

Maverick: Discovering Exceptional Facts from Knowledge Graphs

12/03/19

Paper published in Proc. ACM SIGMOD International Conference on Management of Data, 2018.

Presented by: Juan Carrillo

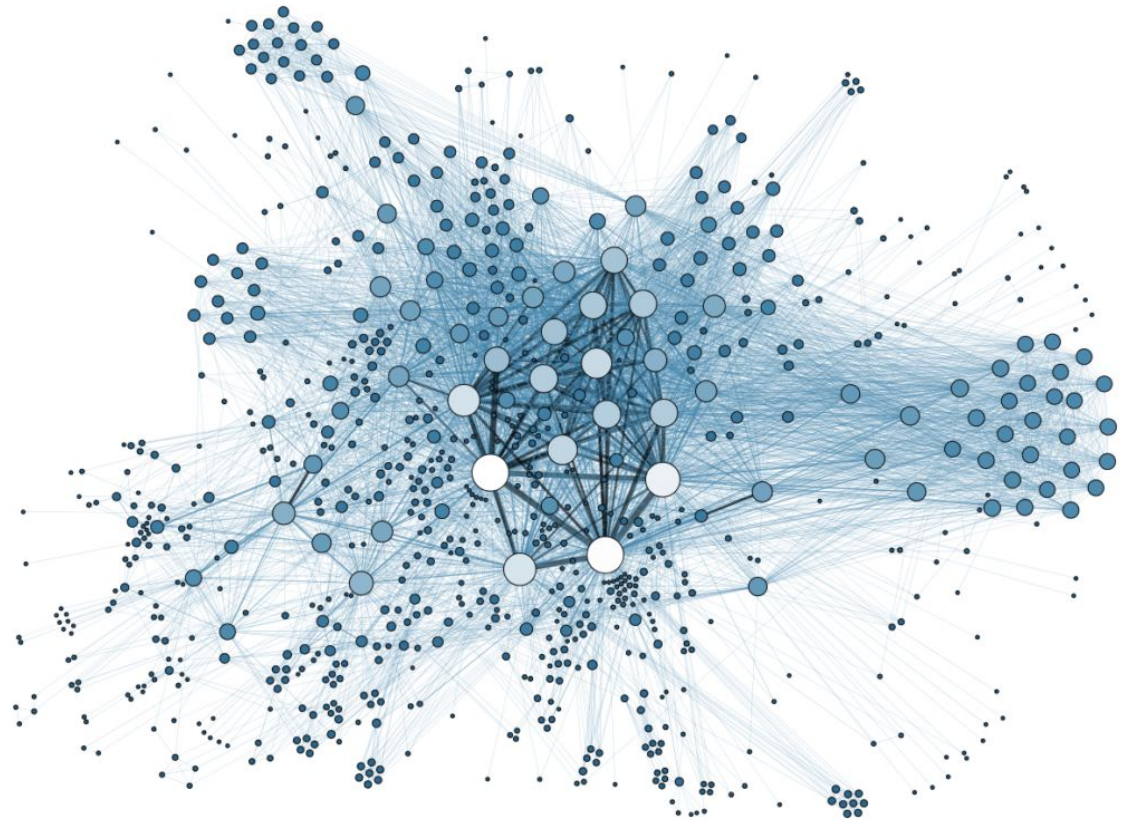
Candidate for MSc. in Computer Software

Department of Electrical & Computer Engineering

University of Waterloo



UNIVERSITY OF
WATERLOO



Agenda

1. Introduction
2. Maverick core features
3. Experiments
4. Conclusions
5. Discussion

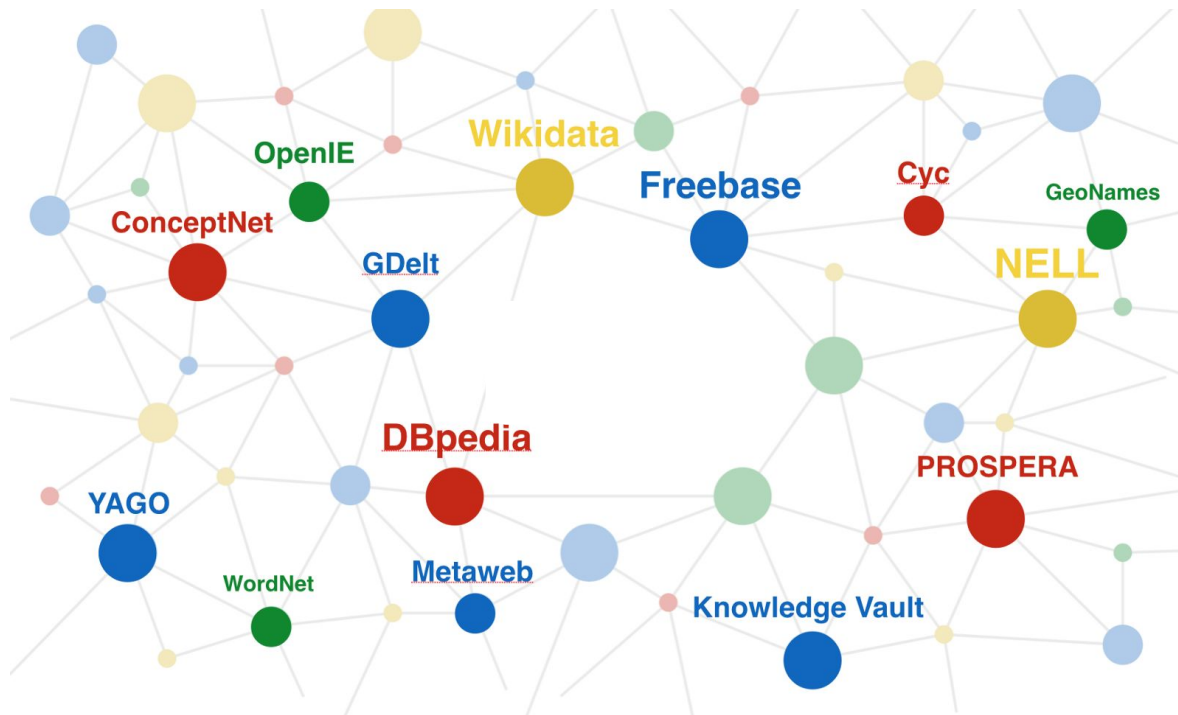




1 Introduction

1. Introduction

From knowledge graphs to exceptional facts



Denzel Washington

Denzel Hayes Washington Jr. is an American actor, director, and producer. He has received three Golden Globe awards, a Tony Award, and two Academy Awards: Best Supporting Actor for the historical war drama film *Glory* and Best Actor for his role as a corrupt cop in the crime thriller *Training Day*.

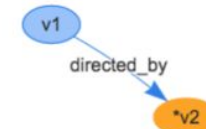


Exceptional Facts

Among all the 95486 film directors, *Denzel Washington* is one of 4665 who appeared in a film.

Isolation: 0.8571

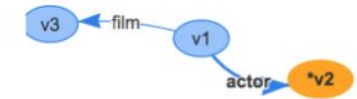
Context



Among all the 60602 film actors, *Denzel Washington* is the only one who served as one of executive producers of Film (*Chasing the Dream*) and Film (*Safe House*).

Outlierness: 0.9792

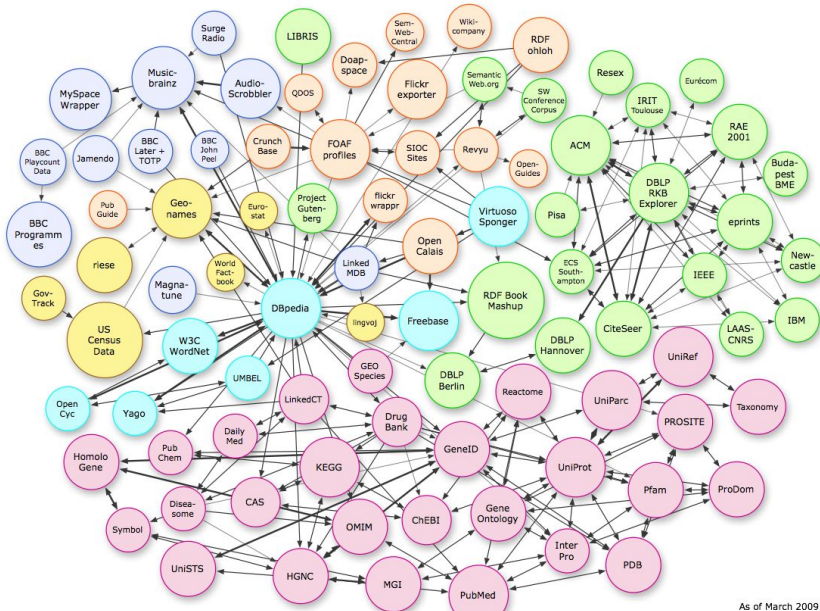
Context



1. Introduction

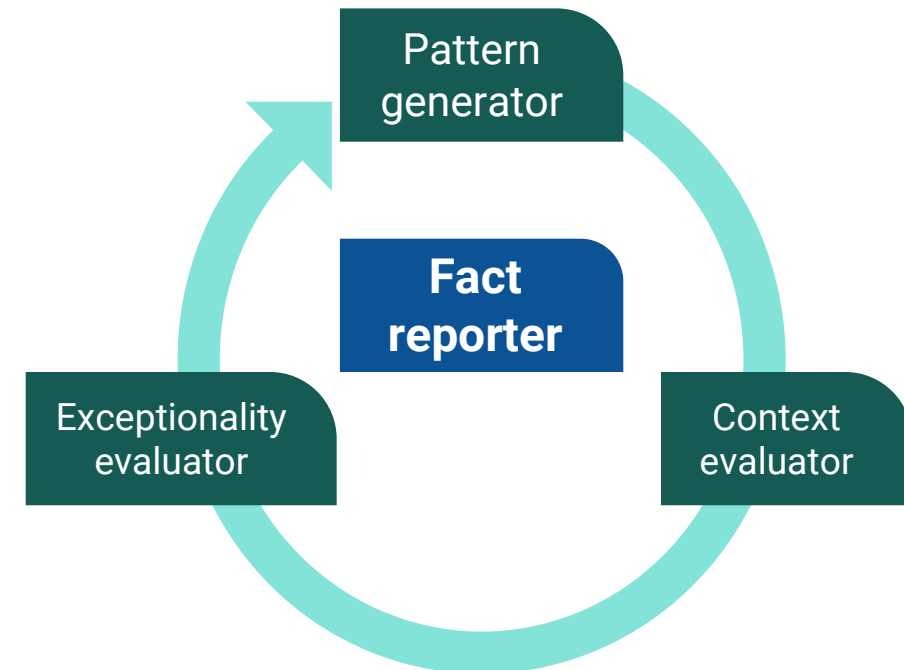
The problem, and the Maverick approach

Knowledge graphs (Linked Data)



 Manually designed queries

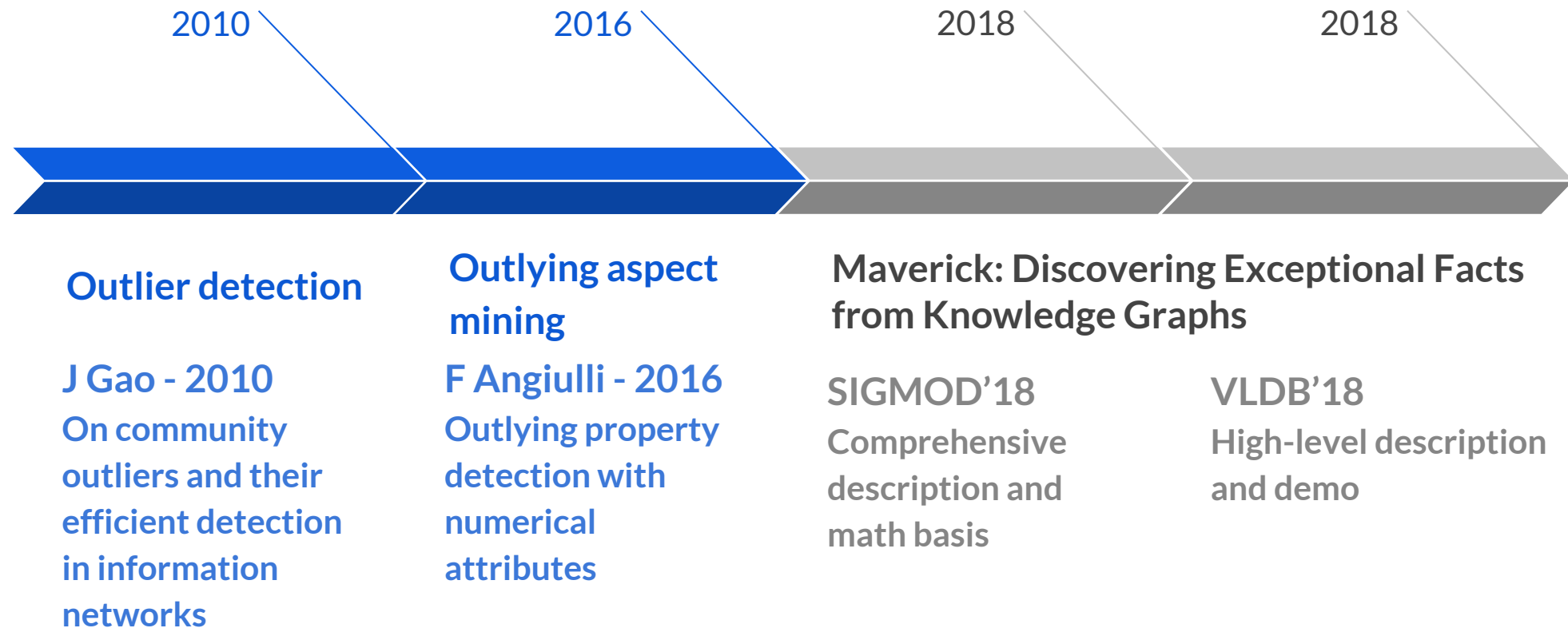
Maverick approach



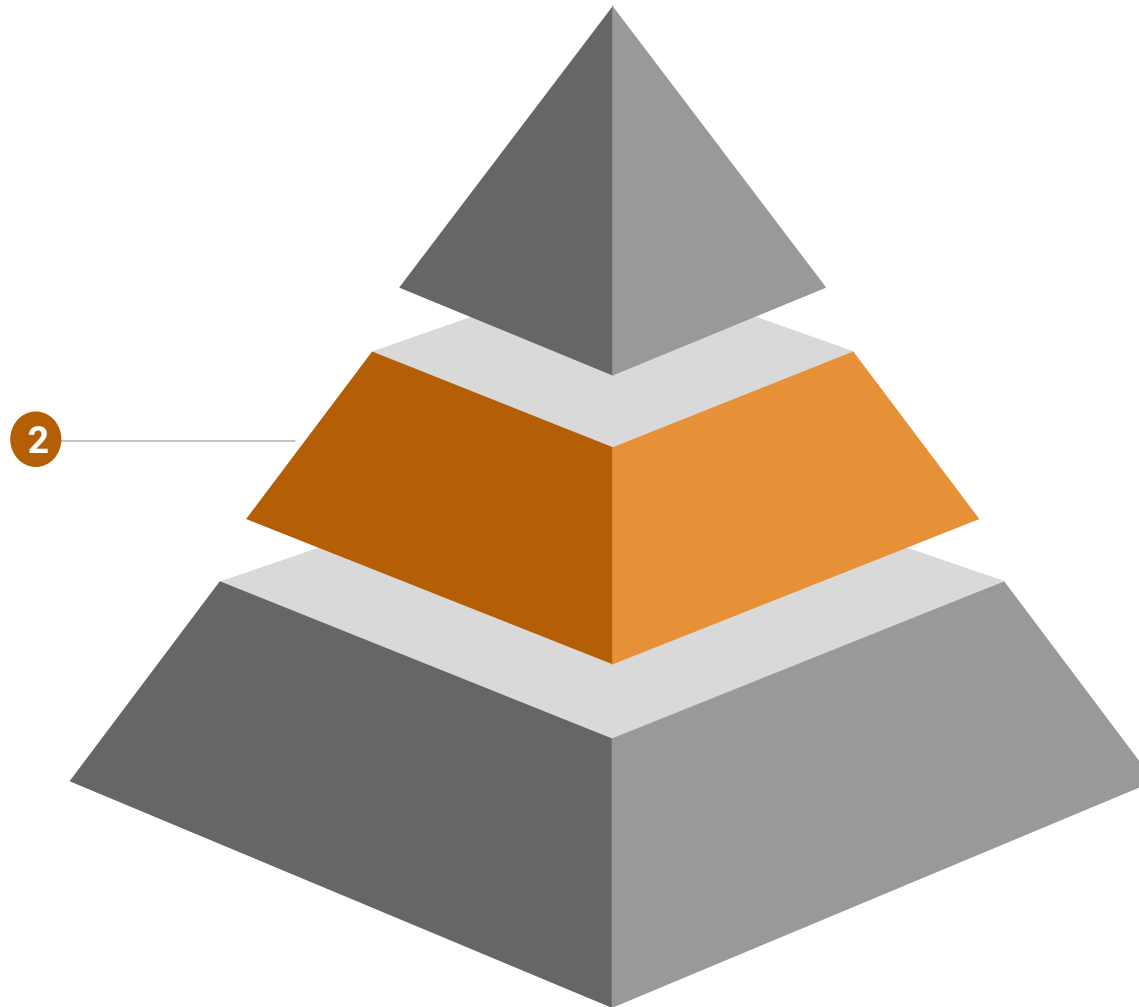
 Automated detection of exceptional facts

1. Introduction

Related background

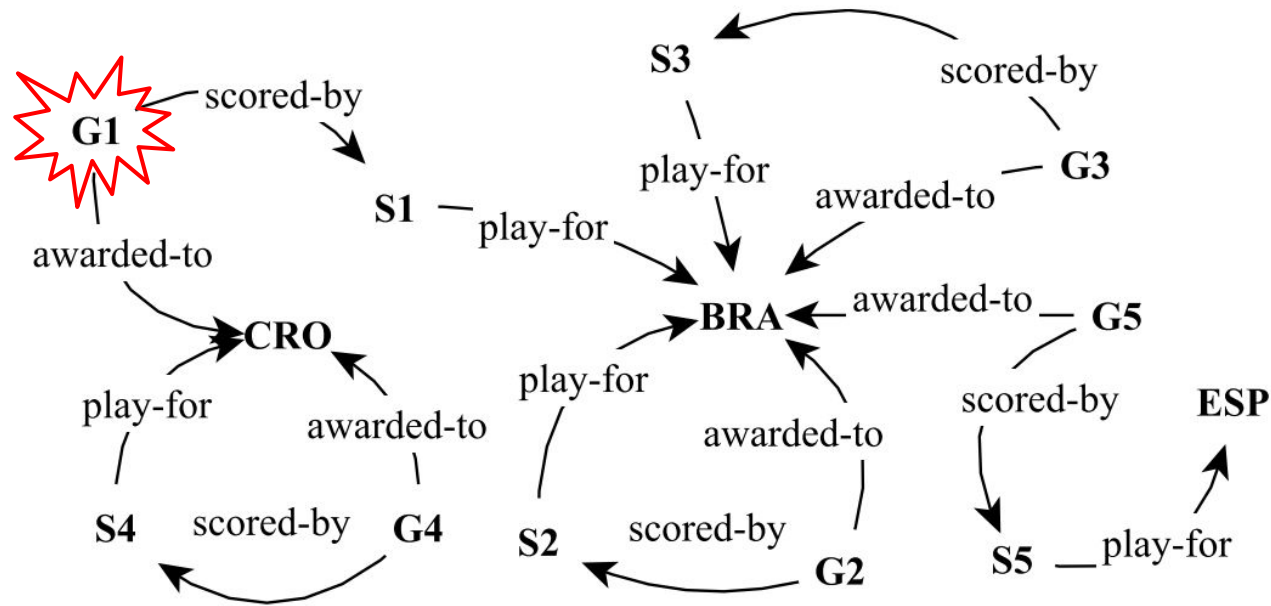


Maverick core features



2. Maverick core features

Entity, context, pattern



$?g \xrightarrow{\text{scored-by}} ?s \xrightarrow{\text{play-for}} \text{BRA}$

a Pattern P_1

$G2 \xrightarrow{\text{scored-by}} S2 \xrightarrow{\text{play-for}} \text{BRA}$

c Match M_2

$G1 \xrightarrow{\text{scored-by}} S1 \xrightarrow{\text{play-for}} \text{BRA}$

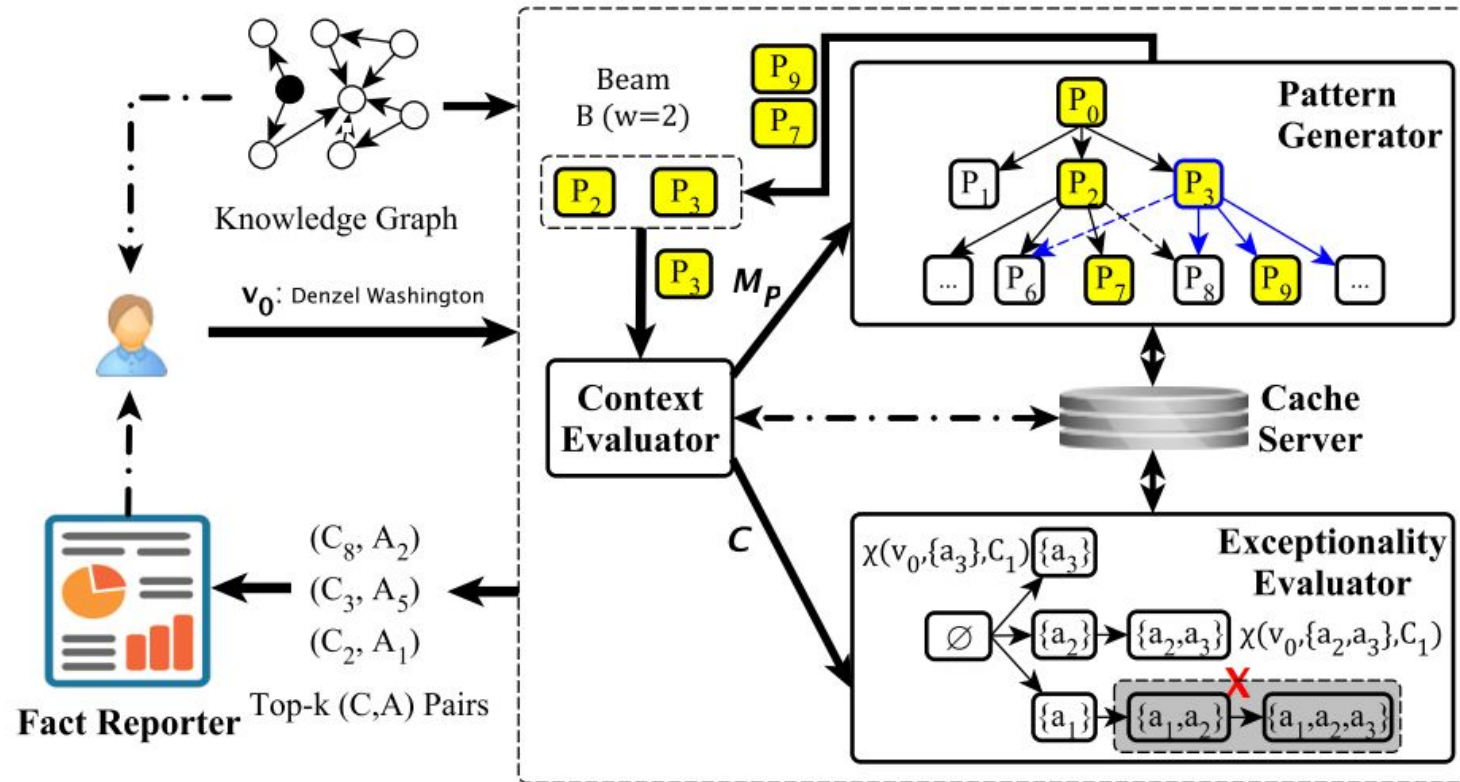
b Match M_1

$G3 \xrightarrow{\text{scored-by}} S3 \xrightarrow{\text{play-for}} \text{BRA}$

d Match M_3

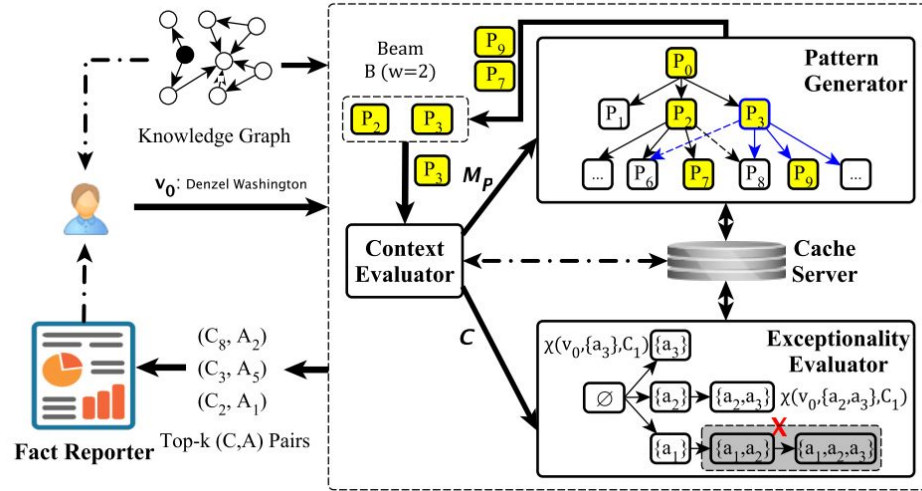
2. Maverick core features

The overall framework



2. Maverick core features

Main Algorithm



```

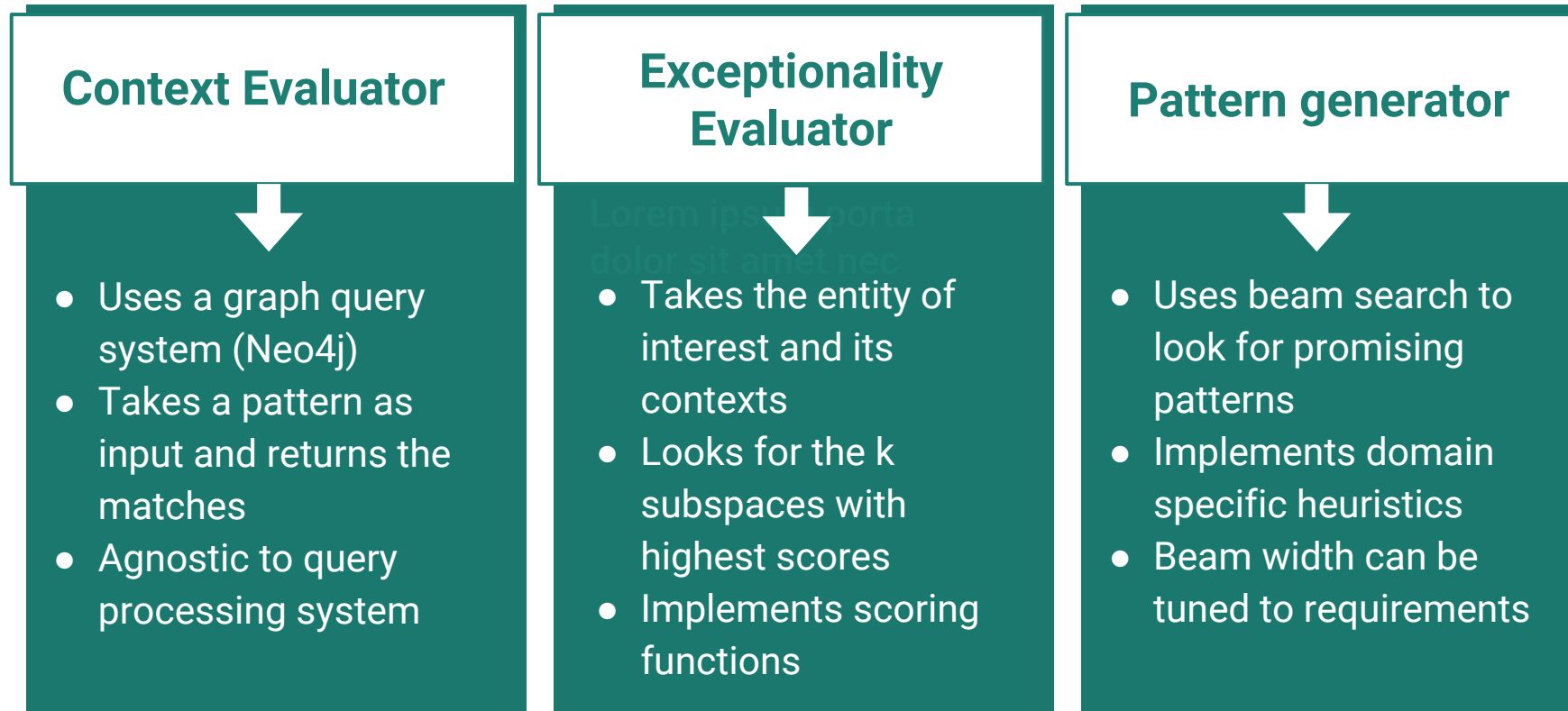
1  FACT-DISCOVER ( $G, v_0, \chi, k, w$ )
   Input:  $G$  : the knowledge graph;  $v_0 \in V_G$  : the entity of interest;
            $\chi$  : the exceptionality scoring function;  $k$  : the size of
           output;  $w$  : the beam width
   Output:  $H$  :  $k$  most exceptional context-subspace pairs

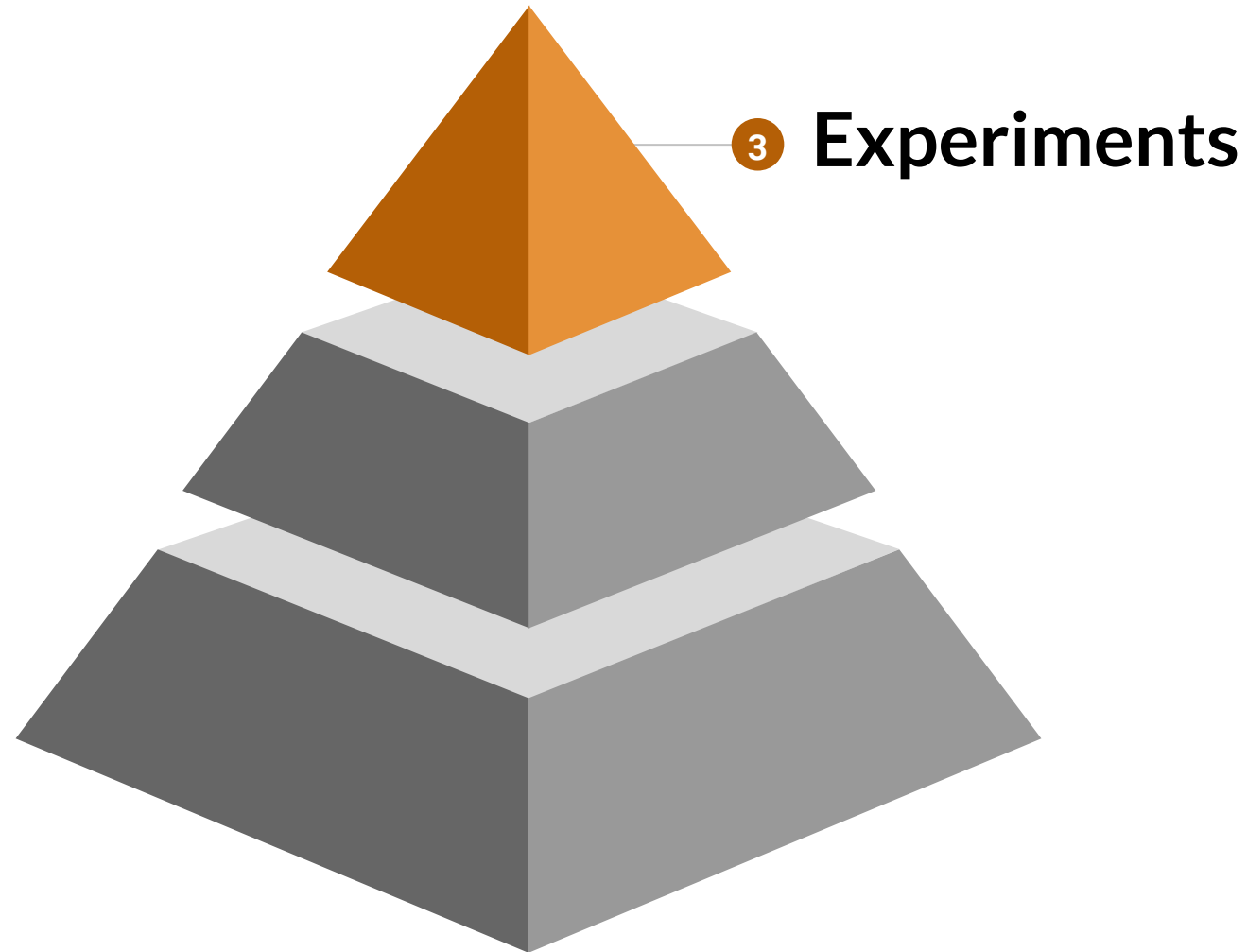
2   $P_0 \leftarrow (V_{P_0} = \{x_0\}, E_{P_0} = \emptyset)$ ;    // Initial state.  $x_0$  is a variable.
3   $B \leftarrow \{P_0\}$ ;                                // Beam.
4   $i \leftarrow 1$ ;                                    // Iteration number.
5  while  $B \neq \emptyset$  and  $i \leq \text{MAX\_ITERATION}$  do
6       $i \leftarrow i + 1$ ;  $B_{tmp} \leftarrow \emptyset$ ;
7      foreach  $P \in B$  do
8          // Obtain contexts of  $v_0$  and matches to  $P$ . (Section 3.1)
9           $C_{v_0}^P, M_P \leftarrow \text{CONTEXT-EVALUATOR}(P, v_0, G)$ ;
10         foreach  $C \in C_{v_0}^P$  do
11             // Exceptionality Evaluation. (Section 4)
12              $\mathcal{A} \leftarrow \text{EXCEPTIONALITY-EVALUATOR}(v_0, C, k, \chi)$ ;
13             foreach  $A \in \mathcal{A}$  do  $H \leftarrow H \cup \{(C, A)\}$ ;
14             // Find  $\mathcal{Y}$  – the children of  $P$ . (Section 5)
15              $\mathcal{Y} \leftarrow \text{PATTERN-GENERATOR}(v_0, P, M_P, w, G)$ ;
16              $B_{tmp} \leftarrow B_{tmp} \cup \mathcal{Y}$ ;
17          $B \leftarrow \text{top-}w \text{ of } B_{tmp} \text{ based on heuristics } h$ ;    // Section 5.4
18 return top- $k$  pairs in  $H$  based on exceptionality scores;

```

2. Maverick core features

Description of components





3. Experiments

Experimental setup

Single node: 16-core, 32GB RAM

Methods compared

- Beam-Rdm
- Beam-Opt
- Beam-Conv
- Breadth-First



Datasets

WCGoals

49.078 nodes, 158.114 edges, 13 different edge labels, and 11 entity types.

OscarWinners

42.148 nodes, 63.187 edges, 24 distinct edge labels, and 13 entity types.

3. Experiments

Efficiency

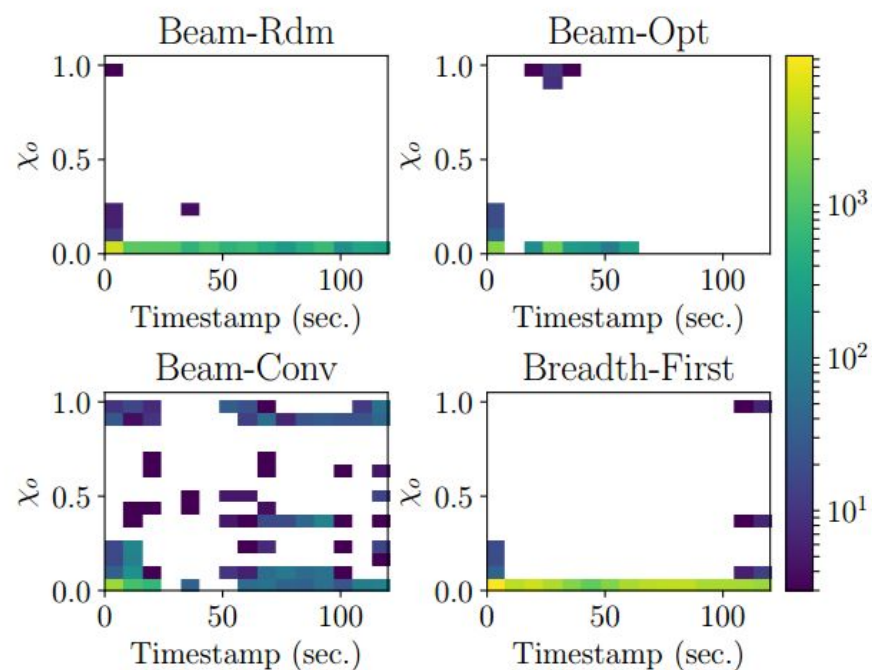
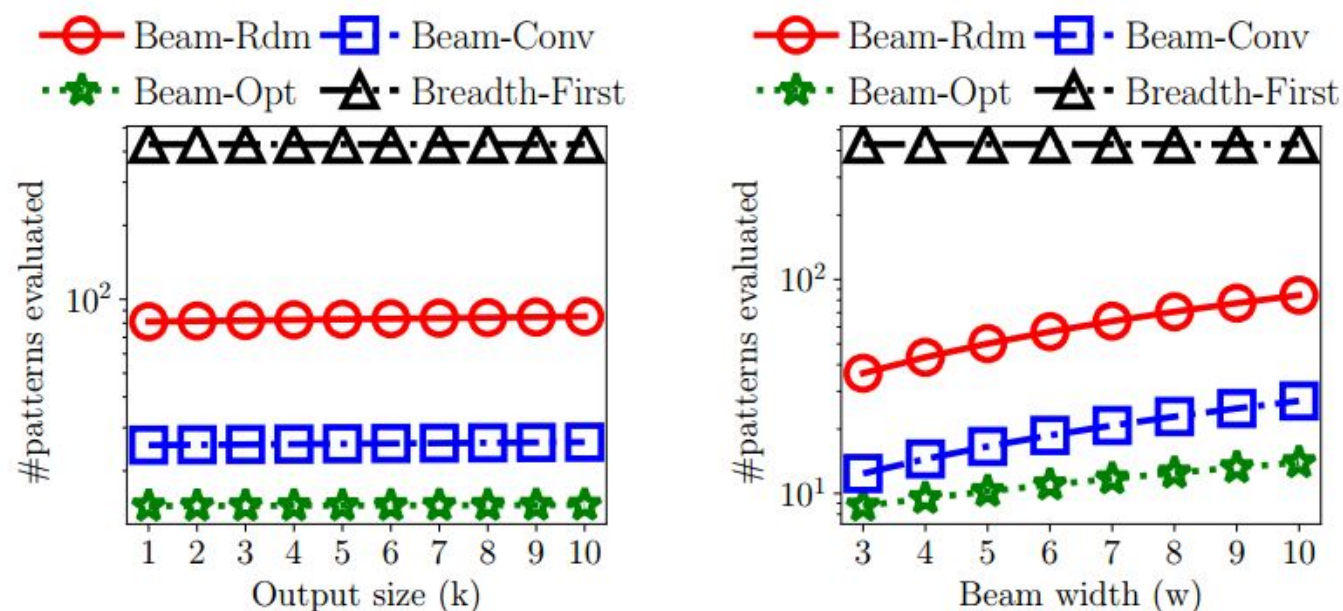


Figure 7: The heat map of exceptionalities scores (χ_o) and timestamps of all the discovered context-subspace pairs during 2-minute runs for 10 entities of interest (v_0) in WCGoals ($k = 10$, $w = 10$).

3. Experiments

Efficiency



a Varying k , fixing $w = 10$. b Varying w , fixing $k = 10$.

Figure 8: Effect of k and w on the number of evaluated patterns.

3. Experiments

Effectiveness

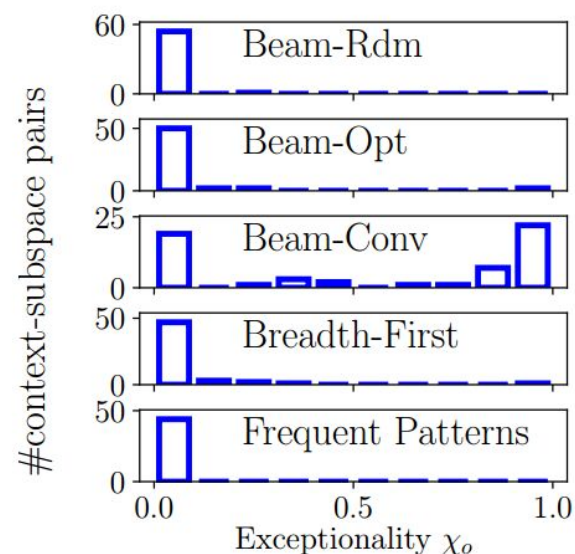


Figure 13: Score distributions of top-10 context-subspace pairs for 10 entities, 10 2-minute runs per entity.

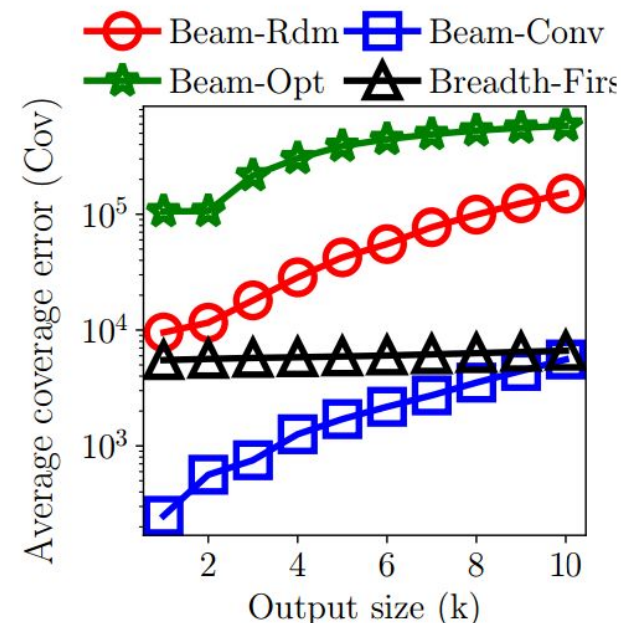


Figure 14: Average coverage error on 10 entities. Beam width 10.

Conclusions ⁴



4. Conclusions

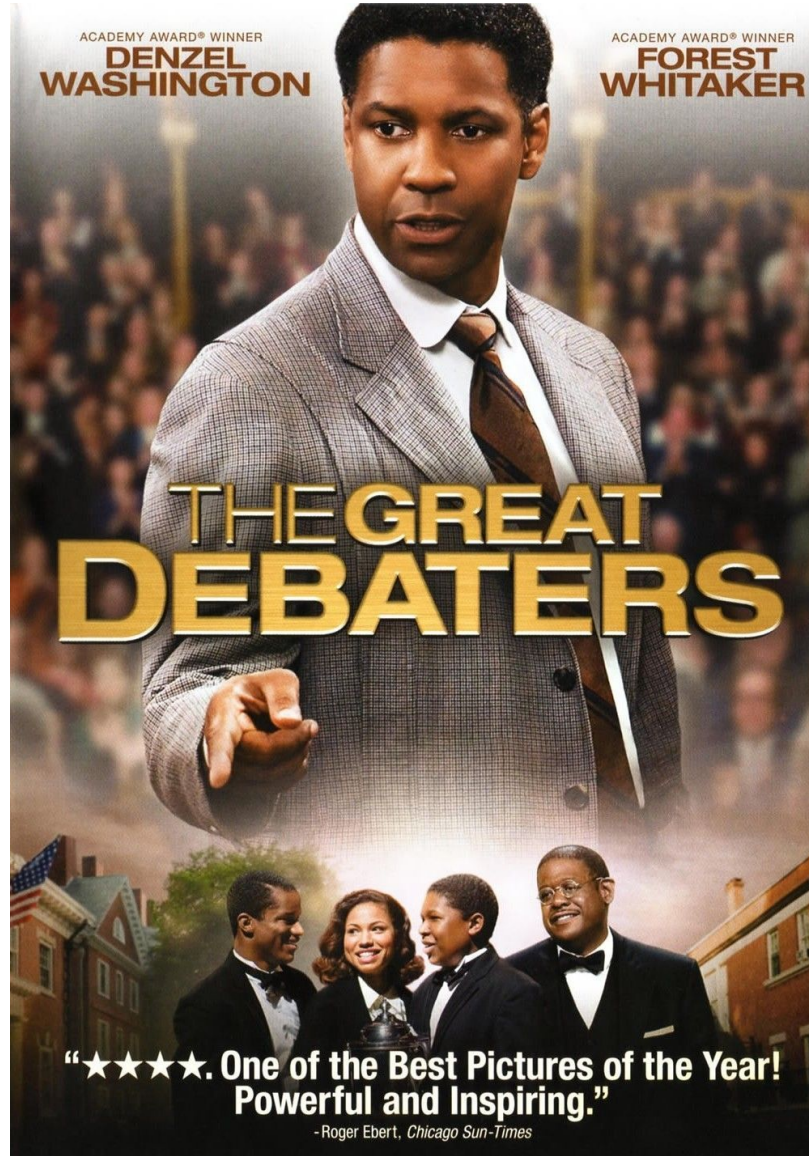
Takeaways and paper contributions

- ✓ The authors model an exceptional fact as a context-pattern pair on a knowledge graph
- ✓ Exponential complexity of search is handled using beam search
- ✓ The framework is adaptable to domain specific requirements

Thanks for your attention



Discussion 5



5. Discussion

Research

1. What other heuristics could be proposed in addition to the two presented in the paper? Design requirements for a third heuristic?
2. How Maverick would perform over a completely different dataset? Different proportions among nodes, edges, edge labels, and entity types.
3. What if we add attributes to the nodes and edges? Constraints
4. How to adapt Maverick to work over multiple/linked knowledge graphs?

Industry

5. What is an example of an application over Google knowledge graph?