

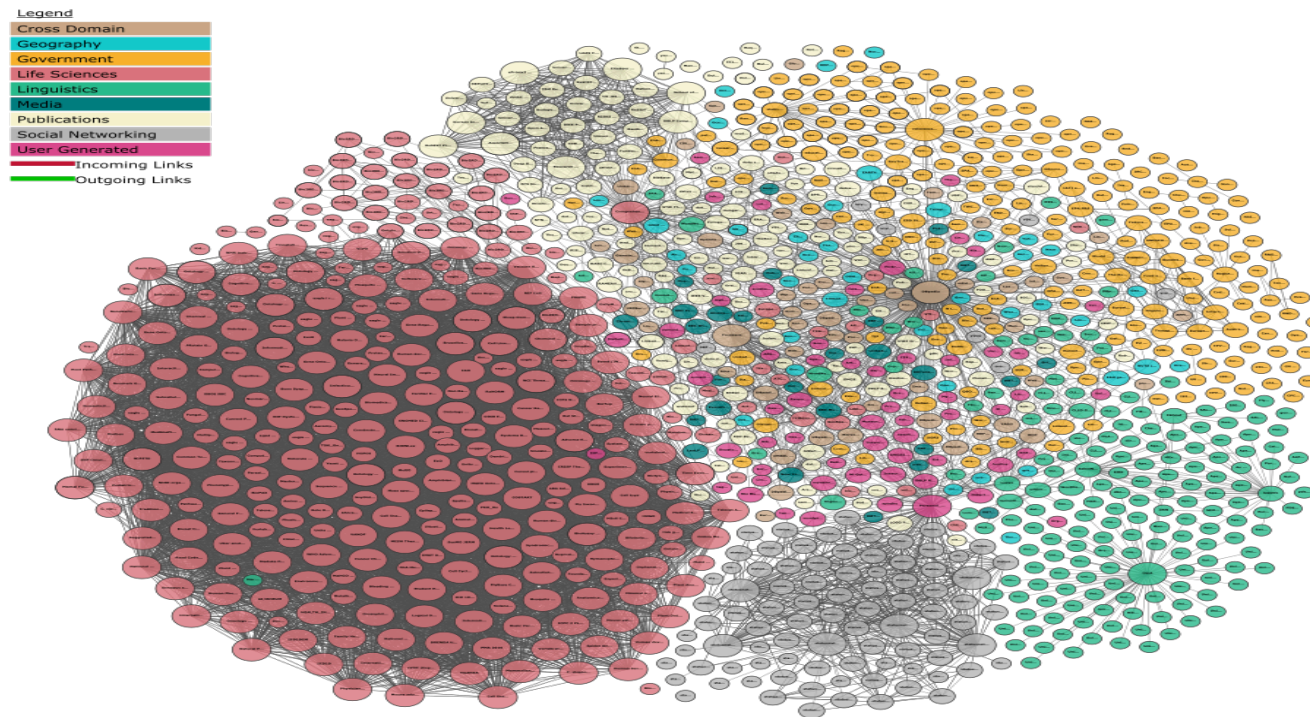


SEMANTIC LABELING

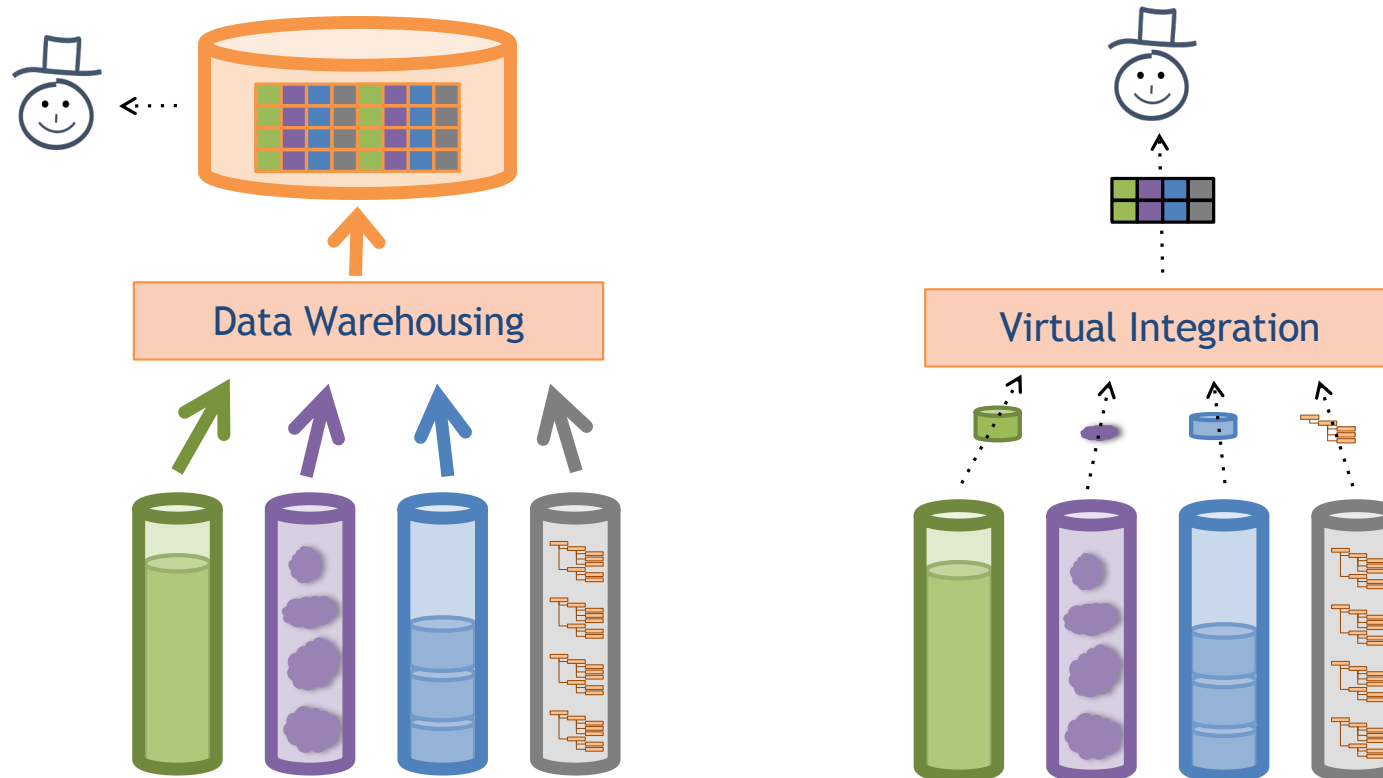
DSCI 558: Building Knowledge Graphs
Craig Knoblock

Based on slides by Pedro Szekely, Minh Pham, & S.K. Ramnandan

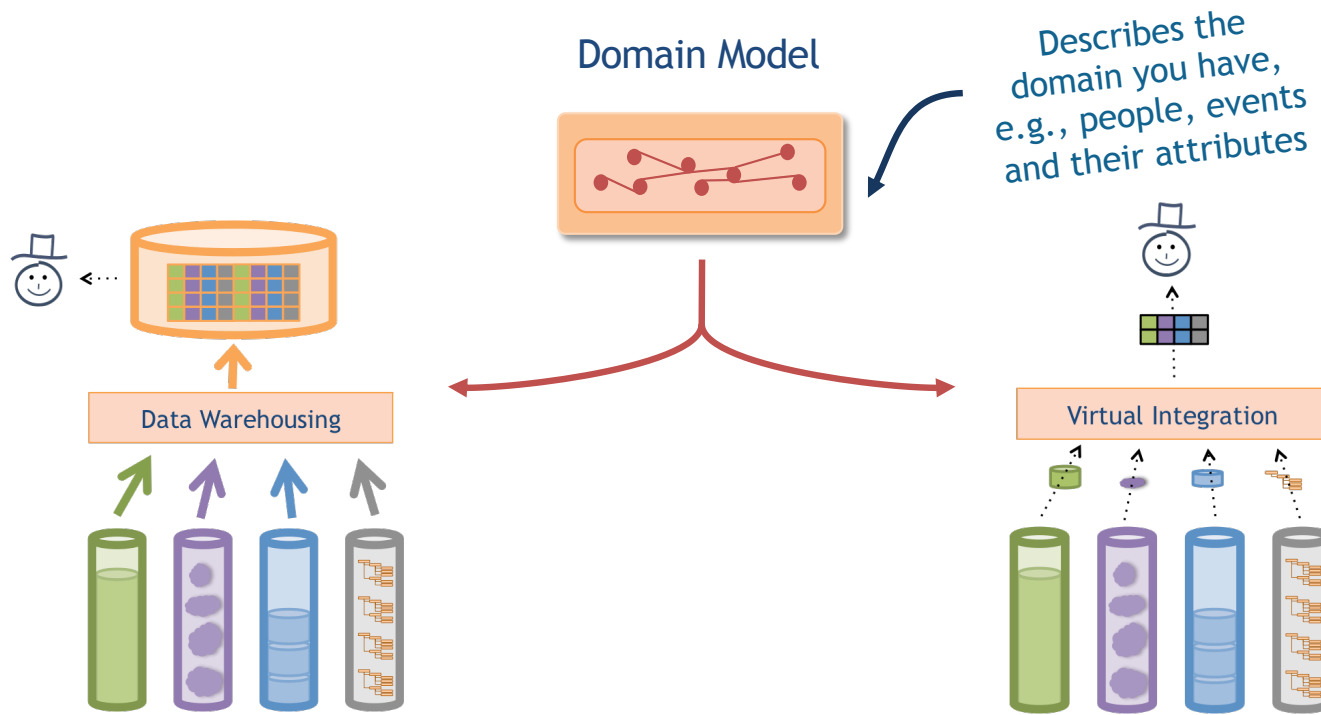
Introduction



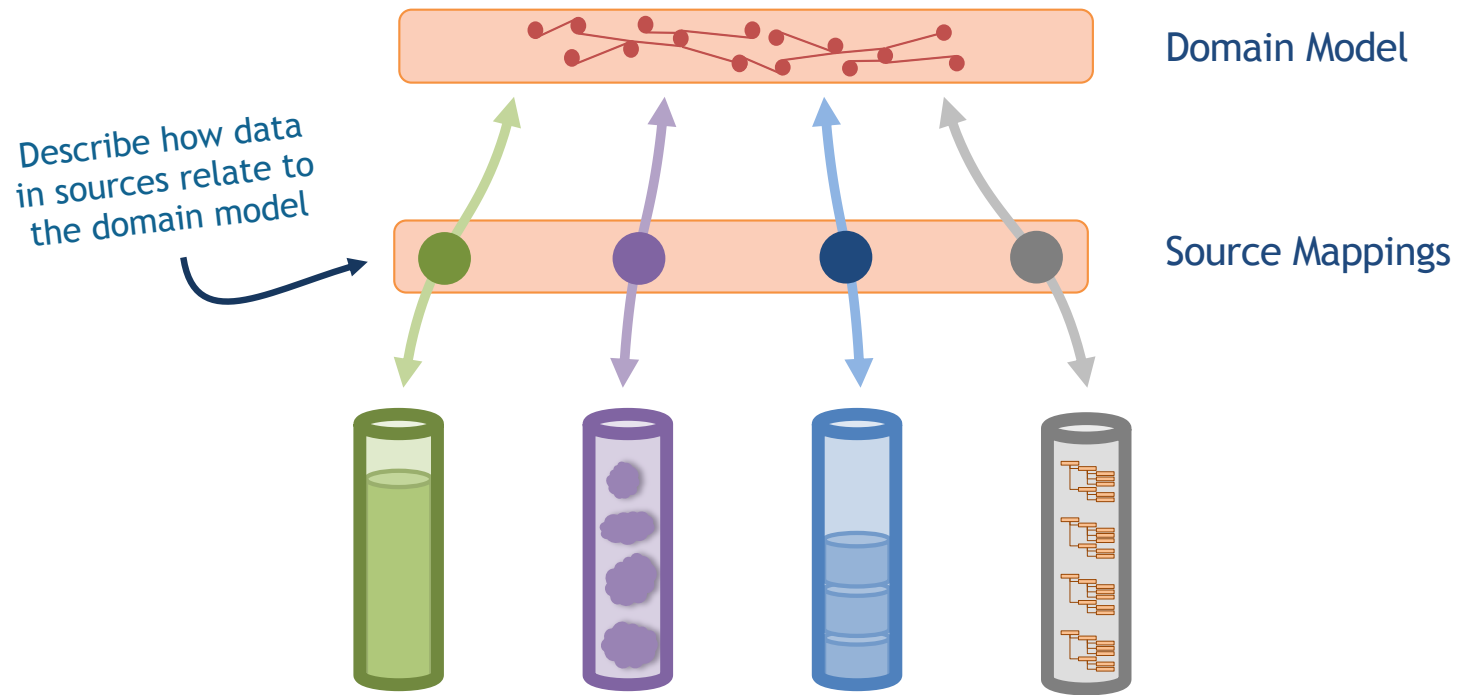
Data Integration Approaches



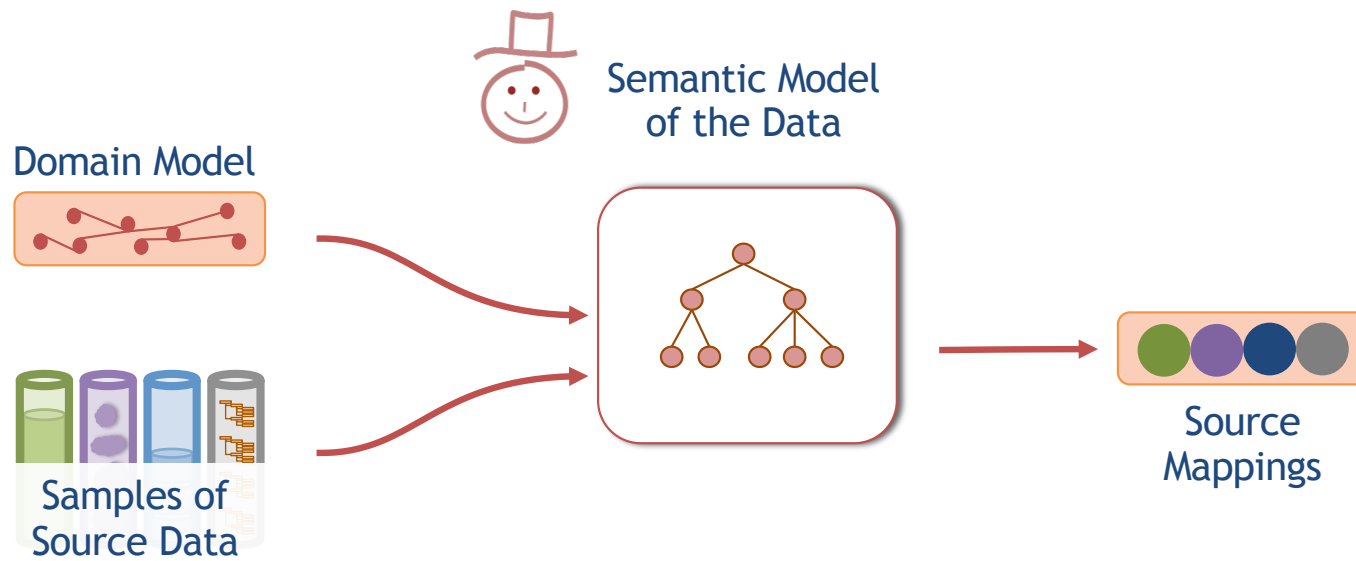
Domain Model



Key Ingredient: Source Mappings



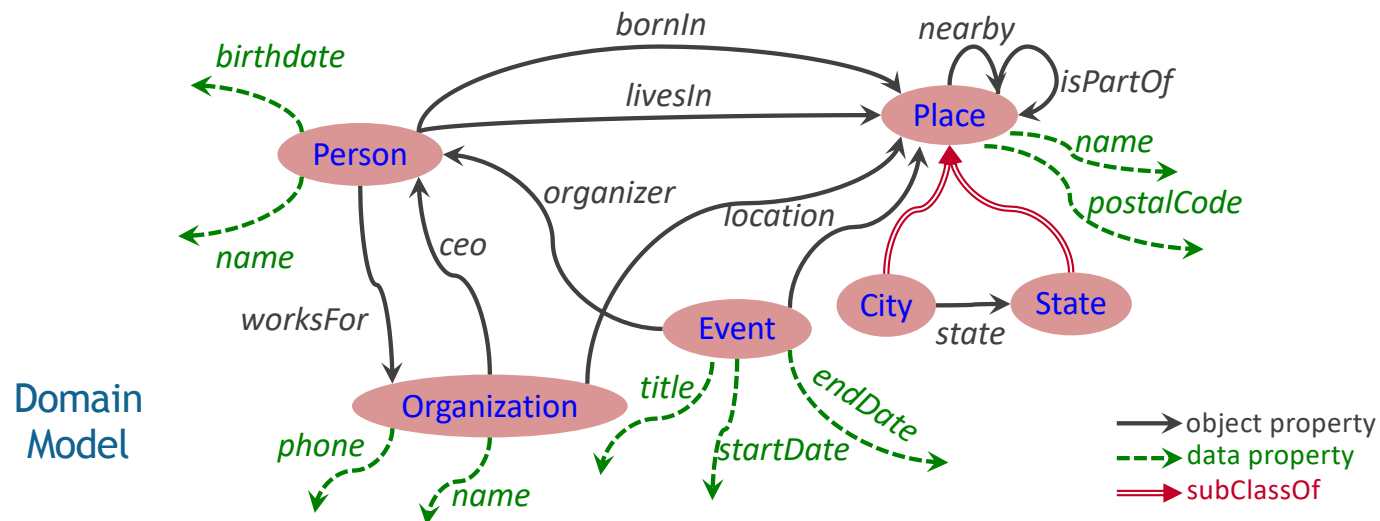
Automatic Source Modeling



What is a Semantic Model?

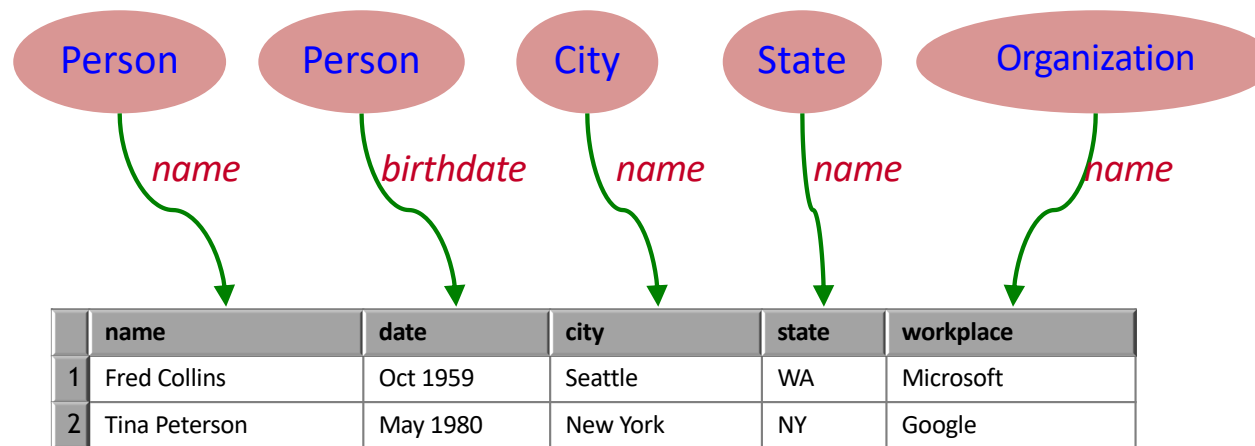
Source

	name	date	city	state	workplace
1	Fred Collins	Oct 1959	Seattle	WA	Microsoft
2	Tina Peterson	May 1980	New York	NY	Google

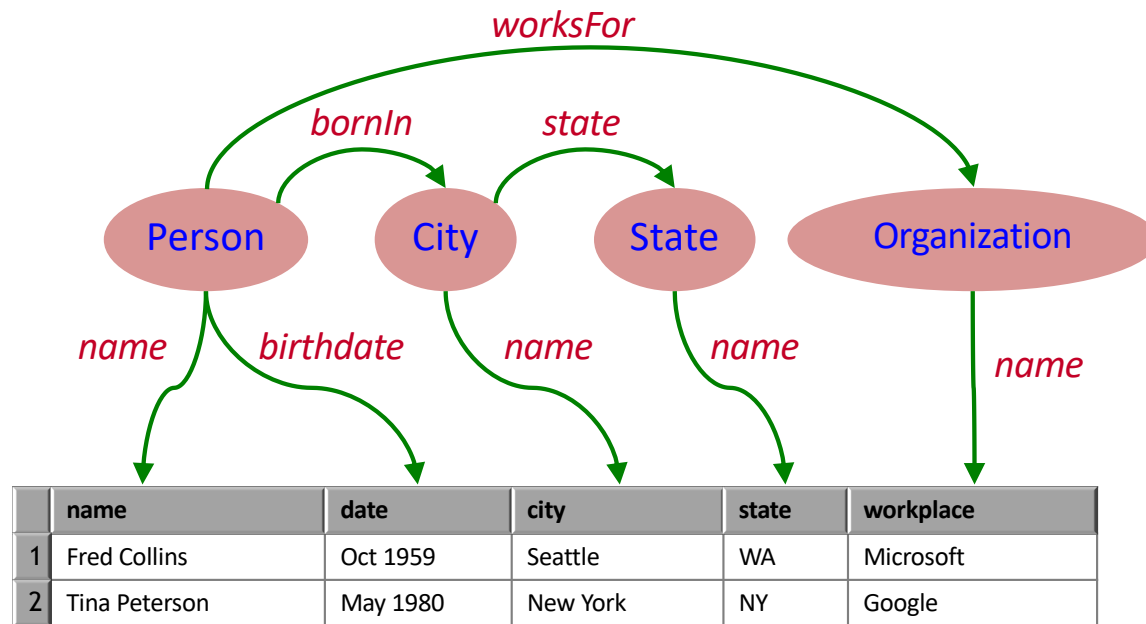


Describe sources using classes & relationships in an ontology

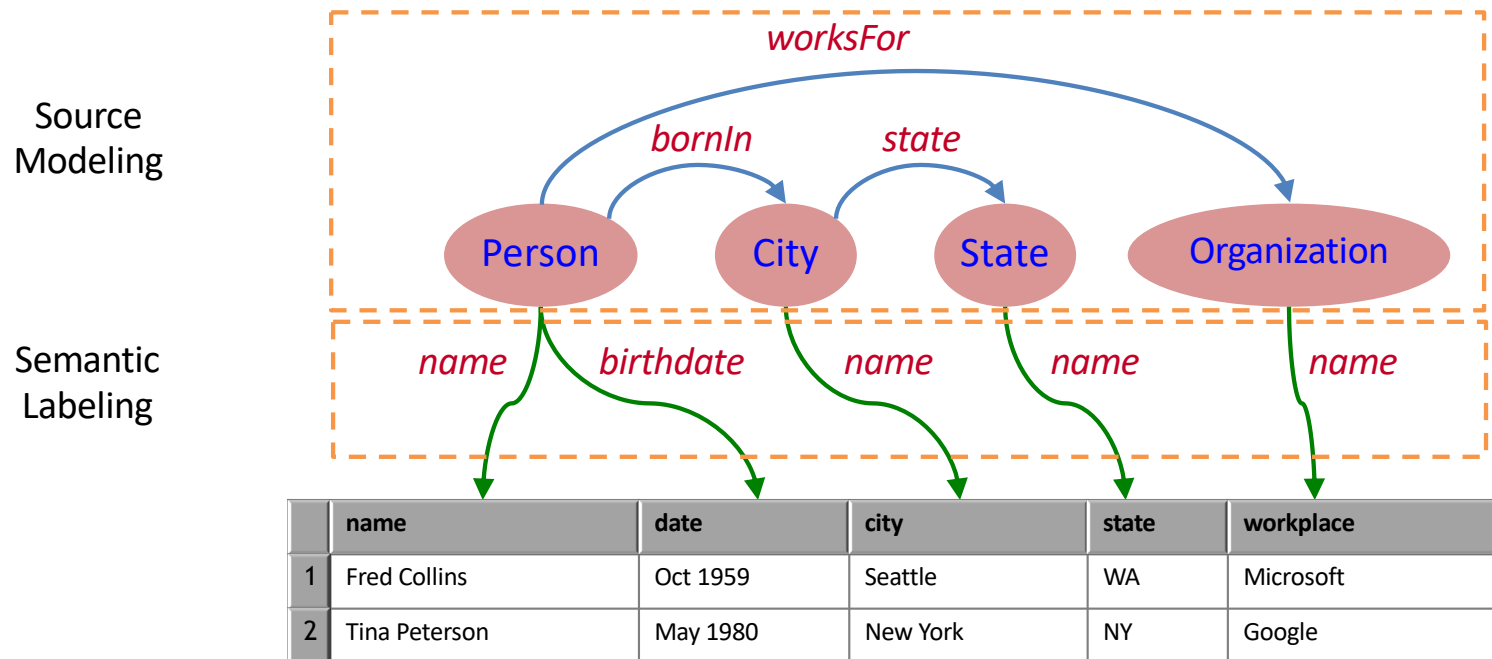
Semantic Types



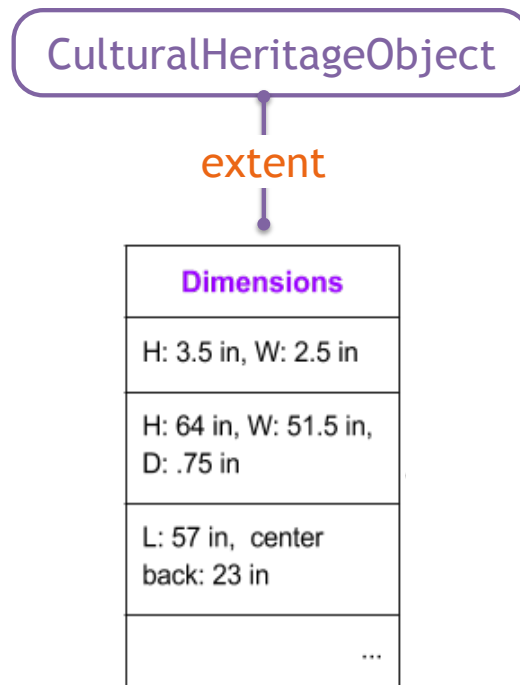
Relationships



Source Modeling Problems



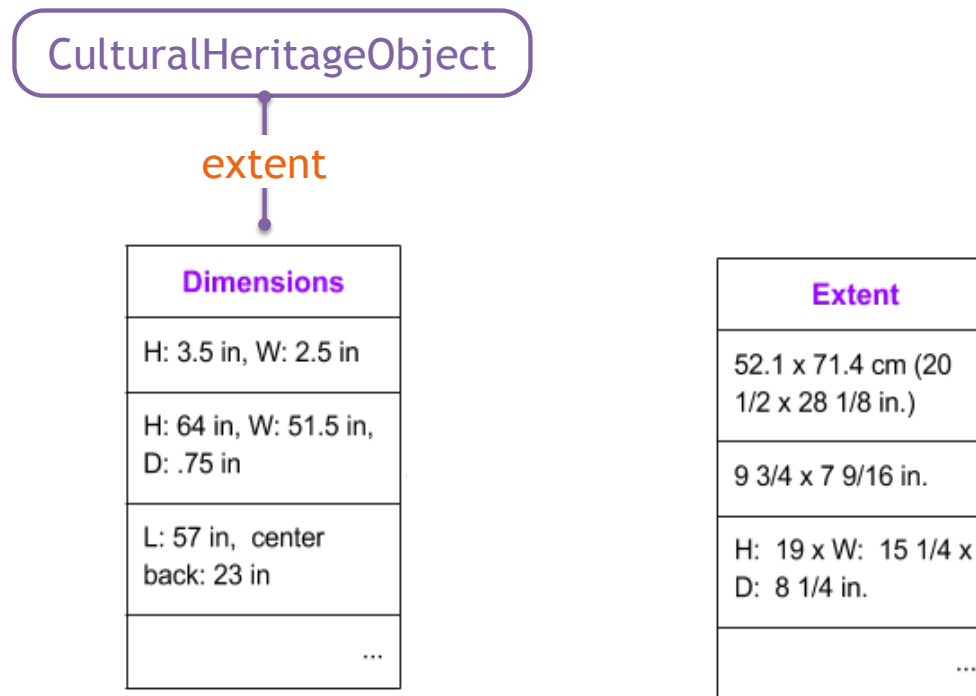
Learning Semantic Types



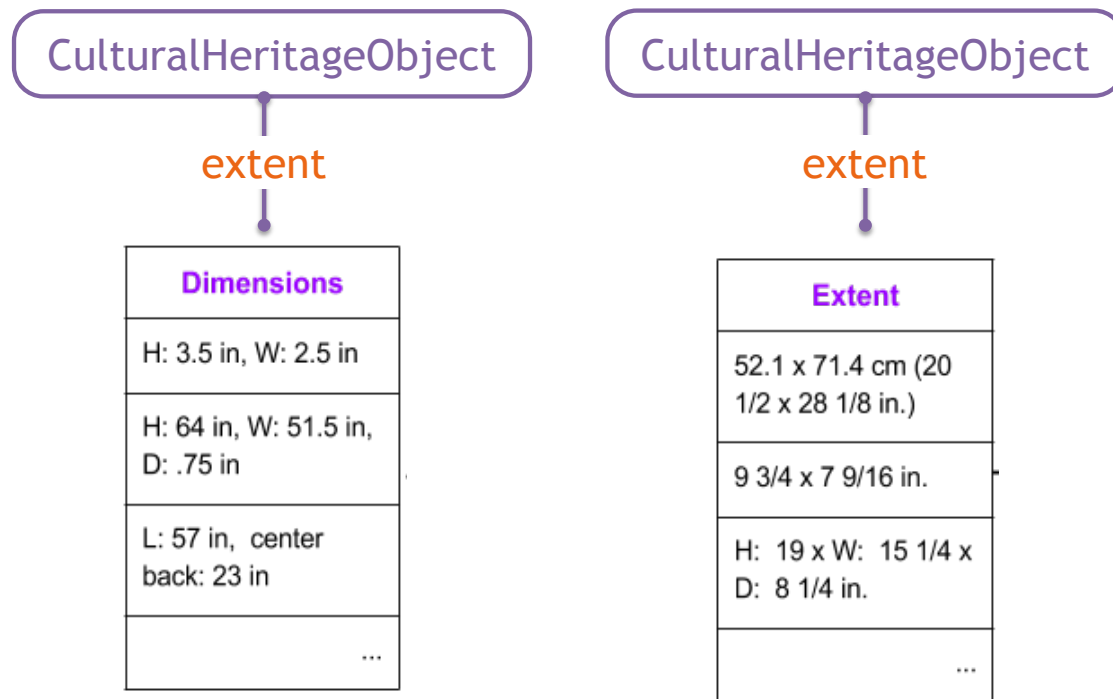
1- User specifies

2- System learns

Learning Semantic Types



Learning Semantic Types



Requirements

- Learn from a small number of examples
- Work on both textual and numeric values
- Learn quickly and highly scalable to large number of semantic types



RULE-BASED APPROACH

Assigning Semantic Labels to Data Sources

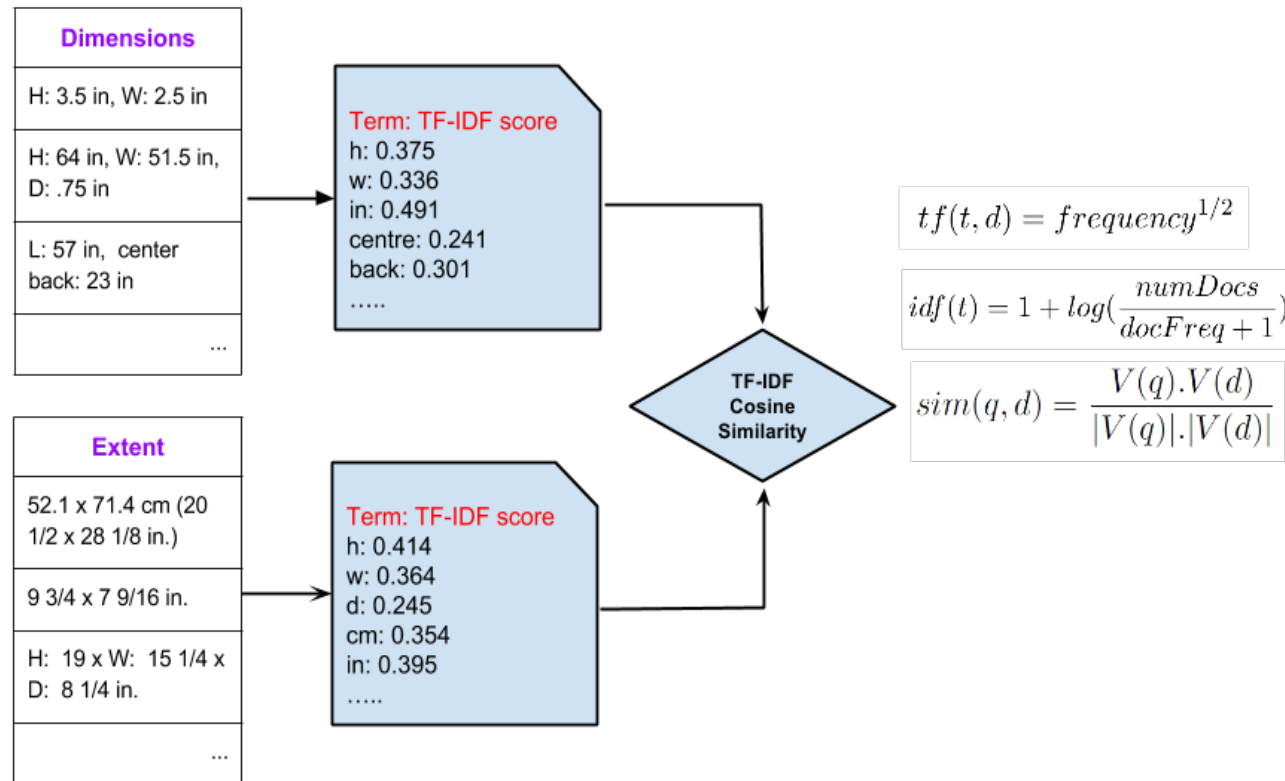
Ramnandan, S.K.; Mittal, A.; Knoblock, C. A.; and Szekely, P.

Approach for Textual Data

- **Document**: each column of data
- **Label**: each semantic type
- Use **Apache Lucene** to index the labeled documents
- Compute **TF/IDF vectors** for documents
- Compare documents using **Cosine Similarity** between TF/IDF vectors

Dimensions
H: 3.5 in, W: 2.5 in
H: 64 in, W: 51.5 in, D: .75 in
L: 57 in, center back: 23 in
...

Approach for Textual Data

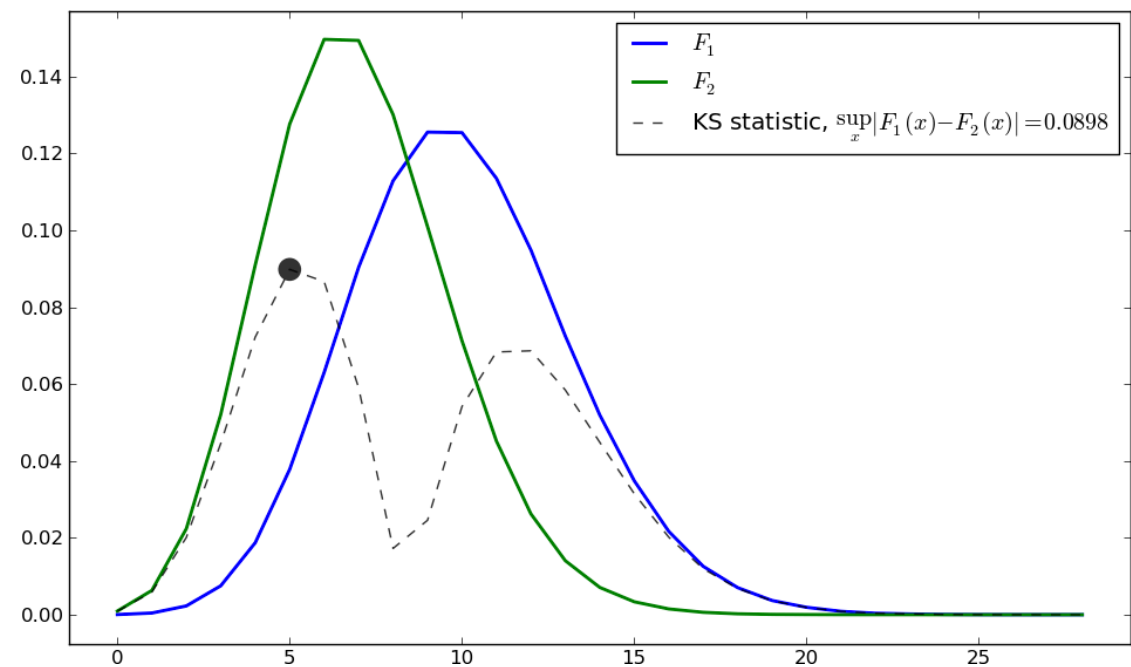
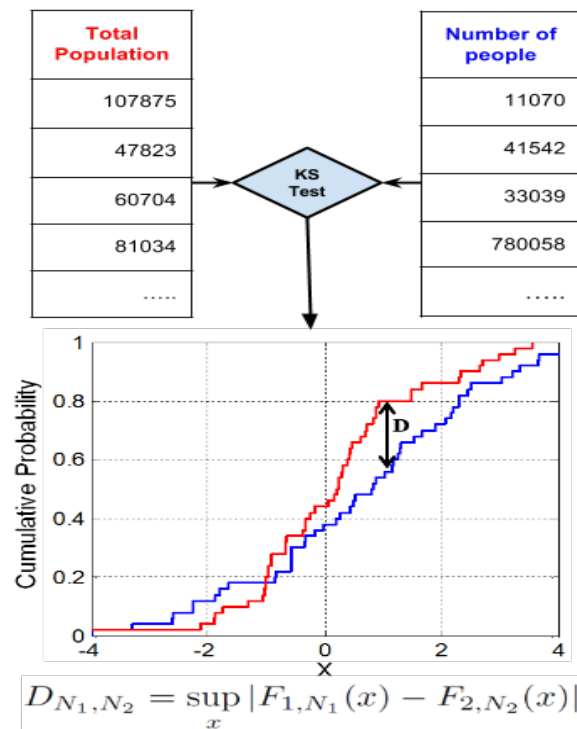


Approach for Numeric Data

- **Distribution** of values in different semantic types is different, e.g., temperature vs. population
- Use **Statistical Hypothesis Testing** to see which distribution fits best
- Welch's T-test, Mann-Whitney U-test and **Kolmogorov-Smirnov Test**

Total Population	Number of people
107875	11070
47823	41542
60704	33039
81034	780058
.....

Approach for Numeric Data

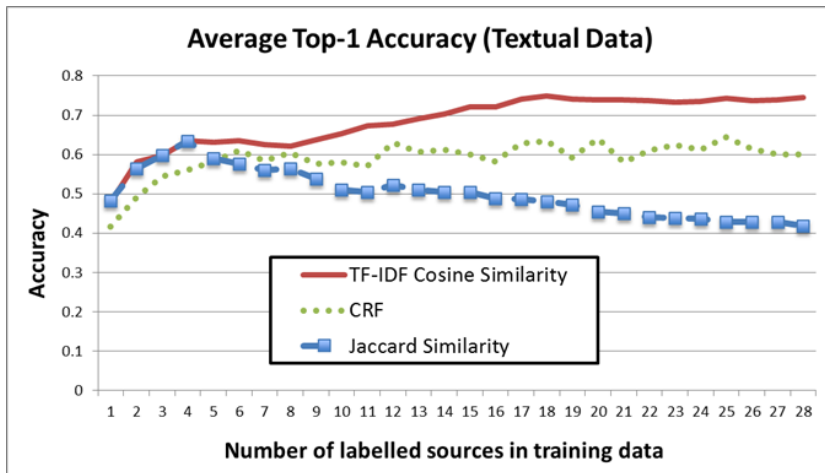


Combined Approach

- **Training**
 - Add new example data as training for either textual or numeric types
 - If ambiguous, train as both textual and numeric
- **Testing**
 - If textual, apply tf/idf
 - If numeric apply KS-test
 - If ambiguous and at least 70% numeric apply KS-test, otherwise tf/idf

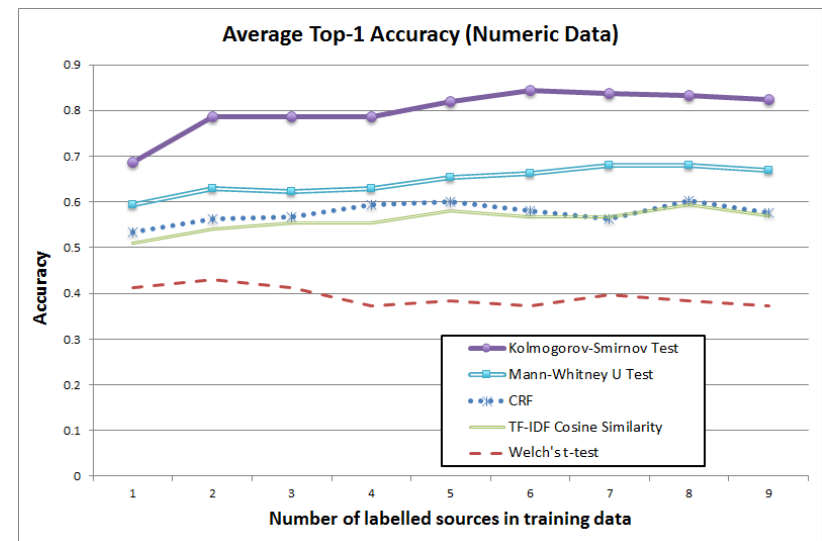
Return Top-k suggestions based on the confidence scores

Evaluation of Semantic Typing



Combined approach achieves 97% accuracy on the top-4 accuracy

Reduced the training time from 110s to 0.45s



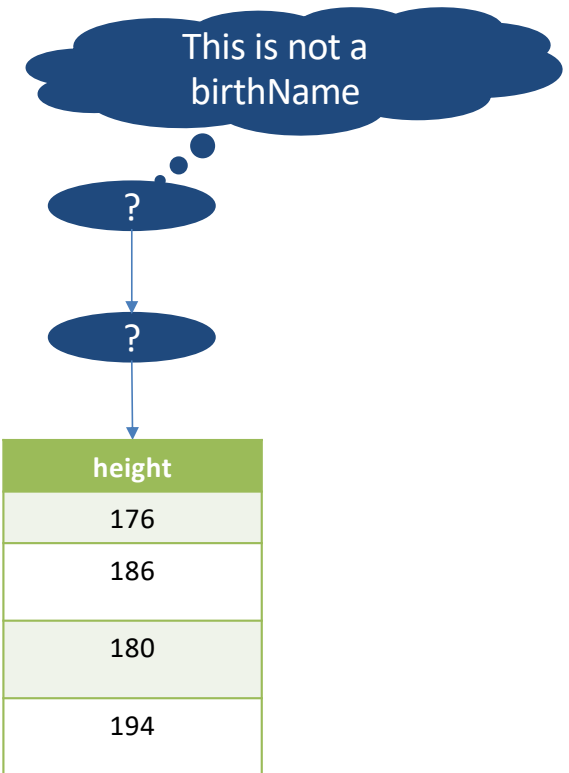
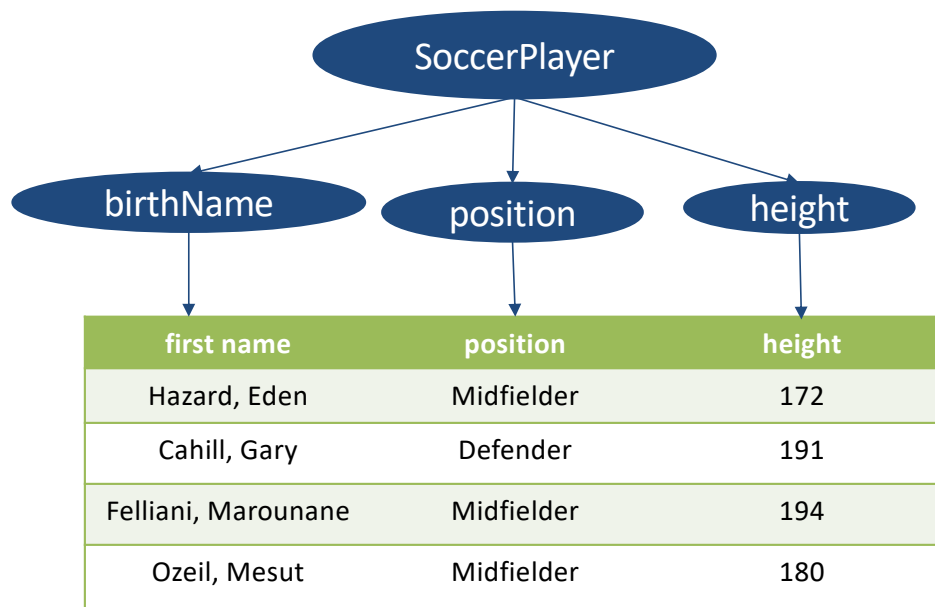


LEARNING-BASED APPROACH: CLASSIFICATION

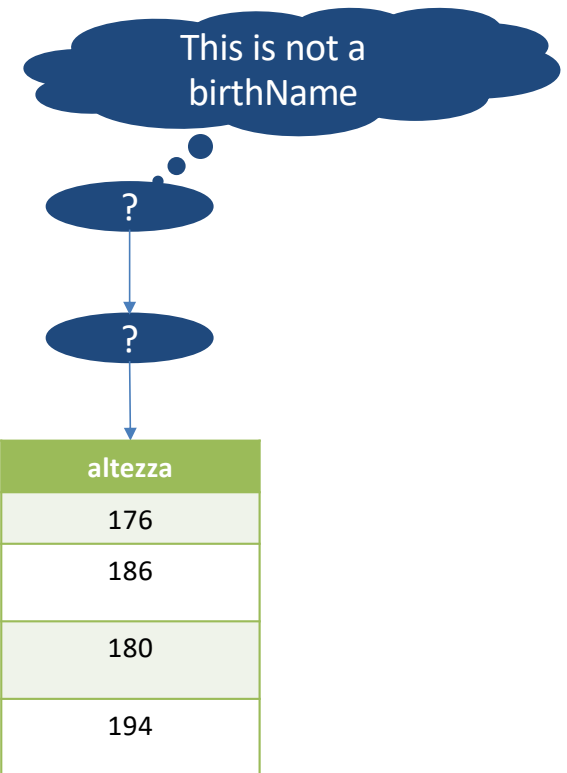
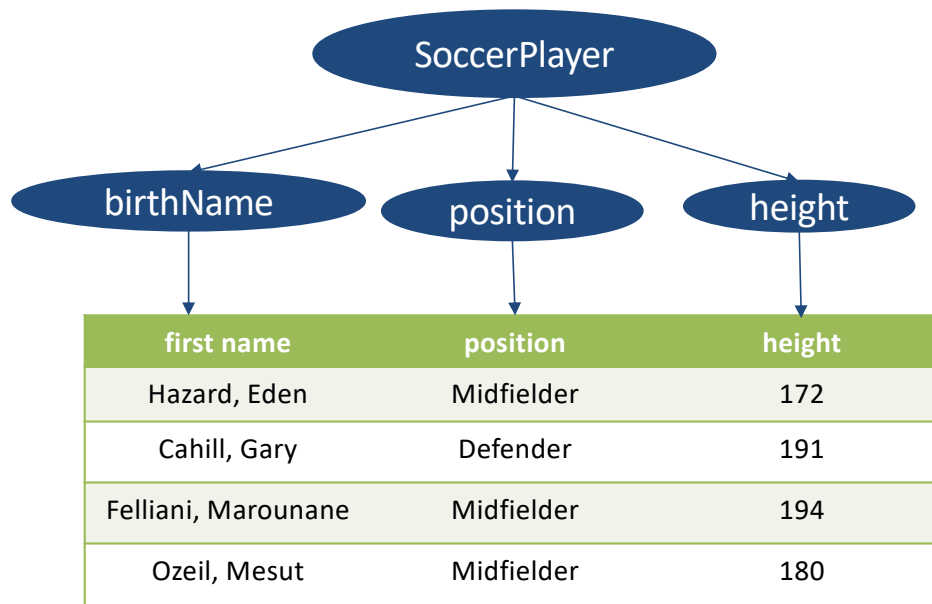
Semantic labeling: a domain-independent approach

Minh Pham, Suresh Alse, Craig Knoblock, Pedro Szekely

General idea

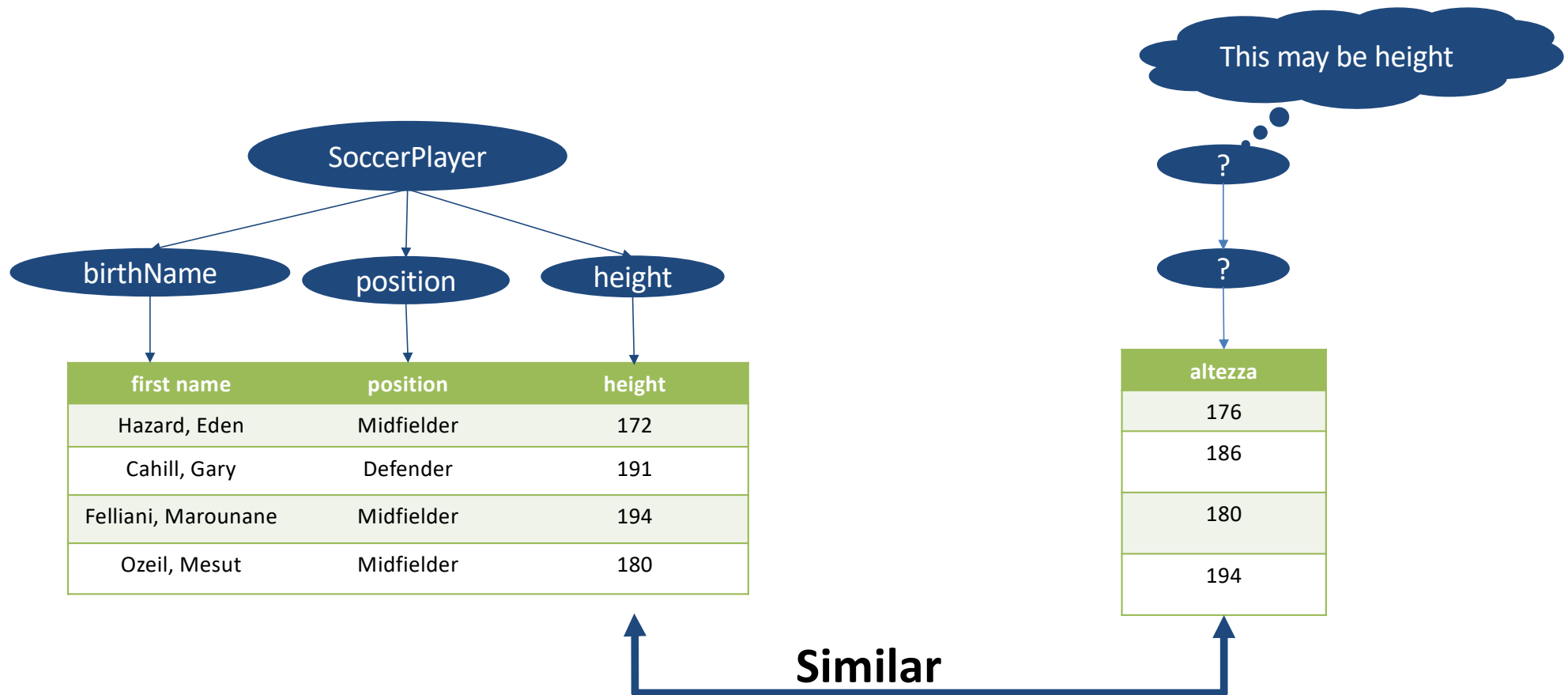


General idea

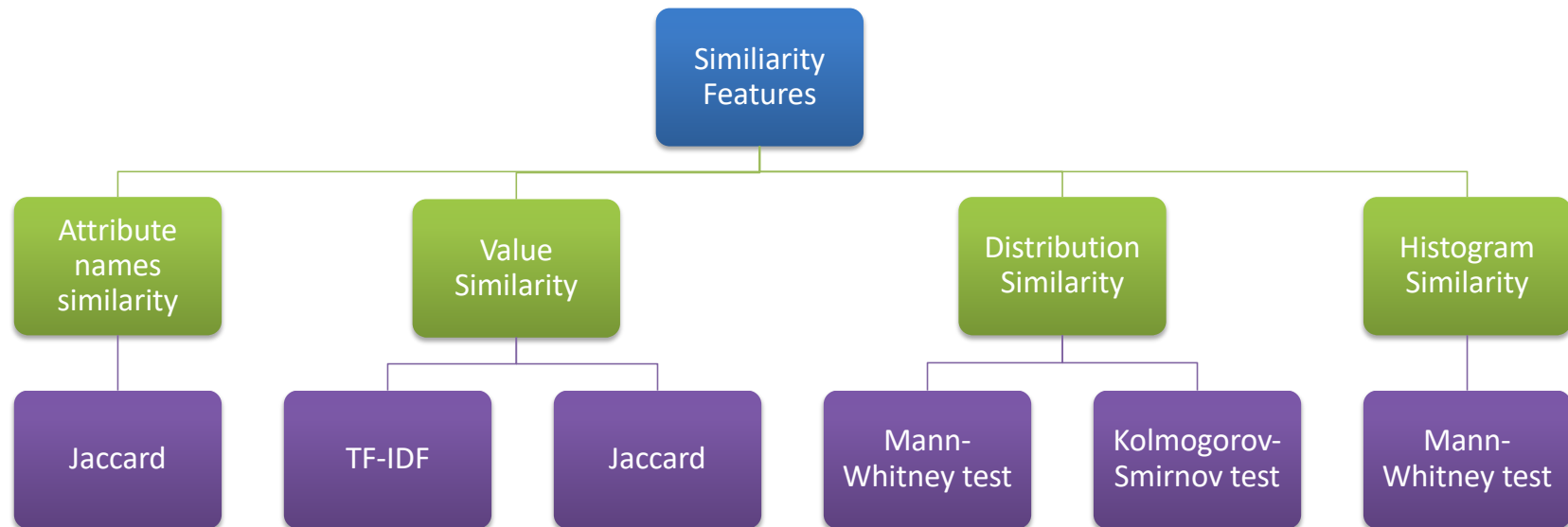


Different

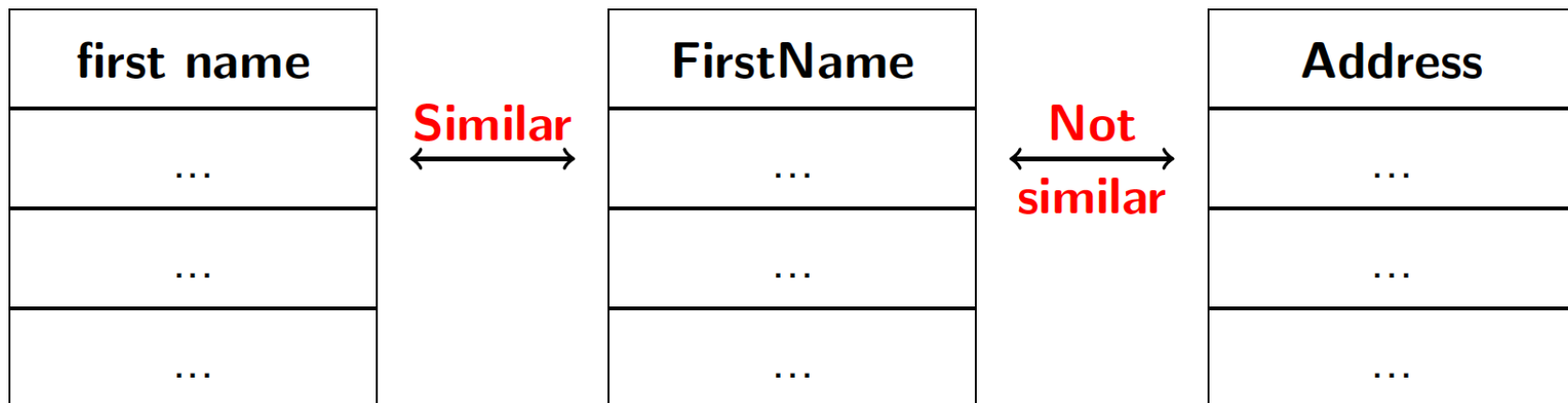
General idea



Similarity features



Attribute name similarity



Value similarity



Value similarity

# game played		# goal scored
4		3
...		...
18		11
23		22

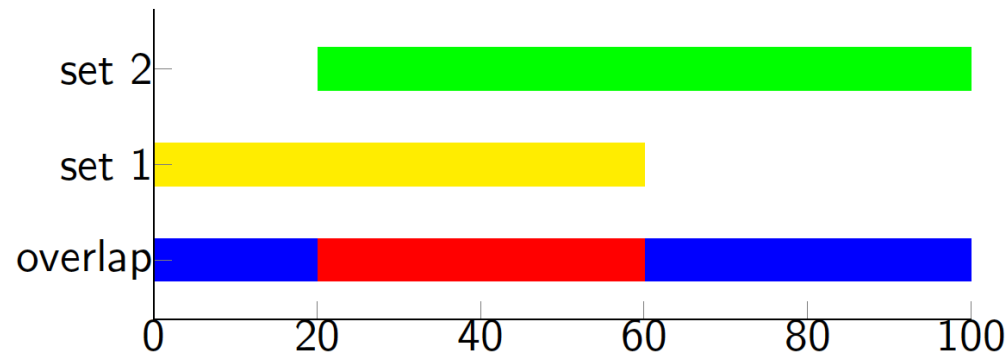
**Overlapping values
is not enough**

Value range similarity

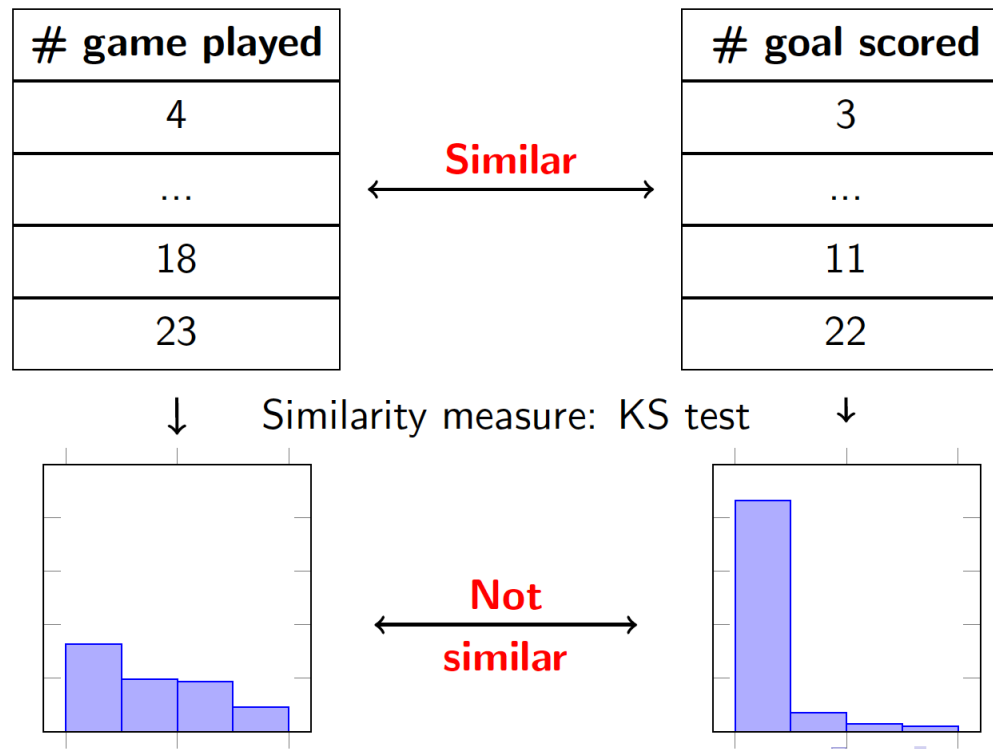
Numeric Jaccard Similarity

Given 2 numeric sets of values A, B ranged in $[a_s, a_e]$ and $[b_s, b_e]$:

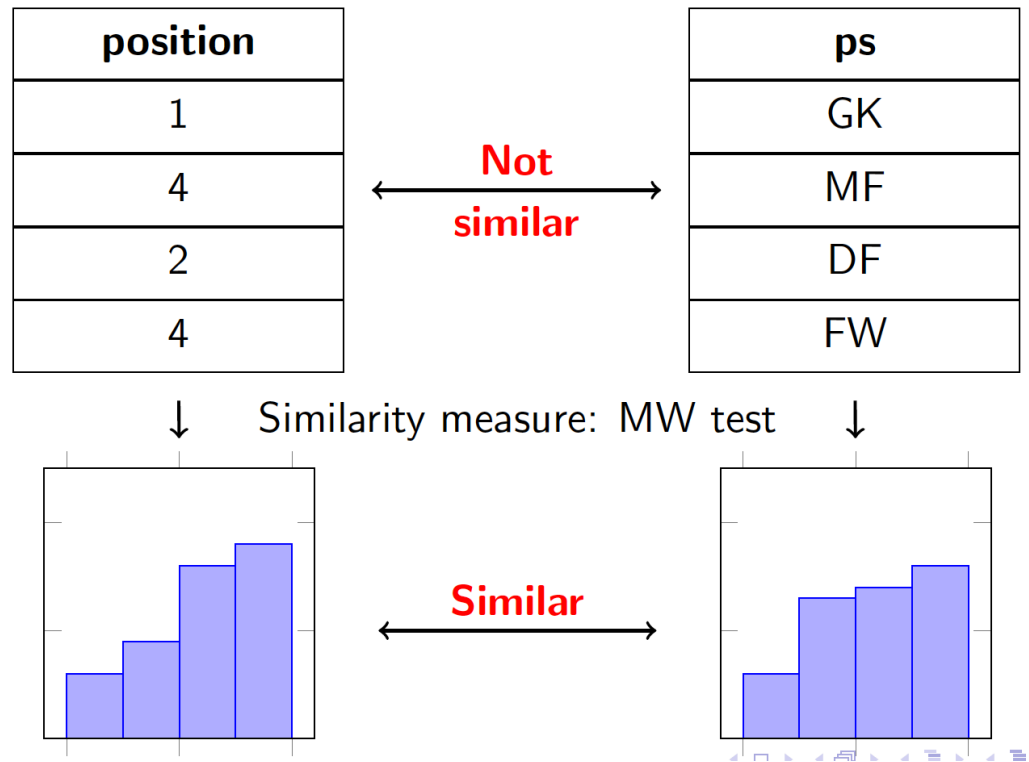
$$\text{numJaccardSim}(A, B) = \frac{|[a_s, a_e] \cap [b_s, b_e]|}{|[a_s, a_e] \cup [b_s, b_e]|}$$



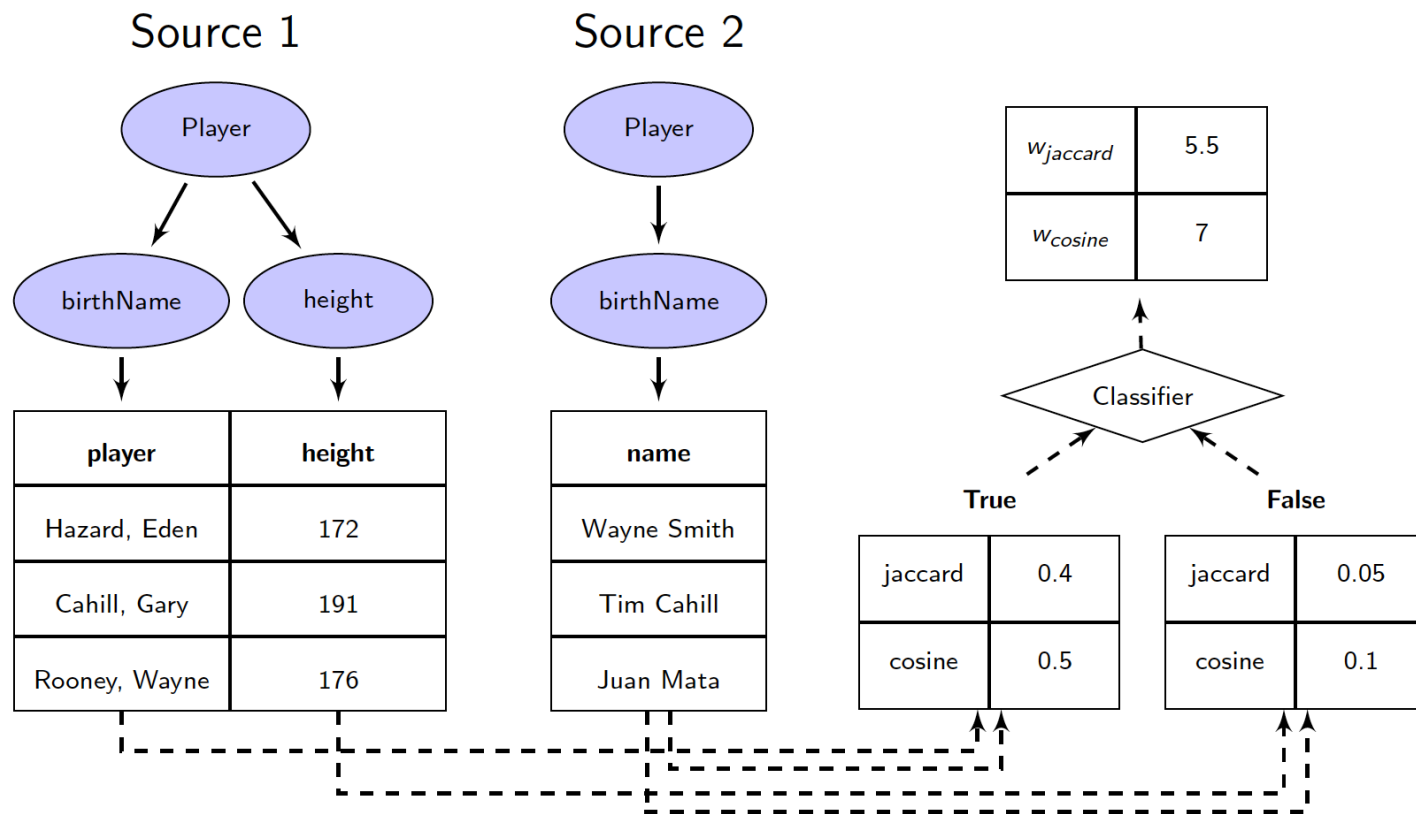
Distribution similarity



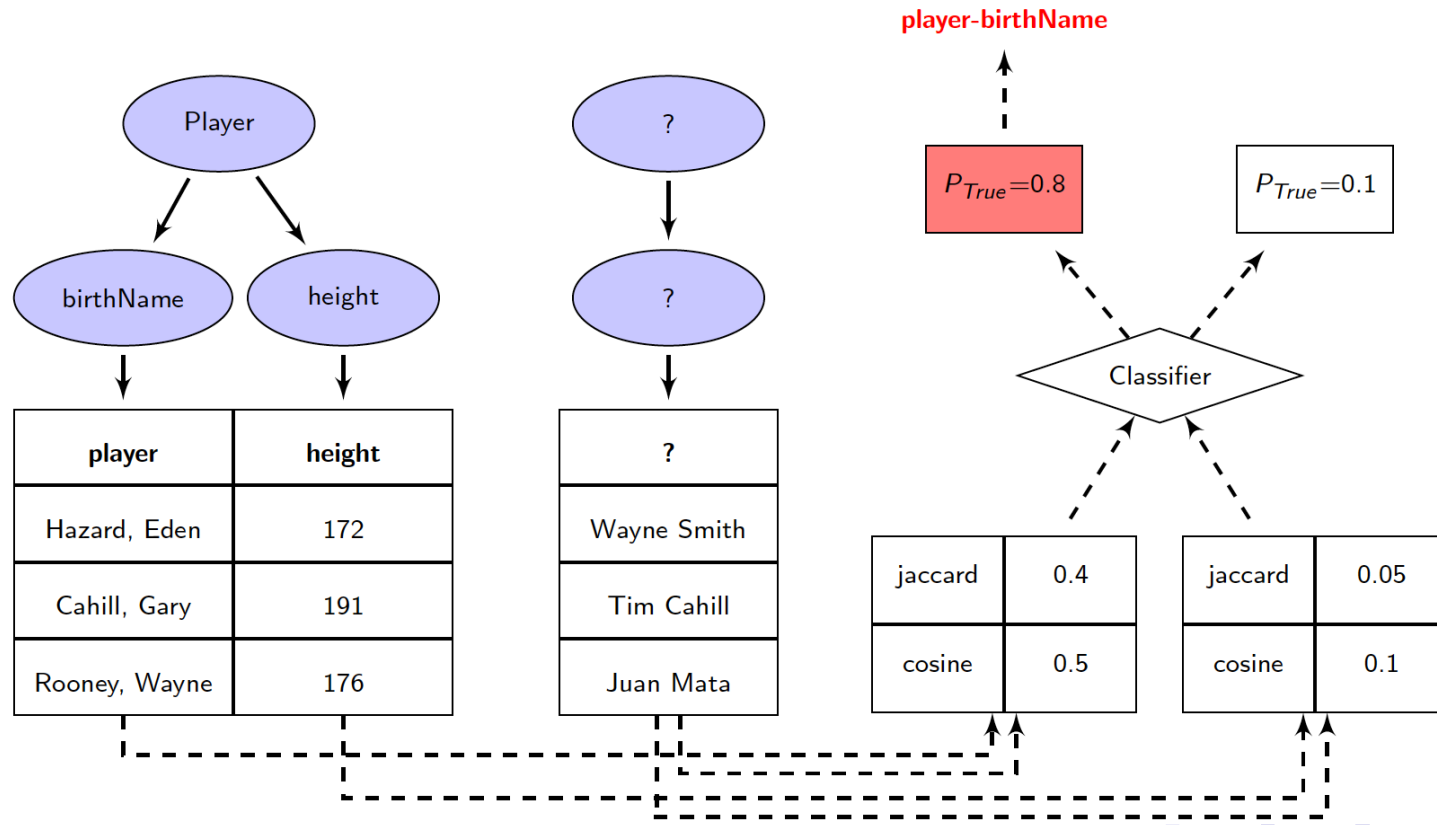
Histogram similarity



Training machine learning model



Predicting new attribute



Evaluation

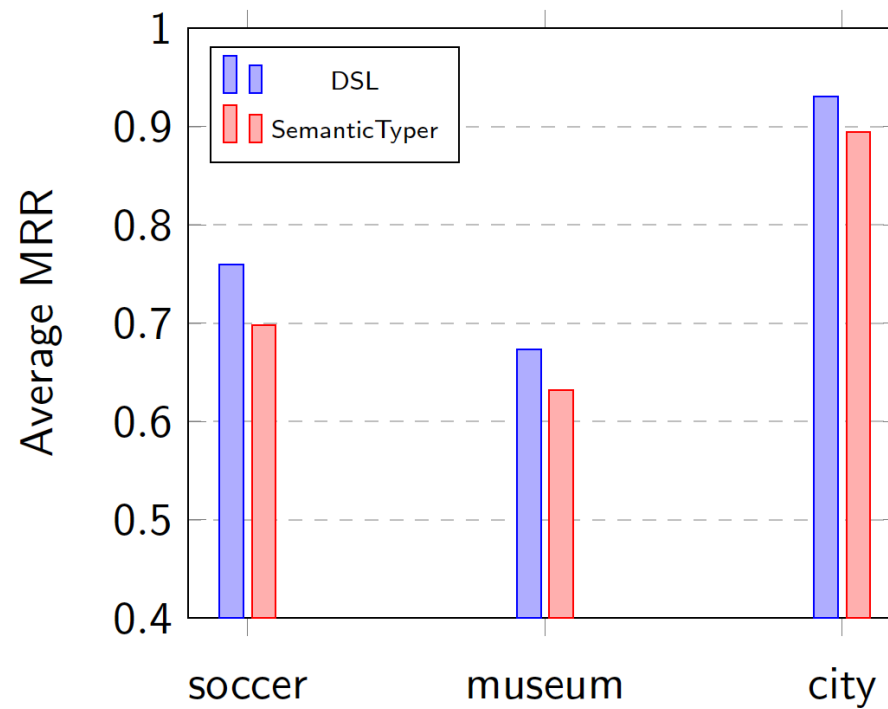
Data sets:

Domain data	# sources	# semantic types	# attributes
soccer	12	14	97
museum	29	20	217
city	10	52	520
weather	4	11	44
T2D Gold	1748	7983	?

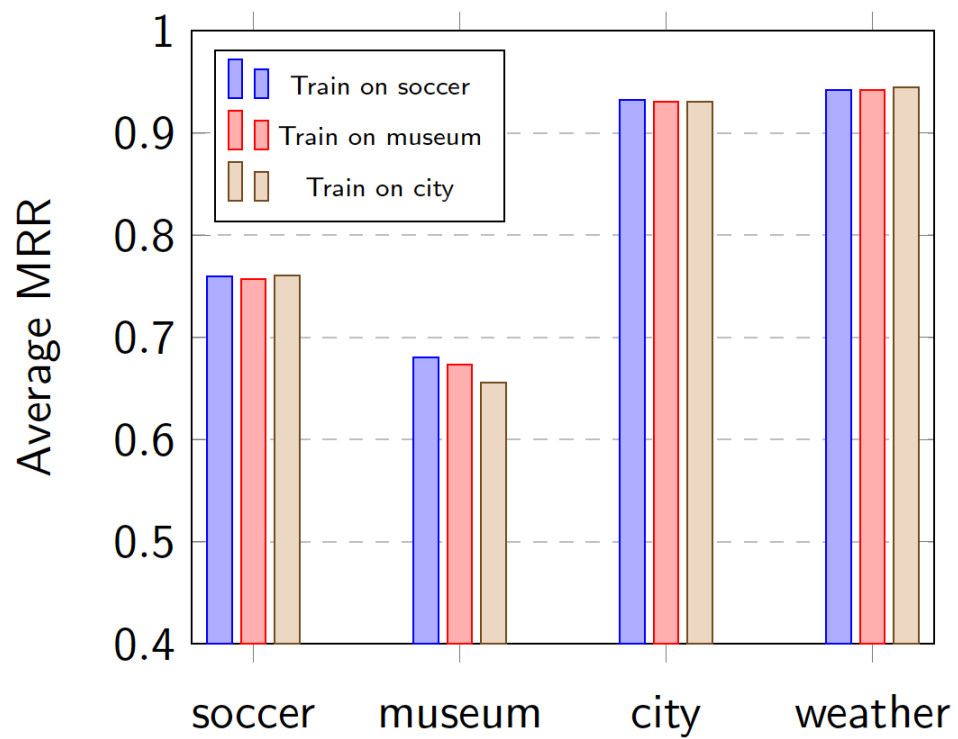
Measurements: Mean Reciprocal Rank (MRR)

Evaluating systems: DSL (our approach), SemanticTyper (Ramnandan et al, 2015), T2K (Ritze et al, 2015)

Evaluation



Evaluation



Conclusion

	Rule-based Approach	Learning-based Classification
Pros	<ul style="list-style-type: none">+ fast+ scalable+ easy to implement	<ul style="list-style-type: none">+ fast+ scalable+ easy to extend+ works in lots of domains
Cons	<ul style="list-style-type: none">+ requires heuristics+ difficult to extend+ may not work in every domain	<ul style="list-style-type: none">+ no use of relationship+ need to train machine learning on a general domain