# PIM and Various Computational Memory Solutions

MEMORY
FOREST

**2022. 04. 02**

**SK Hynix, Euicheol Lim**

We Do Technology | SK hynix

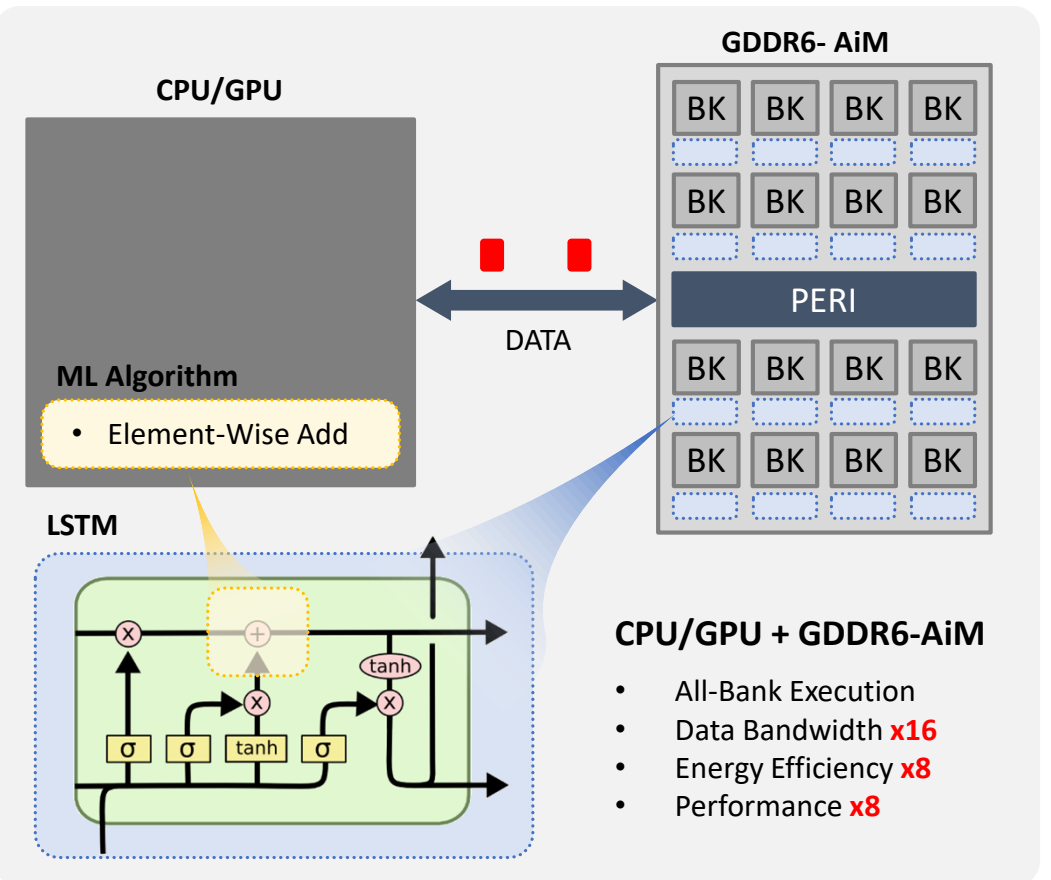# CONTENTS

SK hynix

MEMORY
FOR EST

# GDDR6-AiM Overview

## Definition

GDDR6-AiM is a **GDDR6**-based Accelerator-in-Memory (**AiM**) device targeted for memory-intense Machine Learning algorithm (**RNN, LSTM, MLP**) inference acceleration by offloading certain mathematical operations (**MAC, Activation Function, Element-Wise Multiplication**) from the host (**CPU, GPU, FPGA**).
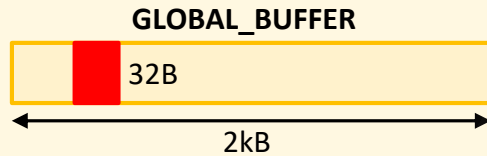
**Conventional System**



CPU/GPU

DRAM

**ML Algorithm**
- MAC
- Activation Function
- Element-Wise Mul
- Element-Wise Add

DATA

BK BK BK BK
BK BK BK BK
PERI
BK BK BK BK
BK BK BK BK

**Machine Learning Algorithm: LSTM**

**In-Memory Accelerated System**

CPU/GPU

GDDR6- AiM

**ML Algorithm**
- Element-Wise Add

DATA

LSTM

**CPU/GPU + GDDR6-AiM**
- All-Bank Execution
- Data Bandwidth **x16**
- Energy Efficiency **x8**
- Performance **x8**

# GDDR6-AiM Overview

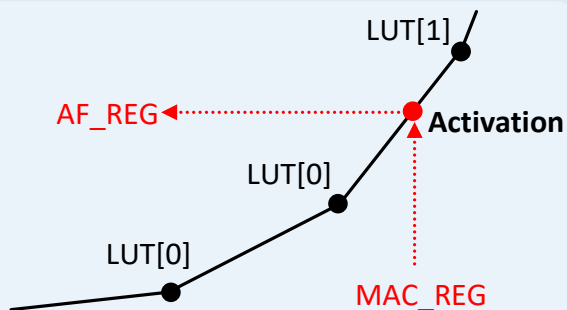## Global Buffer

- Supplementary 2 kB SRAM buffer.
- Provides vector data for MAC.
- Supports 32B WRITE operations.

**GLOBAL_BUFFER**

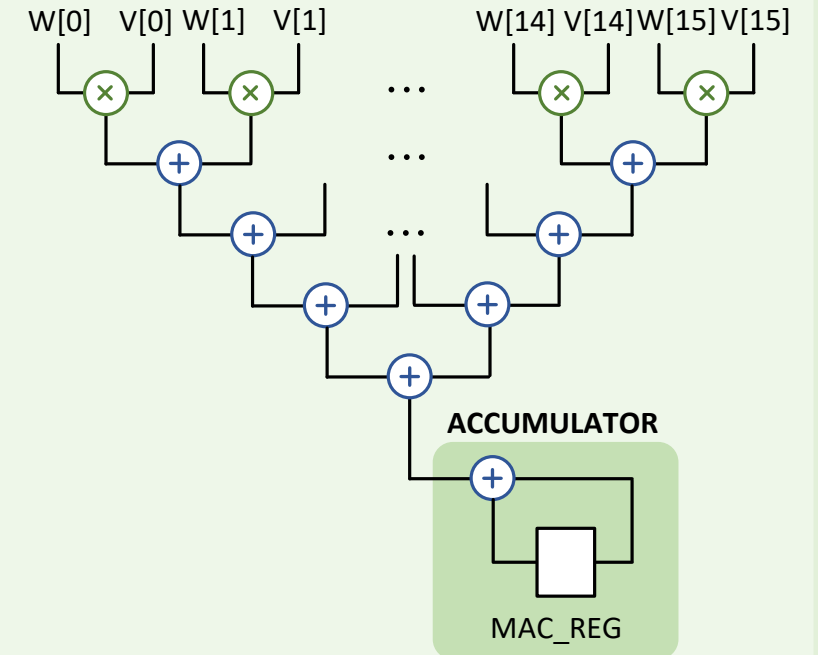| | 32B | |
|---|---|---|

← 2kB →

## Activation Module

- Performs **Activation Function (AF)** computation by linearly Interpolating pre-stored AF template data using MAC calculation results.
- Activation results are stored in a dedicated **AF REG** set and can be later accessed by the user.



LUT[1]
AF_REG ← Activation
LUT[0]
LUT[0]
MAC_REG

| BK0 | BK3 | BK4 | BK7 |
|---|---|---|---|
| MAC | MAC | MAC | MAC |
| Activation | Activation | Activation | Activation |
| Activation | Activation | Activation | Activation |
| MAC | MAC | MAC | MAC |
| BK1 | BK2 | BK5 | BK6 |

| GLOBAL BUFFER | PERI |
|---|---|

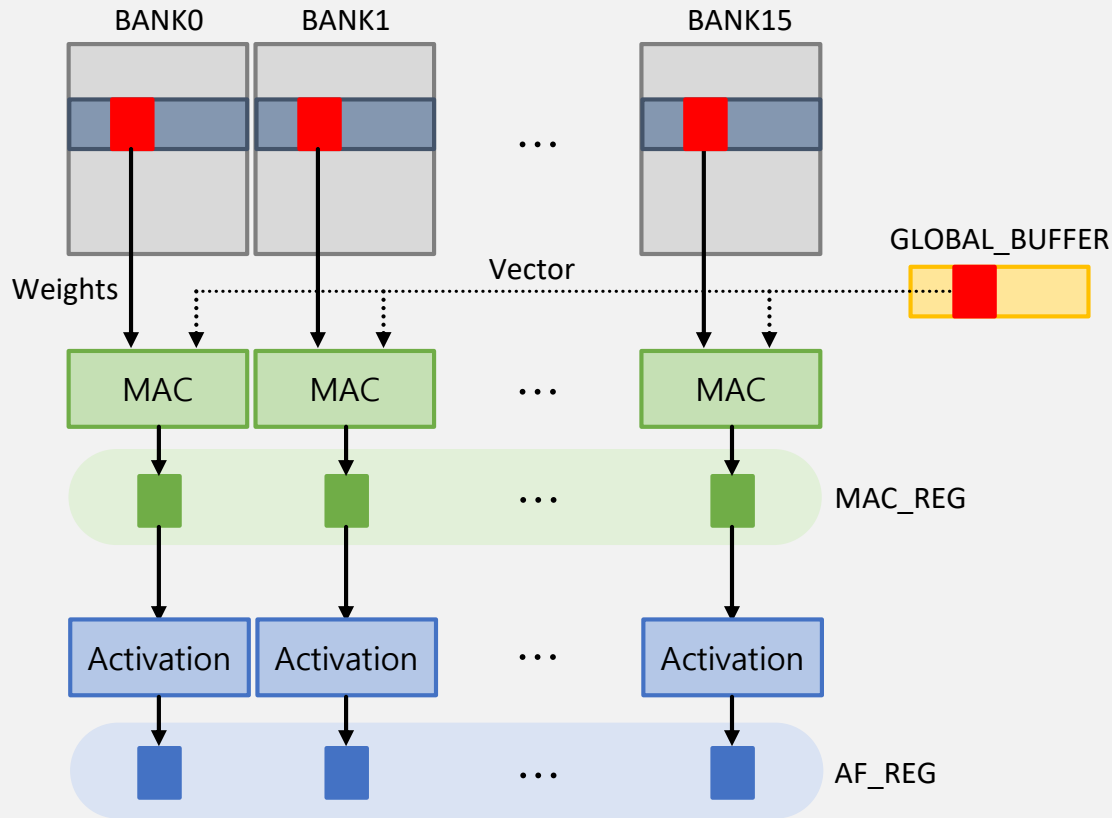| BK8 | BK11 | BK12 | BK15 |
|---|---|---|---|
| MAC | MAC | MAC | MAC |
| Activation | Activation | Activation | Activation |
| Activation | Activation | Activation | Activation |
| MAC | MAC | MAC | MAC |
| BK9 | BK10 | BK13 | BK14 |

## Multiply-And-Accumulate (MAC)

- Performs MAC operation on **sixteen** bfloat16 weight and vector elements (corresponds to a single DRAM column access, i.e. 32 Bytes).
- Computation results are stored in a dedicated **MAC_REG** set and can be later accessed by the user.
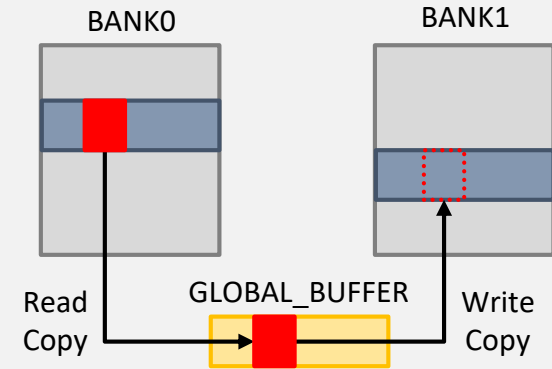


W[0] V[0] W[1] V[1] ... W[14] V[14] W[15] V[15]

**ACCUMULATOR**

MAC_REG

MEMORY FOREST

# GDDR6-AiM Operations
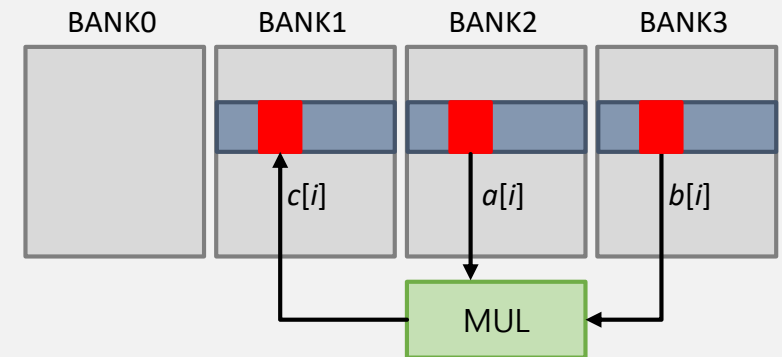
## All-Bank Operations (MAC, Activation Function)



- **MAC** and **Activation Function** operations can be performed in all banks in parallel.
- **Weight** data is sourced from **Banks**; **Vector** data is sourced from the **Global Buffer**.
- **MAC** results are stored in latches collectively referred to as **MAC_REG**.
- **Activation Function** are stored in latches collectively referred to as **AF_REG**.

## In-Channel COPY



- **Global Buffer** acts as FIFO register.
- **Read Copy** fills the FIFO, **Write Copy** transfers FIFO contents to a bank.
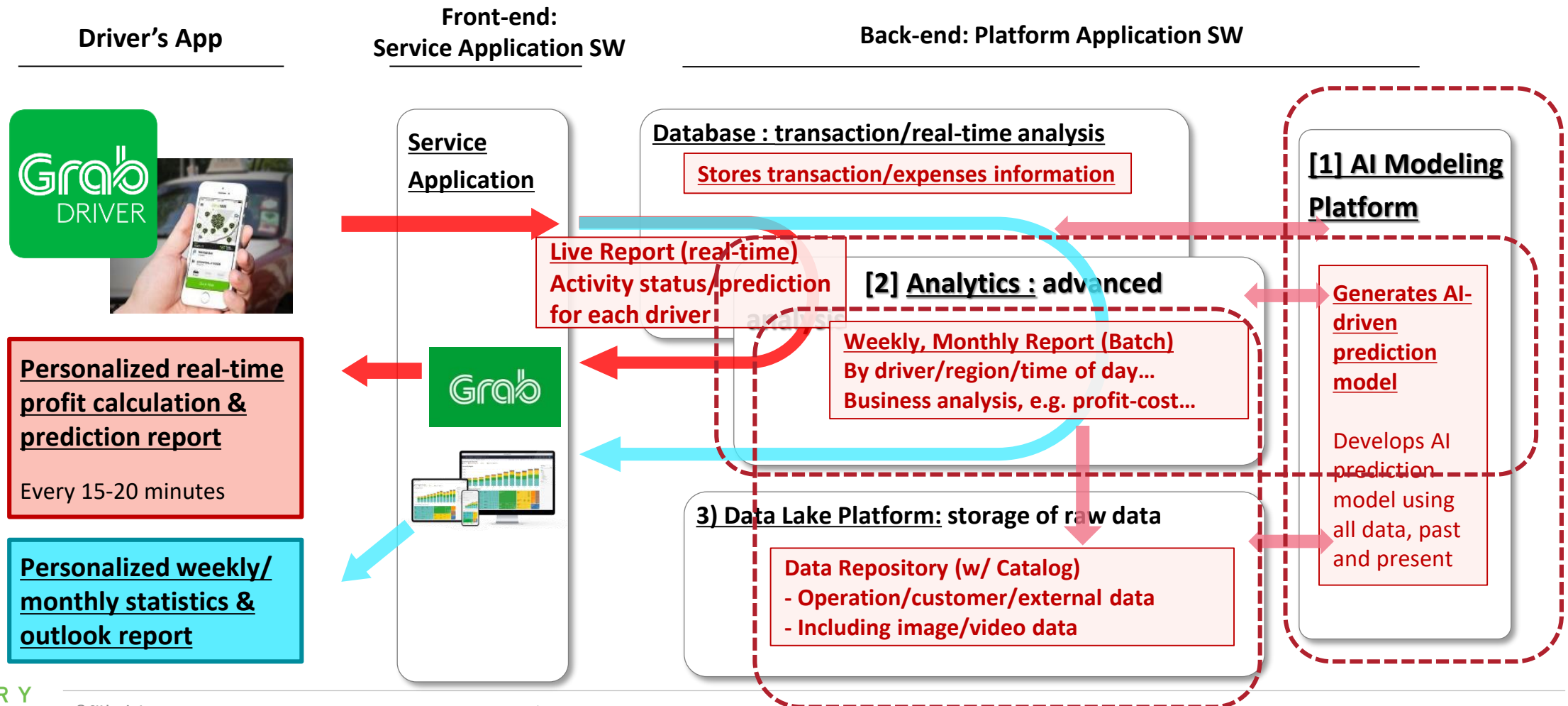
## Element-Wise Multiplication



- Underlying expression: $c[i] = a[i] \cdot b[i]$
- One operation per **Bank Group** can be performed in parallel.

# CONTENTS

SK hynix

MEMORY
FOREST

# AI Service Business Case - Grab Driver

- **Adopt data/analytics/AI platforms due to growing needs for real-time/large-scale/AI-based analysis.**

- **[1] AI Modeling enables prediction, [2] DB/Analytics stores and analyzes key biz data**

**Driver's App**

**Front-end: Service Application SW**

**Back-end: Platform Application SW**

**Service Application**

**Database : transaction/real-time analysis**

**Stores transaction/expenses information**

**[1] AI Modeling Platform**

**Live Report (real-time) Activity status/prediction for each driver**

**[2] Analytics : advanced**

**Generates AI-driven prediction model**

**Weekly, Monthly Report (Batch) By driver/region/time of day… Business analysis, e.g. profit-cost…**

**Personalized real-time profit calculation & prediction report**

Every 15-20 minutes

Develops AI prediction model using all data, past and present

**3) Data Lake Platform: storage of raw data**

**Data Repository (w/ Catalog) - Operation/customer/external data - Including image/video data**

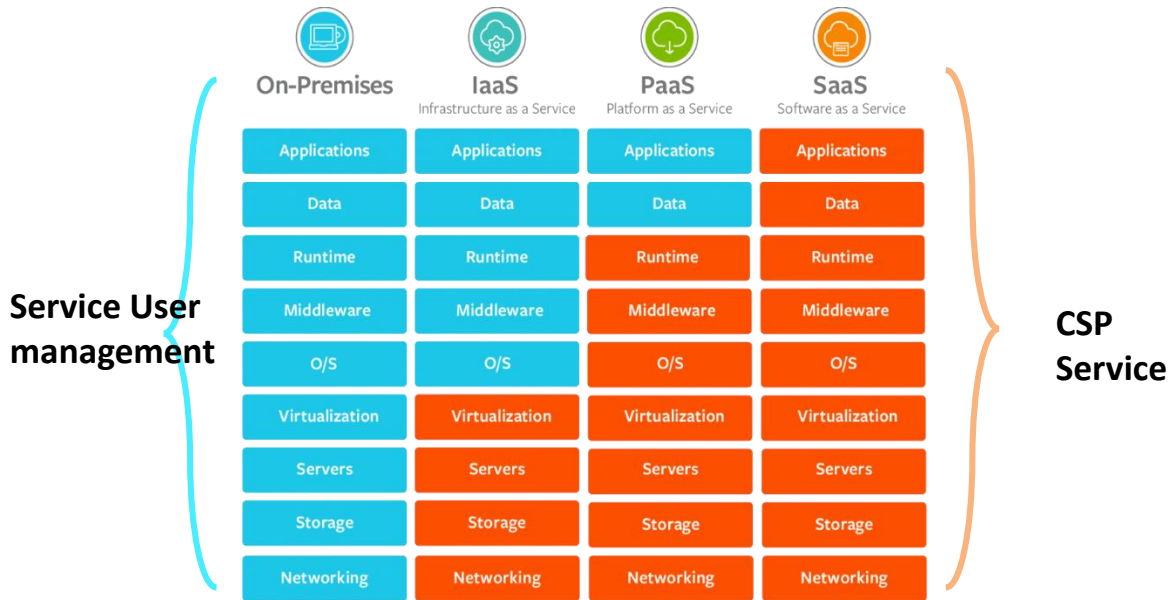**Personalized weekly/ monthly statistics & outlook report**

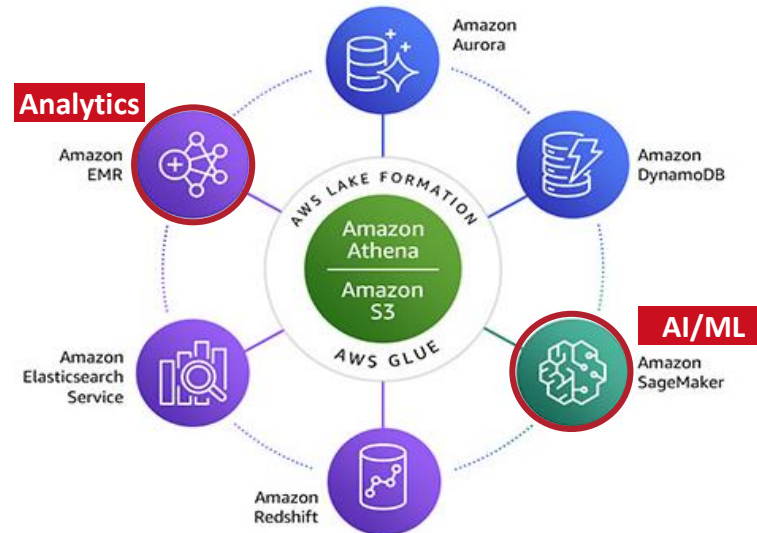# On-premise vs. Cloud (IaaS, PaaS, SaaS)

- On-premise enterprise → Center of gravity shifted to Cloud (IaaS, PaaS, SaaS) from 2010
- The cloud service is enhancing user convenience and strengthening AI/ML and Analytics PaaS

## IaaS, PaaS & SaaS

- **IaaS:** HW resource provision → platform built by Service User
- **PaaS:** + Platform provision → Application development by Service User
- **SaaS:** + Application provision → Service built by Service User



## AI/ML, Analytics PaaS based on Data lake (AWS)



- **[CSP]** Opportunities for service/infrastructure optimization (cost reduction)
- **[Service User]** IT Cost/Business Reduction
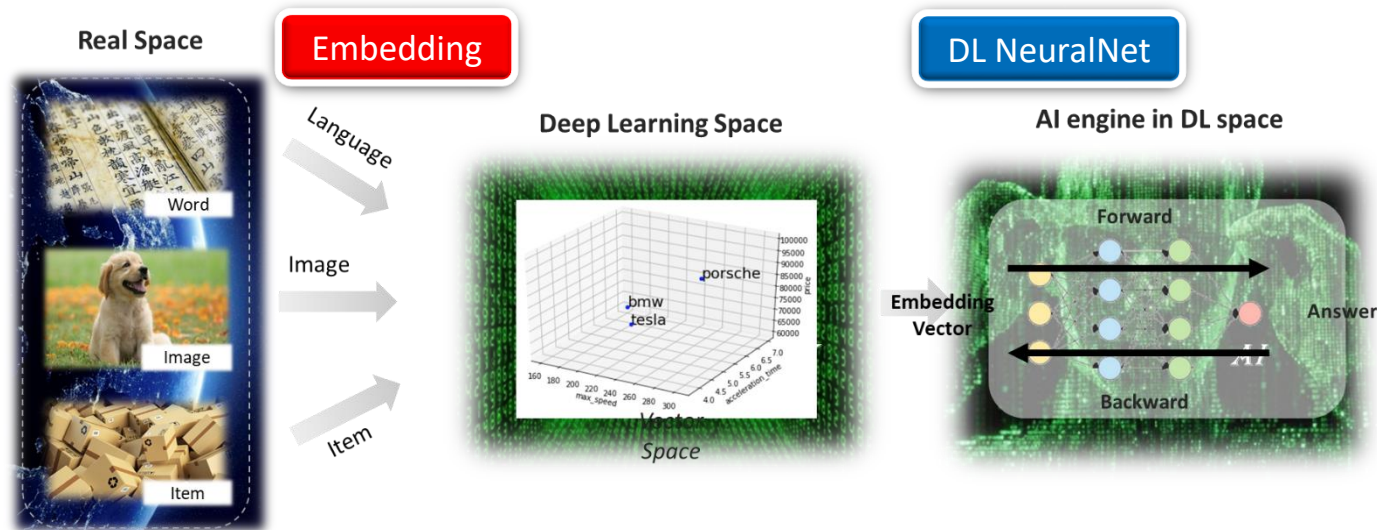- **CSPs understand domain specific customer requirements and workloads**

- **Embedding**

  - To translate Real Space Data into Deep Learning Space Data

  - Memory intensive function

- **DL NeuralNet (Transformer/MLP…)**

  - To do Neural Network in Deep Learning Space

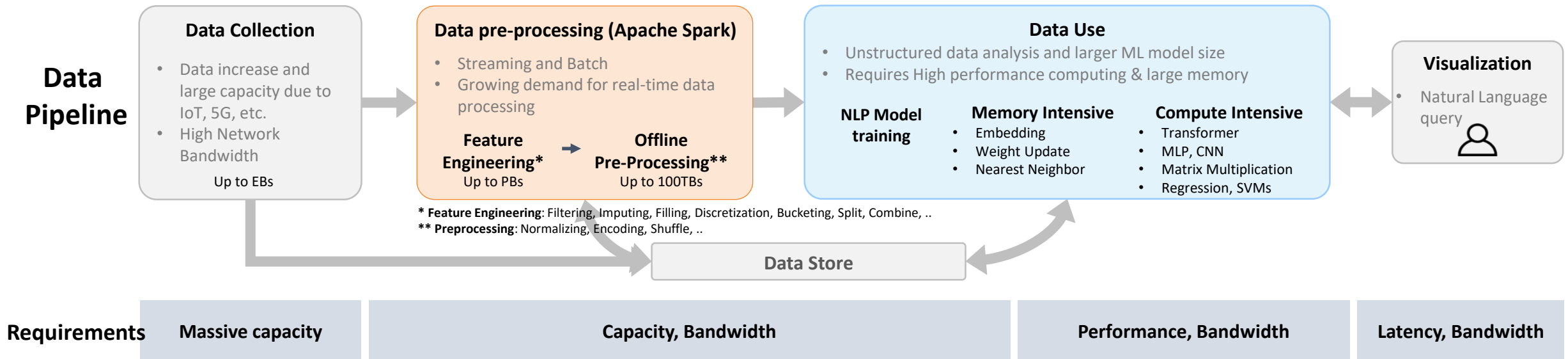  - Memory intensive + Computing intensive function



- **Embedding intensive AI Service: Recommendation**

- **DL NN intensive AI Service: NLP, Vision, …**

# Data pipeline for AI service

- **Real-time data analytics and AI systems are built as pipelines for data processing using AI.**

- **In the future, architectural convergence between AI-Analytics is expected.**

**Data Pipeline**

**Data Collection**
- Data increase and large capacity due to IoT, 5G, etc.
- High Network Bandwidth

Up to EBs

**Data pre-processing (Apache Spark)**
- Streaming and Batch
- Growing demand for real-time data processing

**Feature Engineering*** → **Offline Pre-Processing****
Up to PBs      Up to 100TBs

\* **Feature Engineering**: Filtering, Imputing, Filling, Discretization, Bucketing, Split, Combine, ..
\*\* **Preprocessing**: Normalizing, Encoding, Shuffle, ..

**Data Use**
- Unstructured data analysis and larger ML model size
- Requires High performance computing & large memory

**NLP Model training**

**Memory Intensive**
- Embedding
- Weight Update
- Nearest Neighbor

**Compute Intensive**
- Transformer
- MLP, CNN
- Matrix Multiplication
- Regression, SVMs

**Visualization**
- Natural Language query

**Data Store**

**Requirements**

| Massive capacity | Capacity, Bandwidth | Performance, Bandwidth | Latency, Bandwidth |
|---|---|---|---|

**System ① : Real time Data Analytics**

**System ② : Massive AI**

**Data Analytics**

**Embedding**     **DL NeuralNet**

# CONTENTS

# Back to the Basic – Computational Memory

- **Computational Memory Concept**
  - By performing some host operations on the memory side, energy efficiency and performance are improved
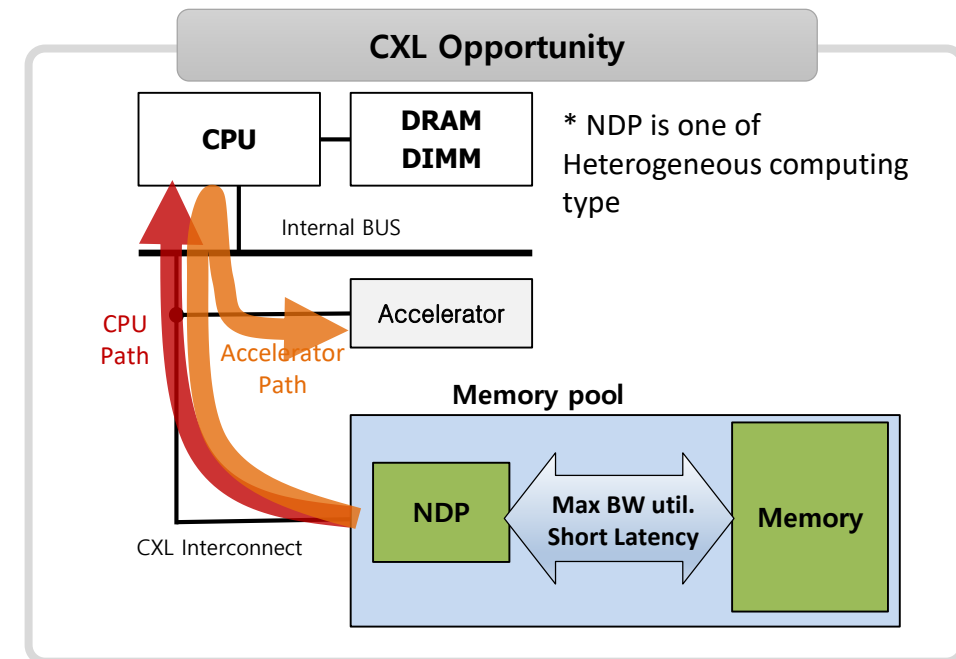


**Conventional**

Host
C

data
bus

Memory
M M M M
M M M M
M M M M
M M M M

**1. Energy efficiency improvement**

By minimizing data movement between host and memory

C : Computing     M : Memory

**2. Performance improvement**

By utilizing the higher BW inside the memory

**CM(Computational Memory )**

Host
C

data
bus

C

bus

Memory
M M M M
M M M M
M M M M
M M M M

***Offloading conditions***

1. ***Low OI (operational intensity)***
2. ***Massive parallelism***
3. ***Data reduction***

Die level computational memory
→ PIM, GDDR6-AiM

: Computing

: Memory

- **CXL interconnect that connects Heterogeneous Computing Elements will become the center of Server System**

  - CXL is an interconnect to support high speed connection between Host Processor and Accelerators/Memory Device.

  - CXL supports 3 protocols based on PCIe Gen5.0

    - 1) CXL.io, 2) CXL.cache, 3) CXL.mem

- **Opportunity: Value added Memory solution available**

  - Unlike conventional DIMMs, CXL-connected memory protocol enables hand-shaking communication, enabling additional functions on memory (ex, DRAM cache, Data processing engine...)

- **With the advent of Memory-intensive Killer Application (AI) and Memory Semantic Interconnect (CXL), research and deploy of CXL memory-based Computational Memory is expected to accelerate.**



CXL Opportunity

CPU — DRAM DIMM

Internal BUS

Accelerator

* NDP is one of Heterogeneous computing type

CPU Path

Accelerator Path

Memory pool

NDP ← Max BW util. Short Latency → Memory

CXL Interconnect

# Card level CM - CMS (CXL base Computational Memory Solution)

- **Higher Performance by fully utilizing Memory BW + energy saving by data reduction + low cost high capacity**
  - Computing core can efficiently handle data-intensive workloads by fully utilizing memory bandwidth in card
  - Data reduction in the cards can significantly improve energy consumption by data movement
  - Cost-effective scalability makes the system to easily scale-up and out without having to pay for expensive servers just to increase the number of memory channels

**CMS-augmented Server**

CXL/PCIe

Computing Core that fully utilizes the internal memory bandwidth

Core

Mem  Mem  Mem  Mem

Mem  Mem  Mem  Mem

Massive memory capacity

Extremely high internal memory bandwidth

Cost-effective scalability

**Die level CM - PIM**

- Bank level parallelism
- Applied within one memory device die
- Small memory capacity per processing node
- Need to define Host interface and standardization

**Card level CM - CMS**

- Channel level parallelism
- Applied across multiple memory devices
- Larger memory capacity per processing node than PIM
- CXL interface available

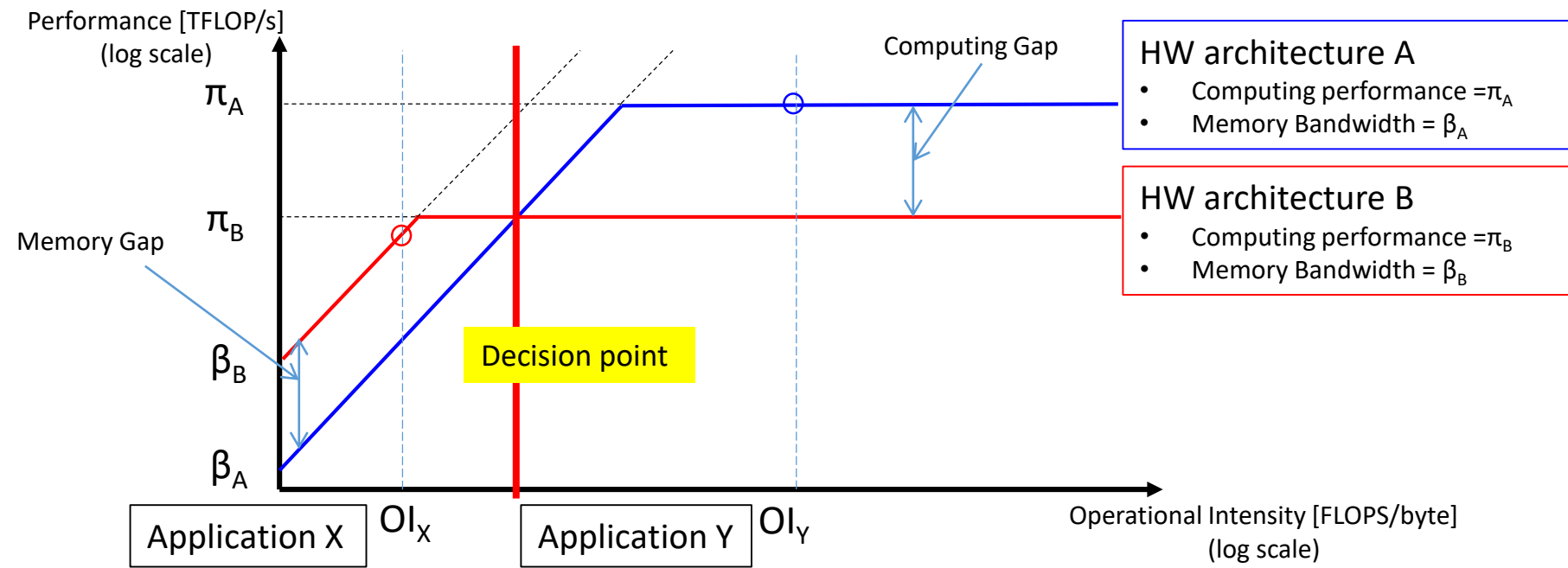**Storage level CM - CSD**

- NAND flash level parallelism
- Applied across multiple NAND flash devices
- Larger memory capacity per processing node than CMS
- Block interface → KV interface

Workload analysis for actual use case

SW framework support such as Compiler, API, Library, device driver

- **Methodology that can analyze HW architecture suitable for a specific processing algorithm**
  - W (Work) : # of operations performed by a given application
  - Q (Memory Traffic) : # of bytes of memory transfers incurred during execution of application
  - OI (Operational Intensity) = W/Q : # of operations per byte of memory traffic.
  - P (Attainable performance) = min ($\pi$, $\beta$ x OI) : In given HW, $\pi$ is max processing performance, $\beta$ is the max bandwidth

- **Representative Data Analytics functions have overall low operational intensity characteristics → Memory bound**

- **The embedding operation is also a memory bound operation with very low operational intensity.**

- **So, if these operations are operated in the Computational Memory, performance and power gains can be obtained.**

### Roofline Analysis

# Workload Analysis – DL Neural network

- **Matrix multiplication, which has been the main target of PIM, is losing its memory-intensive characteristics as the batch increases and algorithm evolves.**

- **There is still an opportunity for offloading of memory-intensive functions regardless of batch size increase such as layer normalization and any kind of function for data itself.**

## Batch vs. OI

- **GEMV (related w/ weight)**

  small batch → memory intensive, large batch → computing intensive

- **Normalization, Optimizer (related w/ data)**

  Even in large batch → memory intensive



GEMV

Vectors    Matrix (Reused)

Layer Normalization

batch / dim / mean / std

Batch vs. OI

Optimizer / Norm. / batch / GEMV / OI

## Low OI Layers in Transformer

- **In Transformer, Memory intensive operation takes significant portion of the total execution time**

- **There is also memory intensive operation in Attention layer (softmax, biases, dropout)**
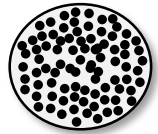


Add & Normalize

Feed Forward    Feed Forward

Add & Normalize

Self-Attention

**Single Transformer Block**  *Source: https://arxiv.org/abs/2007.00072

LayerNorm( + )

**Memory intensive layers**

Softmax( )

© SK hynix Inc. This material is proprietary of SK hynix Inc. and subject to change without notice./ Confidential

# Workload Analysis – Workload density

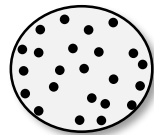**Workload Density – How dense is the data to be processed in the memory**

→ **If the workload is sparse compared to the memory capacity per PE, data reduction per PE is reduced and frequent data movement between PEs is induced.**

High workload density

Deep learning NN
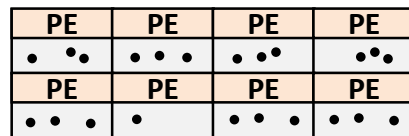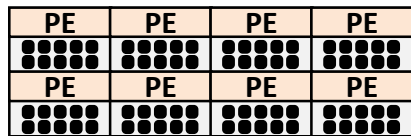(MLP, CNN, Transformer…)
Data Analytics

Low workload density
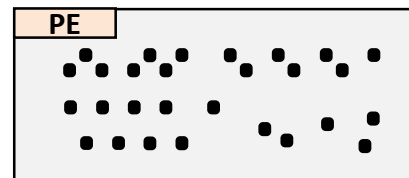
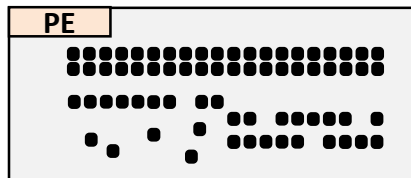Embedding (DLRM)
Graphs
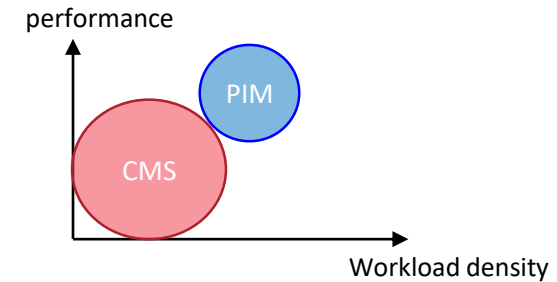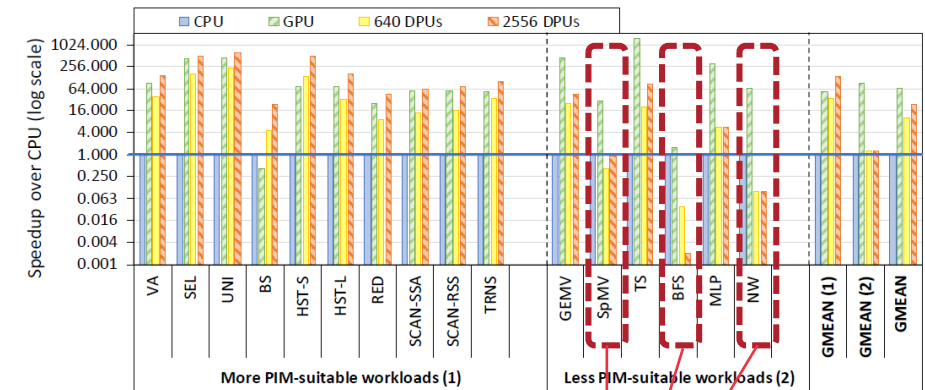Big Data/Data Analytics

PIM
Die level

DRAM Die

CMS
Card level

Memory Card

**The higher workload density, the more appropriate for PIM**

Gómez-Luna, Juan, et al. "Benchmarking a New Paradigm: An Experimental Analysis of a Real Processing-in-Memory Architecture."(2021)
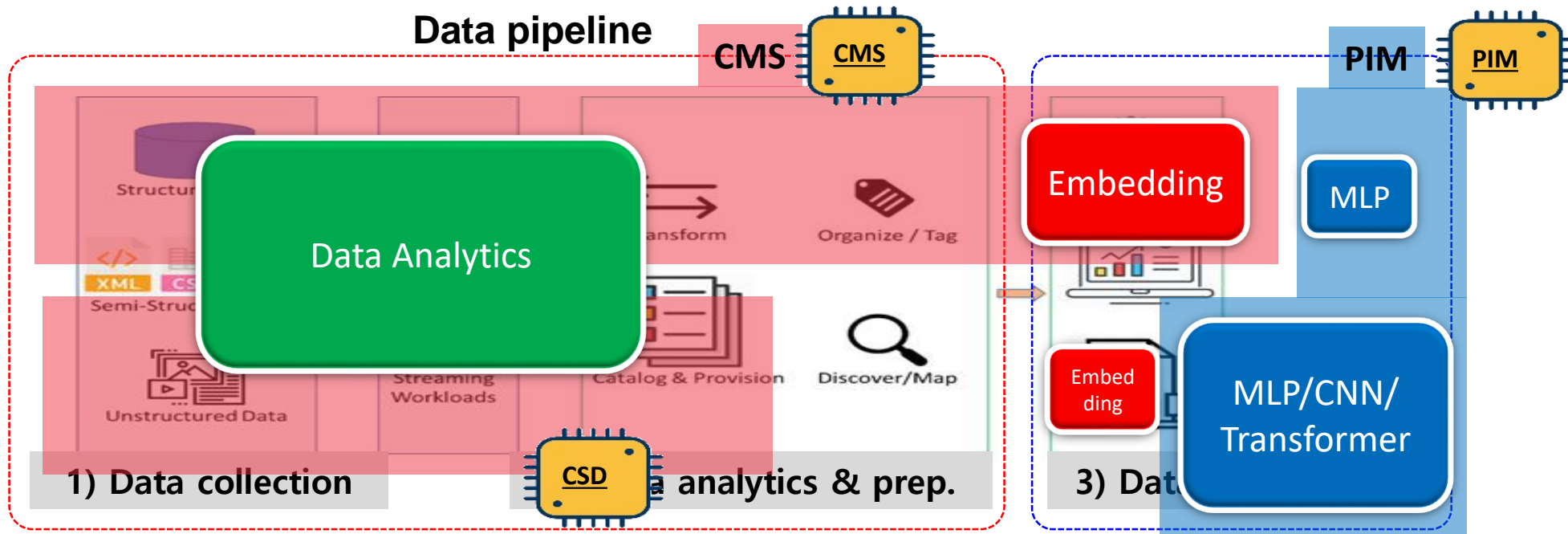
Sparse Data, large data movement among processing node.

# CONTENTS

MEMORY
FOR EST

# Summary

● **Value addition memory base solution can be deployed in whole data pipeline of AI/Data Analytics system**

- Die level CM : PIM

- Card level CM : CMS

- Storage level CM : CSD



**Data pipeline**

CMS · **CMS**   PIM · **PIM**

Data Analytics

Embedding   MLP

Embed ding   MLP/CNN/ Transformer

Structur...   ...ansform   Organize / Tag

XML   C...

Semi-Struc...

Unstructured Data   Streaming Workloads   Catalog & Provision   Discover/Map

1) Data collection   CSD   ...a analytics & prep.   3) Dat...

# *End of Document*