



# **LDBC Social Network Benchmark and Graphalytics**

Gábor Szárnyas  
(CWI Amsterdam, LDBC)

16th LDBC TUC meeting | 2023-06-24 | Seattle

# The Linked Data Benchmark Council (LDBC): Driving competition and collaboration in the graph data management space

Gábor Szárnyas<sup>1\*</sup>, Brad Bebee<sup>2</sup>, Altan Birler<sup>3</sup>, Alin Deutsch<sup>4,5</sup>, George Fletcher<sup>6</sup>, Henry A. Gabb<sup>7</sup>, Denise Gosnell<sup>2</sup>, Alastair Green<sup>8</sup>, Zhihui Guo<sup>9</sup>, Keith W. Hare<sup>8</sup>, Jan Hidders<sup>10</sup>, Alexandru Iosup<sup>11</sup>, Atanas Kiryakov<sup>12</sup>, Tomas Kovatchev<sup>12</sup>, Xincheng Li<sup>13</sup>, Leonid Libkin<sup>14</sup>, Heng Lin<sup>9</sup>, Xiaojian Luo<sup>15</sup>, Arnau Prat-Pérez<sup>16</sup>, David Püroja<sup>1</sup>, Shipeng Qi<sup>9</sup>, Oskar van Rest<sup>17</sup>, Benjamin A. Steer<sup>18</sup>, Dávid Szakállas<sup>19</sup>, Bing Tong<sup>20</sup>, Jack Waudby<sup>21</sup>, Mingxi Wu<sup>5</sup>, Bin Yang<sup>13</sup>, Wenyuan Yu<sup>15</sup>, Chen Zhang<sup>20</sup>, Jason Zhang<sup>13</sup>, Yan Zhou<sup>20</sup>, and Peter Boncz<sup>1</sup>

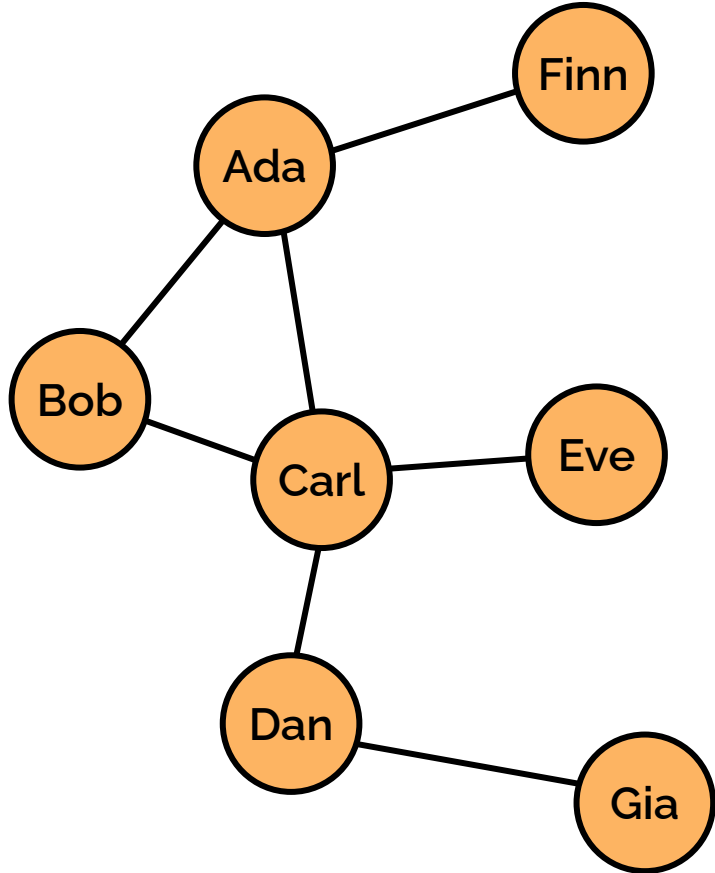
<sup>1</sup> CWI, the Netherlands, <sup>2</sup> Amazon Web Services, <sup>3</sup> Technische Universität München, Germany, <sup>4</sup> UC San Diego, <sup>5</sup> TigerGraph, <sup>6</sup> TU Eindhoven, <sup>7</sup> Intel Corporation, <sup>8</sup> JCC Consulting, <sup>9</sup> Ant Group, <sup>10</sup> Birkbeck, University of London, <sup>11</sup> VU Amsterdam, the Netherlands, <sup>12</sup> Ontotext AD, <sup>13</sup> Ultipa, <sup>14</sup> University of Edinburgh; RelationalAI; ENS, PSL University, <sup>15</sup> Alibaba Damo Academy, <sup>16</sup> *work done while at UPC Barcelona and Sparsity*, <sup>17</sup> Oracle, USA, <sup>18</sup> Pometry Ltd., <sup>19</sup> *individual contributor*, <sup>20</sup> CreateLink, <sup>21</sup> Newcastle University, School of Computing

\* Corresponding author, [gabor.szarnyas@ldbncouncil.org](mailto:gabor.szarnyas@ldbncouncil.org)

**Abstract.** Graph data management is instrumental for several use cases such as recommendation, root cause analysis, financial fraud detection, and enterprise knowledge representation. Efficiently supporting these use cases yields a number of unique requirements, including the need for a concise query language and graph-aware query optimization techniques. The goal of the Linked Data Benchmark Council (LDBC) is to design a set of standard benchmarks that capture representative categories of graph data management problems, making the performance of systems comparable and facilitating competition among vendors. LDBC also conducts research on graph schemas and graph query languages. This paper introduces the LDBC organization and its work over the last decade.

# **LDDB Graphalytics**





The Graphalytics data sets consist of **untyped, unattributed graphs**, which are *either directed or undirected* and *optionally have edge weights*

|                  |
|------------------|
| LDBC SNB Datagen |
| Graph500         |
| Twitter          |
| Friendster       |
| Patents          |
| wiki-Talk        |
| ...              |

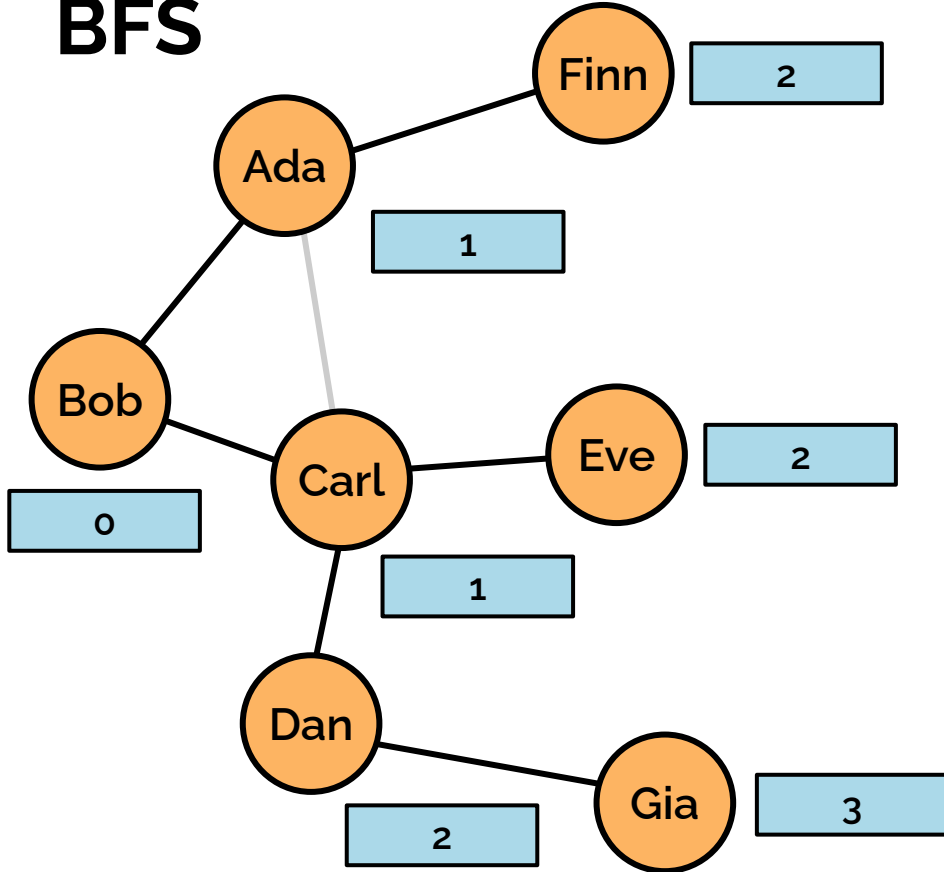
# Largest graphs

| <b>graph</b>     | <b> V </b> | <b> E </b> |
|------------------|------------|------------|
| datagen-9_3-zf   | 555M       | 1.3B       |
| datagen-sf10k-fb | 100M       | 18.8B      |
| graph500-30      | 450M       | 34.0B      |

# Algorithms



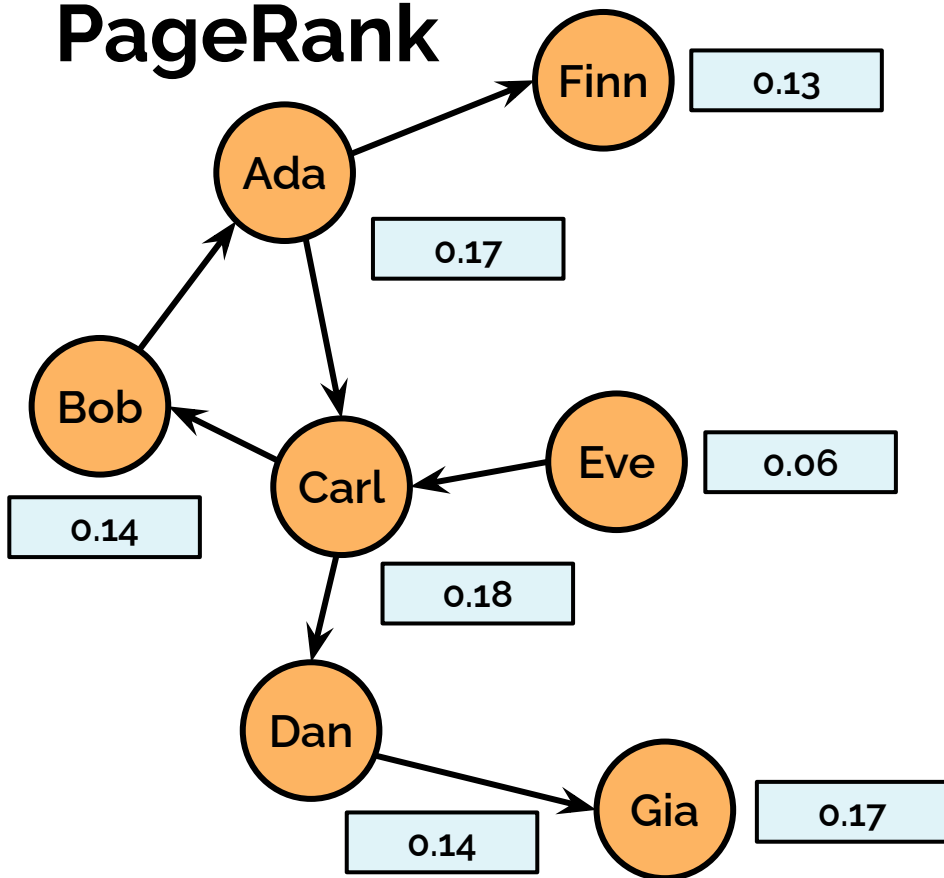
# BFS



**Breadth-first search**(*source*: “Bob”)

Assign the level of traversal for each vertex starting from the source (level = 0).

# PageRank



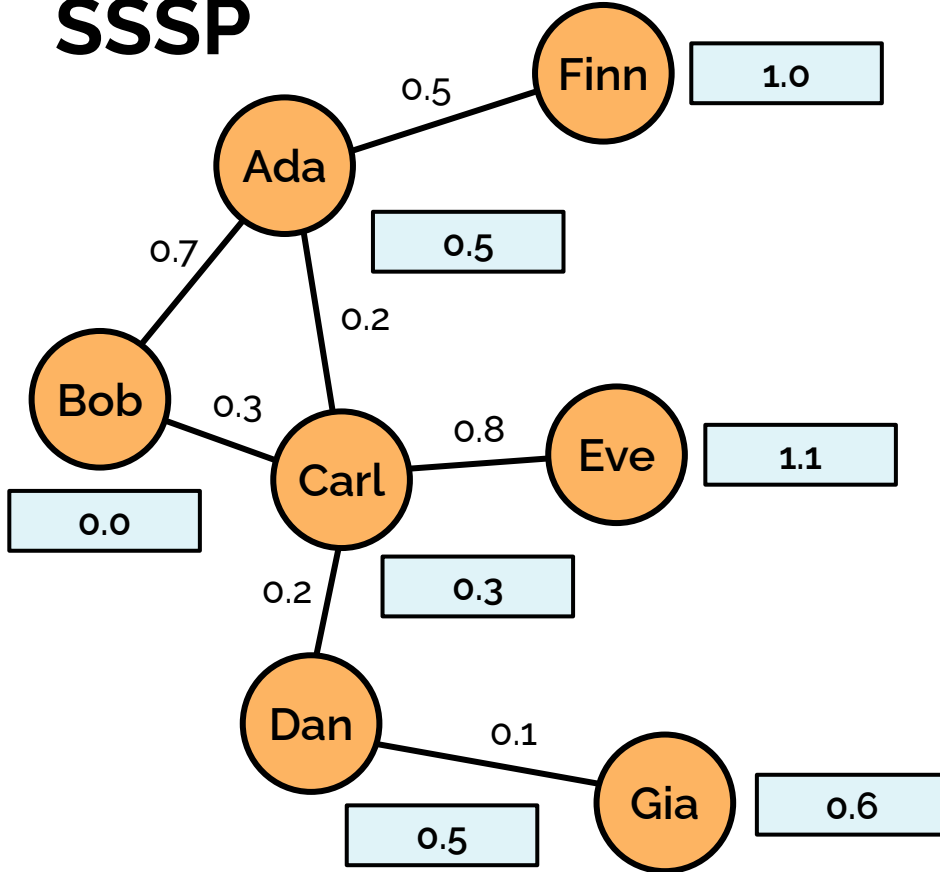
**PageRank**(damping factor: 0.85, iterations: 5)

$$PR_i(v) = \underbrace{\frac{1-d}{|V|}}_{\text{teleport}} + d \cdot \underbrace{\sum_{u \in N_{\text{in}}(v)} \frac{PR_{i-1}(u)}{|N_{\text{out}}(u)|}}_{\text{importance}} + \underbrace{\frac{d}{|V|} \cdot \sum_{w \in D} PR_{i-1}(w)}_{\text{redistributed from dangling}}$$

The PageRank variant in Graphalytics redistributes the PageRank values from sinks among all vertices to avoid “leaking” the PageRank out of the network.



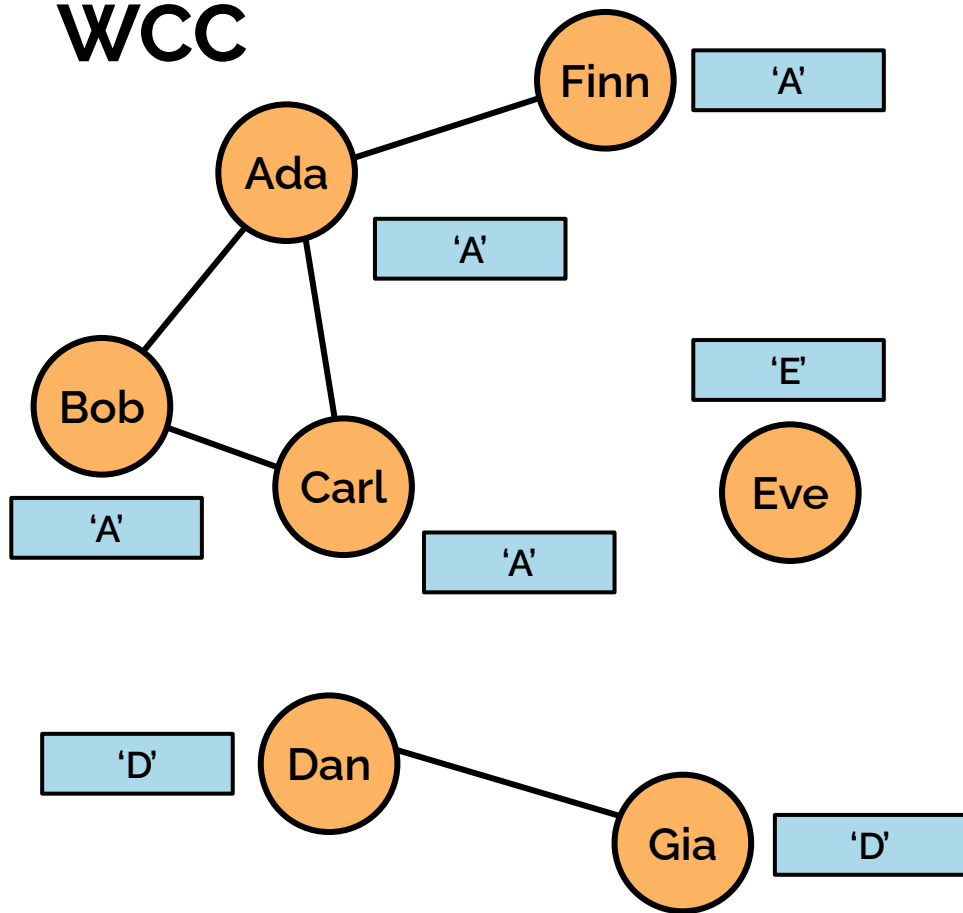
# SSSP



Single-source shortest paths(*source*: "Bob")

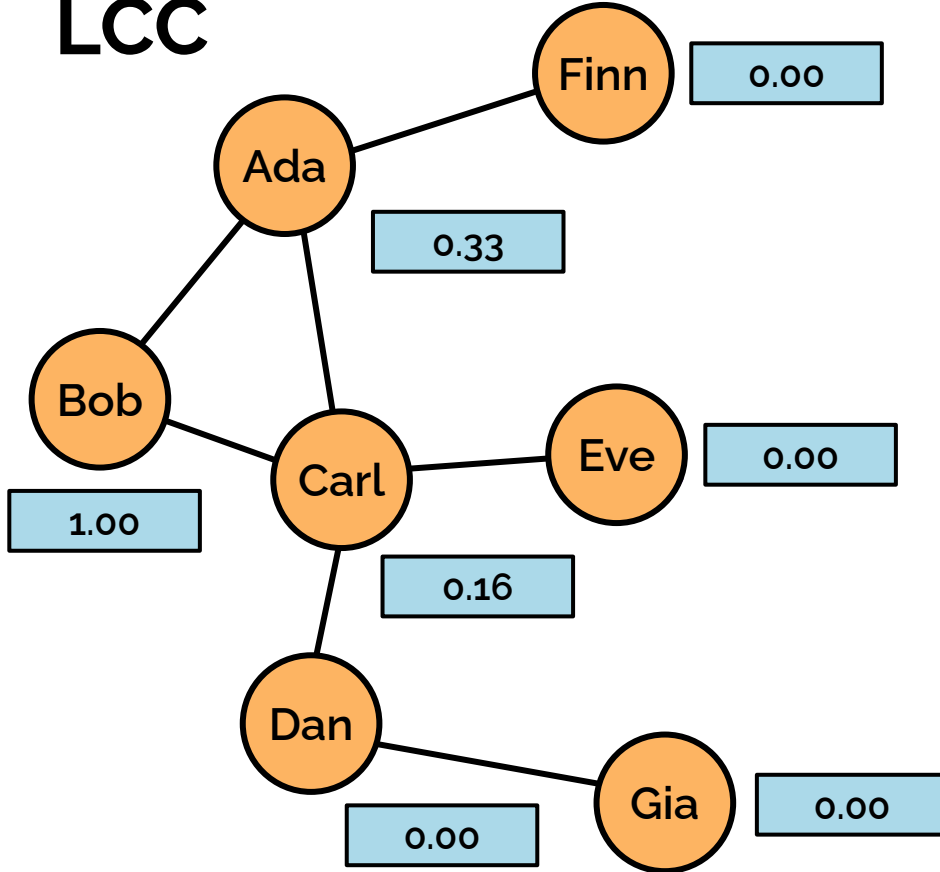
This is the only algorithm that uses edge weights. Many implementations use the delta-stepping SSSP algorithm. These are allowed to specify the delta value for each graph.

# WCC



Weakly connected components

# LCC



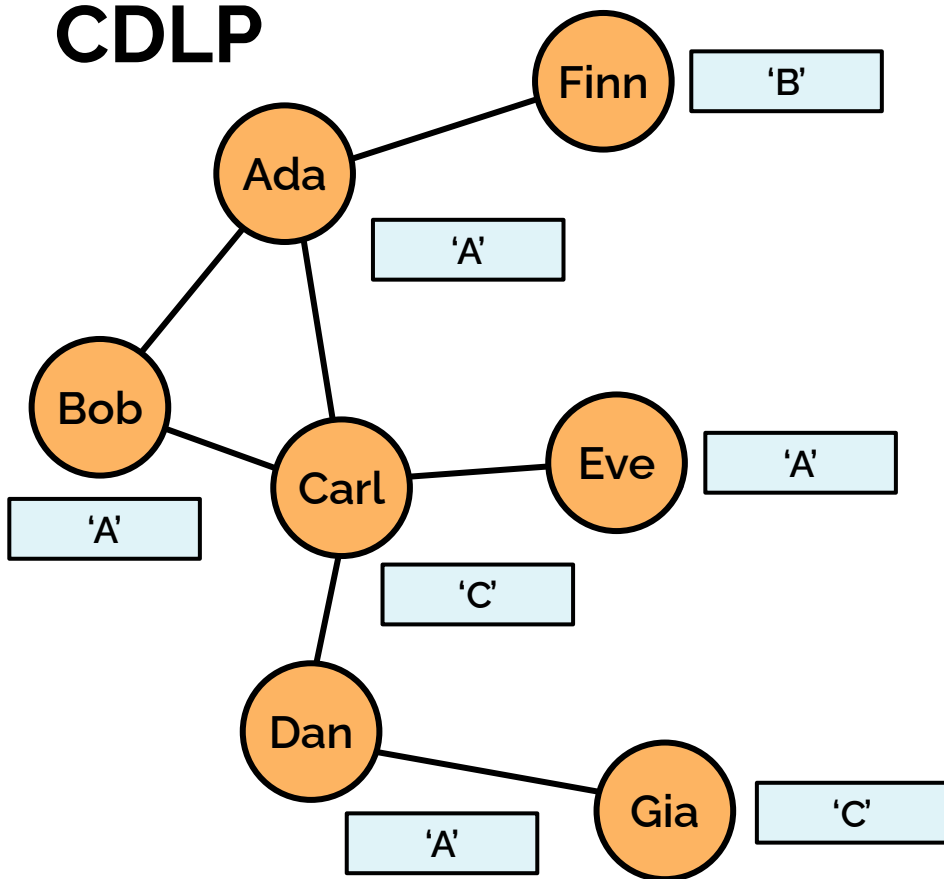
Local clustering coefficient

$$LCC(v) = \begin{cases} 0 & \text{If } |N(v)| \leq 1 \\ \frac{|\{(u,w) | u,w \in N(v) \wedge (u,w) \in E\}|}{|\{(u,w) | u,w \in N(v)\}|} & \text{Otherwise} \end{cases}$$

For each vertex, LCC is #triangles / #wedges.

This algorithm is very similar to triangle count.

# CDLP



Community detection using LP (*iterations: 2*)

$$L_i(v) = \min \left( \arg \max_l \left[ \left| \{u \in N_{\text{in}}(v) \mid L_{i-1}(u) = l\} \right| + \left| \{u \in N_{\text{out}}(v) \mid L_{i-1}(u) = l\} \right| \right] \right)$$

In each iteration, the next label of a vertex is selected as *the minimum mode value among the labels of the neighbours*.

# Graphalytics algorithms

All 6 algorithms:

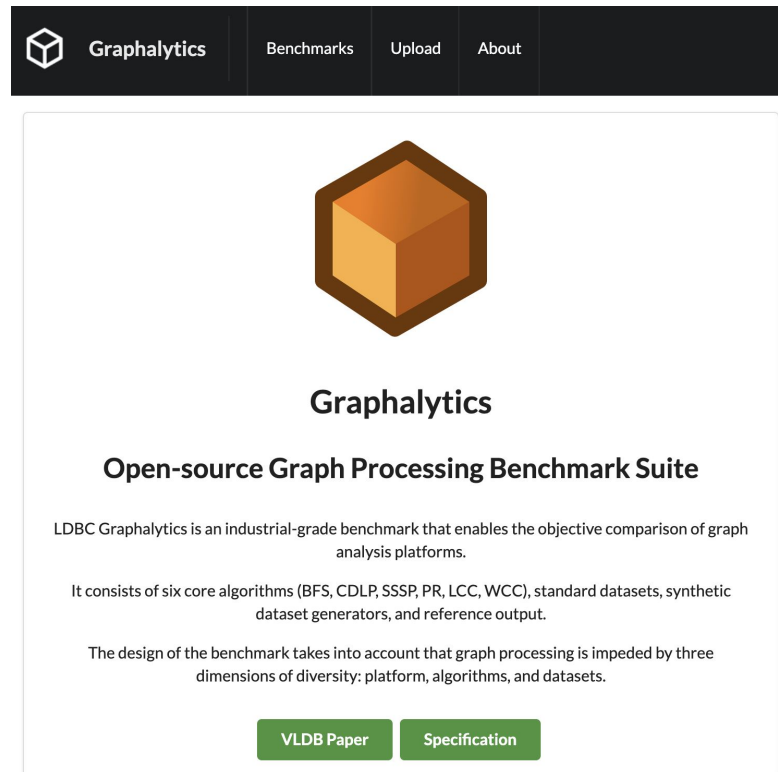
- have **directed and undirected variants**
- are **deterministic**

Validation uses different matching strategies:

- Exact match (BFS, CDLP)
- Epsilon match – relative tolerance of 0.01% (LCC, PR, SSSP)
- Equivalence match – same equivalence classes (WCC)

# Competition site is now open

<https://graphalytics.ldbcouncil.org/>



The screenshot shows the top navigation bar of the Graphalytics website, which is dark-themed with white text. The navigation items are 'Graphalytics' (with a cube icon), 'Benchmarks', 'Upload', and 'About'. Below the navigation bar is a large white box containing the main content. At the top of this box is a large orange 3D cube icon. Below the icon is the title 'Graphalytics' in bold black font, followed by the subtitle 'Open-source Graph Processing Benchmark Suite' in bold black font. The main text describes the benchmark as an industrial-grade tool for comparing graph analysis platforms, lists six core algorithms (BFS, CDLP, SSSP, PR, LCC, WCC), and mentions standard datasets, synthetic dataset generators, and reference output. It also notes that the design considers three dimensions of diversity: platform, algorithms, and datasets. At the bottom of the white box are two green buttons: 'VLDB Paper' and 'Specification'.

Graphalytics

**Open-source Graph Processing Benchmark Suite**

LDDBC Graphalytics is an industrial-grade benchmark that enables the objective comparison of graph analysis platforms.

It consists of six core algorithms (BFS, CDLP, SSSP, PR, LCC, WCC), standard datasets, synthetic dataset generators, and reference output.

The design of the benchmark takes into account that graph processing is impeded by three dimensions of diversity: platform, algorithms, and datasets.

[VLDB Paper](#) [Specification](#)

# **LDDB Social Network Benchmark**



Data set

Queries

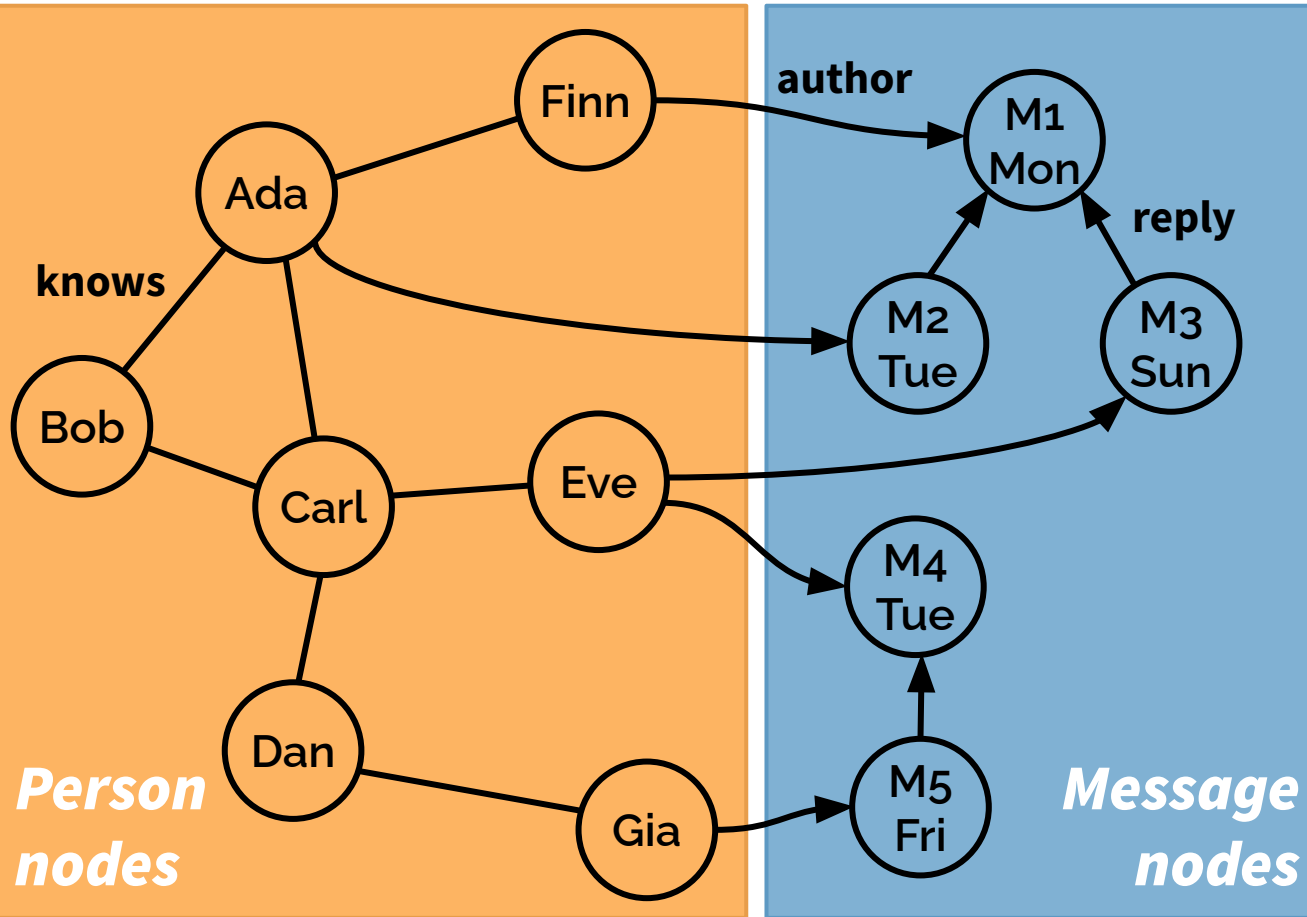
Updates



Data set

Queries

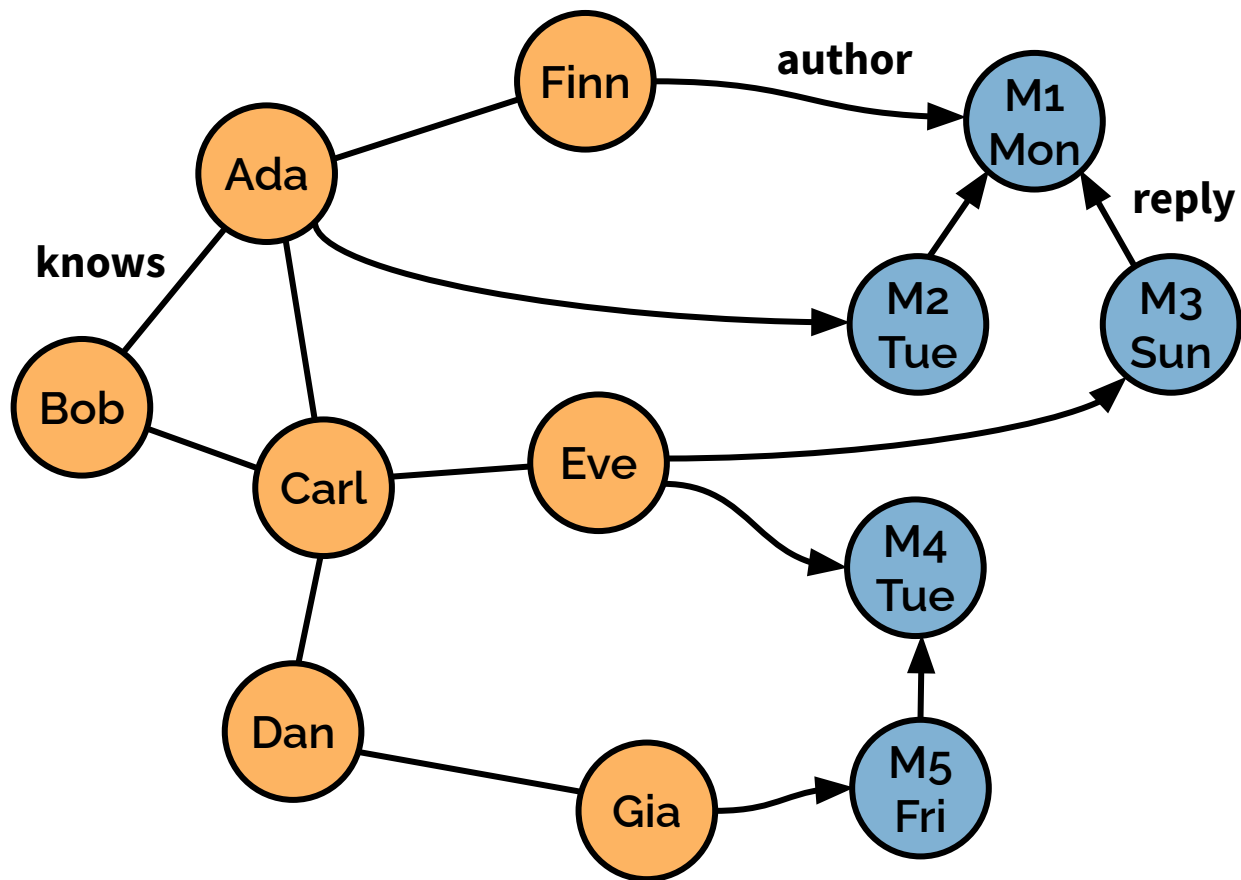
Updates



Data set

Queries

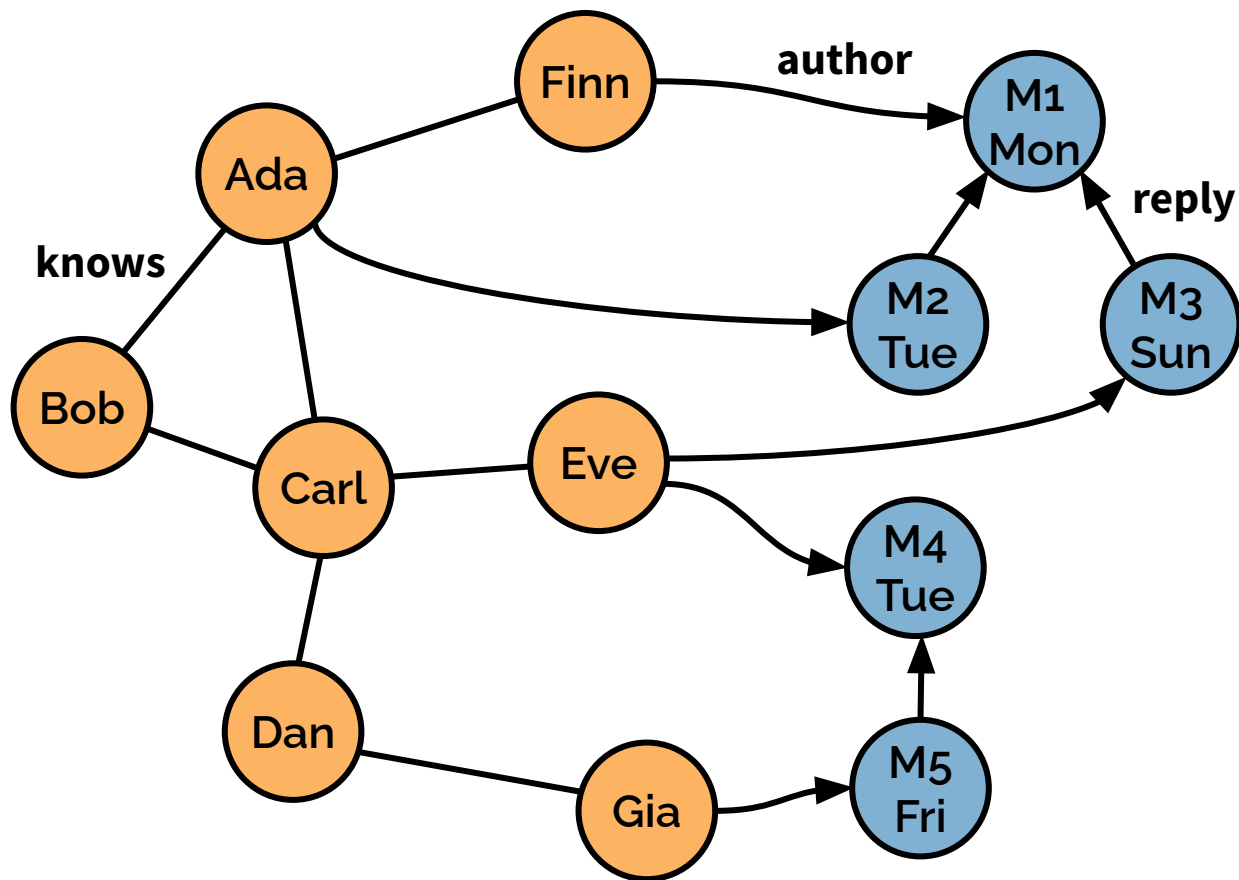
Updates



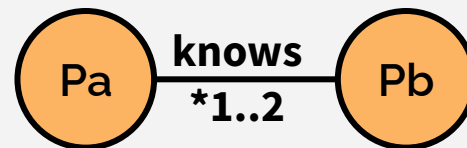
Data set

Queries

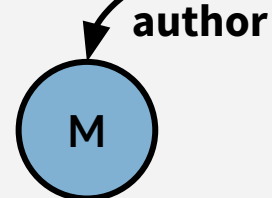
Updates



**Q9(\$name, \$day)**



*name = \$name*

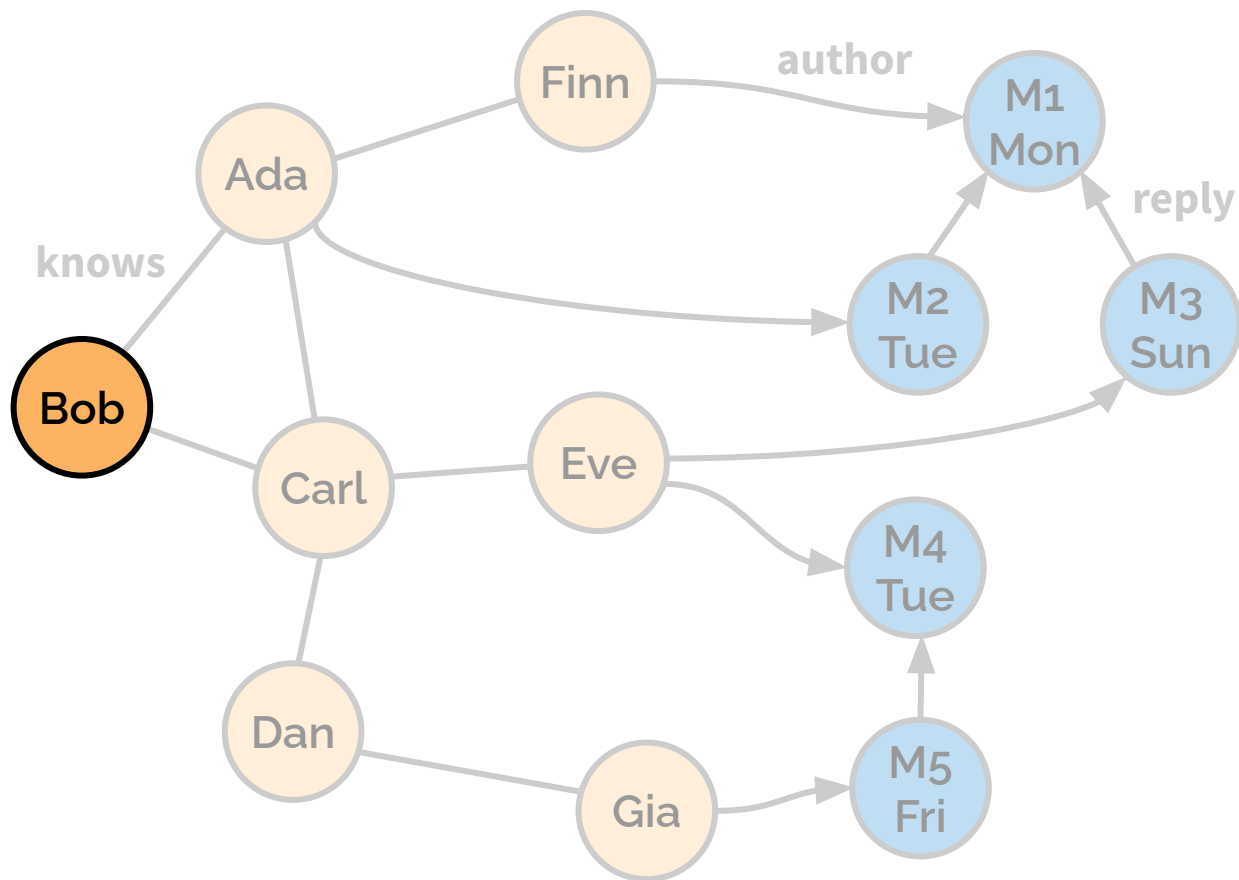


*creation date < \$day*

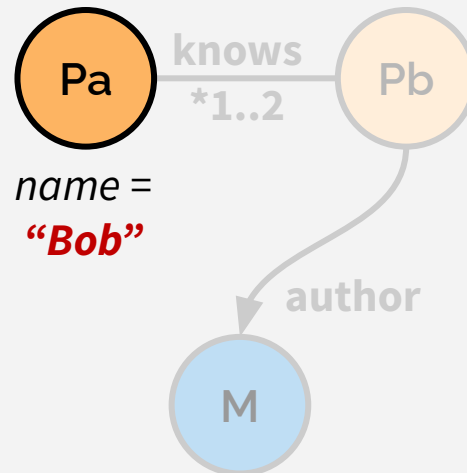
Data set

Queries

Updates



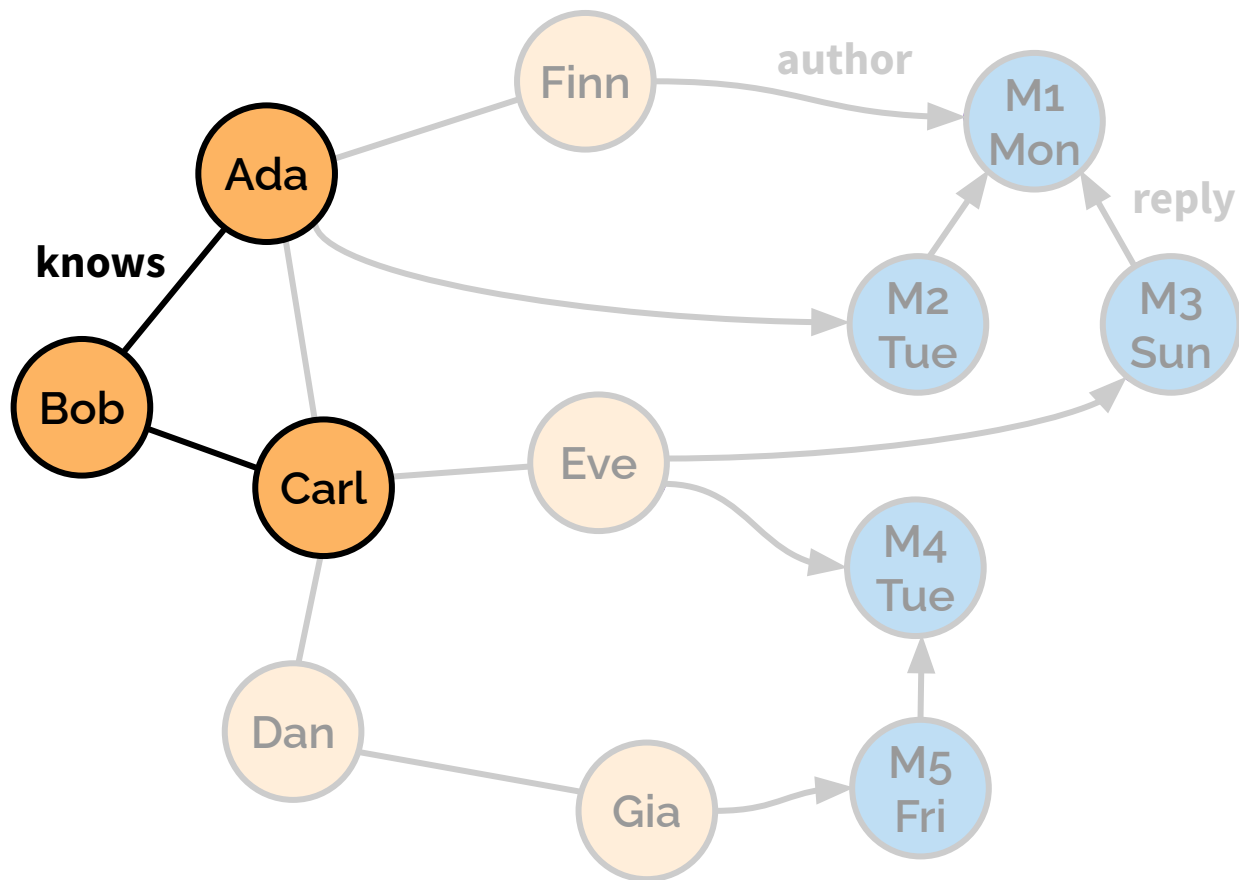
Q9(**“Bob”**, **“Sat”**)



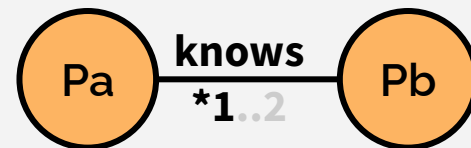
Data set

Queries

Updates



Q9(**“Bob”**, **“Sat”**)



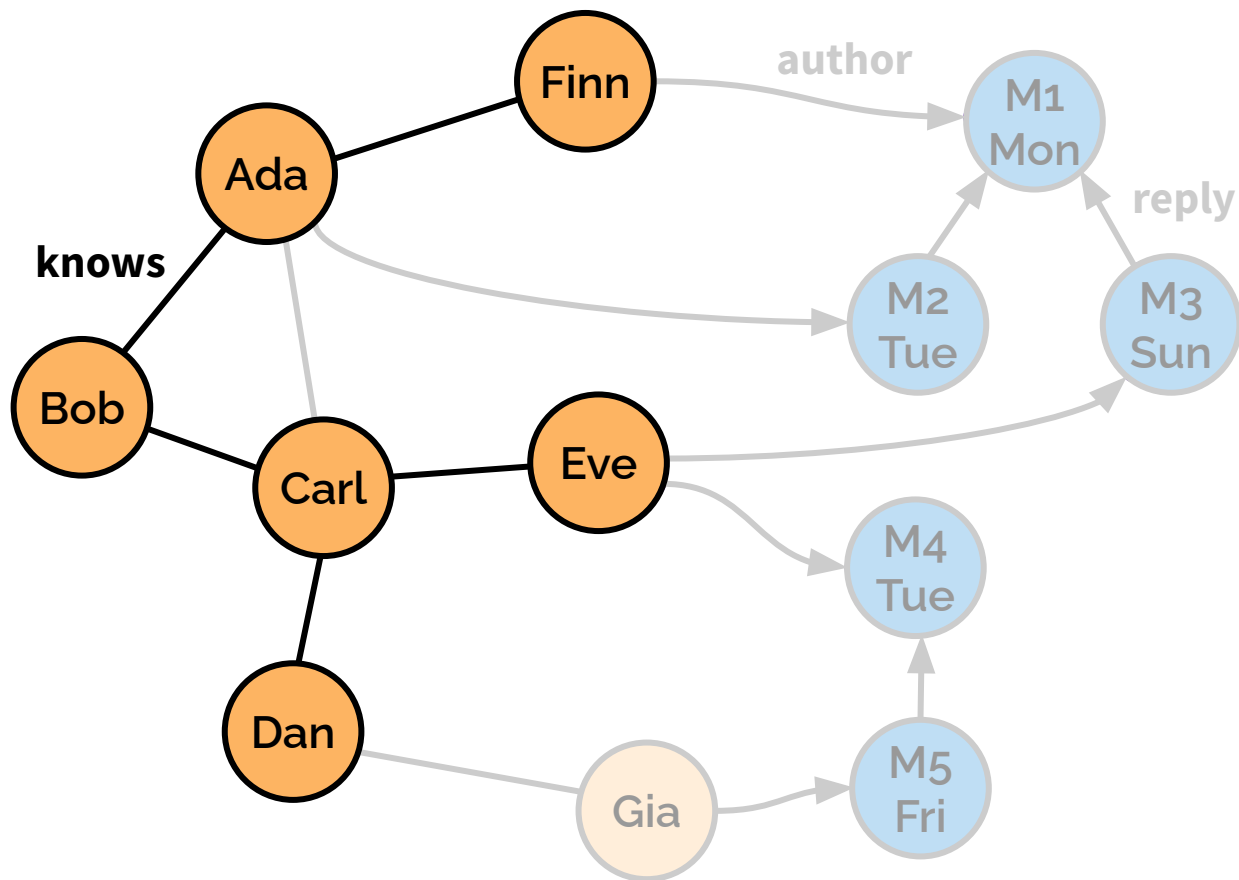
name =  
**“Bob”**

creation date < **“Sat”**

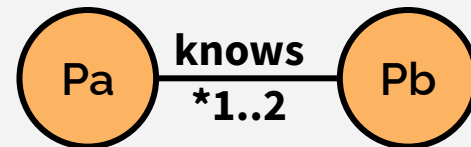
Data set

Queries

Updates



Q9(**"Bob"**, **"Sat"**)



name =  
**"Bob"**

author

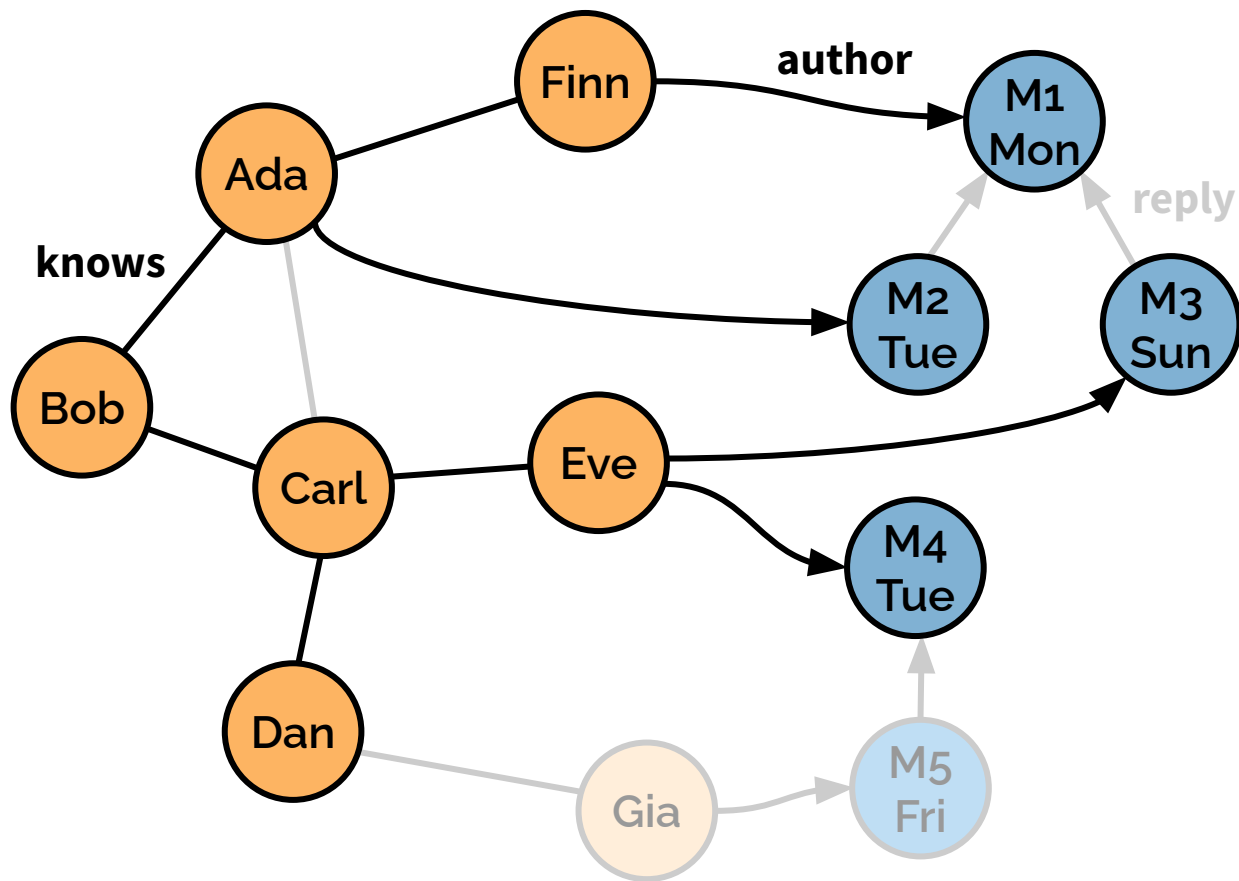


creation date < **"Sat"**

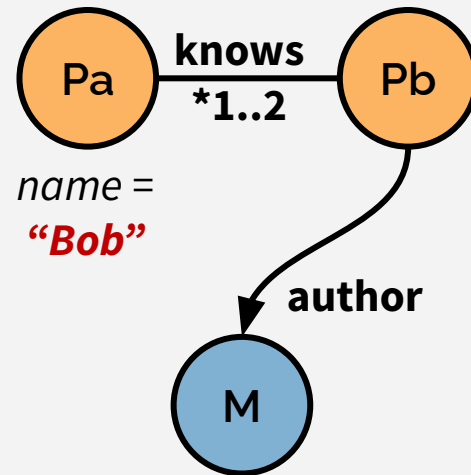
Data set

Queries

Updates



Q9(**"Bob"**, **"Sat"**)



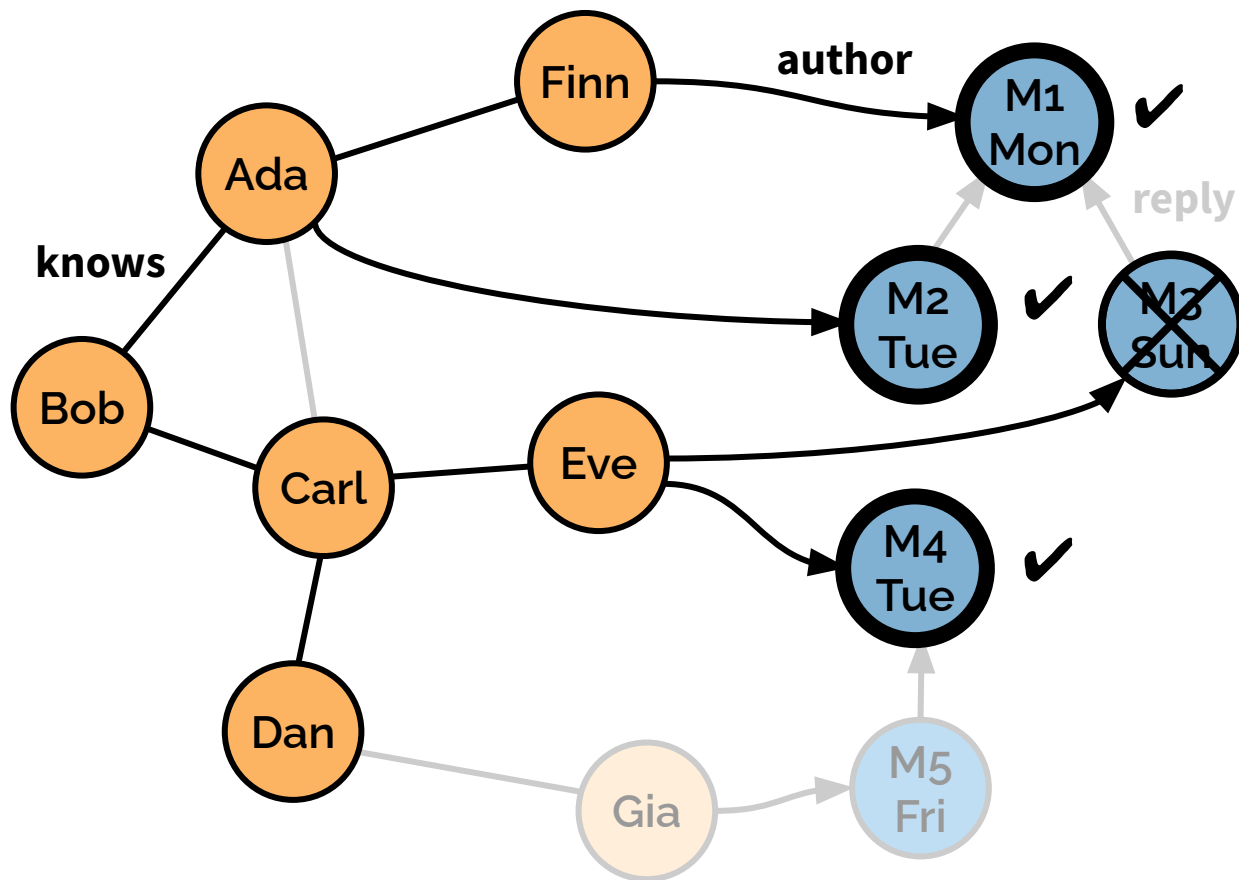
name =  
**"Bob"**

creation date < **"Sat"**

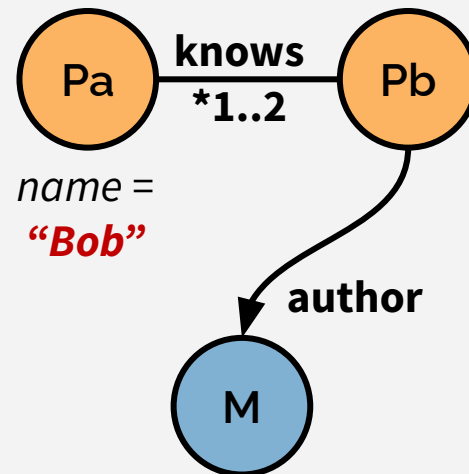
Data set

Queries

Updates



Q9("Bob", "Sat")

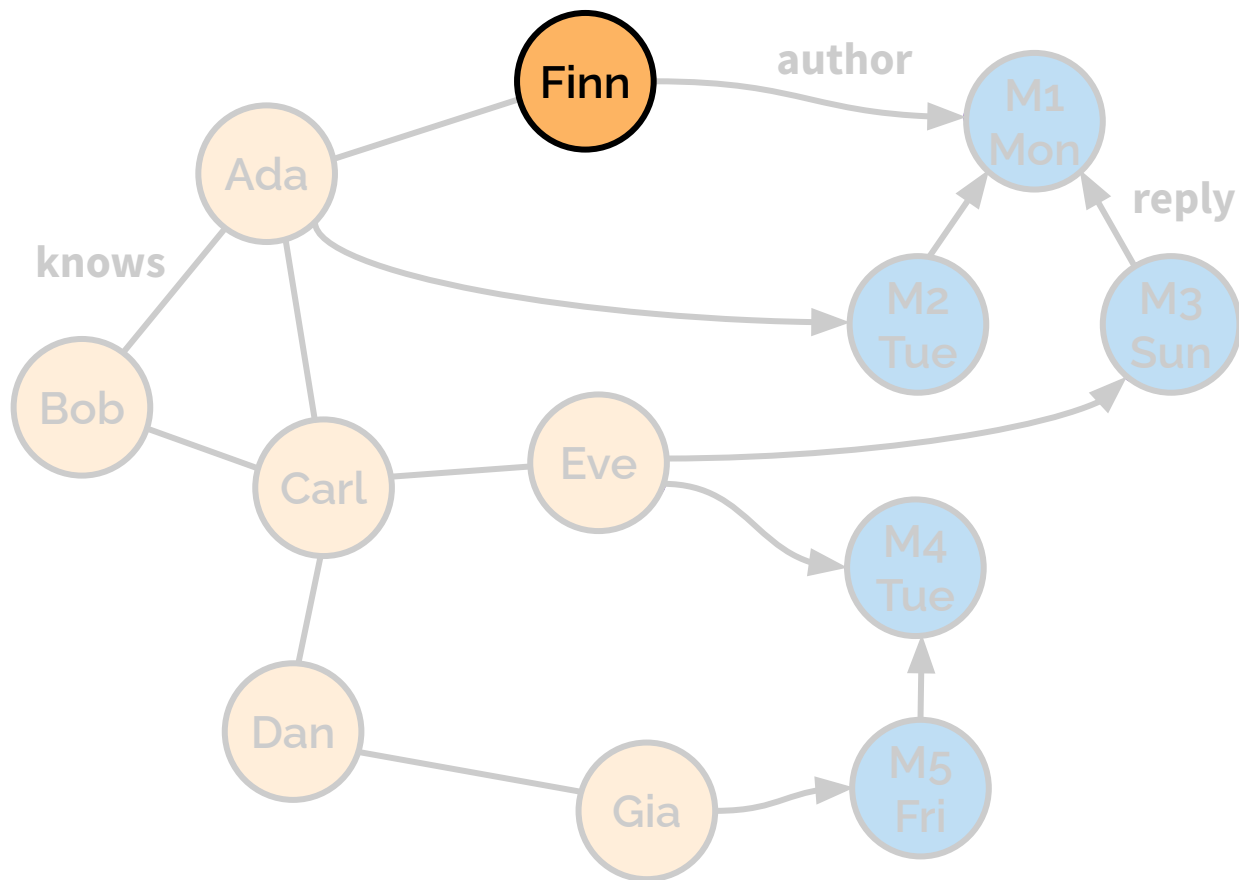




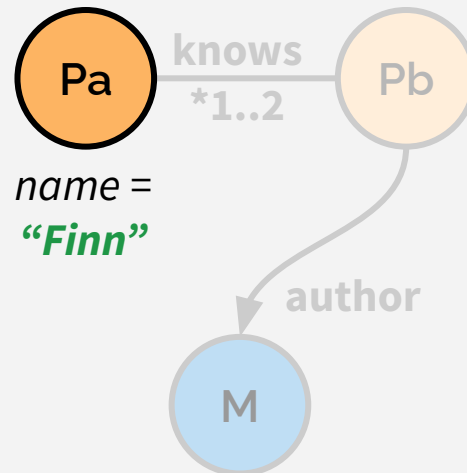
Data set

Queries

Updates



Q9(“Finn”, “Wed”)



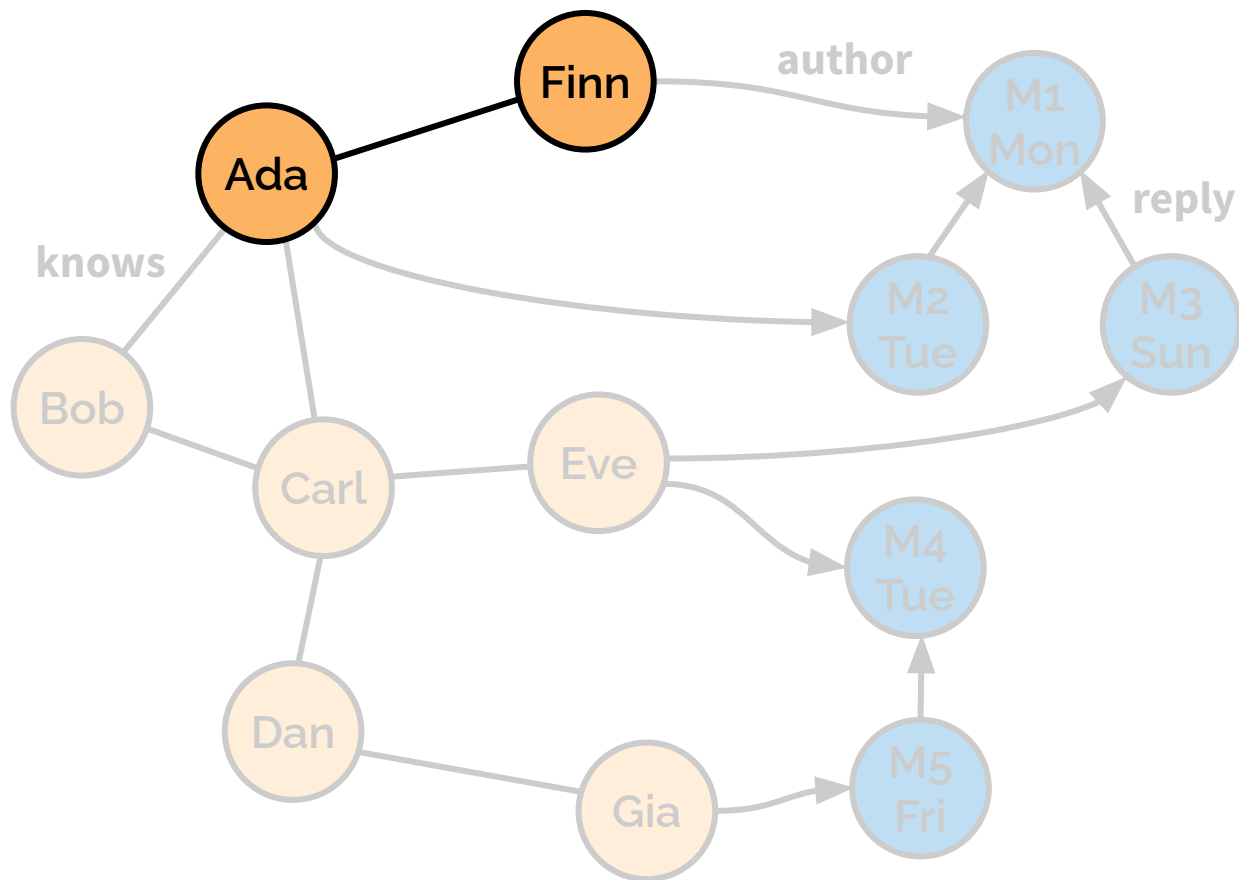
name =  
“Finn”

creation date < “Wed”

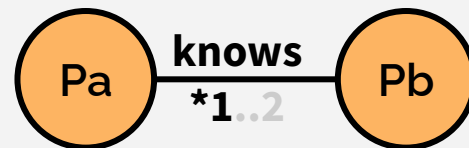
Data set

Queries

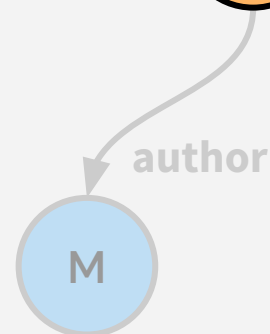
Updates



Q9(“Finn”, “Wed”)



name =  
“Finn”

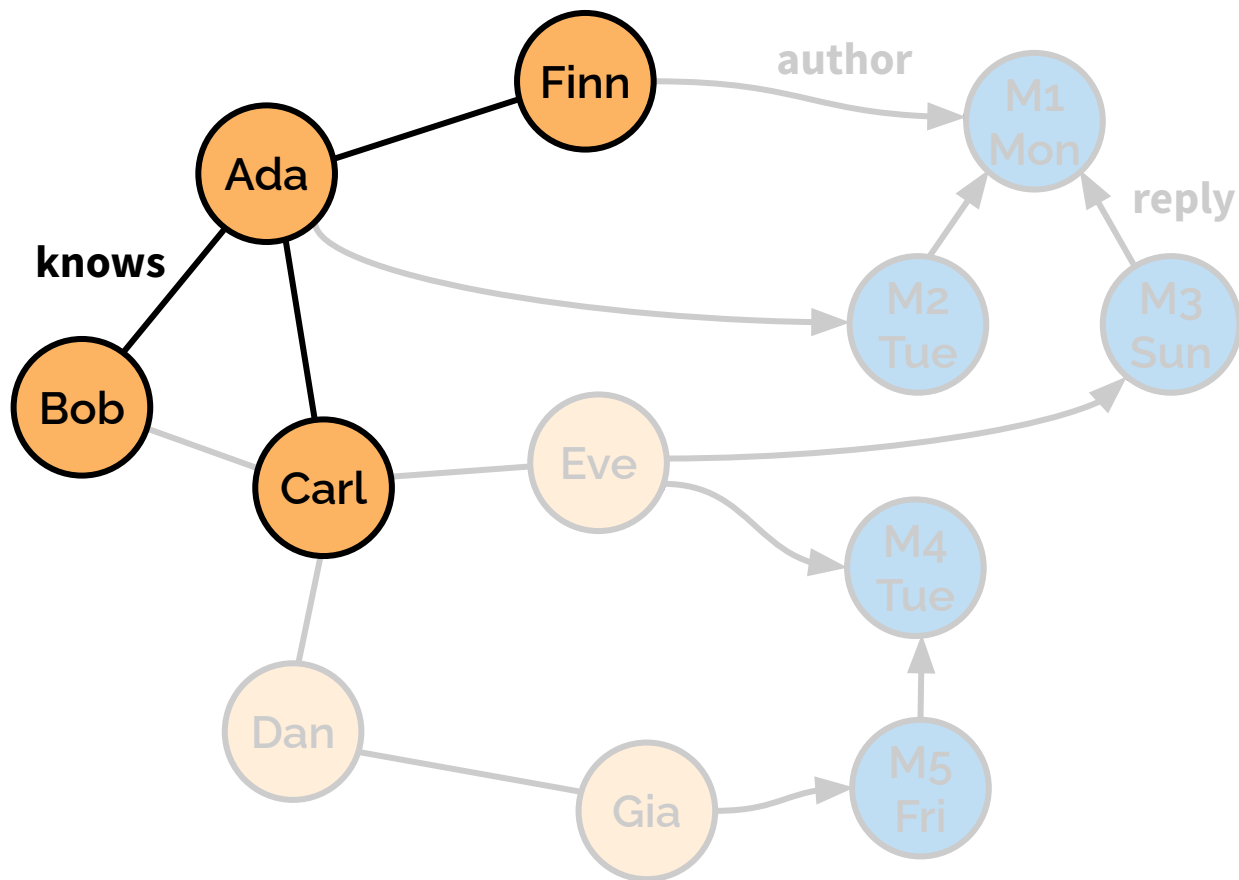


creation date < “Wed”

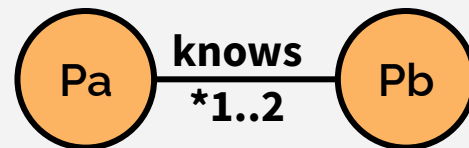
Data set

Queries

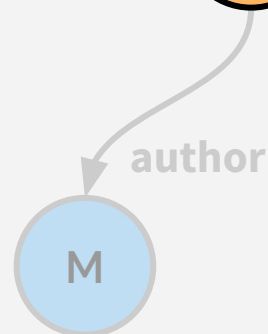
Updates



Q9(“Finn”, “Wed”)



name =  
“Finn”

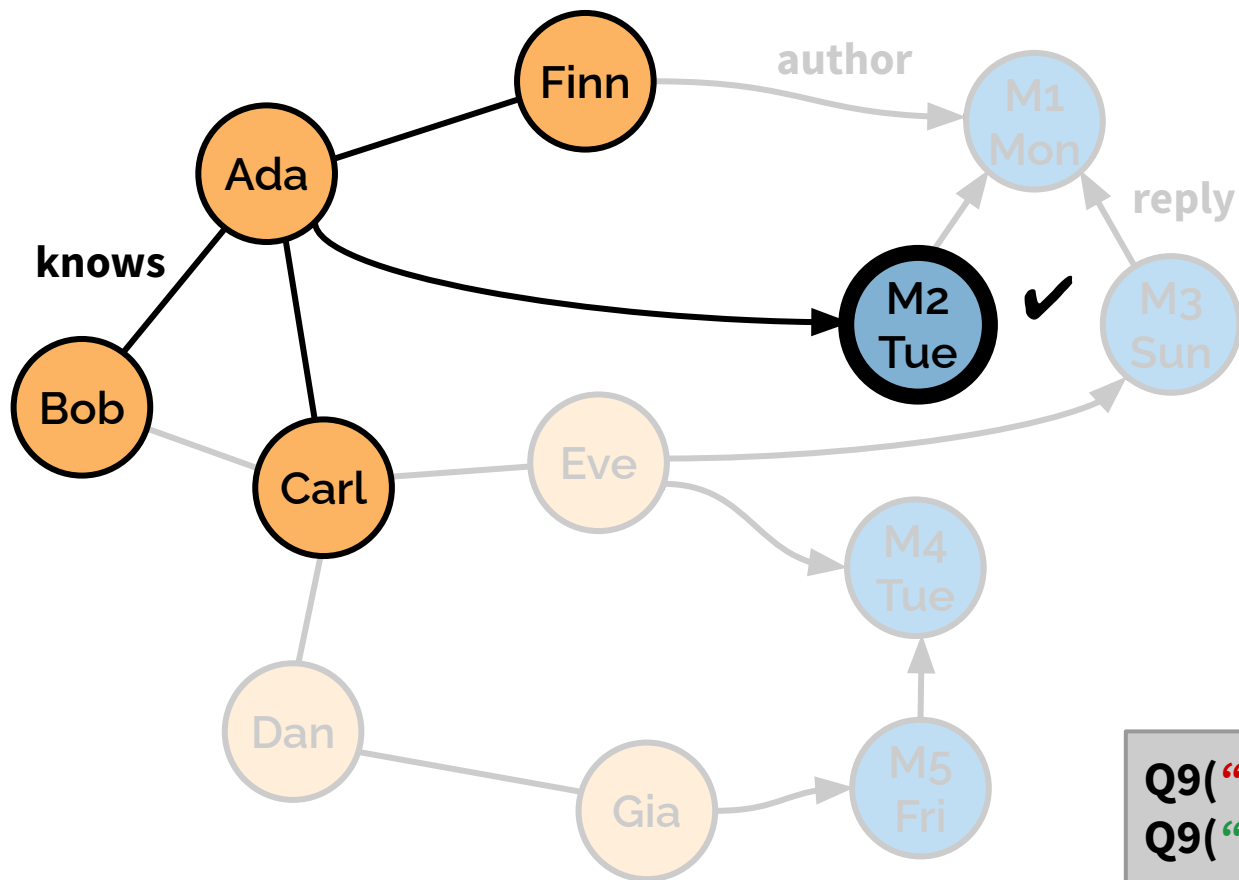


creation date < “Wed”

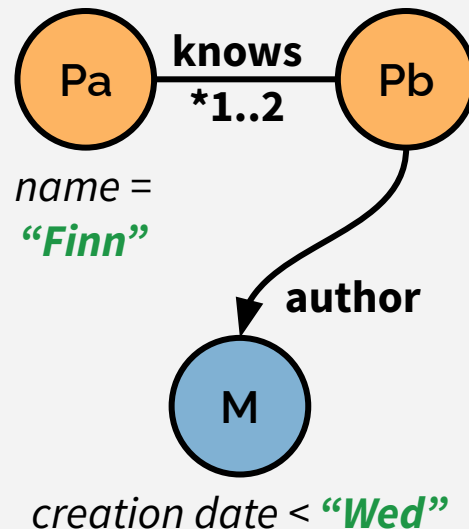
Data set

Queries

Updates



Q9("Finn", "Wed")



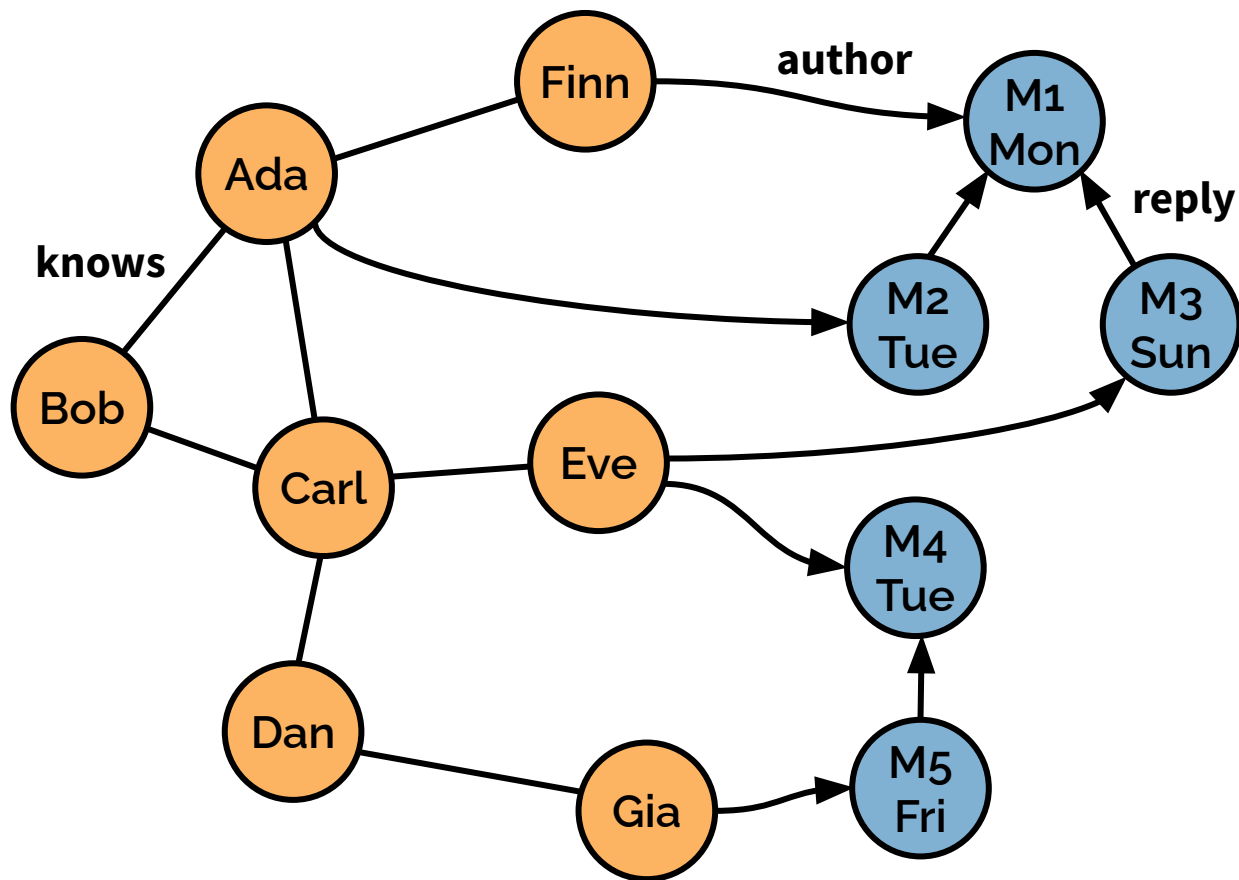
Q9("Bob", "Sat"): 10 nodes

Q9("Finn", "Wed"): 5 nodes

Data set

Queries

Updates



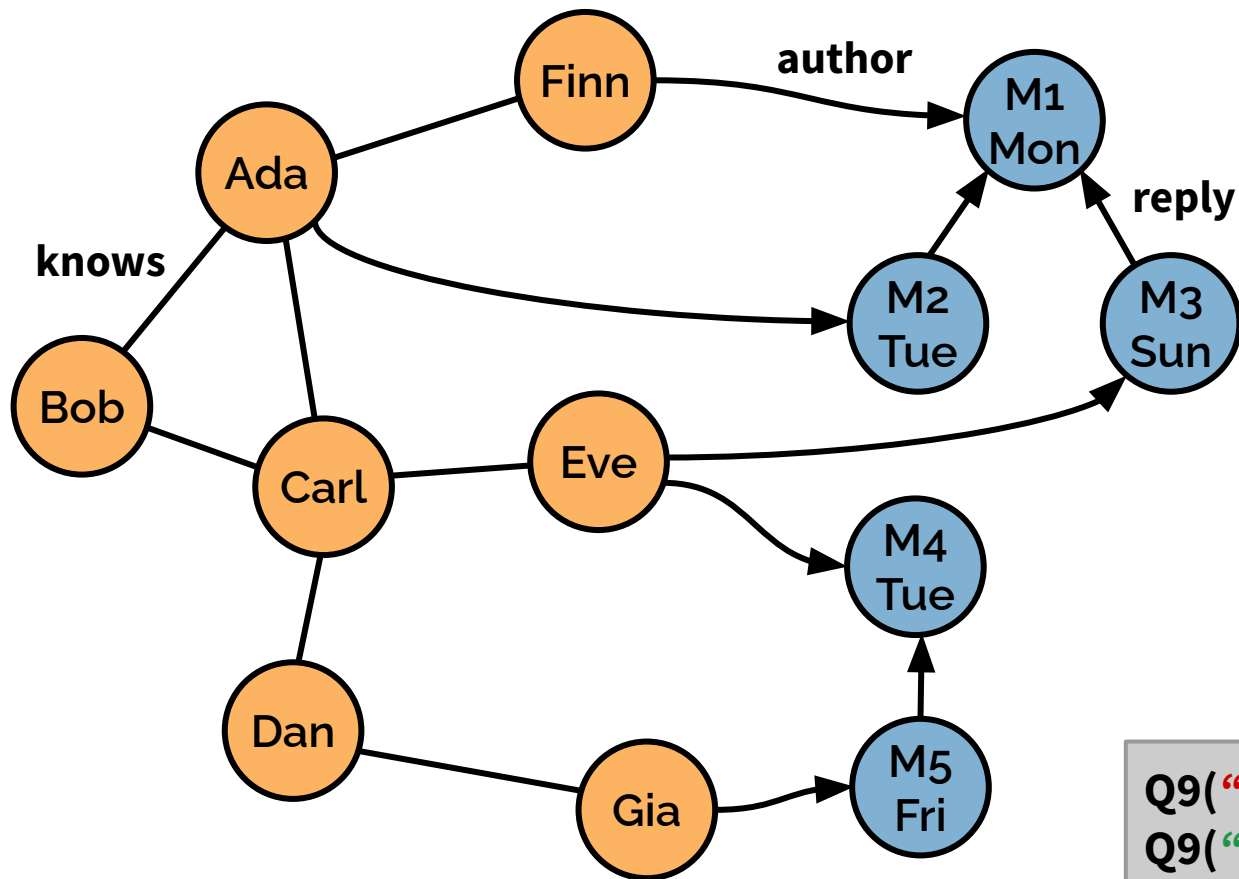
Factor table

| name        | 1-hop friends | 2-hop friends |
|-------------|---------------|---------------|
| <b>Bob</b>  | 2             | 5             |
| Carl        | 4             | 4             |
| Ada         | 3             | 4             |
| Dan         | 2             | 3             |
| Eve         | 1             | 3             |
| <b>Finn</b> | 1             | 2             |
| Gia         | 1             | 1             |

Data set

Queries

Updates



Factor table

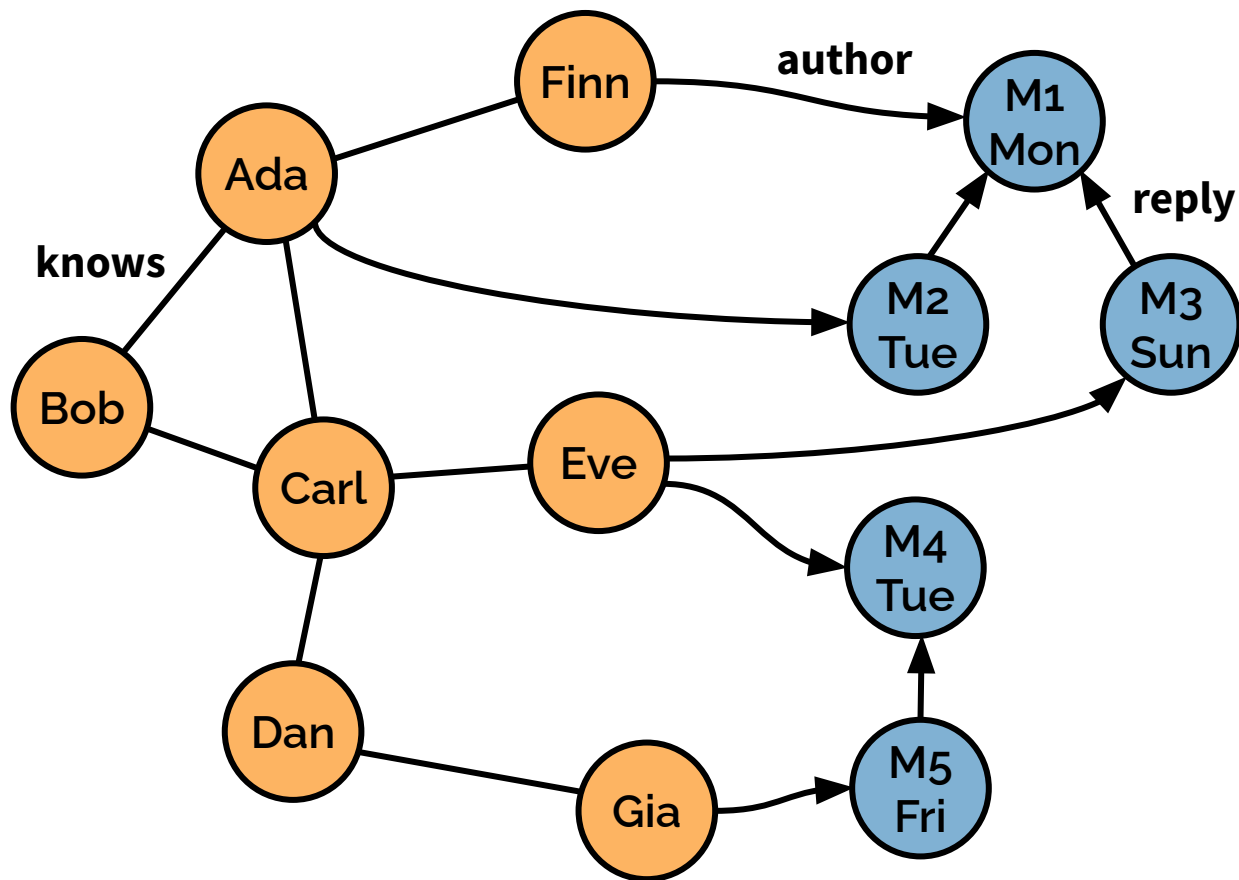
| name        | 1-hop friends | 2-hop friends |
|-------------|---------------|---------------|
| <b>Bob</b>  | 2             | <b>5</b>      |
| Carl        | 4             | 4             |
| Ada         | 3             | 4             |
| Dan         | 2             | 3             |
| Eve         | 1             | 3             |
| <b>Finn</b> | <b>1</b>      | <b>2</b>      |
|             |               | 1             |

Q9("Bob", "Sat"): 10 nodes  
 Q9("Finn", "Wed"): 5 nodes

Data set

Queries

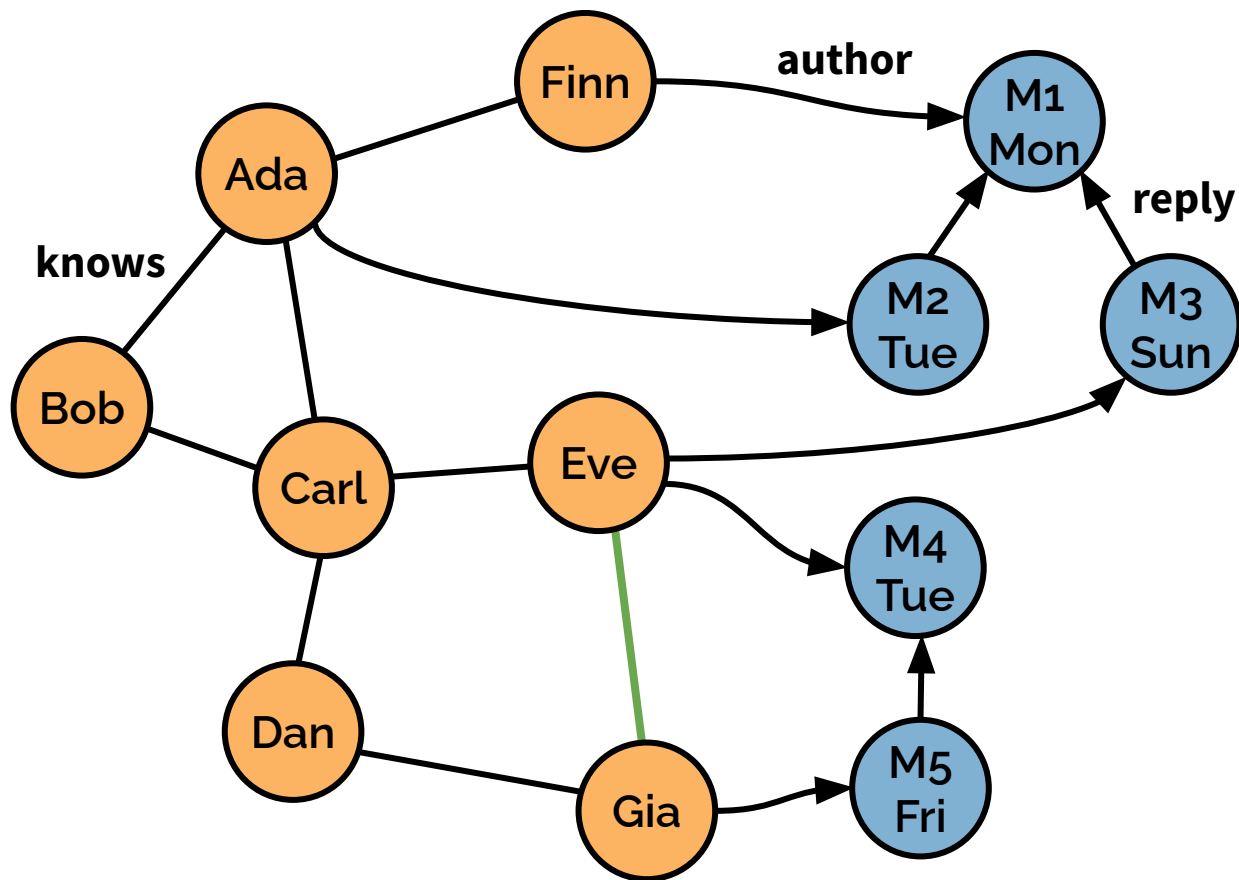
Updates



Data set

Queries

Updates



Updates

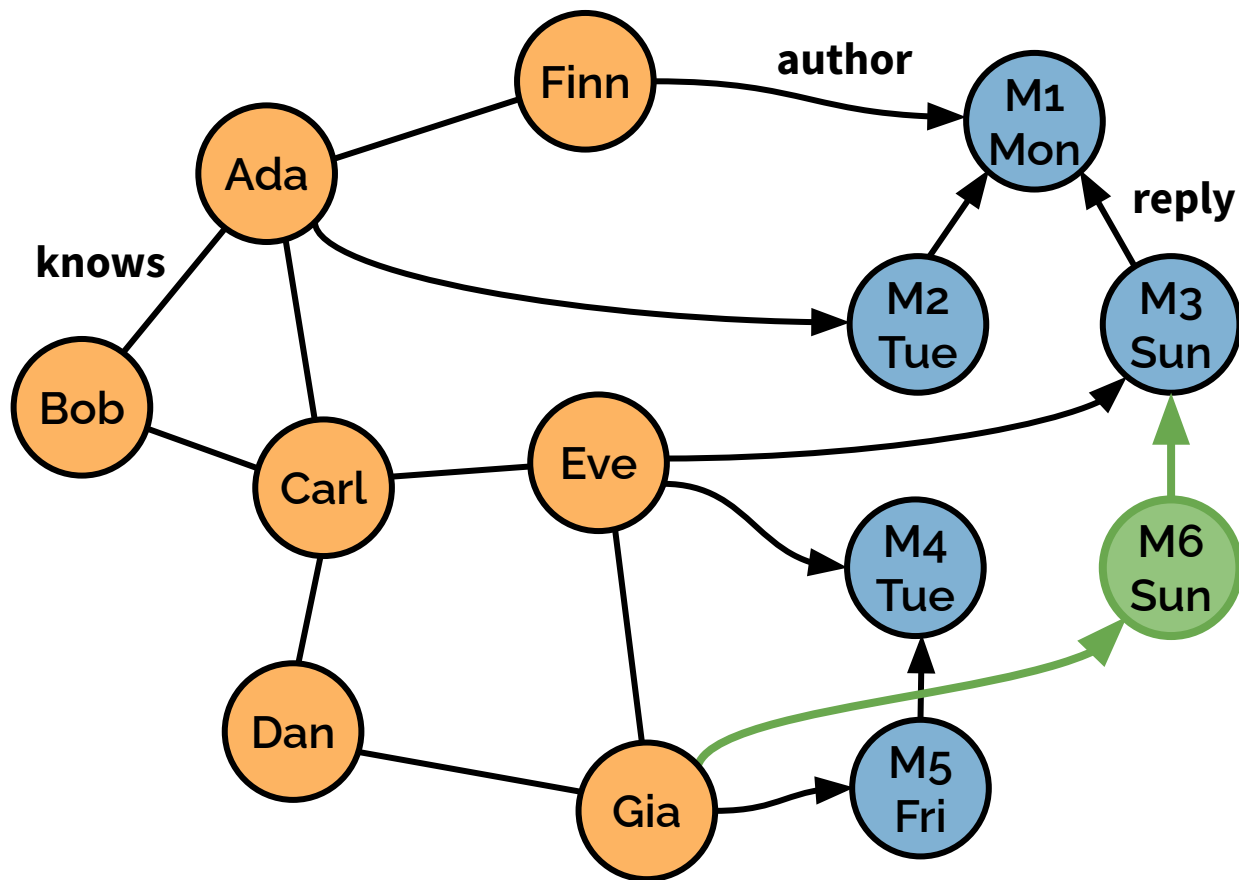
+ knows("Eve", "Gia")



Data set

Queries

Updates



Updates

+ knows("Eve", "Gia")

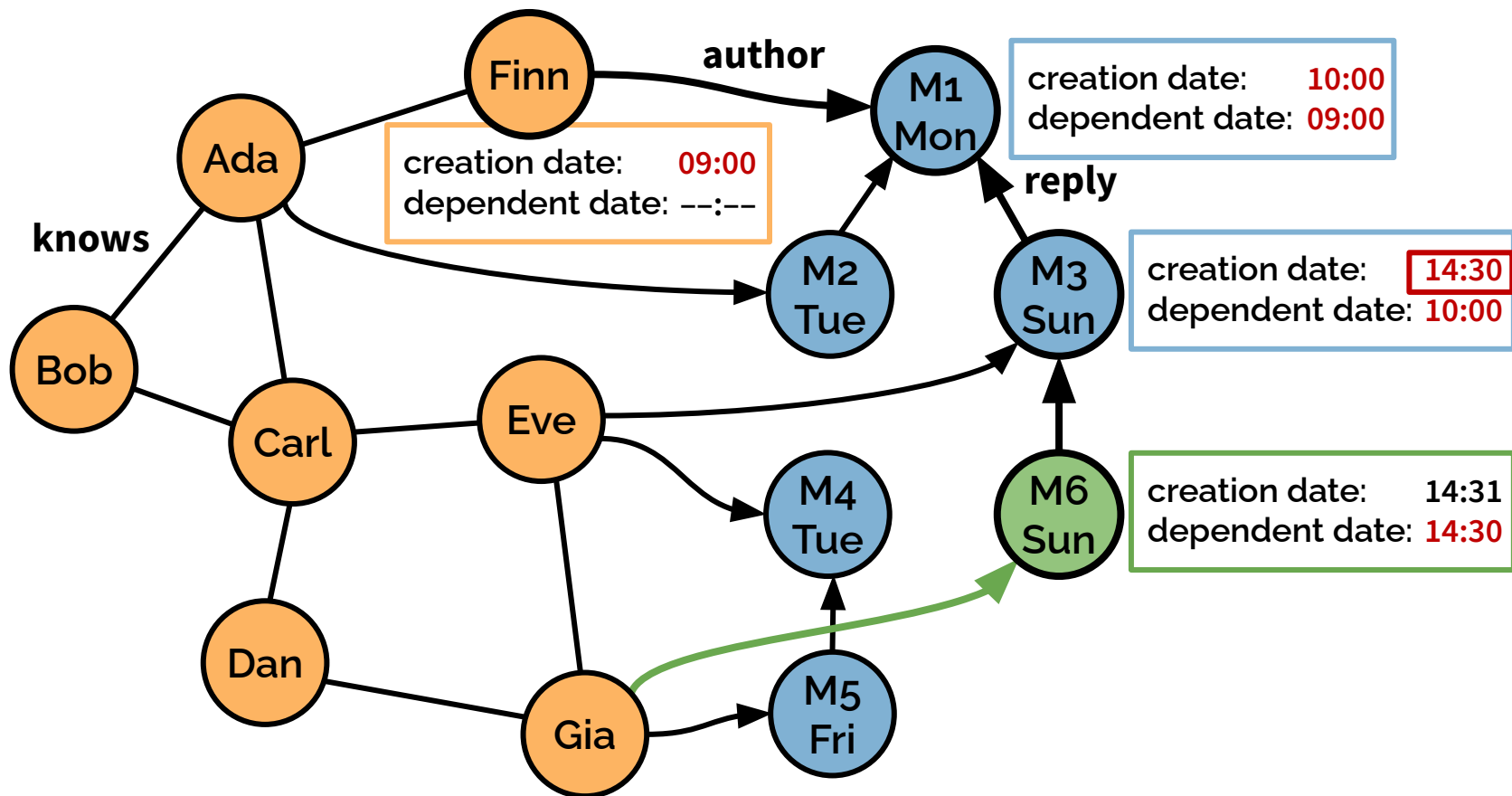
+ Comment("Gia", "M3")

When is this operation executable?

Data set

Queries

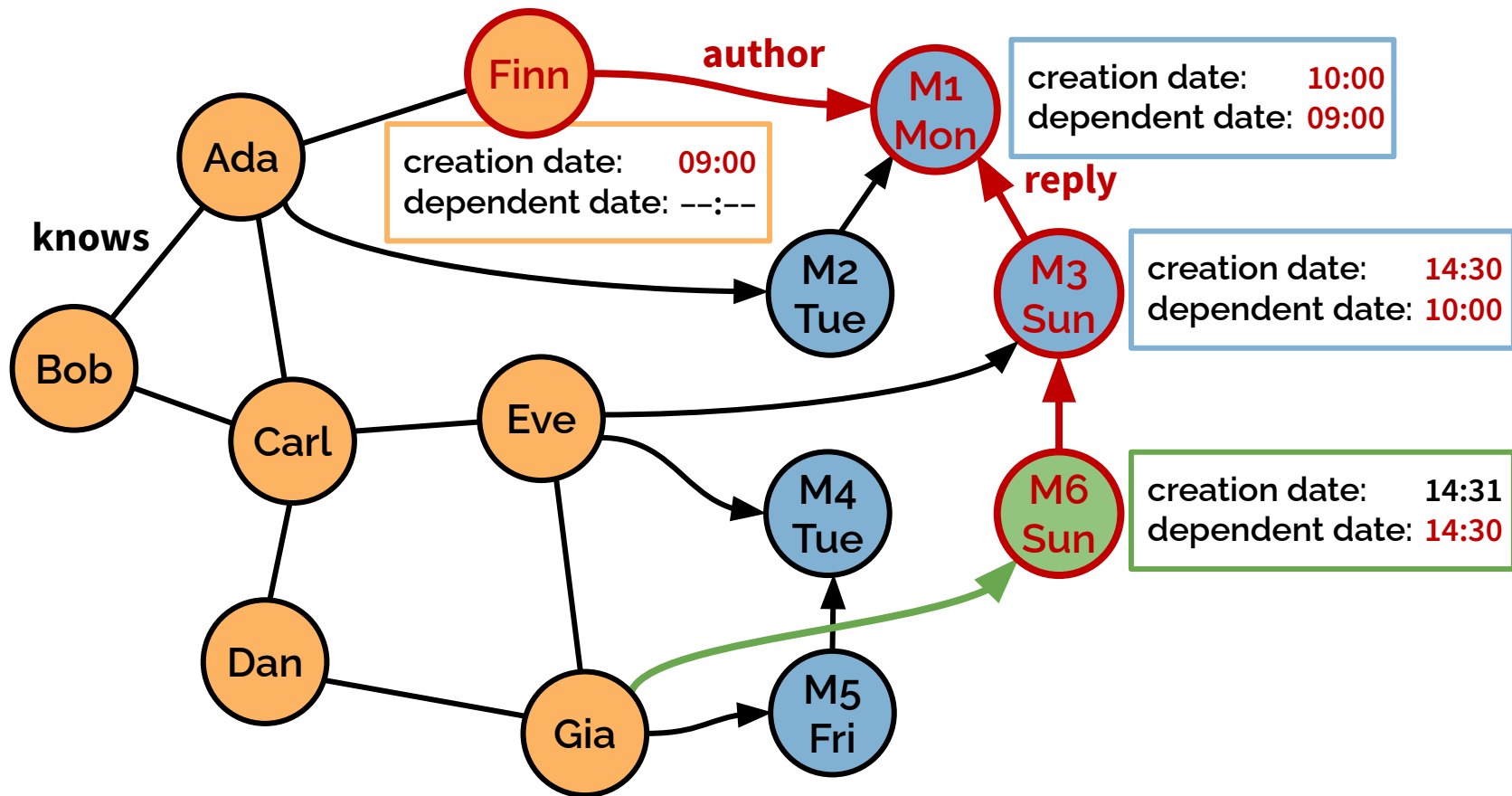
Updates



Data set

Queries

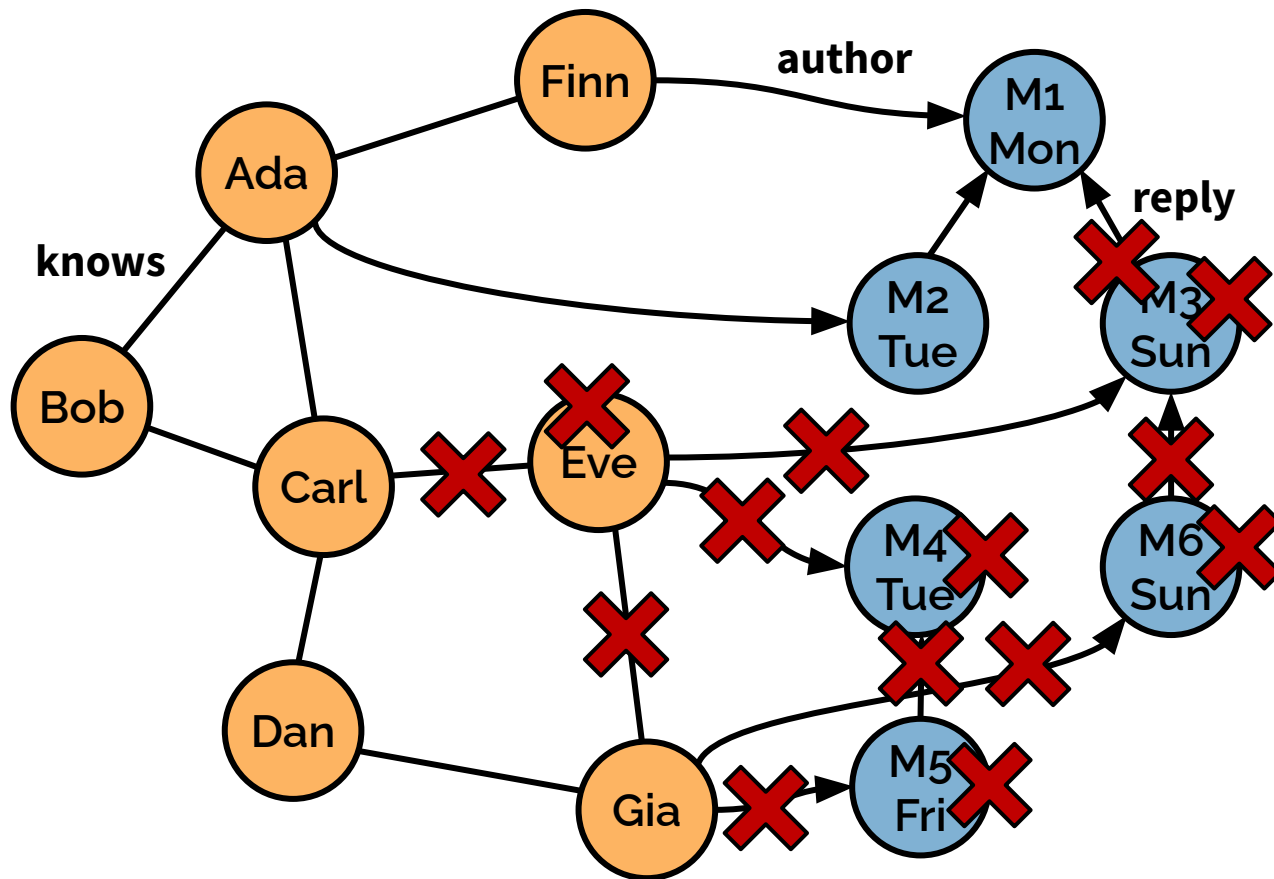
Updates



Data set

Queries

Updates



Updates

+ knows("Eve", "Gia")

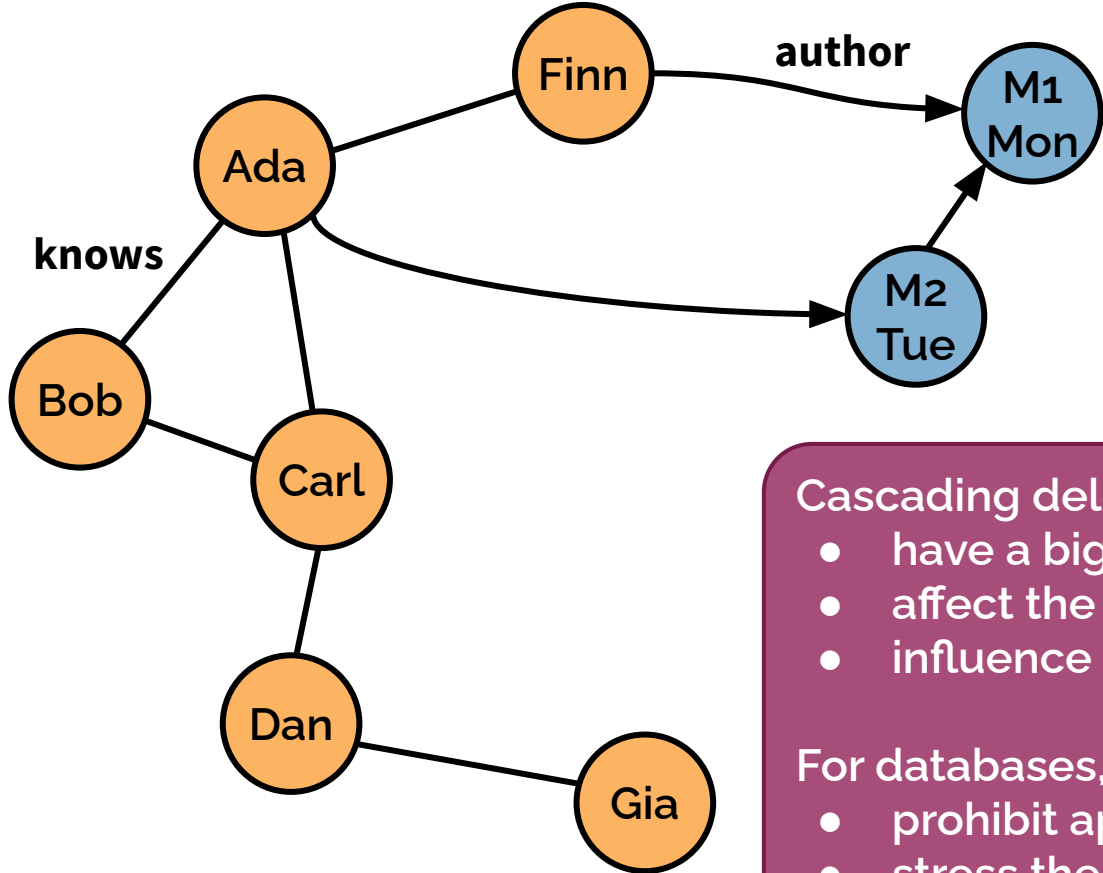
+ Comment("Gia", "M3")

- Person("Eve")

Data set

Queries

Updates



Updates

+ knows("Eve", "Gia")

+ Comment("Gia", "M3")

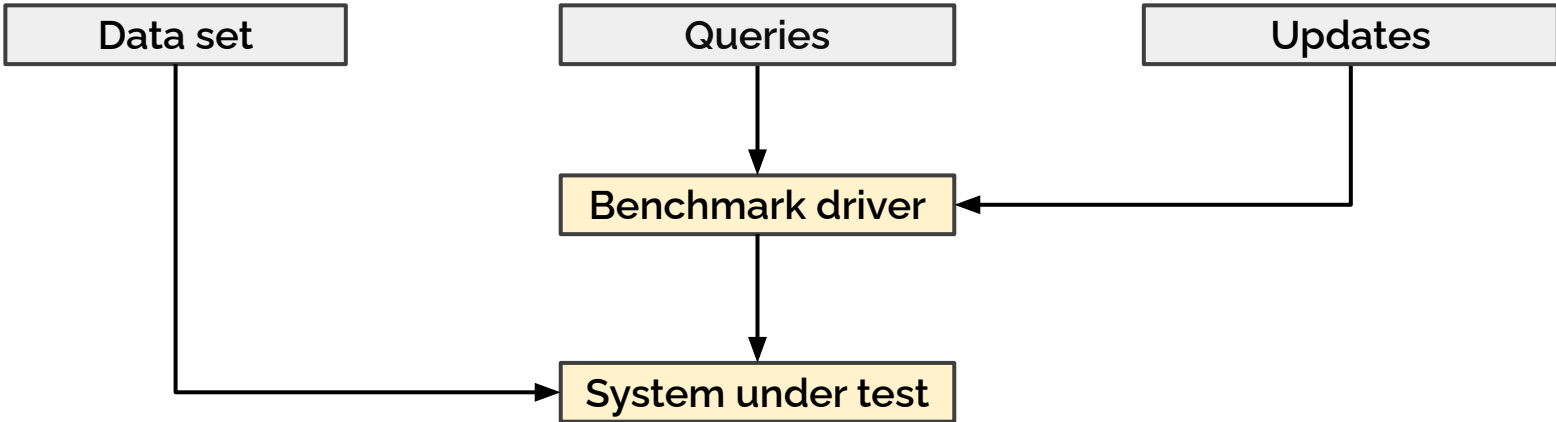
- Person("Eve")

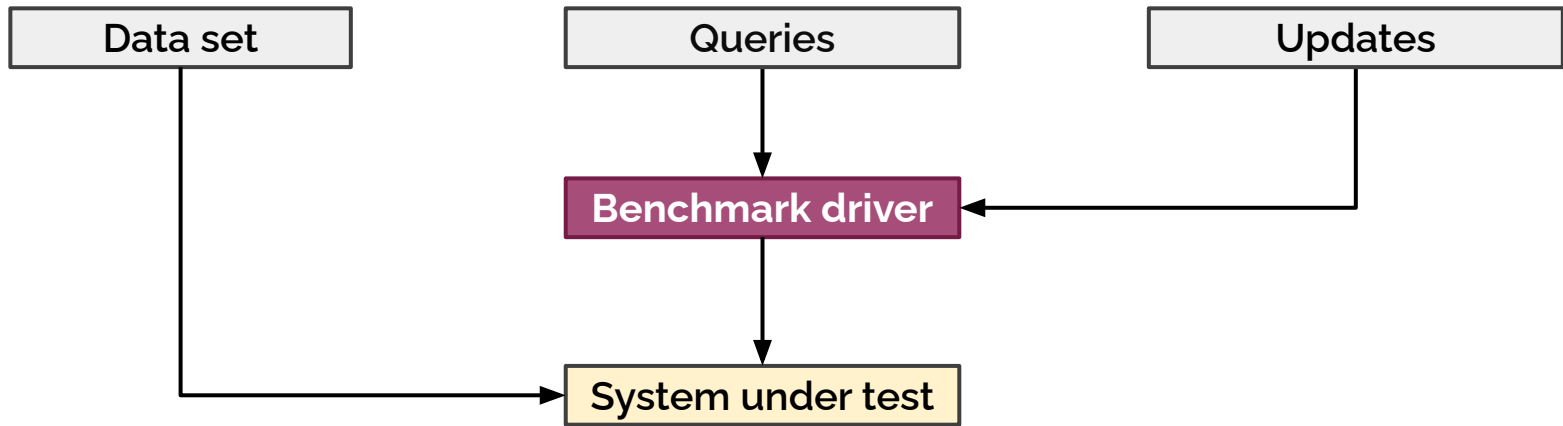
Cascading deletes remove lots of entities:

- have a big impact on the data distribution
- affect the executability of operations
- influence parameter selection

For databases, deletes:

- prohibit append-only data structures
- stress the garbage collector





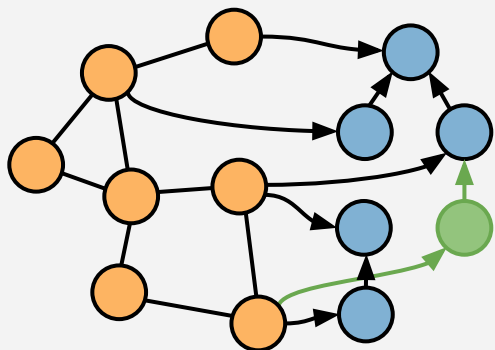
- **Schedules operations to be executable**  
(hard: needs careful parameter selection and dependency tracking)
- **Runs queries and updates concurrently**  
(hard: needs partitioned updates)
- **Collects benchmark results and performs validation**  
(very hard due to concurrent updates: we perform it sequentially)

# **SNB Workloads**

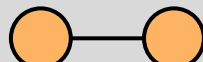




## SNB Interactive v1 (2015)



Q9(\$name, \$day)



name =  
\$name

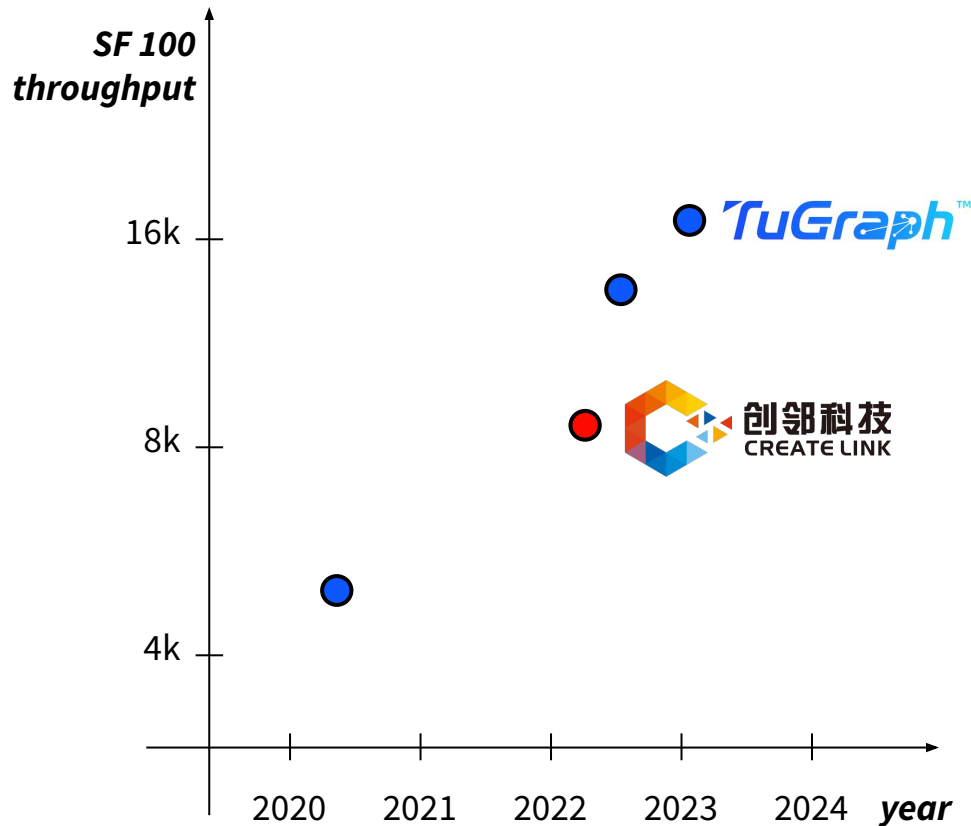
creation date <  
\$day

Queries start in 1-2 person nodes

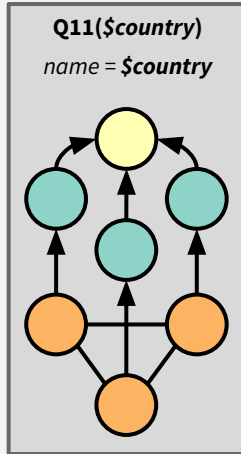
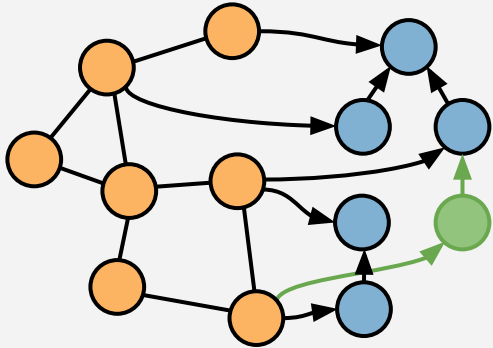
Concurrent inserts (no deletes)

Goal: high throughput (ops/s)

## Results on the 100GB data set



## SNB Business Intelligence (2023)



Queries touch on large portions of the data

Both bulk and concurrent updates allowed

Goal: high throughput & low query runtimes

## Audited results



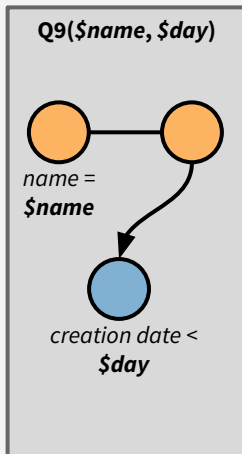
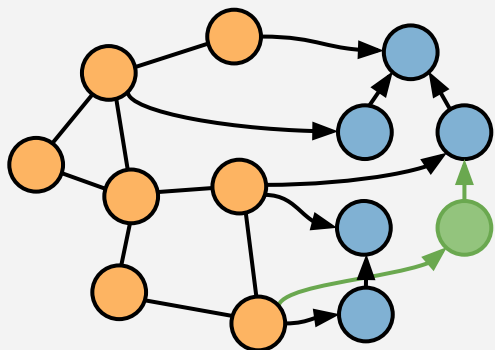
Results for 100GB, 1TB, and 10TB

10TB:

- Power@SF: 89,444
- Throughput@SF: 30,990

More results expected in 2023

## SNB Interactive v2



Queries start in 1-2 person nodes

Concurrent inserts and deletes

Goal: high throughput (ops/s)

## Features backported from BI:

- delete operations
- larger scale factors up to SF30,000
- cheapest path query

## New parameter generation features:

- temporal bucketing for each day
- path curation

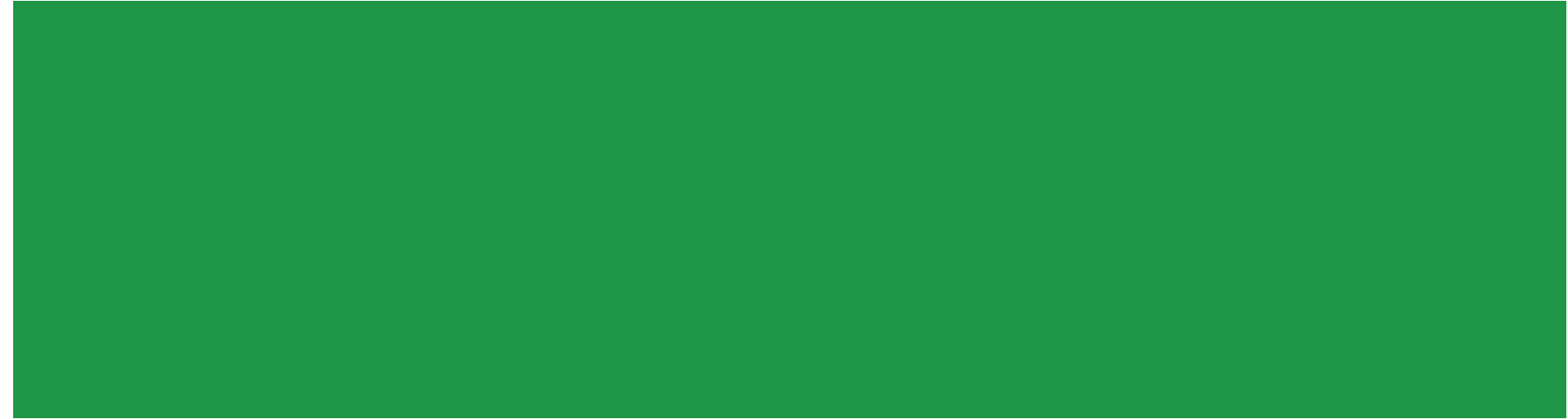
# The LDBC Social Network Benchmark Interactive Workload v2: A Transactional Graph Query Benchmark with Deep Delete Operations

David Püroja<sup>1</sup>, Jack Waudby<sup>2</sup>, Peter Boncz<sup>1</sup>, and Gábor Szárnyas<sup>1</sup>

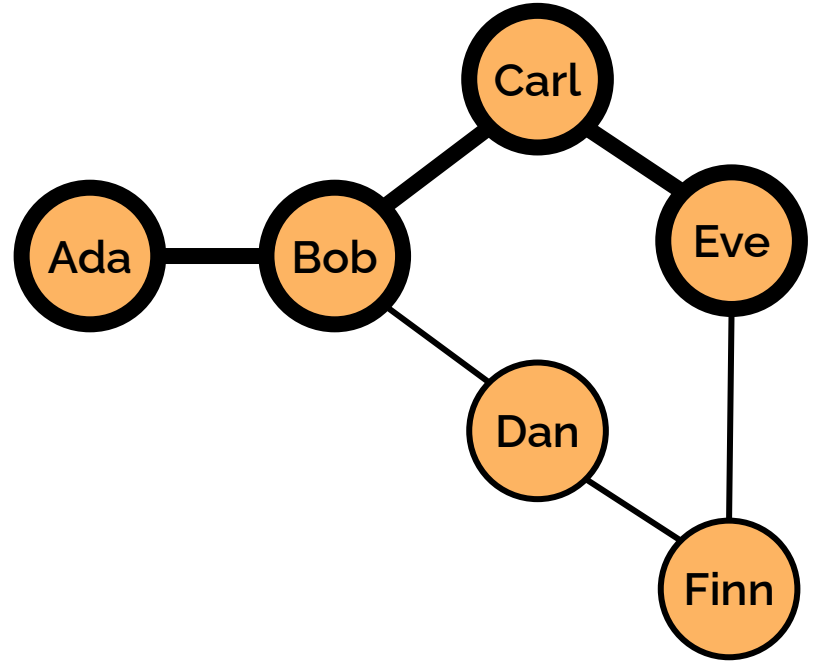
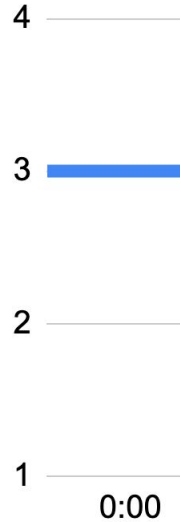
<sup>1</sup> CWI, the Netherlands, <sup>2</sup> Newcastle University, School of Computing  
david.puroja@ldbouncil.org, j.waudby2@newcastle.ac.uk, boncz@cw.nl,  
gabor.szarnyas@ldbouncil.org

**Abstract.** The LDBC Social Network Benchmark’s Interactive workload captures an OLTP scenario operating on a correlated social network graph. It consists of complex graph queries executed concurrently with a stream of updates operation. Since its initial release in 2015, the Interactive workload has become the de facto industry standard for benchmarking transactional graph data management systems. As graph systems have matured and the community’s understanding of graph processing features has evolved, we initiated the renewal of this benchmark. This paper describes the Interactive v2 workload with several new features: delete operations, a cheapest path-finding query, support for larger data sets, and a novel temporal parameter curation algorithm that ensures stable runtimes for path queries.

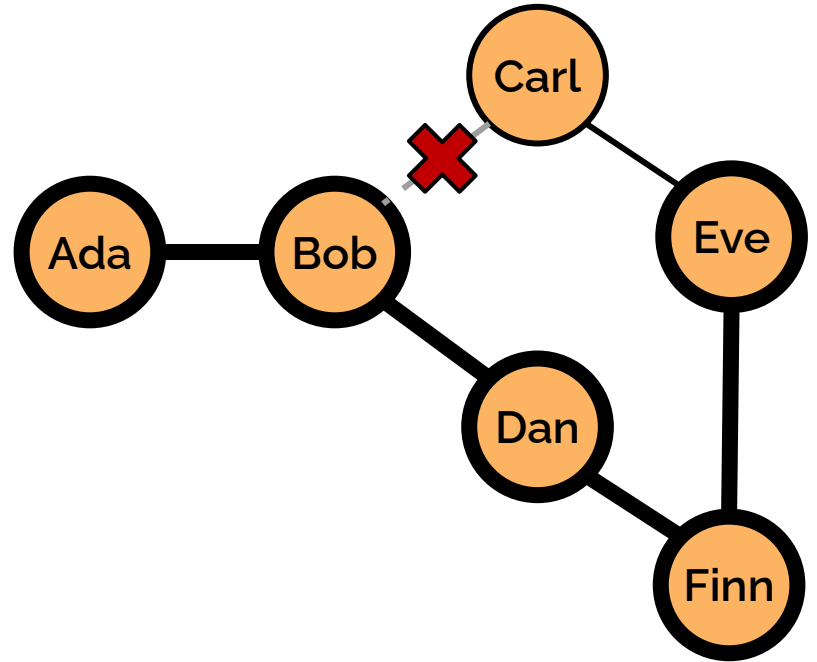
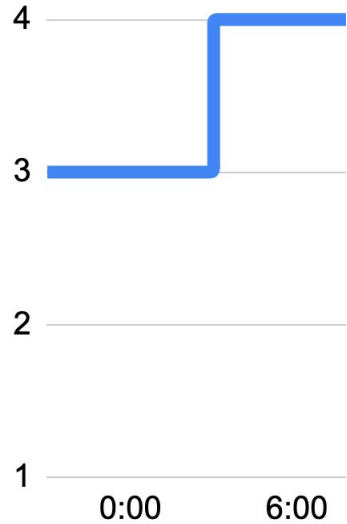
# Path curation



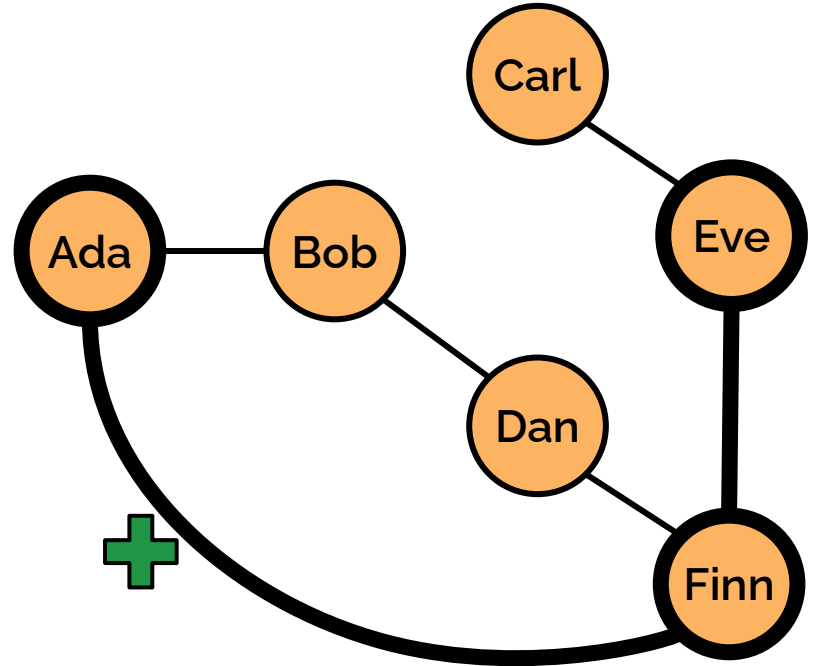
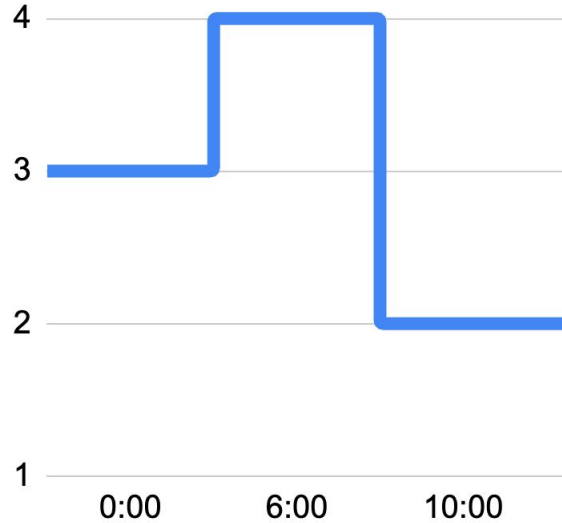
# Shortest distance from “Ada” to “Eve”



# Shortest distance from “Ada” to “Eve”

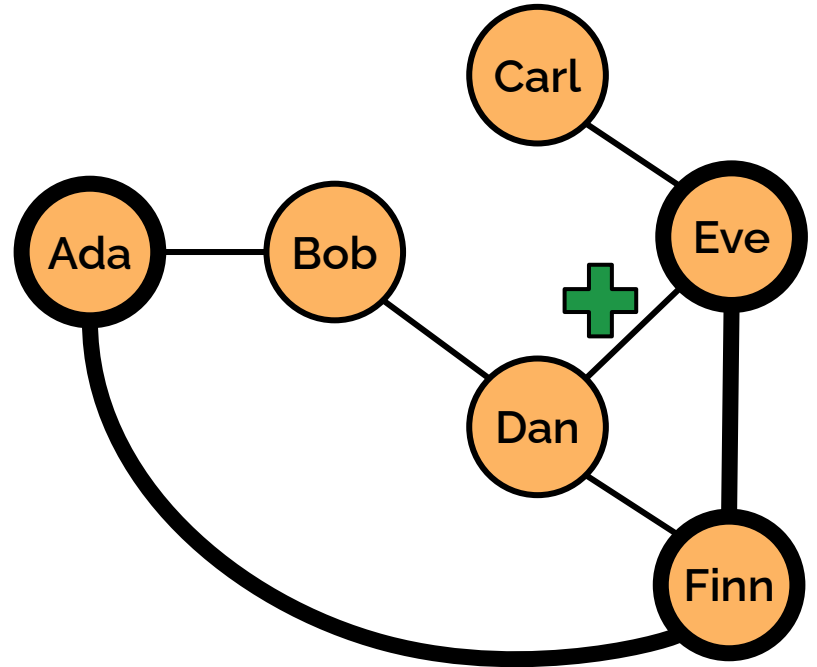
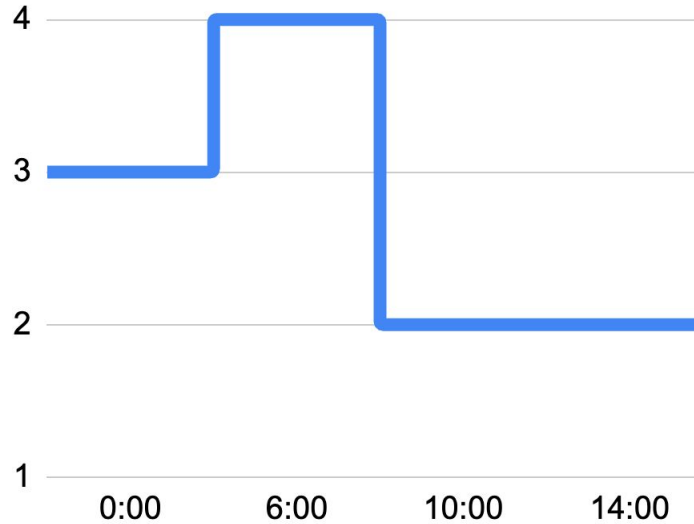


# Shortest distance from “Ada” to “Eve”

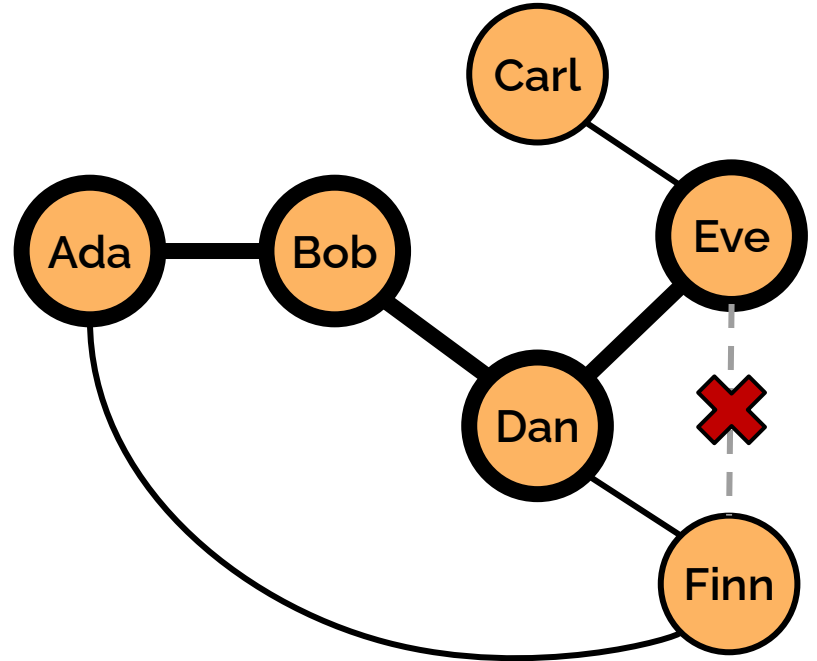
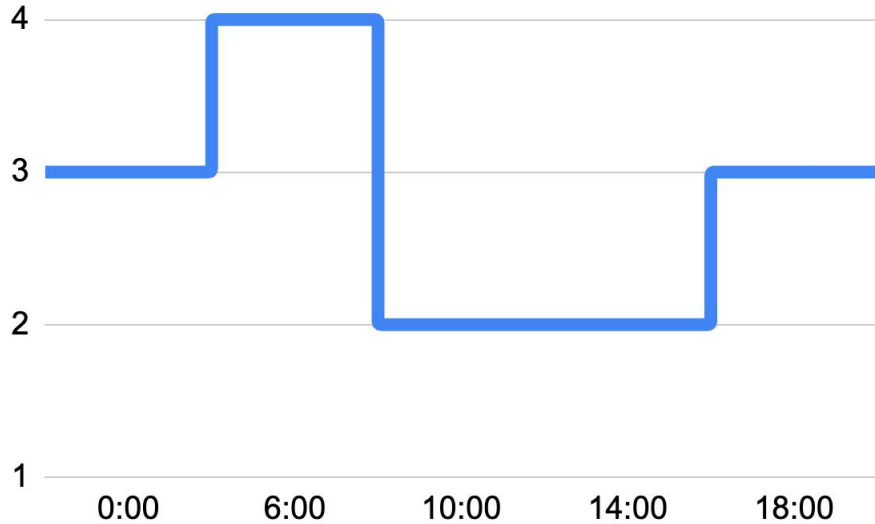




# Shortest distance from “Ada” to “Eve”



# Shortest distance from “Ada” to “Eve”



The shortest path distance changes multiple times during the day.

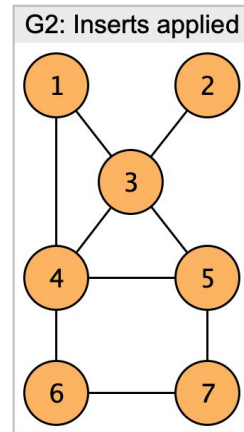
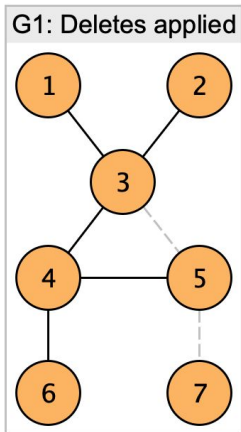
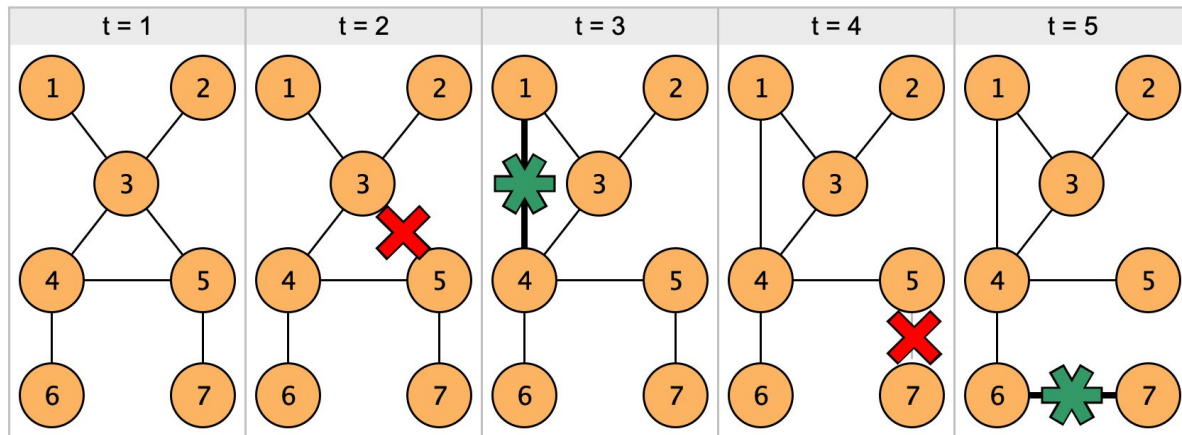
# Path curation with temporal bucketing

For each day, we construct:

**G1** – deletes but no inserts, setting an *upper* bound

**G2** – inserts but no deletes, setting a *lower* bound

$lower \leq actual\ length \leq upper$



Pairs of nodes yielding 3-hop paths in G1 and G2:

- ~~● — 1 to 5~~
- ~~● — 1 to 6~~
- ~~● — 2 to 5~~
- 2 to 6

# Is path curation sufficient?

Not yet:

- We also have to consider the degree distribution of the source–target nodes.

# Is path curation sufficient?

Not yet:

- We also have to consider the degree distribution of the source–target nodes.

Actually:

- For “perfect” parameter curation, we would need to run the entire workload with many parameter candidates and only keep ones which showed a similar behaviour

# Is path curation sufficient?

Not yet:

- We also have to consider the degree distribution of the source–target nodes.

Actually:

- For “perfect” parameter curation, we would need to run the entire workload with many parameter candidates and only keep ones which showed a similar behaviour

The real question:

- Is it worth spending more effort on optimizing the parameter curation?

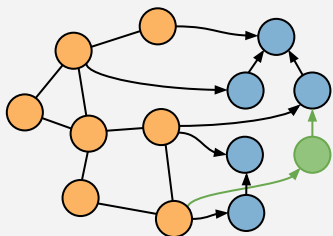
# I'm leaving academia

- Moving to DuckDB Labs (CWI spin-off in Amsterdam)
- Staying involved with LDBC at ~1 day / month

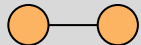


**DuckDB Labs**

## SNB Interactive



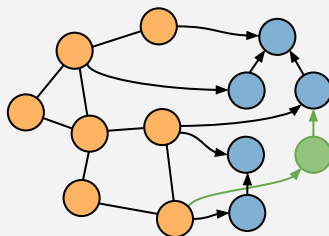
Q9(\$name, \$day)



name =  
\$name

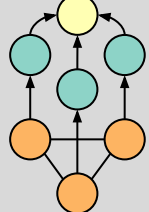
creation date <  
\$day

## SNB Business Intelligence

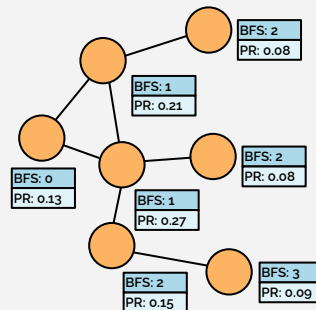


Q11(\$country)

name = \$country



## Graphalytics



Algorithms

|     |      |
|-----|------|
| BFS | CDLP |
| PR  | SSSP |
| LCC | WCC  |

Data sets

|            |
|------------|
| LDBC SNB   |
| Graph500   |
| Twitter    |
| Friendster |
| Patents    |
| wiki-Talk  |

## Semantic Publishing Benchmark

Target: RDF/SPARQL

Domain: Media/publishing industry

Inferencing & continuous updates

## Financial Benchmark

Target: Distributed systems

Domain: Financial fraud detection

Strict latency bound (20 ms)

## Future benchmark ideas

GNNs

Graph mining

Graph streaming



***LDBC*** 

*The graph & RDF  
benchmark reference*