



# Agenda

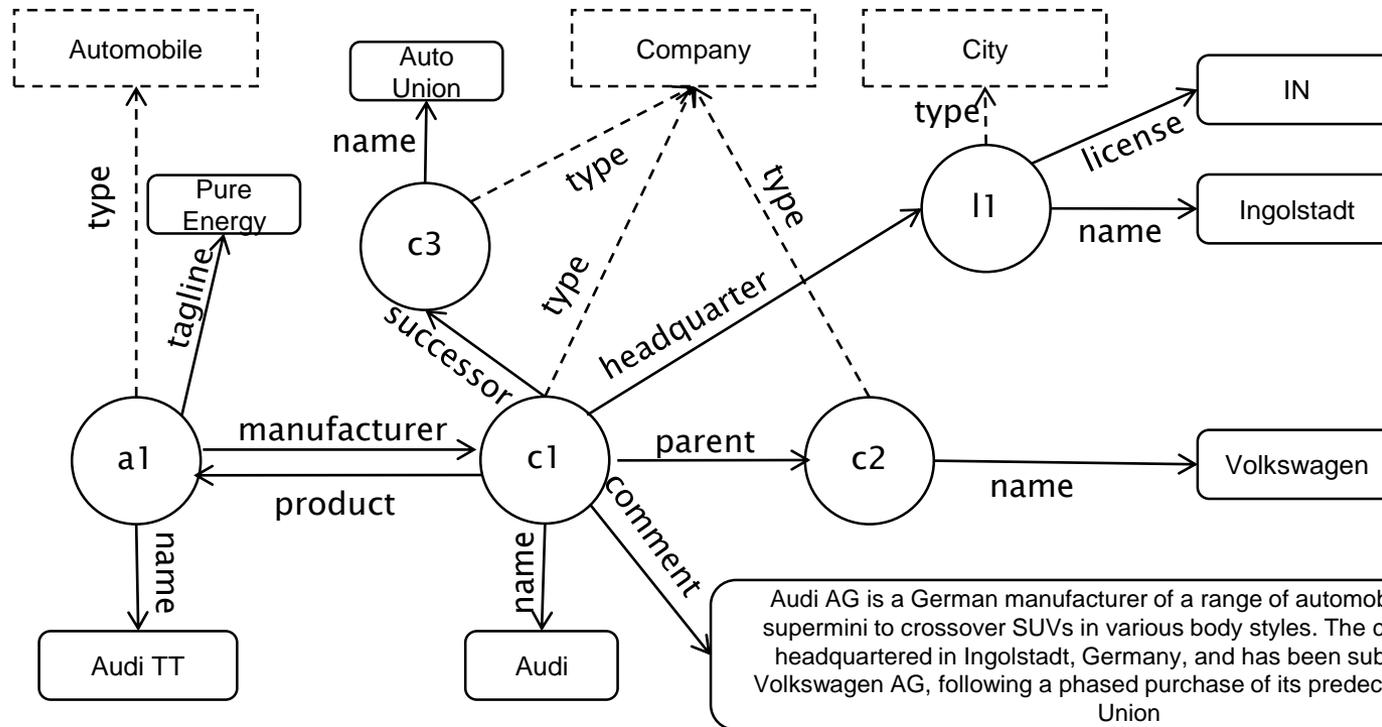
- Introduction
  - Motivation
  - RDF Data
  - State-of-the-Art
- Topical Relational Model (TRM)
  - Overview
  - Template-based Construction
  - TRM Output
- Evaluation
- Conclusion

# Agenda

- **Introduction**
  - **Motivation**
  - **RDF Data**
  - **State-of-the-Art**
- **Topical Relational Model (TRM)**
  - **Overview**
  - **Template-based Construction**
  - **TRM Output**
  - **Applications**
- **Evaluation**
- **Conclusion**

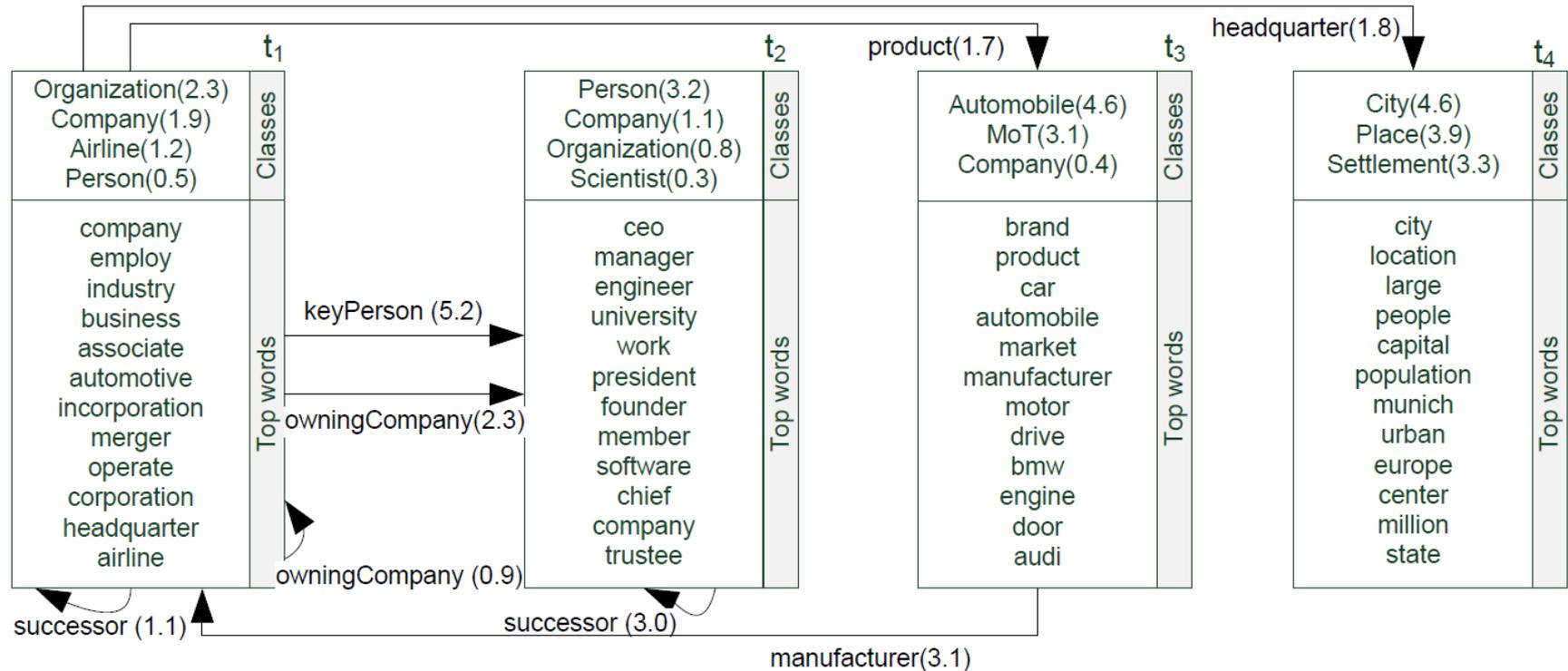
# Motivation

- Bridging the gap between structured and unstructured data:
  - Text-rich structured data is emerging on the Web in the form of RDF
  - Characteristics of RDF data
    - structured data in the form of entities, classes, and relations
    - unstructured data associated with attributes
- Topic Models
  - Model topics as distribution over words
  - Various models considering relational data
    - E.g., social networks, citations, Web links
  - Unique characteristics of RDF are not taken into account
- Statistical Relational Learning
  - Learning probabilistic models from relational data
  - Handling text is a problem



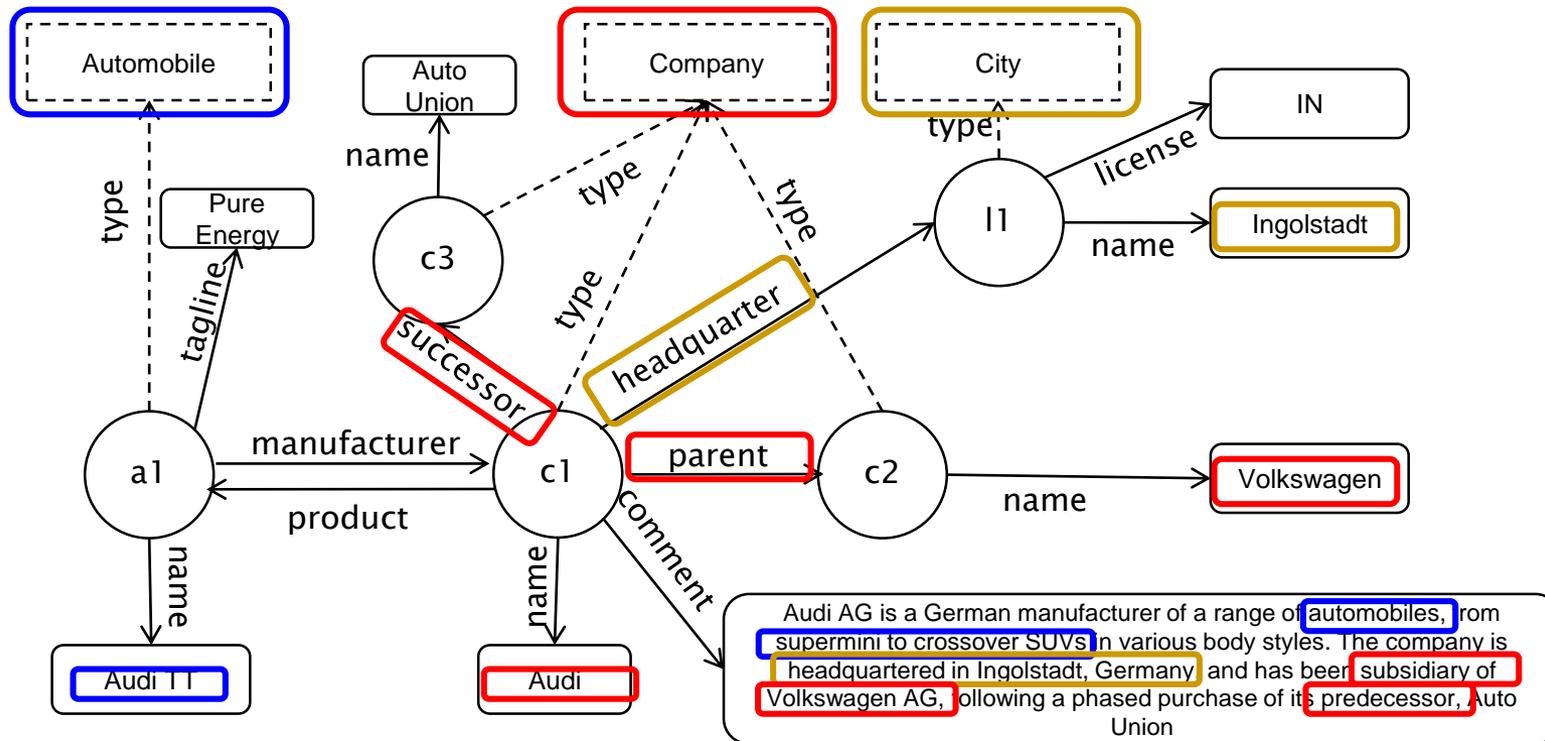
- RDF graph: Vertices denote **resources** and text **values**, connected via **relations** and **attributes**
- **Semi-structured:** Containing both structured and unstructured data

# Topics to be learned

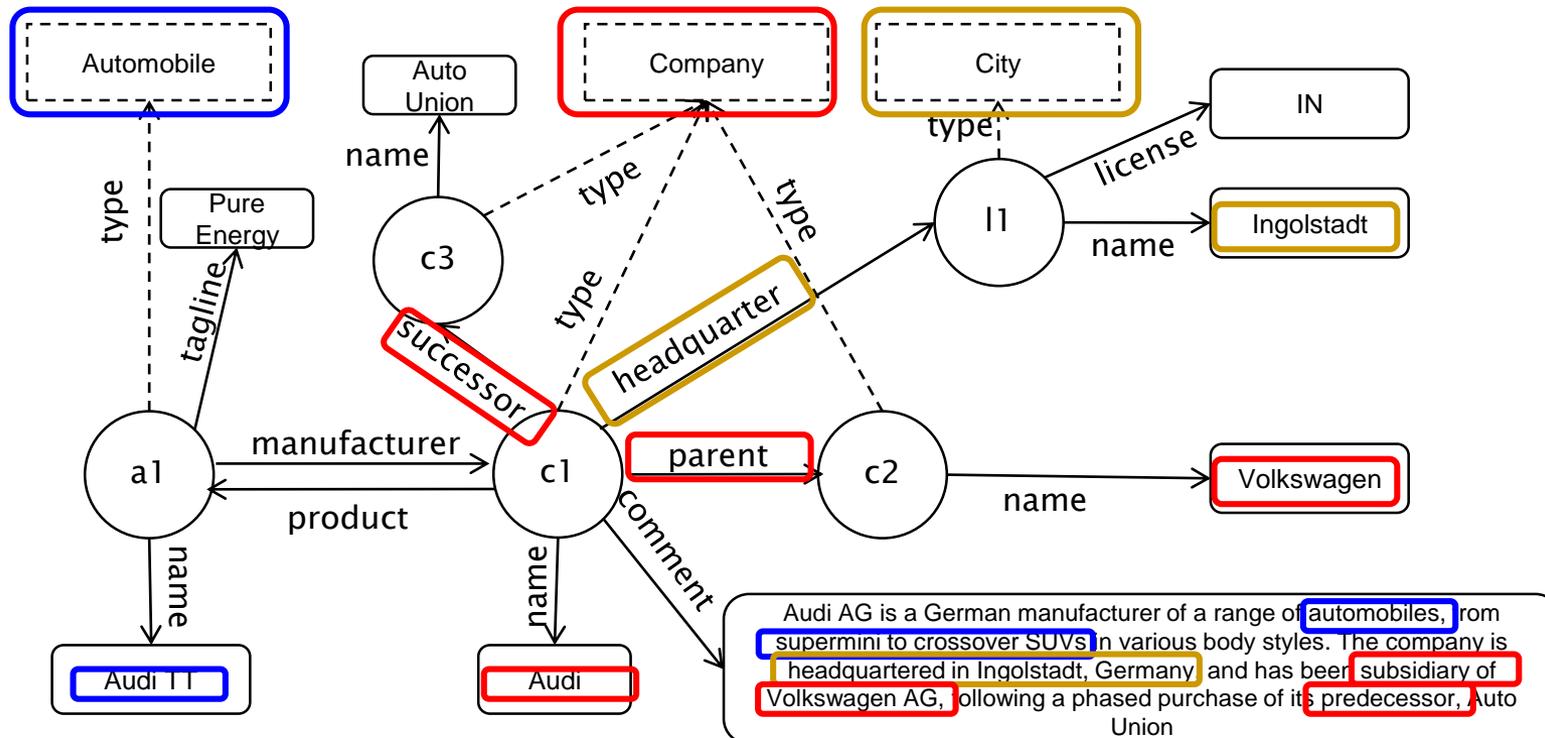


- Topics have to cover words, classes and relations

# RDF Data



- RDF graph: Vertices denote **resources** and text **values**, connected via **relations** and **attributes**
- **Semi-structured:** Containing both structured and unstructured data
- **Highly heterogeneous:** Many classes and relations each of which has varying effects on different topics
- **Sparseness:** Different structure elements have sparse correlations to the topics in the text



- RDF graph: Vertices denote **resources** and text **values**, connected via **relations** and **attributes**
- **Semi-structured**: Containing both structured and unstructured data
- **Highly heterogeneous**: Many classes and relations each of which has varying effects on different topics
- **Sparseness**: Different structure elements have sparse correlations to the topics in the text
- **Challenges**:
  - How to learn topics which better reflect RDF structures
  - How to learn correlations between topics and structure taking heterogeneity and sparseness into account

# State-of-the-art

- Supervised Topic Models
  - sLDA: word co-occurrence and multi-label documents
  - Type-LDA: focus on relation extraction in NLP
- Topic Models for homogeneous networks
  - ReLDA, Pairwise Link-LDA: words and links in hypertext
  - Nubbi, Author-Topic models, Multi-relational TMs
- Topic Models for heterogeneous networks
  - Topic Model with Biased Propagation (TMBP)
    - Limits in capturing complex correlations between structure and topics
    - Sparsity not well addressed
- Statistical Relational Learning
  - MLN: FOL rules for text
  - C<sup>3</sup>: SVM based on textual features (e.g. Jaccard)

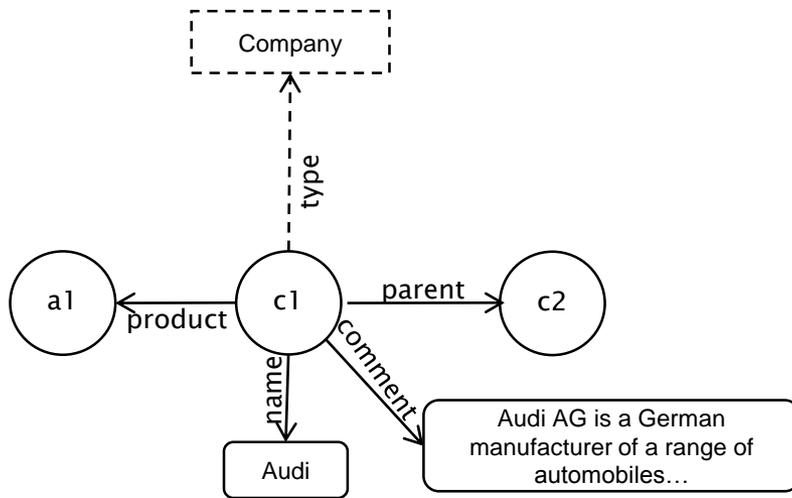
# Agenda

- Introduction
  - Motivation
  - RDF Data
  - State-of-the-Art
- **Topical Relational Model (TRM)**
  - **Overview**
  - **Template-based Construction**
  - **TRM Output**
- Evaluation
- Conclusion

# Topical Relational Models

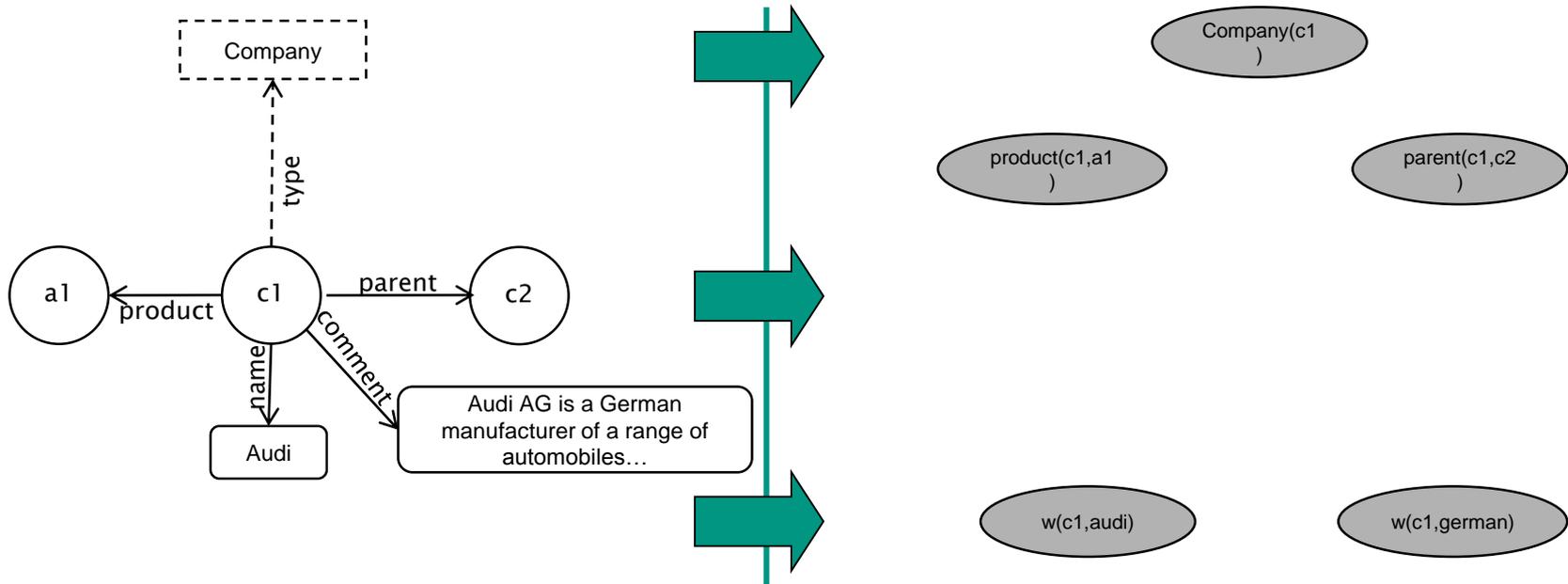
- A model to learn a set of topics and their correlations to the classes, relations and texts from RDF
- Topics as a low-dimensional representation of text in an RDF graph
  - Biased towards the structure
    - e.g. words such as *'employ'* or *'merger'* correlate with classes such as *Company*
- Correlations between the class (relation) and the topics are captured as a vector (matrix) of weights
  - e.g. the topics which assign high probabilities to words *'employ'* or *'merger'* have high weights in the vector of the class *Company*
- Learned using a variational Bayesian EM
  - Model is intractable for exact inference

# Template-based Construction



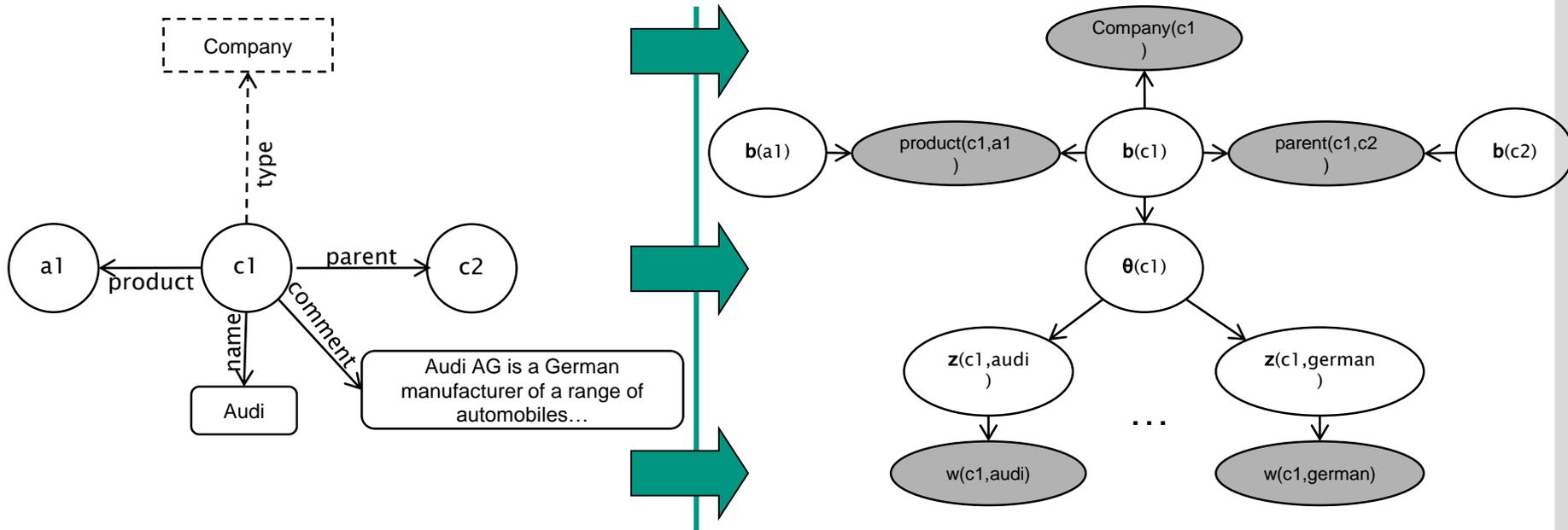
- TRM uses an RDF graph as a template to construct a ground BN

# Template-based Construction



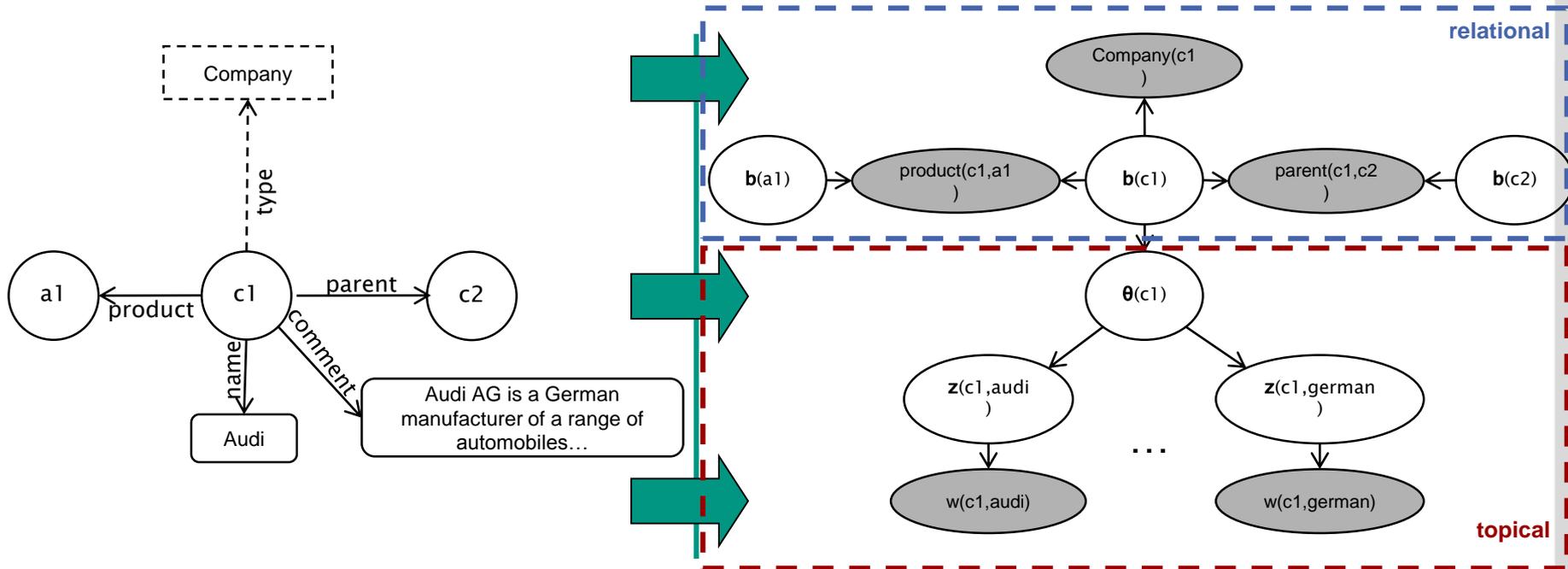
- TRM uses an RDF graph as a template to construct a ground BN
- **Three types of templates for observed variables:** Classes, relations, and words
- Random variables are instantiated for each resource from the template

# Template-based Construction



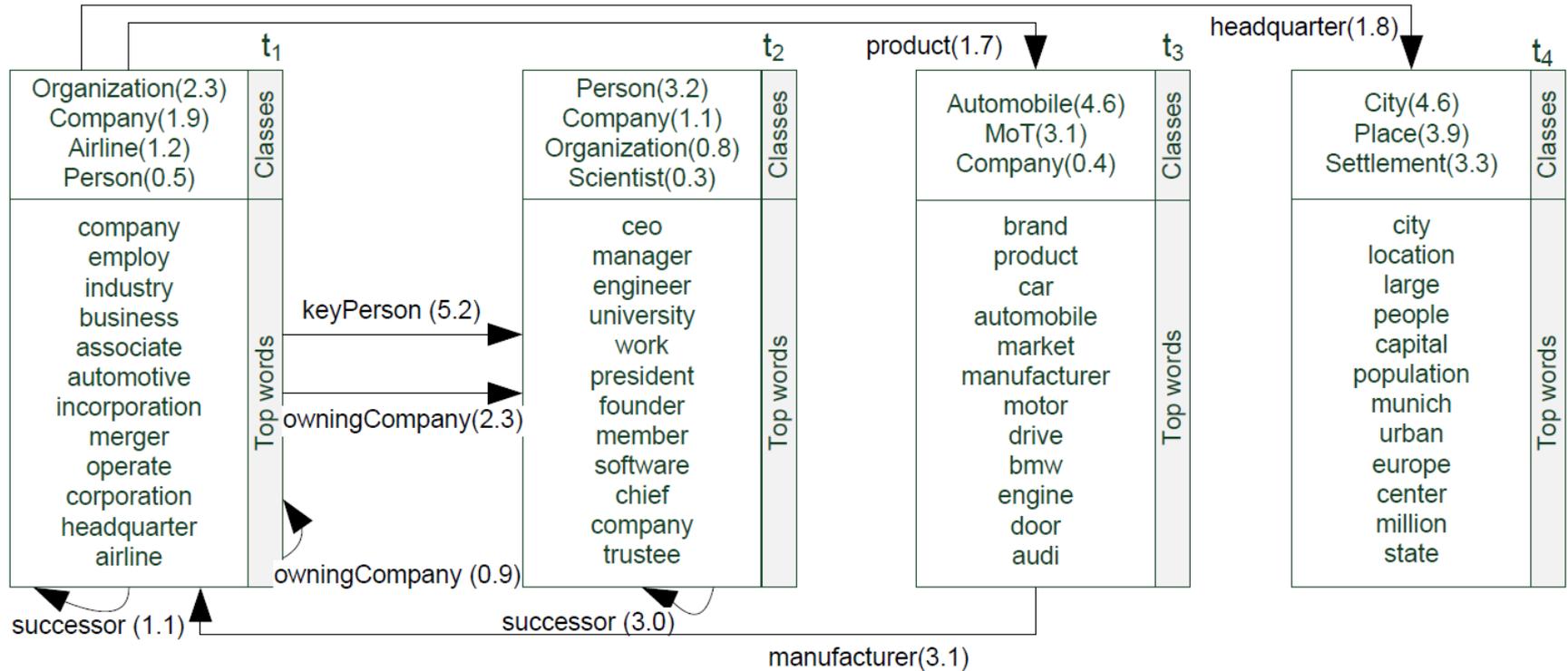
- TRM uses an RDF graph as a template to construct a ground BN
- **Three types of templates for observed variables:** Classes, relations, and words
- Random variables are instantiated for each resource from the template
- Hidden topic-related variables:
  - Topic indicator vector  $\mathbf{b}$
  - Topic proportions  $\theta$
  - Topic-word assignments  $\mathbf{z}$

# Template-based Construction

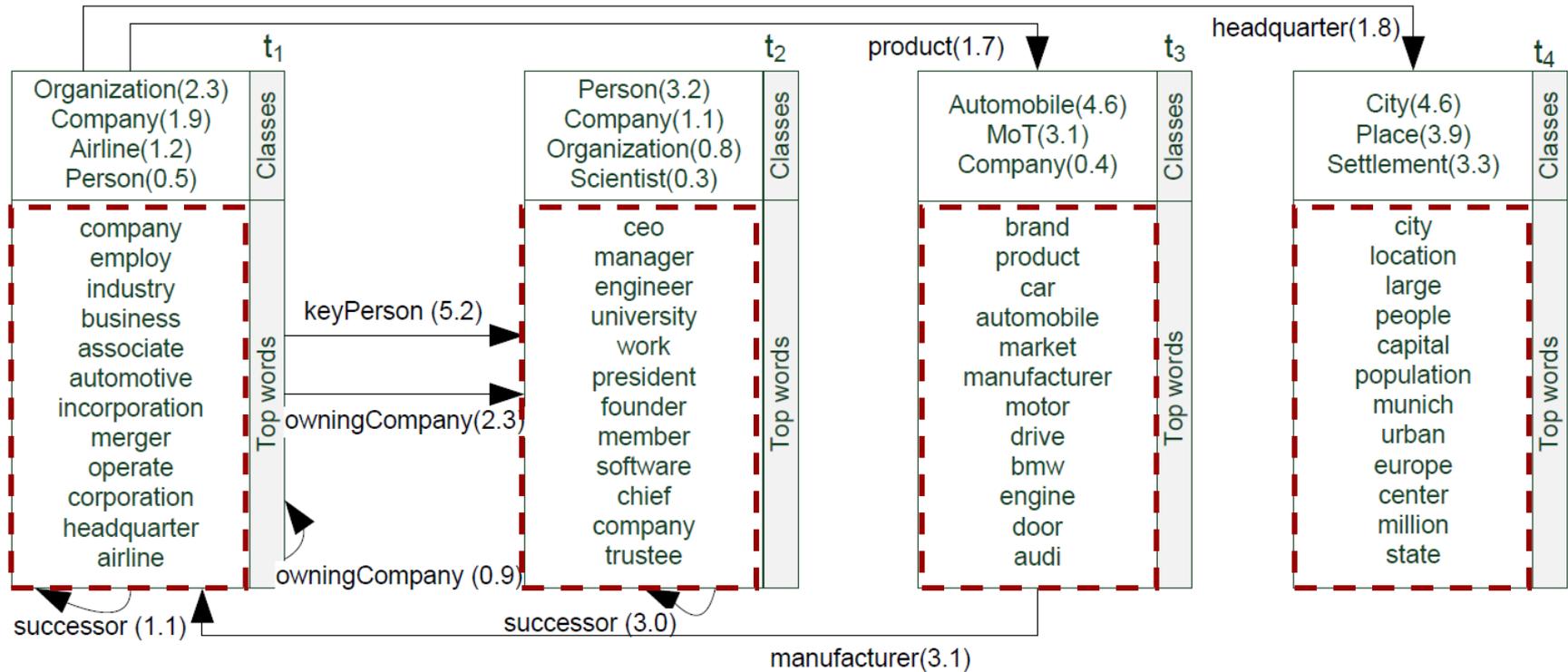


- TRM uses an RDF graph as a template to construct a ground BN
- **Three types of templates for observed variables:** Classes, relations, and words
- Random variables are instantiated for each resource from the template
- Hidden topic-related variables:
  - Topic indicator vector  $\mathbf{b}$
  - Topic proportions  $\theta$
  - Topic-word assignments  $\mathbf{z}$
- The ground BN can be considered as two parts (topical, relational) managed by  $\mathbf{b}$  vector

# TRM Output



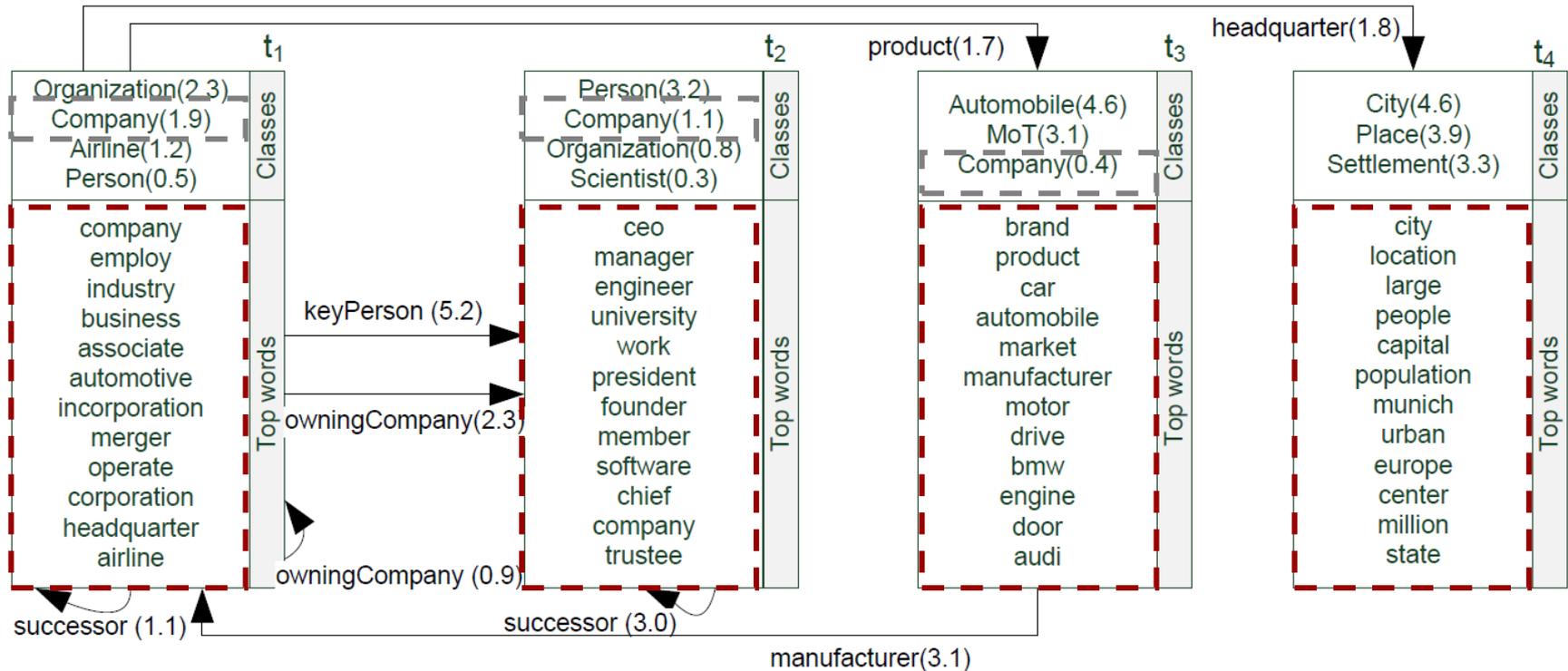
# TRM Output



-TRM discovers the topics biased towards the structure

-Topics are multinomial distributions over the words in the vocabulary

# TRM Output

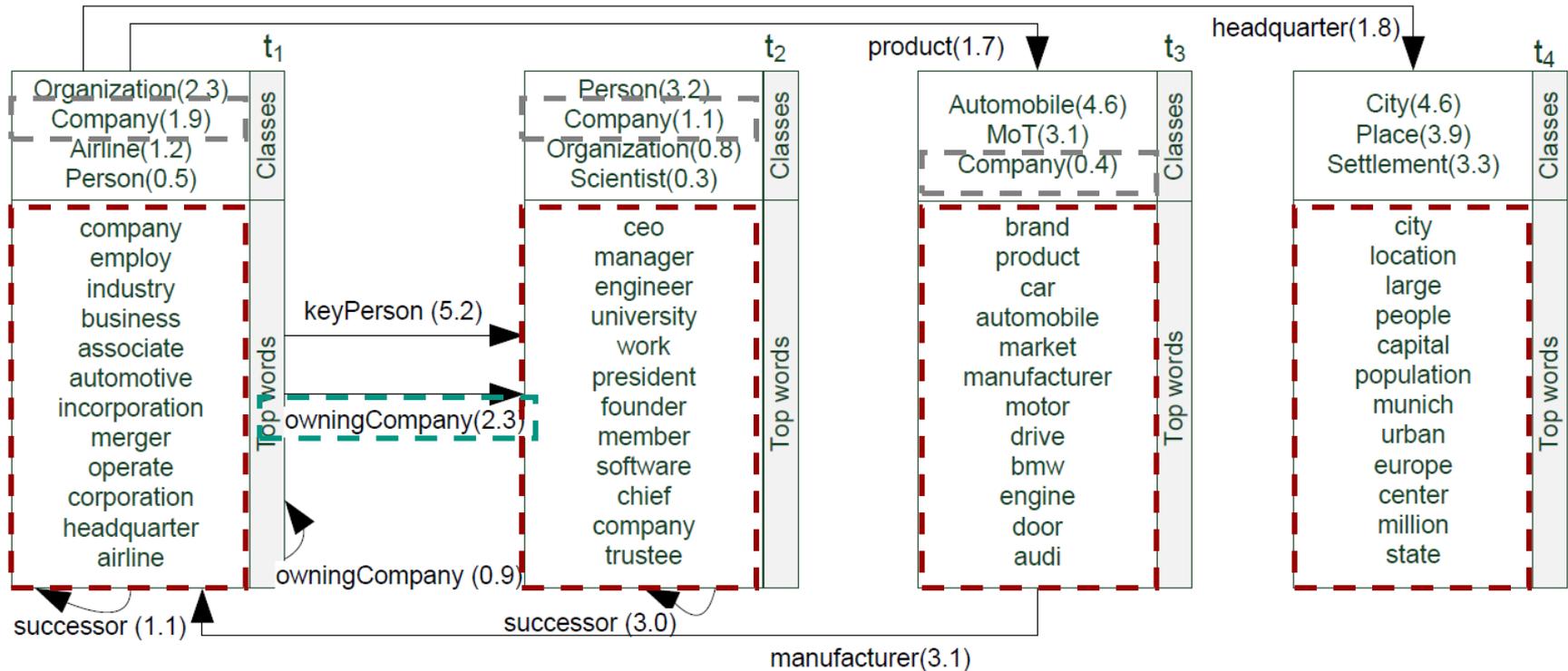


-TRM discovers the topics biased towards the structure

-Topics are multinomial distributions over the words in the vocabulary

-TRM captures the correlations between those topics and the classes via the weights in its class-topic parameter

# TRM Output



-TRM discovers the topics biased towards the structure

-Topics are multinomial distributions over the words in the vocabulary

-TRM captures the correlations between those topics and the classes via the weights in its class-topic parameter

-The weight of observing a relation between any two topics is encoded in its relation-topic parameter (by taking the direction of the relation into account)

# Agenda

- Introduction
  - Motivation
  - RDF Data
  - State-of-the-Art
- Topical Relational Model (TRM)
  - Overview
  - Template-based Construction
  - TRM Output
- **Evaluation**
- Conclusion

# Evaluation

## ■ Datasets:

- **DBLP:** Authors, conferences and their relations to papers
  - 28,569 papers, 28,702 authors and 20 conferences
  - The abstract and title of the papers are treated as textual data
- **DBpedia:** Movie domain
  - 20,094 entities described by 112 distinct classes and 49 different types of relations
  - All attribute values are treated as textual data

## ■ Tasks:

- **Link Prediction**
- **Object Clustering**

# Evaluation

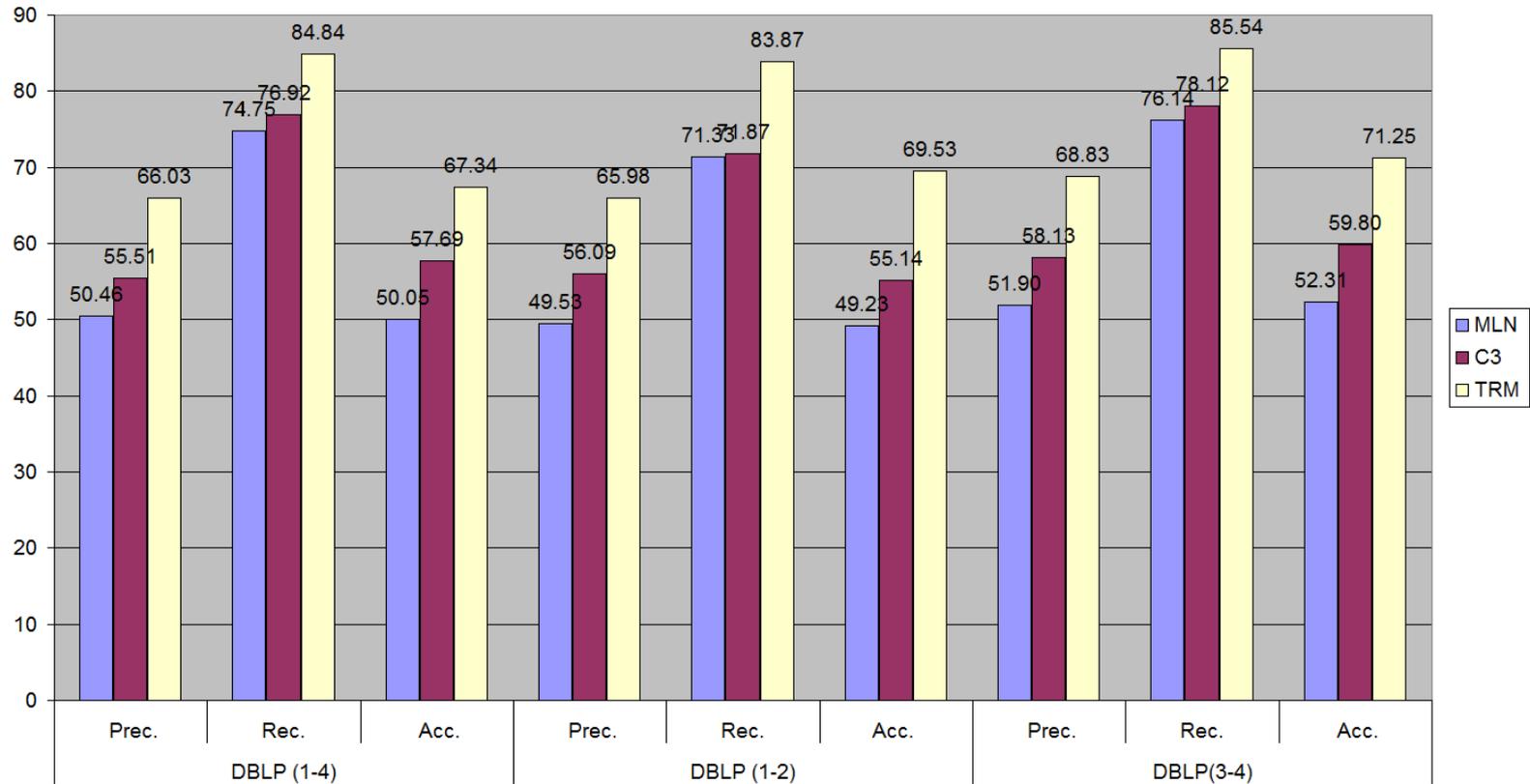
## ■ Link Prediction (LP)

- Predicting **author** relations between papers and authors in DBLP
- *Predicting **starring*** relations between movies and actors in DBpedia

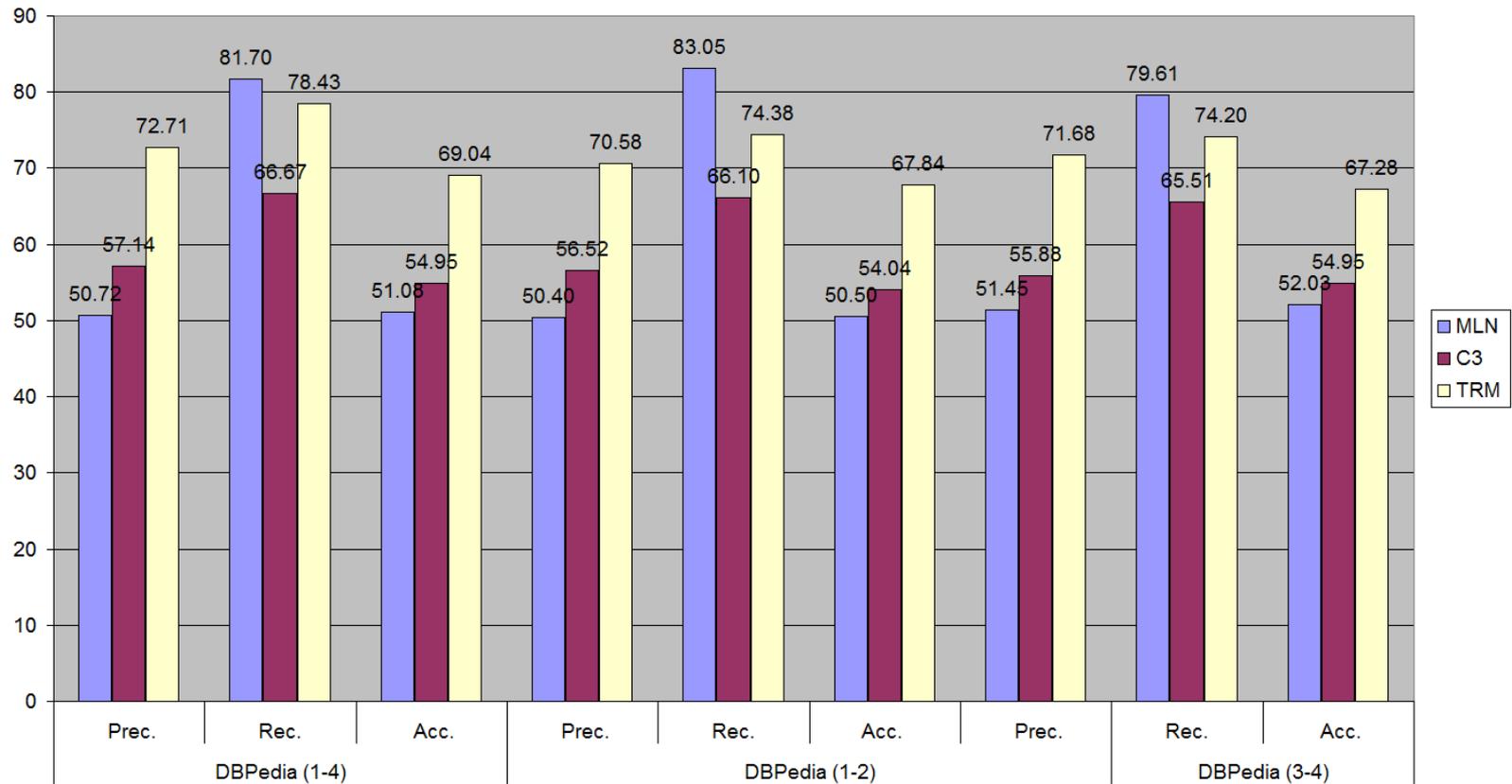
## ■ Baselines:

- MLN [Richardson and Domingos, 2006]
- C<sup>3</sup> [Namata et al., 2011]

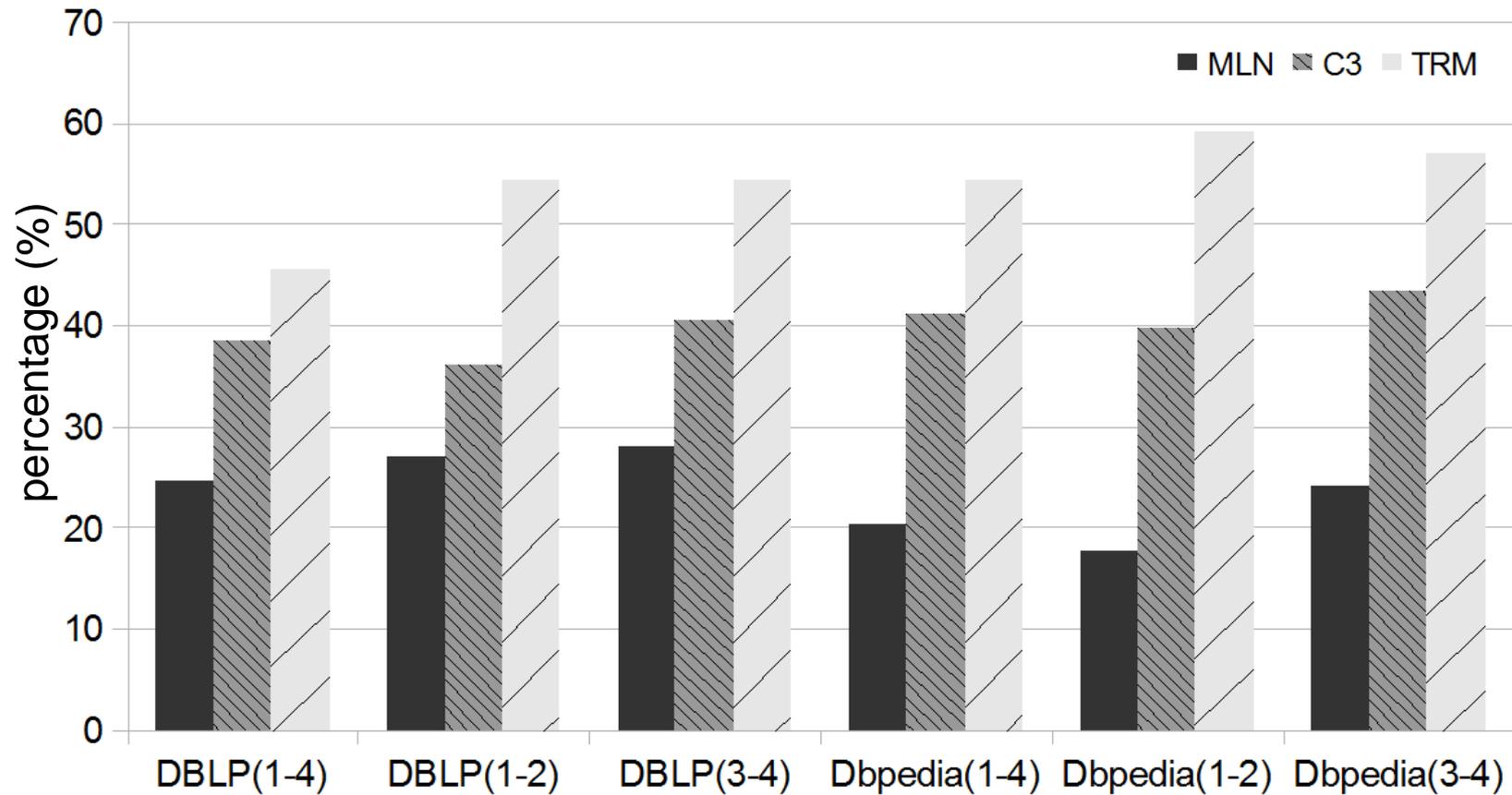
## Precision, recall and accuracy results for link prediction on DBLP



Precision, recall and accuracy results for link prediction on DBPedia



# Evaluation



True negative rate for DBLP and DBpedia

# Agenda

- Introduction
  - Motivation
  - RDF Data
  - State-of-the-Art
- Topical Relational Model (TRM)
  - Overview
  - Template-based Construction
  - TRM Output
- Evaluation
- **Conclusion**

# Conclusion

- TRM provides an effective model for text-rich RDF data
- TRM captures dependencies between words in textual and structured data
- Compared to existing TM approaches, TRM is more effective in exploiting structure information
  
- Application Opportunities
  - Link prediction, object clustering
  - Estimating the result size of hybrid SPARQL queries
  - Keyword search on RDF data: Estimating the most probable connections between entities given the query
  
- Future work:
  - Extension of TRM to richer generative models, such as time-varying and hierarchical topic models

- **Selectivity Estimation on RDF [Wagner et al. 2014]**
  - Estimating the result size of a hybrid SPARQL query, i.e. query with structural and textual predicates
  - TRM provides a uniform data synopsis to summarize the text-rich structured data
    - estimation is based on using TRM output of topic distributions, class-topic and relation-topic parameters
- **Keyword Search on RDF Data**
  - State-of-the-art: Keywords are mapped to RDF elements and connections are discovered among those elements over the schema [Tran et al., ICDE 2009]
  - TRM provides a probabilistic schema capturing the most probable connections given the query

# Evaluation

## ■ Object Clustering

### ■ **DBPedia:**

Clustering movie entities into different types

- e.g. American, British, German movies

### ■ **DBLP:**

Use the six labels in DBLP representing various computer science fields as clusters

- cluster paper, author and venue entities

## ■ Baselines:

- LDA, TMBP-RW, TMBP-Reg

# Evaluation

Accuracy and Normalized Mutual Information (NMI) results for object clustering

