



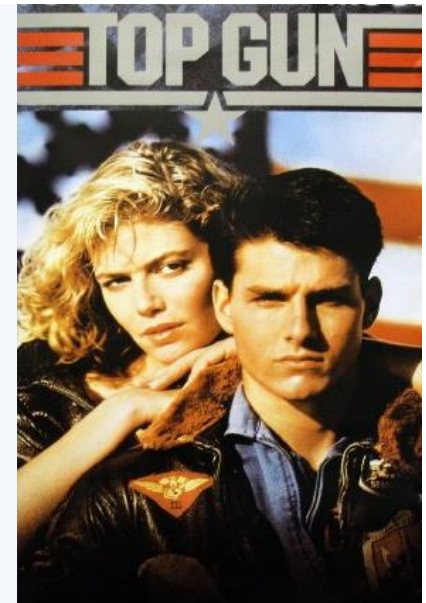
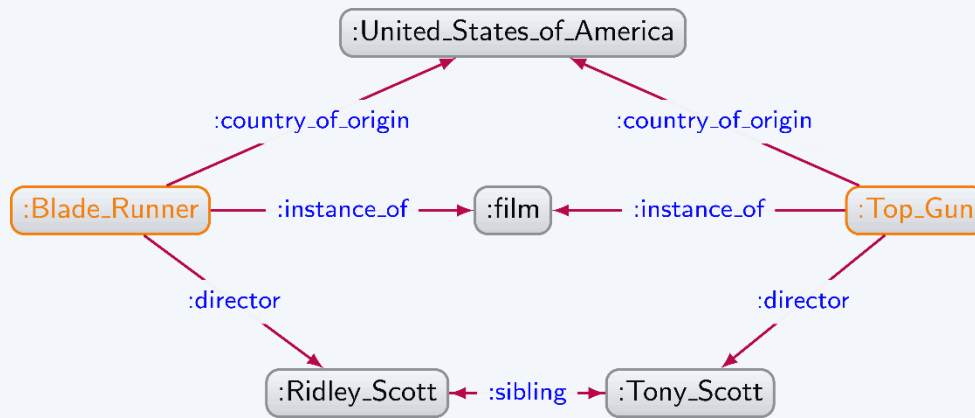
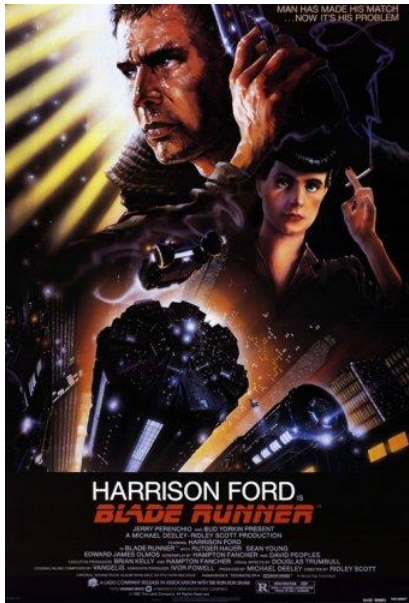
# WISP: WEIGHTED SHORTEST PATHS FOR RDF GRAPHS

Gonzalo Tartari, Aidan Hogan  
DCC, Universidad de Chile

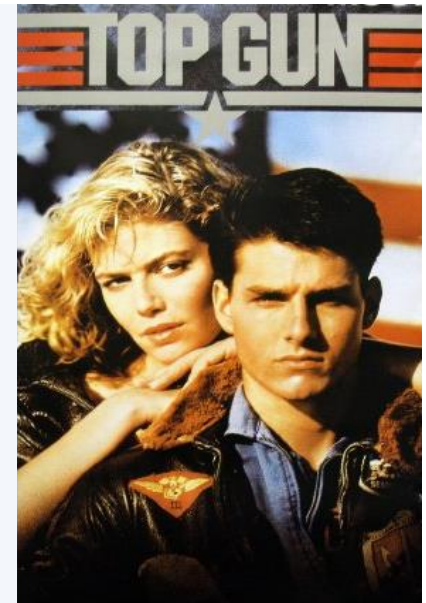
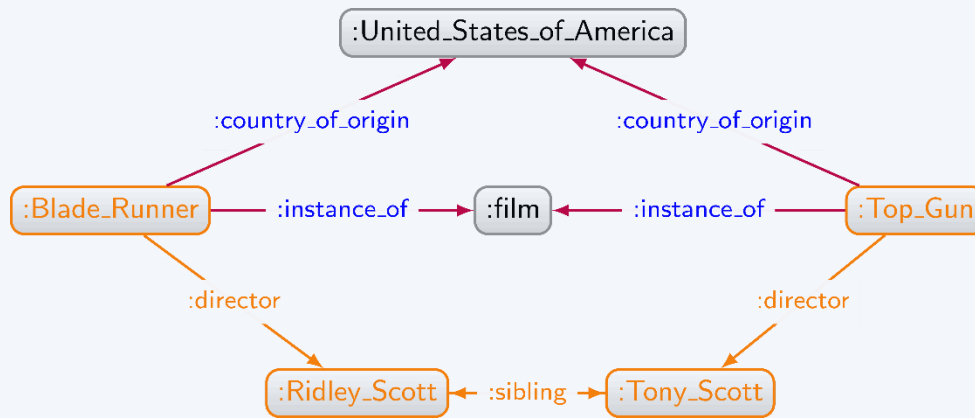
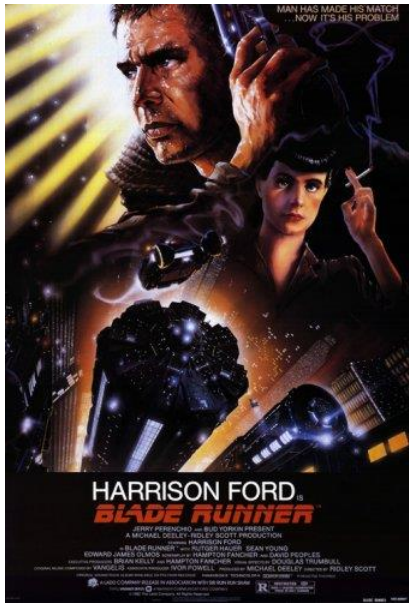


Fundamentos  
de los datos

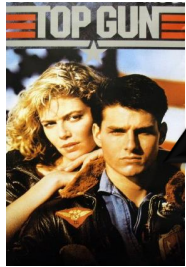
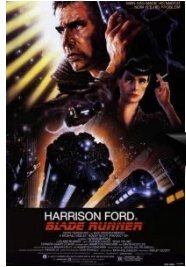
# "INTERESTING" PATHS = SHORTEST PATHS?



# "INTERESTING" PATHS $\neq$ SHORTEST PATHS!



# (MANY OF THE) EXISTING APPROACHES



Enumerate  
Paths

Score  
Paths

Order/Filter  
Paths

Output

# (MANY OF THE) EXISTING APPROACHES



- $n!$  simple paths (but only a few of interest)

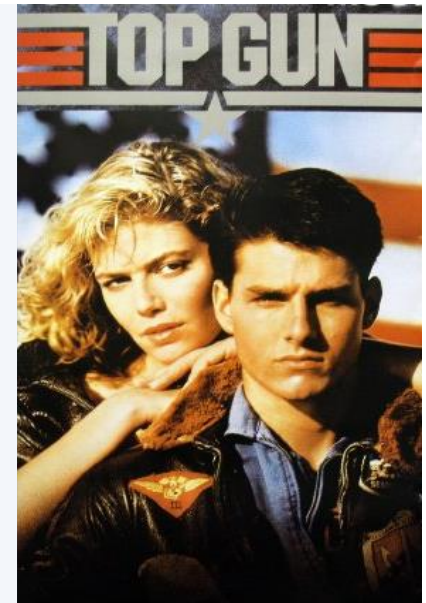
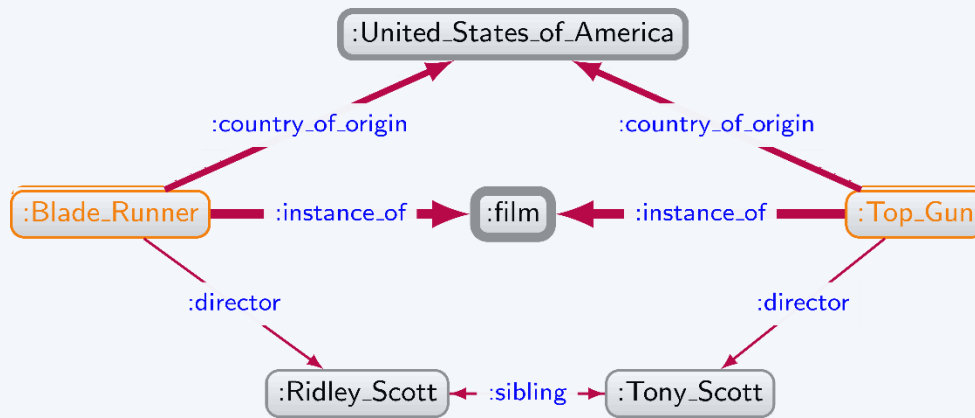
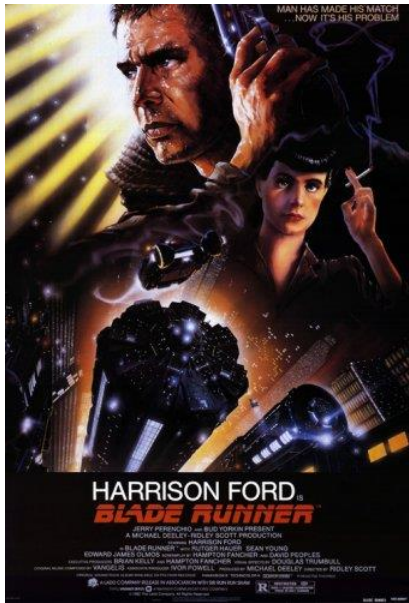
# (MANY OF THE) EXISTING APPROACHES



- $n!$  simple paths (but only a few of interest)
  - Paths thus often bounded in length

# OUR APPROACH: WEIGHT GRAPHS

- *Idea:* “Common” nodes/edges get weighted higher

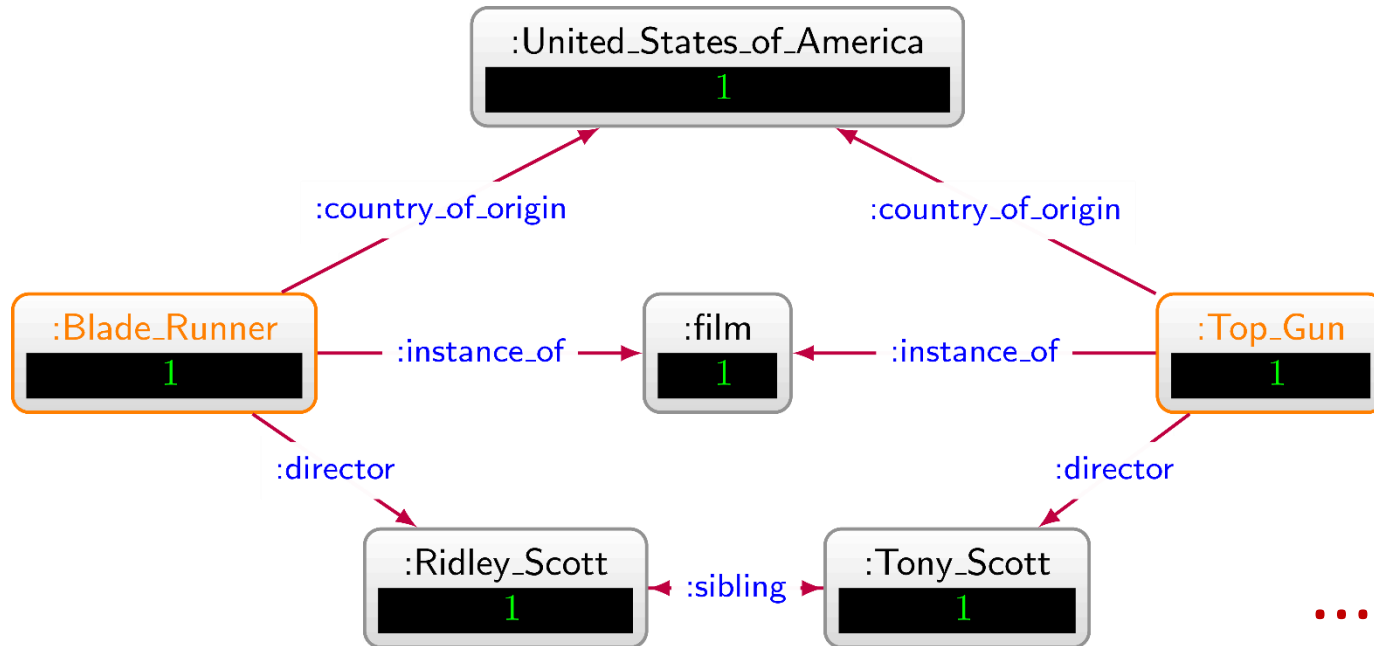


- *Benefit:* Can find “interesting” paths directly

# WEIGHTING GRAPHS: NODES

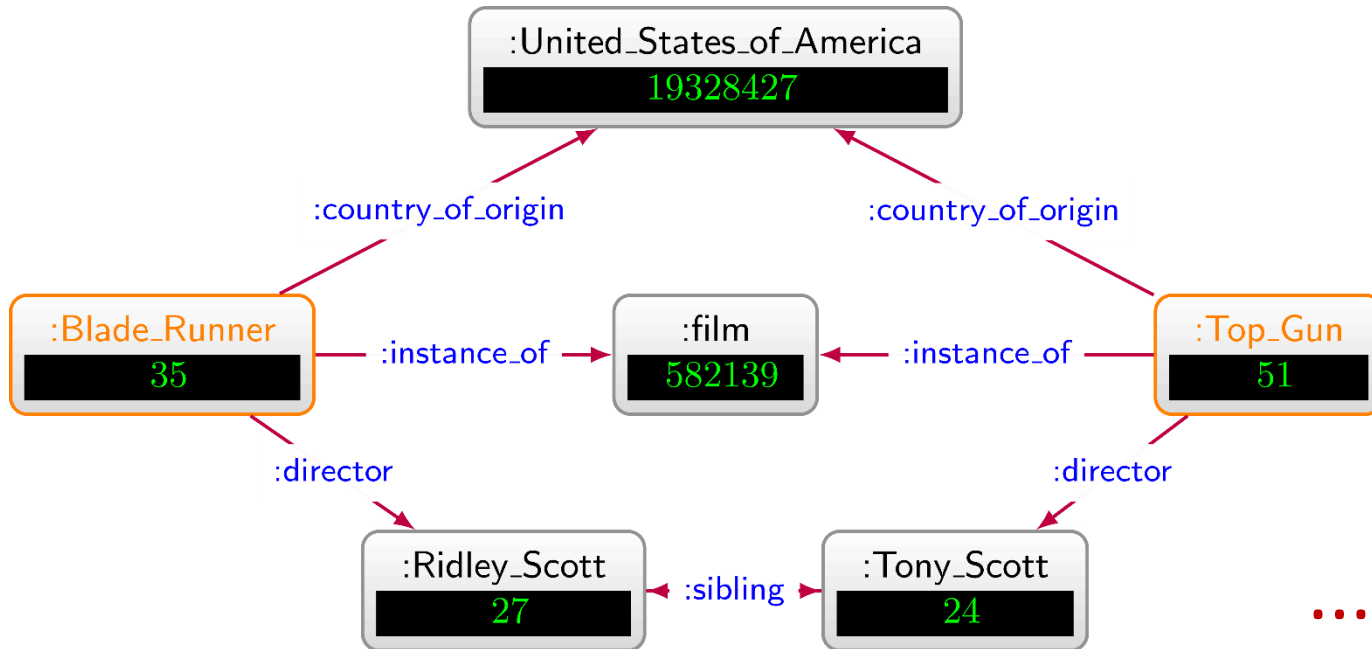


# NODE WEIGHTS: LENGTH (BASELINE)



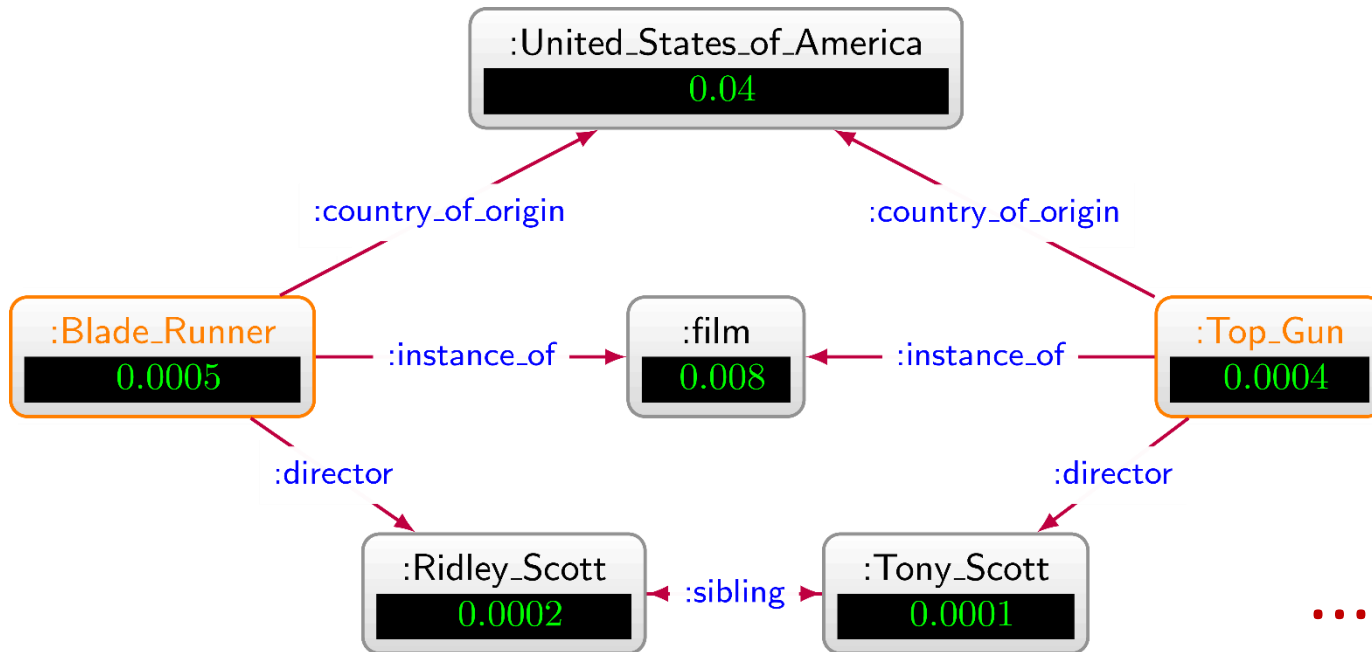
path					score
:country	:U.S.	:country <sup>-</sup>			1
:instance	:film	:instance <sup>-</sup>			1
:director	:RScott	:sibling	:TScott	:director <sup>-</sup>	2
:director	:RScott	:sibling <sup>-</sup>	:TScott	:director <sup>-</sup>	2

# NODE WEIGHTS: DEGREE



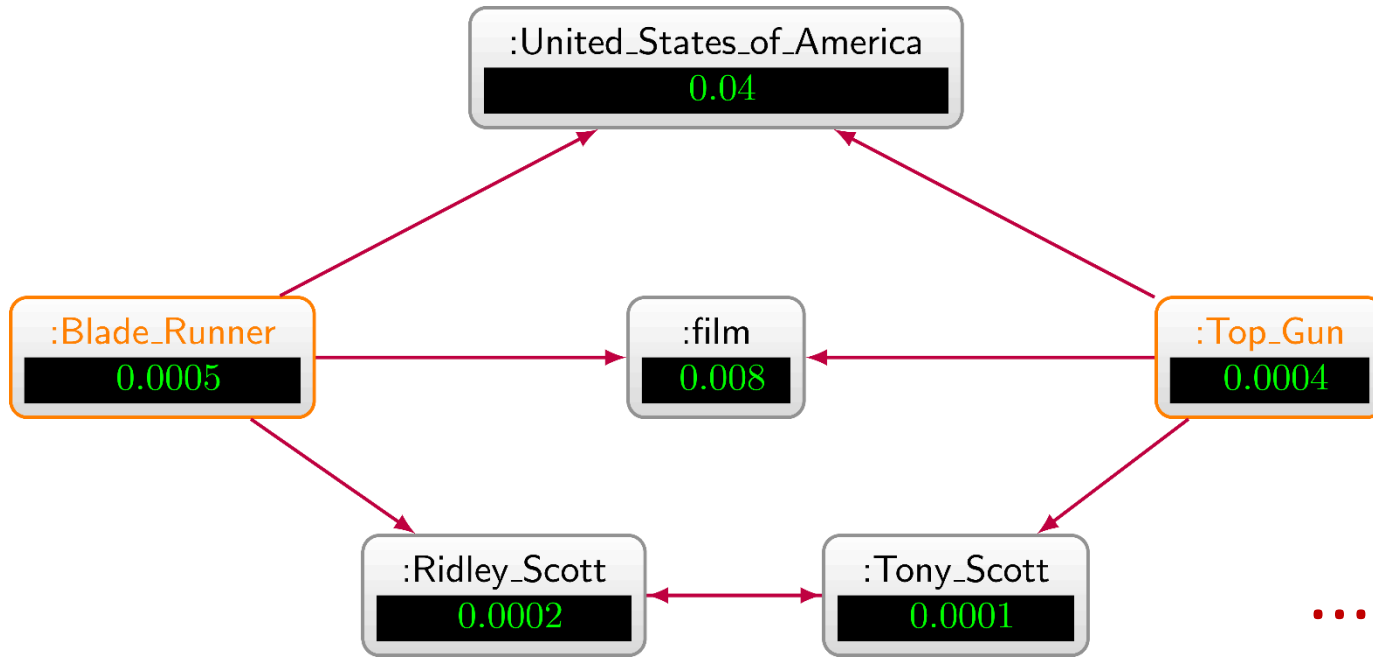
path					score
:director	:RScott	:sibling	:TScott	:director <sup>-</sup>	51
:director	:RScott	:sibling <sup>-</sup>	:TScott	:director <sup>-</sup>	51
:instance	:film	:instance <sup>-</sup>			582139
:country	:U.S.	:country <sup>-</sup>			19328427

# NODE WEIGHTS: PAGERANK



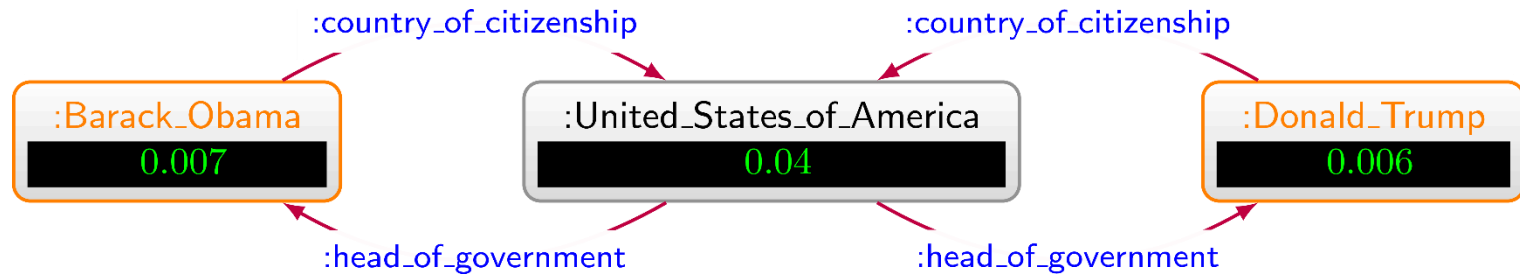
path	score
<code>:director</code> :RScott <code>:sibling</code> :TScott <code>:director</code> <sup>-</sup>	0.0003
<code>:director</code> :RScott <code>:sibling</code> <sup>-</sup> :TScott <code>:director</code> <sup>-</sup>	0.0003
<code>:instance</code> :film <code>:instance</code> <sup>-</sup>	0.008
<code>:country</code> :U.S. <code>:country</code> <sup>-</sup>	0.04

# ASIDE: PAGERANK / DIRECTED GRAPH USED



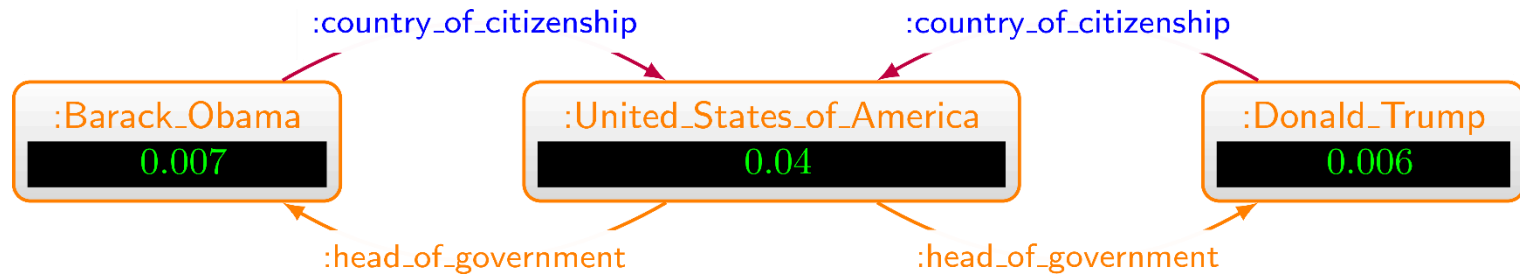
# WEIGHTING GRAPHS: EDGES

# WEIGHTING WITH ONLY NODES



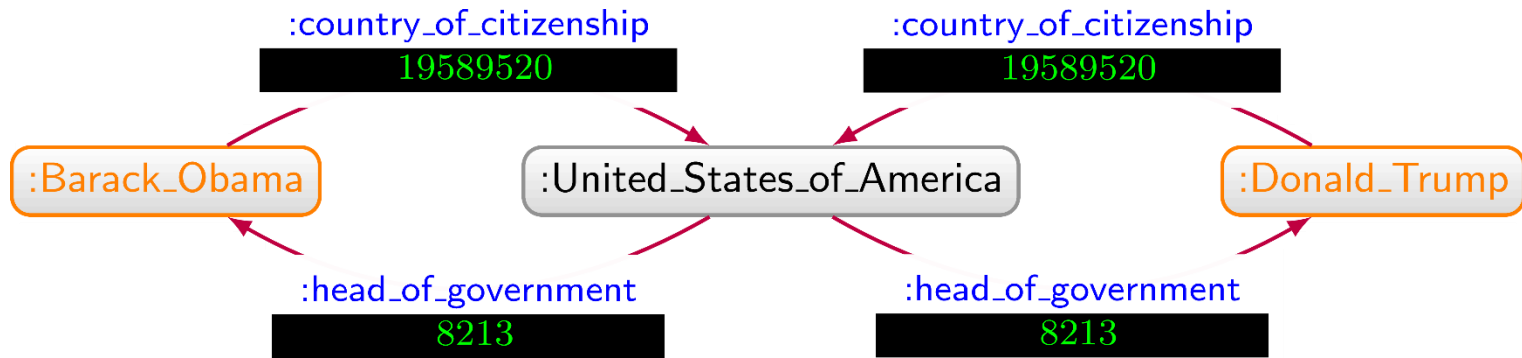
path	score
<b>:country</b> :U.S. <b>:country</b> <sup>-</sup>	0.04
<b>:country</b> :U.S. <b>:head</b> <sup>-</sup>	0.04
<b>:head</b> :U.S. <b>:country</b> <sup>-</sup>	0.04
<b>:head</b> :U.S. <b>:head</b> <sup>-</sup>	0.04

# WEIGHTING WITH ONLY NODES



path	score
<b>:country</b> :U.S. <b>:country</b> <sup>-</sup>	0.04
<b>:country</b> :U.S. <b>:head</b> <sup>-</sup>	0.04
<b>:head</b> :U.S. <b>:country</b> <sup>-</sup>	0.04
<b>:head</b> :U.S. <b>:head</b> <sup>-</sup>	<b>0.04</b>

# EDGE WEIGHTS: FREQUENCY

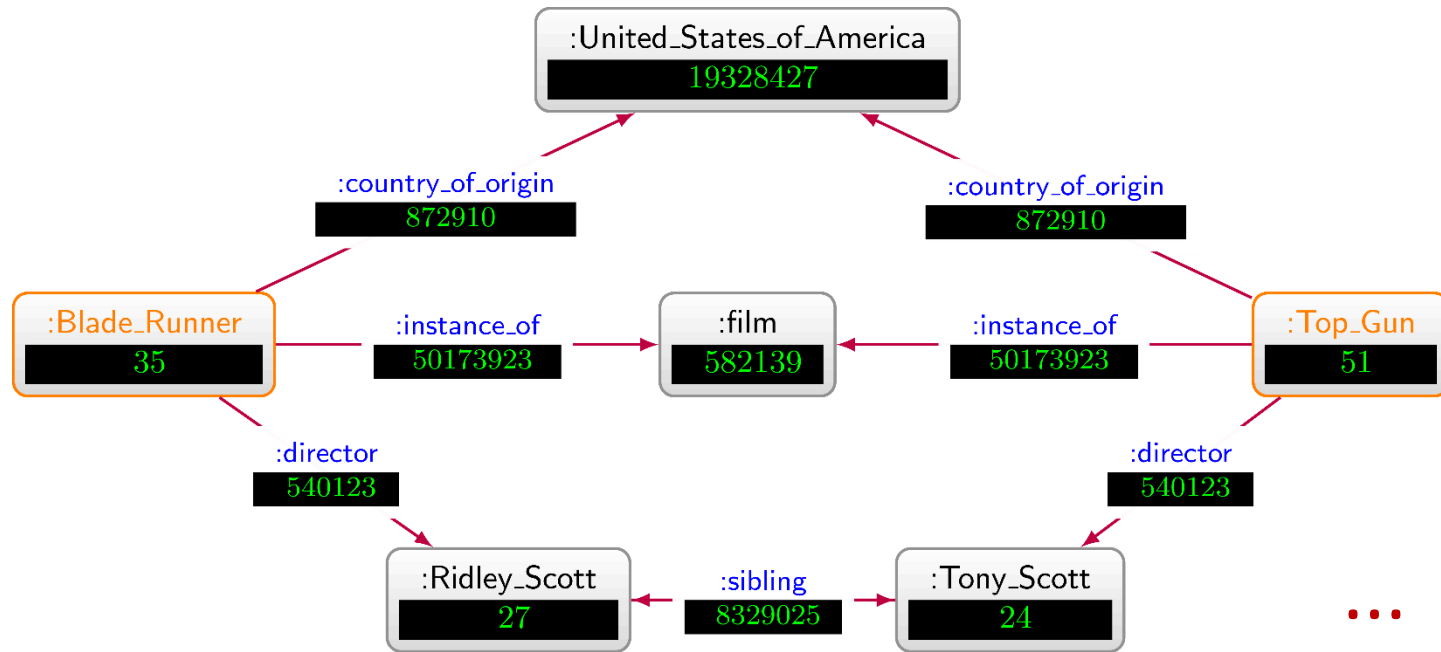


path			score
:head	:U.S.	:head <sup>-</sup>	16426
:country	:U.S.	:head <sup>-</sup>	19597733
:head	:U.S.	:country <sup>-</sup>	19597733
:country	:U.S.	:country <sup>-</sup>	39179040



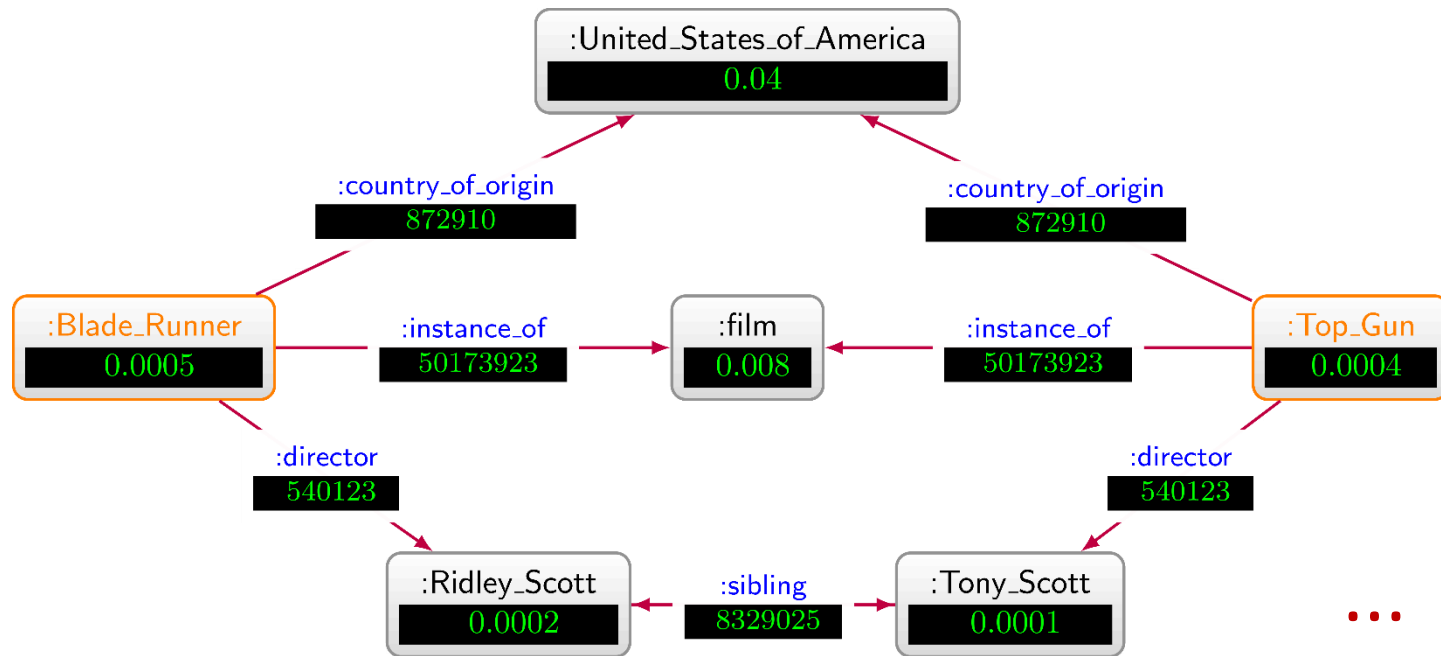
# WEIGHTING GRAPHS: NODES + EDGES

# NODE + EDGE WEIGHTS: DEGREE + FREQUENCY



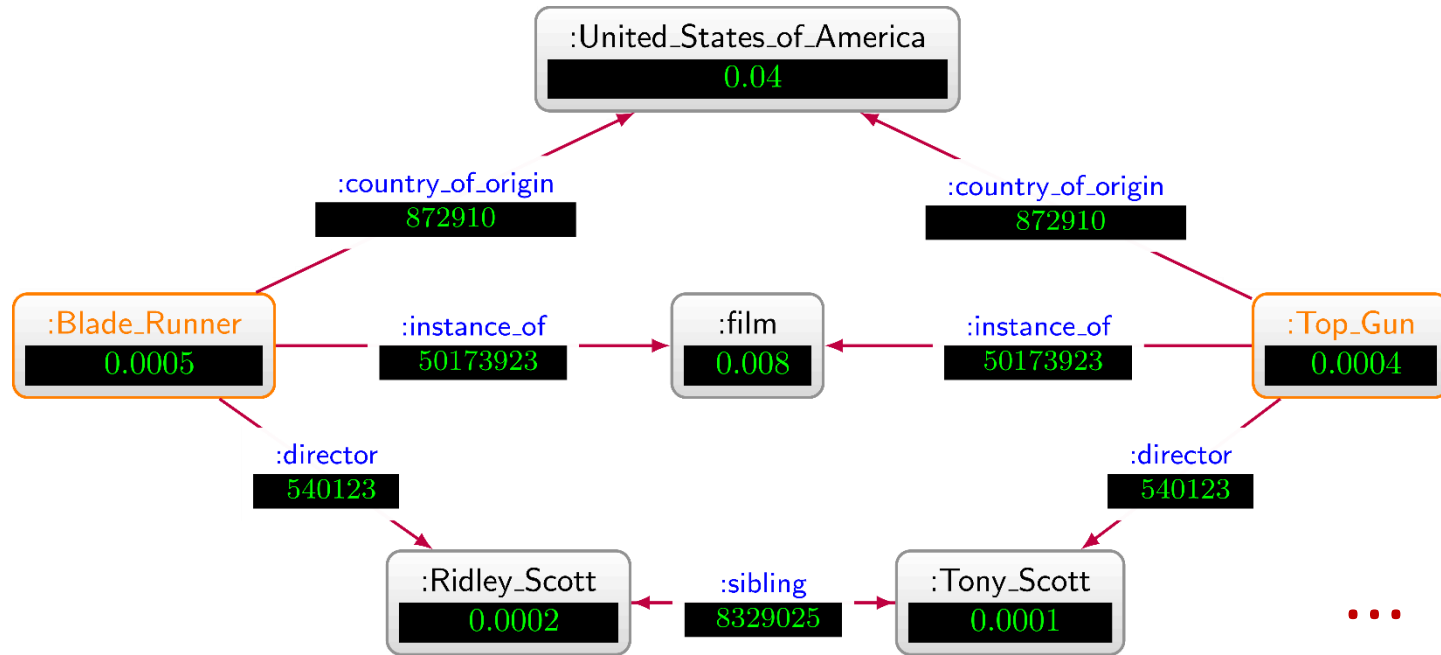
path	score
:director :RScott :sibling :TScott :director <sup>-</sup>	9409322
:director :RScott :sibling <sup>-</sup> :TScott :director <sup>-</sup>	9409322
:country :U.S. :country <sup>-</sup>	21074247
:instance :film :instance <sup>-</sup>	100929984

# NODE + EDGE WEIGHTS: PAGERANK + FREQUENCY



path					score
:country	:U.S.	:country <sup>-</sup>			1745820.04
:director	:RScott	:sibling	:TScott	:director <sup>-</sup>	9409271.0003
:director	:RScott	:sibling <sup>-</sup>	:TScott	:director <sup>-</sup>	9409271.0003
:instance	:film	:instance <sup>-</sup>			100347846.008

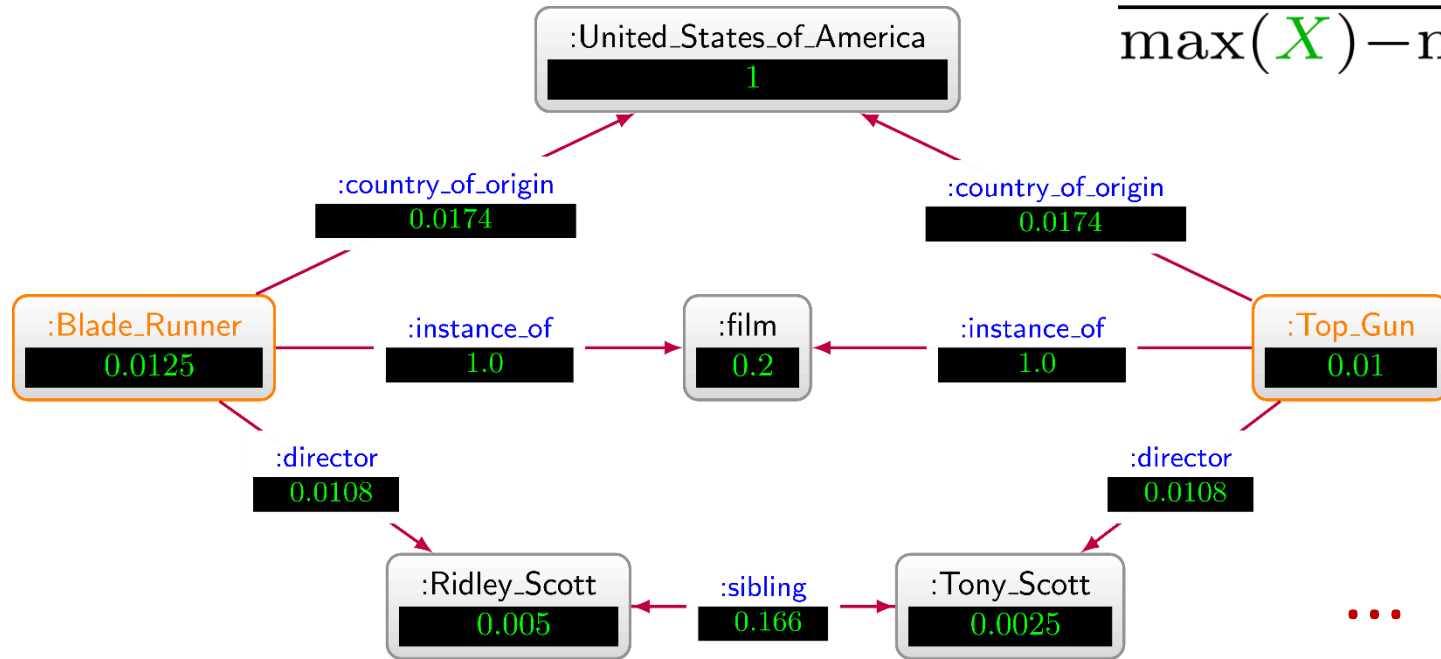
# NODE + EDGE WEIGHTS: PAGERANK + FREQUENCY



path					score
:country	:U.S.	:country <sup>-</sup>			1745820.04
:director	:RScott	:sibling	:TScott	:director <sup>-</sup>	9409271.0003
:director	:RScott	:sibling <sup>-</sup>	:TScott	:director <sup>-</sup>	9409271.0003
:instance	:film	:instance <sup>-</sup>			100347846.008

# NODE + EDGE WEIGHTS: [0,1] NORMALISATION

$$\frac{x - \min(X)}{\max(X) - \min(X)}$$

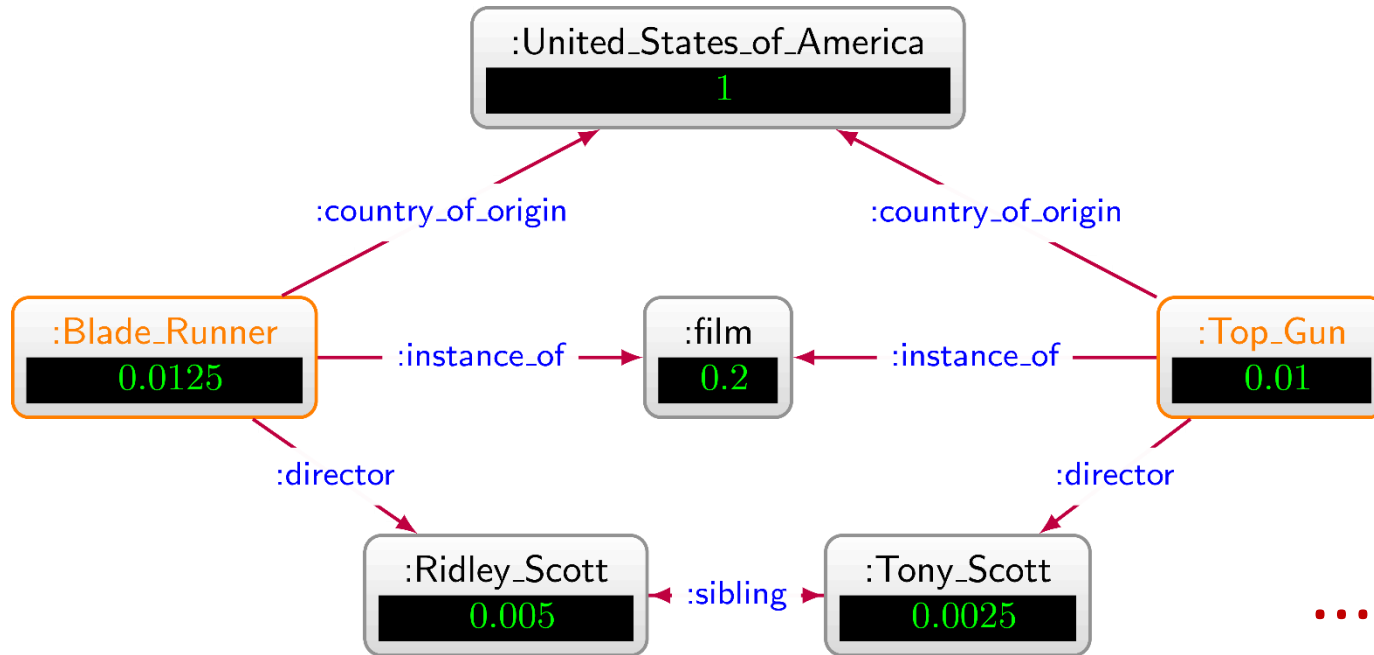


path	score
:director :RScott :sibling :TScott :director <sup>-</sup>	0.195033
:director :RScott :sibling <sup>-</sup> :TScott :director <sup>-</sup>	0.195033
:country :U.S. :country <sup>-</sup>	1.02153
:instance :film :instance <sup>-</sup>	2.2

# HYBRID NODE WEIGHTS

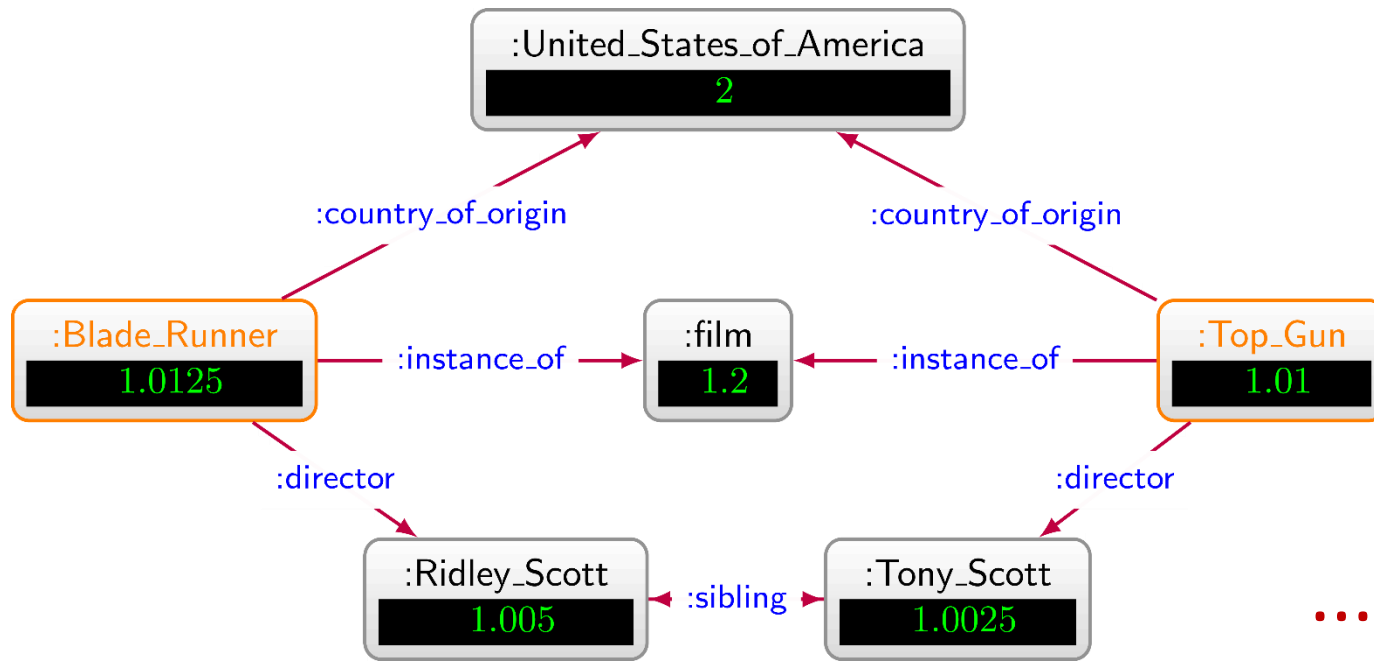
# NODE WEIGHTS: PAGERANK

Visiting one high-centrality node = Visiting thousands of low-centrality nodes



path	score
<code>:director :RScott :sibling :TScott :director</code>	0.0075
<code>:director :RScott :sibling :TScott :director</code>	0.0075
<code>:instance :film :instance</code>	0.02
<code>:country :U.S. :country</code>	1

# HYBRID NODE WEIGHTS: PAGERANK + LENGTH



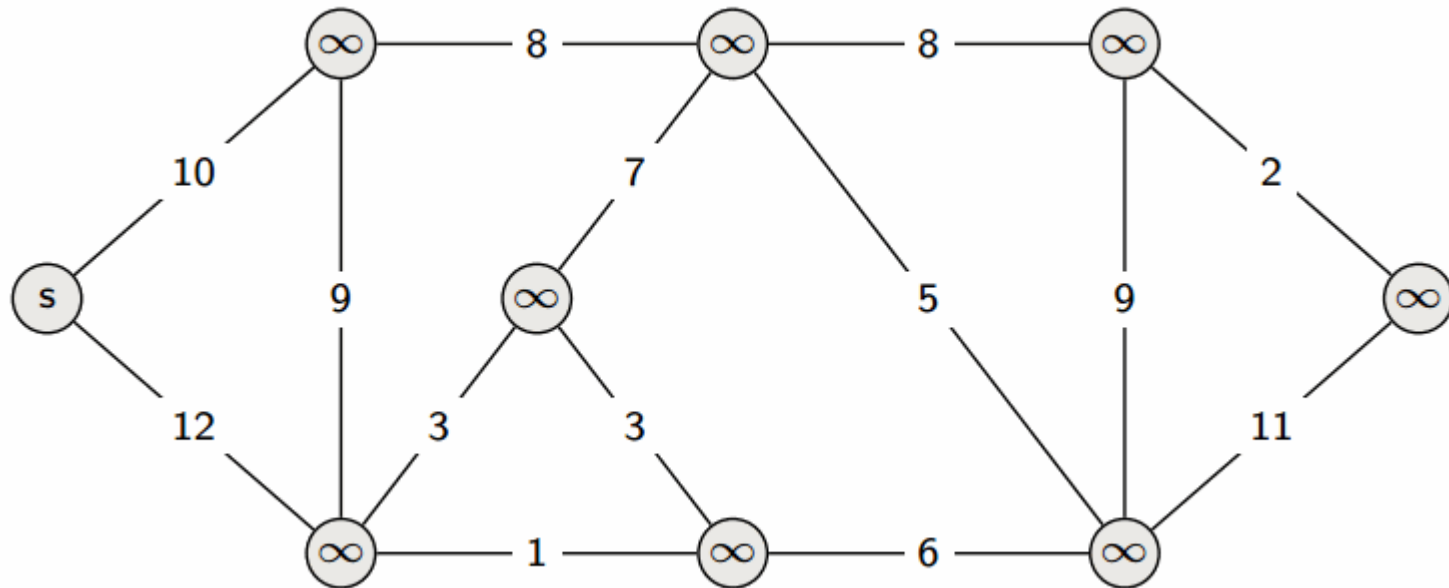
path					score
:instance	:film	:instance <sup>-</sup>			1.2
:country	:U.S.	:country <sup>-</sup>			2
:director	:RScott	:sibling	:TScott	:director <sup>-</sup>	2.0075
:director	:RScott	:sibling <sup>-</sup>	:TScott	:director <sup>-</sup>	2.0075



# IMPLEMENTATION

# WEIGHTED SHORTEST-PATH IMPLEMENTATION

- Dijkstra's algorithm:
  - Worst case:  $O(|E| + |V| \cdot \log|V|)$



# EXPERIMENTS

# QUESTIONS

- **Performance:**
  - How are the runtimes?
  - How is the scalability?
- **Weighting schemes:**
  - How similar are paths for different weightings?
  - Does weighting help find interesting paths?
  - Which weighting finds the most interesting paths?

# DATASET: WIKIDATA

- Truthy dump: 2017-06-07
  - 25 million nodes (Q-IRIs only)
  - 90 million edges



# DATASET: WIKIDATA SLICES

<b>Dataset</b>	1.6 M	3.2 M	6.4 M	12.8 M	FULL
<b>Nodes</b>	1,227,382	2,507,582	5,303,322	10,343,129	25,081,334
<b>Edges</b>	6,603,412	12,160,436	22,008,446	36,404,534	89,878,092



**WIKIDATA**

# MACHINE

- 2 x Intel Xeon Quad Core @1.9GHz
- 32 GB of RAM



# WEIGHTING SCHEMES

- Node
  - Degree (D)
  - PageRank (P)
  - Length (L)
- Node + Edge
  - Degree + Edge Frequency (DE)
  - PageRank + Edge Frequency (PE)
- Hybrid Node + Edge
  - Degree + Length + Edge Frequency (DEL)
  - PageRank + Length + Edge Frequency (PEL)



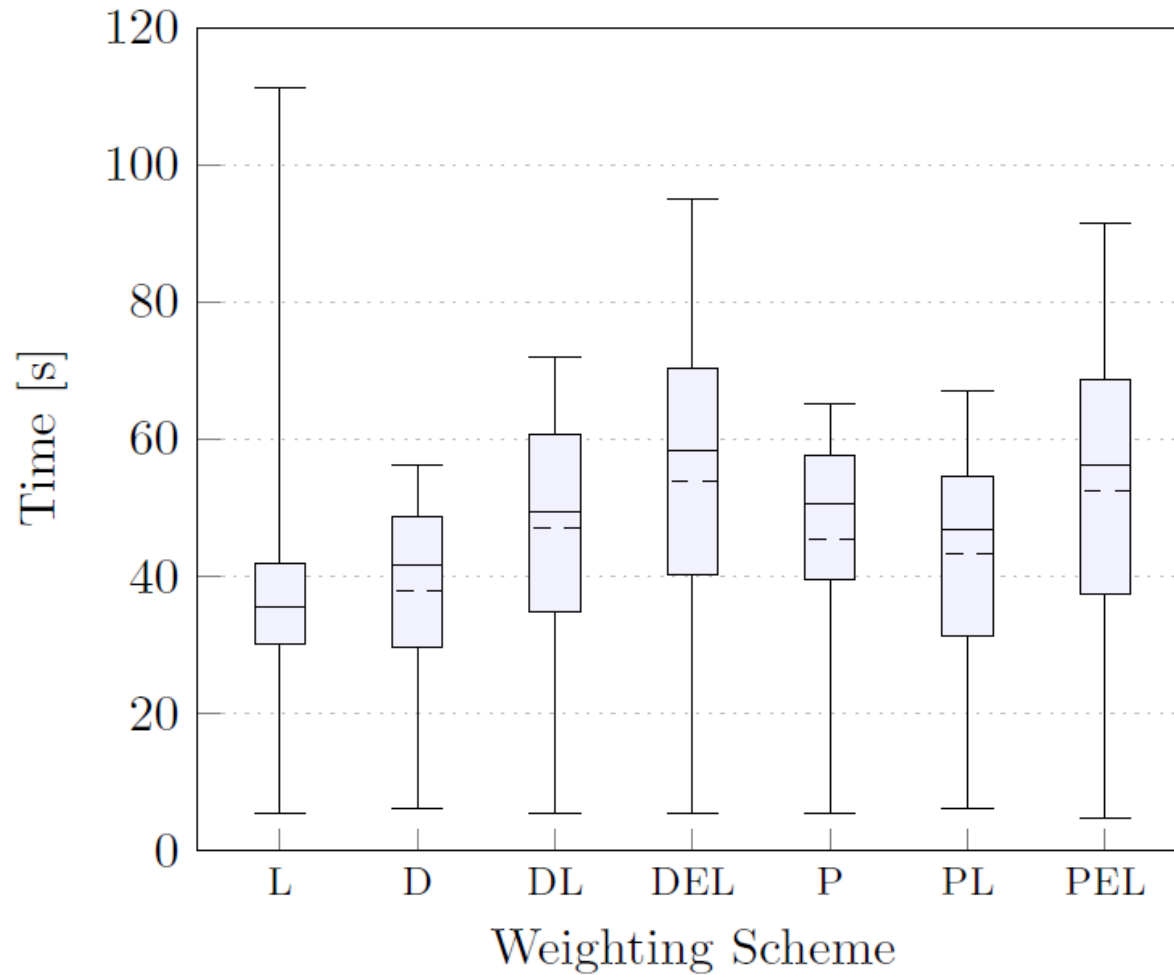
PERFORMANCE

# QUERIES (NODE PAIRS)

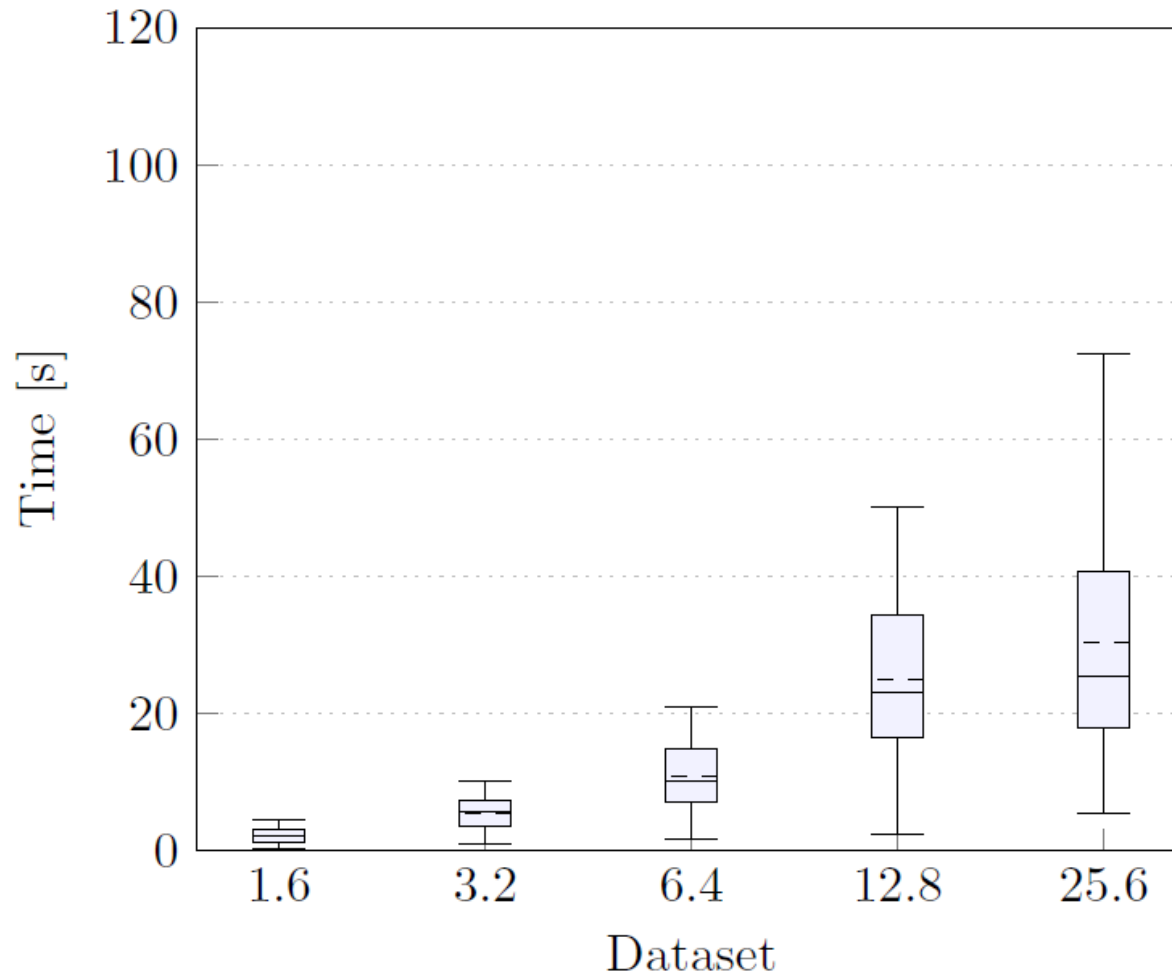
- **Queries:** 100 node pairs randomly sampled
  - From smallest slice (Q code < 100000)
  - From each slice independently
- **Task:** Return one (best) path



# PERFORMANCE RESULTS (FULL DATASET)

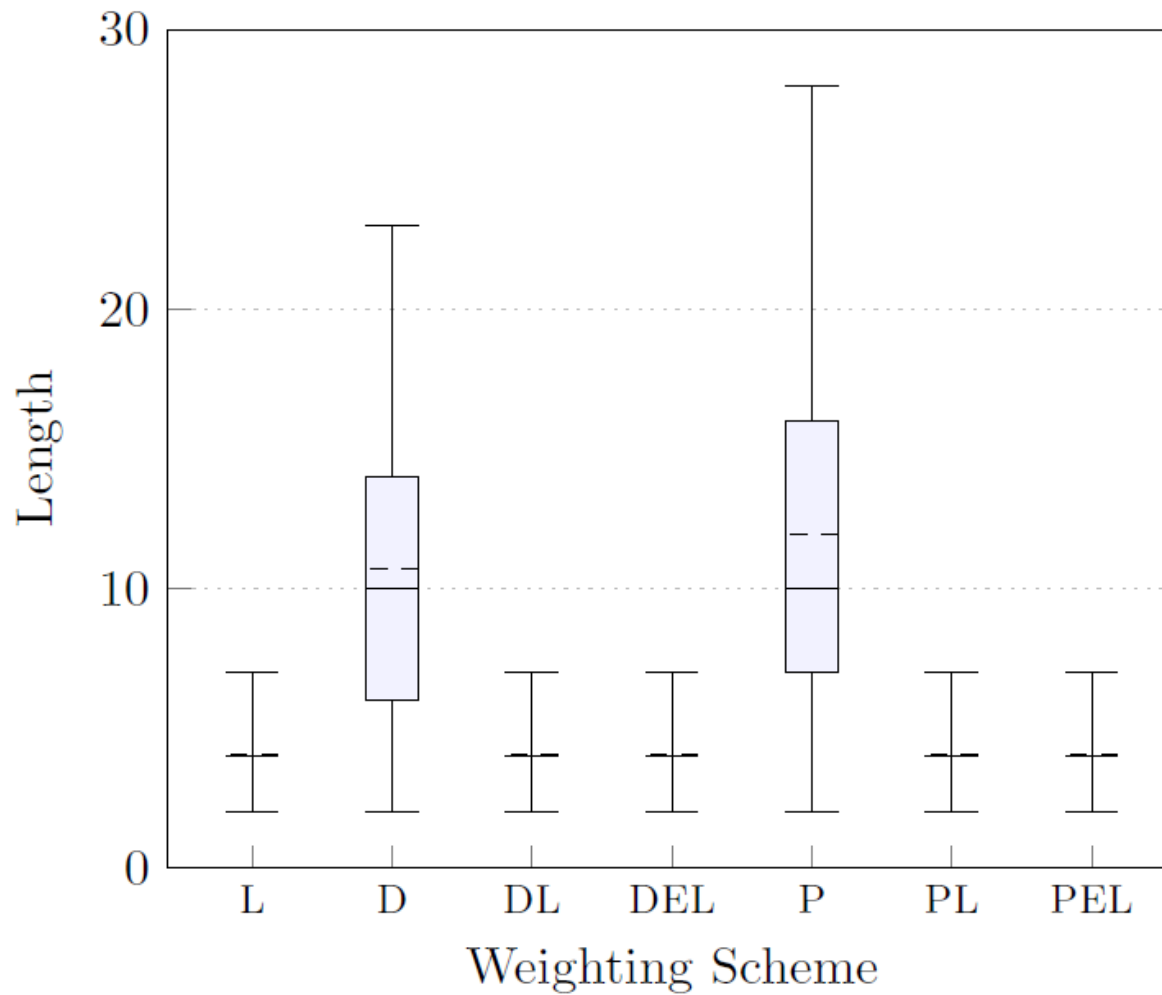


# PERFORMANCE RESULTS (DEL | VARIOUS SCALES)

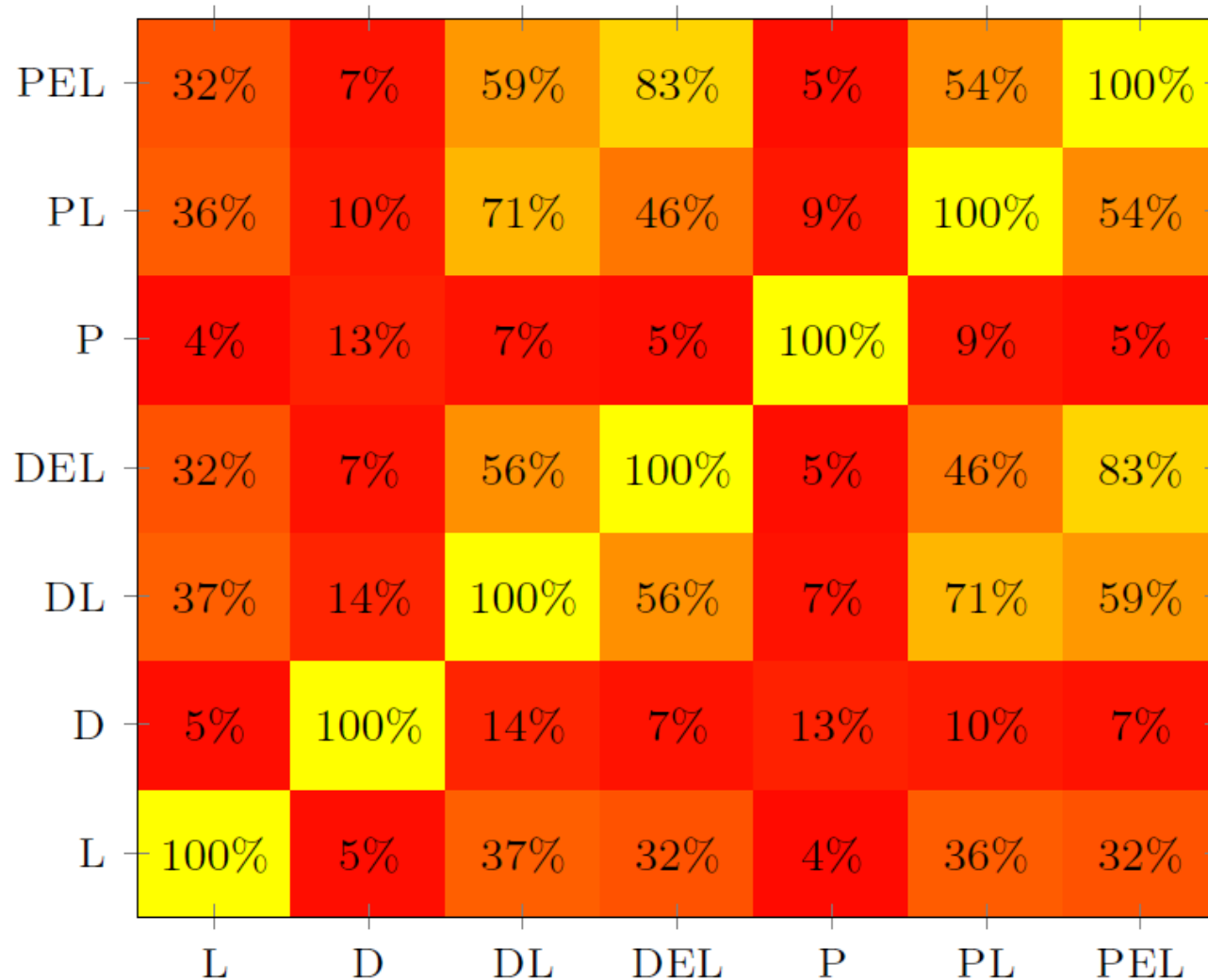


# COMPARISON OF WEIGHTING SCHEMES

# COMPARISON OF PATH LENGTH (FULL DATASET)



# HOW MANY PAIRS GIVE THE SAME PATH? (FULL DATASET)



# USER STUDY



# QUERIES: SAME TYPE

<b>Node</b>	<b>Q code</b>	<b>Node</b>	<b>Q code</b>
George W. Bush	Q207	Barack Obama	Q76
Bolivia	Q750	Portugal	Q45
Elvis Presley	Q303	Deep Purple	Q101505
Lionel Messi	Q615	Stephen Curry	Q352159
Eiffel Tower	Q243	Pisa Tower	Q39054
Alan Turing	Q7251	Leonhard Euler	Q7604
Samsung	Q20716	Apple Inc.	Q312
Ford Motor Company	Q44294	Volkswagen	Q246
Al Pacino	Q41163	Bill Murray	Q29250
CNN	Q48340	BBC	Q9531

# QUERIES: DIFFERENT TYPES

<b>Node</b>	<b>Q code</b>	<b>Node</b>	<b>Q code</b>
Donald Trump	Q22686	Pablo Neruda	Q34189
Steam (videogame platform)	Q337535	Barcelona Football Club	Q7156
Everest (mountain)	Q513	FIFA	Q253414
Lord of the Rings	Q15228	Apollo 11 (space mission)	Q43653
RMS Titanic	Q25173	Saturn	Q193
Netflix	Q907311	Napoleon Bonaparte	Q517
Second World War	Q362	Plato	Q859
The Beatles	Q1299	Star Wars	Q462
Odyssey	Q35160	Facebook (web site)	Q355
X Rays	Q34777	Amazon.com	Q3884

# USER STUDY

- 10 students
- 1.6 M dataset
- Shown all paths for one query together
- Scores: 1 (very poor) - 7 (very good)
- 79 complete evaluations
  - 4 evaluations per query (node pair)
  - 553 scores



# LOWEST-RATED PATH



mean score 1.25 ( {1,1,1,2} )

# HIGHEST-RATED PATH

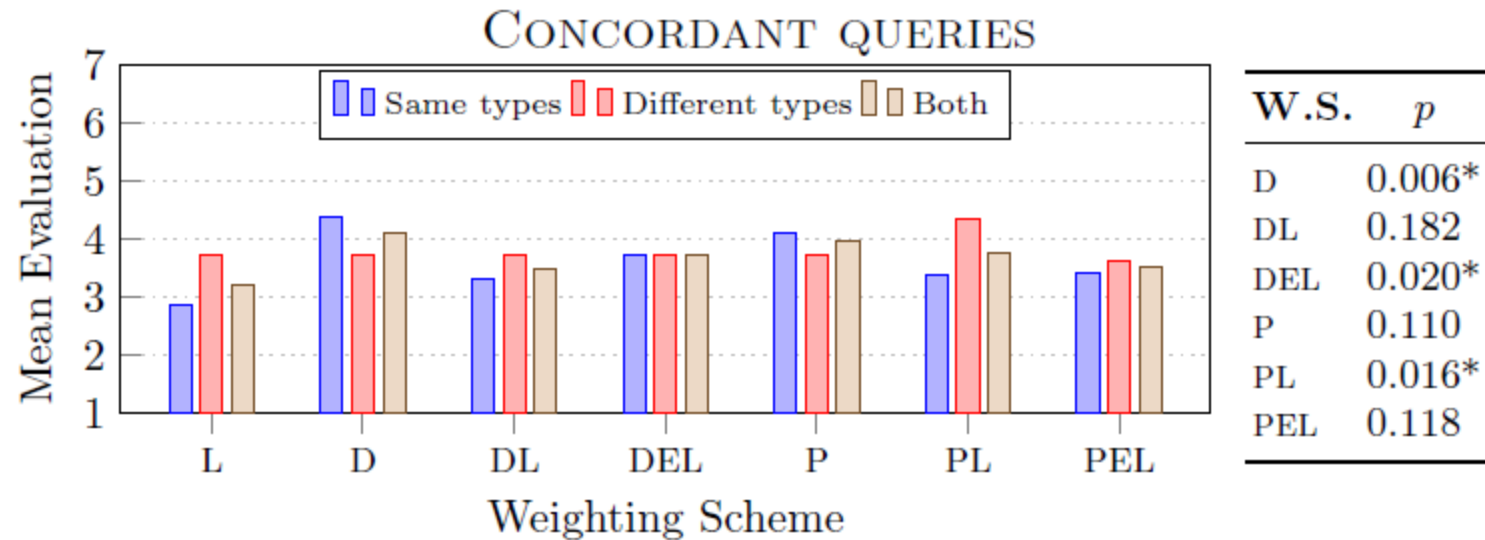
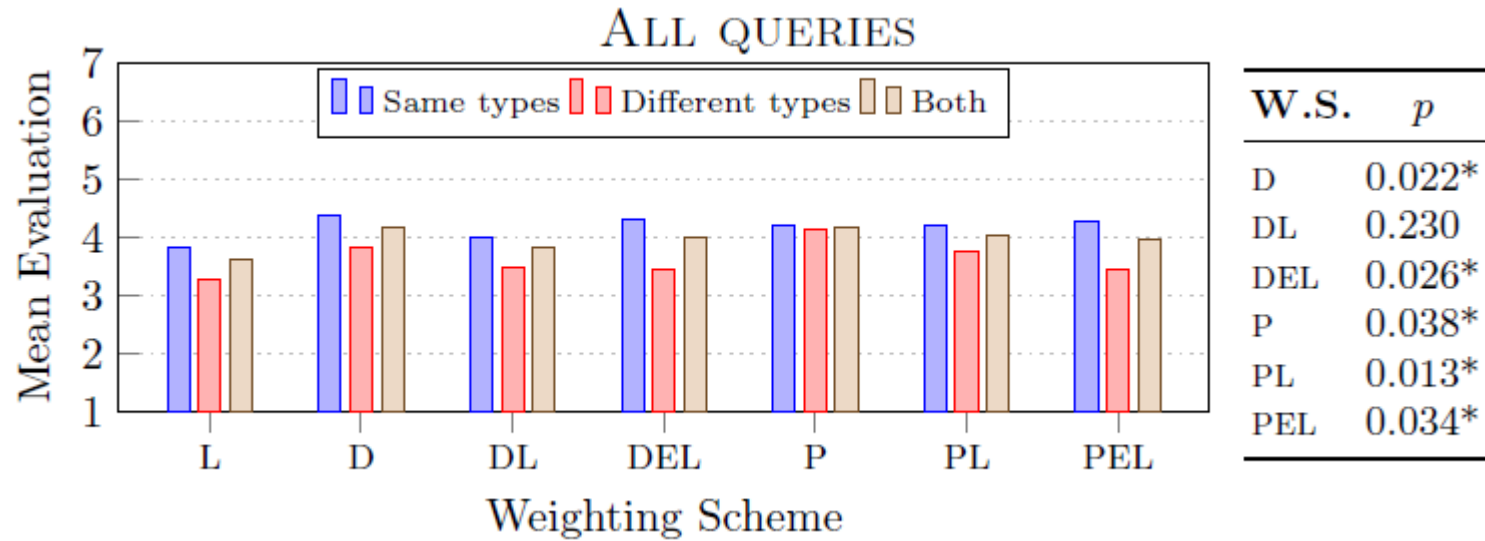


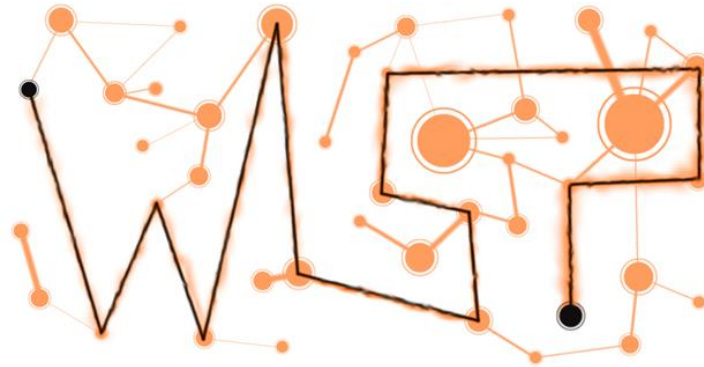
mean score 6.0 ( {5,7} )

# INTER-RATER AGREEMENT

- Kendall's  $\tau$  correlation (ordinal scales)
  - $\tau = 0.201$
  - Slight, positive agreement
- Two sets of results
  - All
    - $\tau = 0.201$ , 20 queries, 79 evaluations
  - Concordant
    - Queries with positive  $\tau$  correlation only
    - $\tau = 0.552$ , 8 queries, 27 evaluations

# USER STUDY: COMPARISON OF WEIGHTINGS





<http://wisp.dcc.uchile.cl>

DEMO



# WISP DEMO

Pathfinder [Dijkstra's Algorithm]

MULTIPLE SEARCH

Nodes weight  
PageRank weight

From:

Adolf Hitler

Entity ID

352

To:

Mahatma Gandhi

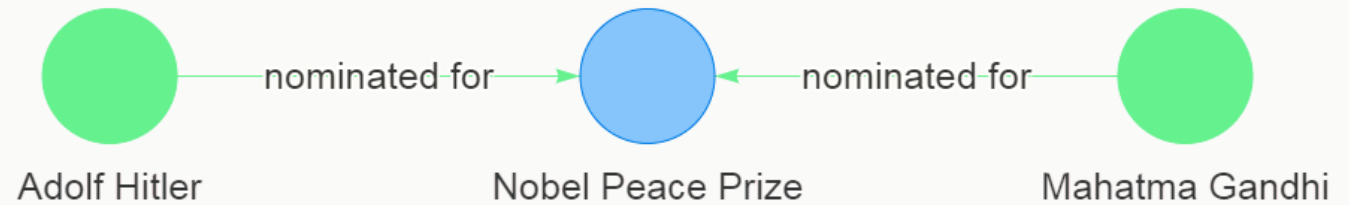
Entity ID

1001

SEARCH

Search time: 5952 ms

Use weighted Edges



# CONCLUSIONS

# CONCLUSIONS

- **Performance:**
  - How are the runtimes?
    - A few seconds (1.6 m) to a few minutes (full dataset)
  - How is the scalability?
    - Linear (roughly)
- **Weighting schemes:**
  - How similar are paths for different weightings?
    - DEL | PEL similar; others not so much
  - Does weighting help find interesting paths?
    - Yes!
  - Which weighting finds the most interesting paths?
    - No clear winner (PL best in most cases)

# FUTURE WORK

- Top-*k* queries
- Explore more weightings
- Normalisation / combinations
- Performance? (Parallelism? Approximation?)
- **iiiEvaluation!!!**