

CR²: Community-aware Compressed Regular Representation for Graph Processing on a GPU

Shinnung Jeong¹, Sungjun Cho², Yongwoo Lee¹, Hyunjun Park¹, Seonyeong Heo³, Gwangsun Kim²,
Yongsok Kim¹, and Hanjun Kim¹

*Yonsei University*¹, *Postech*², *Kyung Hee University*³

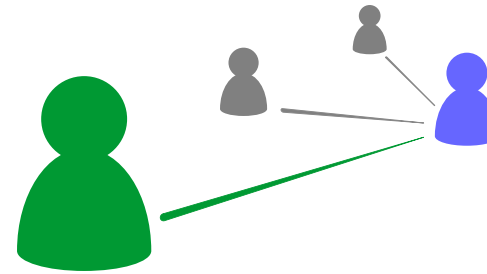


Graphs are One of the Important Data Structures to Abstract and Analyze Real-World Data



Google 검색 또는 URL 입력

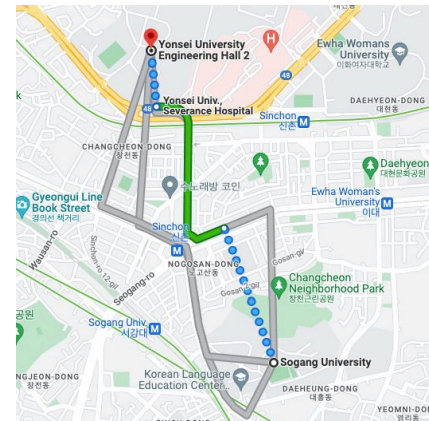
Web Search



Social Network



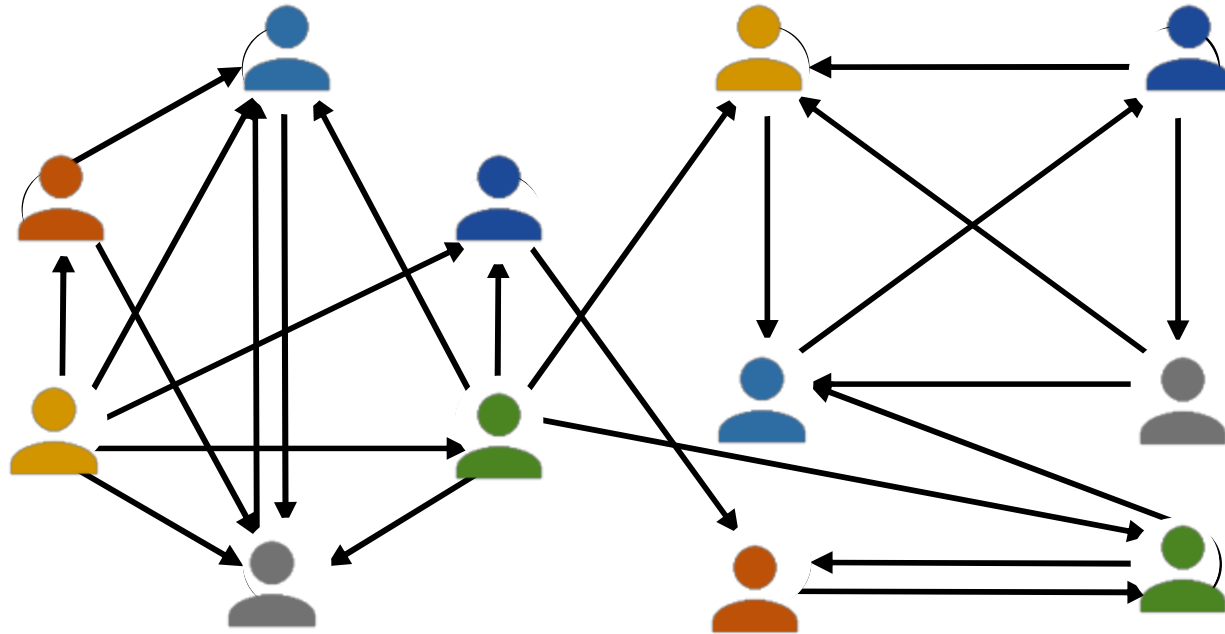
Neuroscience



Path Finding

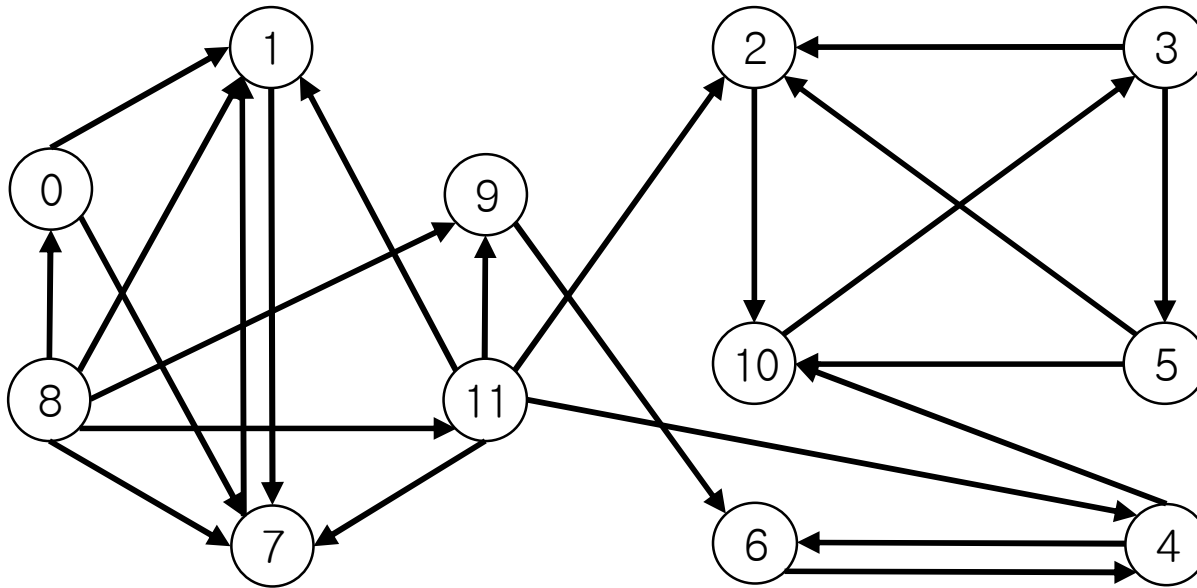
What is Graph?

- An abstract data structure with vertices and their pairs(edges)
- Graph = (Vertex, Edge)



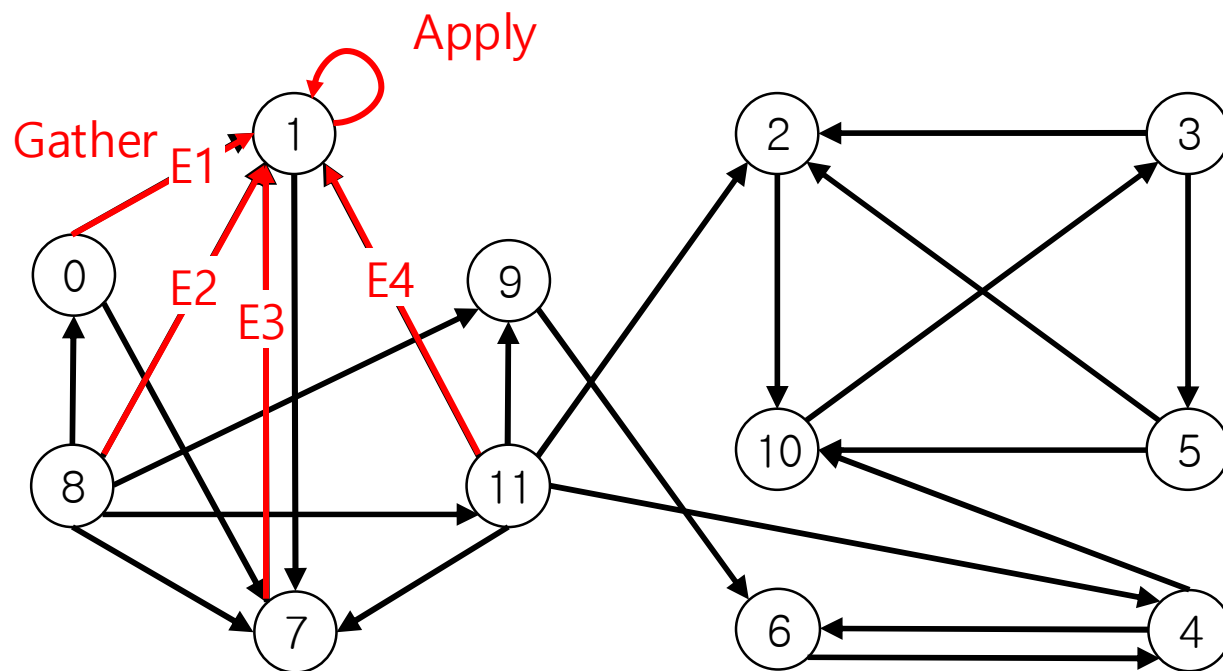
What is Graph?

- An abstract data structure with vertices and their pairs(edges)
- Graph = (Vertex, Edge)



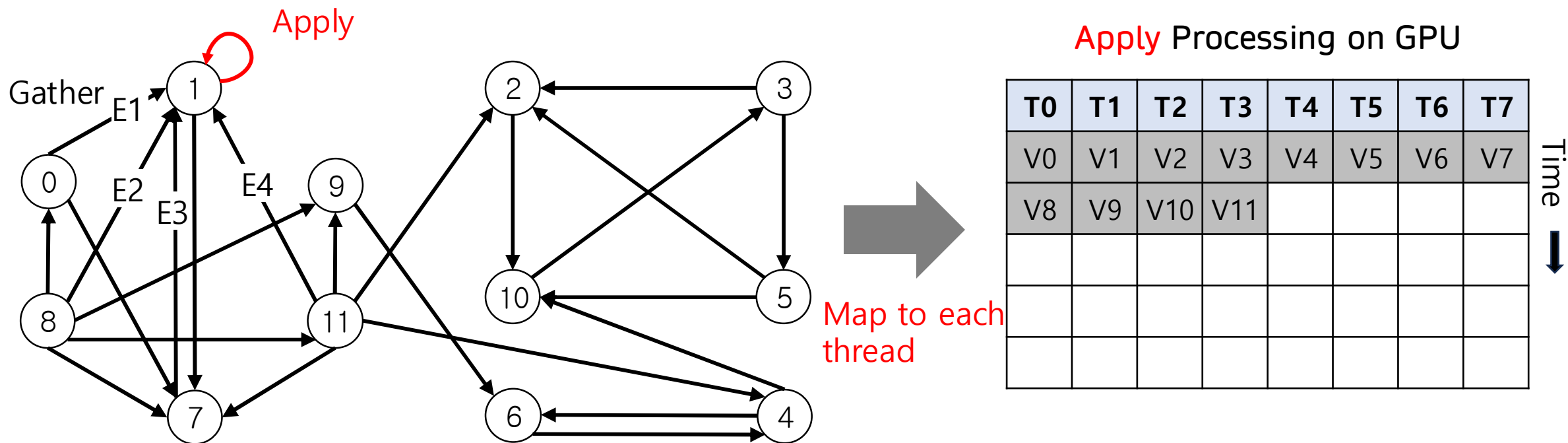
What is the Graph Processing?

- Update vertex data by **gathering** neighbor edge data
- **Apply** computation on vertex data after finishing gathering
- Perform the same algorithm to each vertex and to each edge



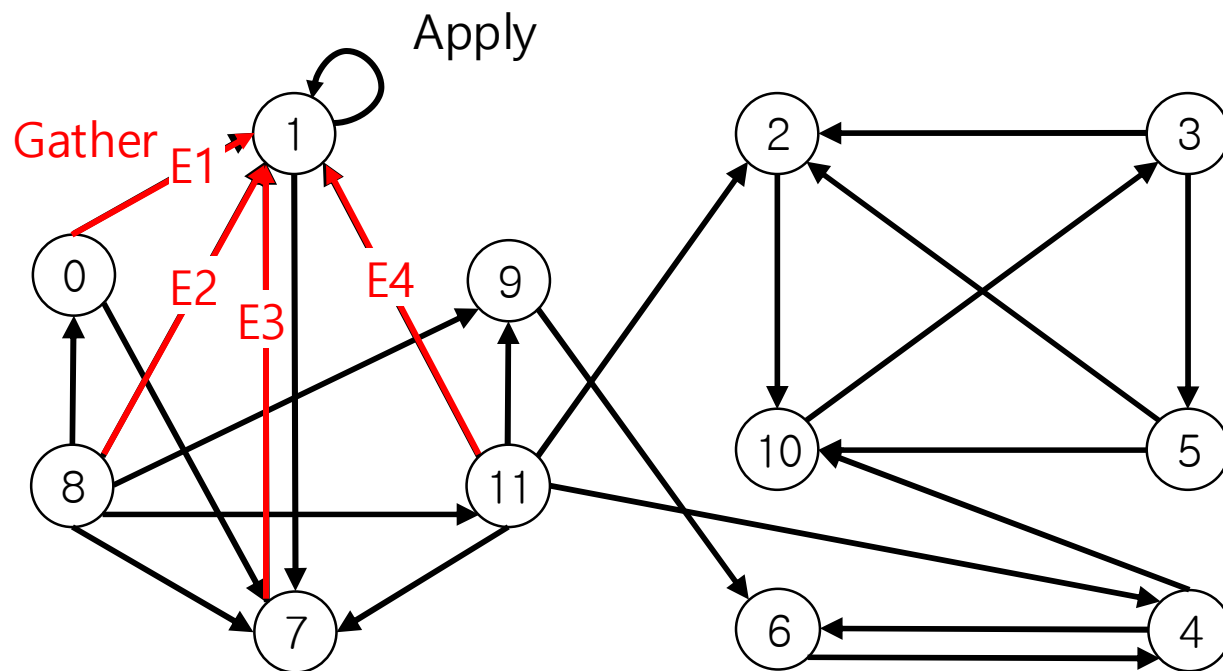
GPU is a Promising Platform for Graph Processing

- Real-world graphs become larger and apply the same algorithms to each vertex and edges
- GPU is designed for Single Instruction Multiple Data(SIMD) processing



However, Graph Processing on GPU is Still Challenging!

- The characteristics of real-world graph are not fit for GPU, such as skewness and sparsity
- Cause irregular distribution inducing **workload imbalance** and reducing resource utilization



Gather Processing on GPU

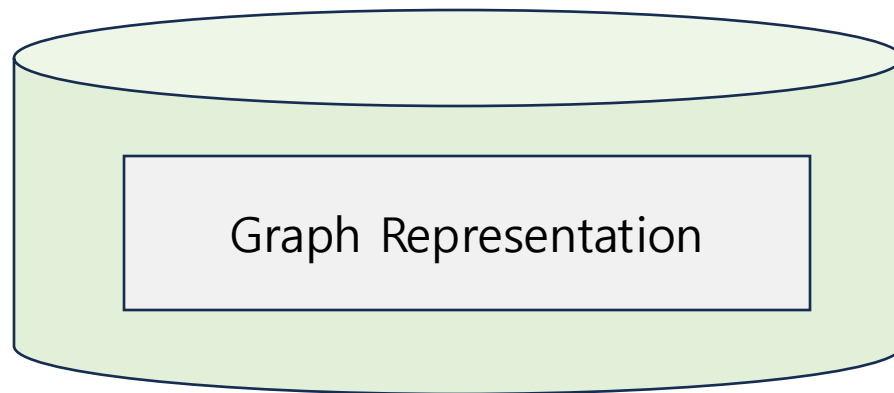
T0	T1	T2	T3	T4	T5	T6	T7
V0	V1	V2	V3	V4	V5	V6	V7
E0	E1	E5	E8	E9	E11	E12	E14
	E2	E6		E10		E13	E15
	E3	E7					E16
	E4						

Time ↓

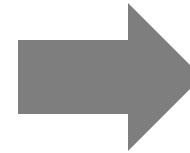
Map each vertex and its neighboring edges to each thread

Changing Graph Representation for GPU can be a Solution

- The graph representation defines how to store a given graph in GPU memory
- Since the graph is large, the graph representation affect to memory usage and perf



Read graph data
from memory



Gather Processing on GPU

T0	T1	T2	T3	T4	T5	T6	T7
V0	V1	V2	V3	V4	V5	V6	V7
E0	E1	E5	E8	E9	E11	E12	E14
	E2	E6		E10		E13	E15
	E3	E7					E16
	E4						

Time ↓

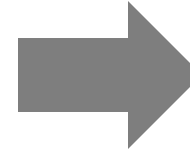
Changing Graph Representation for GPU can be a Solution

- CSR can reduce memory usage and enable neighbor list friendly access

dest id	0	1	2	3	4	5	6	7	8	9	10	11	end											
offset	0	1	5	8	9	11	12	14	18	18	20	23	24											
edge list	8	0	7	8	11	3	5	11	10	6	11	3	4	9	0	1	8	11	8	11	2	4	5	8

Compressed Sparse Row (CSR) ^[1]

Neighbor list access friendly representation



Gather Processing on GPU

T0	T1	T2	T3	T4	T5	T6	T7
V0	V1	V2	V3	V4	V5	V6	V7
E0	E1	E5	E8	E9	E11	E12	E14
	E2	E6		E10		E13	E15
	E3	E7					E16
	E4						

Time ↓

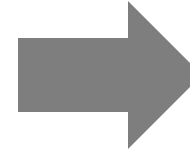
Changing Graph Representation for GPU can be a Solution

- CSR can reduce memory usage and enable neighbor list friendly access
- Graph representation affects **workload imbalance, locality, and memory access pattern**

dest id	0	1	2	3	4	5	6	7	8	9	10	11	end											
offset	0	1	5	8	9	11	12	14	18	18	20	23	24											
edge list	8	0	7	8	11	3	5	11	10	6	11	3	4	9	0	1	8	11	8	11	2	4	5	8

Compressed Sparse Row (CSR) ^[1]

Neighbor list access friendly representation



Gather Processing on GPU

T0	T1	T2	T3	T4	T5	T6	T7
V0	V1	V2	V3	V4	V5	V6	V7
E0	E1	E5	E8	E9	E11	E12	E14
	E2	E6		E10		E13	E15
	E3	E7					E16
	E4						

Time ↓

Existing Graph Representation for GPU: G-shards

Shard 0						
src id	0	3	3	0	1	2
dest id	1	2	4	7	7	10
src value	v_1	v_2	v_4	v_7	v_7	v_{10}

Shard 1						
src id	7	4	6	5	4	5
dest id	1	2	5	6	10	10
src value	v_1	v_2	v_5	v_6	v_{10}	v_{10}

Shard 3												
src id	8	8	11	11	10	11	9	8	11	8	11	8
dest id	0	1	1	2	3	5	6	7	7	9	9	11
src value	v_0	v_1	v_1	v_2	v_3	v_5	v_6	v_7	v_7	v_9	v_9	v_{11}

Rep	Edge Rep	Value	
		Vertex	Edge
CSR	$ E + V + 1$	$ V $	$ E $
G-Shards	$2 E $	$ V + E $	$ E $

G-Shards^[2]

- Achieve balanced execution
- Improve vertex value access locality
- **Cause more memory usage** than CSR

[2] R. Gupta F. Khorasani et al. 2014. CuSha: Vertex-centric Graph Processing on GPUs. HPCA'14

Existing Graph Representation for GPU: Tigr

dest id	0	1	1	2	2	3	4	5	6	7	7	8	9	10	10	11	end							
offset	0	1	3	5	7	8	9	11	12	14	16	18	18	20	22	23	24							
edge list	8	0	7	8	11	3	5	11	10	6	11	3	4	9	0	1	8	11	8	11	2	4	5	8

Tigr^[3]

- Reduce imbalance by vertex split
- Still remain balancing opportunity
- Cause more memory usage
- Do not address vertex value access locality

Rep	Edge Rep	Value	
		Vertex	Edge
CSR	$ E + V + 1$	$ V $	$ E $
G-Shards	$2 E $	$ V + E $	$ E $
Tigr	$ E + (2 + Ft) V $	$ V $	$ E $

[3] Q. Junqiao N. Sabet et al. 2018. Tigr: Transforming Irregular Graphs for GPU-Friendly Graph Processing. ASPLOS'18

Existing Graph Representation for GPU: Tigr

dest id	0	1	1	2	2	3	4	5	6	7	7	8	9	10	10	11	end
offset	0	1	3	5	7	8	9	11	12	14	16	18	18	20	22	23	24

Existing graph representation for GPUs fails to **balance workload** and to improve **locality** while **reducing memory size**

- Still remain workload imbalance
- Cause more memory usage
- Do not address vertex value access locality

Tigr	$ E + (2 + Ft) V $	$ V $	$ E $
------	---------------------	-------	-------

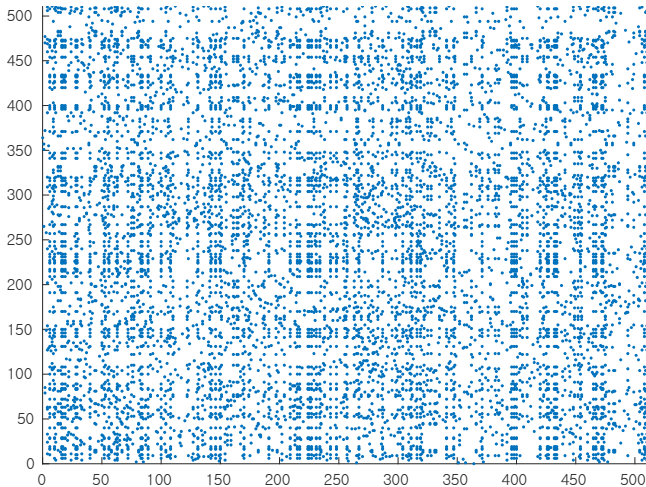
[3] Q. Junqiao N. Sabet et al. 2018. Tigr: Transforming Irregular Graphs for GPU-Friendly Graph Processing. ASPLOS'18

Need a **New Graph Representation!**

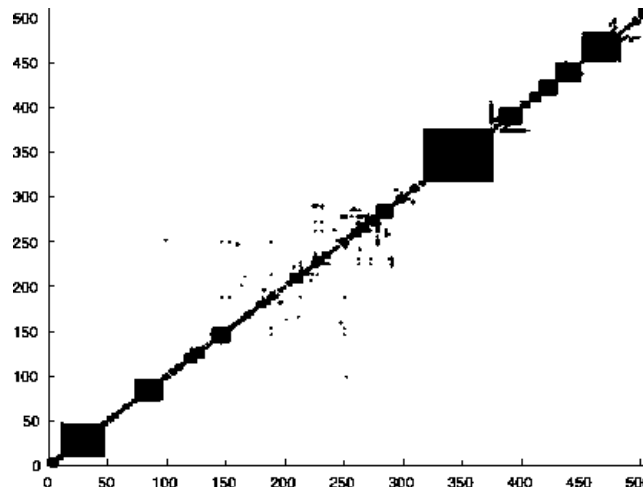
- Leverage the **high locality**
- **Minimize memory usage** to support larger graphs
- **Align with GPU architecture** to maximize performance

Observation 1: The Reordered Graph Has a Community on a Diagonal

- Most vertices are densely connected within certain groups
- With reordering algorithms, real-world graphs can be reordered in diagonal lines



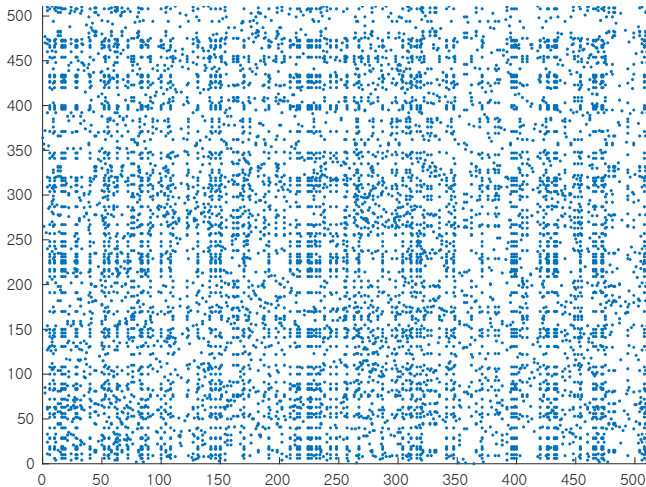
Random Order



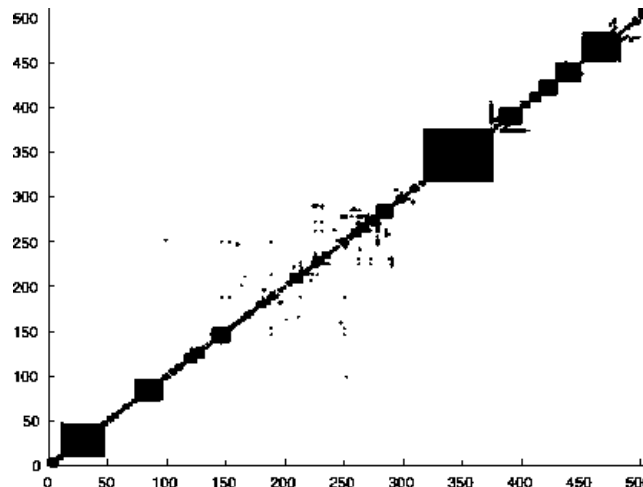
Community-based

Solution 1: Community-aware Clustered Graph Representation

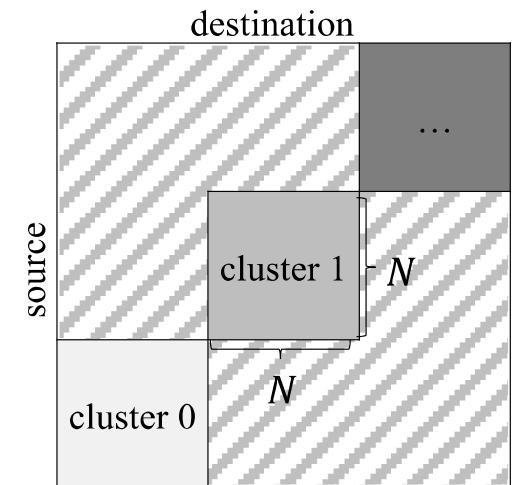
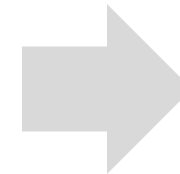
- Extract communities of real-world graph as clusters
 - Execute edges in the same cluster on the same GPU core
- ➔ CR² can **exploit the vertex locality**



Random Order



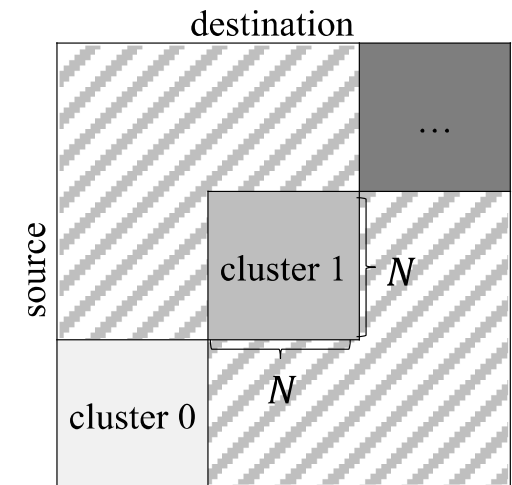
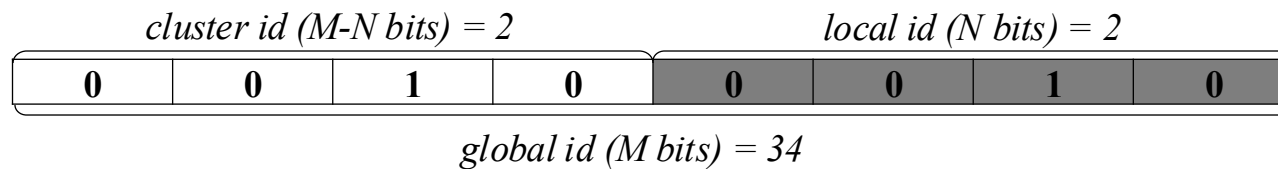
Community-based



Clustering in CR²

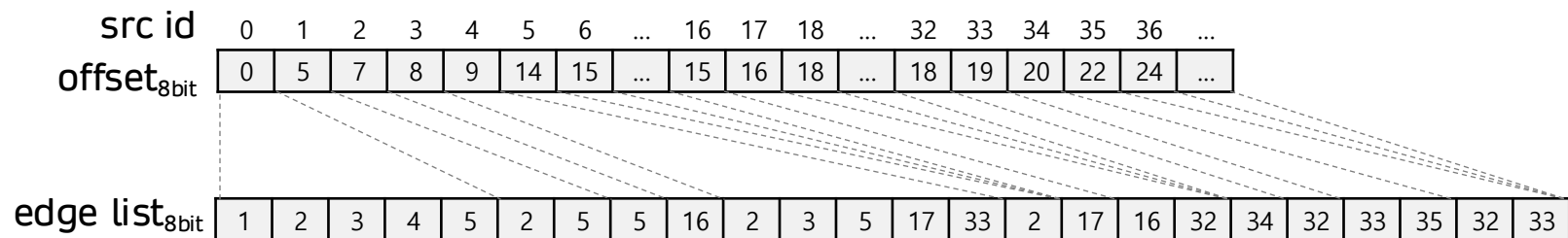
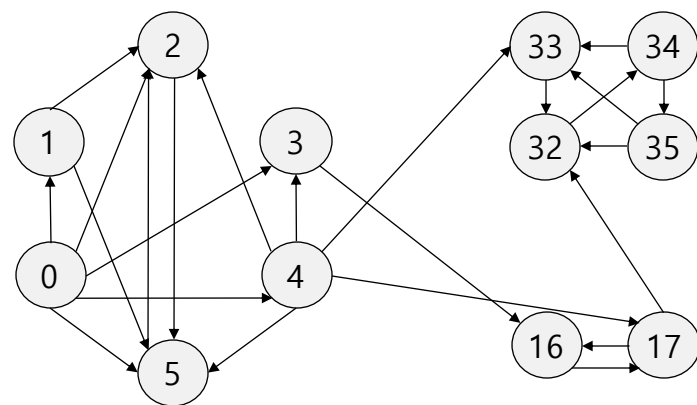
Solution 1: Community-aware Clustered Graph Representation

- Extract communities of real-world graph as clusters
- Execute edges in the same cluster on the same GPU core
- ➔ CR² can **exploit the vertex locality**
- ➔ **Reduce memory usage** by representing vertex id with local id



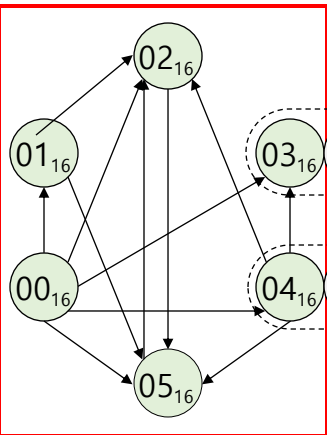
Clustering in CR²

Solution 1: Community-aware Clustered Graph Representation

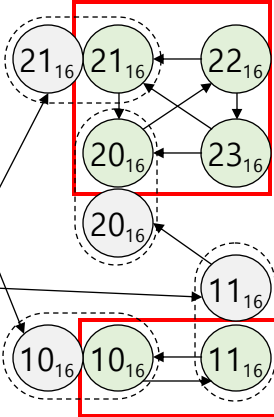


Make clusters to improve locality
Reduce memory usage by using local id

Cluster-0



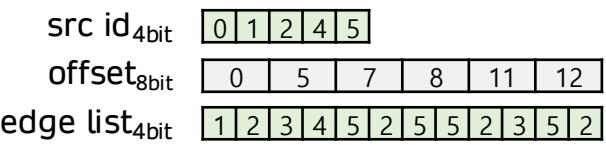
Cluster-2



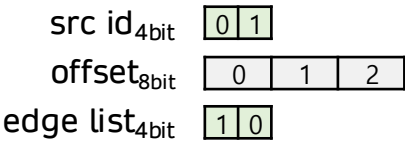
Cluster-1



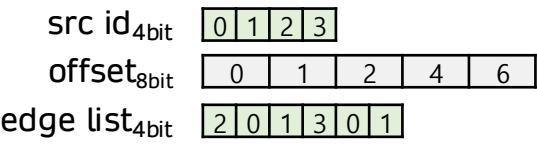
Cluster-0 graph



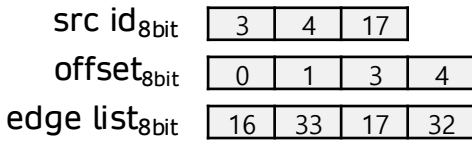
Cluster-1 graph



Cluster-2 graph



Inter-cluster graph



Seperated subgraph for inter cluster

Solution 1: Community-aware Clustered Graph Representation

Cluster-0 graph

src id_{4bit}

0	1	2	4	5
---	---	---	---	---

offset_{8bit}

0	5	7	8	11	12
---	---	---	---	----	----

edge list_{4bit}

1	2	3	4	5	2	5	5	2	3	5	2
---	---	---	---	---	---	---	---	---	---	---	---

Cluster-1 graph

src id_{4bit}

0	1
---	---

offset_{8bit}

0	1	2
---	---	---

edge list_{4bit}

1	0
---	---

Cluster-2 graph

src id_{4bit}

0	1	2	3
---	---	---	---

offset_{8bit}

0	1	2	4	6
---	---	---	---	---

edge list_{4bit}

2	0	1	3	0	1
---	---	---	---	---	---

Inter-cluster graph

src id_{8bit}

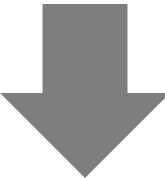
3	4	17
---	---	----

offset_{8bit}

0	1	3	4
---	---	---	---

edge list_{8bit}

16	33	17	32
----	----	----	----



Gather Processing
Execute edges in the same cluster on the same GPU core

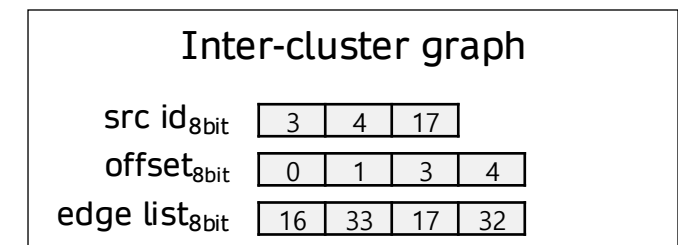
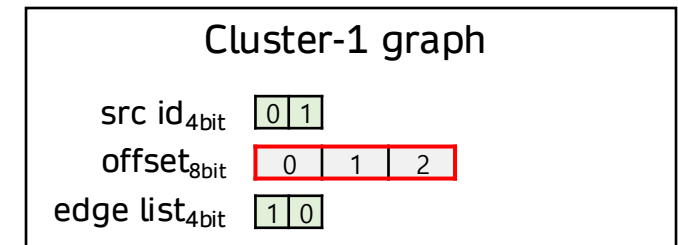
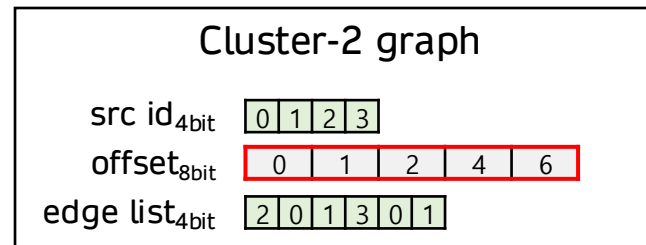
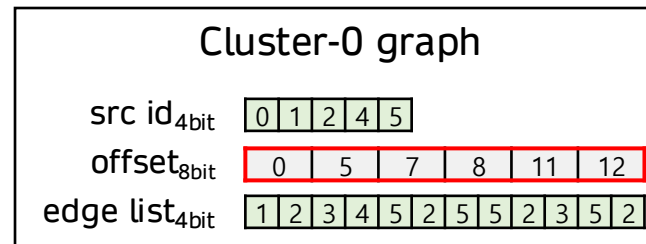
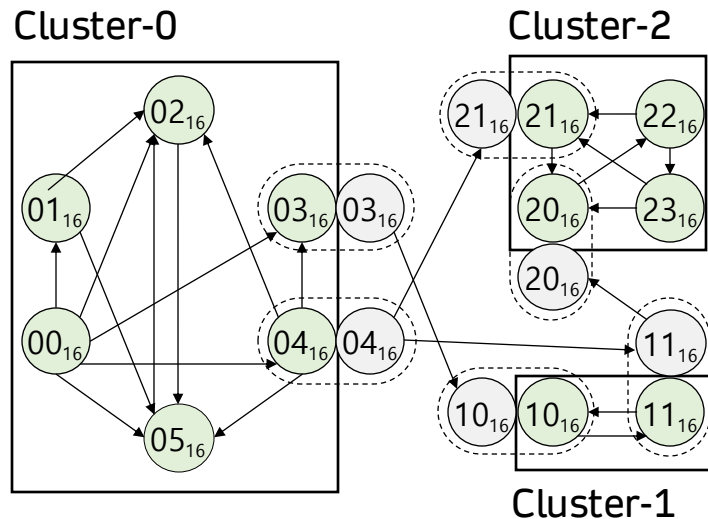
Cluster-0								Cluster-1								Cluster-2								Inter-cluster							
T0	T1	T2	T3	T4	T5	T6	T7	T0	T1	T2	T3	T4	T5	T6	T7	T0	T1	T2	T3	T4	T5	T6	T7	T0	T1	T2	T3	T4	T5	T6	T7
0	1	2	4	5				0	1							0	1	2	3					3	4	17					
0	5	7	8	11				0	1							0	1	2	4					0	1	3					
5	7	8	11	12				1	2							1	2	4	6					1	3	4					
1	2	5	2	2				1	0							2	0	1	0					16	33	32					
2	5		3															3	1						17						
3			5																												
4																															
5																															

- Improve locality
- Reduce memory usage with local id

- Improve locality
- Reduce memory usage with local id
- But, **still suffered from workload imbalance**

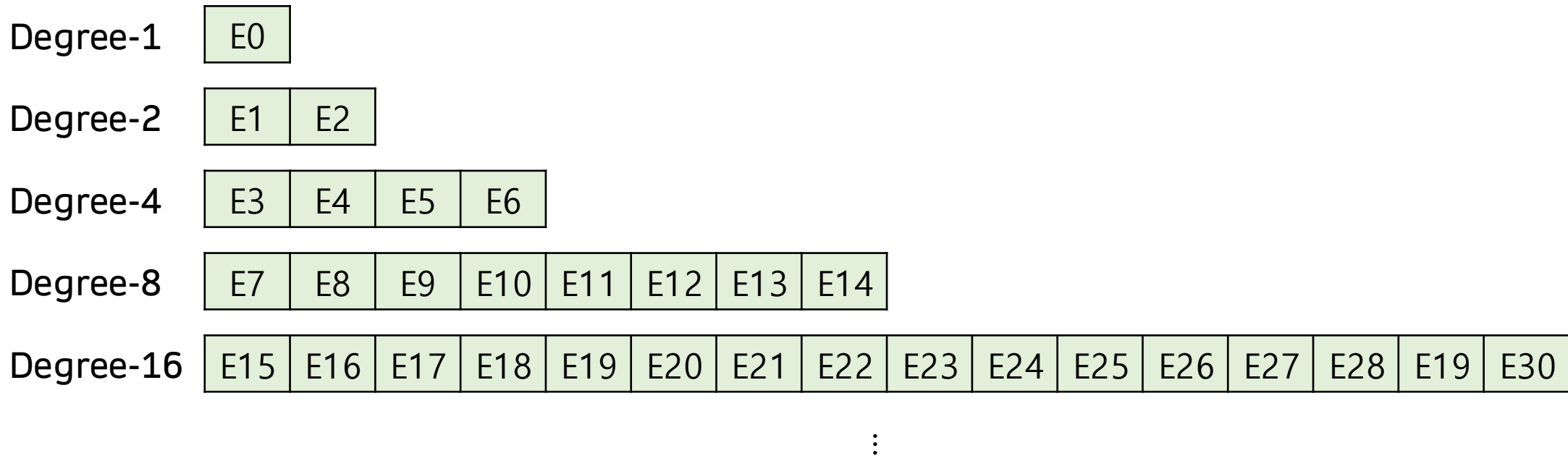
Observation 2: Degree Regulation can Remove Offset Array

- Offset array indicates the number of neighboring edges
- If graph representation has a fixed number of edges per vertex,
 - Can remove offset array
 - Can balance workload

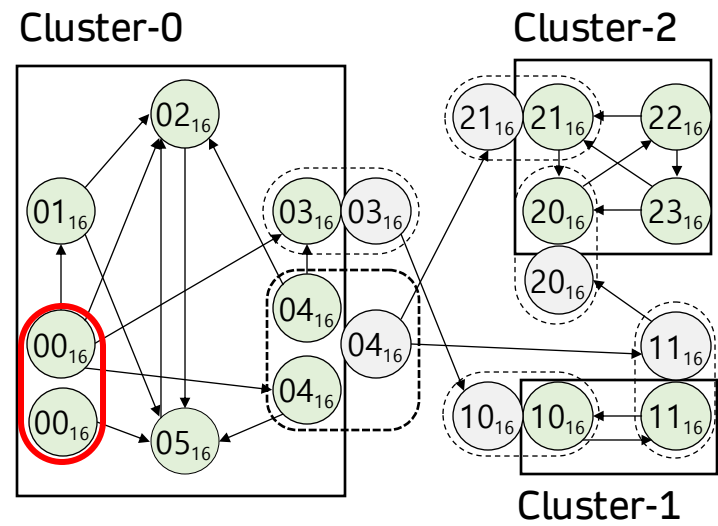


Solution 2: Degree-ordered Subgraph

- Partition the graph into multiple degree-ordered subgraphs in which all the vertices have the **Degree-n regularized number of edges**
- Enable **fine-grained workload balancing** across GPU warps with less memory usage



Solution 2: Degree-ordered Subgraph



Cluster-0 graph

src id_{4bit}

0012445

offset_{8bit}

04578101112

edge list_{4bit}

123452552352

Cluster-1 graph

src id_{4bit}

0

1

offset_{8bit}

0

1

2

edge list_{4bit}

1

0

Cluster-2 graph

src id_{4bit}

0

1

2

3

offset_{8bit}

0

1

2

4

6

edge list_{4bit}

2

0

1

3

0

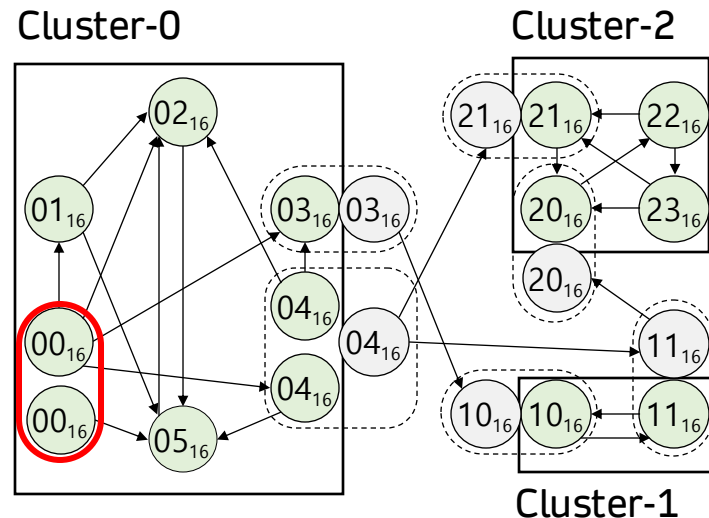
1

Inter-cluster graph

src id _{8bit}	3	4	17	
offset _{8bit}	0	1	3	4
edge list _{8bit}	16	33	17	32

Vertex-split graph & representation
 Split vertex to make it a power-of-two

Solution: CR² Representation



Cluster-0 graph

src id_{4bit}

0012445

offset_{8bit}

04578101112

edge list_{4bit}

123452552352

Cluster-1 graph

src id_{4bit}

0

1

offset_{8bit}

0

1

2

edge list_{4bit}

1

0

Cluster-2 graph

src id_{4bit}

0

1

2

3

offset_{8bit}

0

1

2

4

6

edge list_{4bit}

2

0

1

3

0

1

Inter-cluster graph

src id _{8bit}	3	4	17	
offset _{8bit}	0	1	3	4
edge list _{8bit}	16	33	17	32



Remove offset array & Generate Degree-n Subgraph

Cluster-0 degree-4

src id_{4bit}

0

edge list_{4bit}

1234

Cluster-0 degree-2

src id_{4bit}

1

4

edge list_{4bit}

2

5

2

3

Cluster-0 degree-1

src id _{4bit}	0	2	4	5
edge list _{4bit}	5	5	5	2

Cluster-1 degree-1

src id _{4bit}	0	1
edge list _{4bit}	1	0

Cluster-0 degree-2

src id_{4bit}

23

edge list_{4bit}

1301

Cluster-2 degree-1

src id _{4bit}	0	1
edge list _{4bit}	2	0

Inter-Cluster degree-1

src id_{8bit}

4

edge list_{8bit}

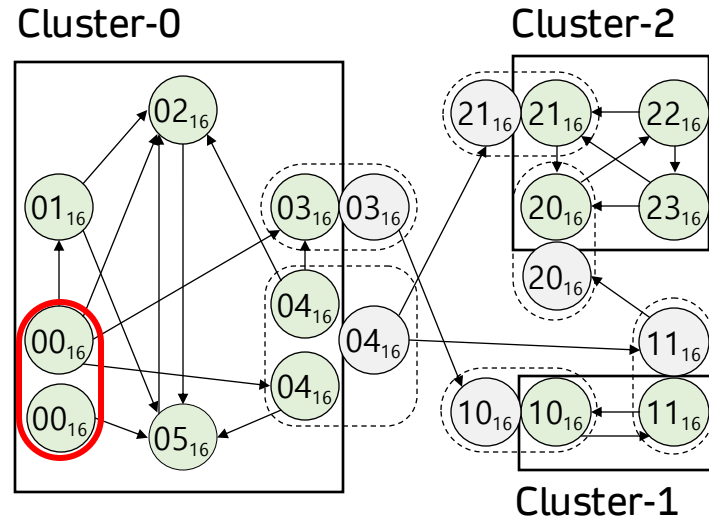
33

17

Inter-Cluster degree-1

src id _{8bit}	3	17
edge list _{8bit}	16	32

Solution: CR² Representation



Cluster-0 graph

src id_{4bit}

0012445

offset_{8bit}

04578101112

edge list_{4bit}

123452552352

Cluster-1 graph				
src id _{4bit}	0	1		
offset _{8bit}	0		1	2
edge list _{4bit}	1	0		

Cluster-2 graph

src id_{4bit}

0	1	2	3
---	---	---	---

offset_{8bit}

0	1	2	4	6
---	---	---	---	---

edge list_{4bit}

2	0	1	3	0	1
---	---	---	---	---	---

Inter-cluster graph				
src id _{8bit}	3	4	17	
offset _{8bit}	0	1	3	4
edge list _{8bit}	16	33	17	32



Remove offset array & Generate Degree-n Subgraph

Cluster-0 degree-4	
src id _{4bit}	0
edge list _{4bit}	1 2 3 4

Cluster-0 degree-2	
src id _{4bit}	1 4
edge list _{4bit}	2 5 2 3

Cluster-0 degree-1	
src id _{4bit}	0 2 4 5
edge list _{4bit}	5 5 5 2

Cluster-1 degree-1	
src id _{4bit}	0 1
edge list _{4bit}	1 0

Cluster-0 degree-2	
src id _{4bit}	2 3
edge list _{4bit}	1 3 0 1

Cluster-2 degree-1	
src id _{4bit}	0 1
edge list _{4bit}	2 0

Inter-Cluster degree-1	
src id _{8bit}	4
edge list _{8bit}	33 17

Inter-Cluster degree-1	
src id _{8bit}	3 17
edge list _{8bit}	16 32

- High locality
- Minimize memory usage
- Align with GPU architecture

Solution: CR² Representation

Cluster-0 degree-4 src id _{4bit} 0 edge list _{4bit} 1 2 3 4	Cluster-0 degree-2 src id _{4bit} 1 4 edge list _{4bit} 2 5 2 3	Cluster-0 degree-1 src id _{4bit} 0 2 4 5 edge list _{4bit} 5 5 5 2	Cluster-1 degree-1 src id _{4bit} 0 1 edge list _{4bit} 1 0
Cluster-0 degree-2 src id _{4bit} 2 3 edge list _{4bit} 1 3 0 1	Cluster-2 degree-1 src id _{4bit} 0 1 edge list _{4bit} 2 0	Inter-Cluster degree-1 src id _{8bit} 4 edge list _{8bit} 33 17	Inter-Cluster degree-1 src id _{8bit} 3 17 edge list _{8bit} 16 32

↓
Gather Processing

Intra-Cluster Execution

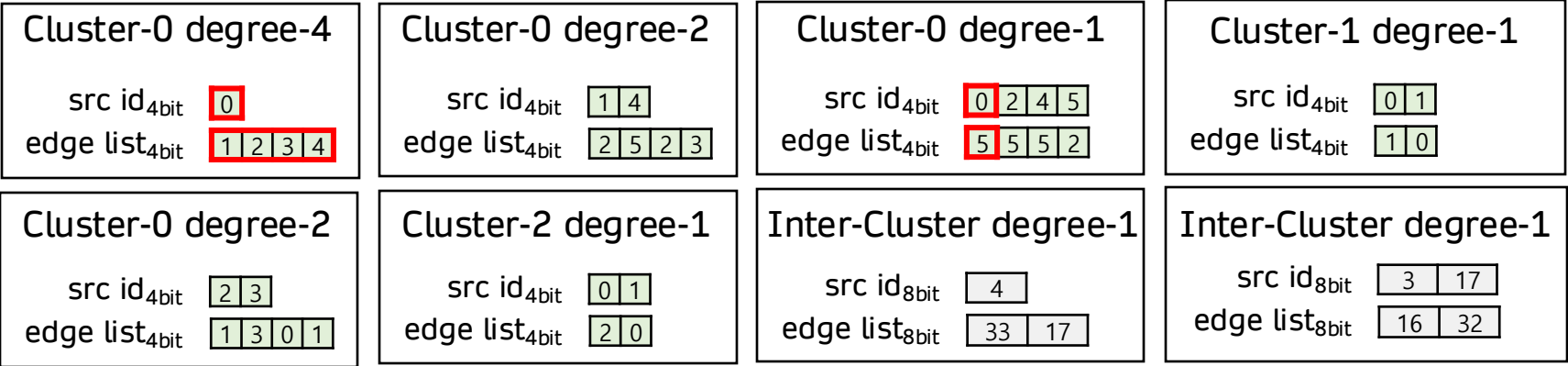
Degree-4 Execution								Degree-2 Execution								Degree-1 Execution							
T0	T1	T2	T3	T4	T5	T6	T7	T0	T1	T2	T3	T4	T5	T6	T7	T0	T1	T2	T3	T4	T5	T6	T7
0	0	0	0					1	1	4	4	2	2	3	3	0	2	4	5	0	1	0	1
1	2	3	4					2	5	2	3	1	3	0	1	5	5	5	2	1	0	2	0

Execute the degree-n subgraphs of each cluster simultaneously

Inter-Cluster Execution

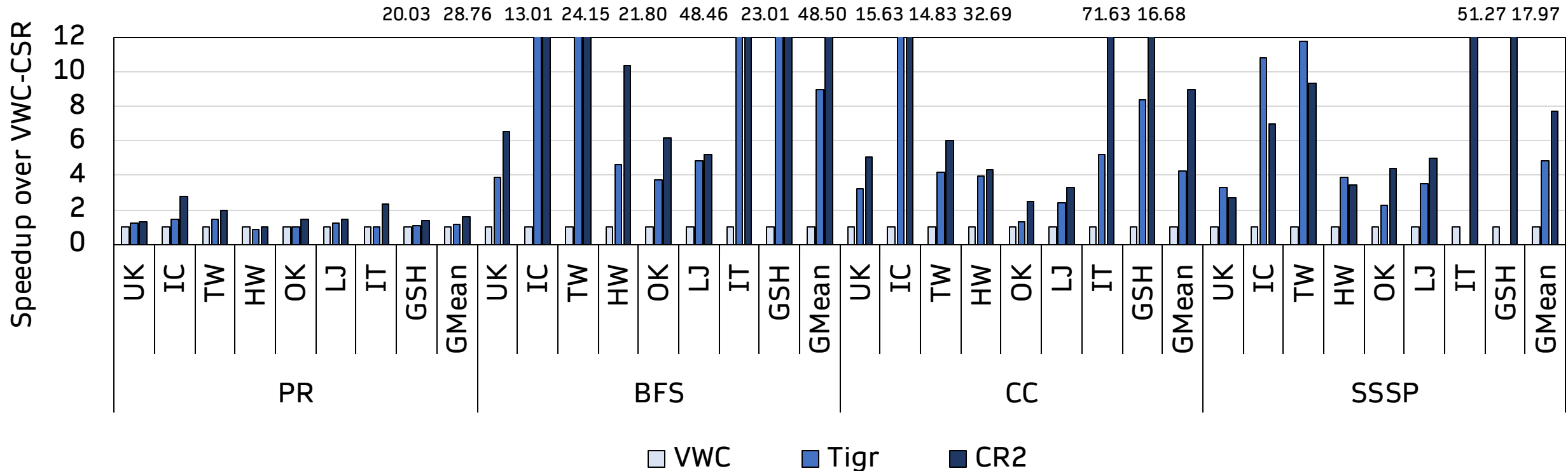
Degree-4 Execution								Degree-2 Execution								Degree-1 Execution							
T0	T1	T2	T3	T4	T5	T6	T7	T0	T1	T2	T3	T4	T5	T6	T7	T0	T1	T2	T3	T4	T5	T6	T7
								4	4							3	17						
								33	17							16	32						

Solution: CR² Representation



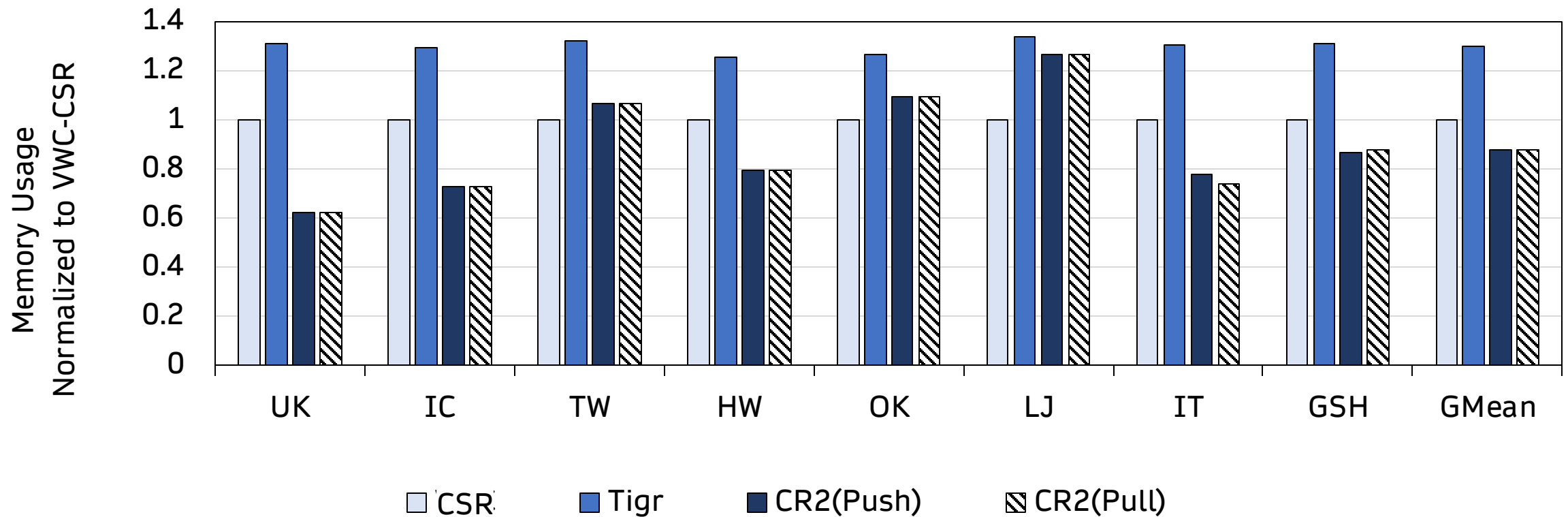
Overall Performance

- Achieve **6.47x** and **1.55x** performance geomean speedup compared to VWC-CSR and Tigr
- Use NVIDIA Geforce RTX 3090
- Use Four graph algorithms: PageRank(PR), Breadth-First Search(BFS), Connected Components(CC), and Single Source Shortest Path(SSSP)



Memory Usage

- Achieve **12.3%** and **32.2%** less memory on geomean average compared with CSR and Tigr



Conclusion

- This work proposes a new graph representation called **CR²**
- Based on two key solution,
 - Community-aware clustered subgraph
 - Degree-ordered subgraphs with vertex degree regulation
- CR² enhance graph processing on GPU by,
 - Improving locality
 - Minimizing memory usage
 - Aligning with GPU architecture
- Achieve 1.53 times performance speedup while using 32.1% less memory on the geomean average compared to the state-of-the-art techniques

Thank you 😊

shin0403@yonsei.ac.kr