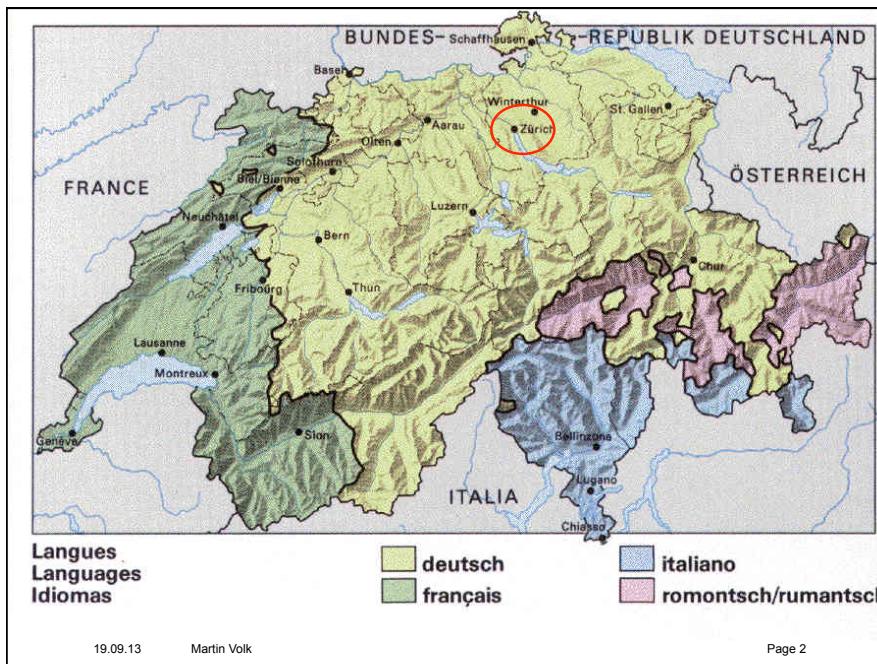


 University of
Zurich^{UZH}
Institute of Computational Linguistics

Multiword Expressions in Parallel Treebanks

Martin Volk
volk@cl.uzh.ch
Institut für Computerlinguistik
Universität Zürich
Parseme COST, Warschau, September 2013



 **University of
Zurich^{UZH}**

Institute of Computational Linguistics

- Biomedical Text Mining (with Novartis)
- Style Checking in Law Texts
- Sentiment Detection
- Digital Humanities (digitizing heritage texts)
- Machine Translation



31. Aug. 2011 Martin Volk, Uni Zurich

3

 **University of
Zurich^{UZH}**

Institute of Computational Linguist

Teaching

Bachelor and Master
in Computational
Linguistics

 **Universität
Zürich^{UZH}**

Master of Arts
Multilingual Text Analysis
Multilinguale Textanalyse
Analyse Multilingue de Texte

The University of Zurich offers an innovative specialized Master in Comparative Corpus Linguistics, combining Computers and Linguistics.

An interdisciplinary program by
The English Department
The Institute of German Studies
The Institute of Romance Studies
The Institute of Computational
Linguistics

Start: September 2012
Application period:
1 Dec. 2011–30 April 2012
Further information:
www.mta.uzh.ch
mta@cl.uzh.ch



19.09.13

 University of
Zurich^{UZH}

Institute of Computational Linguistics

Multiword / Multitoken Expressions

In this talk:

- Names
 - Person, Location, Organization, Event
- Dates and Numbers
- Support/Light verb constructions
- DE verbs with separated prefixes / EN particle verbs
- ES and FR Multiword prepositions

Excluded:

- Idioms, Collocations, ...

19.09.13 Martin Volk Page 5

 University of
Zurich^{UZH}

Institute of Computational Linguistics

Multiword / Multitoken Expressions

Contiguous ... vs.

- Names → *John F. Kennedy*

... Discontiguous

- DE verbs with separated prefixes → *fängt das Spiel an*

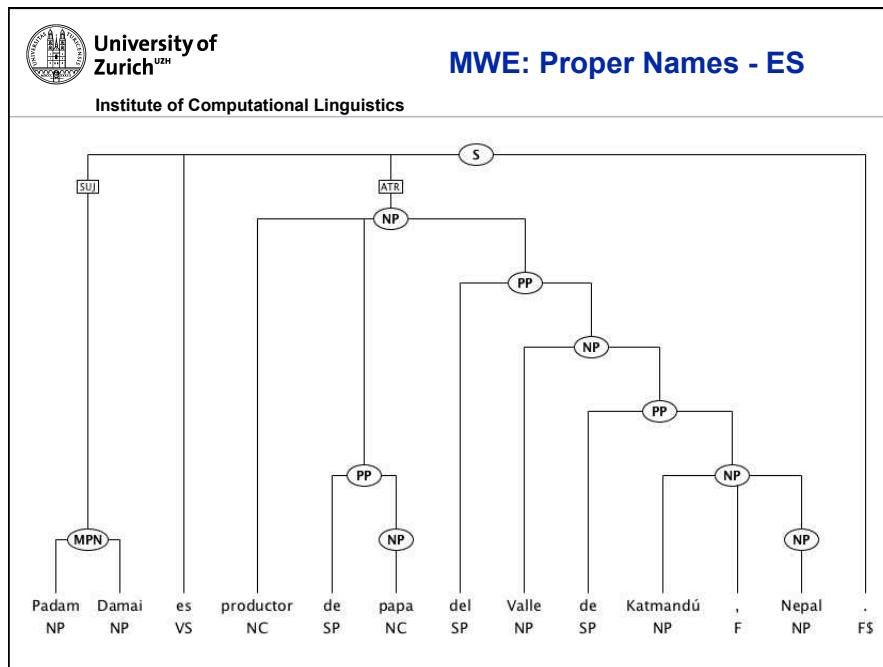
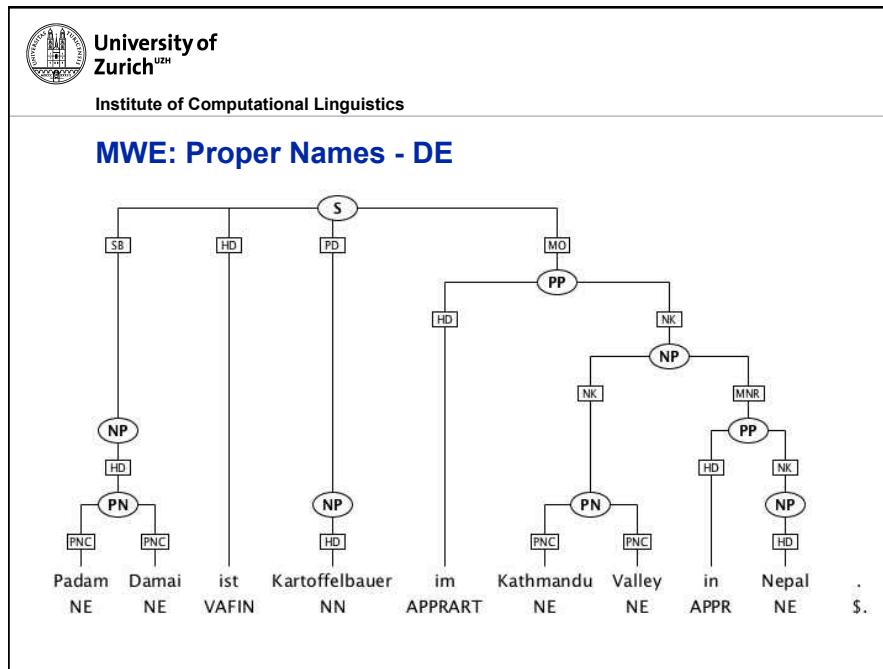
Static ... vs.

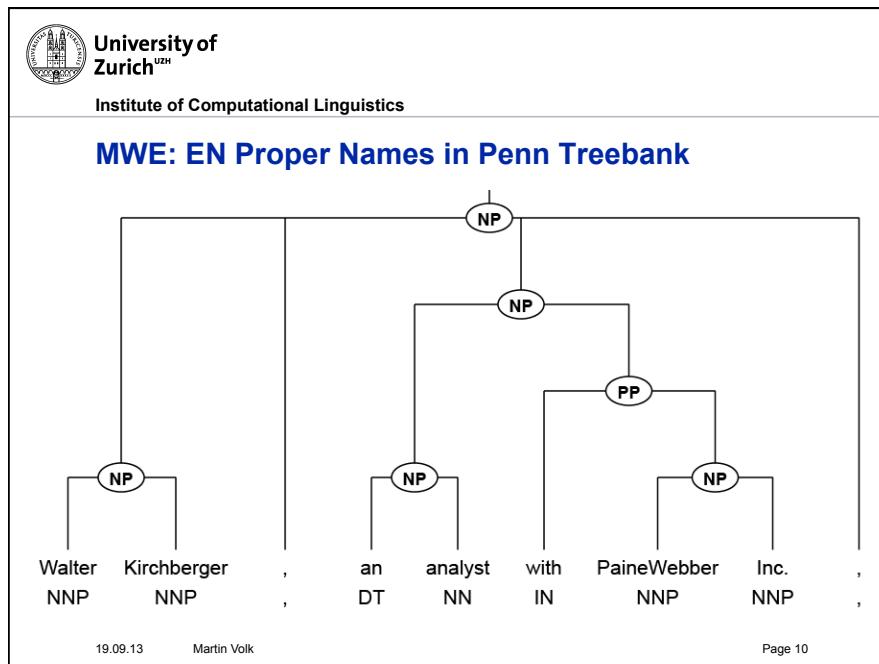
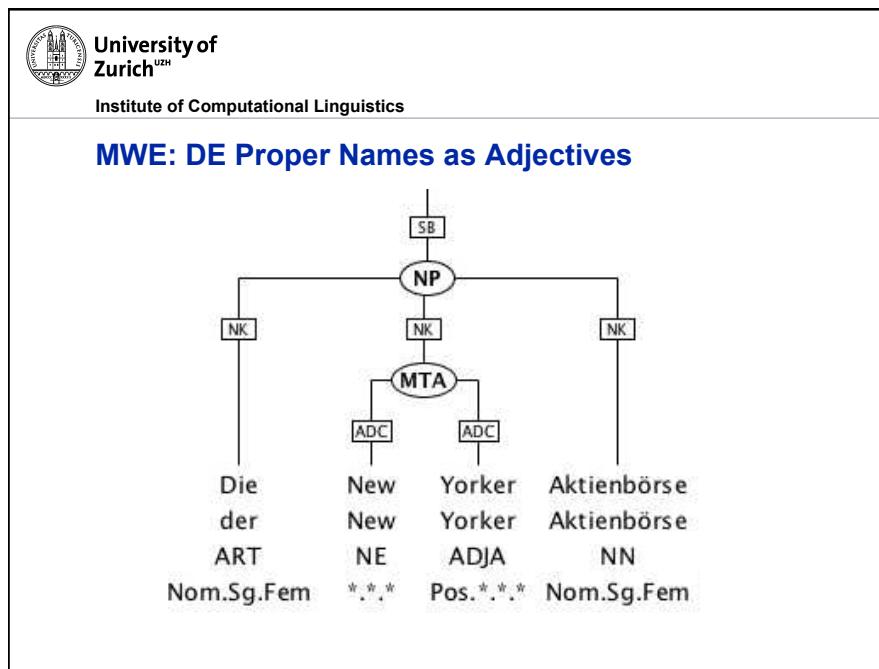
- Dates → *15. June 2005*

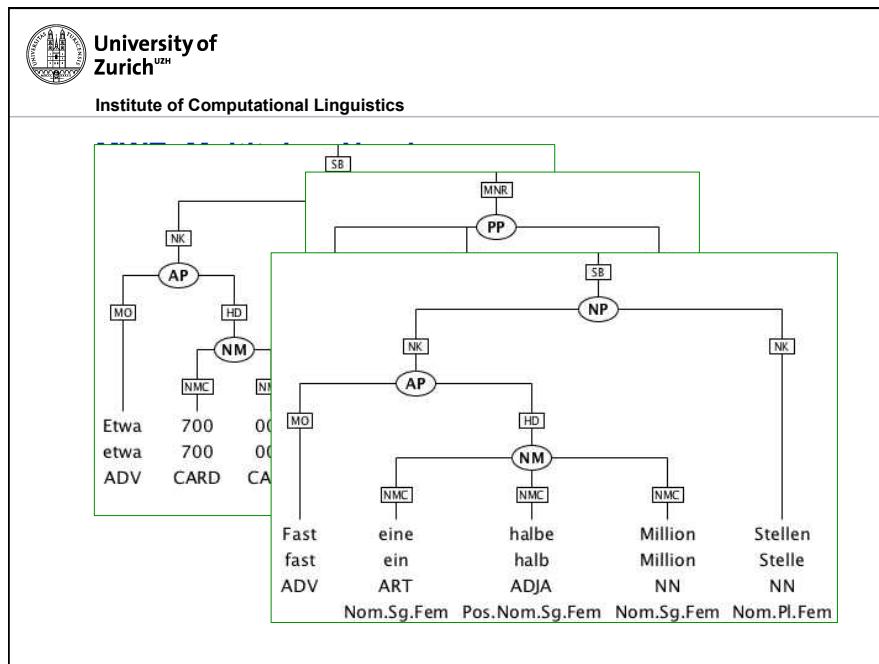
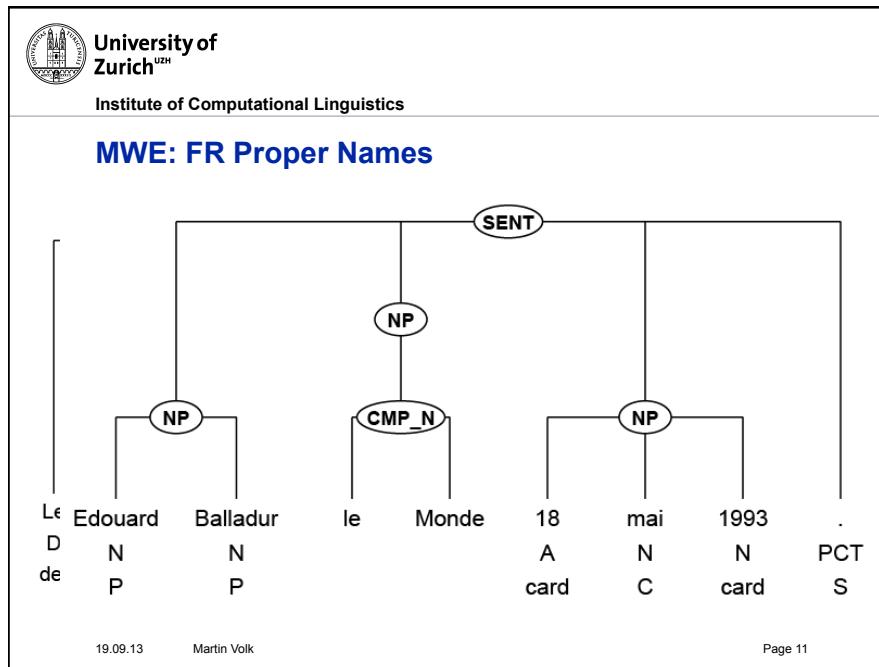
... Dynamic / Inflectable / Modifiable

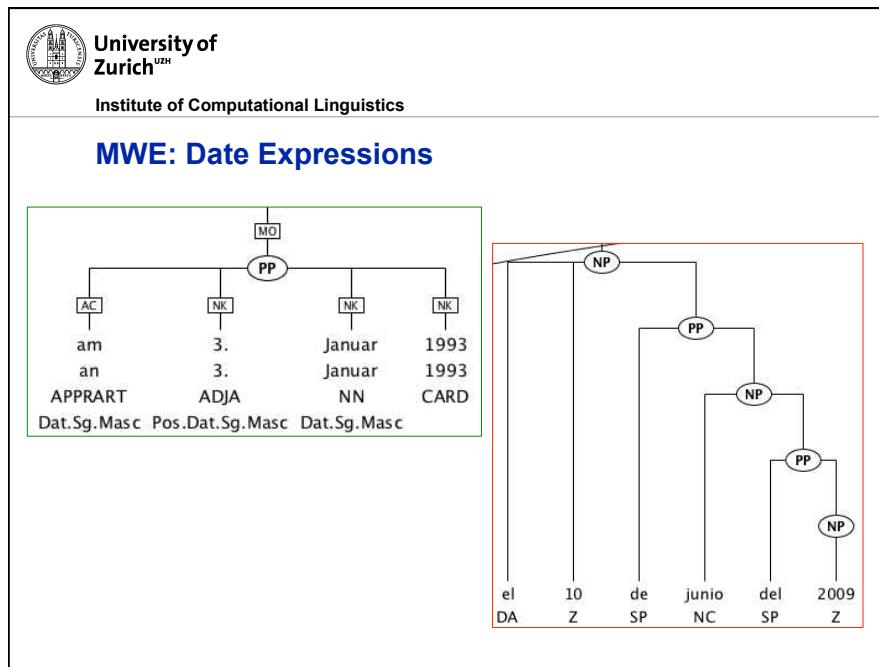
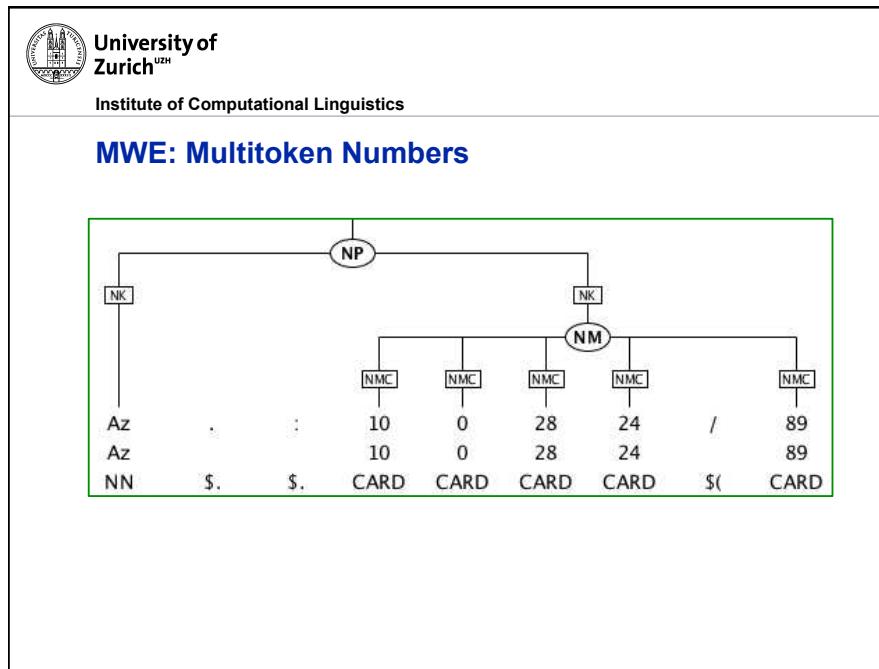
- Support verb constructions → *zur Verfügung stellen*
- Geo names → *refuge solitaire de Sponda*

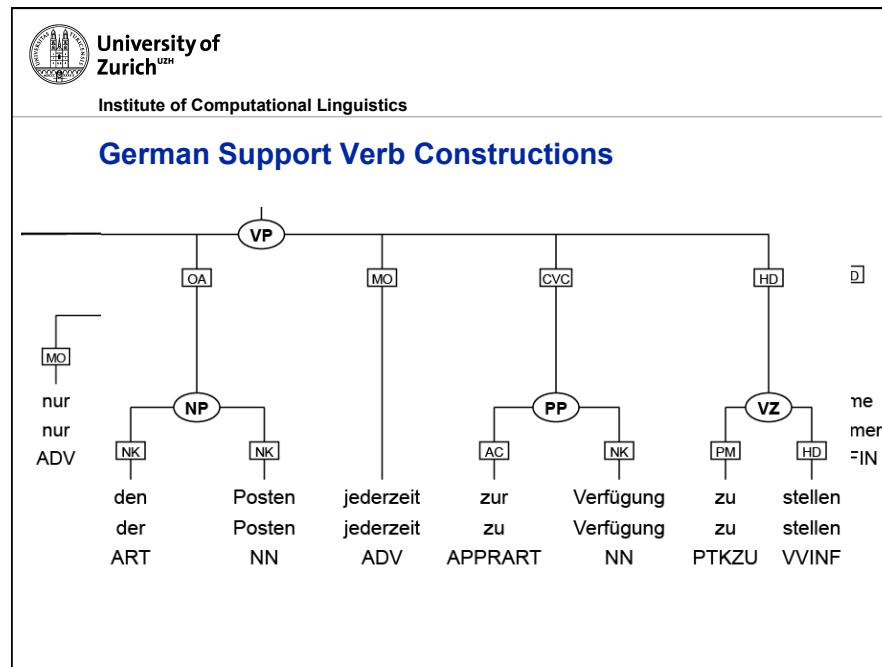
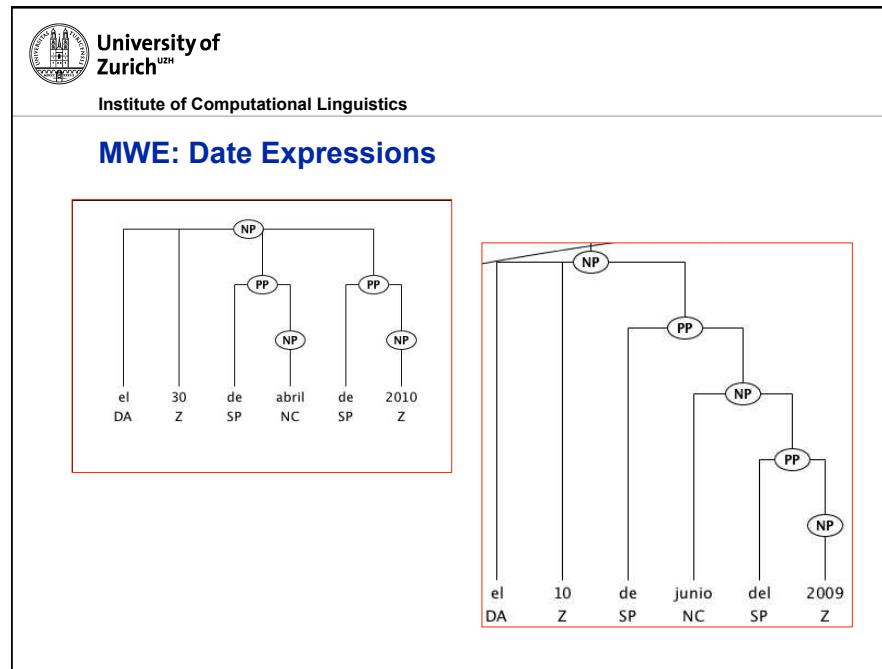
19.09.13 Martin Volk Page 6

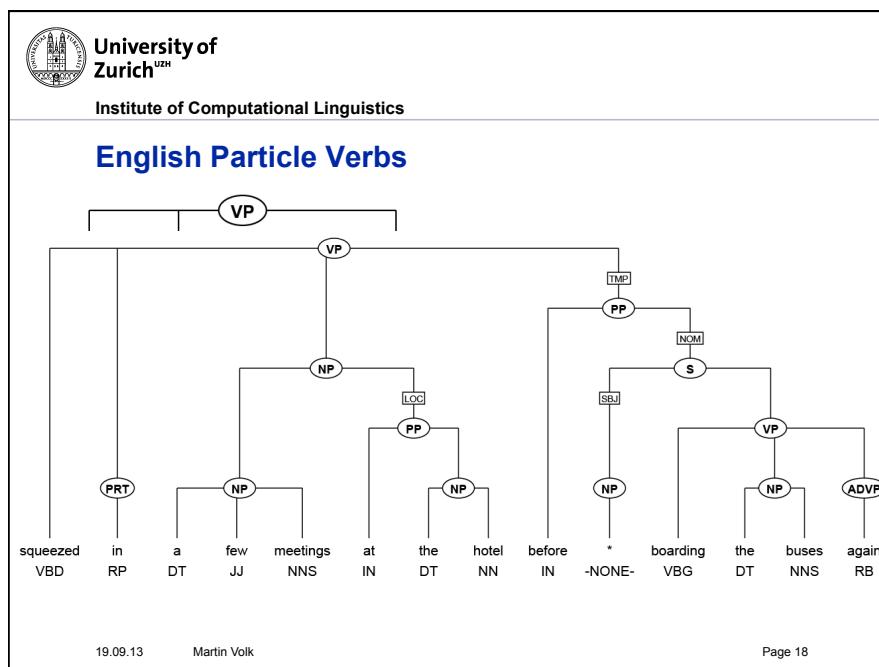
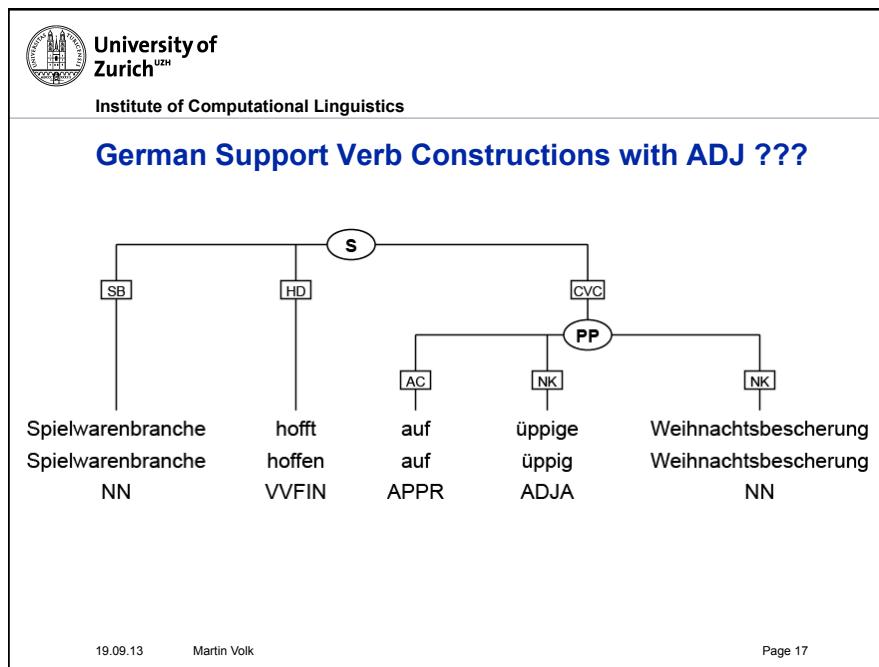












 University of
Zurich^{UZH}

Institute of Computational Linguistics

English Particle Verbs

... in the Penn Treebank (WSJ 0-12): 24.618 trees = 629'798 tokens
PoS = particle (= RP): 1621 occurrences
category = particle (= PRT): 1608 occurrences (= 6.5% of all trees)

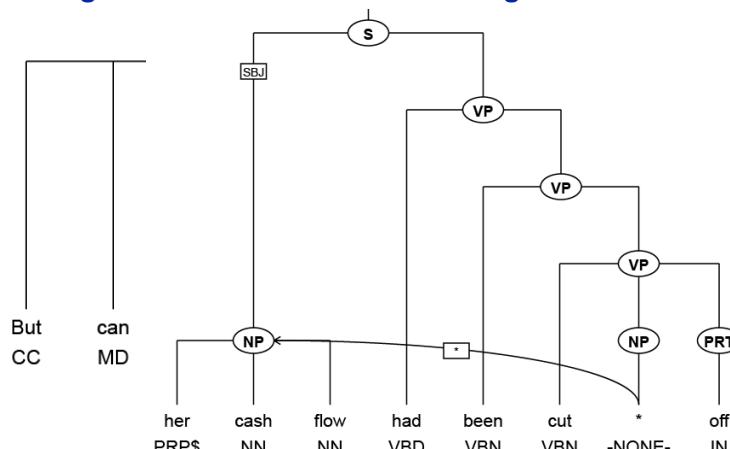
- out of which 212 Verb + PRT pairs are discontiguous (13%)
- out of which 175 occur with 1 intervening token
 - out of which 109 occur with the **empty** token intervening !!!
 - while the remaining 37 have two or more intervening tokens

19.09.13 Martin Volk Page 19

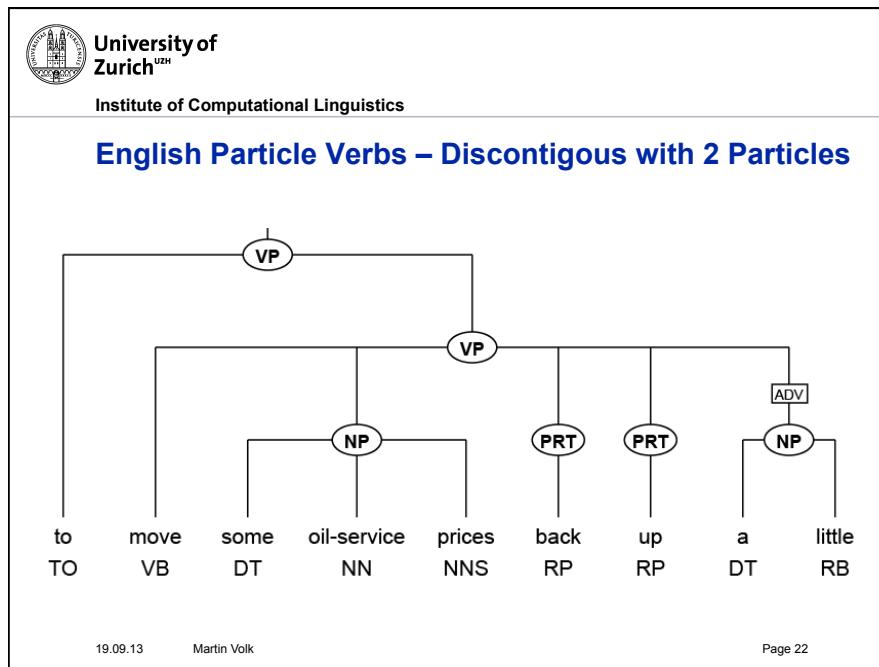
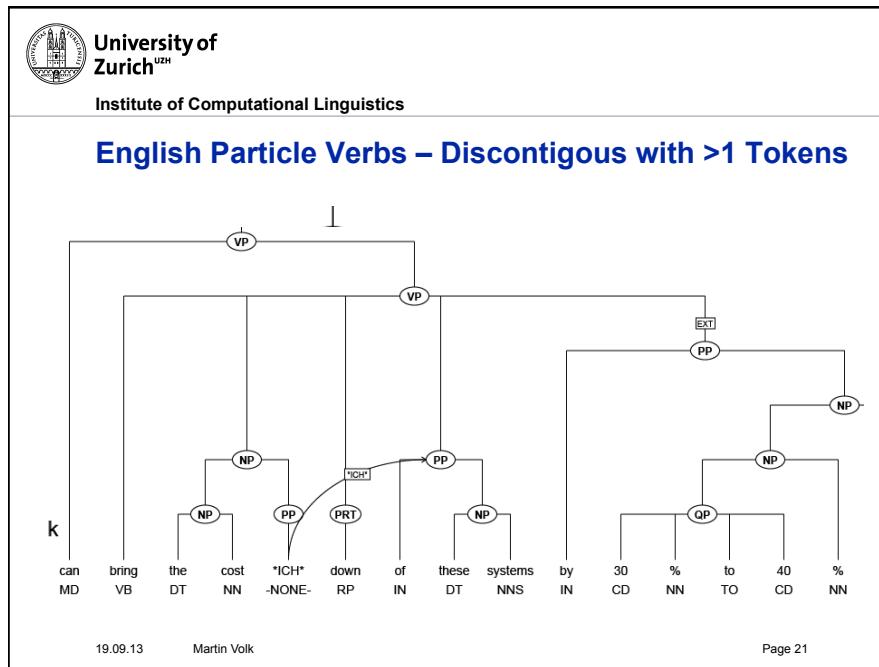
 University of
Zurich^{UZH}

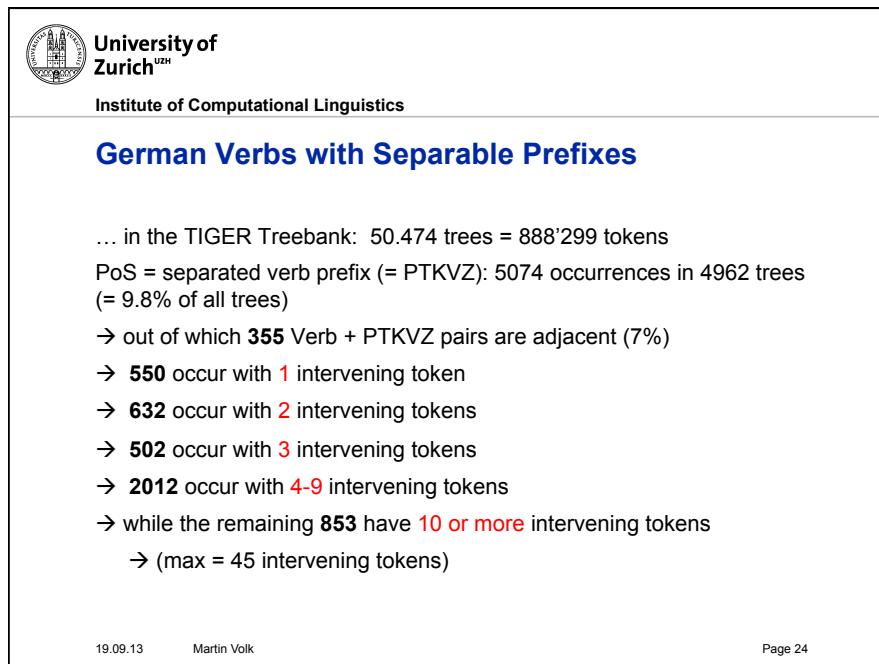
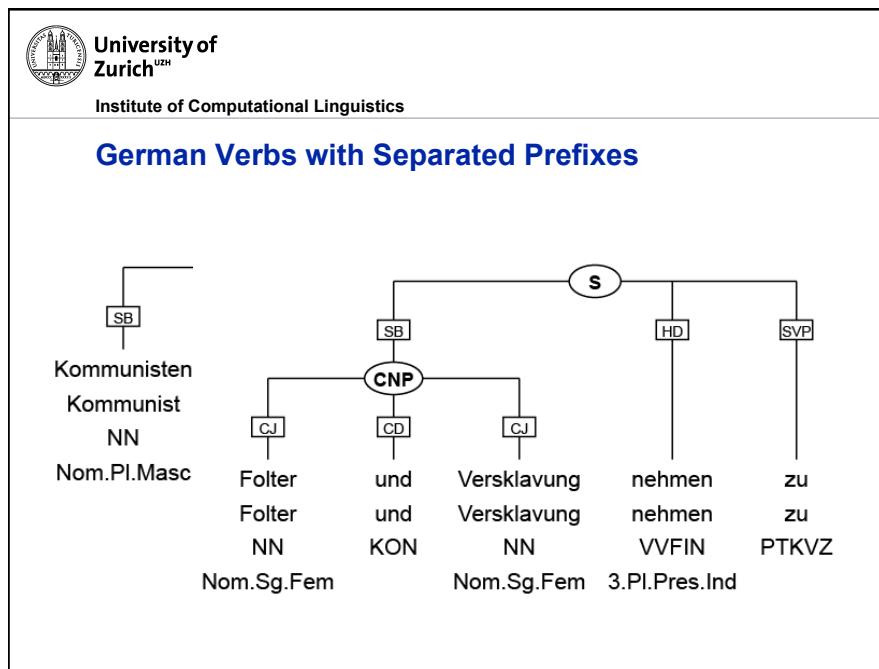
Institute of Computational Linguistics

English Particle Verbs – Discontiguous with 1 Token



19.09.13 Martin Volk Page 20





 University of
Zurich^{UZH}

Institute of Computational Linguistics

German Verb with Separable Prefixe (dist = 45)

Diese Staatsidee wurzelt in bürgerlich-demokratischen Emanzipationsvorstellungen und **weist** dem Staat eine aktive Rolle in der Steuerung wirtschaftlicher und gesellschaftlicher Abläufe in bezug auf eine – über den “ sozialstaatlichen Reparaturauftrag “ hinausweisende – Förderung einer größeren Gleichheit der Lebenschancen in bezug auf eine lebenslange Einkommenssicherung , Gesundheit , Kindererziehung , Wohnen und Bildung **zu**.

19.09.13 Martin Volk Page 25

 University of
Zurich

Institute of Computational Linguistics

German Verbs with Separable Prefixes

 University of
Zurich^{UZH}

Institute of Computational Linguistics

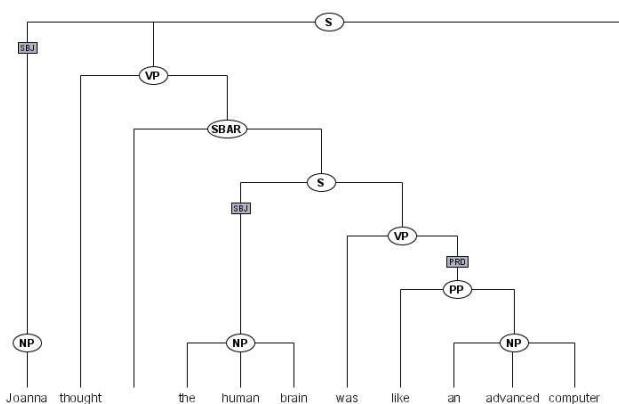
Parallel Treebanks

19.09.13 Martin Volk Page 27

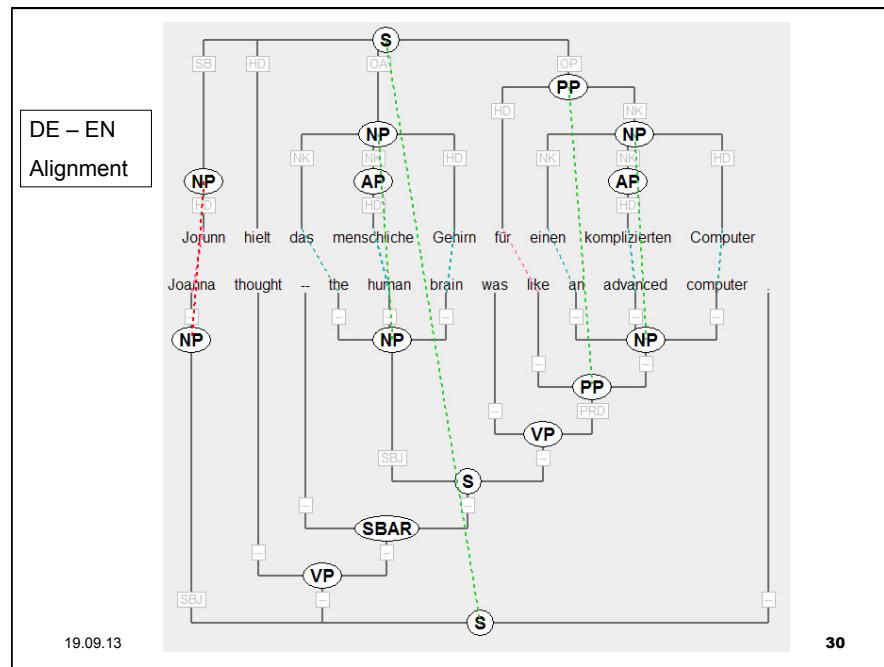
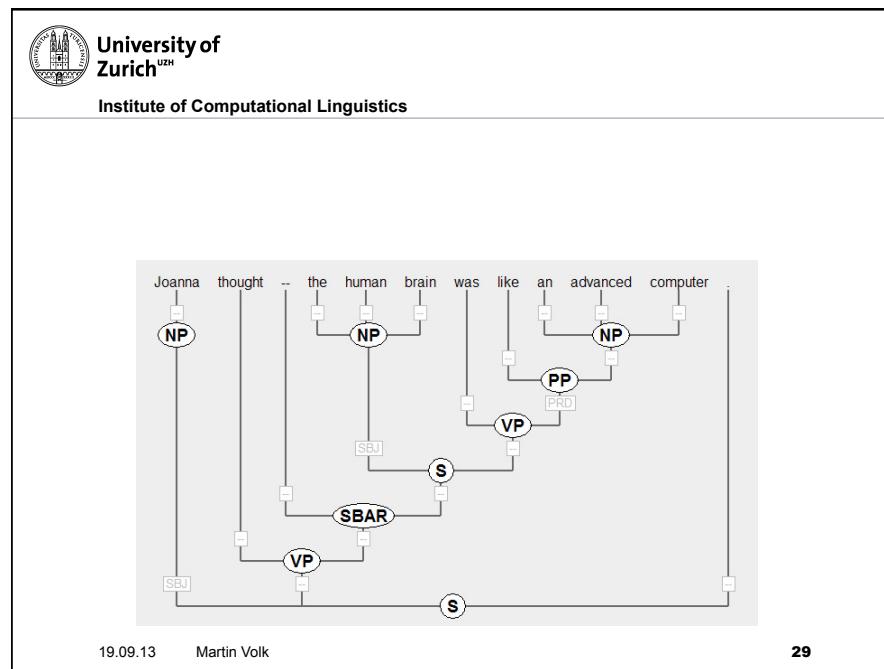
 University of
Zurich^{UZH}

Institute of Computational Linguistics

English Syntax Tree



19.09.13 Martin Volk 28





SMULTRON

= Stockholm MULtilingual TReebank

- 1500 sentences in 3 languages (DE-EN-SV)
 - 500 from Jostein Gaarder's *Sophie's World*
 - 500 from Economy texts
 - ABB Quarterly report
 - Rainforest Alliance: Banana Certification Program
 - SEB Annual report
 - 500 from DVD player manual

19.09.13

Martin Volk

31



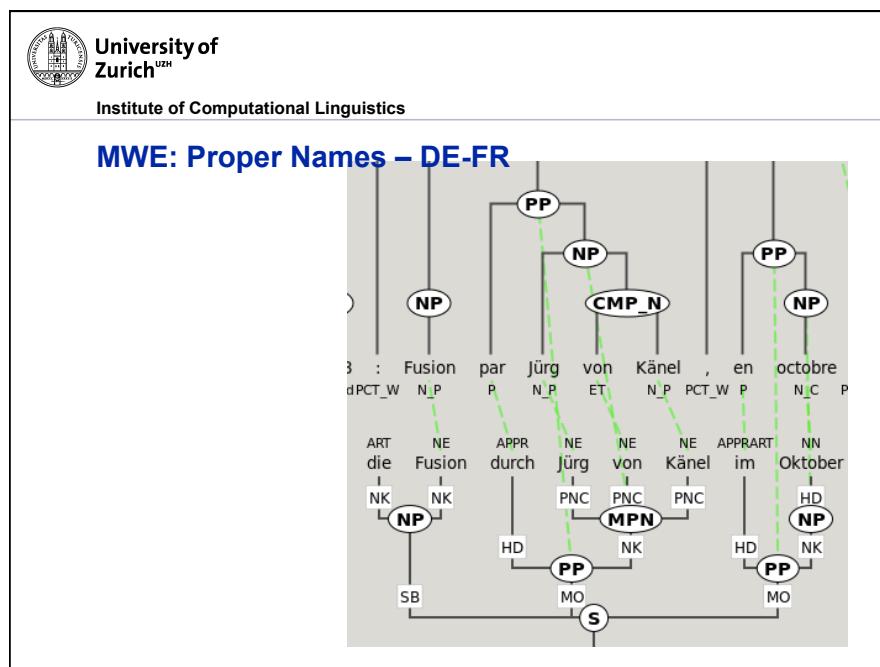
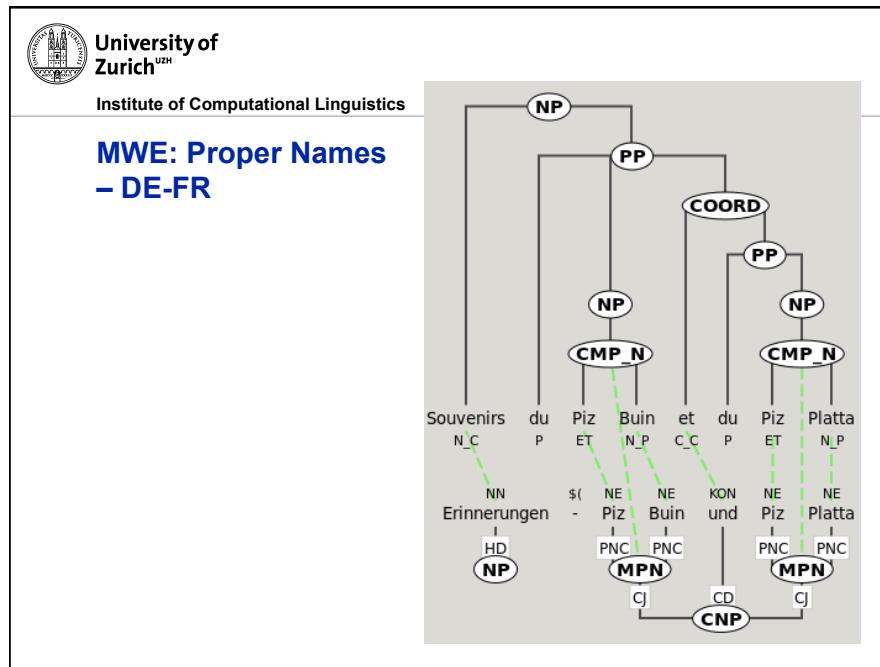
SMULTRON

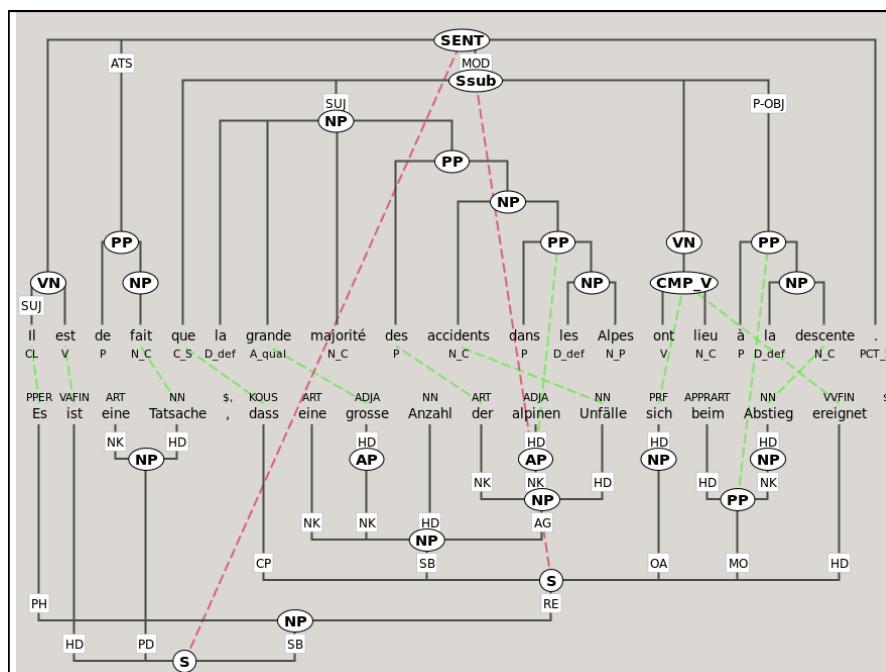
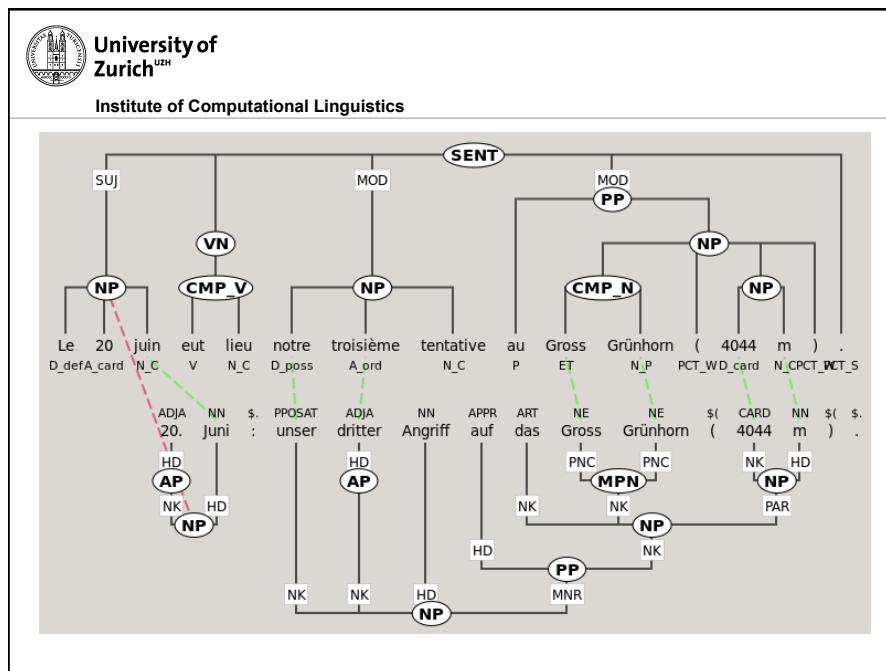
- Plus: 1000 sentences in 2 languages (DE-FR)
 - Swiss Alpine Club mountaineering texts
- Plus: 4000 sentences in 2 languages (DE-ES)
 - agriculture, education, biography, etc.
- In preparation (ES - Quechua)
 - 500 trees for Quechua

19.09.13

Martin Volk

32





 University of
Zurich^{UZH}

Institute of Computational Linguistics

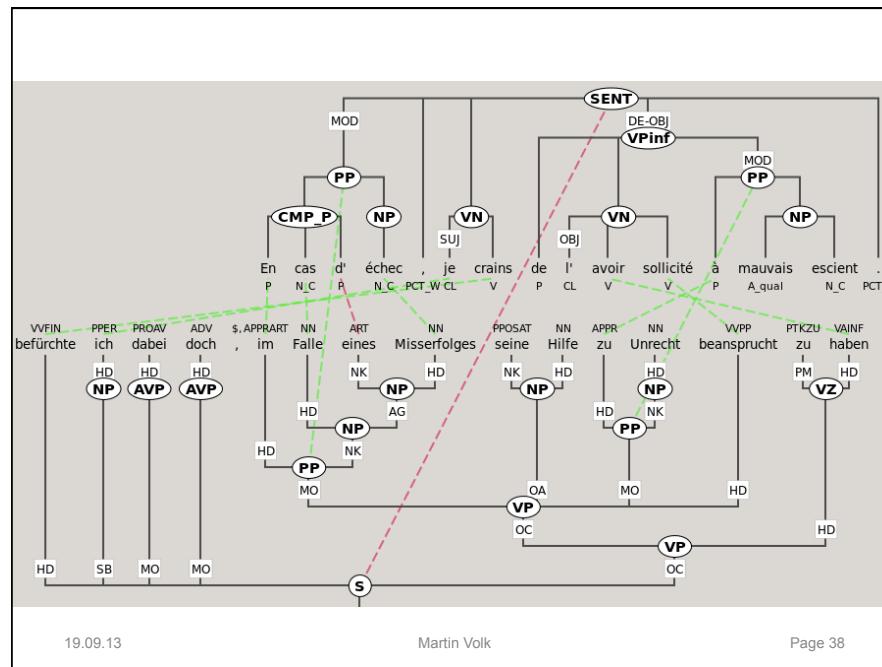
French MW prepositions

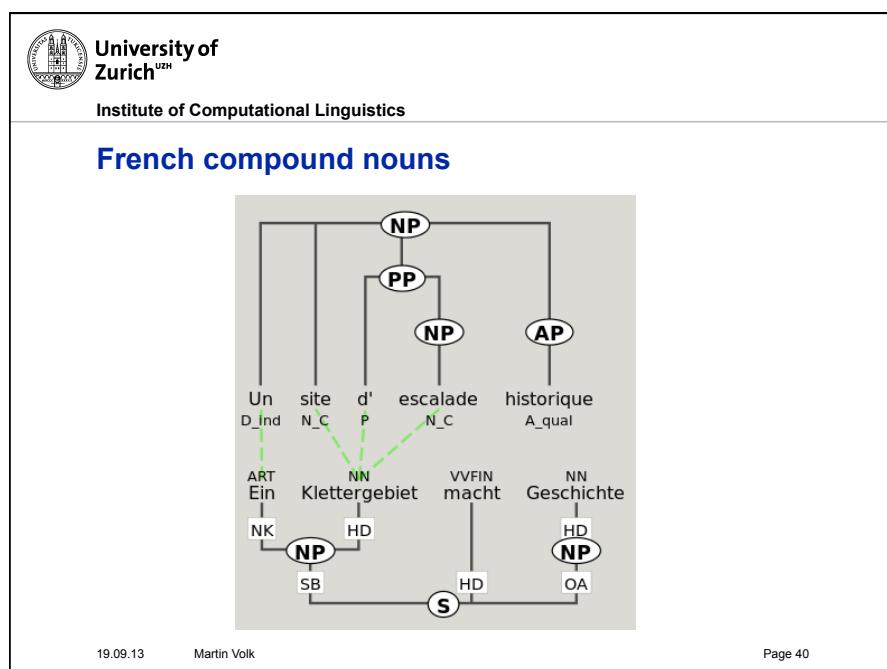
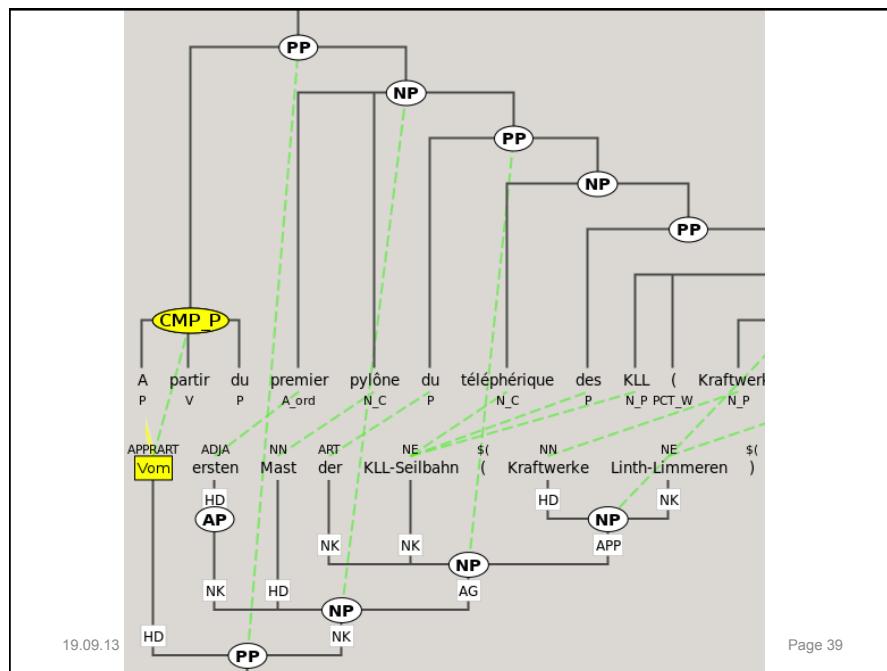
Given: The FR Le Monde Treebank with 20,500 trees

MW prepositions

- consisting of 2 tokens: 1825 occurrences
près de, jusqu' à, plus de, lors de, afin de, ...
- consisting of 3 tokens: 1799 occurrences
par rapport à, en matière de, en raison de, ...
- consisting of 4 tokens: 1443 occurrences
à la fin de, à la tête de, sous le nom de, au _ cours de, ...
- consisting of 5 tokens: 105 occurrences
de l' autre côté de, au _ - dessous de, ...
- consisting of 6 tokens: 74 occurrences
vis – à – vis de, de part et d' autre de, c' est - à - dire, ...

19.09.13 Martin Volk Page 37





University of Zurich^{UZH}

Institute of Computational Linguistics

Spanish MW prepositions

19.09.13 Martin Volk Page 41

University of Zurich^{UZH}

Institute of Computational Linguistics

Summary

MWEs in Treebanks

- cover a wide range of phenomena
 - MW proper names
 - MW numbers
 - Date and time expressions
 - Support verb constructions
 - ES and FR MW prepositions
 - DE verbs with separated prefixes, EN particle verbs
- are handled very differently in treebanks across languages
- are under-researched

 University of
Zurich^{UZH}

Institute of Computational Linguistics

Suggestions for WG 4 Tasks

Perform contrastive state-of-the-art survey

... for as many treebanks as possible

Develop recommendations for

... retrospective annotation of MWEs in existing treebanks
... future annotation of MWEs in treebanks
... and other types of corpora

19.09.13 Martin Volk Page 43

 University of
Zurich^{UZH}

Institute of Computational Linguistics

Thank you!

19.09.13 Martin Volk Page 44