



USENIX ATC 2018





CGraph: A Correlations-aware Approach for Efficient Concurrent Iterative Graph Processing

Yu Zhang¹, Xiaofei Liao¹, Hai Jin¹, Lin Gu¹, Ligang He², Bingsheng He³, Haikun Liu¹

- ¹ SCTS/CGCL, Huazhong University of Science and Technology, China
 - ² University of Warwick, UK
 - ³ National University of Singapore, Singapore







Background and Challenges





What is CGP Job



Many concurrent graph processing jobs are daily executed over the same graph (or its different snapshots) to provide various information for different products





What is CGP Job







What is CGP Job

















CGCL

Challenges: Data Access Problems in the CGP Jobs



The average execution time of each job is significantly prolonged as the number of jobs increases due to higher data access cost





Challenges: An Example



- The CGP jobs access the shared graph partitions in an individual manner along different graph paths
- The processing time of each partition is various for different jobs







Observations:

-Spatial correlation

-Temporal correlation







Observations:

-Spatial correlation: The intersections of the set of graph partitions to be handled by different CGP jobs in each iteration are large (more than 75% of all active partitions on average).

-Temporal correlation







Observations:

-Spatial correlation: The intersections of the set of graph partitions to be handled by different CGP jobs in each iteration are large (more than 75% of all active partitions on average).

-Temporal correlation: Some graph partitions may be accessed by multiple CGP jobs (may be more than 16 jobs) within a short time duration.







Observations:

-Spatial correlation: The intersections of the set of graph partitions to be handled by different CGP jobs in each iteration are large (more than 75% of all active partitions on average).

-Temporal correlation: Some graph partitions may be accessed by multiple CGP jobs (may be more than 16 jobs) within a short time duration.

Develop a solution for efficient use of cache/memory and the data access channel to achieve a higher throughput by fully exploiting the spatial/temporal correlations





Motivations: An Example



Spatial Correlations

 Load the shared partitions for the related jobs along a common order to provide opportunity to consolidate the accesses to the shared graph structure and store a single copy of the shared data in the cache to serve multiple CGP jobs at the same time.

Temporal Correlations





Motivations: An Example



Spatial Correlations

 Load the shared partitions for the related jobs along a common order to provide opportunity to consolidate the accesses to the shared graph structure and store a single copy of the shared data in the cache to serve multiple CGP jobs at the same time.

Temporal Correlations

 Take into account the temporal correlations, e.g., the usage frequency of the graph partitions, when loading them into the cache







Related Work





Existing Graph Processing Systems

Single graph processing



(/	((/		
GraphChi	X-Stream	GridGraph	NXgraph	CLIP	





Existing Graph Processing Systems

Single graph processing









Existing Graph Processing Systems









Our Approach:

A Correlations-aware Data-centric Execution Model





Main Goals

Minimize the redundant accessing and storing cost of the shared graph structure data (occupies more than 70% of the total memory of each job) by fully exploiting the spatial/temporal correlations between the CGP jobs





> Traditional approach:

 $D_1 = (V_1, S_1, E_1, W_1)$ $D_2 = (V_2, S_2, E_2, W_2)$ $D_J = (V_J, S_J, E_J, W_J)$

Most graph structure data G=(V, E, W) is the same for different CGP jobs





> Traditional approach:

 $D_{1} = (V_{1}, S_{1}, E_{1}, W_{1})$ $D_{2} = (V_{2}, S_{2}, E_{2}, W_{2})$ \vdots $D_{J} = (V_{J}, S_{J}, E_{J}, W_{J})$

Most graph structure data G=(V, E, W) is the same for different CGP jobs

Load-Trigger-Pushing (denoted by LTP) model:

D = (V, S, E, W) $G = (V, E, W), \text{ where } G = \bigcup_i G^i \qquad S_1, S_2, \dots, S_J$





- > Load-Trigger-Pushing (denoted by LTP) model:
 - Graph Loading: $L^i \leftarrow L(G^i, \cup_{j \in J} S^i_j)$







Load-Trigger-Pushing (denoted by LTP) model:

- Graph Loading: $L^i \leftarrow L(G^i, \cup_{j \in J} S^i_j)$
- Trigger and Parallel Execution: $S^{i_{new}} \leftarrow \cup_{j \in J} T_j(G^i, S^i_j)$







Load-Trigger-Pushing (denoted by LTP) model:

- Graph Loading: $L^i \leftarrow L(G^i, \cup_{j \in J} S^i_j)$
- Trigger and Parallel Execution: $S^{i_{new}} \leftarrow \cup_{j \in J} T_j(G^i, S^i_j)$
- State Pushing: $S_j^{new} = \bigcup_i S_j^{i_{new}}$







Illustration of Our LTP Model







Implementations: Graph Storage for Multiple CGP Jobs

PageRank Job				I	SSSP Job			
Vertex ID	Value	Vertex ID	Value		Vertex ID	Value	Vertex ID	Value
v1	0.2	v3	0.05		v1	1.2	v3	8
v2	0.1	v4	0.1	•••	v2	0	v4	×
v3	0.25	v5	0.3		v3	2.9	v5	∞
Private Table Partitions				İ	Private Table Partitions			
					Information Associated with Its Edges			
Vertex ID	Edge List	Flag	Master Loca	ation	Informatio	on Associ	ated with Its	Edges
Vertex ID v1	Edge List v3	FlagMaster	Master Loca Partition	ation	Informatio	on Associ	ated with Its	Edges
Vertex ID v1 v2	Edge List v3 v1, v3	FlagMasterMaster	Master Loca Partition Partition	a tion 1 1 1 1	Informatio	on Associ 1 1.2,	ated with Its .1 2.9	Edges
Vertex ID v1 v2 v3	Edge List v3 v1, v3 Ø	FlagMasterMasterMaster	Master Loca Partition Partition Partition	ration n 1 n 1 n 1	Informatio	on Associ 1. 1.2, ¢	ated with Its .1 2.9 Ø	Edges
Vertex ID v1 v2 v3 Vertex ID	Edge List v3 v1, v3 Ø Edge List	FlagMasterMasterMasterFlag	Master Loca Partition Partition Master Loca	ation 1 1 1 1 1 ation	Informatio	on Associ 1. 1.2, 0 0 Associ	ated with Its .1 2.9 Ø ated with Its	Edges
Vertex ID v1 v2 v3 Vertex ID v3	Edge List v3 v1, v3 Ø Edge List v5	FlagMasterMasterMasterFlagMirror	Master Loca Partition Partition Master Loca Partition	ation 1 1 1 1 xation 1	Informatio	on Associ 1. 1.2, 0 Associ 1.	ated with Its .1 2.9 Ø ated with Its .5	Edges Edges
Vertex ID v1 v2 v3 Vertex ID v3 v4	Edge List v3 v1, v3 Ø Edge List v3, v5	FlagMasterMasterMasterFlagMirrorMaster	Master Local Partition Partition Master Local Partition Partition	ation 1 1 1 1 1 ation ation 1	Informatio	on Associ 1.2, 0 Associ 0 Associ 1. 0.9,	ated with Its .1 2.9 3 ated with Its .5 2.5	Edges







Implementations: Details to Store Evolving Graph Structure







Implementations: Load of Partitions



(b) J2 has been submitted



(c) J3 has been submitted







Implementations: Load of Partitions



A core-subgraph based scheduling algorithm can be used to maximize the utilization ratio of each partition loaded into the cache





Implementations: Parallel Processing of Graph Partition







Implementations: Parallel Processing of Graph Partition







Implementations: Vertex State Synchronization





Non-optimized:

Synchronization from **Mirrors** to **Master**

P1:v3->P2:v3 P1:v6->P3:v6 P1:v4->P2:v4







P1:v6->P3:v6

I

Optimized:

Synchronization from **Master** to **Mirrors** Non-optimized:

P2:v3->P1:v3 P2:v3->P3:v3 P2:v4->P1:v4 P2:v4->P3:v4



P2:v3->P1:v3 P2:v4->P1:v4 P2:v3->P3:v3 P2:v4->P3:v4







Performance Evaluation





Experimental setup

Machine information -CPU: 4-way 8-core Intel Xeon CPU E5-2670; each CPU has 20 MB LLC -Main Memory: 64 GB

Typical graph algorithms -PageRank, SSSP, SCC, BFS

Data sets

Properties of data sets

Data sets	Vertices	Edges	Sizes
Twitter	41.7 M	1.4 B	17.5 GB
Friendster	65 M	1.8 B	22.7 GB
uk2007	105.9 M	3.7 B	46.2 GB
uk-union	133.6 M	5.5 B	68.3 GB
hyperlink14	1.7 B	64.4 B	480.0 GB

























Conclusions





Conclusions

What **CGraph** brings in graph processing

- Analysis of temporal/spatial correlations in concurrent graph processing
- A novel data-centric LTP model for concurrent graph processing
- A core-subgraph based scheduling scheme

Future work

- How to further optimize the approach for evolving graph analysis
- How to ensure QoS for some real-time CGP jobs
- How to extend it to a distributed platform and also heterogeneous platform consisting of GPU, FPGA and even ASIC for higher throughput.





Thanks!

Service Computing Technology and System Lab., MoE (SCTS)

Cluster and Grid Computing Lab., Hubei Province (CGCL)

