# Erroneous entities: how to capture?

▪ Multi-relational graphs: a labeled graph with attributes on nodes



**Graph G:** a football database
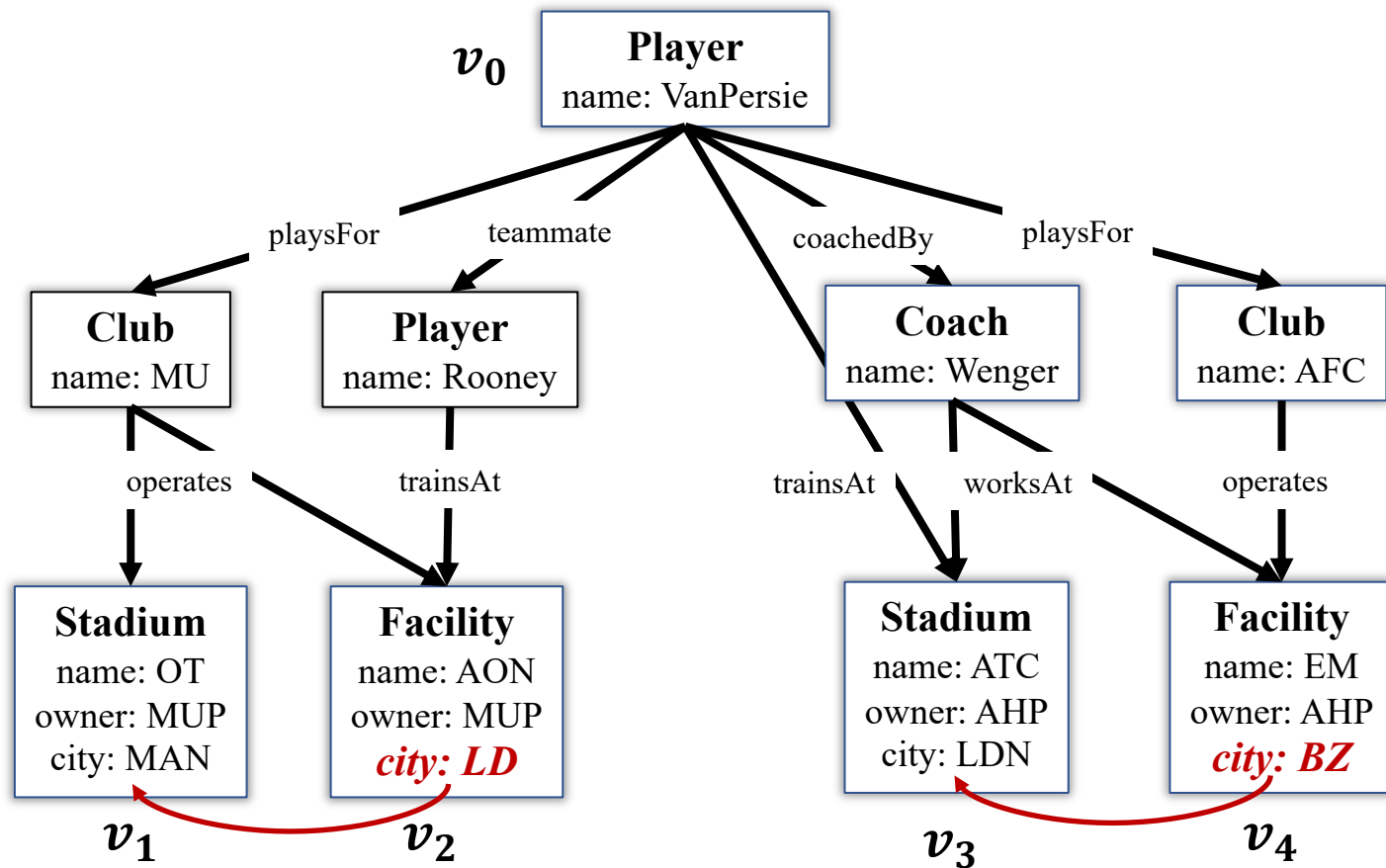
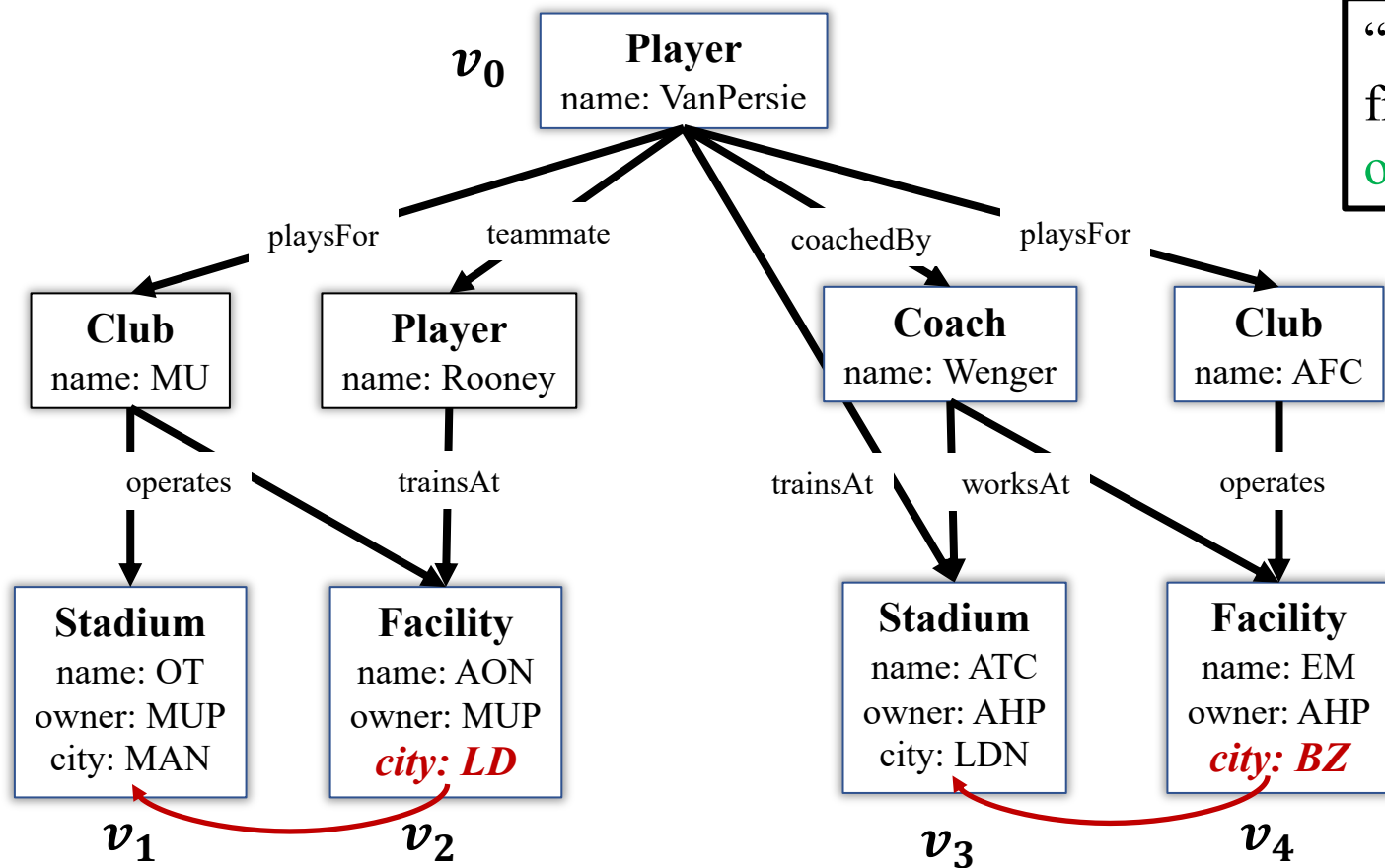# Erroneous entities: how to capture?

- Multi-relational graphs: a labeled graph with attributes on nodes
- Entity errors: incorrect node attributes



**Graph G:** a football database

# Erroneous entities: how to capture?

- Multi-relational graphs: a labeled graph with attributes on nodes
- Entity errors: incorrect node attributes
- Semantics: *relevant paths from a center node*



"For stadium and facility *relevant* to player ($v_0$) from Premier League, if they have the same owner, then they should locate at the same city."

**Graph G:** a football database

# Regular path queries

- Regular expressions: $R = l \,\big|\, l^{\leq k} \,\big|\, R \cdot R \,\big|\, R \cup R$



**Graph G:** a football database

# Regular path queries

- Regular expressions: $R = l \,\big|\, l^{\leq k} \,\big|\, R \cdot R \,|\, R \cup R$

- Paths from **Player** to **Stadium**
- $R_1 = (\text{playsFor} \cdot \text{operates}) \cup (\text{coachedBy} \cdot \text{worksAt})$
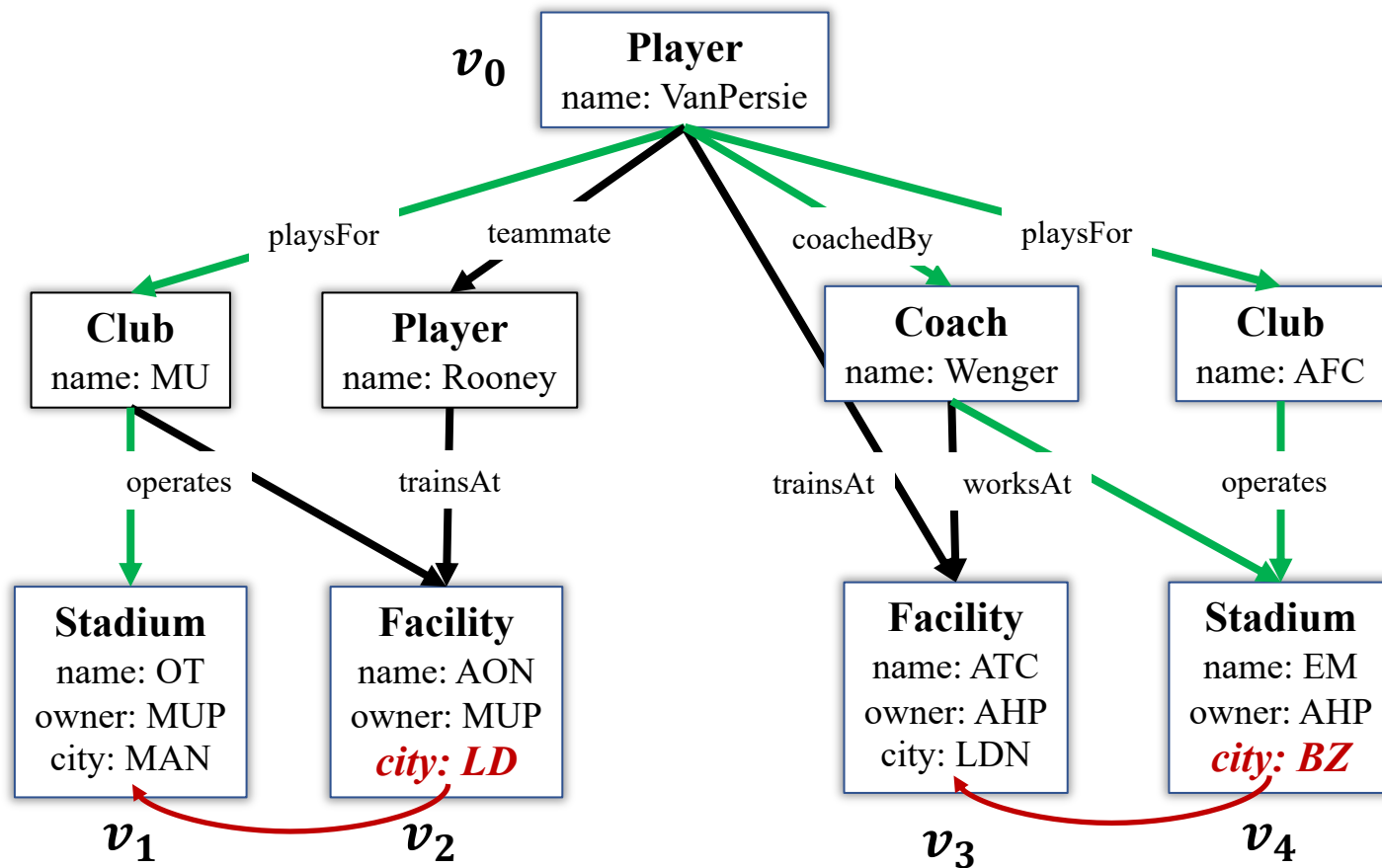


**Graph G:** a football database

# Regular path queries

- Regular expressions: $R = l \,|\, l^{\leq k} \,|\, R \cdot R \,|\, R \cup R$

- Paths from **Player** to **Stadium**
- $R_1 = (\text{playsFor} \cdot \text{operates}) \cup (\text{coachedBy} \cdot \text{worksAt})$

- Paths from **Player** to **Facility**
- $R_2 = (\text{playsFor} \cdot \text{operates}) \cup (\text{teammate}^{\leq 1} \cdot \text{trainsAt})$



**Graph G:** a football database

# Contributions

Graph $G$, StarFDs $\Sigma$
($G$ does not satisfy $\Sigma$) → Error detection → Repair → Repair $G'$
($G'$ satisfies $\Sigma$)

# Contributions

StarFDs: star functional dependencies
new constraints for graphs

Entity repair problem: minimum
editing cost, NP-hard and APX-hard

**StarRepair framework**

Graph $G$, StarFDs $\Sigma$
($G$ does not satisfy $\Sigma$) $\rightarrow$ Error detection $\rightarrow$ Repair $\rightarrow$ Repair $G'$
($G'$ satisfies $\Sigma$)

**Feasible framework** with provable
guarantees whenever possible

# Contributions

**StarFDs:** star functional dependencies
new constraints for graphs

**Entity repair problem:** minimum
editing cost, NP-hard and APX-hard

**StarRepair framework**

Graph $G$, StarFDs $\Sigma$
($G$ does not satisfy $\Sigma$) → Error detection → Repair → Repair $G'$
($G'$ satisfies $\Sigma$)

**Feasible framework** with provable
guarantees whenever possible

**Repair workflow**

Is approximable?

Yes ← → No

Is optimal repairable?      Heuristic solution

Yes ← → No

Optimal solution      Approximation solution

# Star constraints

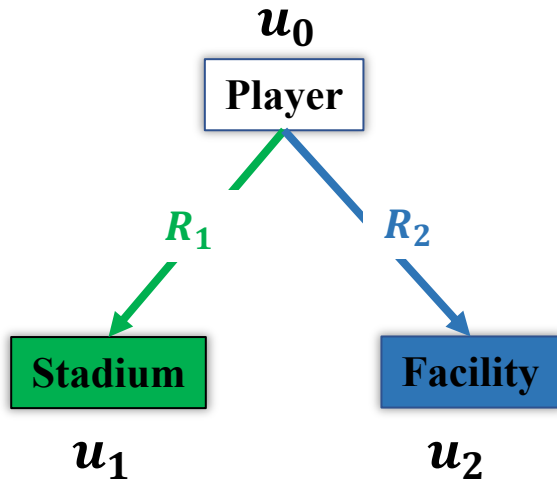- **StarFDs:** $\varphi = (P(u_o), X \rightarrow Y)$

- Star pattern $P(u_o)$:

- Value constraints: $X \rightarrow Y$

# Star constraints

- **StarFDs:** $\varphi = (P(u_o), X \rightarrow Y)$

- Star pattern $P(u_o)$:
  - A two-level tree with center node $u_o$
  - Each branch is a regular expression



$R_1 = (\text{playsFor} \cdot \text{operates}) \cup (\text{coachedBy} \cdot \text{worksAt})$

$R_2 = (\text{playsFor} \cdot \text{operates}) \cup (\text{teammate}^{\leq 1} \cdot \text{trainsAt})$

- Value constraints: $X \rightarrow Y$

# Star constraints

- **StarFDs:** $\varphi = (P(u_o), X \rightarrow Y)$

- Star pattern $P(u_o)$:
  - A two-level tree with center node $u_o$
  - Each branch is a regular expression



$R_1 = (\text{playsFor} \cdot \text{operates}) \cup (\text{coachedBy} \cdot \text{worksAt})$

$R_2 = (\text{playsFor} \cdot \text{operates}) \cup (\text{teammate}^{\leq 1} \cdot \text{trainsAt})$
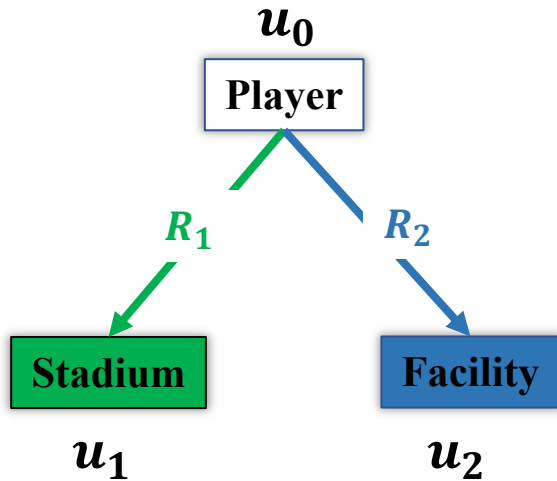
- Value constraints: $X \rightarrow Y$
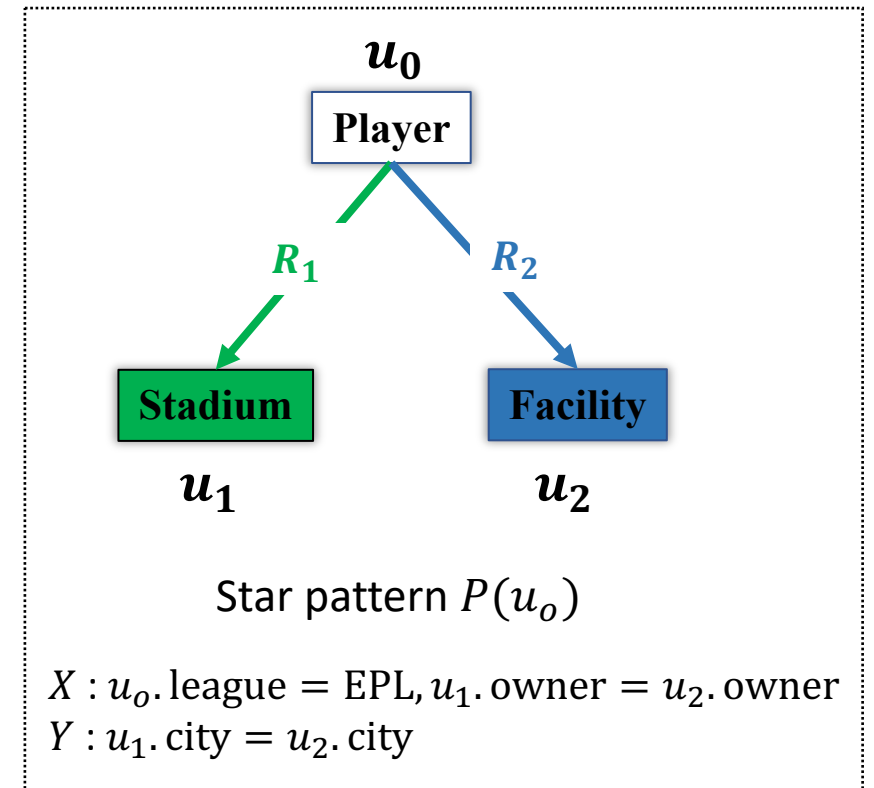  - $X$ and $Y$ are two sets of literals
  - Literals: $u.A = c$, or $u.A = u'.A'$

$X : u_o.\text{league} = \text{EPL}, u_1.\text{owner} = u_2.\text{owner}$

$Y : u_1.\text{city} = u_2.\text{city}$

# Star constraints

- Matching semantics: maximum set matched by star pattern



Star pattern $P(u_o)$

$X : u_o.\text{league} = \text{EPL}, u_1.\text{owner} = u_2.\text{owner}$
$Y : u_1.\text{city} = u_2.\text{city}$

# Star constraints

■ Matching semantics: maximum set matched by star pattern

$u_0$ matches $v_0$
$u_1$ matches $v_1$ and $v_4$
$u_2$ matches $v_2$ and $v_3$



Star pattern $P(u_o)$

$X : u_o.\text{league} = \text{EPL}, u_1.\text{owner} = u_2.\text{owner}$
$Y : u_1.\text{city} = u_2.\text{city}$

# Star constraints

- Matching semantics: maximum set matched by star pattern
- Inconsistencies $I$: matches that $X$ holds but $Y$ does not hold

$u_0$ matches $v_0$
$u_1$ matches $v_1$ and $v_4$
$u_2$ matches $v_2$ and $v_3$



Star pattern $P(u_o)$

$X : u_o.\text{league} = \text{EPL}, u_1.\text{owner} = u_2.\text{owner}$
$Y : u_1.\text{city} = u_2.\text{city}$

# Summary of results
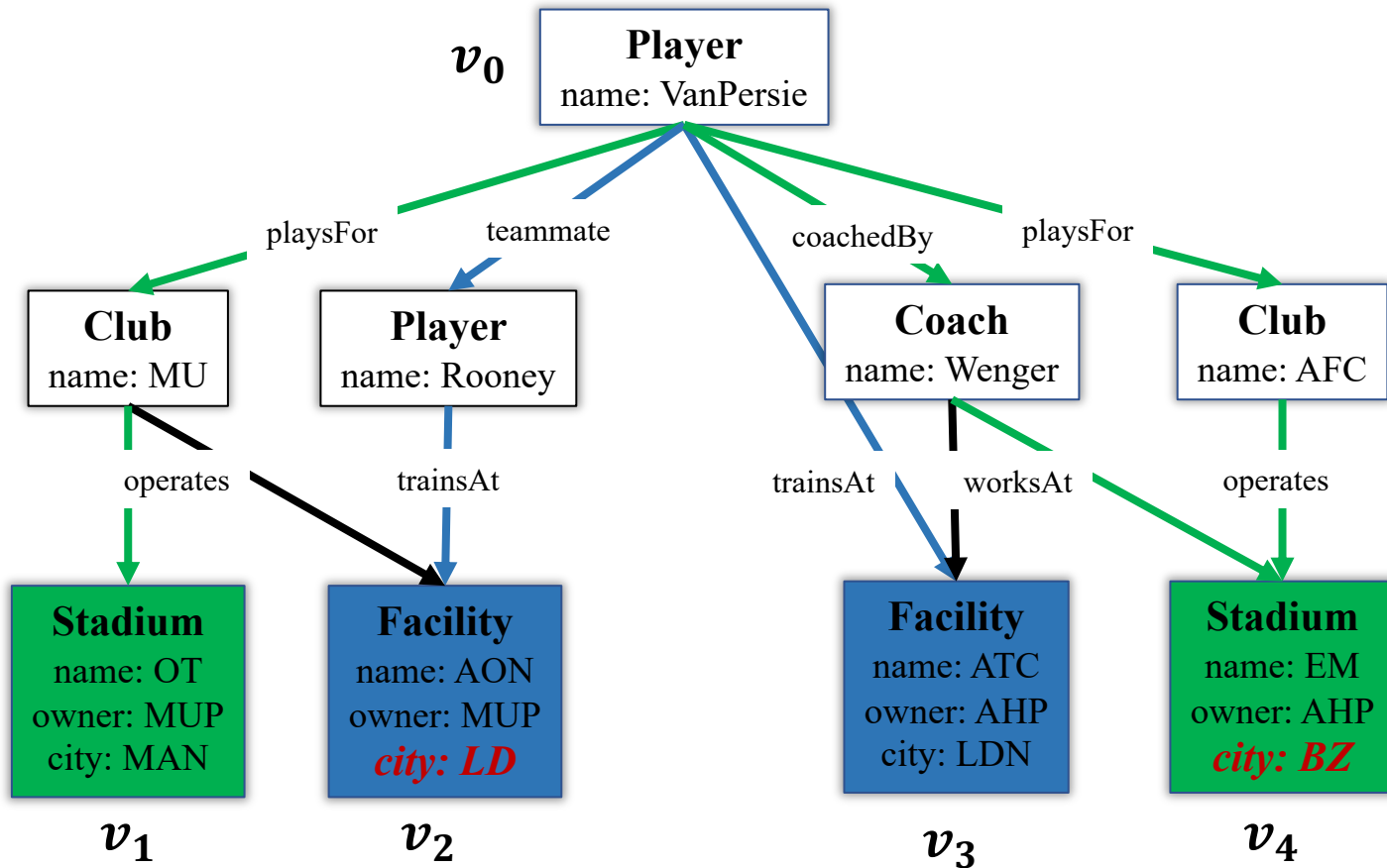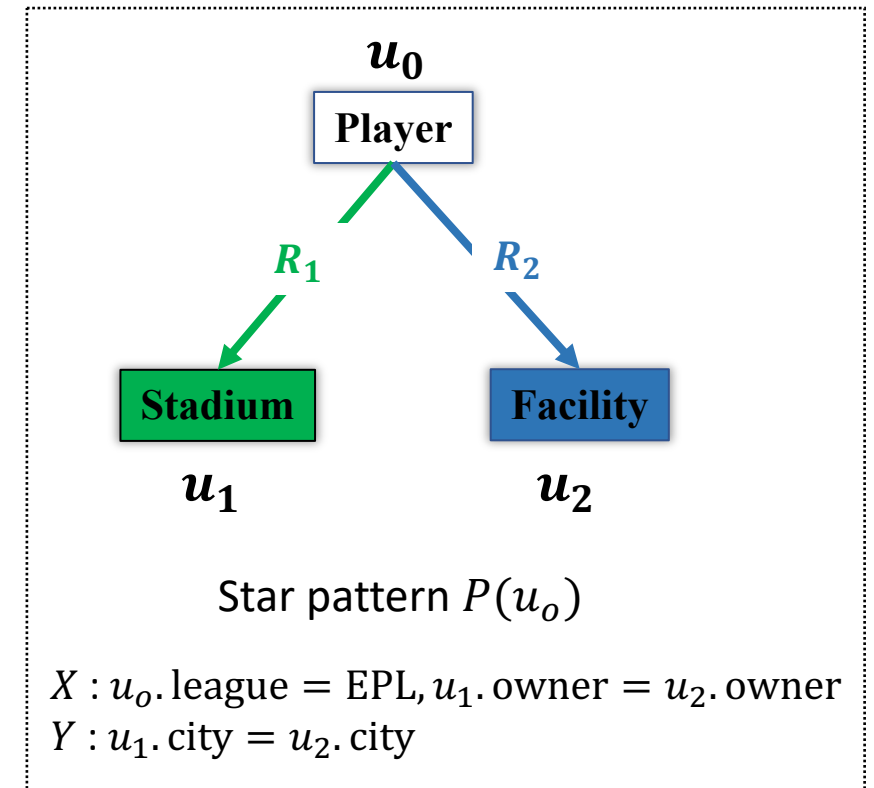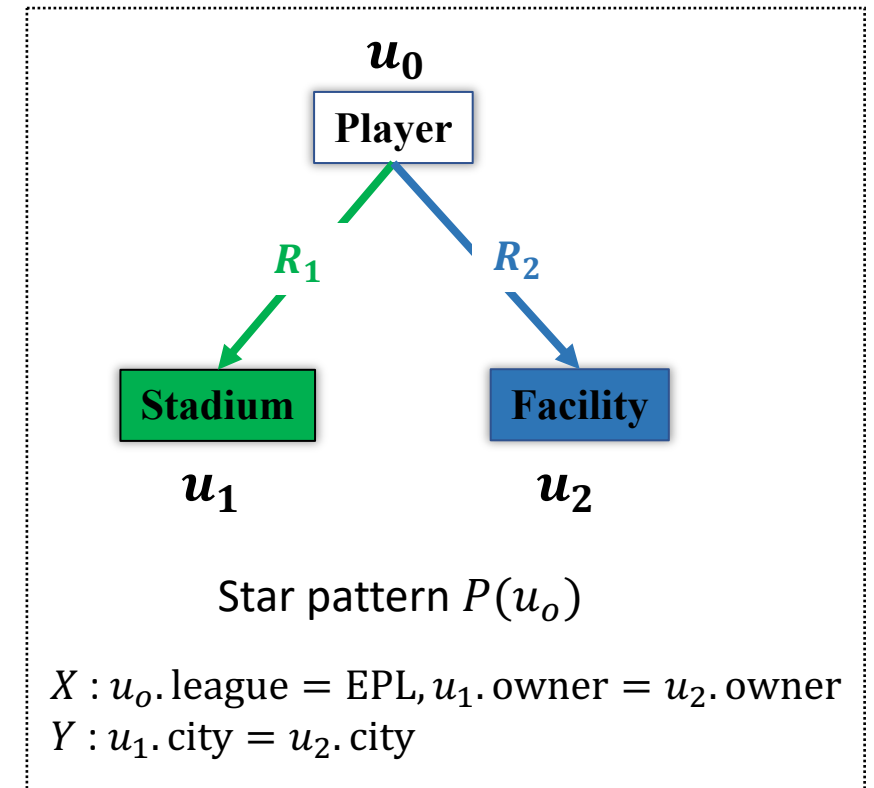
| Problem | Description | Hardness | Solution |
|---------|-------------|----------|----------|
| Satisfiability | **Input:** $\Sigma$<br>decide whether there exists $G$ that satisfies $\Sigma$ | NP-complete | |
| Implication | **Input:** $\Sigma$ and $\varphi$<br>decide whether for all $G$ satisfy $\Sigma$, they satisfy $\varphi$ | coNP-hard | |
| Error detection (validation) | **Input:** $G$ and $\Sigma$<br>**Output:** all inconsistencies $I$ | PTIME | Evaluate regular path queries and validate values<br>- time complexity: $O(|\Sigma||V| + |V|(|V| + |E|))$ |
| Repair | **Input:** $\Sigma$ and $G$ that does not satisfy $\Sigma$<br>**Ouput:** $G'$ that satisfies $\Sigma$ with least repair cost | NP-hard<br>APX-hard | Approximable cases (PTIME checkable)<br>- time complexity $O(|I||\Sigma|^2 + |I|(|I||\Sigma|^2 + |I||\Sigma|))$<br>- approximation ratio: $|I||\Sigma|^2$ |
| | | | Optimal cases<br>- time complexity $O(|I||\Sigma|)$ |
| | | | Heuristic cases<br>- time complexity $O(|I||\Sigma|^2 + |I|(|I||\Sigma|^2 + |I||\Sigma|))$<br>- bounded repairable: cost $\leq |I|$ |

- Notations    $G$: graph      $V$: nodes      $E$: edges
                $\Sigma$: a set of StarFDs    $\varphi$: a single StarFD    $I$: all inconsistencies.

# Updates and repairs

- Updates $O$: operators $o = (v. A, a, c)$ with editing cost $\qquad$ $\text{cost}(O) = \sum_{o \in O} \text{cost}(o)$
- Repair $O$: applying $O$ to $G$, such that obtain $G'$ that satisfies $\Sigma$

# Updates and repairs

- Updates $O$: operators $o = (v.A, a, c)$ with editing cost $\quad \text{cost}(O) = \sum_{o \in O} \text{cost}(o)$
- Repair $O$: applying $O$ to $G$, such that obtain $G'$ that satisfies $\Sigma$



Two repairs:
$O_1 = \{(v_2.\text{city}, \text{LD}, \text{MAN}), (v_4.\text{city}, \text{BZ}, \text{LDN})\}$
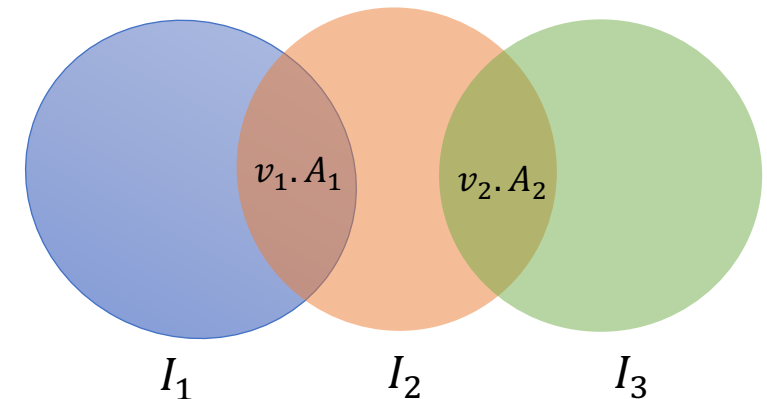$O_2 = \{(v_2.\text{owner}, \text{MUP}, \text{CFG}), (v_4.\text{owner}, \text{EM}, \text{ENIC})\}$

# Entity repair problem

- **Input:** StarFDs $\Sigma$, and graph $G$ does not satisfy $\Sigma$
- **Output:** a repair $O$, such that
  - obtain $G'$ that satisfies $\Sigma$
  - cost$(O) \leq$ cost$(O')$ for any $O'$

# Entity repair problem

- **Input:** StarFDs $\Sigma$, and graph $G$ does not satisfy $\Sigma$
- **Output:** a repair $O$, such that
  - obtain $G'$ that satisfies $\Sigma$
  - $\text{cost}(O) \leq \text{cost}(O')$ for any $O'$

- Solution overview
  - Connected components (CCs): inconsistencies connected at shared node attributes
  - Isolated CCs: no new inconsistency is introduced when a CC is repaired

# Entity repair problem

- **Input:** StarFDs $\Sigma$, and graph $G$ does not satisfy $\Sigma$
- **Output:** a repair $O$, such that
  - obtain $G'$ that satisfies $\Sigma$
  - $\mathrm{cost}(O) \leq \mathrm{cost}(O')$ for any $O'$

- Solution overview
  - Connected components (CCs): inconsistencies connected at shared node attributes
  - Isolated CCs: no new inconsistency is introduced when a CC is repaired

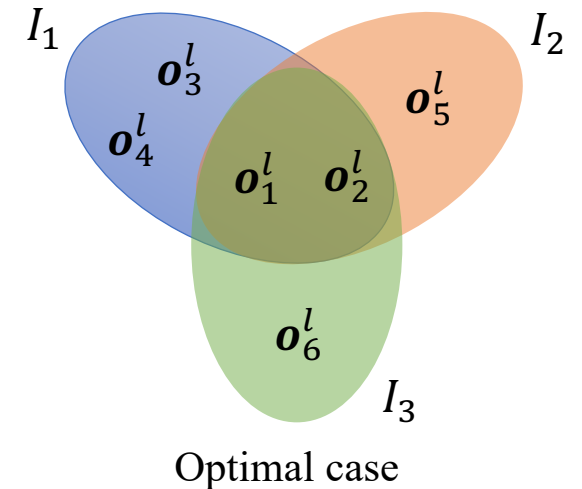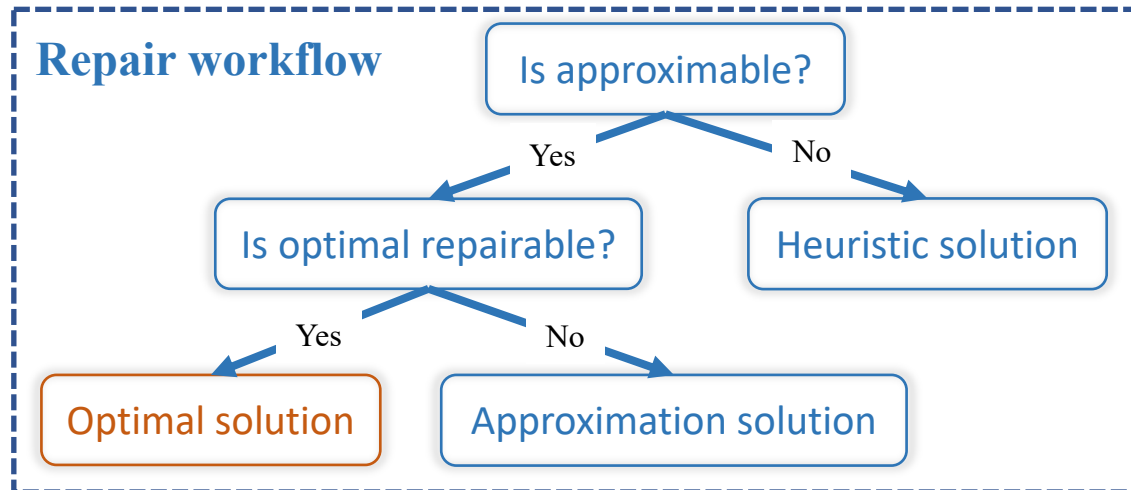Isolated CCs have approximate solutions

**Repair workflow**

Is approximable?

Yes → Is optimal repairable?

No → Heuristic solution

Is optimal repairable?
- Yes → Optimal solution
- No → Approximation solution

$v_1.A_2$  $v_2.A_2$

$I_1$   $I_2$   $I_3$

# Optimal case

- Updates $o^l$: flip the condition of a literal $l$ in $X \cup Y$

- Optimal solution: hyper star structure
  - Select the $o^*$ with least cost in center
  - Select one $o$ with least cost in each petal, and induce $O$
  - If $\text{cost}(o^*) \leq \text{cost}(O)$, return $o^*$; otherwise, return $O$

Example:
- $o^* = o_1^l$
- $O = o_3^l \cup o_5^l \cup o_6^l$
- Return $o^*$ that has less cost



**Repair workflow**

Is approximable?
- Yes → Is optimal repairable?
  - Yes → Optimal solution
  - No → Approximation solution
- No → Heuristic solution

$I_1$

$o_3^l$

$o_4^l$

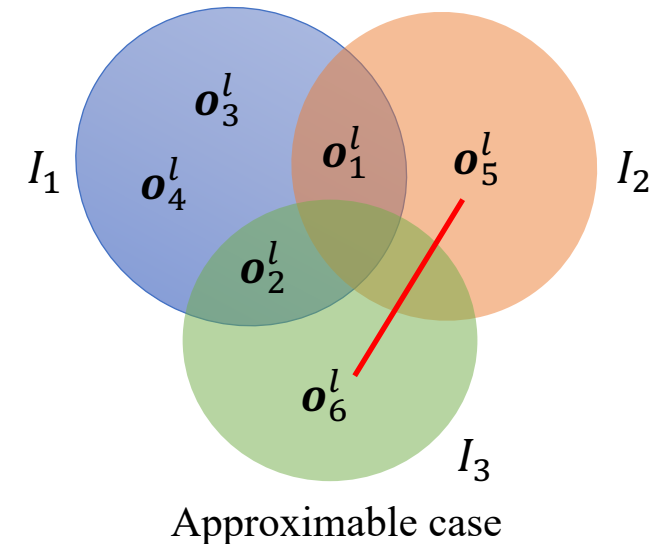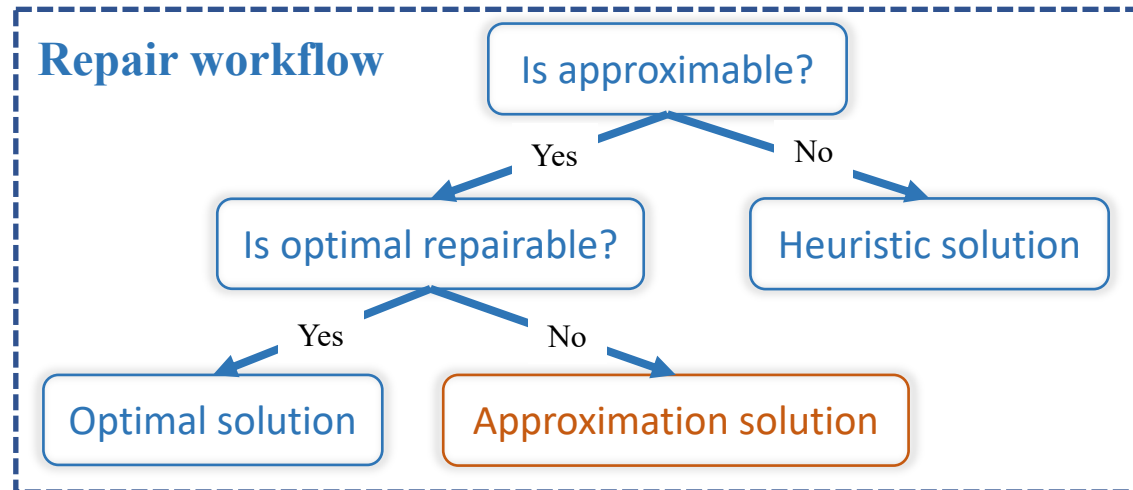$o_1^l$   $o_2^l$

$o_5^l$   $I_2$

$o_6^l$

$I_3$

Optimal case

# Approximable case

- Updates $o^l$: flip the condition of a literal $l$ in $X \cup Y$

- Approximation solution:
  - Hypergraph vertex cover without forbidden pairs
  - Forbidden pairs

    $o_5^l = \{(v_2.\text{owner, MUP, CFG}), (v_4.\text{owner, EM, ENIC})\}$
    $o_6^l = \{(v_2.\text{owner, MUP, FSG}), (v_4.\text{owner, EM, ENIC})\}$

Example:
- Return $\boldsymbol{O} = o_2^l \cup o_5^l$
- $o_6^l$ is pruned

**Repair workflow**

Is approximable?
- Yes → Is optimal repairable?
  - Yes → Optimal solution
  - No → Approximation solution
- No → Heuristic solution

$I_1$   $o_3^l$   $o_1^l$   $o_5^l$   $I_2$

$o_4^l$   $o_2^l$   $o_6^l$   $I_3$

Approximable case

# Heuristic case

- Updates $o^l$: flip the condition of a literal $l$ in $X \cup Y$

- Heuristic solution (for *non-isolated* CC):
  - Select CC introducing fewest inconsistencies
  - Invoke approximation/optimal solution
  - Re-detect inconsistencies
  - Repeat until incur a cost bound

Repair $CC_1$ consisting of $I_1$, $I_2$, and $I_3$

**Repair workflow**

Is approximable?

Yes      No

Is optimal repairable?      Heuristic solution

Yes      No

Optimal solution      Approximation solution
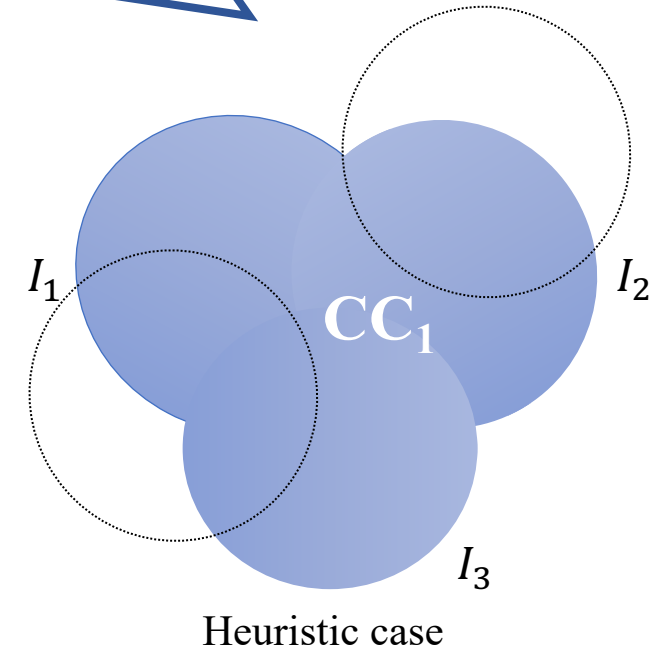
$CC_1$

$I_1$      $I_2$

$I_3$

Heuristic case

# Heuristic case

- Updates $o^l$: flip the condition of a literal $l$ in $X \cup Y$

- Heuristic solution (for *non-isolated* CC):
  - Select CC introducing fewest inconsistencies
  - Invoke approximation/optimal solution
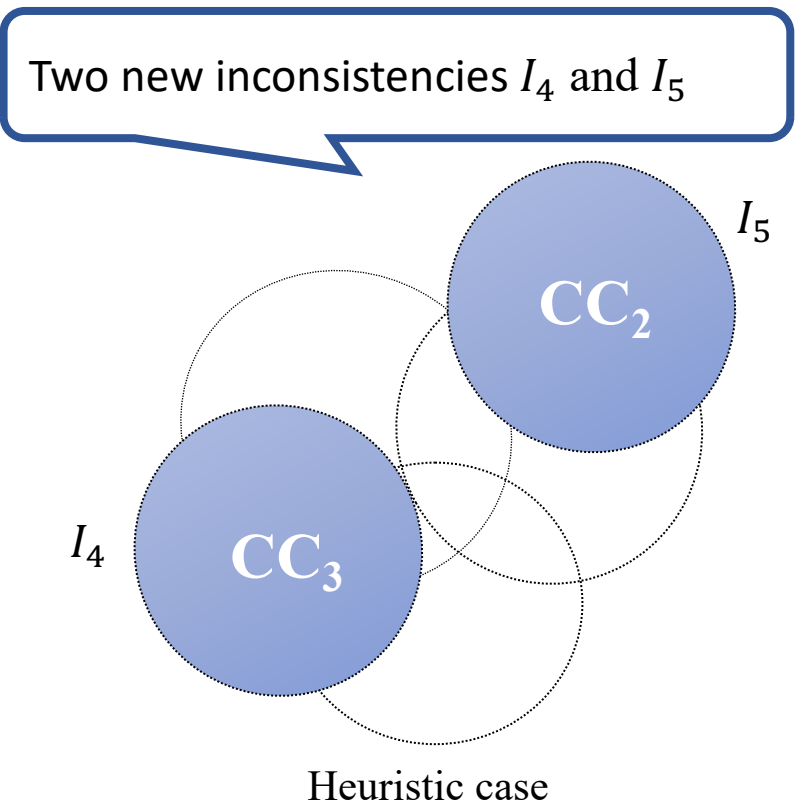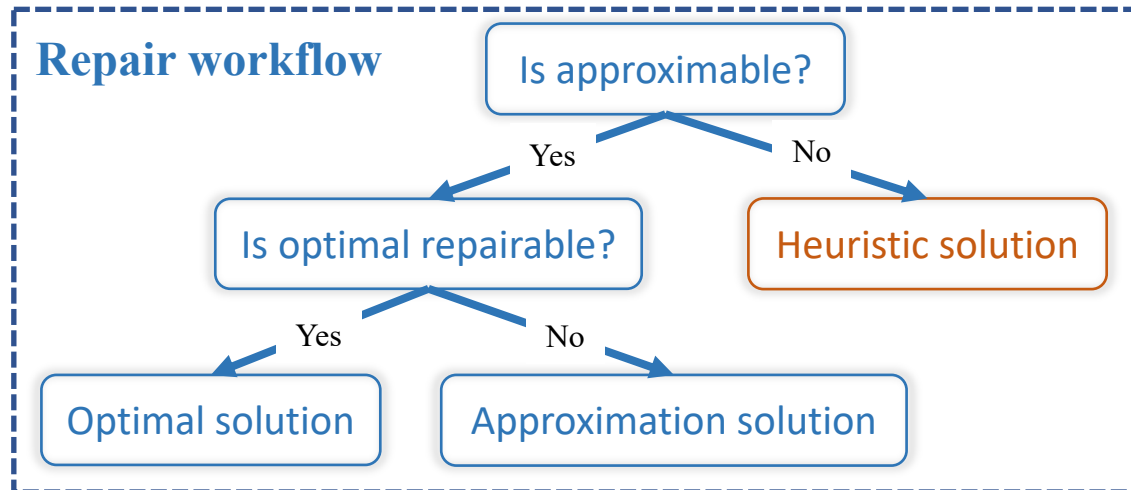  - Re-detect inconsistencies
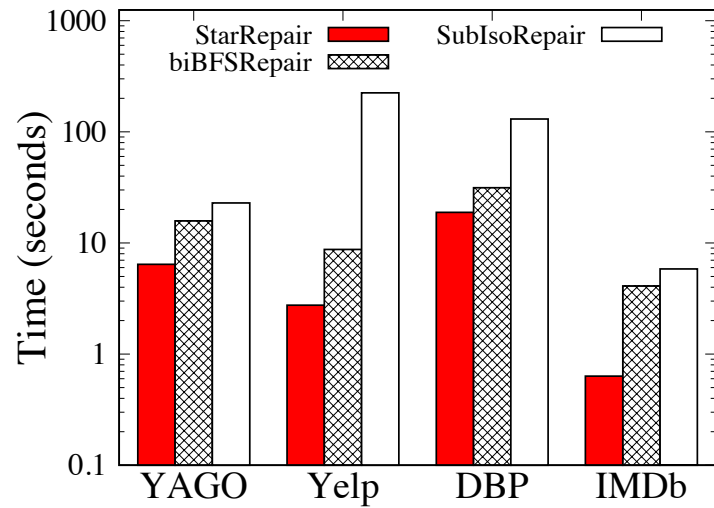  - Repeat until incur a cost bound



Two new inconsistencies $I_4$ and $I_5$

Heuristic case

# Experiment settings

- Datasets

| Data | Description | # of nodes | # of edges | avg. # of attributes per node |
|------|-------------|------------|------------|-------------------------------|
| Yago | Knowledge graph | 2.1M | 4.0M | 3 |
| DBPedia | Knowledge graph | 2.2M | 7.4M | 4 |
| Yelp | Business reviews | 1.5M | 1.6M | 5 |
| IMDb | Movie network | 5.9M | 3.2M | 3 |

- Error generation: adopt silver standard and an error generation benchmark (Arocena et al. 2015)

- StarFD generation: discovered from silver standard (first star patterns and then value constraints)

- Algorithms:

  - **StarRepair:**   use bidirectional search for regular path queries with incremental error detection

  - **biBFSRepair:**  use bidirectional search *without incremental error detection*

  - **SubIsoRepair:** *use subgraph isomorphism* as matching semantics with incremental error detection

# Experiment results

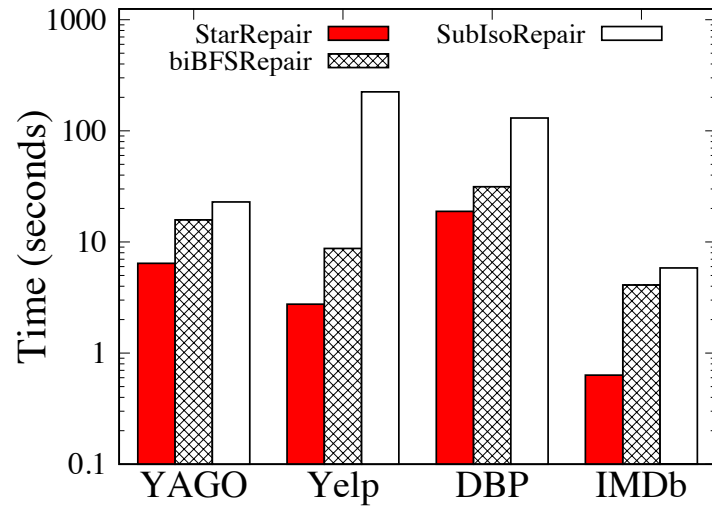■ StarFD repairs: efficiency and effectiveness



StarRepair outperforms biBFSRepair and SubIsoRepair by 3.4 and 7.1 times respectively

■ Case study

# Experiment results

- StarFD repairs: efficiency and effectiveness



StarRepair outperforms biBFSRepair and SubIsoRepair by 3.4 and 7.1 times respectively

StarRepair outperforms SubIsoRepair by 10% in f-score (9% in precision and 14% in recall)
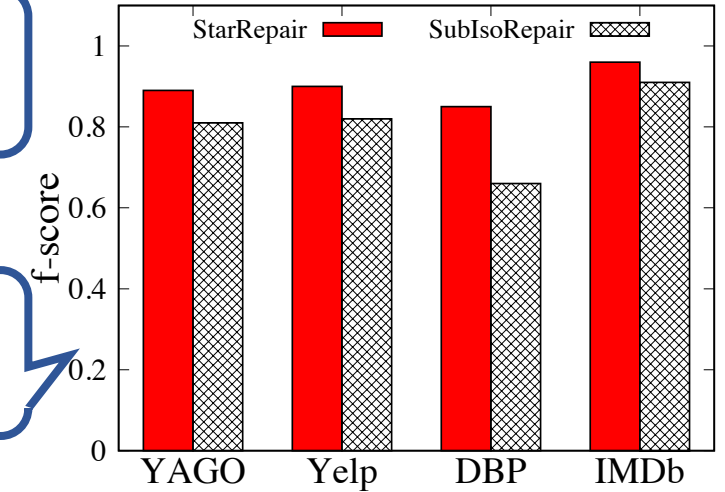
- Case study

# Experiment results

- StarFD repairs: efficiency and effectiveness



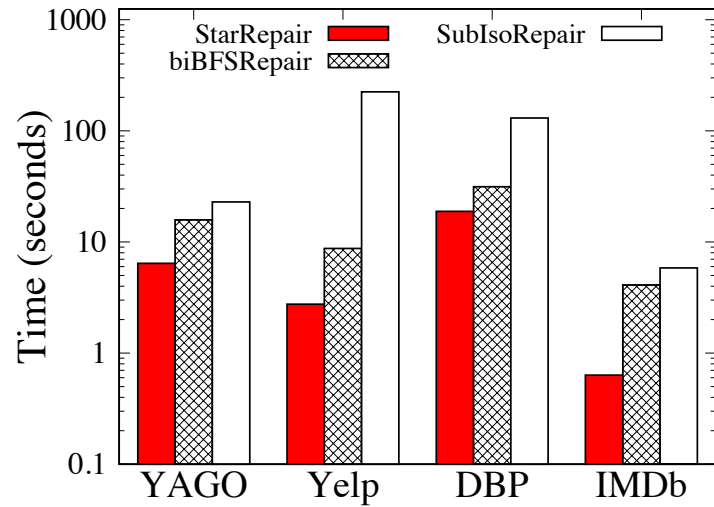StarRepair outperforms biBFSRepair and SubIsoRepair by 3.4 and 7.1 times respectively

StarRepair outperforms SubIsoRepair by 10% in f-score (9% in precision and 14% in recall)
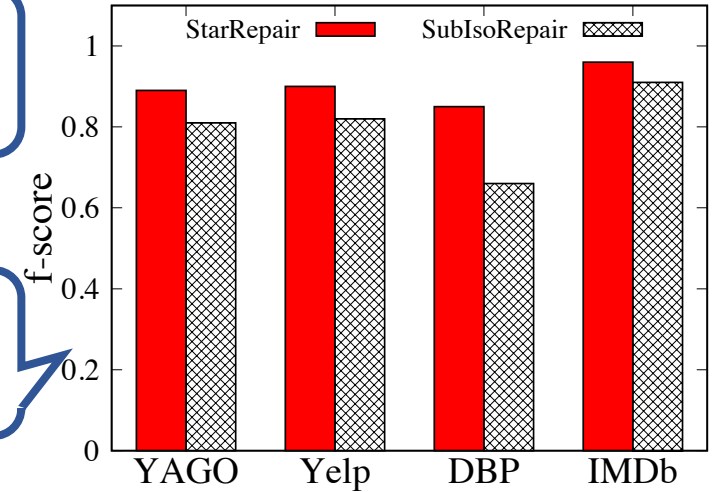
- Case study

**StarFD:** If a person $u_0$ is a politician or president of U.S., and is married to another person $u_1$, then $u_1$'s child is $u_0$'s child.

We found more than 100 such errors in Yago.



$u_0$

Person

$R_1$ = presidentOf
∪ politicianOf

$R_2$ = marriedTo

Country
$u_1$

Person
$u_2$

$u_0$

Person
name: G.W. Bush
child: B. Obama

presidentOf

marriedTo

Country
name: U.S.
$u_1$

Person
name: Laura Bush
child: Barbara Bush
$u_2$

# Compare with GFDs (Fan et al. 2016)

- StarFDs: star functional dependencies
    - Definition: $\varphi = (P(u_o), X \rightarrow Y)$

- GFDs: graph functional dependencies
    - Definition: $\varphi = (P, X \rightarrow Y)$

| Problem | StarFDs | GFDs |
|---|---|---|
| Semantic | star patterns with regex queries | subgraph isomorphism |
| Satisfiability | NP-complete | coNP-complete |
| Implication | coNP-hard | NP-complete |
| Error detection (validation) | PTIME | coNP-complete |

# Repairing Entities using Star Constraints in Multi-relational Graphs

| Problem | Description | Hardness | Solution |
|---------|-------------|----------|----------|
| Satisfiability | **Input:** Σ<br>decide whether there exists $G$ that satisfies Σ | NP-complete | |
| Implication | **Input:** Σ and $\varphi$<br>decide whether for all $G$ satisfy Σ, they satisfy $\varphi$ | coNP-hard | |
| Error detection (validation) | **Input:** $G$ and Σ<br>**Output:** all inconsistencies $I$ | PTIME | Evaluate regular path queries and validate values<br>- time complexity: $O(|\Sigma||V| + |V|(|V| + |E|))$ |
| Repair | **Input:** Σ and $G$ that does not satisfy Σ<br>**Ouput:** $G'$ that satisfies Σ with least repair cost | NP-hard<br>APX-hard | Approximable cases (PTIME checkable)<br>- time complexity $O(|I||\Sigma|^2 + |I|(|I||\Sigma|^2 + |I||\Sigma|))$<br>- approximation ratio: $|I||\Sigma|^2$ |
| | | | Optimal cases<br>- time complexity $O(|I||\Sigma|)$ |
| | | | Heuristic cases<br>- time complexity $O(|I||\Sigma|^2 + |I|(|I||\Sigma|^2 + |I||\Sigma|))$<br>- bounded repairable: cost $\leq |I|$ |

▪ Notations    $G$: graph      $V$: nodes      $E$: edges

                       Σ: a set of StarFDs    $\varphi$: a single StarFD    $I$: all inconsistencies.

# Thank you!



**Kronos:** Lightweight Knowledge-based Event Analysis in Cyber-Physical Data Streams
To appear in Demo Session