

The Future of Graph Computing





Graph Ecosystem : What we have achieved so far

Algorithms :

- Tons of new algorithms ... (ICML/ KDD / SIGMOD...)
- Approximate algorithm, Lower time/space complexity

High Performance Computing

- Parallel and Distributed Processing
- Leveraging new accelerators and storage layers

Graph Analytics Library

- GraphX / Spark, Apache Giraph, ScaleGraph, etc

Graph Database

- IBM System G, Neo4J, Titan, ... etc

Benchmarking

- Graph500, LDBC, etc

Statistical Survey on Graph-Related Papers

• Top-Tier Conferences over 5 years (2011-2015) in 3 research fields

- **HPC**: Supercomputing, IPDPS
- •DB: ICDE, SIGMOD, VLDB
- ML: ICML, SDM, KDD, ICDM

Total



ML (ICML, SDM, KDD, ICDM)



HPC (Supercomputing and IPDPS)



DB



Graph500 Trend : Innovative Graph Algorithms and Implementation brings us 7x performance



Insights from our Graph500 Challenge

1024 nodes (12,288 cores)

We learned that algorithmic and implementation innovation greatly accelerated the performance on the same platforms . What's the next to do ?

GTEPS x13.4 **GTEPS x7** 1600 50000 1400 40000 1200 1000 30000 800 20000 600 400 10000 200 0 0 2012/11 2013/06 2013/11 2014/06 2014/11 2015/06 2015/11 2011/11 2012/06 2012/11 2013/06 2013/11 2014/06 2014/11 Tokyo Tech – TSUBAME 2.5 **Riken K Computer**

GRAPH 500

82,944 nodes (663,552 cores)

Even If you optimize popular frameworks based on JVM, it is still slow ;(. How do we solve this dilemma ?

ScaleGraph vs. GraphX/Spark

Graph analytics involves more IO-bounded operation considering efficient data communication among compute nodes, memory access, workload imbalance, etc.

450 400 350 300 (ع) 250 ية 200 ScaleGraph 150 GraphX/Spark 100 50 0 2 8 16 1 4 Nodes

Weak Scaling (Scale 18), PageRank (30 Steps)

My Observation

- General computational model for incremental graph analytics
 - (vs. Pregel, GMI-V for batch-analytics)
 - A bunch of work on incremental algorithms (incremental PageRank/Community Detection)
 - \rightarrow Could be generalized ?
 - Incremental Pregel ? Other method ?
- How would we align with other areas such as database and machine learning / data mining area ?
 - We just follow the outcome from the ML field ?

Financial Rick Prediction using System G

Risk Factors

- Adjunct Relationship
- Guarantee Relationship
 - Adjunct and Guarantee
- Transaction Relationship
 - Stockholder Relationship

EgoNet Relationship

Enterprise Type





Graph Analysis Machine Learning Cognitive Reasoning

Feature Pattern of Clients with high risk





Graph Database Manage and retrieve Linked Big Data





The Graph 500 List

November 2013

No.	Rank	Machine	Installation Site	<u>Number</u> of nodes	Number of cores	Problem scale	GTEPS
1	1	DOE/NNSA/LLNL Sequoia (IBM - BlueGene/Q, Power BQC 16C 1.60 GHz)	Lawrence Livermore National Laboratory	65536	1048576	40	15363
2	2	DOE/SC/Argonne National Laboratory Mira (IBM - BlueGene/Q, Power BQC 16C 1.60 GHz)	Argonne National Laboratory	49152	786432	40	14328
3	3	JUQUEEN (IBM - BlueGene/Q, Power BQC 16C 1.60 GHz)	Forschungszentrum Juelich (FZJ)	16384	262144	38	5848
4	4	K computer (Fujitsu - Custom supercomputer)	RIKEN Advanced Institute for Computational	65536	524288	40	5524.12

Home Complete Results Benchmarks Green Graph 500 Log In

IBM System G: Graph Computing Framework



© 2016 IBM Corporation

IEM

Graph Analytics Project with Financial Institute

■ **Project Duration** : 1-2 months (1 – 2 people fully involved with the project)

Graph Size and Hardware:

- Property graph with more than 10 attributes
- Graph Size :
 - First phase : Millions of vertices and edges
 - <u>Second phase</u>: Billion-scale graph with more fine-grained time-series transaction data
- Hardware : 60 CPU cores in total with shared-memory machine

Some Graph Analytics

- Cycle detection with various condition on properties
- Egonet extraction
- Prediction based on Machine Learning



Graph Database / Other data stores

Missing Graph Component from Industry Point of View IBM &

1. Scalable Graph Database with trillion-scale time-series data

- -Seamless programming model with regular graph processing
- -Leveraging Memory Hierarchy

17

2. More powerful Graph Query language for temporal graphs

-Gremlin, Sparql, Cypher... are not enough

3. <u>Multi-layer and Multiplex networks</u>

Multi-modal data source brings requirements for multi-layer networks

4. <u>How to leverage Accelerators for high performance graph</u> processing ?

- -Vertex-centric model (Pregel) on GPGPUs / NVLink
- –**Domain-specific language** for Pregel that allows you to write your Graph algorithms on CPU, or GPU, or hybrid engines.

IBM. 🍯

5. <u>General Incremental Processing Model for Dynamic Graphs</u>

- -Vertex-centric model is still attractive since it brings better parallelism
- -Pregel + Spark RDD-like paradigm might be better combination .
- -But like the modularity-optimization based paradigm, it needs global synchronization
- 6. Workflow language for ETL, Graph processing, and other machine learning
- 7. Deep learning for Static/Dynamic Graphs
 - -How could we apply CNN / RNN for static / dynamic graphs ?
 - Applications: Fraud detection, Non performing loan, and anti-money laundering



- Design A : More intuitive way for expressing graphs
 - Add_edge(A, B, tk, 100)
- Design B : Separating time-series data from connectivity
 - Users can define graph model as Design A, but the system could convert Design A to Design B automatically.
 - Eab.addTransaction(tk, 100)





Getting a snapshot of Temporal Graph

Give me a snapshot at 2015/11/13 8:30 from temporal graph





Getting a snapshot of Temporal Graph

• How can efficiently store a temporal graph efficiently ?

- -1) Space efficiency for storing many snapshots
- -2) Or read efficiency to access particular snapshot

Versioning Issue

21

One year ago – 2015/11/13



Leveraging Memory Hierarchy

Vertex property (time, space) Latest data should be in DRAM

- -frequently accessed vertices such as Superhubs
- Edge property (time, space)





Discussion ! Thank You