

# Large-scale graph recessing with emerging storage evices

Huibing Dong













Algorithm	Description
PageRank	Ranking nodes based on their incoming/outgoing edges
Breath-First Search	Graph traversal based on breath
Shortest Path	Finding the path between 2 nodes where the sum of weights is minimal
Pattern Matching	Finding certain structures (e.g. path, star)
Triangle Count	Counting the number of triangles
Connected Component	Finding the subgraphs in which any two vertices are connected



Algorithm	Description
PageRank	Ranking nodes based on their incoming/outgoing edges
Breath-First Search	Graph traversal based on breath
Shortest Path	Finding the path between 2 nodes where the sum of weights is minima
Pattern Matching	Finding certain structures (e.g. path, star)
Triangle Count	Counting the number of triangles
Connected Component	Finding the subgraphs in which any two vertices are connected

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

$$Iterative \qquad Neighbors matter$$







Irregular access



Expensive DRAM

## Shared-memory

Memory Single machine

Limited graph size

2013-Ligra, 2014-GraphX, 2015-Chaos Distributed

Memory

Clusters

Costly

2010-Pregel, 2010-Graphlab, 2012-PowerGraph, 2014-GraphX, 2014-Blogel, 2015-Chaos, 2016-Gemini

### External-memory

Memory + Storage Single machine Larger size + cost efficient

2012-GraphChi, 2013-X-stream, 2013-TurboGraph, 2015-GridGraph, 2015-FlashGraph, 2017-Graphene, 2017-Mosaic, 2018-GraFBoost, 2019-GraphOne, 2019-Lumos

\*Blue: Semi-external systems

Others architecture

Graph-optimized database: Neo4j

# **PLATFORM FAMILY**





#### **Graph Processing Platforms**

Large scale graph processing systems: survey and an experimental evaluation, Omar Batarfi et. al, Cluster Computing'15

# **EXTERNAL/OUT-OF-CORE SYSTEM**





# **GRAPH REPRESENTATION**









Graph

Adjacency list

#### Adjacency matrix

✓ Sparse matrix

 $\checkmark$  Much less storage space needed

Storage format: Compressed Sparse Column (CSC) & Compressed Sparse Row (CSR) files

#### Semi-external systems

- Vertex data
  - ✤ In the main memory
  - Fine-grained accesses, byte-addressable
- Edge data
  - On the secondary storage
  - ✤ Coarser accesses

### External systems

- Even the vertex data itself is too large
- Both Vertex & Edge data on the secondary storage





Continue to scale Large capacity with lower latency





# Programming Model Vertex-centric Edge-centric \*IO-centric

## **Execution Model**

Bulk synchronous

Asynchronous

# **PAPER LIST**



## **External systems**

OSDI'12	GraphChi	Parallel Sliding Window	8GB DRAM + 256GB SSD + 750GB HD
KDD'13	Turbograph		12GB DRAM + 2x 512GB SSD
SOAP'13	X-Stream	Edge-centric	64GB DRAM + 2x 200GB
ICDE'15	Venus		
ATC'15	GridGraph	2D partition	30.5GB DRAM + 2x 2TB HDD/1800 GB SSD
FAST'15	FlashGraph	semi-external; merges I/O	512GB DRAM + 15x OCZ Vertex 4 SSD
SC'16	G-Store		
FAST'17	Graphene	Fine-grained IO	128GB DRAM + 16x 500GB SSD
EuroSys'17	Mosaic		
ISCA'18	GraFBoost		48GB DRAM + Xilinx VC707 FPGA + 2x 512GB SSD
FAST'19	GraphOne		512GB DRAM + 512GB SSD
ATC' 19	Lumos	Dependency-Driven	64GB DRAM + 4x 1.9TB SSD

# **PAPER LIST**



## GraphChi

Parallel Sliding Window

- ✤ one sub-graph at a time
- ✤ 3 phases



## X-stream

#### Edge-centric



## GridGraph



## Graphene

- Semi-external
- ✤ Merge edges I/O requests

FlashGraph





••		•••		•	
••••		•••	••••		•••
••••					

# **POTENTIAL RESEARCH PROBLEMS**

- Programming model optimization
- Execution model support
- Partitioning
- Serializing
- Emerging Storage devices selection:
  - Zone-named space SSD
    - ZAC/ZBC
  - Open-channel SSD
  - Hybrid devices
  - Storage tiers