

# How and Why is An Answer (Still) Correct? Maintaining Provenance in Dynamic Knowledge Graphs

---

Garima Gaur<sup>¶,a</sup>, Arnab Bhattacharya<sup>¶</sup>, Srikanta Bedathur<sup>†</sup>  
garimag@cse.iitk.ac.in, arnabb@cse.iitk.ac.in, srikanta@cse.iitd.ac.in

<sup>¶</sup> Indian Institute of Technology, Kanpur, India

<sup>†</sup> Indian Institute of Technology, Delhi, India

<sup>a</sup> Thanks SIGIR for covering  
conference registration cost



- Knowledge Graph (KG): collection of *facts*



- Fact extractors extracting information from various sources

---

<sup>1</sup><https://en.wikipedia.org/wiki/Wikipedia:Statistics>

- Knowledge Graph (KG): collection of *facts*



- Fact extractors extracting information from various sources
- Dynamic KGs
  - NELL is continuously at work since 2010
  - 1.9 Wikipedia edits/second <sup>1</sup>

---

<sup>1</sup><https://en.wikipedia.org/wiki/Wikipedia:Statistics>

# Dynamic data $\implies$ Evolving answer

- List democrats who are running for US president 2020

**On 30th Jan 2019**

---

Andrew Yang  
Tulsi Gabbard  
John Delaney  
Julián Castr  
Kamala Harris

---

**On 28th Feb 2019**

---

Andrew Yang  
:  
Elizabeth Warren  
Amy Klobuchar  
Bernie Sanders

---

**On 30th March 2020**

---

Bernie Sanders  
Joe Biden

---

---

<sup>2</sup>“Provenance Semirings”, PODS, 2007

## Dynamic data $\implies$ Evolving answer

- List democrats who are running for US president 2020

<u>On 30th Jan 2019</u>	<u>On 28th Feb 2019</u>	<u>On 30th March 2020</u>
Andrew Yang	Andrew Yang	
Tulsi Gabbard	⋮	
John Delaney	Elizabeth Warren	Bernie Sanders
Julián Castr	Amy Klobuchar	Joe Biden
Kamala Harris	Bernie Sanders	

- Important to propagate changes in *facts* down to the precomputed (“standing”) queries

---

<sup>2</sup>“Provenance Semirings”, PODS, 2007

## Dynamic data $\implies$ Evolving answer

- List democrats who are running for US president 2020

On 30th Jan 2019

Andrew Yang  
Tulsi Gabbard  
John Delaney  
Julián Castr  
Kamala Harris

On 28th Feb 2019

Andrew Yang  
⋮  
Elizabeth Warren  
Amy Klobuchar  
Bernie Sanders

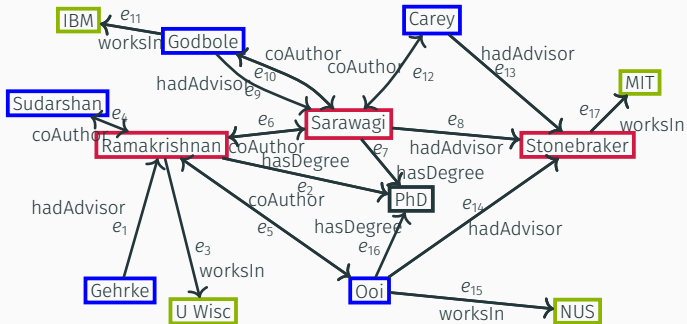
On 30th March 2020

Bernie Sanders  
Joe Biden

- Important to propagate changes in *facts* down to the precomputed (“standing”) queries
- Need a mechanism to keep track of extraction process and the source of information
- How provenance<sup>2</sup> captures how a query answer is generated
- Encode provenance as a polynomial – monomial corresponds to derivation

---

<sup>2</sup>“Provenance Semirings”, PODS, 2007



**Figure 1:** KG encoding information about phd students, their advisors and collaborators

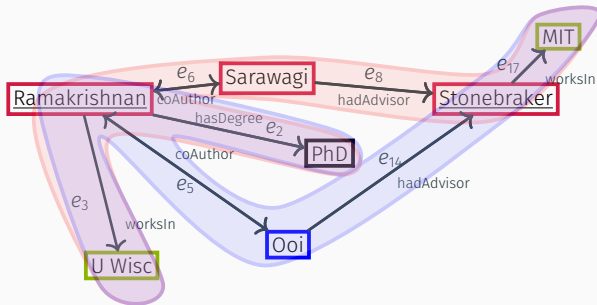
## Provenance Polynomial

- *Find pairs of advisors and collaborators of their students such that the collaborator has a PhD and works in an institute*



# Provenance Polynomial

- Find pairs of advisors and collaborators of their students such that the collaborator has a PhD and works in an institute



- Two derivations of answer  $\langle \text{Stonebraker}, \text{Ramakrishnan} \rangle$  :
  - Red-colored subgraph :  $\{e_2, e_3, e_6, e_8, e_{17}\}$
  - Blue-colored subgraph :  $\{e_2, e_3, e_5, e_{14}, e_{17}\}$
- Resultant polynomial is  $e_2.e_3.e_6.e_8.e_{17} + e_2.e_3.e_5.e_{14}.e_{17}$

## Query result maintenance under edge update

Given a knowledge graph  $G(V, E)$  and a set of standing queries  $Q = \{Q_1, Q_2, \dots, Q_n\}$ , maintain result along with their provenance of a subset  $Q' \subseteq Q$  such that  $Q_i, \forall Q_i \in Q'$ , gets affected on the deletion or insertion of an edge  $e_d, e_d \in E$

Query re-computation is impractical due to KG size!

- Framework **HUKA** which incrementally maintains the query result and its provenance under edge insertion or deletion.

HUKA – maintaining How  
provenance under Updates to  
Knowledge grAph

---

## Handling Edge Insertion: Primary Idea

Shifting focus from exact matches of a query pattern to its partial matches

# Handling Edge Insertion: Primary Idea

Shifting focus from exact matches of a query pattern to its partial matches

## Potential Match (PM)

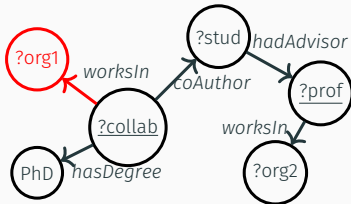
Any subgraph  $S$  of the knowledge graph  $G$  which can become an exact match of a query  $Q$  after a *single* edge insertion is called a **potential match**.

# Handling Edge Insertion: Primary Idea

Shifting focus from exact matches of a query pattern to its partial matches

## Potential Match (PM)

Any subgraph  $S$  of the knowledge graph  $G$  which can become an exact match of a query  $Q$  after a *single* edge insertion is called a **potential match**.

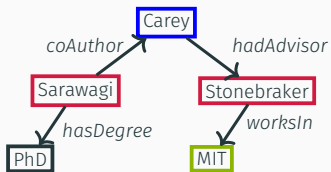
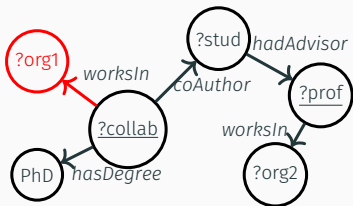


# Handling Edge Insertion: Primary Idea

Shifting focus from exact matches of a query pattern to its partial matches

## Potential Match (PM)

Any subgraph  $S$  of the knowledge graph  $G$  which can become an exact match of a query  $Q$  after a *single* edge insertion is called a **potential match**.



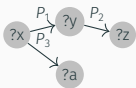
- Incremental insertion handling approach – *one* edge at a time
- Addressing three sub-problems:
  1. Pre-compute potential matches (PM) of each query
  2. After insertion, efficiently identify *transformed* PM
  3. Maintain PM to ensure correctness while handling subsequent updates



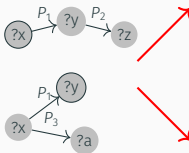
- Incremental insertion handling approach – *one* edge at a time
- Addressing three sub-problems:
  1. Pre-compute potential matches (PM) of each query
  2. After insertion, efficiently identify *transformed* PM
  3. Maintain PM to ensure correctness while handling subsequent updates
- HUKA operates in two phases
  - **Query Registration**
  - **Update Processing**

# Query Registration

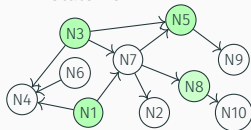
0. Submit standing query



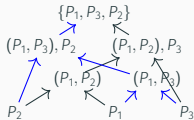
1. Subquery Construction



2 Annotate KG

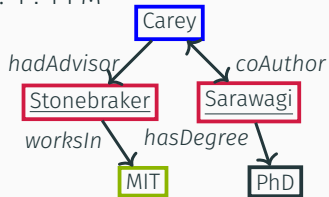


3 Execution Plan

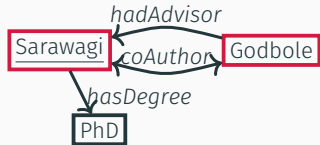


# Task 1: PM computation

$G_1$ : 1 : 1 PM



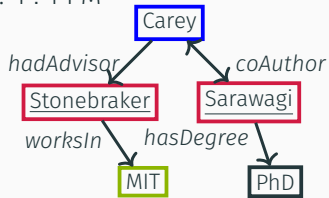
$G_2$ : 1 : M PM



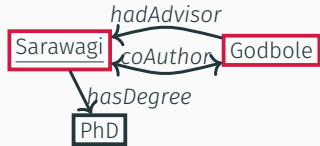
- Insert  $\langle \text{Sarawagi}, \text{worksIn}, \text{IITB} \rangle$ :  $G_1$  and  $G_2$  becomes exact match

# Task 1: PM computation

$G_1$ : 1 : 1 PM



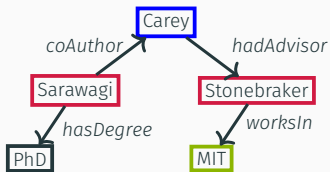
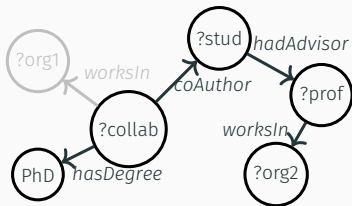
$G_2$ : 1 : M PM



- Insert  $\langle \text{Sarawagi}, \text{worksIn}, \text{IITB} \rangle$ :  $G_1$  and  $G_2$  becomes exact match
- Unmatched triple patterns:
  - $G_1$ :  $\langle ?collab, \text{worksIn}, ?org1 \rangle$
  - $G_2$ :  $\langle ?collab, \text{worksIn}, ?org1 \rangle$  and  $\langle ?prof, \text{worksIn}, ?org2 \rangle$
- Types of potential matches:
  - 1 : 1 PM: New edge matches to *single* triple pattern
  - 1 : M PM: New edge satisfies *multiple* triple constraints

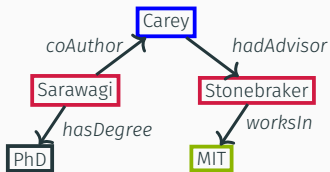
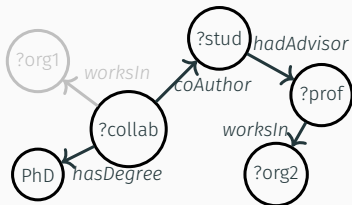
# Task1: PM Computation

- 1 : 1 PM (**pre-computed**): Satisfies subqueries with one less triple pattern



# Task1: PM Computation

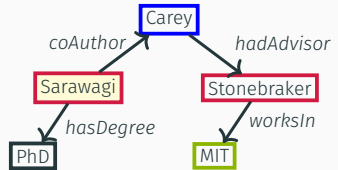
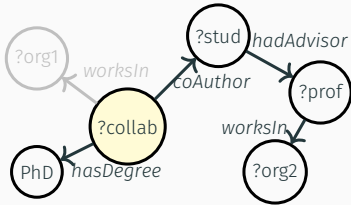
- 1 : 1 PM (**pre-computed**): Satisfies subqueries with one less triple pattern



- 1 : M PM (**lazily computed**): On appropriate (*expected*) edge insertion,
  - If new edge satisfies all the unmatched triple patterns
  - PM directly becomes an exact match
  - An exact match also a partial match – satisfies all subqueries

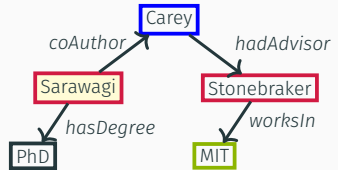
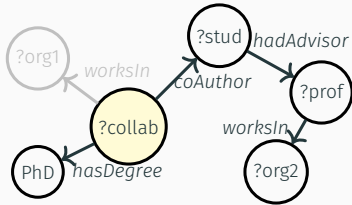
## Task 2: KG annotation

- Efficiently check if the new edge has converted a PM to an exact match
- **Connection points**: PM node expecting an edge



## Task 2: KG annotation

- Efficiently check if the new edge has converted a PM to an exact match
- **Connection points**: PM node expecting an edge



- Annotate all the connection points – avoids materializing subquery results
- Annotation – expected edge and provenance polynomial of corresponding PM



## Task 3: PM maintenance

- **Local Plan:** For *each* subquery
  - AND-OR tree<sup>3</sup> – all possible execution plans
  - Best plan selection based on graph data specific cardinality estimator<sup>4</sup>
  - Collects node signatures – *characteristic set* (CS)

$$CS(u) = \{P \mid \langle u, P, v \rangle\}$$

- Cardinality estimation based on the frequency of a CS

---

<sup>3</sup>“Materialized View Selection and Maintenance Using Multi-query Optimization”, SIGMOD, 2001

<sup>4</sup>“Characteristic sets: Accurate cardinality estimation for RDF queries with multiple joins”, ICDE, 2011

## Task 3: PM maintenance

- **Local Plan:** For *each* subquery
  - AND-OR tree<sup>3</sup> – all possible execution plans
  - Best plan selection based on graph data specific cardinality estimator<sup>4</sup>
  - Collects node signatures – *characteristic set* (CS)

$$CS(u) = \{P \mid \langle u, P, v \rangle\}$$

- Cardinality estimation based on the frequency of a CS
- **Global Plan:** For *all* subqueries
  - Merging best local plans of all subqueries of standing queries
  - Promotes re-usability – share intermediate expression computation

---

<sup>3</sup>“Materialized View Selection and Maintenance Using Multi-query Optimization”, SIGMOD, 2001

<sup>4</sup>“Characteristic sets: Accurate cardinality estimation for RDF queries with multiple joins”, ICDE, 2011

# Local Plan Construction

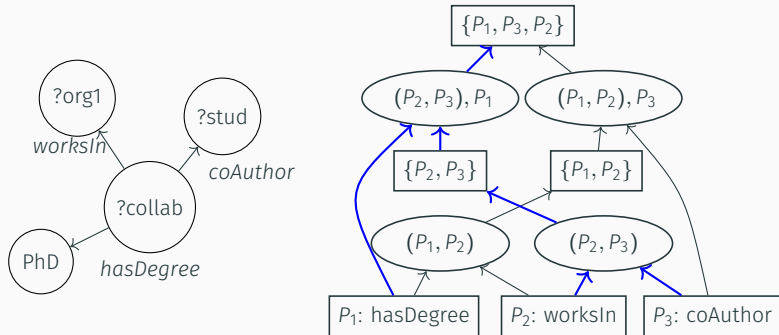


Figure 2: Subquery and its AND-OR tree (Boxes  $\equiv$  OR; Ellipses  $\equiv$  AND)

- Greedily choose the best plan – traversing bottom-up

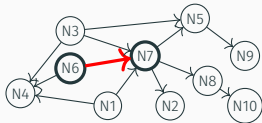
$$\sum_{\{P_1, P_2\} \subset CS_i} \text{Freq}(CS_i) > \sum_{\{P_2, P_3\} \subset CS_i} \text{Freq}(CS_i)$$

- Global plan is a combination of best local plans

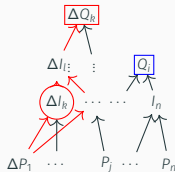
# Update Processing

- Insert:  $\langle N6, P_1, N7 \rangle$

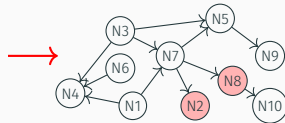
1. Examine incident vertices



2. Find new PM



3.1 Annotate new CP



- HUKA also supports result maintenance under *fact* deletion
- Inverted indexes to support deletion and insertion together

# Experimental Results

---

- Statistics of datasets

Dataset	Vertices	Edges	Predicates	Queries	Avg. Query Size	Subqueries
YAGO2	8.8M	23M	78	4	6.25	26
DBpedia	32M	117M	53K	215	3.90	879

- **Query Set:**

- YAGO2: Benchmark queries used to evaluate RDF-3X<sup>5</sup>;
- DBpedia: real world queries over DBpedia available from the USEWOD 2014.

- **Workload Configuration:** Randomly generated with controlled ratio of deletion to insertion operations.

---

<sup>5</sup>“The RDF-3X engine for scalable management of RDF data”, VLDB, 2010

# Efficiency Comparison

- Baselines against HUKA

Dataset	HUKA <sup>6</sup>	GProM <sup>7</sup>	ProvSQL <sup>8</sup>	Neo4j
YAGO2	0.119 s	25.121 s	75.657 s	5.709 s
DBpedia	1.252 s	5.217 s	6.870 s	99.318 s

- Varying workload impact

Dataset	Deletion -Heavy	Deletion -Moderate	Balanced	Insertion -Moderate	Insertion -Heavy
YAGO2	0.062 s	0.091 s	0.126 s	0.146 s	0.169 s
DBpedia	0.943 s	1.056 s	1.315 s	1.403 s	1.475 s

<sup>6</sup>Code available at <https://github.com/gaugarima/HUKA>

<sup>7</sup>"GProM-a swiss army knife for your provenance needs", IEEE Data Engineering Bulletin, 2018

<sup>8</sup>'ProvSQL: provenance and probability management in postgresQL', VLDB, 2018

- First provenance-aware query result maintenance solution
- **HUKA** – an end-to-end framework to support maintenance of query result and its how provenance
- Seamlessly handles both insertion and deletion update operations



Thank you!  
Questions?