

Dependency-Driven Disk-based Graph Processing

Keval Vora

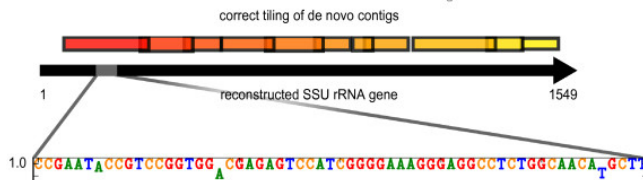
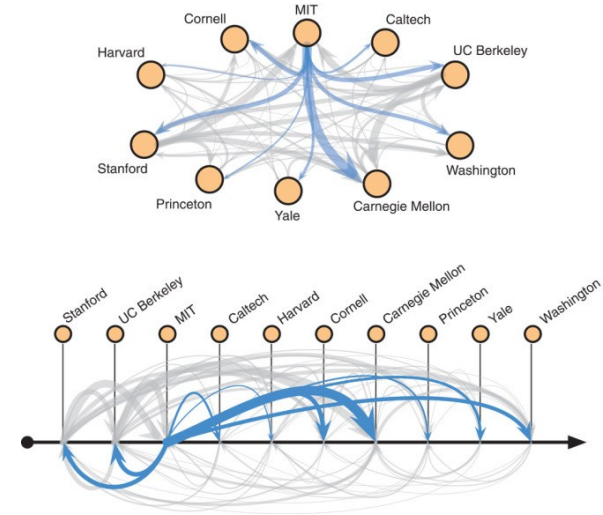
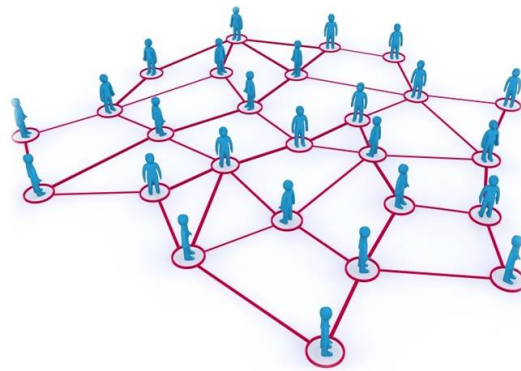
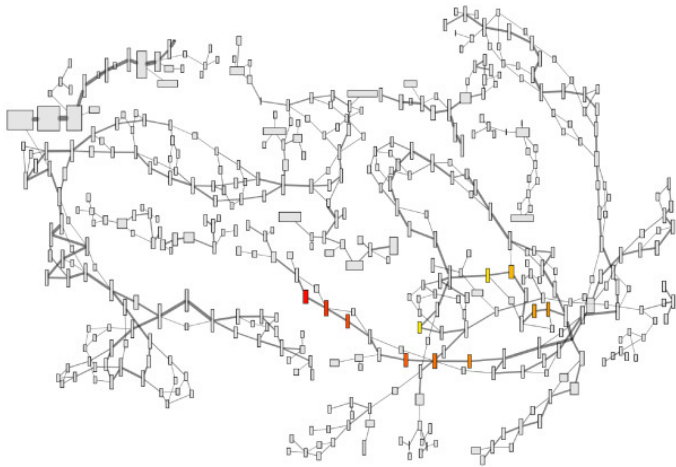
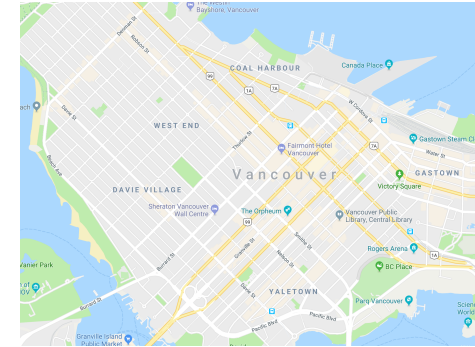
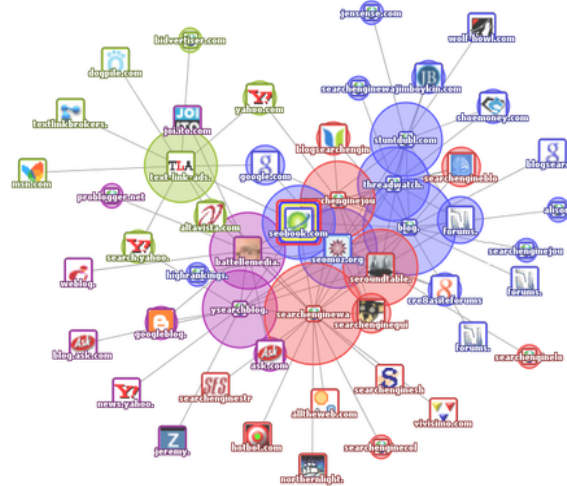
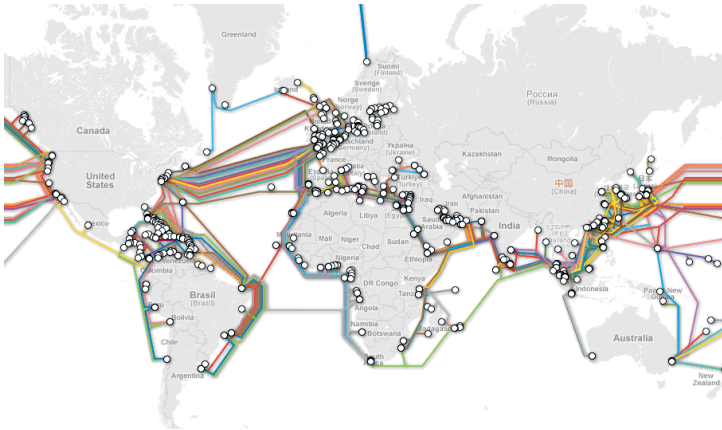
Simon Fraser University

USENIX ATC'19 - Renton, WA

July 11, 2019

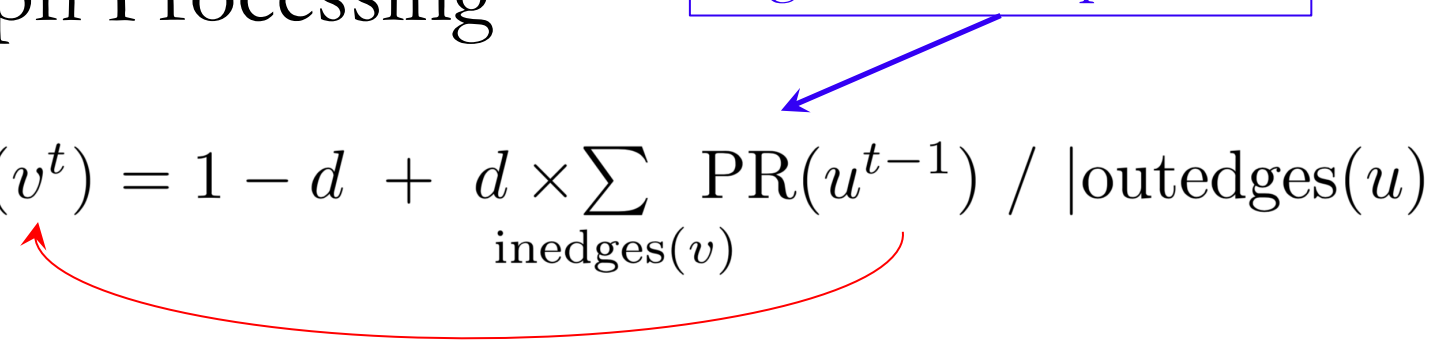
<https://github.com/pdclab/lumos>

Graph Processing

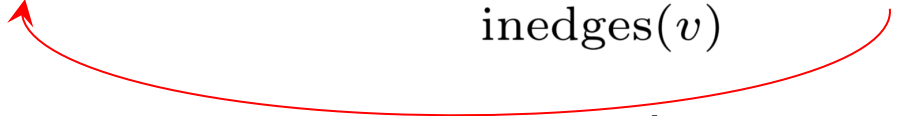


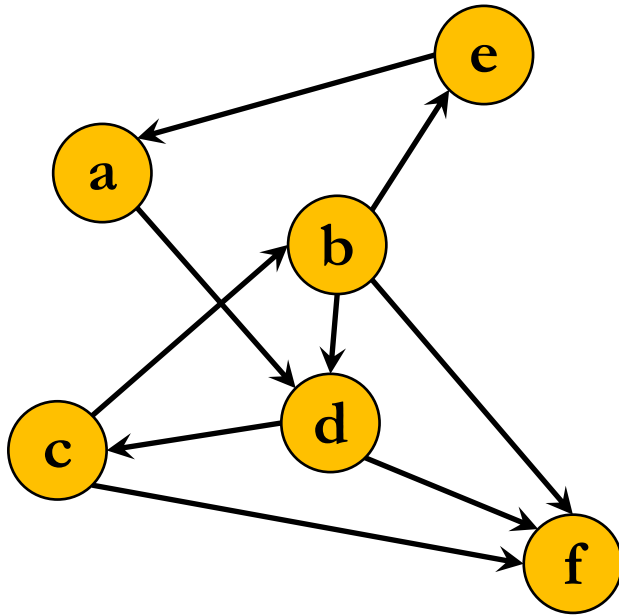
Iterative Graph Processing

PageRank computation

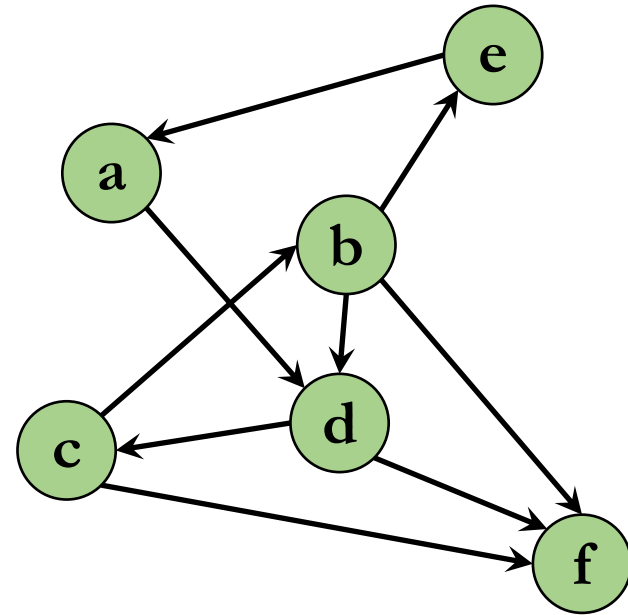
$$\text{PR}(v^t) = 1 - d + d \times \sum_{\text{inedges}(v)} \text{PR}(u^{t-1}) / |\text{outedges}(u)|$$


Iterative Graph Processing

$$\text{PR}(v^t) = 1 - d + d \times \sum_{u \in \text{inedges}(v)} \text{PR}(u^{t-1}) / |\text{outedges}(u)|$$


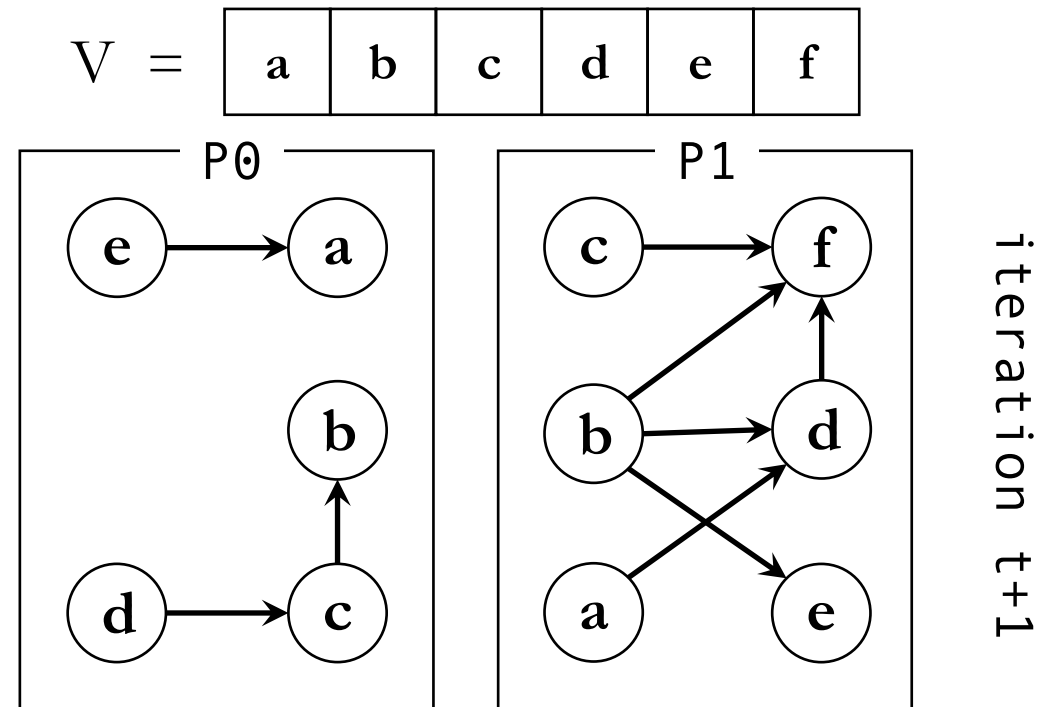
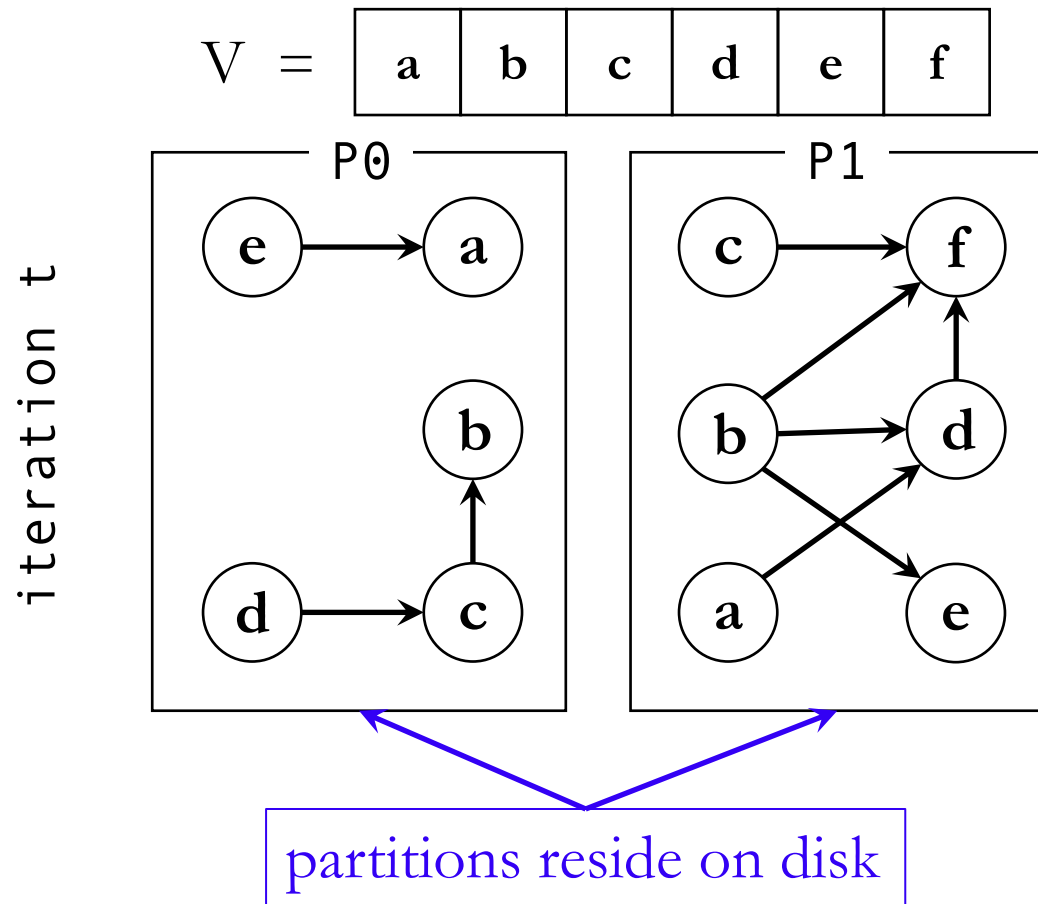


iteration t

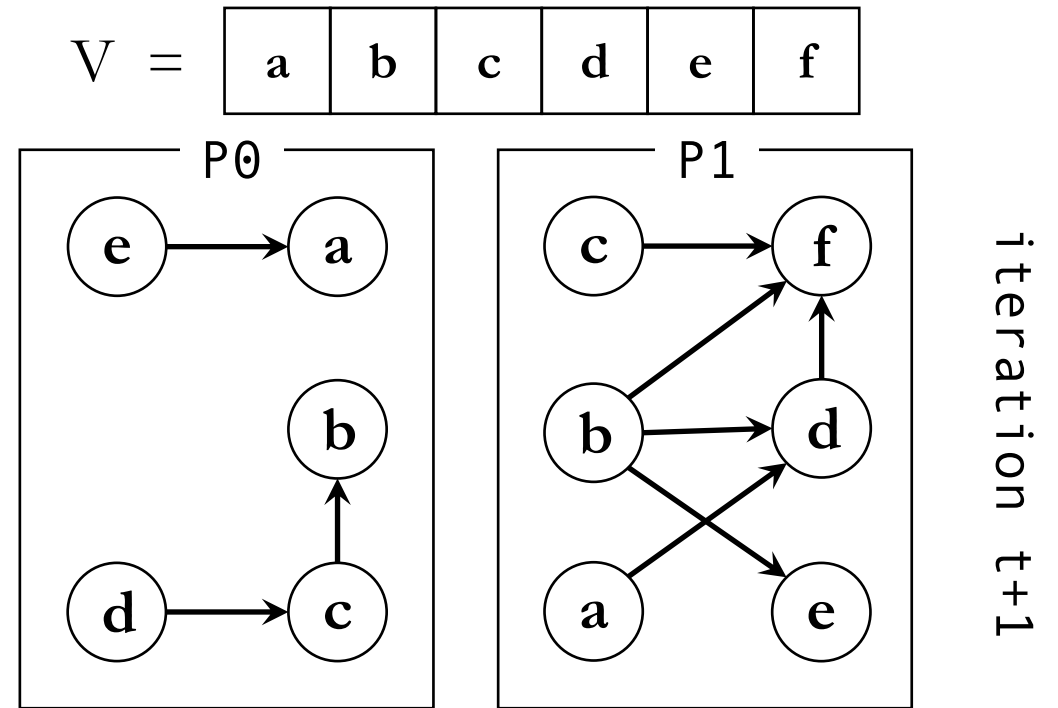
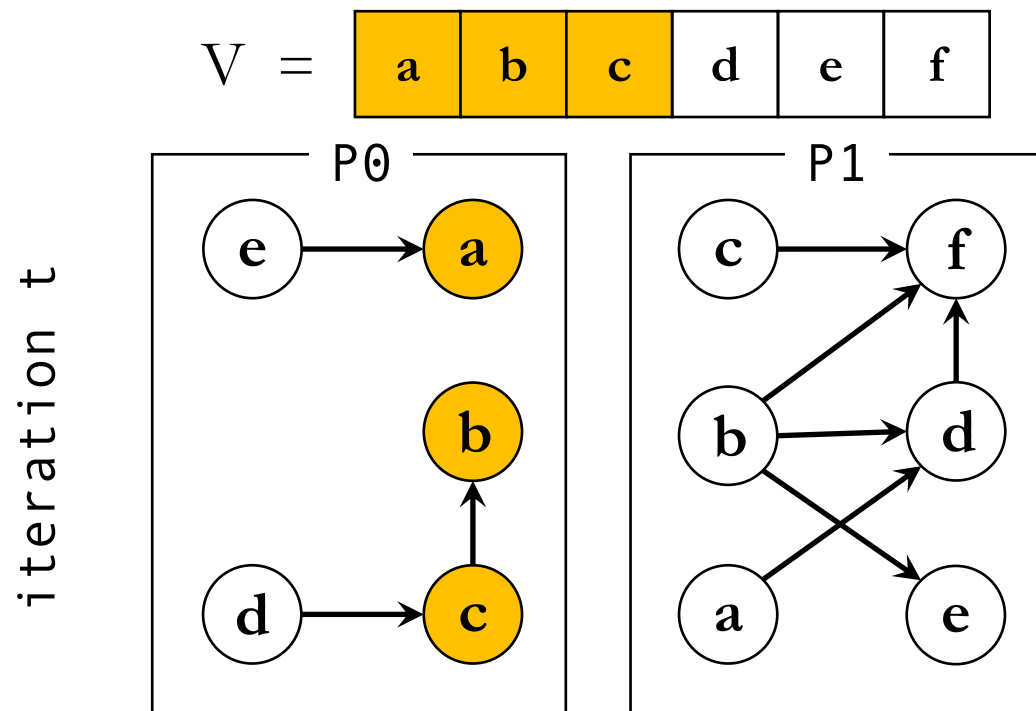


iteration t+1

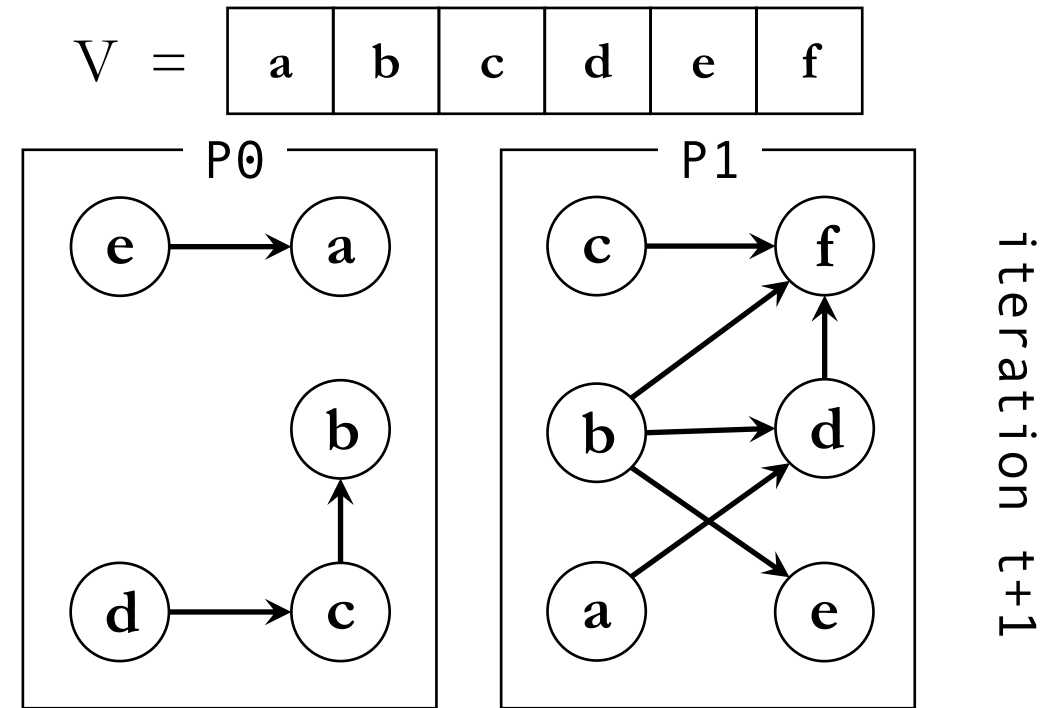
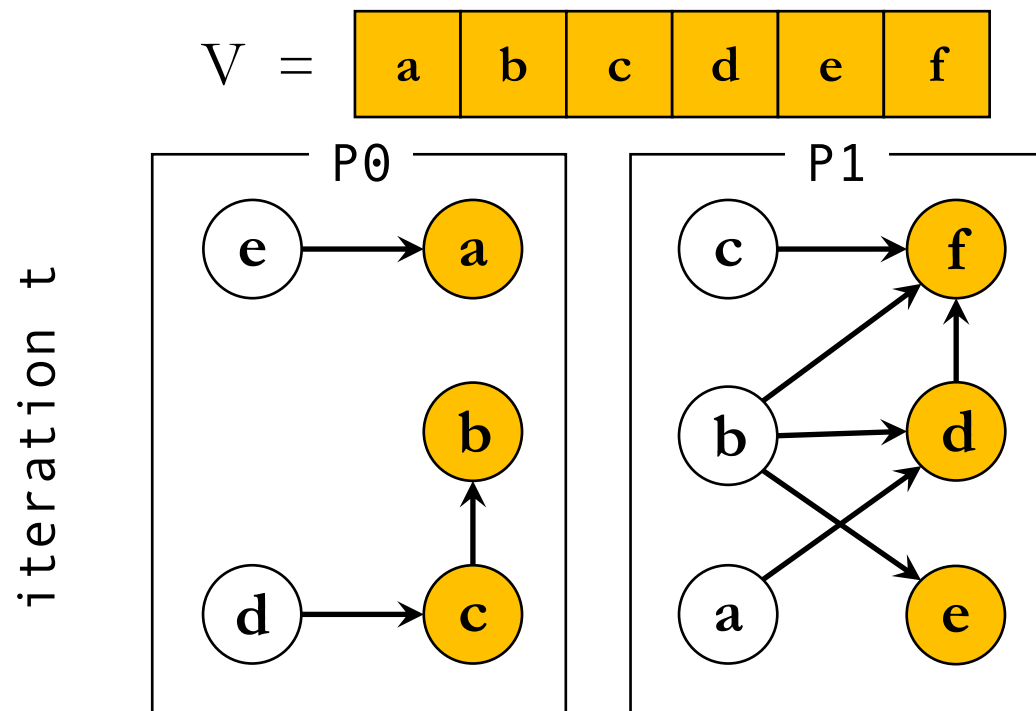
Out-of-Core Graph Processing



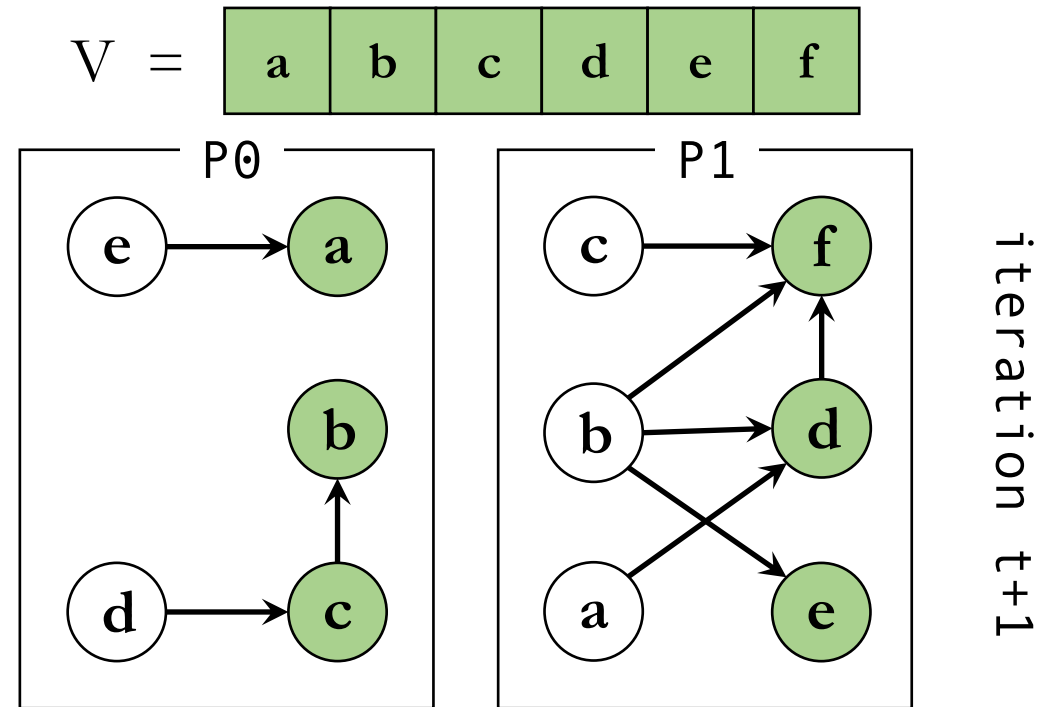
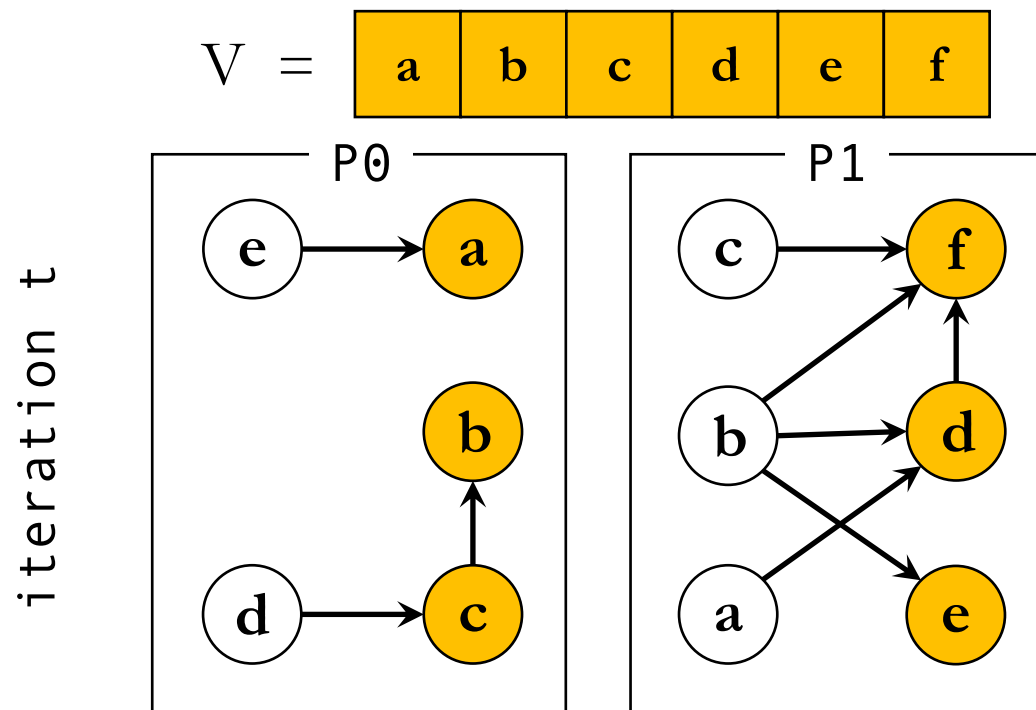
Out-of-Core Graph Processing



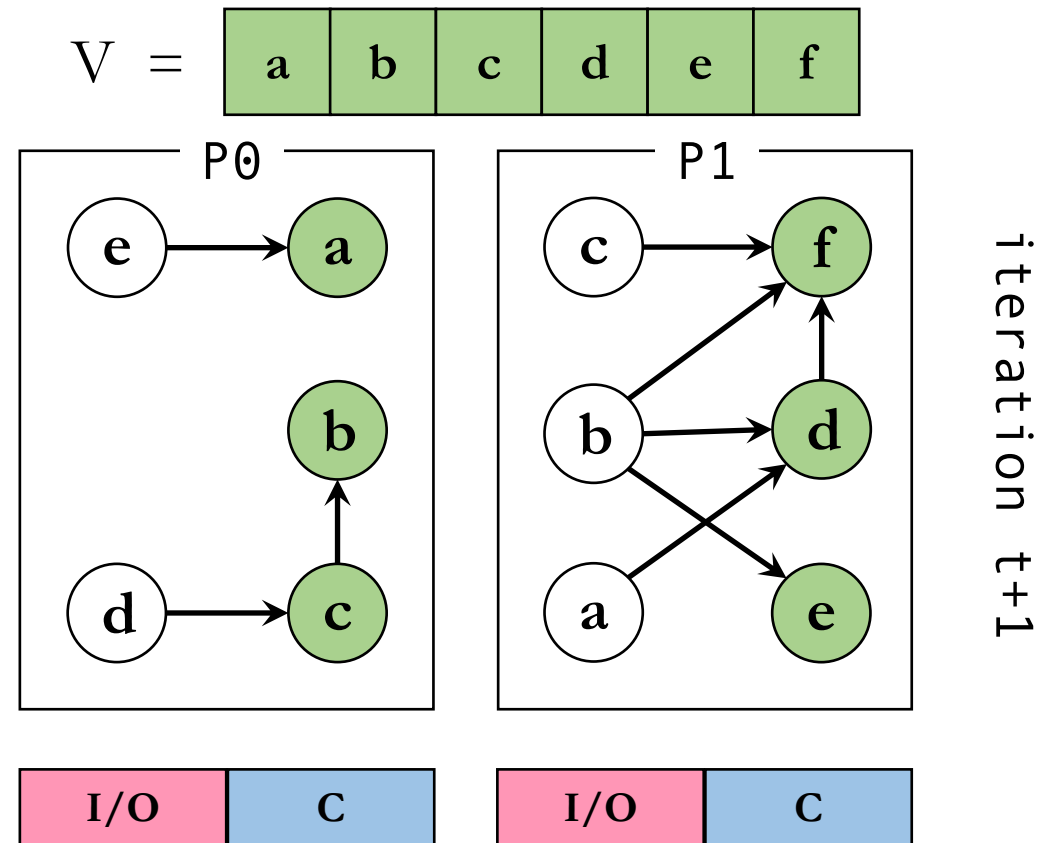
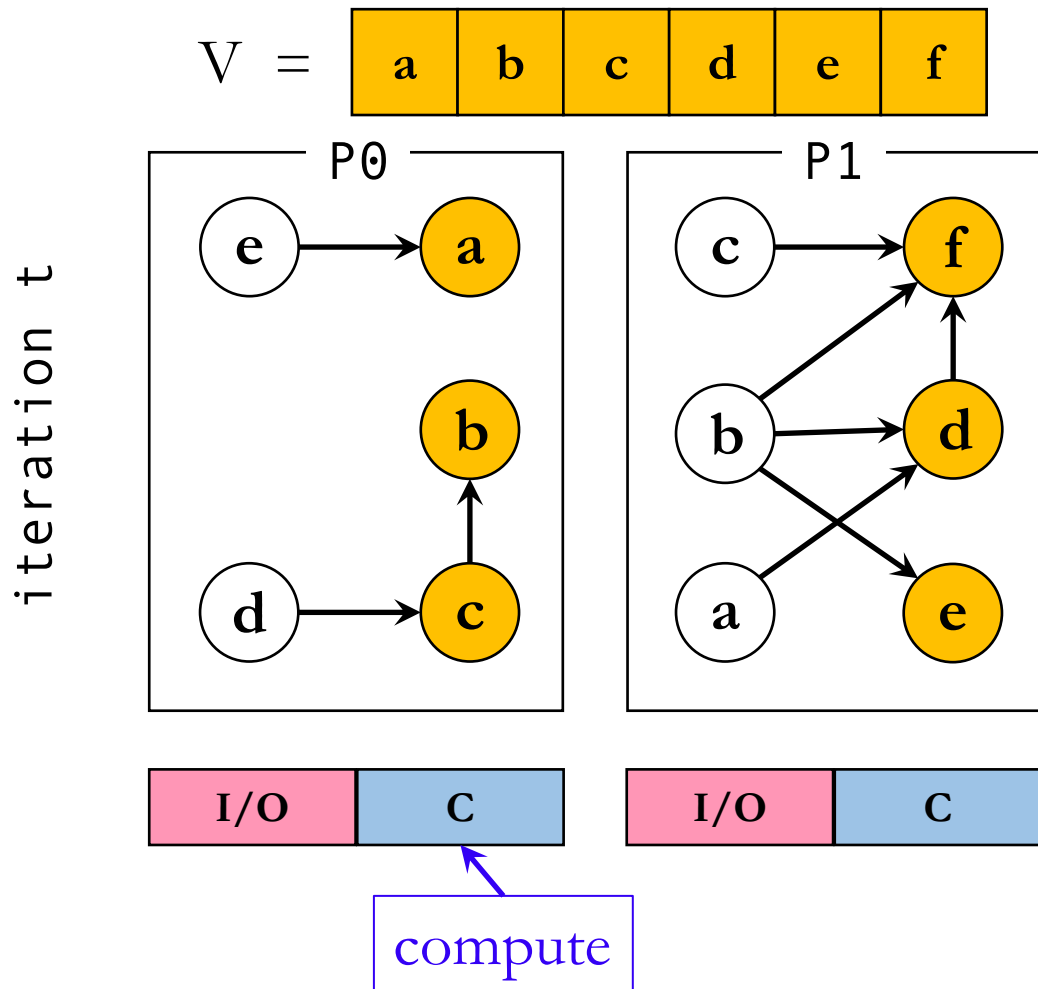
Out-of-Core Graph Processing



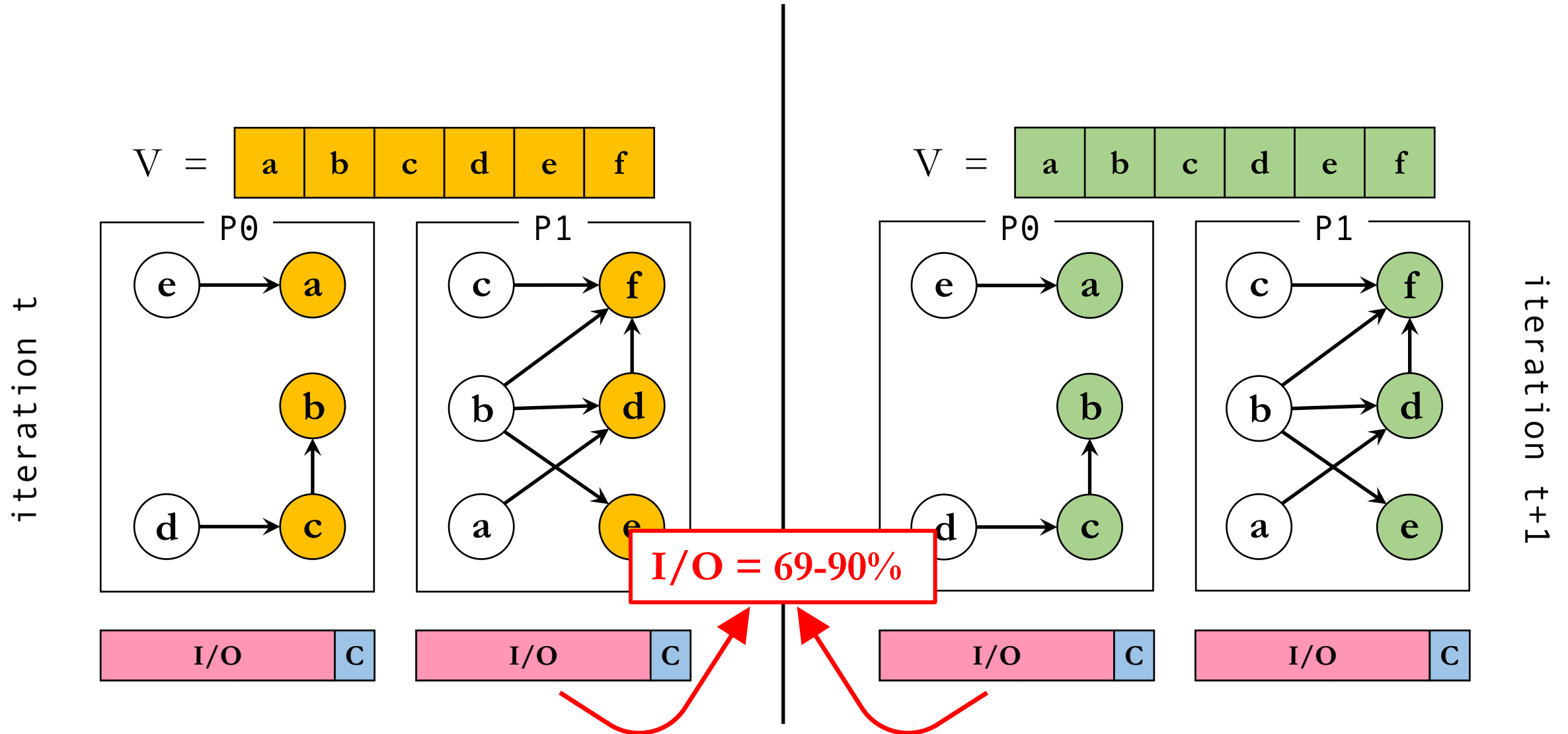
Out-of-Core Graph Processing



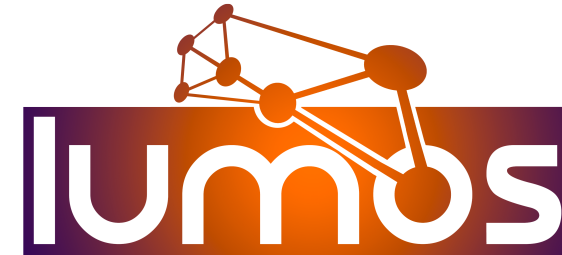
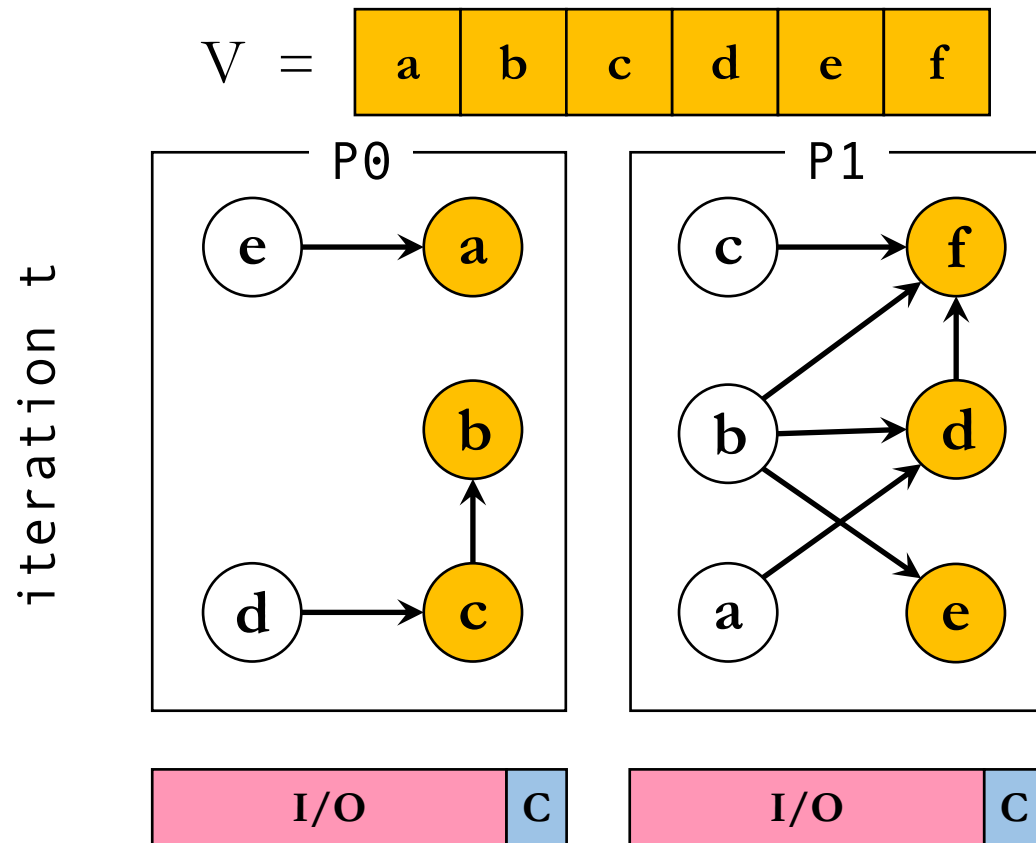
Out-of-Core Graph Processing



Out-of-Core Graph Processing



Lumos



- **Out-of-Order** processing
 - Amortize I/O across multiple iterations
- Guarantee **Bulk Synchronous Parallel** Semantics
- **Dependency-Driven Cross-Iteration** value propagation
 - *Safely* compute future values

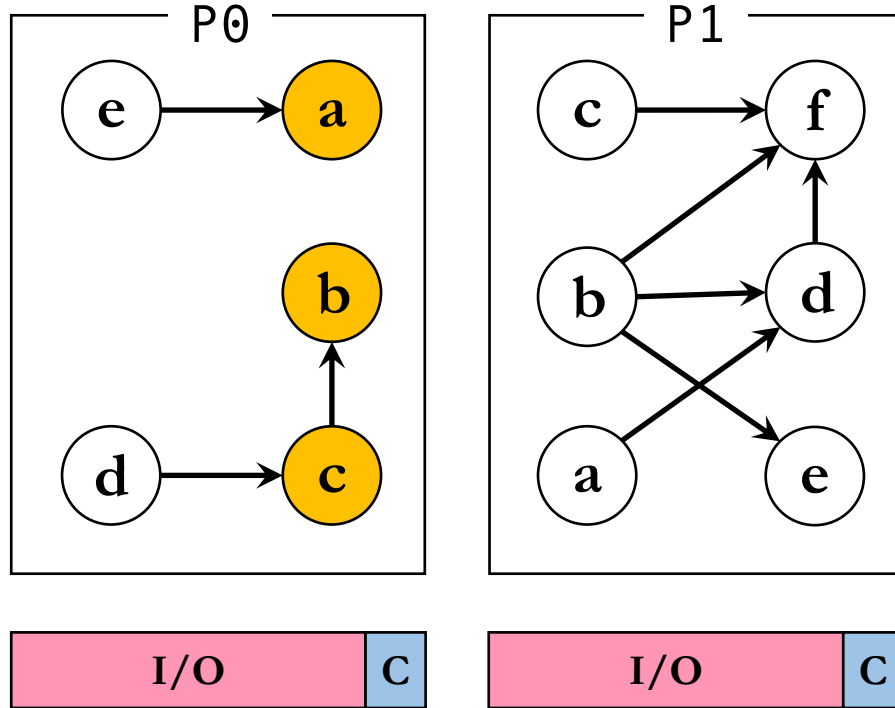
Computing Future Values

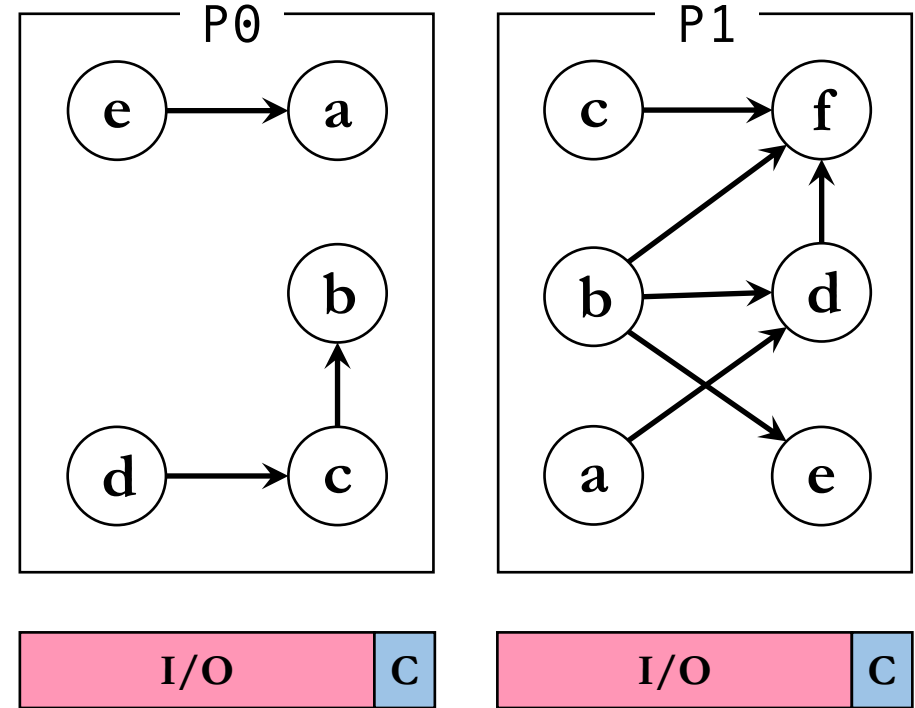
future vertex values

$$V^F = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$

$$V = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$

iteration t

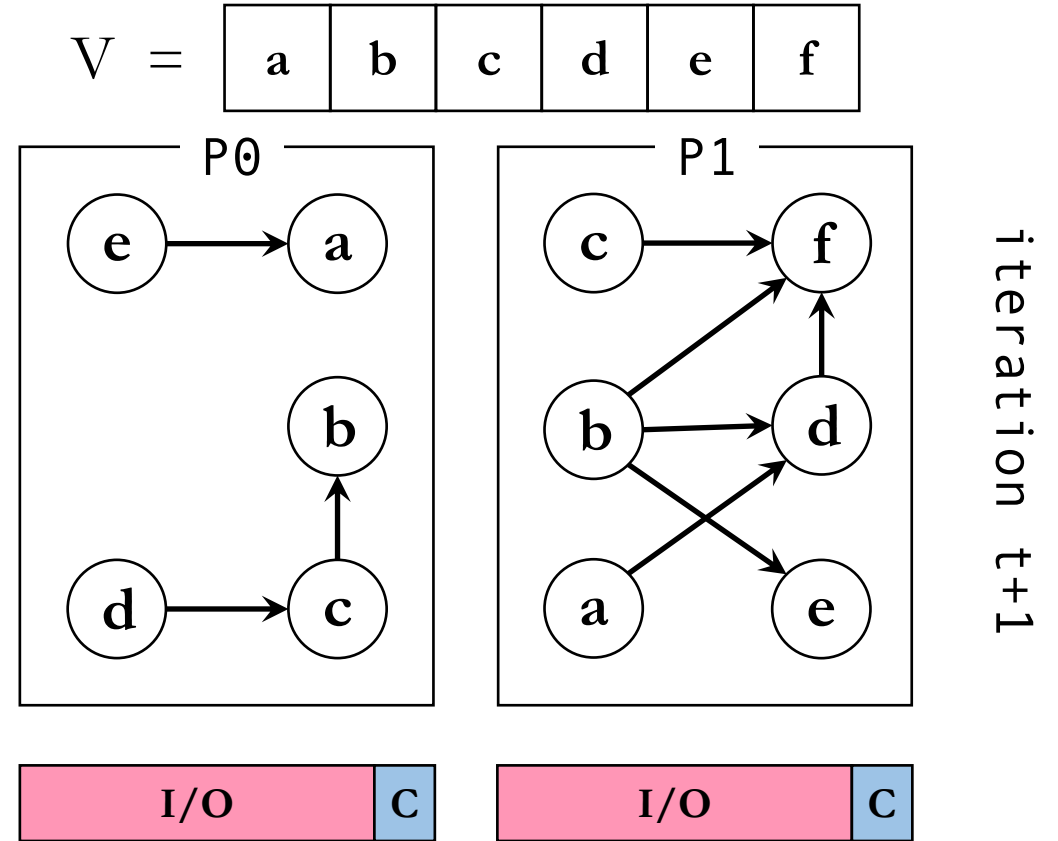
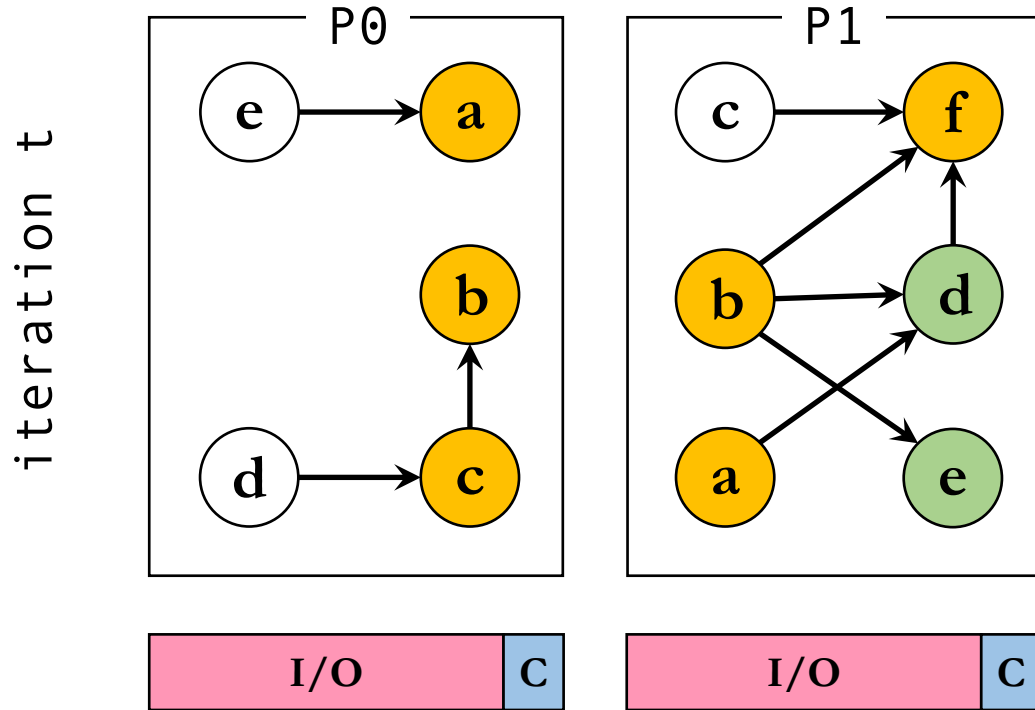


$$V = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$


iteration $t+1$

Computing Future Values

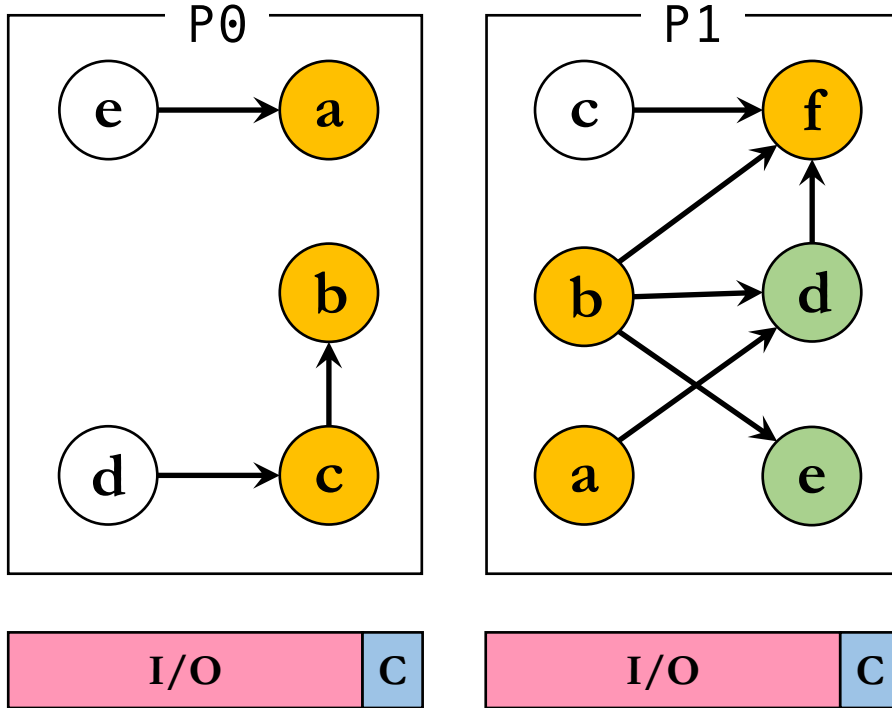
$$V^F = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$

$$V = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$


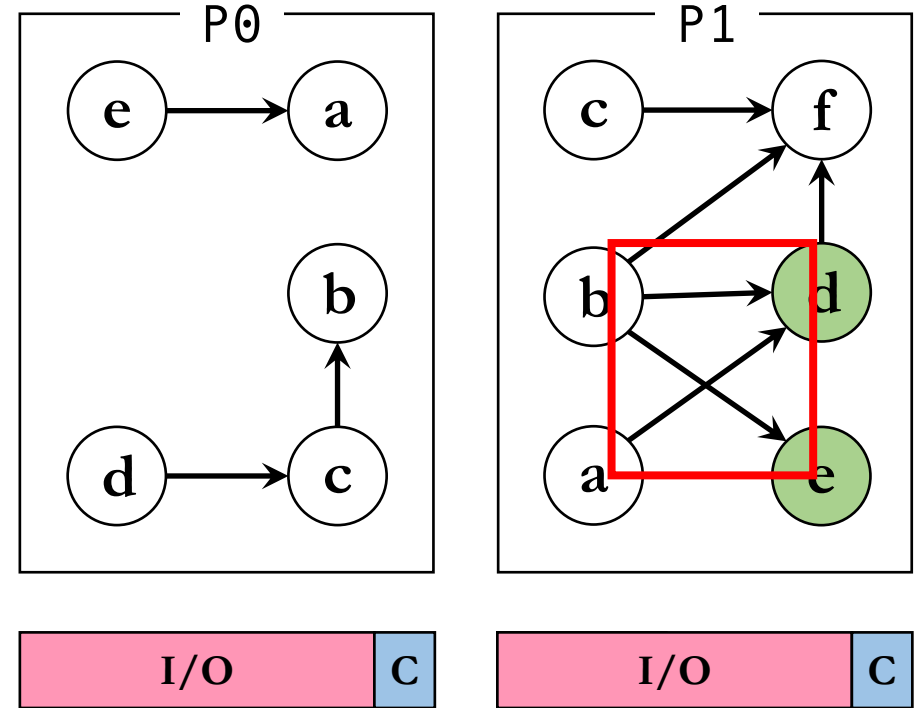
Computing Future Values

$$V^F = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$

$$V = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$

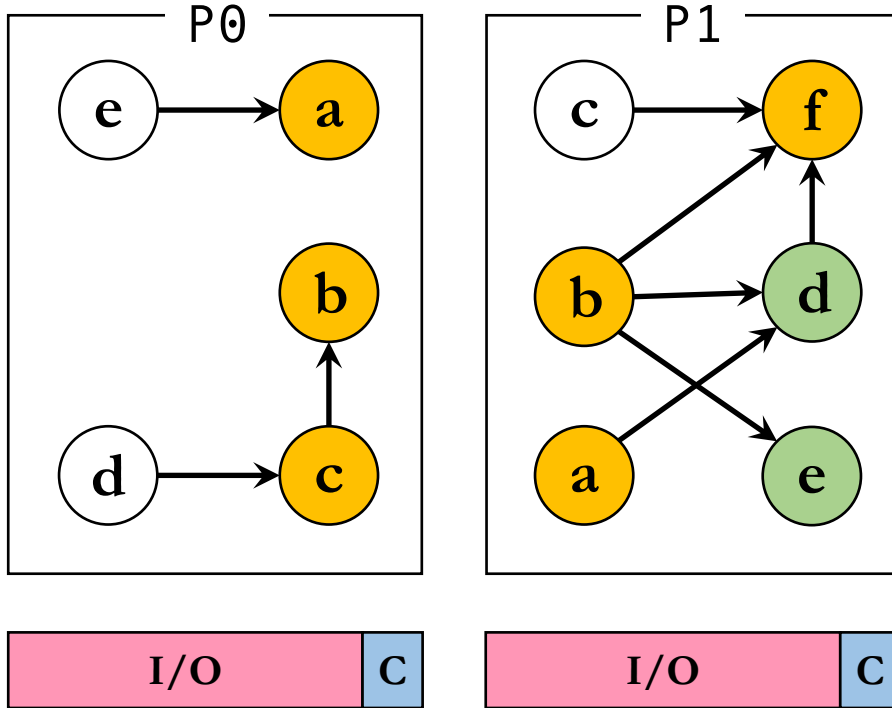
iteration t 

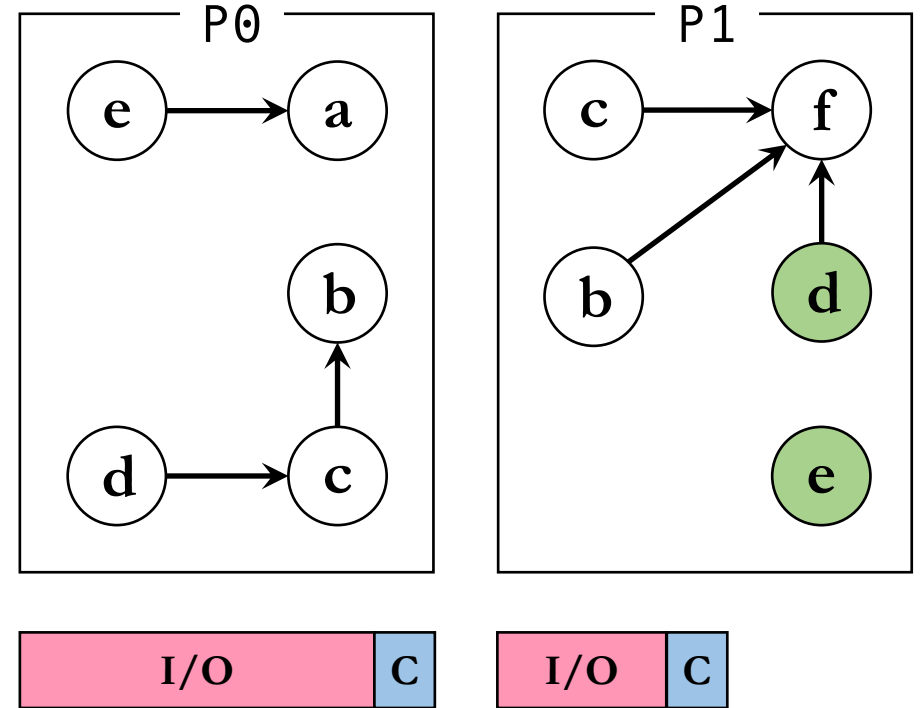
$$V = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$

iteration $t+1$

Computing Future Values

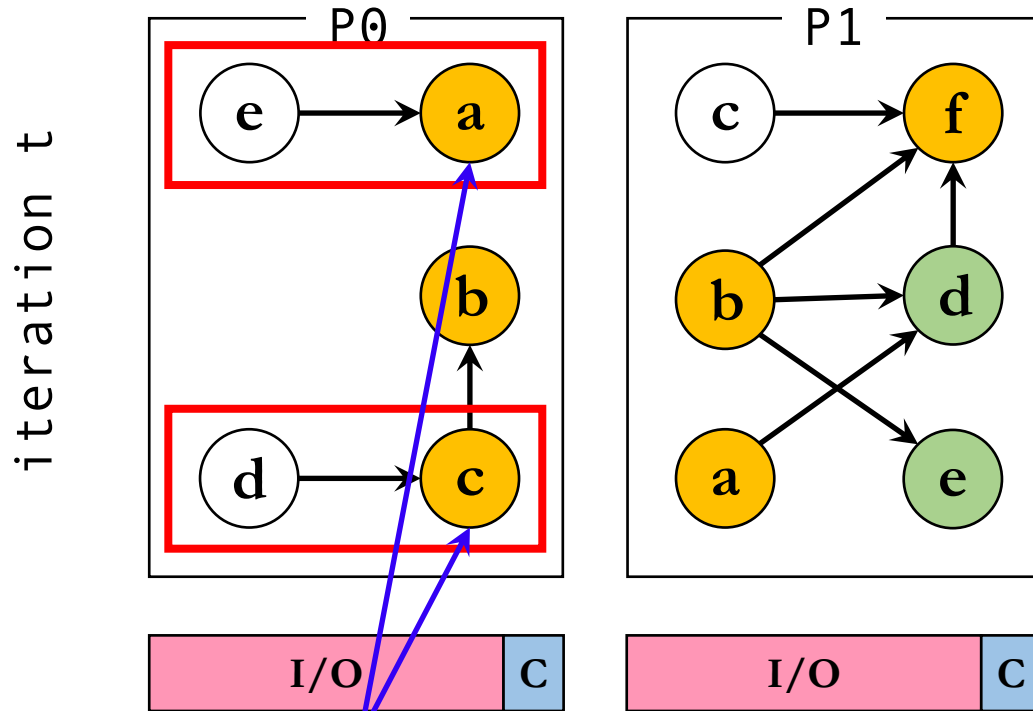
$$V^F = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$

$$V = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$
iteration t 

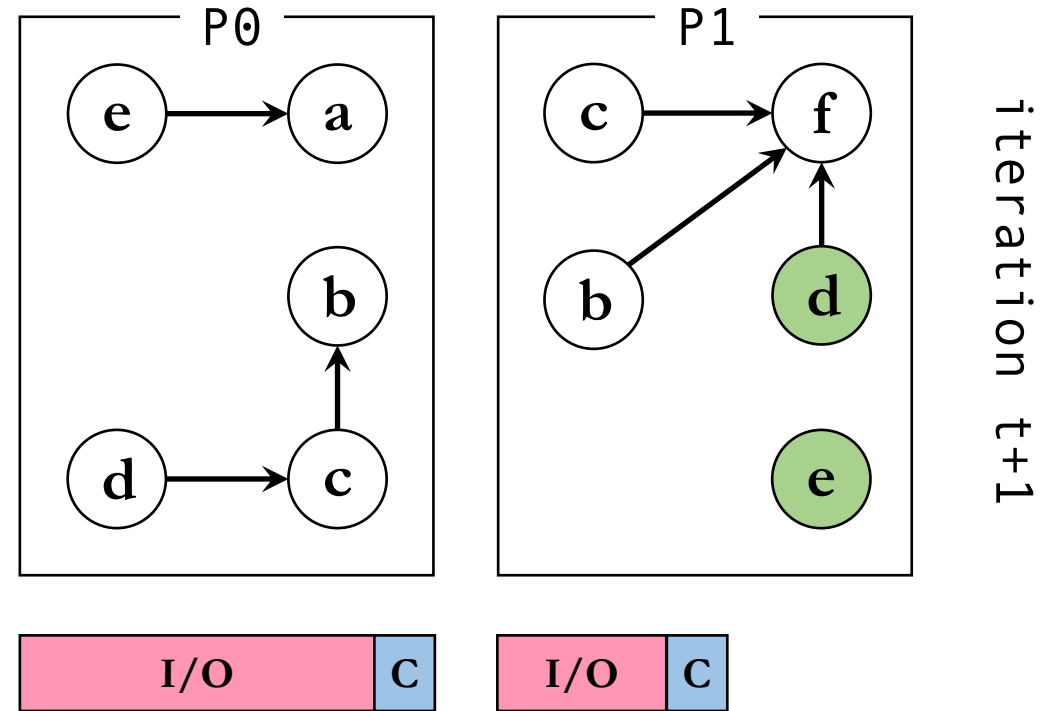
$$V = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$
iteration $t+1$

Computing Future Values

$$V^F = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$

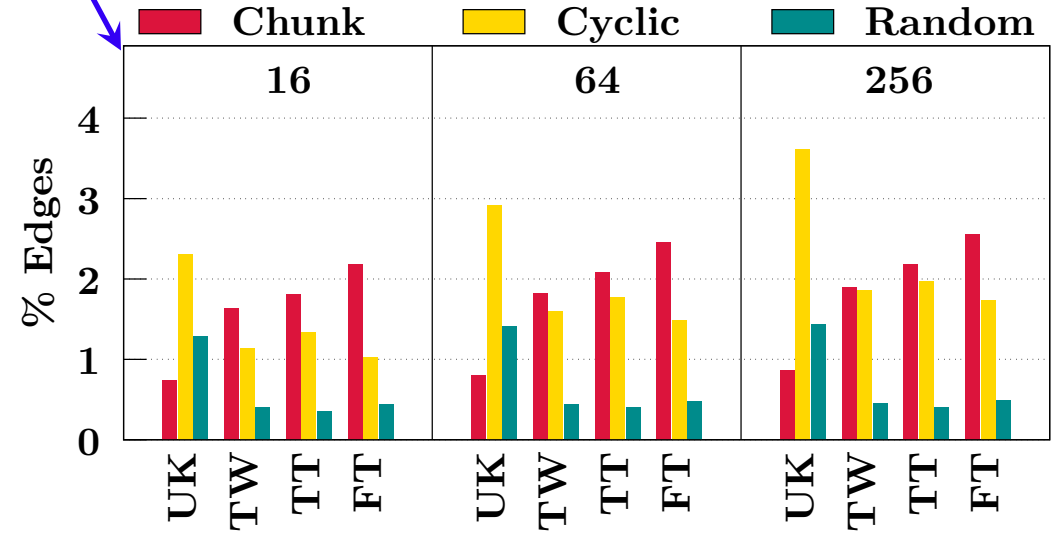
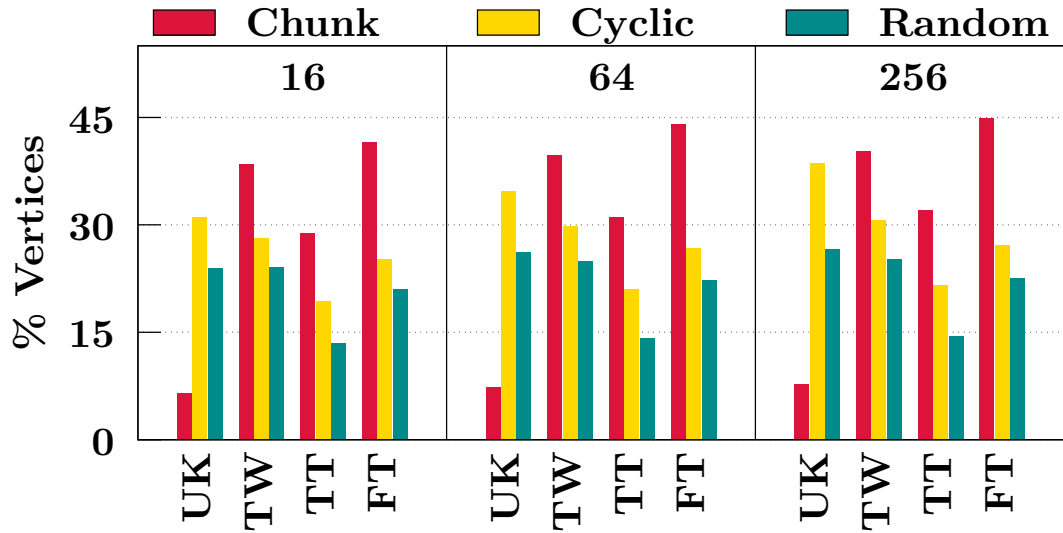
$$V = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$


future dependencies not met

$$V = \begin{array}{|c|c|c|c|c|c|} \hline a & b & c & d & e & f \\ \hline \end{array}$$


Computing Future Values

30-45% vertices compute future values, and they contribute for less than 4% edges



	Edges	Vertices
UK	1B	39.5M
TW	1.5B	41.7M
TT	2B	52.6M
FT	2.5B	68.3M

Graph Computation

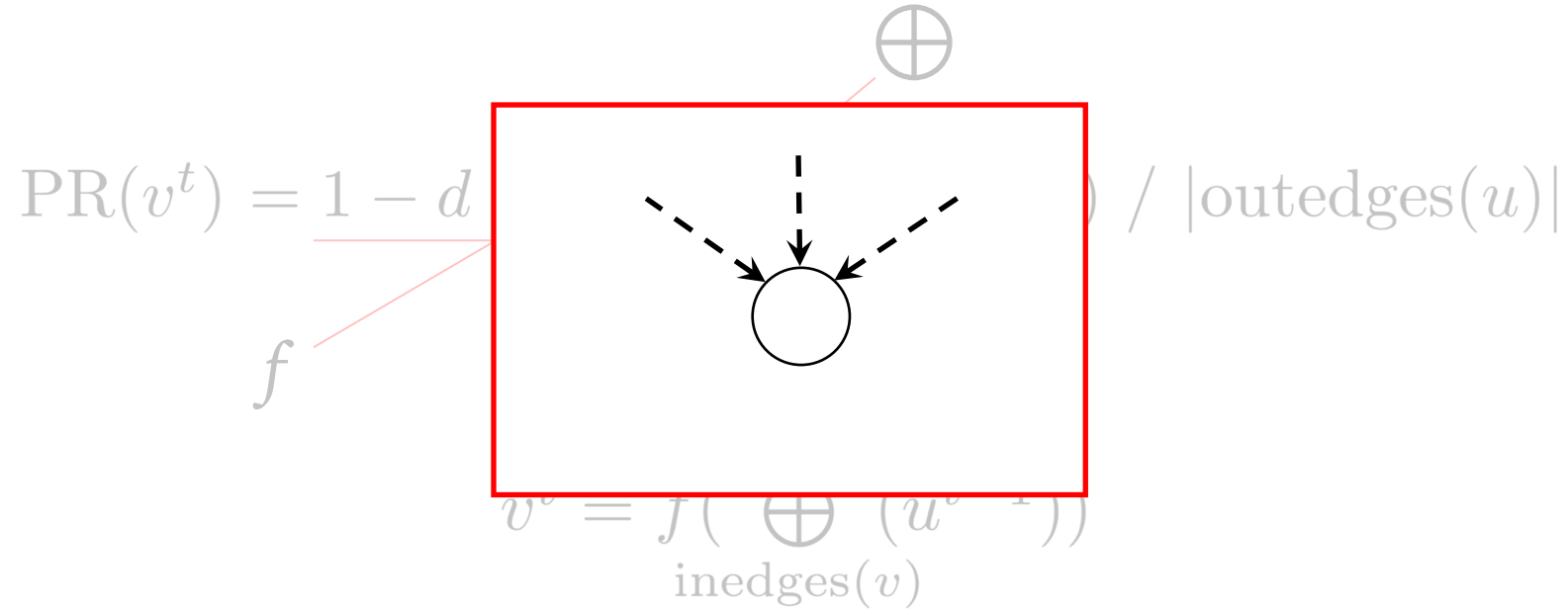
$$\text{PR}(v^t) = \frac{1 - d + d \times \sum_{u \in \text{inedges}(v)} \text{PR}(u^{t-1})}{|\text{outedges}(v)|}$$

$$v^t = f\left(\bigoplus_{u \in \text{inedges}(v)} (u^{t-1})\right)$$

Partial Aggregation

$$\text{inedges}(v) = X \cup Y \implies \bigoplus_{u \in \text{inedges}(v)} (u^t) = \bigoplus \left(\bigoplus_X (u^t), \bigoplus_Y (u^t) \right)$$

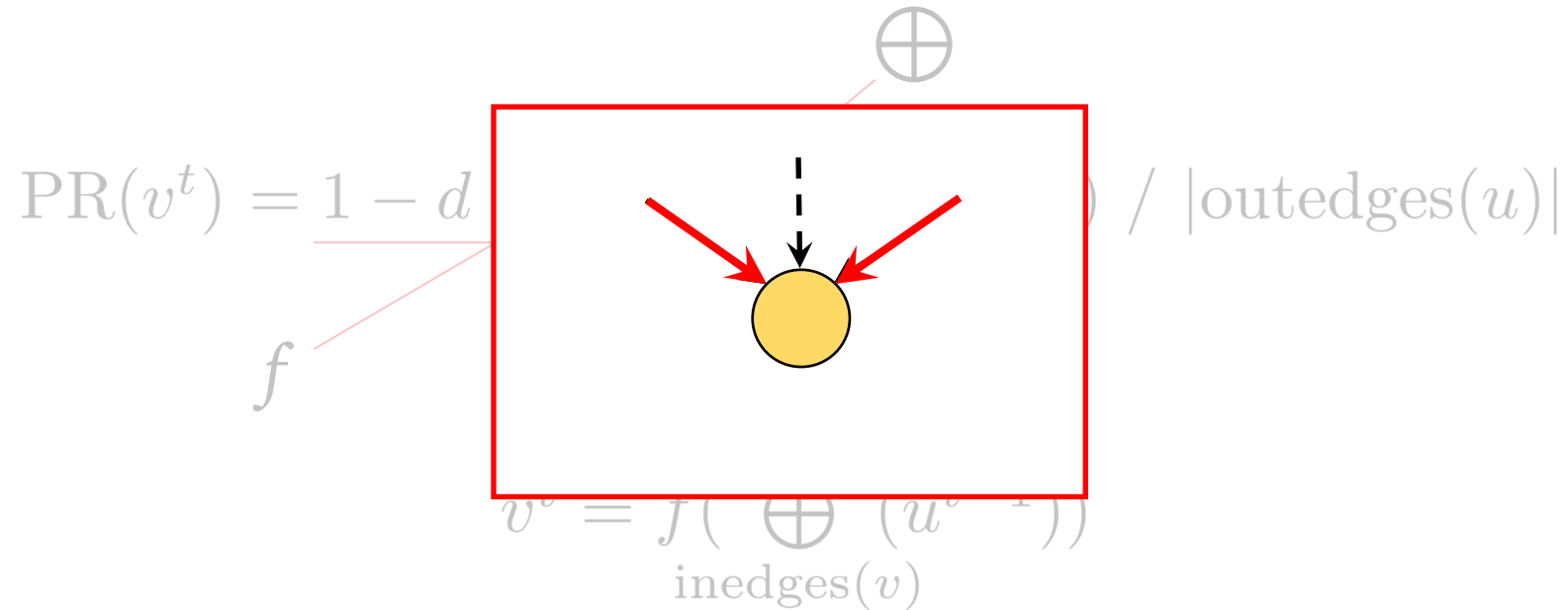
Graph Computation



Partial Aggregation

$$inedges(v) = X \cup Y \implies \bigoplus_{inedges(v)} (u^t) = \bigoplus \left(\bigoplus_X (u^t), \bigoplus_Y (u^t) \right)$$

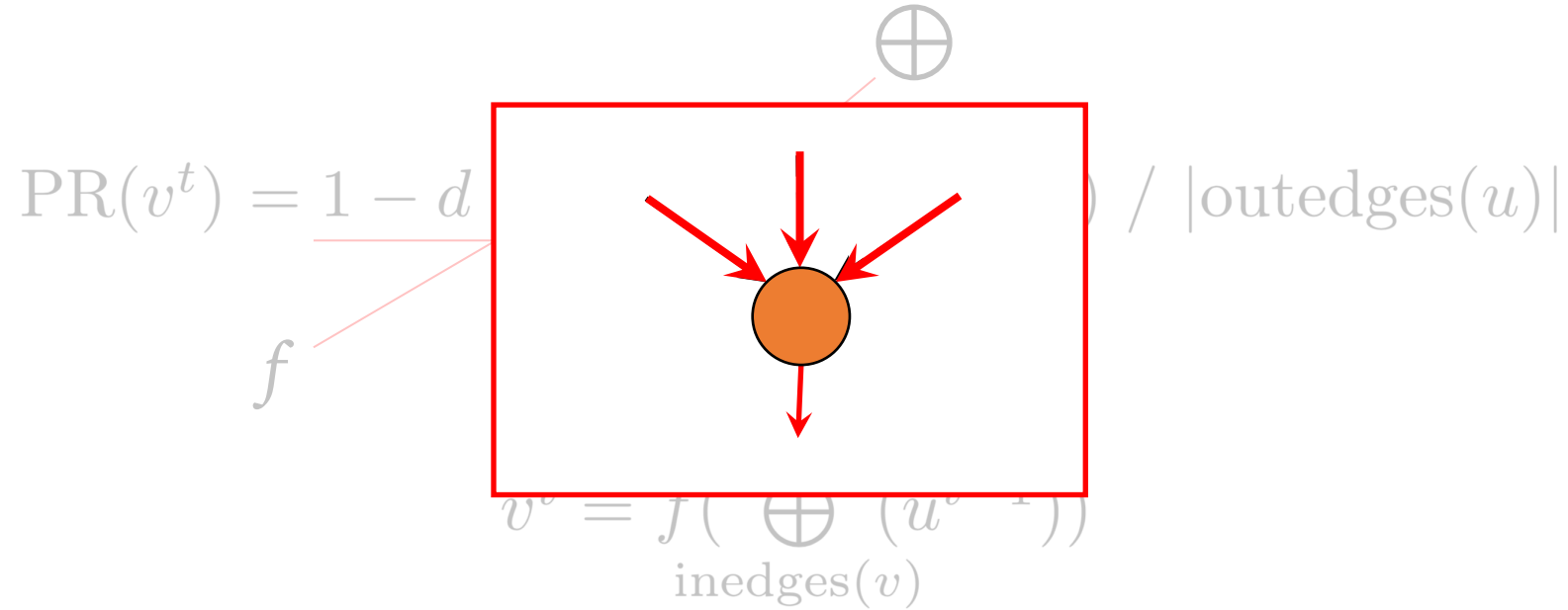
Graph Computation



Partial Aggregation

$$inedges(v) = X \cup Y \implies \bigoplus_{inedges(v)} (u^t) = \bigoplus \left(\bigoplus_X (u^t), \bigoplus_Y (u^t) \right)$$

Graph Computation

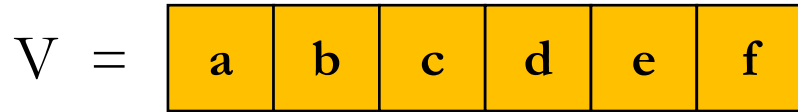
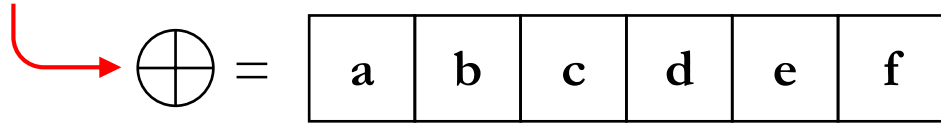


Partial Aggregation

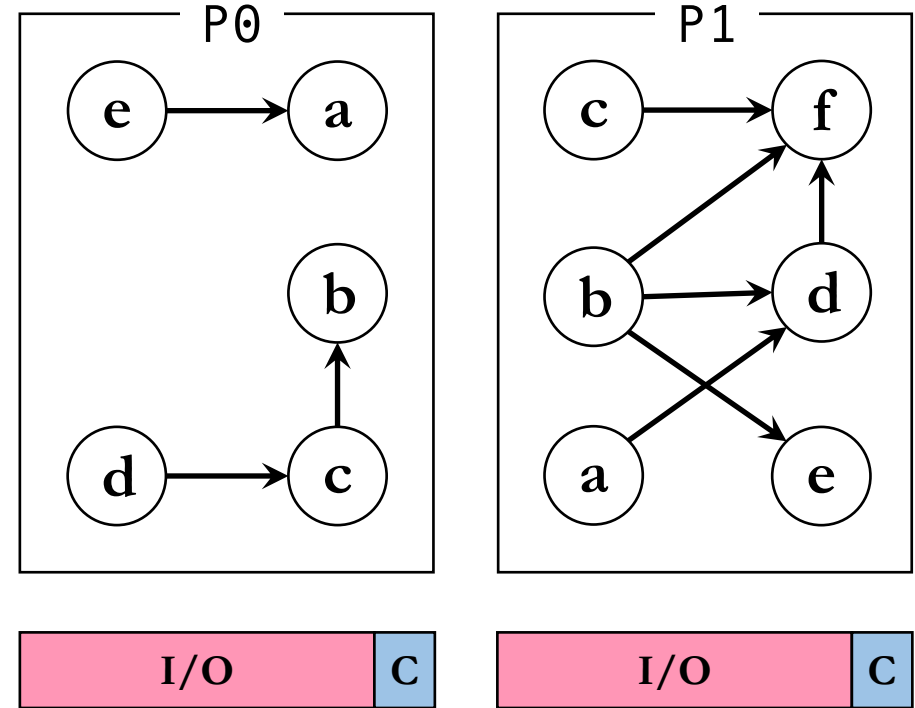
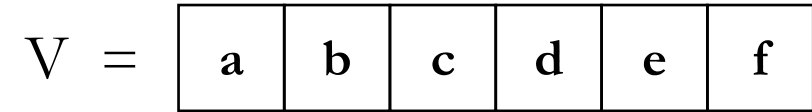
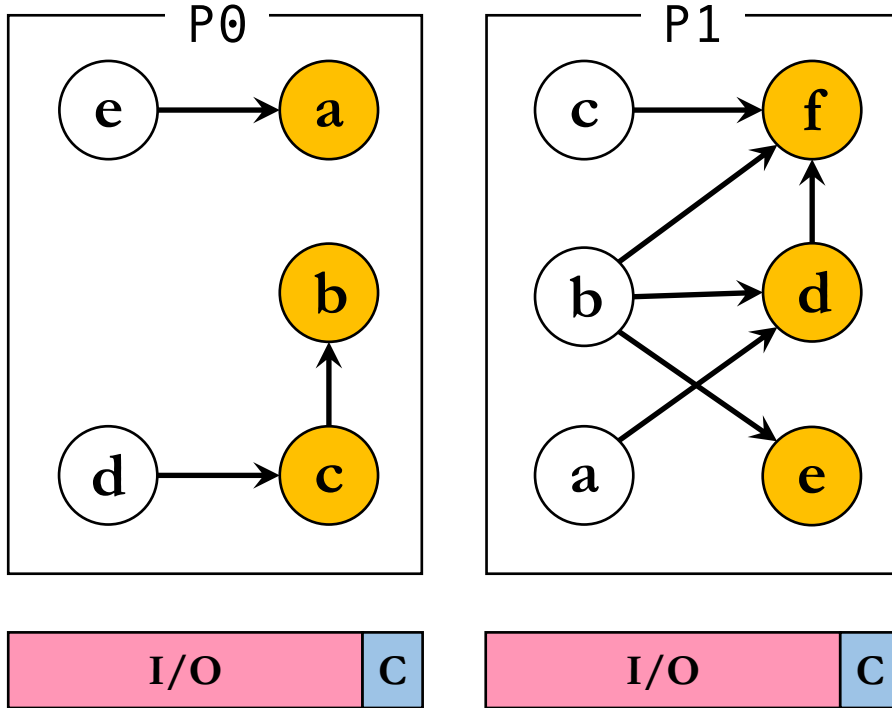
$$inedges(v) = X \cup Y \implies \bigoplus_{inedges(v)} (u^t) = \bigoplus \left(\bigoplus_X (u^t), \bigoplus_Y (u^t) \right)$$

Cross-Iteration Value Propagation

Partial Aggregation



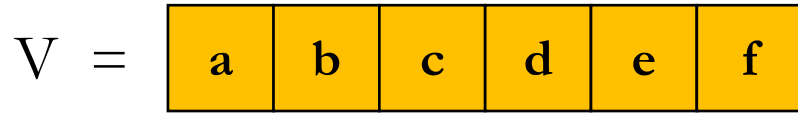
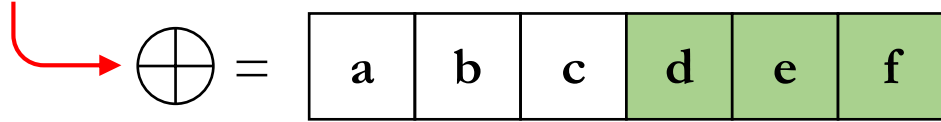
iteration t



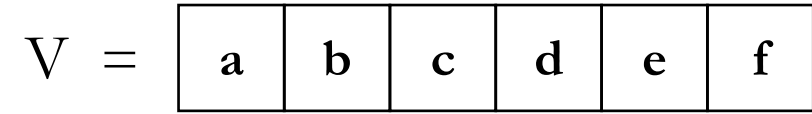
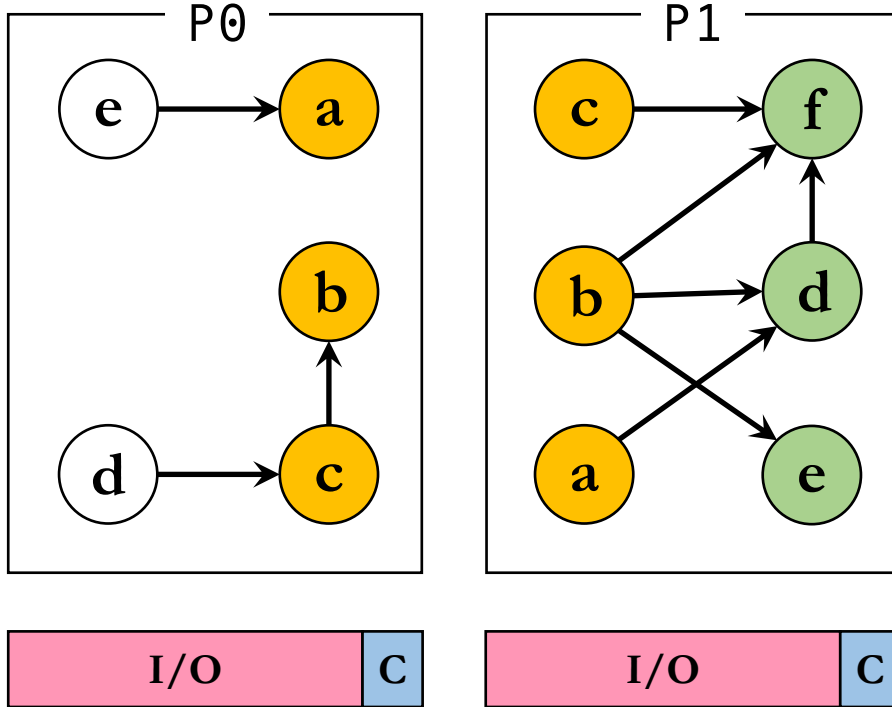
iteration t+1

Cross-Iteration Value Propagation

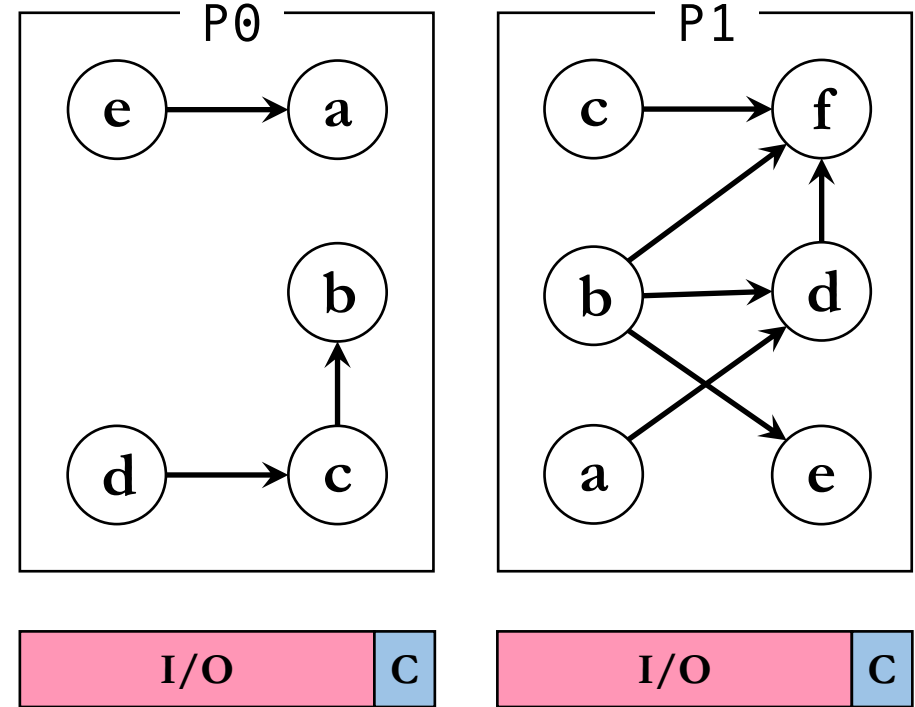
Partial Aggregation



iteration t

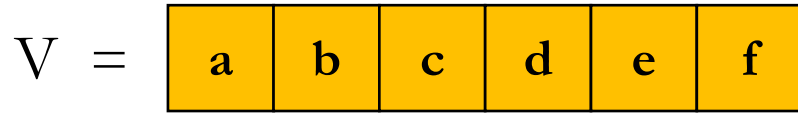
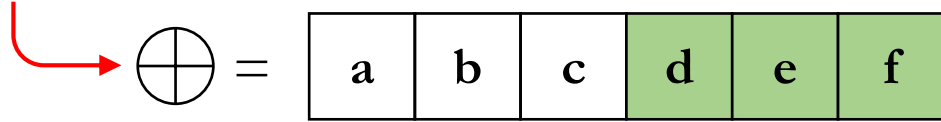


iteration t+1

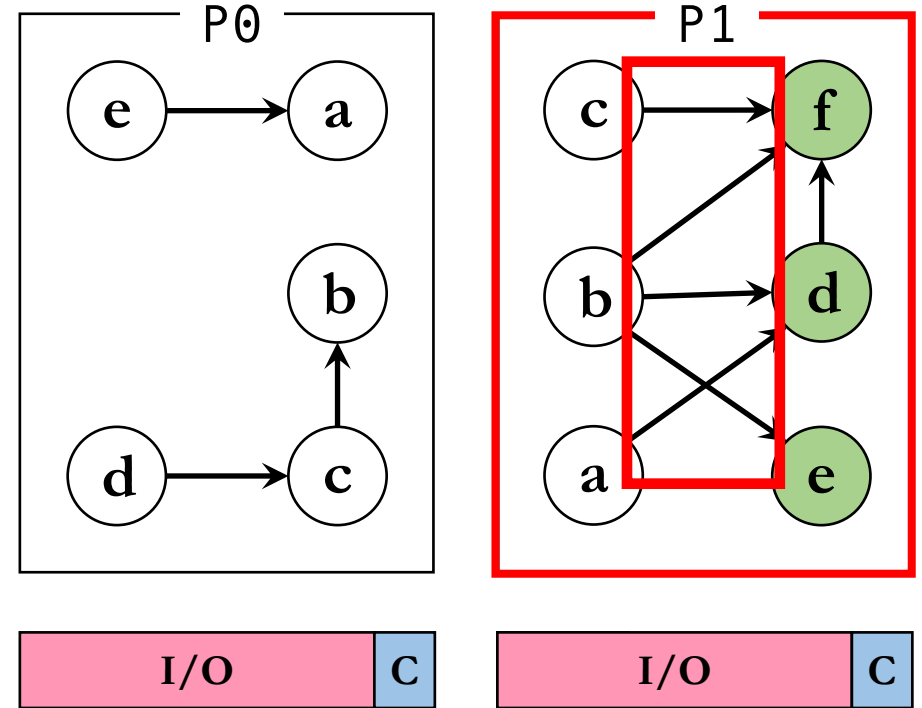
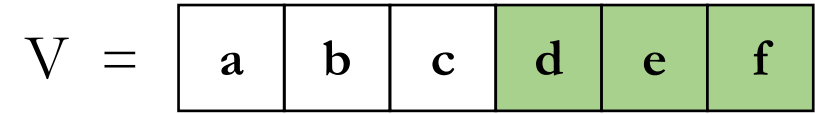
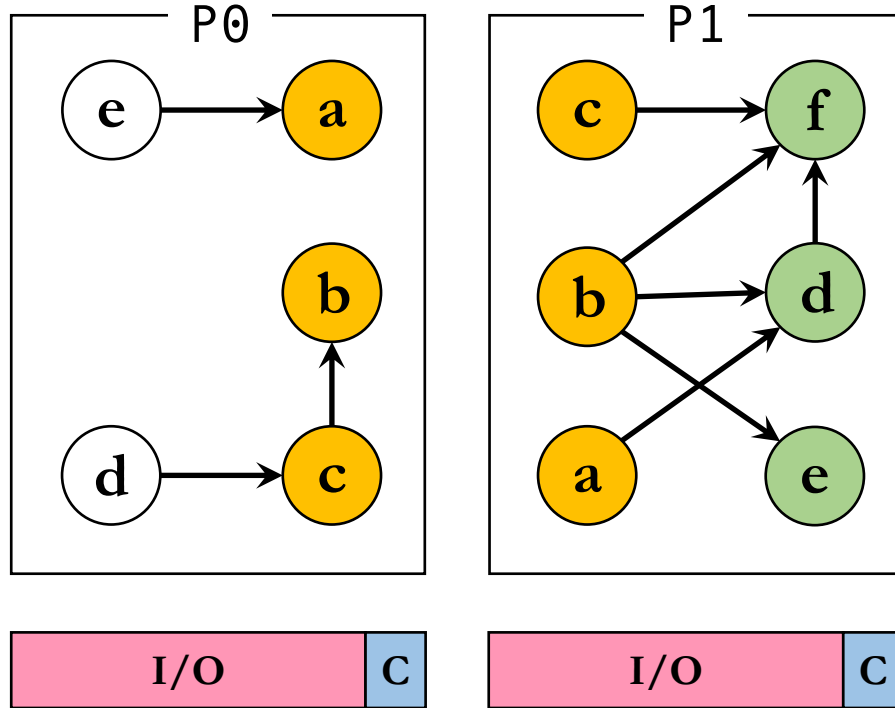


Cross-Iteration Value Propagation

Partial Aggregation



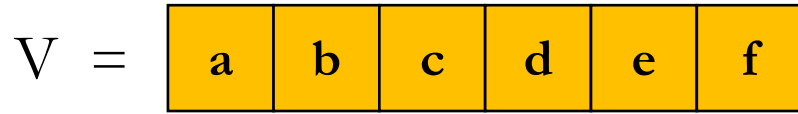
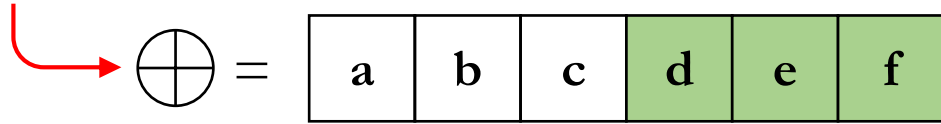
iteration t



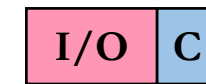
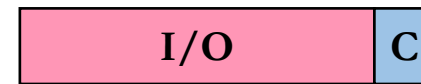
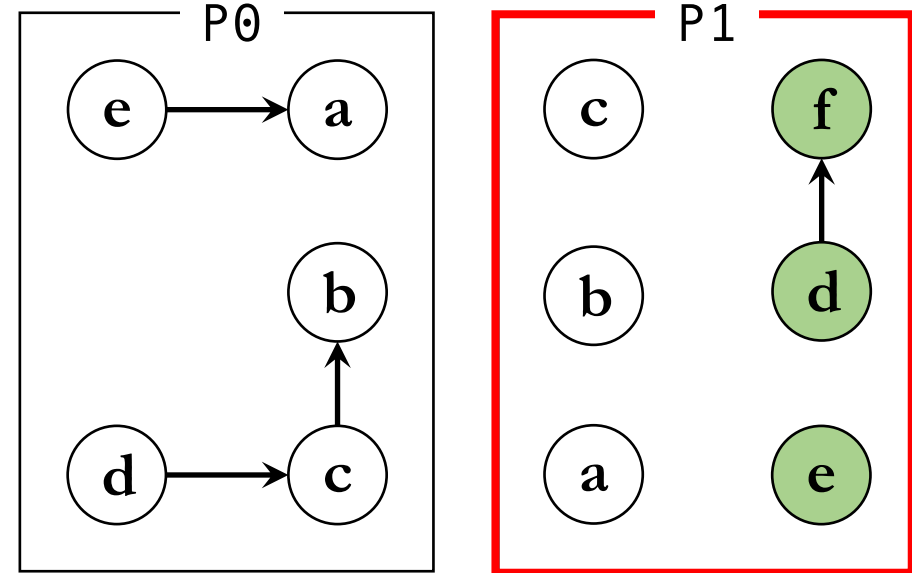
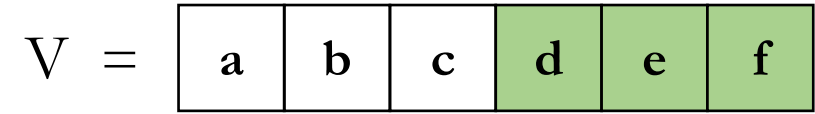
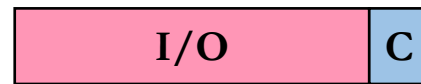
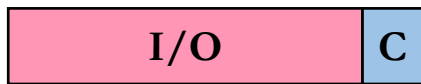
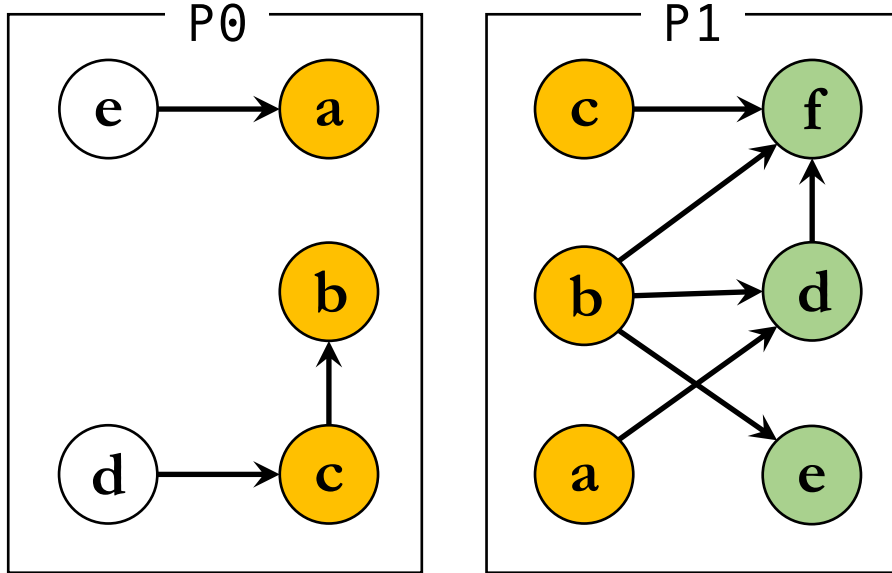
iteration t+1

Cross-Iteration Value Propagation

Partial Aggregation

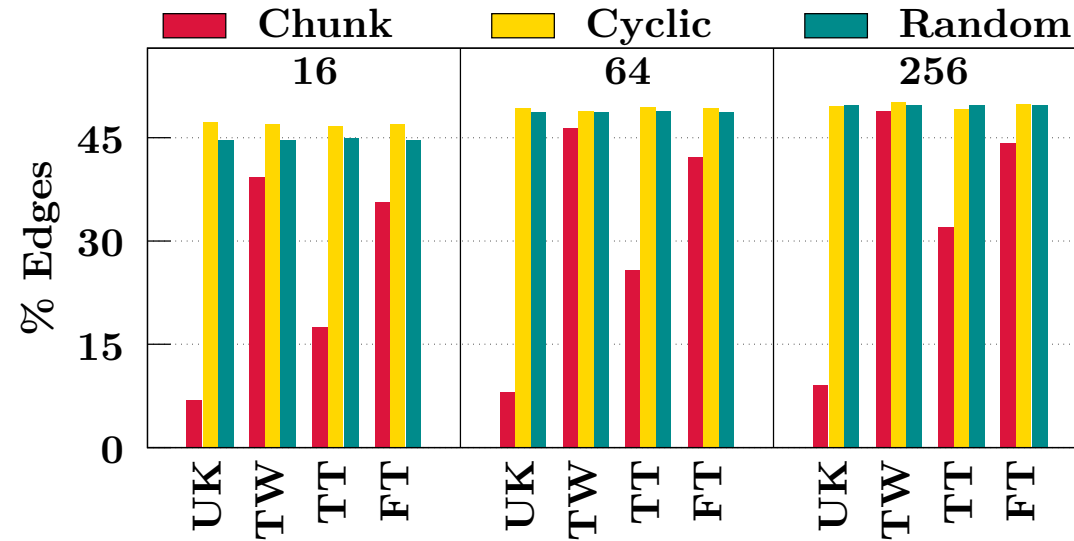


iteration t



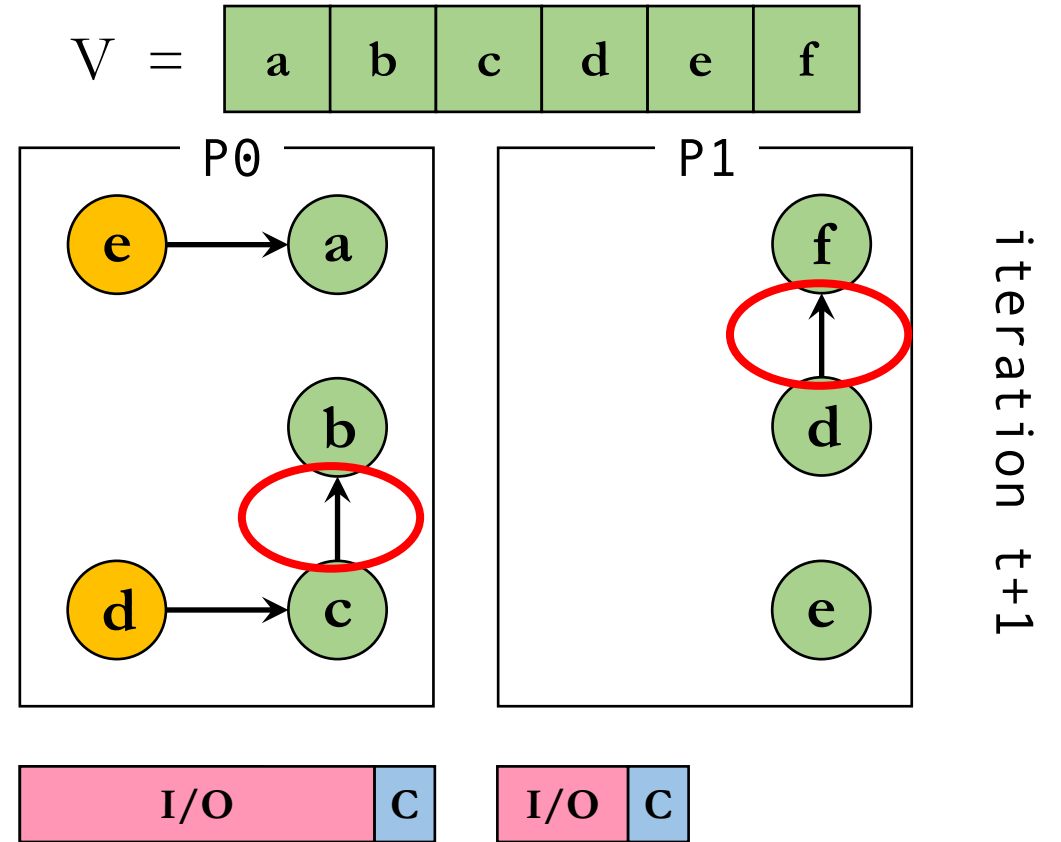
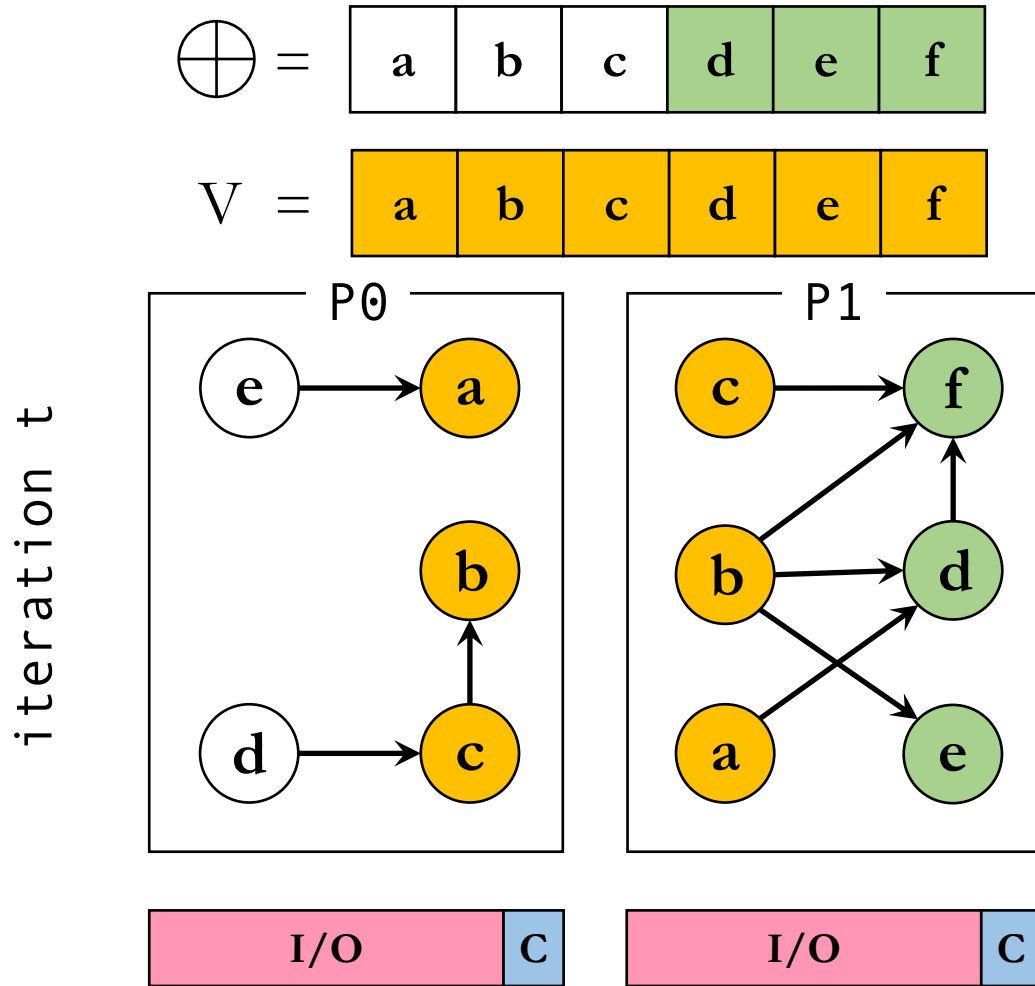
iteration t+1

Cross-Iteration Value Propagation

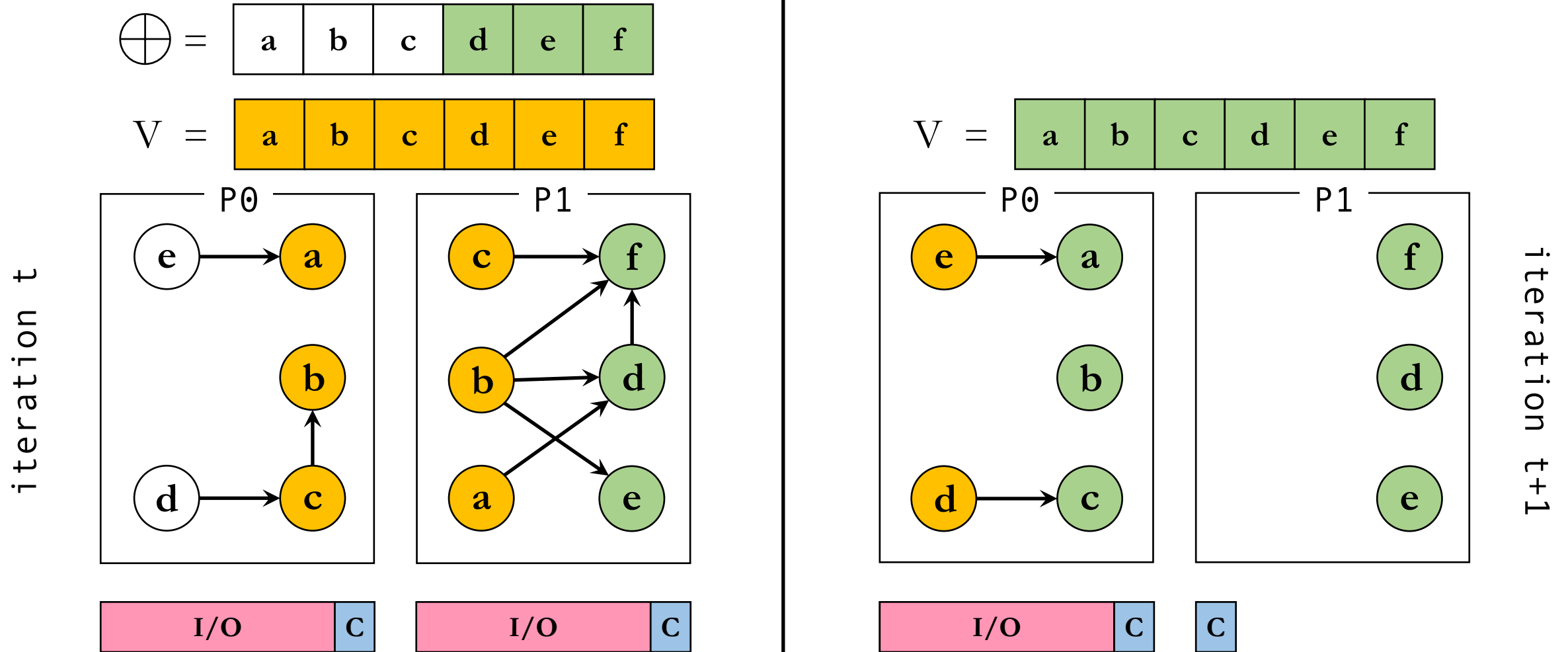


future values computed
across >45% edges

Intra-Partition Propagation



Intra-Partition Propagation

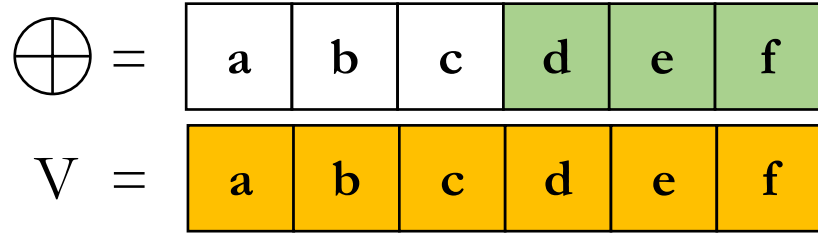


Lumos: Cross-Iteration Value Propagation

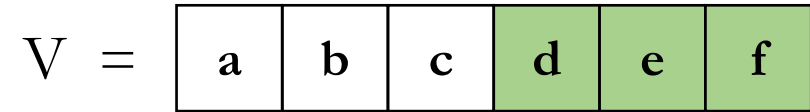
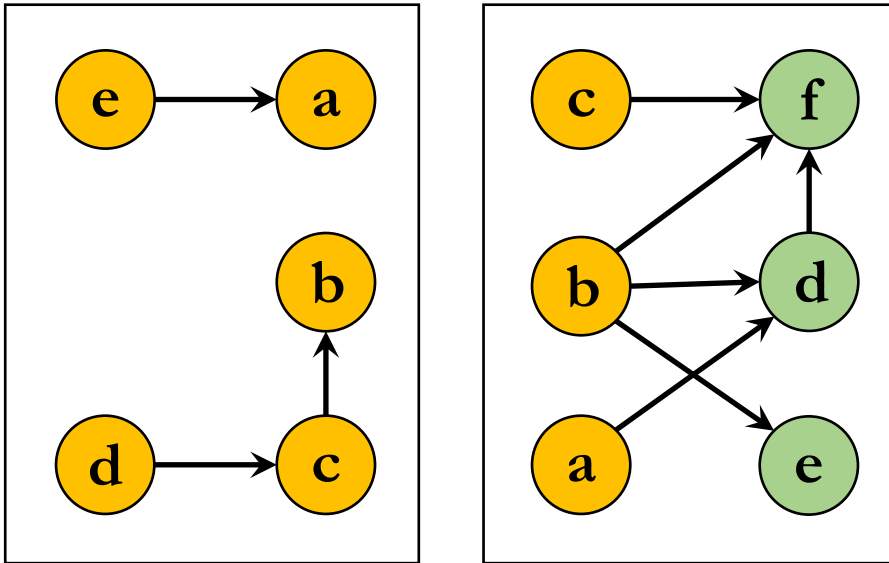
- Propagation across **multiple future iterations**
- Synergy with **Dynamic Shards** [ATC'16]
- Detailed **comparison with out-of-order** techniques
- Generalized **graph layout & partitioning strategies**
- Cross-iteration propagation across **different partition sizes**
- Interplay with **selective scheduling**
- Lumos for **Asynchronous Algorithms**

Lumos for Asynchronous Algorithms

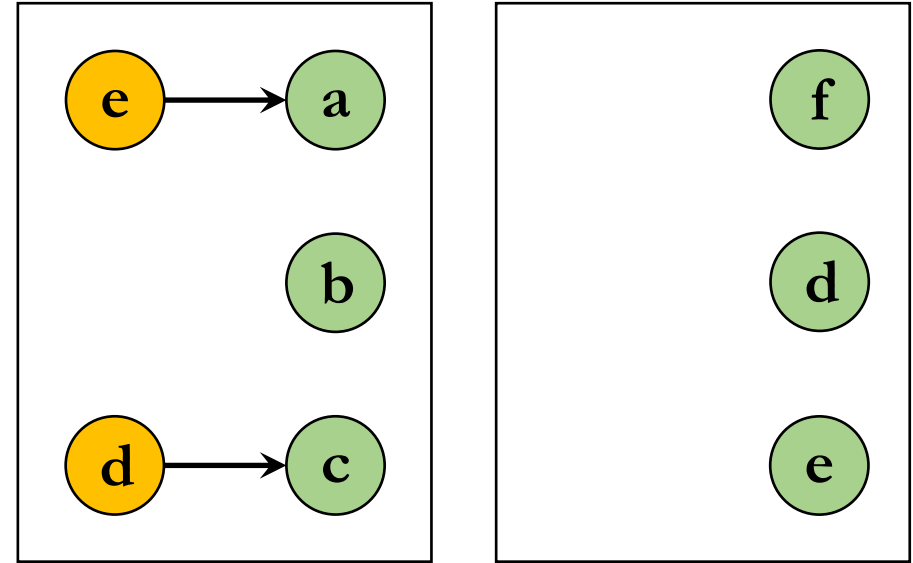
$$SSSP(v^{t+1}) = \min_{inedges(v)} SSSP(u^t) + \text{weight}(u, v)$$



iteration t

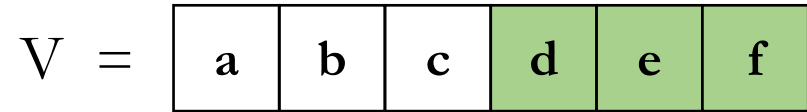
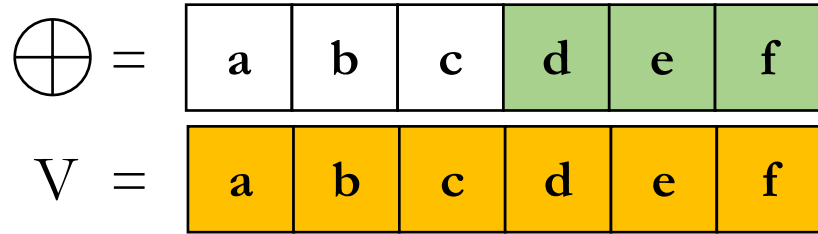


iteration t+1

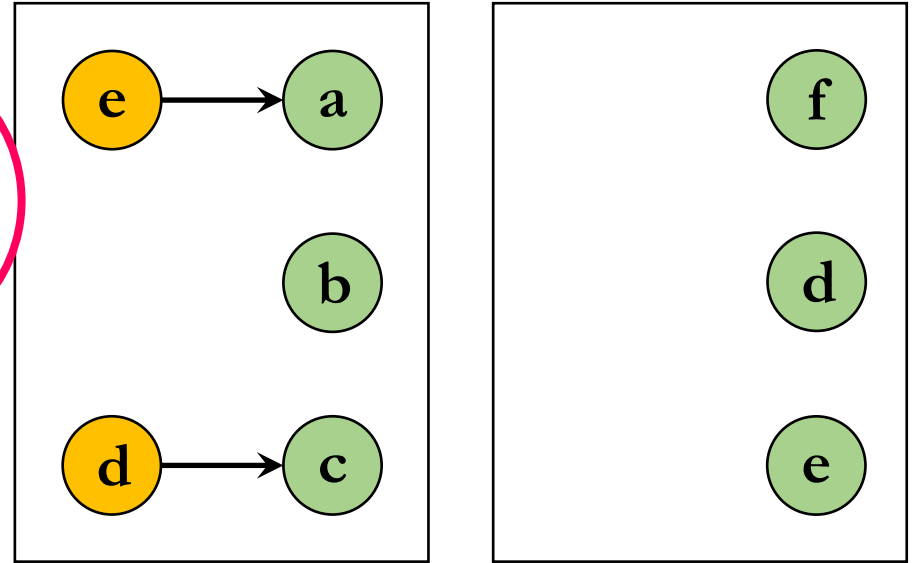
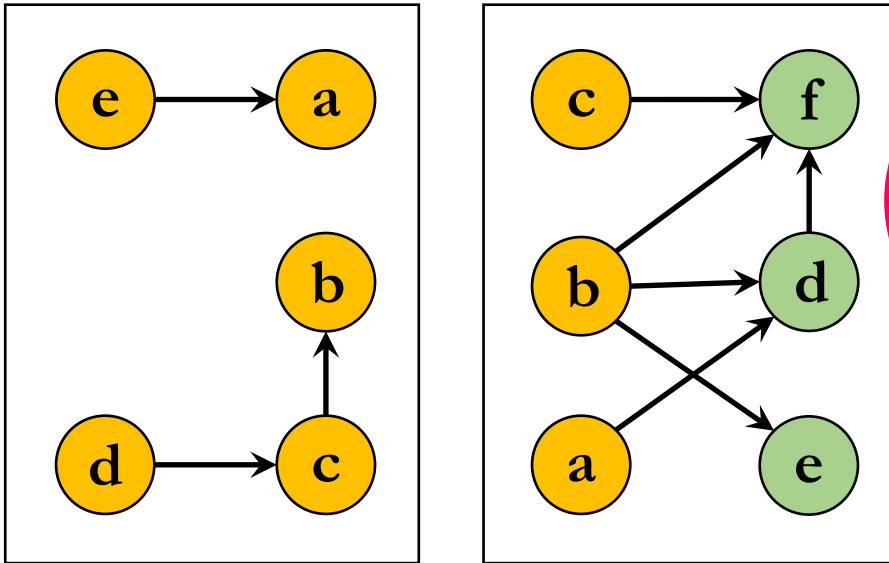


Lumos for Asynchronous Algorithms

$$SSSP(v^{t+1}) = \min_{u \in \text{inedges}(v)} SSSP(u^t) + \text{weight}(u, v)$$



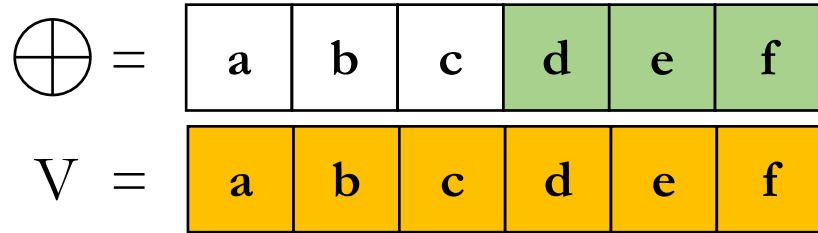
iteration t



iteration t+1

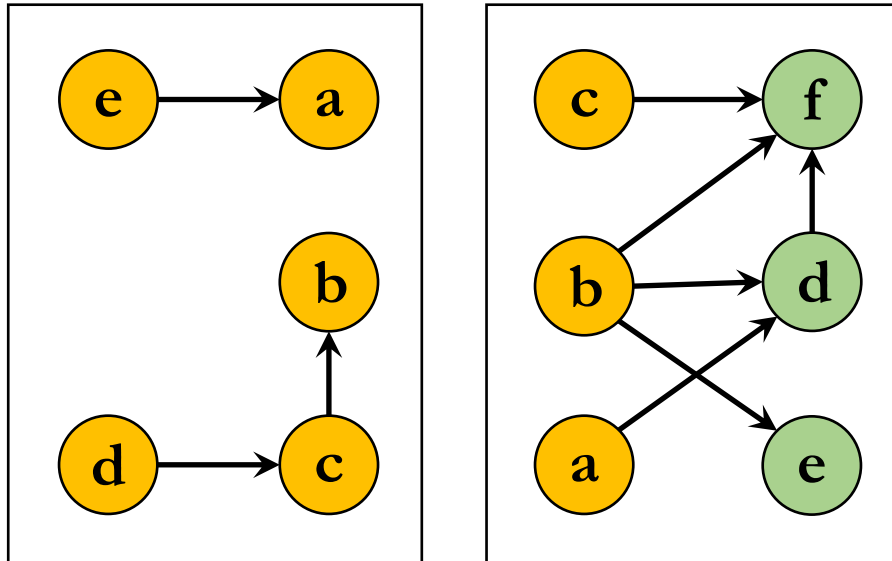
Lumos for Asynchronous Algorithms

$$\text{SSSP}(v^{t+1}) = \min_{\text{inedges}(v)} \text{SSSP}(u^t) + \text{weight}(u, v)$$



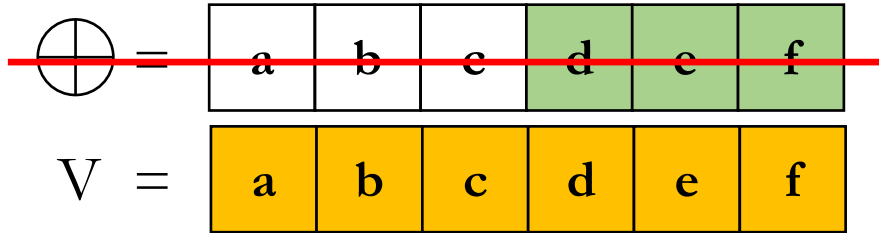
- Monotonic selection (KickStarter [ASPLOS'17])

iteration t



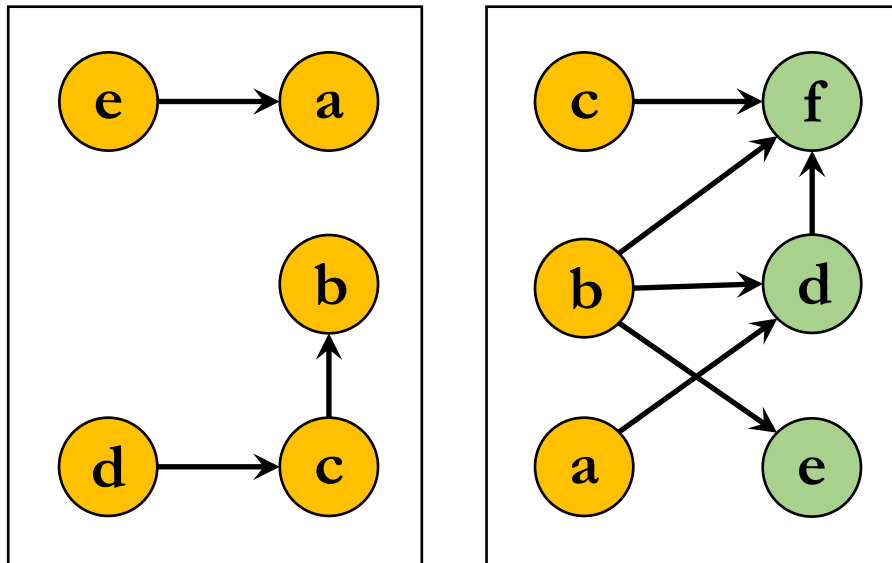
Lumos for Asynchronous Algorithms

$$\text{SSSP}(v^{\overline{t+1}}) = \min_{\text{inedges}(v)} \text{SSSP}(u^{\overline{t}}) + \text{weight}(u, v)$$



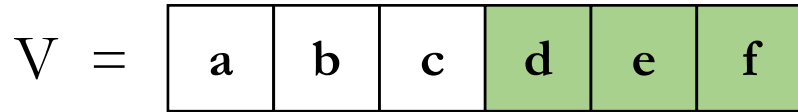
- Monotonic selection (KickStarter [ASPLOS'17])
- Partial aggregations merged directly

iteration t

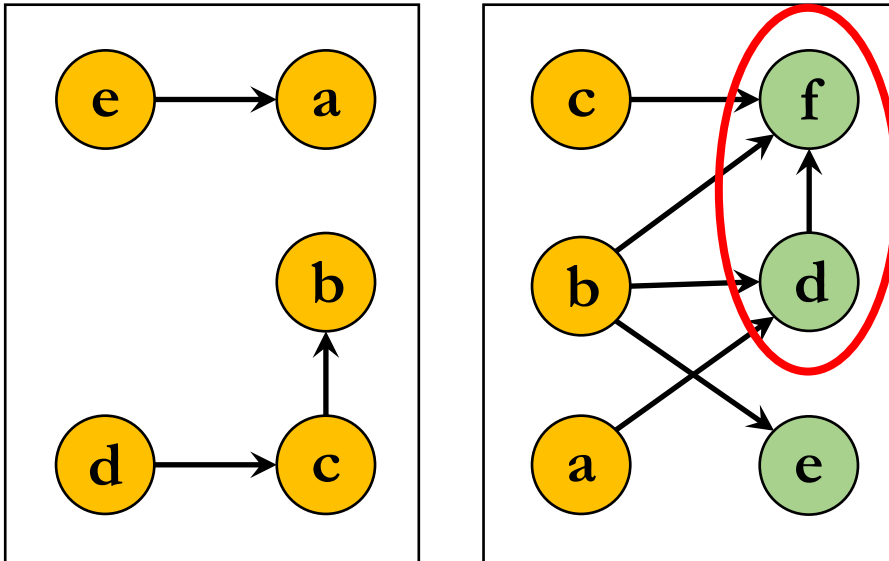


Lumos for Asynchronous Algorithms

$$\text{SSSP}(v) = \min_{\text{inedges}(v)} \text{SSSP}(u) + \text{weight}(u, v)$$



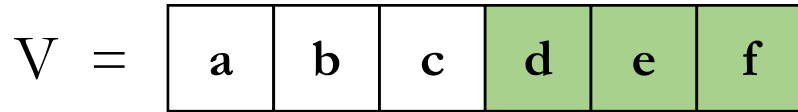
iteration t



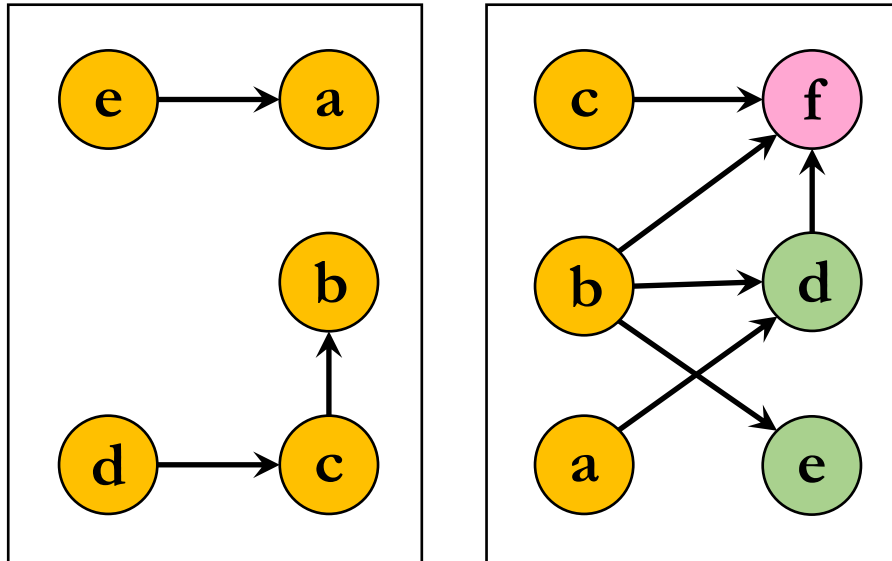
- Monotonic selection (KickStarter [ASPLOS'17])
- Partial aggregations merged directly
- Partial values safe to be propagated
 - Any change can be propagated immediately

Lumos for Asynchronous Algorithms

$$\text{SSSP}(v) = \min_{\text{inedges}(v)} \text{SSSP}(u) + \text{weight}(u, v)$$



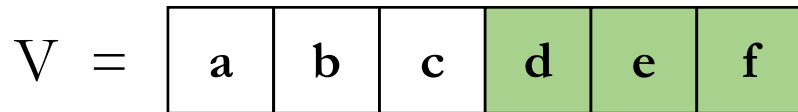
iteration t



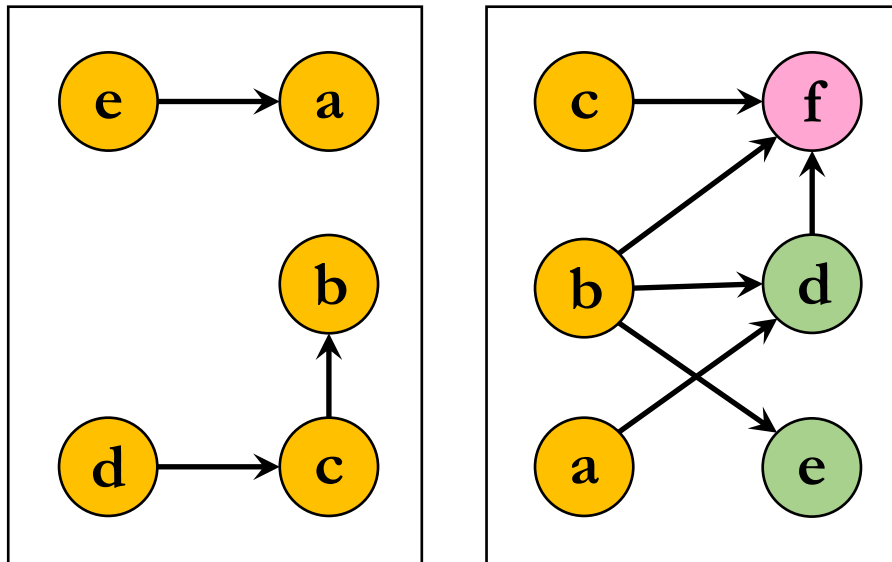
- Monotonic selection (KickStarter [ASPLOS'17])
- Partial aggregations merged directly
- Partial values safe to be propagated
 - Any change can be propagated immediately

Lumos for Asynchronous Algorithms

$$\text{SSSP}(v) = \min_{\text{inedges}(v)} \text{SSSP}(u) + \text{weight}(u, v)$$



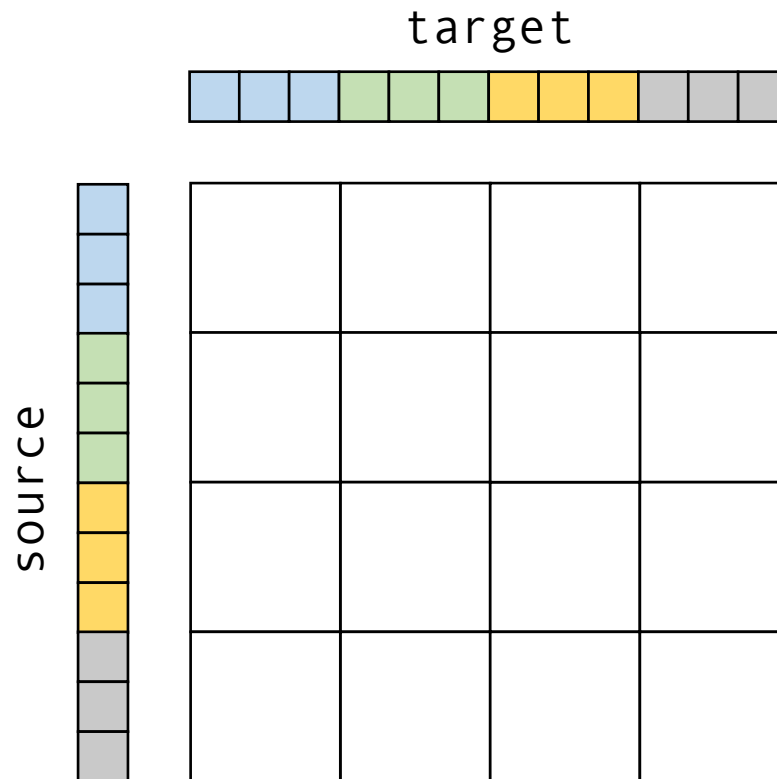
iteration t



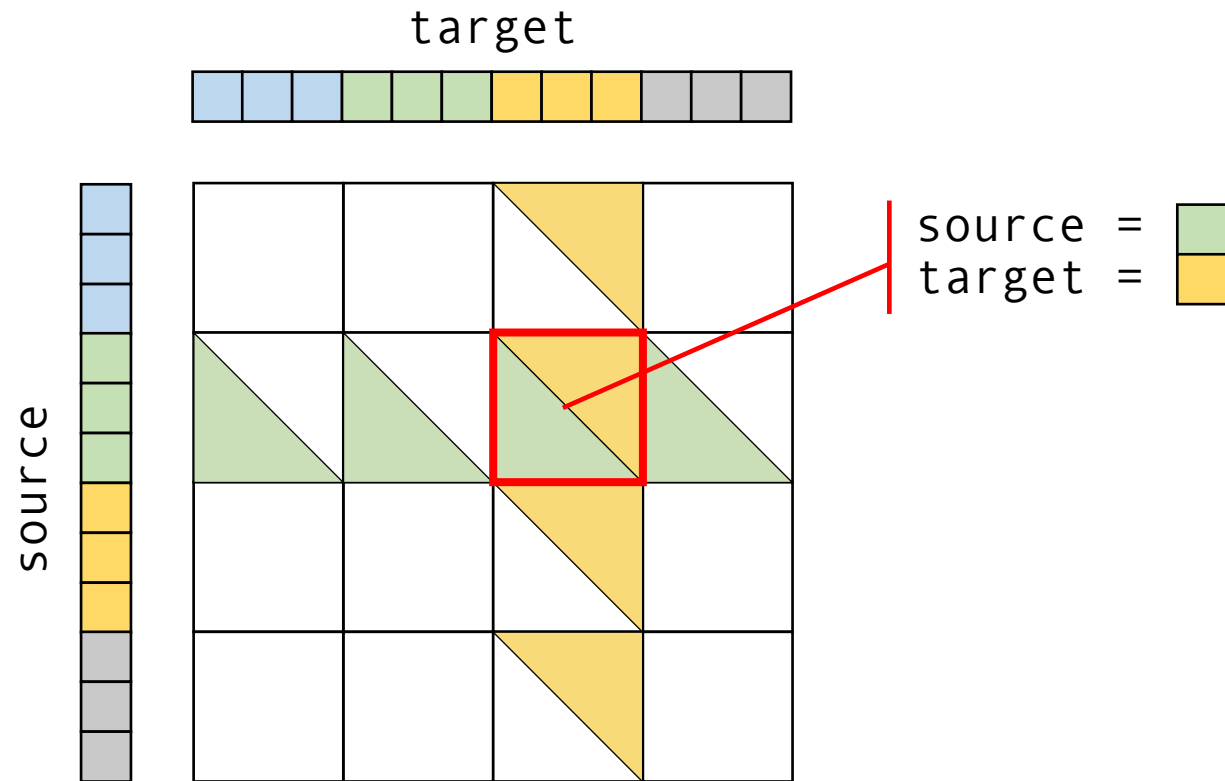
- Monotonic selection (KickStarter [ASPLOS'17])
- Partial aggregations merged directly
- Partial values safe to be propagated
 - Any change can be propagated immediately
- Intra-partition propagation based on degree of asynchrony (ASPIRE [OOPSLA'14])

Lumos System

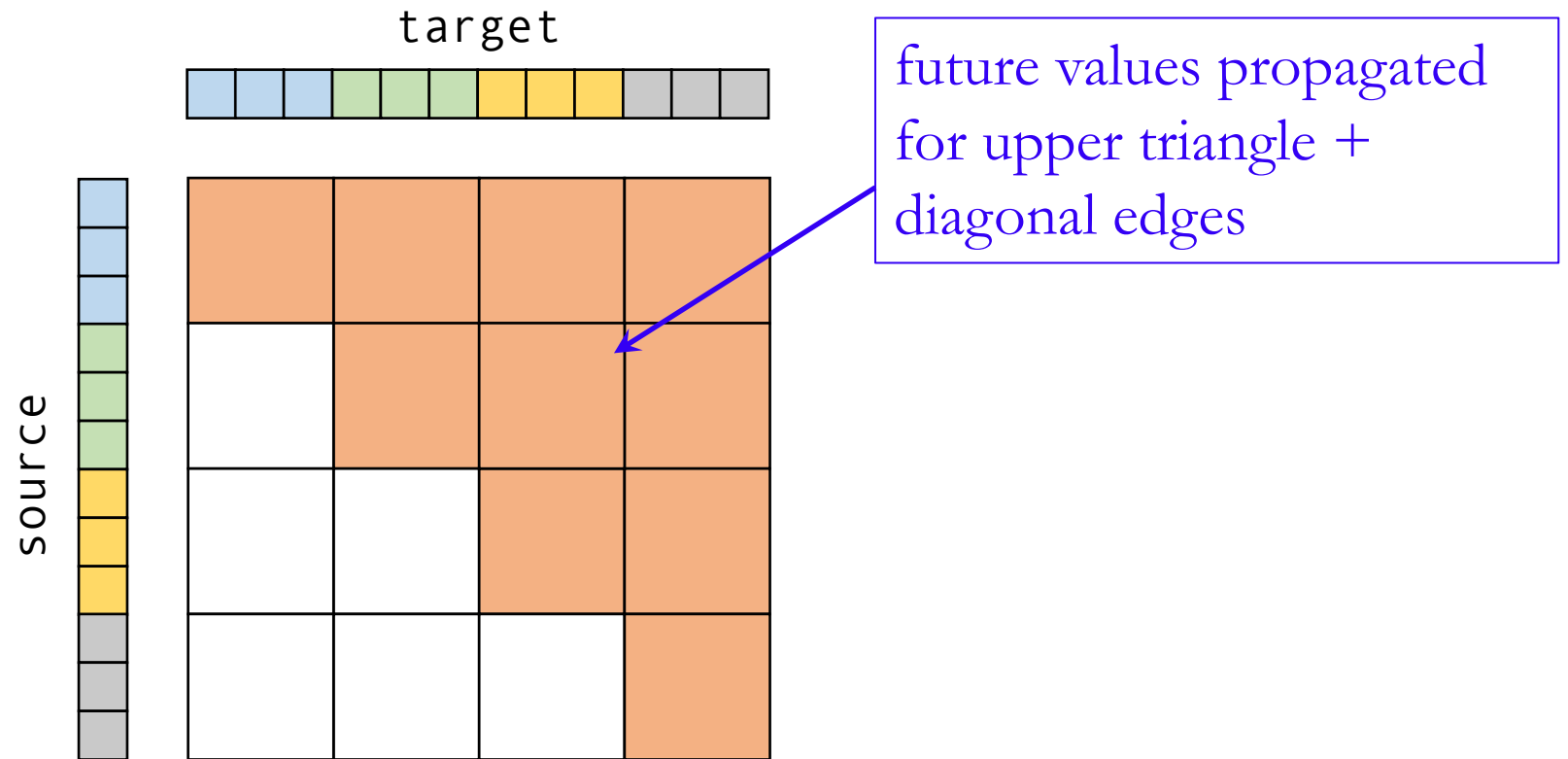
- Grid layout based on GridGraph [ATC'15]



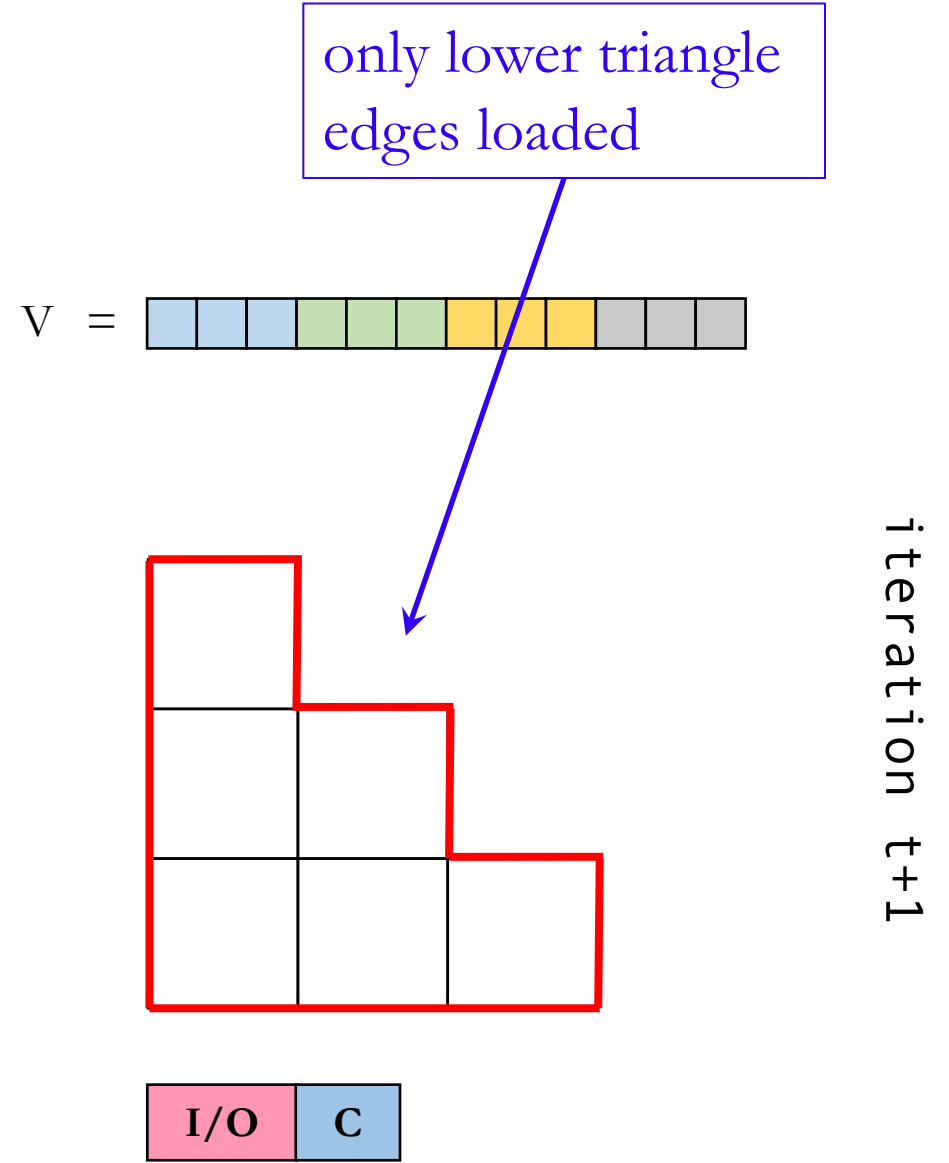
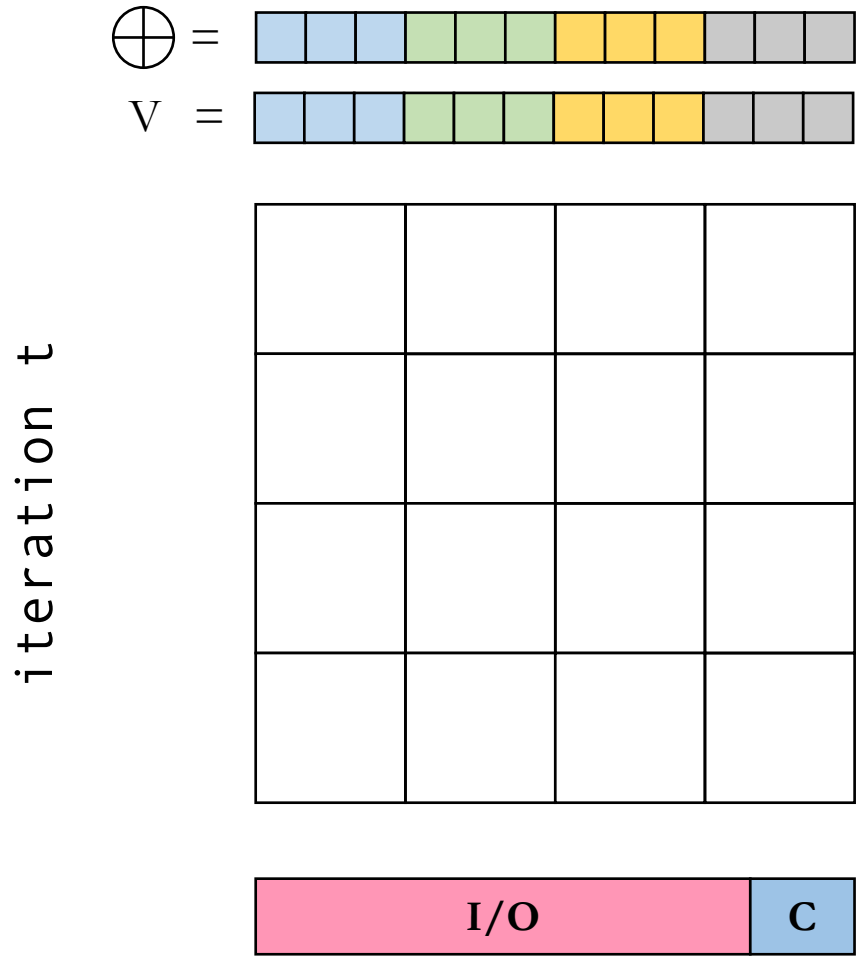
Lumos System



Lumos System

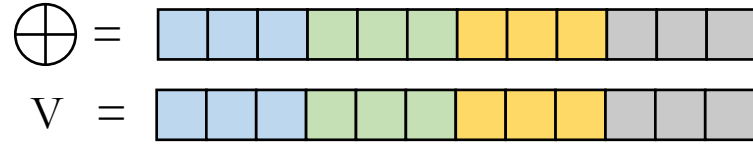


Lumos System

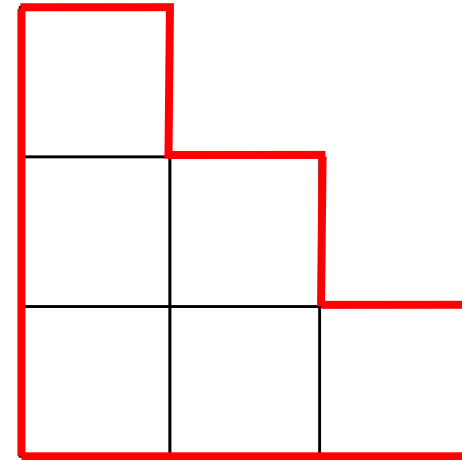
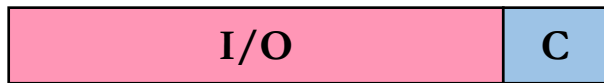
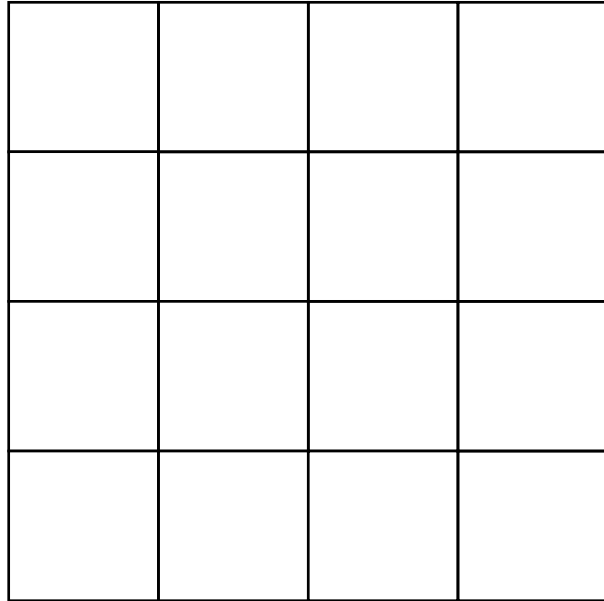


Lumos System

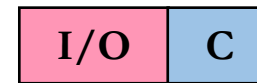
```
function PROCESSPRIMARY(PROPAGATE, CROSSPROPAGATE, COMPUTE)  
function PROCESSSECONDARY(PROPAGATE, COMPUTE)
```



iteration t



iteration t+1

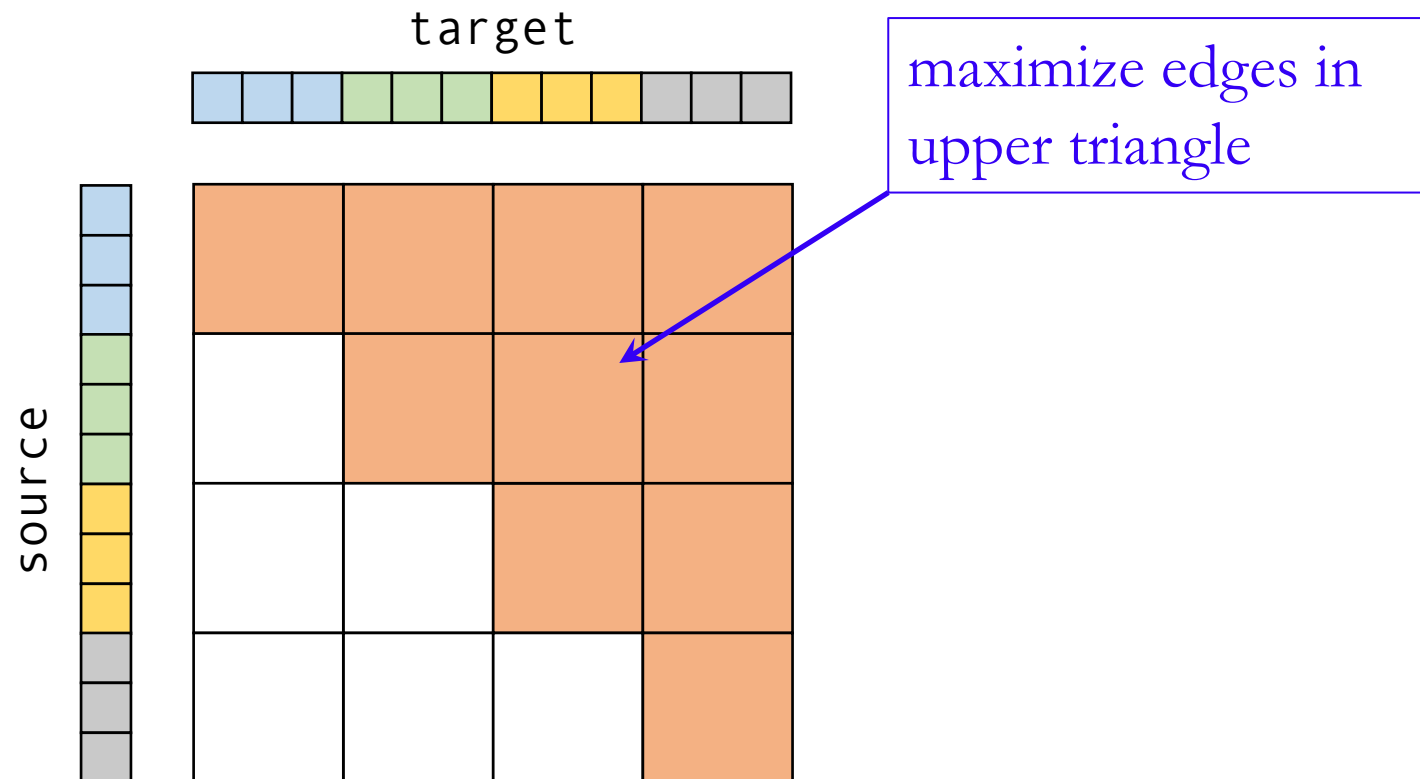


Light-weight Degree-Aware Partitioning

HOF: Highest Out-Degree First

HIL: Highest In-Degree Last

HRF: Highest Out-Degree to In-Degree Ratio First

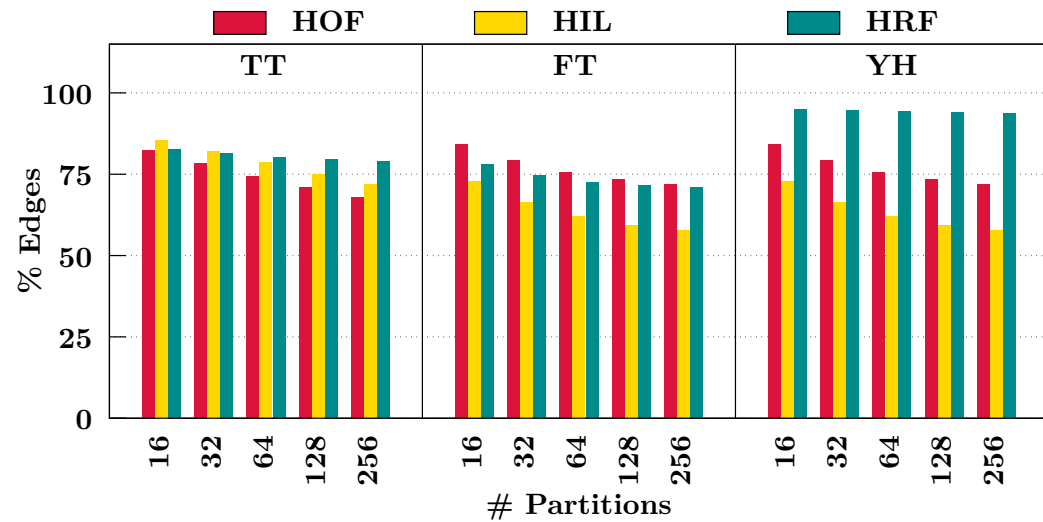


Light-weight Degree-Aware Partitioning

HOF: Highest Out-Degree First

HIL: Highest In-Degree Last

HRF: Highest Out-Degree to In-Degree Ratio First



cross-iteration
propagation across
72-93% edges

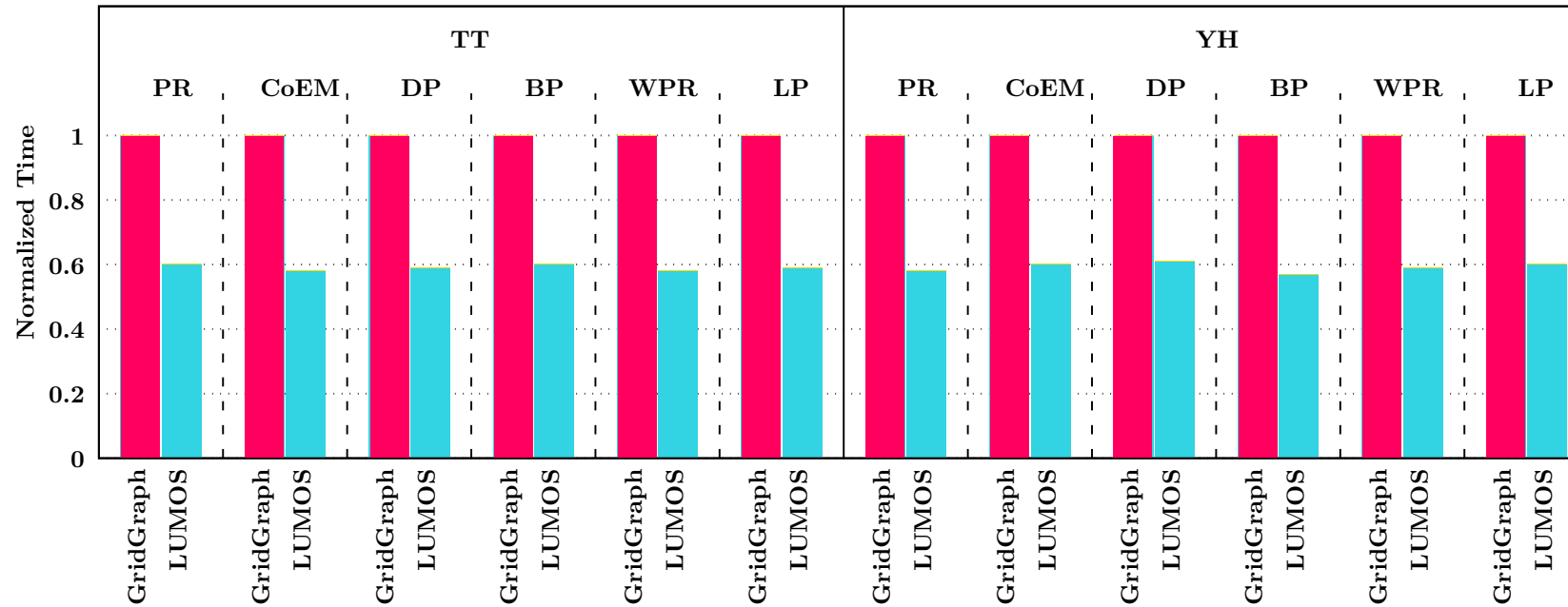
	Edges	Vertices
TT	2B	52.6M
FT	2.5B	68.3M
YH	6.6B	1.4B

Experimental Setup

- Graph algorithms
 - PageRank (weighted and unweighted), Belief Propagation, Co-Training Expectation Maximization, Dispersion, Label Propagation
- Performance on single disk
 - h1.2xlarge: 278MB/sec HDD read bandwidth
- Scaling I/O
 - d2.4xlarge: 195-768MB/sec read bandwidth over 1-4 HDDs
 - i3.8xlarge: 1.2-4.1GB/sec read bandwidth over 1-4 SSDs

Performance

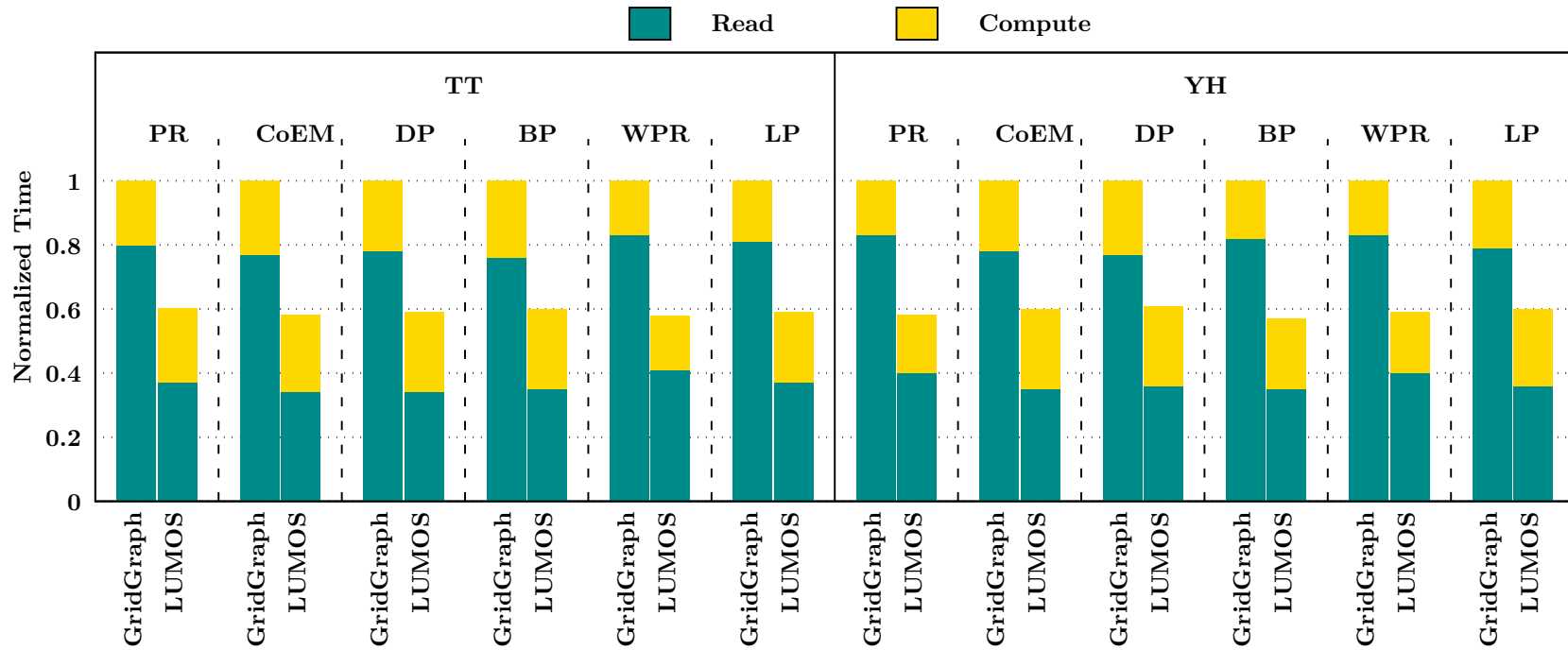
h1.2xlarge: 278MB/sec HDD read bandwidth



	Edges	Vertices
TT	2B	52.6M
YH	6.6B	1.4B

Performance

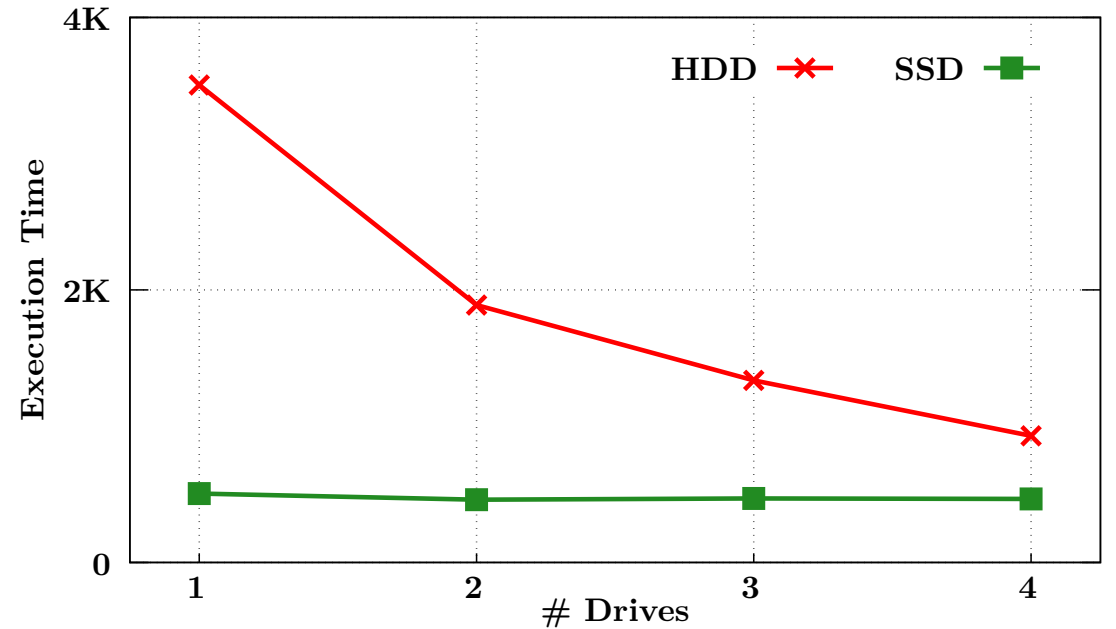
h1.2xlarge: 278MB/sec HDD read bandwidth



	Edges	Vertices
TT	2B	52.6M
YH	6.6B	1.4B

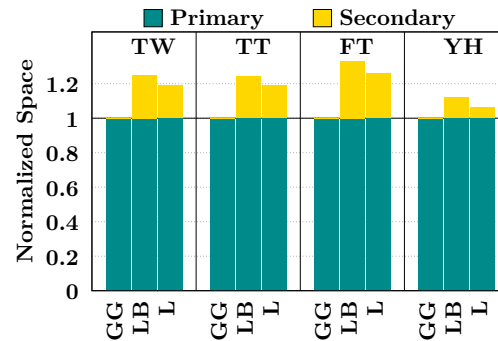
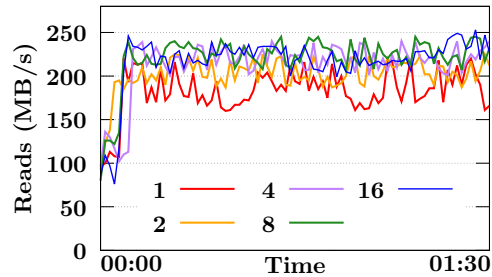
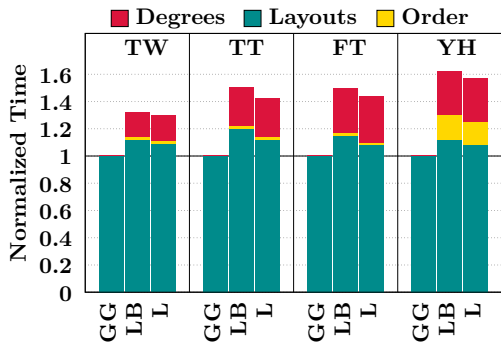
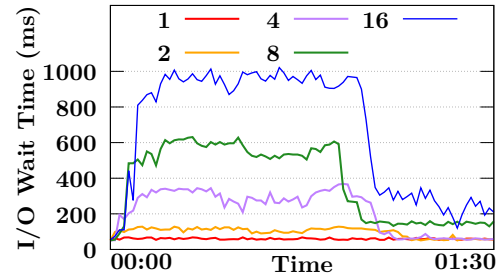
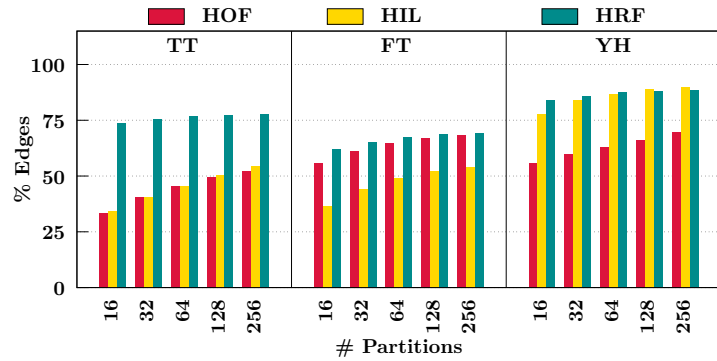
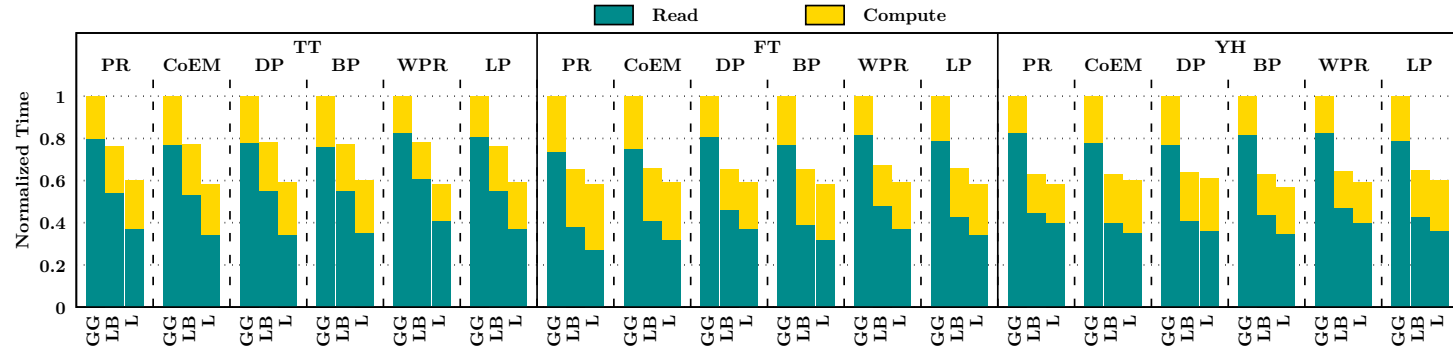
Scaling I/O

	Edges	Vertices
RMAT29	8.6B	537M



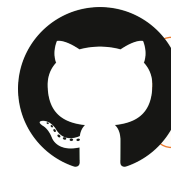
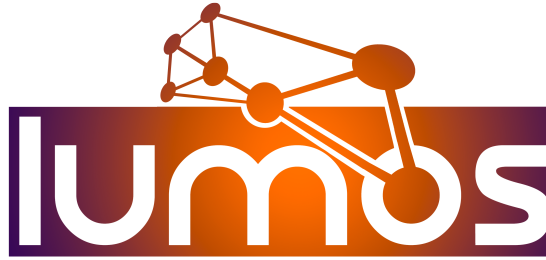
	Device Type	Single Drive	RAID-0 with k drives		
			k = 2	k = 3	k = 4
d2.4xlarge	HDD	195MB/s	368MB/s	590MB/s	768MB/s
i3.8xlarge	SSD	1.2GB/s	3.8GB/s	4.1GB/s	3.9GB/s

Detailed Evaluation of Lumos



	Version	TT	FT	YH
PR	GridGraph	737	1008	3223
	LUMOS-BASE	563	659	2027
	LUMOS	439	583	1885
	× LUMOS	1.68×	1.73×	1.71×
CoEM	GridGraph	1119	1554	5082
	LUMOS-BASE	861	1029	3216
	LUMOS	651	914	3043
	× LUMOS	1.72×	1.70×	1.67×
DP	GridGraph	846	1032	3484
	LUMOS-BASE	656	675	2219
	LUMOS	498	611	2111
	× LUMOS	1.70×	1.69×	1.65×
BP	GridGraph	2498	3782	13769
	LUMOS-BASE	1921	2456	8660
	LUMOS	1487	2212	7913
	× LUMOS	1.68×	1.71×	1.74×
WPR	GridGraph	984	1302	4330
	LUMOS-BASE	769	874	2758
	LUMOS	569	770	2547
	× LUMOS	1.73×	1.69×	1.70×
LP	GridGraph	1054	1421	4583
	LUMOS-BASE	805	935	2976
	LUMOS	624	826	2728
	× LUMOS	1.69×	1.72×	1.68×

Conclusion



github.com/pdclab/lumos

- **Dependency-Driven Cross-iteration** value propagation
 - **Out-of-Order** processing to reduce I/O
 - Guarantees **Bulk Synchronous Parallel** semantics
- Generic technique that fundamentally eliminates the performance barrier