

Relation Extraction
and
Machine Learning for IE

Feiyu Xu

feiyu@dfki.de

Language Technology-Lab
DFKI, Saarbrücken

Relation in IE

Information Extraction is ...

a technology that is futuristic from the user's point of view in the current information-driven world.

Rather than indicating which documents need to be read by a user, it extracts pieces of information that are salient to the user's needs ...

provided by NIST:

[http://www-nlpir.nist.gov/related_projects/muc/]

Information Extraction: A Pragmatic Approach

- Identify the types of entities that are relevant to a particular task
- Identify the range of facts that one is interested in for those entities
- Ignore everything else

IE from Research Papers

A Critical Evaluation of Commensurable Abduction Models for Semantic Interpretation (1990) (Correct) (5 citations)

Peter Norvig Robert Wilensky University of California, Berkeley Computer...
Thirteenth International Conference on Computational Linguistics, Volume 3

Download: norvig.com/coling.ps
Cached: [PS.gz](#) [PS](#) [PDF](#) [DjVu](#) [Image](#) [Update](#) [Help](#)

From: norvig.com/resume (more)
Home: [R.Wilensky](#) [HPSearch](#) (Correct)

NEC ResearchIndex [Bookmark](#) [Context](#) [Related](#)

[\(Enter summary\)](#) Rate this article: 1 2 3 4 5 (best)
[Comment on this article](#)

Abstract: this paper we critically evaluate three recent abductive interpretation models, those of Charniak and Goldman (1989); Hobbs, Stickel, Martin and Edwards (1988); and Ng and Mooney (1990). These three models add the important property of commensurability: all types of evidence are represented in a common currency that can be compared and combined. While commensurability is a desirable property, and there is a clear need for a way to compare alternate explanations, it appears that a single scalar measure is not enough to account for all types of processing. We present other problems for the abductive approach, and some tentative solutions. ([Update](#))

Context of citations to this paper: [More](#)

.... (break slight modification of the one given in [Ng and Mooney, 1990] The new definition remedies the anomaly reported in [Norvig and Wilensky, 1990] of occasionally preferring spurious interpretations of greater depths. Table 1: Empirical Results Comparing Coherence and...

.... costs as probabilities, specifically within the context of using abduction for text interpretation, are discussed in [Norvig and Wilensky \(1990\)](#). The use of abduction in disambiguation is discussed in Kay et al. 1990) We will assume the following: 13) a. Only literals...

Cited by: [More](#)

[Translation Mismatch in a Hybrid MT System - Gawron \(1999\) \(Correct\)](#)
[Abduction and Mismatch in Machine Translation - Gawron \(1999\) \(Correct\)](#)
[Interpretation as Abduction - Hobbs, Stickel, Appelt, Martin \(1990\) \(Correct\)](#)

Active bibliography (related documents): [More](#) [All](#)

0.1: [Critiquing: Effective Decision Support in Time-Critical Domains - Gertner \(1995\) \(Correct\)](#)
0.1: [Decision Analytic Networks in Artificial Intelligence - Matzkevich, Abramson \(1995\) \(Correct\)](#)
0.1: [A Probabilistic Network of Resolutions - DeRose, Liu \(1992\) \(Correct\)](#)

Extracting Job Openings from the Web: Semi-Structured Data

foodscience.com-Job2

JobTitle: Ice Cream Guru
Employer: foodscience.com
JobCategory: Travel/Hospitality
JobFunction: Food Services
JobLocation: Upper Midwest
Contact Phone: 800-488-2611
DateExtracted: January 8, 2001
Source: www.foodscience.com/jobs_midwest.html
OtherCompanyJobs: foodscience.com-Job1

Ice Cream Guru

If you dream of cold creamy chocolate or gooey gooey cookie, there's a great opportunity for you to maintain and expand this major corporation's high-end ice cream brand. Will be based in the Upper Midwest for about a year. After that, California here I come! Requires a BS in Food Science or dairy, plus ice cream formulation experience. Will consider entry level with an MS and an internship.
Contact Susan: e-mail
1-800-488-2611

On the Notion *Relation Extraction*

Relation Extraction is the cover term for those Information Extraction tasks in which instances of semantic relations are detected in natural language texts.

Types of Information Extraction in LT

- Topic Extraction
- Term Extraction
- Named Entity Extraction
- Binary Relation Extraction
- N-ary Relation Extraction
- Event Extraction
- Answer Extraction
- Opinion Extraction
- Sentiment Extraction

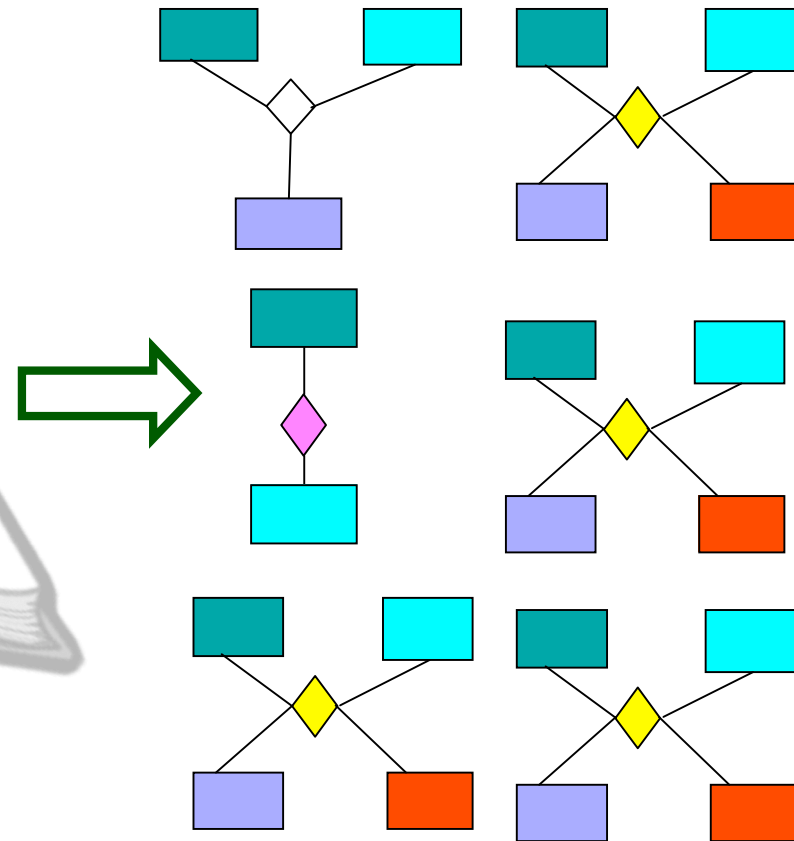
Types of Information Extraction in LT

- Topic Extraction
- Term Extraction
- Named Entity Extraction
- **Binary Relation Extraction**
- **N-ary Relation Extraction**
- **Event Extraction**
- **Answer Extraction**
- **Opinion Extraction**
- **Sentiment Extraction**

Types of Relation Extraction



Relation Extraction is a demanding sub-area of Information Extraction



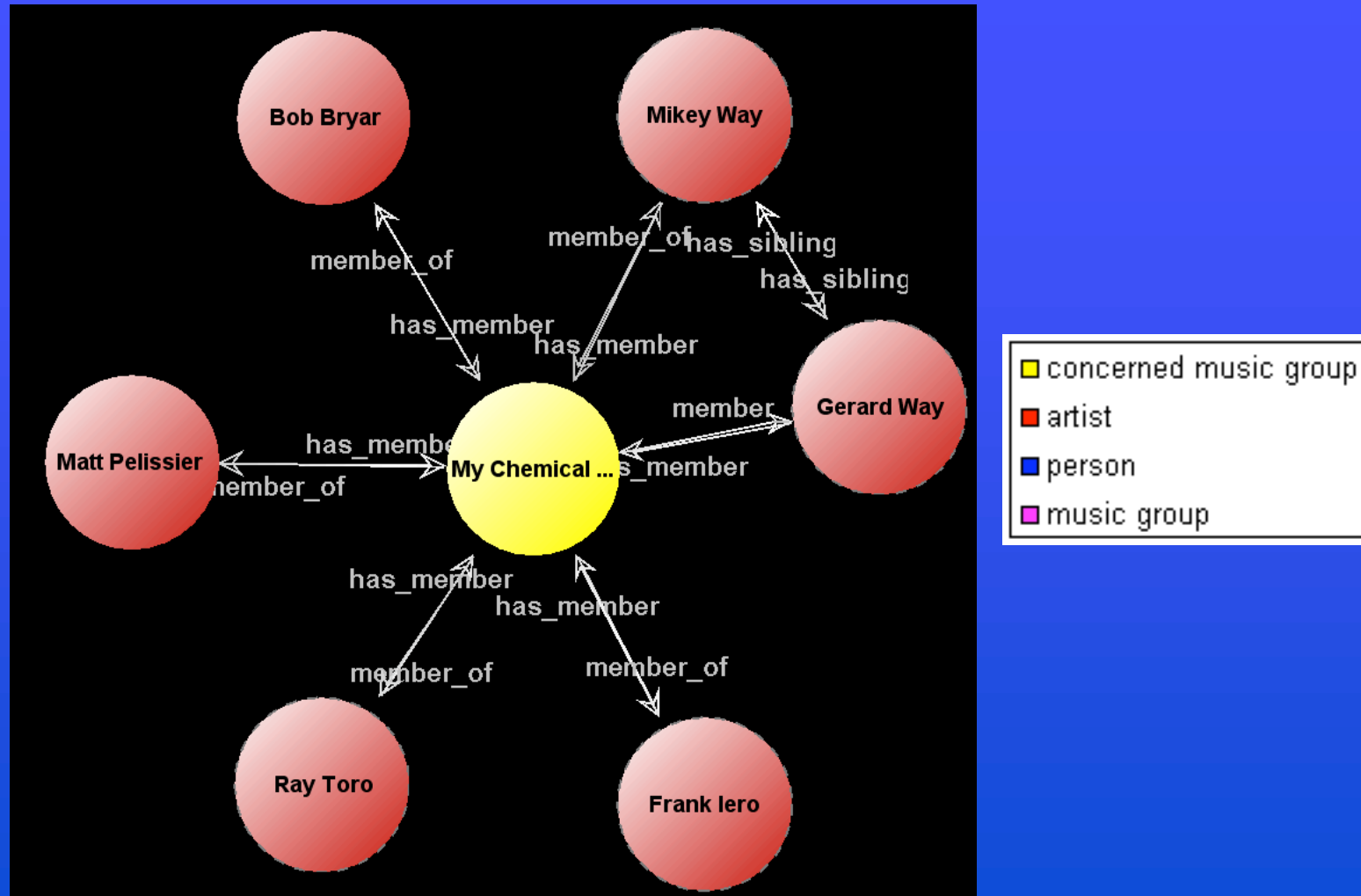
Example of Binary Social Relations

Social Network of "Madonna" (Depth = 1)



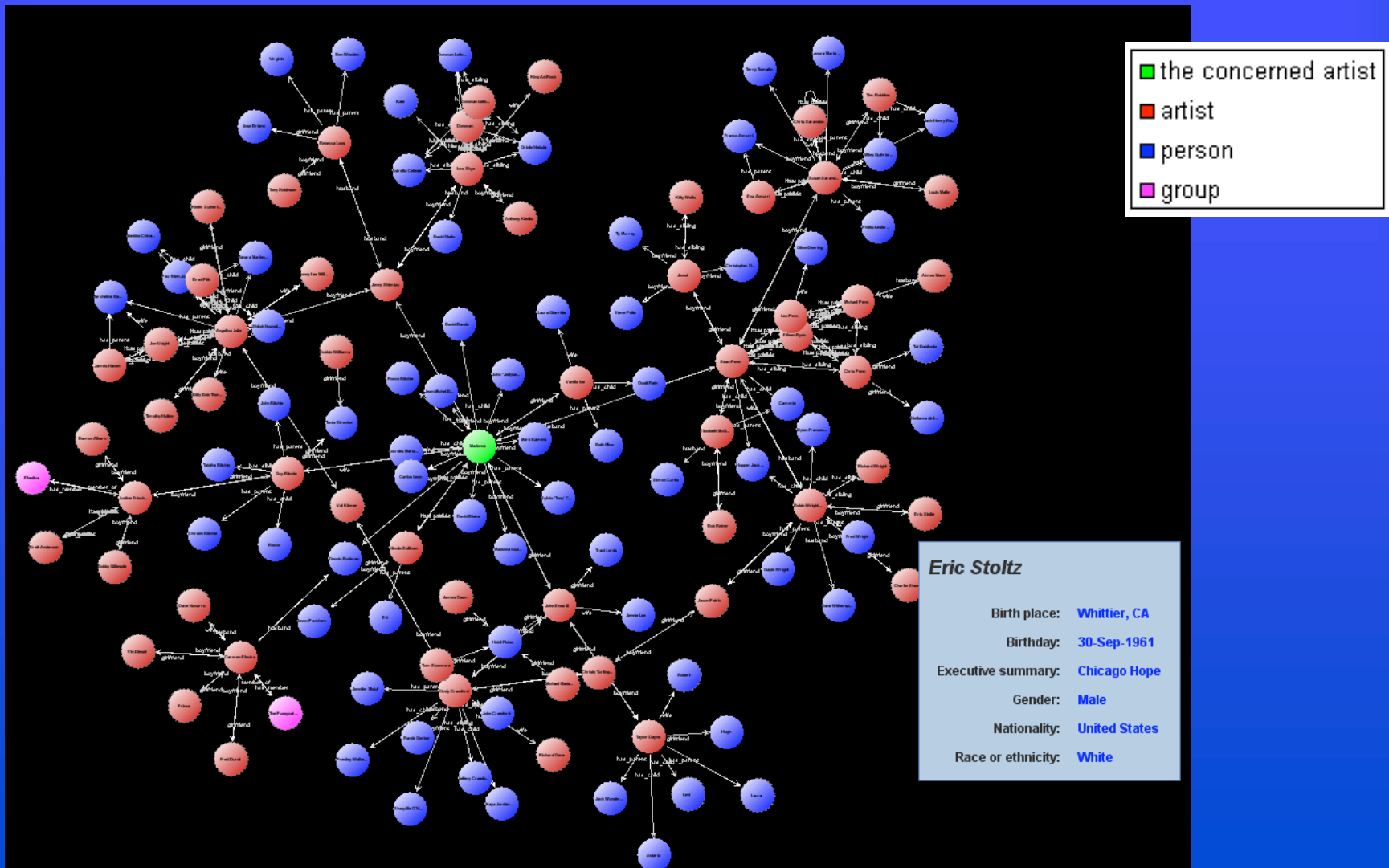
Examples of Binary Relations

Social Network of "My Chemical Romance" (Depth = 1)



Examples

Social Network of “Madonna” (Depth = 3)



Relation about Person, Title and Organization

October 14, 2002, 4:00 a.m. PT

For years, **Microsoft Corporation CEO Bill Gates** railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, **Microsoft** claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. **Gates** himself says **Microsoft** will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said **Bill Veghte**, a **Microsoft VP**. "That's a super-important shift for us in terms of code access."

Richard Stallman, **founder** of the **Free Software Foundation**, countered saying...

*	Microsoft Corporation CEO Bill Gates	}
*	Microsoft Gates Microsoft	
*	Bill Veghte Microsoft VP	}
*	Richard Stallman founder Free Software Foundation	

NAME	TITLE	ORGANIZATION
Bill Gates	CEO	Microsoft
Bill Veghte	VP	Microsoft
RichardStallman	founder	Free Soft..

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.

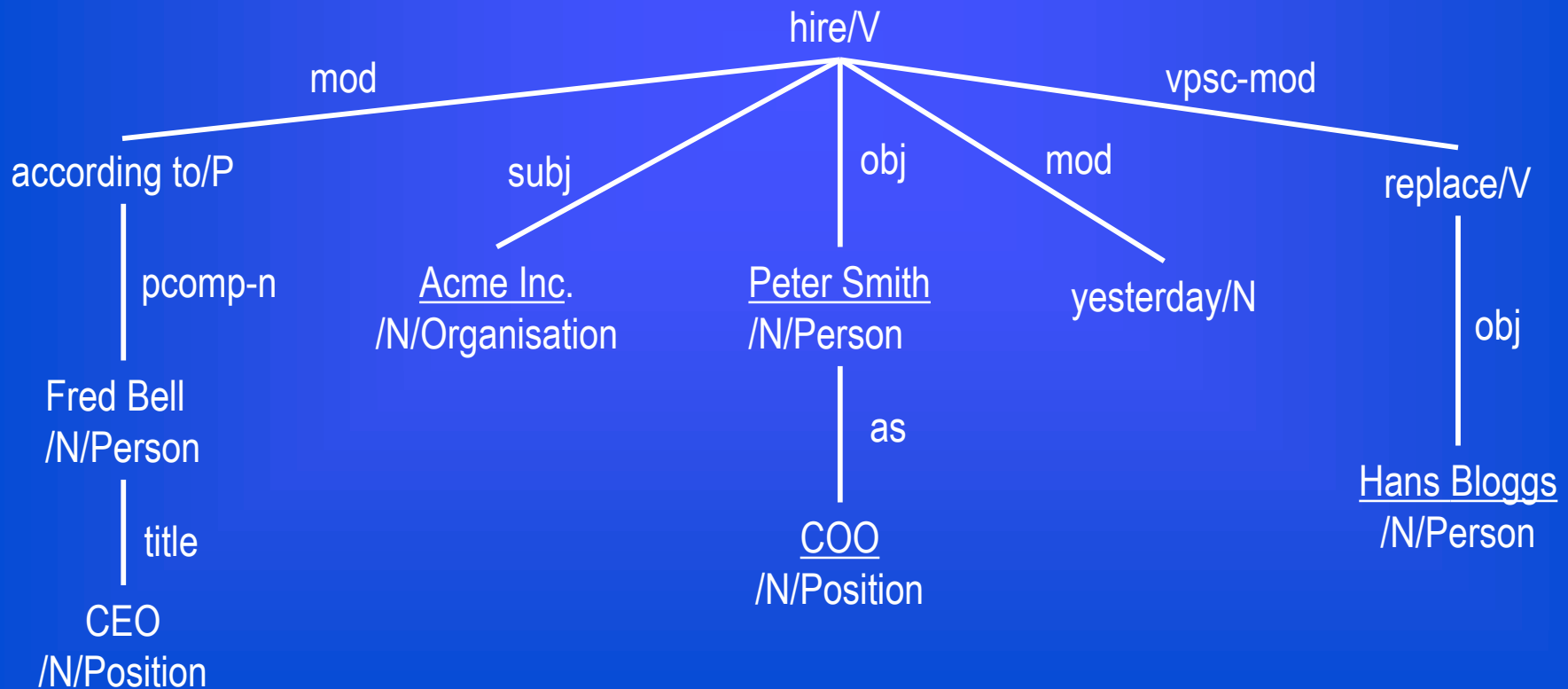
<person_in, person_out, position, organisation>

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

<person_in, person_out, position, organisation>

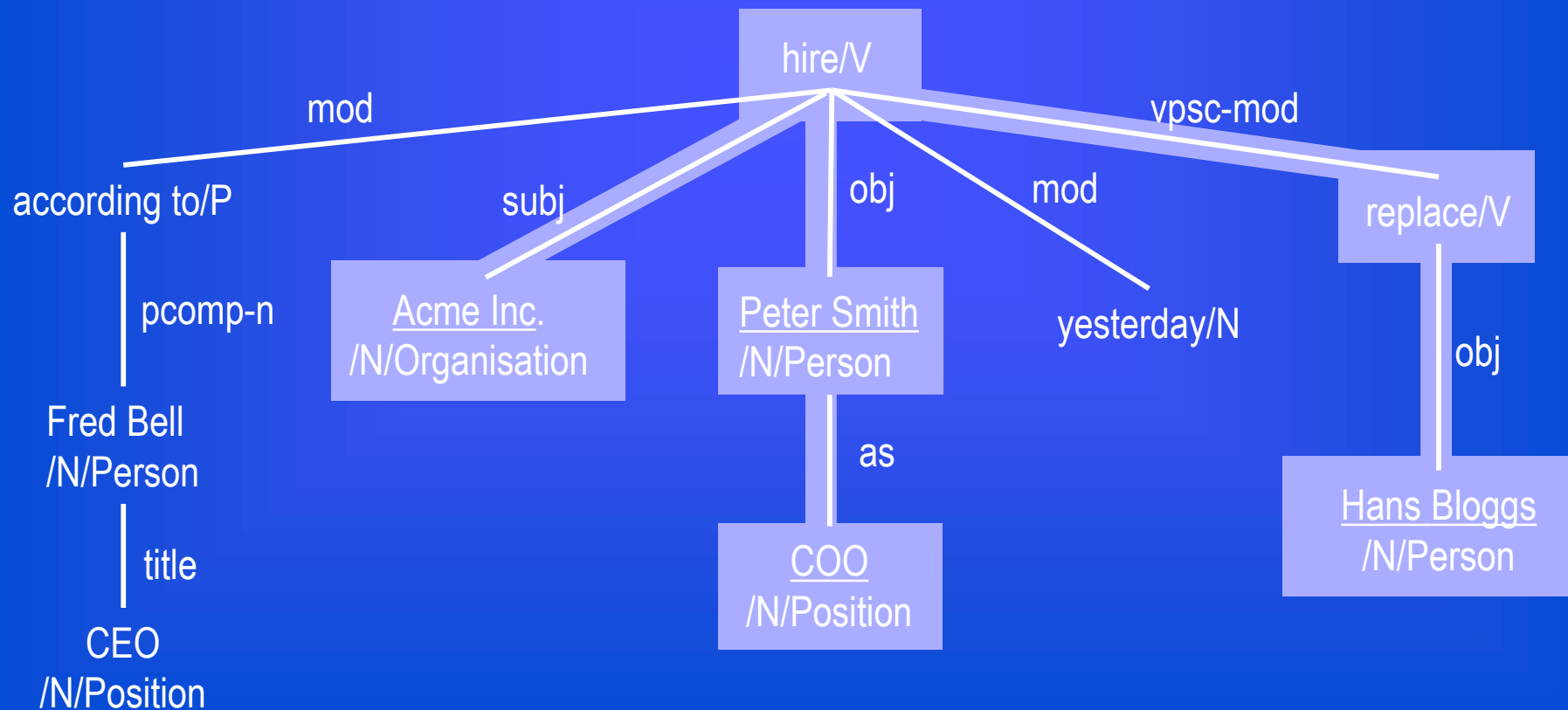
According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<person_in, person_out, position, organisation>

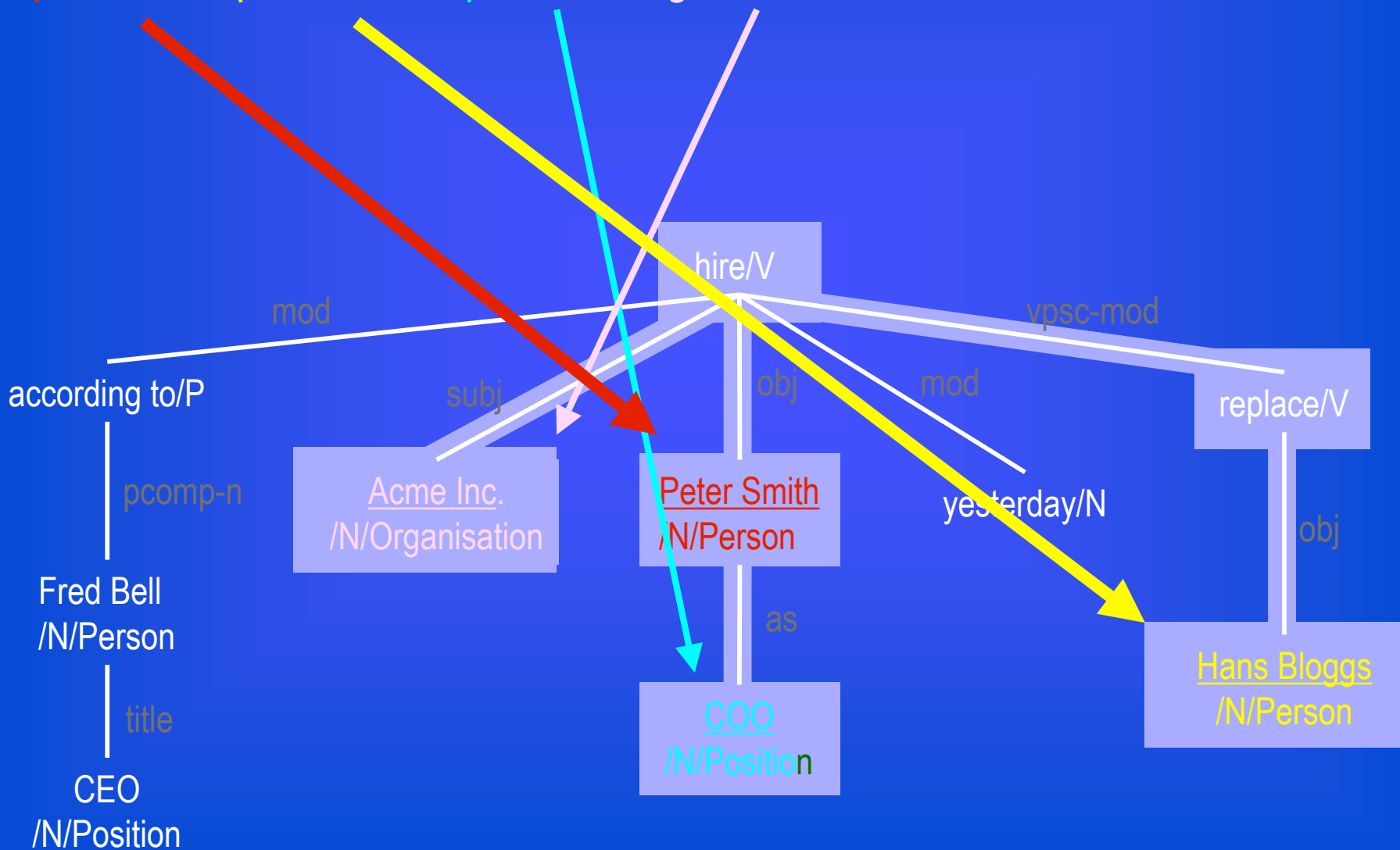


According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<person_in, person_out, position, organisation>



<person_in, person_out, position, organisation>



A Brief History of IE

Message Understanding Conferences

[MUC-7 98]

- U.S. Government sponsored conferences with the intention to coordinate multiple research groups seeking to improve IE and IR technologies (since 1987)
- defined several generic types of information extraction tasks (MUC Competition)
- MUC 1-2 focused on automated analysis of military messages containing textual information
- MUC 3-7 focused on information extraction from newswire articles
 - terrorist events
 - international joint-ventures
 - management succession event

Evaluation of IE systems in MUC

- Participants receive description of the scenario along with the annotated *training corpus* in order to adapt their systems to the new scenario (1 to 6 months)
- Participants receive new set of documents (*test corpus*) and use their systems to extract information from these documents and return the results to the conference organizer
- The results are compared to the manually filled set of templates (*answer key*)

Evaluation of IE systems in MUC

- precision and recall measures were adopted from the information retrieval research community

$$recall = \frac{N_{correct}}{N_{key}} \quad precision = \frac{N_{correct}}{N_{correct} + N_{incorrect}}$$

$$F = \frac{(\beta^2 + 1) \times precision \times recall}{\beta^2 \times precision + recall}$$

- Sometimes an F -measure is used as a combined recall-precision score

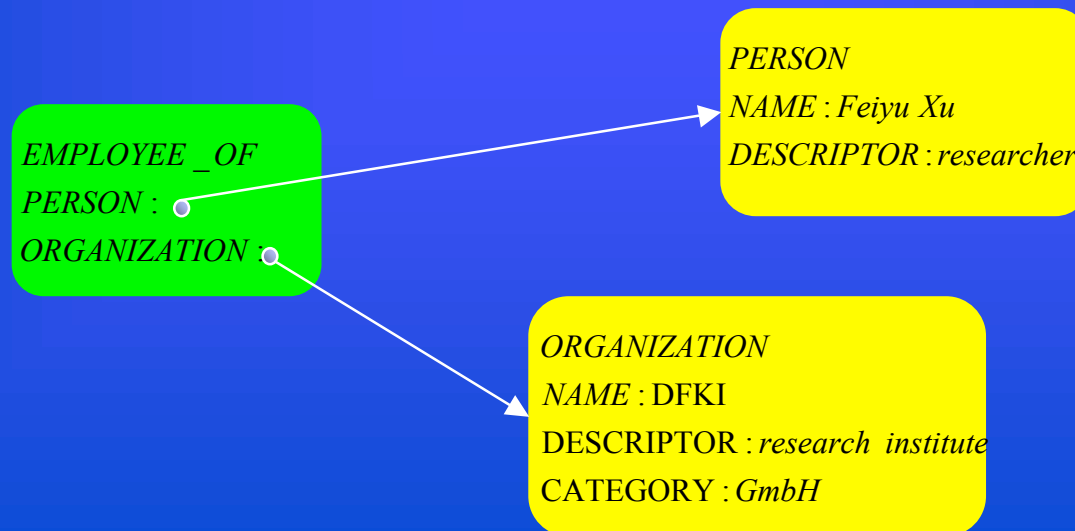
Generic IE tasks for MUC-7

- (NE) Named Entity Recognition Task requires the identification and classification of named entities
 - organizations
 - locations
 - persons
 - dates, times, percentages and monetary expressions
- (TE) Template Element Task requires the filling of small scale templates for specified classes of entities in the texts
 - Attributes of entities are slot fills (identifying the entities beyond the name level)
 - Example: Persons with slots such as name (plus name variants), title, nationality, description as supplied in the text, and subtype.

“Capitan Denis Gillespie, the comander of Carrier Air Wing 11”

Generic IE tasks for MUC-7

- (TR) Template Relation Task requires filling a two slot template representing a binary relation with pointers to template elements standing in the relation, which were previously identified in the TE task
 - subsidiary relationship between two companies (employee_of, product_of, location_of)



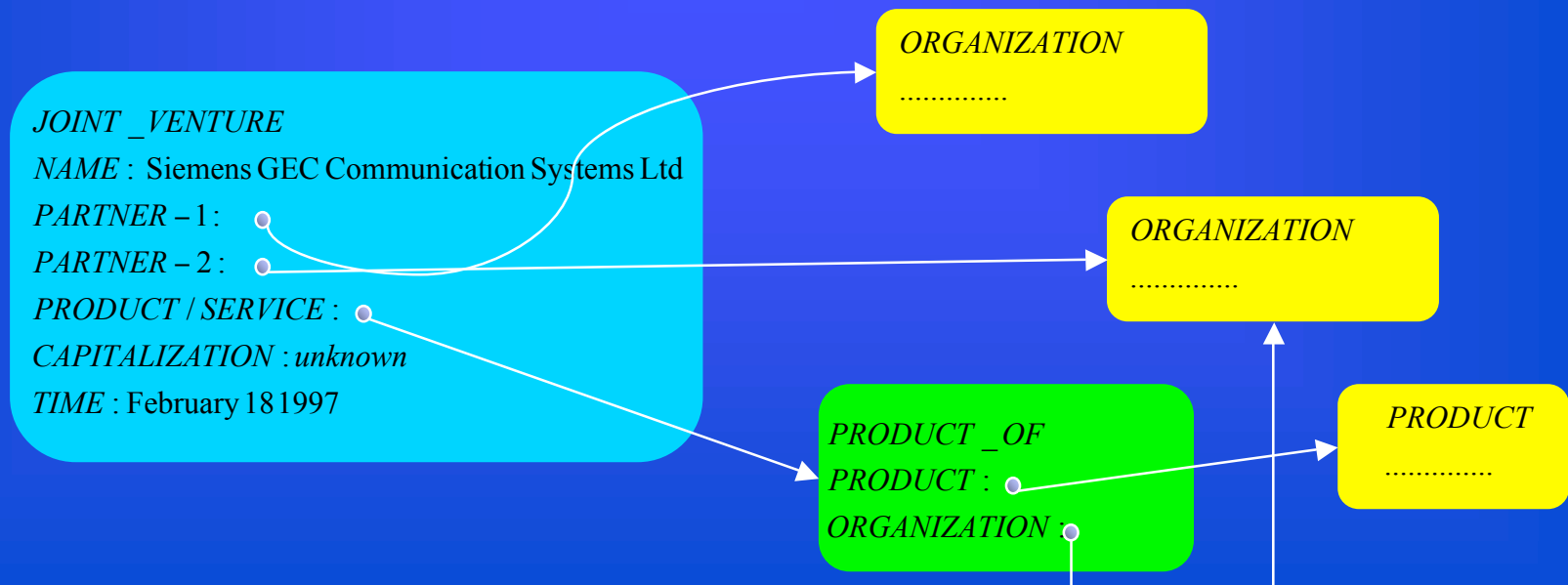
Generic IE tasks for MUC-7

- (CO) Coreference Resolution requires the identification of expressions in the text that refer to the same object, set or activity
 - variant forms of name expressions
 - definite noun phrases and their antecedents
 - pronouns and their antecedents

“**The U.K. satellite television broadcaster** said its subscriber base grew 17.5 percent during the past year to 5.35 million”

Generic IE tasks for MUC-7

- (ST) Scenario Template requires filling a template structure with extracted information involving several relations or events of interest
 - intended to be the MUC approximation to a real-world information extraction problem
 - identification of partners, products, profits and capitalization of joint ventures



Tasks evaluated in MUC 3-7

[Chinchor, 98]

EVALITASK	NE	CO	RE	TR	ST
MUC-3					YES
MUC-4					YES
MUC-5					YES
MUC-6	YES	YES	YES		YES
MUC-7	YES	YES	YES	YES	YES

Development Steps within IE Communities

- from attempts to use the methods of full text understanding to shallow text processing;
- from pure knowledge-based hand-coded systems to (semi-) automatic systems using machine learning methods;
- from complex domain-dependent event extraction to standardized domain-independent elementary entity identification, simple semantic relation and event extraction.

The ACE Program

- “Automated Content Extraction” since 1999
- Develop core information extraction technology by focusing on extracting specific semantic entities and relations over a very wide range of texts.
- Corpora: Newswire and broadcast transcripts, but broad range of topics and genres.
 - Third person reports
 - Interviews
 - Editorials
 - Topics: foreign relations, significant events, human interest, sports, weather
- Discourage highly domain- and genre-dependent solutions

Components of a Semantic Model

- Entities - Individuals in the world *that are mentioned in a text*
 - Simple entities: singular objects
 - Collective entities: sets of objects of the same type *where the set is explicitly mentioned in the text*
- Relations – Properties that hold of tuples of entities.
- Complex Relations – Relations that hold among entities and relations
- Attributes – one place relations are attributes or individual properties

Components of a Semantic Model

- Temporal points and intervals
- Relations may be timeless or bound to time intervals
- Events – A particular kind of simple or complex relation among entities involving a change in relation state at the end of a time interval.

Relations in Time

- timeless attribute: $\text{gender}(x)$
- time-dependent attribute: $\text{age}(x)$
- timeless two-place relation: $\text{father}(x, y)$
- time-dependent two-place relation: $\text{boss}(x, y)$

Relations vs. Features or Roles in AVMs

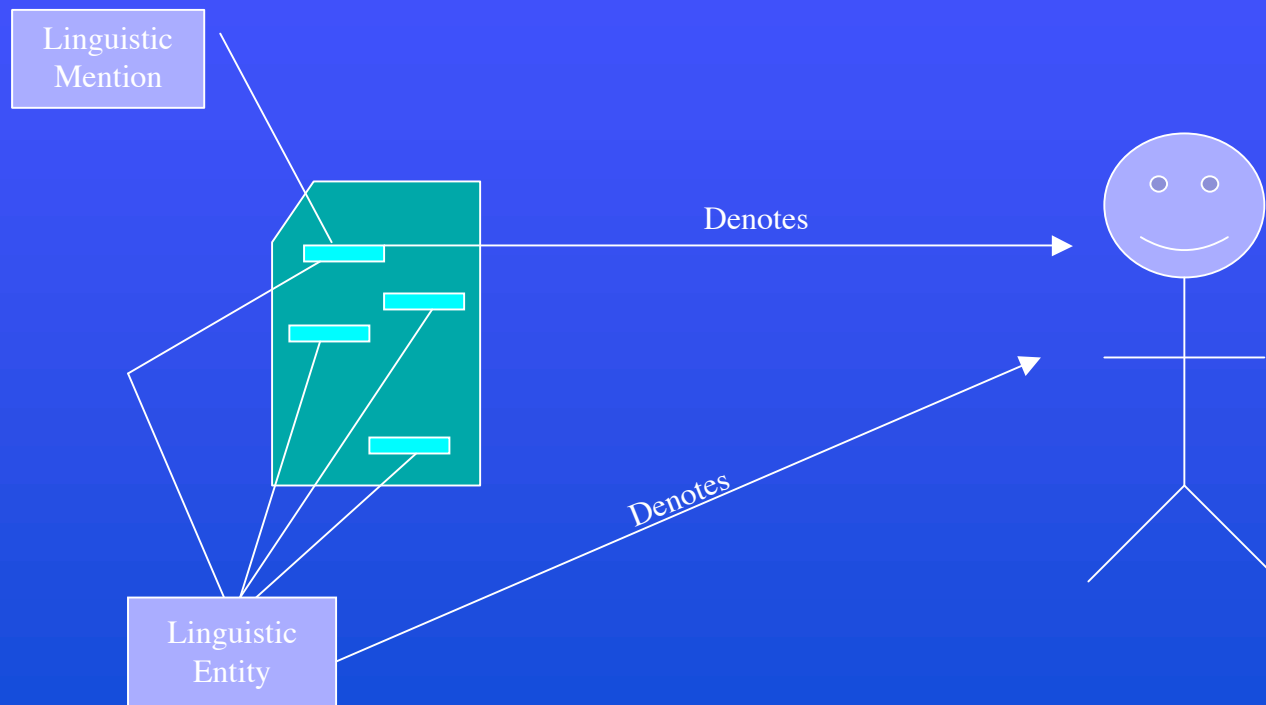
- Several two place relations between an entity x and other entities y_i can be bundled as properties of x . In this case, the relations are called roles (or attributes) and any pair $\langle \text{relation} : y_i \rangle$ is called a role assignment (or a feature).
- name $\langle x, CR \rangle$

name: Condoleezza Rice
office: National Security Advisor
age: 49
gender: female

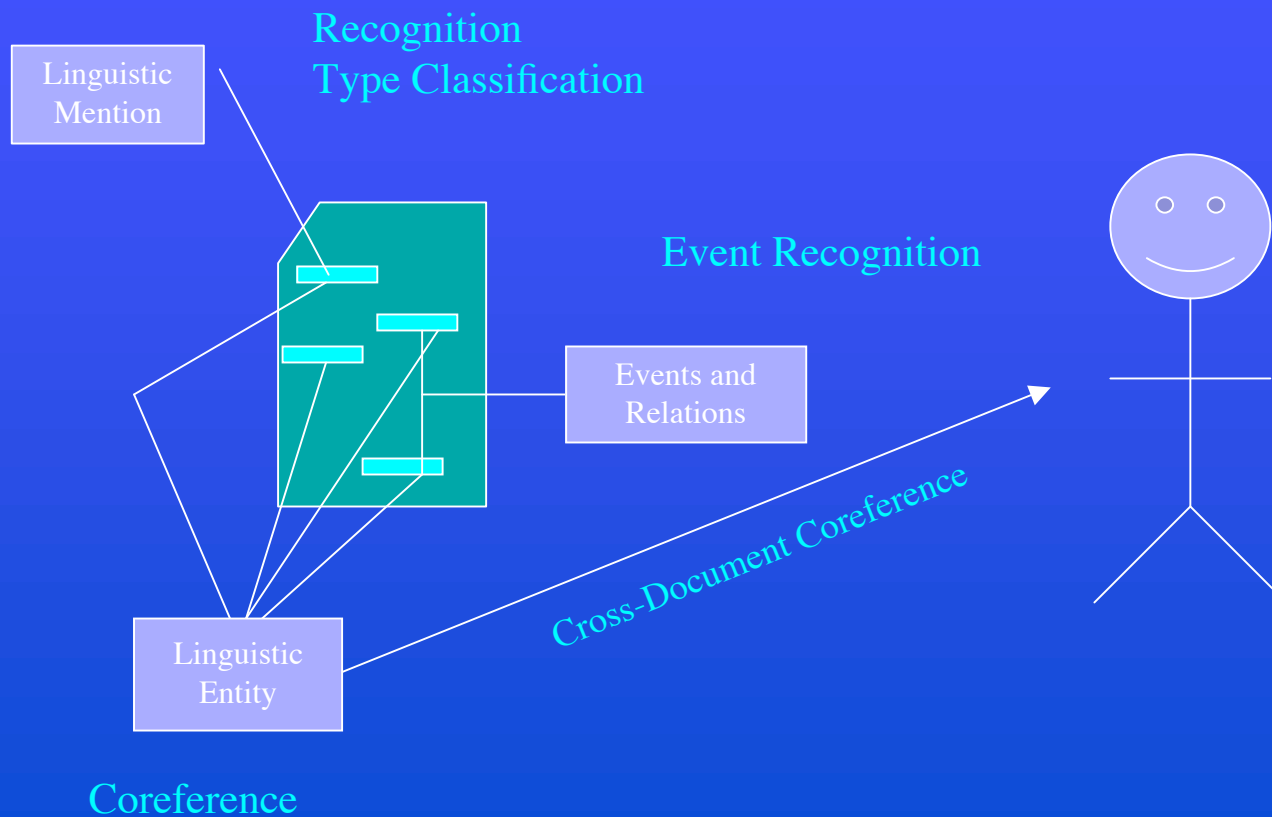
Semantic Analysis: Relating Language to the Model

- Linguistic Mention
 - A particular linguistic phrase
 - Denotes a particular entity, relation, or event
 - A noun phrase, name, or possessive pronoun
 - A verb, nominalization, compound nominal, or other linguistic construct relating other linguistic mentions
- Linguistic Entity
 - Equivalence class of mentions with same meaning
 - Coreferring noun phrases
 - Relations and events derived from different mentions, but conveying the same meaning

Language and World Model



NLP Tasks in an Extraction System



Example

1. Three of the Nobel Prizes for Chemistry during the first decade were awarded for pioneering work in organic chemistry.
2. In 1902 Emil Fischer (1852-1919), then in Berlin, was given the prize for his work on sugar and purine syntheses.
3. Another major influence from organic chemistry was the development of the chemical industry, and a chief contributor here was Fischer's teacher, Adolf von Baeyer (1835-1917) in Munich, who was awarded the prize in 1905.

Anaphora in Texts

He/The scientist won the 2005 Nobel Prize for Peace on Friday for his efforts to limit the spread of atomic weapons.



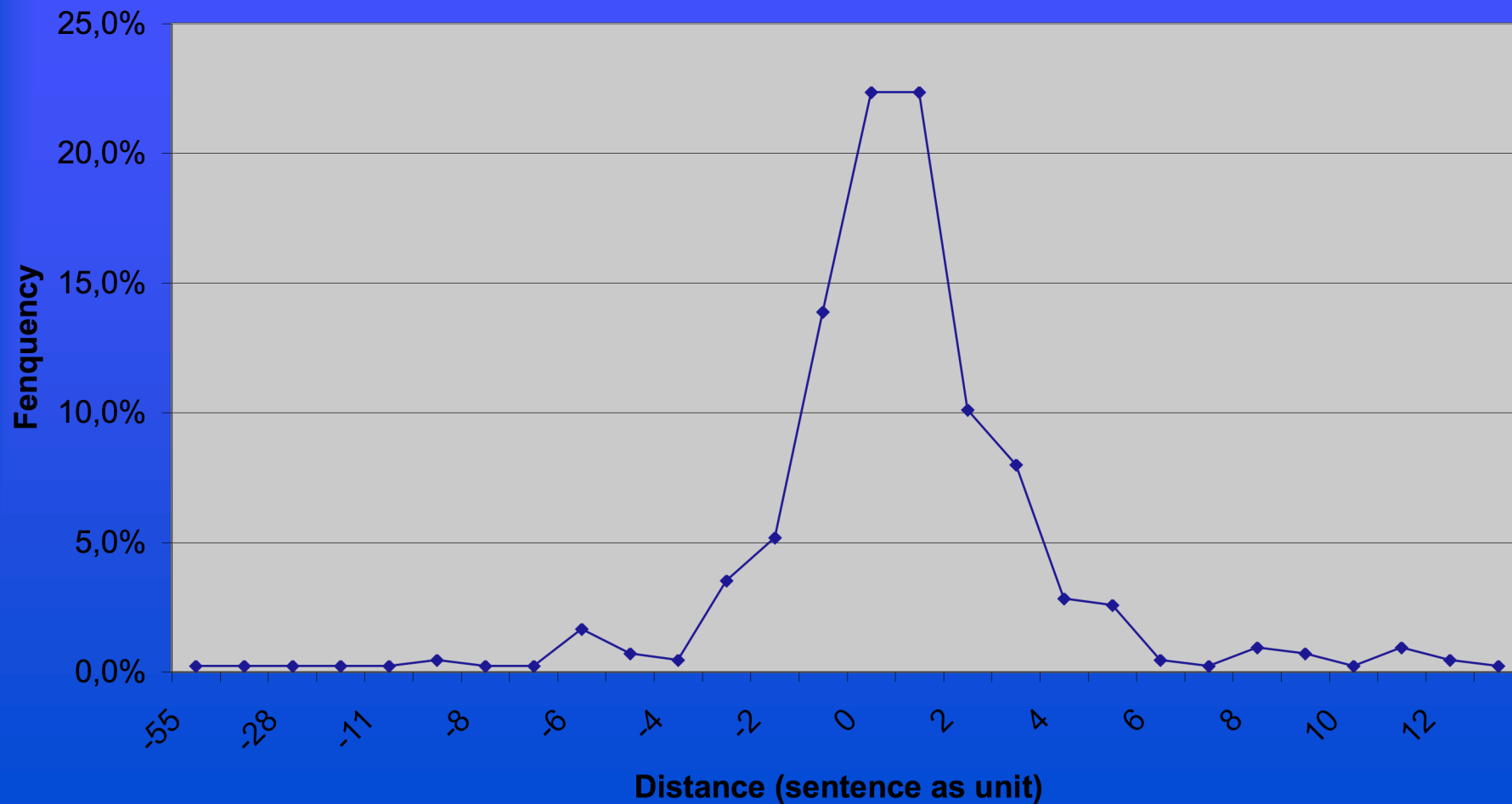
<?PERSON, Nobel, Peace, 2005>

Coreference Relations and Indicators

- Complex linguistic phenomena, influenced by lexical, syntactic, semantic and discourse constraints
- The indicators shared by many approaches are
 - Distance: coreference expressions are often close to each other in the surface structure;
 - Syntactic: pronominal resolution constraints within sentence
 - Semantic: same or compatible semantic category, agreement in number, gender and person;
 - Discourse: parallelism, repetition, apposition, name alias.

Receny Indicator in Nobel Prize Domain

- News reports from New York Times, online BBC and CCN (18.4 MB, 3328 documents)



1. Two Americans have won the 2002 Nobel Prize in Economic Sciences.
2. The two scientists, Daniel Kahneman and Vernon L. Smith, received the honour on Wednesday for their work using psychological research and laboratory experiments in economic analysis.

1. Egypt honours its Nobel Prize chemist.
2. President Hosni Mubarak of Egypt has awarded the country's most prestigious prize - the Nile Necklace - to the Egyptian-born chemist Ahmed Zewail.

Repetition and Elaboration

- Cohesion indicator *repetition* is often used as indicator for semantic similarity and semantic consistency, e.g.,
 - „two Americans“ and „two scientists“
 - „chemist“ and „chemist“
- Elaboration phenomena are normal in newspaper texts

S1 is an Elaboration of S0 if a proposition P follows from the assertions of both S0 and S1, but S1 contains a property of one of the elements of P that is not in S0 (Hobbs, 1979)

Relation Argument as a Complex Semantic Object

- A complex noun phrase contains often more than one property about an argument: e.g.

Egyptian-born chemist Ahmed Zewail

- Relevant properties of a winner in Nobel Prize domain
 - Nationality/origin/inhabitant: e.g., two Americans, the Egyptian-born, a Dutch
 - Profession/occupation: e.g., novelist, chemist, scientist, researcher
 - Title/position: e.g., professor, president
 - Domain description: e.g., recipient, winner, Nobel Laureate
 - General description: e.g., the man, a woman, the team

„two Americans“

sentence_id : i

*number : [type : plural
amount : 2]*

definite : indef

grammarrole : subject

semantics : [nationality : american]

„two scientists“

$$\left[\begin{array}{l} \textit{sentence_id} : i + 1 \\ \textit{number} : \left[\begin{array}{l} \textit{type} : \textit{plural} \\ \textit{amount} : 2 \end{array} \right] \\ \textit{definite} : \textit{def} \\ \textit{grammarrole} : \textit{subject} \\ \textit{semantics} : \left[\textit{profession} : \textit{scientist} \right] \\ \textit{names} : \langle \textit{name1} \quad \textit{name2} \rangle \end{array} \right]$$

Unification of
„two Americans“ and „two scientists“

$$\left[\begin{array}{l} \textit{number} : \left[\begin{array}{l} \textit{type} : \textit{plural} \\ \textit{amount} : 2 \end{array} \right] \\ \textit{semantics} : \left[\begin{array}{l} \textit{nationality} : \textit{american} \\ \textit{profession} : \textit{scientist} \end{array} \right] \\ \textit{names} : \langle \textit{name1} \quad \textit{name2} \rangle \end{array} \right]$$

The Basic Semantic Tasks of an IE System

- Recognition of linguistic entities
- Classification of linguistic entities into semantic types
- Identification of coreference equivalence classes of linguistic entities
- Identifying the actual individuals that are mentioned in an article
 - Associating linguistic entities with predefined individuals (e.g. a database, or knowledge base)
 - Forming equivalence classes of linguistic entities from different documents.

The ACE Ontology

- **Persons**
 - A natural kind, and hence self-evident
- **Organizations**
 - Should have some persistent existence that transcends a mere set of individuals
- **Locations**
 - Geographic places with no associated governments
- **Facilities**
 - Objects from the domain of civil engineering
- **Geopolitical Entities**
 - Geographic places with associated governments

Why GPEs

- An ontological problem: certain entities have attributes of physical objects in some contexts, organizations in some contexts, and collections of people in others
- Sometimes it is difficult to impossible to determine which aspect is intended
- It appears that in some contexts, the same phrase plays different roles in different clauses

Aspects of GPEs

- Physical
 - San Francisco has a mild climate
- Organization
 - The United States is seeking a solution to the North Korean problem.
- Population
 - France makes a lot of good wine.

Types of Linguistic Mentions

- Name mentions
 - The mention uses a proper name to refer to the entity
- Nominal mentions
 - The mention is a noun phrase whose head is a common noun
- Pronominal mentions
 - The mention is a headless noun phrase, or a noun phrase whose head is a pronoun, or a possessive pronoun

Explicit and Implicit Relations

- Many relations are true in the world. Reasonable knowledge bases used by extraction systems will include many of these relations. Semantic analysis requires focusing on certain ones that are directly motivated by the text.
- Example:
 - Baltimore is in Maryland, which is in United States.
 - “Baltimore, MD”
 - Text mentions Baltimore and United States. Is there a relation between Baltimore and United States?

Another Example

- *Prime Minister Tony Blair attempted to convince the British Parliament of the necessity of intervening in Iraq.*
- Is there a role relation specifying Tony Blair as prime minister of Britain?
- A test: a relation is implicit in the text if the text provides convincing evidence that the relation actually holds.

Explicit Relations

- Explicit relations are expressed by certain surface linguistic forms
 - Copular predication - Clinton was the president.
 - Prepositional Phrase - The CEO of Microsoft...
 - Prenominal modification - The American envoy...
 - Possessive - Microsoft's chief scientist...
 - SVO relations - Clinton arrived in Tel Aviv...
 - Nominalizations - Anan's visit to Baghdad...
 - Apposition - Tony Blair, Britain's prime minister...

Types of ACE Relations

- **ROLE** - relates a person to an organization or a geopolitical entity
 - Subtypes: member, owner, affiliate, client, citizen
- **PART** - generalized containment
 - Subtypes: subsidiary, physical part-of, set membership
- **AT** - permanent and transient locations
 - Subtypes: located, based-in, residence
- **SOC** - social relations among persons
 - Subtypes: parent, sibling, spouse, grandparent, associate

Event Types (preliminary)

- Movement
 - Travel, visit, move, arrive, depart ...
- Transfer
 - Give, take, steal, buy, sell...
- Creation/Discovery
 - Birth, make, discover, learn, invent...
- Destruction
 - die, destroy, wound, kill, damage...

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

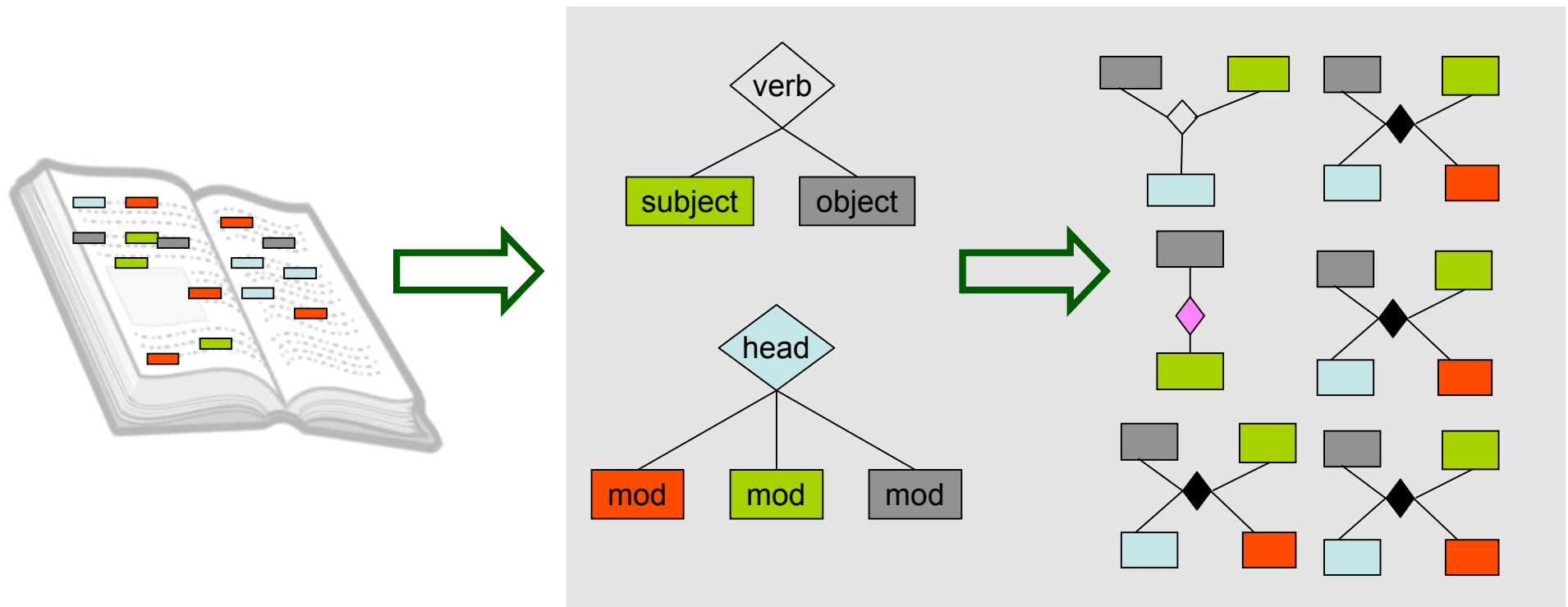
- Document structure
 - Free text
 - Semi-structured
 - Structured
 - Richness of the annotation
 - Shallow NLP
 - Deep NLP
 - Complexity of the template filling rules
 - Single slot
 - Multi slot
 - Amount of data
- Degree of automation
 - Semi-automatic
 - Supervised
 - Semi-Supervised
 - Unsupervised
 - Human interaction/contribution
 - Evaluation/validation
 - during learning loop
 - Performance: recall and precision

Documents

- Unstructured (Free) Text
 - Regular sentences and paragraphs
 - Linguistic techniques, e.g., NLP
- Structured Text
 - Itemized information
 - Uniform syntactic clues, e.g., table understanding
- Semi-structured Text
 - Ungrammatical, telegraphic (e.g., missing attributes, multi-value attributes, ...)
 - Specialized programs, e.g., wrappers

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

Outline

- State of the art
- Domain Adaptive Relation Extraction Framework (DARE)
- Experiments and evaluations
- Performance analysis and discussion
- Conclusion and future work

Outline

- State of the art
- Domain Adaptive Relation Extraction Framework (DARE)
- Experiments and evaluations
- Performance analysis and discussion
- Conclusion and future work

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.

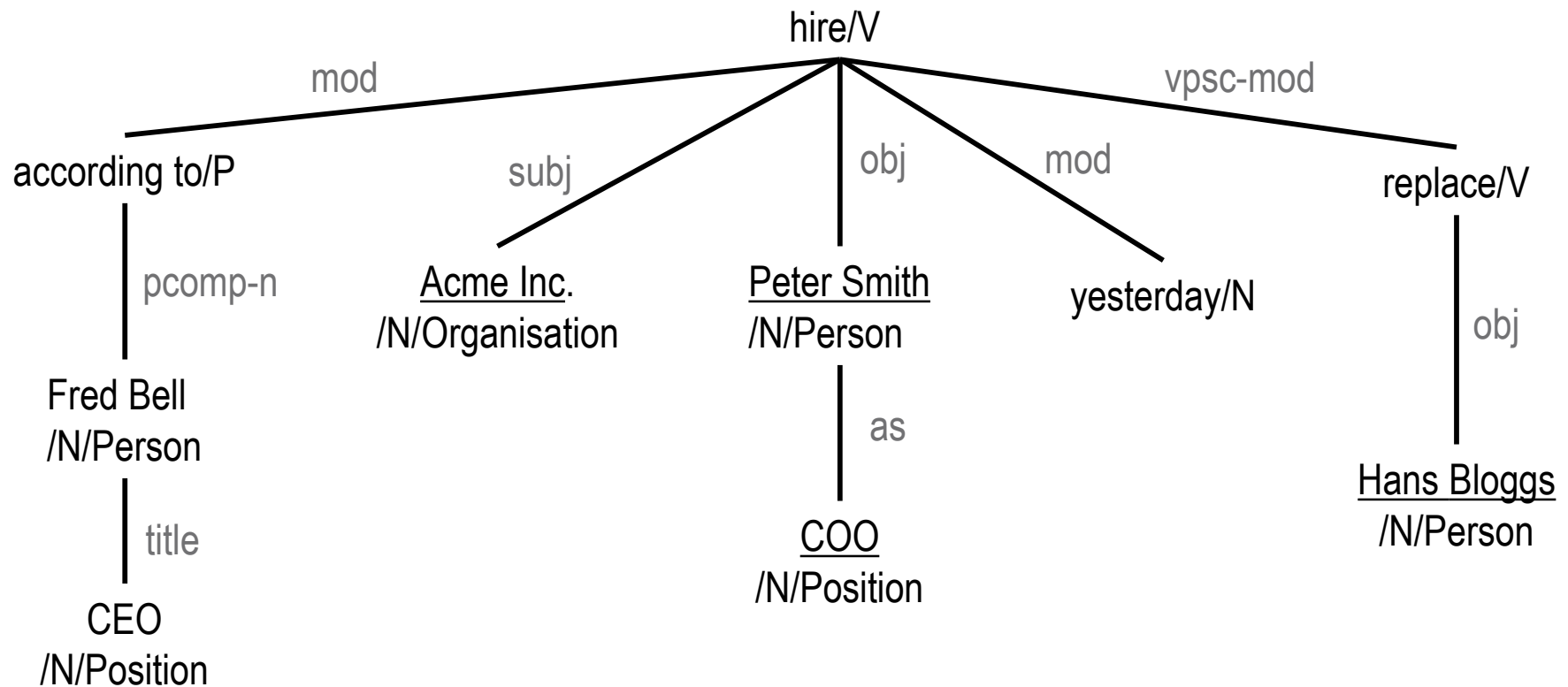
<person_in, person_out, position, organisation>

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

<person_in, person_out, position, organisation>

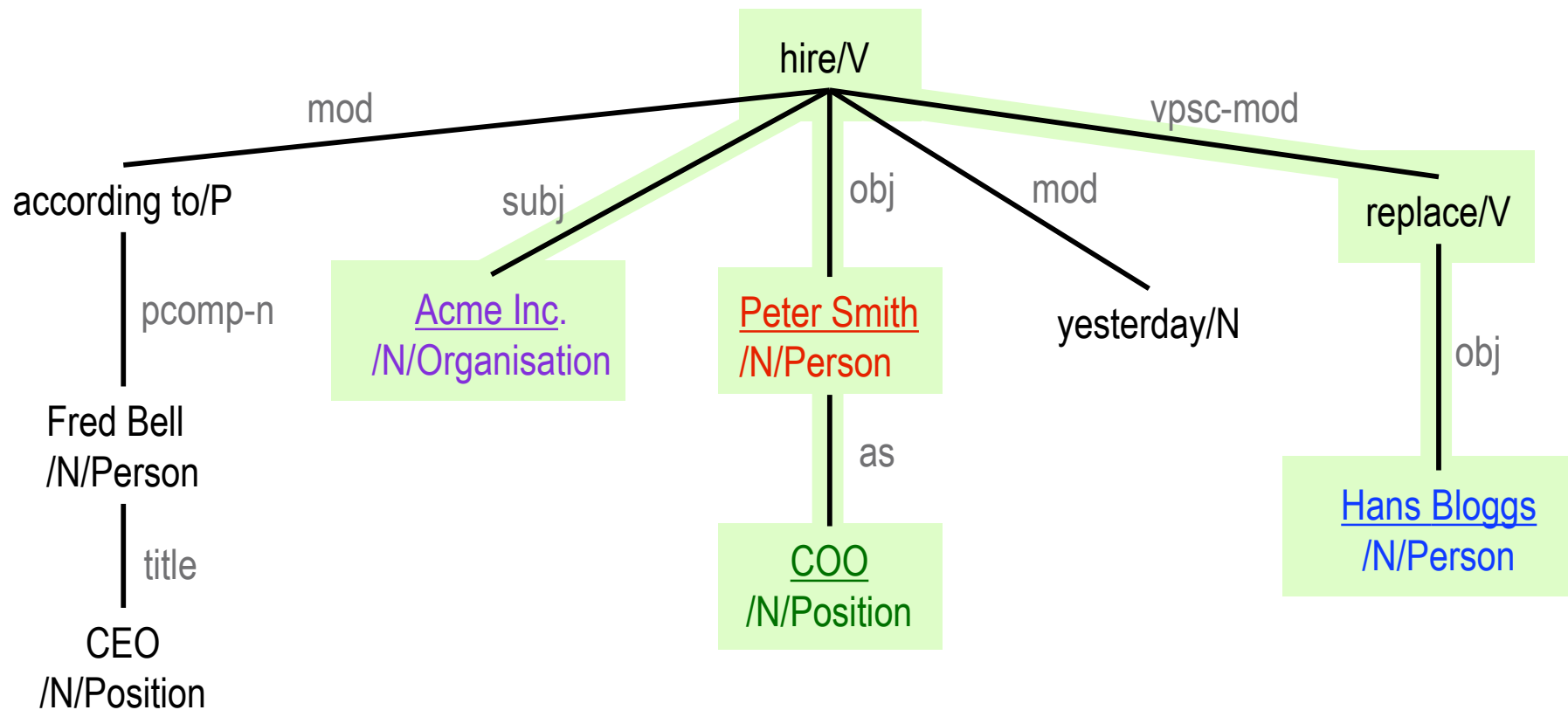
According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<person_in, person_out, position, organisation>

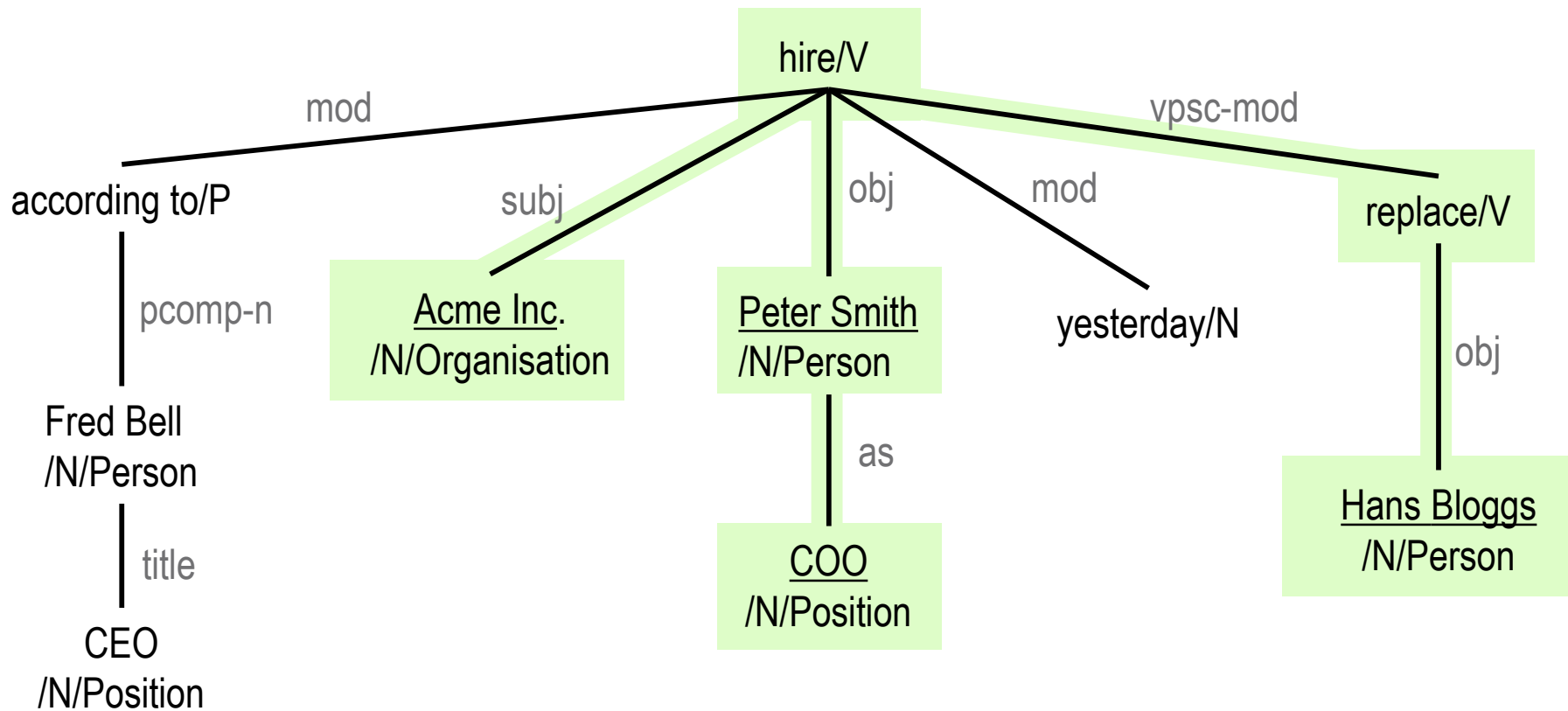


According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<person_in, person_out, position, organisation>



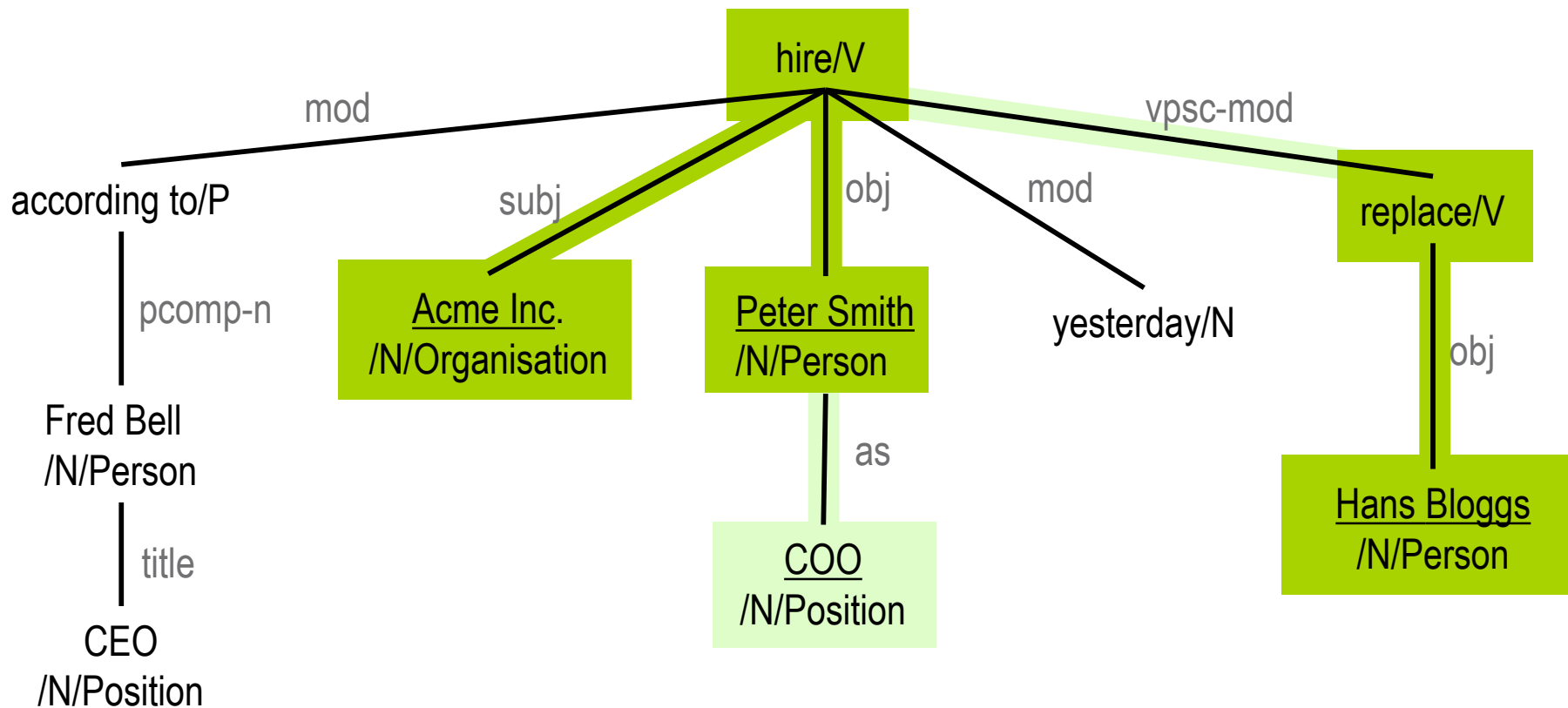
Ideal Target Pattern



Previous Work: SVO Model

Yangarber (2001)

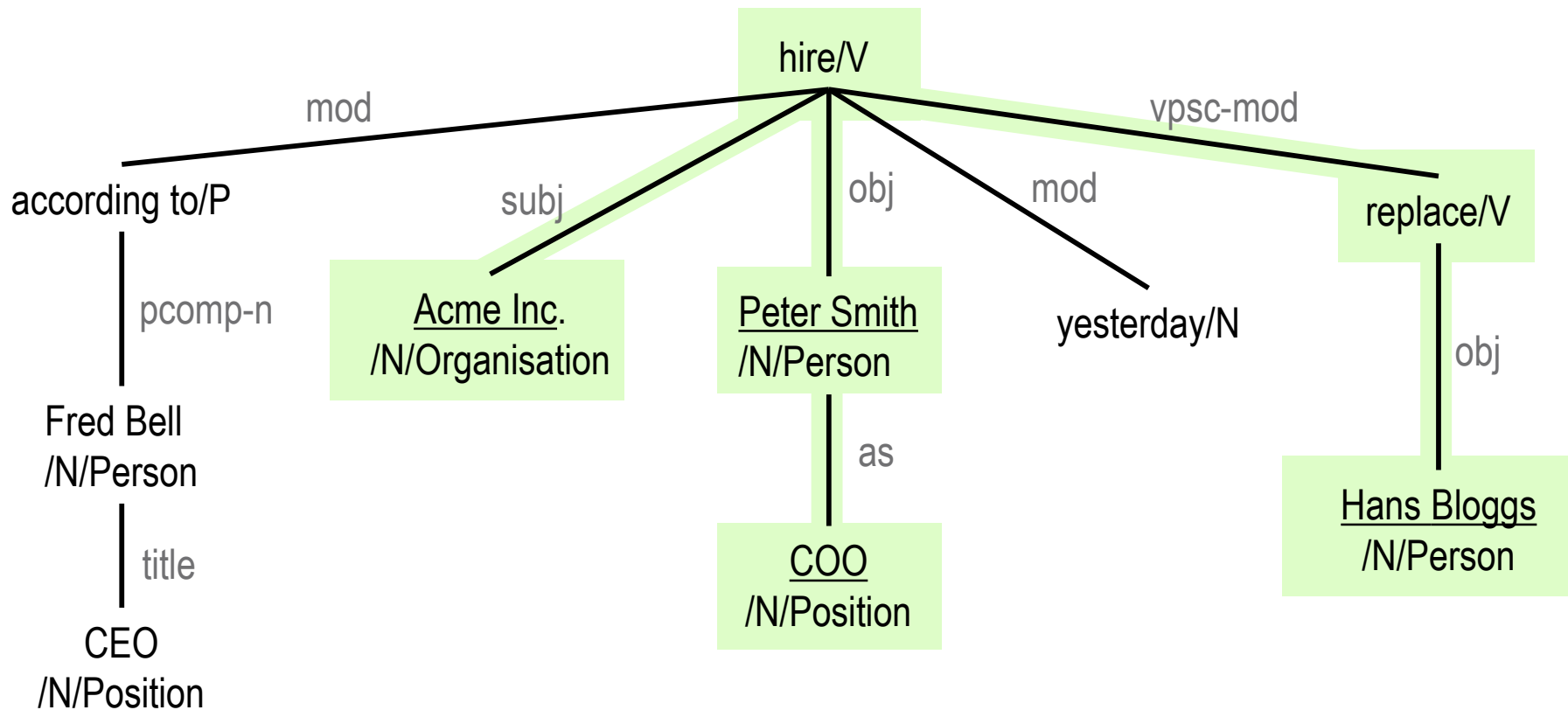
- Verb centered
- Direct relations between subject-verb-object
- Complex NP can not be extracted, e.g., the person and position relation
- The linguistic relations among patterns are not considered, e.g., hire and replace



Previous Work: Chain Model

Sudo et al. (2001)

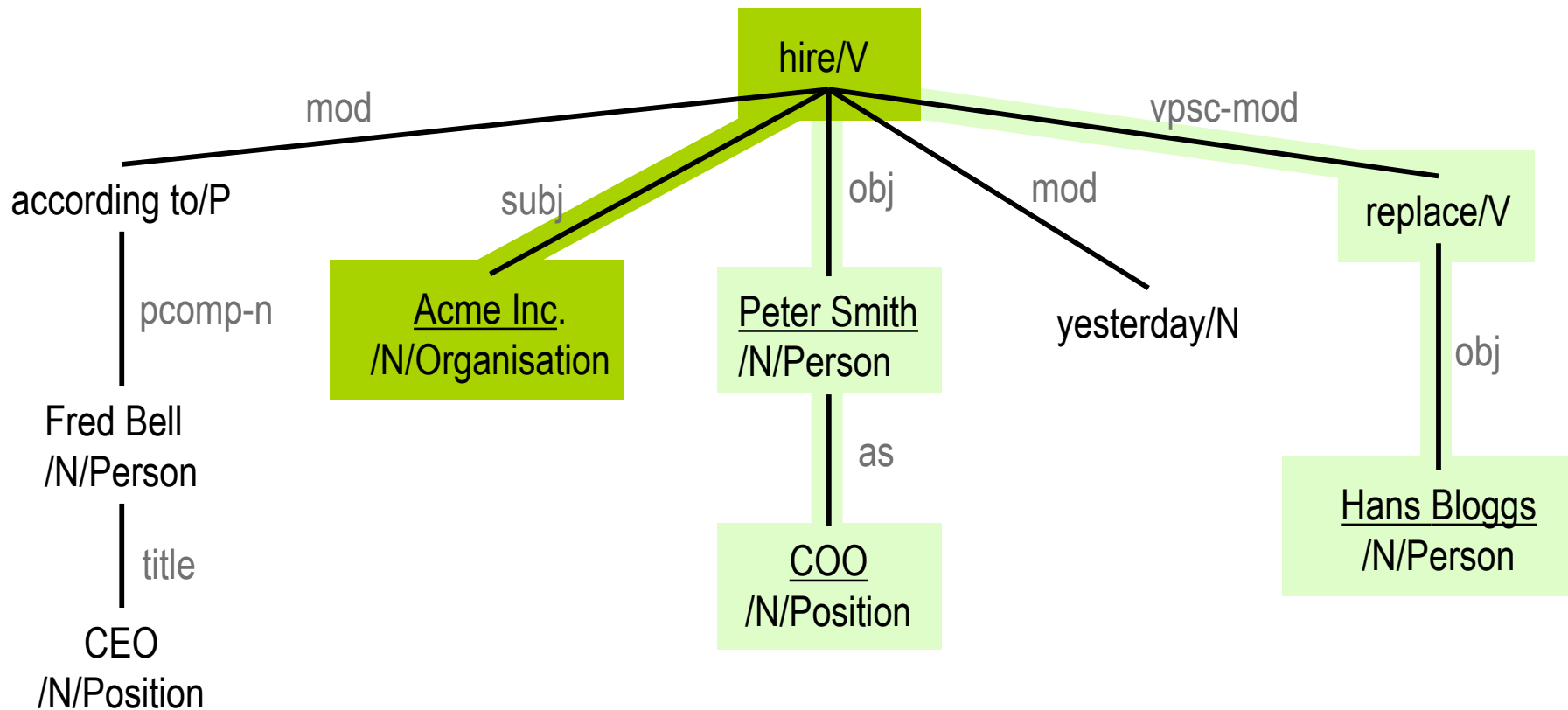
- ▣ Verb centered
- ▣ A single syntactic path dominated by a verb containing at least one relevant named entity concept



Previous Work: Chain Model

Sudo et al. (2001)

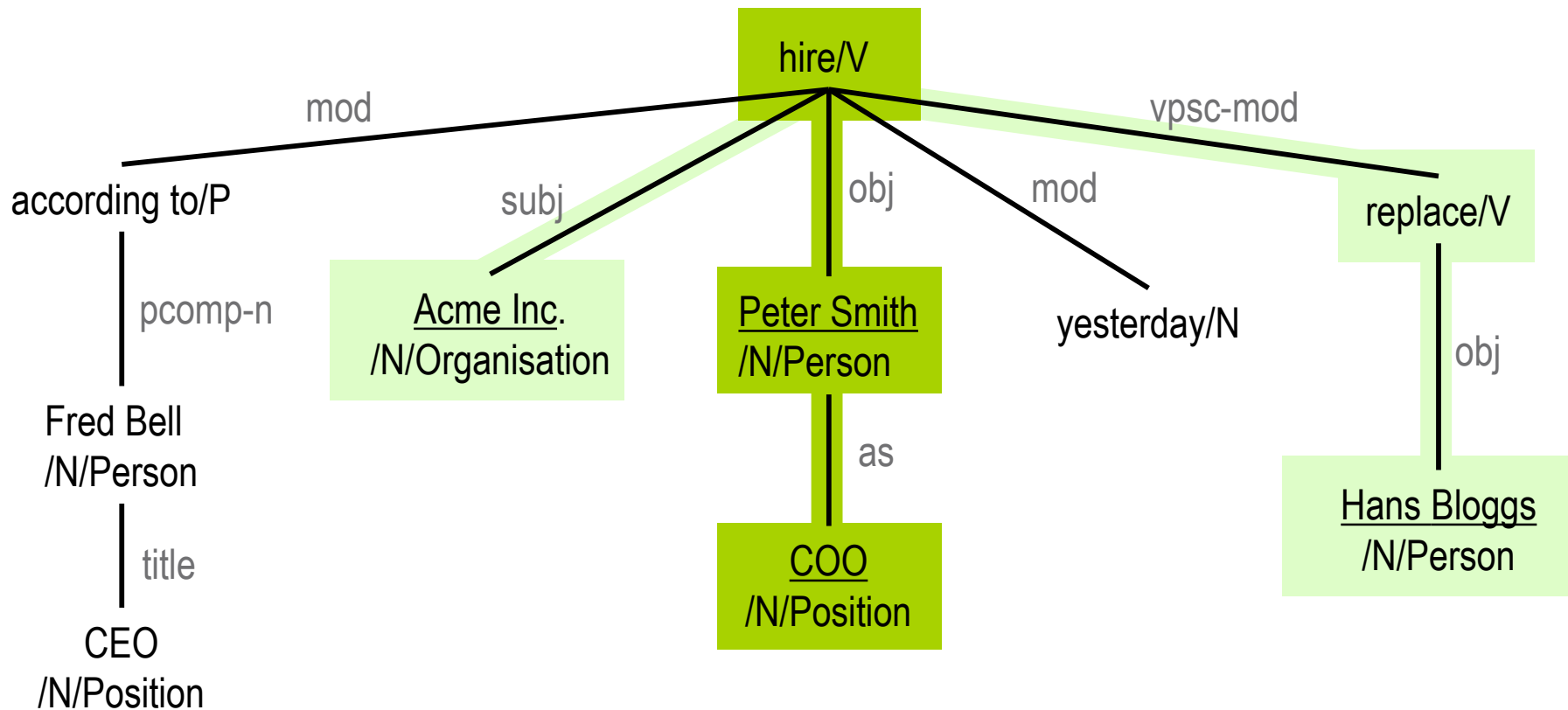
- Verb centered
- A single syntactic path dominated by a verb containing at least one relevant named entity concept



Previous Work: Chain Model

Sudo et al. (2001)

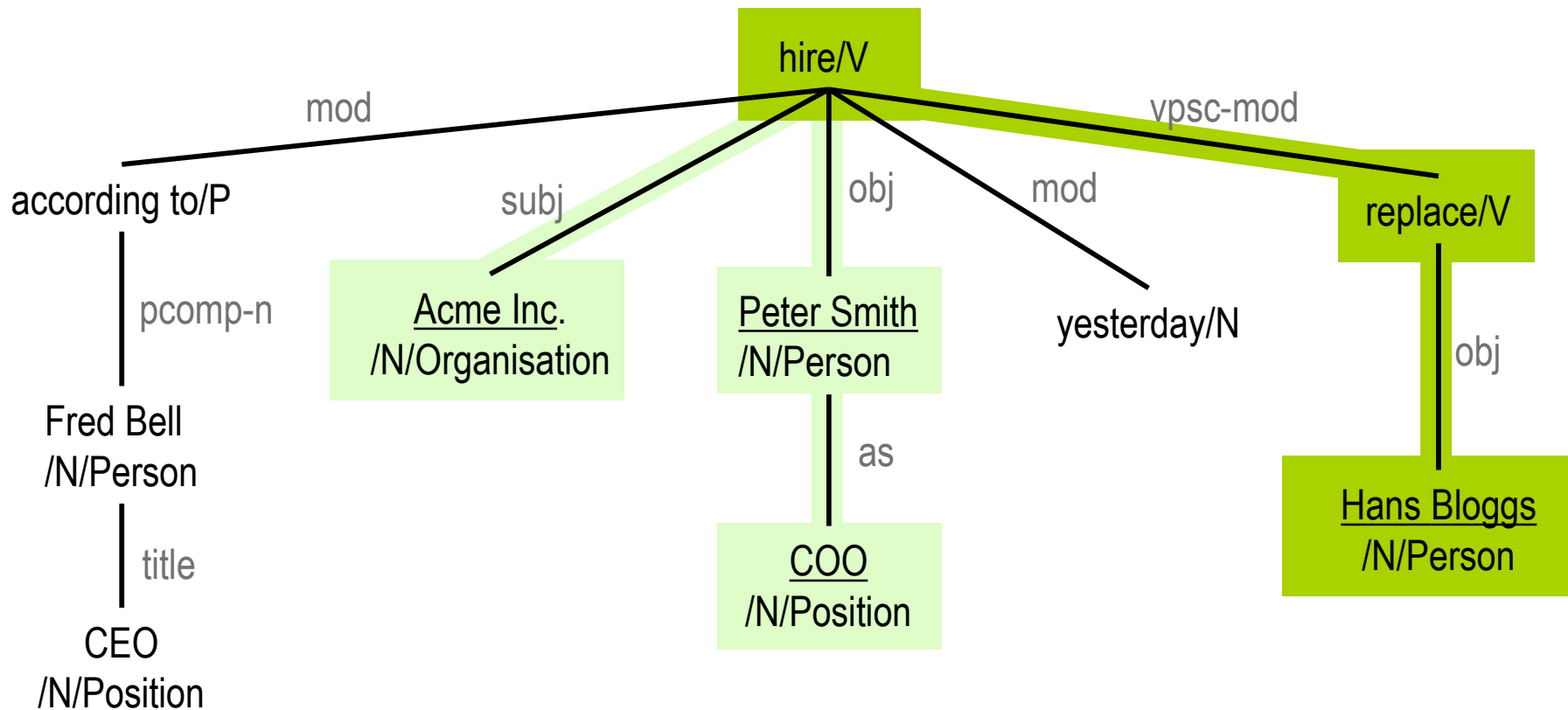
- Verb centered
- A single syntactic path dominated by a verb containing at least one relevant named entity concept



Previous Work: Chain Model

Sudo et al. (2001)

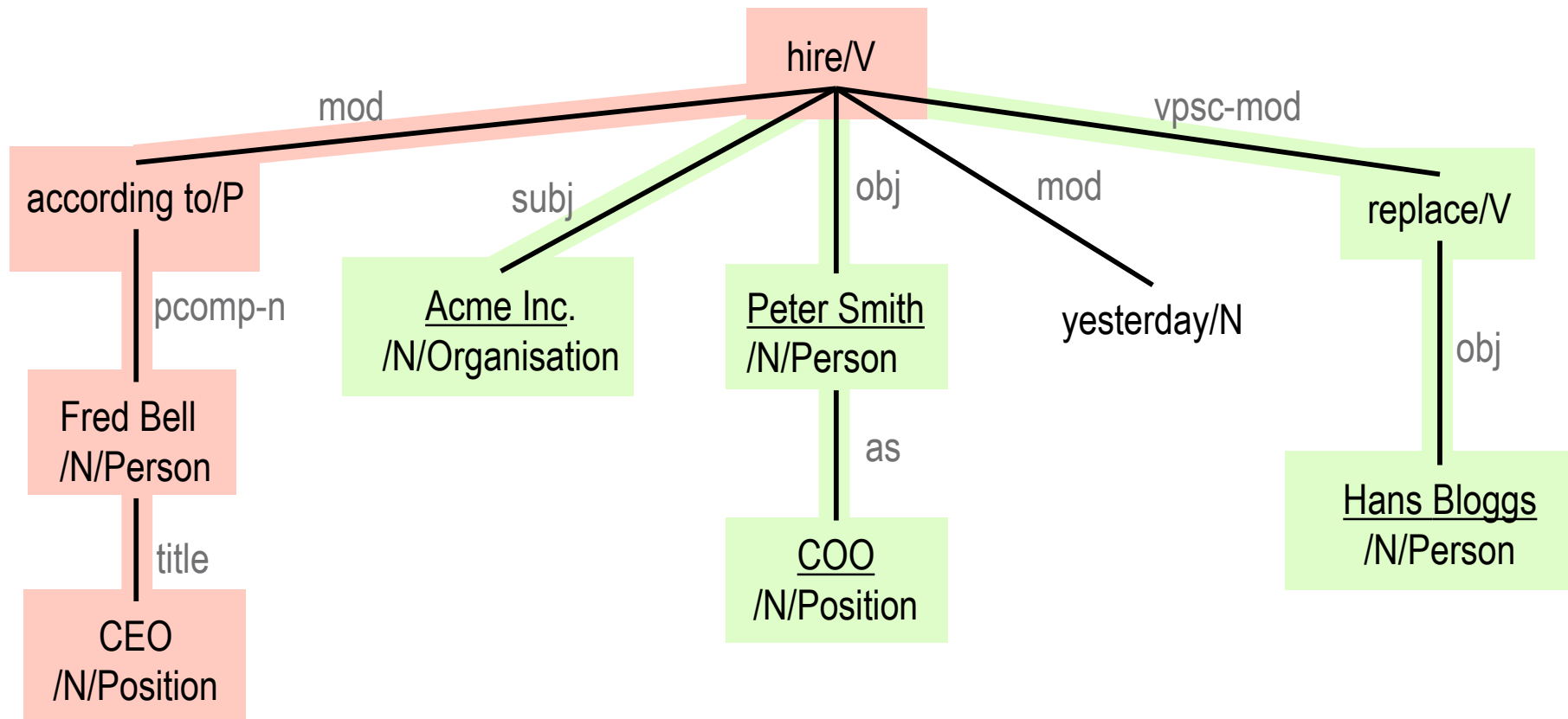
- Verb centered
- A single syntactic path dominated by a verb containing at least one relevant named entity concept



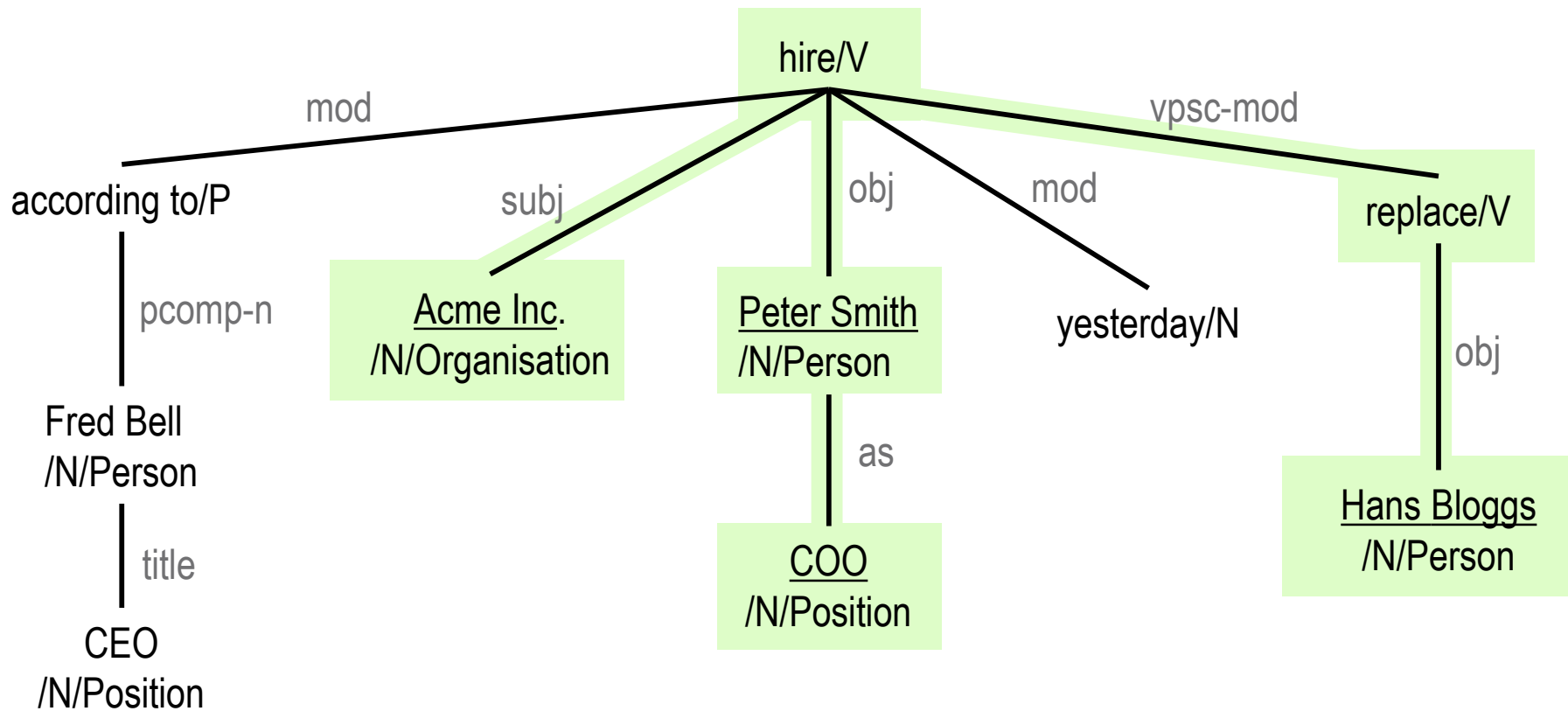
Previous Work: Chain Model

Sudo et al. (2001)

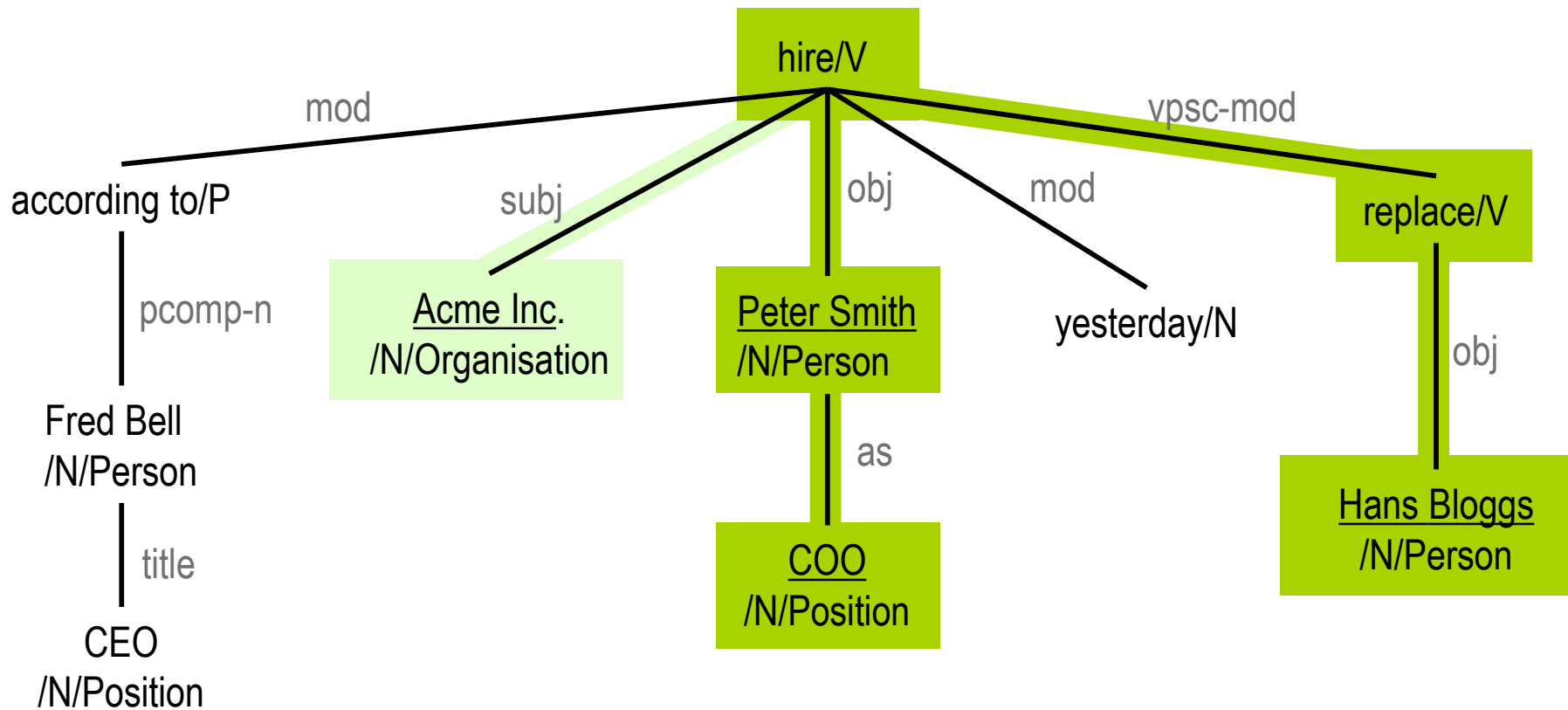
- Verb centered
- A single syntactic path dominated by a verb containing at least one relevant named entity concept



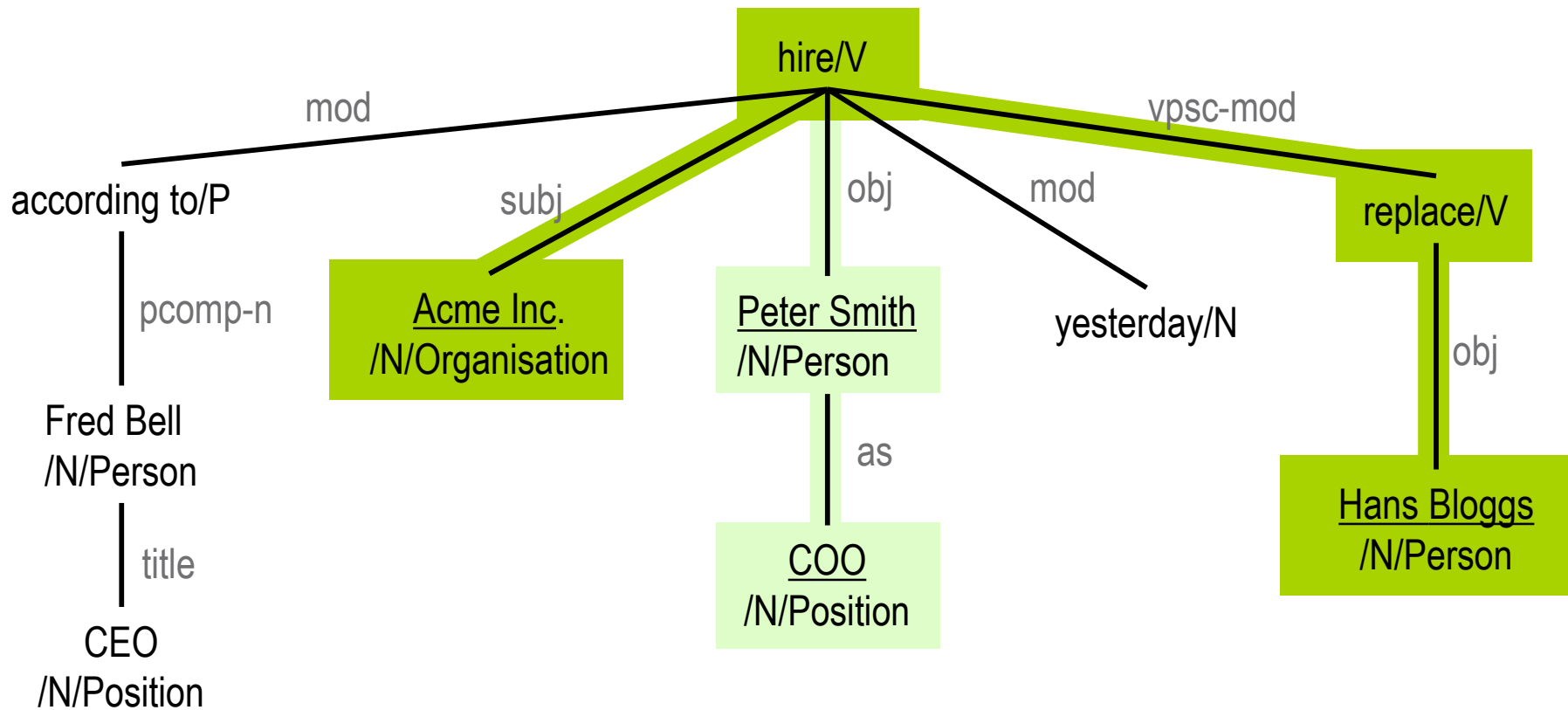
- verb centered
- pairs of chains instead of single paths



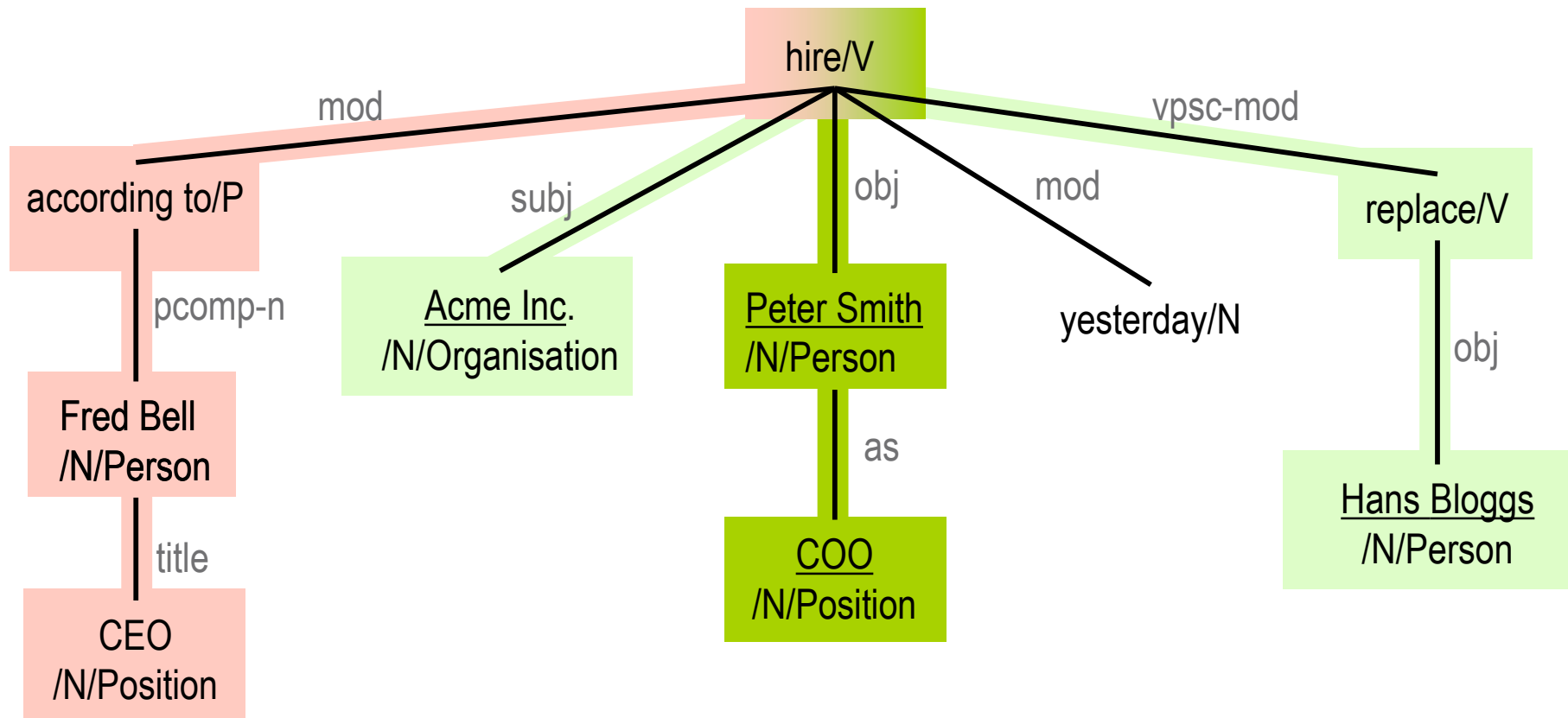
- verb centered
- pairs of chains instead of single paths



- ▣ verb centered
- ▣ pairs of chains instead of single paths



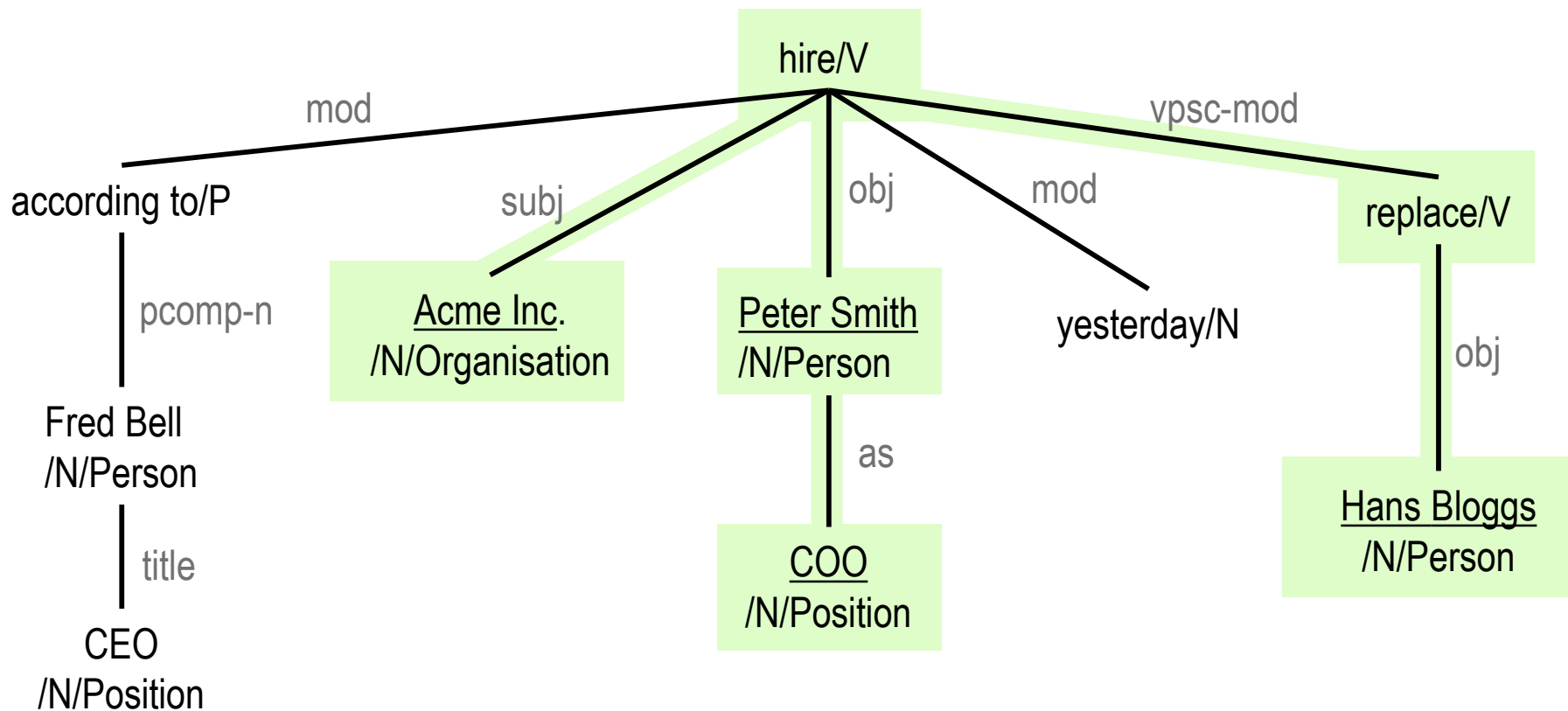
- verb centered
- pairs of chains instead of single paths



Previous Work: Subtree-Model

Sudo et al. (2003)

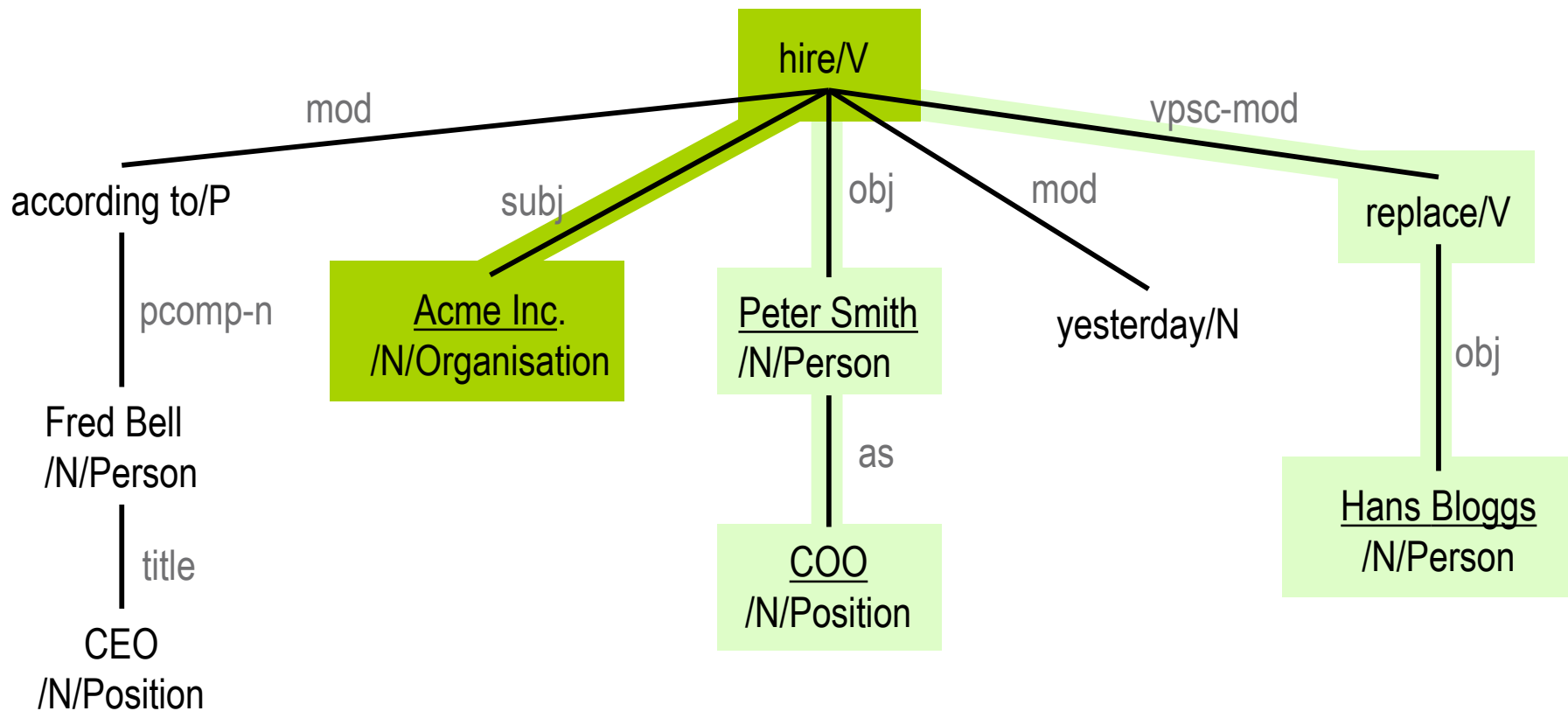
- verb centered
- All chains dominated by a verb, which contain at least one relevant named entity and their combinations



Previous Work: Subtree-Model

Sudo et al. (2003)

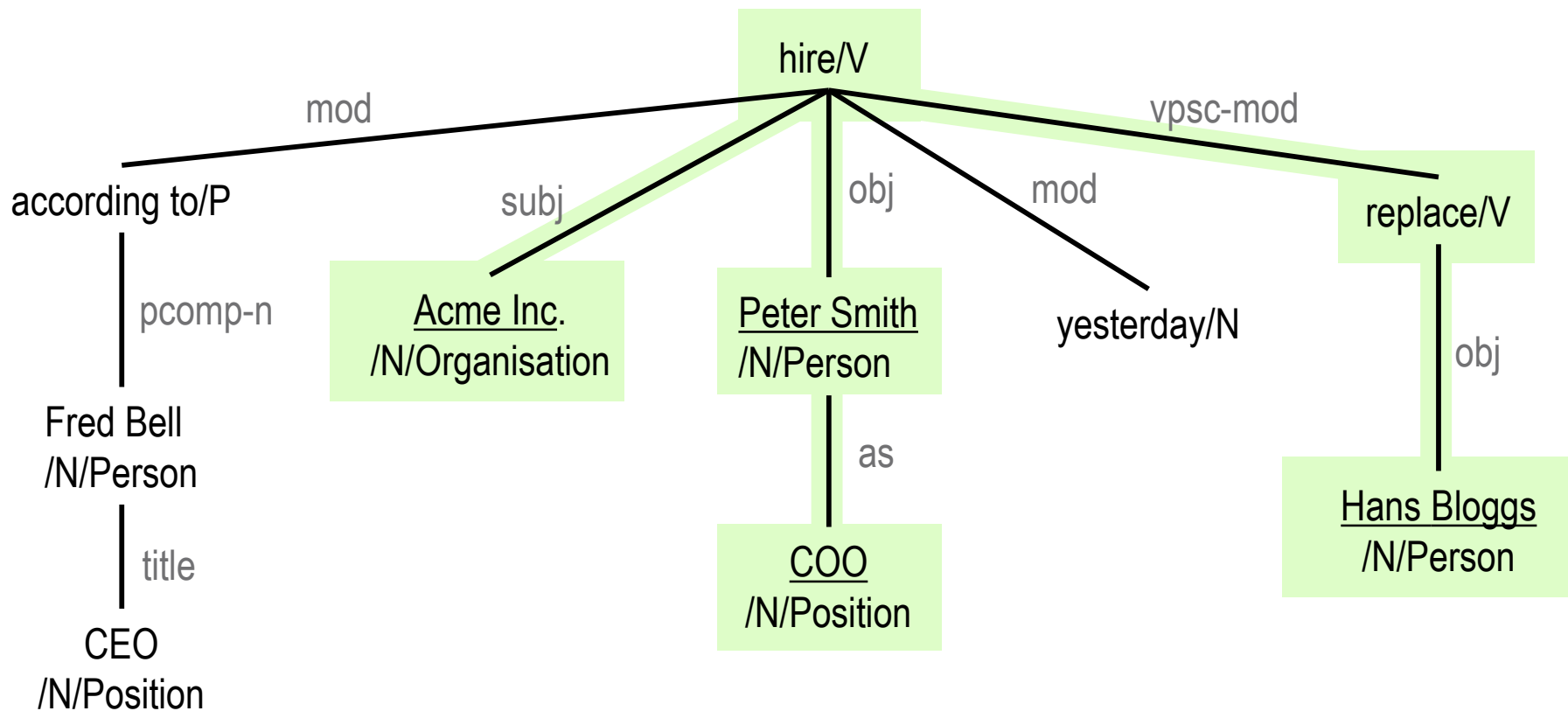
- verb centered
- All chains dominated by a verb, which contain at least one relevant named entity and their combinations



Previous Work: Subtree-Model

Sudo et al. (2003)

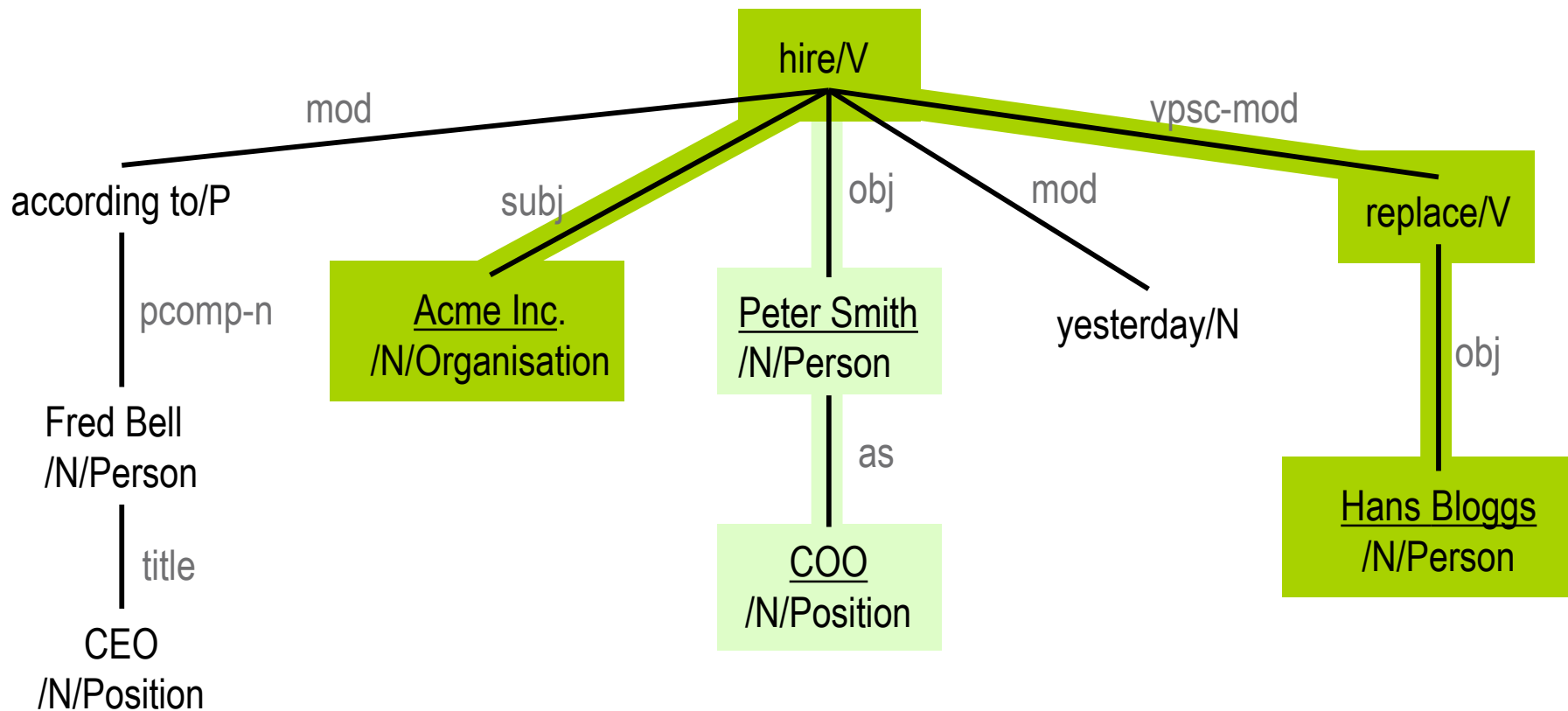
- ▣ verb centered
- ▣ All chains dominated by a verb, which contain at least one relevant named entity and their combinations



Previous Work: Subtree-Model

Sudo et al. (2003)

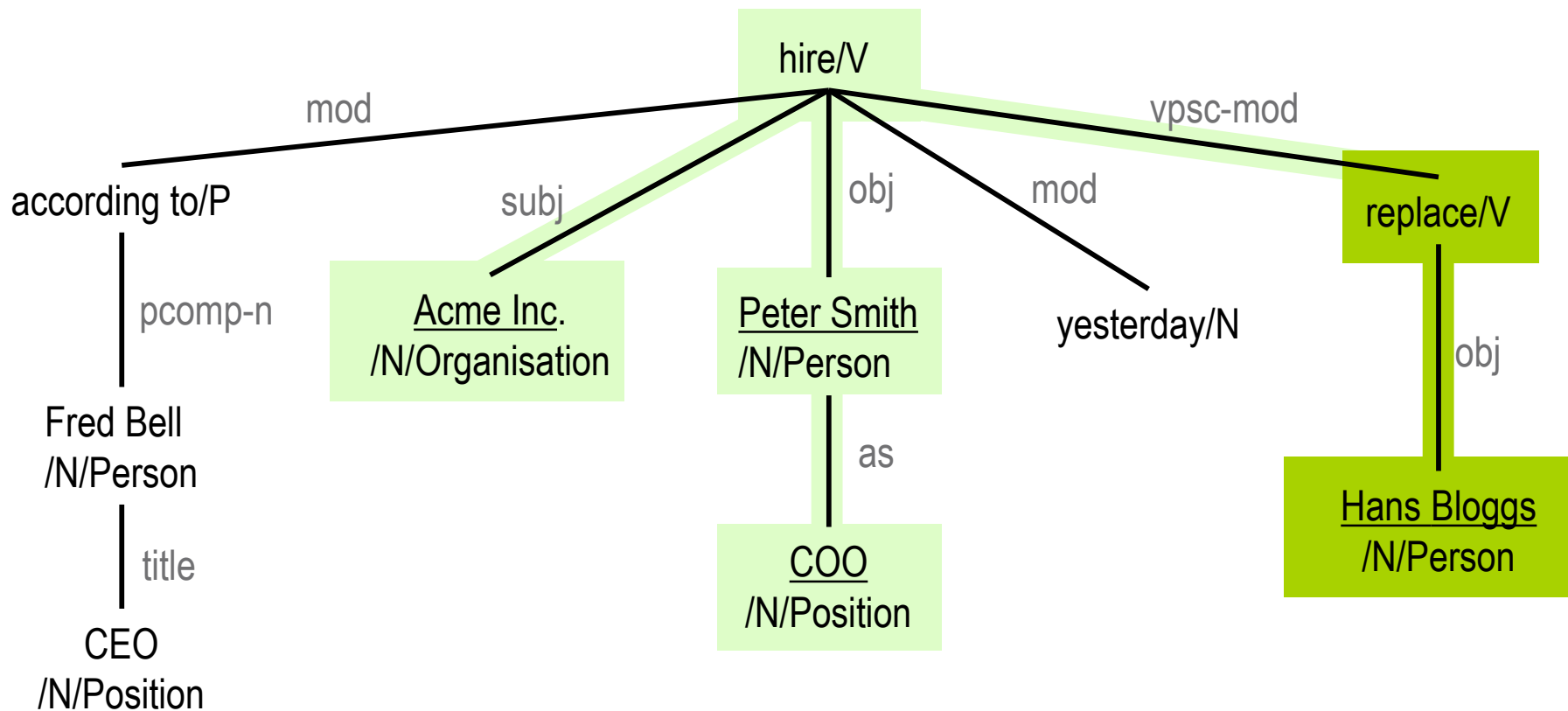
- verb centered
- All chains dominated by a verb, which contain at least one relevant named entity and their combinations



Previous Work: Subtree-Model

Sudo et al. (2003)

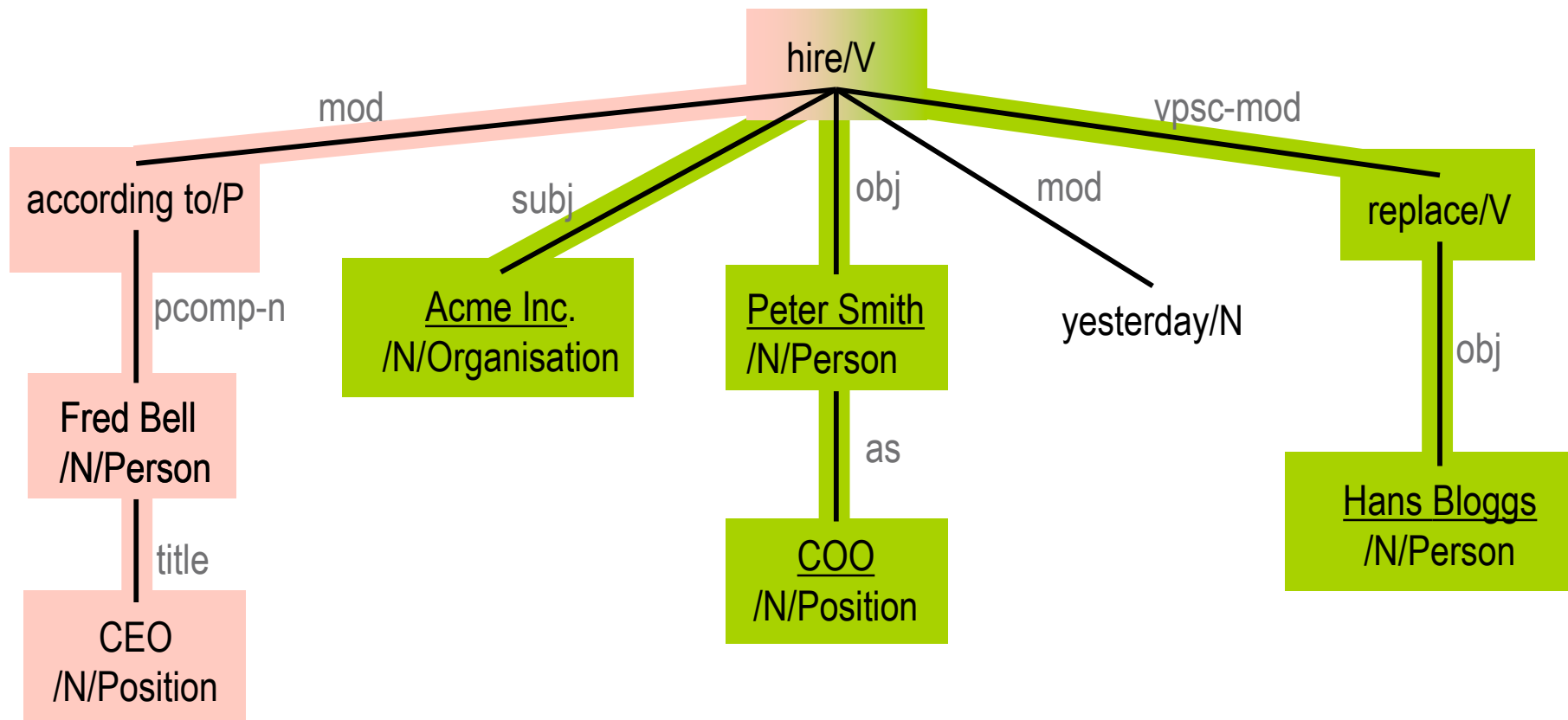
- ▣ verb centered
- ▣ All chains dominated by a verb, which contain at least one relevant named entity and their combinations



Previous Work: Subtree-Model

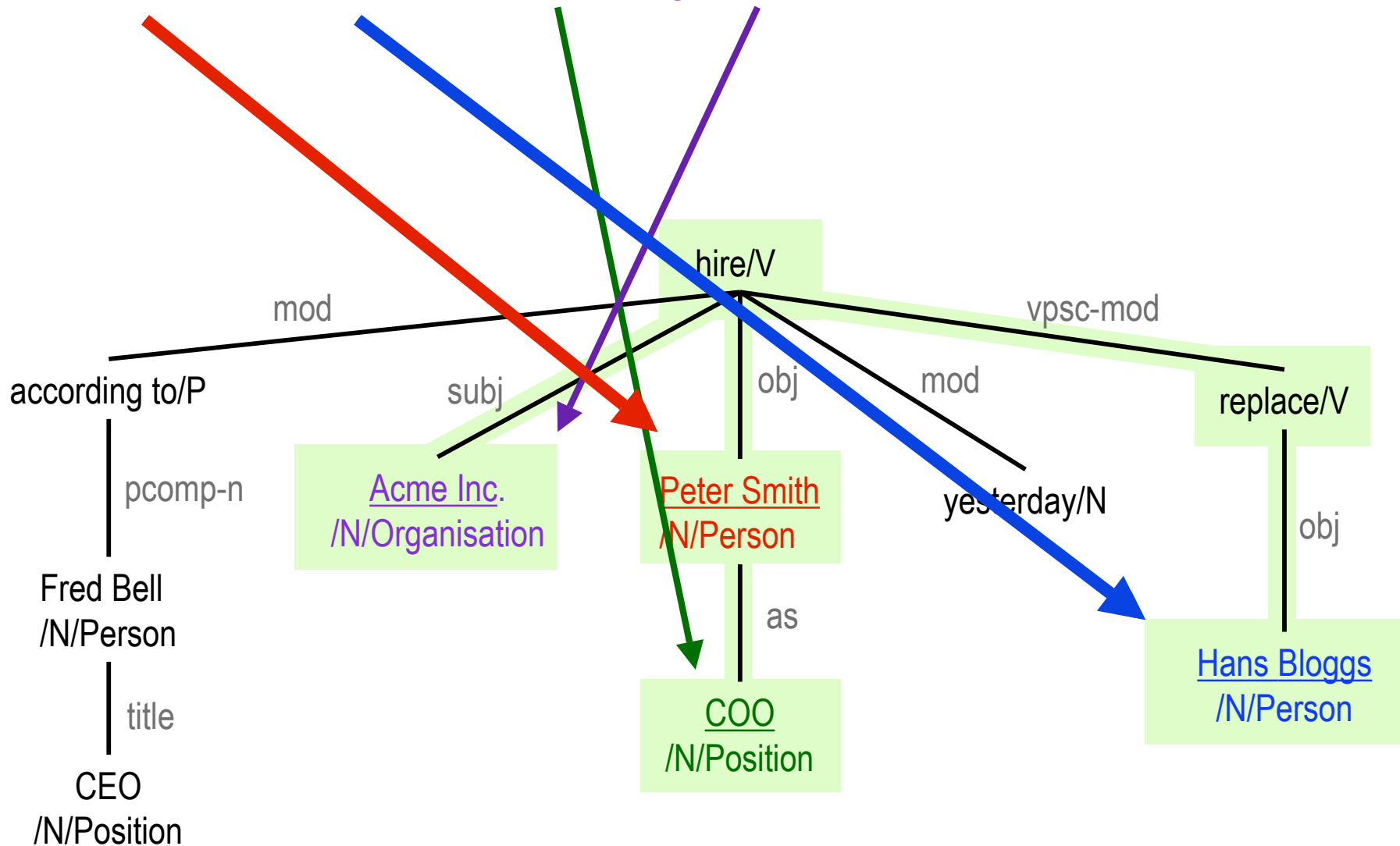
Sudo et al. (2003)

- 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>



- State of the art
- Domain Adaptive Relation Extraction Framework (**DARE**)
- Experiments and evaluations
- Performance analysis and discussion
- Conclusion and future work

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

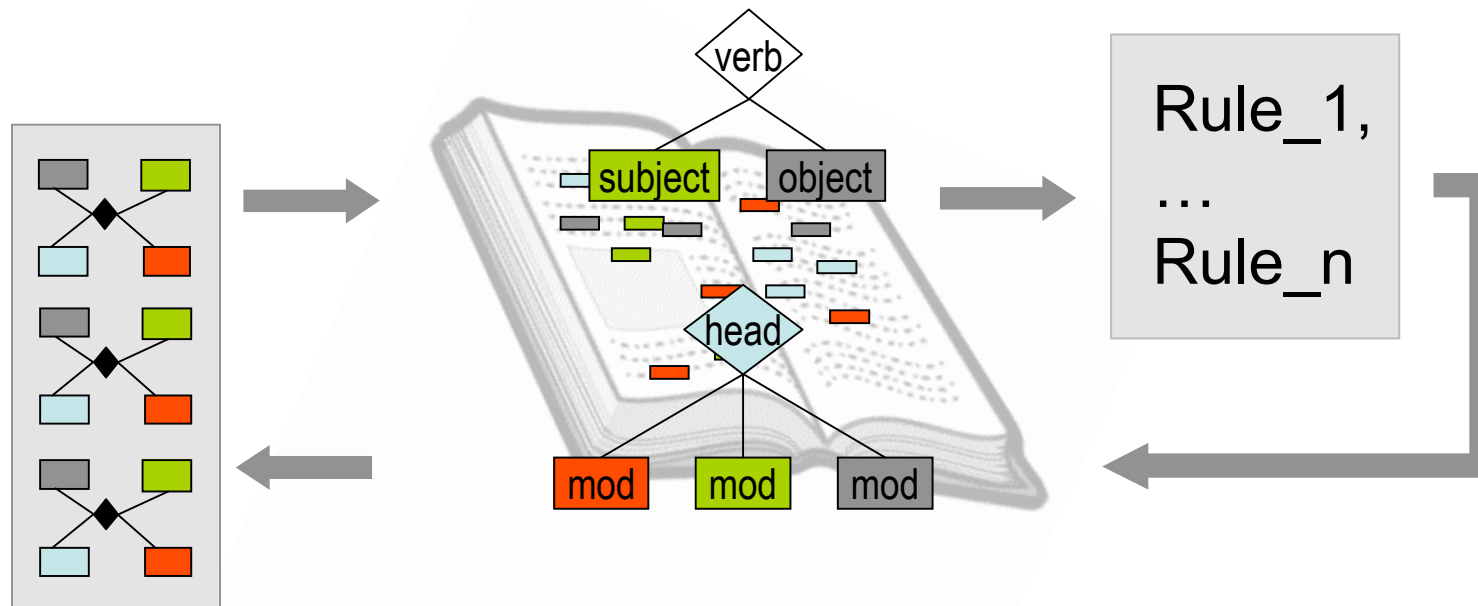
- Samples of target relation instances serve as semantic seed
- Systematic treatment of n-ary relations and their projections
- Exploitation of relation projections for pattern discovery
- Bottom-up compositional pattern discovery
- A recursive linguistic rule representation
- Rules contain semantic 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

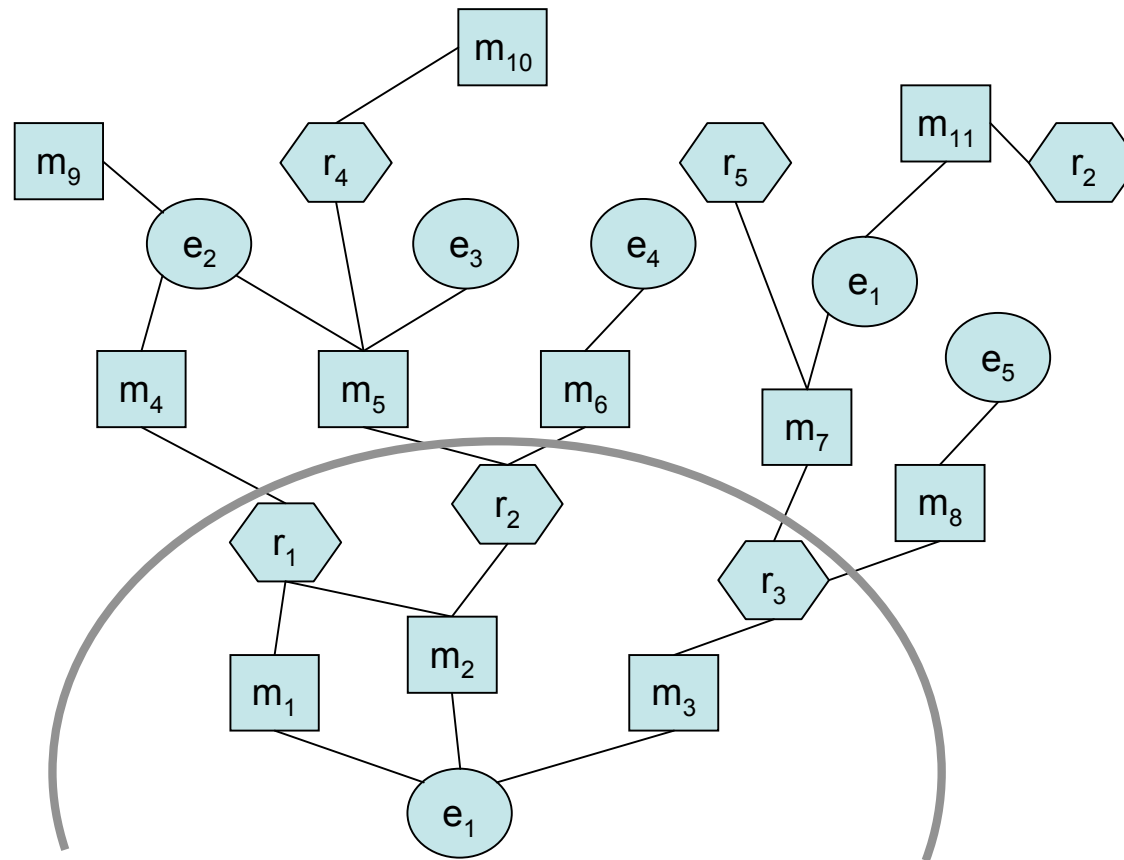
- binary relations only, no projections, no linguistic analysis

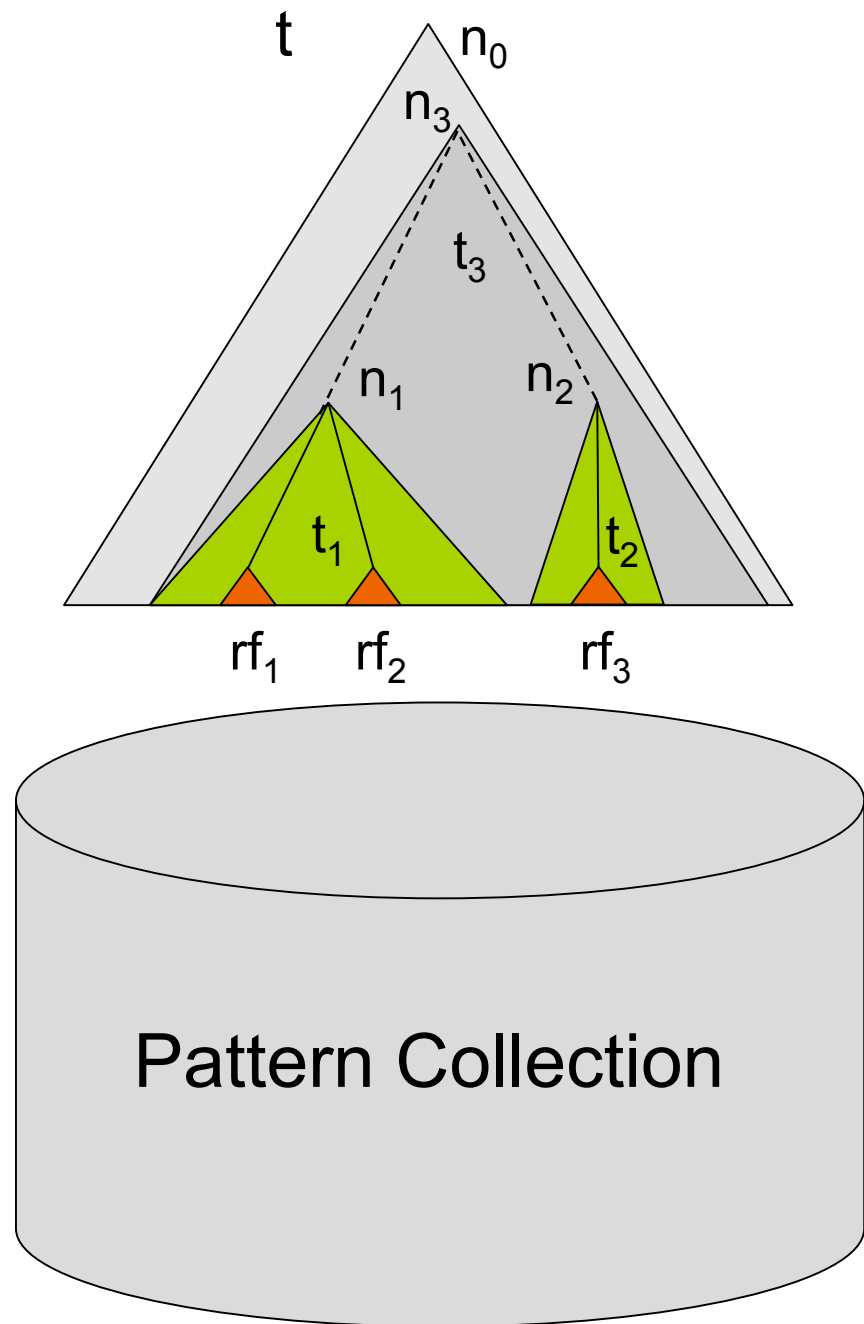
□ DARE

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

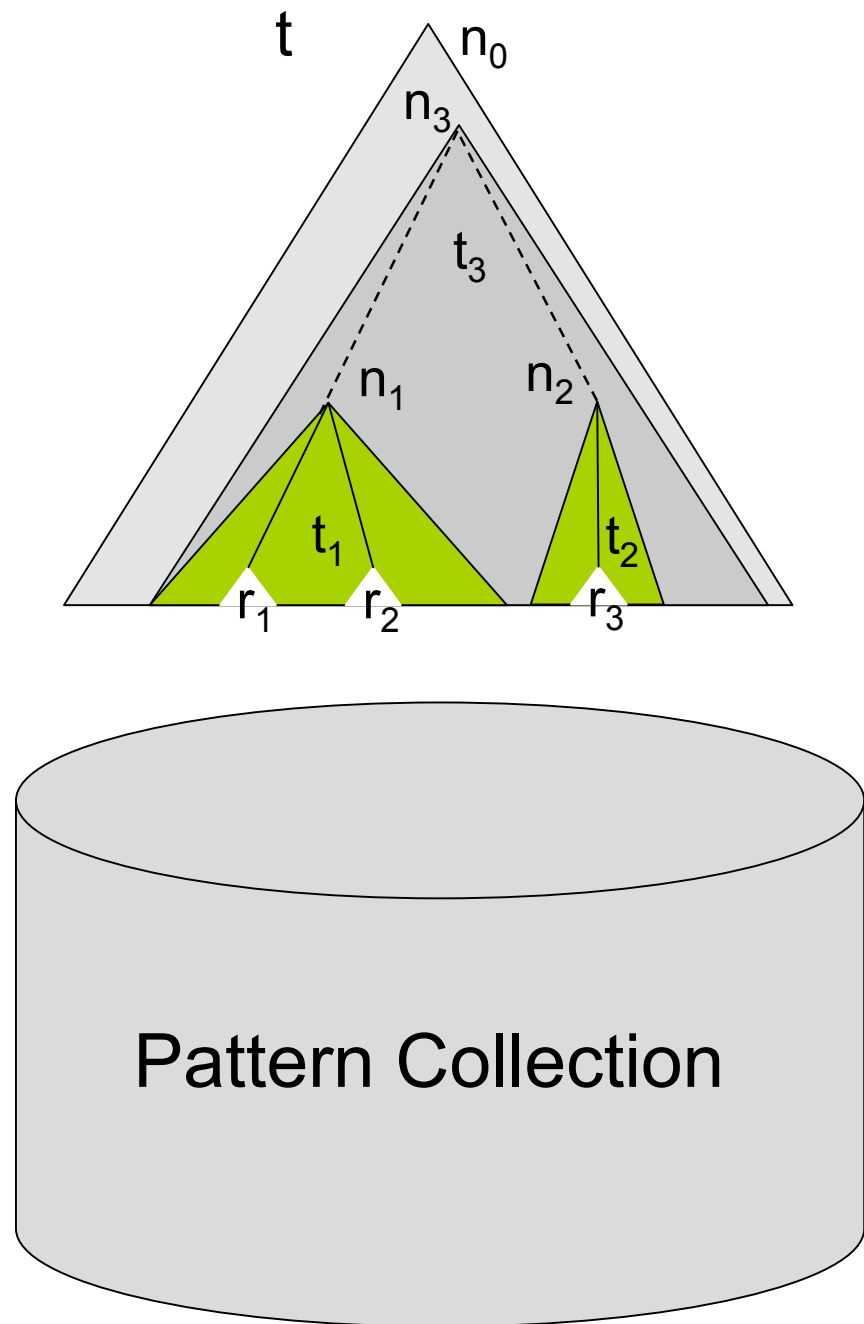
(in the experiments I use MINIPAR by Dekan Lin 1999)

Start of Bootstrapping (simplified)

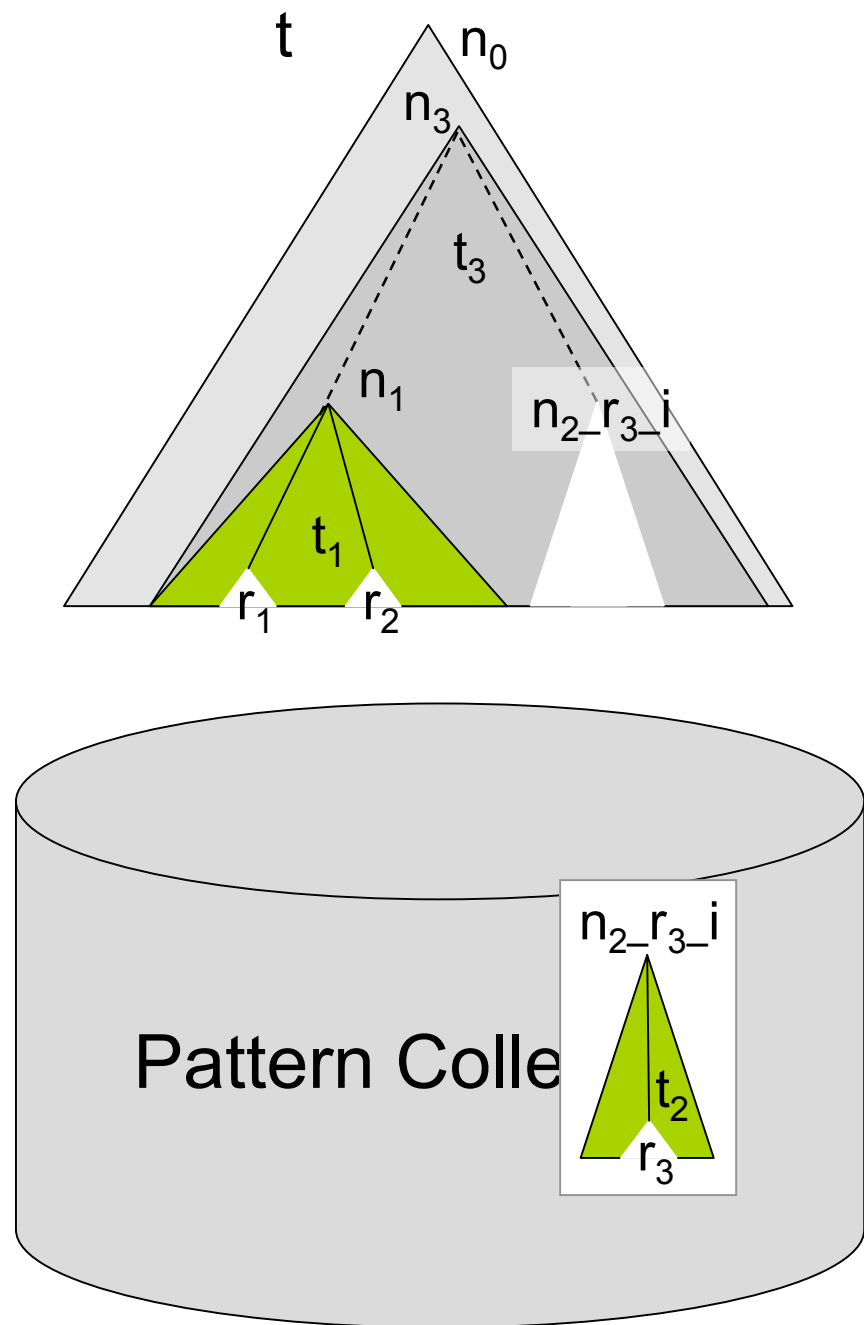




0. replace all nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;



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;



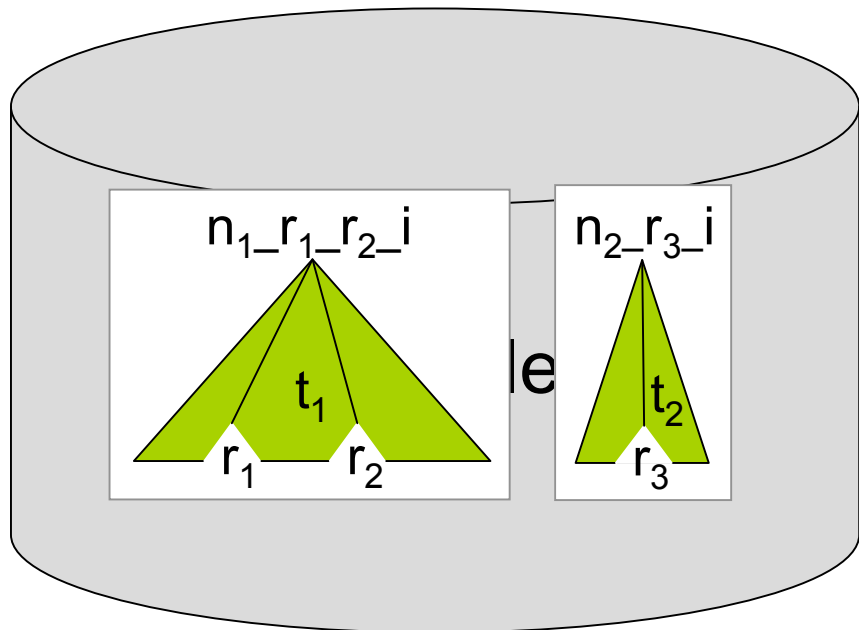
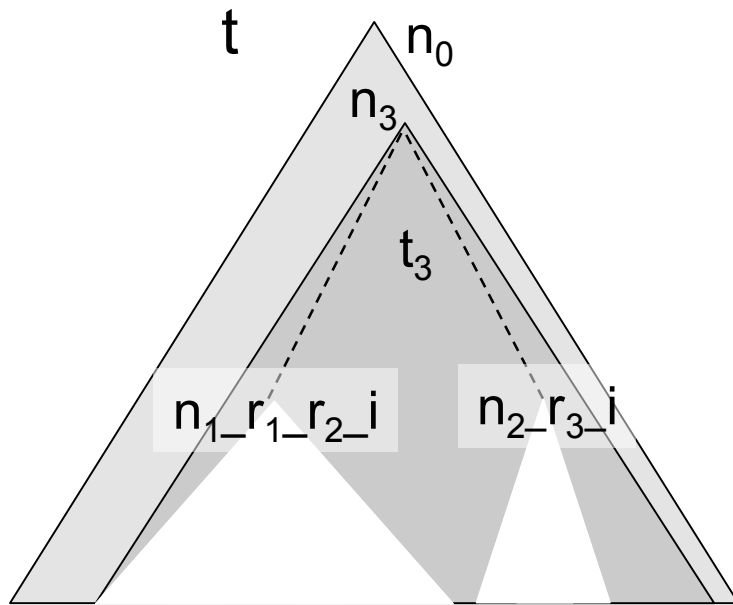
0. replace all nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;

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.



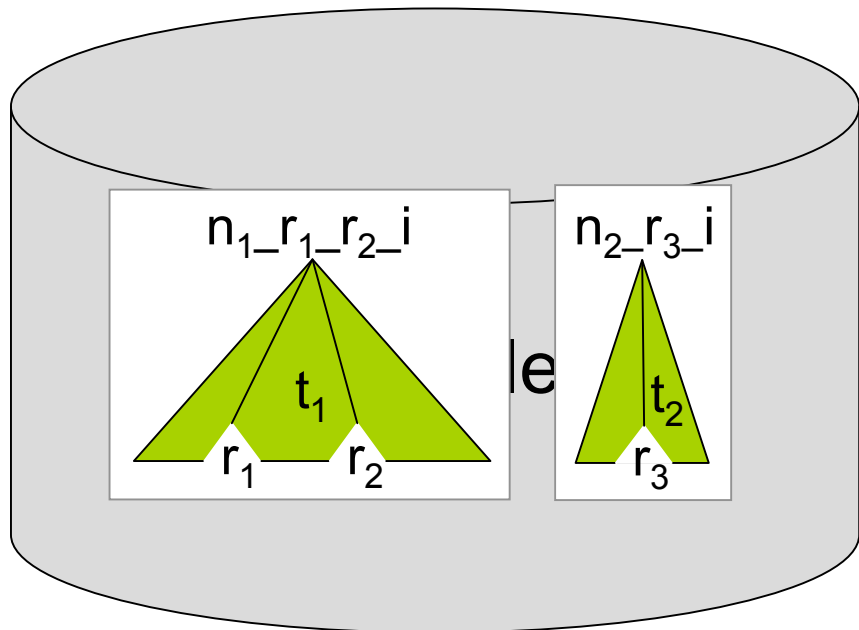
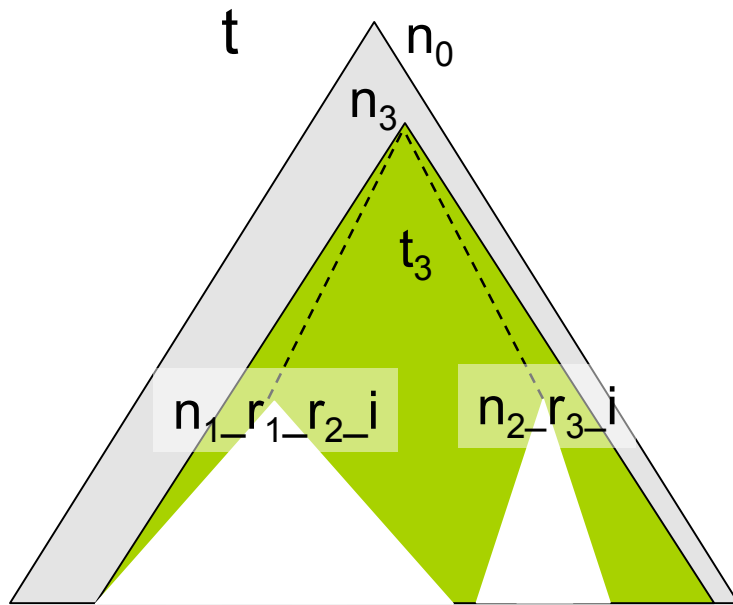
0. replace all nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;

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.



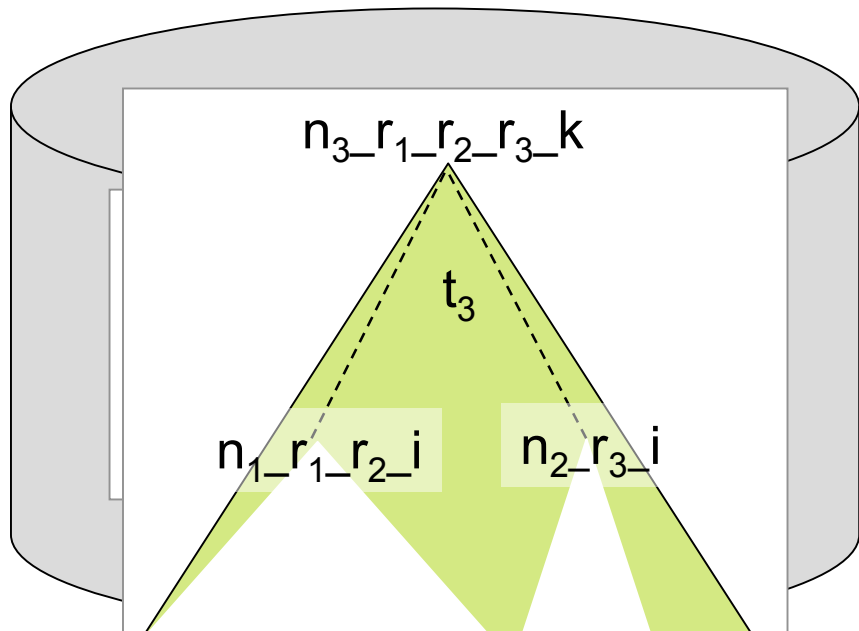
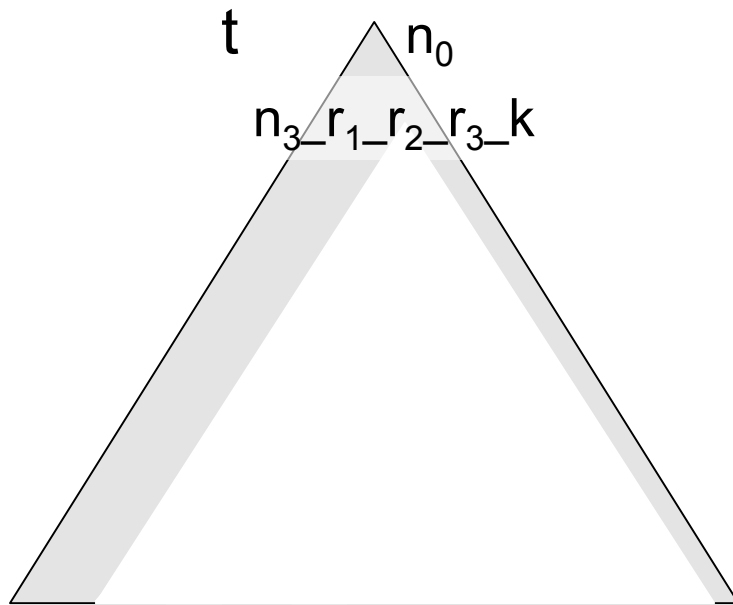
0. replace all nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;

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.



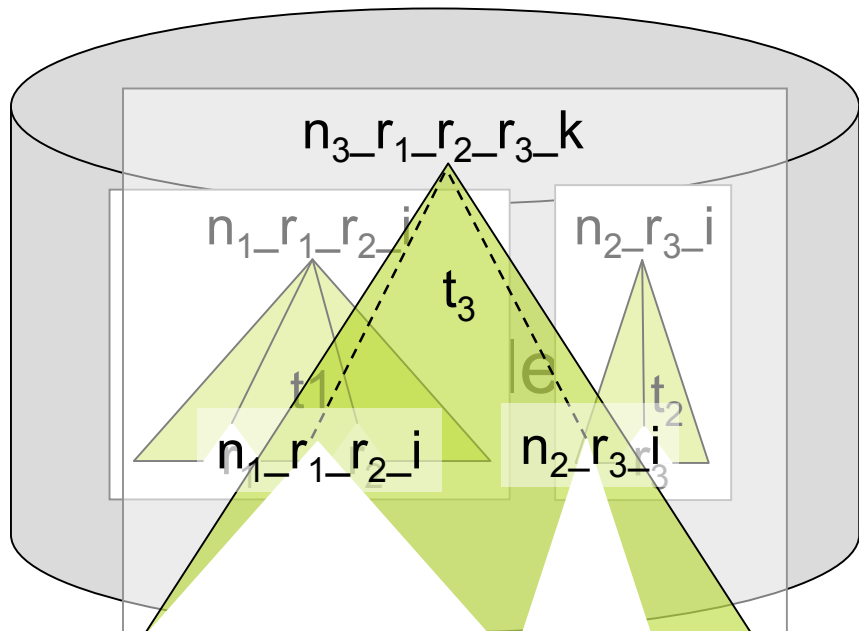
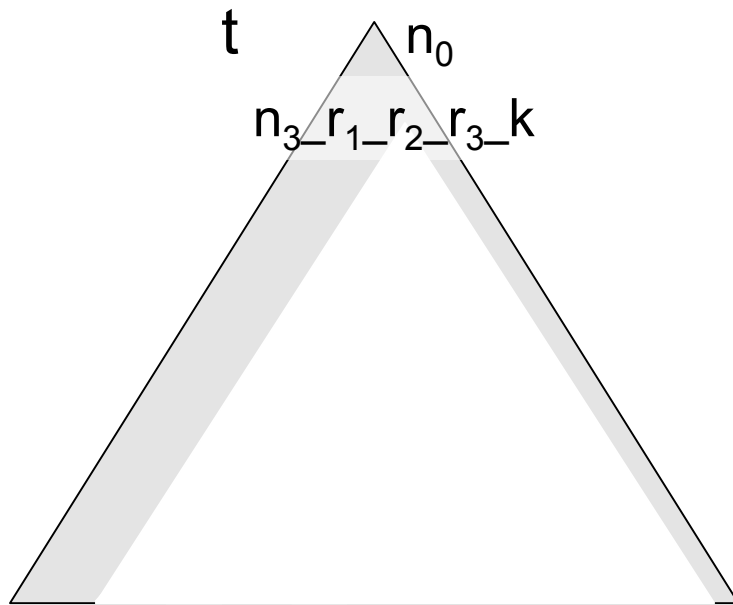
0. replace all nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;

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

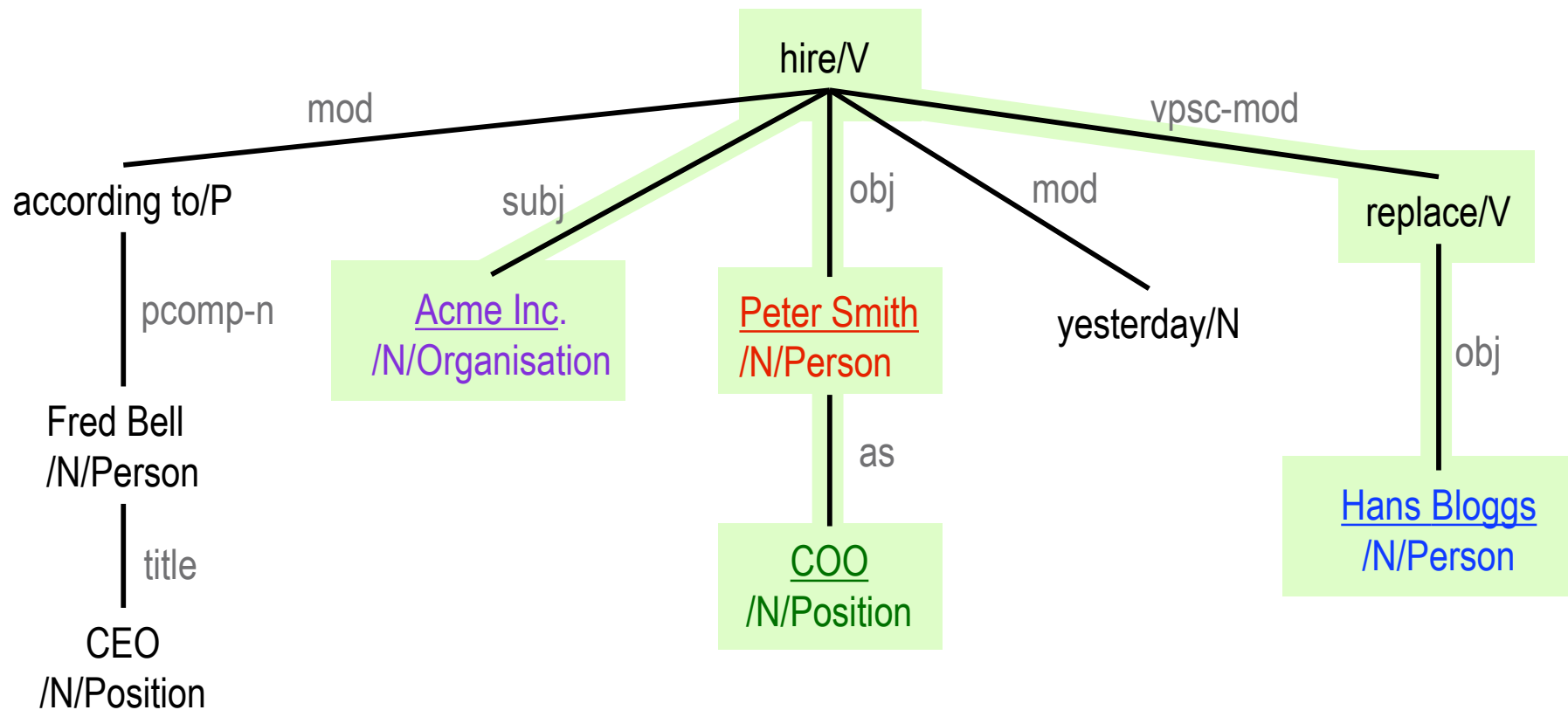
1. identify the set of the lowest non-terminal nodes N_1 in t that dominate i arguments (possibly among other nodes).

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.

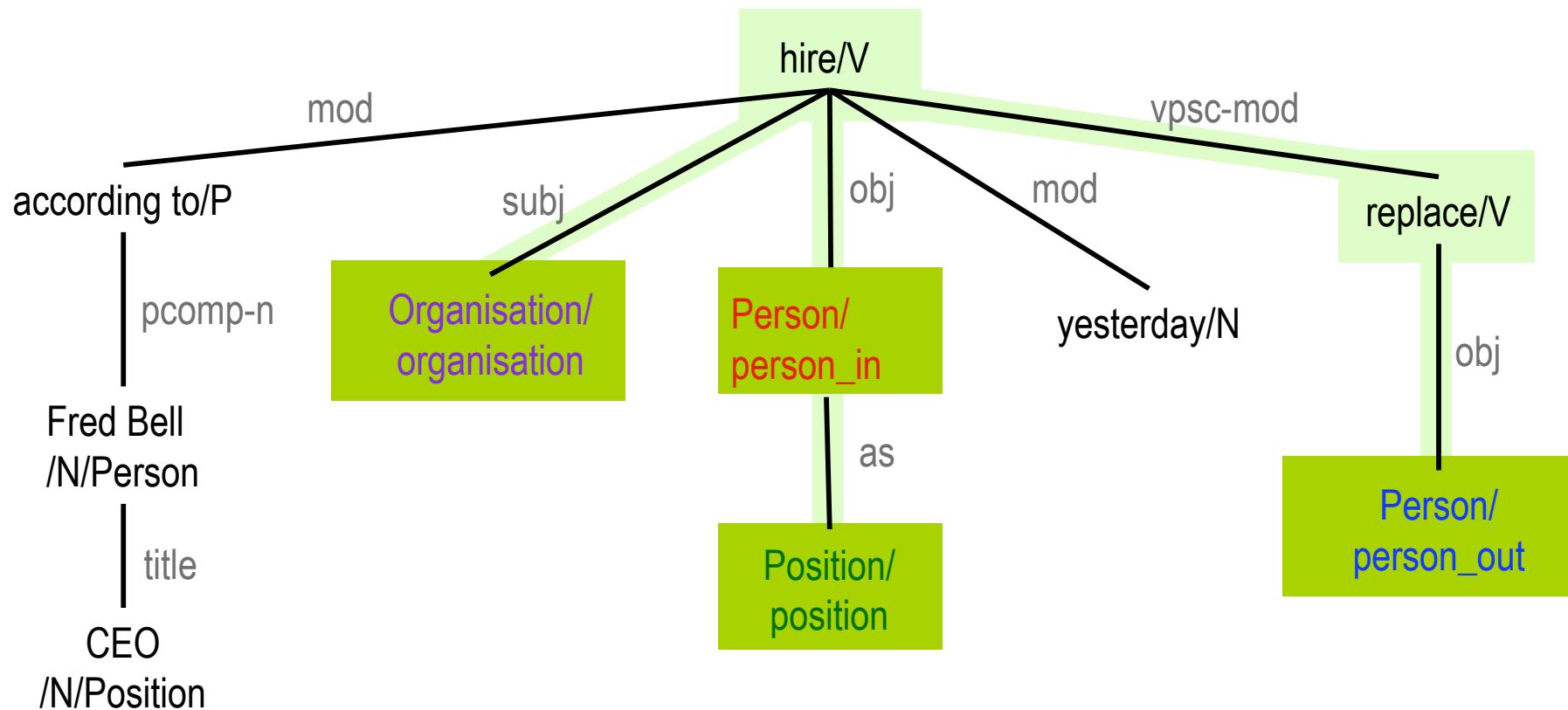
According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<Peter Smith/person_in, Hans Bloggs/person_out, COO /position, Acme Inc. /organisation>



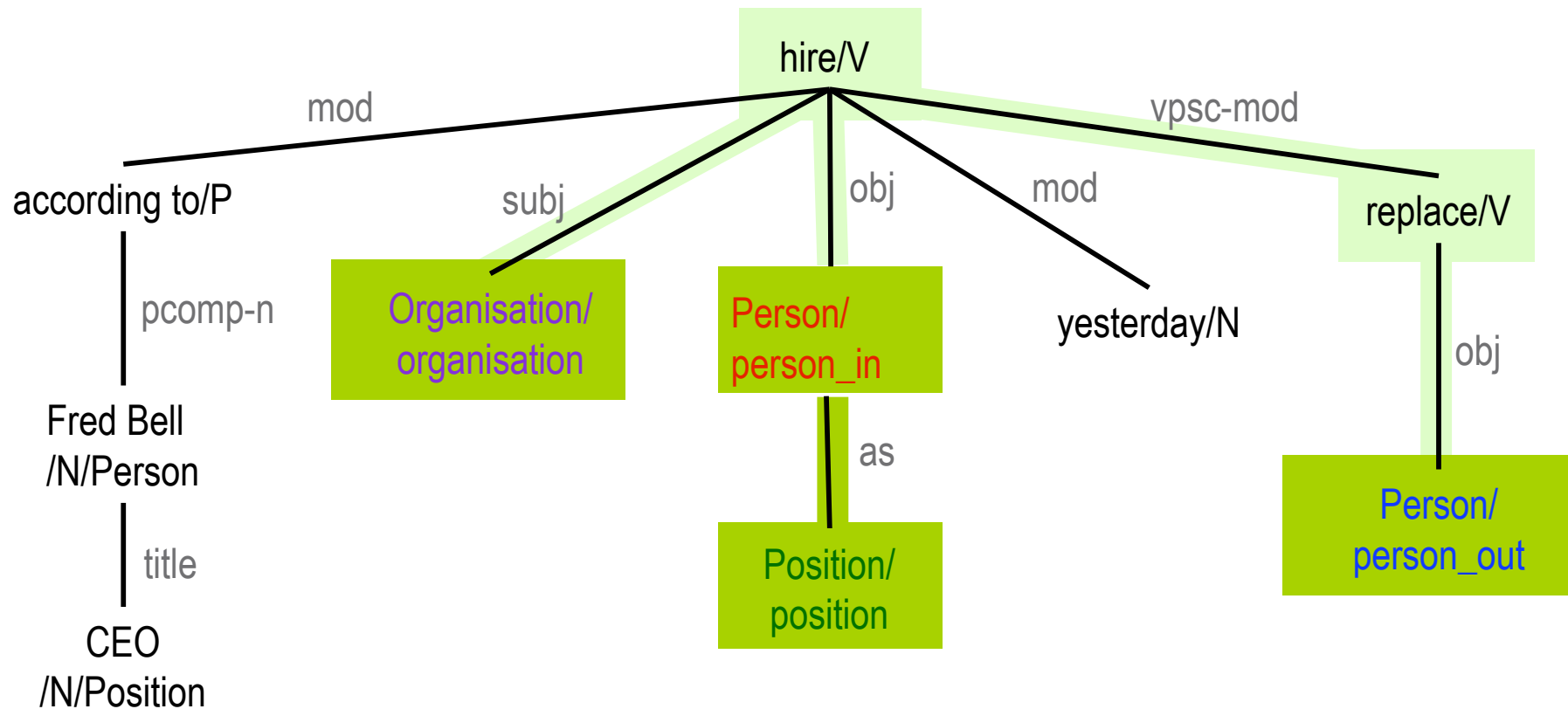
According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<Peter Smith/person_in, Hans Bloggs/person_out, COO /position, Acme Inc. /organisation>



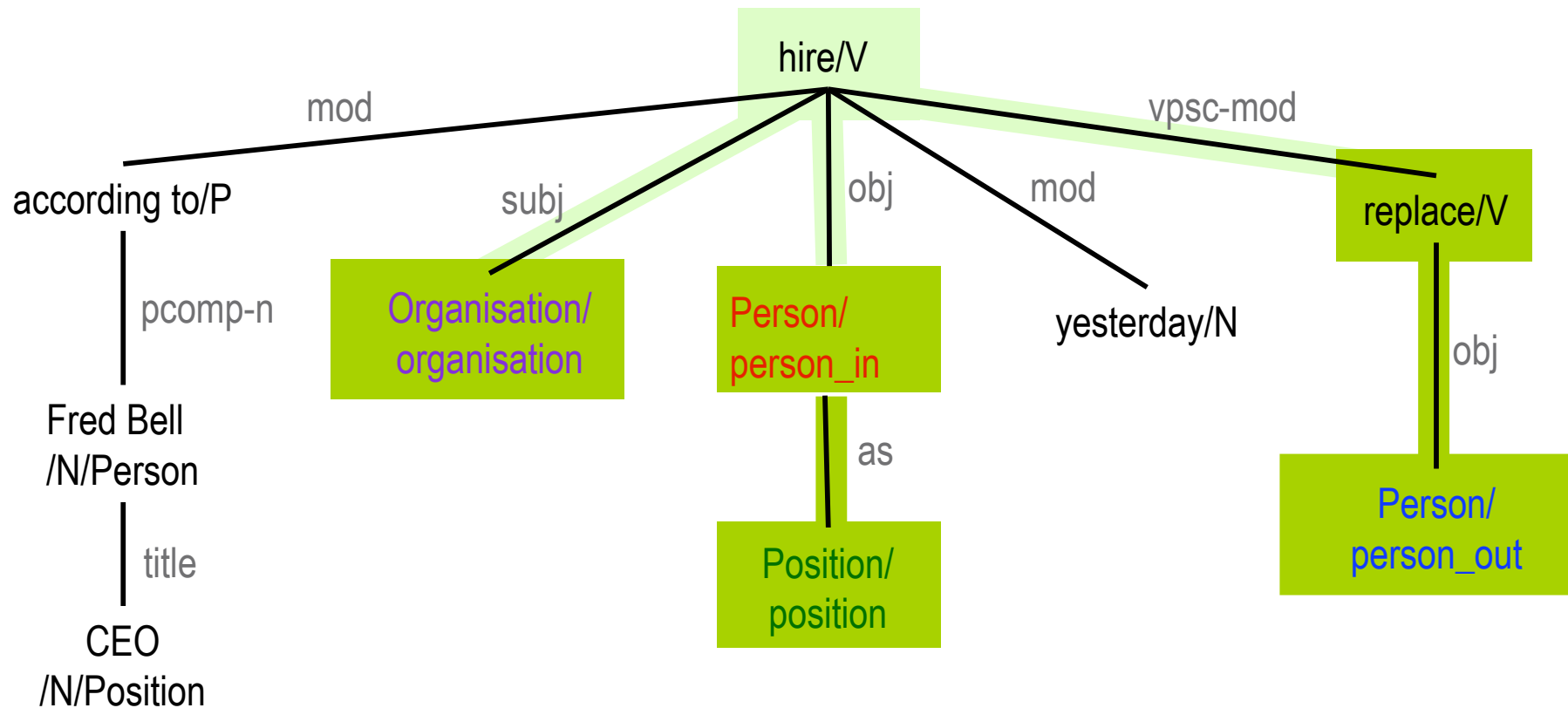
According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<Peter Smith/person_in, Hans Bloggs/person_out, COO /position, Acme Inc. /organisation>



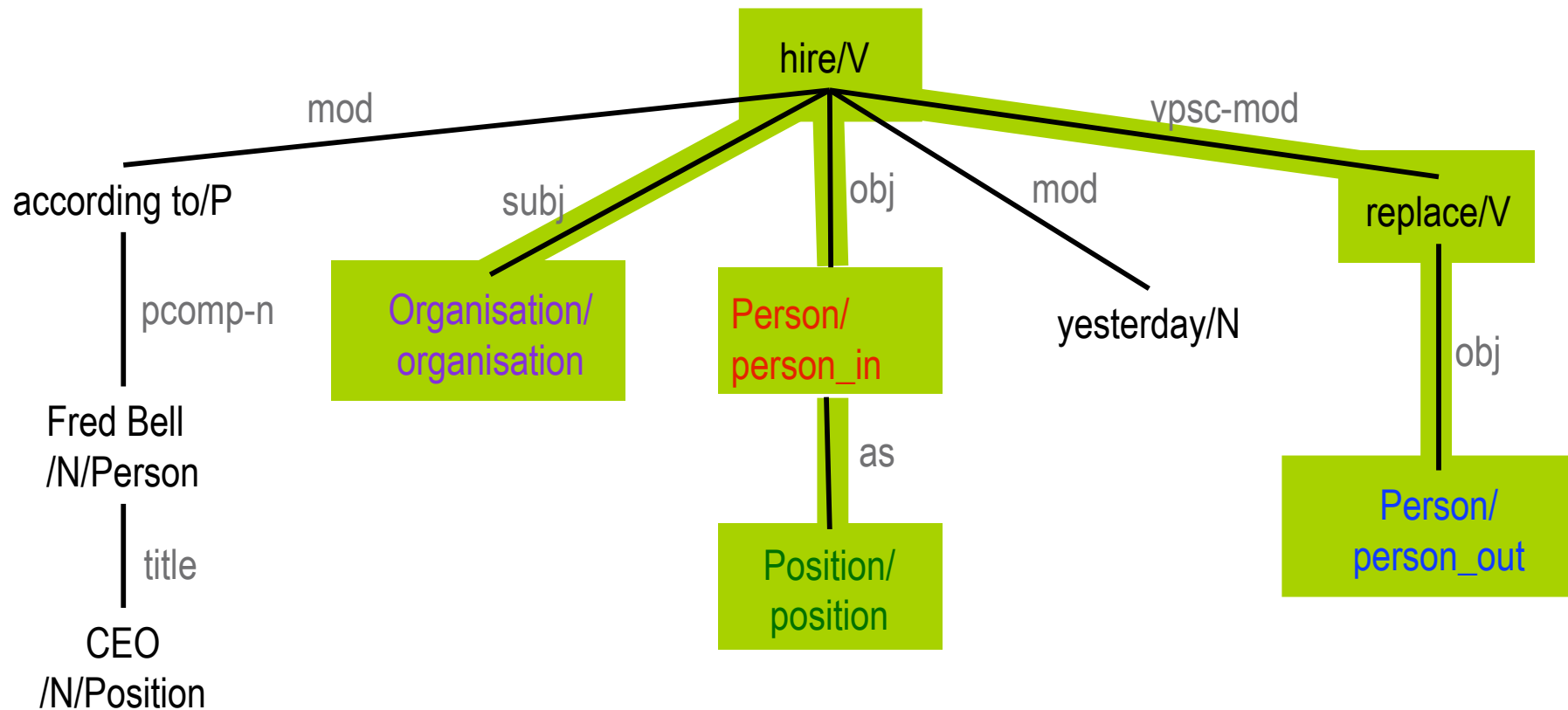
According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<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>



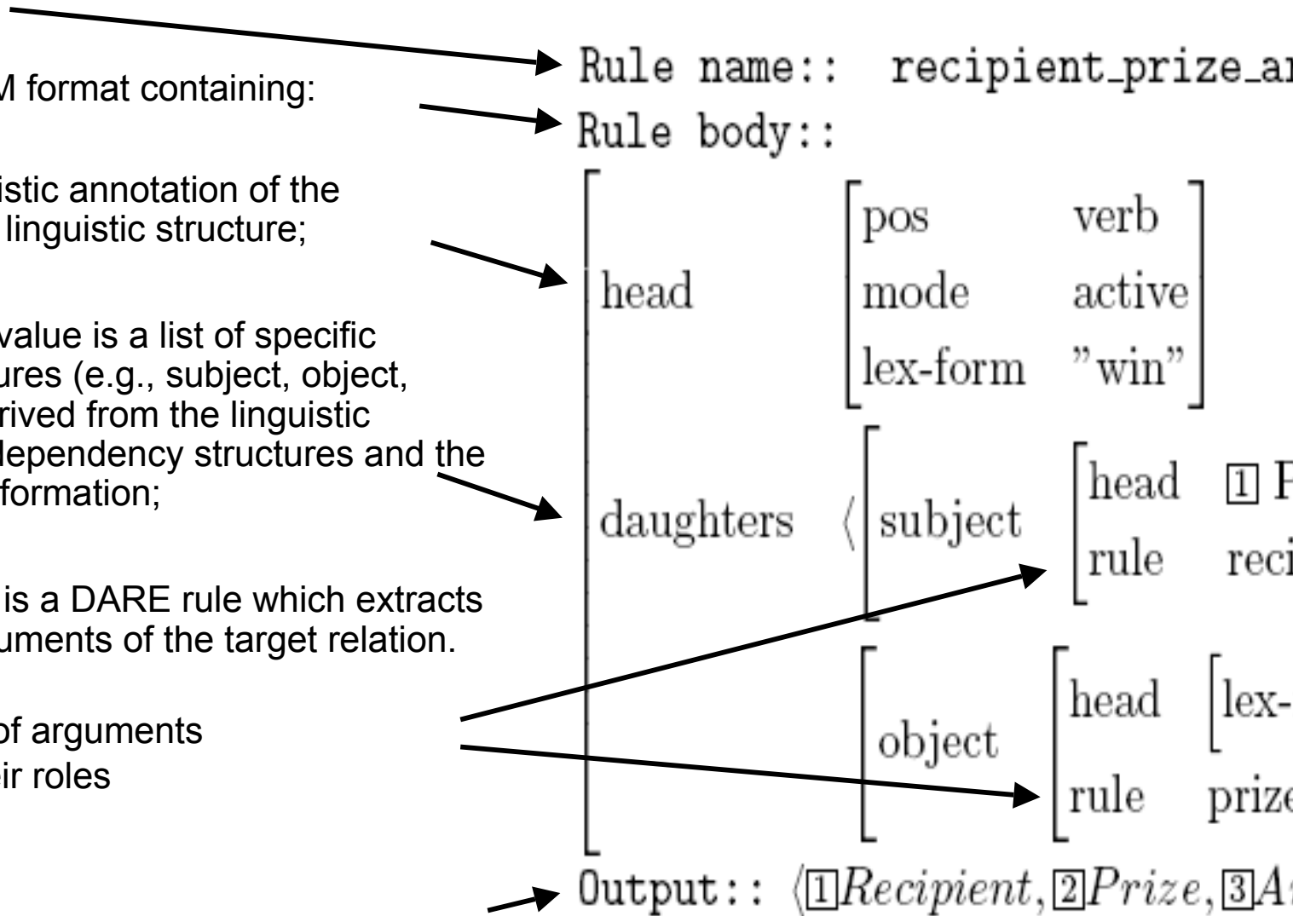
DARE Rule Components

1. rule name: r_i

2. rule body: in AVM format containing:

- **head**: the linguistic annotation of the top node of the linguistic structure;
- **daughters**: its value is a list of specific linguistic structures (e.g., subject, object, head, mod), derived from the linguistic analysis, e.g., dependency structures and the named entity information;
- **rules**: its value is a DARE rule which extracts a subset of arguments of the target relation.

3. **Output**: n-tupel of arguments with their roles



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>

- State of the art
- Domain **A**daptive **R**elation **E**xtraction **F**ramework (**DARE**)
- Experiments and evaluations
- Performance analysis and discussion
- Conclusion and future work

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

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

Evaluation Against Ideal Tables

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

Management Succession Domain

| Initial Seed # | Precision | Recall |
|----------------|--------------|--------------|
| 1 | 12.6% | 7.0% |
| 1 | 15.1% | 21.8% |
| 20 | 48.4% | 34.2% |
| 55 | 62.0% | 48.0% |

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%

- State of the art
- Domain **A**daptive **R**elation **E**xtraction **F**ramework (**DARE**)
- Experiments and evaluations
- Performance analysis and discussion
- Conclusion and future work

The Dream

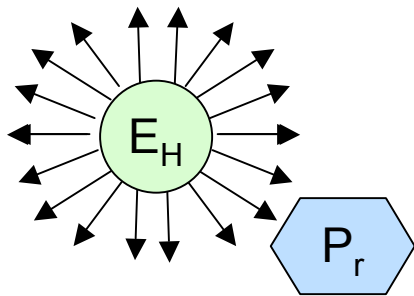
- Wouldn't it be wonderful if we could always automatically learn most or all relevant patterns of some relation from one single semantic instance!
- Or at least find all event instances.
- This sounds too good to be true!

Research Questions

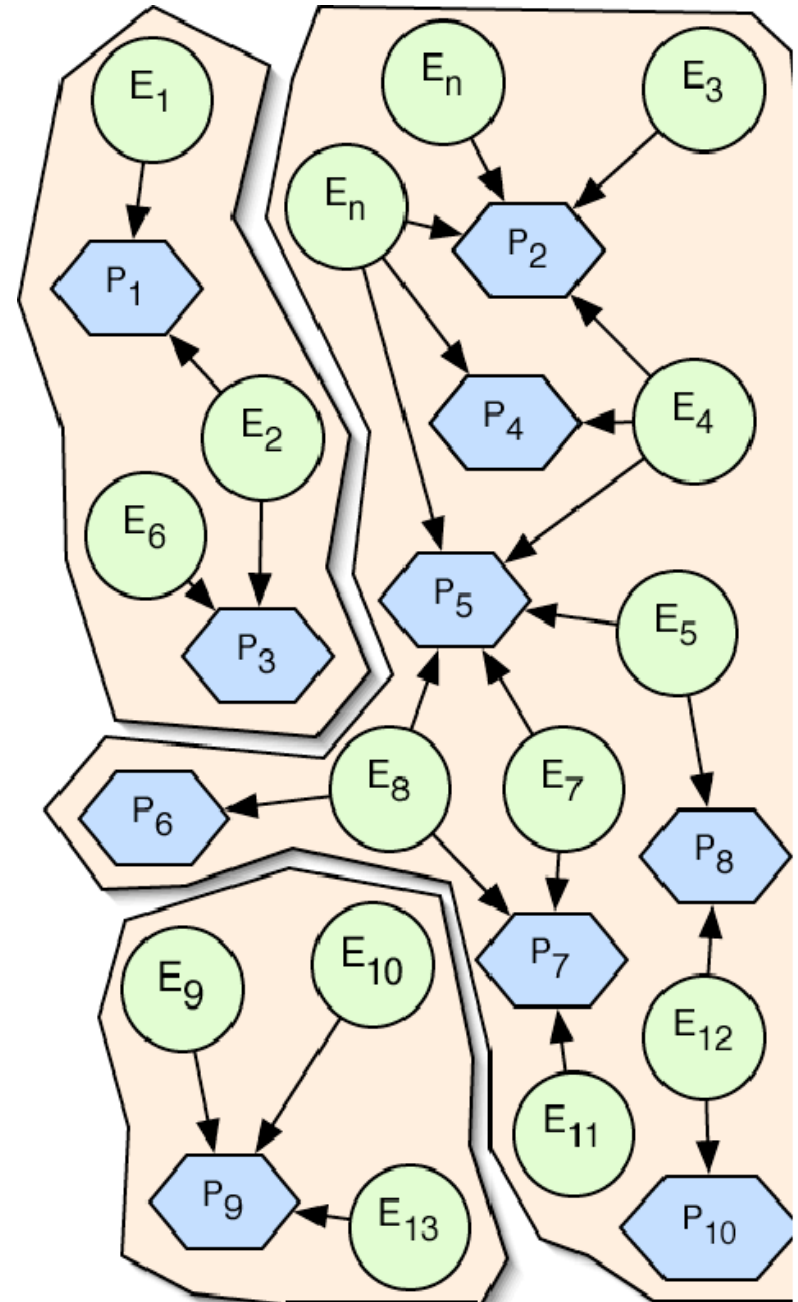
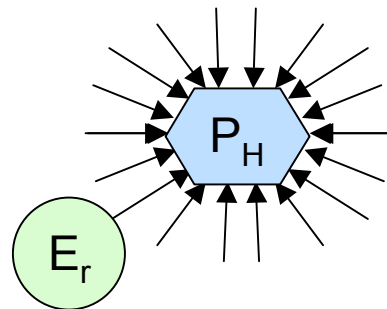
As scientists we want to know

- 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

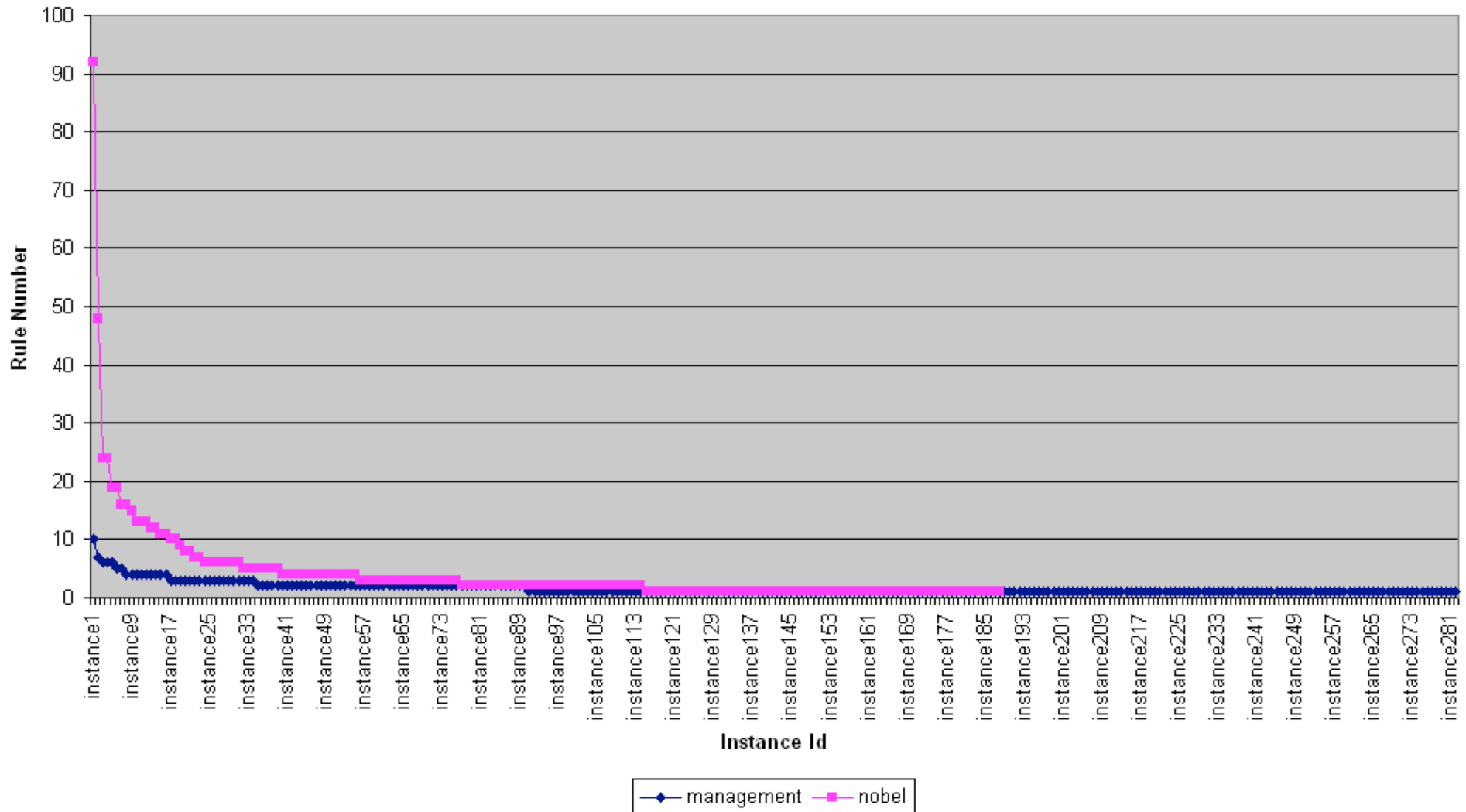


Popular patterns help to reach rarely mentioned events

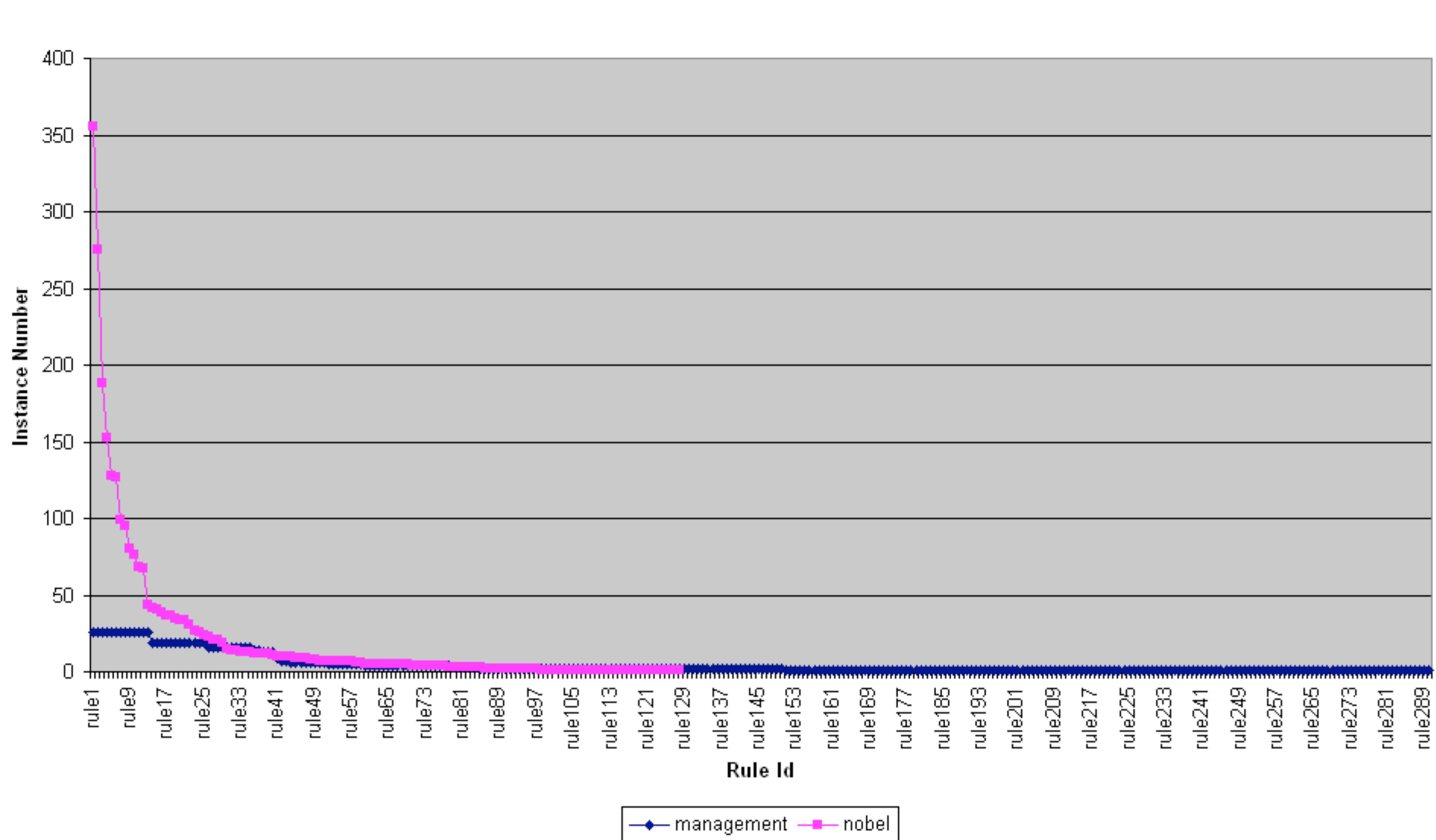


Instance to Pattern

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)

- State of the art
- Domain **A**daptive **R**elation **E**xtraction **F**ramework (**DARE**)
- Experiments and evaluations
- Performance analysis and discussion
- Conclusion and future work

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:
 - reusable generic linguistic knowledge
 - + raw data
 - + a few examples (seed)
 - = domain specific relation extraction grammar

- 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.

Future Work

- Improvement of recall
 - Extension of learning data
 - Bridging the islands by new additional data
 - Use of a related domain, e.g, Nobel Prize for other prizes
 - Improvement of rule generalization
 - Intersentential extraction

- Improvement of precision
 - Negative rules (domain independent and domain specific)
 - Integration of high-precision NLP analysis (HPSG)

References

1. N. Kushmerick. Wrapper induction: Efficiency and Expressiveness, Artificial Intelligence, 2000.
2. I. Muslea. Extraction Patterns for Information Extraction. AAI-99 Workshop on Machine Learning for Information Extraction.
3. Riloff, E. and R. Jones. Learning Dictionaries for Information Extraction by Multi-Level Bootstrapping. In Proceedings of the Sixteenth National Conference on Artificial Intelligence (AAAI-99) , 1999, pp. 474-479.
4. R. Yangarber, R. Grishman, P. Tapanainen and S. Huttunen. Automatic Acquisition of Domain Knowledge for Information Extraction. In Proceedings of the 18th International Conference on Computational Linguistics: COLING-2000, Saarbrücken.
5. F. Xu, H. Uszkoreit and Hong Li. Automatic Event and Relation Detection with Seeds of Varying Complexity. In Proceedings of AAAI 2006 Workshop Event Extraction and Synthesis, Boston, July, 2006.
6. F. Xu, D Kurz, J Piskorski, S Schmeier. A Domain Adaptive Approach to Automatic Acquisition of Domain Relevant Terms and their Relations with Bootstrapping. In Proceedings of LREC 2002.
7. W. Drozdowski, H.U. Krieger, J. Piskorski, U. Schäfer and F. Xu. Shallow Processing with Unification and Typed Feature Structures -- Foundations and Applications. In KI (Artificial Intelligence) journal 2004.
8. Feiyu Xu, Hans Uszkoreit, Hong Li. A Seed-driven Bottom-up Machine Learning Framework for Extracting Relations of Various Complexity. In Proceedings of ACL 2007, Prague
9. <http://www.dfki.de/~neumann/ie-essli04.html>