

GridGraph: Large-Scale Graph Processing on a Single Machine Using 2-Level Hierarchical Partitioning

Xiaowei ZHU

Tsinghua University

Widely-Used Graph Processing





Existing Solutions

- Shared memory
 - Single-node & in-memory
 - Ligra, Galois, Polymer
- Distributed
 - Multi-node & in-memory
 - GraphLab, GraphX, PowerLyra
- Out-of-core
 - Single-node & disk-based
 - GraphChi, X-Stream, TurboGraph

\bigcirc

Existing Solutions

Shared memory

- Single-node & in-memory
- Ligra, Galois, Polymer

• Distributed

- Multi-node & in-memory
- GraphLab, GraphX, PowerLyra

• Out-of-core

- Single-node & disk-based
- GraphChi, X-Stream, TurboGraph

Large-scale Limited capability to big graphs

Irregular structure Imbalance of computation and communication

Inevitable random access Expensive disk random access

\bigcirc

Existing Solutions

- Shared memory
 - Single-node & in-memory
 - Ligra, Galois, Polymer
- Distributed
 - Multi-node & in-memory
 - GraphLab, GraphX, PowerLyra
- Out-of-core
 - Single-node & disk-based
 - GraphChi, X-Stream, TurboGraph Most cost effective!

Large-scale Limited capability to big graphs

Irregular structure Balance of computation and communication

Inevitable random access Expensive disk random access



Methodology

- How to handle graphs that is much larger than memory capacity?
 - Partition!

System	Unit
GraphChi	Shards
TurboGraph	Page
X-Stream	Streaming Partitions
PathGraph	Tree-Based Partitions
FlashGraph	Page
GridGraph	Chunks & Blocks



State-of-the-Art Methodology

- X-Stream
 - Access edges <u>sequentially</u> from <u>disks</u>
 - Access vertices **randomly** inside **memory**
 - Guarantee locality of vertex accesses by partitioning



State-of-the-Art Methodology

- X-Stream
 - Access edges <u>sequentially</u> from <u>slow</u> memory
 - Access vertices **randomly** inside **fast** memory
 - Guarantee locality of vertex accesses by partitioning





9

Edge-Centric Scatter-Gather



Scatter:

for each streaming partition

load <u>source</u> vertex chunk of edges into fast memory stream edges

append to several updates

Gather:

for each streaming partition

load destination vertex chunk of updates into fast

memory

stream <u>updates</u>

apply to vertices

X-Stream: Edge-centric Graph Processing using Streaming Partitions, A. Roy et al., SOSP 2013

\bigcirc

Motivation





Basic Idea



Answer: Guarantee the locality of both <u>source</u> and <u>destination</u> vertices when streaming <u>edges</u>!

Streaming-Apply: for each streaming edge block load <u>source</u> and <u>destination</u> vertex chunk of edges into memory stream <u>edges</u> read from <u>source</u> vertices write to <u>destination</u> vertices

\mathbf{C}

Solution

- Grid representation
 - Dual sliding windows
 - Selective scheduling
- 2-level hierarchical partitioning

\bigcirc

Grid Representation

- Vertices partitioned into P equalized chunks
- Edges partitioned into P × P blocks
 - Row ⇔ source
 - Column ⇔ destination



Streaming-Apply Processing Model

- Stream edges block by block
 - Each block corresponding to two vertex chunks
 - Source chunk + destination chunk
 - Fit into memory
- Difference with scatter-gather
 - 2 phases \rightarrow 1 phase
 - Updates are applied on-the-fly



\bigcirc

Dual Sliding Windows

- Access edge blocks in column-oriented order
 - From left to right
 - Destination window slides as column moves
 - From top to bottom
 - Source window slides as row moves
 - Optimize write amount
 - 1 pass over the destination vertices







P=2	PR	Deg	NewP	R ₩
	V	V	0 0	00
	1 1	2 2	(1, 2) (2, 1)	(1, 3) (2, 4)
	1	1 2	(3, 2) (4, 2)	(4, 3)

Object	I/O Amt.
Edges	0
Src. vertex	0
Dest. vertex	0

1, 2, 0

Initialize

$$PageRank_{i} = (1 - d) + d * \sum_{j \in N_{in}(i)} \frac{PageRank_{j}}{OutDegree_{j}}$$





P=2	PR	Πρα	NewP	R ₩
	₽	↓	0.5 0.5	00
	1 1	2 2	(1, 2) (2, 1)	(1, 3) (2, 4)
	1 1	1 2	(3, 2) (4, 2)	(4, 3)

Object	I/O Amt.
Edges	$0 \rightarrow 2$
Src. vertex	$0 \rightarrow 2$
Dest. vertex	$0 \rightarrow 2$

Stream Block (1, 1)

Cache source vertex chunk 1 in <u>memory</u> (<u>miss</u>); Cache destination vertex chunk 1 in <u>memory</u> (<u>miss</u>); Read edges (1, 2), (2, 1) from <u>disk</u>;





P=2	PR	Deg	NewP	R ₩
	$\mathbf{\Psi}$	↓	0.5 2	00
	1 1	2 2	(1, 2) (2, 1)	(1, 3) (2, 4)
	1 1	1 2	(3, 2) (4, 2)	(4, 3)

Object	I/O Amt.
Edges	2 → 4
Src. vertex	$2 \rightarrow 4$
Dest. vertex	2

)

Stream Block (2, 1)

Cache source vertex chunk 2 in <u>memory</u> (<u>miss</u>); Cache destination vertex chunk 1 in <u>memory</u> (<u>hit</u>); Read edges (3, 2), (4, 2) from <u>disk</u>;





P=2	PR	Deg	NewP	R ₩
	V	V	0.5 2	0.5 0.5
	1 1	2 2	(1, 2) (2, 1)	(1, 3) (2, 4)
	1 1	1 2	(3, 2) (4, 2)	(4, 3)

Object	I/O Amt.
Edges	4 → 6
Src. vertex	4 → 6
Dest. vertex	2 → 6

1, 2, 0.5

Stream Block (1, 2)

Cache source vertex chunk 1 in <u>memory</u> (<u>miss</u>); Cache destination vertex chunk 2 in <u>memory</u> (<u>miss</u>); Write back destination vertex chunk 1 to <u>disk</u>; Read edges (1, 3), (2, 4) from <u>disk</u>;





P=2	PR	Deg	NewP	R ♥
	V	↓	0.5 2	1 0.5
	1 1	2 2	(1, 2) (2, 1)	(1, 3) (2, 4)
	1 1	1 2	(3, 2) (4, 2)	(4, 3)

Object	I/O Amt.
Edges	$6 \rightarrow 7$
Src. vertex	$6 \rightarrow 8$
Dest. vertex	6

1, 2, 0.5

Stream Block (2, 2)

Cache source vertex chunk 2 in <u>memory</u> (<u>miss</u>); Cache destination vertex chunk 2 in <u>memory</u> (<u>hit</u>); Read edges (4, 3) from <u>disk</u>;





P=2	PR	Deg	NewP	R ♥
	¥	↓	0.5 2	1 0.5
	1 1	2 2	(1, 2) (2, 1)	(1, 3) (2, 4)
	1 1	1 2	(3, 2) (4, 2)	(4, 3)

Object	I/O Amt.
Edges	7
Src. vertex	8
Dest. vertex	6 → 8

1, 2, 0.5

Iteration 1 finishes

Write back destination vertex chunk 2 to disk;



I/O Access Amount

• For 1 iteration $E + (2 + P) \times V$

1 pass over the edges (read)

P pass over the source vertices (read)

1 pass over the destination vertices (read+write)

Implication: P should be the minimum value that enables needed vertex data to be fit into memory.



I/O Access Amount



1 pass over the destination vertices (read+write)

Implication: P should be the minimum value that enables needed vertex data to be fit into memory.



Memory Access Amount



1 pass over the destination vertices (read+write)

Implication: P should be the minimum value that enables needed vertex data to be fit into memory.



- Skip blocks with no active edges
 - Very simple but important optimization
 - Effective for lots of algorithms
 - BFS, WCC, ...















BFS from 1 with P = 2



	Parent	1	-	1	1		2	
BFS finishes		(1, 2)	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		4			
		(3, 2 (4, 2	2) 2)	(4	., 3)	- A		vcce

4+7+3=14 Access 14 edges in all

Impact of P on Selective Scheduling



Р	1	2	4
Edge accesses	21(=7+7+7)	14=(4+7+3)	7=(2+3+2)

Effect becomes better with more fine-grained partitioning.

Implication: A larger value of P is preferred.



Dilemma on Selection of P

Coarse-grained

Ρ

Fewer accesses on vertices

Poorer locality Less selective scheduling **Fine-grained**

Better locality More selective scheduling

More accesses on vertices

small

Dilemma from Memory Hierarchy

- Different selections of P
 - Disk Memory hierarchy
 - Fit hot vertex data into memory
 - Memory Cache hierarchy
 - Fit hot vertex data into cache
 - Disk Memory Cache





2-Level Hierarchical Partitioning

• Apply a Q × Q partitioning over the P × P grid



P = number of partitions, M = memory capacity, C = LLC capacity, V = size of vertices $_{33}$



Programming Interface

StreamVertices(F_v, F)

• StreamEdges(F_e, F)

Algorithm 3 PageRank **function** CONTRIBUTE(*e*) Accum(&NewPR[e.dest], $\frac{PR[e.source]}{Deg[e.source]}$ end function **function** COMPUTE(*v*) $NewPR[v] = 1 - d + d \times NewPR[v]$ **return** |NewPR[v] - PR[v]|end function d = 0.85 $PR = \{1, ..., 1\}$ Converged = 0while ¬*Converged* do *NewPR* = $\{0, ..., 0\}$ StreamEdges(Contribute) *Diff* = <u>StreamVertices</u>(Compute) Swap(*PR*, *NewPR*) Converged = $\frac{Diff}{V} \leq Threshold$ end while

\bigcirc

Evaluation

- Test environment
 - AWS EC2 i2.xlarge
 - 4 hyperthread cores
 - 30.5GB memory
 - 1 × 800GB SSD
 - AWS EC2 d2.xlarge
 - 4 hyperthread cores
 - 30.5GB memory
 - 3 × 2TB HDD

- Applications
 - BFS, WCC, SpMV,
 PageRank

Dataset	V	E	Data size	Р
LiveJournal	4.85M	69.0M	527 MB	4
Twitter	61.6M	1.47B	11 GB	32
UK	106M	3.74B	28 GB	64
Yahoo	1.41B	6.64B	50 GB	512







Yahoo

Runtime(S)	BFS	WCC	SpMV	PageRank
GraphChi	-	114162	2676	13076
X-Stream	-	-	1076	9957
GridGraph	16815	3602	263.1	4719

"-" indicates failing to finish in 48 hours

i2.xlarge, memory limited to 8GB



Disk Bandwidth Usage



I/O throughput of a 10-minute interval running PageRank on Yahoo graph



Effect of Dual Sliding Windows



PageRank on Yahoo



Effect of Selective Scheduling



WCC on Twitter



Impact of P on In-Memory Performance



PageRank on Twitter Edges cached in 30.5GB memory Optimal around P=32 where needed vertex data can be fit into L3 cache.



Impact of Q on Out-of-Core Performance





Comparison with Distributed Systems

- PowerGraph, GraphX *
 - 16 × m2.4xlarge, \$15.98/h
- GridGraph
 - i2.4xlarge, \$3.41/h
 - 4 SSDs
 - 1.8GB/s disk bandwidth



* GraphX: Graph Processing in a Distributed Dataflow Framework, JE Gonzalez et al., OSDI 2014



Conclusion

- GridGraph
 - Dual sliding windows
 - Reduce I/O amount, especially writes
 - Selective scheduling
 - Reduce unnecessary I/O
 - 2-level hierarchical grid partitioning
 - Applicable to 3-level (cache-memory-disk) hierarchy



Thanks!