



HANS USZKOREIT & FEIYU XU 07

# Bootstrapping Relation Extraction Grammars from Semantic Seeds

Hans Uszkoreit & Feiyu Xu

DFKI and  
Saarland University

Thanks to Li Hong for her contributions to the experiments



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007



HANS USZKOREIT & FEIYU XU 07

## ★ Task and motivation

## ★ A new approach to seed-based learning for relation extraction

- Learning extraction rules for various complexity
- Experiments and evaluation

## ★ Scientific questions, insights and conclusion

- Seed-based learning in small and big worlds
- Lessons learned and outlook



German Research Center for Artificial Intelligence GmbH

T-FaNT ★ TOKYO UNIVERSITY ★ 13 MARCH 2007



HANS USZKOREIT & FEIYU XU 07

## Challenge

- ★ Development of a generic strategy for extracting relations/events of various complexity from large collections of open-domain free texts

## Central Motivation

- ★ Enable inexpensive adaptation to new relation extraction tasks/domains



German Research Center for Artificial Intelligence GmbH

T-FaNT ★ TOKYO UNIVERSITY ★ 13 MARCH 2007

## Existing Unsupervised or Minimally Supervised IE Approaches



HANS USZKOREIT & FEIYU XU 07

- ☆ Lack of expressiveness (Stevenson and Greenwood, 2006)
  - Restricted to a certain linguistic representation, mainly verb-centered constructions  
e.g., subject verb object construction (Yangarber, 2003)  
*subject(company)-verb("appoint")-object(person)*
  - other linguistic constructions can not be discovered: e.g., apposition, compound NP  
*the 2005 Nobel Peace Prize*
- ☆ Lack of semantic richness (Riloff, 1996; Agichtein and Gravano, 2000; Yangarber, 2003, Greenwood and Stevenson, 2006)
  - Pattern rules cannot assign semantic roles to the arguments  
*subject(person)-verb("succeed")-object(person)*
- ☆ No good method to select pattern rules, in order to deal with large number of tree patterns (Sudo et al., 2003)
- ☆ No systematic way to handle relations and their projections
  - do not consider the linguistic interaction between relations and their projections, which is important for scalability and reusability of rules



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

## Two Approaches to Seed Construction by Bootstrapping



HANS USZKOREIT & FEIYU XU 07

- ★ Pattern-oriented (e.g., ExDisco (Yangarber 2001))
  - too closely bound to the linguistic representation of the seed, e.g.,  
*subject(company) v("appoint") object(person)*
  - An event can be expressed by more than one pattern and by various linguistic constructions
  
- ★ Relation and event instances as seeds (e.g., DIPRE (Brin 1998) and Snowball (Agichtein and Gravano 2000), (Xu et al. 2006))
  - domain independence: it can be applied to all relation and event instances
  - flexibility of the relation and event complexity: it allows n-ary relations and events
  - processing independence: the seeds can lead to patterns in different processing modules, thus also supporting hybrid systems, voting approaches etc.
  - Not limited to a sentence as an extraction unit



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

## Our Approach: *DARE* (1)



HANS USZKOREIT & FEIYU XU 07

- ★ seed-driven and bottom-up rule learning in a bootstrapping framework
  - starting from sample relation instances as seeds
    - complexity of the seed instance defines the complexity of the target relation
  - – pattern discovery is bottom-up and compositional, i.e., complex patterns are derived from simple patterns for relation projections
  - – bottom-up compression method to cluster and generalize rules
  - – only subtrees containing seed arguments are pattern candidates
  - pattern rule ranking and filtering method considers two aspects of a pattern
    - its domain relevance and
    - the trustworthiness of its origin



German Research Center for Artificial Intelligence GmbH

T-FaNT ★ TOKYO UNIVERSITY ★ 13 MARCH 2007



HANS USZKOREIT & FEIYU XU 07

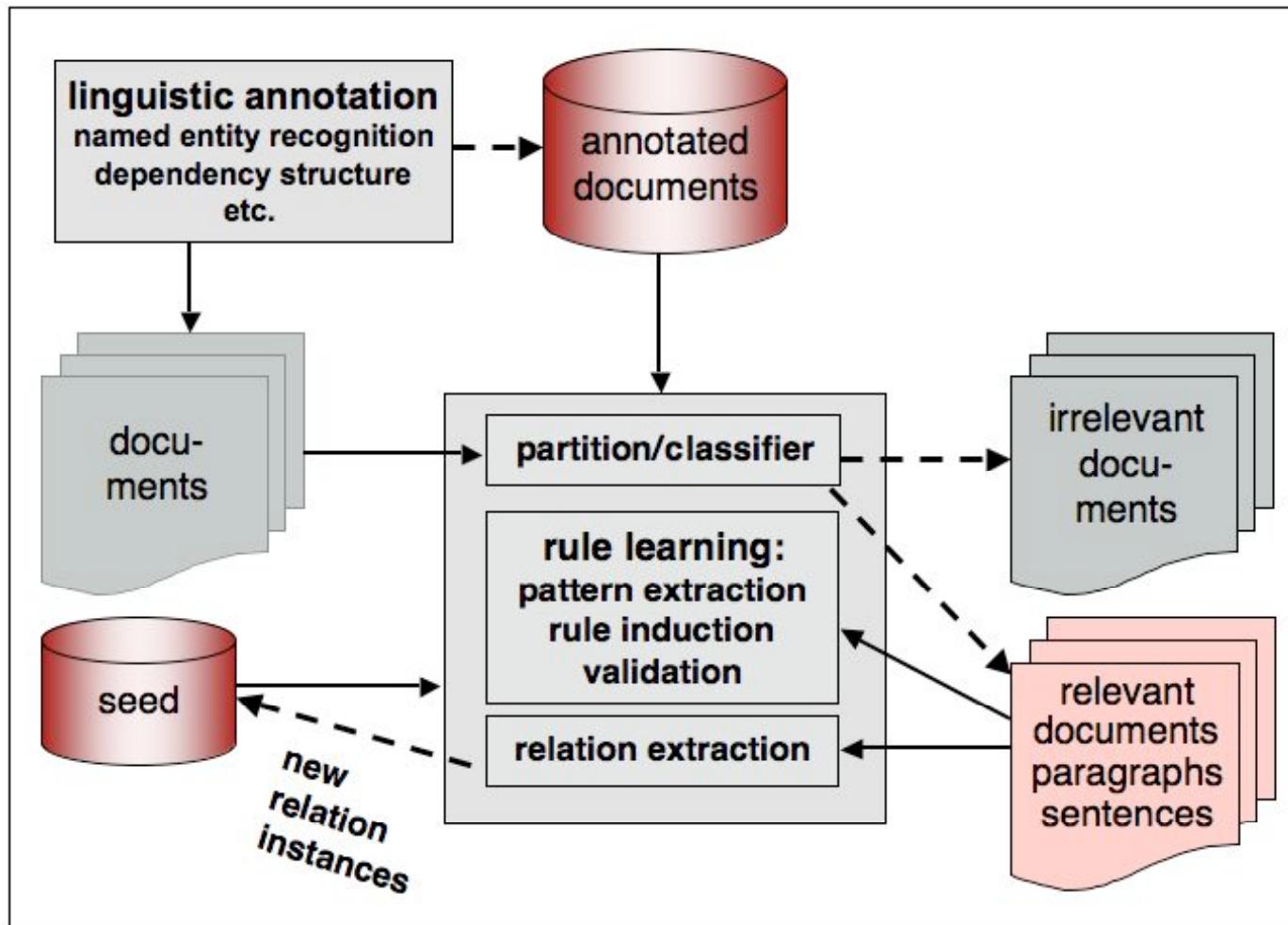
### → Compositional rule representation model

- support the bottom-up rule composition
- expressive enough for the representation of rules for various complexity
- precise assignment of semantic roles to the slot arguments
- reflects the precise linguistic relationship among the relation arguments and reduces the template merging task in the later phase
- the rules for the subset of arguments (projections) may be reused for other relation extraction tasks.

# DARE System Architecture



HANS USZKOREIT & FEIYU XU 07



# Algorithm



HANS USZKOREIT & FEIYU XU 07

1. Given
  - A large corpus of un-annotated and un-classified documents
  - A trusted set of relation or event instances, initially chosen ad hoc by the user, the seed, normally, one or two.
2. NLP annotation
  - Annotate the relevant documents with named entities and dependency structures
3. Partition
  - Apply seeds to the documents and divide them into relevant and irrelevant documents  
A document is relevant, if its text fragments contain a minimal number of relation arguments of a seed
  - Paragraph/sentence retrieval
4. Rule learning
  - Extract patterns
  - Rule induction/compression
  - Rule validation
5. Apply induced rules to the same document set
6. Rank new seeds
7. Stop if no new rules and seeds can be found, else repeat 3-6



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007



HANS USZKOREIT & FEIYU XU 07

## ★ Target relation

$\langle \text{recipient}, \text{prize}, \text{area}, \text{year} \rangle$

## ★ Example

*Mohamed ElBaradei won the 2005 Nobel Peace Prize on Friday for his efforts to limit the spread of atomic weapons.*



German Research Center for Artificial Intelligence GmbH

T-FaNT ★ TOKYO UNIVERSITY ★ 13 MARCH 2007

## Example Rules



HANS USZKOREIT & FEIYU XU 07

Rule name:: prize\_area\_year\_1

Rule:: ( $\exists$  year)( $\exists$  prizename) ( $\exists$  areaname) 'Prize'

Output::  $\langle \exists Prize, \exists Area, \exists Year \rangle$

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

Rule:: [verb [mode active  
wordform "win"]  
subject [recipient  $\exists$  Person  
rule recipient\_1::  $\langle \exists Person \rangle$ ]  
object [prize  $\exists$  Prize  
area  $\exists$  Area  
year  $\exists$  Year  
rule prize\_area\_year\_1::  $\langle \exists Prize, \exists Area, \exists Year \rangle$ ]]]

Output::  $\langle \exists Recipient, \exists Prize, \exists Area, \exists Year \rangle$



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

# Rule Interaction



HANS USZKOREIT & FEIYU XU 07

*Mohamed ElBaradei won the 2005 Nobel Peace Prize on Friday  
for his efforts to limit the spread of atomic weapons*

☆ **prize\_area\_year\_1:**

extracts a ternary projection instance  $\langle \text{prize}, \text{area}, \text{year} \rangle$  from  
a noun phrase compound

☆ **recipient\_prize\_area\_year\_1:**

triggers **prize\_area\_year\_1** in its object argument and extracts  
all four arguments.



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

# Rule Components



HANS USZKOREIT & FEIYU XU 07

1. rule name:  $r_i$ ;
2. output: a set  $A$  containing the  $n$  arguments of the  $n$ -ary relation, labelled with their argument roles;
3. rule body in AVM format containing:
  - a possibly empty set  $R_i$  of **DARE** rules, each of which extracts some proper subset of  $A$ ;
  - a possibly empty set of constraints  $C_i$  defining which functional arguments in  $r_i$  call which rules in  $R_i$ ;
  - rule-specific linguistic labels (e.g., **subject**, **object**, **head**, **mod**), derived from the linguistic analysis.



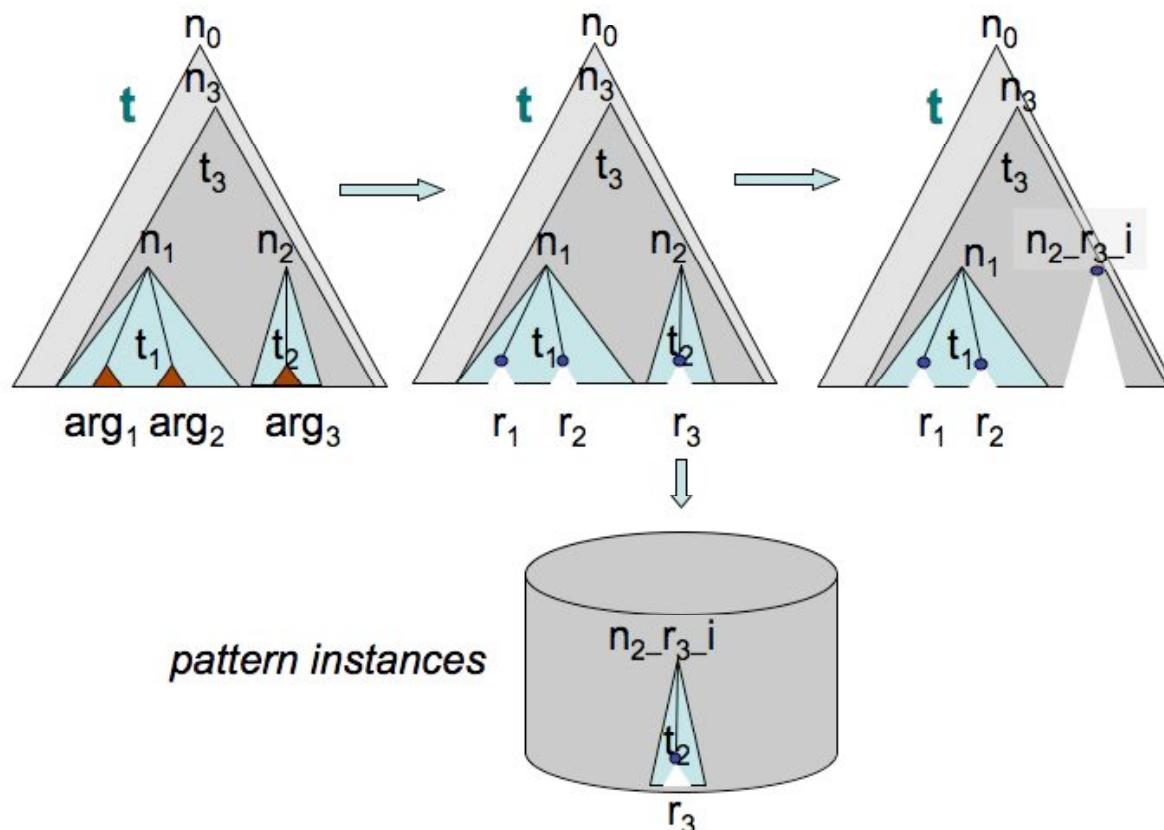
German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

# Pattern Extraction Step 1



HANS USZKOREIT & FEIYU XU 07



1. replace all terminal nodes that are instantiated with the seed arguments by new nodes. Label these new nodes with the seed argument roles and their entity classes;
2. identify the lowest nonterminal nodes  $N_1$  in  $t$  that dominate at most one argument (possibly among other nodes).
3. substitute  $N_1$  by nodes labelled with the seed argument roles and their entity classes
4. prune the subtrees dominated by  $N_1$  from  $t$  and add these subtrees into  $P$ . These subtrees are assigned the argument role information and a unique id.



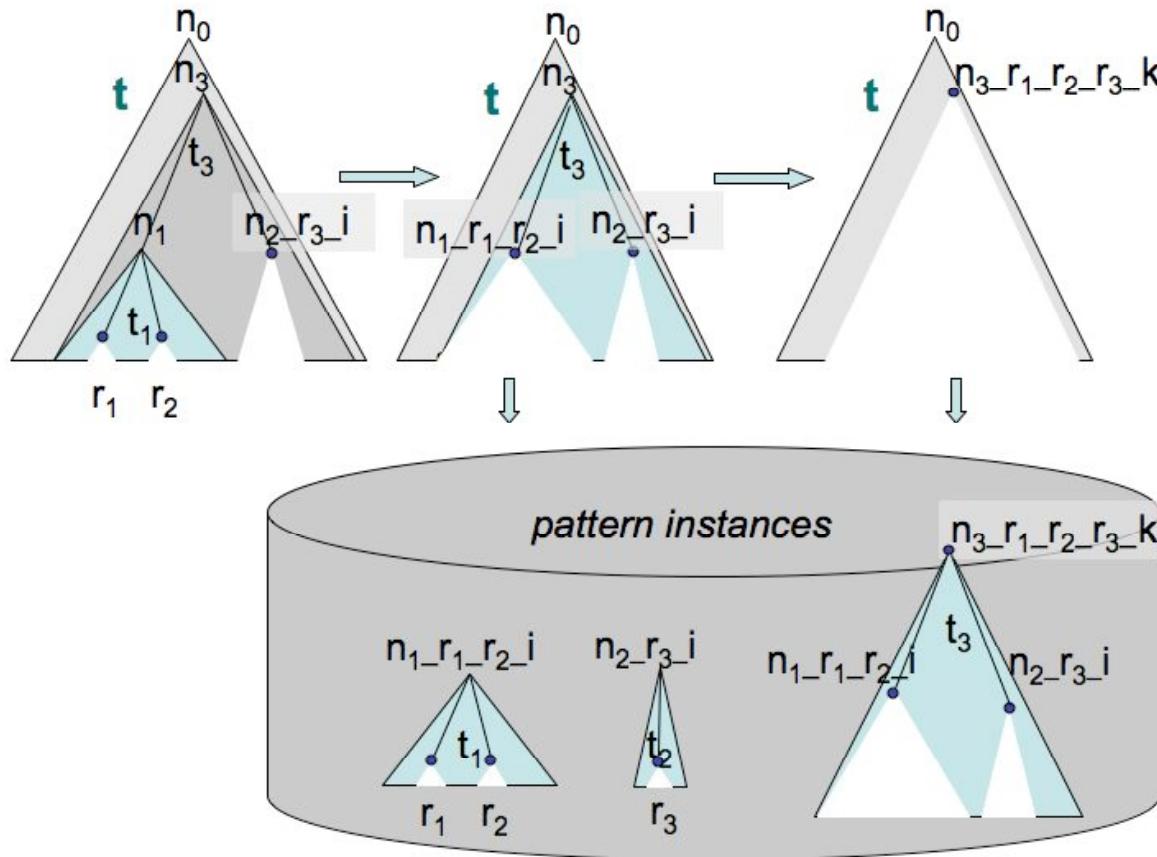
German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

# Pattern Extraction Step 2



HANS USZKOREIT & FEIYU XU 07



For  $i=2$  to  $n$

1. find the lowest nodes  $N_i$  in  $t$  that dominate in addition to other children only  $i$  seed arguments;
2. substitute  $N_i$  by nodes labelled with the  $i$  seed argument role combination information (e.g.,  $r_i-r_i$ ) and with a unique id.
3. prune the subtrees  $T_i$  dominated by  $N_i$  in  $t$ ;
4. add  $T_i$  together with the argument, role combination information and the unique id to  $P$

# Event Instance as Seed



HANS USZKOREIT & FEIYU XU 07

Here a relation-seed is a quadruple of 4 entity types

*event* &

Prize Name : *prize\_name*  
Prize Area : *area\_name*  
Recipient List : list of *person*  
Year: *year*

Examples in xml

```
<seed id="1">
<prize name="Nobel"/>
<year>1999</year>
<area name="chemistry"/>
<recipient>
  <person>
    <name>Ahmed H. Zewail</name>
    <surname>Zewail</surname>
    <gname>Ahmed</gname>
    <gname>H</gname>
  </person>
</recipient>
</seed>
```



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

# Sentence Analysis and Pattern Identification

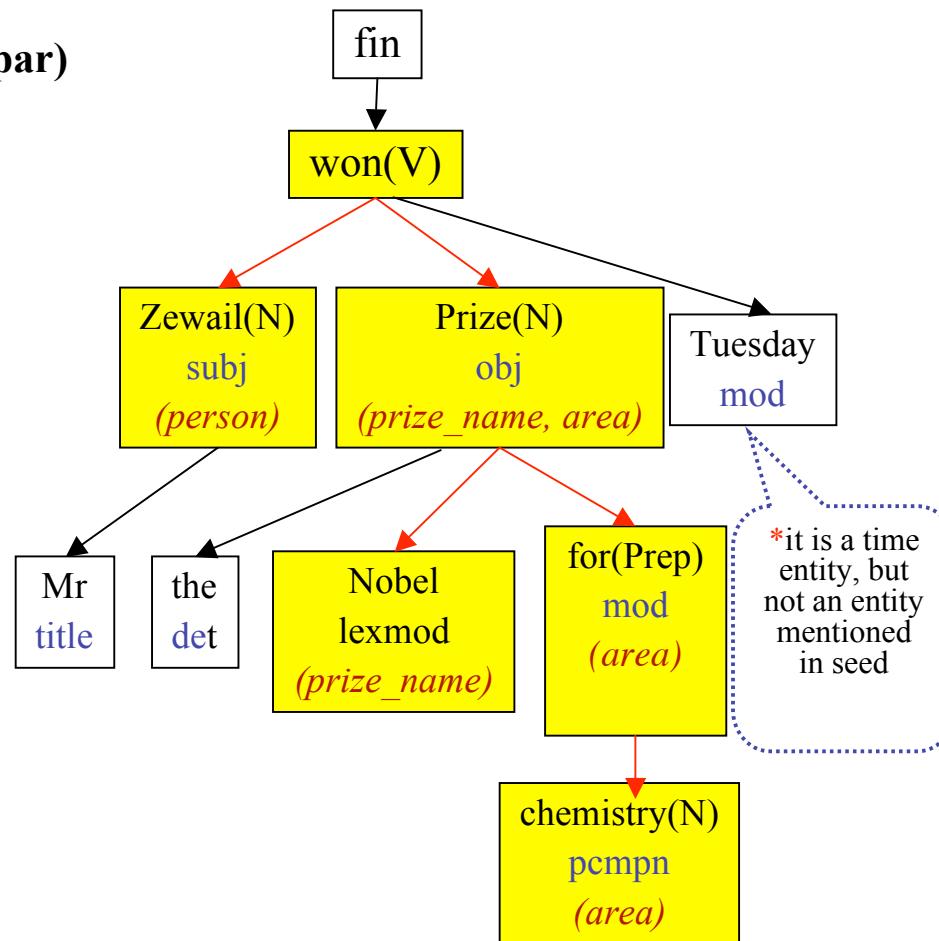


HANS USZKOREIT & FEIYU XU 07

Seed: *(Nobel, chemistry, [Ahmed H. Zewail], 1999)*

Sentence: *Mr. Zewail won the Nobel Prize for chemistry Tuesday.*

Parse Tree (SProUT + Minipar)



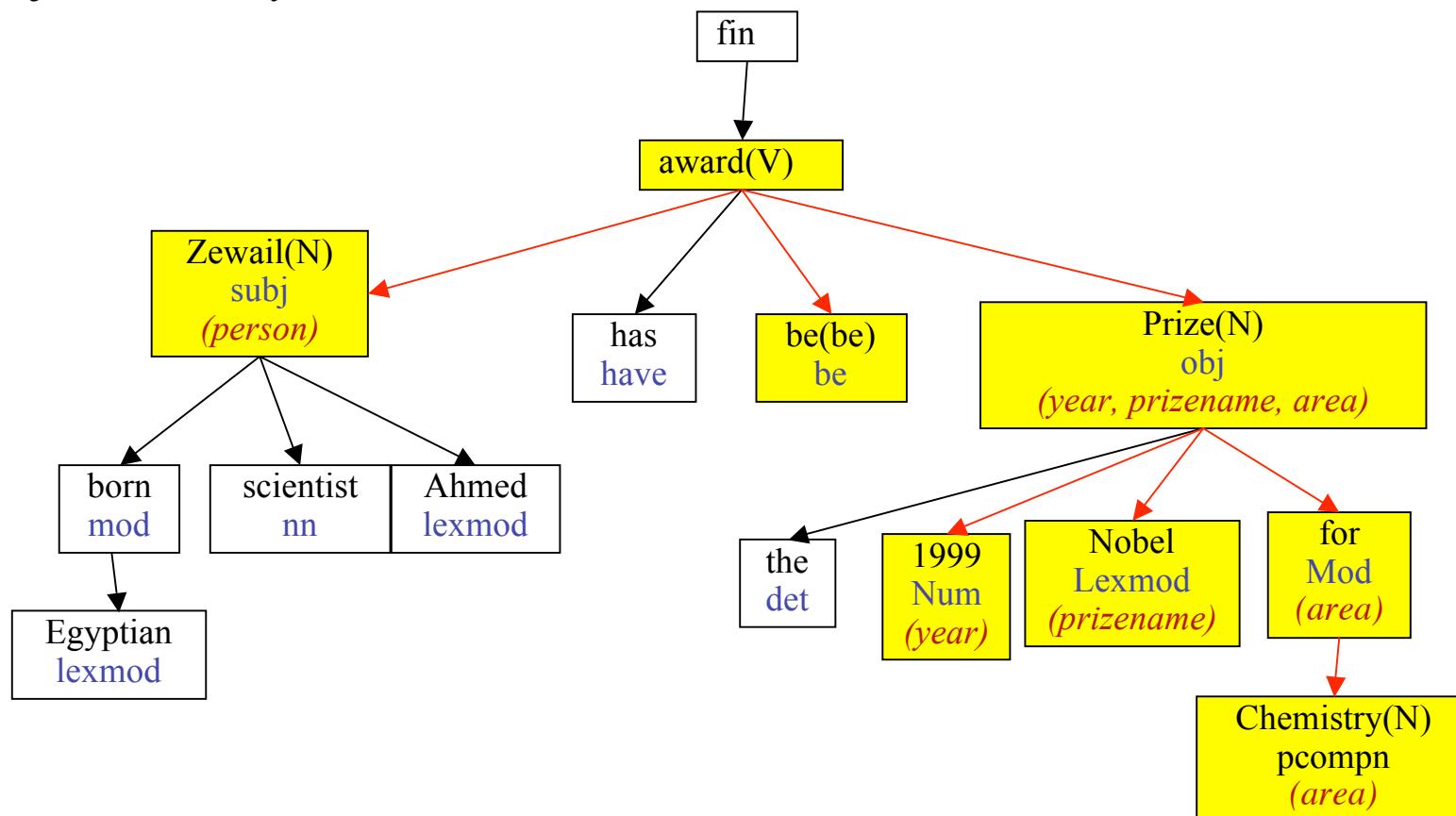
## Example 2



HANS USZKOREIT & FEIYU XU 07

Seed: *(Nobel, chemistry, [Ahmed H. Zewail], 1999)*

Sentence: *Egyptian-born scientist Ahmed Zewail has been awarded the 1999 Nobel Prize for Chemistry.*





- ☆ Which kind of sentences could represent an event?

complexity	matched sentence	event sentence	Relevant sentences in %
4-ary	36	34	94.0
3-ary	110	96	87.0
2-ary	495	18	3.6

Table 1. distribution of the seed complexity

## Distribution of Relation Projections



HANS USZKOREIT & FEIYU XU 07

combination (3-ary, 2-ary)	matched sentence	event sentence	relevant sentences in %
person, prize, area	103	91	82%
person, prize, time	0	0	0%
person, area, year	1	1	100%
prize, area, year	6	4	68%
person, prize	40	15	37%
person, area	123	0	0%
person, year	8	3	37%
prize, area	286	0	0%
prize, year	25	0	0%
area, year	12	0	0%

Table 2. distribution of entity combinations



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

# Experiments



HANS USZKOREIT & FEIYU XU 07

- ☆ Two domains
  - Nobel Prize award: *<recipient, prize, area, year>*
  - management succession: *<Person\_In, Person\_Out, Position, Organisation>*
- ☆ Test data sets

Data Set Name	Files	Volume
Nobel Prize A (1999-2005)	2296	12,6 MB
Nobel Prize B (1981-1998)	1032	5,8 MB
MUC-6	199	1 MB



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007



## ★ Conditions and Problems

- Complete list of Nobel Prize award events from online portal Nobel-e-Museum
- No gold-standard evaluation corpus available

## ★ Solution

- our system is successful if we capture one instance of the relation tuple or its projections, namely, one mentioning of a Nobel Prize award event. (Agichtein and Gravano, 2000)
- construction of so-called *Ideal* tables that reflexe an approximation of the maximal detectable relation instances
  - The Ideal tables contain all Nobel Prize winners that co-occur with the word “Nobel” in the test corpus and integrate the additional information from the Nobel-e-Museum

# Evaluation Against Ideal Tables



HANS USZKOREIT & FEIYU XU 07

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



German Research Center for Artificial Intelligence GmbH

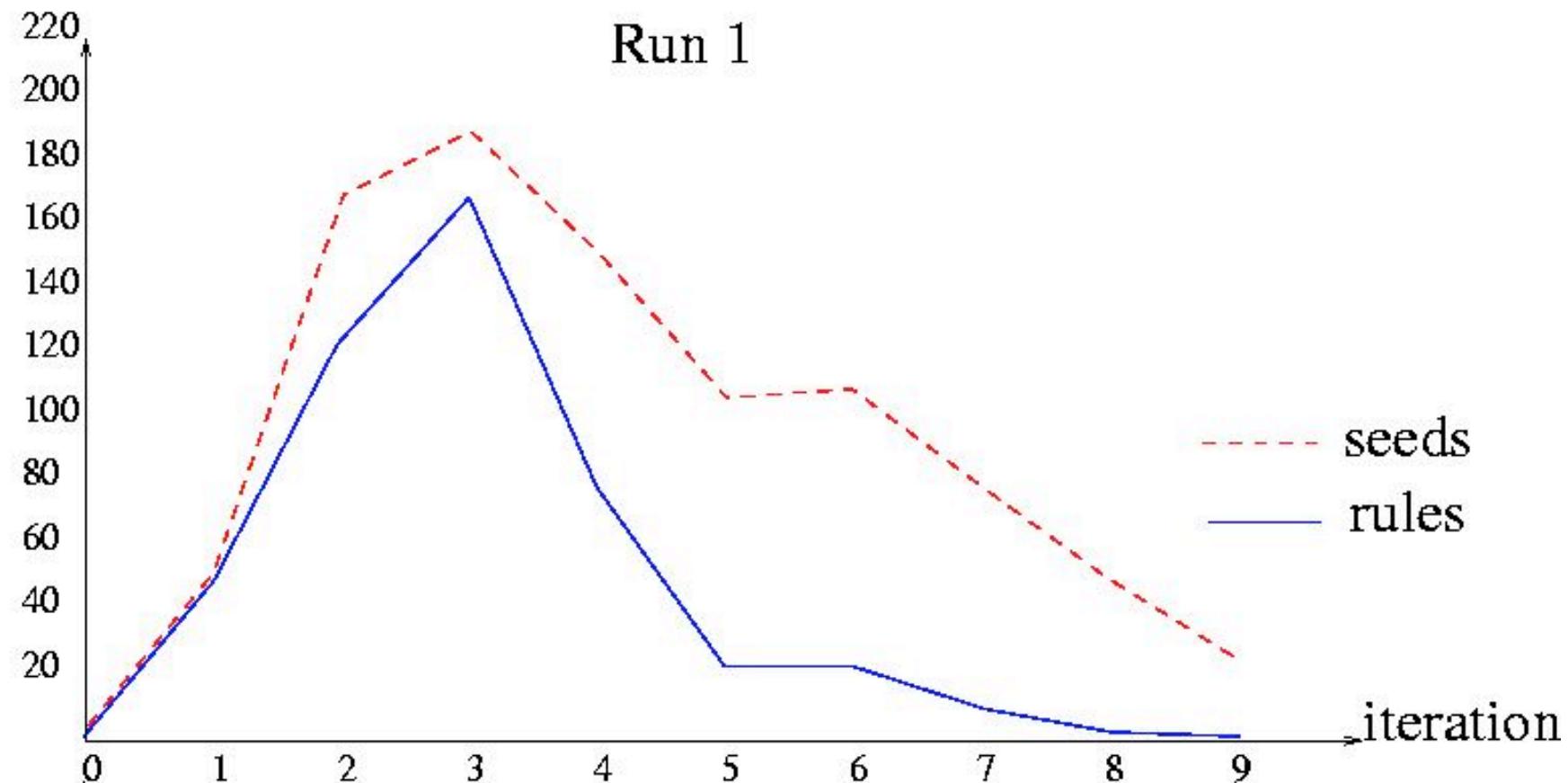
T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

## Iteration Behavior (Seed vs. Rule)



HANS USZKOREIT & FEIYU XU 07

seeds/rules



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

## The Dream



HANS USZKOREIT & FEIYU XU 07

- ☆ 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. (IDEAL Tables or Completeness)
- ☆ This sounds too good to be true!



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

## Research Questions



HANS USZKOREIT & FEIYU XU 07

★ 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?



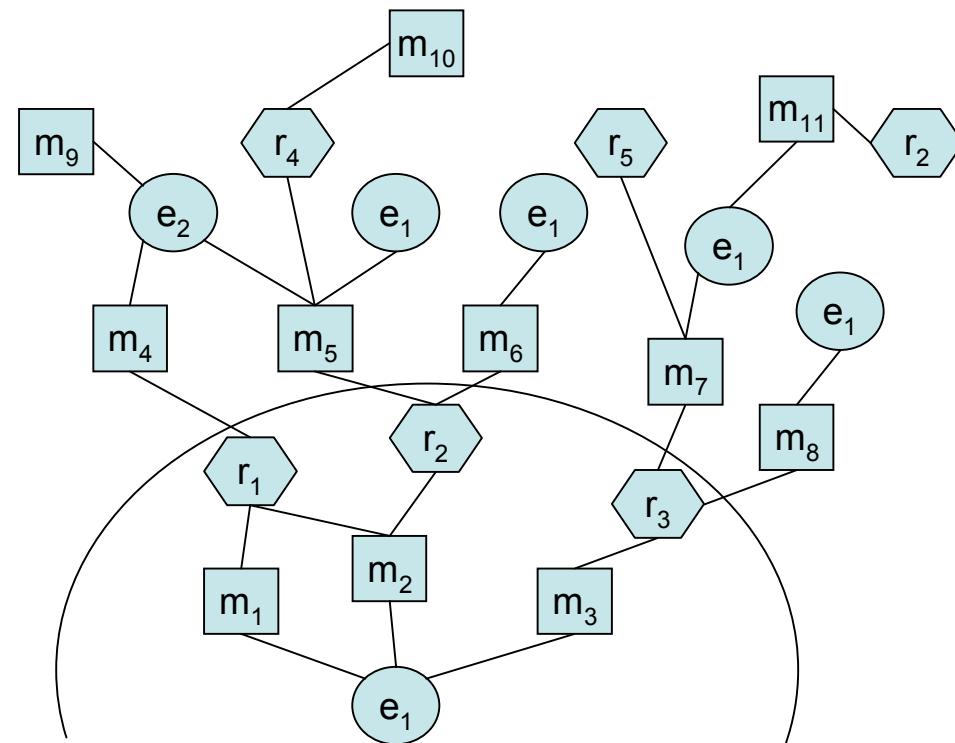
German Research Center for Artificial Intelligence GmbH

T-FaNT ★ TOKYO UNIVERSITY ★ 13 MARCH 2007

# Start of Bootstrapping (simplified)



HANS USZKOREIT & FEIYU XU 07



## Questions



HANS USZKOREIT & FEIYU XU 07

Can we reach all events in the graph?

By how many steps?  
From any event instance?



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

## Abstraction

bipartite graph

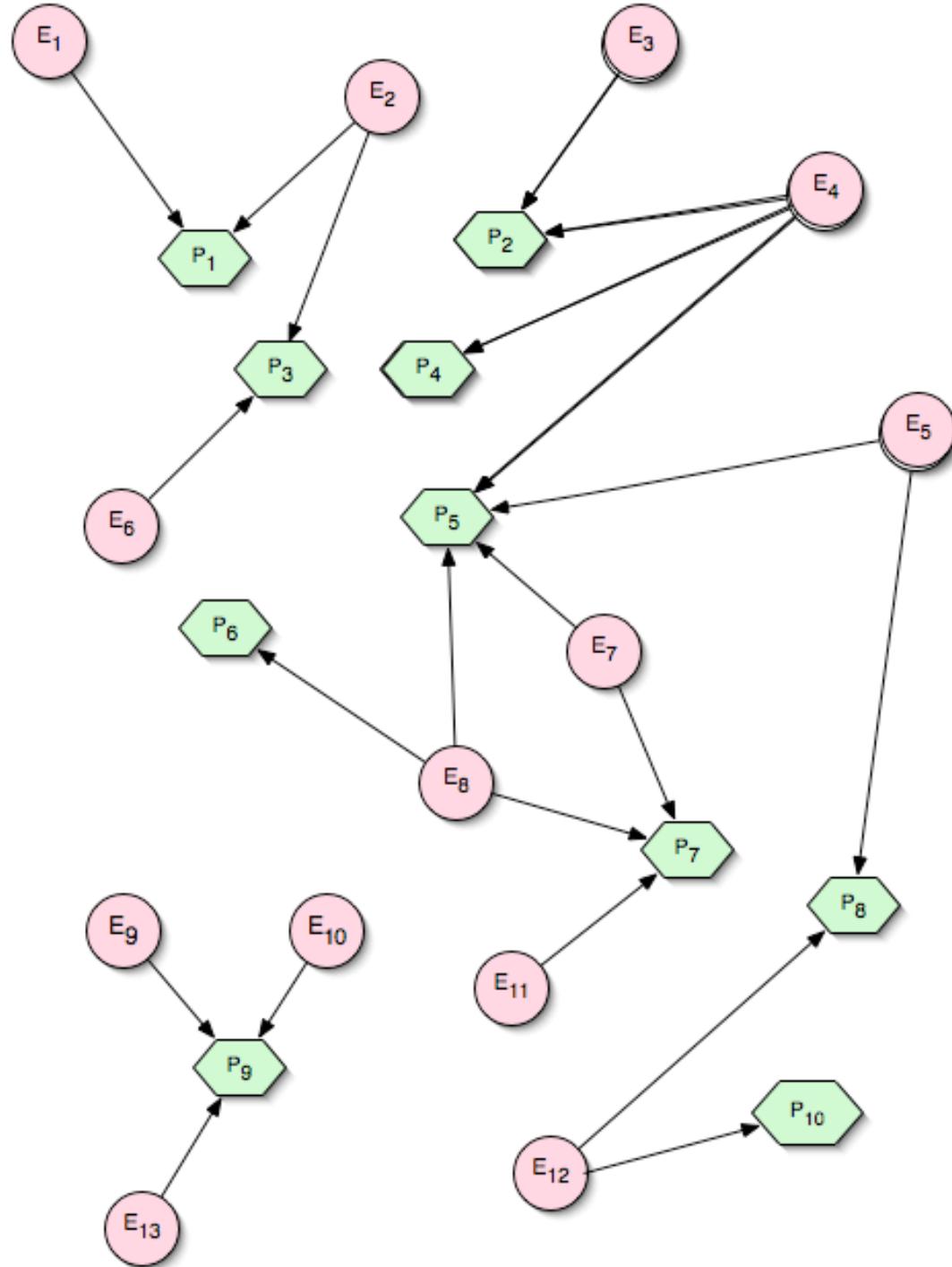
two types of vertices

$E_i$  = event instance

$P_j$  = linguistic pattern

relevant properties:

- ★ two degree distributions
- ★ connectedness
- ★ average and maximum path lengths between events



## Two Distributions



HANS USZKOREIT & FEIYU XU 07

1. Distributions of Pattern in Texts
2. Distribution of Mentionings to Relation Instances



German Research Center for Artificial Intelligence GmbH

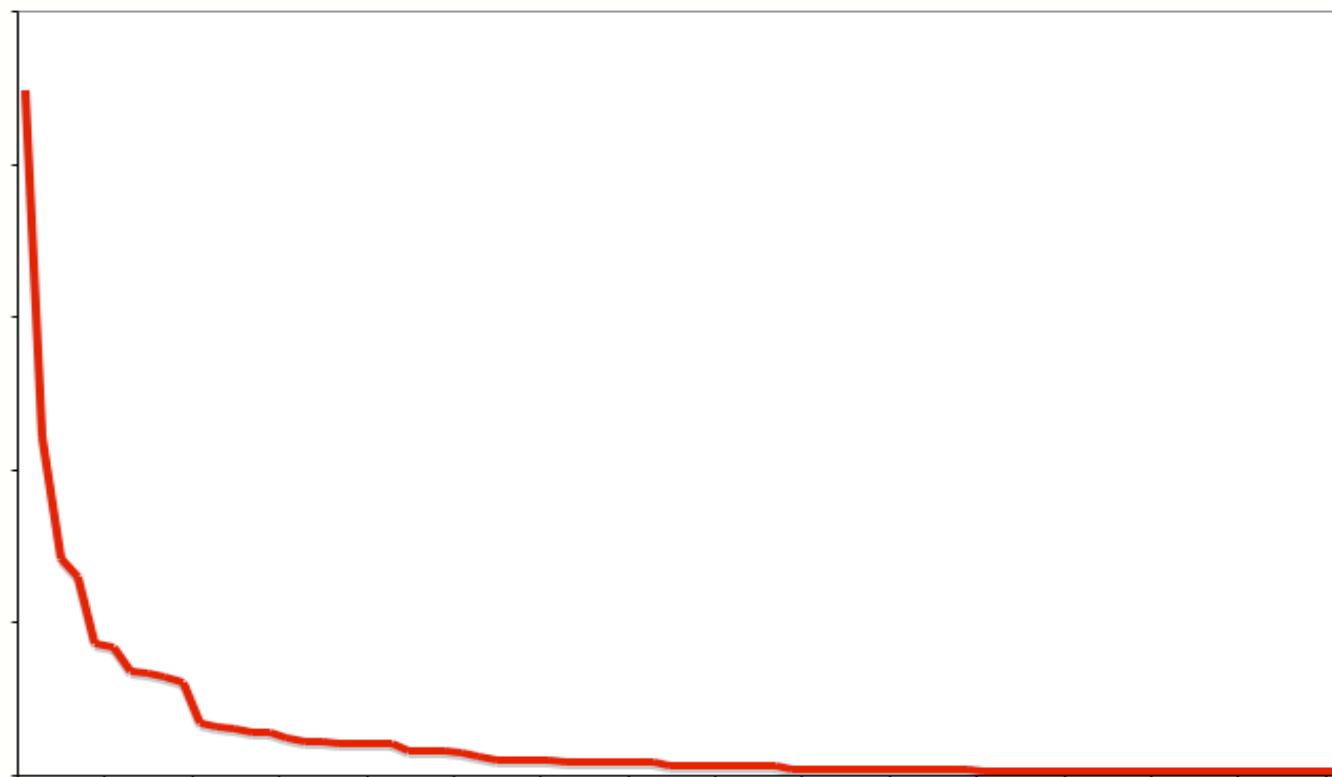
T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

## Two Distributions



HANS USZKOREIT & FEIYU XU 07

General distribution of patterns in texts probably follows Church's Conjecture: Zipf distribution (a heavy-tailed skewed distribution)



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

# Distribution of Mentionings to Events



HANS USZKOREIT & FEIYU XU 07

- ☆ Distribution of mentionings to relation instances (events) differs from one task to the other.
- ☆ The distribution reflects the redundancy in textual coverage of events.
- ☆ Distribution depends on text selection, e.g. number of sources (newspapers, authors, time period)

example 1: several periodicals report on Nobel Prize events

example 2: one periodical reports on management succession events



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007



HANS USZKOREIT & FEIYU XU 07

- ★ Degree Distribution gives the probability distribution of degrees in a complex network

$$p(k) = \sum_{v \in V \mid \deg(v)=k} 1$$

- ★ scale-free networks

$$P(k) \sim k^{-\gamma}$$

Zipf-like distribution (heavy-tailed skewed distribution) of degrees



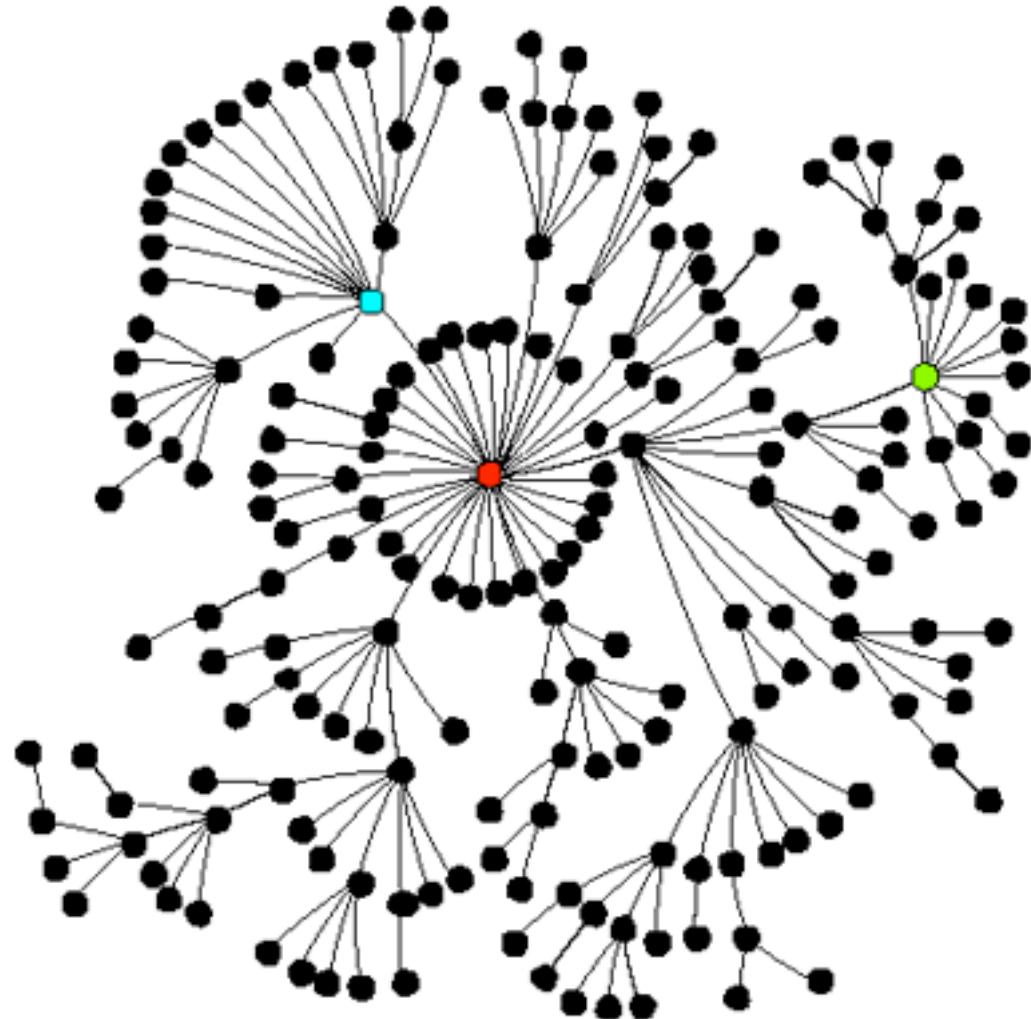
German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

## Example of Scale-Free Nets



HANS USZKOREIT & FEIYU XU 07



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

# Small-World Property



HANS USZKOREIT & FEIYU XU 07

Networks exhibiting the small-world property

- social networks (max path-length 5-7)
- co-authorship networks (Erdős number)
- Internet
- WWW
- air traffic route maps (max. 3 hops)

Networks that do not exhibit the small-world property

- road networks
- railway networks
- kinship networks



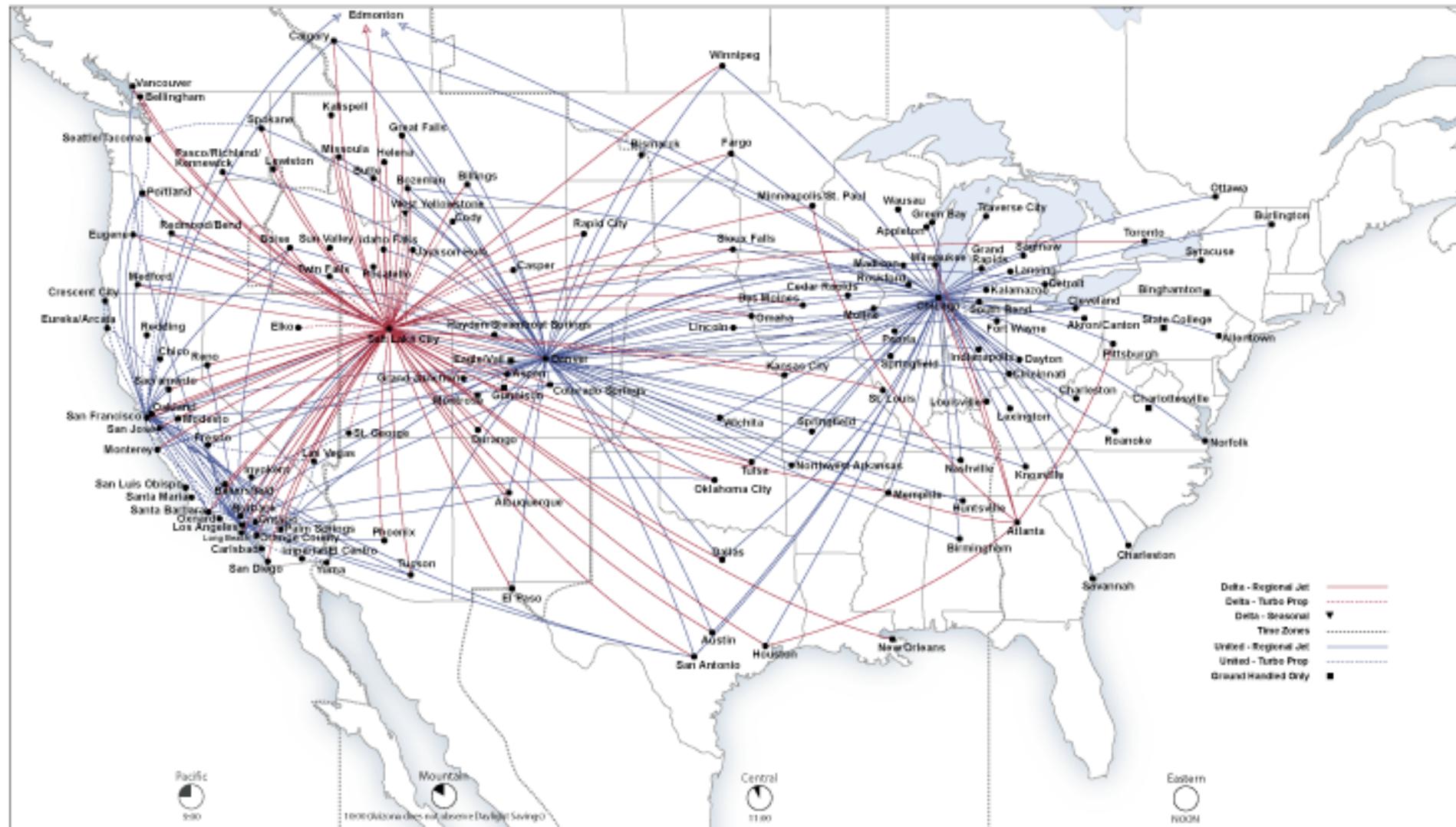
German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

# Airline Route Networks



HANS USZKOREIT & FEIYU XU 07



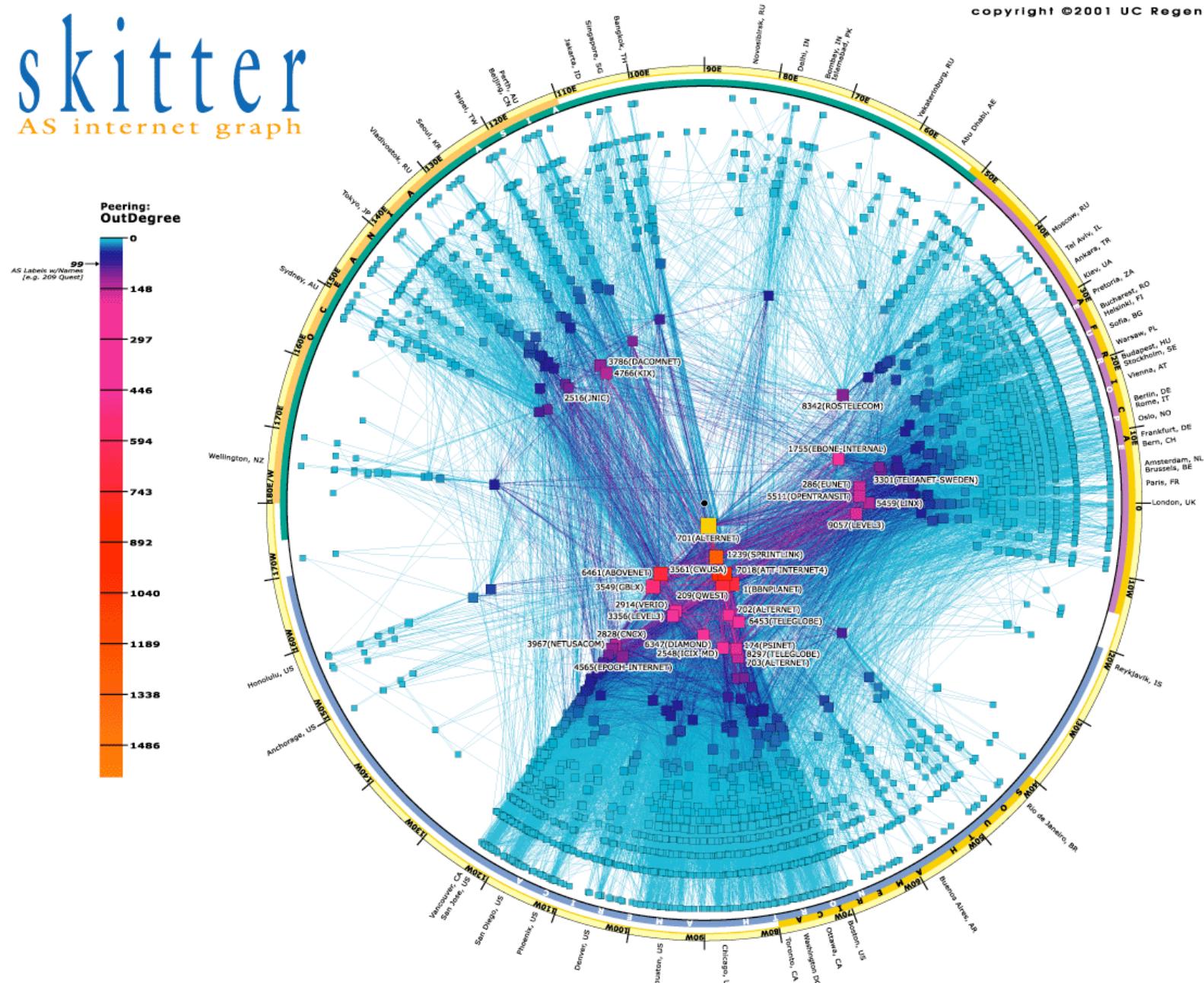
German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

# skitter

AS internet graph

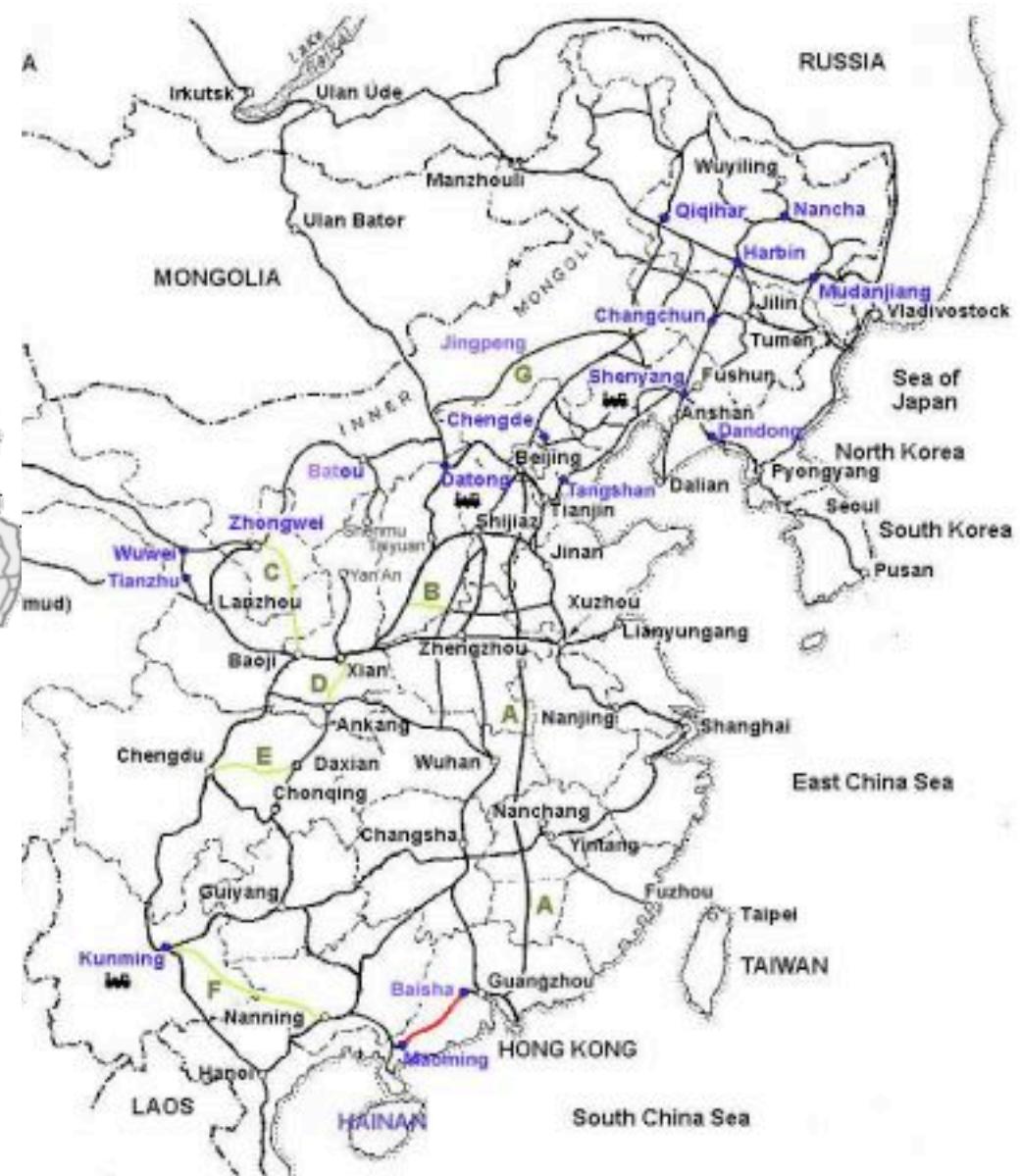
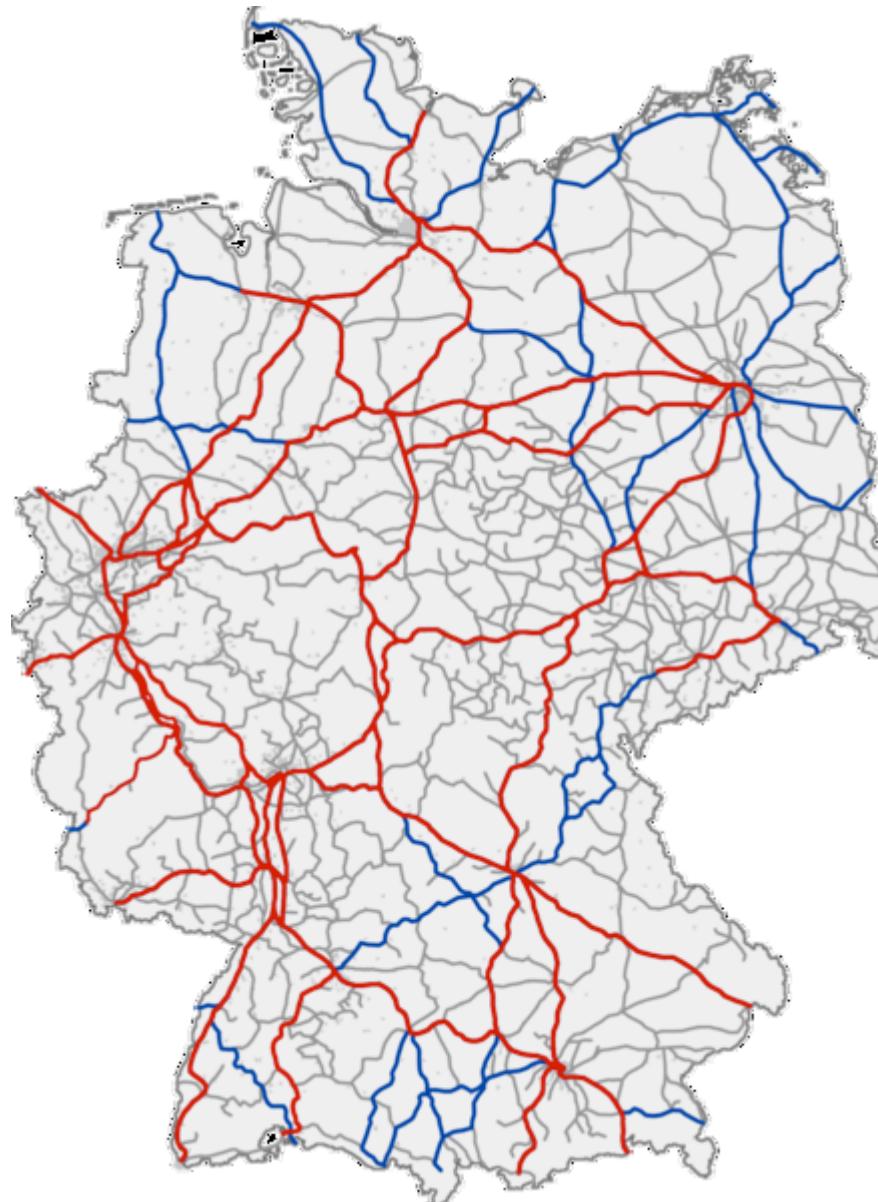
copyright ©2001 UC Regents. all rights reserved.



cooperative association for Internet data analysis   O san diego supercomputer center   O university of california, san diego

9500 gillman drive, mc0505   O la jolla, ca 92093-0505   O tel. 858-534-5000   O http://www.caida.org/





German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

# Small Worlds for Bootstrapping



HANS USZKOREIT & FEIYU XU 07

- ★ If both distributions follow a skewed distribution and if the distributions are independent from each other, then we get a scale-free network in the broader sense of the term.
  
- ★ For each type of vertices we get strong hubs. This leads to very short paths (for most connections).



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

# Degrees of Small for Small Worlds



HANS USZKOREIT & FEIYU XU 07

- ★ However, there are degrees of the small-world property.
- ★ Small World Networks are further optimized if there are forces beyond probability that cause hubs to be directly connected.



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

## Approaches to Solve the Problem



HANS USZKOREIT & FEIYU XU 07

- ☆ Enlarging the domain

Pulitzer Prize --> all Prizes

- ☆ selecting Carrier Domains (parallel learning domains)

Pulitzer Prize --> Nobel Prize

Ernst Winter Preis --> Nobel Prize

Fritz Winter Preis --> Nobel Prize



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

## Other Discovered Award Events



HANS USZKOREIT & FEIYU XU 07

Academy Award

actor % (Cannes Film Festival's Best Actor award)

American Library Association Caldecott Award

American Society

award

Blitzker

Emmy

feature % (feature photography award)

first % (the first Caldecott Medal)

Francesca Primus Prize

gold % (gold medal)

Livingston Award

National Book Award

Newbery Medal

Oscar

P.G.A

PEN/Faulkner Award

prize

reporting % (the investigative reporting award)

Tony

Tony Award

U.S. Open

But also:

nomination

\$1 million

\$29,000

about \$226,000

praise

acclaim

discovery

doctorate

election



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007



HANS USZKOREIT & FEIYU XU 07

☆ enlarging the text base for finding seeds and patterns

- New York Times MUC data --> general press corpora
- New York Times MUC data --> WWW

☆ enlarging the text base for finding new seeds

- New York Times MUC data --> WWW
- German Press Data --> English Press Data

## Summary



HANS USZKOREIT & FEIYU XU 07

- ☆ Our approach works with semantic seeds.
- ☆ It learns rules for an n-ary relation and its projections.
- ☆ Rules mark the slot-filler with their roles.



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

## Conclusions and Outlook



HANS USZKOREIT & FEIYU XU 07

- ☆ For some relation extraction tasks, the semantic seed based bootstrapping approach works surprisingly well.
- ☆ For others, it still works to some degree.
- ☆ Our deeper understanding of the problem helps us to select or prepare data for effective learning.



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007

## Next Steps



HANS USZKOREIT & FEIYU XU 07

- ☆ Go beyond the sentence.
- ☆ Investigate properties of relations w.r.t. data.
- ☆ Try to describe them as graph properties.
- ☆ Try out auxiliary data sets (such as the Web).
- ☆ Try out deep processing: extract patterns from RMRS with extended ERG (first tests by Zhang Yi 80% coverage for Nobel prize sentences, 61% for management succession)



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007



HANS USZKOREIT & FEIYU XU 07

The material presented here has been submitted for publication.  
An earlier stage of the results was published in:

Feiyu Xu, Hans Uszkoreit & Hong Li: Automatic Event and Relation  
Detection with Seeds of Varying Complexity. In Proceedings of the AAAI  
2006 Workshop Event Extraction and Synthesis, Boston, July, 2006.



German Research Center for Artificial Intelligence GmbH

T-FaNT ☆ TOKYO UNIVERSITY ☆ 13 MARCH 2007