

---

# Relation Extraction

## Part 2

Feiyu Xu



---

# Machine Learning for Relation Extraction



# Motivations of ML

---

- ⊙ Porting to new domains or applications is expensive
- ⊙ Current technology requires IE experts
  - ⊙ Expertise difficult to find on the market
  - ⊙ SME cannot afford IE experts
- ⊙ Machine learning approaches
  - ⊙ Domain portability is relatively straightforward
  - ⊙ System expertise is not required for customization
  - ⊙ “Data driven” rule acquisition ensures full coverage of examples

## Problems

---

- ⊙ Training data may not exist, and may be very expensive to acquire
- ⊙ Large volume of training data may be required
- ⊙ Changes to specifications may require reannotation of large quantities of training data
- ⊙ Understanding and control of a domain adaptive system is not always easy for non-experts

# Parameters of IE Real-World Tasks

---

- ⊙ Document structure
  - ⊙ Free text
  - ⊙ Semi-structured
  - ⊙ Structured
- ⊙ Linguistic annotation
  - ⊙ Shallow NLP
  - ⊙ Deep NLP
- ⊙ Complexity and specificity of relation
  - ⊙ Unary
  - ⊙ N-ary
- ⊙ Depth of extraction
  - ⊙ Recognition
  - ⊙ Classification
  - ⊙ Semantic role labelling
- ⊙ Degree of automation
  - ⊙ Semi-automatic
  - ⊙ Supervised
  - ⊙ Semi-Supervised
  - ⊙ Minimally-Supervised
  - ⊙ Distant Supervision
  - ⊙ Unsupervised
- ⊙ Human interaction/contribution
- ⊙ Data properties
  - ⊙ Domain relevance
  - ⊙ Redundancy
  - ⊙ Connectivity
- ⊙ Evaluation/validation
  - ⊙ With/without gold standard
  - ⊙ Performance: recall & precision
  - ⊙ Interaction among parameters

---

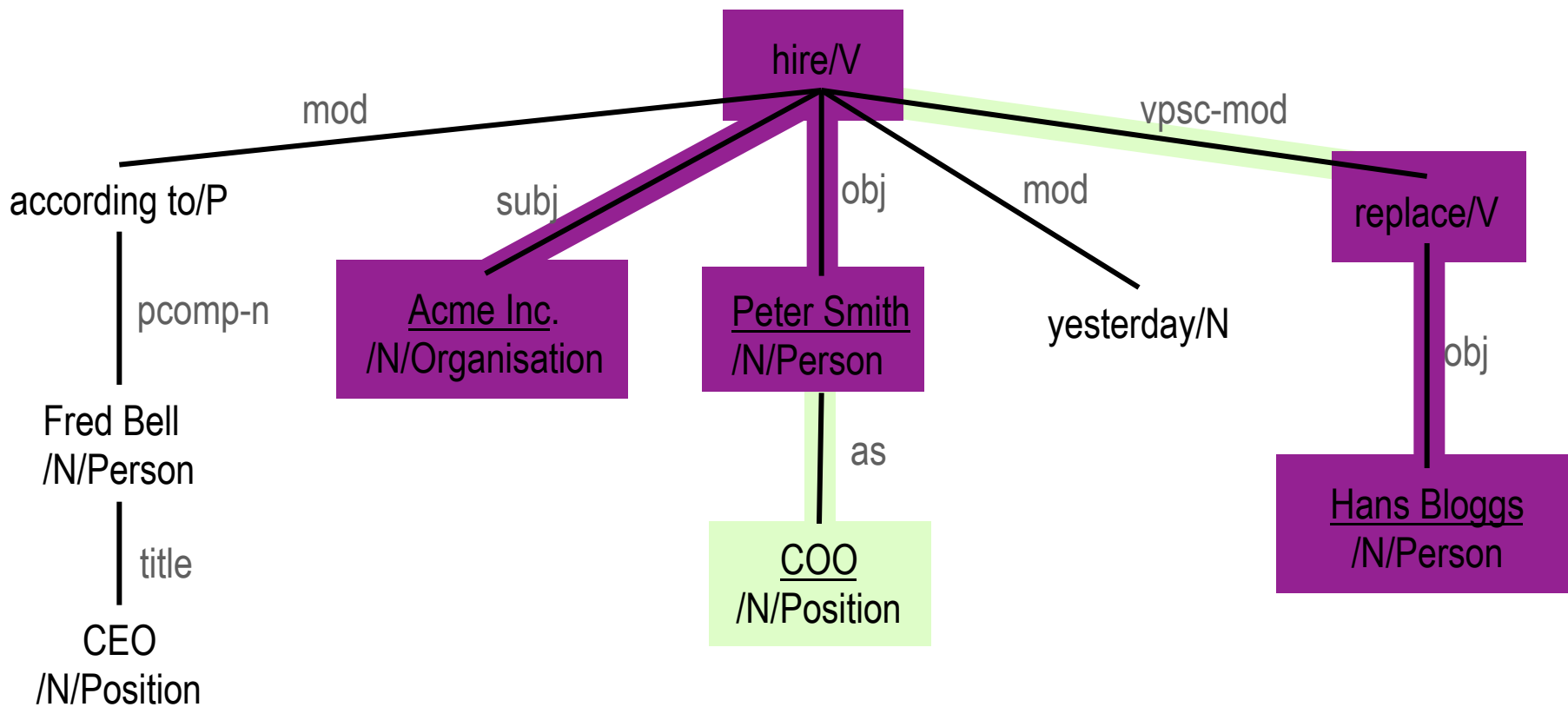
# State of the Art



# Binary Relation Only Approaches

---

- ⊙ Extraction of binary relations only such as
  - ⊙ *author-book*
  - ⊙ *company-location*
- ⊙ Do not employ the existing syntactic and semantic structures among  $n > 2$  arguments and rely on a later component to merge binary relations into complex relations
- ⊙ Approaches
  - ⊙ Ravichandran and Hovy, 2002
  - ⊙ Pantel et al., 2004
  - ⊙ Pasca et al., 2006a; Pasca et al., 2006b



# Surface-oriented Rule Representation

---

- ◎ Shallow linguistic analyses
- ◎ These formalisms are robust and efficient but only handle binary relations.
- ◎ Work best for relations whose arguments usually co-occur in close proximity within a sentence and whose mentions exhibit limited linguistic variation
- ◎ Approaches
  - ◎ Pasca et al., 2006a; Pasca et al., 2006b;
  - ◎ Kozareva et al., 2008; Hovy et al., 2009; Kozareva and Hovy, 2010

## Minimally Supervised (Bootstrapping)

---

- ⊙ Based on iterative learning with limited initial knowledge
  - ⊙ Start a small number of initial examples of relation instances (or patterns)
  - ⊙ Label the free texts during iterations (e.g., Agichtein and Gravano, 2000; Yangarber et al., 2000; Ravichandran and Hovy, 2002; Stevenson and Greenwood, 2005).
- ⊙ Often suffer from semantic drift or the propagation of errors occurring during iterations
- ⊙ Performance depends on data properties

## Distant Supervision

---

- ◎ A massive seed-based, one step version of bootstrapping
- ◎ Rely on a large amount of trustworthy facts
- ◎ Their performance does not hinge on corpus data properties such as redundancy
- ◎ Approaches
  - ◎ (Mintz et al., 2009)
  - ◎ Others: (Wu and Weld, 2007); (Wu et al., 2008); (Weld et al., 2008); (Hoffmann et al., 2010); (Xu et al., 2011); (Nguyen and Moschitti, 2011)

- ⊙ They do not target given relations.
- ⊙ They are very useful for applications continuously faced with new relation or event types, e.g., online social media monitoring
- ⊙ However, the results of these systems cannot be directly taken for filling knowledge databases, because the semantics of the new relations including the roles of the entities remains unknown.
- ⊙ Example: TextRunner (Banko et al., 2007; Yates et al., 2007)

## Reality in IE Projects

---

- ◎ Our IE users are often not domain experts
- ◎ IE experts have to develop methods and strategies for
  - ◎ Prospecting a domain
  - ◎ Proposing relevant relations
  - ◎ Finding relevant and suitable data

---

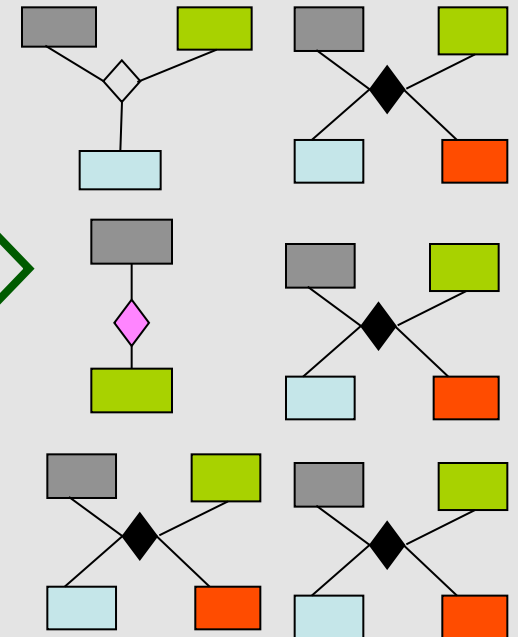
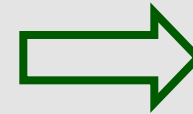
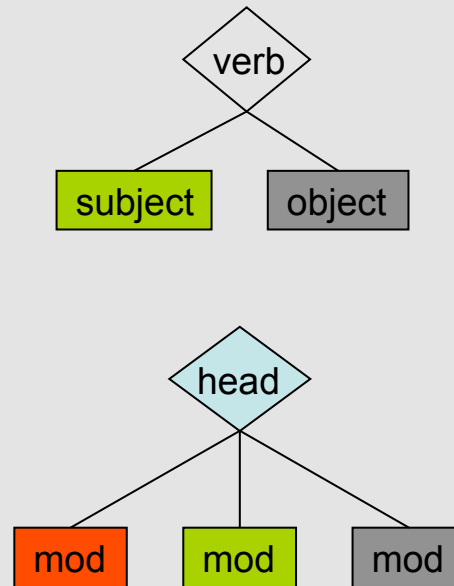
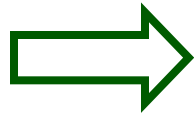
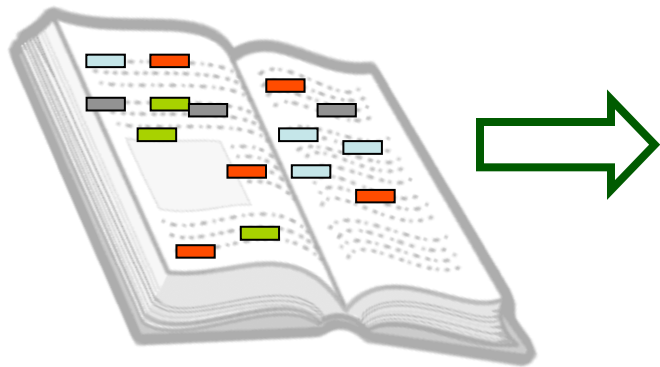
# DARE: Domain Adaptive Relation Extraction

<http://dare.dfki.de>



# Research Goal

Development of a general framework for automatically learning mappings between linguistic analyses and target semantic relations, with minimal human intervention.



# Challenges

- Easy adaptation to new relation types with varied complexity
- Automatic learning without annotated corpus
- Exhaustive discovery of relevant linguistic patterns
- Integration of semantic role information into linguistic patterns

# Example

A relation extraction task in the domain *management succession* (MUC-6)

< person\_in, person\_out, position, organisation >

- *person\_in*: the person who obtained the position
- *person\_out*: the person who left the position
- *position*: the job position that the two persons were involved in
- *organisation*: the organisation where the position was located

According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

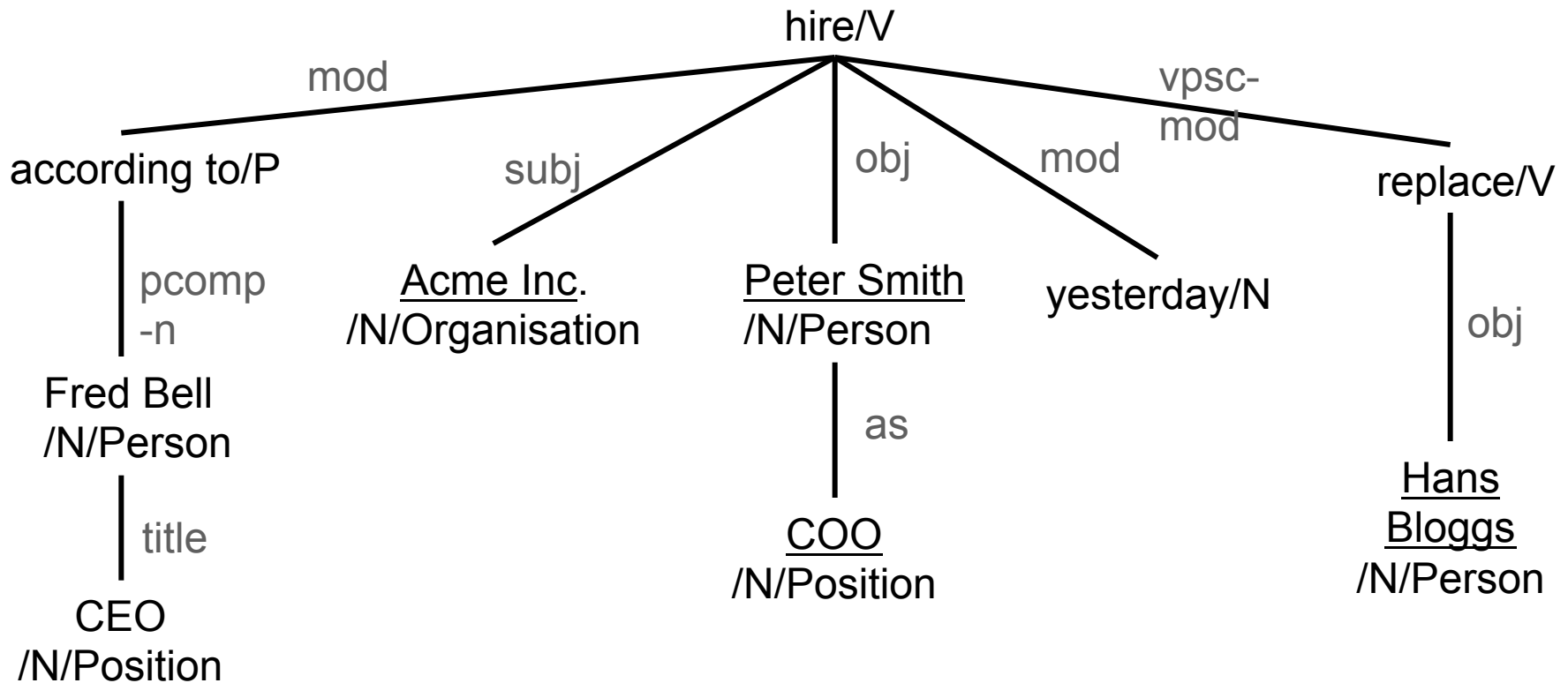
<person\_in, person\_out, position, organisation>

According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<person\_in, person\_out, position, organisation>

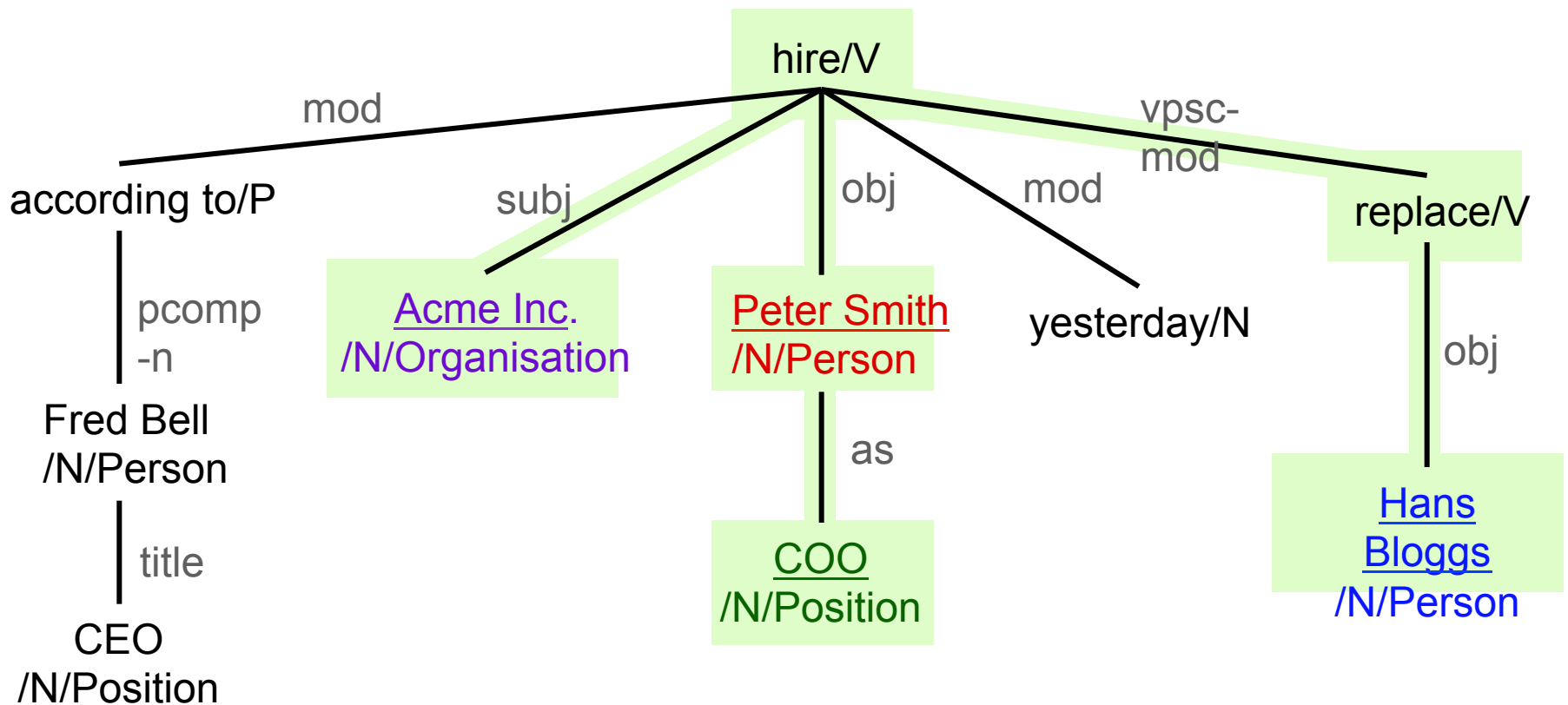
According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<person\_in, person\_out, position, organisation>

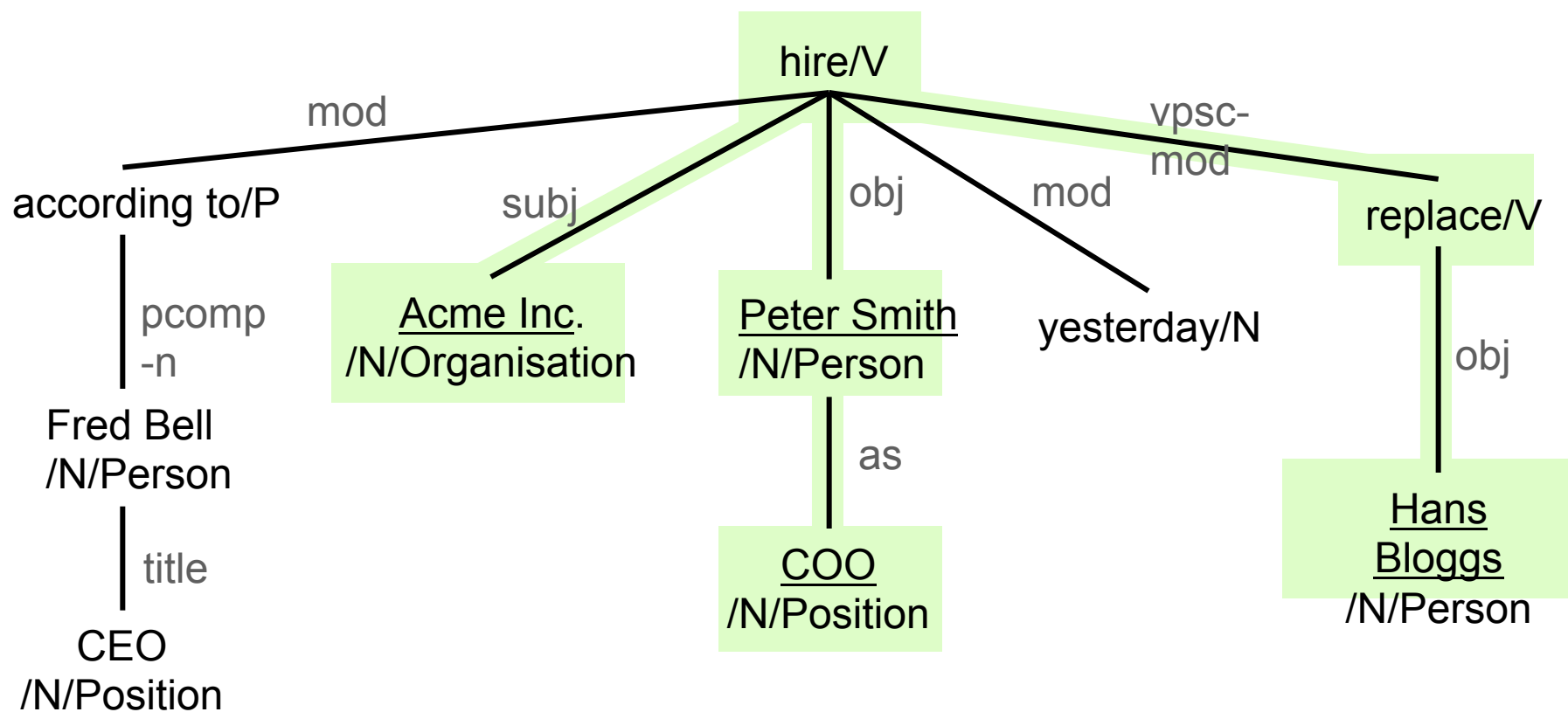


According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

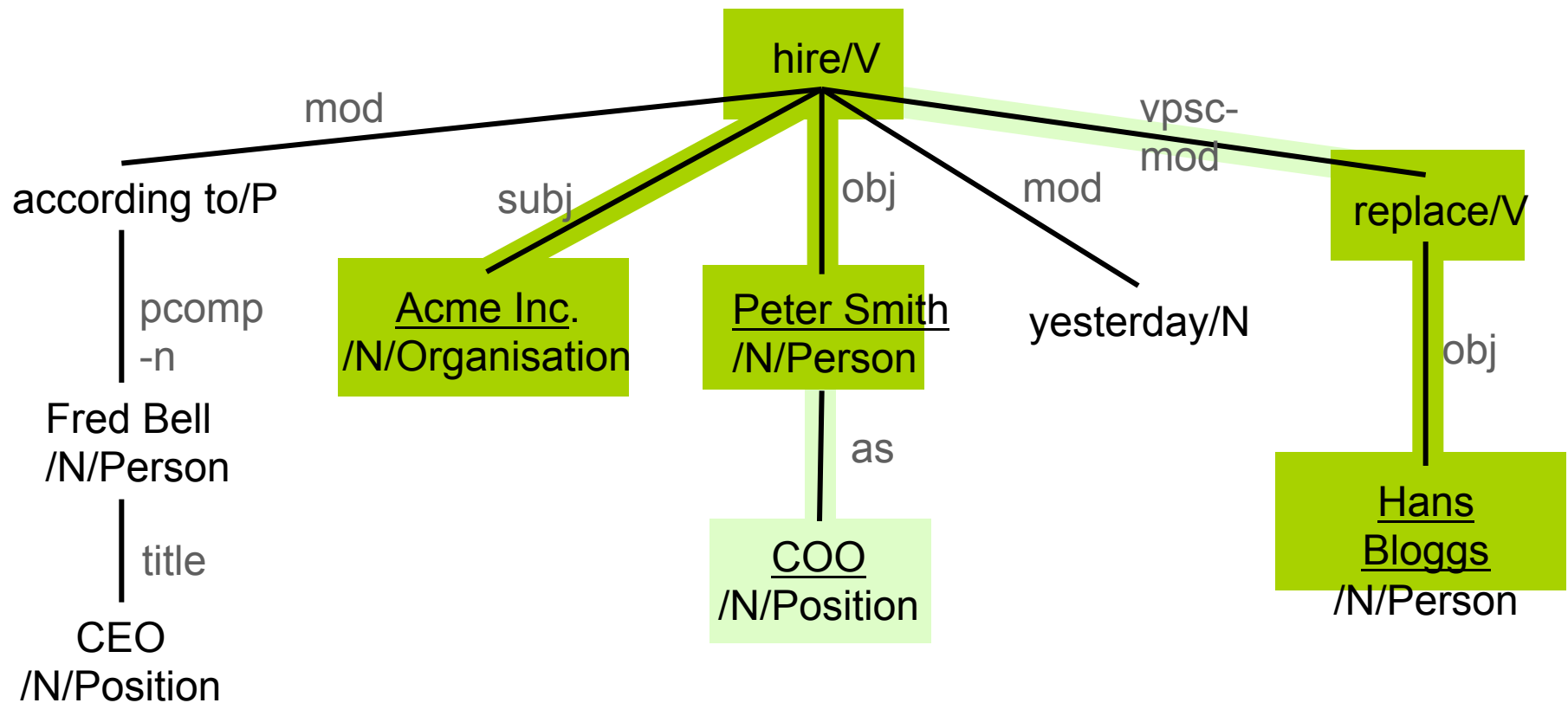
<person\_in, person\_out, position, organisation>



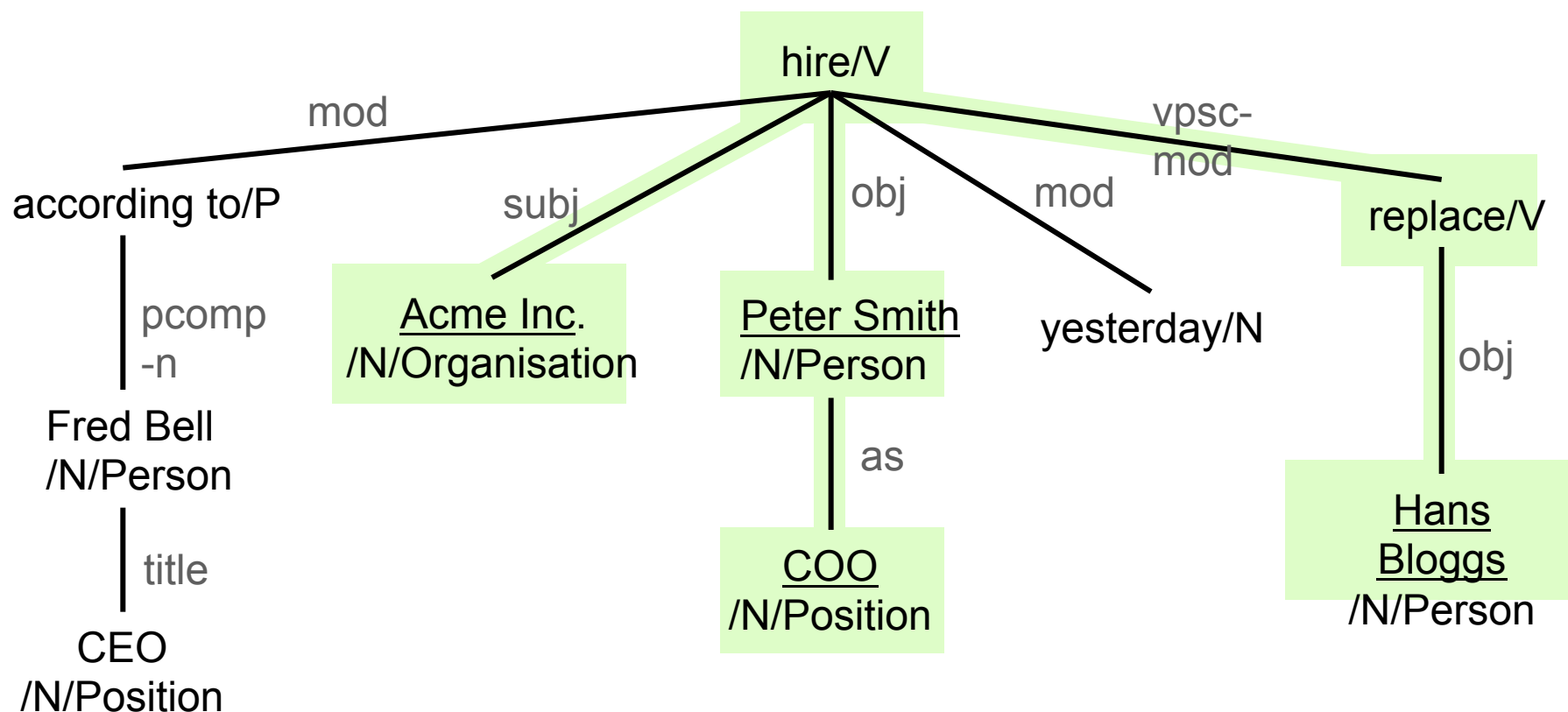
# Ideal Target Pattern



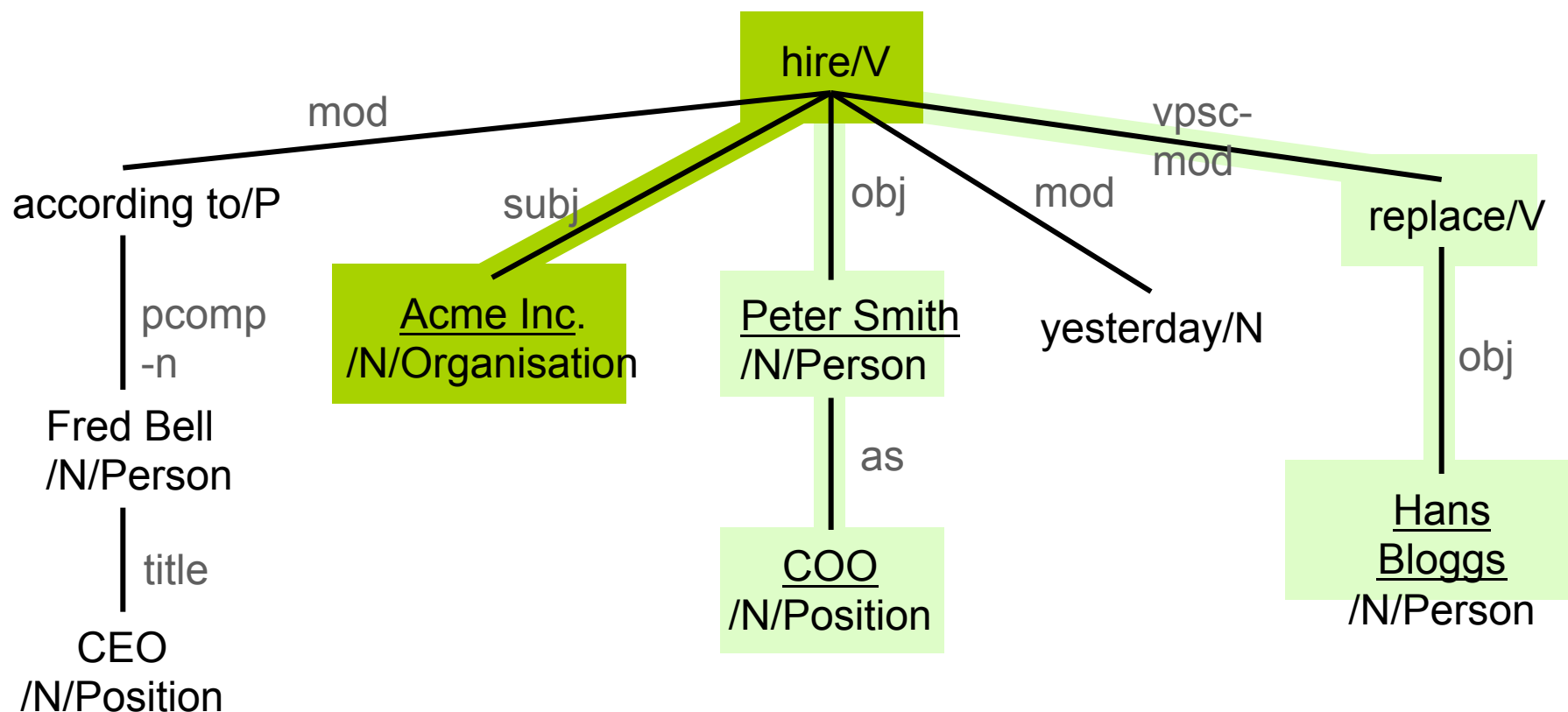
- Verb centered
- Direct relations between subject-verb-object
- Complex NP can not be extracted, e.g., the person and position relation
- The linguistic relations among patterns are not considered, e.g., hire and replace



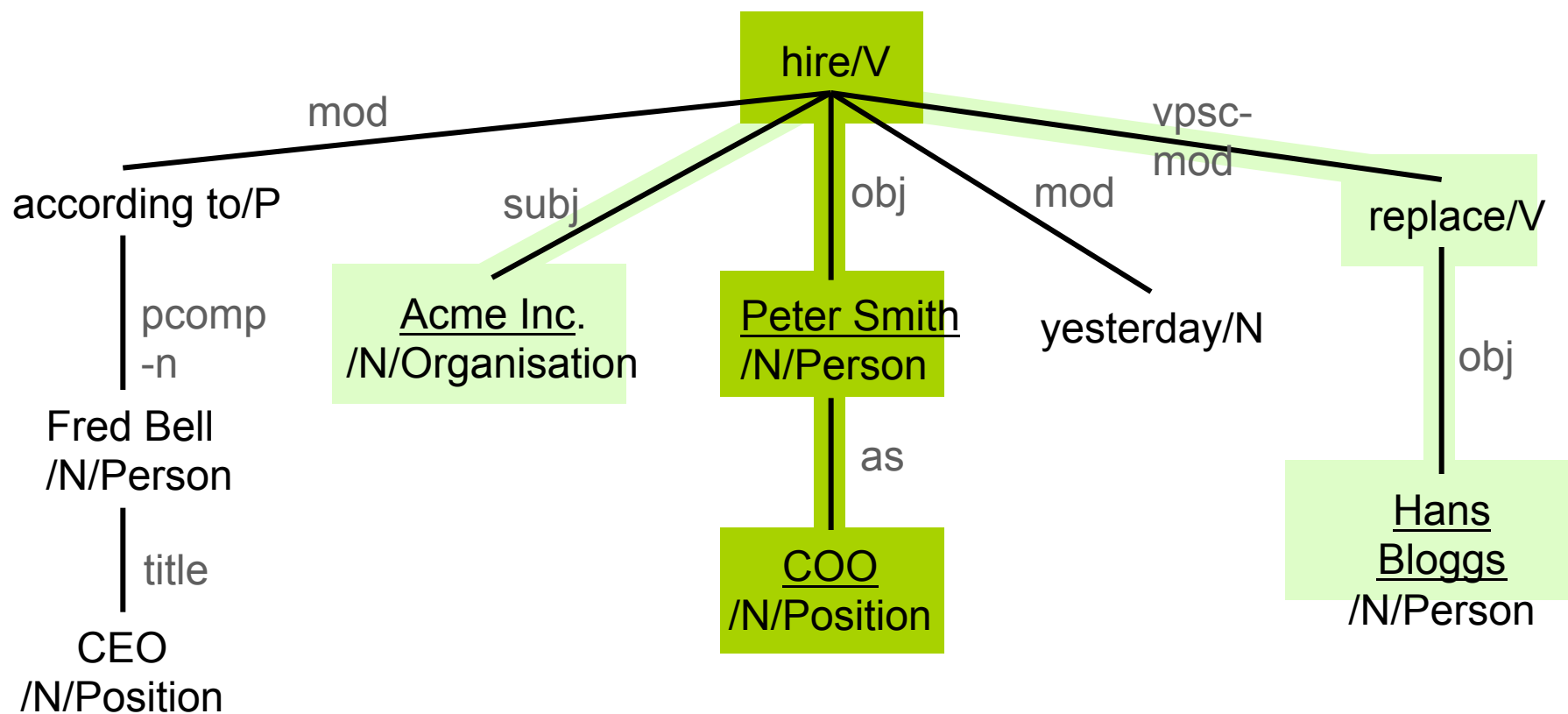
- Verb centered
- A single syntactic path dominated by a verb containing at least one relevant named entity concept



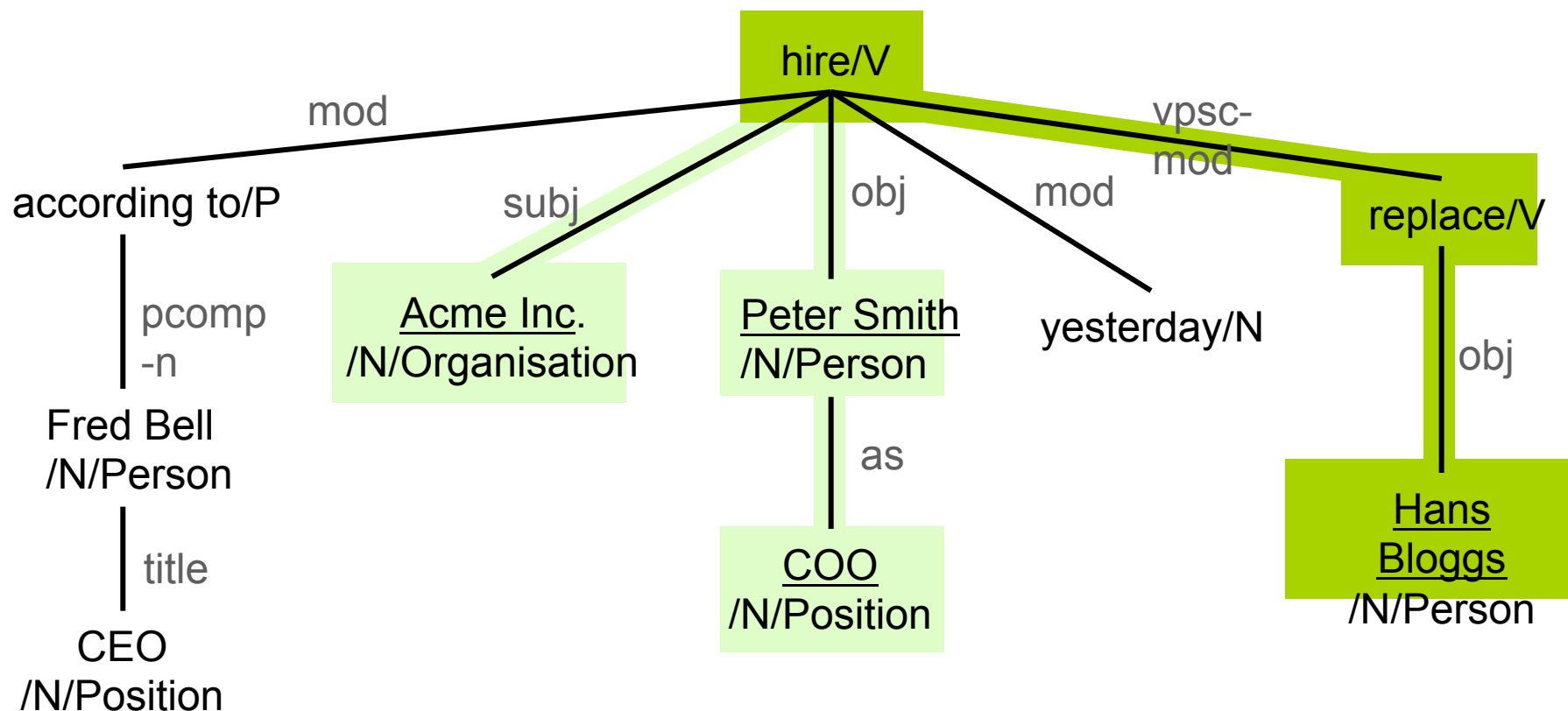
- Verb centered
- A single syntactic path dominated by a verb containing at least one relevant named entity concept



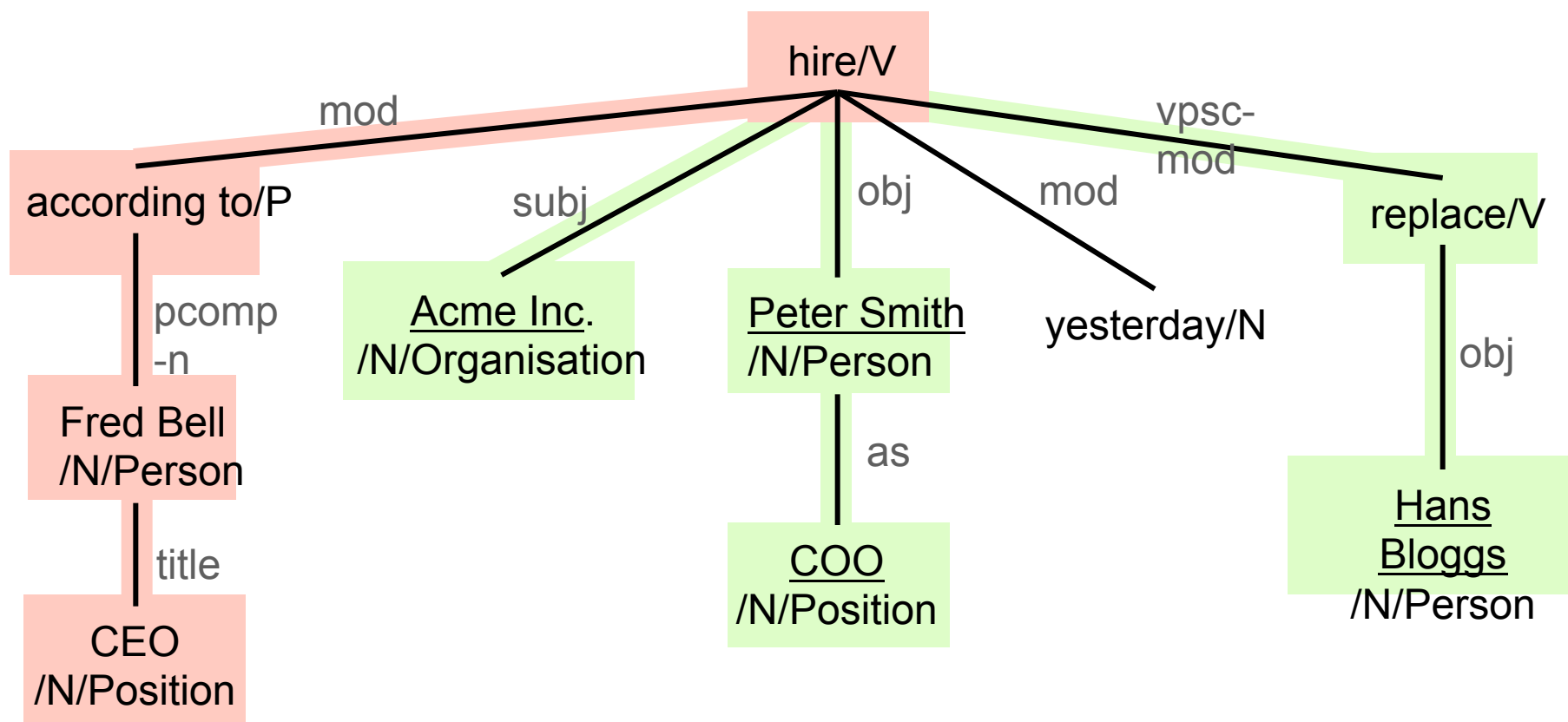
- Verb centered
- A single syntactic path dominated by a verb containing at least one relevant named entity concept



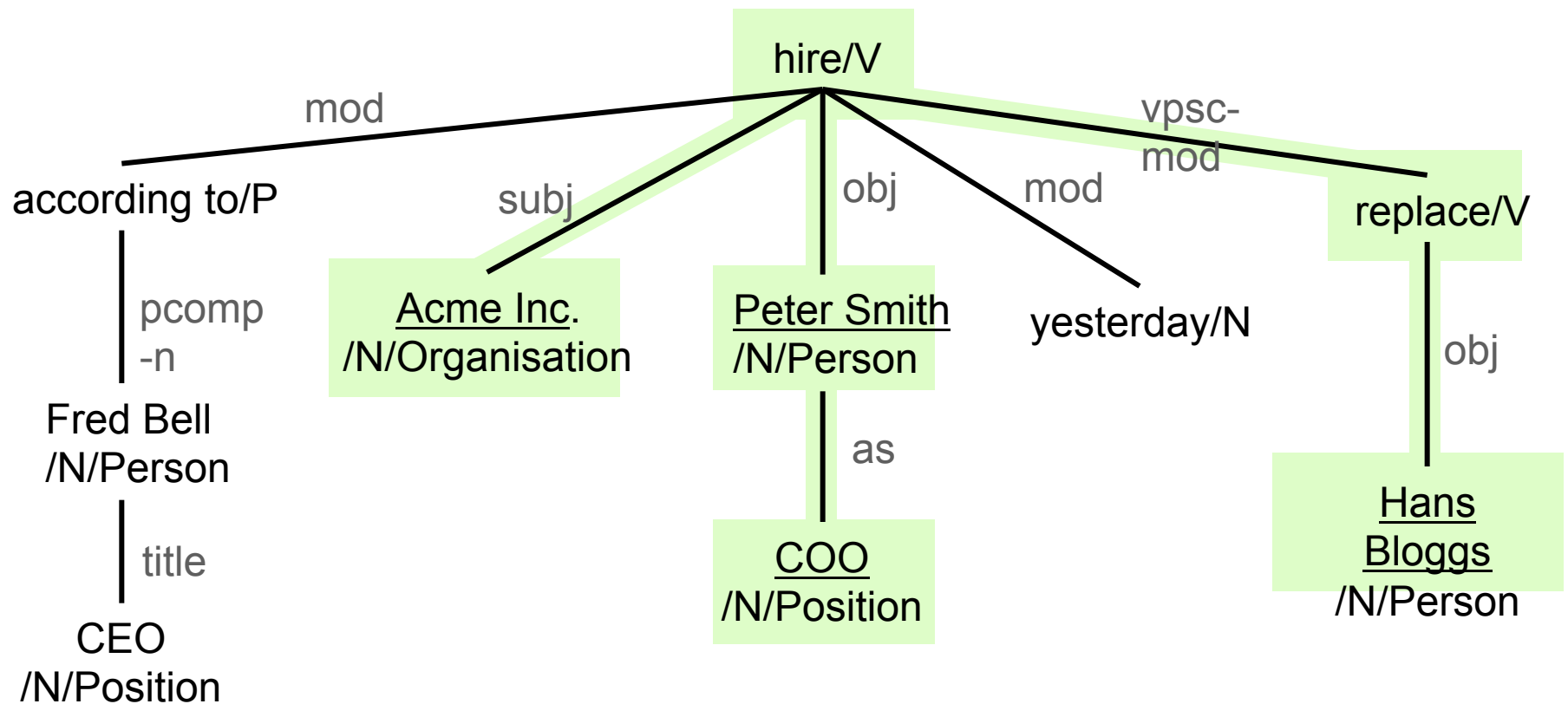
- Verb centered
- A single syntactic path dominated by a verb containing at least one relevant named entity concept



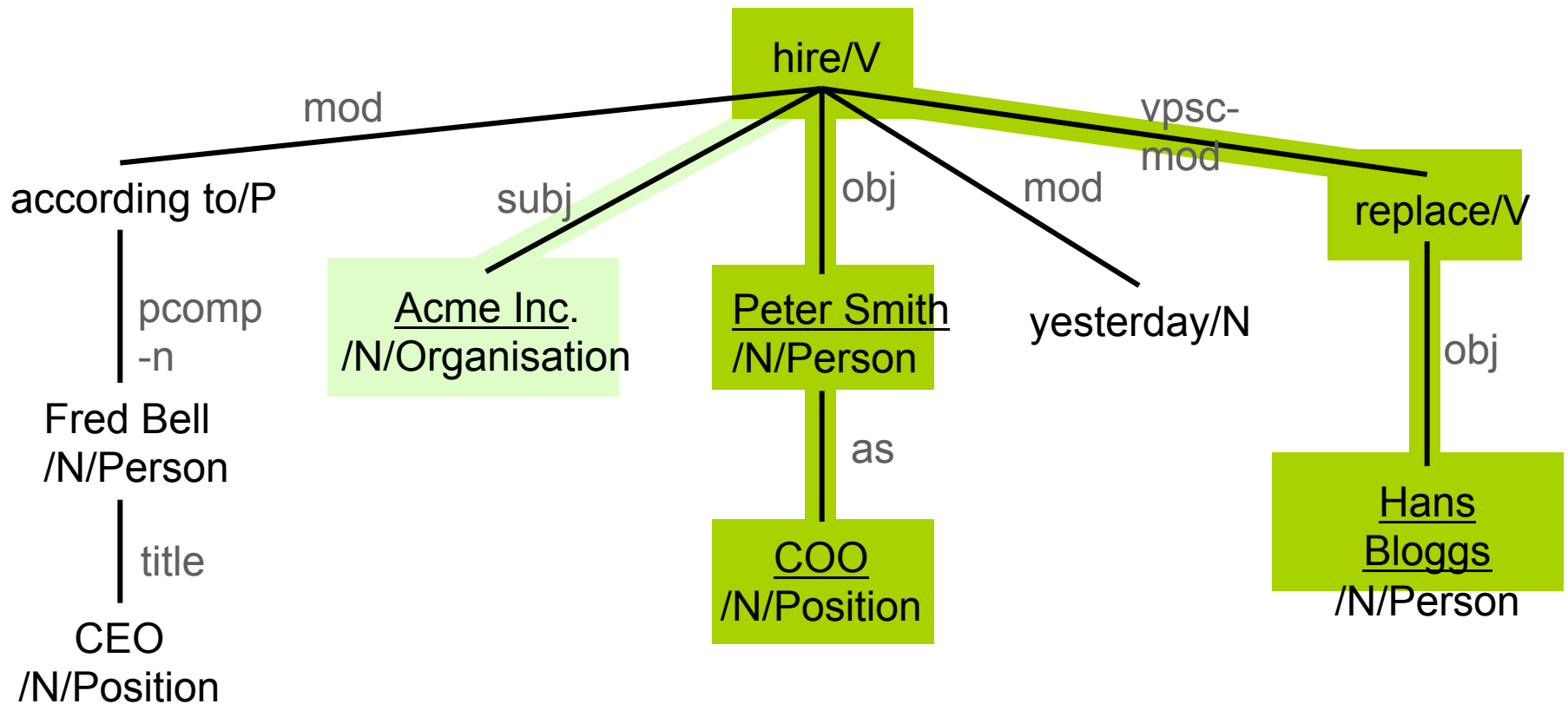
- Verb centered
- A single syntactic path dominated by a verb containing at least one relevant named entity concept



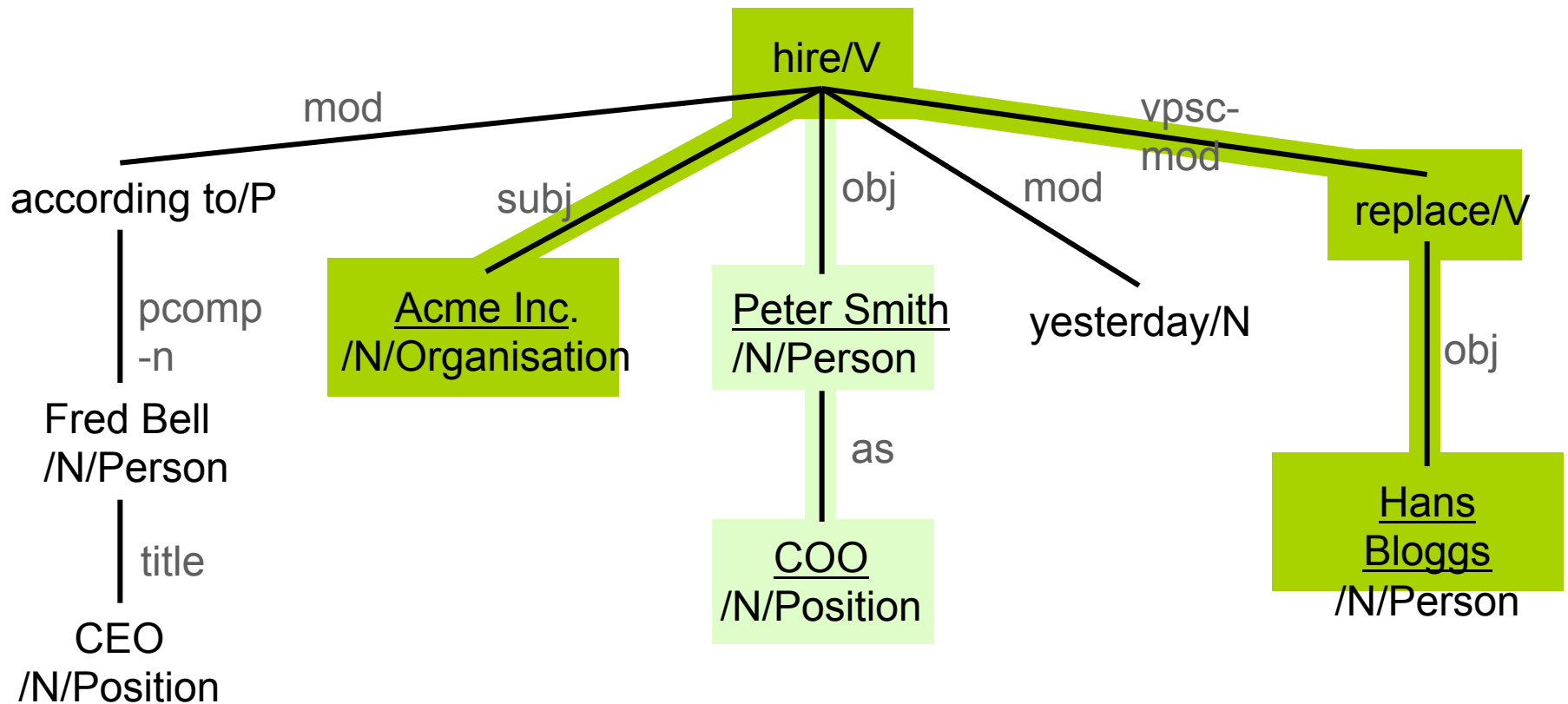
- verb centered
- pairs of chains instead of single paths



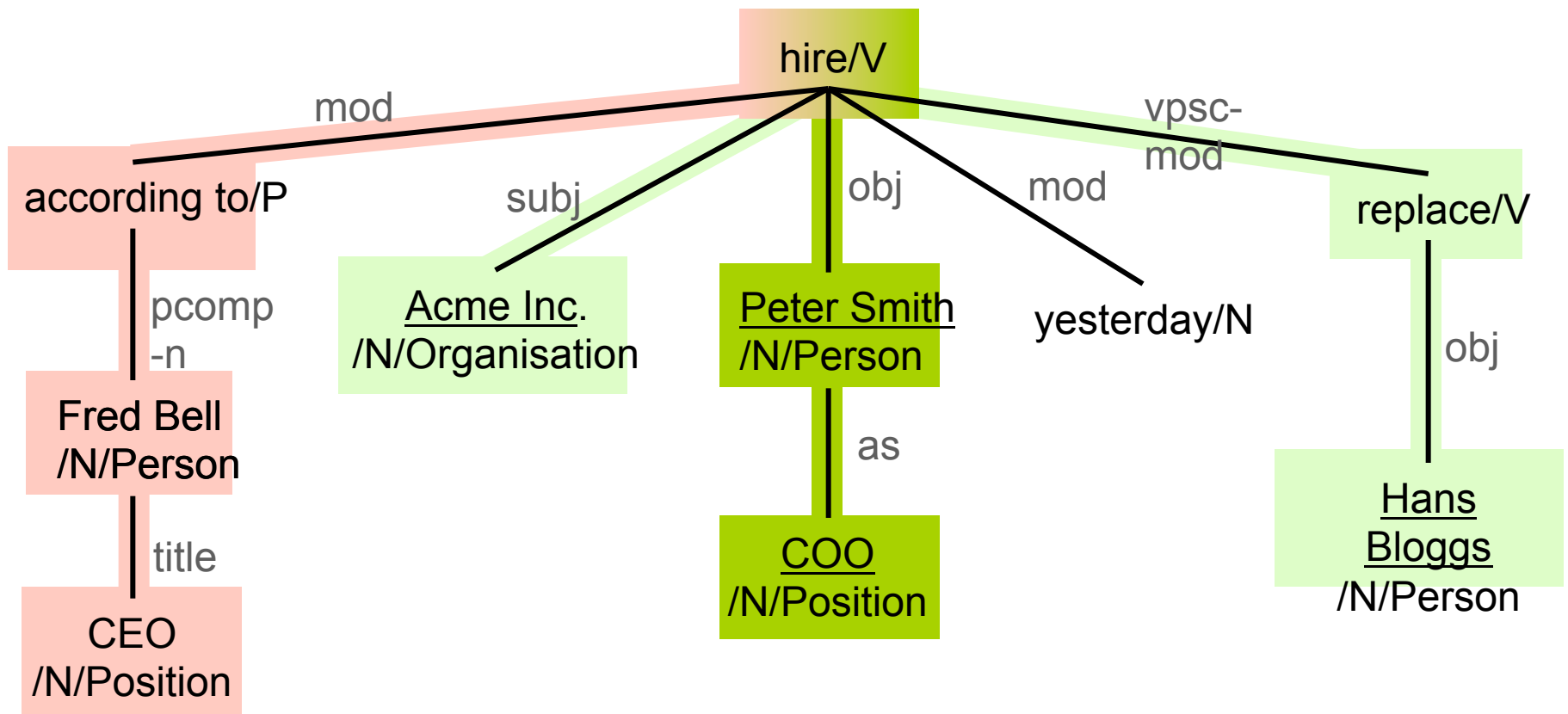
- verb centered
- pairs of chains instead of single paths



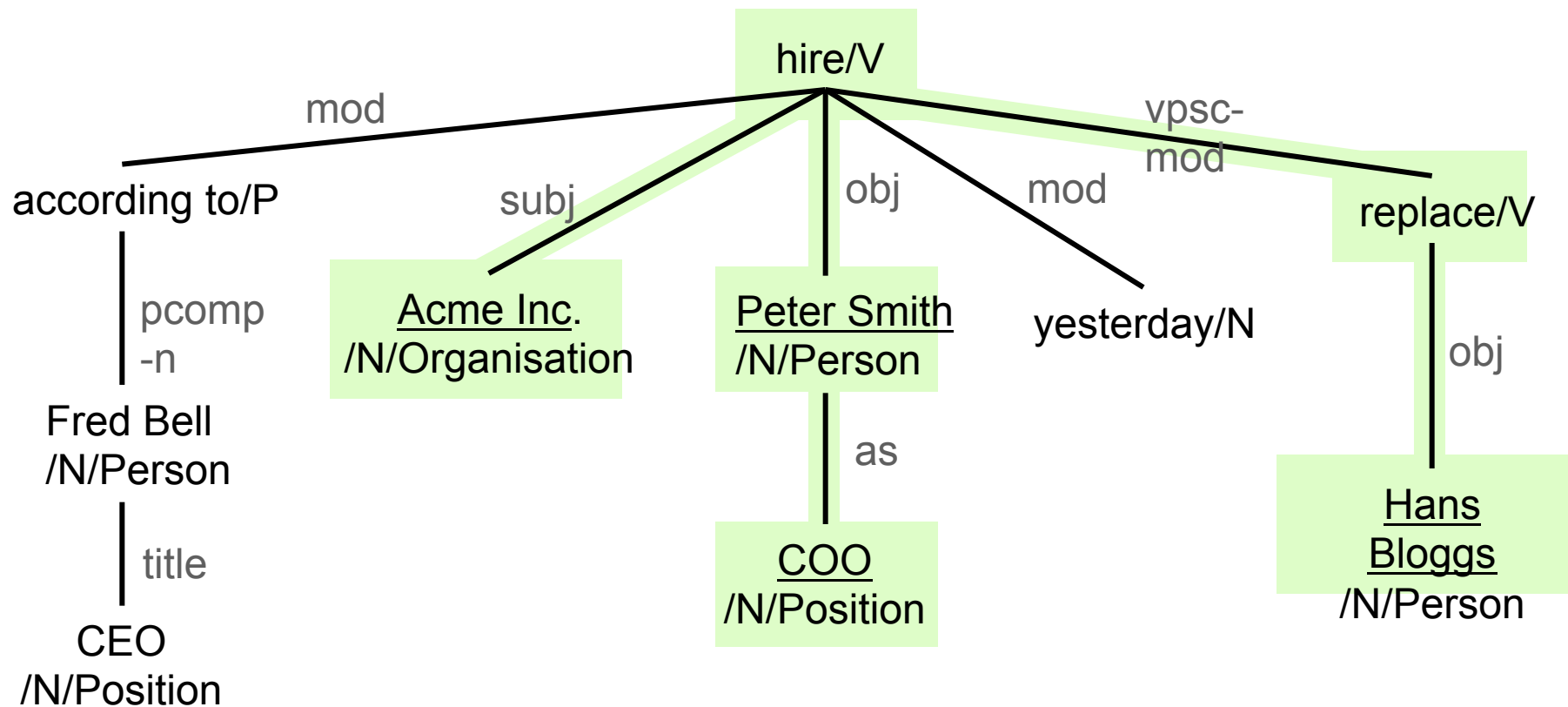
- verb centered
- pairs of chains instead of single paths



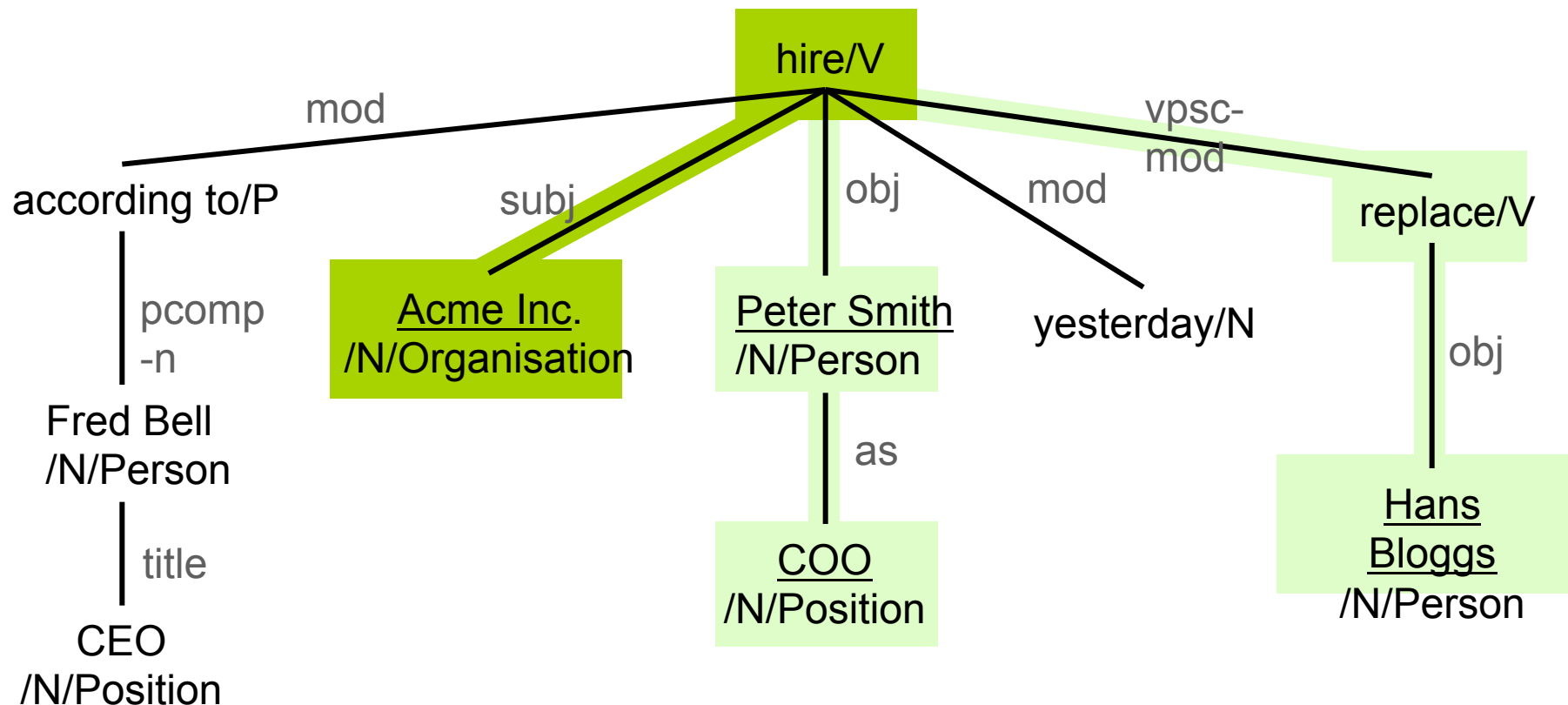
- verb centered
- pairs of chains instead of single paths



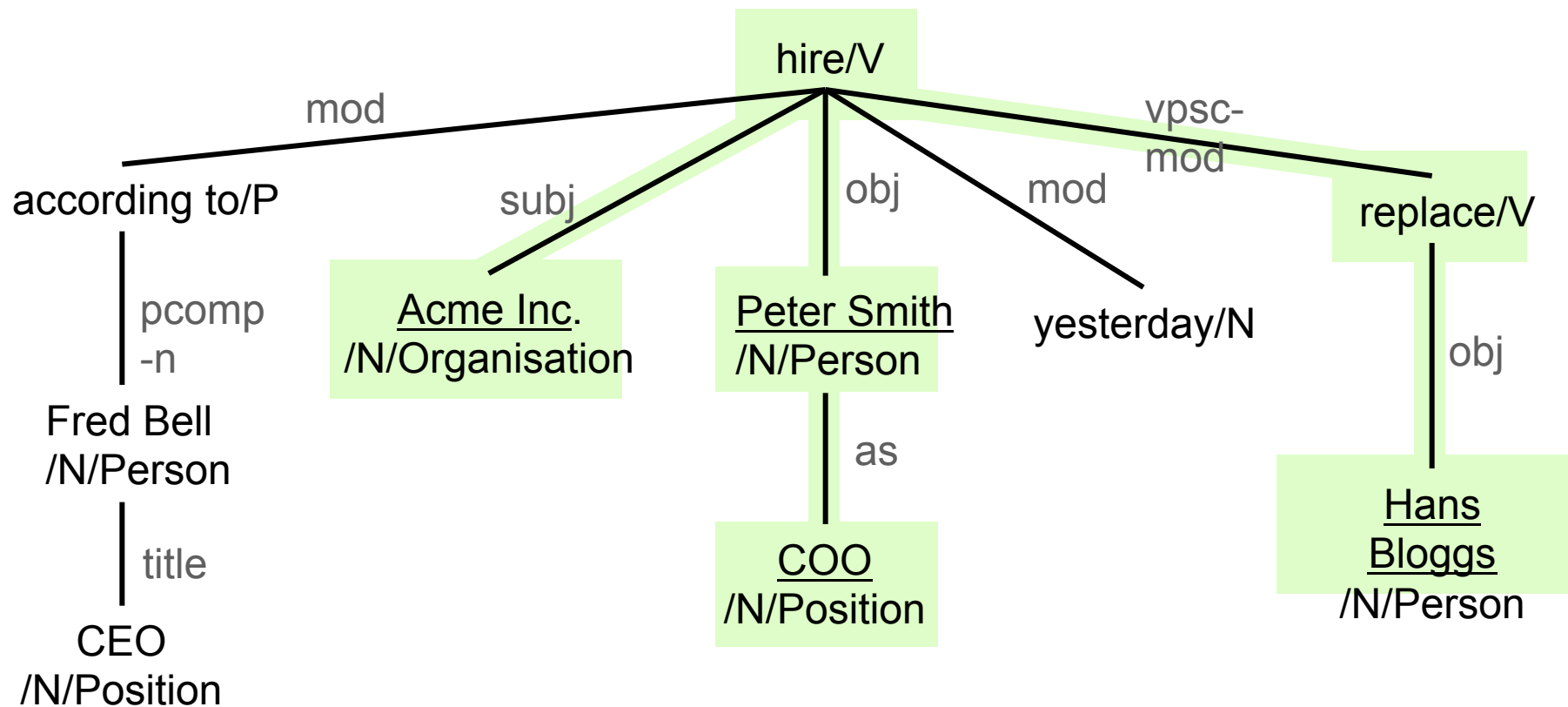
- verb centered
- All chains dominated by a verb, which contain at least one relevant named entity and their combinations



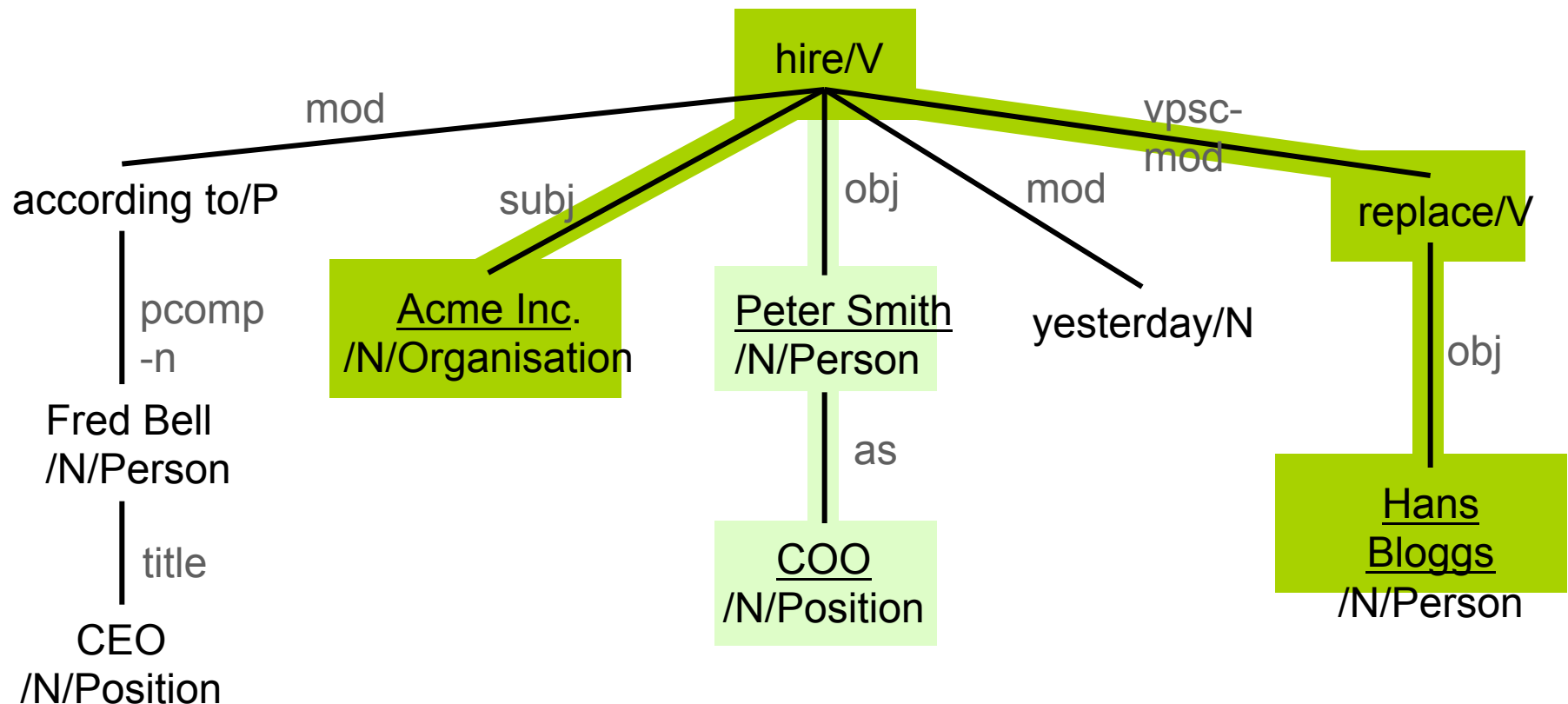
- verb centered
- All chains dominated by a verb, which contain at least one relevant named entity and their combinations



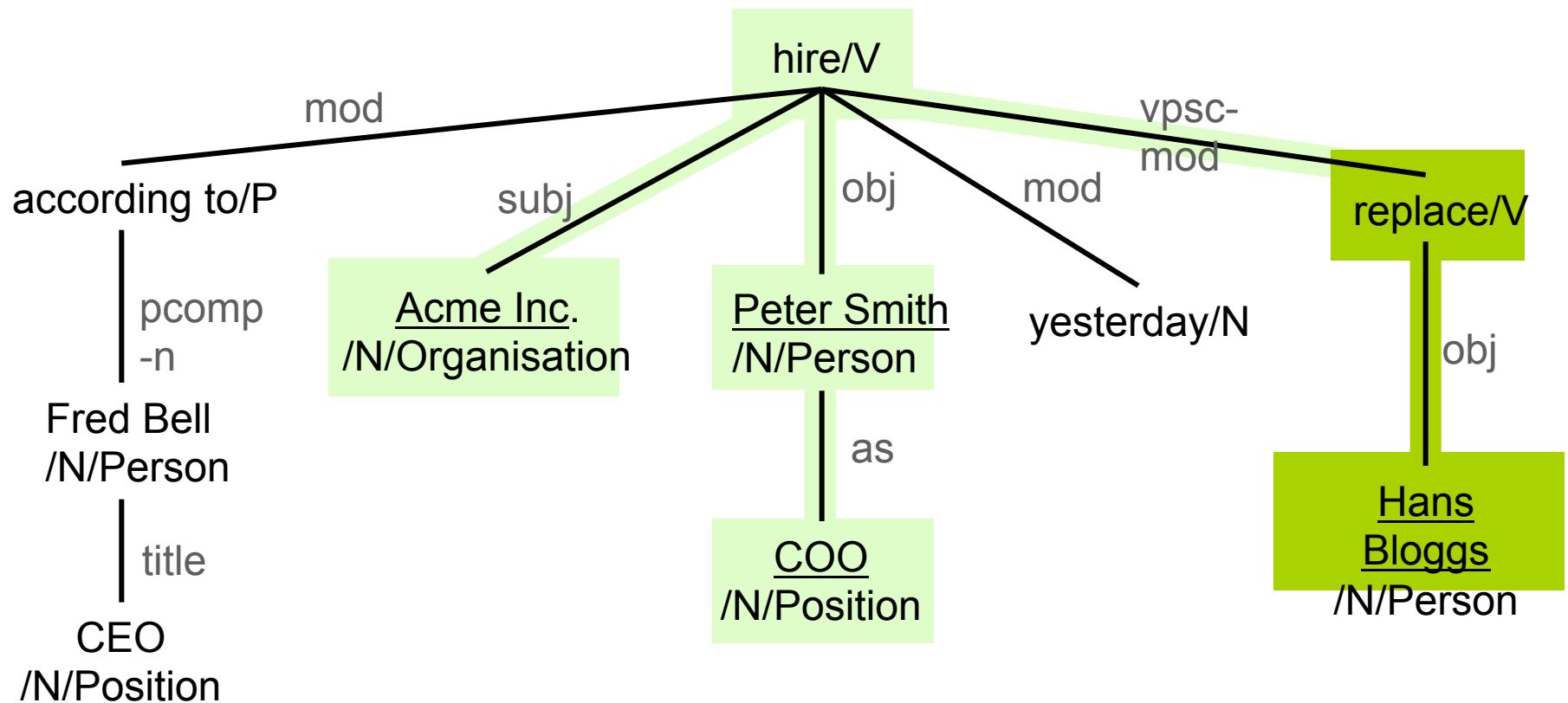
- verb centered
- All chains dominated by a verb, which contain at least one relevant named entity and their combinations



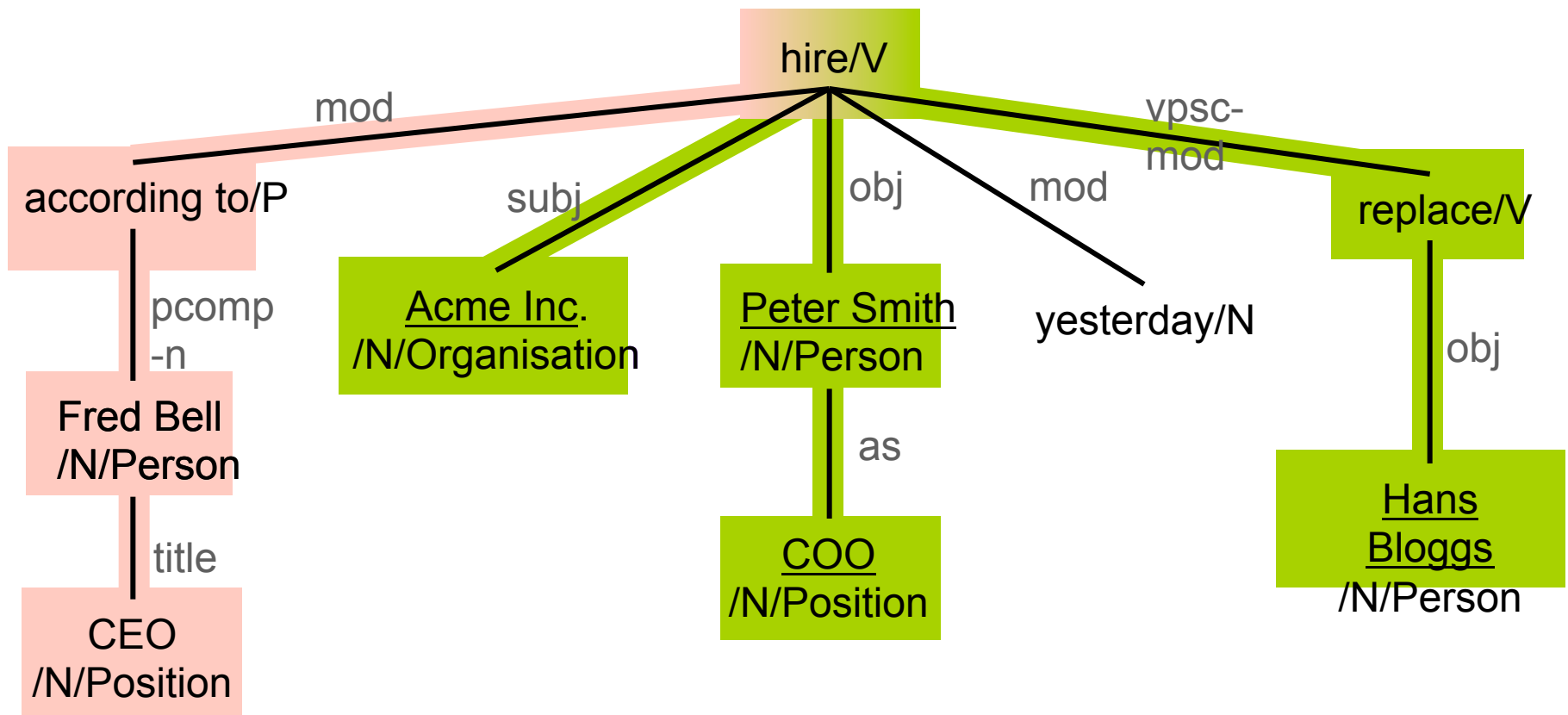
- verb centered
- All chains dominated by a verb, which contain at least one relevant named entity and their combinations



- verb centered
- All chains dominated by a verb, which contain at least one relevant named entity and their combinations

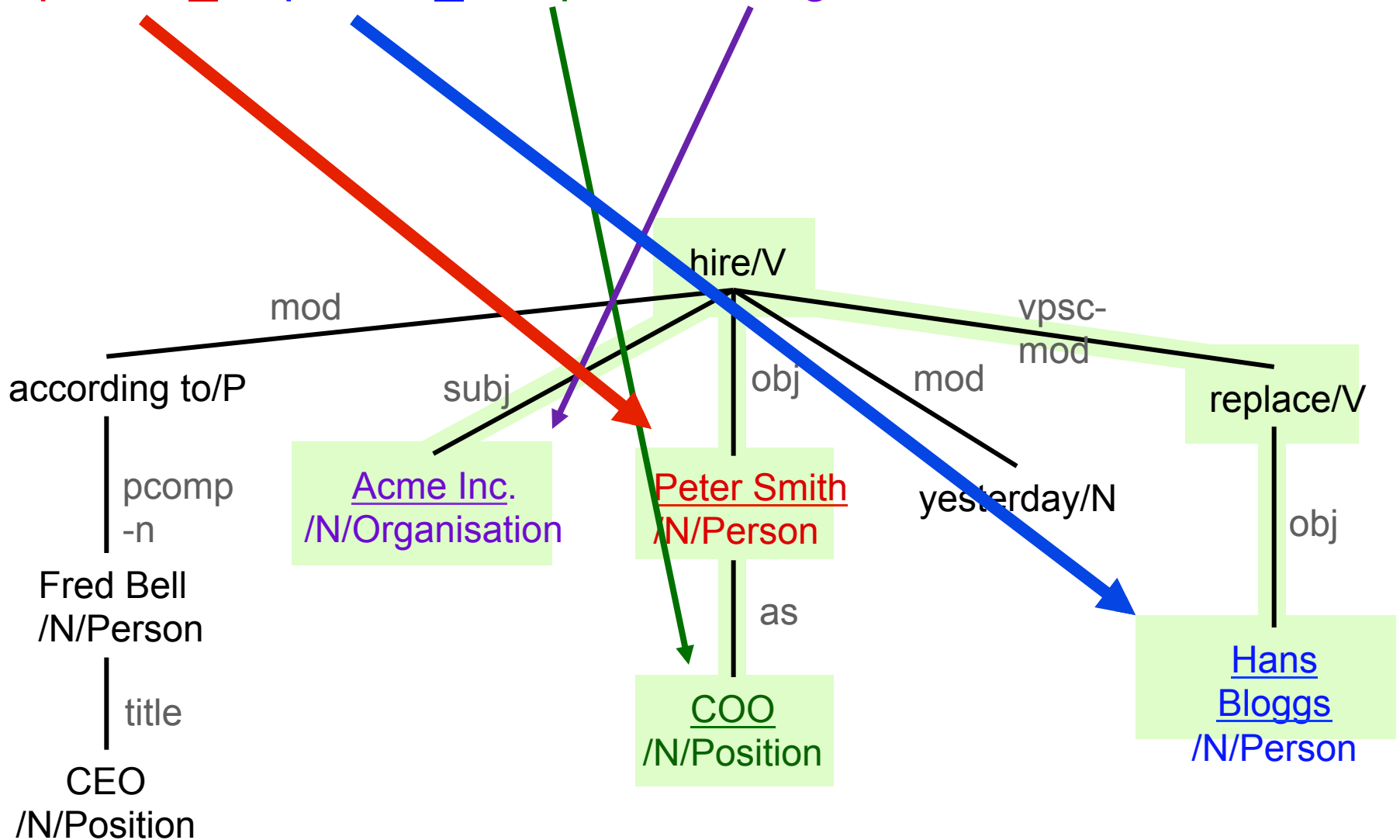


- verb centered
- All chains dominated by a verb, which contain at least one relevant named entity and their combinations



None of the existing models links the detected slot-filling candidates with their respective semantic roles

<person\_in, person\_out, position, organisation>



# Properties of DARE

- ❑ Samples of target relation instances serve as semantic seed
- ❑ Systematic treatment of n-ary relations and their projections
- ❑ Exploitation of relation projections for pattern discovery
- ❑ Bottom-up compositional pattern discovery
- ❑ A recursive linguistic rule representation
- ❑ Rules contain semantic roles w.r.t. to target relation
- ❑ Bottom-up compression method to generalize rules
- ❑ Filtering of rule candidates by “domain relevance”

# Novel Properties of DARE

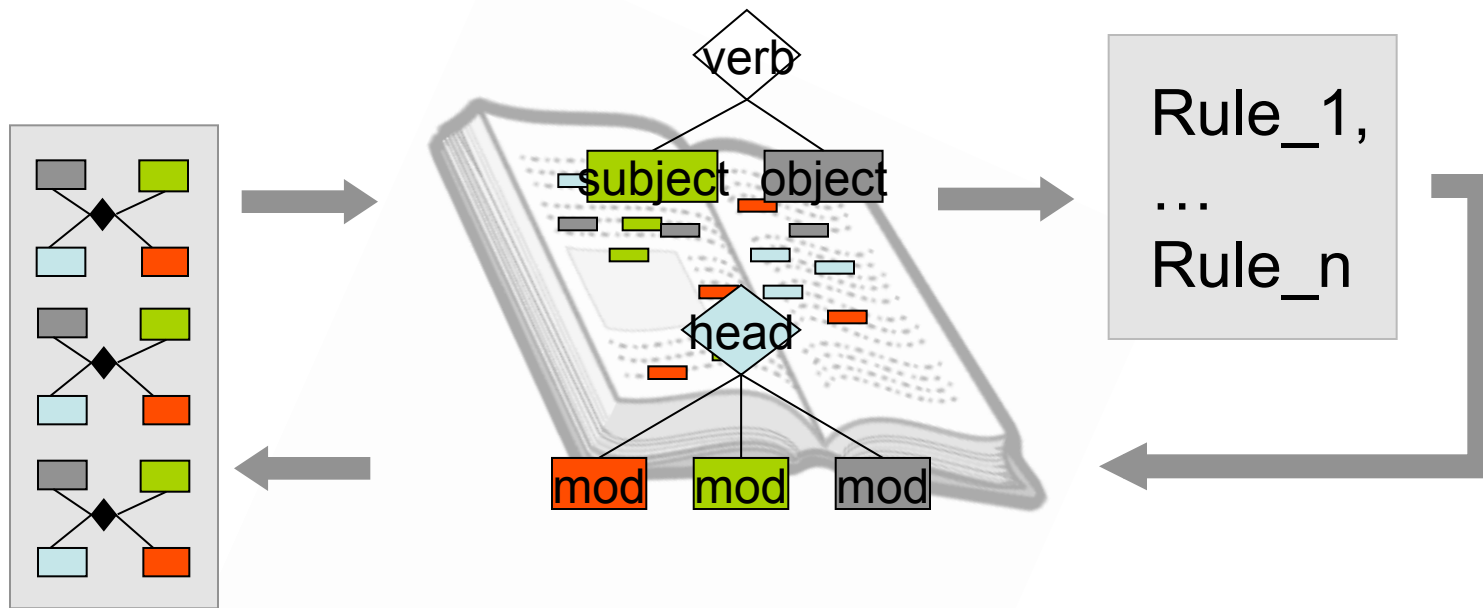
- Samples of target relation instances serve as semantic seed
- Systematic treatment of n-ary relations and their projections
- Exploitation of relation projections for pattern discovery
- Bottom-up compositional pattern discovery
- A recursive linguistic rule representation
- Rules contain semantic roles w.r.t. to target relation
- Bottom-up compression method to generalize rules
- Filtering of rule candidates by “domain relevance”

# Bootstrapping Relation Extraction with Semantic Seed

Adapted from

DIPRE (Brin, 1998) and Snowball (Agichtein & Gravano, 2000)

but extended and enriched with linguistic analysis



# Bootstrapping Relation Extraction with Semantic Seed

- DIPRE and Snowball

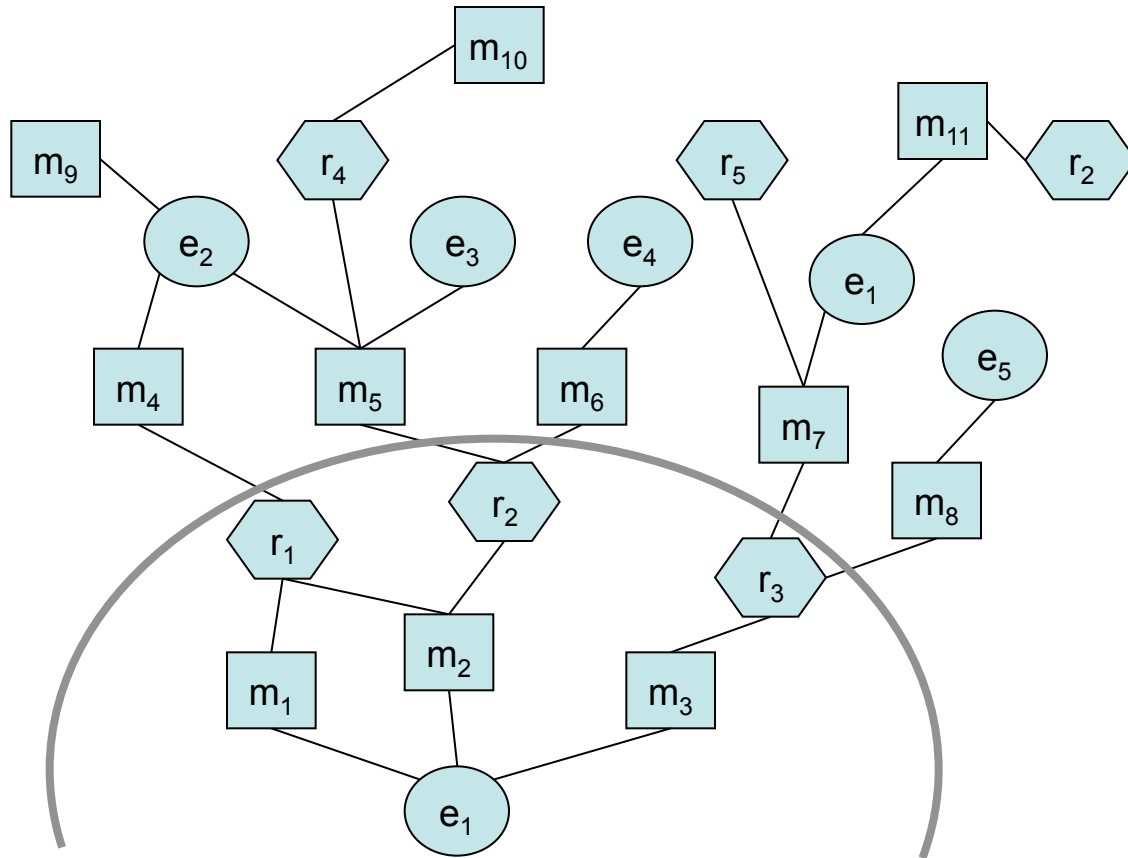
- binary relations only, no projections, no linguistic analysis

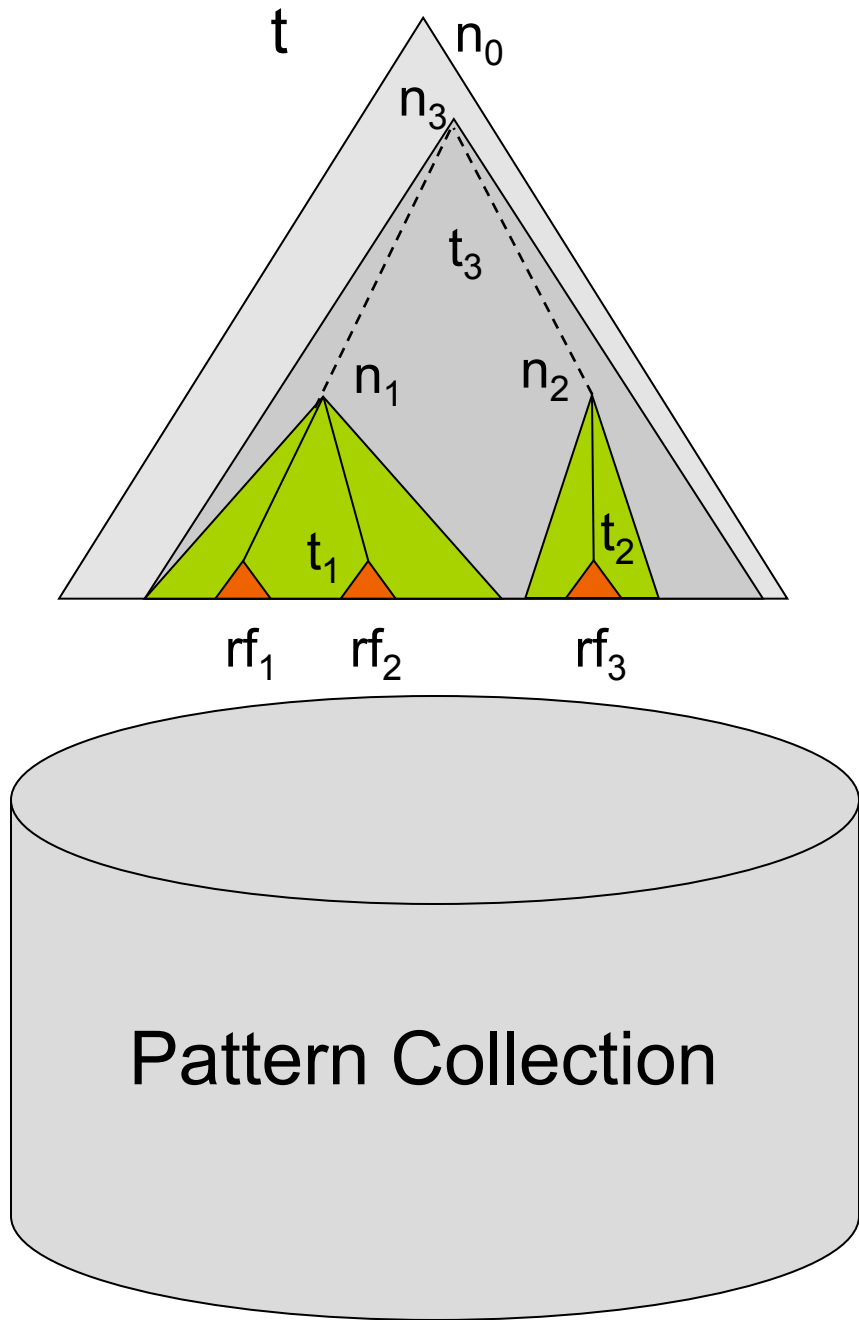
- DARE

- n-ary relations and their projections, deep linguistic analysis

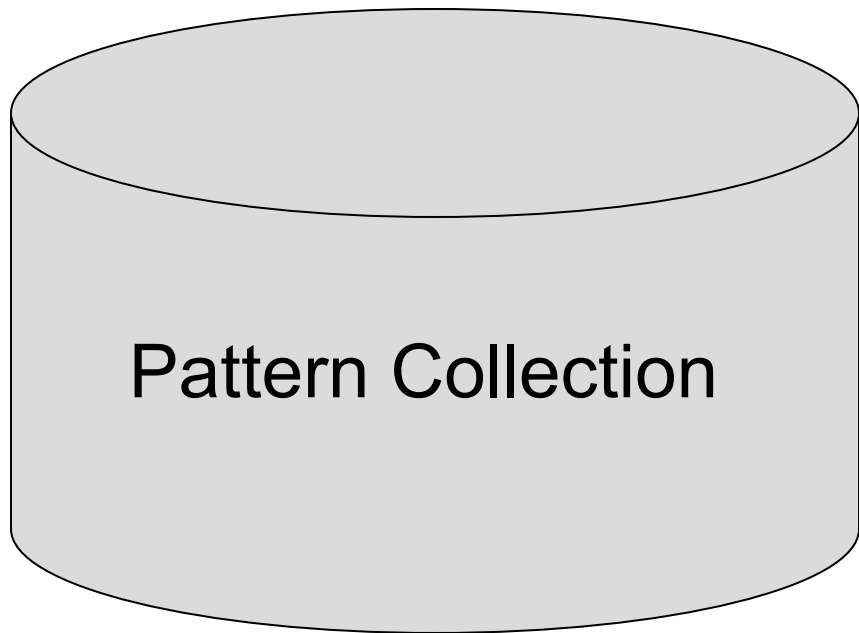
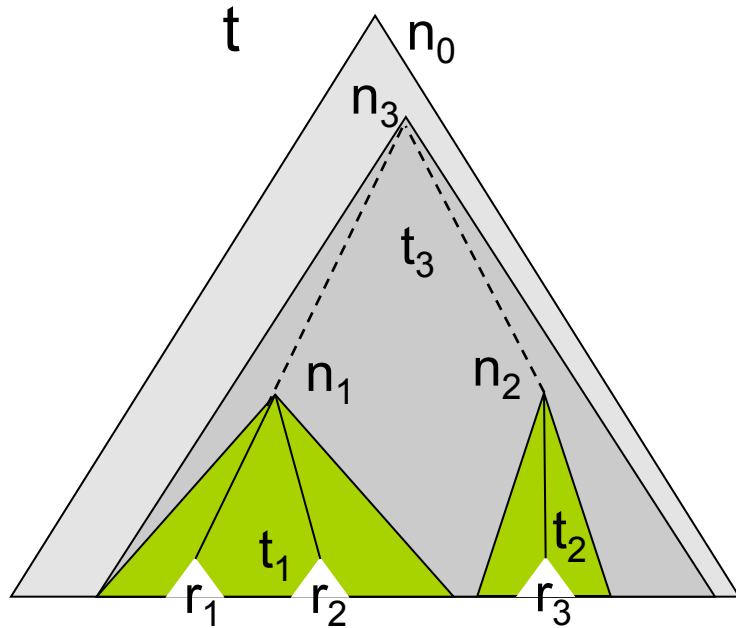
(in the experiments I use MINIPAR by Dekan Lin 1999)

# Start of Bootstrapping (simplified)

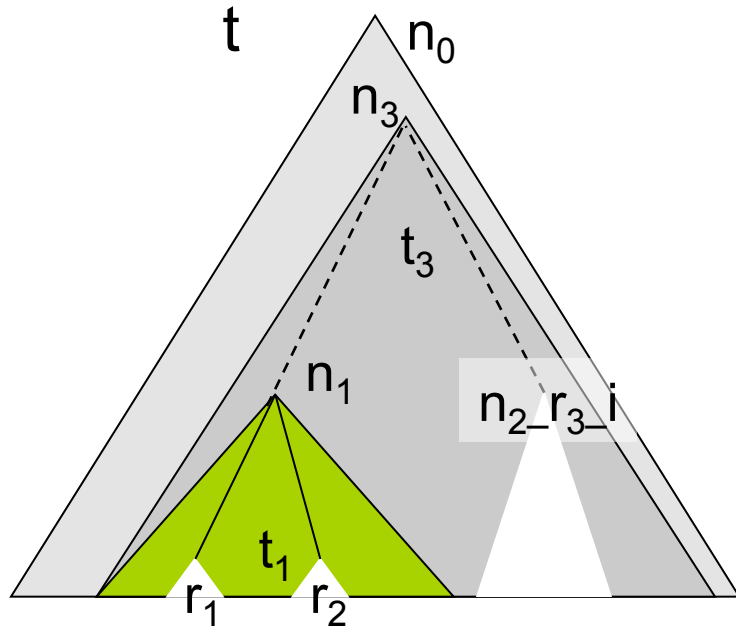




0. replace all nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;



0. replace all nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;



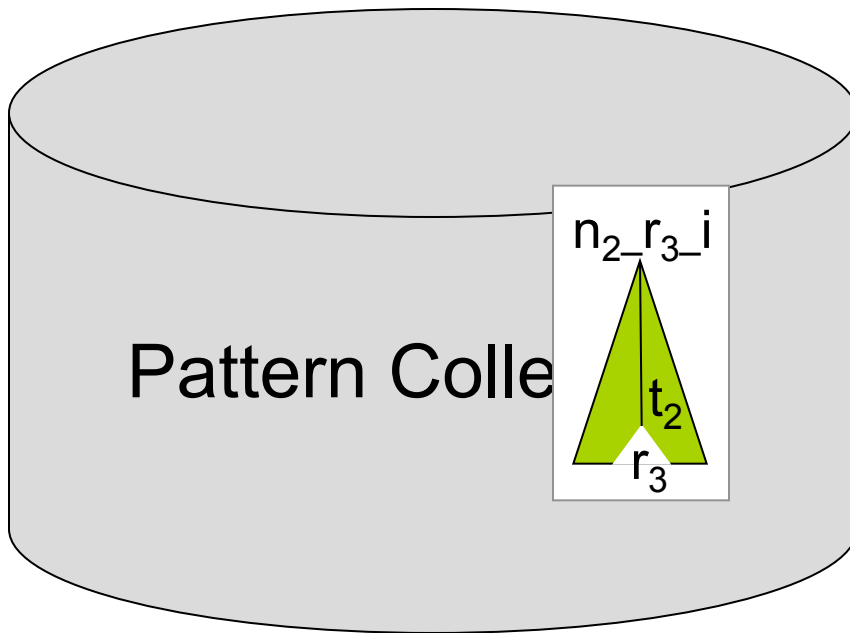
0. replace all nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;

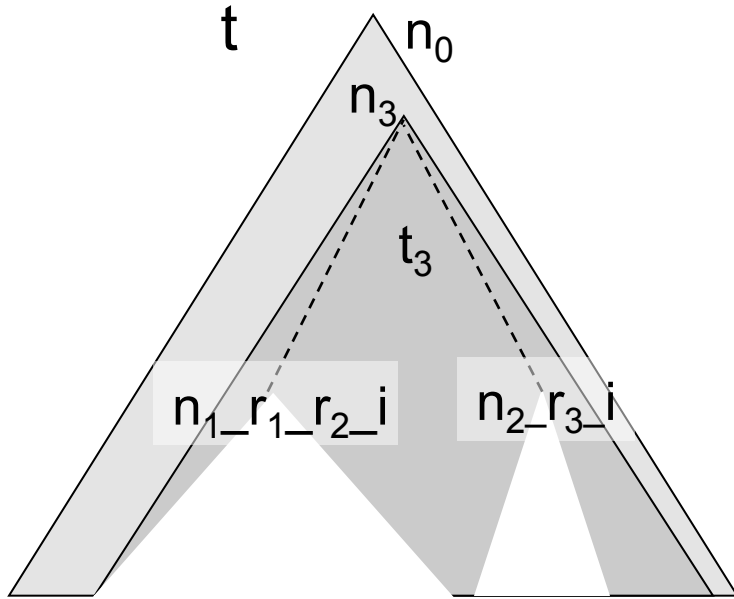
for  $i=1$  to  $n$

1. identify the set of the lowest non-terminal nodes  $N_1$  in  $t$  that dominate  $i$  arguments (possibly among other nodes).

2. substitute  $N_1$  by nodes labelled with the seed argument roles and their entity classes

3. prune the subtrees dominated by  $N_1$  from  $t$  and add these subtrees into the pattern collection. These subtrees are assigned the argument role information and a unique id.





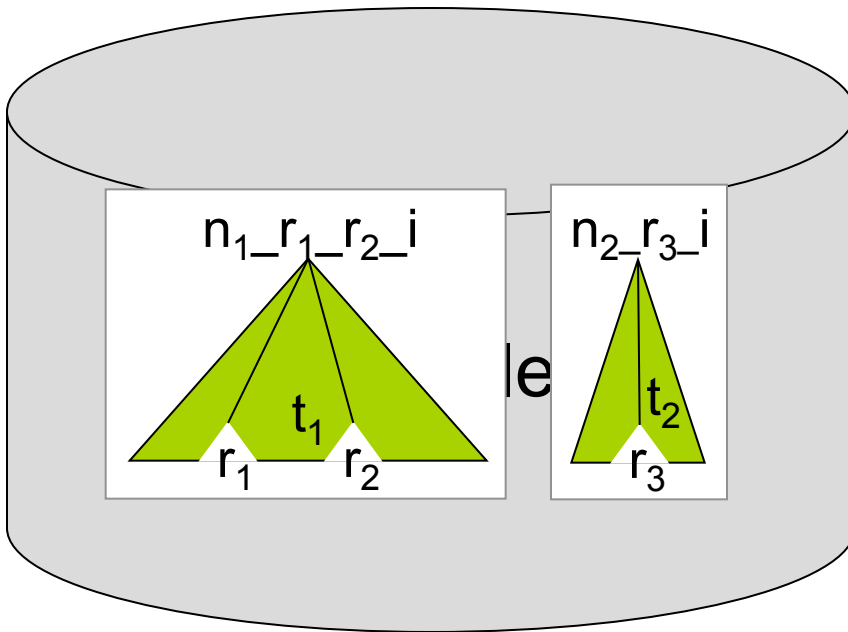
0. replace all nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;

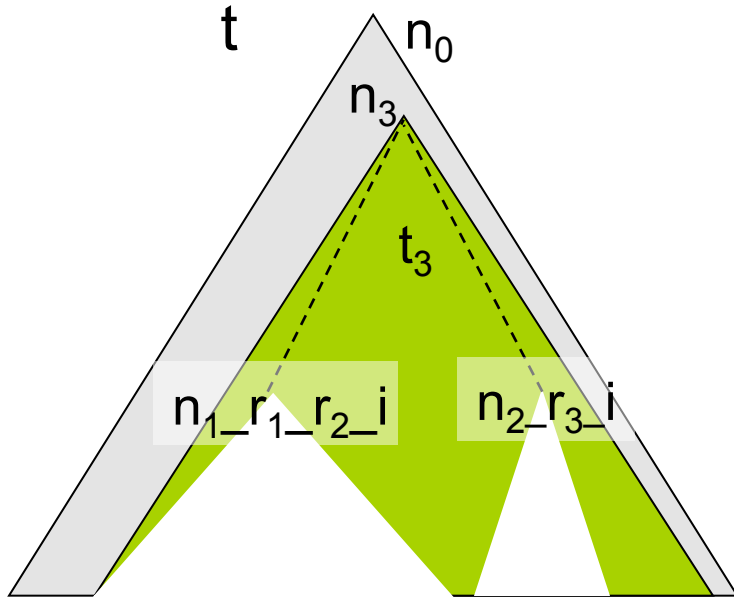
for  $i=1$  to  $n$

1. identify the set of the lowest non-terminal nodes  $N_1$  in  $t$  that dominate  $i$  arguments (possibly among other nodes).

2. substitute  $N_1$  by nodes labelled with the seed argument roles and their entity classes

3. prune the subtrees dominated by  $N_1$  from  $t$  and add these subtrees into the pattern collection. These subtrees are assigned the argument role information and a unique id.





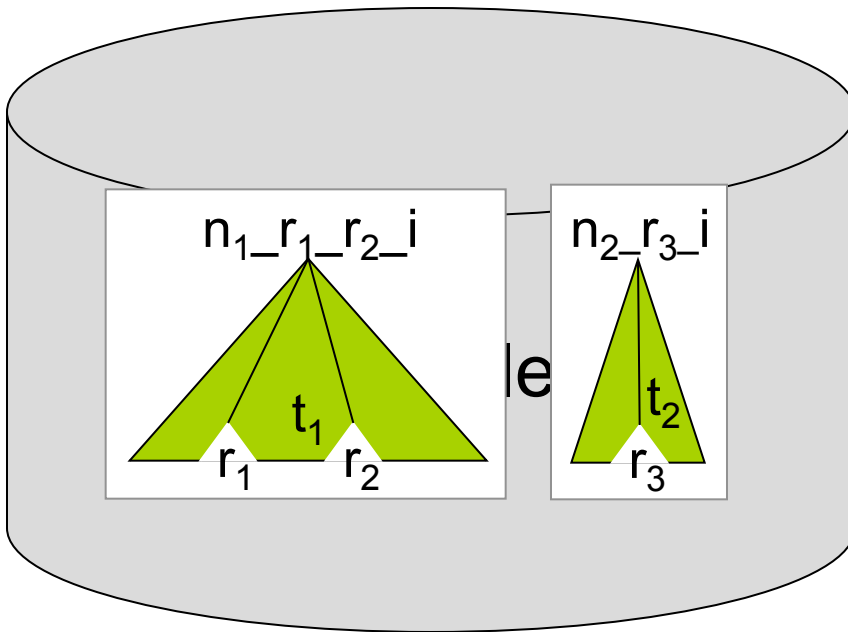
0. replace all nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;

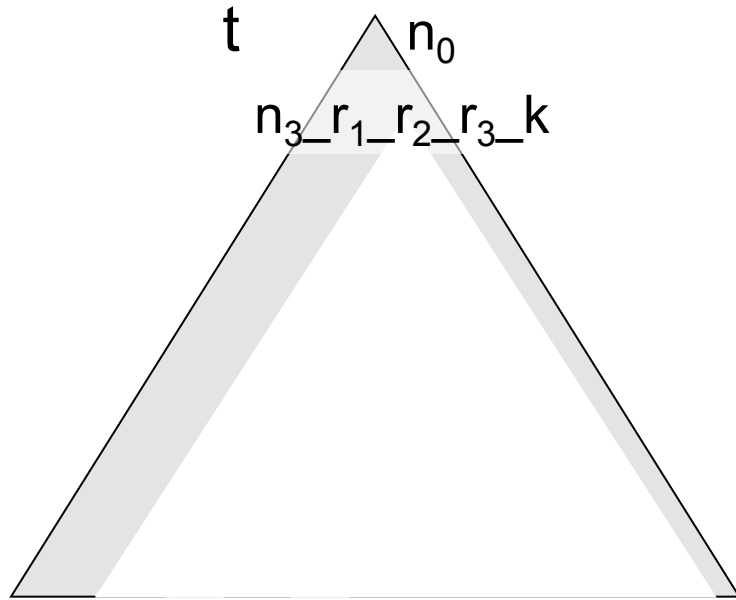
for  $i=1$  to  $n$

1. identify the set of the lowest non-terminal nodes  $N_1$  in  $t$  that dominate  $i$  arguments (possibly among other nodes).

2. substitute  $N_1$  by nodes labelled with the seed argument roles and their entity classes

3. prune the subtrees dominated by  $N_1$  from  $t$  and add these subtrees into the pattern collection. These subtrees are assigned the argument role information and a unique id.





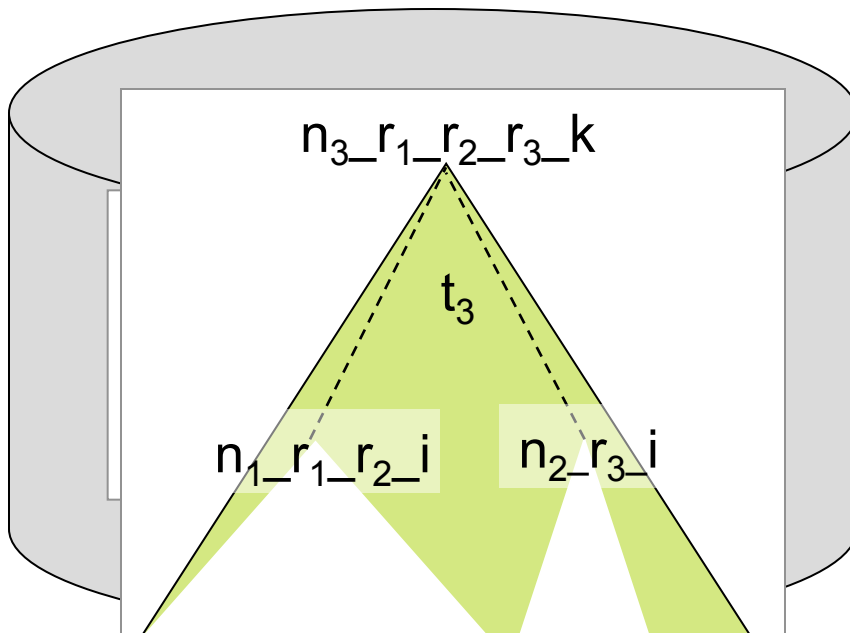
0. replace all nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;

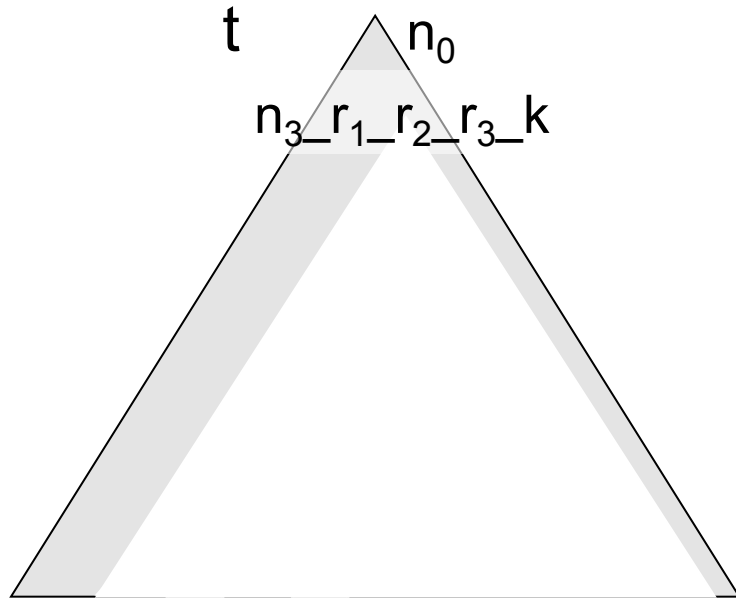
for  $i=1$  to  $n$

1. identify the set of the lowest non-terminal nodes  $N_1$  in  $t$  that dominate  $i$  arguments (possibly among other nodes).

2. substitute  $N_1$  by nodes labelled with the seed argument roles and their entity classes

3. prune the subtrees dominated by  $N_1$  from  $t$  and add these subtrees into the pattern collection. These subtrees are assigned the argument role information and a unique id.





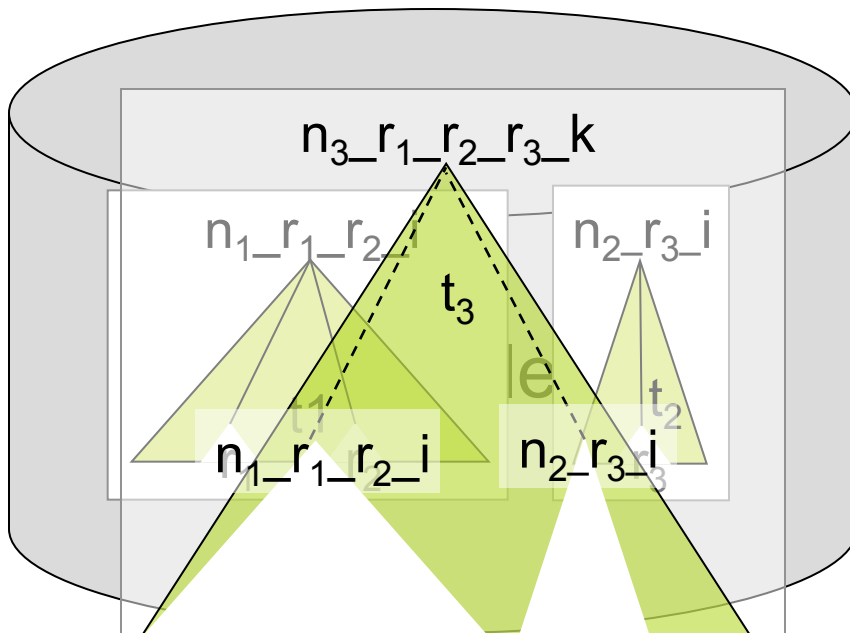
0. replace all nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;

for  $i=1$  to  $n$

1. identify the set of the lowest non-terminal nodes  $N_1$  in  $t$  that dominate  $i$  arguments (possibly among other nodes).

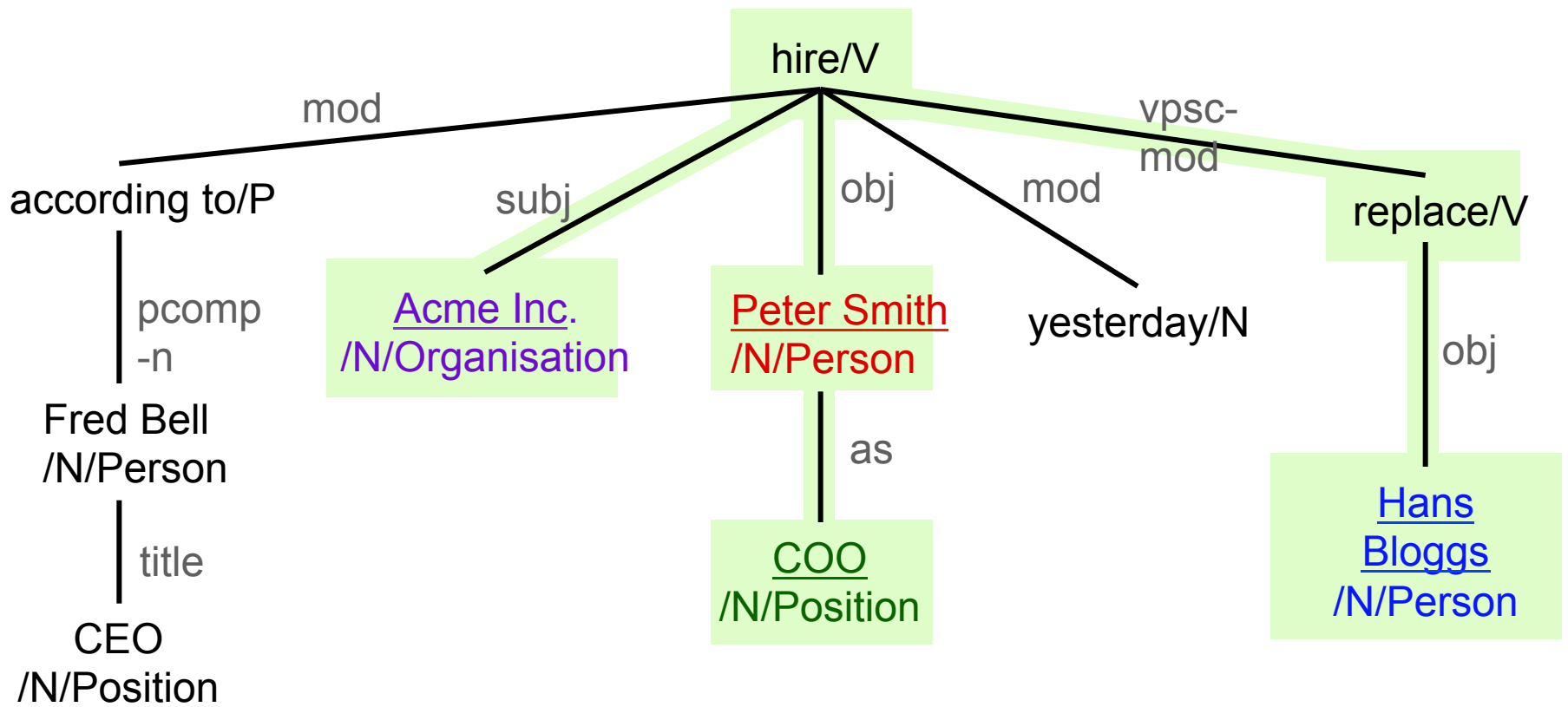
2. substitute  $N_1$  by nodes labelled with the seed argument roles and their entity classes

3. prune the subtrees dominated by  $N_1$  from  $t$  and add these subtrees into the pattern collection. These subtrees are assigned the argument role information and a unique id.



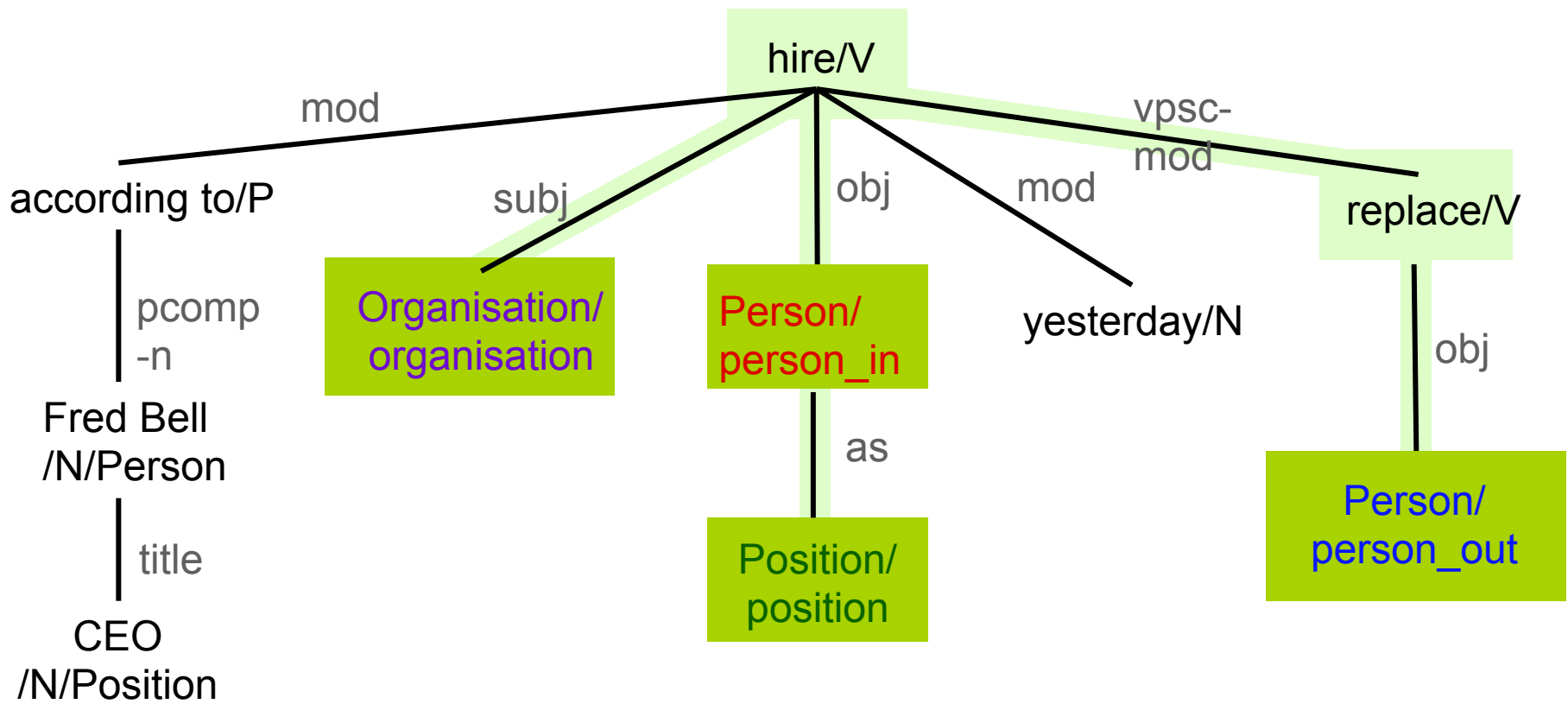
According to CEO Fred Bell, [Acme Inc.](#) hired [Peter Smith](#) as [COO](#) yesterday, replacing [Hans Bloggs](#).

<[Peter Smith](#)/person\_in, [Hans Bloggs](#)/person\_out, [COO](#) /position, [Acme Inc.](#) /organisation>



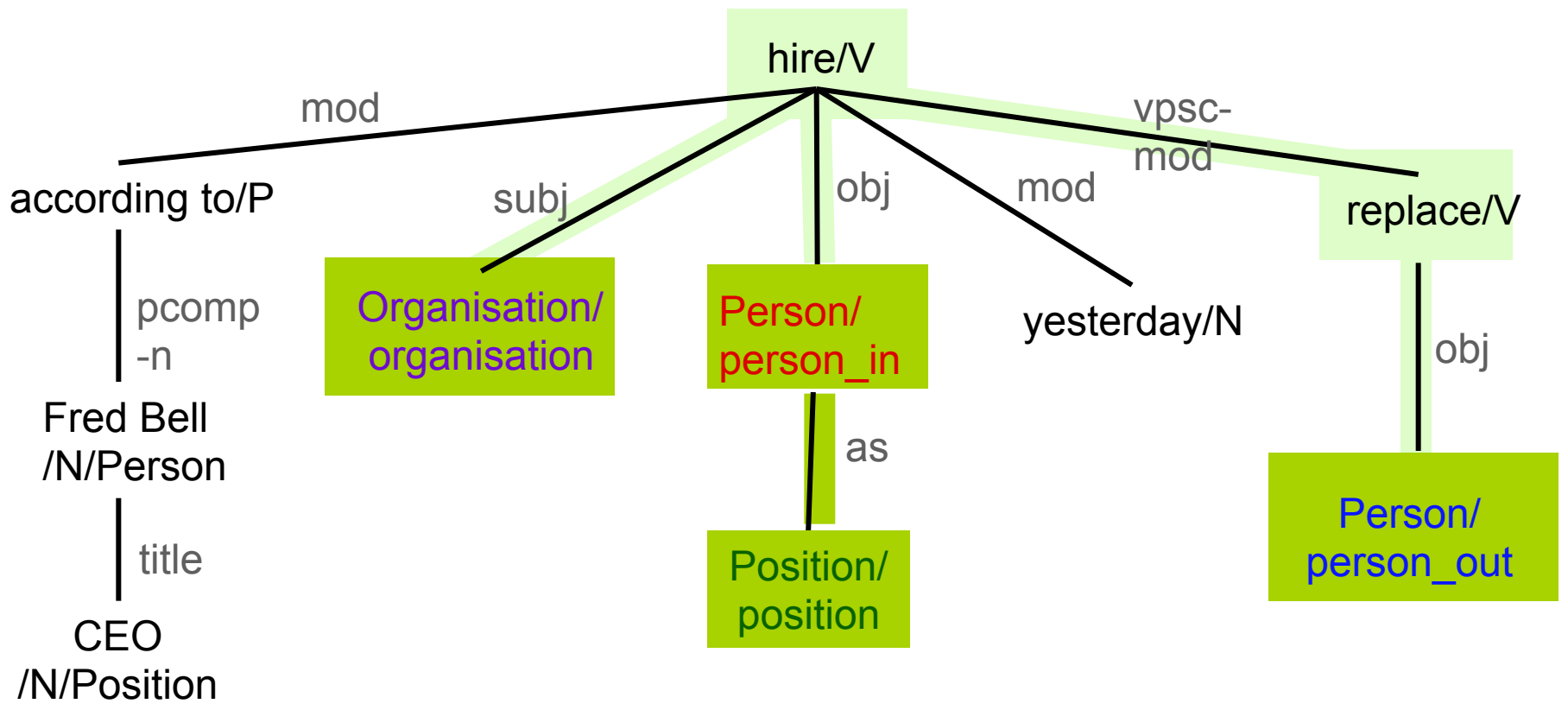
According to CEO Fred Bell, [Acme Inc.](#) hired [Peter Smith](#) as [COO](#) yesterday, replacing [Hans Bloggs](#).

<[Peter Smith](#)/person\_in, [Hans Bloggs](#)/person\_out, [COO](#) /position, [Acme Inc.](#) /organisation>



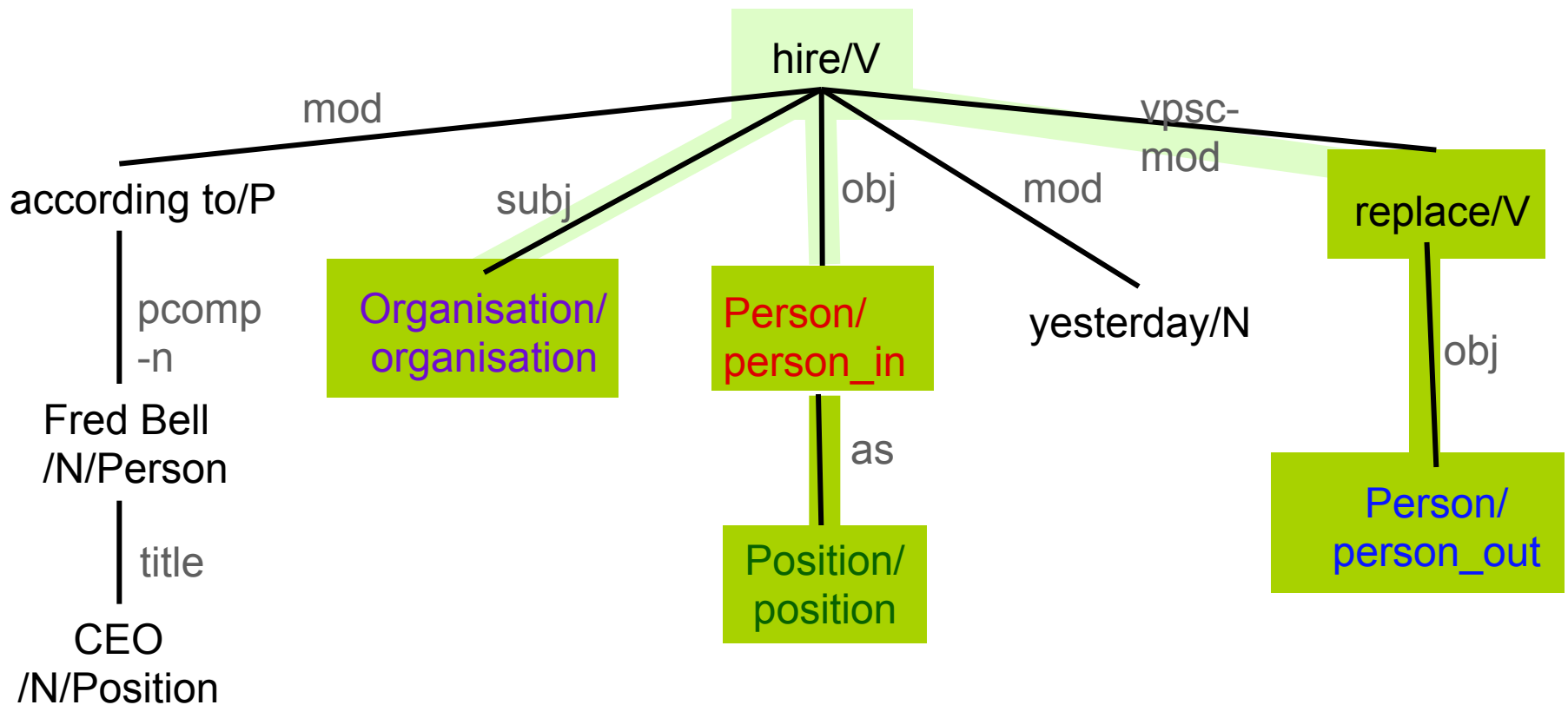
According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<Peter Smith/person\_in, Hans Bloggs/person\_out, COO /position, Acme Inc. /organisation>



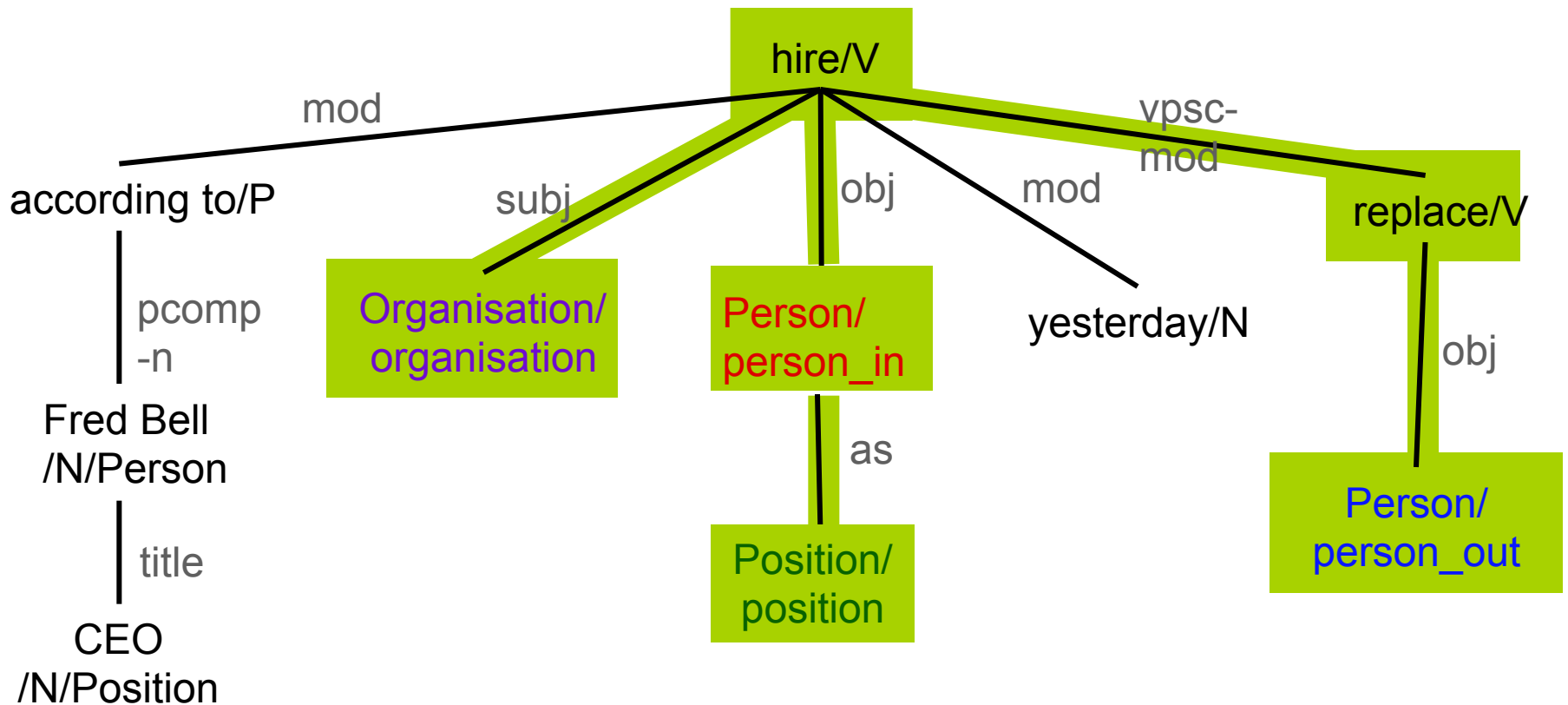
According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<Peter Smith/person\_in, Hans Bloggs/person\_out, COO /position, Acme Inc. /organisation>

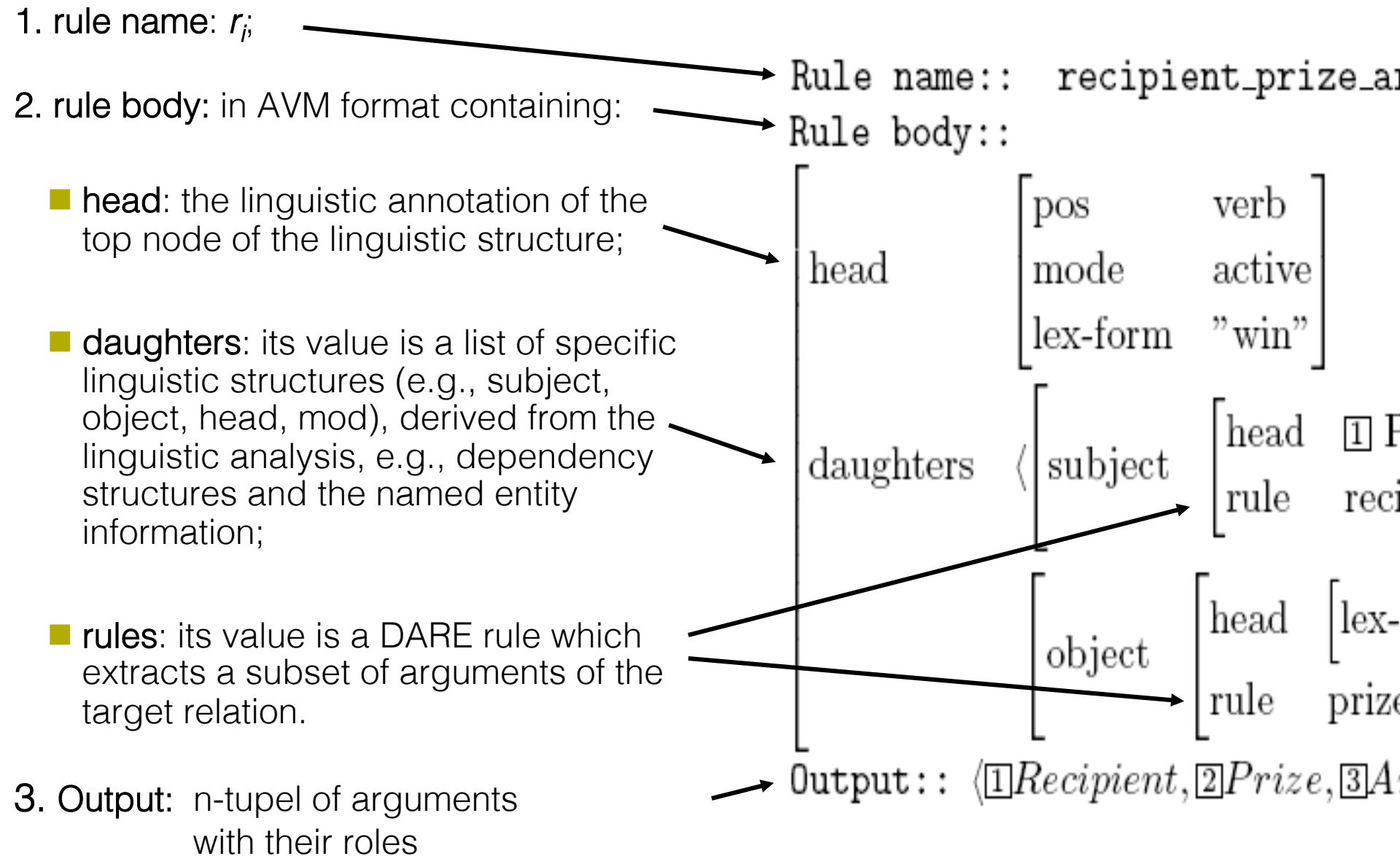


According to CEO Fred Bell, [Acme Inc.](#) hired [Peter Smith](#) as [COO](#) yesterday, replacing [Hans Bloggs](#).

<[Peter Smith](#)/person\_in, [Hans Bloggs](#)/person\_out, [COO](#) /position, [Acme Inc.](#) /organisation>



# DARE Rule Components



# DARE Rule Components

Rule name:: recipient\_prize\_area\_year\_1

Rule body::

```
[ head [ pos      verb
        mode     active
        lex-form  "win" ]
  daughters < [ subject [ head [1] Person
                          rule  recipient_1:: <[1]Person > ] ],
              [ object [ head [ lex-form  "prize" ]
                          rule  prize_area_year_1:: <[2]Prize, [3]Area, [4]Year > ] ] ] >
```

Output:: <[1]Recipient, [2]Prize, [3]Area, [4]Year >

# prize\_area\_year\_1

Rule name:: prize\_area\_year\_1

Rule body::

head	[	pos	noun	]							
		lex-form	"prize"								
daughters	<	[	lex-mod	[	head	[	3	Year	]]	,	
			[	lex-mod	[	head	[	1	Prize	]]	,
			[	lex-mod	[	head	[	2	Area	]]	>

Output:: <[1]Prize,[2]Area,[3]Year>

# Two Domains

- Award Events (start with subdomain Nobel Prizes)

reasons: good news coverage  
complete list of all award events  
good starting point for other award domains

- Management Succession Events

reason: comparison with previous work

# Experiments

- Two domains
  - Nobel Prize Awards: <recipient, prize, area, year>
  - Management Succession: <person\_in, person\_out, position, organisation>
- Test data sets

<b>Data Set Name</b>	<b>Doc Number</b>	<b>Data Amount</b>
<b>Nobel Prize A (1981-1998)</b>	<b>1032</b>	<b>5.8 MB</b>
<b>Nobel Prize B (1999-2005)</b>	<b>2296</b>	<b>12.6 MB</b>
<b>Nobel Prize A+B</b>	<b>3328</b>	<b>18.4 MB</b>
<b>MUC-6</b>	<b>199</b>	<b>1MB</b>

# Evaluation Against Ideal Tables

Data Set	Seed	Precision	Recall
Nobel Prize A	<[Sen, Amartya], nobel, economics, 1998>	<b>87.3%</b>	<b>31.0%</b>
Nobel Prize A	<[Arias, Oscar], nobel, peace, 1987>	<b>83.8%</b>	<b>32.0%</b>
Nobel Prize B	<[Zewail, Ahmed H], nobel, chemistry, 1999>	<b>71.6%</b>	<b>50.7%</b>
A+B	<[Zewail, Ahmed H], nobel, chemistry, 1999>	<b>80.6%</b>	<b>62.9%</b>

# Management Succession Domain

Initial Seed #	Precision	Recall
<b>1</b>	<b>12.6%</b>	<b>7.0%</b>
<b>1</b>	<b>15.1%</b>	<b>21.8%</b>
<b>20</b>	<b>48.4%</b>	<b>34.2%</b>
<b>55</b>	<b>62.0%</b>	<b>48.0%</b>

# Comparison

Our result with 20 seeds (after 4 iterations)

- precision: 48.4%
- recall: 34.2%

compares well with the best result reported so far by (Greenwood and Stevenson, 2006) with the linked chain model starting with 7 hand-crafted patterns (after 190 iterations)

- precision: 43.4%
- recall: 26.5%

# Reusability of Rules

## □ Prize award patterns

- Detection of other Prizes such as *Pulitzer Prize*, *Turner Prize*
- Precision: 86.2%

## □ Management succession

- Domain independent binary pattern rules:  
*Person-Organisation*, *Person-Position*
- Evaluation of top 100 relation instances  
Precision: 98%

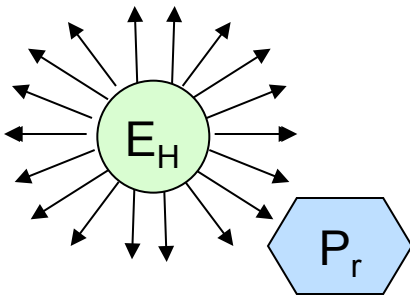
# Research Questions

As scientists we want to know

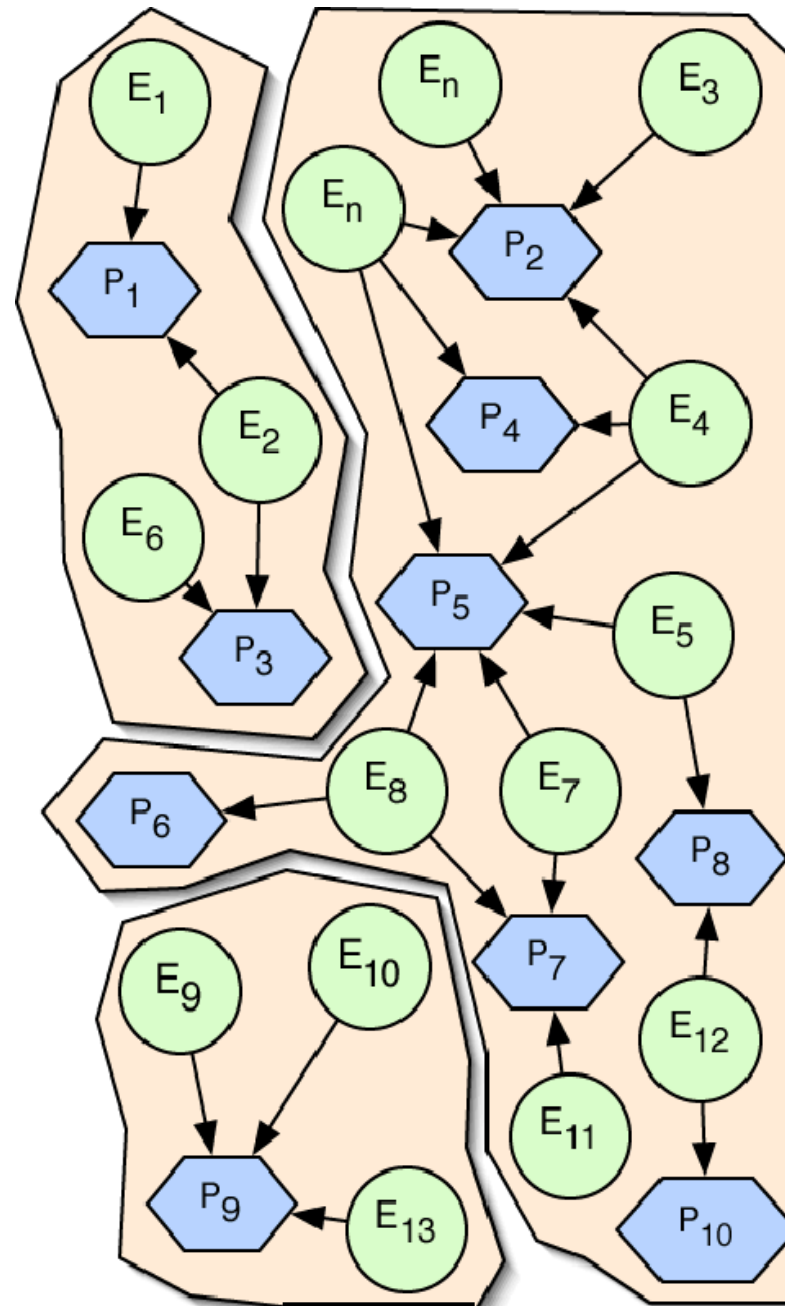
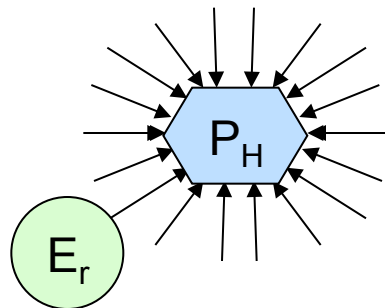
- Why does it work for some tasks?
- Why doesn't it work for all tasks?
- How can we estimate the suitability of domains?
- How can we deal with less suitable domains?

Careful analysis confirmed the following assumption:  
redundancy, both on patterns and event mentions, helps.

Frequently reported events make rare patterns reachable

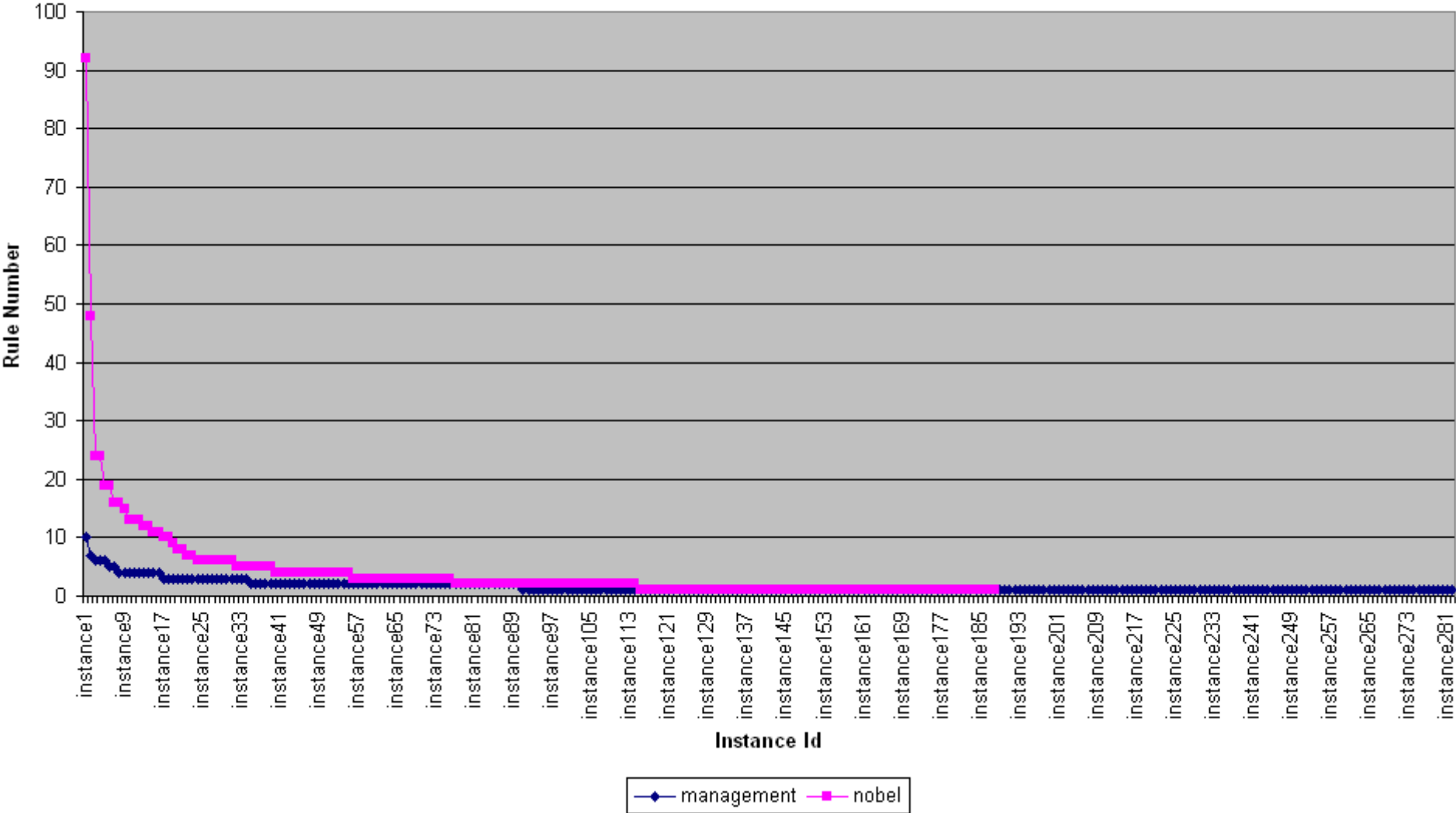


Popular patterns help to reach rarely mentioned events

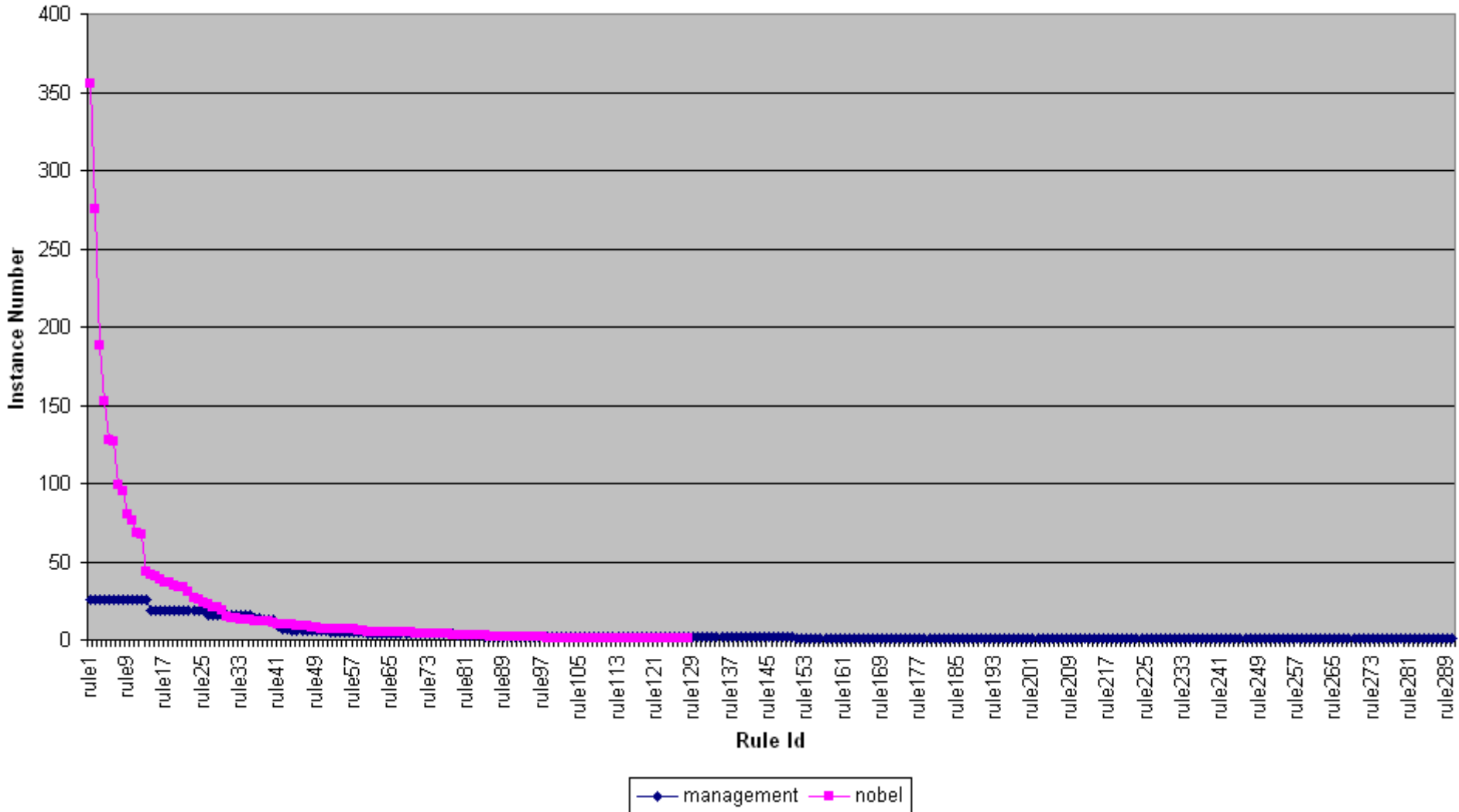


# Instance to Pattern

## Nobel Prize vs. Management Succession



# Rule to Instances (Nobel Prize vs. Management Succession)



# Insights

- Results from graph theory help to understand the requirements on data.

**Example:** small world property

- For data sets with continents and islands, we can sometimes exploit additional data or auxiliary domains to bridge the islands by learning rare patterns.

**Example:** use of Nobel prize domain for learning patterns for events concerning less popular prizes (many other prizes could be detected)

# Conclusion

- DARE is the first approach to combine the idea of bootstrapping IE systems with a linguistic grammar
- This can be illustrated by a simple formula:  
$$\begin{array}{l} \text{reusable generic linguistic knowledge} \\ + \text{ raw data} \\ + \text{ a few examples (seed)} \\ \hline = \text{ domain specific relation extraction grammar} \end{array}$$
- In addition to the obvious practical advantages, the approach offers theoretical benefits: It supports a view of IE as a systematic gradual approximation of language understanding.

---

# Domain-Adaptation for Big Text Data Analytics



# Text Analytics for Big Textual Data

---

- ◎ Three main features of big data
  - ◎ *Volume*: large-scale in volume
  - ◎ *Variety*: with respect to heterogeneous domains and formats
  - ◎ *Velocity*: because of its rapid and steady growing.
- Requirements of text analytics technologies for big data
  - *efficient*
  - *robust*
  - *scalable*
  - *domain-adaptive*

# Domain Adaptation is Essential for Big Data!

---

- Among the three big data features, **variety** and **velocity** are even more challenging than the sheer size **volume**

## Reasons:

---

- ⊙ New domains have been constantly emerging, rapidly growing in size.
- ⊙ Domains can differ in
  - **topics** (e.g., medicine, chemistry or mechanics)
  - **genres** (e.g., news, novels, blogs, scientific publications or patents)
  - **targets** (e.g., different relations such as marriage, person-parent relation, disease-symptom relation)
  - **data internal properties** (e.g., size or redundancy or connectivity).
- ⊙ Systems, methods or strategies developed or trained for so-called general purpose or one specific domain can often not be directly taken over by other domains, because
  - each domain needs its own domain knowledge and
  - each application data has its own special properties.

# Relevant Strategies for Domain Adaptation

---

- ③ Minimally dependent on the labeled training data
  - ④ Minimally or weakly supervised machine learning methods
- ③ Strategies for
  - confidence estimation of automatically learned information and knowledge
  - filtering of irrelevant and wrong information
- ③ Domain adaptation of generic systems for specific applications

## Our solutions (1)

---

- ③ minimally supervised and distantly supervised automatic learning of domain-specific grammar-based pattern rules for n-ary RE: DARE and Web-DARE Systems
  - ③ Feiyu Xu, Hans Uszkoreit, Hong Li, "A Seed-driven Bottom-up Machine Learning Framework for Extracting Relations of Various Complexity (2007)". In ACL 2007.
  - ③ Hans Uszkoreit, Feiyu Xu, Hong Li. "Analysis and Improvement of Minimally Supervised Machine Learning for Relation Extraction". In NLDB 2009.
  - ③ Sebastian Krause, Hong Li, Hans Uszkoreit, Feiyu Xu, "Large-Scale Learning of Relation-Extraction Rules with Distant Supervision from the Web". In Proceedings of the 11th International Semantic Web Conference (ISWC 2012).

## Our solutions (2)

---

### ◎ Various filtering and confidence estimation methods for high-performance and large-scale relation extraction

◎ Sebastian Krause, Hong Li, Hans Uszkoreit, Feiyu Xu, "Large-Scale Learning of Relation-Extraction Rules with Distant Supervision from the Web". In Proceedings of the 11th International Semantic Web Conference (ISWC 2012)

◎ Andrea Moro, Hong Li, Sebastian Krause, Feiyu Xu, Roberto Navigli, Hans Uszkoreit, "Semantic rule filtering for web-scale relation extraction". In Proceeding of International Semantic Web Conference (ISWC 2013).

◎ Feiyu Xu, Hans Uszkoreit, Sebastian Krause, Hong Li. Boosting Relation Extraction with Limited Closed-World Knowledge. COLING 2009, Poster.

### ⊙ Automatic adaptation and improvement of generic parsing results for specific domains

⊙ Peter Adolphs, Feiyu Xu, Hans Uszkoreit, Hong Li, "Dependency Graphs as a Generic Interface between Parsers and Relation Extraction Rule Learning". In Proceedings of KI 2011, pp. 50-62, 2011.

⊙ Feiyu Xu, Hong Li, Yi Zhang, Hans Uszkoreit, Sebastian Krause, "Parse reranking for domain-adaptive relation extraction". Journal of Logic and Computation, doi: 10.1093/logcom/exs055, Oxford University Press, 2012.

### ⊙ Automatic generation of domain-specific linguistic knowledge resources

⊙ Hans Uszkoreit and Feiyu Xu, "From Strings to Things, SAR-Graphs: A New Type of Resource for Connecting Knowledge and Language". In Proceedings of 1st International Workshop on NLP and DBpedia volume 1064, Sydney, NSW, Australia, CEUR Workshop Proceedings, 10/2013

⊙ Open source: [sargraph.dfki.de](http://sargraph.dfki.de)

---

# Web-DARE

## Distant-supervised Web-scale RE



## Web-DARE: Distant Supervision based RE

---

- ③ Large number of RE rules are automatically learned by using Freebase as seed knowledge and Web as training corpus
- ③ Goal:
  - covering most linguistic variants for expressing a relation
  - thus solving the notorious long-tail problem of real-world NL applications

## Data Set

---

- ◎ rules learned for 39 relations
  - ◎ n-ary relations  $n \geq 2$
- ◎ three domains: business, awards and people
- ◎ 2.8 million relation instances retrieved from Freebase as seed
- ◎ 20 million web documents as training corpus

# Example in Nobel Prize Award Domain

---

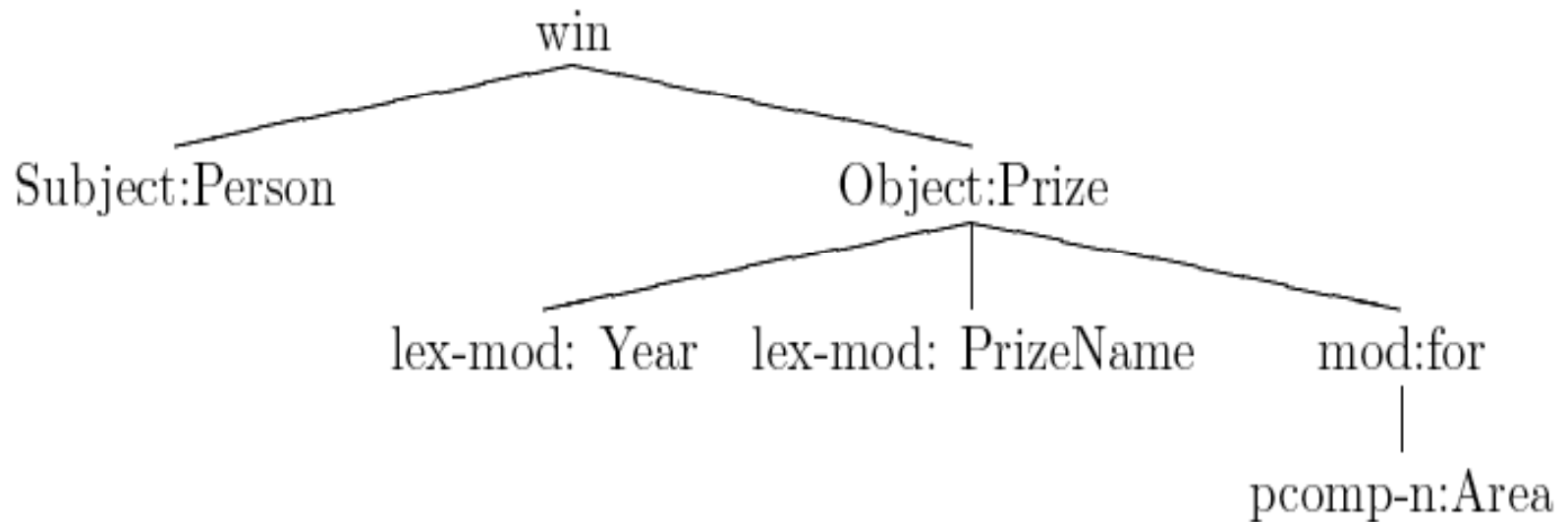
- Seed example

*<Mohamed ElBaradei/Person, Nobel/Prize, Peace/Area, 2005/Year>*

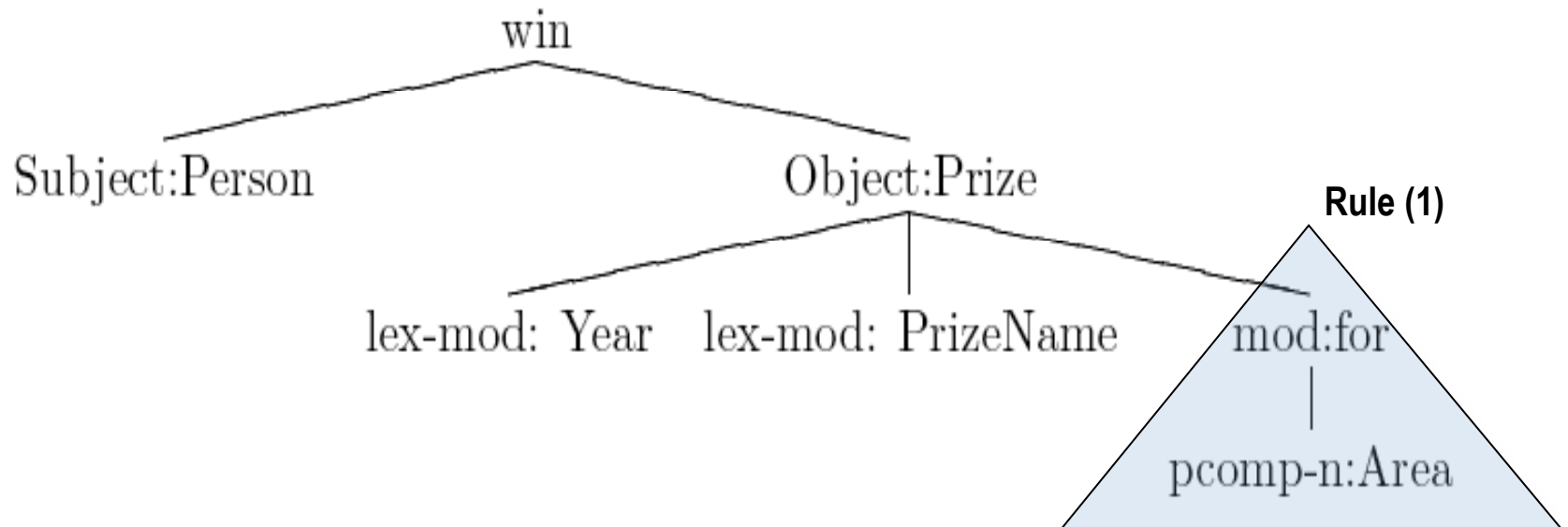
- Sentence matched with the seed

Mohamed ElBaradei won the 2005 Nobel Prize for Peace on Friday ...

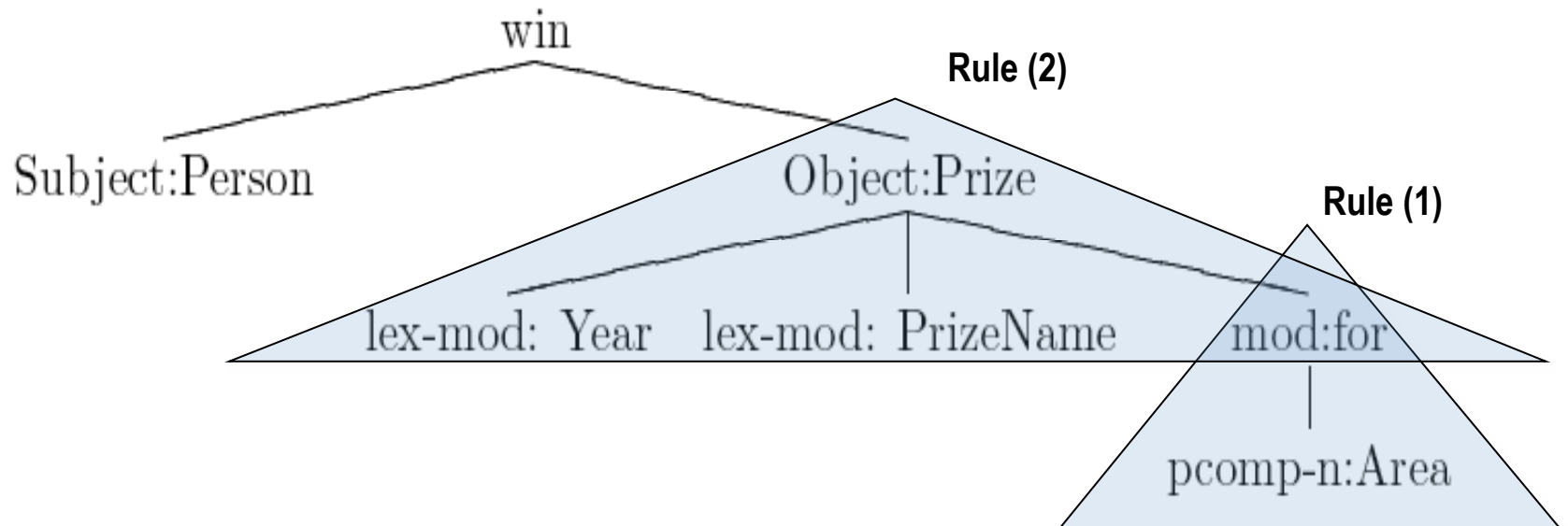
# Dependency Parse Result



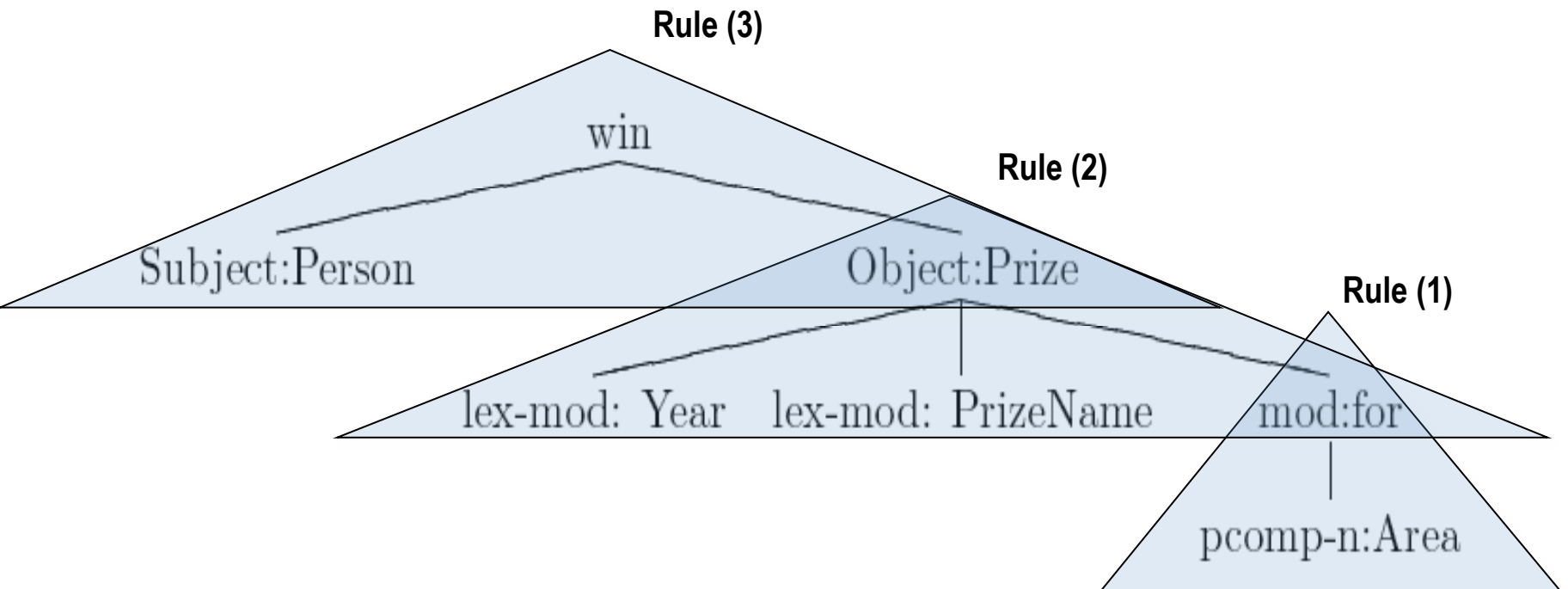
# Bottom Up Rule Learning



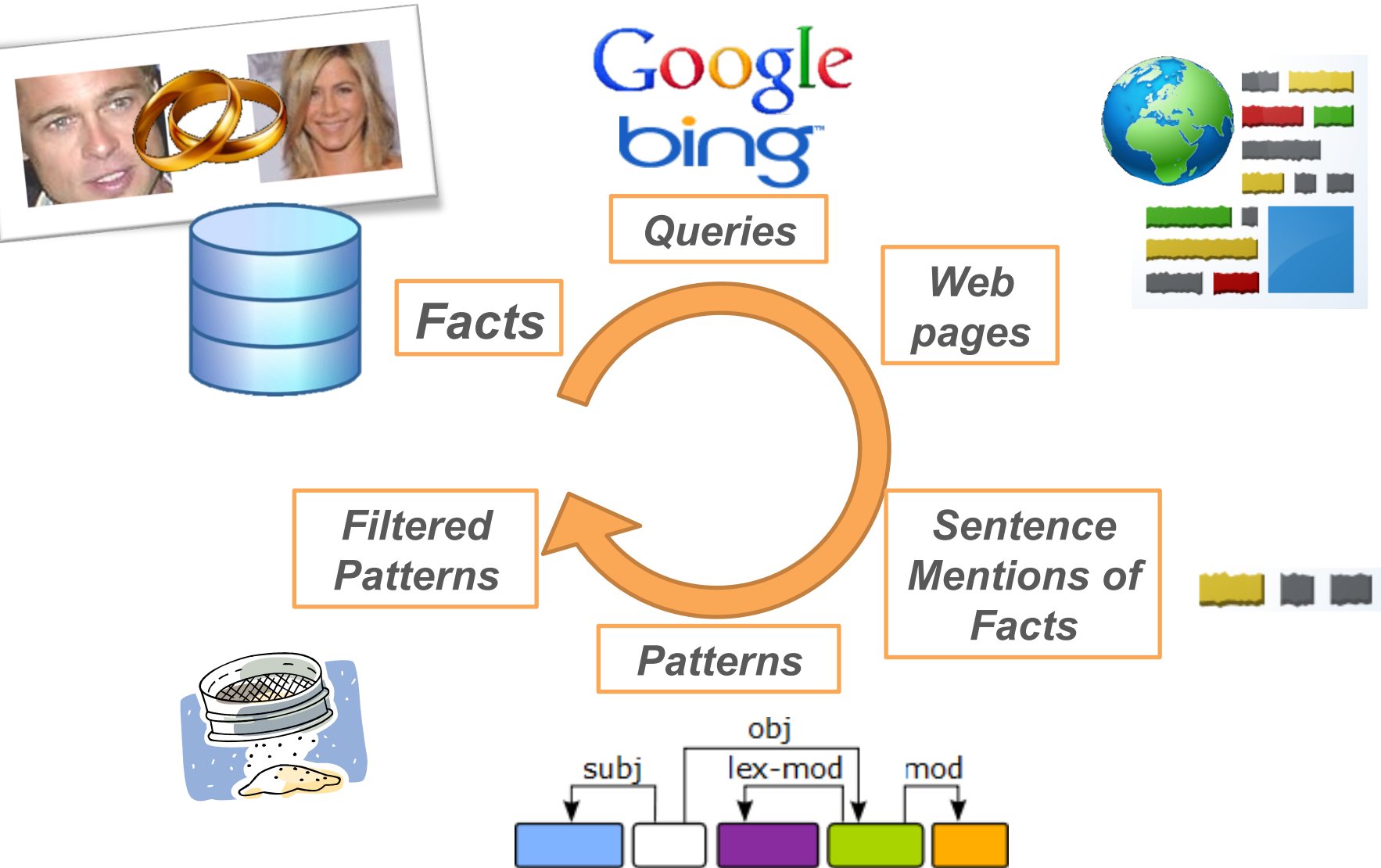
# Bottom Up Rule Learning



# Bottom Up Rule Learning



# Web-DARE Architecture



# Some Statistics of Web-DARE Rules

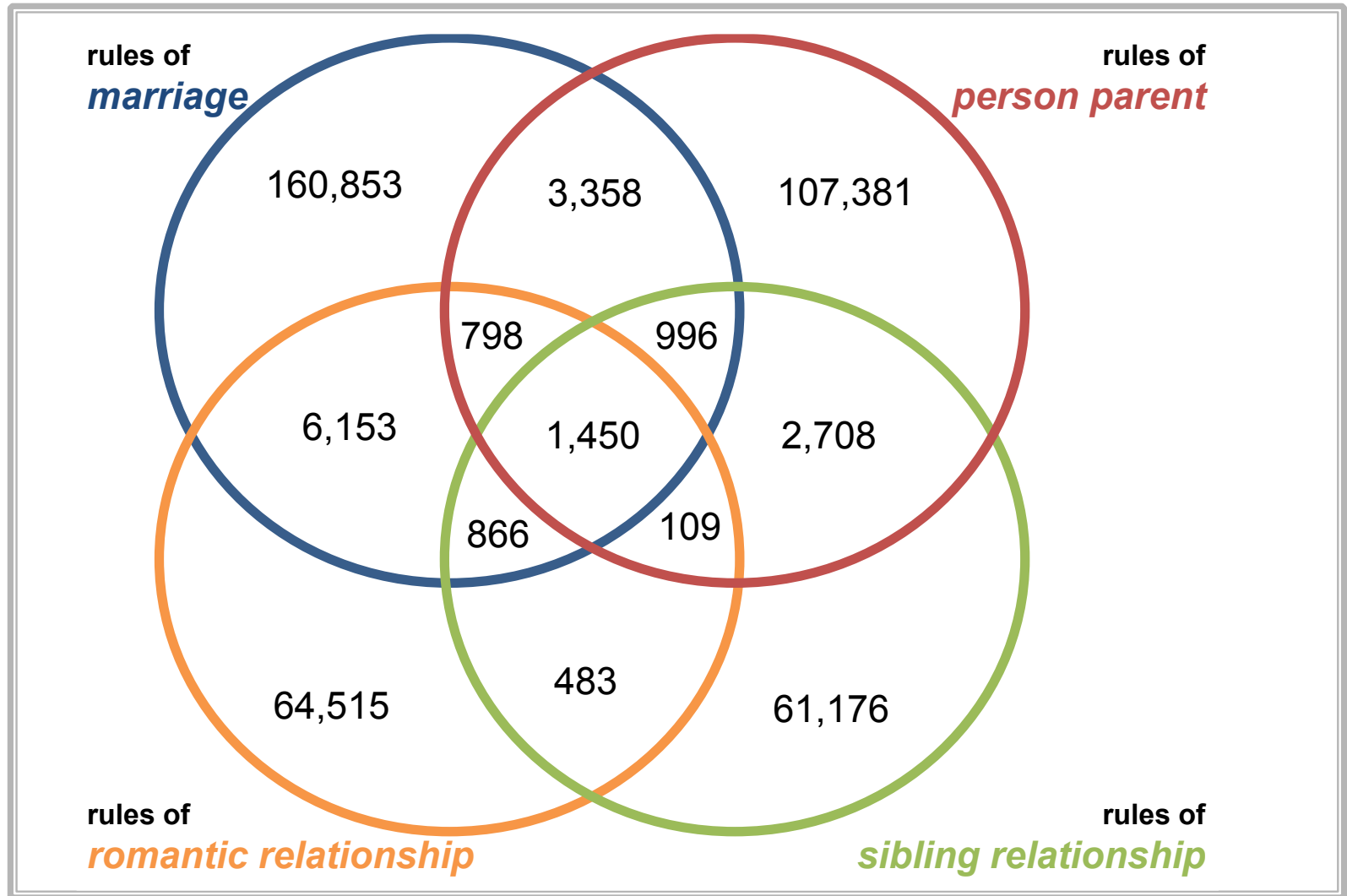
Relation	# Sentences used	# Sentences w/ a learned rule	# Rules	# Rules w/o duplicates
<i>award nomination</i>	13,966	13,149	23,987	7,800
<i>award honor</i>	50,550	49,001	106,550	40,578
<i>hall of fame induction</i>	31,244	28,278	44,920	17,450
<i>organization relationship</i>	46,331	42,824	60,379	28,903
<i>acquisition</i>	63,967	60,903	96,747	50,544
<i>organization merger</i>	2,996	1,521	3,243	1,758
<i>company name change</i>	9,433	9,132	15,619	6,910
<i>spin off</i>	5,247	5,094	8,319	4,798
<i>marriage</i>	342,895	335,313	557,478	176,949
<i>sibling relationship</i>	167,611	160,893	255,788	69,596
<i>romantic relationship</i>	155,335	152,878	229,393	74,895
<i>person parent</i>	192,610	186,834	390,878	119,238
<b>average of 39 relations</b>	66,545	66,509	109,435	41,620

## Problems of Large-Scale Approach

---

- ◎ Very low precision
  - ◎ a lot of noisy rules
  - ◎ many rules are learned from more than one relation

# Euler Diagram for Four People-Relations



---

# Various Filtering Strategies for High-Performance Web-Scale RE



# Frequency-Driven Rule Filters

◎ Merged Filter:

$$valid_m^{\mathcal{R}}(r) = valid_{freq}^{\mathcal{R}}(r) \wedge valid_{inter}^{\mathcal{R}}(r)$$

- 1) **absolute frequency filtering:** a threshold to exclude rules with low occurrence

# Rule Frequency Driven Filters

## ⊙ Merged Filter:

$$valid_m^{\mathcal{R}}(r) = valid_{freq}^{\mathcal{R}}(r) \wedge valid_{inter}^{\mathcal{R}}(r)$$

- 1) **absolute frequency filtering:** a threshold to exclude rules with low occurrence
- 2) **inter-relation filter (Overlap Filter – FO Filter):**
  - based on mutual exclusiveness of relations with similar entity-type signatures.
  - a rule is only valid for a relation, if its relative frequency is higher than any other relations with similar entity type signatures.

$$valid_{inter}^{\mathcal{R}}(r) = \begin{cases} true & \text{if } \forall \mathcal{R}' \in \mathbb{R} \setminus \{\mathcal{R}\} : rf_{r, \mathcal{R}} > rf_{r, \mathcal{R}'} \\ false & \text{otherwise} \end{cases}$$

# Weakness of Filtering with Frequency

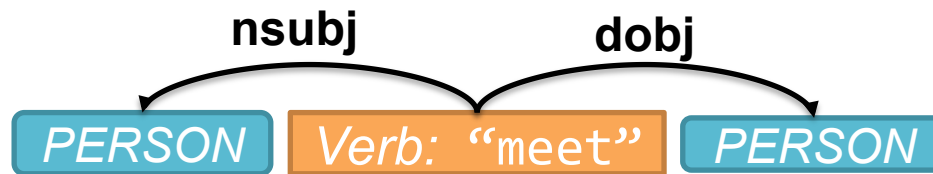
---

- ◎ Undetected low-quality patterns:
  - high frequency in target relation, low frequency in coupled relations

# Weakness of Filtering with Rule Frequency

---

- ⊙ Undetected low-quality patterns:
  - ⊙ high frequency in target relation, low frequency in coupled relations



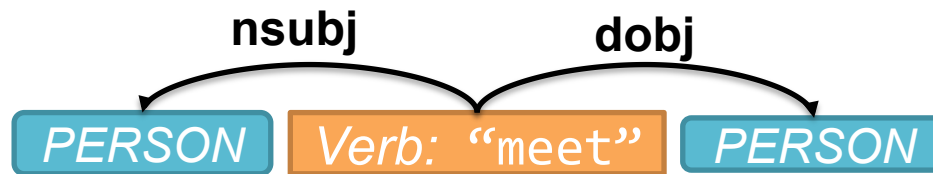
X ?



✓ ?

# Weakness of Filtering with Rule Frequency

- ⊙ Undetected low-quality patterns:
  - ⊙ high frequency in target relation, low frequency in coupled relations



X ?

- ⊙ Erroneously-deleted good patterns:
  - ⊙ infrequent patterns



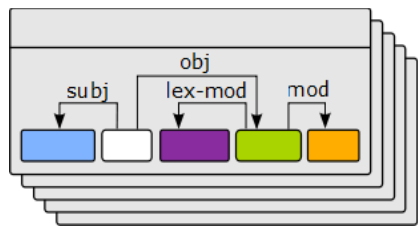
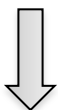
✓ ?



# Lexical Semantics can help!



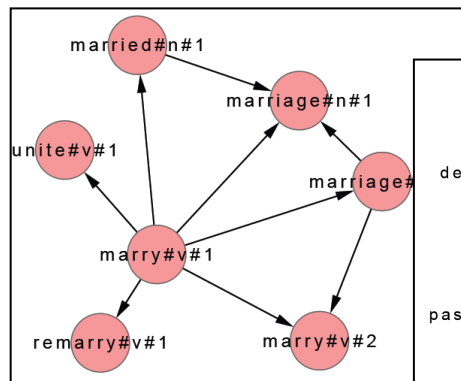
World Wide Web



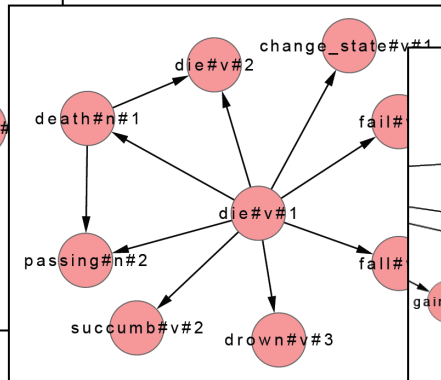
✓ ? ✗ ?

Candidate RE Patterns

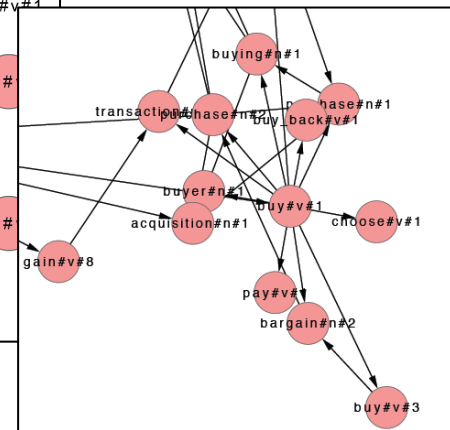
marriage



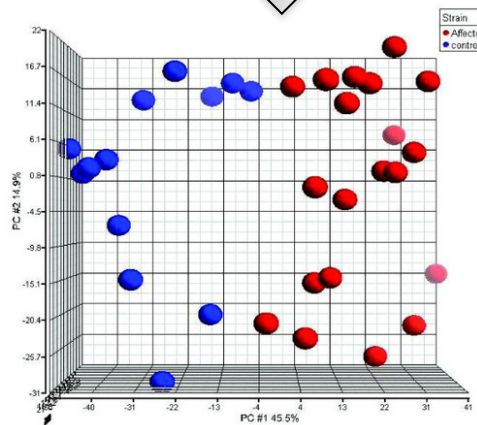
person-death



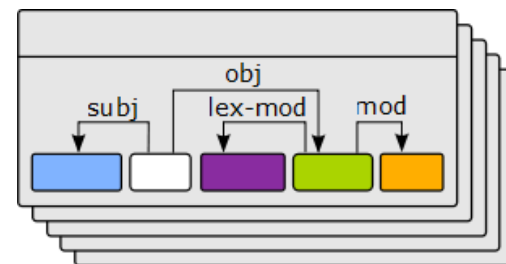
acquisition



Relation-specific lexical semantic graphs



Unsupervised Classification

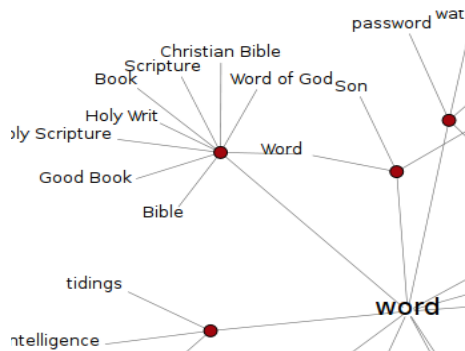


✓ !!!

High-quality RE Patterns

# Automatic learning of relation-specific lexical semantic network

Generic Lexical Semantic Network (BabelNet)

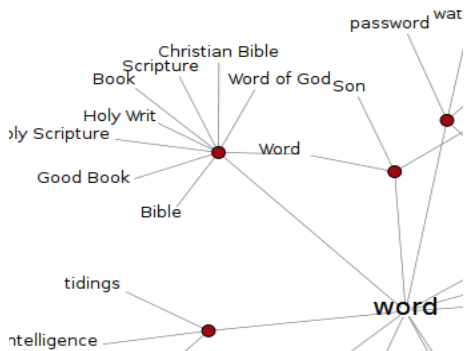


automatically learned unfiltered RE rules and their mentions



# Automatic learning of relation-specific lexical semantic network

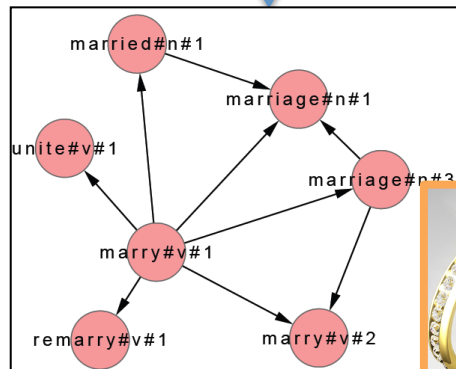
Generic Lexical Semantic Network (BabelNet)



automatically learned unfiltered RE rules and their mentions

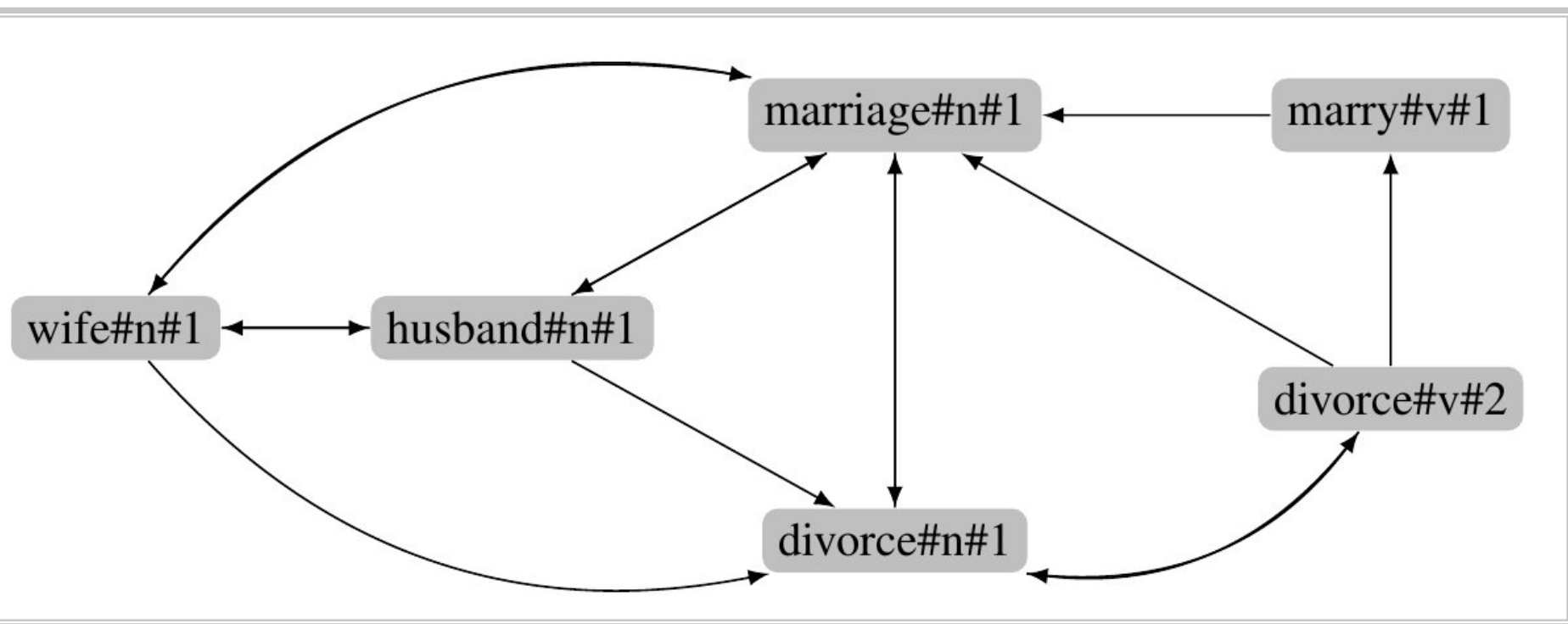


Word Sense Disambiguation

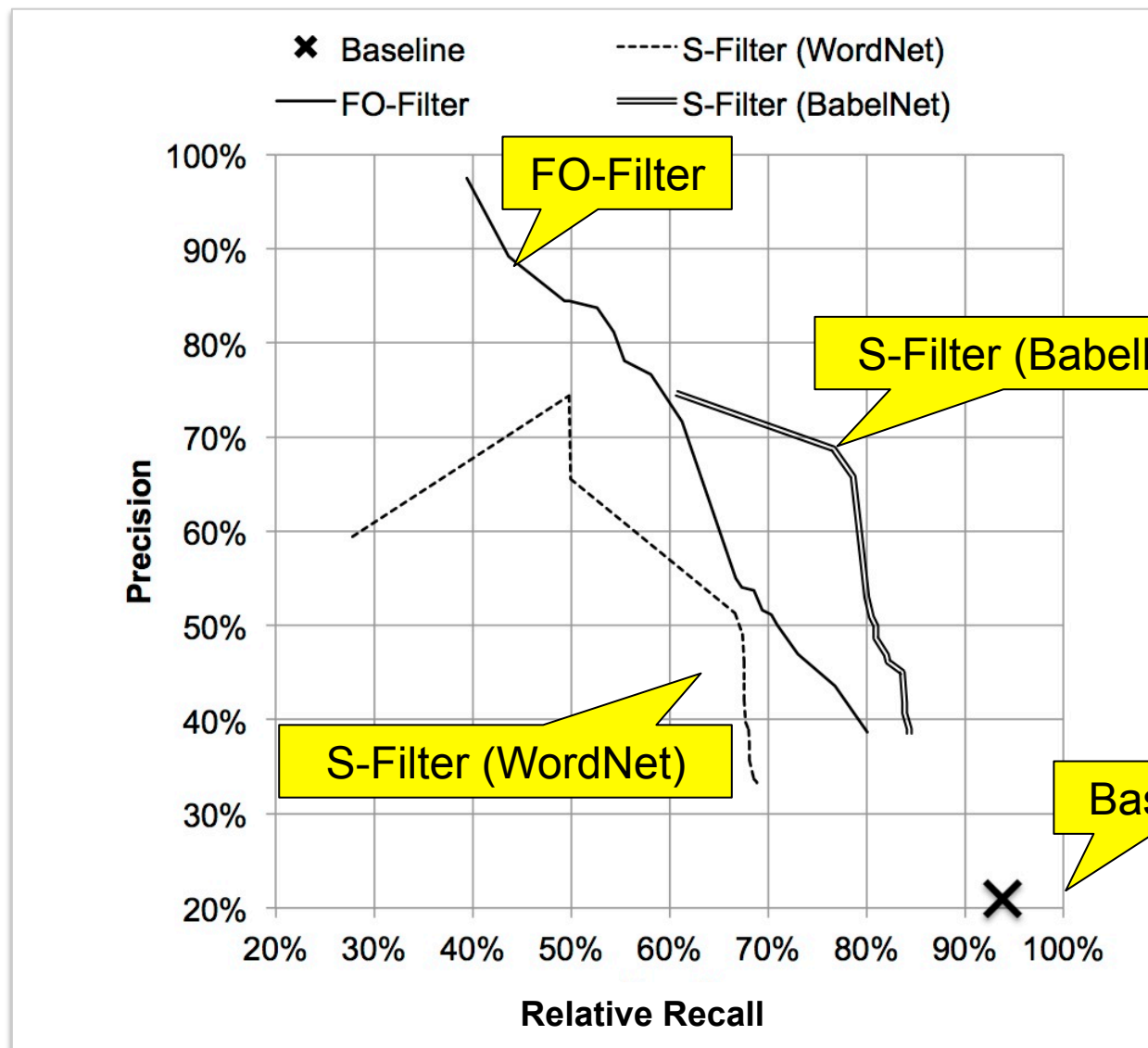


# The Relation-Specific Semantic Graph

An excerpt of the semantic graph for the relation *marriage*



# Extrinsic Eval. - Web-DARE



---

# Parse-Reranking for Domain-adaptive RE



## Error Types of Extracted Wrong Instances

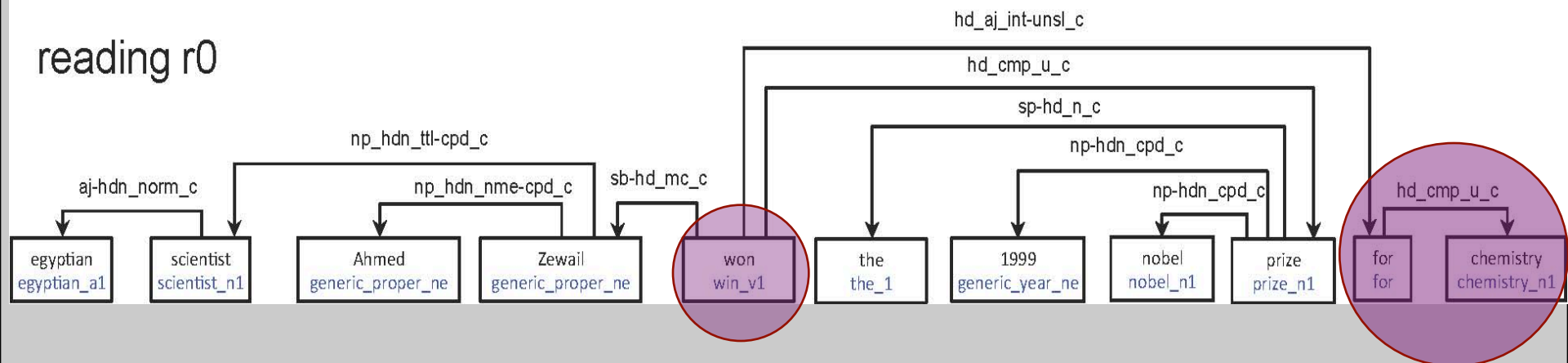
---

Content	Modality	Named Entity Recognition (NER)	Parsing	NER & Parsing	DARE Rules
11.8%	17.6%	5.9%	38.2%	11.8%	14.7%

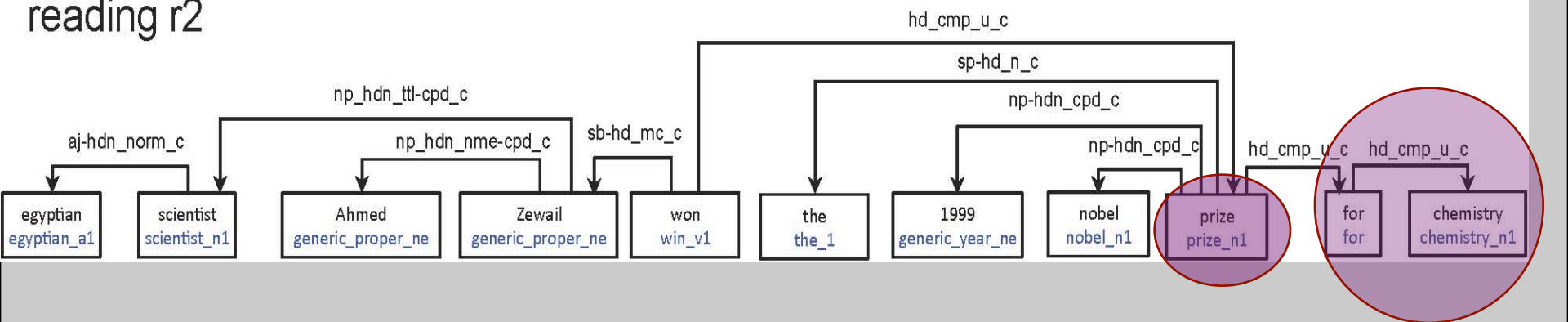
# HPSG Parses: PP Attachment

*Egyptian scientist Ahmed Zewail won the 1999 Nobel Prize for chemistry*

reading r0



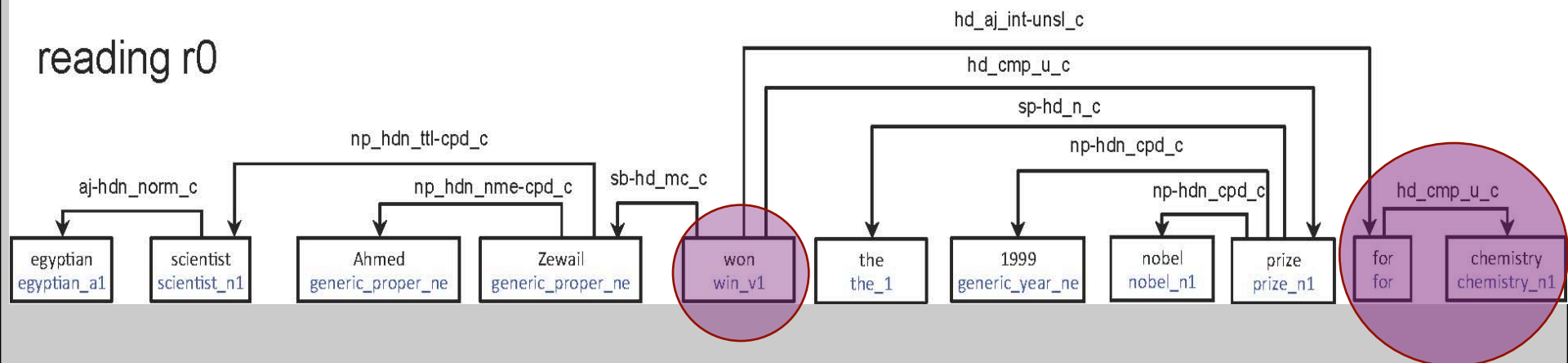
reading r2



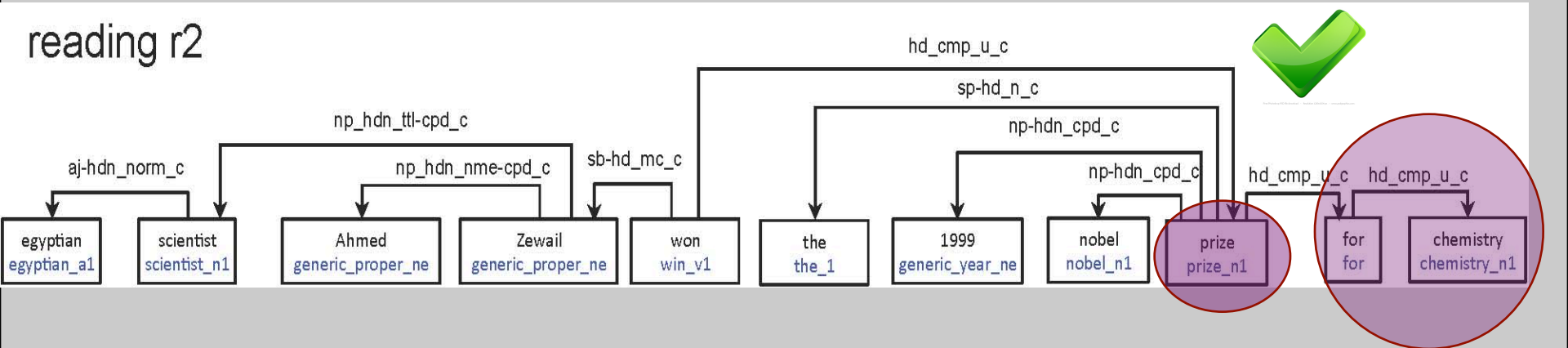
# HPSG Parses: PP Attachment

*Egyptian scientist Ahmed Zewail won the 1999 Nobel Prize for chemistry*

reading r0

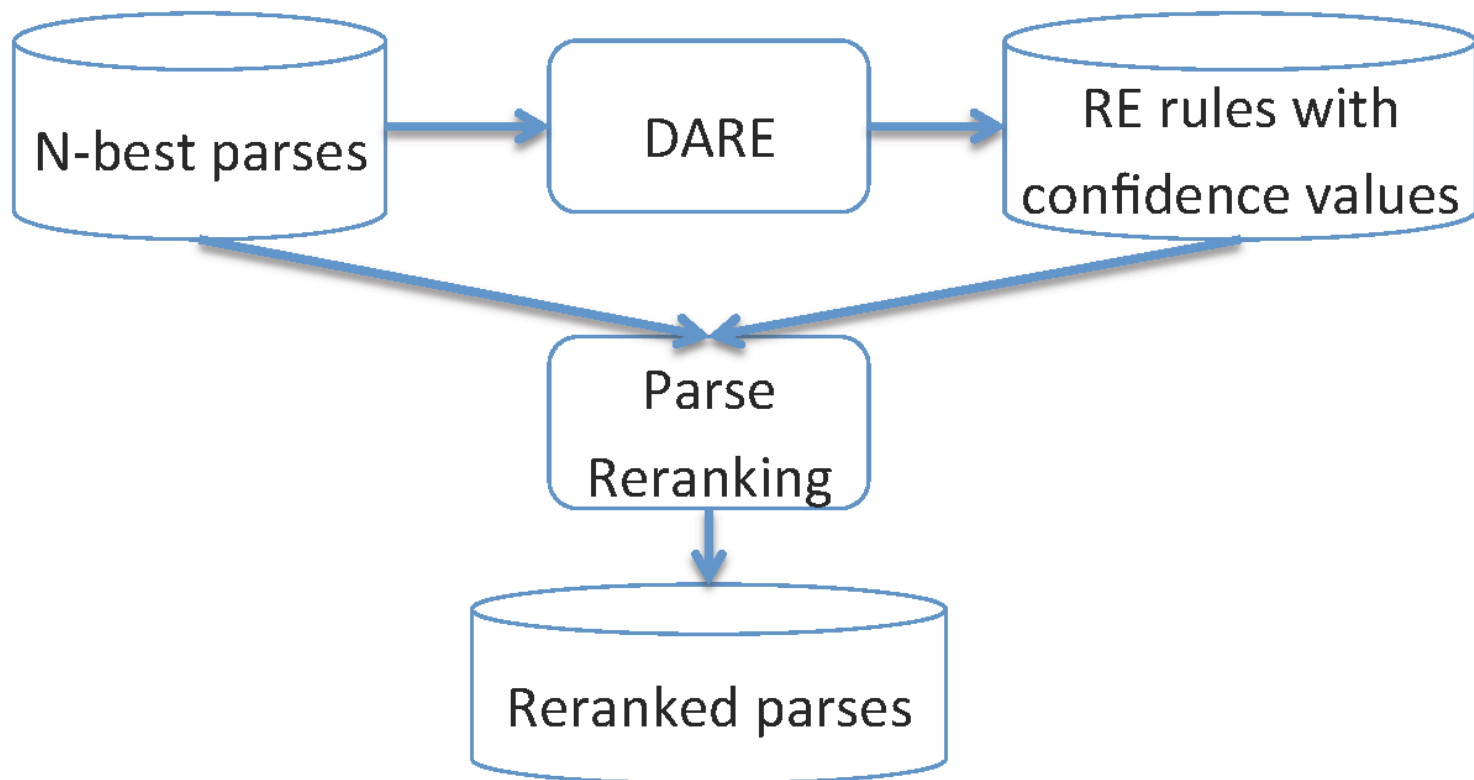


reading r2



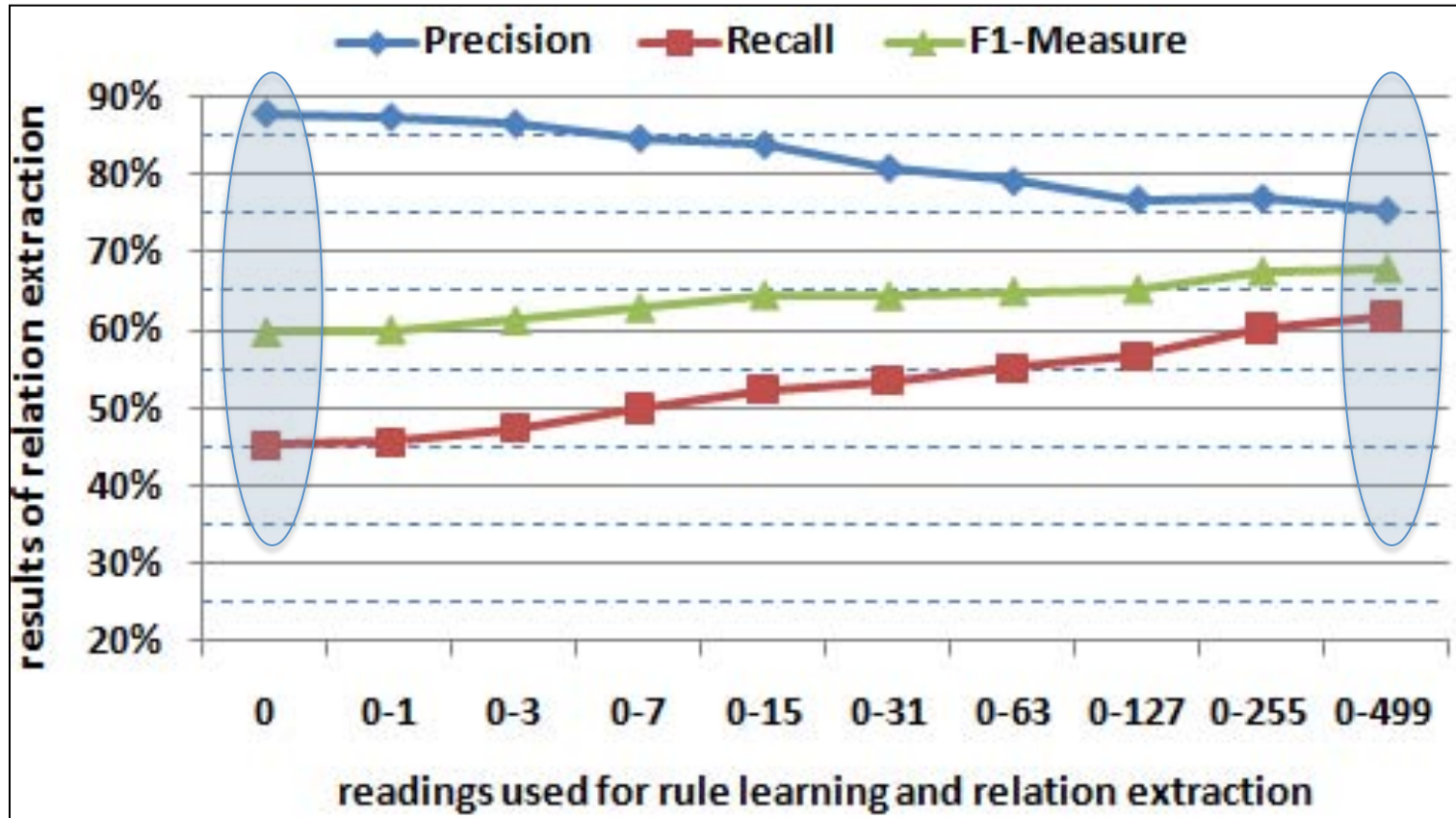
# Reranking Architecture

---



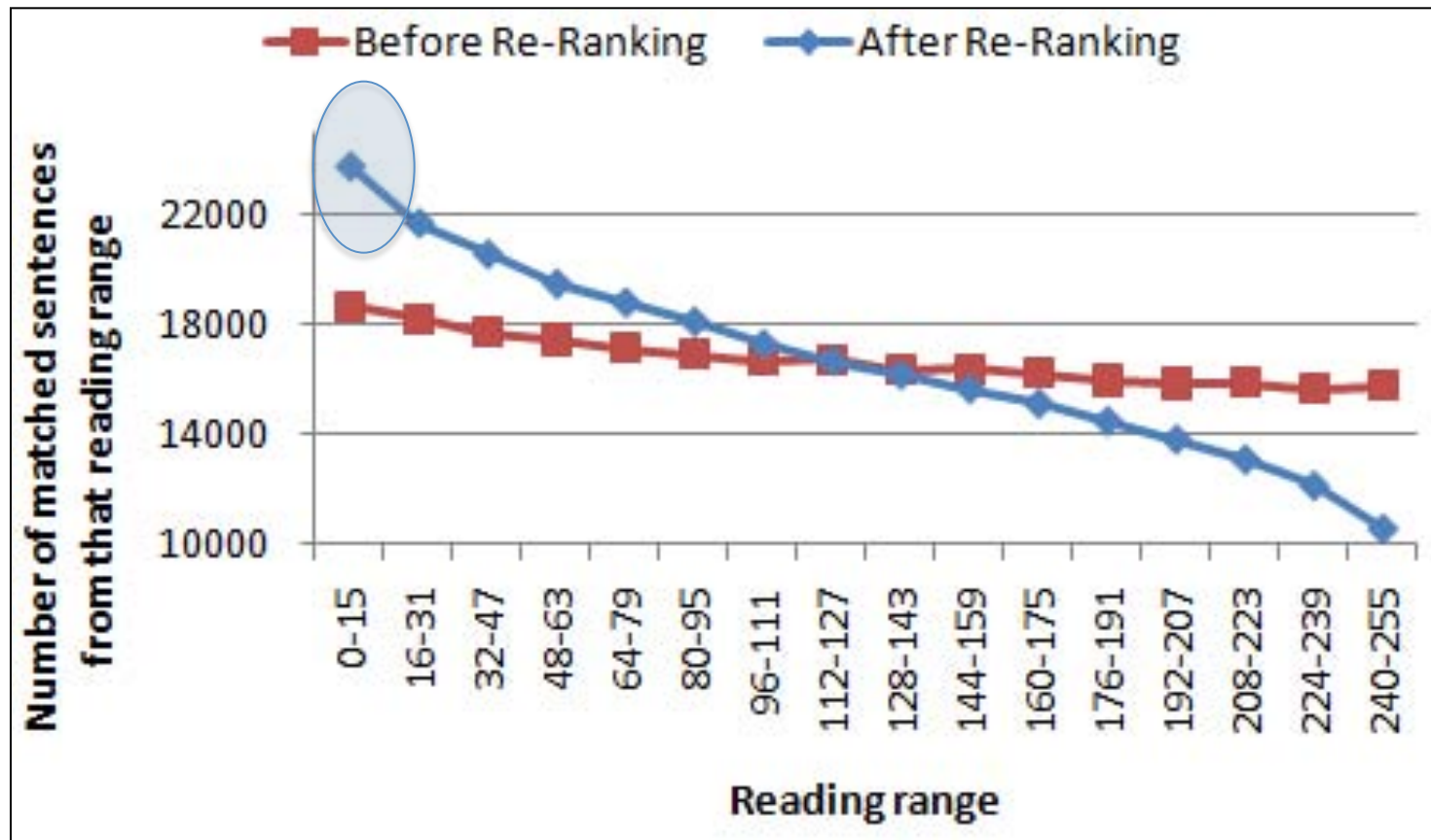
## Baseline: before Re-ranking

- ◎ Best reading: high precision, low recall, low F-measure
- **500 readings**: lower precision, higher recall, higher F-measure



## After Re-Ranking:

- ⊙ Re-ranked top readings match more sentence mentions containing RE instances
- ⊙ Improvements of Recall and F-Measure



## Conclusion

---

- ◎ The performance of large-scale RE for each application is dependent on the performance of domain-adaptation methods
- Three original contributions (among others):
  - Extension of relation extraction to n-ary relations
  - Semantic filtering with large lexical knowledge bases
  - Parser improvement for the specific RE task by reranking
- For our work we received a Google Focused Research Award



Google  
Focused Research Awards