# **Big Graph Data Processing**

Wook-Shin Han POSTECH Apr. 11, 2014



- Big (Graph) Data
- Big Graph Processing Techniques
- Conclusion

#### Welcome to Graph World



**Protein interactions** 

3

## **Big Graphs in Real World**

	8.	the last	i.		
2	a de			1.	
• • •	and the				3.
• %					-
000	12		-		
o					

Facebook	
Twitter	
LinkedIn	
Last.FM	
LiveJournal	
del.icio.us	

**1Billion** 640 Million 60 Million 40 Million 25 Million 5.3 Million

# of nodes (vertices)

# of edges
140 Billion
10 Billion
0.9 Billion
2 Billion
2 Billion
0.7 Billion

#### Outline

- My Background
- Big (Graph) Data
- Big Graph Data Processing
- Conclusion

### Big Graph Processing [VLDB10, SIGMOD11, VLDB13, SIGMOD13, KDD13, SIGMOD14]

#### iGraph vI.0

•<u>Han, W.</u>, Lee, J., Duc, P., and Yu, J., "iGraph: A Framework for Comparisons of Diskbased Graph Indexing Techniques," In *VLDB* 2010. (invited to the VLDB Journal as **best** of VLDB 2010 papers)

•<u>Han, W.</u>, Duc, P., Lee, J., Kasperovics, R., and Yu, J., "iGraph in Action: Performance Analysis of Disk-Based Graph Indexing Techniques," In *SIGMOD* 2011.

#### iGraph v2.0 + Turbo<sub>ISO</sub>

•Lee, J., <u>Han, W.</u>, Kasperovics, R., and Lee, J., "An In-depth Comparison of Subgraph Isomorphism Algorithms in Graph Databases," In *VLDB* 2013.

•<u>Han, W.</u>, Lee, J., and Lee, J., "Turbo<sub>ISO</sub>: Towards Ultra-Fast Subgraph Isomorphism Search in Graph Databases," In SIGMOD 2013.

#### **TurboGraph**

<u>Han, W.</u> et al., "TurboGraph: A Fast Parallel Graph Engine Handling Billion-scale Graphs in a Single PC," In KDD 2013. (oral presentation)
Kim, J., <u>Han, W.</u>, et al., "OPT: A New Framework for Overlapped and Parallel Triangulation in Large-scale Graphs," In SIGMOD 2014

# Subgraph Isomorphism (Subgraph Matching)

- One of the most important graph queries
- Find all subgraphs that match a query graph



#### **Many Important Applications**



#### Social network search



RDF query processing





#### Motif search

Put	Chen Ibchei	n Sket m.ncb	tcher '	V2.4 · .nih.g	Chro ov/ed	me it2/in	dex.h	tml?si	miles:	-C1%	BDCC(%3DCC%3DC1)C(CC)C.C.C.C.C
Bro	adbar	nd 🔻	]	SM	LES		•	C1(	=CC(	=CC=	C1)CC(C)CC(=O)O)[N+](=O)[O-]
N	BM	Udo	Cin	Sty	Del	Qry	÷	0	ŧ	*	
-	=	=	-	11		×	*	SIA	D/A	S/D	0, *****************
Δ		Ô	$\bigcirc$	$\bigcirc$	O	Ô		⊕	Θ	0	
	~	Y	~~	$\sim$	L	+	сно	со <sub>2</sub> н	NO2	SO3H	
н		?	?	¥						He	
Li	Be				в	С	Ν	0	F	Ne	
Na	Mg				AI	Si	Ρ	s	CI	Ar	
К	Ca	Sc	Sc	•	Ga	Ge	As	Se	Br	Kr	
Rb	Sr	Y	Y	•	In	Sn	Sb	Te	1	Xe	0 0
Cs	Ba	Lu	Lu	•	Π	Pb	Bi	Po	At	Rn	
Exp	ort	MD	L Mol	file	-				Do	me	
Hydr	ogen	Kee	p Asl	s		•			Help		
Imp	oort	[II]9	일 선택	목 선	[백된	파일	없음				

# Chemical compound search

# Problems in Existing Indexing Methods: Motivation for iGraph

- Serious problems in existing experiments
  - Compared indexes are implemented in different code bases
- Elapsed times reported can vary depending on implementation skills
- Number of disk I/Os not used.
- Small database (≤ 20 Mbytes)
  - All files are cached in the OS file system cache

# iGraph vI.0 [VLDBI0, SIGMODII]

- First common framework for disk-based graph indexes
- Supports both mining-based and non-mining based indexes
- Selected as a best track paper in VLDB 2010
- Open source: http://www.igraph.or.kr
  - Has been used in 26 countries

#### Universities using iGraph vI.0



# iGraph v2.0 [VLDBI3]

- Focuses on the subgraph isomorphism algorithm
- The first generic framework that allows implementation of any subgraph isomorphism algorithm by extending this framework
- Provides in-depth analysis and comparison of the state-of-the art algorithms

# Turbo<sub>ISO</sub> [SIGMODI3]

• A new subgraph isomorphism algorithm

#### **Three things to remember**

- Candidate region exploration
- Neighborhood equivalence class (NEC)
- Comb/Perm strategy

# Review of Existing Subgraph Isomorphism Algorithm

- Backtracking algorithm
  - Find solutions by incrementing partial solutions or discarding them when they cannot be completed



matching order:  $\langle u_1, u_2, u_3 \rangle$ 





#### **Related Work**

- Exact search
  - Non signature-based
    - Ullmann [JACM1976]

An In-depth Comparison of Subgraph Isomorphism Algorithms in Graph Databases [VLDB2013]

- Signature-based
  - GraphQL [SIGMOD2008]
  - GADDI [EDBT2009]
  - Spath [VLDB2010]
- Similarity search
  - TALE [ICDE2008] NESS [SIGMOD2011], NEMA [VLDB2013], ...

#### Importance of matching order

matching order  $O_1$ :  $\langle \mathbf{u}_1, \mathbf{u}_3, \mathbf{u}_2, \mathbf{u}_4 \rangle$ 



#### Importance of matching order

matching order  $O_1$ :  $\langle \mathbf{u}_1, \mathbf{u}_3, \mathbf{u}_2, \mathbf{u}_4 \rangle$ 



**Bad** matching order  $O_1$ : needs **500,001** matches

#### Importance of matching order

matching order  $O_2$ :  $\langle \mathbf{u}_1, \mathbf{u}_4, \mathbf{u}_2, \mathbf{u}_3 \rangle$ 



# Motivation I: One good matching order is not enough!

 $O_1 (= \langle u_1, u_3, u_2, u_4 \rangle)$ :500,001 matches for  $g_1 + 51$  matches for  $g_2$  $O_2 (= \langle u_1, u_4, u_2, u_3 \rangle)$ :51 matches for  $g_1 + 500,001$  matches for  $g_2$ 











#### **Motivation 2: Useless Permutations**



matching order:  $\langle u_1, u_2, ..., u_7 \rangle$ 



call tree T 25

# Neighborhood Equivalence Class (NEC)

- Each query vertex in the same NEC has
  - the same label
  - the same *adjacent* query vertices



#### Key Idea 2: Comb/Perm strategy with NEC

#### Avoid useless permutations



## **Contributions of Turbo**ISO

- Up to four orders of magnitude performance improvement over the state of the art method\*
- Candidate region exploration
  - Provides good and robust matching order
  - Completely solves the notorious matching order problem
- Neighborhood equivalence class and Comb/Perm strategy
  - Avoid useless permutations

#### **Comparison with STW**

- Comparison with STW\*
- Graph: WordNet

\*Su

De



ÞB\_2012.

#### **Best Characteristics of Turbo**ISO

- Ultrafast and Robust performance
  - Online analytics is possible!
- Parameter-free
- Little index maintenance cost
  - Supports very large evolving graphs!
- Easily parallelizable

TurboGraph: A Fast Parallel Graph Engine Handling Billion-scale Graphs in a Single PC [KDD13]

#### **Motivation**

Distributed System approach





Single machine approach

GraphChi [OSDI'12]



DBMS approach VERY SLOW for mining???

Can we exploit nice concepts in DBMSs without losing performance?

#### **Comparison with other engines**

Query	Input graph	Vertices/Edges	Representative perfomance result	GraphChi [3] on a single PC	TurboGraph on a single PC
Pagerank (5 iterations)	Twitter	41M/1.46B	Spark[1](50 machines), 8.1 min	19.35 min	2.74 min
Pagerank (1 iteration)	Twitter	41M/1.46B	Gbase[2](100 machines), 13.5 min	6.42 min	0.62 min
Pagerank (1 iteration)	YahooWeb	1.41B/6.63B	Gbase[2](100 machines), 13.0 min	20.90 min	3.71 min

# Why is TurboGraph Ultra-fast?

- Fast Graph Storage Engine
- Full parallelism
  - Multi-core parallelism
  - Flash SSD IO parallelism
    - Reading 400~500 Mbytes/sec from commodity SSDs
    - 97K IOPS (High-performance Random Read)
- Full overlap
  - CPU processing and I/O processing
  - I/O latency can be hidden!



#### **Three things to remember**

• Efficient disk/in-memory graph storage

• Pin-and-slide model

Handling general vectors (see the paper)

## **Challenges for Graph Storage**

- Adjacency list vs. adjacency matrix
- Two types of graph operations in disk-based graphs
  - Graph traversal (unique in graphs)
  - Bitmap operations during computation

# Disk-based representation in TurboGraph

- Slotted page of I Mbyte size
  - Page contains records corresponding to adjacency lists
  - RID consists of a page ID and a slot number
- Vertex IDs or RIDs in adjacency list
  - Vertex ID approach
    - Good for bitmap operation
    - Bad for graph traversal
      - requires a potentially LARGE mapping table!
  - RID approach
    - Good for graph traversal
    - Seems to be bad for bitmap operation??

# **RID (mapping) table**

- Each entry corresponds to a page (not a single RID)
  - Size is very small
- Each entry stores the starting vertex ID in the page
- Translation of RID (pageID, slotNo) to vertex ID
  - RIDTable[pageID].startVertex + slotNo
  - Can be done in O(I)

#### **In-memory Data structures**

- Buffer pool
- Mapping table from page ID to frame ID
  - Hash-table based mapping incurs significant performance overhead for graph traversal!
  - TurboGraph uses a page table approach!
- A data structure for handling large adjacency list (see paper)

#### Example





### **Core operations in buffer pool**

- PINPAGE(pid)/UNPINPAGE(pid)
  - support large adjacency lists
- PINCOMPUTEUNPIN(pid, RIDList, uo)
  - Prepins an available frame
  - Issues an asynchronous I/O request to the FlashSSD
  - On completion of the I/O, a callback thread processes the vertices in the RIDList by invoking the user-defined function uo.Compute
  - After processing all vertices in RIDList, unpin the page

## Supported Query Power: Matrix-vector multiplication

- G = (V,E), X (column vector)
- M(G)<sub>i</sub>: *i*-th column vector of G
- Column view:  $Y = \sum_{i=1}^{|V|} M(G)_i \times X_i$ 
  - Applications can define their own multiplication and summation semantics (the user-defined function Compute can generalize both)
  - $M(G)_i$  is represented as the adjacency list of  $v_i$
- We can restrict the computation to just a subset of vertices  $v_{I[1]} \sim v_{I[k]}$

# Column-view of matrix-vector multiplication in TurboGraph

Algorithm 1 Matrix-Vector-Multiplication(G = (V, E), X, I, Y)

- 1: for i = 1 to |I| do
- 2:  $Compute(v_{I[i]}.adj, X_{I[i]}, Y)$
- 3: end for

#### **Pin-and-Slide Model**

- New computing model for efficiently processing the generalized matrix-vector multiplication in the column view
- Utilizing execution thread pool and callback thread pool

## Pin-and-Slide Model (cont'd)

- Given a set V of vertices of interest,
  - Identify the corresponding pages for V
  - **Pin** the pages in the buffer pool
  - Issue parallel asynchronous I/O requests for pages which are not in the buffer
  - Without waiting for the I/O completion, execution threads concurrently process vertices in V that are in the pages pinned
  - Slide the processing window one page at a time as soon as either an execution thread or a callback thread finishes the processing of a page

## Example



I = (0, 1, 1, 1, 0, 1, 1)

- identify pages  $(p_0, p_1, p_2, p_3, p_4)$  for I
- Pin  $p_1$  and  $p_2$ 2.
- 3. Issue asynchronous I/O request for  $p_0$
- Execution threads process  $v_2$ ,  $v_3$ , and  $v_5$  concurrently 4.
- 5. On completion of I/O request for  $p_0$ , callback threads process  $v_1$
- After processing any page, unpin the page and slide execution window 6. i.e., process  $p_3$  and finally process  $p_4$ 46

 $v_6$ 

 $p_4$ 

# **Processing Graph Queries**

- We support graph queries based on matrixvector multiplication
  - Targeted queries processing only part of a graph
  - Global queries processing the whole graph
- Targeted queries
  - BFS, K-step neighbors, Induced subgraph, K-step egonet, K-core, cross-edges etc.
- Global queries
  - PageRank, connected component

## **Experimental setup**

- Datasets
  - LiveJournal (4.8M vertices), Twitter (42M vertices), YahooWeb (1.4B vertices)
- Intel i7 6-core PC with 12 GB RAM
- 512GB SSD (Samsung 840 series)
- Bypass OS cache to guarantee real I/Os
- Main competitor
  - GraphChi

#### **Targeted Queries**



TurboGraph outperforms GraphChi by up to four orders of magnitude.

#### **Global Queries**



TurboGraph outperforms GraphChi by up to

27.69 times for PageRank.144.11 times for Connected Component<sup>+</sup>.

<sup>+</sup>upcoming paper for details and much faster performance

OPT: A New Framework for Overlapped and Parallel Triangulation in Large-scale Graphs [SIGMOD14]

# **Highlights of OPT**

- an overlapped and parallel disk-based triangulation framework for billion-scale graphs
- At the macro level, overlaps the internal triangulation and the external triangulation
- At the micro level, overlaps the CPU and I/O operations
- Achieves almost ideal performance

#### Internal triangles vs external triangles



#### **Running example**







 $\mathbf{p}_1$ 







54

# Load $p_1$ and $p_2$





p<sub>2</sub> p<sub>3</sub>

 $p_1$ 

55

 $p_4$ 

# Load p<sub>3</sub>





p<sub>2</sub> p<sub>3</sub>

 $p_1$ 

 $p_4$ 

# Load p<sub>4</sub>



#### MEMORY



57

#### **Effective disk access order of OPT**



Scan direction

#### **Comparison with GraphChi**

	LJ	ORKUT T	WITTER	UK
OPT	6.39	18.51	469.40	480.918
GraphChi-Tri	85.87	196.95	1850.26	4046.77
GraphChi-Tri/OPT	13.44	10.64	3.94	8.41

For the Yahoo dataset, OPT outperforms GraphChi by 31 times!

#### Conclusions

- We have presented a series of graph processing frameworks for large-scale graphs
  - iGraph 1.0 for graph indexing
  - iGraph 2.0 + Turbo<sub>ISO</sub> for subgraph isomorphism
  - TurboGraph for graph analytics



# Handling general vectors

- Indicator vector can be implemented as a bitmap
- However, what if we want to use general vectors instead?
  - Consider PageRank where we need random accesses to pagerank values and out-degrees in two general vectors

# Main idea of handling general vectors

- Adopt the concept of block-based nested loop join
- a general vector is partitioned into multiple
   chunks such that each chunk fits in memory
- Regard the pages pinned in the current buffer as a block
- Join a block with a chunk of each random vector in-memory until we consume all chunks
- Hide this mechanism as much as possible from users! (see paper)

#### Example

