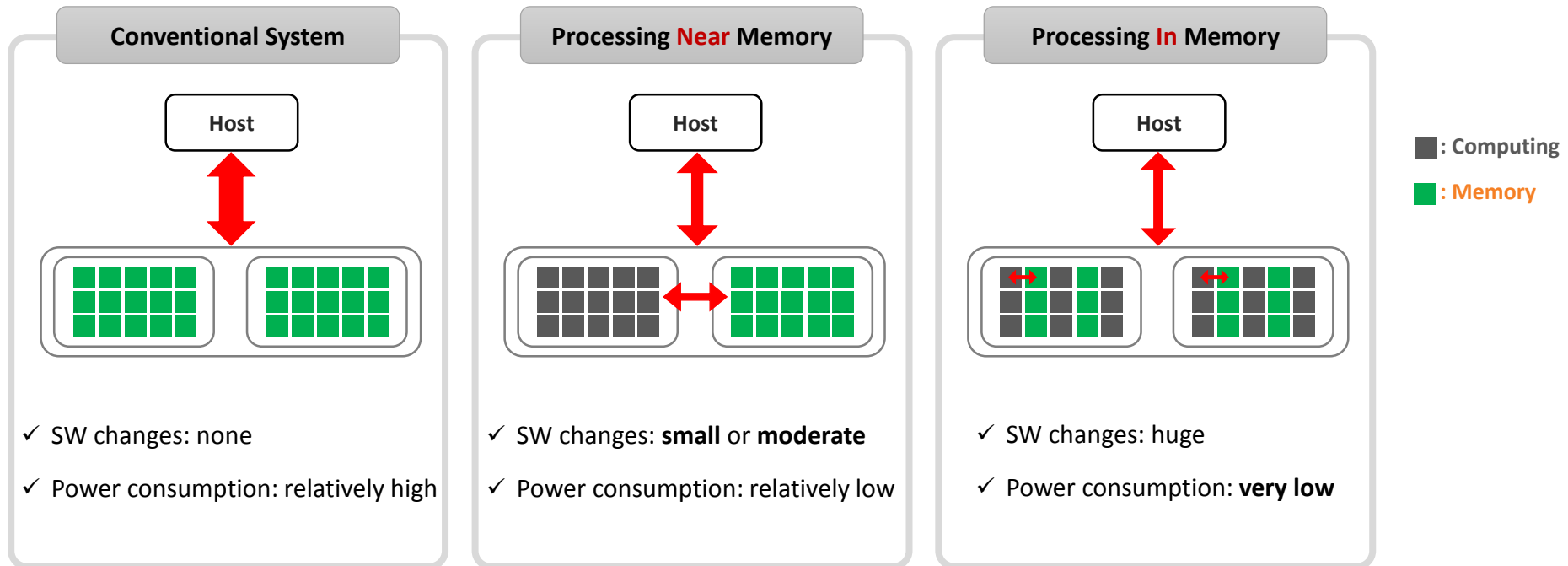
Moving Computation To Data

- Near Data Processing can reduce energy for data transfer by locating compute node where data lives



Near Data Processing Concept at Memory Level – PIM, PNM

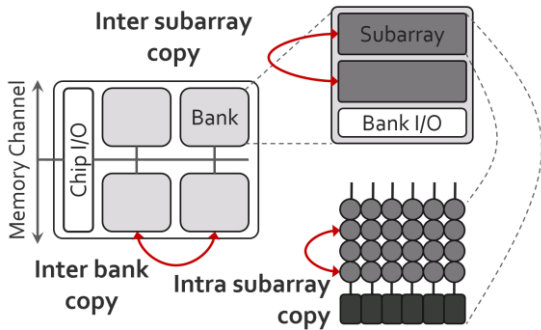
- PIM (Processing In Memory), PNM (Processing Near memory)
 - Pros: Overcome bandwidth limitation between logic and memory devices with lower power consumption
 - Cons: Not backward compatible with legacy software stack, requiring some changes in system software



PIM, PNM Researches

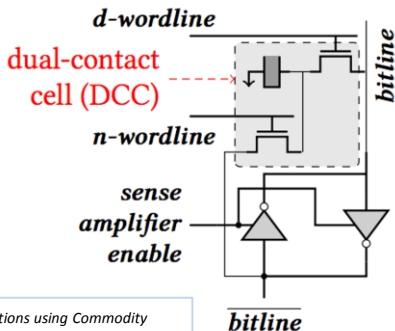
PIM

- **Row data copy**
e.g.) Bulk copy



Seshadri et al., "RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization," MICRO'13

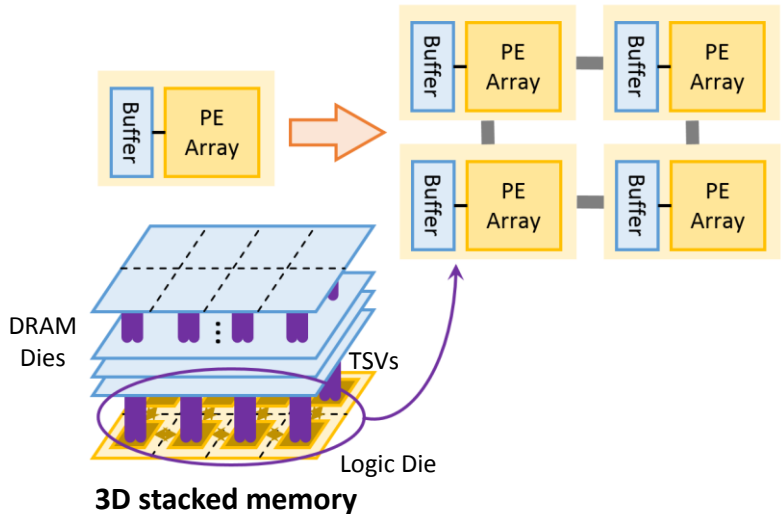
- **Bulk bitwise operations**
e.g.) In-DRAM bit-wise NOT



Seshadri et al., "Ambit: In-Memory Accelerator for Bulk Bitwise Operations using Commodity DRAM Technology," MICRO'17

PNM

- **3D stacked Memory**
e.g.) NN accelerator

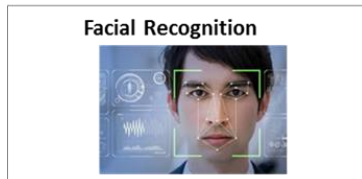
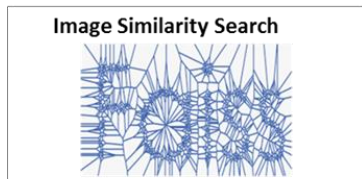


M. Gao et al., "TETRIS: Scalable and Efficient Neural Network Acceleration with 3D Memory," ASPLOS'17

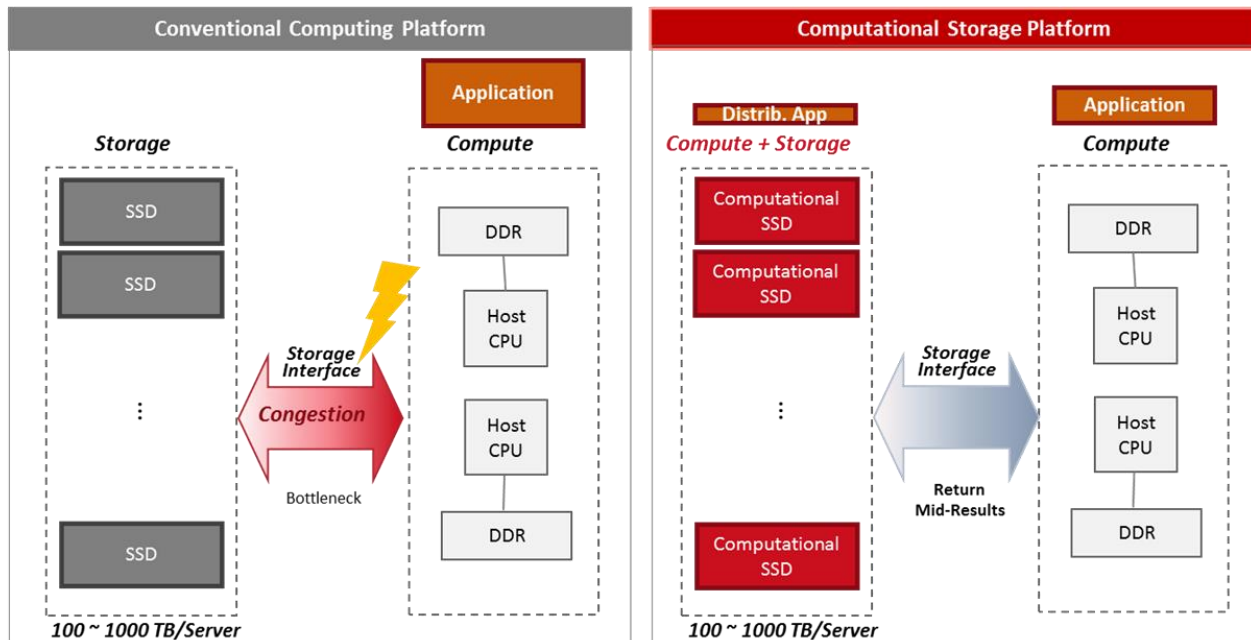
Near Data Processing at Storage Level – In Storage Processing

- As ML/AI based Big Data analytics workload increases I/O congestion between storage and compute node which could cause a performance bottleneck

Emerging Use Cases



ML/AI
Processing
of
Massive
Immersive
Media
Data



4.5hrs*

Use Case Reducing Analytics Time by ~x5

1hr*

* 64x 8TB SSD (Storage System 0.5PB), PCIe 3.0x16, DRAM B/W 2x16GB/s

Researches on In Storage Processing

- Big-data analysis, distributed machine learning are key target workloads for computational SSD researches

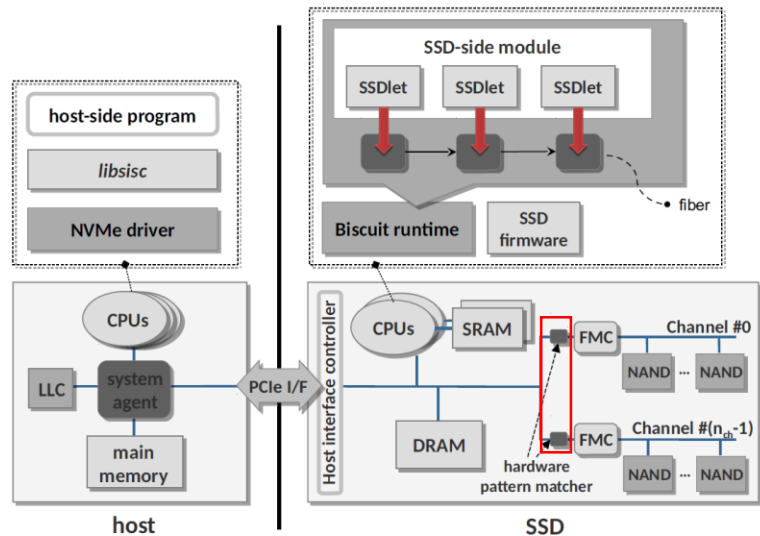


Figure 2. Overall architecture of a Biscuit system.

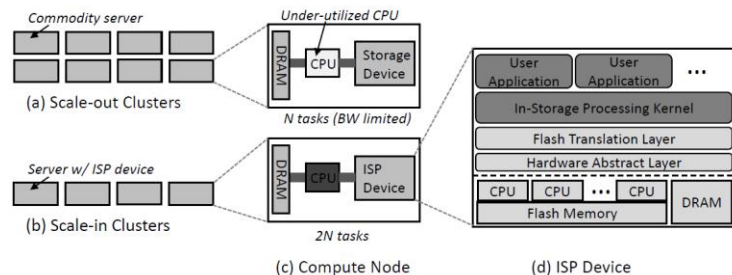
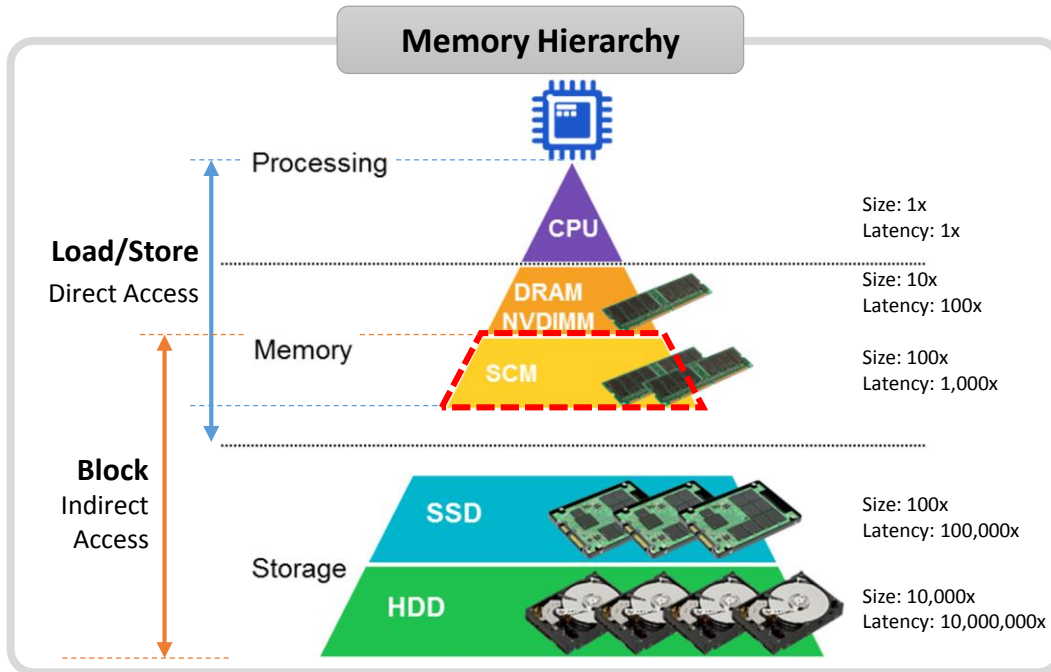


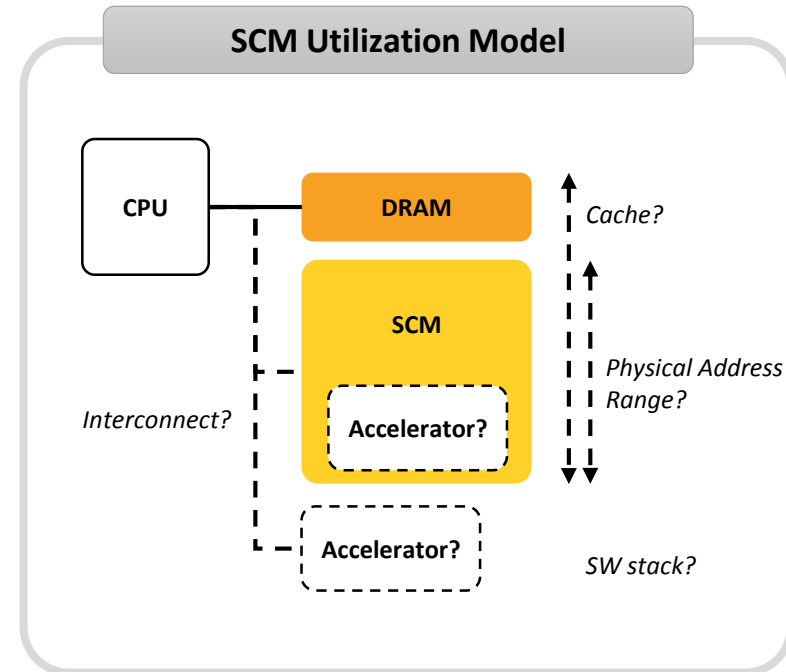
Figure 5: Scale-in clusters with ISP. ISP devices can be implemented utilizing modern SSD architectures.

Near Data Processing for SCM/NVM Layer

- **SCM (Storage Class Memory)** could be useful in increasing the memory capacity: slower but larger, cost efficient compared to DRAM, and persistent
- **SCM utilization models** are actively being researched including the application in the NDP structure



Source: Rambus



Why Near Data Processing Again?

- Emerging applications' strong demands are driving Near Data Processing Architecture
- Lower hanging fruit – 3D, 2.5D stacking technology improves feasibility of Near Data Processing

Strong Demands

General application

- Computing centric workloads
- High locality, data reuse



Emerging big data, AI applications

- Memory-intensive workload
- Low locality, abundant parallelism
- Bandwidth & Energy Constraints
- Irregular access patterns

Feasibility

High cost of integrating compute units within DRAM



3D, 2.5D die stacking

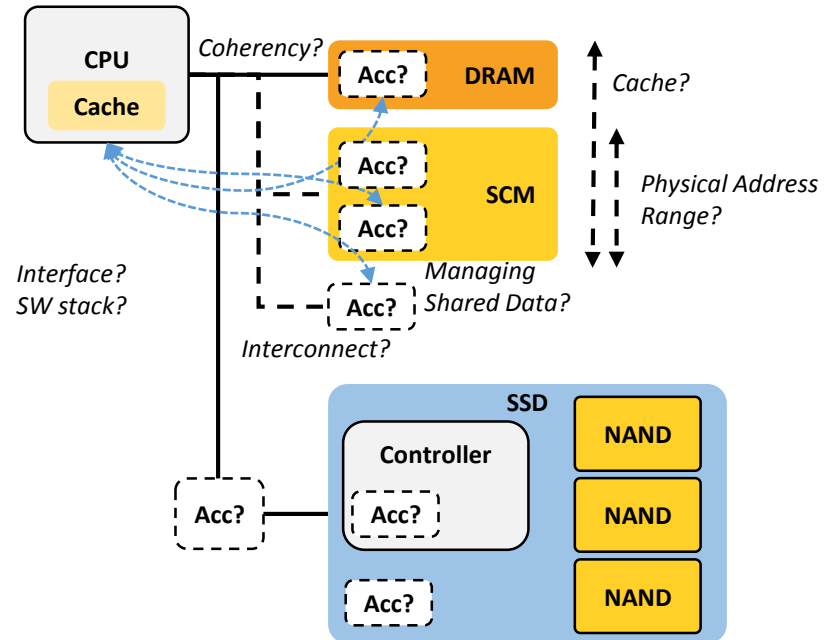
- Low cost
 - High feasibility
- than computing in DRAM

Key challenges in NDP - System

- Many challenges exist in enabling near data processing technology to be widespread

System Issue

- How to design **novel accelerator architecture**?
- How to support **virtual memory space**?
- How to maintain **cache coherence**?
- **Parallel computing** scheme?
 - e.g.) How to manage shared data
- How to design **data structure** for NDP?
- What kind of **interconnect** should be used?

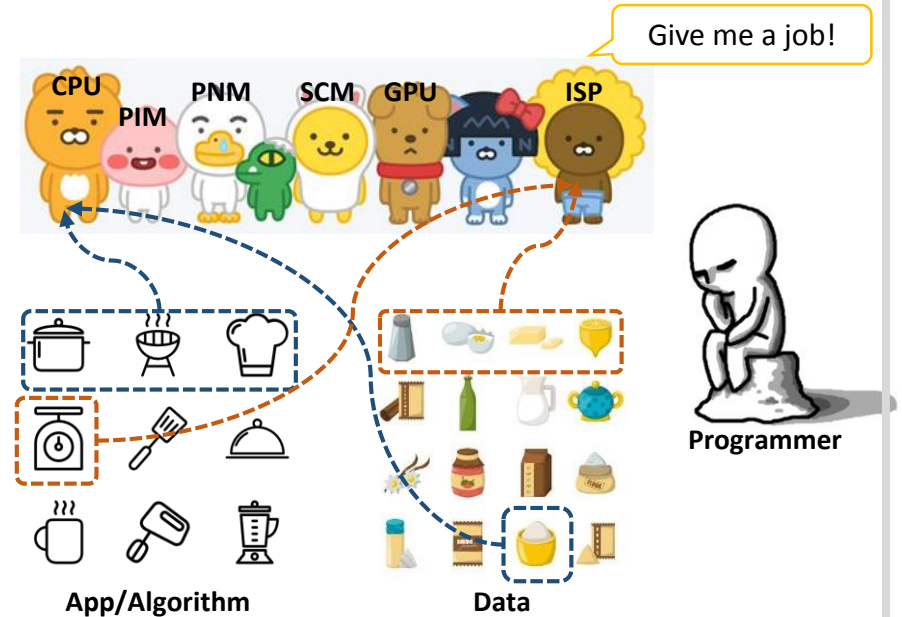


Key Challenges in NDP – Programming Model

- How to establish a programming model is also an important issue

Programming Model

- Which **part** of an application should run in the NDP units?
- How to schedule code? Static? Dynamic?
- Who schedule? Programmer? Compiler? OS?
- Need to **integrate** system software with programming frameworks



Contents



1 ○ **Computing Trend – AI, Big Data**

2 ○ **Architectures for AI**

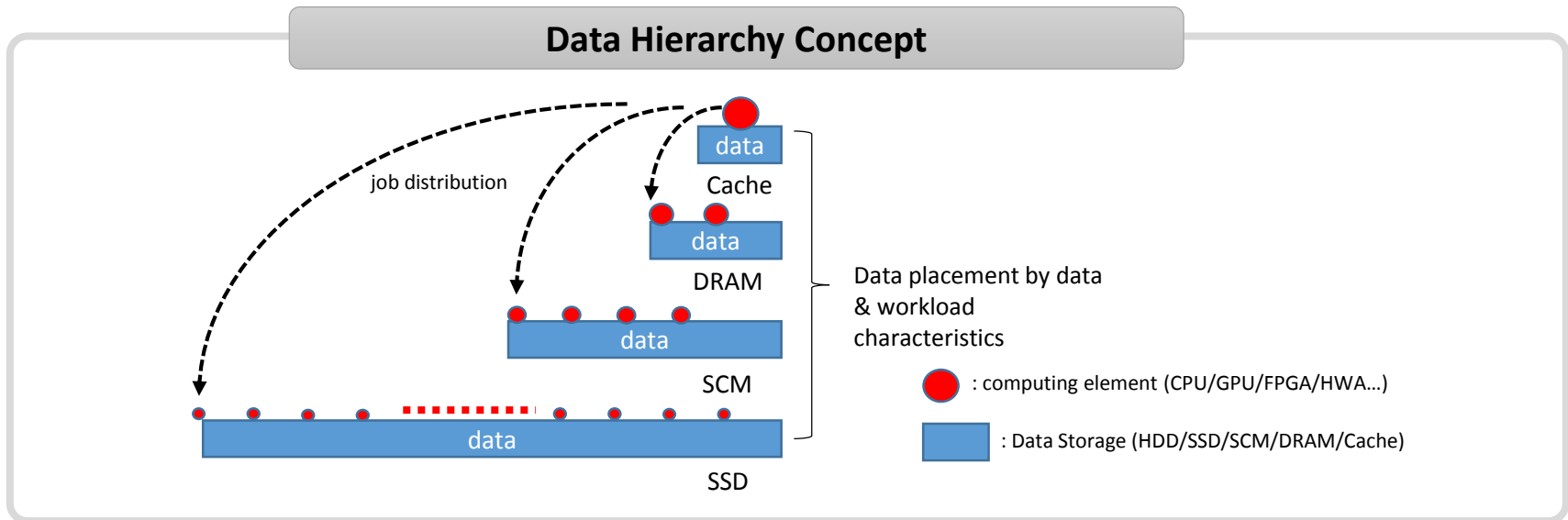
3 ○ **Near Data Processing**

4 ○ **Data Hierarchy – holistic approach for Near Data Processing**

5 ○ **Conclusion**

Data Hierarchy – NDP For All Memory Layers

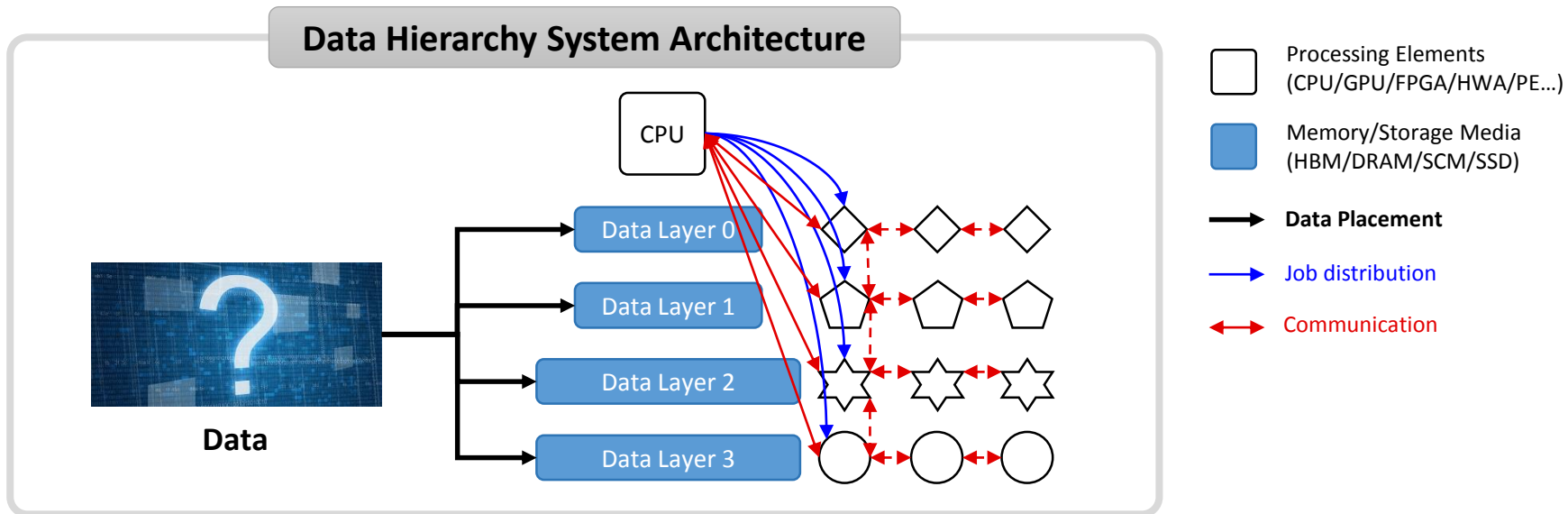
- **Basic Concept – Every Memory layers have own processing element**
 - Minimizing data movement for the **entire memory hierarchy**
 - **Data placed** based on the data & workload’s characteristics
 - **Compute data where data reside**
- We name this system architecture as “**Data Hierarchy**”



Research Topics in Data Hierarchy – Data placement

▪ Data placement

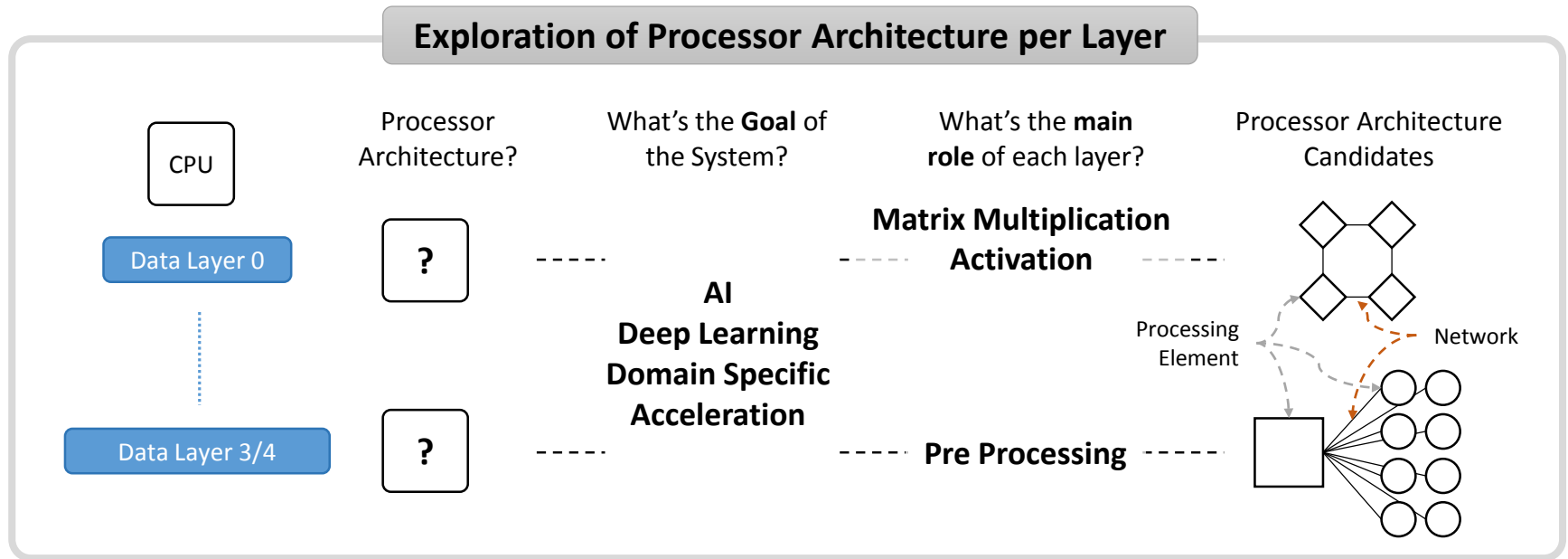
- How to **classify** the characteristics of data?
- How to **place** the data at the appropriate layer
- Whether the characteristics of data will be **determined only by the data itself**, or **by the algorithms** that utilize the data



Research Topics in Data Hierarchy – Processor Architecture

Processor architecture per layer

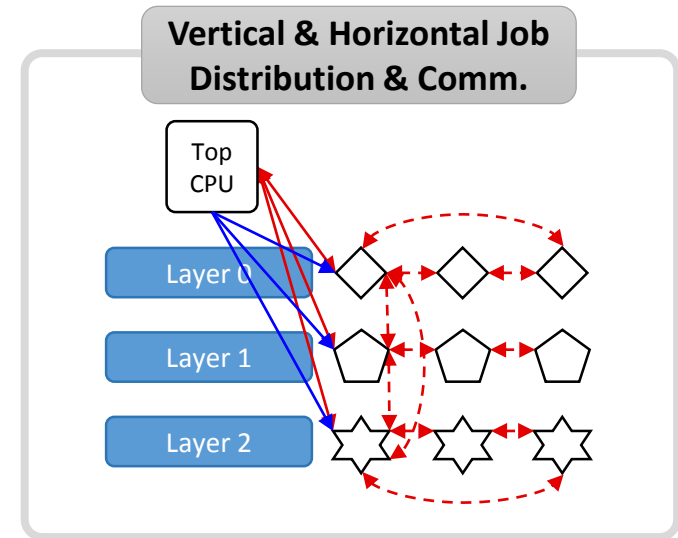
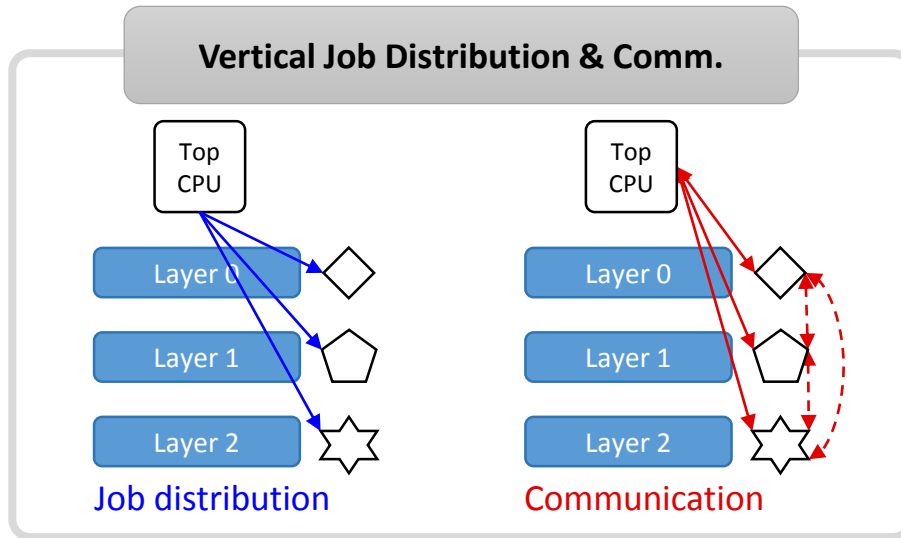
- What is the feasible structure of processing elements for each layer?
- **Domain specific architecture** is the key



Research Topics in Data Hierarchy – Job Distribution & Comm.

Job distribution & Communication

- In terms of job distribution, Data Hierarchy is similar to heterogeneous computing
- Framework is required to assign a job to each layer and aggregate the results



Contents



1	Computing Trend – AI, Big Data
2	Architectures for AI
3	Near Data Processing
4	Data Hierarchy – holistic approach for Near Data Processing
5	Conclusion

Conclusion

- As AI and Big-data became widely spread, system architecture has to be improved to solve energy efficiency and performance scaling issues
- One of the ways to solve the issues is to **place the data in different layers based on characteristics** and **processing the data independently** within that layer → **Data Hierarchy** Concept
- Yes, there are lots of research topics for Data Hierarchy
 - Data placement, data mapping
 - Processor architecture for each layer
 - Job distribution

“Start from simple and domain specific applications, but have to target as general as possible”



THANK YOU

E-mail: euicheol.lim@sk.com



MEMORY SYSTEMS R&D

