



JOHNS HOPKINS
UNIVERSITY
CENTER FOR LANGUAGE
& SPEECH PROCESSING



Penn
UNIVERSITY OF PENNSYLVANIA



human language technology
center of excellence

Large-Scale Paraphrase Extraction & Application

Juri Ganitkevitch

Jonathan Weese, Courtney Napoles, Benjamin Van Durme,
Matt Post, and Chris Callison-Burch

#tbt

Paraphrase Extraction

Text-to-Text Applications

Scaling to PPDB

Contributions

Paraphrases

Differing textual expressions of the same meaning:

cup \leftrightarrow mug

the king's speech \leftrightarrow His Majesty's address

X₁ talks to X₂ \leftrightarrow X₁ converses with X₂

one JJ instance of NP \leftrightarrow a JJ case of NP

..in NLP

Recognition or generation of meaning-constrained relationships in **text**, i.e. automatic...

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

Where do paraphrases
come from?

Data.

Types of Corpora

Parallel: English – English

Comparable: English ~ English

Plain: 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**.

Types of Corpora

Parallel: English – English

Comparable: English ~ English

Plain: English

Bilingual parallel: English – French

A staggering 5 million Americans **have been victims of identity theft** in the last five years, according to a federal trade commission survey out this week.

In the last year alone, 1 million people **have had their identity purloined**.

Types of Corpora

Parallel: English – English

Comparable: English ~ English

Plain: English

Bilingual parallel: English – French

duty | responsibility

Modified by:

additional, administrative,
assigned, assumed,
collective, congressional,
constitutional

Object of:

accept, articulate, assert,
assign, assume, attend
to, avoid, become,
breach

Types of Corpora

Parallel: English – English

Comparable: English ~ English

Plain: English

Bilingual parallel: English – French

Bilingual Parallel Corpora

Sentence-aligned corpora in English
and any foreign language

Strong meaning equivalence signal

Available in large quantities

Pivoting over a Foreign Language

... their

... ihre

plans in the long term

langfristigen Pläne

would ...

würden ...

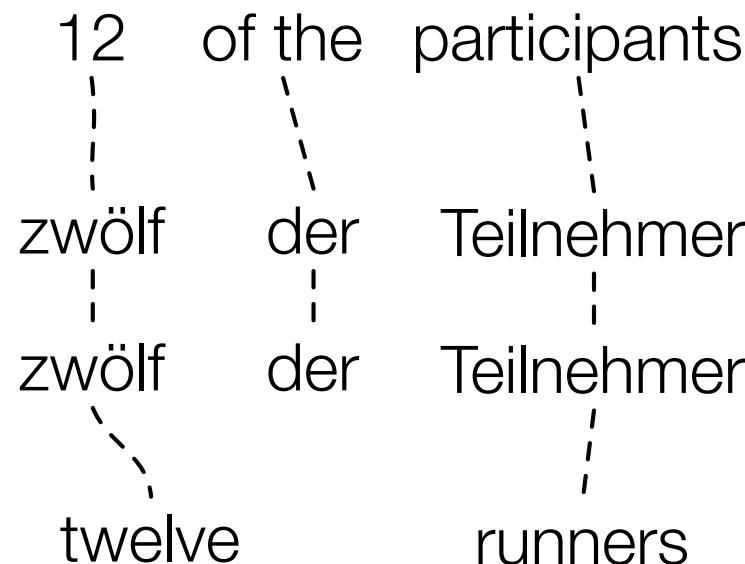
Goal: Expressiveness

Generalization vs. memorization

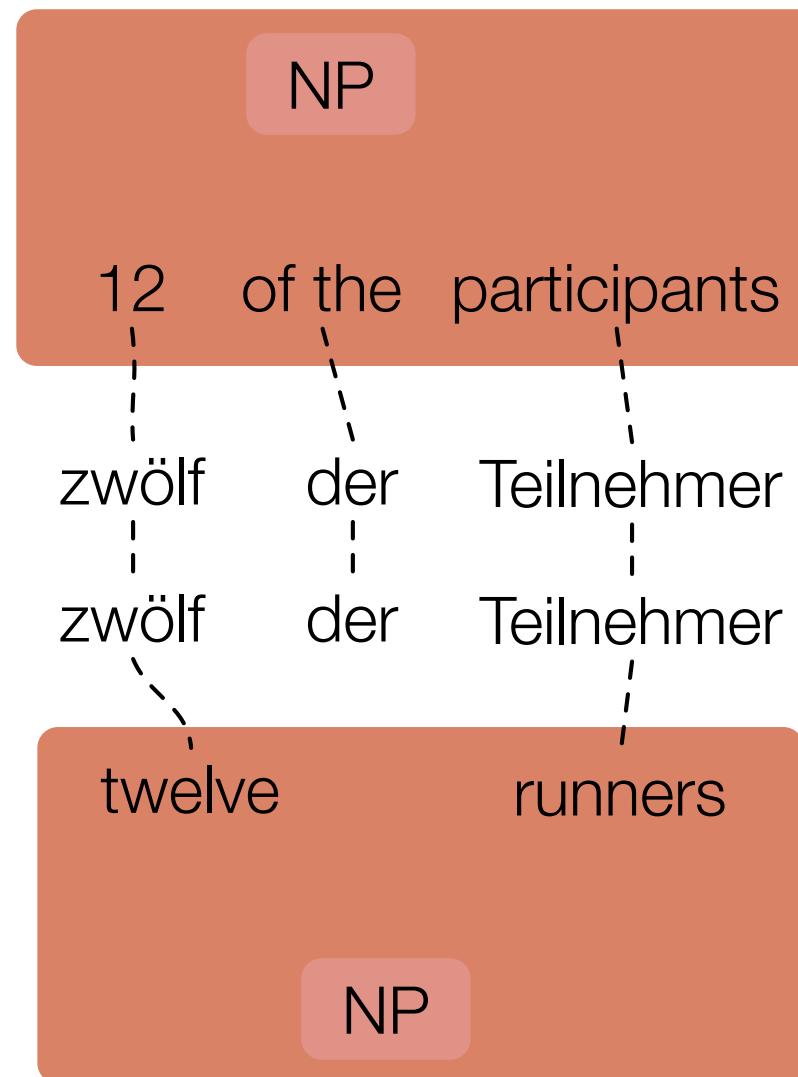
Prefer general patterns over memorized phrases

Use linguistic information

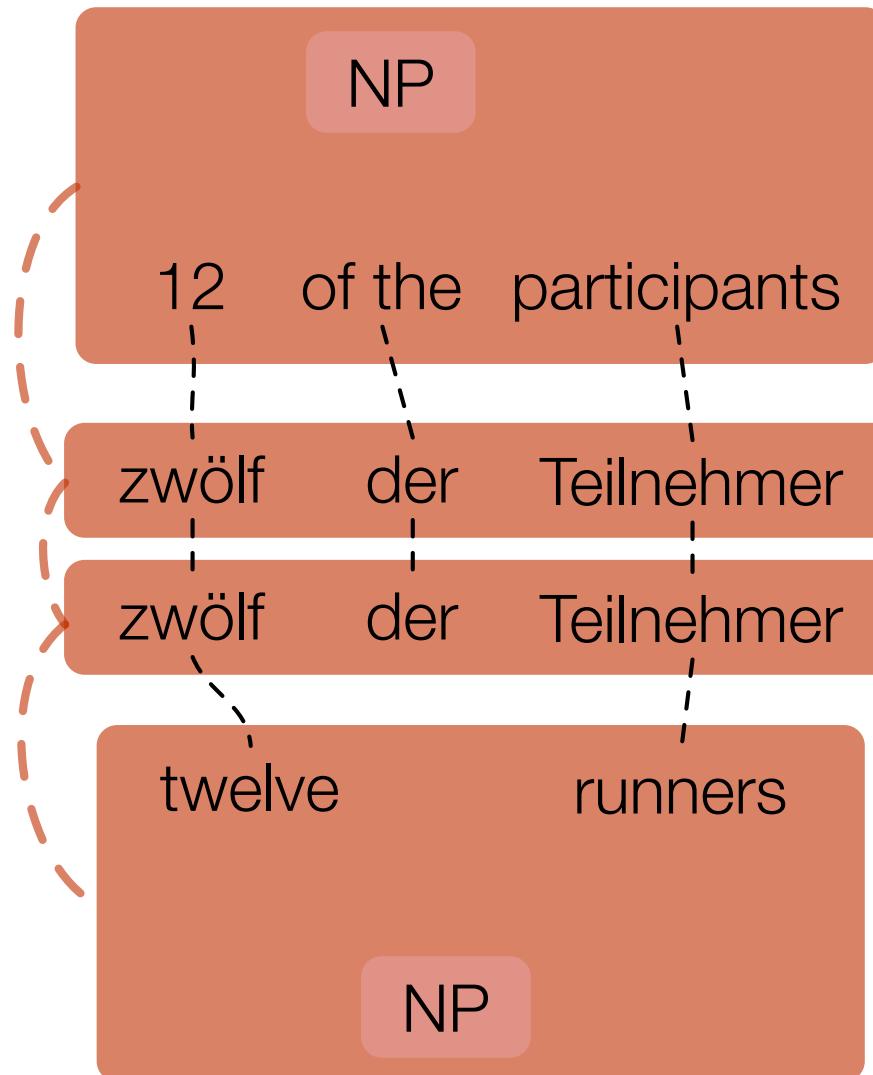
Syntactic Pivoting



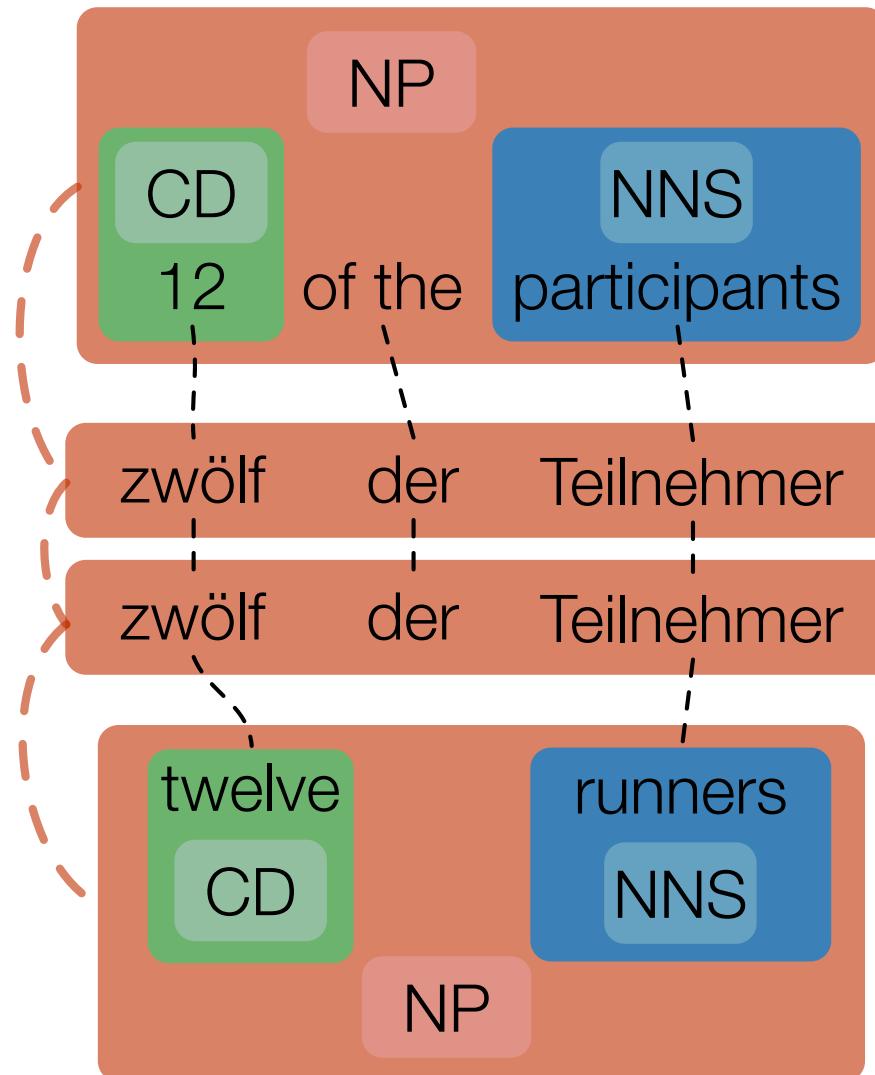
Syntactic Pivoting



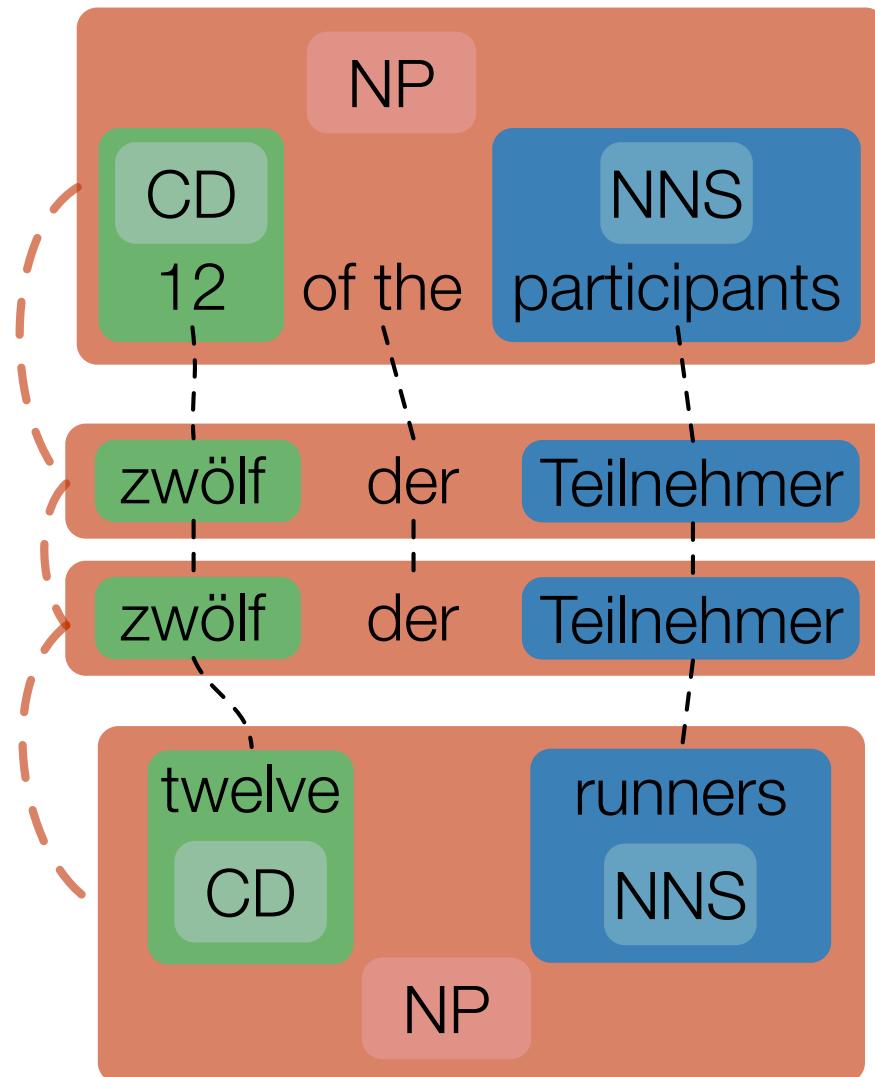
Syntactic Pivoting



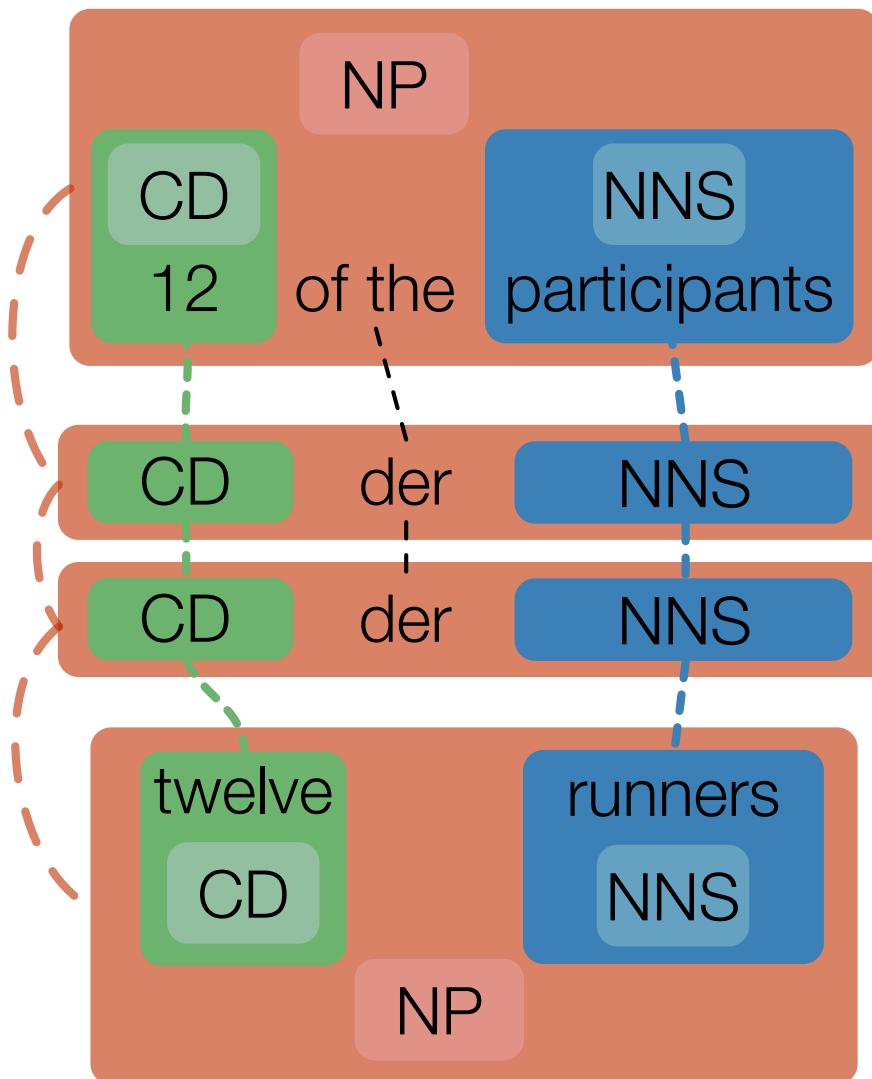
Syntactic Pivoting



Syntactic Pivoting

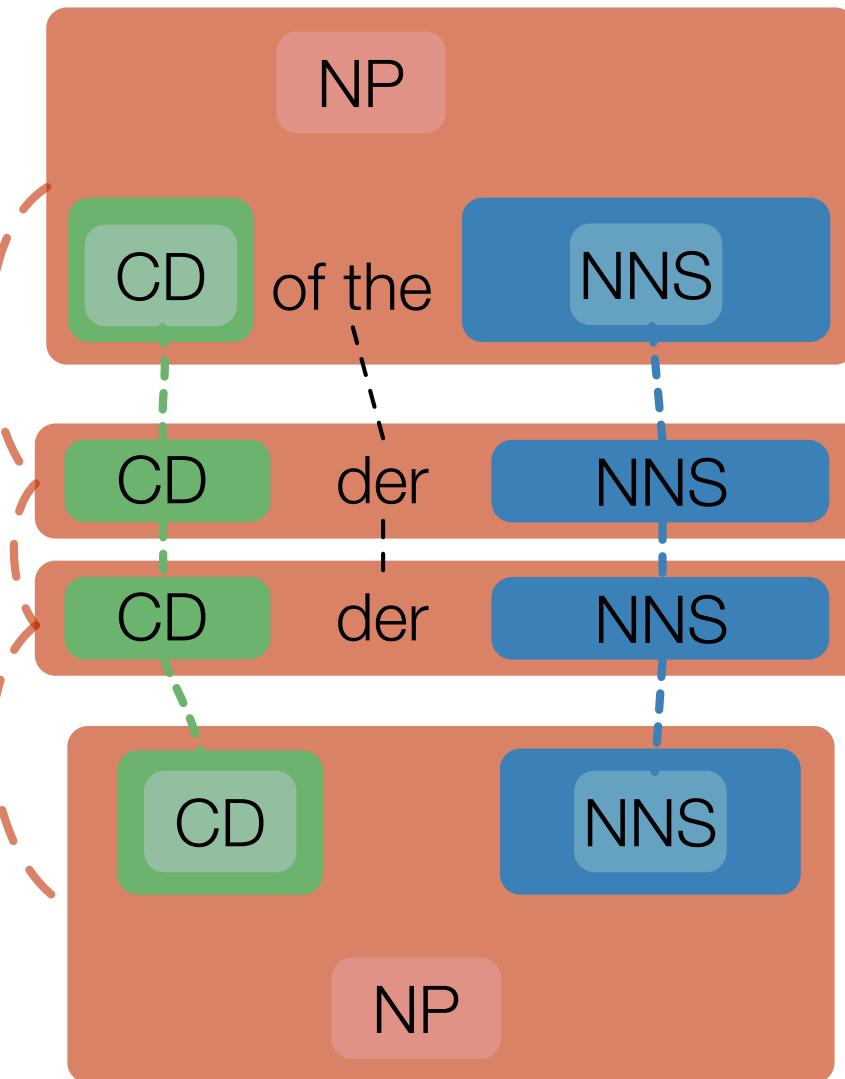


Syntactic Pivoting

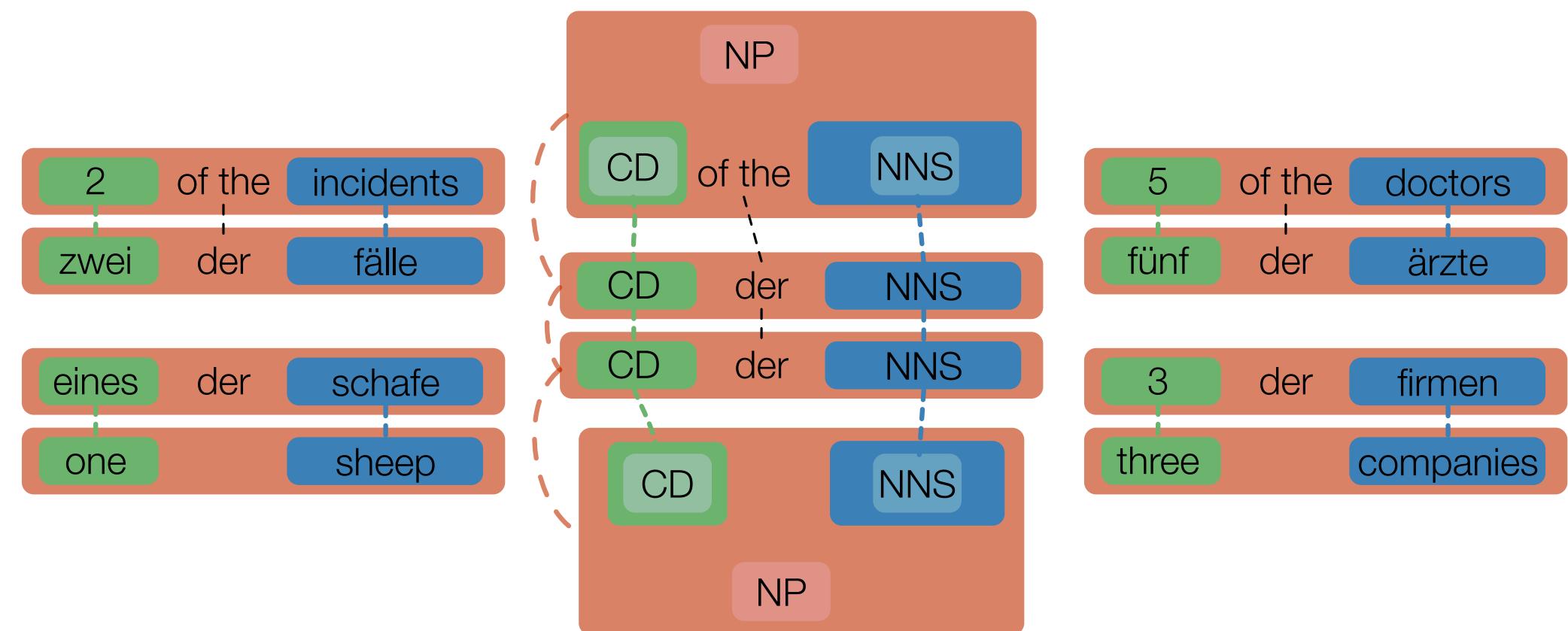


Ganitkevitch et al., 2011

Syntactic Pivoting



Syntactic Pivoting



Possessive rule	NP → the NN of the NNP the NNP's NN NP → the NNS ₁ made by the NNS ₂ the NNS ₂ 's NNS ₁
Dative shift	VP → give NN to NP give NP the NN VP → provide NP ₁ to NP ₂ give NP ₂ NP ₁
Adv./adj. phrase move	S/VP → ADVP they VBP they VBP ADVP S → it is ADJP VP VP is ADJP
Verb particle shift	VP → VB NP up VB up NP
Reduced relative clause	SBAR/S → although PRP VBP that although PRP VBP ADJP → very JJ that S JJ S
Partitive constructions	NP → CD of the NN CD NN NP → 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 → take action ADVP to act ADVP VP → to take a decision PP to decide PP

Paraphrase Quality

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

Pivoting Probability

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

Paraphrase Extraction

Text-to-Text Applications

Scaling to PPDB

Contributions

Text-to-Text



Paraphrases form an SCFG

“Monolingual translation” via Joshua

Sentence compression, text
simplification, poetry generation,
lawyer-to-English, etc.

Sentence Compression

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

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

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

Adaptation Scheme

Tuning Data

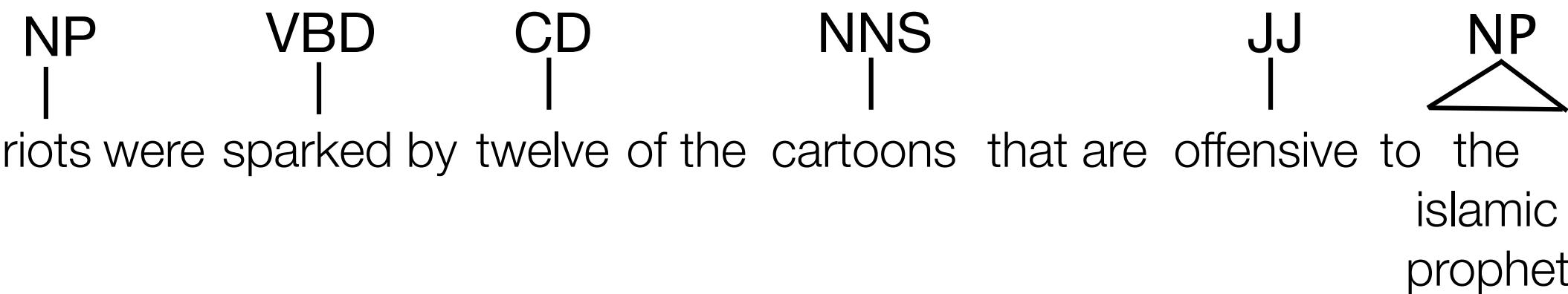
~2k long/short sentence pairs

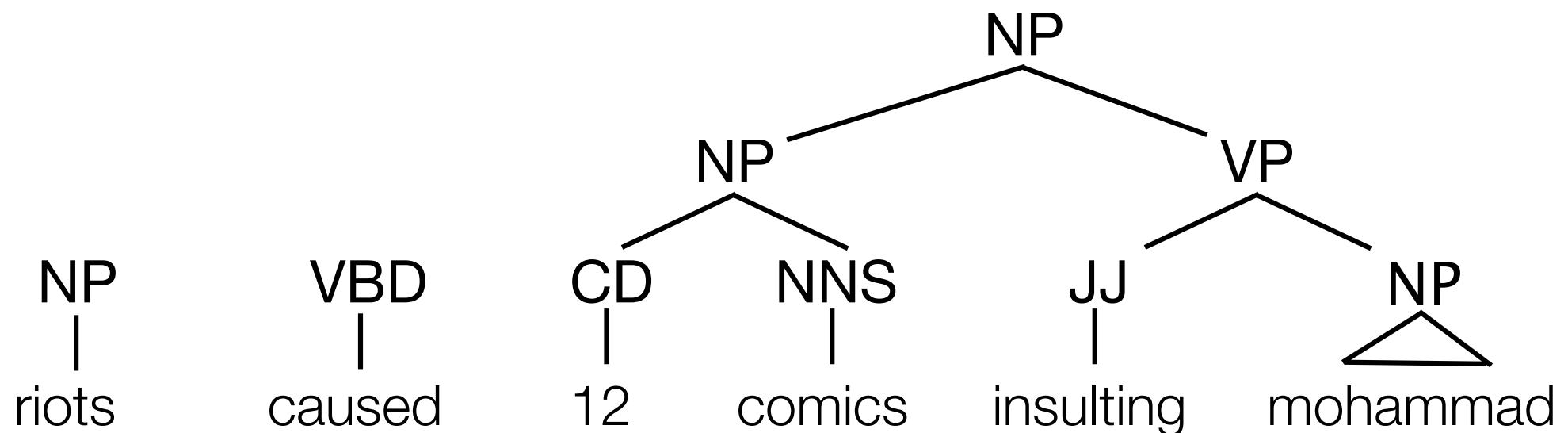
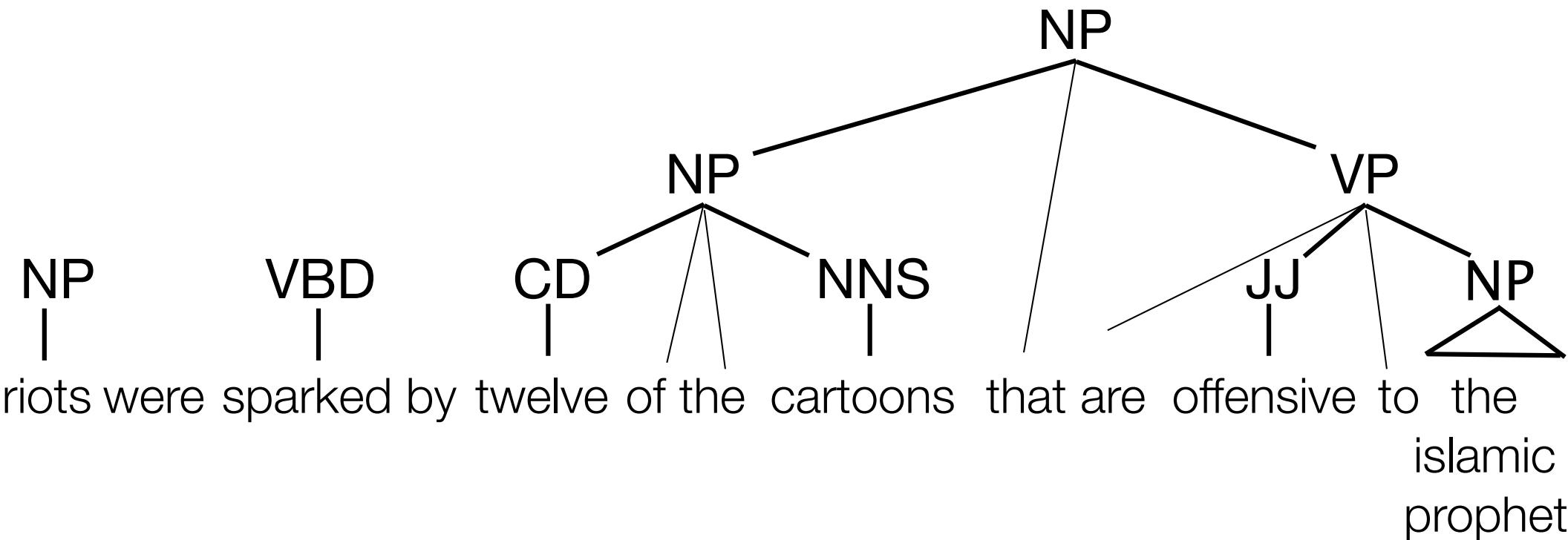
Objective Function

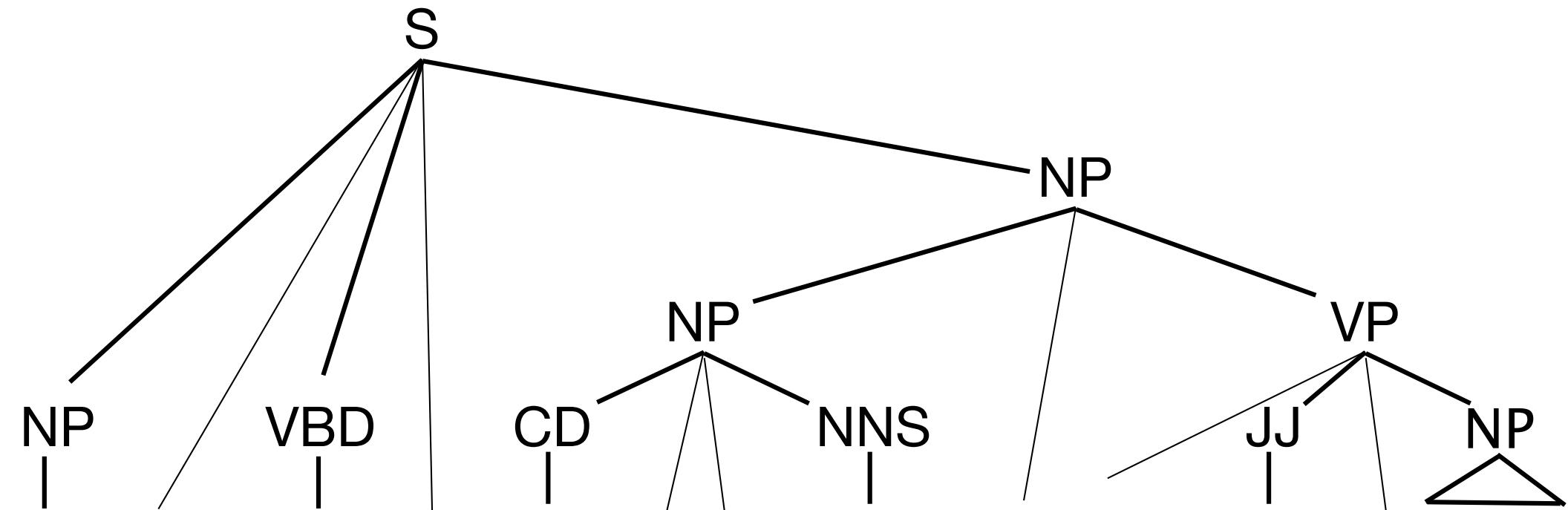
Compression-ratio targeted BLEU

Paraphrase Features

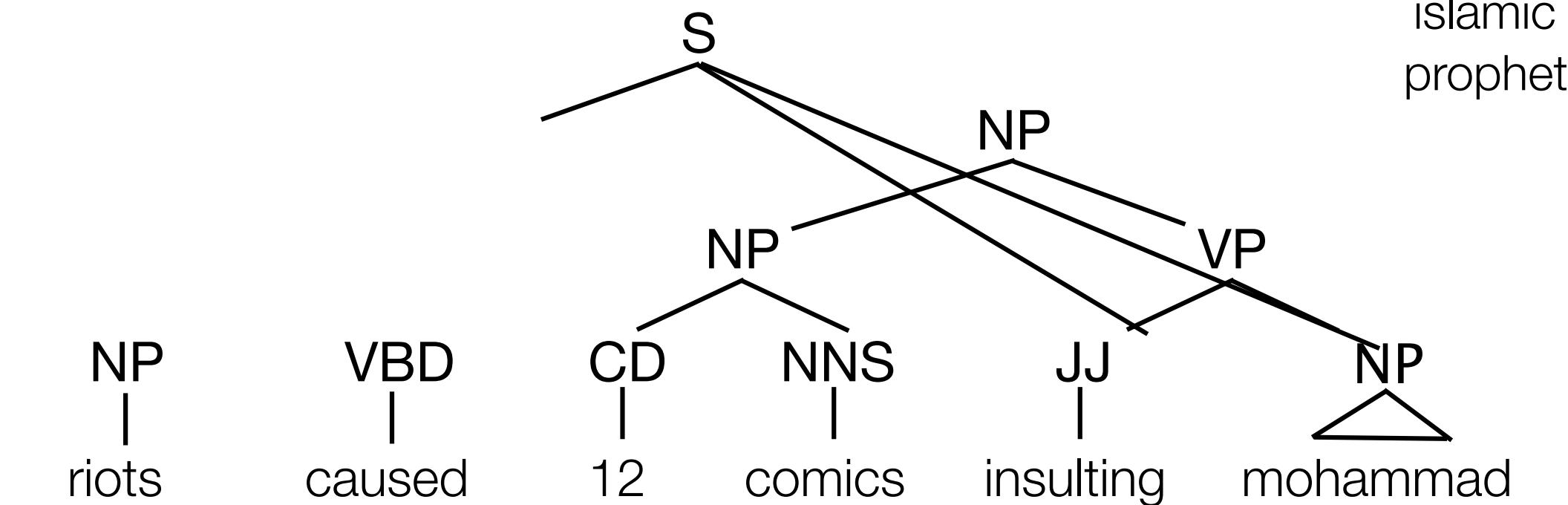
Compression-ratio of paraphrase etc.





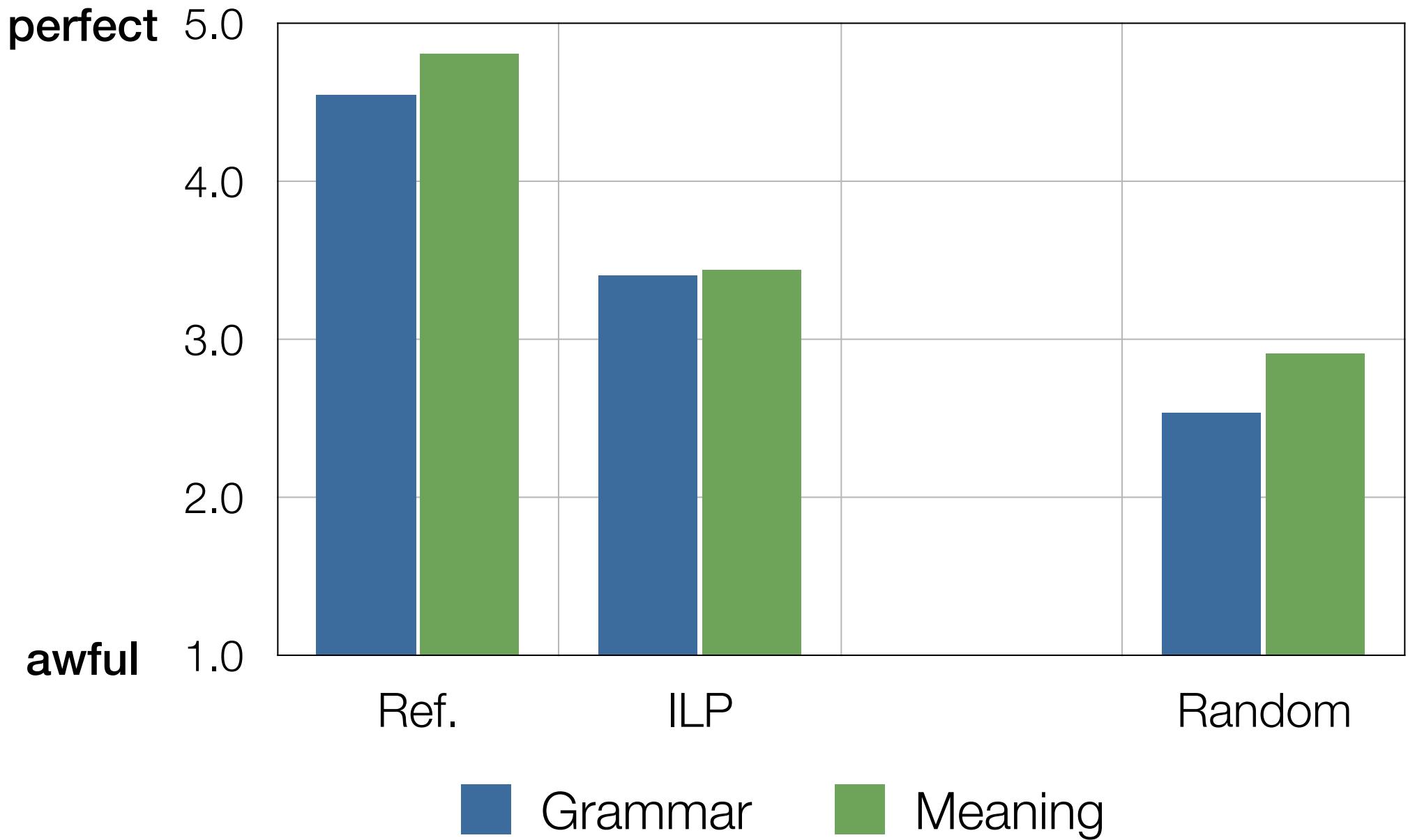


riots were sparked by twelve of the cartoons that are offensive to the
islamic prophet

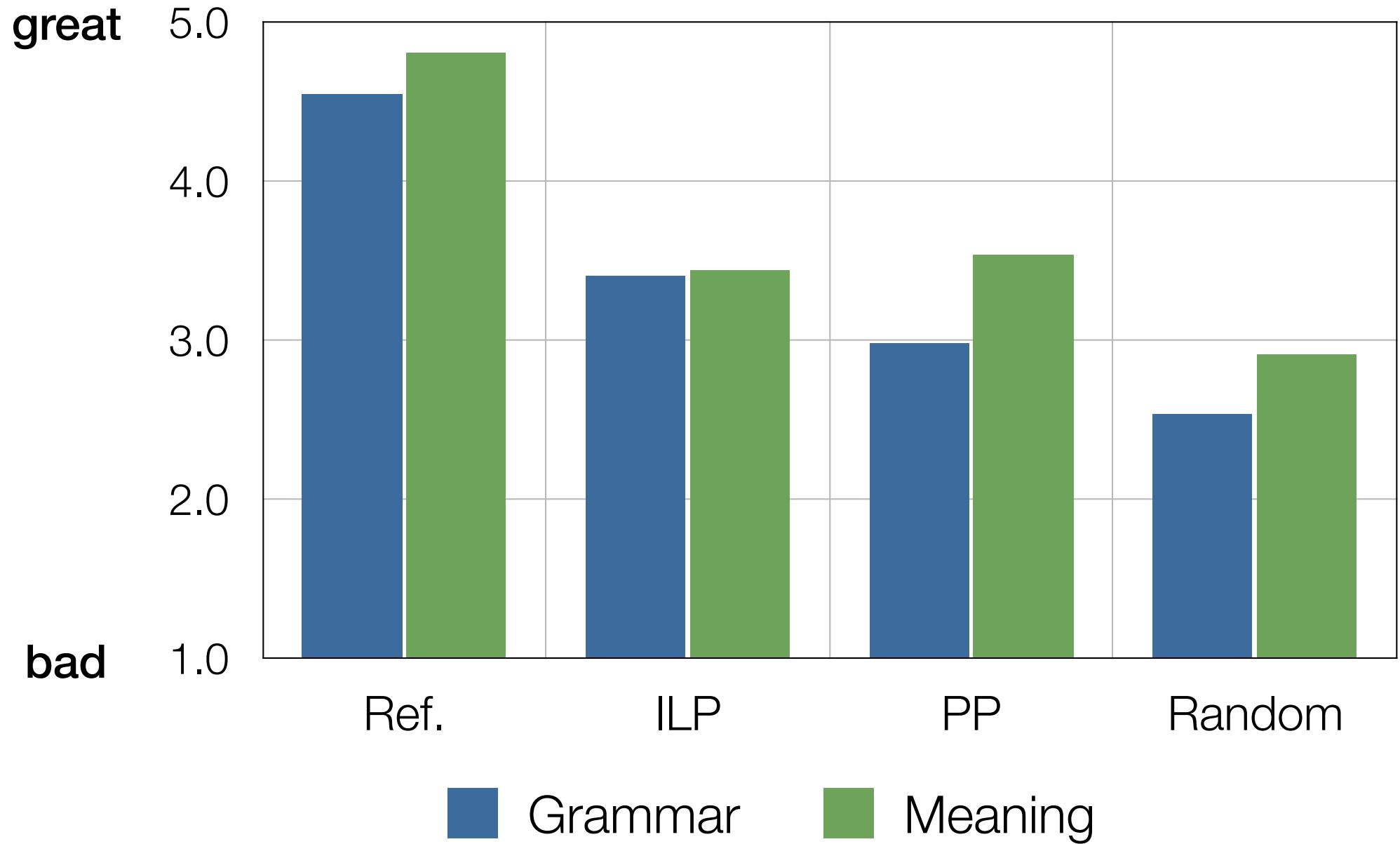


riots caused 12 comics insulting mohammad

Compression Results



Compression Results



Paraphrase Extraction

Text-to-Text Applications

Scaling to PPDB

Contributions

Going Beyond Translation

Monolingual data & distributional similarity:

Orthogonal signal

Rich contextual features possible

Vast amounts of data available

duty | responsibility

Modified by:

additional, administrative,
assigned, assumed,
collective, congressional,
constitutional

Object of:

accept, articulate, assert,
assign, assume, attend
to, avoid, become,
breach

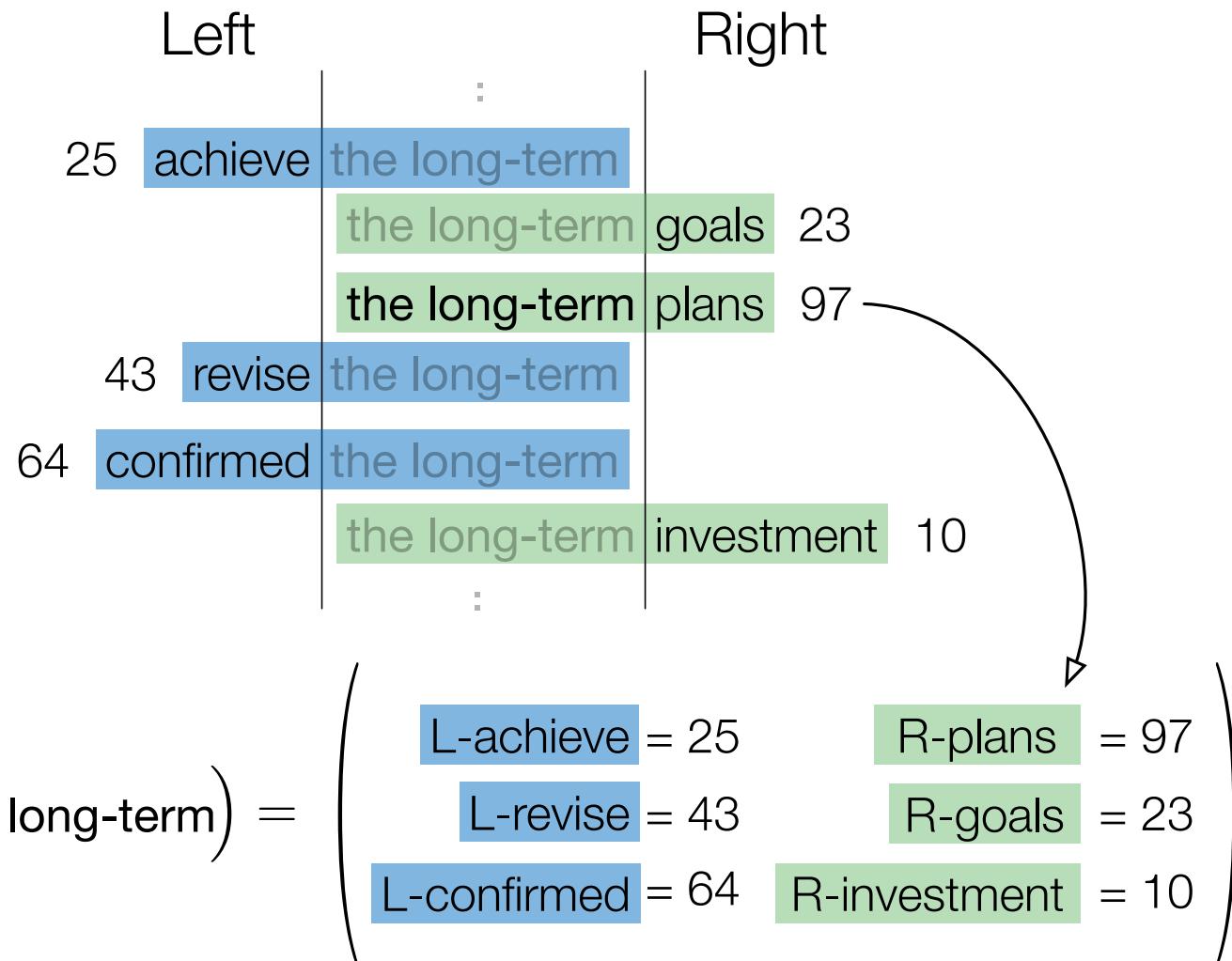
Monolingual Data

Google n -grams

Annotated Gigaword

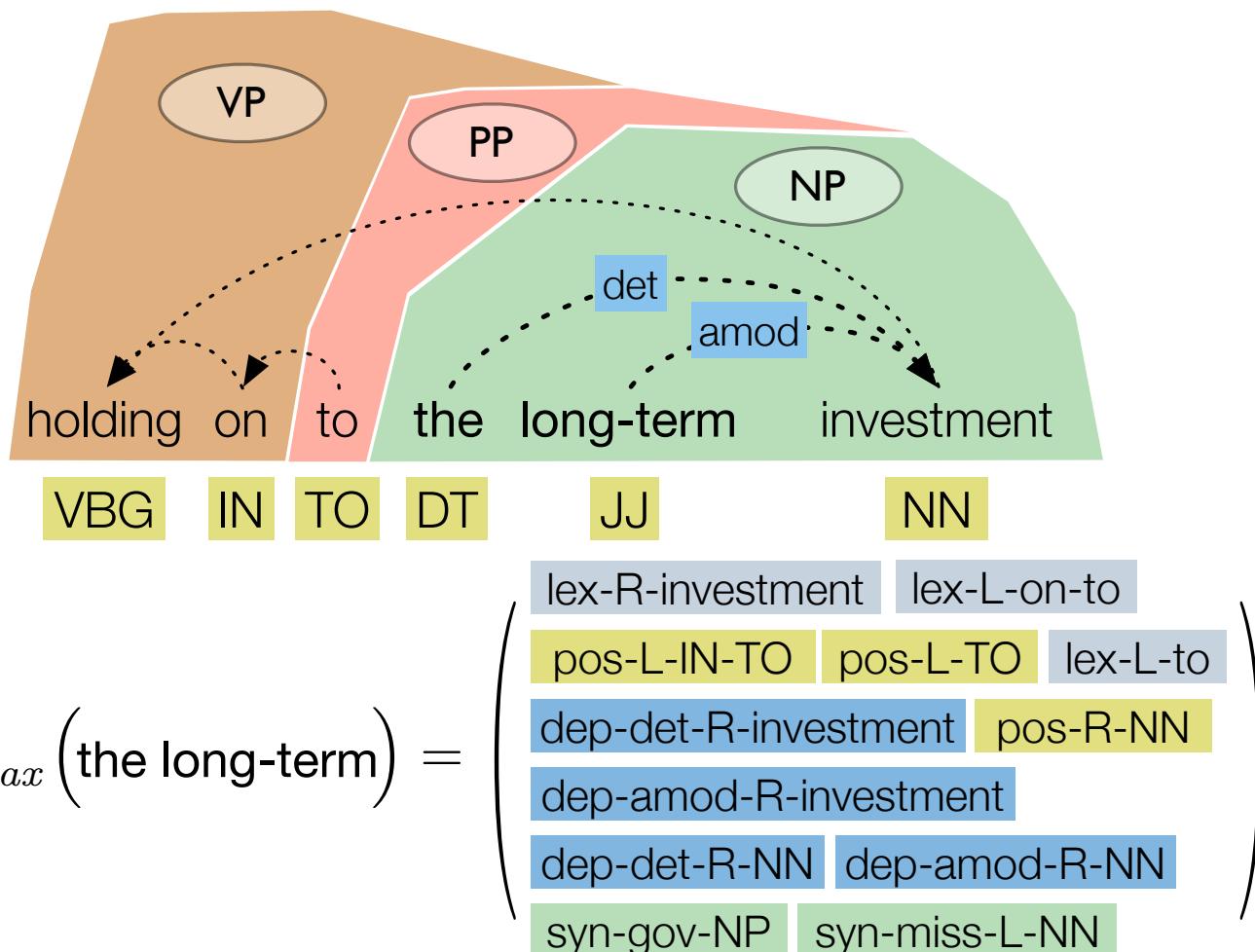
Ganitkevitch et al., 2012

n -gram Context

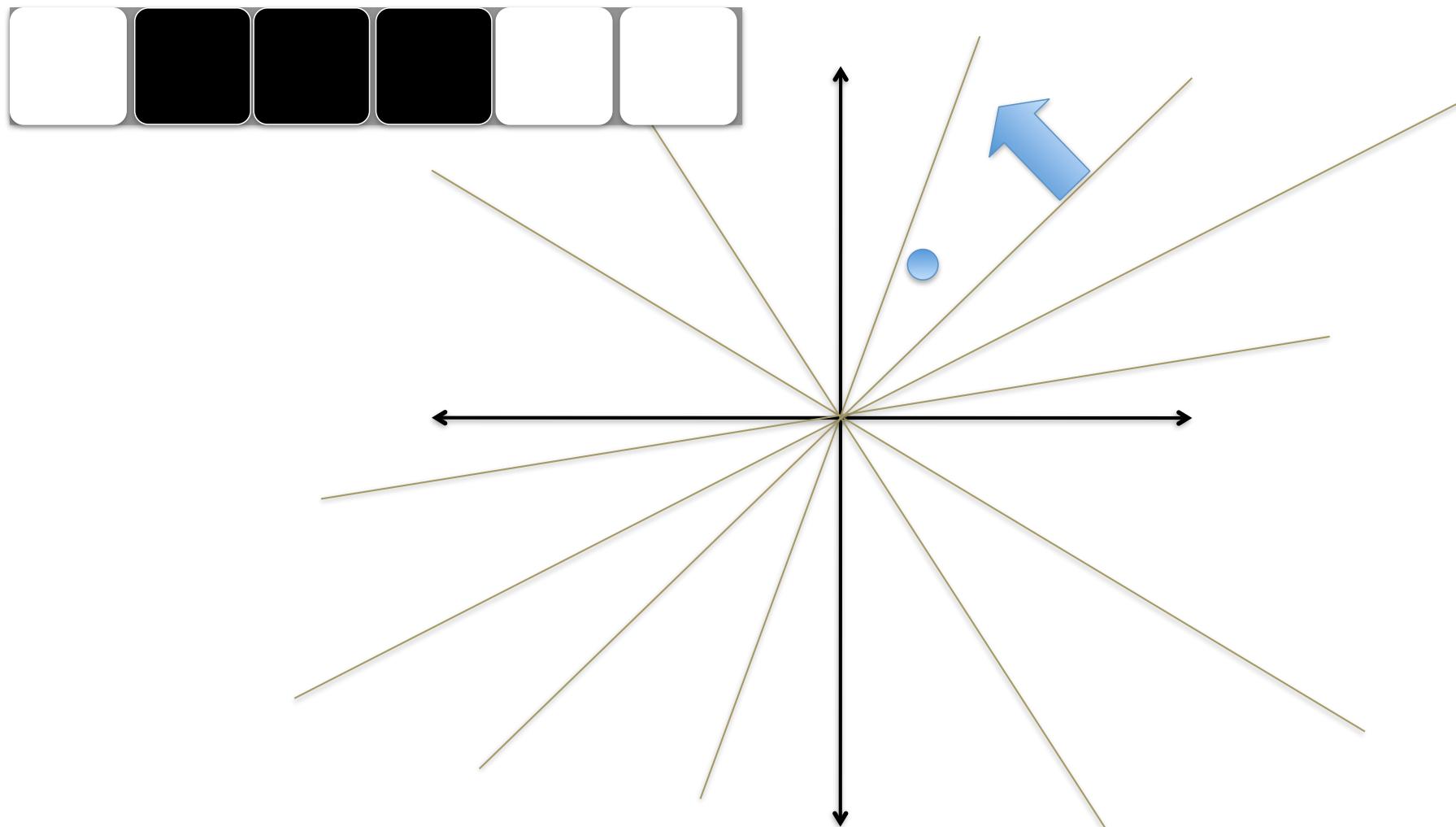


Ganitkevitch et al., 2012

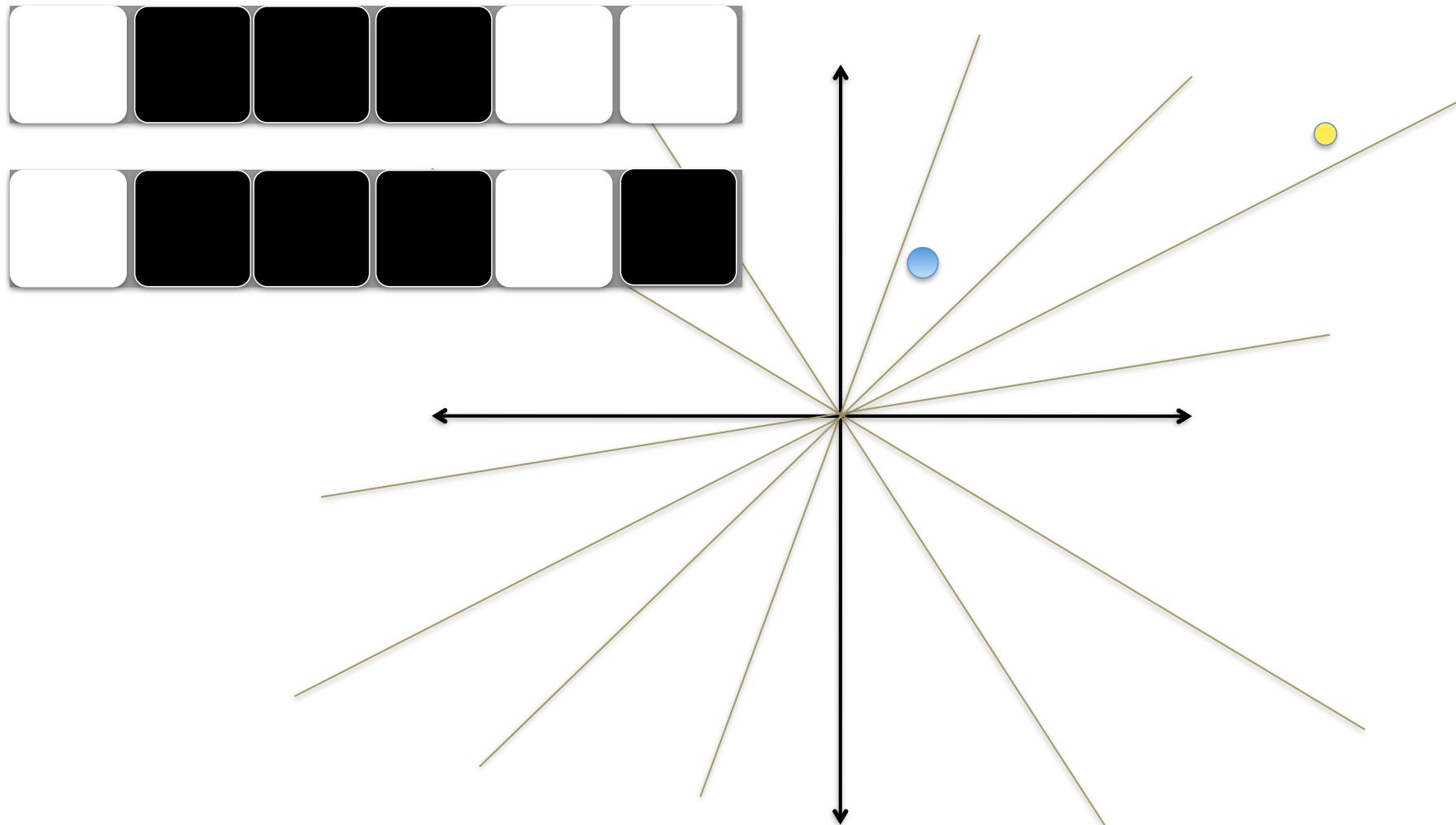
Annotated Context



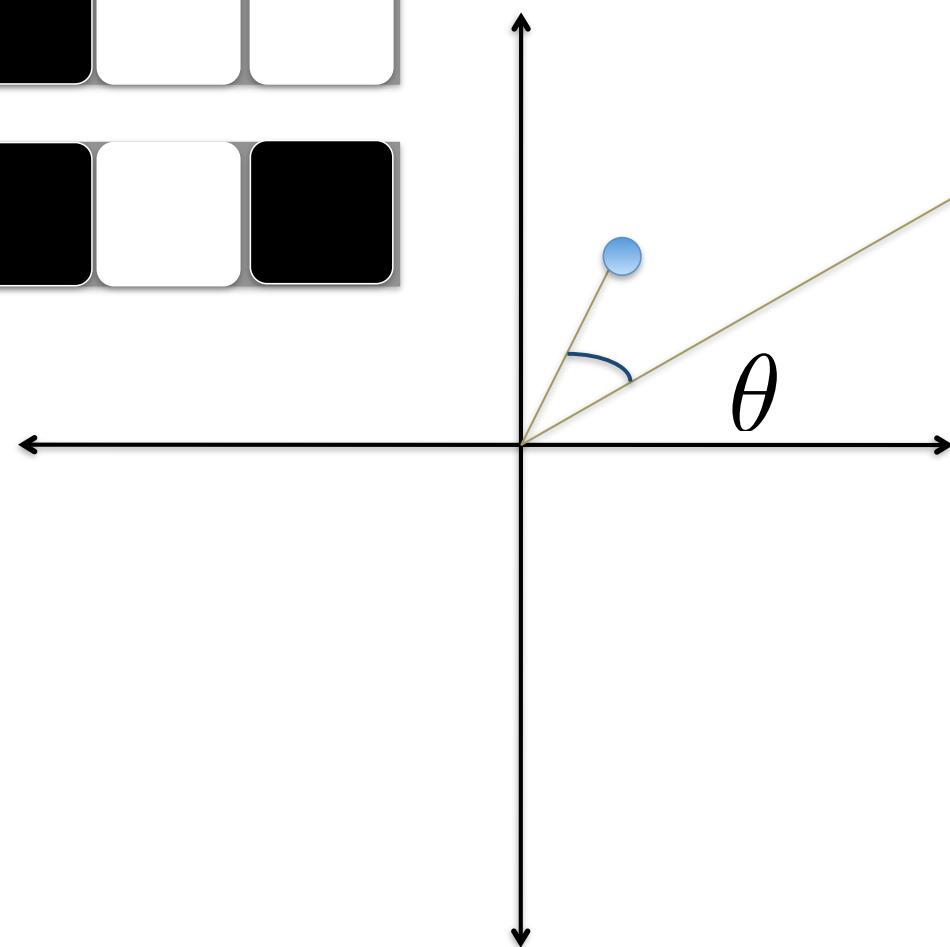
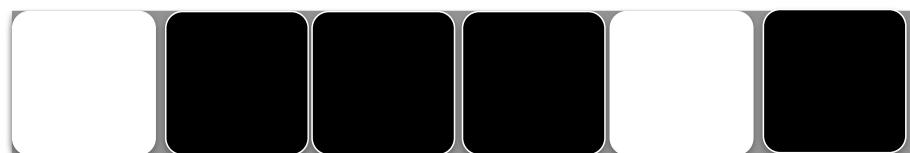
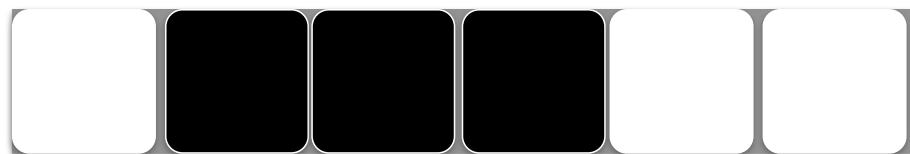
Locality Sensitive Hashing



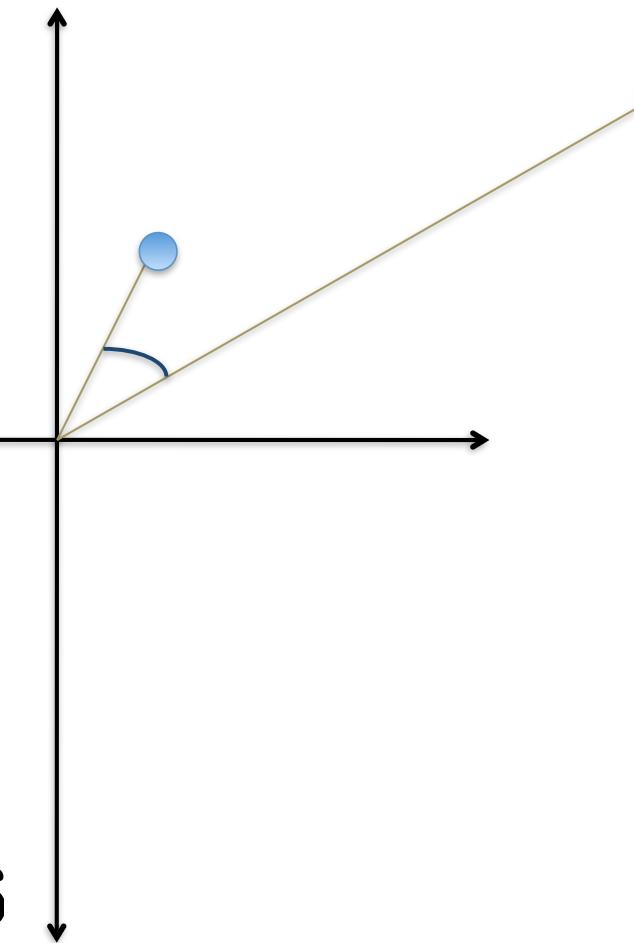
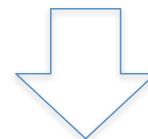
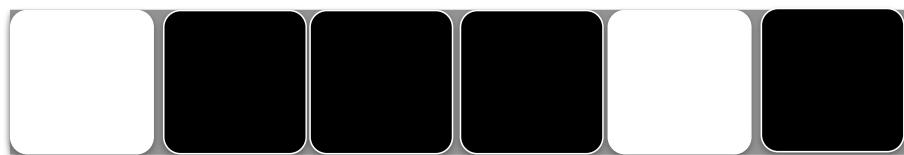
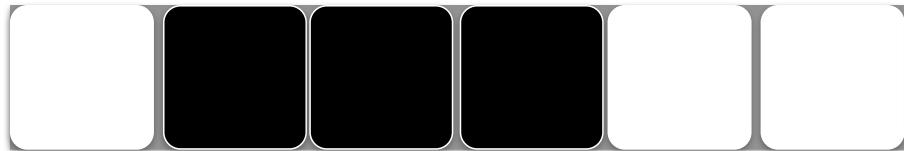
Locality Sensitive Hashing



Locality Sensitive Hashing



Locality Sensitive Hashing



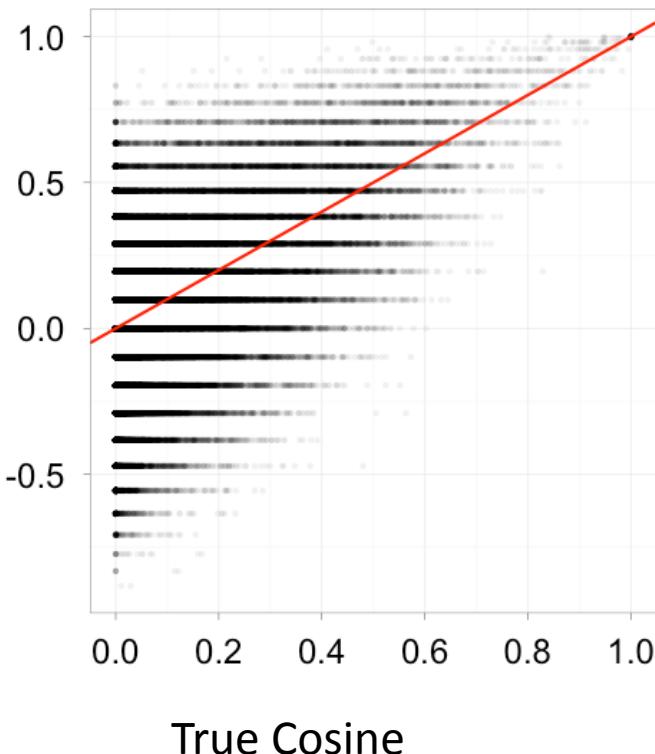
Hamming Distance := $h = 1$

Signature Length := $b = 6$

Accuracy as function of bit length

32 bit signatures

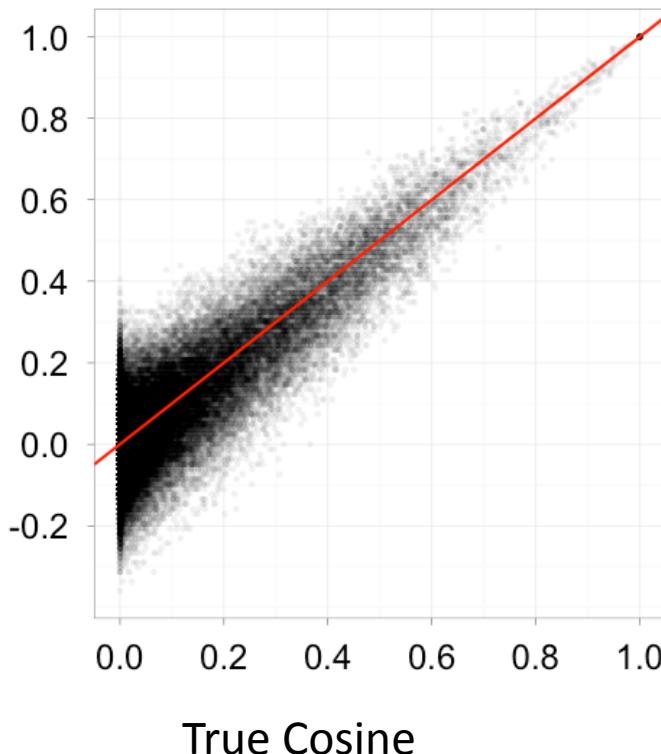
Approximate Cosine



True Cosine

256 bit signatures

Approximate Cosine



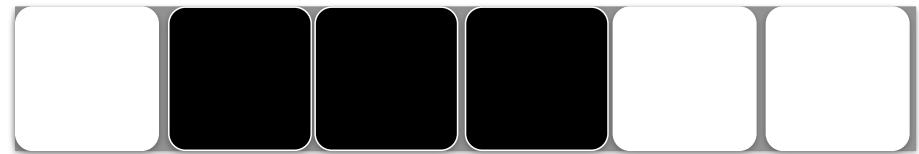
True Cosine

Cheap

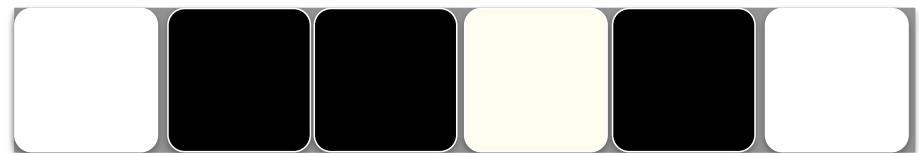
Accurate

Signature Collection

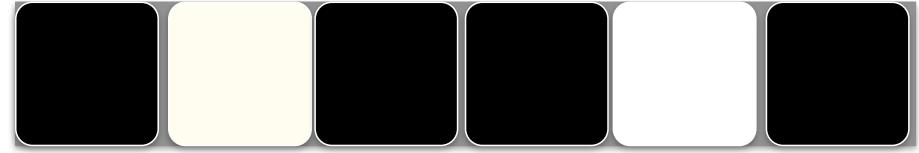
of



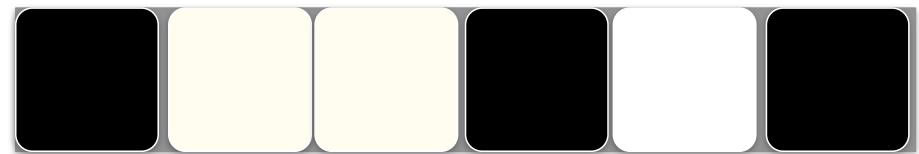
's



the long term



in the long term



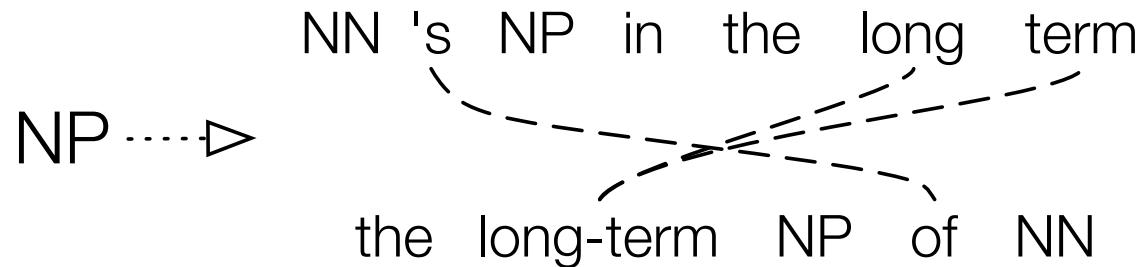
Syntactic Paraphrase Similarity

NN 's NP in the long term

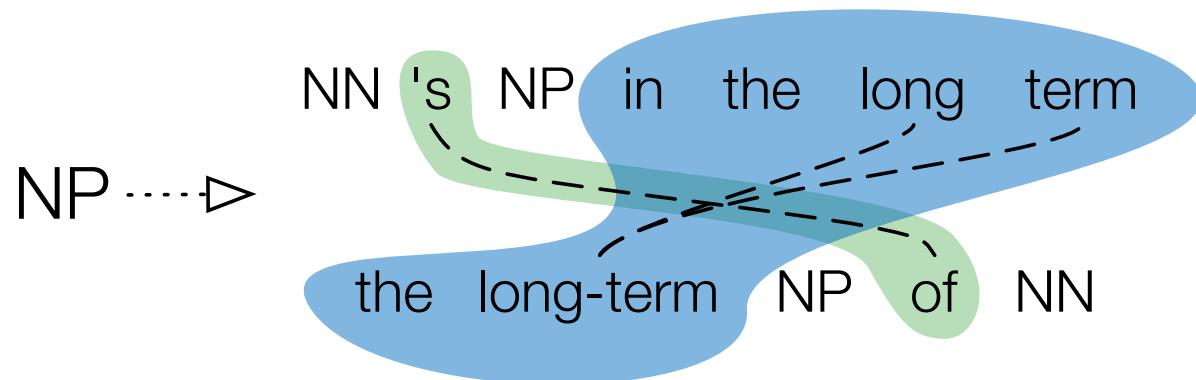
NP ---->

the long-term NP of NN

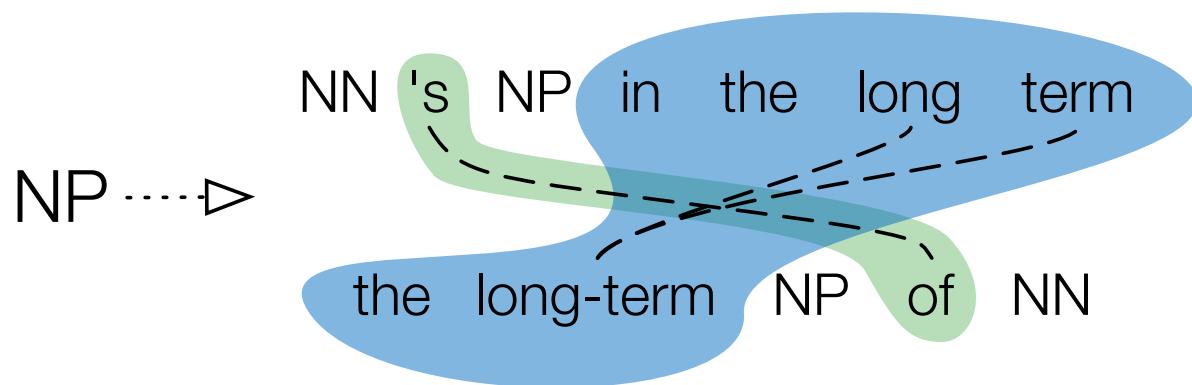
Syntactic Paraphrase Similarity



Syntactic Paraphrase Similarity

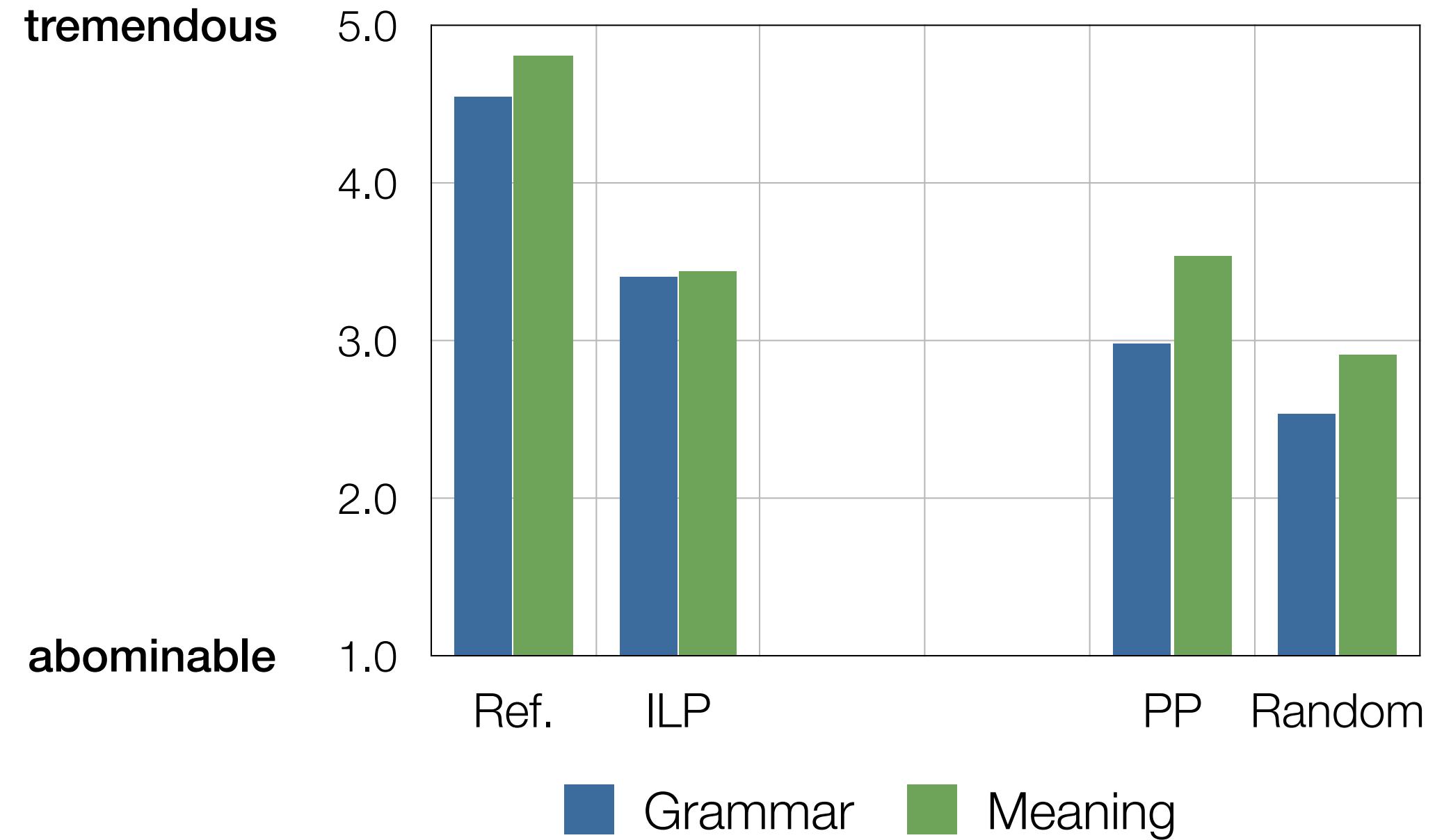


Syntactic Paraphrase Similarity



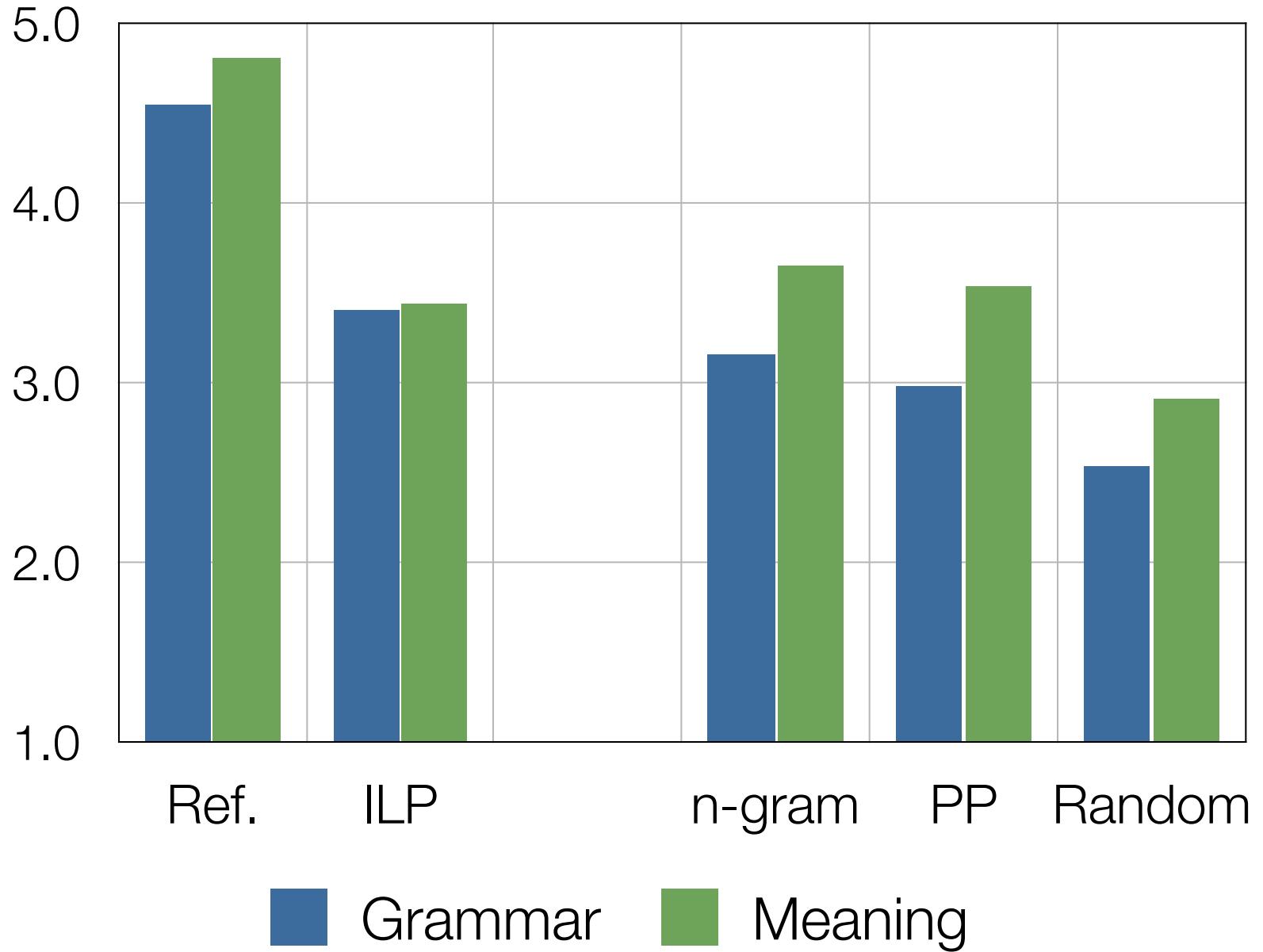
$$sim(\mathbf{r}) = \frac{1}{2} \left(sim\left(\begin{matrix} \text{the long-term} \\ \text{in the long term} \end{matrix} \right) + sim\left(\begin{matrix} \text{'s} \\ \text{of} \end{matrix} \right) \right)$$

Compression Results



Compression Results

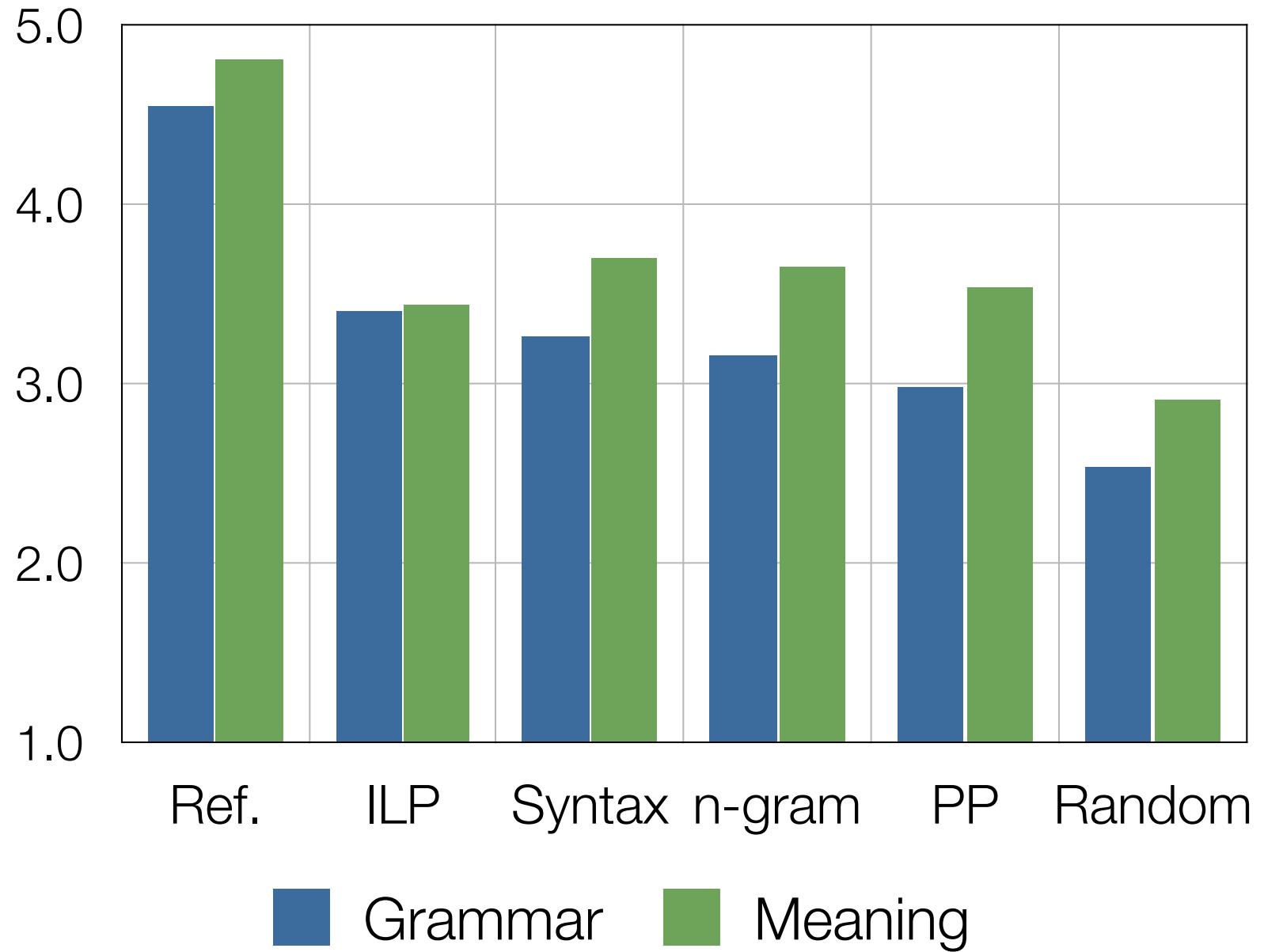
excellent



horrible

Compression Results

amazing



pathetic

Paraphrase Extraction

Text-to-Text Applications

Scaling to PPDB

Contributions

Bilingual Data

Europarl, Fr-En Gigaword, News Commentary,
UN, JRC Acquis, GALE, NIST Urdu, OpenST,
TAUS

106 million sentence pairs

2 billion English words

22 pivot languages

Monolingual Data

Google n -grams

Collection of 1 trillion tokens with counts

Based on vast amounts of text

Annotated Gigaword (Napoles et al., 2012)

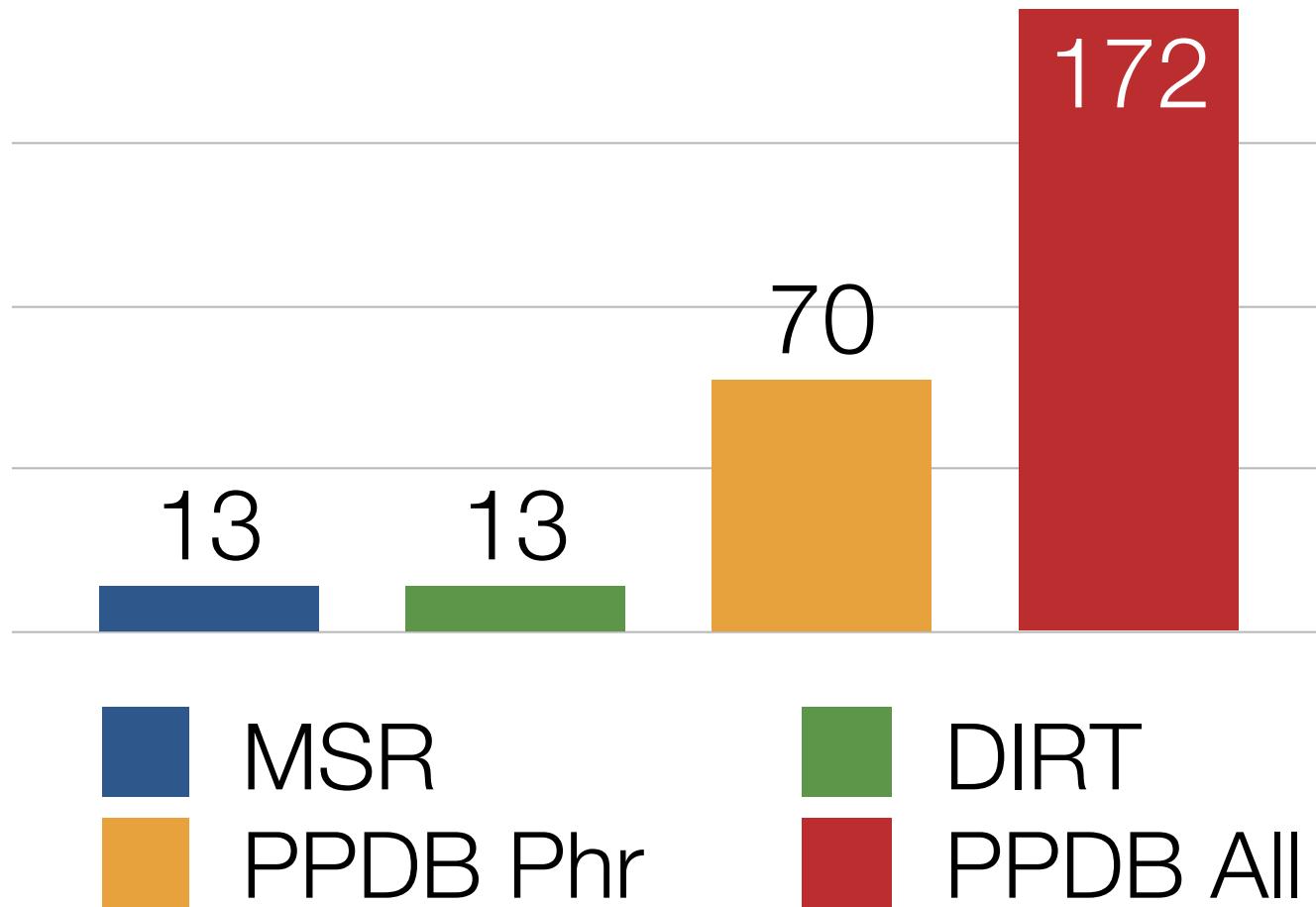
Collection of 4 billion words, parsed and tagged

Ganitkevitch et al., 2013

PPDB:Eng

	Paraphrases
Lexical	7.6 M
One-Many	2.8 M
Phrasal	68.4 M
Syntactic	93.6 M
Total	172.4 M

Other Resources



Ganitkevitch et al., 2013



Machinery

Joshua & Thrax

SCFG-based decoder for translation and paraphrasing

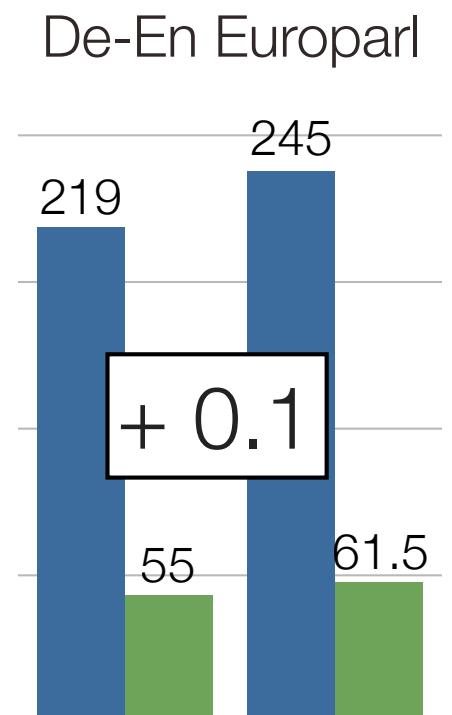
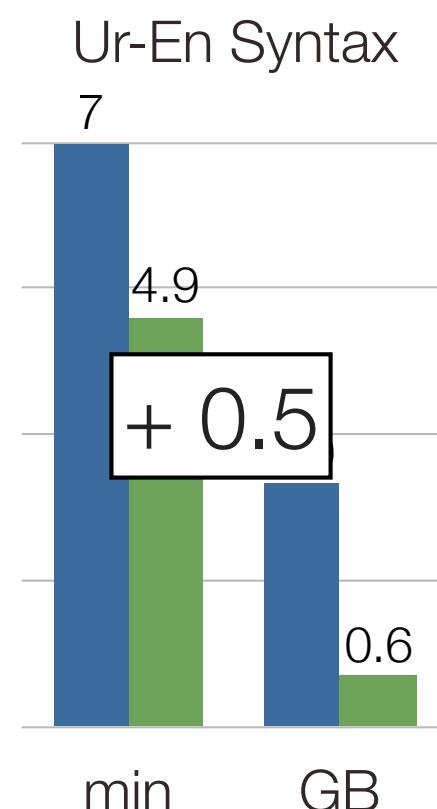
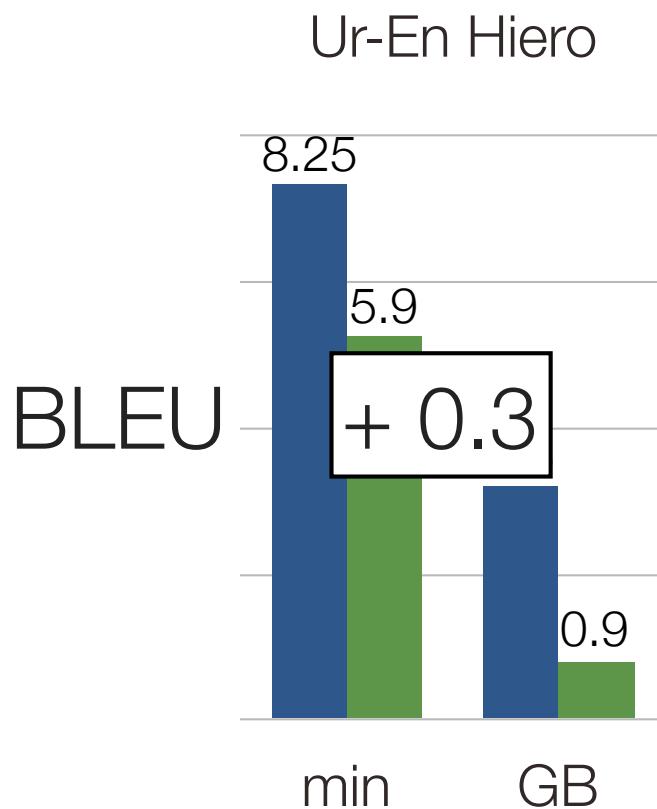
Includes: parallelized extraction of paraphrase grammars & distributional signatures

Jerboa

Toolkit for randomized and streaming algorithms (LSH, Bloom filters, etc.)

Thrax 2.0

1.0 2.0



75% slimmer, up to 300% faster

Ganitkevitch et al., 2012

Extraction on Amazon

PPDB:Spa

100 machines, 4 cores, 15GB RAM

~ 8 hours on spot instances

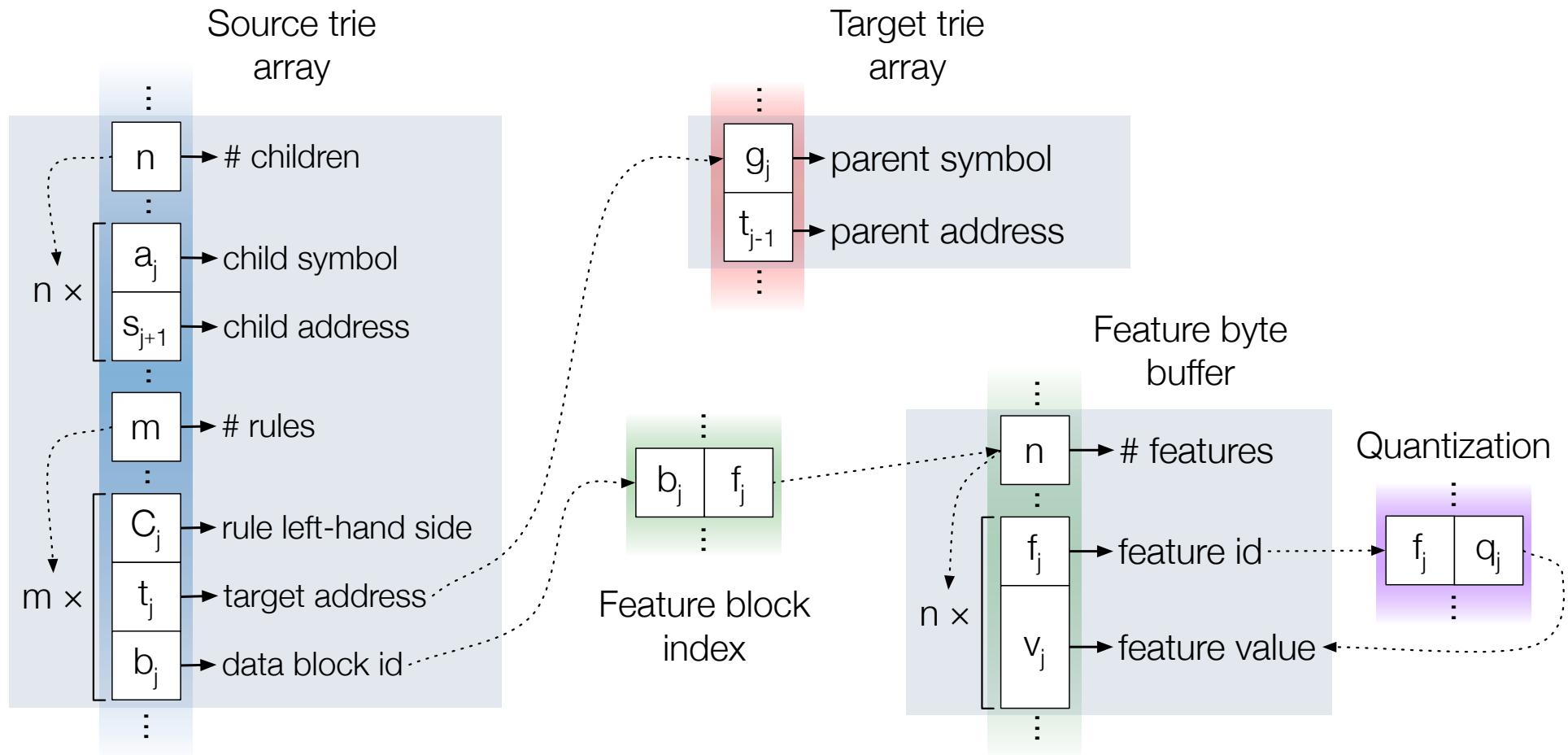
~ \$100

Text-to-Text with PPDB?

Filtering PPDB to a dev & test set:

- 118 million paraphrases
- Needs ~192 GB of memory
- Takes hours to load

Packed Grammars



Ganitkevitch et al., 2012

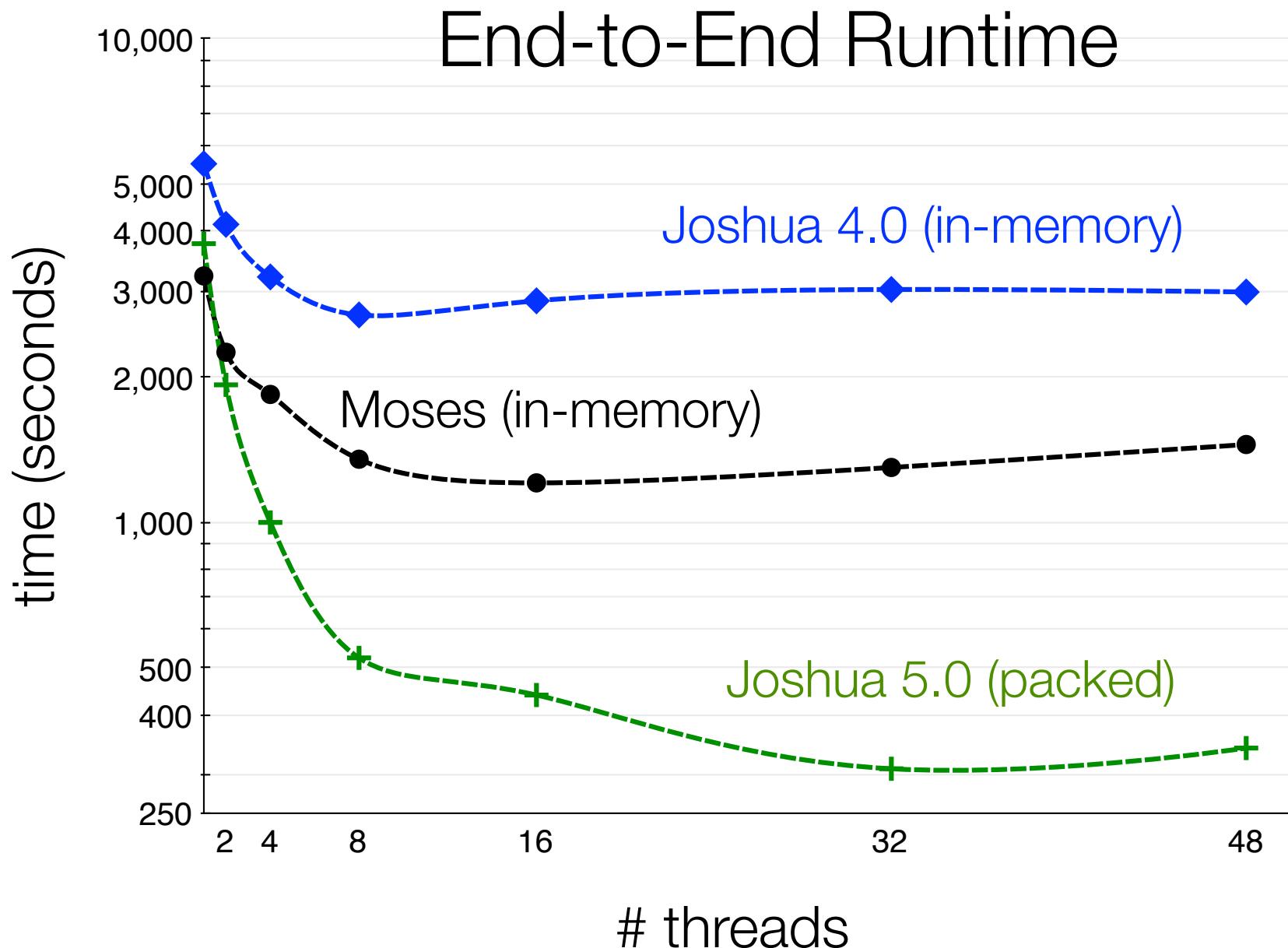
Smaller & Faster

Takes up ~6.5 GB of space (~3% of
unpacked)

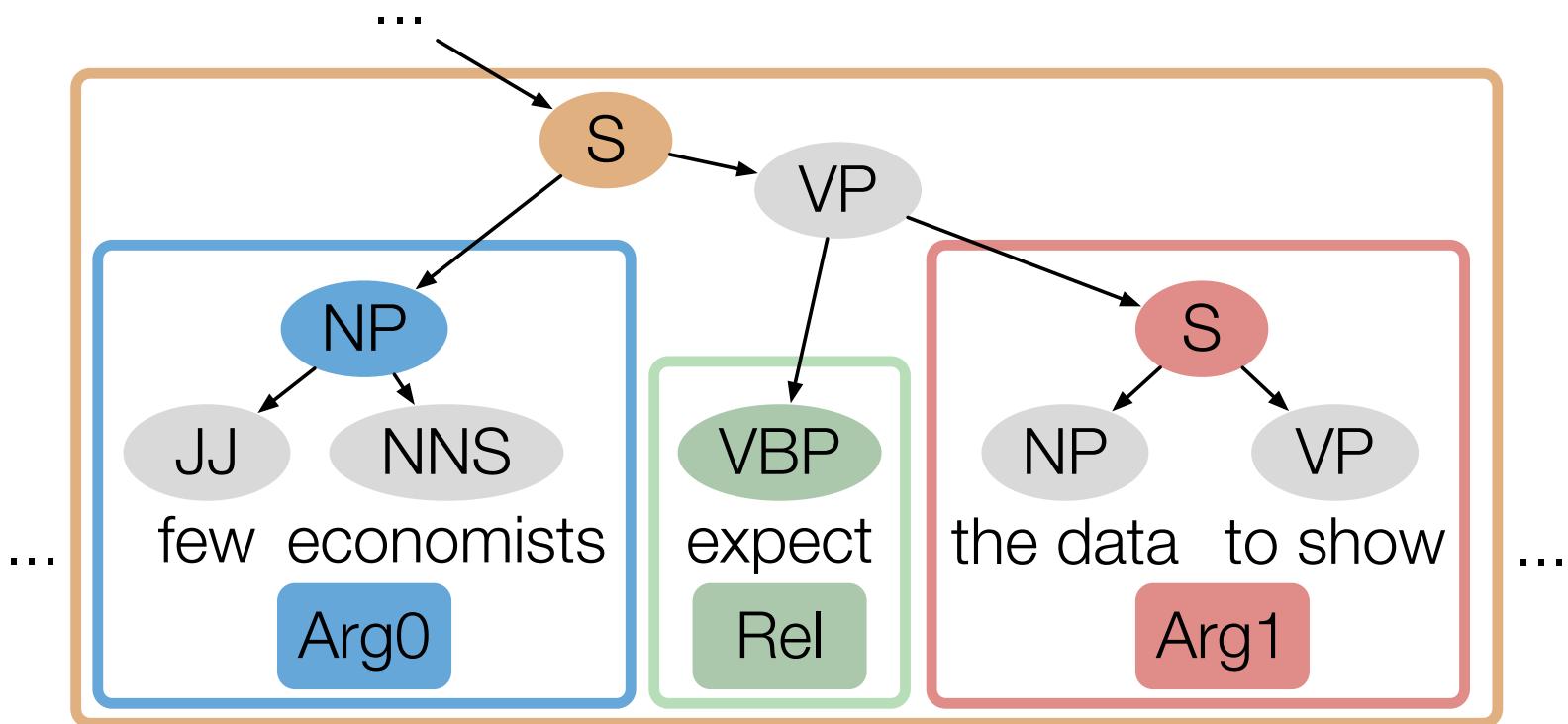
Little over 2 minutes to load

Slightly slower decoding

Joshua



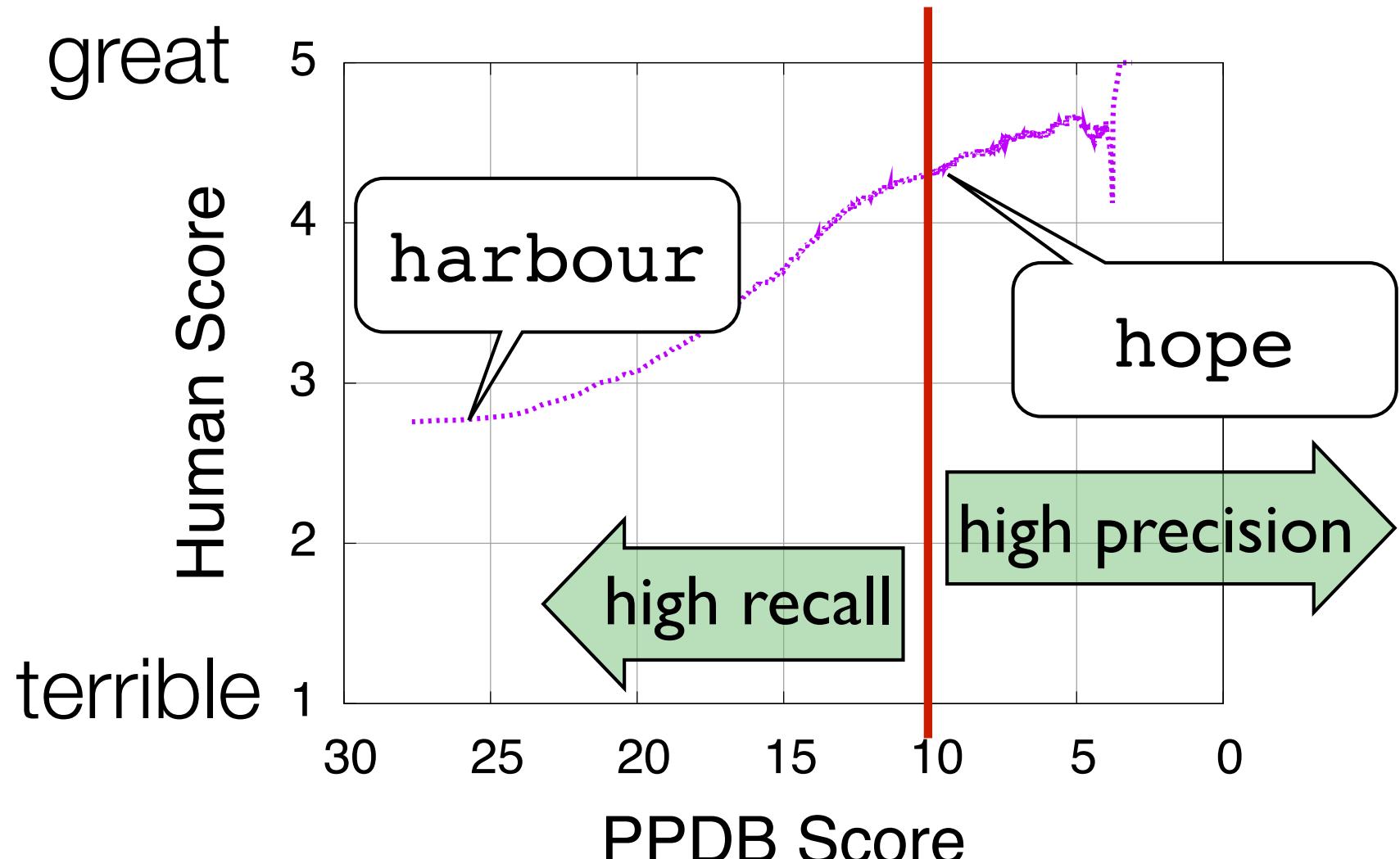
Expanding Propbank



expect

7.40909	[VBP]	await
7.48333	[VBP]	wait
8.55643	[VBP]	hope
9.21413	[VBP]	anticipate
...		
26.7452	[VBP]	grant
26.8336	[VBP]	let
26.8438	[VBP]	save
26.8846	[VBP]	harbour

Do the Scores Work?



Ganitkevitch et al., 2013

Paraphrase Extraction

Text-to-Text Applications

Scaling to PPDB

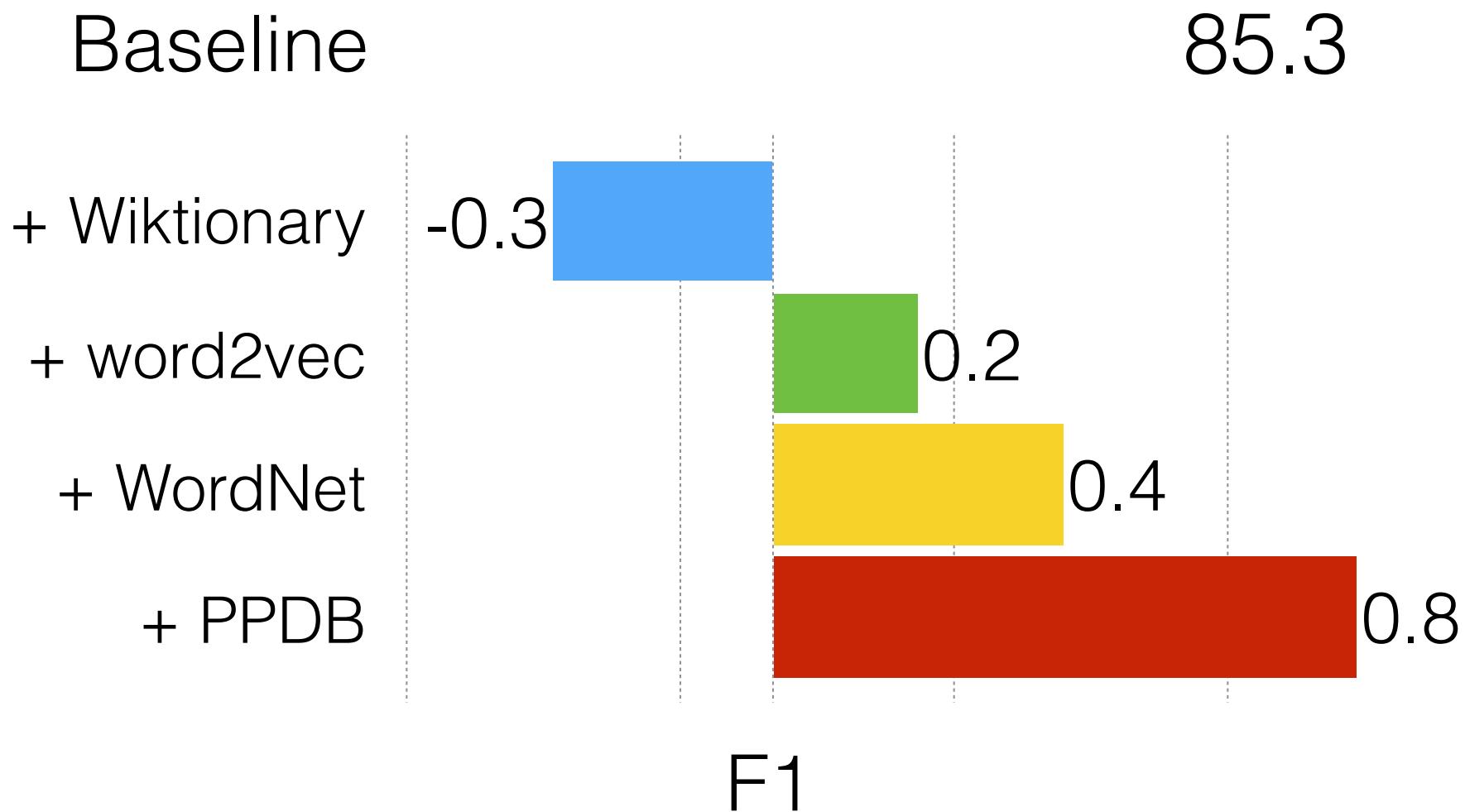
Contributions

Monolingual Alignment

Jacana Aligner (Yao et al., ACL 2013)

- Wiktionary (90k)
- Snow's 400k-WordNet
- word2vec (trained on enwik9)
- PPDB (270k)

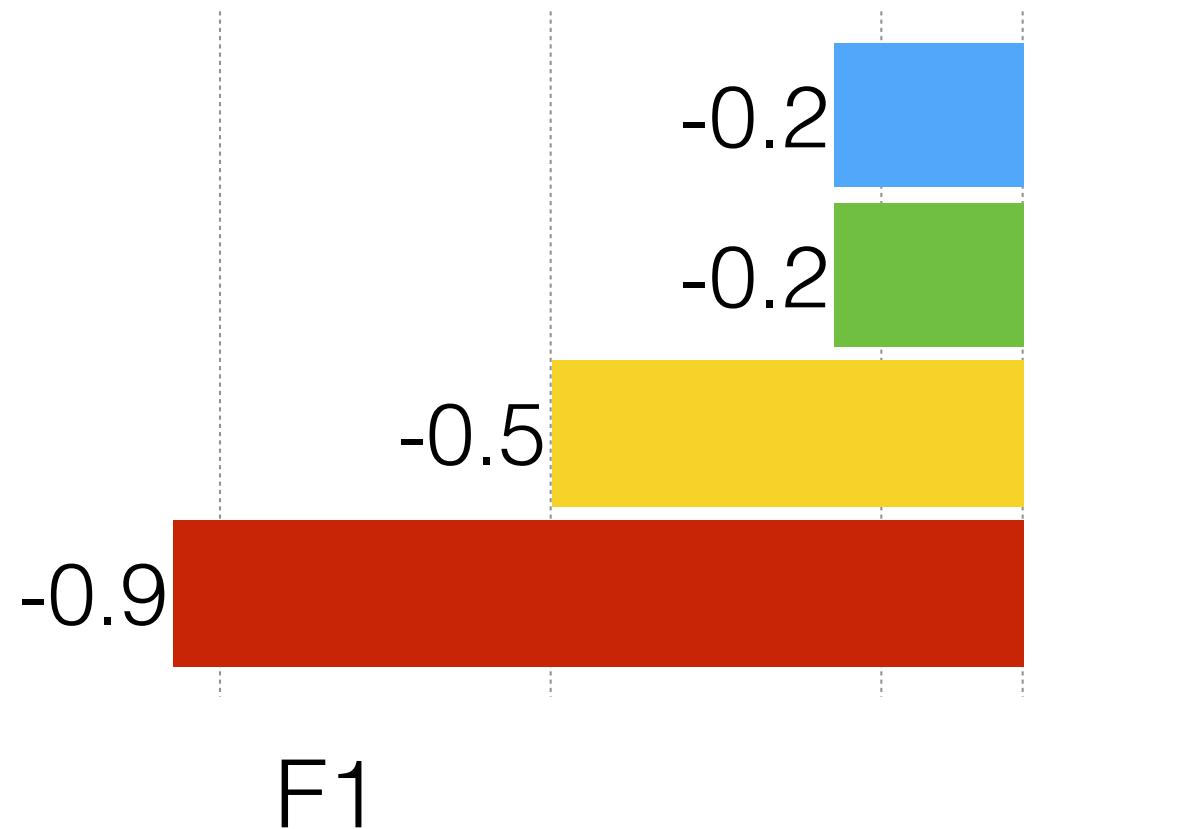
Monolingual Alignment



Monolingual Alignment

Full System 86.5

- Wiktionary
- word2vec
- WordNet
- PPDB



Retrofitting Embeddings

Semantic Relatedness (Faruqui et al., '14)

- WordNet
- FrameNet
- PPDB (XL)

Other Languages

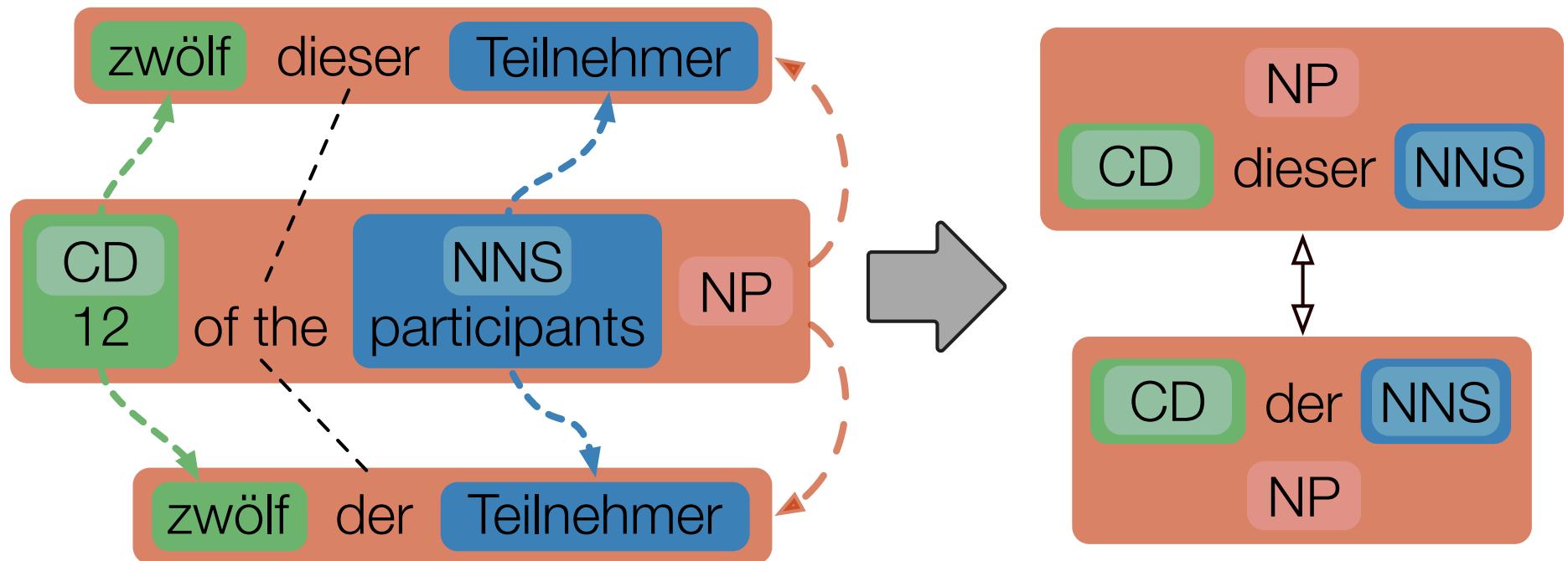
PPDB bilingual data:

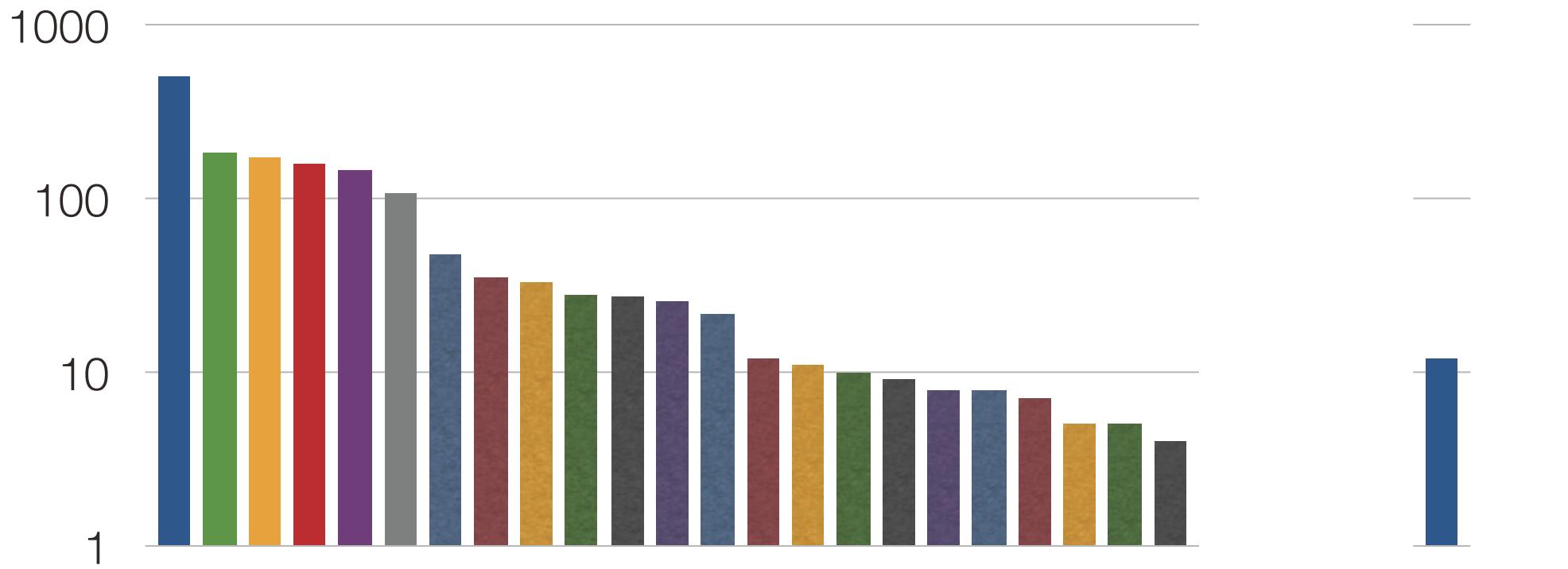
106 million sentence pairs

2 billion English words

22 pivot languages

Projected Labels





Coverage Expansion

Plug-in paraphrasing for NLP:

Query expansion, feature

mapping, input normalization, ...

AT&T

2:47 AM



**“ Siri, tell me some
astounding facts about
potatoes ”**





astounding

astonishing amazing staggering surprising stunning
startling mind-boggling incredible striking
breathtaking impressive unbelievable uncanny
perplexing narcotic stunning shocking appalling
remarkable astonishing extraordinary upside-down
dreadful puzzling awesome miraculous fantastic
meteoric horrific hideous dazzling alarming surprised
horrendous dramatic gruesome overwhelming
frightening phenomenal terrific tremendous
spectacular dire strange narcotics reversal

astounding

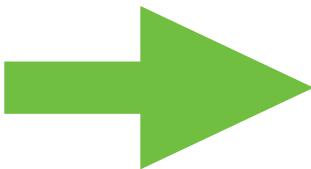
awful unexpected wonderful atrocious precipitous
non-recurring paradoxical odd whopping growth
intimidating terrible unusual horrible awe-inspiring
great embarrassing sensational runaway gorgeous
wacky intriguing unreal horrifying exciting splendid
special unthinkable surprising daunting egregious
heinous marvellous shocking inconceivable
disconcerting explosive notable brilliant noticeable
incredible rapid fabulous excellent superb bizarre
outrageous impossible

astounding

confusing troubling fascinating unlikely noteworthy
sheer massive unprecedented non-routine
exceptional unbelievable memorable amazing extra-
ordinary curious formidable groovy growing peculiar
mighty insane distinguished spectacular crazy weird
wondrous impressed compelling cool sterling stellar
extra suspicious commendable freaky sizeable
admirable immense singular unspeakable
unmatched magnificent conspicuous unconventional
scary outstanding jaw-dropping wretched

Canonicalization

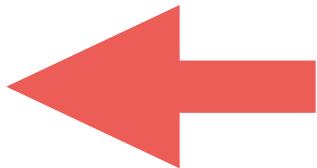
astonishing
amazing
staggering
surprising
stunning
startling
mind-boggling
incredible
striking
breathtaking
impressive
unbelievable
uncanny
perplexing
narcotic
stunning
shocking
appalling
remarkable
astonishing
extraordinary
upside-down
dreadful
puzzling
awesome
miraculous
fantastic
meteoric
horrific
hideous
dazzling
alarming
surprised
horrendous
dramatic
gruesome
overwhelming
frightening
phenomenal



astounding

Expansion

astonishing
amazing
staggering
surprising
stunning
startling
mind-boggling
incredible
striking
breathtaking
impressive
unbelievable
uncanny
perplexing
narcotic
stunning
shocking
appalling
remarkable
astonishing
extraordinary
upside-down
dreadful
puzzling
awesome
miraculous
fantastic
meteoric
horrific
hideous
dazzling
alarming
surprised
horrendous
dramatic
gruesome
overwhelming
frightening
phenomenal



astounding

Authoring



A screenshot of a Facebook search results page. The search bar at the top contains the text "Search for people, places and things". Below the search bar, a post from user "Meg Mitchell" is displayed. The post features a profile picture of Meg Mitchell in a white spacesuit. The text of the post reads: "Quick, everyone give me 10 zillion ways to say "this suggests"" followed by "about 2 weeks ago". Below the post are three interaction buttons: "Like", "Comment", and "Share".



Search for people, places and things



Juri Ganitkevitch Psh. Puny humans. Have some robot help:

there appears to be a	it meant	this would
it seems to	it indicates	they seem to
the data suggest	it seems like	there are signs
this proves	as a result	this makes it
it seems likely	that demonstrates	this calls
it is reported	it would seem	this includes
our results suggest	the results suggest	this means ,
it followed	and that means	this would mean
this illustrates	there seems to	this meant
as a consequence	there are indications	this will require
, therefore	this reflects	we conclude
it follows ,	it was found	there appear to
these findings suggest	what this means is	this would require
it also means	it shows	, it would appear
, as a consequence	the results indicated	it requires
that is to say	he seemed	it follows
the evidence suggests	and it means	this seems to
it is apparent	it would appear ,	it would appear
it does mean	, therefore ,	that means ,
it means	this means that	it sounds
we seem to	it will mean	as a result ,
this represents	evidence suggests	that is ,
there appeared to	we understand	, consequently
, it would seem	it seems as	he seemed to
this means is	that includes	this has meant
, it seems	it therefore appears	therefore to
it mean	it seems ,	it seemed to
, as a result ,	this results	that implies
it seems	it seems to be	he seems to be
apparently ,	this points	, it follows
that indicates	these results suggest	the data indicate
, it appeared	it may be concluded	, this means
the implication is	there seems	that means
it was felt	, it means	these observations suggest
this means	there is some evidence	this creates
there seems to be a		

The Real World™

pi

P A R E N T A L

A D V I S O R Y

s.

E X P L I C I T C O N T E N T

a.

ed

To Summarize

- Scaled rich paraphrase extraction to large bitext corpora
- Integrated distributional signal with bitext-derived paraphrases
- Open source infrastructure to enable text-to-text with PPDB
- Created a resource for multilingual coverage expansion in NLP

thank you for your time uh , thanks
thanks , man thank you for your attention
keep the change you look amazing
anyway , thanks thank you , frank
diet coke **Thank you!** gee , thanks

any objections question time anything else
Questions?
the interrogation your questions
any more questions
oral questions the crux of the problem
budgetary questions humanitarian problems