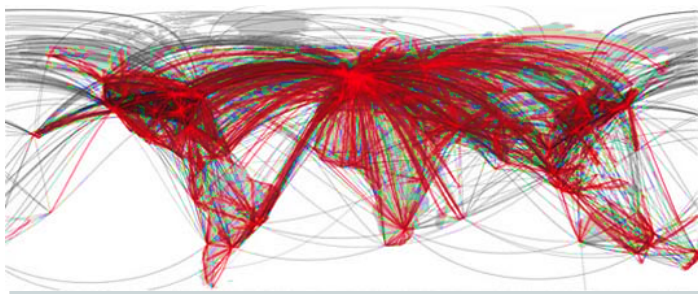# Fast and Concurrent RDF Queries using *RDMA-assisted* GPU Graph Exploration
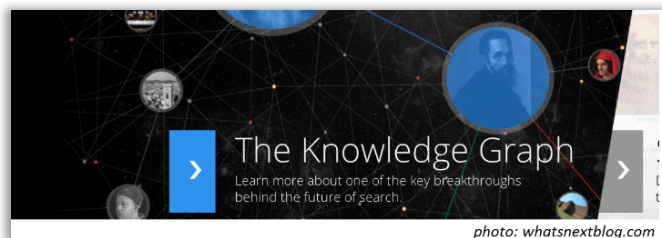
Siyuan Wang, Chang Lou, **Rong Chen**, Haibo Chen

Institute of Parallel and Distributed Systems (IPADS)
Shanghai Jiao Tong University
http://ipads.se.sjtu.edu.cn

# Graphs are Everywhere

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

The Knowledge Graph
Learn more about one of the key breakthroughs behind the future of search.

*photo: whatsnextblog.com*

Jena

TAO

Unicorn

# RDF and SPARQL

## Resource Description Framework (**RDF**)

► Representing linked data on the Web

► Public knowledge bases: DBpedia, Wikidata, PubChemRDF

► Google's knowledge graph

# RDF and SPARQL
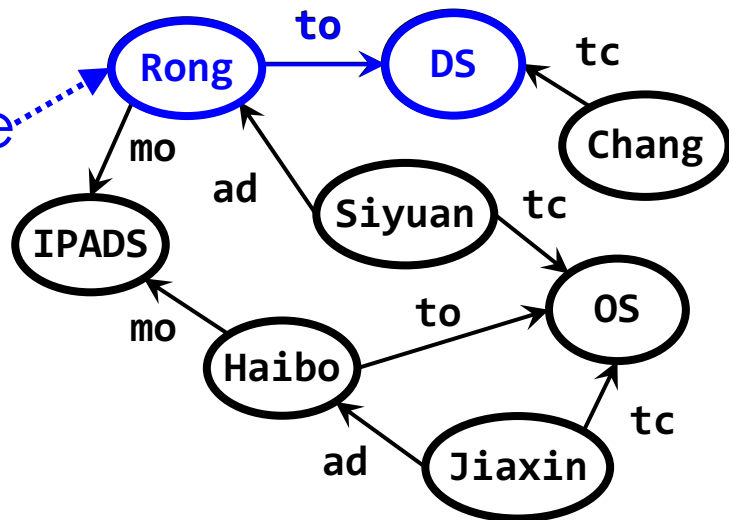
RDF is a graph composed by a set of ⟨**S**ubject, **P**redicate, **O**bject⟩ triples



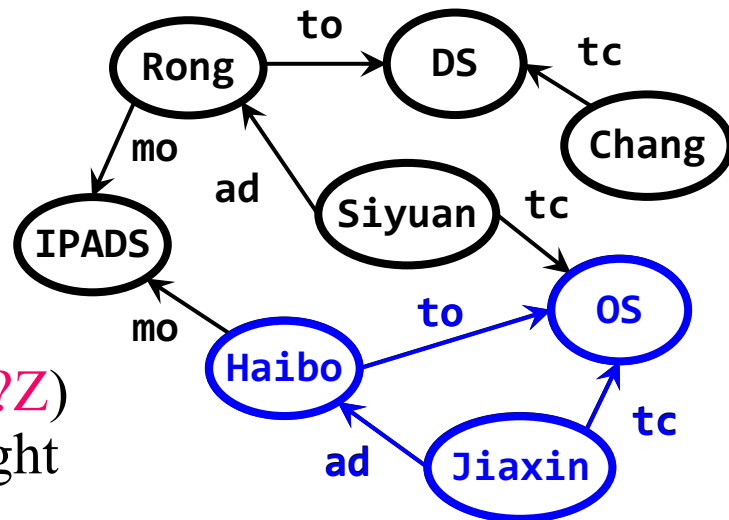| Rong   | to  | DS    |
|--------|-----|-------|
| Rong   | mo  | IPADS |
| Siyuan | ad  | Rong  |
| Siyuan | tc  | OS    |
| Haibo  | to  | OS    |
| Haibo  | mo  | IPADS |
| Jiaxin | ad  | Haibo |
| . . .  |     |       |

triple

4

# RDF and SPARQL

SPARQL is standard query language for RDF

**Triple Pattern**

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

TP1
TP2
TP3

**Variable**

Professor (?X) advises (**ad**) student (?Z) who also takes (**tc**) a course (?Y) taught by (**tc**) the professor
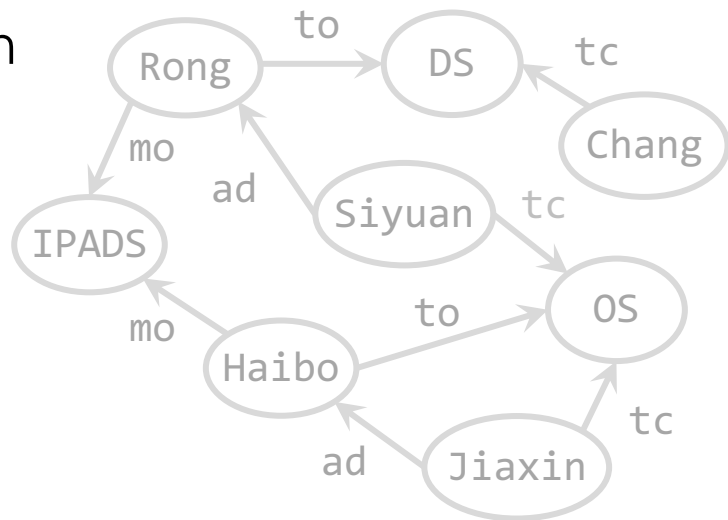
# Queries are Heterogeneous

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

## Heavy Query (Q$_H$)

► Start from a set of vertices

► Explore a large part of graph

## Light Query (Q$_L$)

# Queries are Heterogeneous

```
SELECT ?X ?Y ?Z WHERE {
TP1   ?X teacherof ?Y .
      ?Z takecourse ?Y .
      ?Z adivsor ?X .        }
```

## Heavy Query ($Q_H$)

► Start from a set of vertices

► Explore a large part of graph

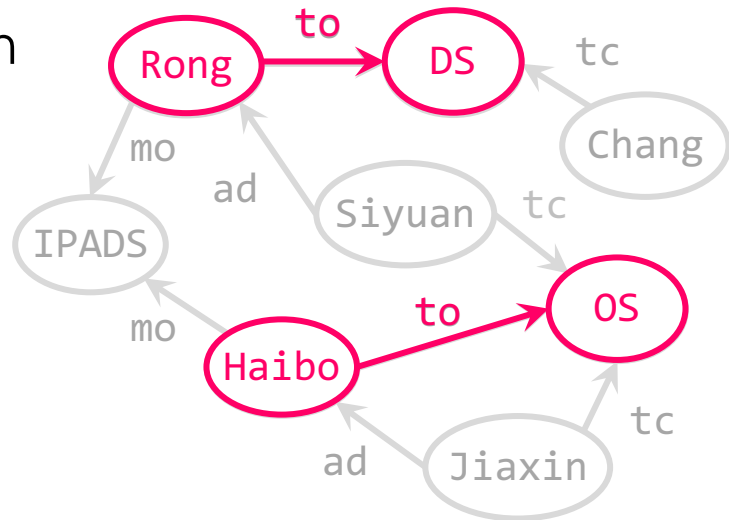## Light Query ($Q_L$)

# Queries are Heterogeneous

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

## Heavy Query ($Q_H$)

► Start from a set of vertices

► Explore a large part of graph

## Light Query ($Q_L$)

# Queries are Heterogeneous

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

## Heavy Query (Q_H)

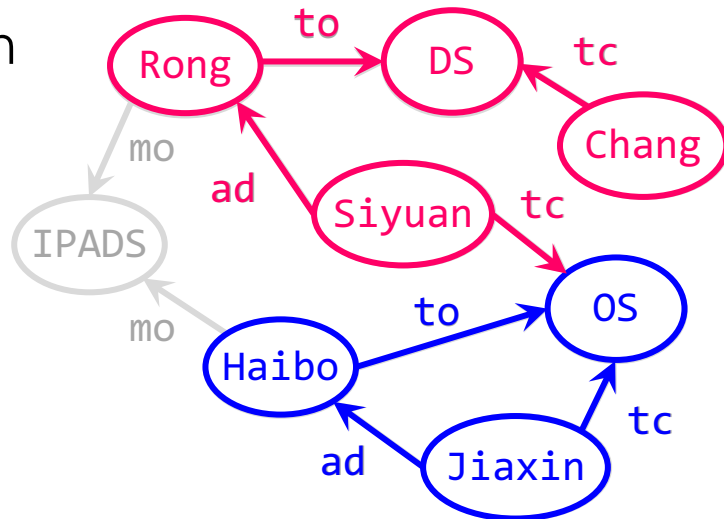► Start from a set of vertices

► Explore a large part of graph

## Light Query (Q_L)

# Queries are Heterogeneous

```
SELECT ?X WHERE {
  ?X advisor Rong .
  ?X takecourse OS . }
```

## Heavy Query ($Q_H$)

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

## Light Query ($Q_L$)

► Start from a given vertex
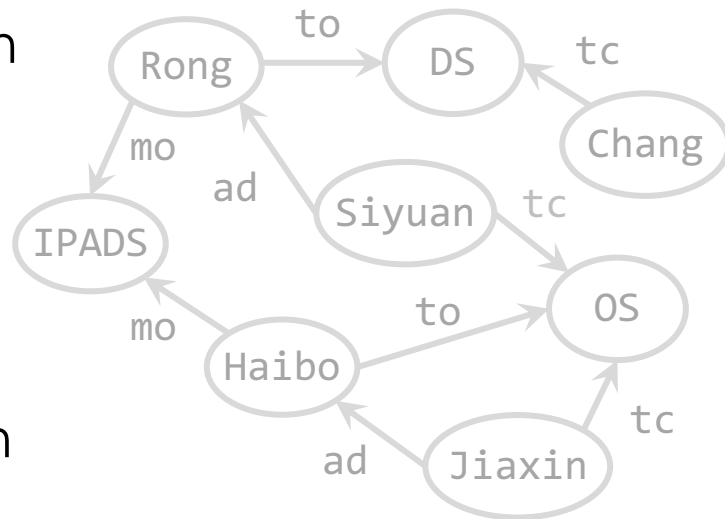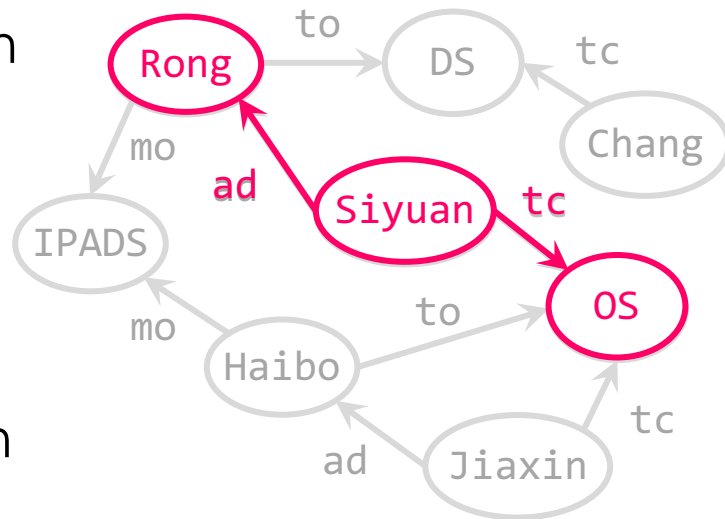► Explore a small part of graph

# Queries are Heterogeneous

```
SELECT ?X WHERE {
  ?X advisor Rong .
  ?X takecourse OS . }
```

## Heavy Query ($Q_H$)

► Start from a set of vertices

► Explore a large part of graph

## Light Query ($Q_L$)

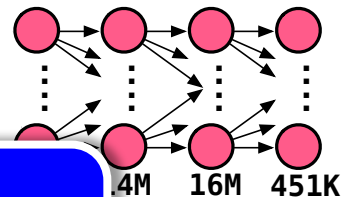► Start from a given vertex

► Explore a small part of graph

# Queries are Heterogeneous

Heavy Query (Q_H)

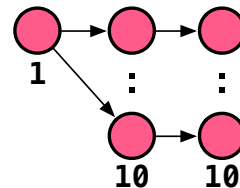► Start from a set of vertices

► Expl

Q7*
390 ms

.4M    16M    451K

**Incompetent to handle heavy queries efficiently**

Light Query (Q_L)

► Start from a given vertex

► Explore a small part of graph

Q5*
0.13 ms

1    :    :

10    10

# Concurrent Workload

**logarithmic scale**

Median Latency (msec)

$10^4$
$10^3$
$10^2$
$10^1$
$10^0$
$10^{-1}$

**Latency**

**Throughput**

$10^0$ $10^1$ $10^2$ $10^3$ $10^4$ $10^5$ $10^6$

Throughput (queries/sec)

**better**

# Concurrent Workload

**logarithmic scale**

**Latency**

**Throughput**

PURE light query workload

Pure:Light

**Thpt: 398K q/s**

**Lat: 0.10 ms**

Median Latency (msec)

Throughput (queries/sec)

# Concurrent Workload

HYBRID light & heavy query workload

Thpt: 10 q/s
Lat: 8,600 ms

logarithmic scale



Hybrid:Heavy    Pure:Light

Lat: 100 ms

Thpt: 398K q/s

Lat: 0.10 ms

Latency

Media

Throughput (queries/sec)

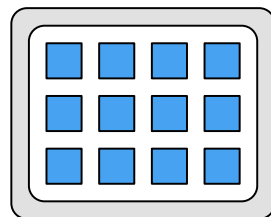$10^0$   $10^1$   $10^2$   $10^3$   $10^4$   $10^5$   $10^6$

Throughput

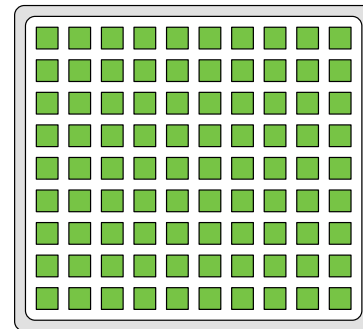**poor** performance when facing **hybrid** workload

# Advanced Hardware

## Heterogeneity

▶ **GPU** has many cores and high memory bandwidth

CPU

GPU

## Fast Communications
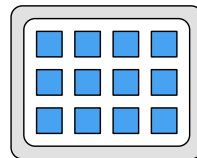
▶ **RDMA** enables direct data transfer between machines

Server

Server

memory

memory

RNIC

RNIC

CPU

CPU

# General Idea

**Heterogeneous**
Workload

**Heterogeneous**
Hardware



Light
Query

CPU

Heavy
Query

GPU

1.7M   14M   16M   451K

# System Overview

Wukong+G : a distributed graph query system that can leverage CPU/GPU to handle hybrid workload

1. GPU-enable graph exploration
2. GPU-friendly RDF store
3. Heterogeneous RDMA Communication

Performance improvement

► Latency: **2.3X** - **9.0X** speedup for **heavy** query
► Throughput: **345K** queries/sec in **hybrid workloads**

# Wukong Architecture

# Wukong+G Architecture

# Query Execution on CPU

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

# Query Execution on CPU

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

22

# Query Execution on CPU

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

**Work Thread**

**Metadata**

**Graph Store**

CPU

Query

TP-1
TP-2
...

Key

Value

Store

CPU DRAM

**Intermediate Result**

History Table

?X ?Y ?Z

# Query Execution on CPU

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

**Work Thread**

**Metadata**

**Graph Store**

Query

Key    Value    ore    RAM

**Observation**: all of traversal paths can be explored **independently**

**Intermediate Result**

History Table

**?X ?Y ?Z**

# Query Execution on CPU

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



**Work Thread**

**Metadata**

**Graph Store**

CPU

Query

TP-1
TP-2
...

Key

Value

Store

CPU DRAM

**Intermediate Result**

History Table

?X ?Y ?Z

25

# Query Execution on GPU

```
SELECT ?X ?Y ?Z WHERE {
```
**TP1**  `?X teacherof ?Y .`
**TP2**  `?Z takecourse ?Y .`
**TP3**  `?Z adivsor ?X .      }`
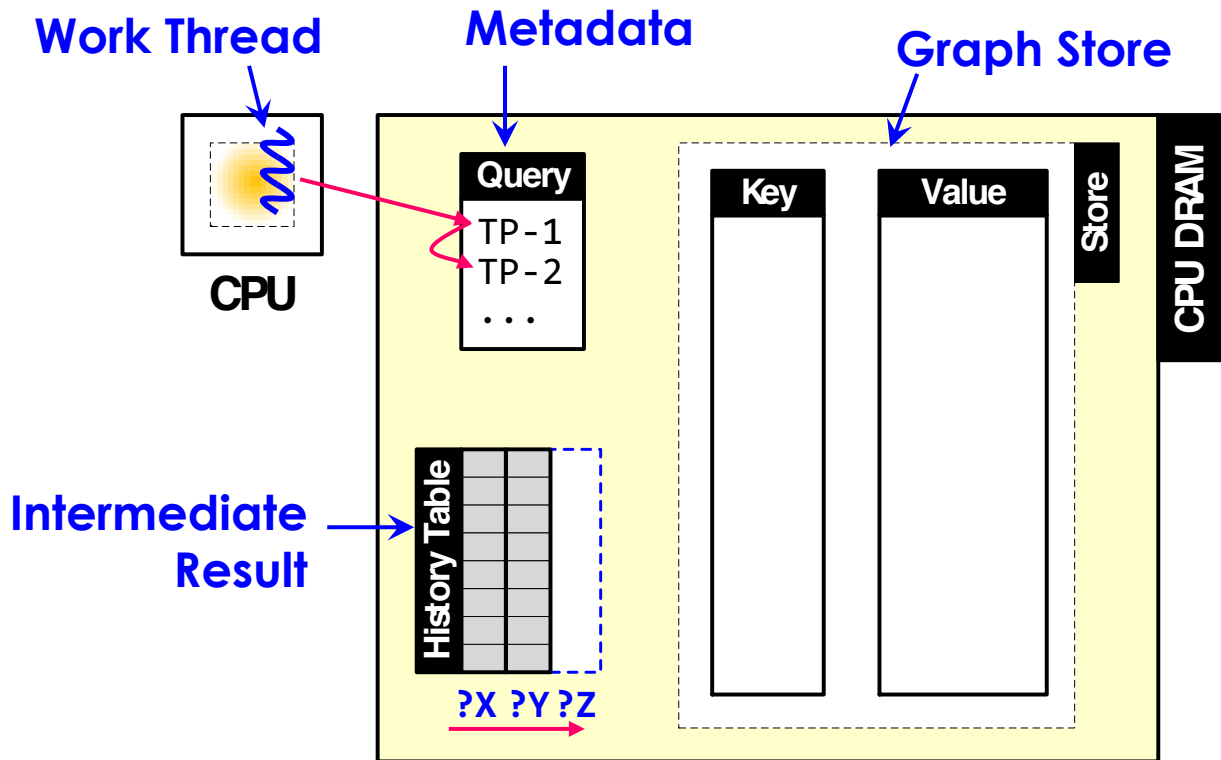
# Query Execution on GPU

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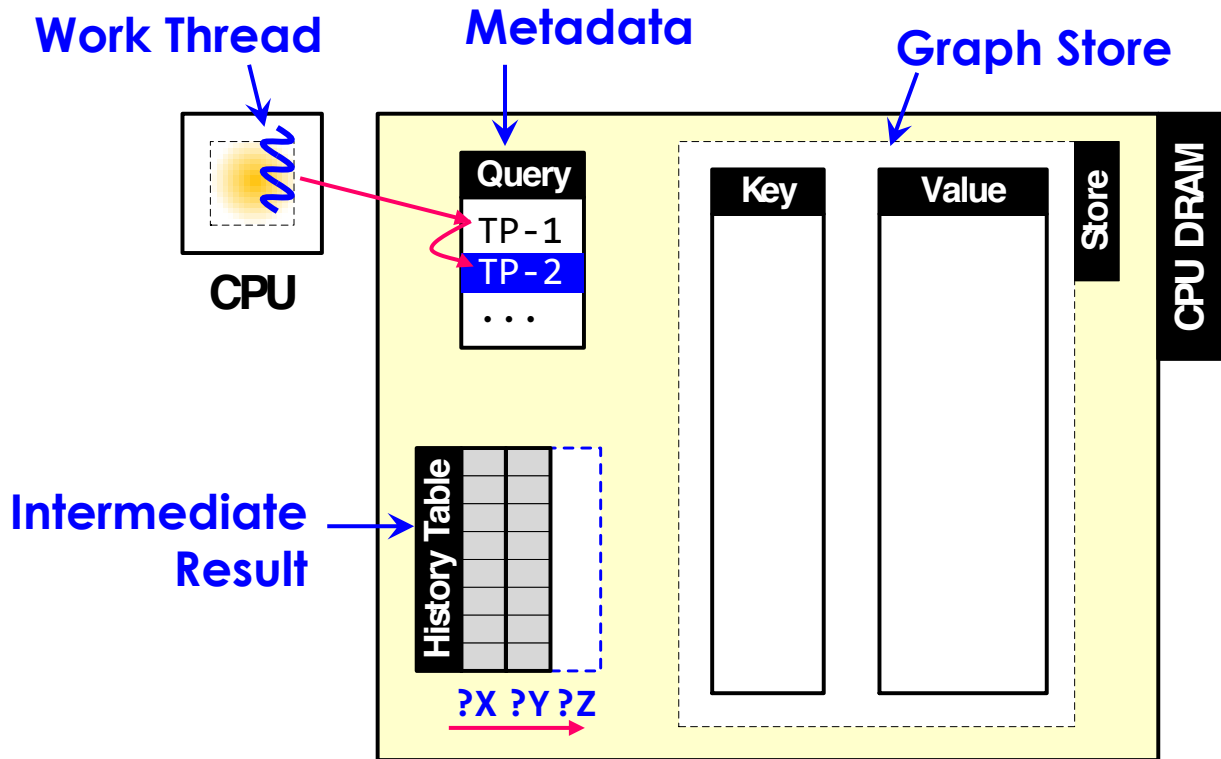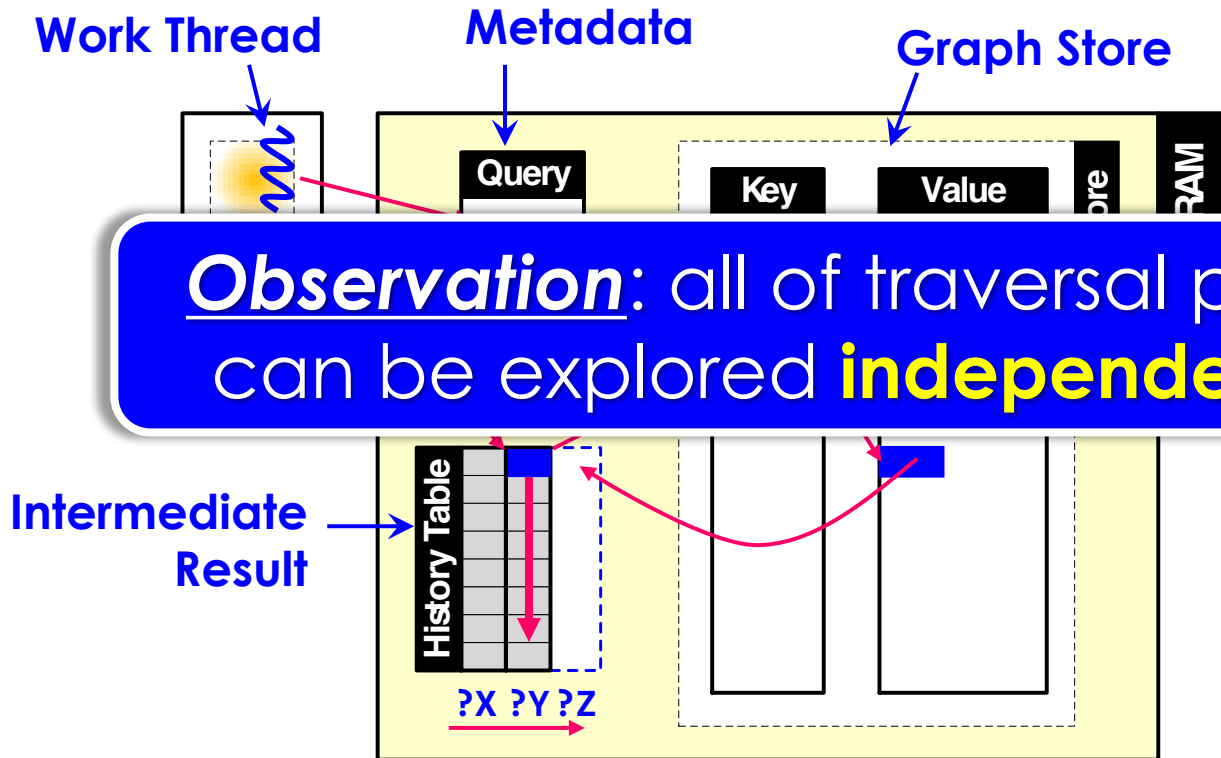
TP1
TP2
TP3



**GPU**

**Work Thread**

**Agent Thread**

**CPU**

History Table

?X ?Y ?Z

Key    Value    Cache

**GPU DRAM**

**Prefecthing**

Query
TP-1
TP-2
...

Key    Value    Store

**CPU DRAM**

27

# Challenges

1. **Small GPU memory**
2. **Limited PCIe bw.**
3. **Long comm. path**

**GPU (2880)**

**CPU (16)**

**68GB/s**

**288GB/s**

**10GB/s**

**10GB/s**

16GB DRAM

256GB DRAM

NETWORK

256GB DRAM

16GB D

Smart data prefetching

GPU-friendly key/value store

Heterogeneous RDMA comm.

Evaluation

# Smart Data Prefetching

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

Time

Entire RDF graph

**X** Out-of-memory

**OB1**: a query only touches **a part** of RDF data

| | Granularity | Footprint | Data Transfer |
|---|---|---|---|
| | Entire graph | 16.3GB | **Failed** |

2560

**OB2**: the **predicate** of TP is commonly **known**

# Smart Data Prefetching

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

**Time**

**Entire RDF graph**

**X** **Out-of-memory**

**Per-query**

**load time**

**compute time**

**Case study**: Q7 on LUBM-2560

| | Memory | Data ... fer |
|---|---|---|
| ...d | | |
| Per-query | **5.6GB** | 5.6GB |

**OB3**: **TPs** of a query will be executed **in sequence**

# Smart Data Prefetching

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

**Time**

**Entire RDF graph**

**X Out-of-memory**

**Per-query**

**load time**

**compute time**

**Per-pattern**

IDX    TP1    TP2    TP3

**Pipeline**

**Case study**: Q7 on LUBM-2560

| Granularity | Memory Footprint | Data Transfer |
|---|---|---|
| Entire graph | 16.3GB | Failed |
| Per-query | 5.6GB | 5.6GB |
| Per-pattern | **2.9GB** | 5.6GB |

32

# Smart Data Prefetching

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



**Entire RDF graph** — ✗ **Out-of-memory**

**Per-query** — load time / compute time

**Per-pattern** — IDX, TP1, TP2, TP3

**Pipeline**

**Per-block**

Time

**Case study**: Q7 on LUBM-2560

| Granularity | Memory Footprint | Data Transfer |
|---|---|---|
| Entire graph | 16.3GB | Failed |
| Per-query | 5.6GB | 5.6GB |
| Per-pattern | 2.9GB | 5.6GB |
| Per-block | 2.9GB | **0.7GB*** |

* evaluated on 6GB GPU memory

33

GPU-enable query processing

GPU-friendly key/value store

Heterogeneous RDMA comm.

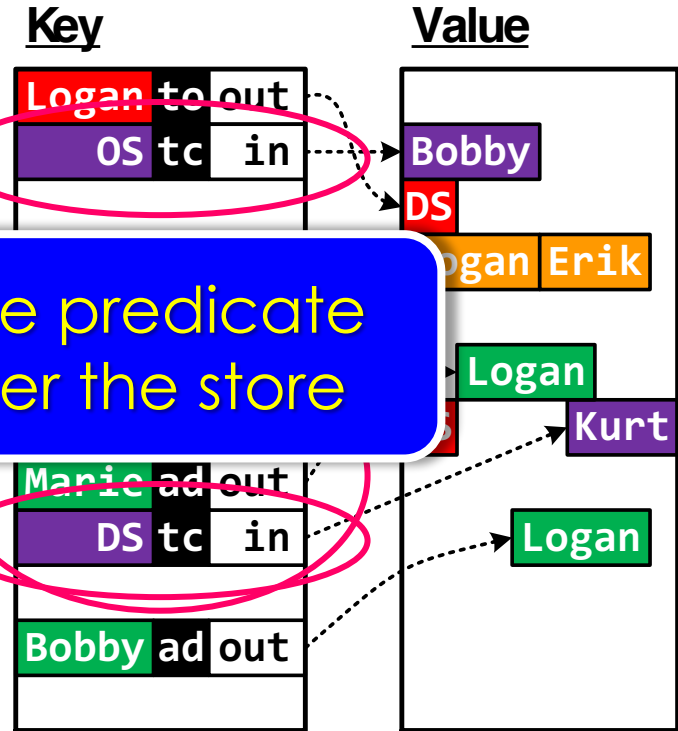Evaluation

# Original RDF Store (Wukong)

▶ **Predicate-based decomposition**

1. Efficient query processing on CPU

   *Hash*

2. Provide possibility to prefetching triples with a certain predicate

**Key**

| Logan | to | out |
| OS | tc | in |

| Maria | ad | out |
| DS | tc | in |

| Bobby | ad | out |

**Value**

Bobby
DS
ogan Erik

Logan
Kurt

Logan

Triples with the same predicate are sprinkled all over the store

# GPU-friendly RDF Store

► **Predicate-based grouping**

SEG_OFFSET( pred, dir )

+ *Hash*( vtx )

% SEG_SZ( pred, dir )



*Segment*

36

# GPU-friendly RDF Store
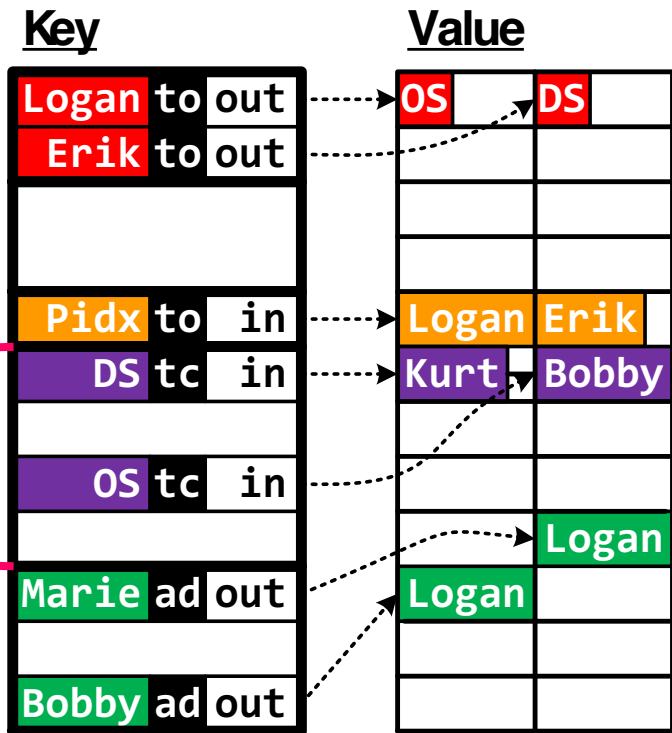
▶ **Predicate-based grouping**

1. *Segment*: prefetching in **batch**
2. *Hashing*: lookup **efficiency**
3. *Occupancy rate of segment*:
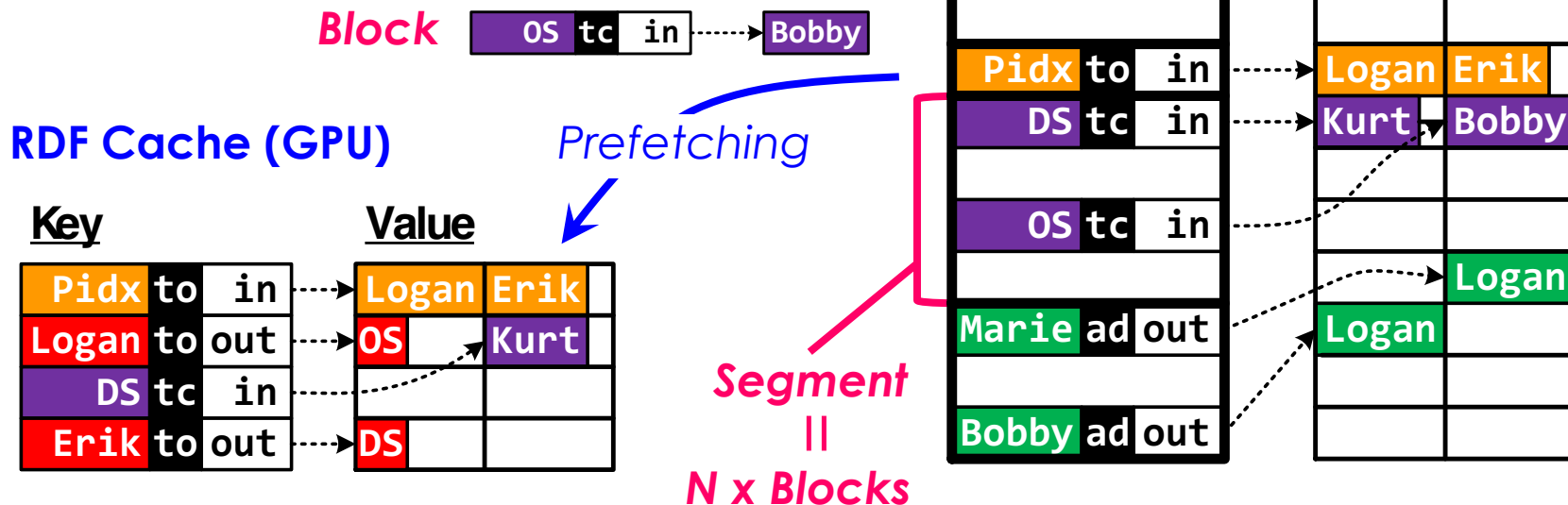
*Prefetching Cost*
vs.
*Lookup Overhead*

Tradeoff



*Segment*

37

# GPU-friendly RDF Store



► Predicate-based grouping
► **Fine-grained swapping**

**RDF Store (CPU)**

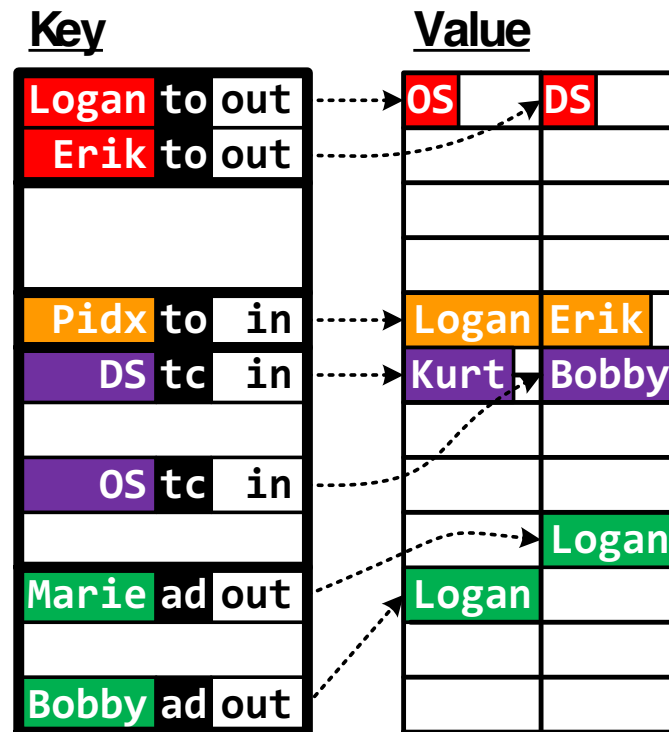**RDF Cache (GPU)**

*Block*

*Prefetching*

*Segment = N x Blocks*

# GPU-friendly RDF Store

- ▶ Predicate-based grouping
- ▶ Fine-grained swapping

- ▶ **Pairwise caching**

- ▶ **Look-ahead replacement**
  (see paper)

**Key**

| | | |
|---|---|---|
| **Logan** | **to** | **out** |
| **Erik** | **to** | **out** |
| | | |
| | | |
| **Pidx** | **to** | **in** |
| **DS** | **tc** | **in** |
| | | |
| **OS** | **tc** | **in** |
| | | |
| **Marie** | **ad** | **out** |
| | | |
| **Bobby** | **ad** | **out** |

**Value**

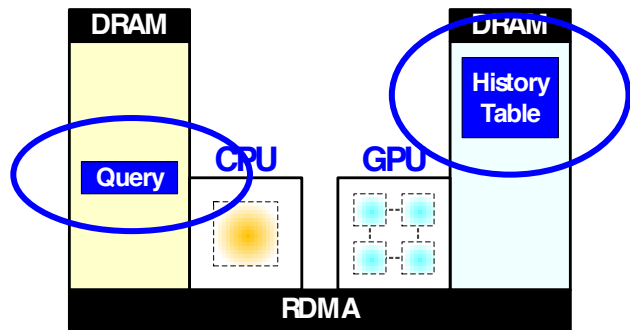| | |
|---|---|
| **OS** | **DS** |
| | |
| | |
| | |
| **Logan** | **Erik** |
| **Kurt** | **Bobby** |
| | |
| | |
| | **Logan** |
| **Logan** | |
| | |
| | |

39

# Agenda

GPU-enable query processing

GPU-friendly key/value store
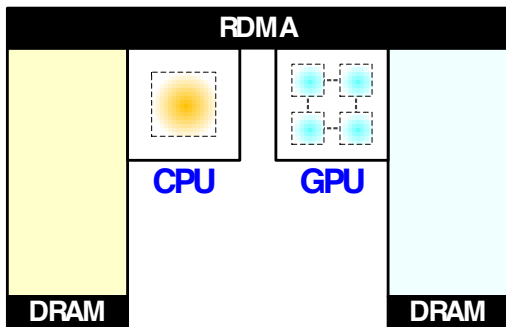
Heterogeneous RDMA comm.

Evaluation
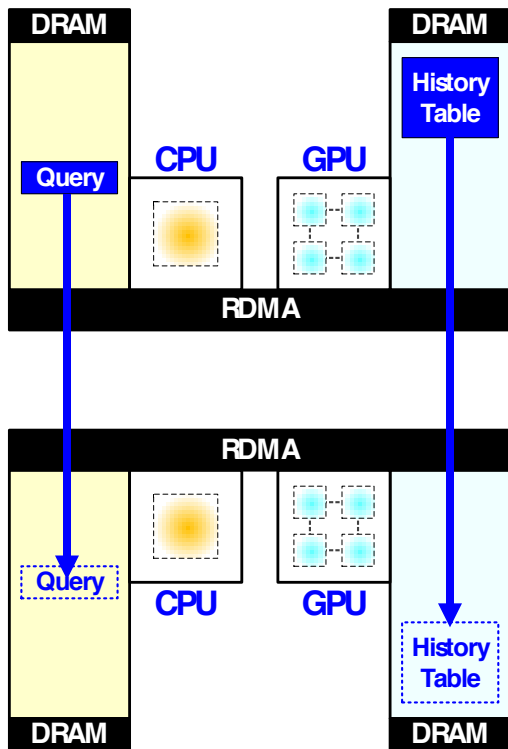
# Heterogeneous RDMA Communication

Metadata: *Query*

Data: *History Table*

# Heterogeneous RDMA Communication



Metadata: *Query*

Data: *History Table*

# Heterogeneous RDMA Communication



*(Native)* **RDMA**

Metadata: *Query*
  ② CPU → CPU (RDMA)


Data: *History Table*
  ① GPU → CPU (PCIe)
  ② CPU → CPU (RDMA)
  ③ CPU → GPU (PCIe)
  ④ GPU → CPU (PCIe)
  ⑤ CPU → CPU (RDMA)

43

# Heterogeneous RDMA Communication



**RDMA** with **GPUDirect**

Metadata: *Query*
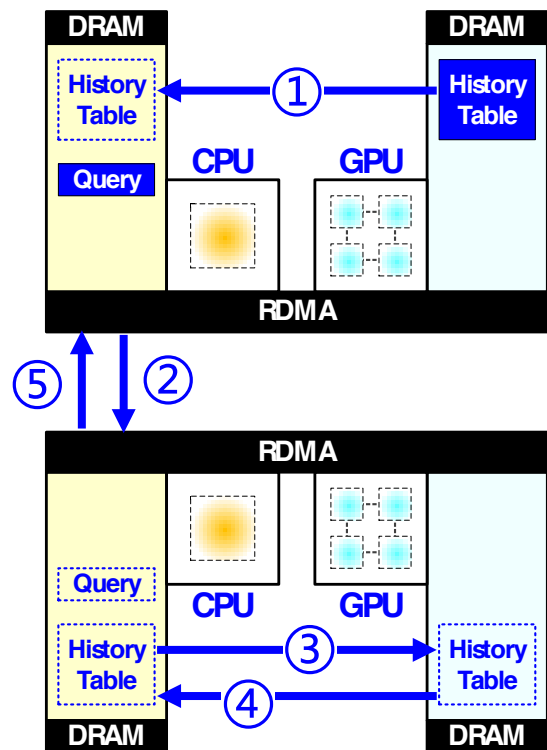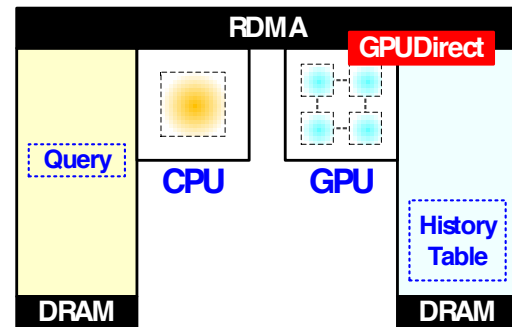
② CPU → CPU (RDMA)

Data: *History Table*

① GPU → CPU (PCIe)

② CPU → CPU (RDMA)

③ CPU → GPU (PCIe)

④ GPU → CPU (PCIe)

⑤ CPU → CPU (RDMA)

# Heterogeneous RDMA Communication



**RDMA** with **GPUDirect**

Metadata: *Query*
① CPU → CPU (RDMA)

Data: *History Table*
① GPU → CPU (PCIe)
① GPU → CPU (RDMA+G)
③ CPU → GPU (PCIe)
④ GPU → CPU (PCIe)
⑤ CPU → CPU (RDMA)

# Heterogeneous RDMA Communication



**RDMA** with **GPUDirect**

Metadata: *Query*
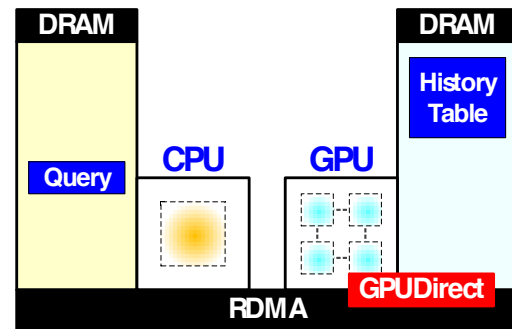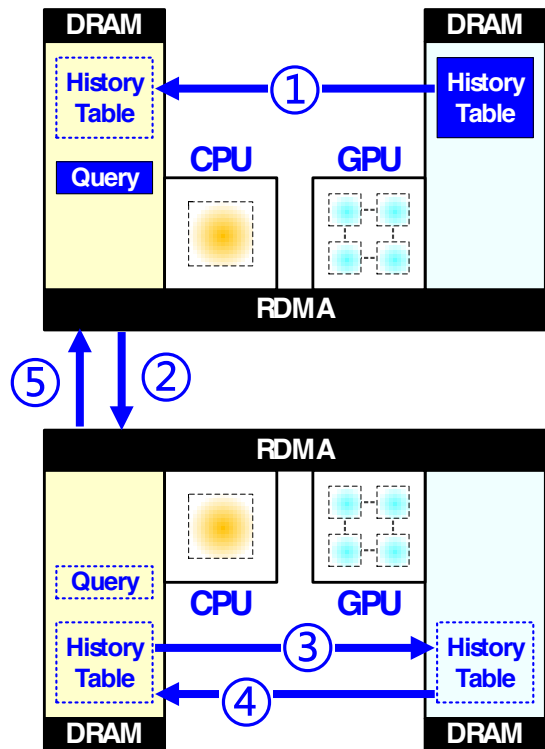
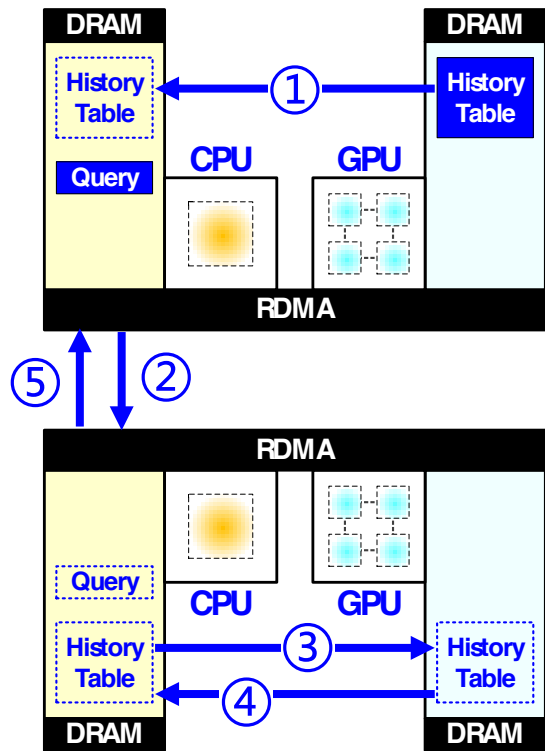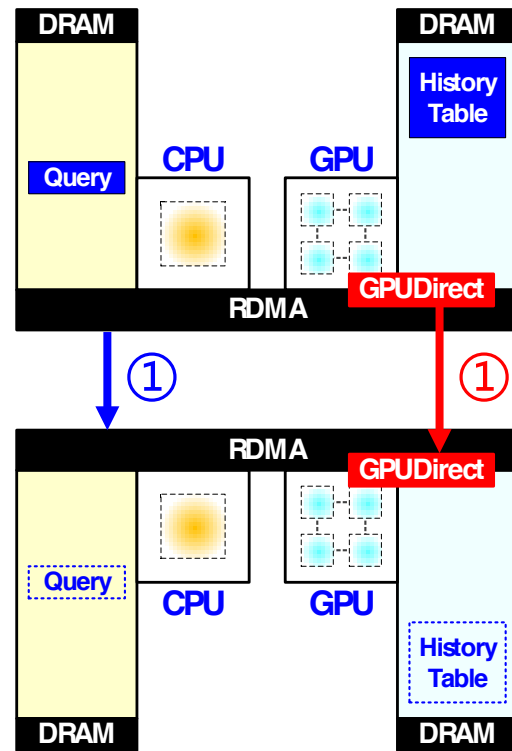① CPU → CPU (RDMA)

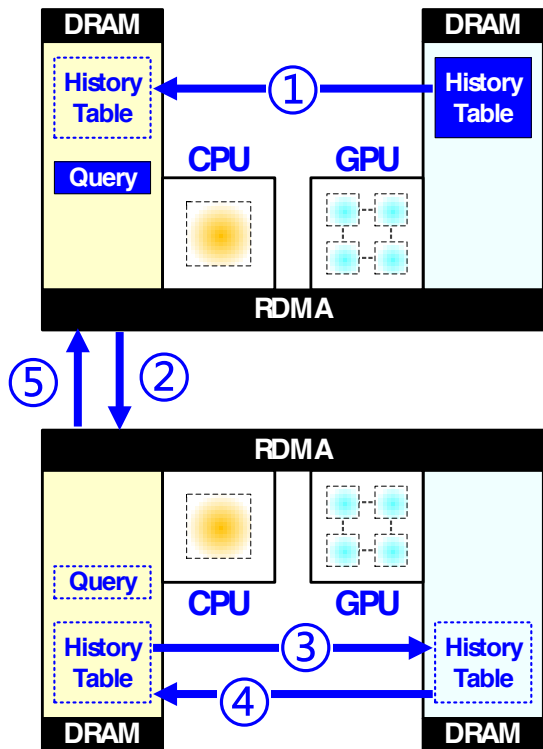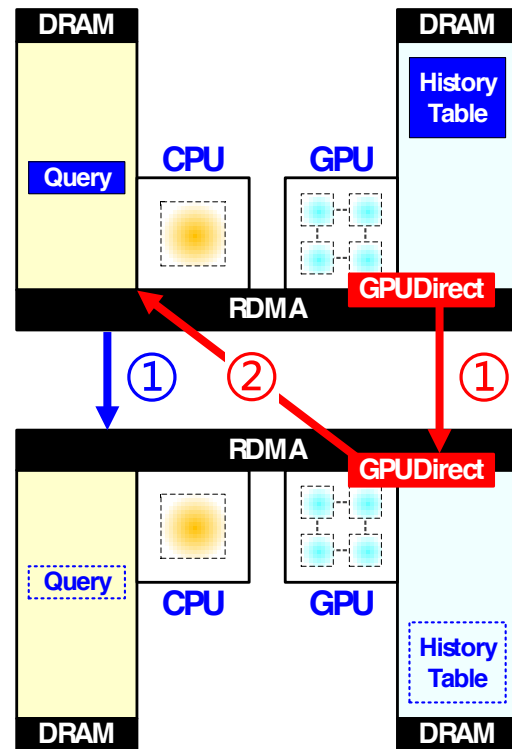Data: *History Table*

① GPU → CPU (PCIe)

① GPU → CPU (RDMA+G)

③ CPU → GPU (PCIe)

④ GPU → CPU (PCIe)

② GPU → CPU (RDMA+G)

GPU-enable query processing

GPU-friendly key/value store

Heterogeneous RDMA comm.

Evaluation

# Evaluation



**Baseline:** state-of-the-art systems

☐ Wukong, TriAD (distributed triple store)

**Platforms: 10** servers on a rack-scale **5**-node cluster

☐ RDMA: Mellanox 56Gbps IB NIC, 40Gbps IB Switch

☐ Two servers run on a single machine

☐ Each server: 12-core Intel Xeon, 128GB DRAM,
    NVIDIA Tesla K40m (2880 cores, 12GB DRAM)

## Benchmarks

☐ Synthetic: LUBM
☐ Real-life: DBPSB, YAGO2

| Dataset | #T | #S | #O | #P | Size† |
|---------|------|------|------|------|------|
| **LUBM-2560** | 352 M | 55 M | 41 M | 17 | 58GB |
| **LUBM-10240** | 1,410 M | 222 M | 165 M | 17 | 230GB |
| **DBPSB** | 15 M | 0.3 M | 5.2 M | 14,128 | 2.8GB |
| **YAGO2** | 190 M | 10.5 M | 54.0 M | 99 | 13GB |

# Single Query Latency (msec)

## Heavy queries (Q1-Q3, Q7)

- ▶ Start from a set of vertices
- ▶ Touch a large part of graph
- ▶ Speedup: **2.3X~9X** vs. Wukong

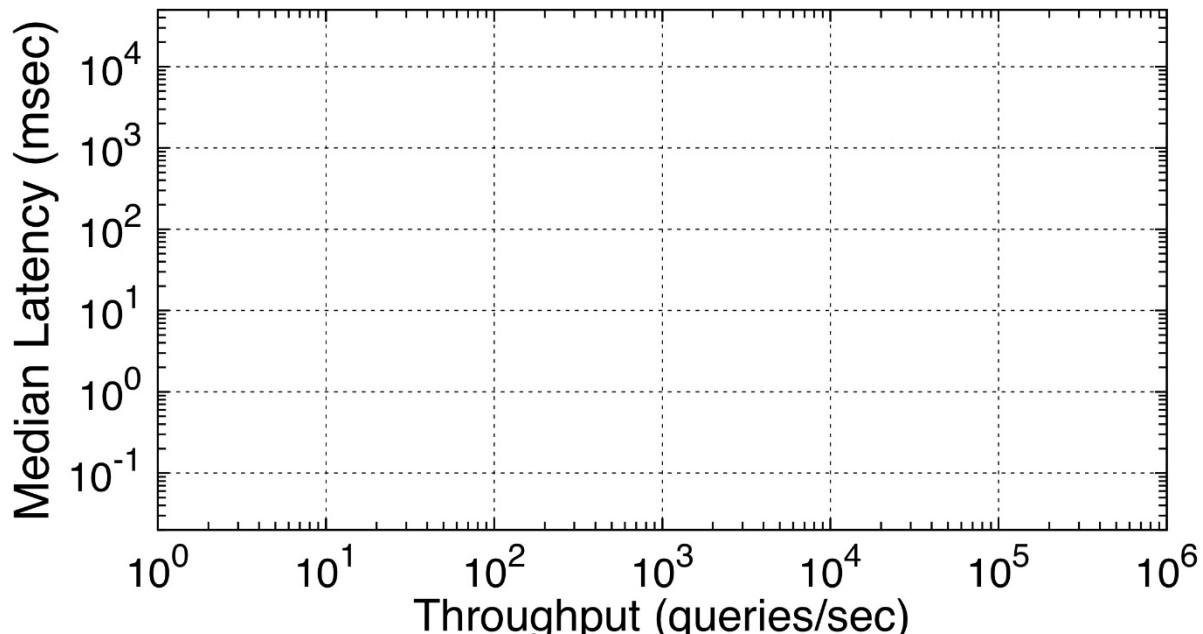## Light queries (Q4-Q6)

- ▶ Start from a given vertex
- ▶ Touch a small part of graph
- ▶ Negligible slowdown

| LUBM-2560 | | Wukong | Wukong+G |
|---|---|---|---|
| **H** | **Q1** (3.6GB) | 992 | **165** |
| | **Q2** (2.4GB) | 138 | **31** |
| | **Q3** (3.6GB) | 340 | **63** |
| | **Q7** (5.6GB) | 828 | **100** |
| | Geo. M | 443 | **75** |
| **L** | **Q4** | **0.13** | 0.16 |
| | **Q5** | **0.09** | 0.11 |
| | **Q6** | **0.49** | 0.51 |
| | Geo. M | **0.18** | 0.21 |

**single server**

| LUBM-10240 | | Wukong | Wukong+G |
|---|---|---|---|
| **H** | **Q1** (14.25GB) | 480 | **211** |
| | **Q2** (9.74GB) | 66 | **12** |
| | **Q3** (14.25GB) | 171 | **19** |
| | **Q7** (22.58GB) | 390 | **100** |
| | Geo. M | 215 | **47** |
| **L** | **Q4** | **0.44** | 0.46 |
| | **Q5** | **0.13** | 0.17 |
| | **Q6** | **0.70** | 0.71 |
| | Geo. M | **0.34** | 0.38 |

**10-server cluster**

49

# Performance of Hybrid Workloads



**WKD** (default)
Heavy/Light: ALL of CPUs **(10)**

**WKI** (Isolation)
Heavy: HALF of CPUs **(5)**
Light: HALF of CPUs **(5)**

**WKG** (+G)
Heavy: CPU **(1)** + GPUs **(1)**
Light: **REST of** CPU **(9)**

<u>Workload</u>: 6 classes of light queries + 4 classes of heavy queries

50

# Performance of Hybrid Workloads



**WKD** (default)
Heavy/Light: ALL of CPUs **(10)**

**WKI** (Isolation)
Heavy: HALF of CPUs **(5)**
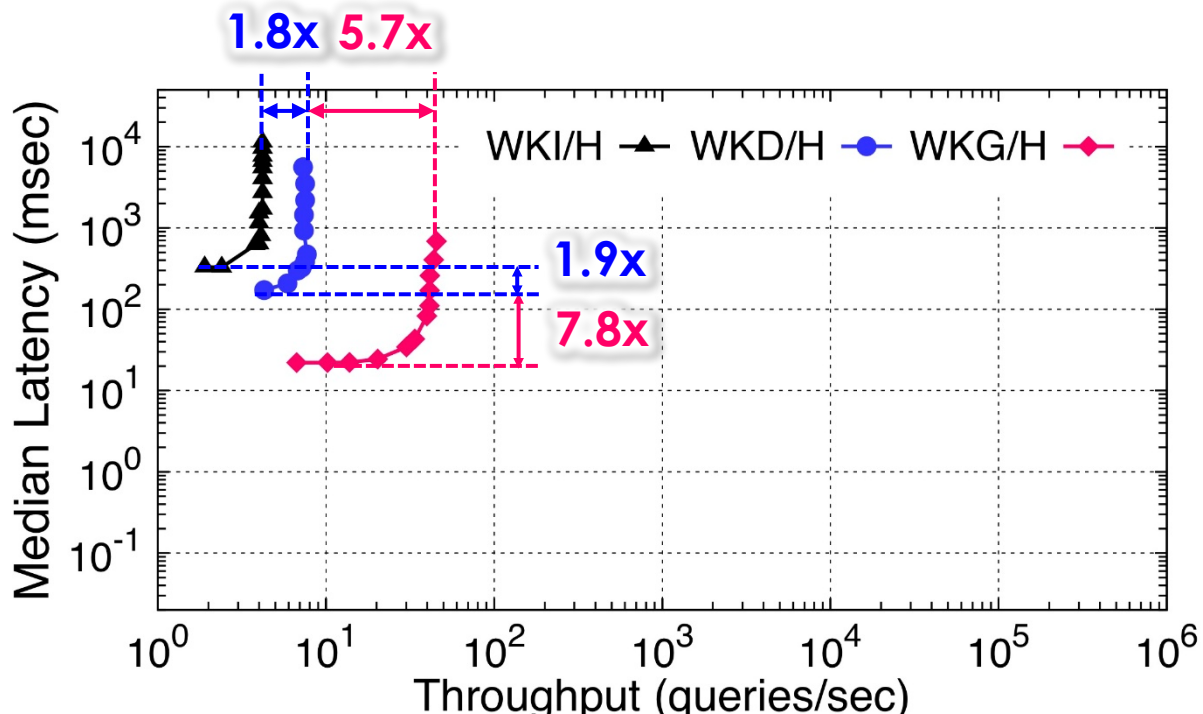Light: HALF of CPUs **(5)**

**WKG** (+G)
Heavy: CPU **(1)** + GPUs **(1)**
Light: **REST of** CPU **(9)**

<u>Workload</u>: 6 classes of light queries + 4 classes of heavy queries

# Performance of Hybrid Workloads



**WKD** (default)
Heavy/Light: ALL of CPUs **(10)**

**WKI** (Isolation)
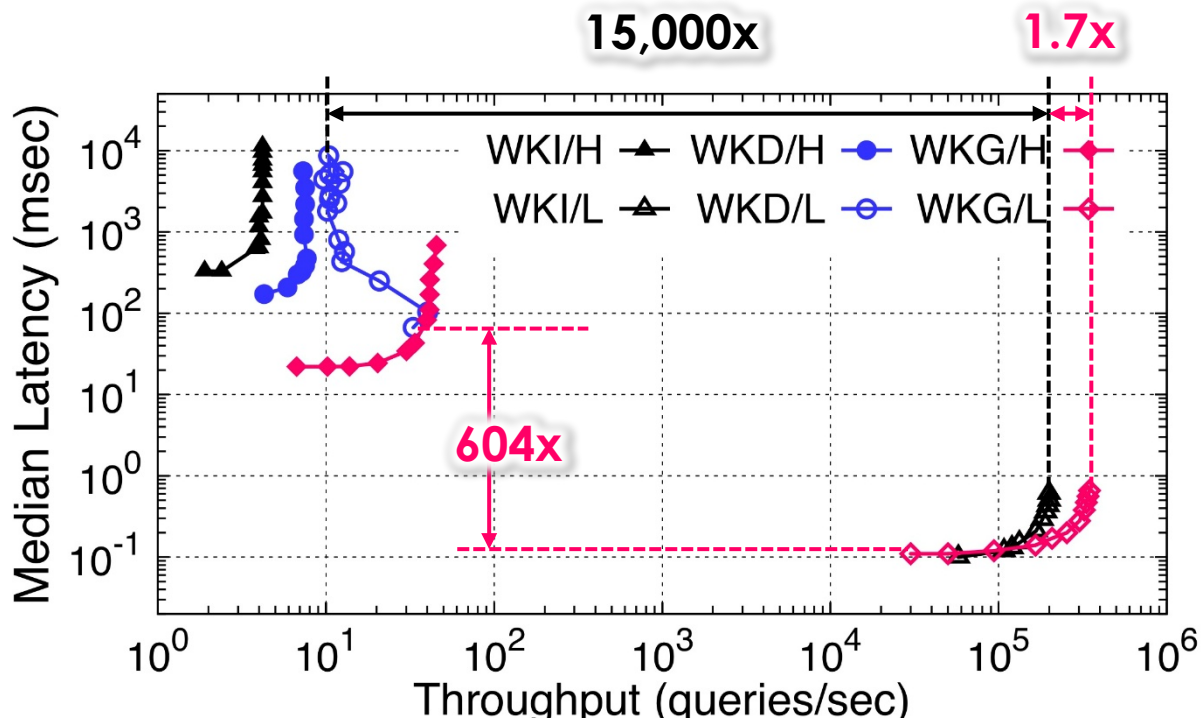Heavy: HALF of CPUs **(5)**
Light: HALF of CPUs **(5)**

**WKG** (+G)
Heavy: CPU **(1)** + GPUs **(1)**
Light: **REST of** CPU **(9)**

Workload: 6 classes of light queries + 4 classes of heavy queries

# Conclusion

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
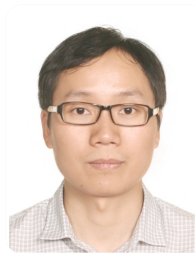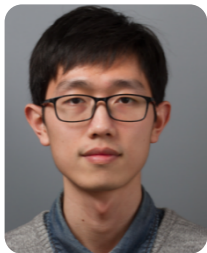one order of magnitude when facing hybrid workloads

**Website**: http://ipads.se.sjtu.edu.cn/projects/wukong

**GitHub**: https://github.com/SJTU-IPADS/wukong

# Thanks

## Wukong+G

**GitHub**: https://github.com/SJTU-IPADS/wukong

Institute of Parallel and Distributed Systems
Shanghai Jiao Tong University

## Questions

# Distinguish Heavy & Light Queries

**Query plan optimizer**

► **Query plan**: the order of triple patterns

► Using a **cost-model** to estimate the execution time of different plans for a given query

► For SPARQL query, cost model is roughly based on **#paths** may be explored

Wukong+G uses a user-defined threshold for #paths to distinguish heavy and light queries

# GPU Memory Size Limitation

**Too large predicate segment**

1. Load one part of segment to GPU memory
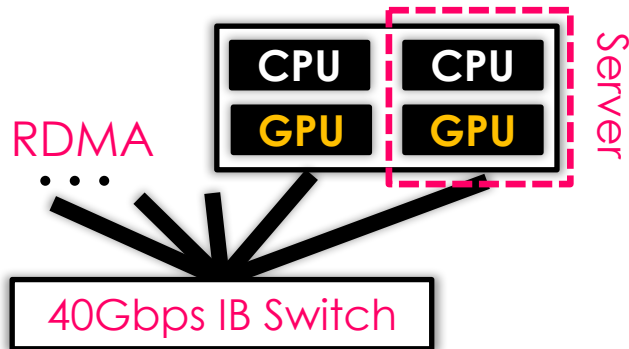2. Do traversal work

Repeat 1 and 2

**Too large intermediate results**

1. Load one part of history table to GPU memory
2. Do traversal work

Repeat 1 and 2

# Multi-GPUs Support

▶ Run a separate server for each GPU card and several co-located CPU cores (usually a socket)

▶ All servers comply with the same communication mechanism via GPUDirect RDMA operations

# Graph Analytics vs. Graph Query

|  | **Graph Analytics** | **Graph Query** |
|---|---|---|
| Graph Model | Property Graph | Semantic (RDF) Graph |
| Working Set | A whole Graph | A small frac. of Graph |
| Processing | Batched & Iterative | Concurrent |
| Metrics | Latency | Latency & Throughput |

# Factor Analysis of Improvement

► Single Server w/ 3GB GPU memory

► LUBM-2560

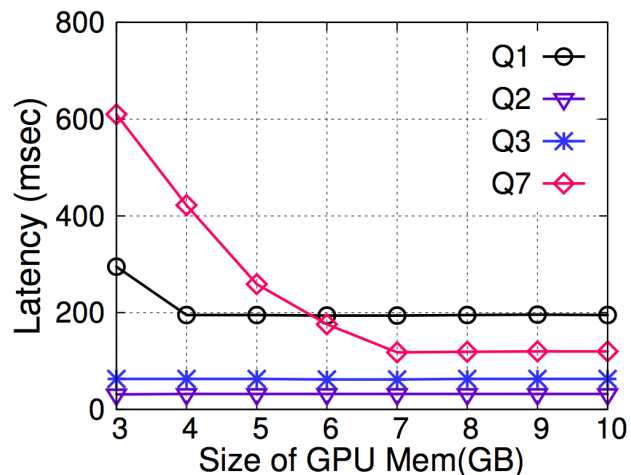| LUBM-2560 | Per-query | Per-parttern | Per-block | Pipeline |
|---|---|---|---|---|
| **Q1** (3.6GB) | x | 743 | 313 | 295 |
| **Q2** (2.4GB) | 284 | 283 | 32 | 31 |
| **Q3** (3.6GB) | x | 309 | 62 | 63 |
| **Q7** (5.6GB) | x | 893 | 622 | 610 |

# RDF Cache of GPU

► LUBM-2560



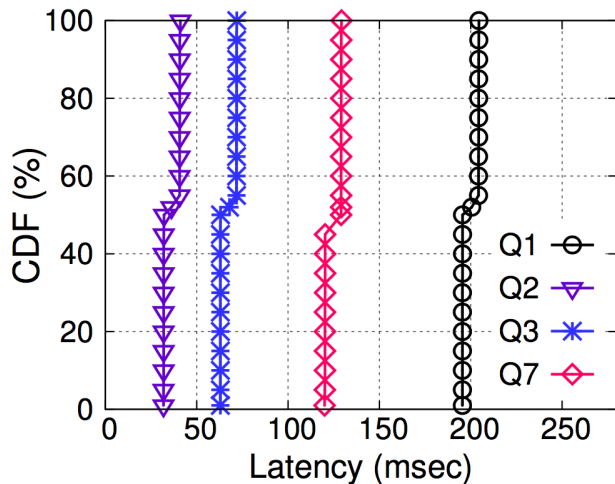Fig. 11: The latency with the increase of GPU memory.

► LUBM-2560

► 10GB GPU Memory



Fig. 12: The CDF of latency for mixed heavy workload.