



**MADSYS**  
THU • HPC

# Squeezing out

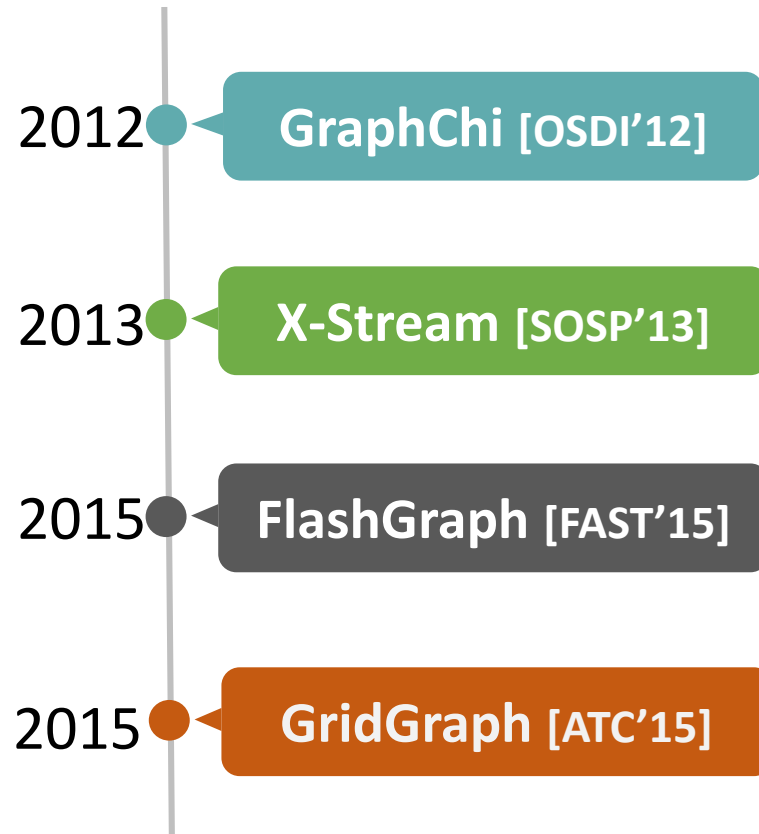
**All the Value of Loaded Data:**

An Out-Of-Core Graph Processing System with Reduced Disk I/O

**Zhiyuan Ai<sup>1</sup>**, Mingxing Zhang<sup>1</sup>, Yongwei Wu<sup>1</sup>,  
Xuehai Qian<sup>2</sup>, Kang Chen<sup>1</sup>, Weimin Zheng<sup>1</sup>

Tsinghua University<sup>1</sup>, University of Southern California<sup>2</sup>

## Single machine graph processing systems (Disk I/O is bottleneck)



**Wrong  
trade-off !!**



# Comparison

**Trade-off**

**Bottleneck**

**Wrong**

	I/O locality	The amount of Disk I/O
Prior systems	Better	Large
CLIP	Not better	<b>Little</b>

**CLIP**



Squeezing out All the Value of Loaded Data  
**Up to tens or even thousands of times speedup**

Prior systems

**CLIP**



**VS**

**Loaded Data  
(In-Memory)**

Squeezing gently  
Many iterations

Squeezing **hard**  
**Few** iterations

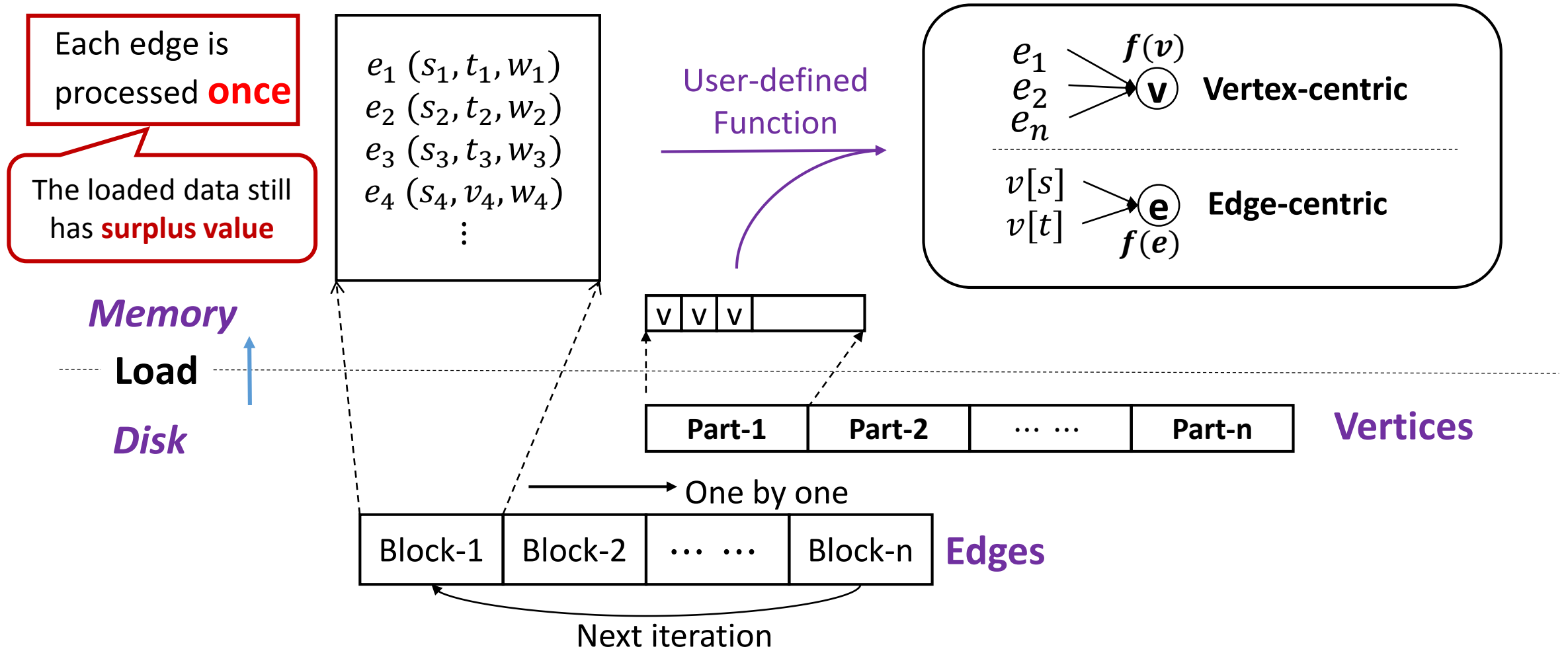
# Outline

1. Background
- 2. CLIP**
3. Evaluation
4. Conclusion



# 2.1 Shortcomings of prior works (1)

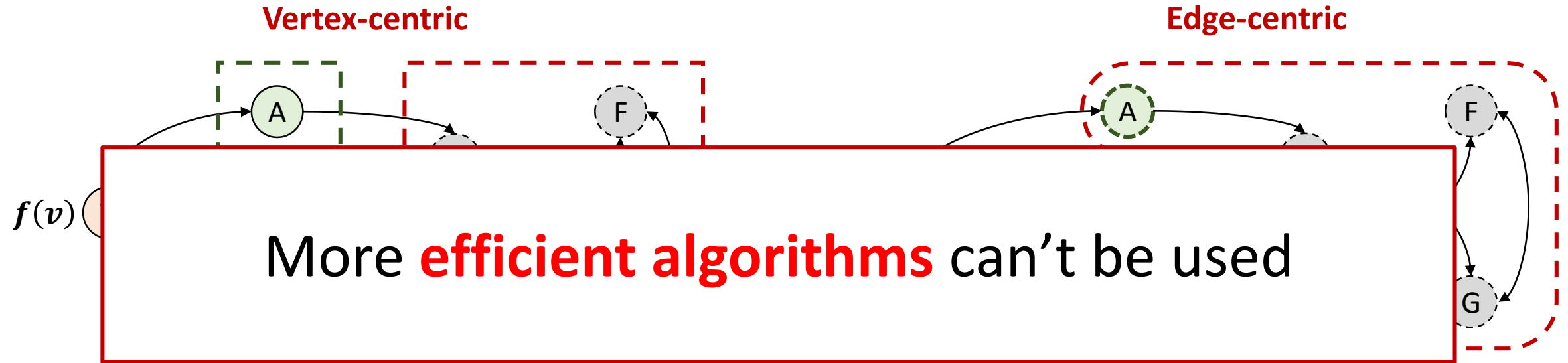
## ① Processing order is **not flexible!!**



# 2.1 Shortcomings of prior works (2)



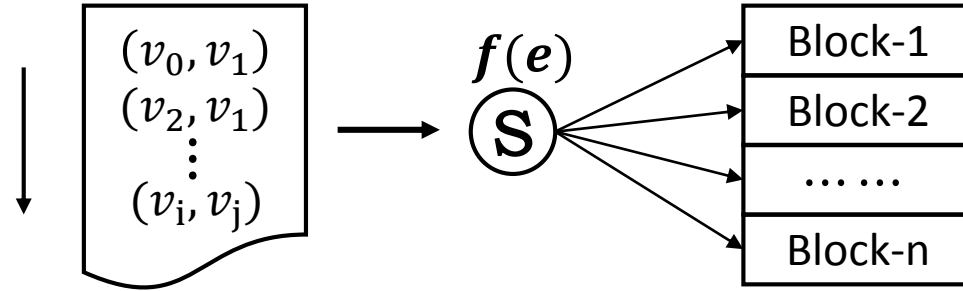
## ② Constrained programming model



```
Update(vertex) begin
  x[] ← read values of in- and out-edges of vertex ;
  vertex.value ← f(x[]);
  foreach edge of vertex do
    edge.value ← g(vertex.value, edge.value);
  end
  Graphchi [OSDI'12]
```

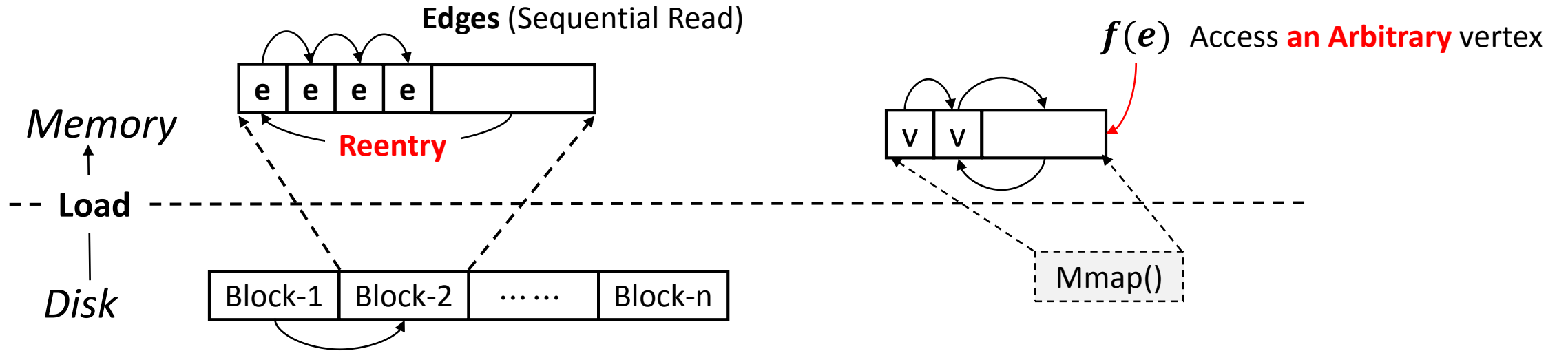
```
function STREAMEDGES( $F_e, F$ )
  Sum = 0
  for each active block do ▷ block with active edges
    for each edge  $\in$  block do
      if  $F(\text{edge.source})$  then
        Sum +=  $F_e(\text{edge})$ 
      end if
  end if
  GridGraph [ATC'15]
```

# 2.2 Our solutions



↑ *Sorting*

↓ *Execution*





**CLIP:** Squeezing out all the value of loaded data

## 1. Reentry of Loaded Data

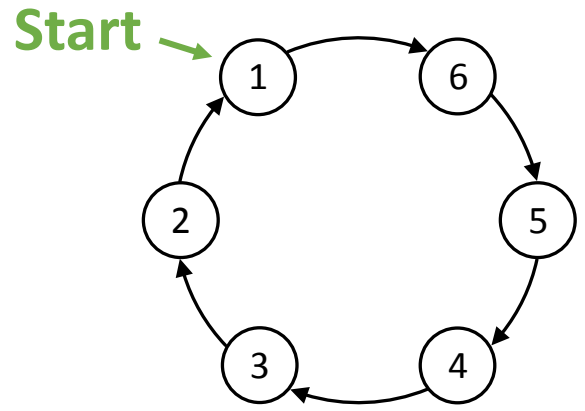
- **Principle:** Processing the loaded data **multiple times** rather than once.
- **Example:** SSSP (Single Source Shortest Path).

## 2. Beyond the Neighborhood

- **Principle:** The user-defined function is allowed to update **an arbitrary** vertex's property.
- **Example:** WCC (Weakly Connected Component).

# 2.2.1 Reentry of Loaded Data

**Example:** Calculating single source shortest path (SSSP)



1 → 6

*Init dist*

0	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$
---	----------	----------	----------	----------	----------

	<i>source</i>	<i>target</i>	<i>weight</i>
e1	1	6	1
e2	2	1	1
e3	3	2	1
e4	4	3	1
e5	5	4	1
e6	6	5	1

Prior system: process once

1 → 6
2 → 1
3 → 2
4 → 3
5 → 4
6 → 5

Disk Memory

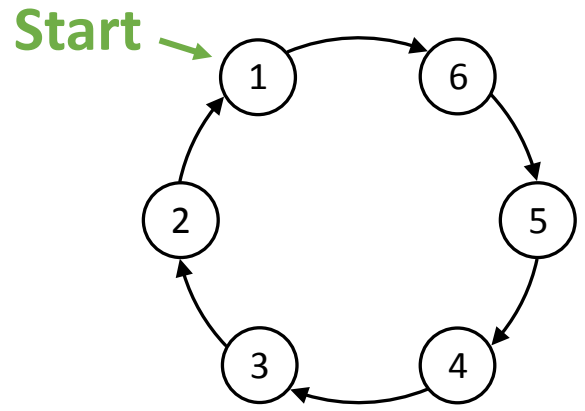
↓ once

Iteration 1

$v_1$	$v_2$	$v_3$	$v_4$	$v_5$	$v_6$
0	$\infty$	$\infty$	$\infty$	2	1

# 2.2.1 Reentry of Loaded Data

**Example:** Calculating single source shortest path (SSSP)



1 → 6

*Init dist*

0	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$
---	----------	----------	----------	----------	----------

	<i>source</i>	<i>target</i>	<i>weight</i>
e1	1	6	1
e2	2	1	1
e3	3	2	1
e4	4	3	1
e5	5	4	1
e6	6	5	1

Prior system: process once

1 → 6
2 → 1
3 → 2
4 → 3
5 → 4
6 → 5

Disk Memory

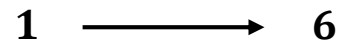
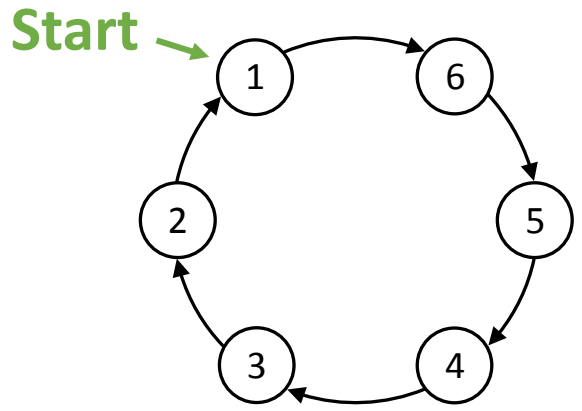
	$v_1$	$v_2$	$v_3$	$v_4$	$v_5$	$v_6$
Iteration 1	0	$\infty$	$\infty$	$\infty$	2	1
Iteration 2	0	$\infty$	$\infty$	3	2	1

once

# 2.2.1 Reentry of Loaded Data



**Example:** Calculating single source shortest path (SSSP)



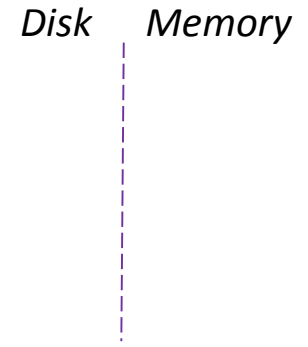
Prior system: process once

*Init dist*

0	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$
---	----------	----------	----------	----------	----------

	<i>source</i>	<i>target</i>	<i>weight</i>
e1	1	6	1
e2	2	1	1
e3	3	2	1
e4	4	3	1
e5	5	4	1
e6	6	5	1

1 → 6
2 → 1
3 → 2
4 → 3
5 → 4
6 → 5

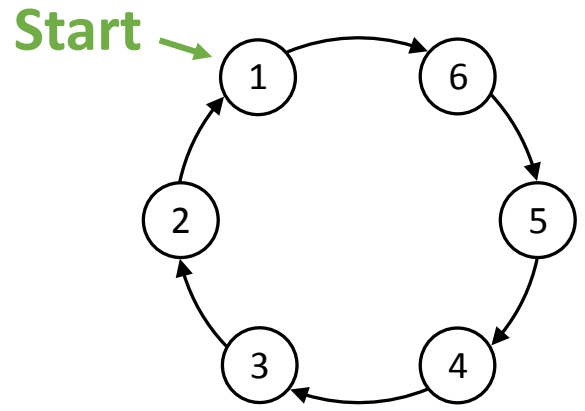


	$v_1$	$v_2$	$v_3$	$v_4$	$v_5$	$v_6$
Iteration 1	0	$\infty$	$\infty$	$\infty$	2	1
Iteration 2	0	$\infty$	$\infty$	3	2	1
Iteration 3	0	$\infty$	4	3	2	1
Iteration 4	0	5	4	3	2	1

Total amount of disk I/O:  $4 * \text{sizeof}(E)$

# 2.2.1 Reentry of Loaded Data—CLIP

**Example:** Calculating single source shortest path (SSSP)



Init dist

0	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$
---	----------	----------	----------	----------	----------

	source	target	weight
e1	1	6	1
e2	2	1	1
e3	3	2	1
e4	4	3	1
e5	5	4	1
e6	6	5	1

CLIP: Process multiple times

1 → 6
2 → 1
3 → 2
4 → 3
5 → 4
6 → 5

Disk Memory

Iteration 1

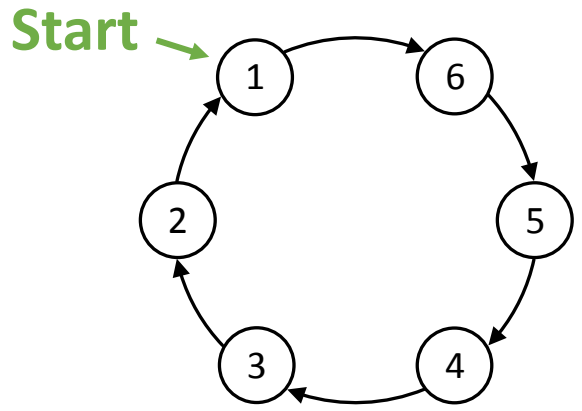
$v_1$	$v_2$	$v_3$	$v_4$	$v_5$	$v_6$
0	$\infty$	$\infty$	3	2	1

Pass two

# 2.2.1 Reentry of Loaded Data—CLIP



**Example:** Calculating single source shortest path (SSSP)



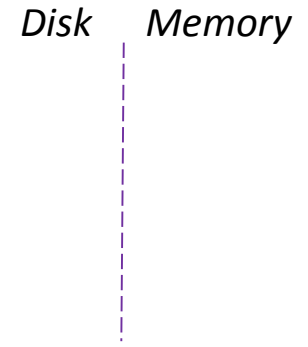
**CLIP: Process multiple times**

*Init dist*

0	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$
---	----------	----------	----------	----------	----------

	<i>source</i>	<i>target</i>	<i>weight</i>
e1	1	6	1
e2	2	1	1
e3	3	2	1
e4	4	3	1
e5	5	4	1
e6	6	5	1

1 → 6
2 → 1
3 → 2
4 → 3
5 → 4
6 → 5



	$v_1$	$v_2$	$v_3$	$v_4$	$v_5$	$v_6$
Iteration 1	0	$\infty$	$\infty$	3	2	1
Iteration 2	0	5	4	3	2	1

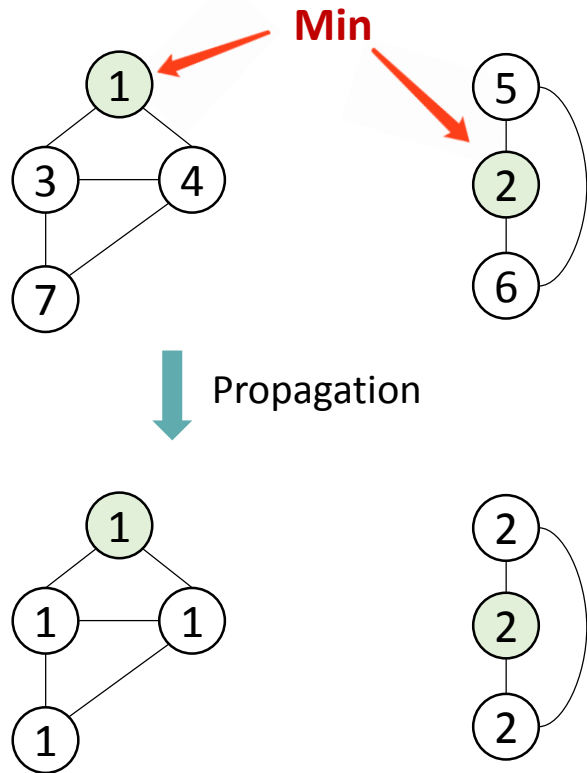
**Total amount of disk I/O:  $2 * \text{sizeof}(E)$**

# 2.2.2 Beyond the Neighborhood



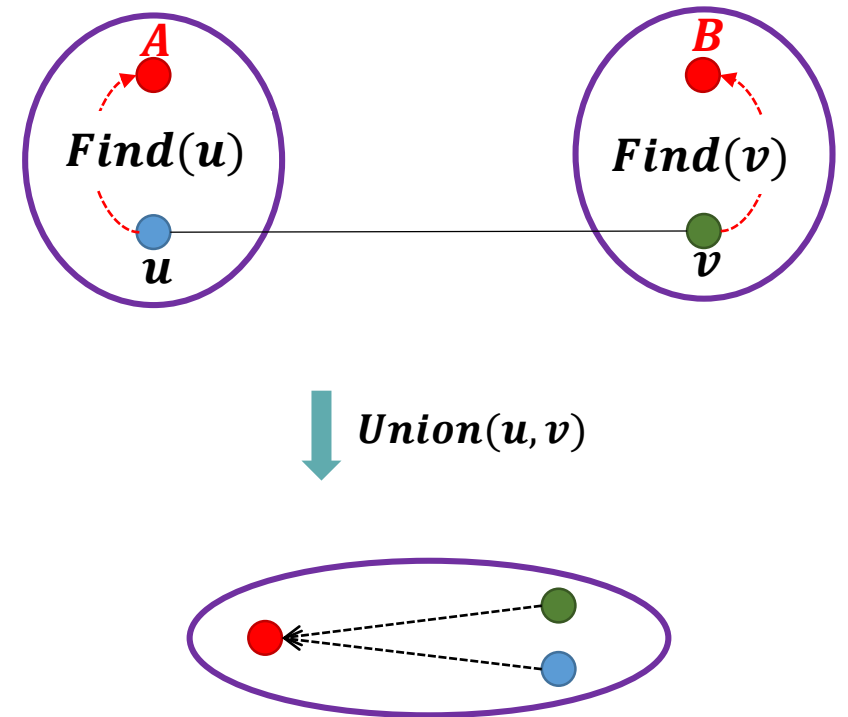
**Example:** Weakly Connected Component (WCC)

1. Label propagation based algorithm (Prior systems)



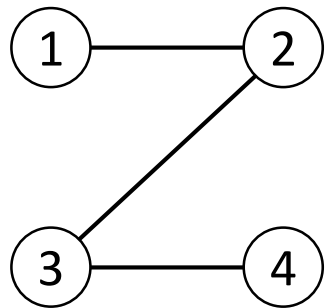
VS

2. Disjoint set algorithm (CLIP)



# 2.2.2 Beyond the Neighborhood

**Example:** Calculating weakly connected component (WCC)



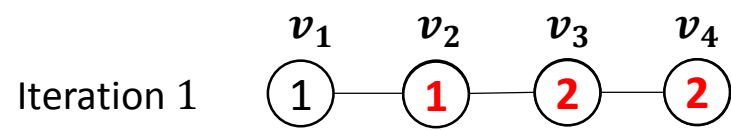
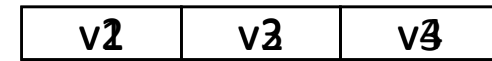
Undirected graph

Prior system: label propagation based

	<i>source</i>	<i>target</i>
e1	2	3
e2	2	1
e3	1	2
e4	4	3
e5	3	4
e6	3	2

2 → 1
2 → 3
1 → 2
4 → 3
3 → 4
3 → 2

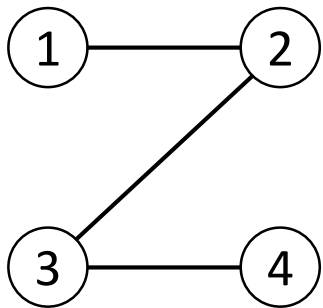
Disk    Memory





# 2.2.2 Beyond the Neighborhood

**Example:** Calculating weakly connected component (WCC)



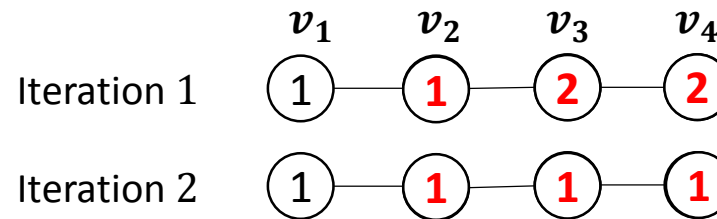
Undirected graph

	<i>source</i>	<i>target</i>
e1	2	3
e2	2	1
e3	1	2
e4	4	3
e5	3	4
e6	3	2

Prior system: label propagation based

2 → 1
2 → 3
1 → 2
4 → 3
3 → 4
3 → 2

Disk    Memory

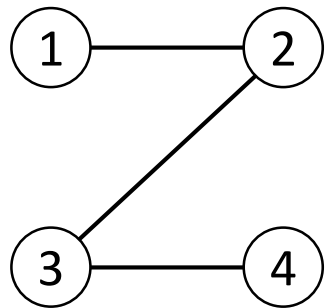


Total amount of disk I/O:  $2 * \text{sizeof}(E)$

# 2.2.2 Beyond the Neighborhood—CLIP



**Example:** Calculating weakly connected component (WCC)



Undirected graph

	<i>source</i>	<i>target</i>
e1	2	3
e2	2	1
e3	1	2
e4	4	3
e5	3	4
e6	3	2

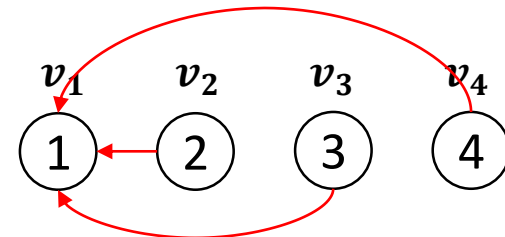
Clip: disjoint set (Can access an arbitrary vertex)

<b>2 → 1</b>
<b>2 → 3</b>
<b>1 → 2</b>
<b>4 → 3</b>
<b>3 → 4</b>
<b>3 → 2</b>

Disk    Memory

v1	v2	v3	v4
----	----	----	----

Iteration 1



Total amount of disk I/O: **1** \* sizeof(E)

# 2.2.3 APIs

## Programming model of C<sub>LIP</sub>

---

<i>Sort</i> ( $\mathcal{F}_s$ )	—	$\mathcal{F}_s :=$ <b>double</b> function( <i>Edge</i> & <i>e</i> )
<i>Exec</i> ( $\mathcal{F}_e$ )	—	$\mathcal{F}_e :=$ <b>void</b> function( <i>Vertexes</i> & <i>v_list</i> , <i>Edge</i> & <i>e</i> )
<i>VMap</i> ( $\mathcal{F}_v$ )	—	$\mathcal{F}_v :=$ <b>void</b> function( <i>Vertexes</i> & <i>v_list</i> , <i>VertexID</i> & <i>vid</i> )

---

All vertices

---

### Algorithm 2 SSSP Algorithm in CLIP.

---

**Functions:**

```
 $\mathcal{F}_v(v\_list, vid) := \{$   
  if vid == start do  
    v_list[vid].dist  $\leftarrow$  0; v_list.setActive(vid, true);  
  else v_list[vid].dist  $\leftarrow$  INF; v_list.setActive(vid, false);  
 $\}$ 
```

```
 $\mathcal{F}_e(v\_list, e) := \{$   
  if v_list[e.dst].dist > v_list[e.src].dist + e.weight do  
    v_list[e.dst].dist  $\leftarrow$  v_list[e.src].dist + e.weight;  
    v_list.setActive(e.dst, true);  
  else v_list.setActive(e.dst, false);  
 $\}$ 
```

**Computation:**

*VMap*( $\mathcal{F}_v$ );  
**Until convergence:**  
*Exec*( $\mathcal{F}_e$ );

Simple

---

### Algorithm 3 WCC Algorithm in CLIP.

---

**Functions:**

```
 $\mathcal{F}_{find}(v\_list, vid) := \{$   
  if v_list[vid].pa == vid do return vid;  
  else return v_list[vid].pa =  $\mathcal{F}_{find}(v\_list, v\_list[vid].pa)$ ;  
 $\}$ 
```

```
 $\mathcal{F}_{union}(v\_list, src, dst) := \{$   
  s  $\leftarrow$   $\mathcal{F}_{find}(v\_list, src)$ ;  
  d  $\leftarrow$   $\mathcal{F}_{find}(v\_list, dst)$ ;  
  if s < d do v_list[d].pa  $\leftarrow$  v_list[s].pa;  
  else if s > d do v_list[s].pa  $\leftarrow$  v_list[d].pa;  
 $\}$ 
```

```
 $\mathcal{F}_e(v\_list, e) := \{ \mathcal{F}_{union}(v\_list, e.src, e.dst); \}$ 
```

```
 $\mathcal{F}_v(v\_list, vid) := \{$   
  v_list[vid].pa  $\leftarrow$  vid; v_list.setActive(vid, true);  
 $\}$ 
```

**Computation:**

*VMap*( $\mathcal{F}_v$ );  
*Exec*( $\mathcal{F}_e$ );

# 2.2.4 Applications

①

- **Reentry of loaded data (Asynchronous applications)**

- **SSSP (Single Source Shortest Path)**

- **BFS (Breadth-first Search)**

Delta-based PageRank, Diameter Approximation, Transitive Closures, Betweenness Centrality .....

②

- **Beyond-neighborhood applications**

- **WCC (Weakly connected component)**

- **MIS (Maximal Independent Set )**

Graph Stream Algorithms, SCC, Coloring, MCST, Triangle counting .....

# Outline

1. Background
2. CLIP
- 3. Evaluation**
4. Conclusion



# 3.1 Setup

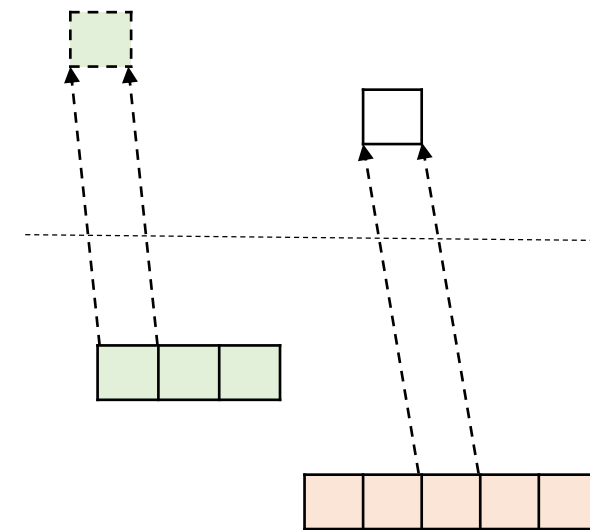
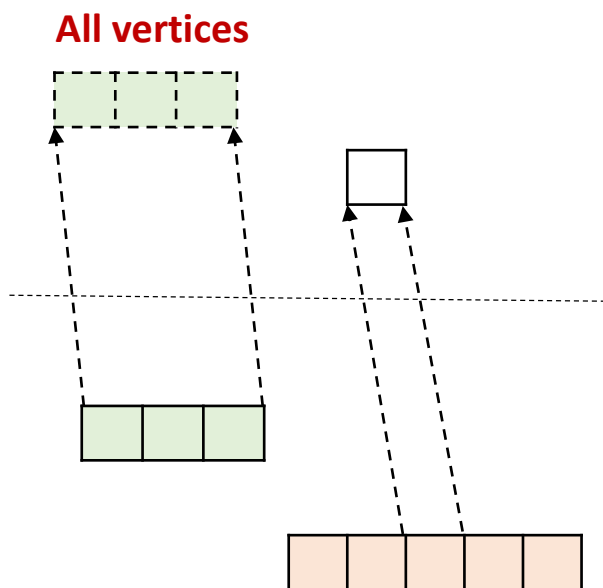
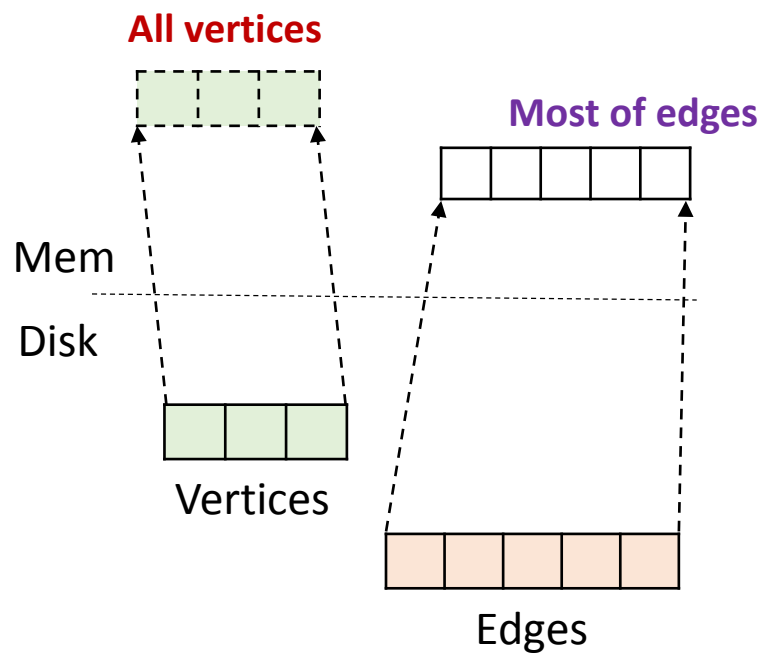


Testing three different scenarios:

All-In-Memory

Semi-External

External



(Limit is 32GB)

## The real-world graph datasets

	Vertexes	Edges	Threshold	
			External	Semi
LiveJournal	4.85M	69.0M	16MB	256MB
Dimacs	23.9M	58.3M	64MB	256MB
Twitter	41.7M	1.47B	128MB	4GB
Friendster	65.6M	1.8B	128MB	4GB
Yahoo	1.4B	6.64B	4GB	8GB

Tips: Dimacs and Yahoo have a large diameter

## Test environment

- **CPU:** Two Intel(R) Xeon(R) CPU E5-2640 v2 @ 2.00GHz (each has 8-cores)
- **DRAM:** 32GB, L3 cache 20MB
- **Disk:** Standard 1TB SSD drive. The average throughput is about 450MB/s for sequential read.

# 3.2 Evaluation—Loaded data reentry



Execution time (in seconds, external / semi-external / all-in-memory)

		LiveJournal	Dimacs	Friendster	Twitter	Yahoo
SSSP	X-Stream	357.9 / 118.4 / 8.45	77212/ 22647/ 853.2	6352 / 3346 / -	4065 / 2255 / -	$\infty$ / $\infty$ / -
	GridGraph	66.42 / 48.1 / 6.97	14618/ 13480/ 889.9	1086 / 784.6 / 85.31	1639 / 1083 / 83.51	77298/ 17432/ -
	CLIP	30.14 / 11.23 / 5.09	3202 / 1981 / 316.1	176.2 / 55.79 / 55.85	1353 / 600.6 / 91.82	18160/ 6932 / -
BFS	X-Stream	91.50 / 22.94 / 4.06	8934 / 6538 / 114.9	2526 / 1084 / -	1421 / 627.4 / -	$\infty$ / $\infty$ / -
	GridGraph	13.20 / 15.4 / 2.49	5199 / 5239 / 406.2	499.6 / 493.7 / 61.54	220.5 / 209.6 / 32.16	35572/ 7403 / -
	CLIP	10.01 / 5.46 / 2.53	1768 / 1059 / 96.12	98.87 / 38.55 / 38.72	141.2 / 110.4 / 44.7	10533/ 3297 / -

The number of iterations

		LiveJornal	Dimacs	Friendster	Twitter	Yahoo
SSSP	X-Stream	45	10790	50	39	>500
	GridGraph	37	10788	27	33	5824
	CLIP	8	934	1	19	484
BFS	X-Stream	16	6263	30	14	>720
	GridGraph	15	6262	29	13	4375
	CLIP	4	574	1	5	243

SSSP speedup: 1.8x-14.06x

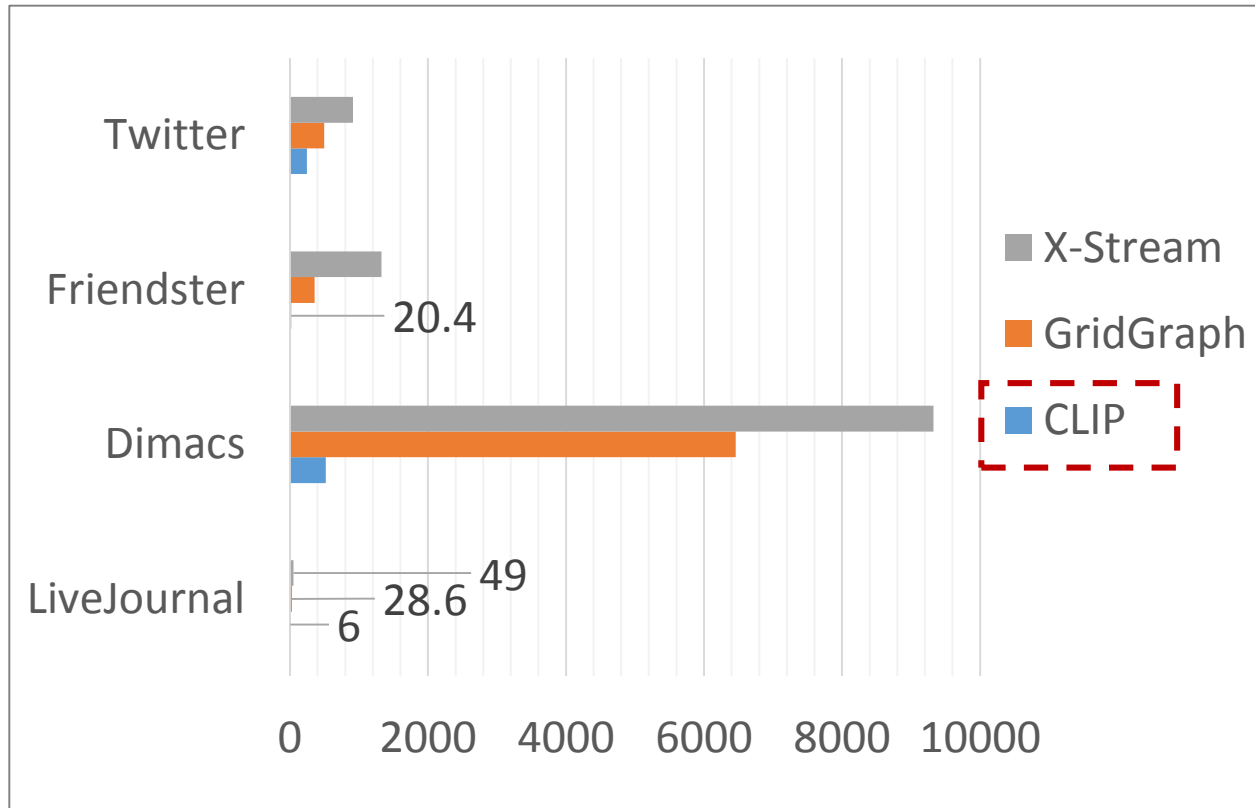
BFS speedup: 1.9x-12.08x



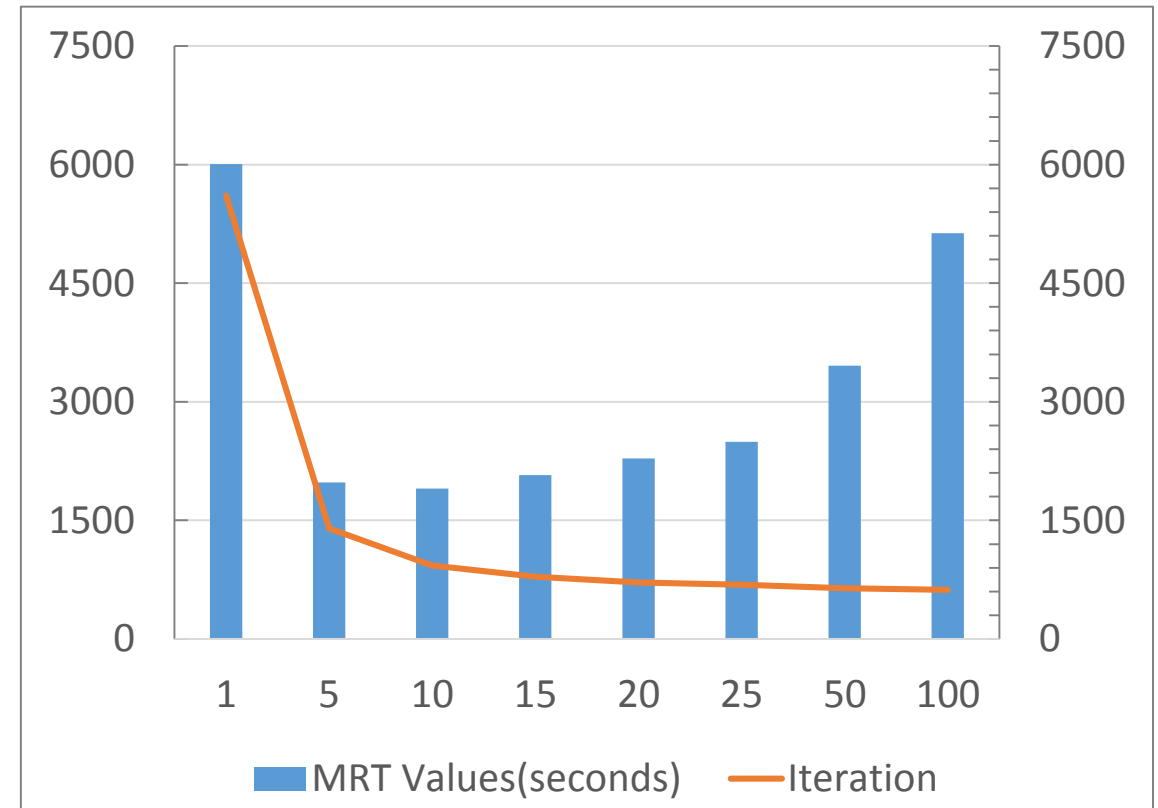
# 3.2 Evaluation—Loaded data reentry



### Total amount of disk I/O for SSSP (GB)



### SSSP on Dimacs graph



- All the values we reported is measured at **“MRT (Maximum Reentry Times) = 5”**
- Heuristically setting MRT to **5-10** is usually enough (less than 4% difference)

# 3.3 Evaluation—Beyond-neighborhood



Execution time (in seconds, external / semi-external / all-in-memory)

		LiveJournal	Dimacs	Friendster	Twitter	Yahoo
WCC	X-Stream	179.5 / 57.77 / 10.25	16633/ 6751 / 185.3	4521 / 2341 / -	1904 / 1194 / -	$\infty$ / $\infty$ / -
	GridGraph	22.32 / 13.8 / 3.57	6547 / 5757 / 422.5	967.5 / 466.6 / 82.95	431.5 / 272.3 / 62.3	19445/ 2916 / -
	CLIP	3.73 / 2.40 / 2.43	2.61 / 1.35 / 1.33	186 / 65.48 / 64.56	132.7 / 49.03 / 48.85	310.6 / 220.9 / -
MIS	X-Stream	422.1 / 152.6 / 13.06	103.4 / 41.42 / 5.95	9880 / 4867 / -	5513 / 3042 / -	$\infty$ / $\infty$ / -
	GridGraph	166.6 / 122.1 / 2.98	46.32 / 39.19 / 14.46	3945 / 3777 / 253.7	2510 / 2473 / 156.1	$\infty$ / $\infty$ / -
	CLIP	6.7 / 2.57 / 2.58	1.6 / 1.17 / 1.21	188.8 / 62.49 / 62.18	90.44 / 49.08 / 49.13	321.5 / 220.2 / -

The number of iterations

		LiveJornal	Dimacs	Friendster	Twitter	Yahoo
WCC	X-Stream	13	6263	24	16	>360
	GridGraph	37	6261	12	10	1368
	CLIP	1	1	1	1	1
MIS	X-Stream	54	31	65	53	>300
	GridGraph	53	31	64	52	>400
	CLIP	1	1	1	1	1



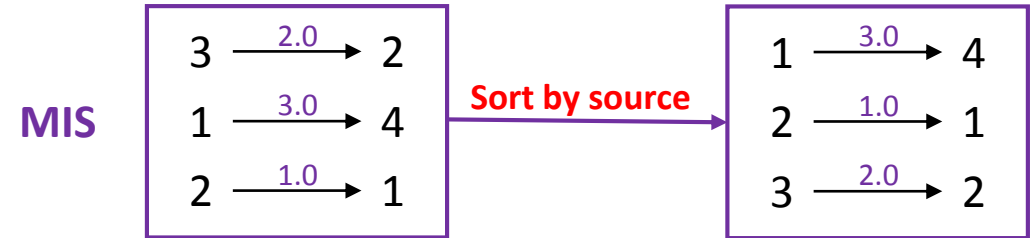
→ Sequential implementation  
 → But only need **one** iteration

WCC speedup: 5.6x-4264x  
 MIS speedup: 34x-60x

# 3.4 Evaluation—Preprocessing

## Preprocessing (Sorting) Evaluation

- Sorting is required by some algorithms and can be **amortized by reusing**



### Preprocessing time (seconds)

	GridGraph	CLIP (by source ID)	CLIP (by weight)
LiveJornal	4.57	5.06	9.06
Dimacs	4.09	5.12	6.81
Friendster	185.5	145.3	228.9
Twitter	160.3	126.2	188.1
Yahoo	1616	1410	1563

### Total time on Friendster (seconds)

		Proc.	Exec.	Total
MIS	X-Stream	0	4867	4867
	GridGraph	185.5	3777	3962.5
	<b>Clip</b>	<b>145.3</b>	<b>62.49</b>	<b>207.79</b>

MIS 19.07x faster than GridGraph

MIS 23.42x faster than X-Stream

# 4 Conclusion

**CLIP:** A single-machine graph processing system

**Focus on squeezing out all the value of loaded data**

## Two techniques:

- **Reentry of Loaded Data**

Processing the loaded data **multiple times** rather than once



Reducing iterations



Reducing the total amount of disk I/O



Reducing execution time

- **Beyond the Neighborhood**

The user-defined function is allowed to update **an arbitrary** vertex's property

Up to **tens or even thousands of** times speedup

**Thank You!!**