



Using Machine Translation Techniques for Paraphrasing

Chris Callison-Burch

April 10, 2014

with Benjamin Van Durme, Juri Ganitkevitch, Ellie Pavlick, Wei Xu,
Courtney Napolis, Xuchen Yao, Peter Clark, Jonny Weese,
Matt Post, Tsz Ping Chan, Rui Wang, Trevor Cohn, Mirella Lapata and
Colin Bannard

Paraphrases

Differing textual expressions of the same meaning:

cup ↔ mug

the king's speech ↔ His Majesty's address

X ↔ X

one JJ instance of NP ↔ a JJ case of NP

Paraphrasing in NLP

Recognition or generation of paraphrases plays a part in...

...information extraction, question answering, entailment recognition, summarization, translation, compression, simplification, natural language generation, etc.

Data-Driven Paraphrasing

Monolingual parallel: English – English

Monolingual comparable: English ~ English

Plain monolingual: English

Bilingual parallel: English – French



What a scene! Seized by the tentacle and **glued to** its suckers, the unfortunate man was **swinging in the air** at the **mercy** of this enormous appendage. He gasped, he choked, he yelled: "Help! Help!" I'll hear his **harrowing plea** the rest of my life!
The **poor fellow** was **done for**.

What a scene! The unhappy man, seized by the tentacle and **fixed to** its suckers, was **balanced in the air** at the **caprice** of this enormous trunk. He rattled in his throat, he was stifled, he cried, "Help! help!" That **heart-rending cry**! I shall hear it all my life.
The **unfortunate man** was **lost**.

Paraphrasing with parallel monolingual data

Barzilay and McKeown (2001) identify paraphrases using identical contexts in aligned sentences:

Emma burst into tears and he tried to comfort her,
saying things to make her smile.

Emma cried and he tried to console her, adorning
his words with puns.

burst into tears = cried and comfort = console

Paraphrasing with comparable texts

Dolan, Quirk, and Brockett (2004) extract sentential paraphrases from newspaper articles published on the same topic and date:

On its way to an extended mission at Saturn, the Cassini probe on Friday makes its closest rendezvous with Saturn's dark moon Phoebe.

The Cassini spacecraft, which is en route to Saturn, is about to make a close pass of the ringed planet's mysterious moon Phoebe.

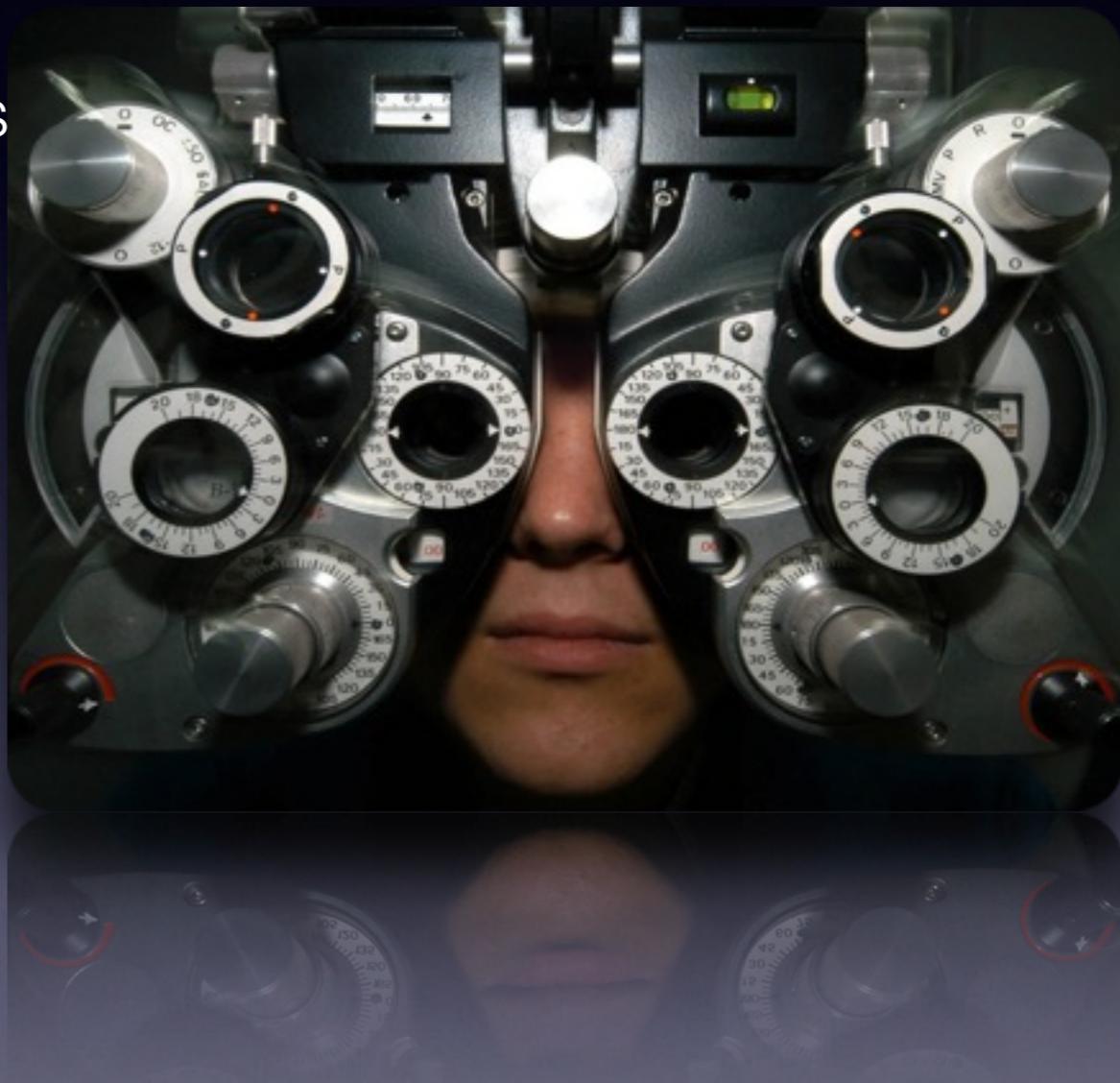
Distributional Hypothesis

If we consider **oculist** and **eye-doctor** we find that, as our corpus of utterances grows, these two occur in almost the same environments. In contrast, there are many sentence environments in which **oculist** occurs but **lawyer** does not...

It is a question of the relative frequency of such environments, and of what we will obtain if we ask an informant to substitute any word he wishes for **oculist** (not asking what words have the same meaning).

These and similar tests all measure the probability of particular environments occurring with particular elements... If A and B have almost identical environments we say that they are synonyms.

–Zellig Harris (1954)



DIRT

Lin and Panel (2001) operationalize the Distributional Hypothesis using dependency relationships to define similar environments.

Duty and responsibility share a similar set of dependency contexts in large volumes of text:

modified by adjectives	objects of verbs
additional, administrative, assigned, assumed, collective, congressional, constitutional ...	assert, assign, assume, attend to, avoid, become, breach ...

My focus: Paraphrasing & Translation

Translation is re-writing a text using words in a different language.

Paraphrasing is translation into the same language.

Inspiration from Statistical Machine Translation

We reuse & adapt:

Training data + alignment algorithms

Models + feature functions

Parameter estimation

Decoder

Bilingual Data

Sentence-aligned parallel corpora in English and any foreign language

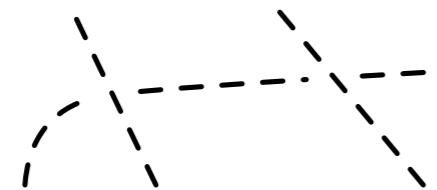
Available in large quantities

Strong meaning equivalence signal

... but different languages.

Bilingual Pivoting

... 5 farmers were



... fünf Landwirte

thrown into jail

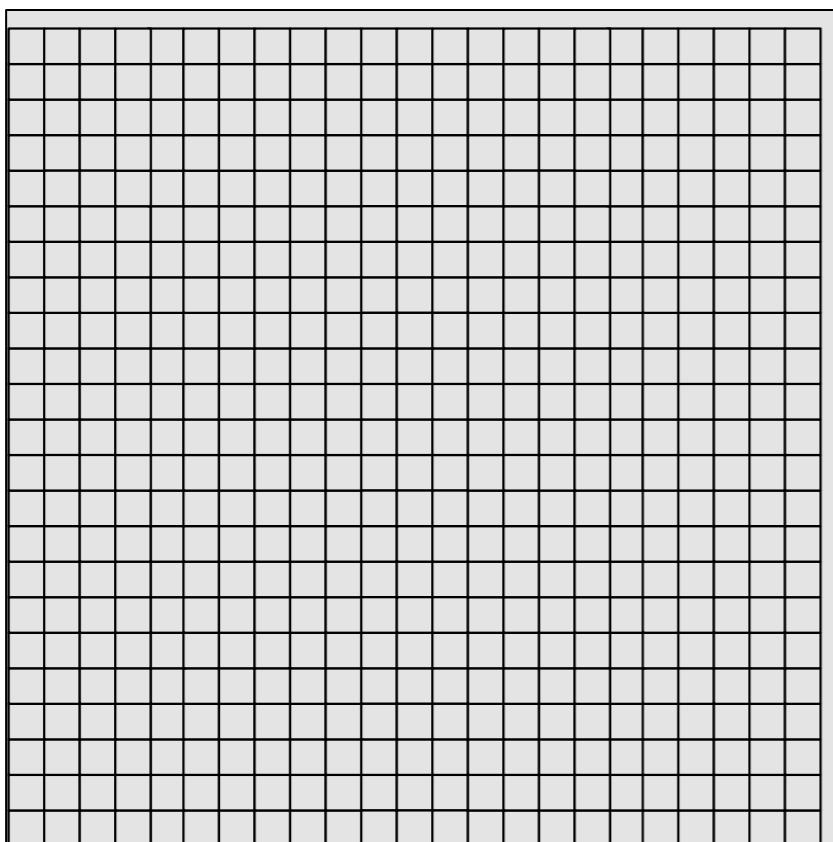
festgenommen

in Ireland ...

, weil ...

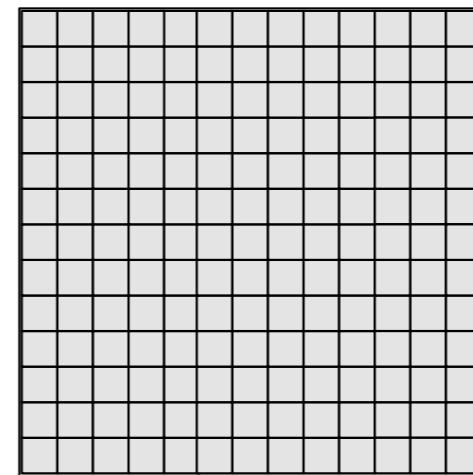
Large, diverse sets of bilingual training data

1000M



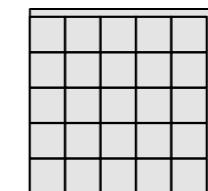
French-English
 10^9 word webcrawl

2 languages @
250M each



DARPA
GALE Program

21 languages @
50-80M each



European
Parliament

Wide range of paraphrases

thrown into jail

arrested

be thrown in prison

arrest

detained

been thrown into jail

cases

imprisoned

being arrested

custody

incarcerated

in jail

maltreated

jailed

in prison

owners

locked up

put in prison for

protection

taken into custody

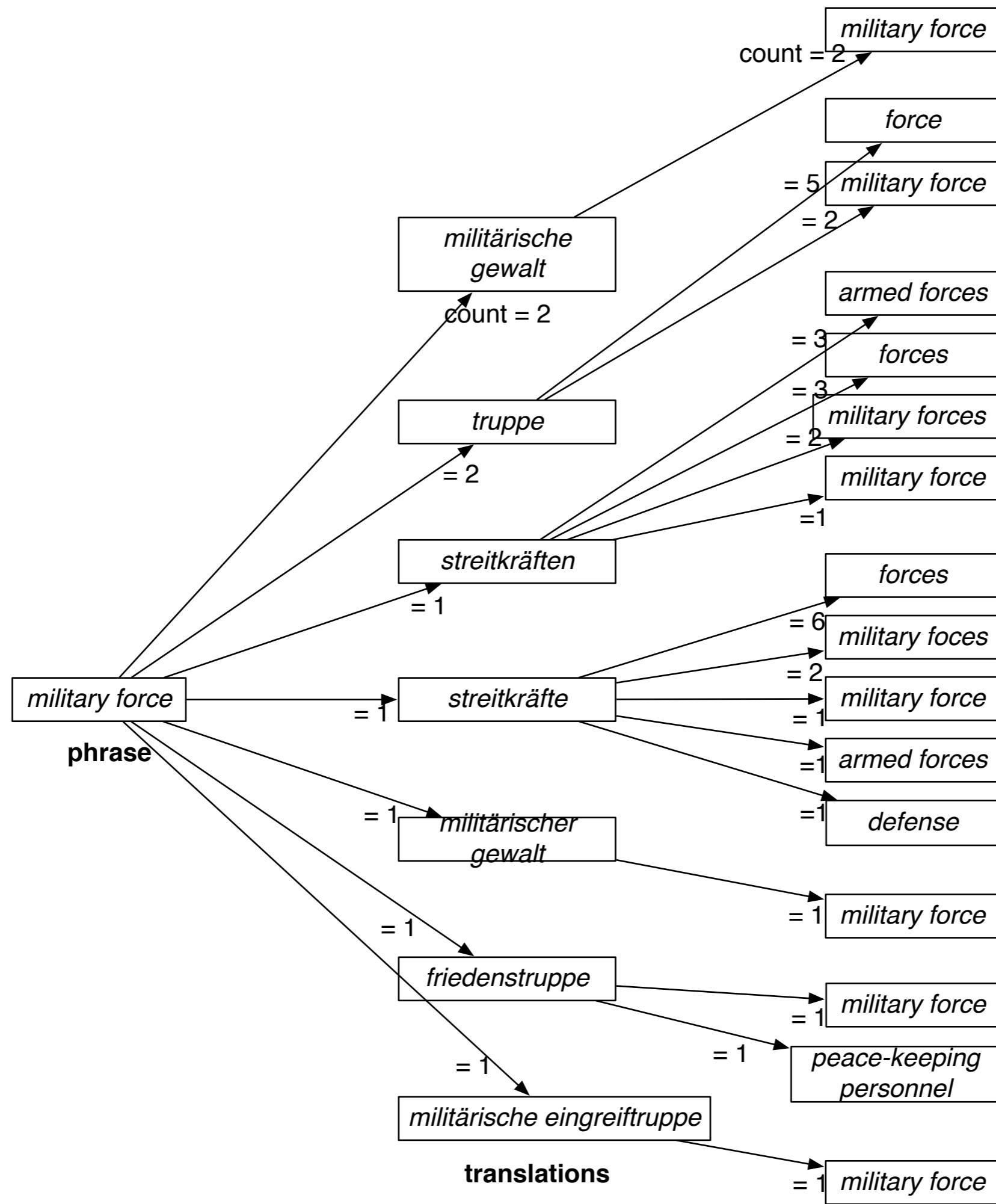
were thrown into jail

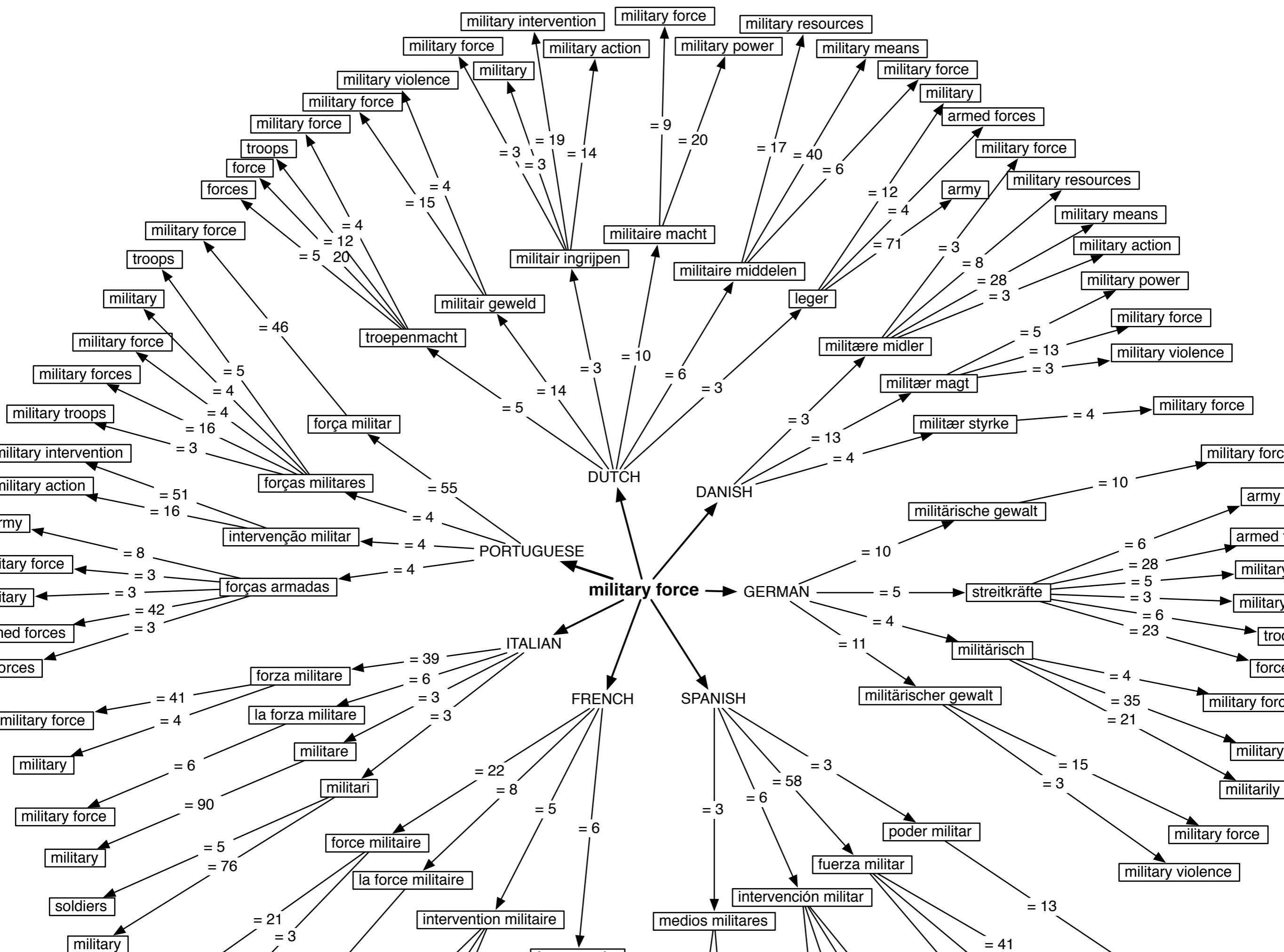
thrown

thrown into prison who are held in detention

Paraphrase Probability

$$\begin{aligned} p(e_2|e_1) &= \sum_f p(e_2, f|e_1) \\ &= \sum_f p(e_2|f, e_1)p(f|e_1) \\ &\approx \sum_f p(e_2|f)p(f|e_1) \end{aligned}$$





Syntactic constraints

thrown into jail

arrested

be thrown in prison

arrest

detained

been thrown into jail

cases

imprisoned

being arrested

custody

incarcerated

in jail

maltreated

jailed

in prison

owners

locked up

put in prison for

protection

taken into custody

were thrown into jail

thrown

thrown into prison

who are held in detention

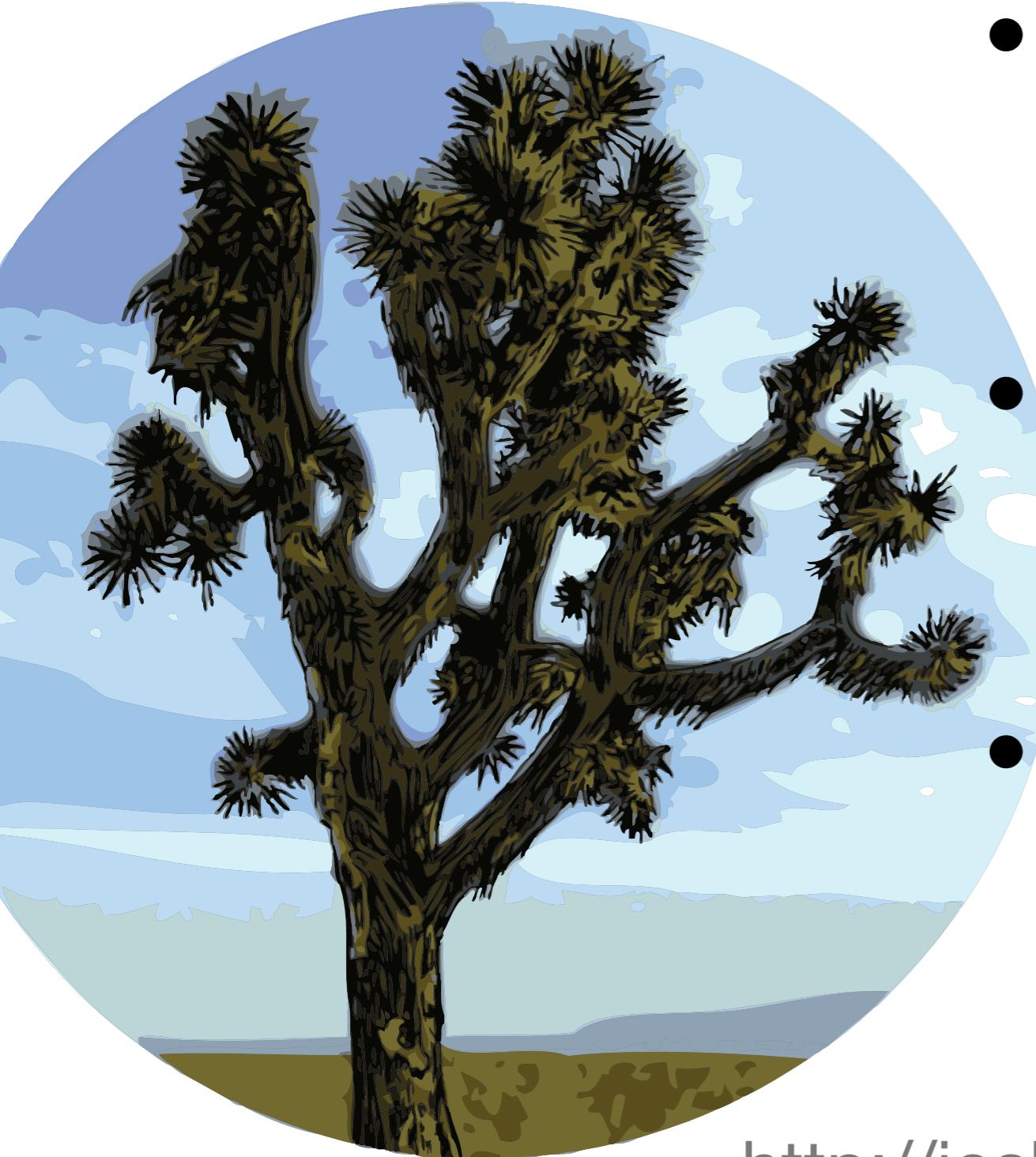
Sentential paraphrases from bitexts?

Bilingual parallel corpora provide an excellent source
of lexical and phrasal paraphrases.

Sentential | structural paraphrases are more
obviously learned from English-English sentence
pairs.

Can we learn structural paraphrases from bitexts?
How should we represent them?

Syntactic MT in the Joshua Decoder



- Synchronous context free grammars generate pairs of corresponding strings
- Can be used to describe translation and re-ordering between languages
- Because Joshua uses SCFGs, it translates sentences by parsing them

Example SCFG for translation

	Urdu	English
$S \rightarrow$	$NP\textcircled{1} \ VP\textcircled{2}$	$NP\textcircled{1} \ VP\textcircled{2}$
$VP \rightarrow$	$PP\textcircled{1} \ VP\textcircled{2}$	$VP\textcircled{2} \ PP\textcircled{1}$
$VP \rightarrow$	$V\textcircled{1} \ AUX\textcircled{2}$	$AUX\textcircled{2} \ V\textcircled{1}$
$PP \rightarrow$	$NP\textcircled{1} \ P\textcircled{2}$	$P\textcircled{2} \ NP\textcircled{1}$
$NP \rightarrow$	<i>hamd ansary</i>	<i>Hamid Ansari</i>
$NP \rightarrow$	<i>na}b sdr</i>	<i>Vice President</i>
$V \rightarrow$	<i>namzd</i>	<i>nominated</i>
$P \rightarrow$	<i>kylye</i>	<i>for</i>
$AUX \rightarrow$	<i>taa</i>	<i>was</i>

NP1
hamd ansary

NP2
na}b sdr

P3
kylye

V4
namzd

AUX5
taa

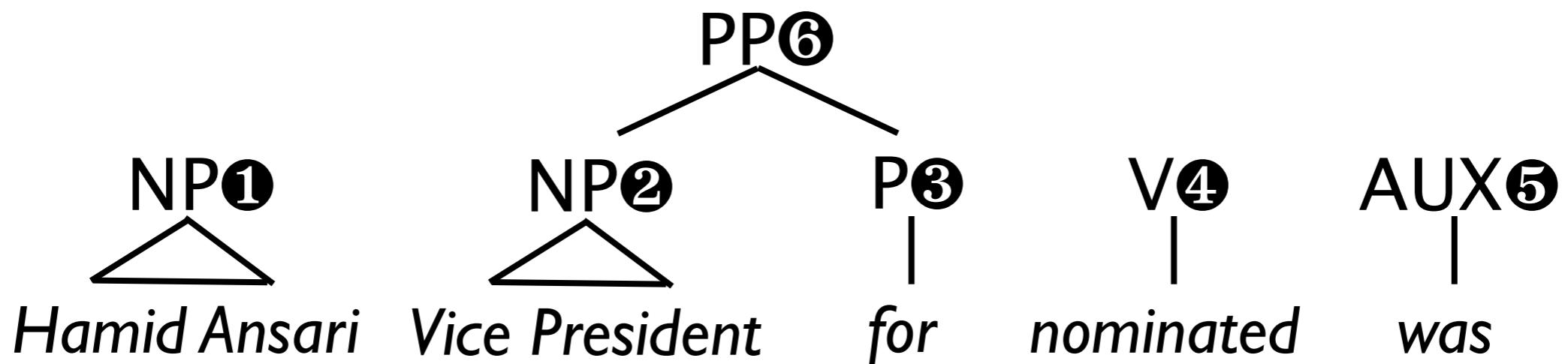
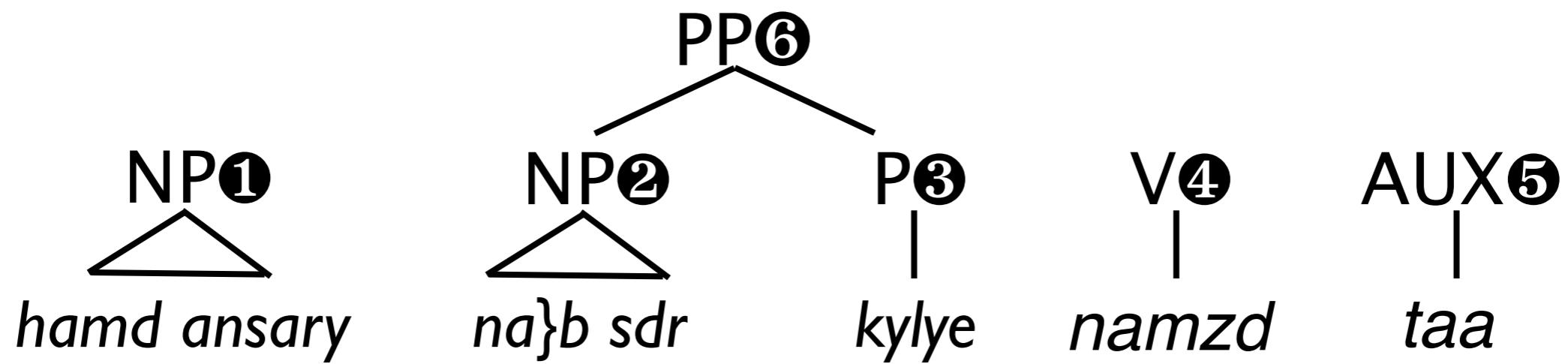
NP1
Hamid Ansari

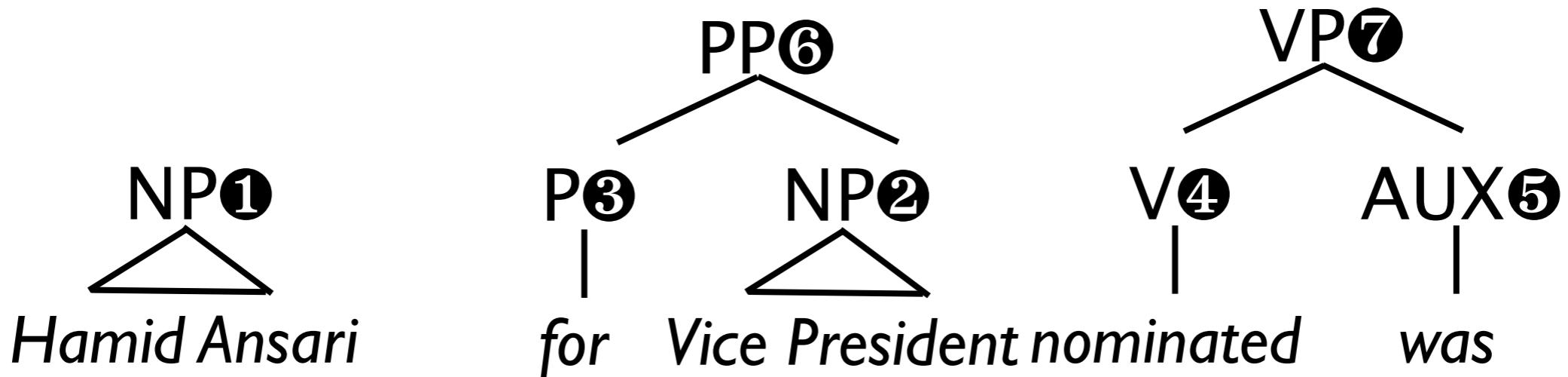
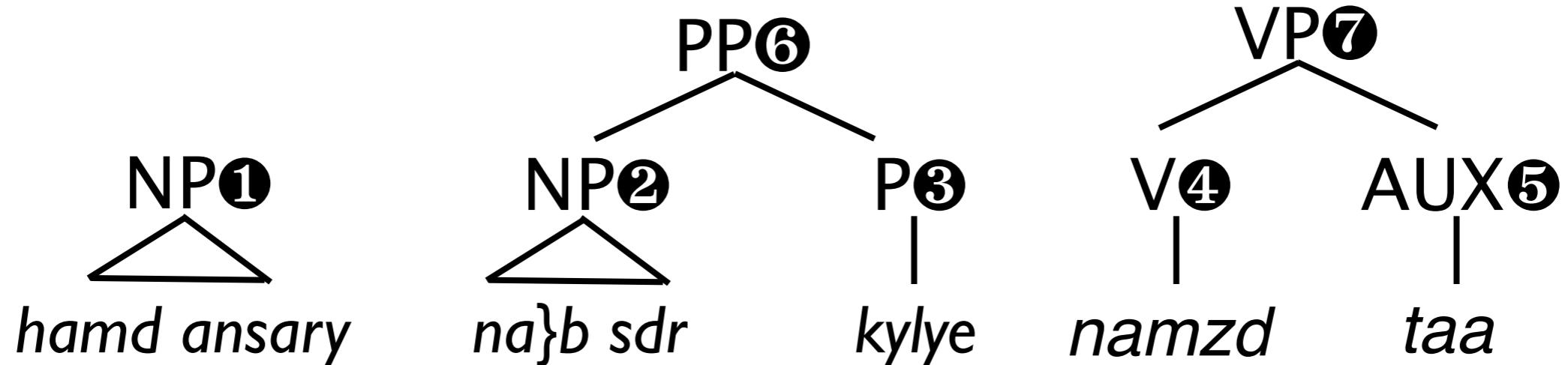
NP2
Vice President

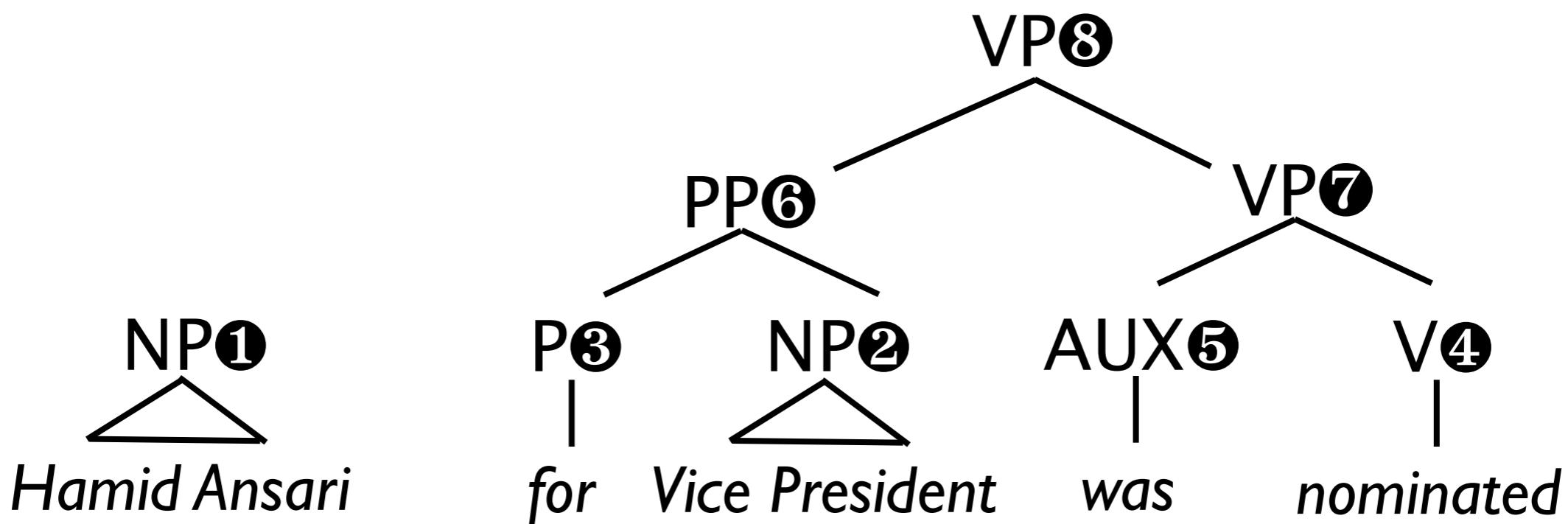
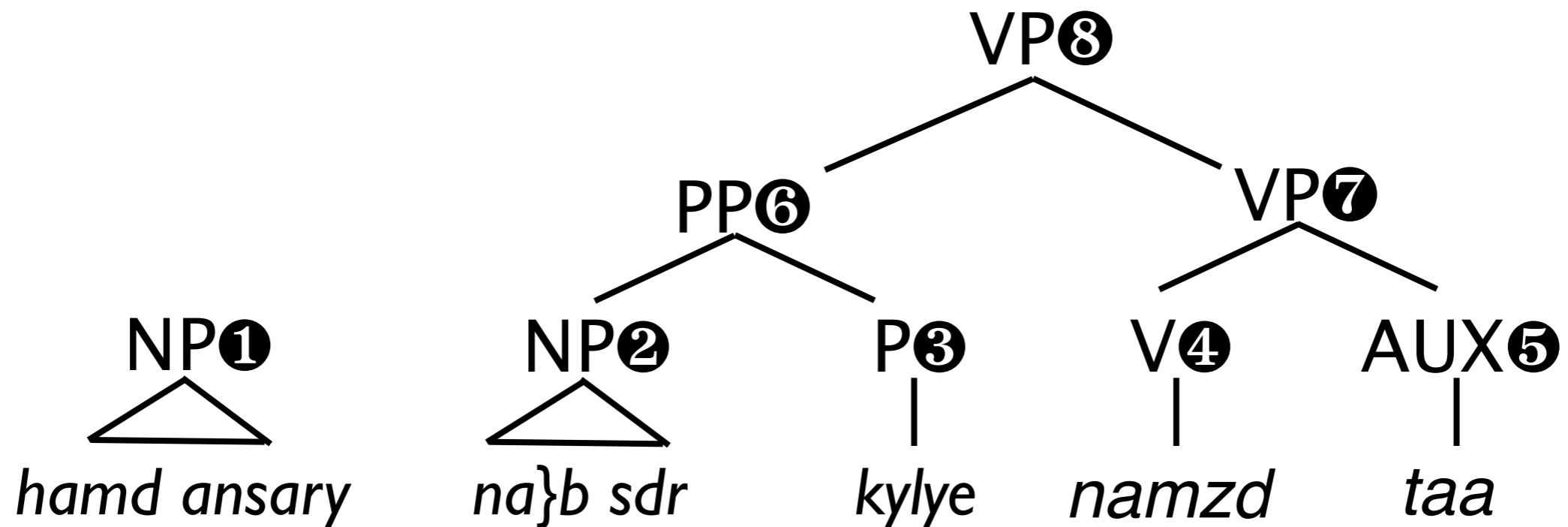
P3
for

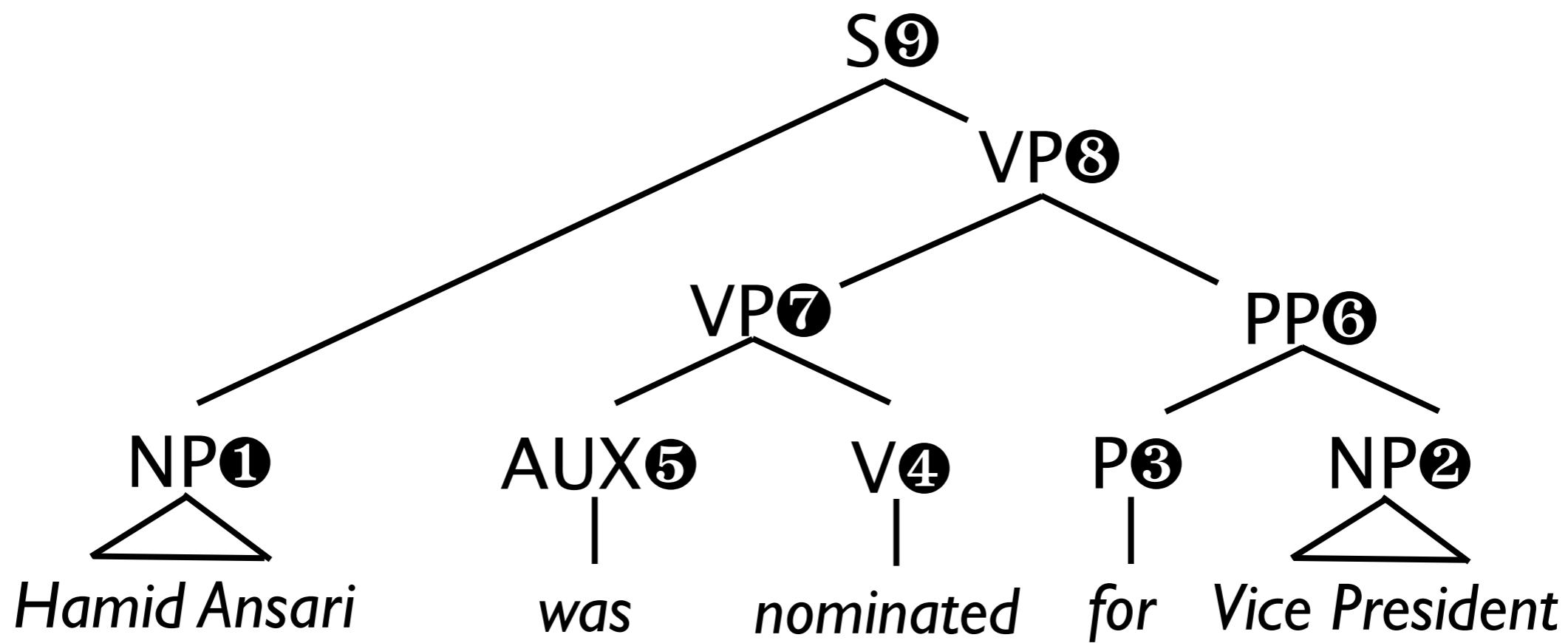
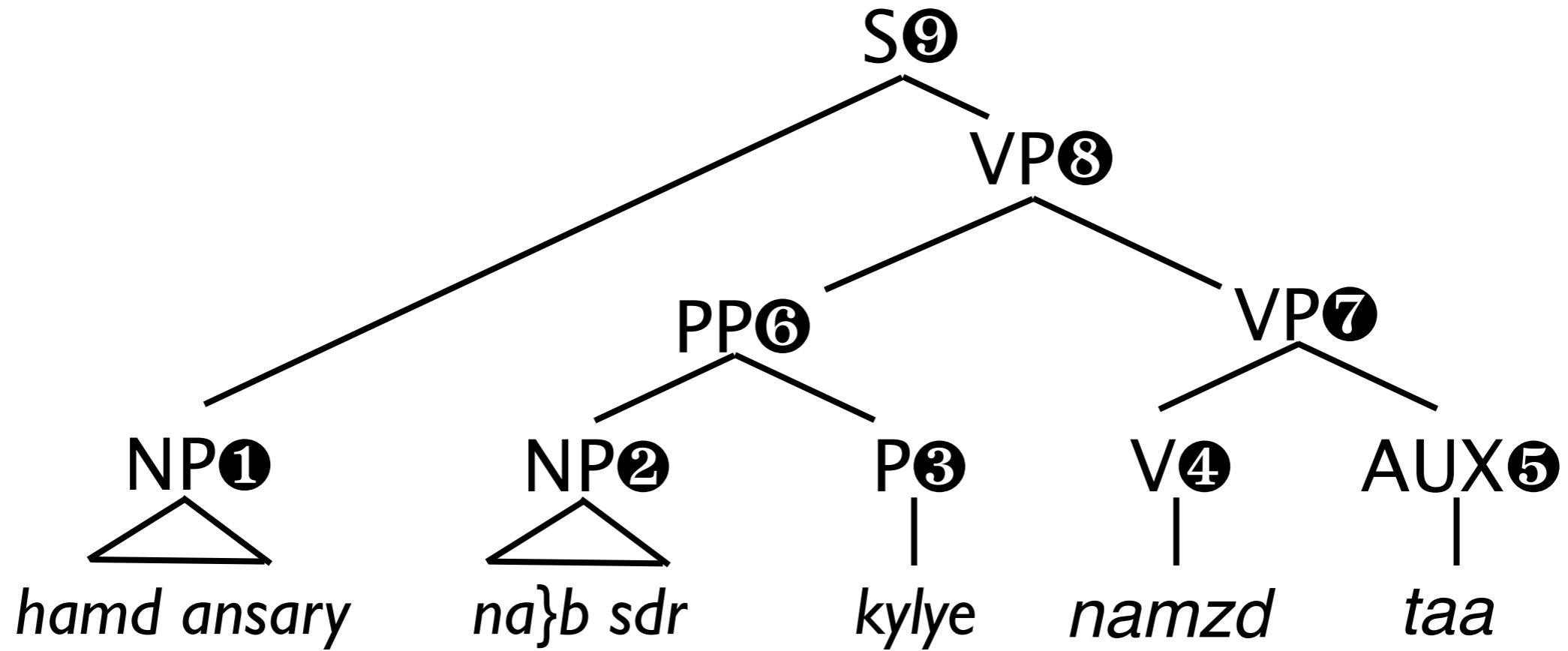
V4
nominated

AUX5
was









SCFGs via Pivoting

- Adapting our syntactic MT models, we learn structural transformations, like the English possessive rule

$$\begin{array}{l} \text{NP} \rightarrow \quad \text{NP}'s \text{ NN} \mid \text{le NN de NP} \\ \\ \text{NP} \rightarrow \quad \text{the NN of NP} \mid \text{le NN de NP} \end{array}$$

combine to

$$\text{NP} \rightarrow \quad \text{NP}'s \text{ NN} \mid \text{the NN of NP}$$

Possessive rule	NP → NP → the NNS	the NN of the NNP the NNP's NN the NNS
Dative shift	VP → VP →	give NN to NP give NP the NN provide NP give NP
Adv. adj. phrase move	S VP → S →	ADVP they VBD they VBD ADVP it is ADJP VP VP is ADJP
Verb particle shift	VP →	VB NP up VB up NP
Reduced relative clause	SBAR S ADJP →	although PRP VBD that although PRP VBD very JJ that S JJ S
Partitive constructions	NP → NP →	CD of the NN CD NN all DT\NP all of the DT\NP
Topicalization	S →	NP, VP. VP, NP.
Passivization	SBAR →	that NP had VBN which was VBN by NP
Light verbs	VP → VP →	take action ADVP to act ADVP to make a decision PP to decide PP

Text-to-Text Generation

T2T involves generating meaning-equivalent text that is *subject to some constraints*:

sentence compression, *shorter*

simplification, *easier to understand*

poetry from prose, *rhyme and meter*

Sentence Compression

Reduce length of a sentence (#chars) while retaining the meaning

$$\text{Compression ratio: } \varphi = \frac{\text{length}_{\text{compression}}}{\text{length}_{\text{original}}}$$

Paraphrasing as a task and problem is of paramount importance to a multitude of applications in the field of NLP.

Sentence Compression

Reduce length of a sentence (#chars) while retaining the meaning

$$\text{Compression ratio: } \varphi = \frac{\text{length}_{\text{compression}}}{\text{length}_{\text{original}}}$$

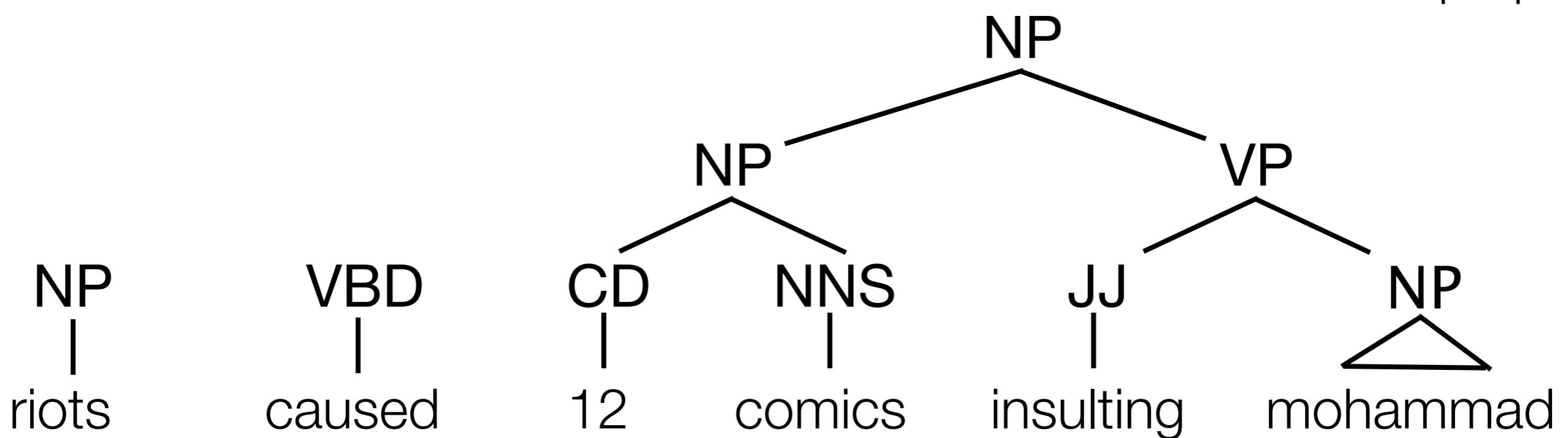
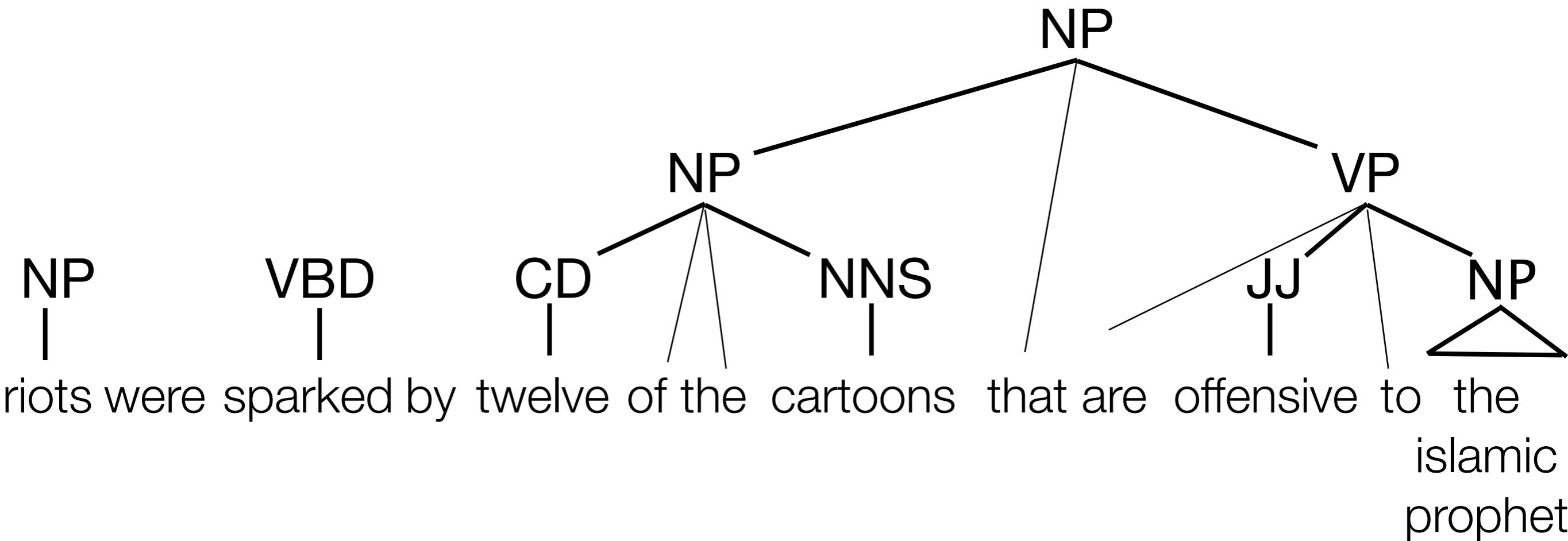
~~Paraphrasing as a task and problem is of paramount importance to a multitude of applications in the field of NLP.~~
is awesome

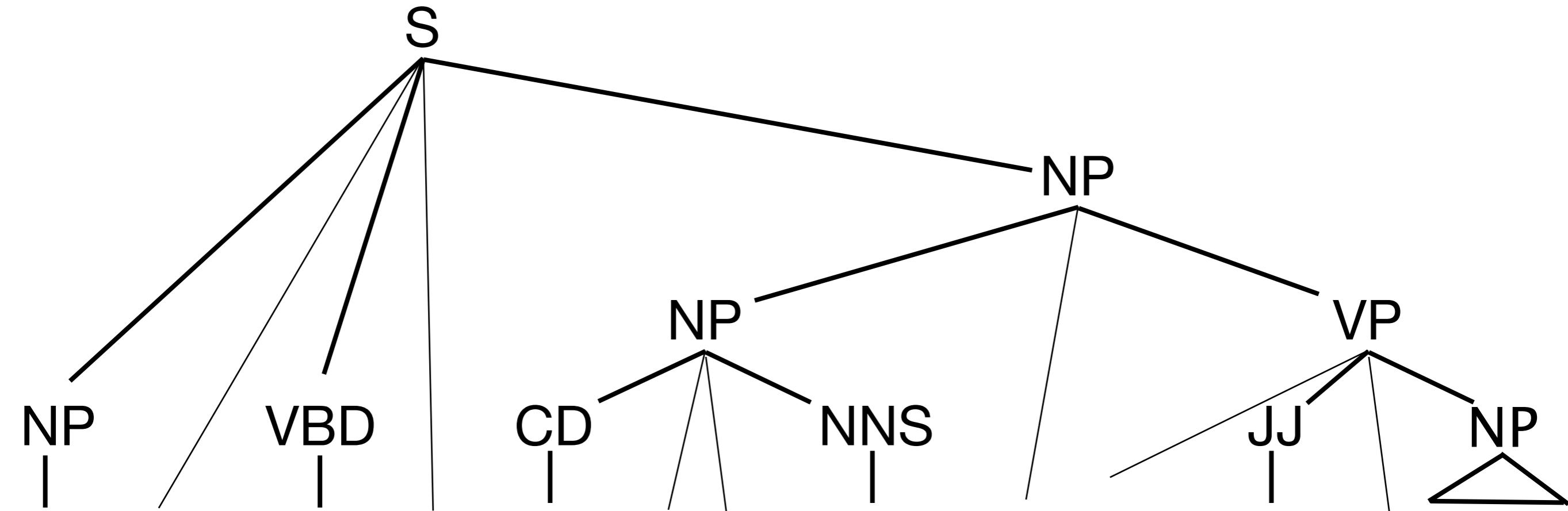
Paraphrase Grammar

	English	English
$S \rightarrow NP①$	were VBD by NP②	NP② VBD NP①
$NP \rightarrow$	NP that VP	NP VP
$VP \rightarrow$	are JJ to NP	JJ NP
$NP \rightarrow$	CD of the NNS	CD NNS
$CD \rightarrow$	twelve	12
$NNS \rightarrow$	cartoons	comics
$JJ \rightarrow$	offensive	insulting
$NP \rightarrow$	the islamic prophet	mohammed
$VBD \rightarrow$	sparked	caused

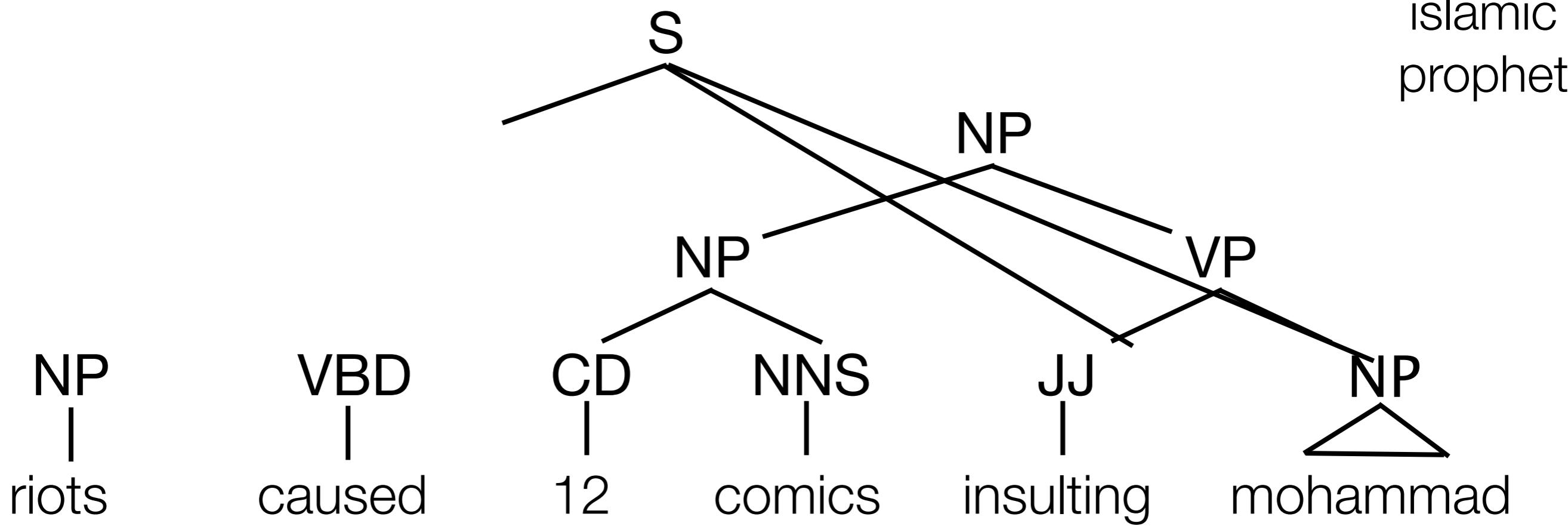
NP VBD CD NNS JJ NP
| | | | | |
riots were sparked by twelve of the cartoons that are offensive to the
islamic prophet

NP VBD CD NNS JJ NP
| | | | | |
riots caused 12 comics insulting mohammad

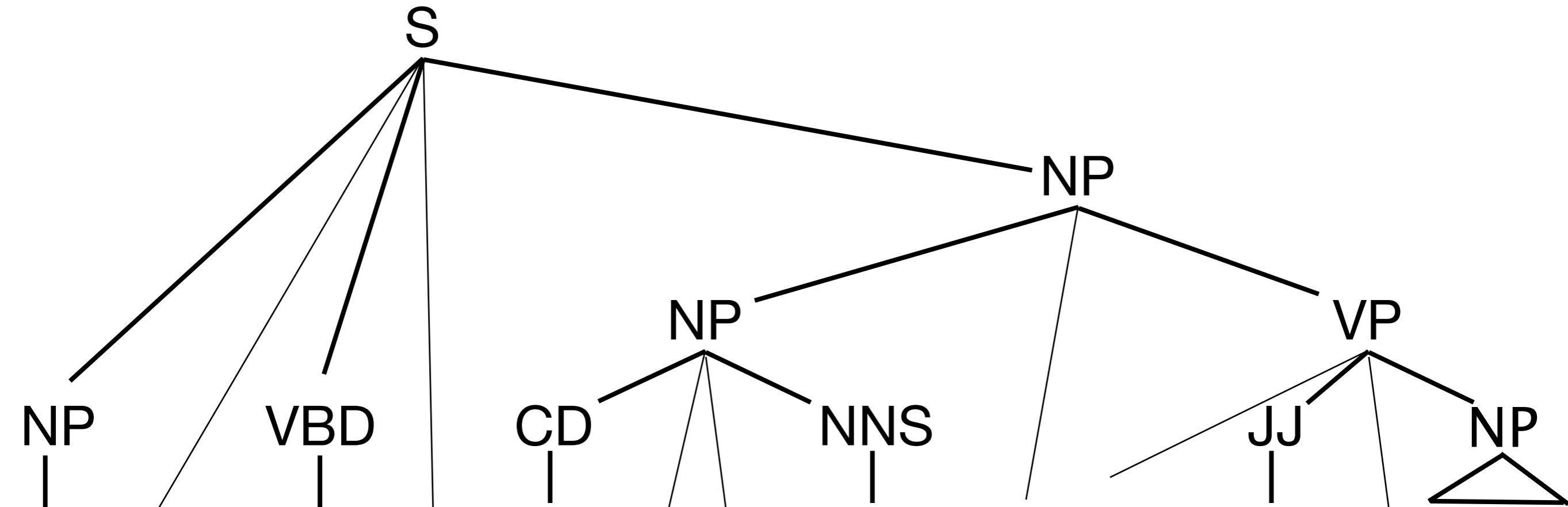




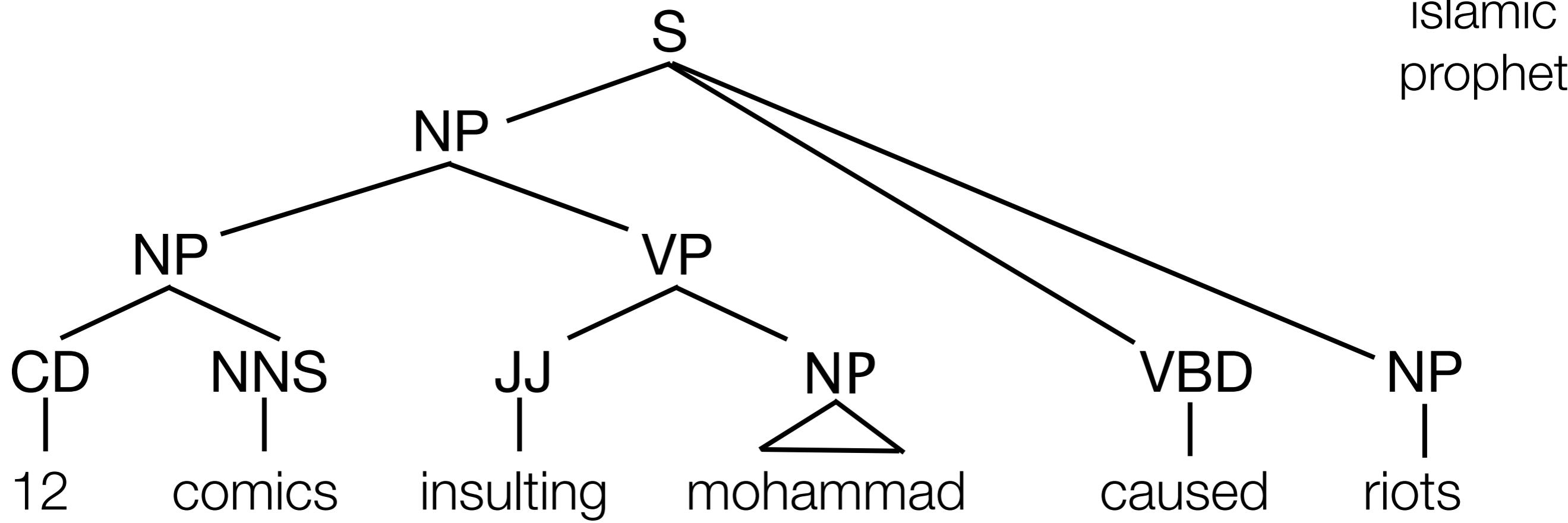
riots were sparked by twelve of the cartoons that are offensive to the
islamic prophet



riots caused 12 comics insulting mohammad



riots were sparked by twelve of the cartoons that are offensive to the
islamic prophet



12 comics insulting mohammad caused riots

Text-to-Text Applications

Claim:

Paraphrasing is suitable to tackle sentential text-to-text tasks, and we can re-use SMT machinery for T2T.

However:

Naive application of MT techniques will not work, need to adapt them

Task Adaptation

SMT	T2T
Naive application of the MT machinery to the task	Task-specific adaptations

- Development data
- Objective function
- Feature set
- Grammar augmentations

Development Data

SMT	T2T
English reference translations that are used to calculate BLEU for SMT.	Selected pairs of reference translations that significantly differ in length.

and he said that the project **will cover** the needs of the region in the long term.

82

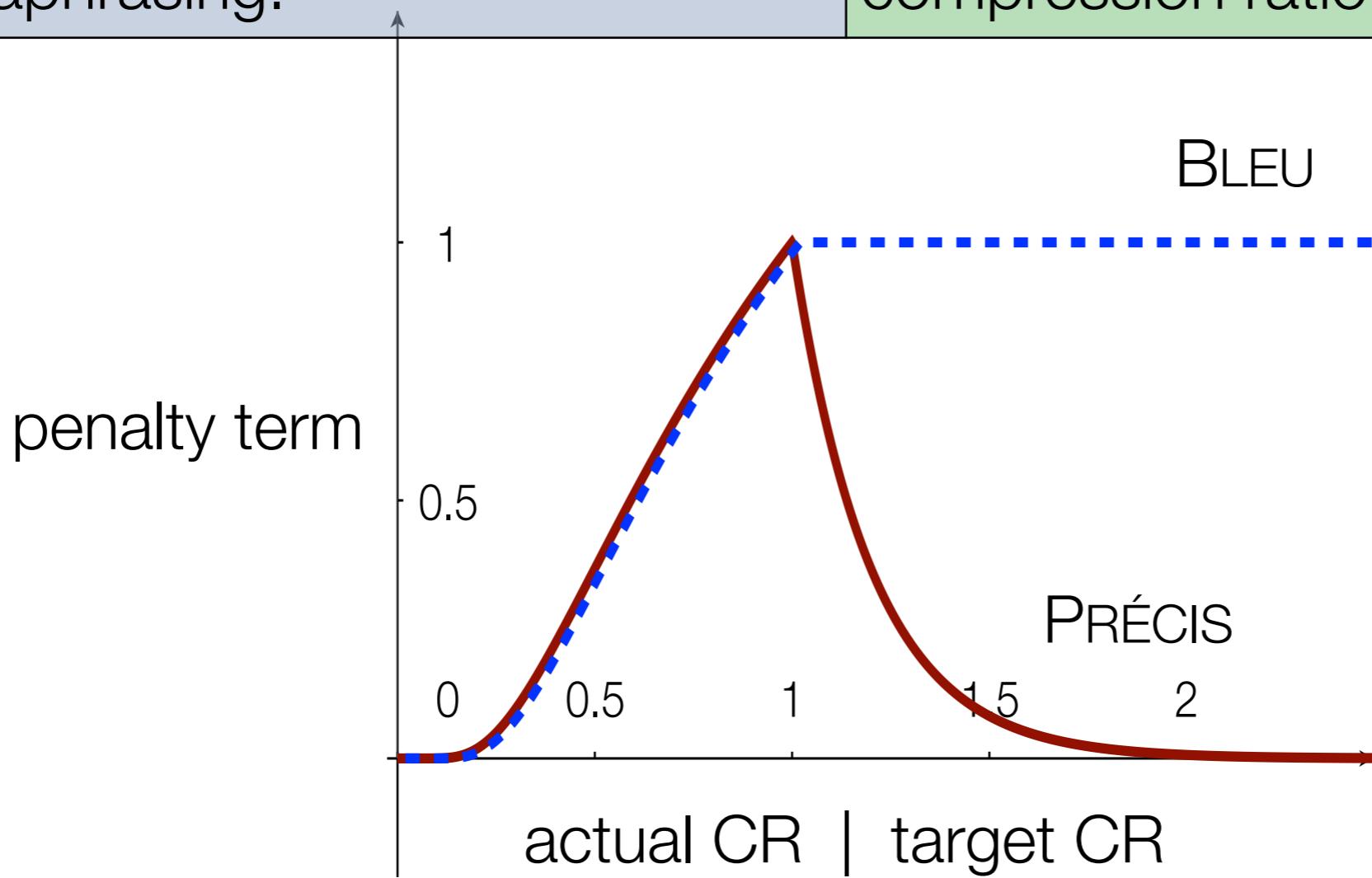
he said the project **includes** all the district's long-term needs.

65

compression ratio = 0.79

Objective Function

SMT	T2T
Optimized for English-to-English BLEU score. Causes self-paraphrasing.	Add a “verbosity penalty” to BLEU that allows a target compression ratio to be set.



Features

SMT	T2T
Phrasal and lexical probabilities quantify general paraphrase quality.	Features counting number of source and target words and the difference between them.

$\text{VP} \rightarrow \text{NP}$ was eaten by NN | NN ate NP

$$p(e_1|e_2) = 0.1 \quad c_{e_1} = 14 \quad c_{e_2} = 5 \quad \log CR = \log \frac{c_{e_1}}{c_{e_2}}$$
$$c_{diff} = -9$$

Augmentations

SMT	T2T
It is not typical for additional task-specific rules to be added in the standard SMT pipeline.	Augment the grammar with deletion rules for specific POS (JJ, RB, DT) allowing for shorter compressions.

$JJ \rightarrow \text{superfluous} \mid \epsilon$

$RB \rightarrow \text{redundantly} \mid \epsilon$

$DT \rightarrow \text{the} \mid \epsilon$

Monolingually-derived Features

SMT	T2T
All features, aside from the LM, are bilingually derived.	Calculate distributional similarity of paraphrase pairs from monolingual data

Orthogonal signal to bilingual pivoting

Even more data available

Incorporated as features in T2T model

Distributional Similarity

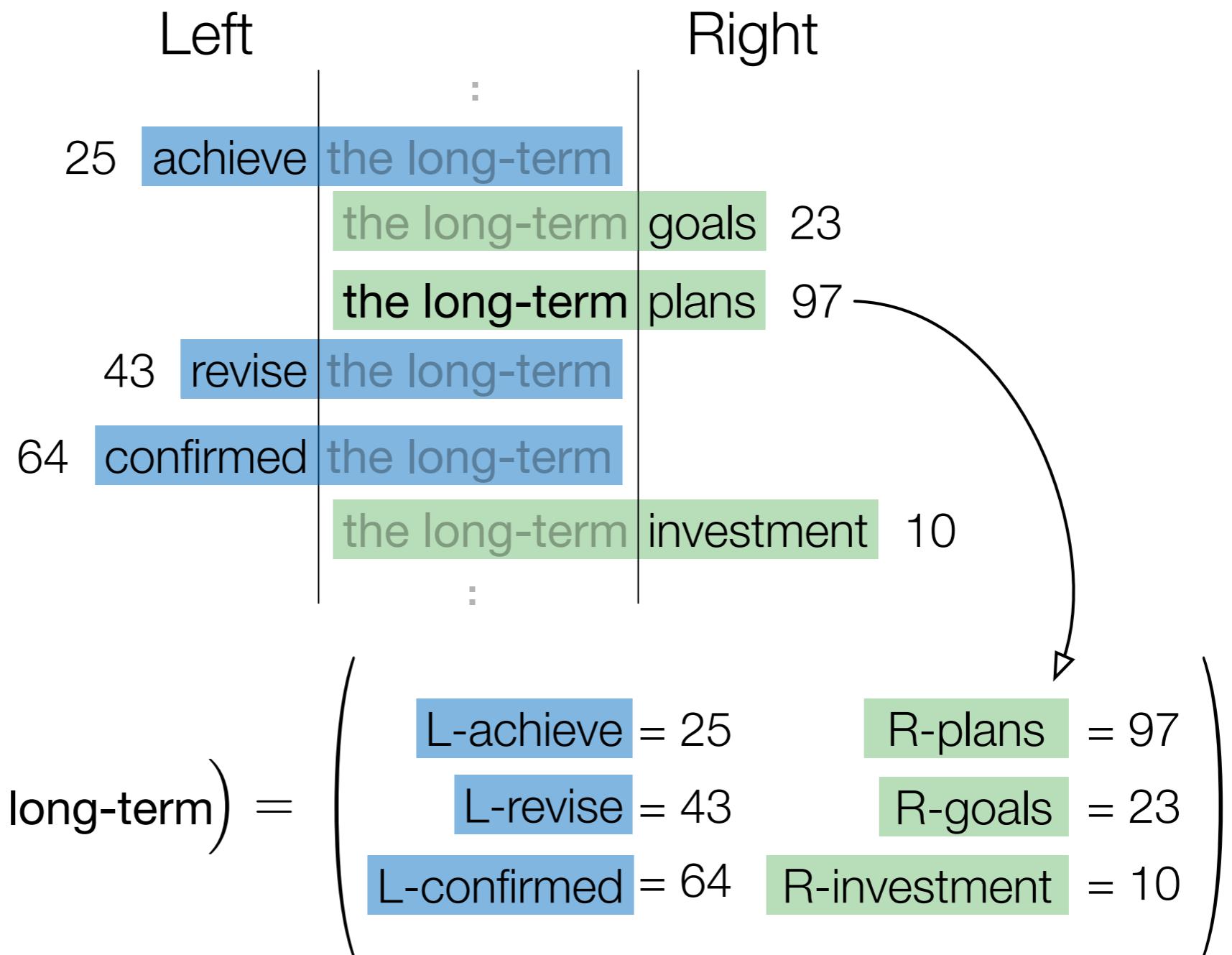
Idea: similar words occur in similar contexts.

Characterize words by their contexts

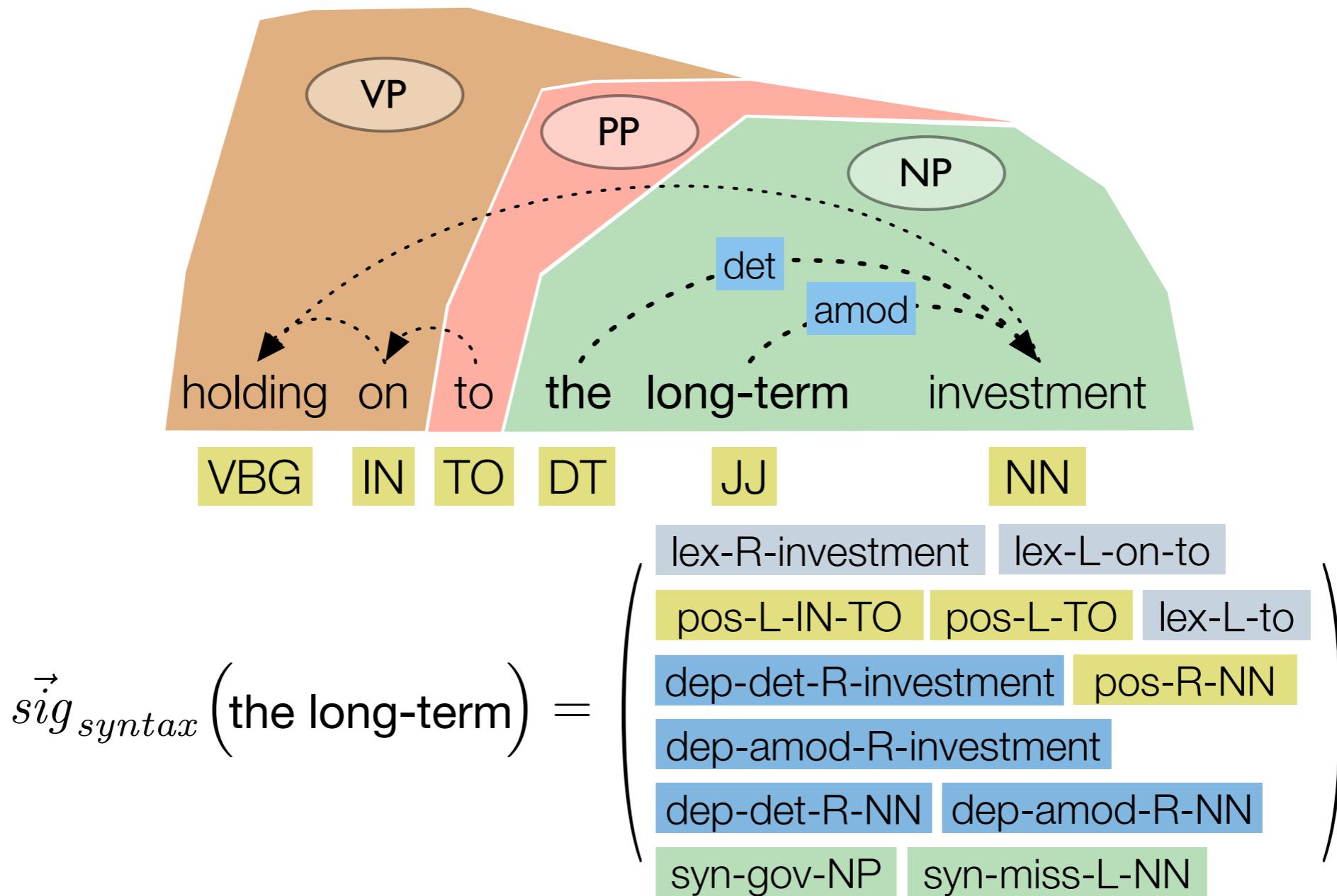
Contexts represented by co-occurrence vectors, similarity quantified by cosine

“Are these paraphrases substitutable?”

n -gram Context



Syntactic Context



Large Monolingual Data Sets

Google n-grams

Collection of 1 trillion tokens with counts

Based on vast amounts of text

Annotated Gigaword (AKBC-WEKEX '12)

Collection of 4 billion words, parsed and tagged

Task-based Evaluation

Evaluated paraphrases in the context of a T2T compression task.

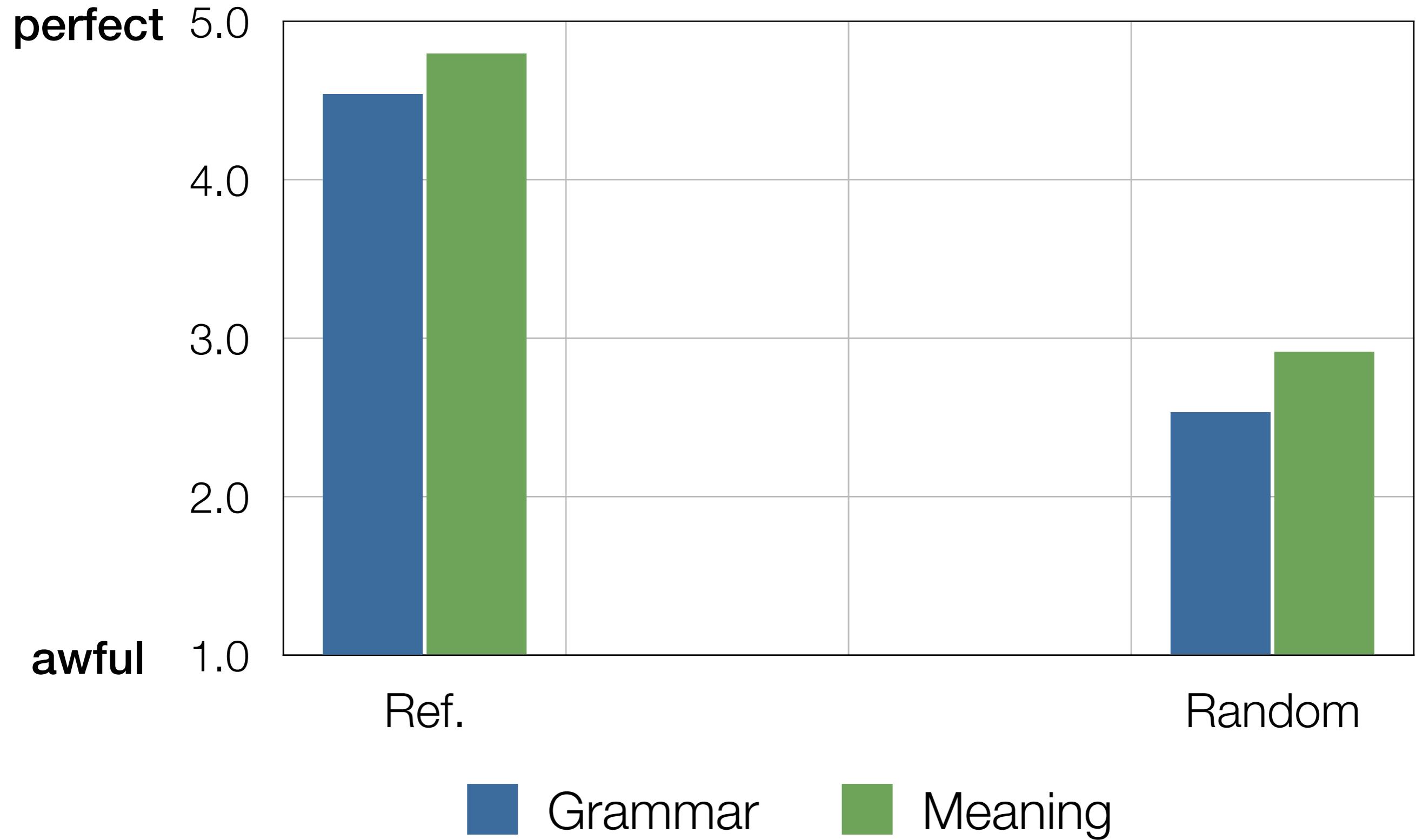
Compared against a state of the art system.

Human assessment (5-point scale):

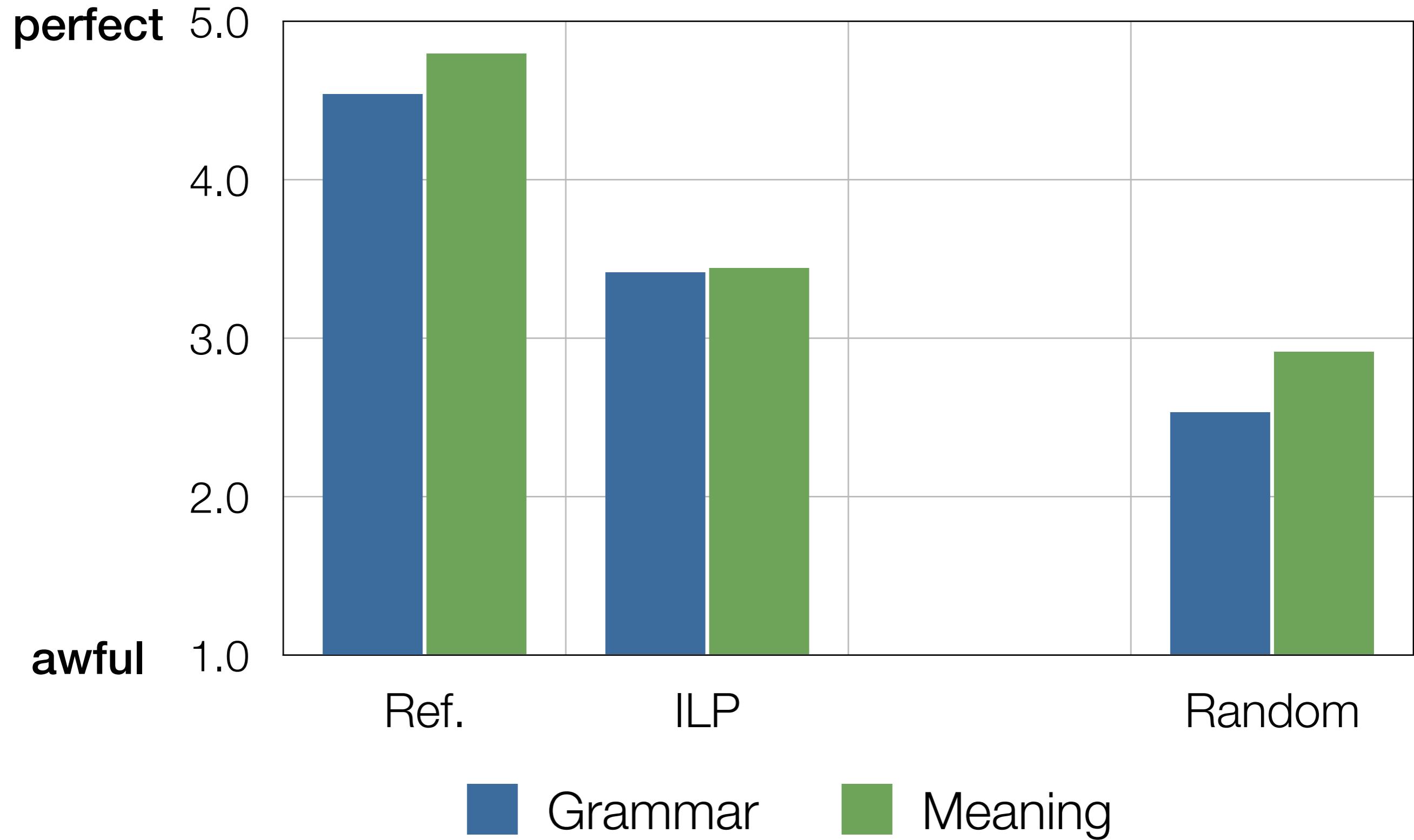
How well do these sentences retain the meaning of original?

How grammatical is the resulting sentence?

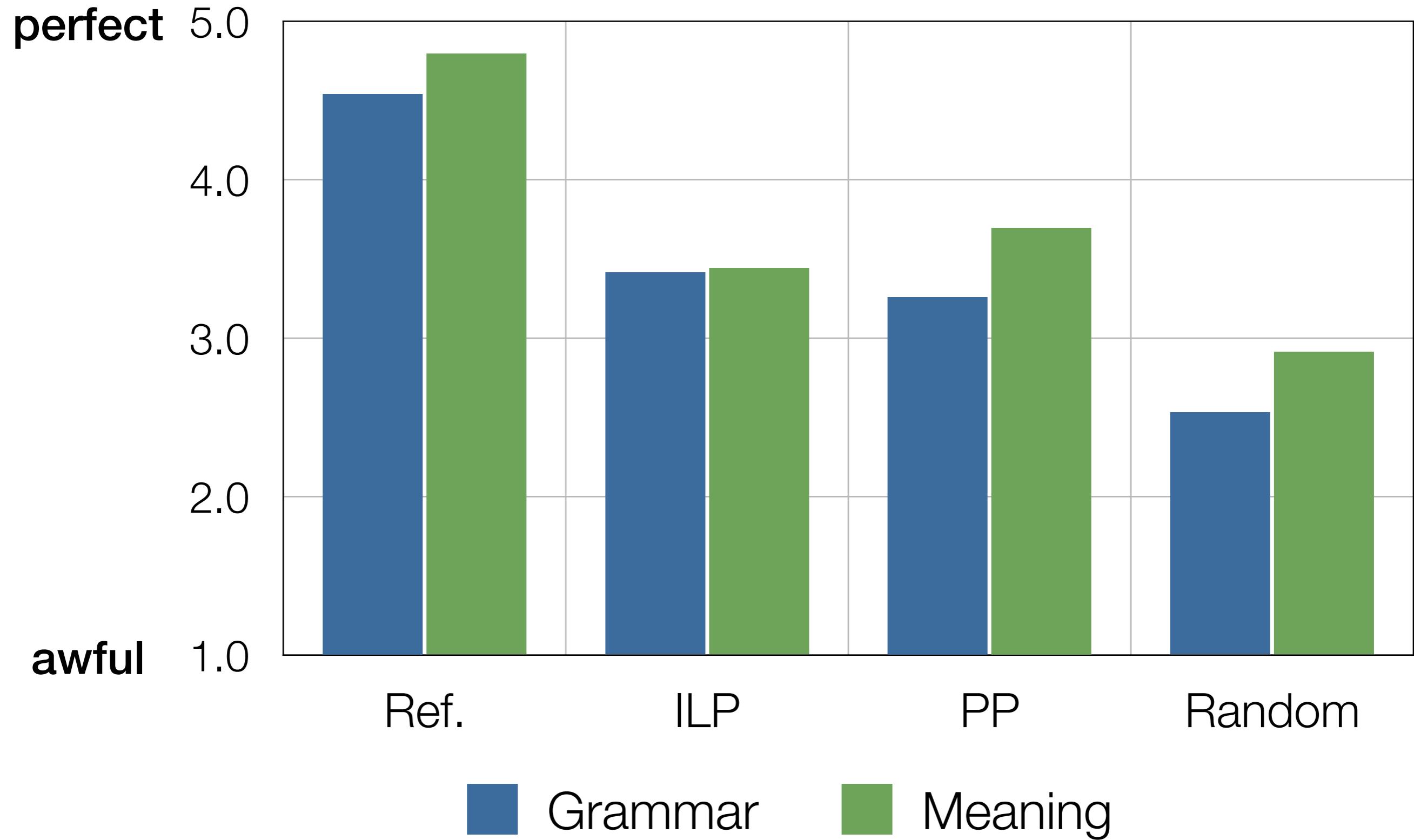
Compression Quality



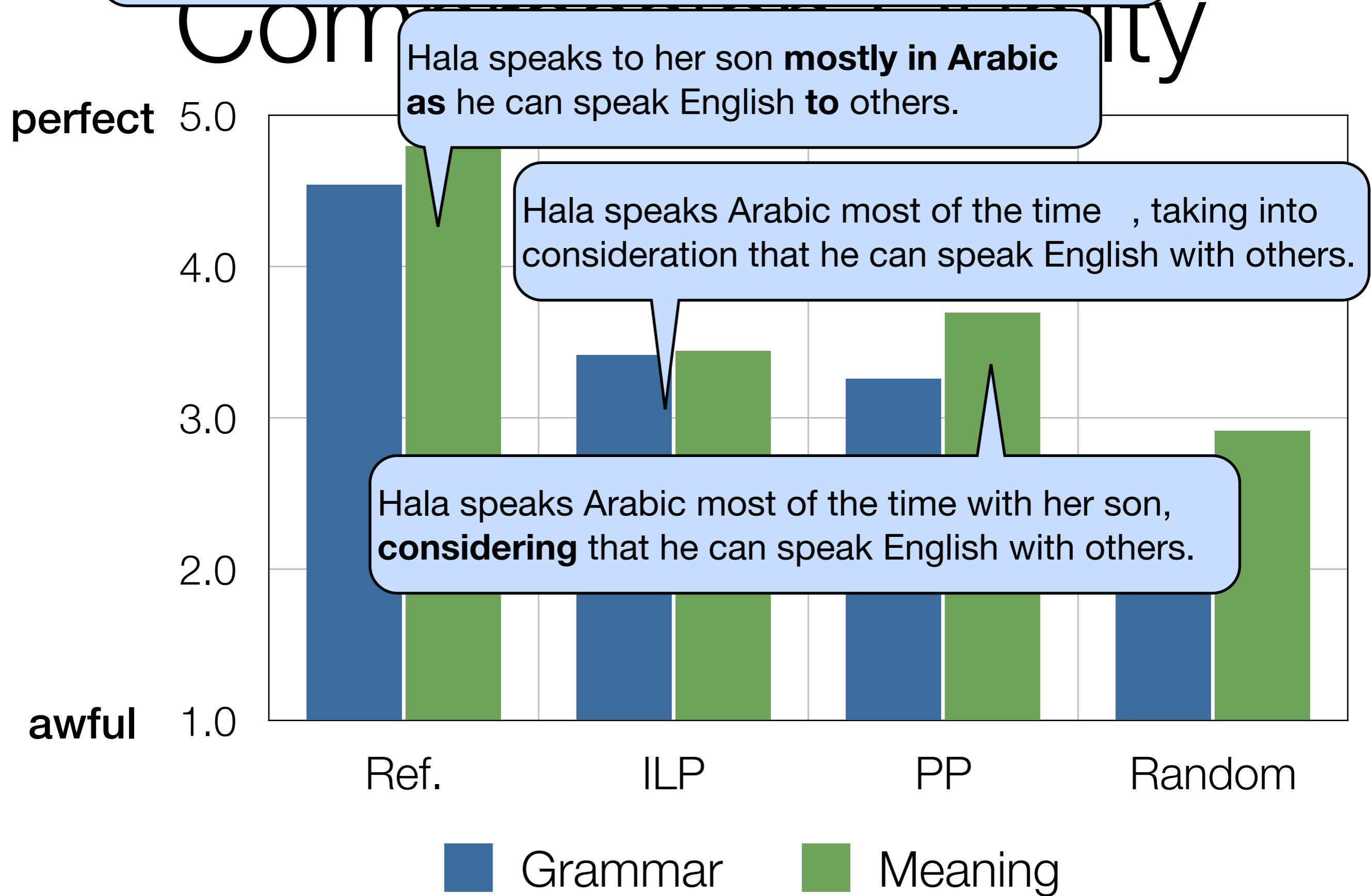
Compression Quality



Compression Quality



Input: Hala speaks Arabic most of the time with her son, taking into consideration that he can speak English with others.



Adaptation in 5 easy steps

Step	SMT to T2T Adaptation
1	Dev data: Collect a set of sentence pairs that reflects the task that you are trying to model
2	Objective function: Create a new objective function that indicates how well the system output the constraints of your task
3	Task-specific features: Add new features to the model that will allow it to score its own output for the task
4	Augment the grammar: Use your domain knowledge to add any rules that would not normally be contained in a paraphrase grammar.
5	Other features: Take advantage of the English to English to add other features that model grammaticality more generally.

Resources

Joshua Decoder



- An open source decoder that synchronous context free grammars to translate
- Implements all algorithms needed for translating with SCFGs
 - grammar extraction
 - chart-parsing
 - n-gram LM integration

<http://joshua-decoder.org>

PPDB: The Paraphrase Database

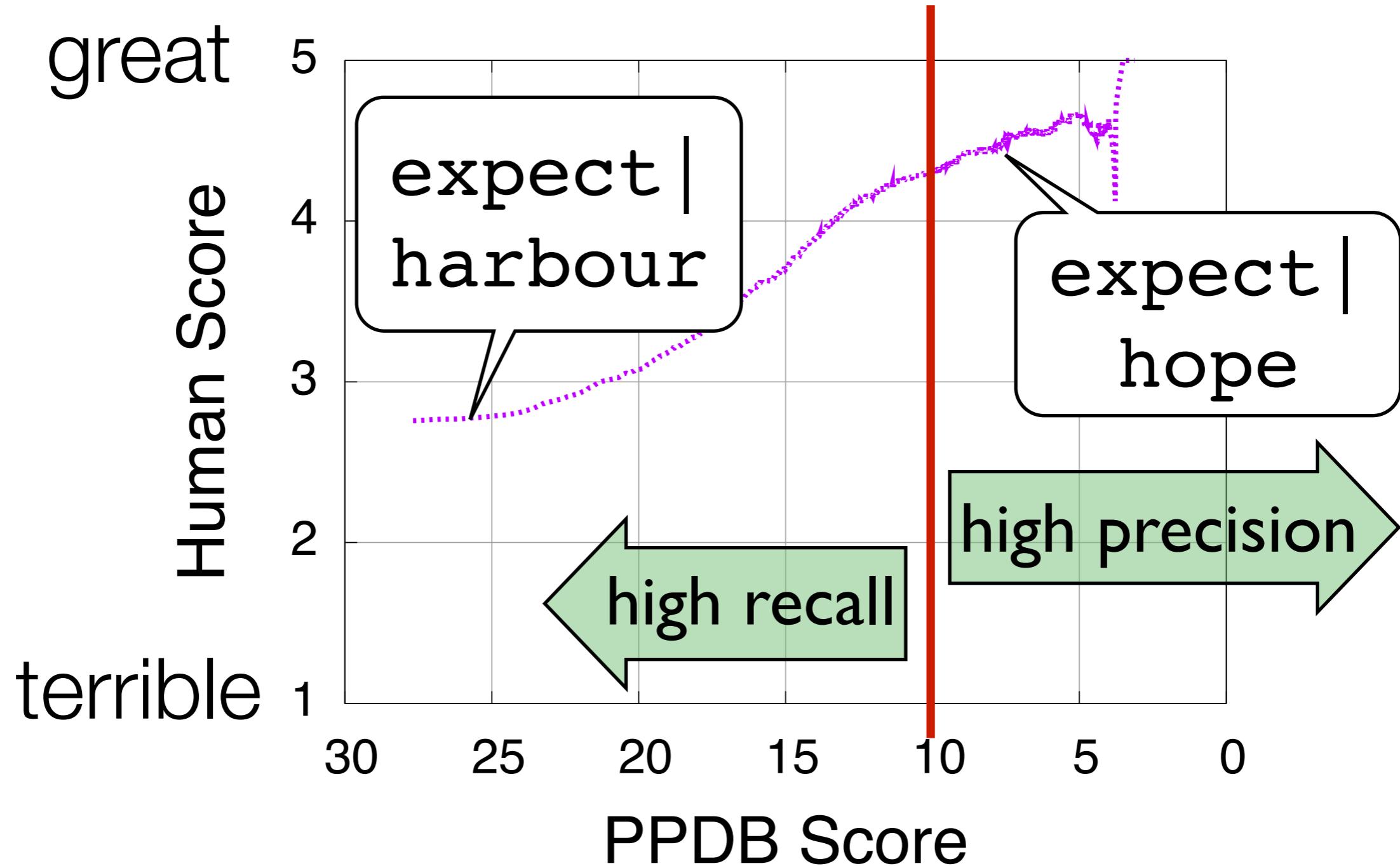
- A huge collection of paraphrases
- Extracted from 106 million sentence pairs,
2 billion English words, 22 pivot languages

	Paraphrases
Lexical	7.6 M
Phrasal	68.4 M
Syntactic	93.6 M
Total	169.6 M

<http://paraphrase.org>

	All	Lexical	Phrasal	Syntactic
S	Paraphrases (424MB, 6.8M rules)	Paraphrases (1.7MB, 31k rules)	Paraphrases (42MB, 637k rules)	Constituent (38MB, 585k rules)
		Identity (16MB, 437k rules)	Identity (170MB, 4.1M rules)	Non-Constituent (343MB, 5.6M rules)
M	Paraphrases (757MB, 11.9M rules)	Paraphrases (1.7MB, 69k rules)	Paraphrases (42MB, 1.2M rules)	Constituent (69MB, 1.0M rules)
		Identity (16MB, 468k rules)	Identity (170MB, 4.3M rules)	Non-Constituent (601MB, 9.6M rules)
L	Paraphrases (1.5GB, 23.5M rules)	Paraphrases (12MB, 198k rules)	Paraphrases (209MB, 3.0M rules)	Constituent (148MB, 2.2M rules)
		Identity (19MB, 503k rules)	Identity (191MB, 4.5M rules)	Non-Constituent (1.2GB, 18.2M rules)
XL	Paraphrases (2.8GB, 43.2M rules)	Paraphrases (33MB, 548k rules)	Paraphrases (486MB, 6.9M rules)	Constituent (300MB, 4.4M rules)
		Identity (20MB, 532k rules)	Identity (198MB, 4.7M rules)	Non-Constituent (2.1GB, 31.4M rules)
XXL	Paraphrases (5.7GB, 86.4M rules)	Paraphrases (125MB, 2.1M rules)	Paraphrases (1.5GB, 20.2M rules)	Constituent (644MB, 9.3M rules)
		Identity (21MB, 559k rules)	Identity (204MB, 4.8M rules)	Non-Constituent (3.6GB, 54.8M rules)
XXXL	Paraphrases (12.2GB, 169M rules)	Paraphrases (451MB, 7.6M rules)	Paraphrases (4.9GB, 68.4M rules)	Constituent (1.1GB, 16.1M rules)
		Identity (22MB, 570k rules)	Identity (207MB, 4.9M rules)	Non-Constituent (5.1GB, 77.4M rules)

Do the Scores Work?



Fun PPDB Examples

munchies ||| hungry

hustle ||| scam

sexiest ||| hottest

**P A R E N T A L
A D V I S O R Y
E X P L I C I T C O N T E N T**

losers

ly ||| indeed

Summary

Extraction & Representation

Extended large-scale paraphrase acquisition
from bitexts to syntactic paraphrases

Generation

Introduced a straightforward and effective
adaptation framework

Extensions beyond SMT

Improved performance by using monolingual
information

Current directions

Domain-specific paraphrasing

What if we want to generate paraphrases for specific domains like biology? Do they vary? How do we ensure ours are appropriate

Paraphrase recognition and entailment

The RTE problem diverges in interesting ways from paraphrasing. We are combining natural language inference and data-driven paraphrasing.

Divide

Parliament

gap

division

split

divided

gulf

dividing

share

divide up

divisions

separate

distinction

rift

difference

Biology

divided

division

dividing

divides

break

split

dispense

multiply

cleave

fracture

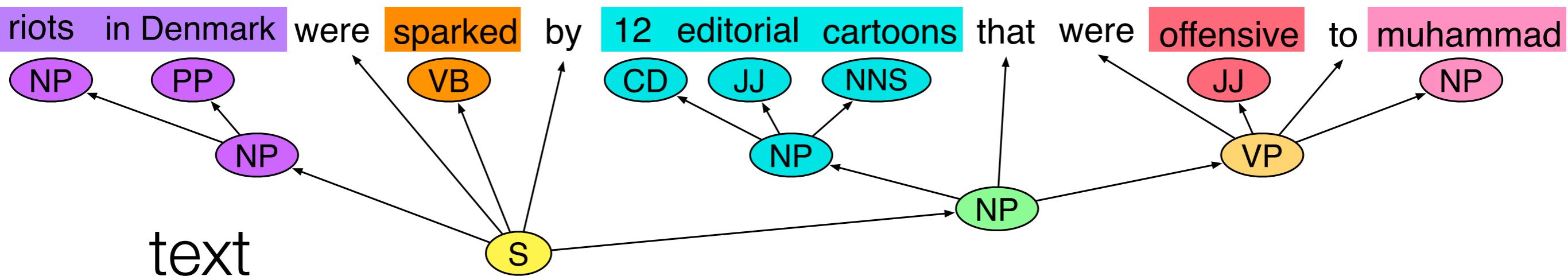
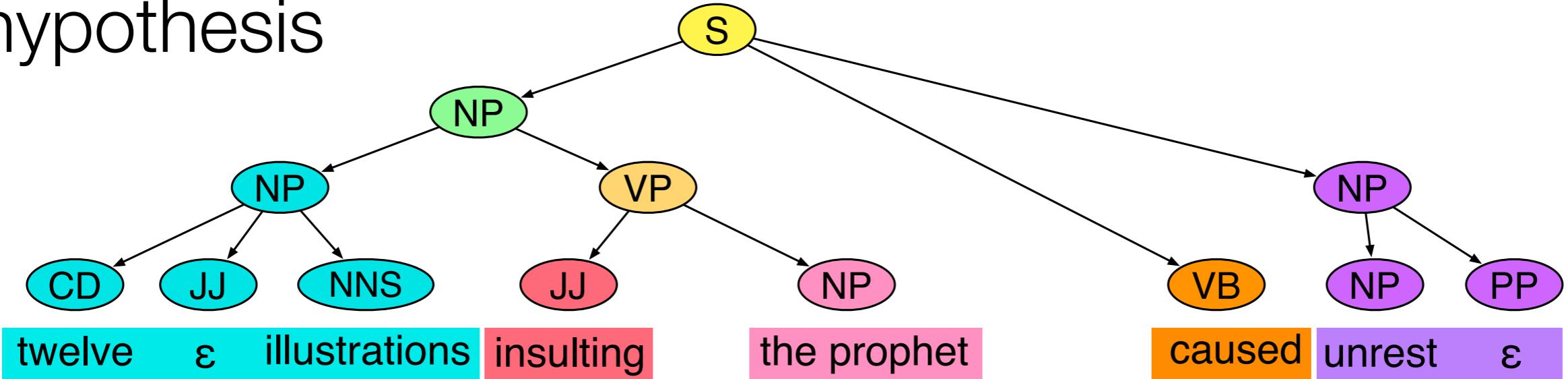
separate

mitotic division

partition

Textual Inference

hypothesis



Attaching a Semantics

twelve	12	equivalence
cartoons	illustrations	forward backward
Riots in Greece → Civil unrest in Europe Civil unrest in Europe → Riots in Greece		
caused	prevented	negation
Europe	the middle East	alternation

thank you for your time

many thanks

here you go anyway , thanks

leave a message

gee , thanks

thanks , man you look amazing

bless you

diet coke

Thank you!

thank you very much

keep the change thank you for your attention

uh , thanks

why , thank you

don't thank me

hey , thanks

thank you , frank

Bibliography

Paraphrasing with Bilingual Parallel Corpora. Colin Bannard and Chris Callison-Burch. ACL 2005.

Improved Statistical Machine Translation Using Paraphrases. Chris Callison-Burch, Philipp Koehn and Miles Osborne, 2006. In Proceedings NAACL-2006.

Paraphrase Substitution for Recognizing Textual Entailment. Wauter Bosma and Chris Callison-Burch. Lecture Notes in Computer Science, 2007.

Paraphrasing and Translation. Chris Callison-Burch, 2007. PhD Thesis, University of Edinburgh.

Syntactic Constraints on Paraphrases Extracted from Parallel Corpora. Chris Callison-Burch. EMNLP 2008.

Constructing Corpora for the Development and Evaluation of Paraphrase Systems. Trevor Cohn, Chris Callison-Burch, Mirella Lapata, 2008. Computational Linguistics: Volume 34, Number 4.

ParaMetric: An Automatic Evaluation Metric for Paraphrasing. Chris Callison-Burch, Trevor Cohn and Mirella Lapata. COLING 2008

Reranking Bilingually Extracted Paraphrases Using Monolingual Distributional Similarity. Charley Chan, Chris Callison-Burch, and Benjamin Van Durme. GEMS 2011.

Learning Sentential Paraphrases from Bilingual Parallel Corpora for Text-to-Text Generation. Juri Ganitkevitch, Chris Callison-Burch, Courtney Napoles, and Benjamin Van Durme. EMNLP 2011.

Monolingual Distributional Similarity for Text-to-Text Generation. Juri Ganitkevitch, Ben Van Durme and Chris Callison-Burch. StarSEM 2012.

PPDB: The Paraphrase Database Juri Ganitkevitch, Ben Van Durme and Chris Callison-Burch. NAACL 2013.

Entailment relations

Hypernym	Synonym	Antonyms	Alternations	Independent
beetle insect	icebox refrigerator	advantage disadvantage	cheese butter	advocacy spokesman
honeybee bee	impasse deadlock	competence incompetence	cliff cave	aircraft sky
fees spending	infirmary hospital	continuity discontinuity	clothing equipment	actor arena
know-how knowledge	insurrection revolt	inflow outflow	clothing housing	actor maker
pond lake	jewel gem	insanity sanity	coating asphalt	actor movie
fertilizer manure	john lavatory	legitimacy illegitimacy	columnist newspaperman	actor singer
actor entertainer	kale cabbage	niece nephew	commentator reporter	actor spokesman
actor performer	labyrinth maze	descendants ancestors	competence productivity	advantage equipment
acquisition buying	laundry washing	husbands wives	compliance enforcement	ambassador delegation