# Maverick: Discovering Exceptional Facts from Knowledge Graphs

**Gensheng Zhang, Damian Jimenez, Chengkai Li**

SIGMOD, Houston, June 14, 2018

UNIVERSITY OF TEXAS ARLINGTON

iDiR
idir.uta.edu

# Exceptional Facts

Denzel Washington followed Sidney Poitier as only the second black to win the Best Actor award.

**Entity of interest**     Denzel Washington

**Context**     Best Actor award winners

**Attributes**     Ethnicity

**Peculiar value**     African American
(only two satisfy)

# Exceptional Facts

**abc NEWS** Denzel Washington followed Sidney Poitier as only the second black to win the Best Actor award.

**Entity of interest** Denzel Washington

Given an entity **x**

**find**

**Context** Best Actor award winners

A context

**Attributes** Ethnicity

A set of attributes (subspace)

**Peculiar value** African American (only two satisfy)

idir.uta.edu 3

# Exceptional Facts

Denzel Washington followed Sidney Poitier as only the second black to win the Best Actor award.

| | | | |
|---|---|---|---|
| **Entity of interest** | Denzel Washington | Given an entity **x** | |
| | | **find** | **such that** |
| **Context** | Best Actor award winners | A context | the context has many entities, including **x** |
| **Attributes** | Ethnicity | A set of attributes (subspace) | |
| **Peculiar value** | African American (only two satisfy) | | **x** bears a peculiar value w.r.t. the subspace (few in the context have the value) |

**4**

# Exceptional Facts

**abc NEWS** Denzel Washington followed Sidney Poitier as only the second black to win the Best Actor award.

**YAHOO! SPORTS** This was Brazil's first own goal in World Cup history.

**Chicago Tribune** Hillary Clinton becomes first female presidential nominee.

idir.uta.edu **5**

# Applications

## Computational Journalism

- Fact-finding
- Fact-checking
  - The first female presidential nominee was Victoria Woodhull, not Hillary Clinton (snopes.com)
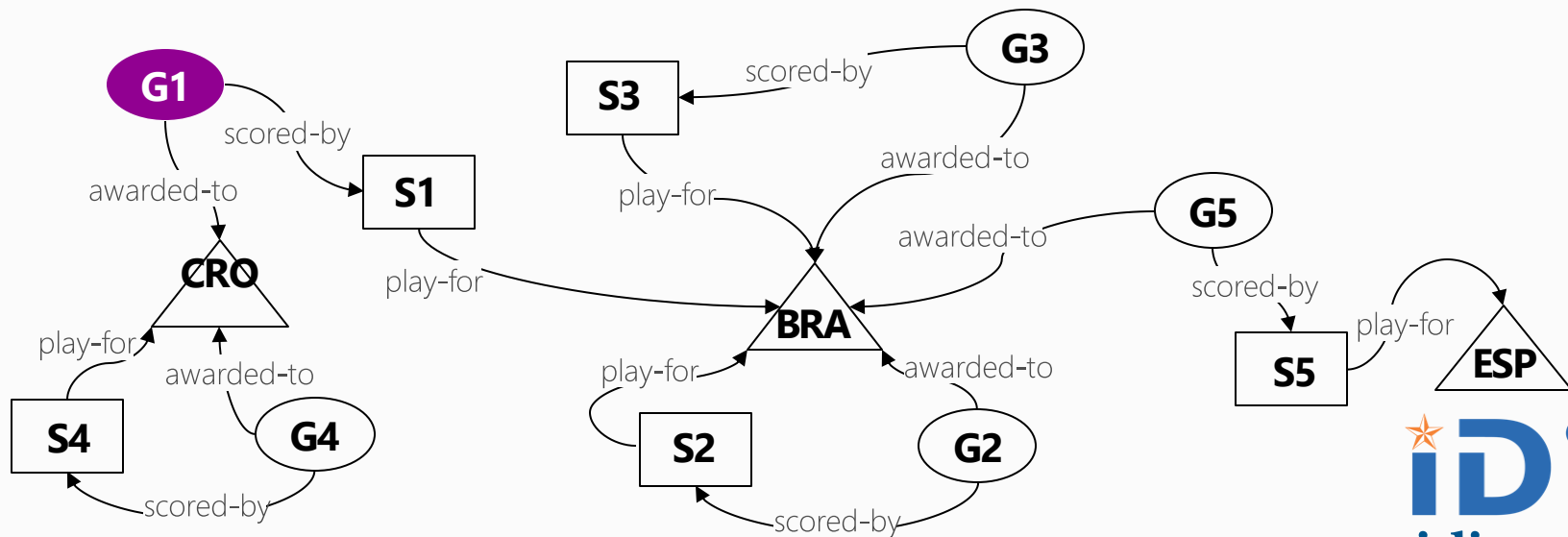
## Data Cleaning

## Recommendation Systems

- Friends, news, and product promotion

**Willis Tower**

Website   Directions

4.4   1,556 Google reviews

Skyscraper in Chicago, Illinois

The Willis Tower, built as and still commonly referred to as Sears Tower, is a 108-story, 1,450-foot skyscraper in Chicago, Illinois, United States. Wikipedia

Hours: Open today   6AM–3PM

**Did you know:** Willis Tower in Chicago is the second-tallest building in the US. wikipedia.org

# Exceptional Facts from Knowledge Graphs

What is exceptional about G1?

# Modeling

Attributes: labels of incoming/outgoing edges

Subspace: a subset of attributes

G1. awarded-to = CRO

# Modeling

Context: entities sharing some common characteristics

Defined by a pattern-variable pair

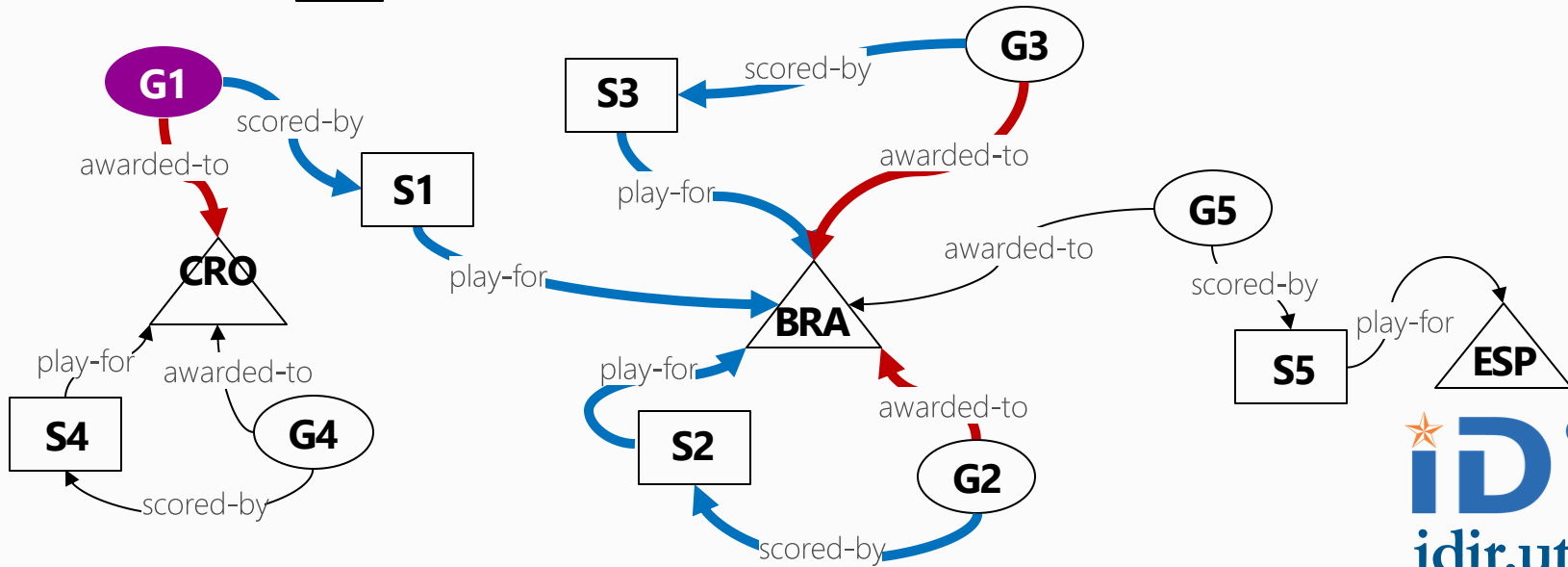?g —awarded-to→ **?s** —play-for→ **BRA** Goals scored by Brazilian players

# Modeling

Context: entities sharing some common characteristics
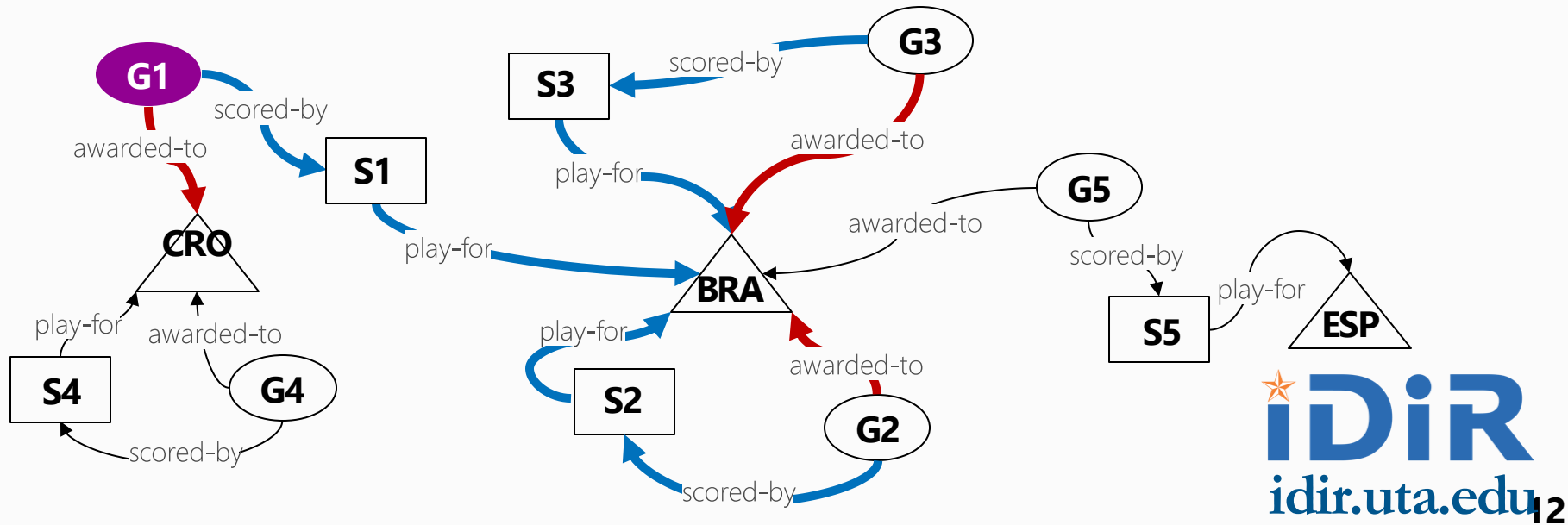Defined by a pattern-variable pair



?g —awarded-to→ ?s —play-for→ BRA    Goals scored by Brazilian players

# Modeling

Context: entities sharing some common characteristics
Defined by a pattern-variable pair



Goals scored by Brazilian players

# Modeling

**What is exceptional about G1?**

Among all the goals scored by BRA players, G1 is the only own goal.

# Problem Formulation

## Input

- Entity of interest $v_0$
- Exceptionality function $\chi$
- Result size $k$

## Output

- Top-$k$ (context, subspace) pairs with regard to $\chi$, in which $v_0$ stands out

# Challenges

- Number of attribute subspaces: $O(2^{|A_{v_0}|})$
- Number of patterns (contexts): $\Omega(2^{|V_G|})$

# Related Work

## Outlier detection

o   Maverick finds conditions that make an object stand out, although the object may not necessarily be an outlier.

## Outlying aspect mining

Challenges in adopting existing algorithms:

o   Many assume a single-table model: a graph can be an extremely large and sparse table

o   Conjunctive queries on a single table ≠ pattern queries

o   Multiple tables: unclear how to handle joins

o   Unclear how to handle set values

# Maverick

# Exceptionality Function $\chi$

$\chi(v, A, C) \in \mathcal{R}$

outlierness $(\chi_o)$ [Angiuli2009TODS], one-of-the-few$(\chi_f)$ [Wu2012KDD], isolation scores $(\chi_i)$ [Liu2008ICDM]

## Upper bound function

Theorem 4.2

$$upper_o(v, A, C) = \sum_{S \in \mathcal{S}_A} (p_S)^2 - \frac{(2\, p_{v.A} + 1) \times |C| - 2}{|C|^2}$$

Theorem 4.3

$$upper_f(v, A, C) = \left| \{u \mid u \in \overline{C_v}, \; p_{u.A}^A > 1/|C|\} \right| \Big/ |C|$$

Theorem 4.4

$$upper_i(v, A, C) = 1 - 2^{-\frac{-\log_2 \frac{1}{|C|}}{-q_{v.A} - sum_{S \in \mathcal{S}_A \setminus \{v.A\}} \, p_S \times \log_2 p_S}}$$

# Pattern Generator (PG)

# Match–based Pattern Generation

o   Construct Partial Order of Valid Patterns

THEOREM 5.4. *Suppose $P'$ is a child of $P \in \mathbb{P}$, i.e., $(P, P') \in E_{\mathbb{P}}$ and thus $P'$ is a valid pattern with matches. Given any match $M'$ to $P'$, there exists a match $M$ to $P$ that is a subgraph of $M'$, i.e., $\forall M' \in \mathcal{M}_{P'}, \exists M \in \mathcal{M}_P \text{ s.t. } V_M \subseteq V_{M'} \text{ and } E_M \subseteq E_{M'}$.*

19

# Datasets and Experiments

## WCGoals

Created based on FIFA.com

11 node types, 13 edge types

49,078 nodes, 158,114 edges

## Film-Award

A subgraph of Freebase

95 node types, 117 edge types

5,437,628 nodes, 10,879,448 edges

# See you in Rio



VLDB2018 demo

**Maverick: A System for Discovering Exceptional Facts from Knowledge Graphs**

idir.uta.edu