Building In-memory Graph Store for <u>Fast</u> and <u>Concurrent</u> Querying with Advanced Hardware Features

RONG CHEN IPADS, Shanghai Jiao Tong University IPADS Workshop, 2018



Joint work w/ Haibo, Jiaxin, Youyang, Yunhao, Youyang, Siyuan, Chang, Ning, Jiaqi, Xiating, Xuehan, Wenhao, and Zhenhan @IPDADS

Research in DS Group @IPADS



WUKONG Project







Graphs are Everywhere



Online **graph query** plays a vital role for searching, mining and reasoning linked data



Knowledge Graphs and Querying



+ F	ollow
-----	-------

Ö

Google knowledge graph has more than 70 billion facts about people, places, things. + language, image, voice translation -

"...Every piece of information that we crawl, index, or search is analyzed in the context of Knowledge Graph."

How a Database of the World's Knowledge Shapes Google's Future MIT Technology Review, January 27, 2014

"...Google users will able to browse through the company's 'knowledge graph,' or its ever-expanding database of information about 'entities' – people, places and things – the 'attributes' of those entities and how different entities are connected to one another."

> What Google's Search Changes Might Mean for You Wall Street Journal, March 14, 2012

Bing's Satori Adds Timeline Data For About 500k Famous People

Isaac Asimov

Matt McGee on February 21, 2014 at 12:59 pm

Abraham Lincoln Abraham Lincoln was the 16th President of the United States, serving from March 1861 until his assassination in April 1865. Lincoln led the United States through its Cult War-its bloodjets war and

its greatest moral, constitutional and politica... 4 en.wikipedia.org

.ived: Feb 12, 1809 - Apr 15, 1865 (age 56) Height: 6' 4" (1.93 m)

Spouse: Mary Todd Lincoln (1842 - 1865 Romance: Ann Rutledge

Children: Robert Todd Lincoln · Tad Lincoln · William Wallace Linco Edward Baker Lincoln

Related movies: Lincoln - Abraham Lincoln: Vampire Hunter



1835: Lincoln's first romantic interest was Ann Rutledge, whom he met when he first moved to New Salem; by 1835, they were in a relationship but not formally engaged.

1836: Admitted to the bar in 1836, he moved to Springfield, Illinois, and began to practice law under John T. Stuart, Mary Todd's cousin. en wikeda org Born: Jan 2, 1920 - Petrovichi, Russia Didei: Apr 6, 1992 - New York, New York Spouse: Janet Asimov (1973 - 1992) - Gertrude Blugeman (1942 -

ac Asimov was an American author an

rofessor of biochemistry at Boston University

ast known for his works of science fiction and for

his popular science books. Asimov was one of the

ost prolific writers of all time, having writte....

Awards: Nebula Award for Best Novel - Hugo Award for Best Novel Hugo Award for Best Novelette +

Education: Columbia University

Parents: Judah Asimov - Anna Rachel Berman Asimov

1941: In September 1941 Astounding published the 32nd story Asimov wrote, "Nightim", which has been described as one of "the most famous science-fiction stories of all time".

1955: Doubleday also published collections of Asimov's short stories beginning with The Martian Way and Other Stories in 1955.

Bing's Knowledge Repository, Satori, Adds More Interactive Content



RDF and **SPARQL**

RDF: <u>R</u>esource <u>D</u>escription <u>F</u>ramework

- Public: DBpedia, PubChemRDF, Bio2RDF
- Google's Knowledge Graph, Bing's Satori, Facebook's, LinkedIn's, Yahoo's, ...



<u>SPARQL</u> Protocol and RDF Query Language



Learning SPARQL @LearningSPARQL · 19 Nov 2017 Google seeks a "Linguist/Ontologist, Google Knowledge Graph" with "Experience with ontology development (RDF(S)/OWL, SPARQL, the Semantic Web or Frame based KR systems)" indeed.com/viewjob?jk=181...



Learning SPARQL @LearningSPARQL · 27 Mar 2013 Satori (Bing's answer to Google's Knowledge Graph) uses SPARQL research.microsoft.com/en-us/projects...

Pub



BIO22RDF

RDF and **SPARQL**

- mo: MemberOf
- ad: ADvisor
- to: TeacherOf
- tc: TakeCourse

RDF is a graph composed by a set of (Subject, Predicate, Object) triples

	- V	¥	
Rong	to	DS	Kong US
Rong	mo	IPADS	
Siyuan	ad	Rong	ad (Siyuan) tc
Siyuan	tc	0S	
Haibo	to	OS	mo to OS
Haibo	mo	IPADS	Haibo
Jiaxin	ad	Haibo	ad Jiavin to
•••			

RDF and **SPARQL**

SPARQL is standard query language for **RDF**



Queries are Heterogeneous

Heavy Query (Q_H)

- Non-selective query
- Start from a set of vertices
- Explore a large part of graph

SELECT ?X ?Y ?Z WHERE {
 ?X teacherof ?Y .
 ?Z takecourse ?Y .
 ?Z adivsor ?X .
 }



Queries are Heterogeneous

Heavy Query (Q_H)

- Non-selective query
- Start from a set of vertices
- Explore a large part of graph

Light Query (Q_L)

- Selective query
- Start from a given vertex
- Explore a small part of graph

SELECT ?Y WHERE {
 ?X memberOf IPADS
 ?X teachOf DS .
 ?Y advisor ?X . }





Triple Store and Scan-Join

Store RDF data as a set of triples in RDBMS



SELECT ?Y WHERE {
TP1 ?X memberOf IPADS
TP2 ?X teachOf DS .
TP3 ?Y advisor ?X . }



SELECT ?Y WHERE {
TP1 ?X memberOf IPADS
TP2 ?X teachOf DS .
TP3 ?Y advisor ?X . }

Triple Store and Scan-Join

Store RDF data as a set of triples in RDBMS







	y (ms)	Latenc	CORES	CONE	LUBM
	LIGHT	HEAVY	CORES	CONF	#T=1.41T
Graph Store	5.65	12,034	144	VLDB	Trinity.RDF
Triple Store	8.72	1,920	96	sigmod	TriAD

[Wukong, OSDI 2016]

RDMA: Remote Direct Memory Access

- High speed, low latency, and low CPU overhead
 - Interface: IPoIB, SEND/RECV Verbs, READ/WRITE (one-sided primitive)
 - Bypass OS kernel: zero copy
 - Round-trip time: one-sided/~1-3µs, verb msg/~7µs, IPoIB/~100 µs



LUBM	CONE	CODES	Latency		
#T=1.41T	CONF	CORES	HEAVY	LIGHT	
Trinity .RDF	VLDB	144	12,034	5.65	Graph Store
TriAD	SIGMOD	96	1,920	8.72	Triple Store
+RDMA	40G IB	96	+4%	+137%	

[Wukong, OSDI 2016]



Systematic Approaches



RDMA-friendly in-memory graph store



2. Differentiate Partitioning







RDMA-enable graph exploration



I. Full-history Pruning 2. Data/Exec Migration 3. Worker-obliger Strategy





Graph Model and Indexing

		SELECT ?X ?Y ?Z WHERE {
	TP1	?X teacherof ?Y .
	TP2	?Z takecourse ?Y .
	TP3	<pre>?Z adivsor ?X .</pre>
$(Rong) \longrightarrow (DS) tc$		
/mo (Chang)		
ad since to		
TPADS Siyuan CC		
to (OS)		
(Haibo)		
tc		SELECT ?X ?Y WHERE {
ad Jiaxin	TP1	?X type Course .
	TP2	<pre>?Y teacherOf ?X . }</pre>



Differentiated (Graph) Partitioning



Differentiated (Graph) Partitioning



Start from normal vertex
Exploit locality

SELECT ?X ?Y WHERE {	
<pre>?X teacherOf ?Y .</pre>	TP1
?Z takerCourse ?Y .	TP2
<pre>?Z advisor ?X . }</pre>	TP3

- Start from index vertex
- Exploit parallelism





Vertex Decomposition



Vertex Decomposition







- Inefficient lookup
- Unnecessary data transfer

Vertex Decomposition



Traditional (e.g., Trinity.RDF) key value Rong→IN=ad:Siyuan;OUT=mo:IPADS,to:DS,..

SELECT ?Y WHERE	{			
RONG teacherOf	?X	•		TP1
?Y takerCourse	?X	0	}	TP2

- Inefficient lookup
- Unnecessary data transfer

Decomposition



Efficient for both local and remote (RDMA) accesses

Systematic Approaches



RDMA-friendly in-memory graph store



I. Model & Indexing 2. Differentiate Partitioning





RDMA-enable graph exploration

I. Full-history Pruning



2. Data/Exec Migration









Migrate Execution or Data

Fork-join (migrate exec)



Send sub-query by RDMA WRITE
 Async exploration w/ full-History
 Exploit high parallelism

In-place (migrate data)



- Fetch data by RDMA READ
- Bypass remote CPU & OS
 Exploit low latency

Worker-obliger Strategy



- Latency-centric work-stealing algorithm
- Ring policy: practical, effective, and easy to impl.

LUBM	CONE	CORES	Latency (m)		
#T=1.41T	CONF	CORES	HEAVY	LIGHT	
Trinity .RDF	VLDB	144	12,034	5.65	
TriAD	sigmod	96	1,920	8.72	
+RDMA	40G IB	96	+4%	+137%*	

^{*} using only 1 core

[Wukong, OSDI 2016]

LUBM	CONE	CODES	Laten	cy (m)
#T=1.41T		COKES	HEAVY	LIGHT
Trinity .RDF	VLDB	144	12,034	5.65
TriAD	sigmod	96	1,920	8.72
+RDMA	40G IB	96	+4%	+137%*
WUKONG	OSDI	96	248	0.40 [*]
			7.4X	21.8X

[•] using only 1 core

[Wukong, OSDI 2016]

LUBM	CONE	CODES	Laten	cy (m)	ТНРТ
#T=1.41T		CORES	HEAVY	LIGHT	(q/s)
Trinity .RDF	VLDB	144	12,034	5.65	≈ 400
TriAD	sigmod	96	1,920	8.72	250
+RDMA	40G IB	96	+4%	+137%*	
WUKONG	OSDI	96	248	0.40*	269 <u>K</u>
			7.4X	21.8X	I,076X
only 1 core				[Wukor	ng, OSDI 2018

Conclusion: Wukong@OSDI16

New hardware technologies open opportunities

Wukong: a distributed in-memory RDF store that leverages RDMA-based graph exploration to support fast and concurrent RDF queries

Achieving orders-of-magnitude lower latency & higher throughput than prior state-of-the-art systems



Website: <u>http://ipads.se.sjtu.edu.cn/projects/wukong</u> GitHub: <u>https://github.com/SJTU-IPADS/wukong</u>

Fast and Concurrent RDF Queries using RDMA-assisted GPU Graph Exploration



Queries are Heterogeneous

Heavy Query (Q_H)

- Start from a set of vertices
- Explore a large part of graph





3000X

 $Q5^*$

- Light Query (Q_L)
 - Start from a given vertex
 - Explore a small part of graph 0.13 ms

* Wukong on a 10-server cluster for LUBM-10240 dataset

Concurrent Workload





Query Execution on GPU

SELECT ?X ?Y ?Z WHERE {TP1?X teacherof ?Y .TP2?Z takecourse ?Y .TP3?Z adivsor ?X .





Systematic Approaches

1. Smart data prefetching

2. GPU-friendly key/value store

3. Heterogeneous RDMA comm.









LUBM	CONE	CODES	Latenc	x y (ms)	THPT	(q/s)
#T=1.41T	CONF	CORES	HEAVY	LIGHT	HEAVY	LIGHT
TriAD	sigmod	100	3,094	7.04		
WUKONG	OSDI	100	215	0.34*	7.7	40



[Wukong+G, ATC 2018]

LUBM	CONE	CORES	Latency (ms)		THPT	(q/s)
#T=1.41T		CORES	HEAVY	LIGHT	HEAVY	LIGHT
TriAD	sigmod	100	3,094	7.04		
WUKONG	OSDI	100	215	0.34 [*]	7.7	40
Wuk+G	ATC	100 10 GPU	47	0.38*	45.4	346K
			7.4X	- %	5.9X	8.650X

[Wukong+G, ATC 2018]

[•] using only 1 core

Conclusion: Wukong+G@ATC18

Hardware heterogeneity opens opportunities for hybrid workloads on graph data

Wukong+G : a distributed RDF query system supports heterogeneous CPU/GPU processing for hybrid queries on graph data

Outperform prior state-of-the-art systems by more than one order of magnitude when facing hybrid workloads



Website: <u>http://ipads.se.sjtu.edu.cn/projects/wukong</u> GitHub: <u>https://github.com/SJTU-IPADS/wukong</u> Sub-millisecond Stateful Stream Querying over Fast-evolving Linked Data



Streaming Data and Querying

Multiple data sources are continuously generating **streaming data** in **high velocity**



Example: Social Networking



Stored Data



<Rong, creates, Feed>. 12:30 <Feed, hash_tag, SOSP>. 12:30 <Yunhao, likes, Feed>. 12:31 <Haibo, likes, Feed>. 12:40

Streaming Data

Example: Social Networking



Stored Data

Streaming Data

Example: Social Networking

Last 30 minutes, which IPADS members created feeds that are liked by other IPADS members?



time Registered by user Triggered by system Triggered by system Triggered by system Canceled by user

Stateful Streaming Query

A stateful streaming query needs to integrate streaming data with stored data



Real-time User Activity

Durable Social Graph

Conventional Approach



Conventional Approach

1. Cross-system Cost

~40% execution time waste on data transformation and transmission

2. Inefficient Query Plan

Semantic gaps between the two systems impair query optimizations

3. Limited Scalability

Stream processing systems dedicate all resources to improve the performance of a single job

Integrated Design



Systematic Approaches

1. Hybrid store: timed & timing

2. Stream index & partitioning

3. Bounded snapshot scalarization







LSBENCH #T=118M R=133K/s	CORES	Wukong Storm	CSPARQL- engine	
GEO. M	24 x1	5.91	757	

LSBENCH	CORES	Wuk	ONG	SPARK	
R=133K/s	CORES	STORM	HERON	STREAMING	
GEO. M	24 x8	6.29	5.85	679	

[Wukong+S, SOSP 2017]

LSBENCH #T=118M R=133K/s	CORES	Wukon Stor	IG CSP. M er	ARQL- ngine	Wυκc	ong+S	
GEO. M	24 x1	5.91 757		0.48ms			
		12.3	X I	,577X			
LSBENCH #T=3.75B R=133K/s	CORES	Wuk Storm	ong Heron	Spark Streaming		Wυκοι	NG+S
GEO. M	24 x8	6.29 3.7X	5.85 2.7X	I,	679 ,476X	0.4	6ms

[Wukong+S, SOSP 2017]

Conclusion: Wukong+S@SOSP17

Wukong+S : distributed querying engine adopting integrated design for stateful stream queries over fast-evolving linked data

Achieving **sub-millisecond** latency and exceeding **one million** queries per second



Website: <u>http://ipads.se.sjtu.edu.cn/projects/wukong</u> GitHub: <u>https://github.com/SJTU-IPADS/wukong</u>

Current Projects on Wukong

Lightweight, non-invasive migration for graph store

Fast and accurate optimizer for graph query

Supporting pipeline workload on linked data

Shard-based Migration



	Tdeal	Shard-based		
	IdedI	Before	After	
Throughput (K ops/sec)	3,204	116	159	
Median/50 th Latency (msec)	0.61	17.15	12.02	
Tail/99 th Latency (msec)	2.45	78.60	63.50	
Remote Access Rate (%)	0	87.9	68.1	
Data Volume Migrated (GB)	-	-	5.xx	

Dataset: Graph500

Workload: 2-hop gueries

- Graph data: poor locality Additional metadata (POS)
- Conflict with READ & WRITE operations

Split Live Migration



Integrate with Location Cache



Lightweight & Non-invasive Migration



Fast and Accurate Optimizer



Wukong							
OPT	OSDI16	Trinity.RDF	Wukong+P				
311	513	325 0.073	311 0.520				
80	80	80 0.011	80 0.012				
225	227	274 0.119	†(225) 0 0.305				
0.030	0.030	0.030 0.021	0.030 0.008				
0.023	0.023	0.023 0.005	0.023 0.005				
0.095	0.095	0.097 0.009	0.095 0.007				
194	309	347 0.077	195 0.290				
4.93	5.67	6.25	(5.31)2.07				

Cardinality Estimation: Type-centric
 Cost model: Mimic-based
 Plan enumeration: budget-aware



Pipeline Graph Processing





	1		-					
	Query			Transform		Preprocessing	Analytics	Total
Lubm-640/ 4-node /s	excutution	result bring back	final aggregate and send back	maping	format			
Wukong + powerlyra	0.640	3.096	9.82		0	2.86	2.64	15.776
Wukong + Gemini	0.640	3.096	9.82	6.28	0.67	2.97	0.25	23.086
Wukong-pipeline	0.462	1.343	0		0	0.702	0.322	2.829
					1			



Website: <u>http://ipads.se.sjtu.edu.cn/projects/wukong</u> GitHub: <u>https://github.com/SJTU-IPADS/wukong</u>

Wukong, short for Sun Wukong, who is known as the Monkey King and is a main character in the Chinese classical novel "Journey to the West". Since Wukong is known for his extremely fast speed (21,675 kilometers in one somersault) and the ability to fork himself to do massive multi-tasking, we term our system as Wukong.

Questions









