

# Graph Neural Networks Acceleration with Adaptive Dataflow Architectures & FPGA

Advanced Processor Technologies Research Group

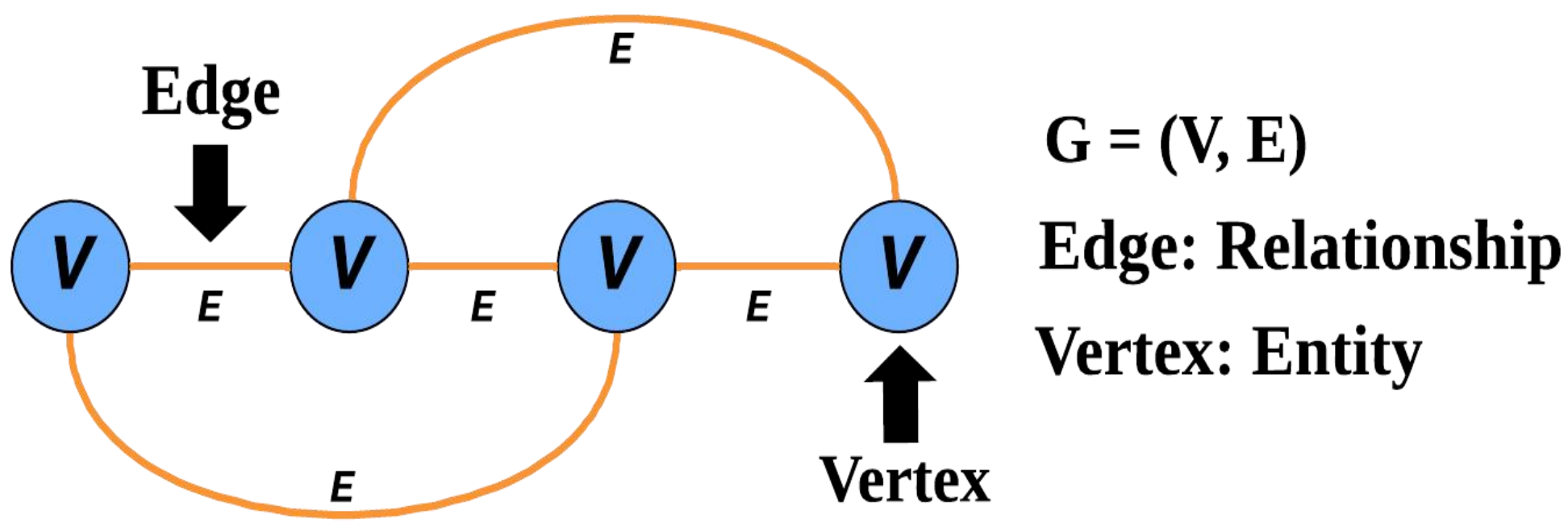
Zheyu Liu

Supervisor: Christos-efthymios Kotselidis

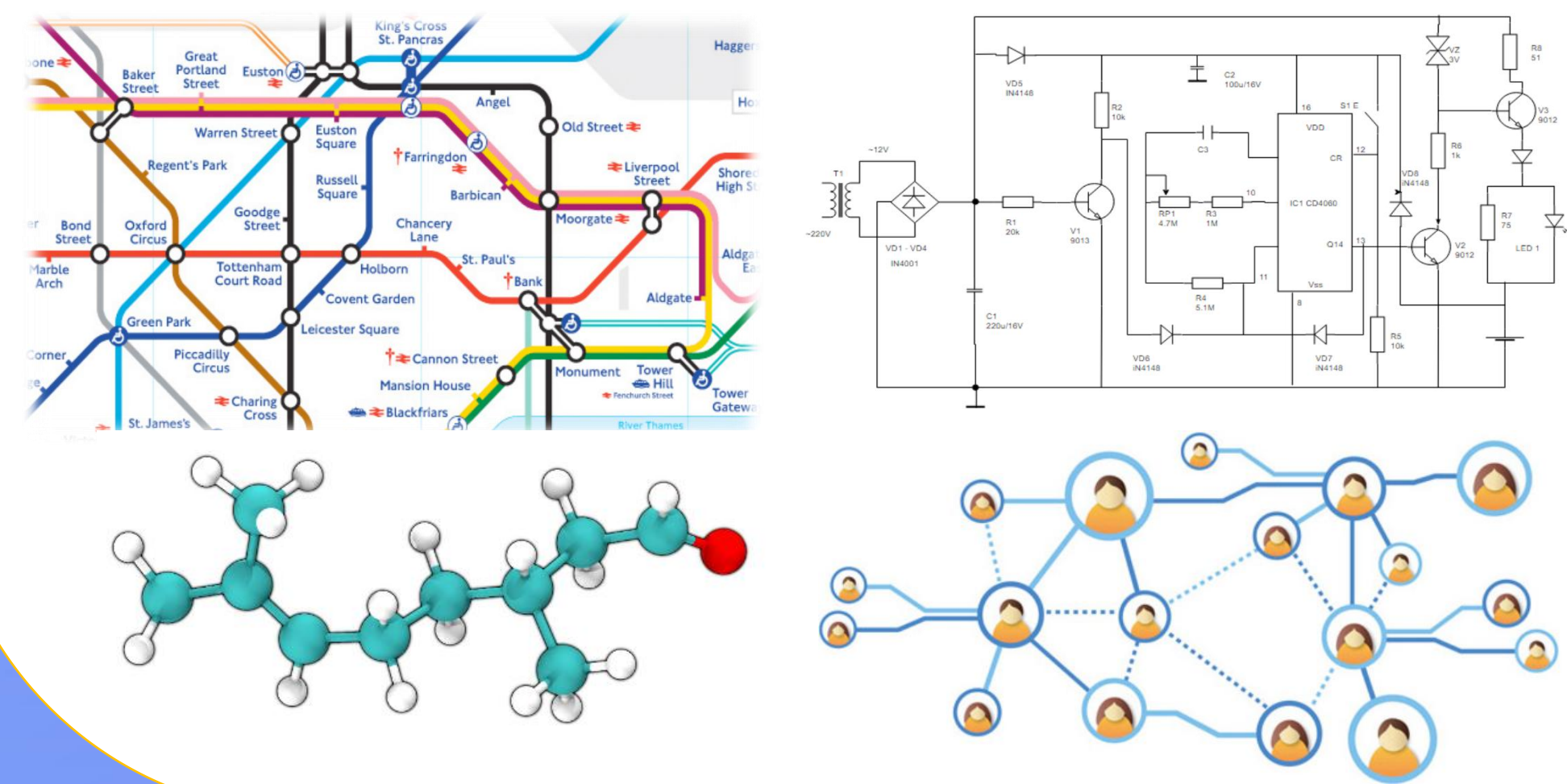
## I. Introduction

### What is Graph?

Graph is a non-Euclidean data structure used to represent complex relationships between entities, consisting of numerous vertices and edges.



### Graph Examples



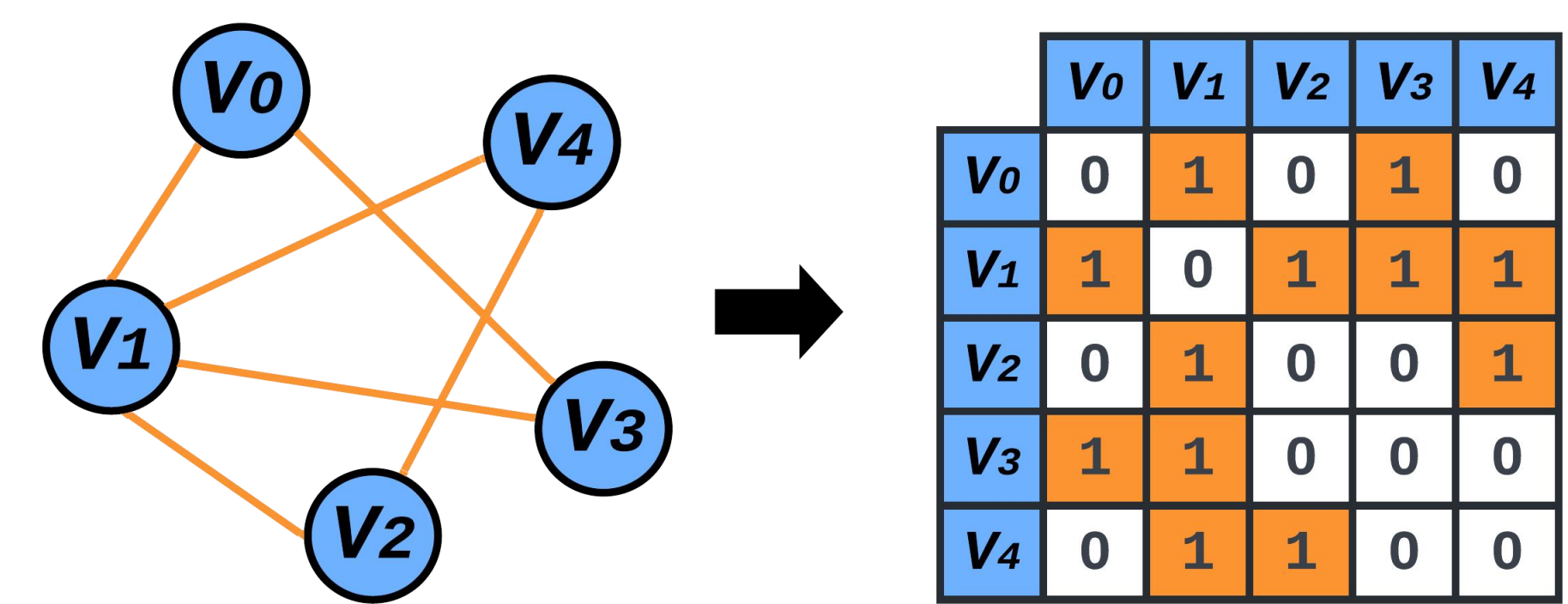
### Graph in Memory

Vertex in a graph contains its unique information, such as a person's height, weight and hobbies. All the features of a vertex can be represented as a numerical **Embedding Vector** with  $N$  elements.

**Embedding Vector**

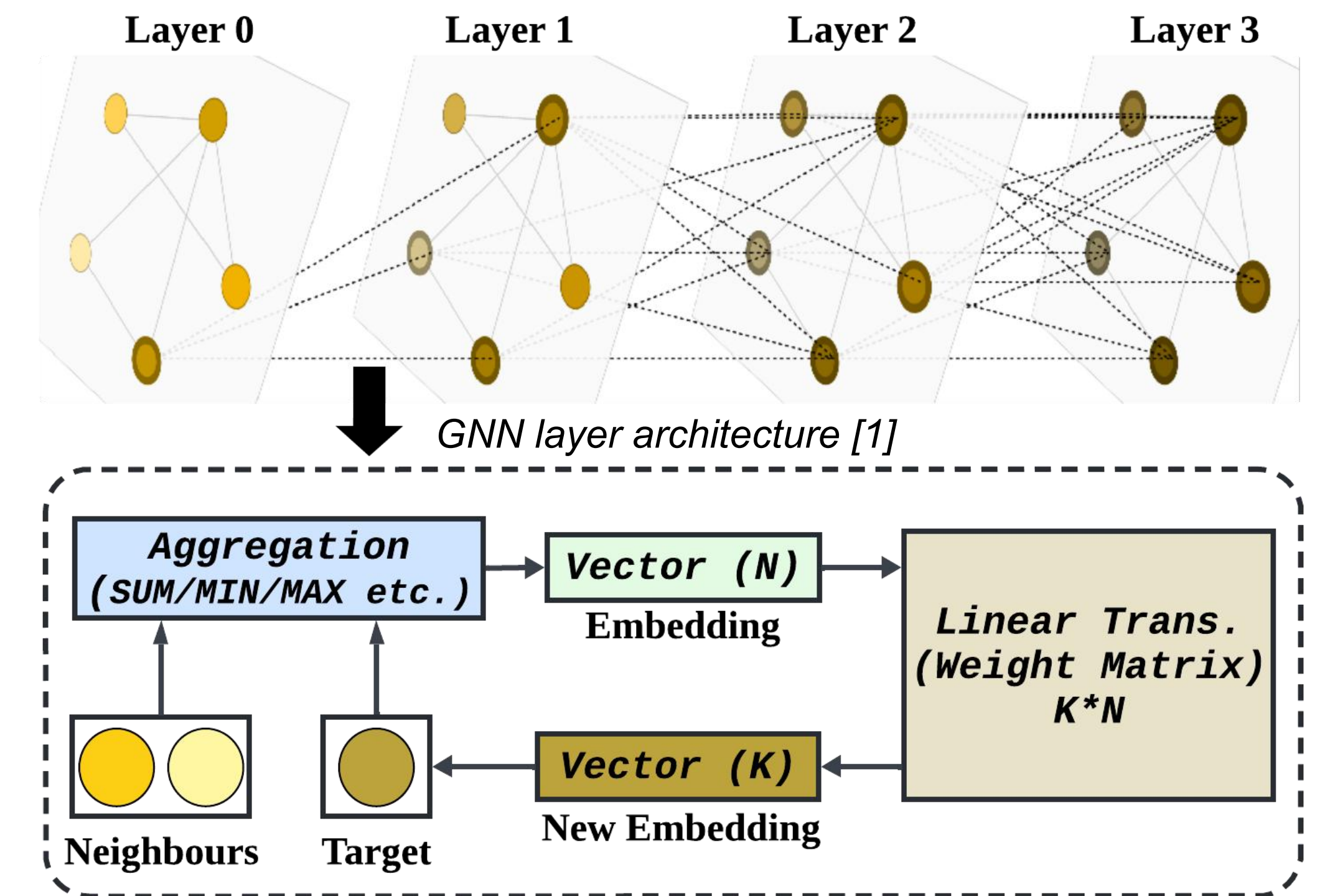
$V_0$	$V_0 F(0)$	$V_0 F(1)$	$V_0 F(2)$	...	$V_0 F(N-1)$
$V_1$	$V_1 F(0)$	$V_1 F(1)$	$V_1 F(2)$	...	$V_1 F(N-1)$
$V_2$	$V_2 F(0)$	$V_2 F(1)$	$V_2 F(2)$	...	$V_2 F(N-1)$
...	...	...	...	...	...
$V_{M-1}$	$V_{M-1} F(0)$	$V_{M-1} F(1)$	$V_{M-1} F(2)$	...	$V_{M-1} F(N-1)$

Connectivity information of vertices is usually stored in form of an **Adjacency Matrix**.



### Apply Neural Network to Graph

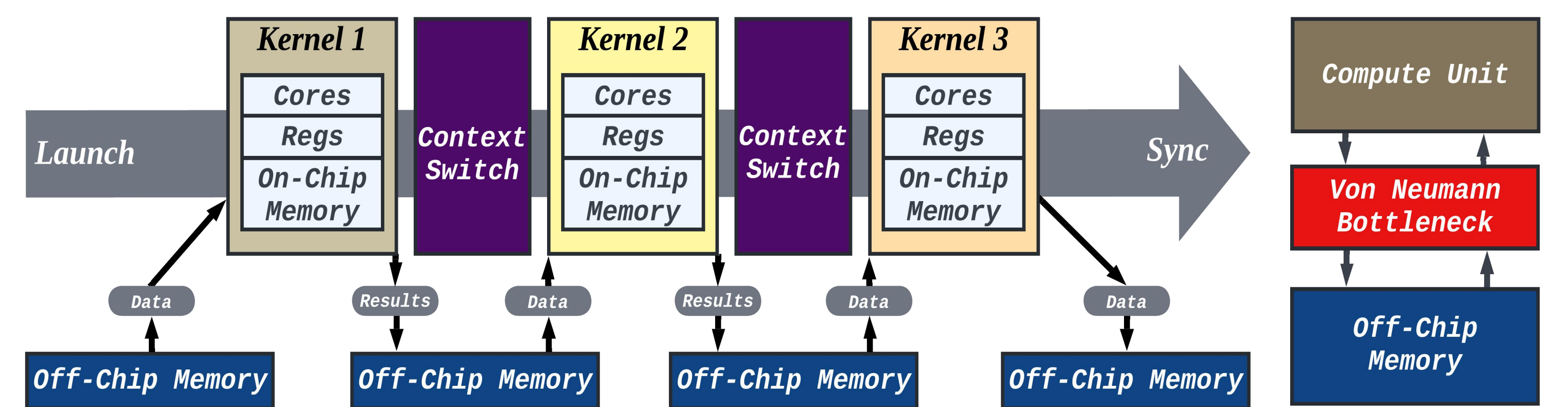
The message passing framework establishes the data flow in a graph neural network (GNN) which primarily comprises **Aggregation** and **Linear Transformation**.



GNNs are showing promising applications in diverse domains such as traffic prediction, fraud detection and drug discovery.

## II. Problems of Running GNNs on GPU

GPUs are equipped with cache and very high bandwidth off-chip memory to alleviate the Von Neumann bottleneck, enabling faster data transfer between the processing units (core) and memory.



Operators	L2 Hit % (A100)
torch.gather	62.75
torch.index_add_	82.09
torch.index_select	70.02
torch.matmul( $\leq 10^8$ el.)	88.58
torch.matmul( $> 10^8$ el.)	89.56
torch.transpose	93.59

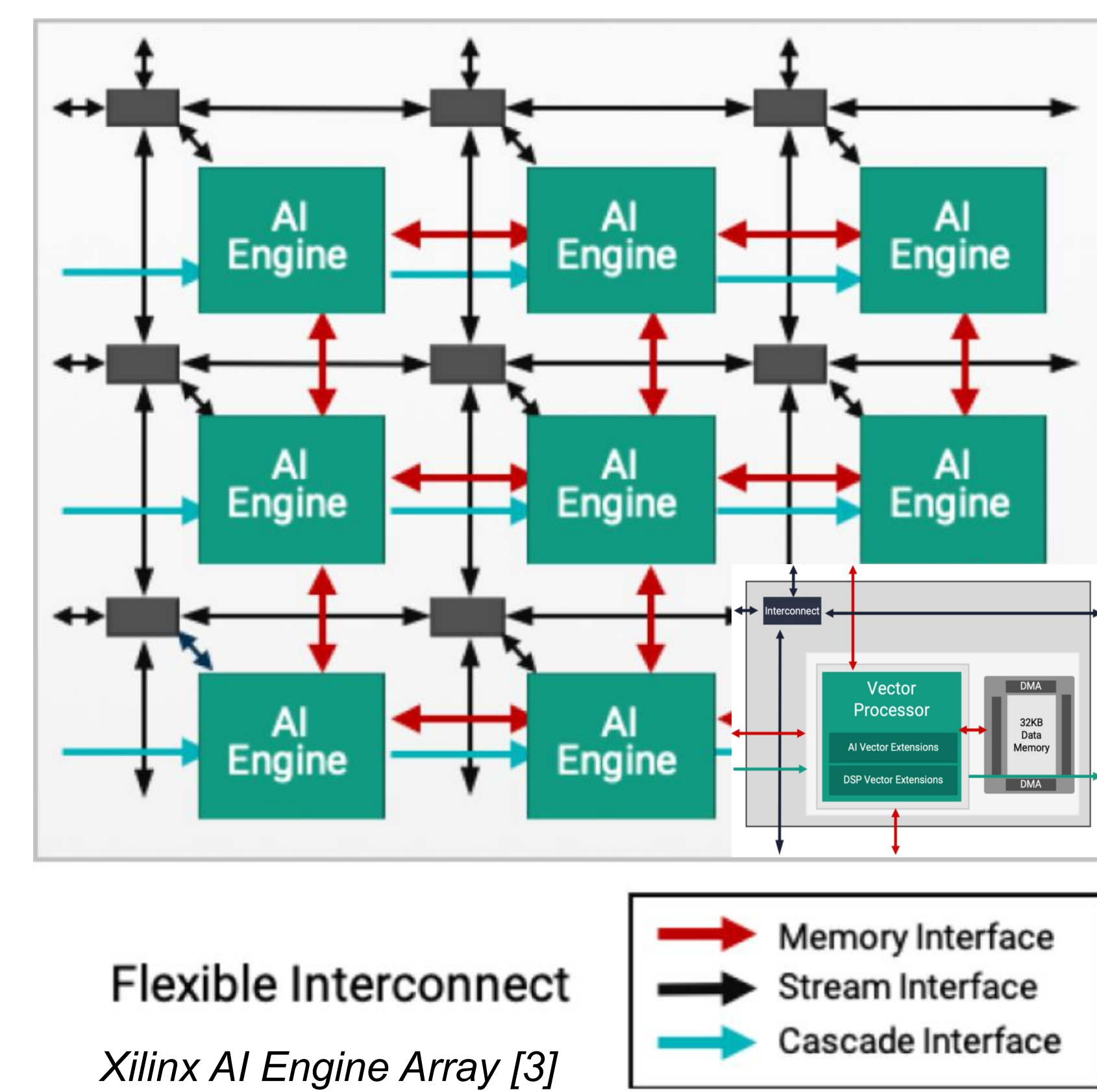
A100 Cache Hit Rate on ML Micro-operations [2]

Due to frequent random memory accesses, the cache hit rate for operations related to GNN is significantly reduced.

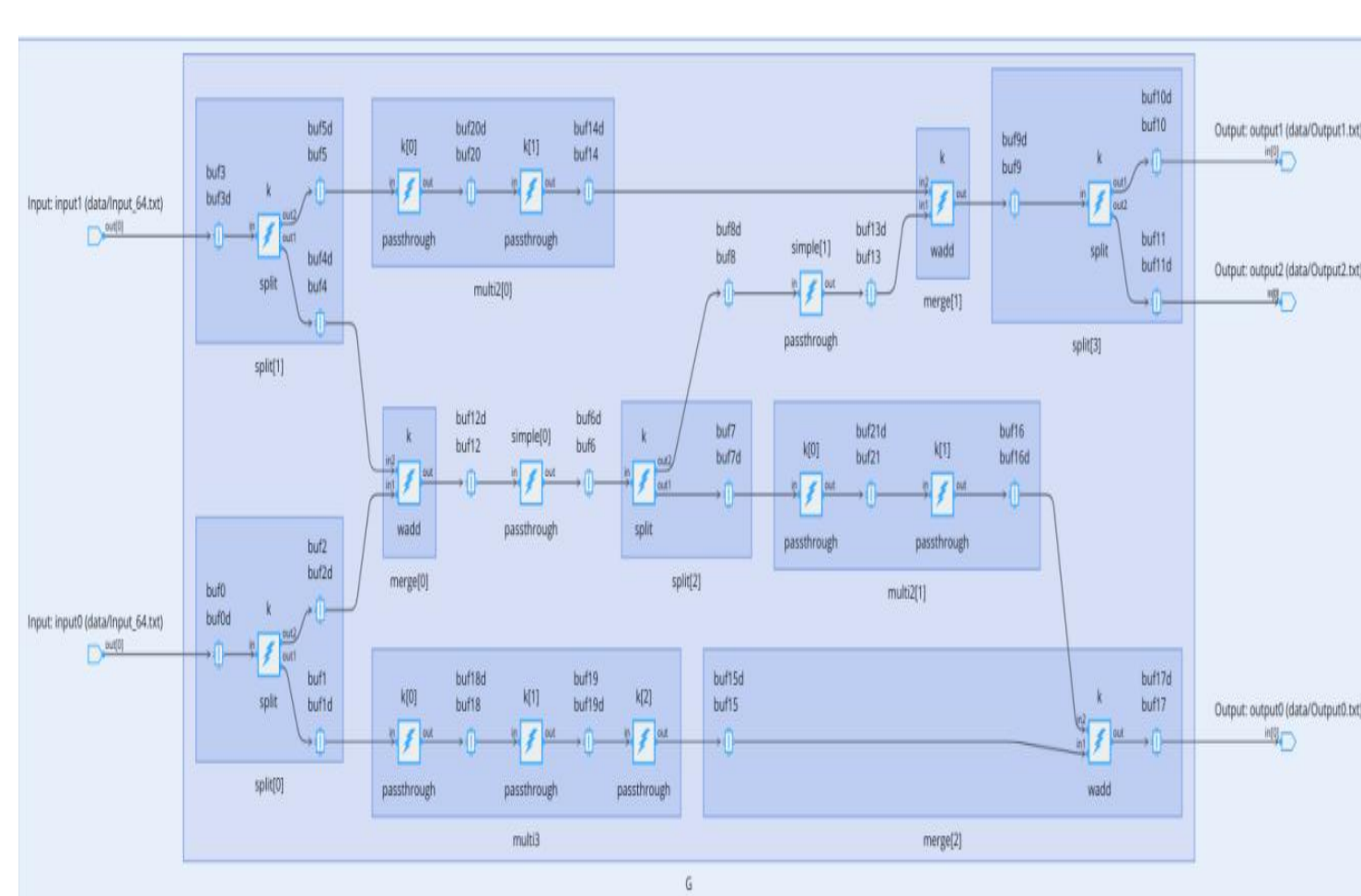
- **Sparse Connection Among Vertices**
- **Irregular Memory Access Pattern**
- **Unbalanced Workload**
- **Complex Vertex Dependency**
- **Hard to Batching**

## III. GNNs on Reconfigurable Architecture Acceleration

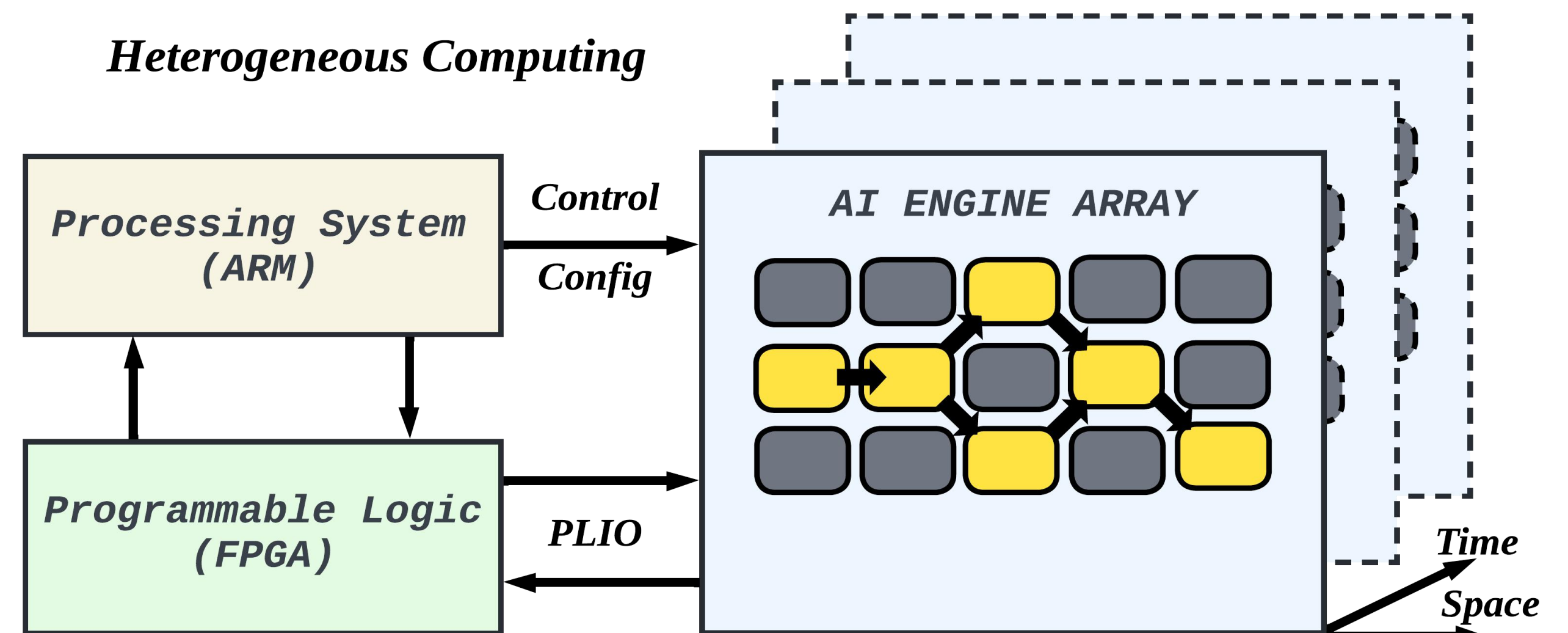
### Xilinx Adaptive Dataflow Architecture



- **Data-Driven Computing**
- **Implicit Parallelism**
- **Distributed On-Chip Memory**
- **Reconfigurable SIMD Pipeline**
- **Graph based Programmability**



### GNN Adaptive Acceleration



FPGA is a fine-grained reconfigurable architecture capable of efficiently performing data preprocessing tasks and implementing specialized hierarchical caches to hide memory access latency.

The PS unfolds GNN computation along both temporal and spatial dimensions, dynamically configuring new computation graphs.