

---

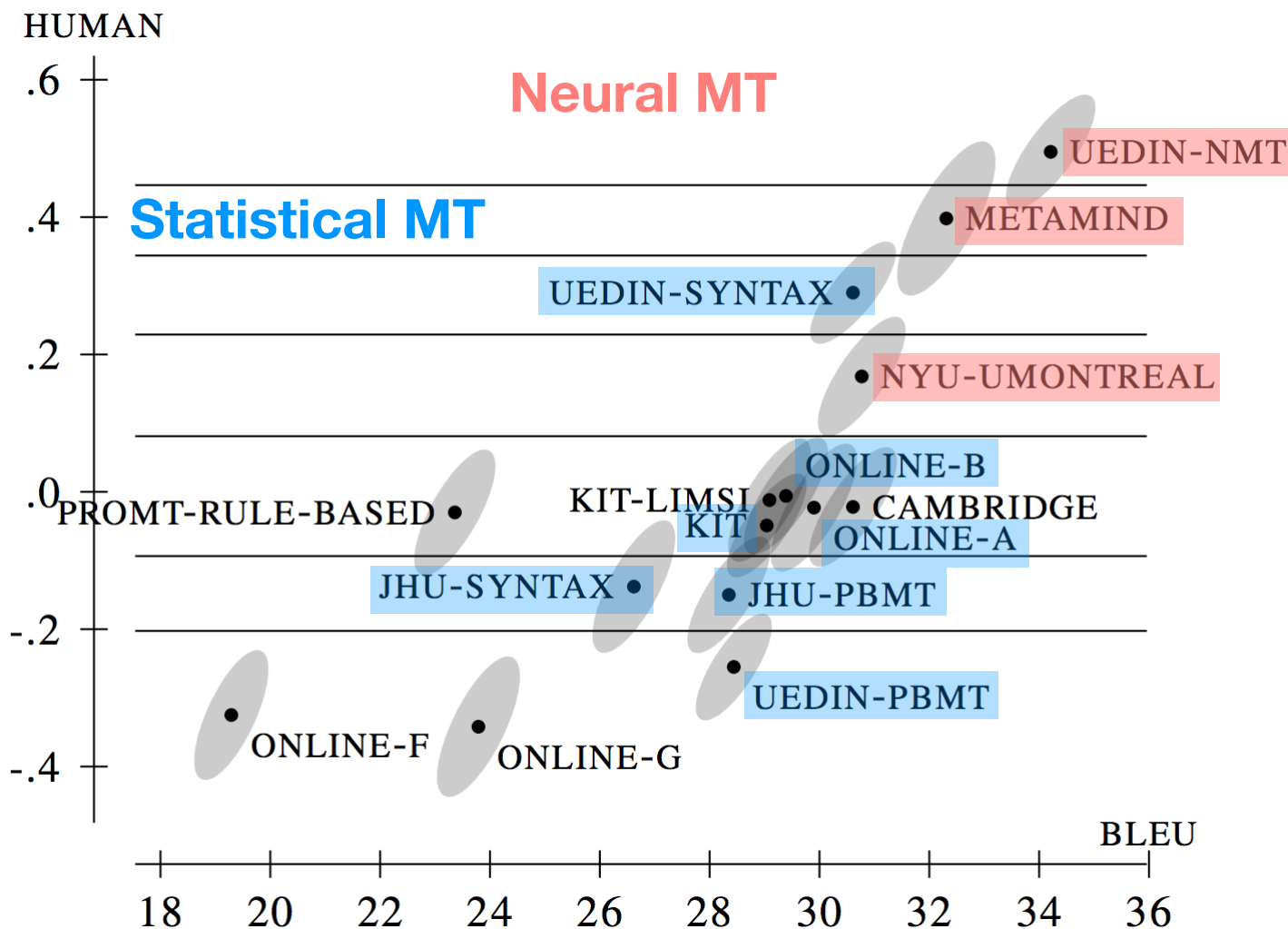
# Current Challenges

Philipp Koehn

3 November 2022

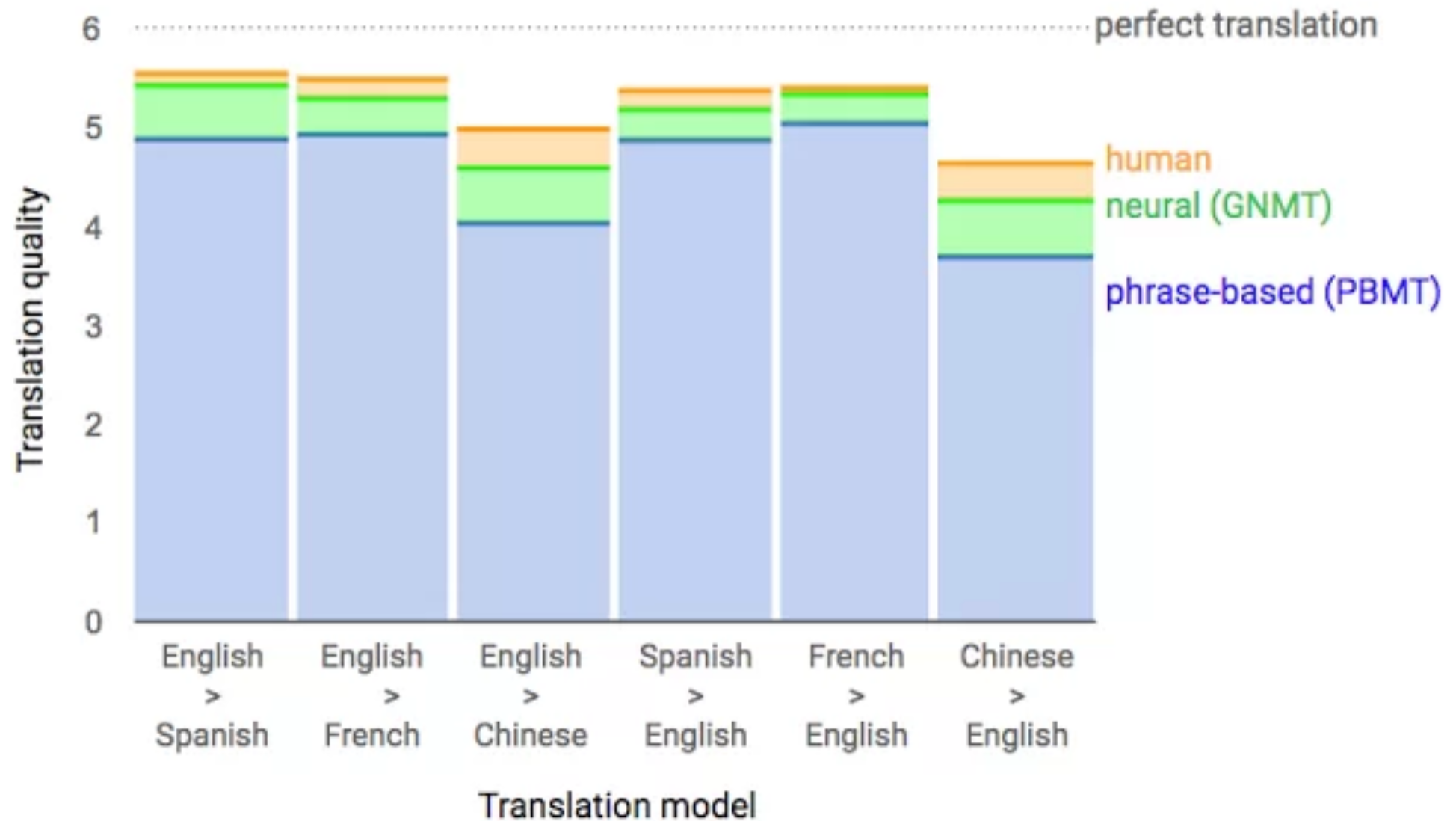


# WMT 2016



(in 2017 barely any statistical machine translation submissions)

# 2017: Google: "Near Human Quality"



# 2018: More Hype



## Microsoft Research Achieves Human Parity For Chinese English Translation

Written by Sue Gee

Wednesday, 21 March 2018

Researchers in Microsoft's labs in Beijing and in Redmond and Washington have developed an AI machine translation system that can translate with the same accuracy as a human from Chinese to English.

## SDL Cracks Russian to English Neural Machine Translation

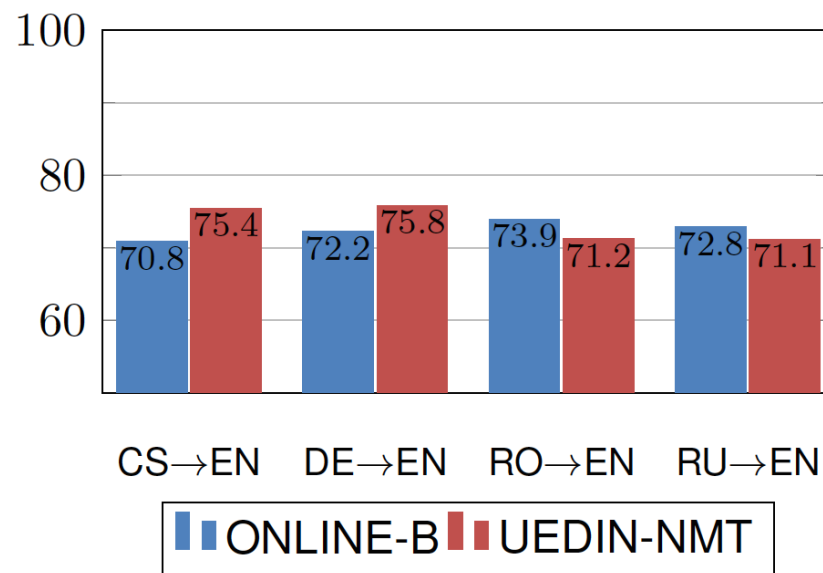
Global Enterprises to Capitalize on Near Perfect Russian to English Machine Translation as SDL Sets New Industry Standard

*“90% of the system’s output labelled as perfect by professional Russian-English translators”*

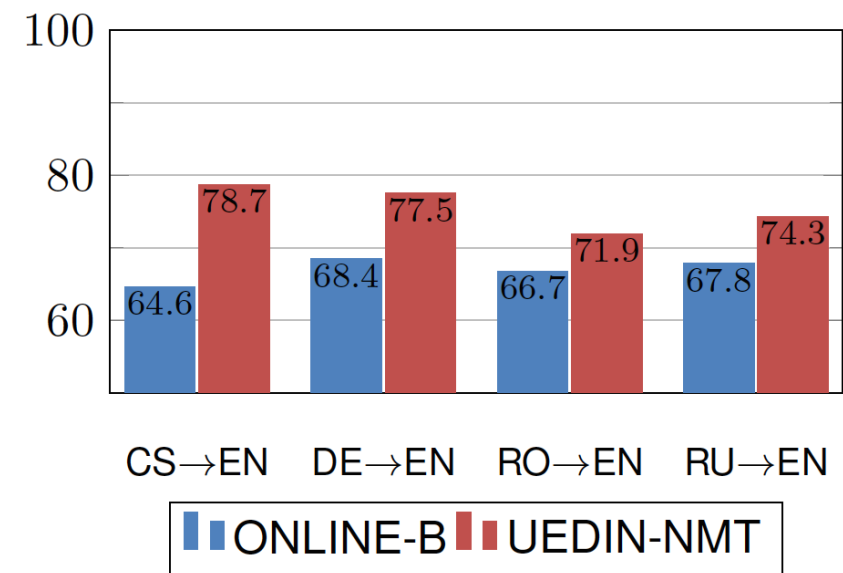
# Just Better Fluency?



Adequacy  
+1%



Fluency  
+13%



(from: Sennrich and Haddow, 2017)

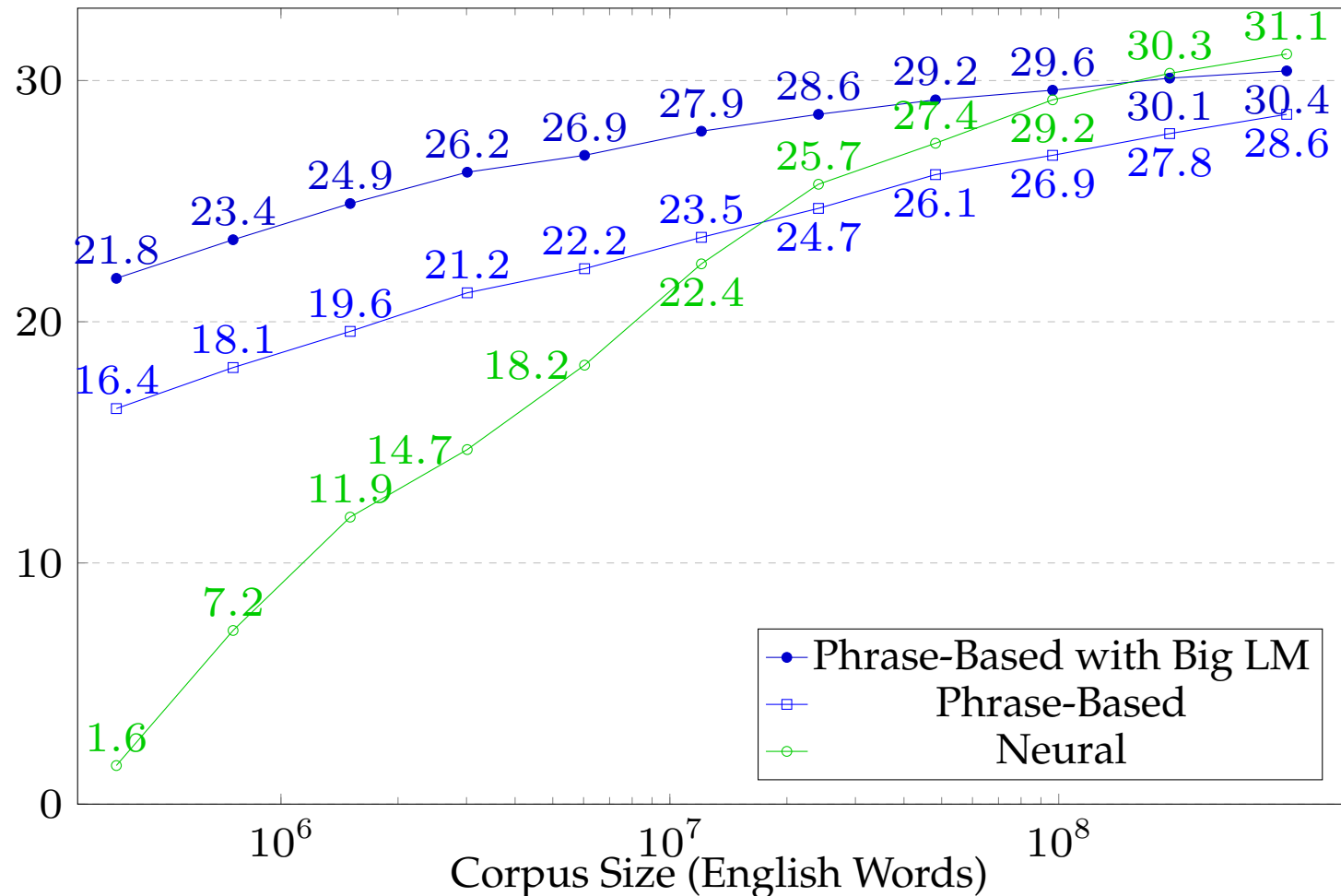
# Challenges



- Lack of training data
- Domain mismatch
- Rare words
- Word alignment
- Beam search
- Noise
- Control over output
- Interpretability

# lack of training data

# Amount of Training Data



English-Spanish systems trained on 0.4 million to 385.7 million words








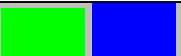
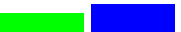




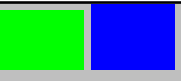

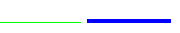






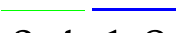
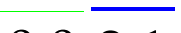
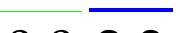







# Translation Examples



Source	A Republican strategy to counter the re-election of Obama
$\frac{1}{1024}$	Un órgano de coordinación para el anuncio de libre determinación
$\frac{1}{512}$	Lista de una estrategia para luchar contra la elección de hojas de Ohio
$\frac{1}{256}$	Explosión realiza una estrategia divisiva de luchar contra las elecciones de autor
$\frac{1}{128}$	Una estrategia republicana para la eliminación de la reelección de Obama
$\frac{1}{64}$	Estrategia siria para contrarrestar la reelección del Obama .
$\frac{1}{32} +$	Una estrategia republicana para contrarrestar la reelección de Obama

# domain mismatch

# Domain Mismatch

System ↓	Law	Medical	IT	Koran	Subtitles
<b>All Data</b>	 30.5 32.8	 45.1 42.2	 35.3 44.7	 17.9 17.9	 26.4 20.8
<b>Law</b>	 31.1 34.4	 12.1 18.2	 3.5 6.9	 1.3 2.2	 2.8 6.0
<b>Medical</b>	 3.9 10.2	 39.4 43.5	 2.0 8.5	 0.6 2.0	 1.4 5.8
<b>IT</b>	 1.9 3.7	 6.5 5.3	 42.1 39.8	 1.8 1.6	 3.9 4.7
<b>Koran</b>	 0.4 1.8	 0.0 2.1	 0.0 2.3	 15.9 18.8	 1.0 5.5
<b>Subtitles</b>	 7.0 9.9	 9.3 17.8	 9.2 13.6	 9.0 8.4	 25.9 22.1

# Translation Examples

Source	Schaue um dich herum.
Ref.	Look around you.
All	NMT: Look around you. SMT: Look around you.
Law	NMT: Sughum gravecorn. SMT: In order to implement dich Schaue .
Medical	NMT: EMEA / MB / 049 / 01-EN-Final Work programme for 2002 SMT: Schaue by dich around .
IT	NMT: Switches to paused. SMT: To Schaue by itself . \t \t
Koran	NMT: Take heed of your own souls. SMT: And you see.
Subtitles	NMT: Look around you. SMT: Look around you .



# rare words

# Rare Words

- More frequent in training  $\rightarrow$  more likely to get right in test
- Let's measure this

- More frequent in training → more likely to get right in test
- Let's measure this
- One problem
  - frequency measured for input words
  - translation correctness measured for output words

# Translation Accuracy for Input Words

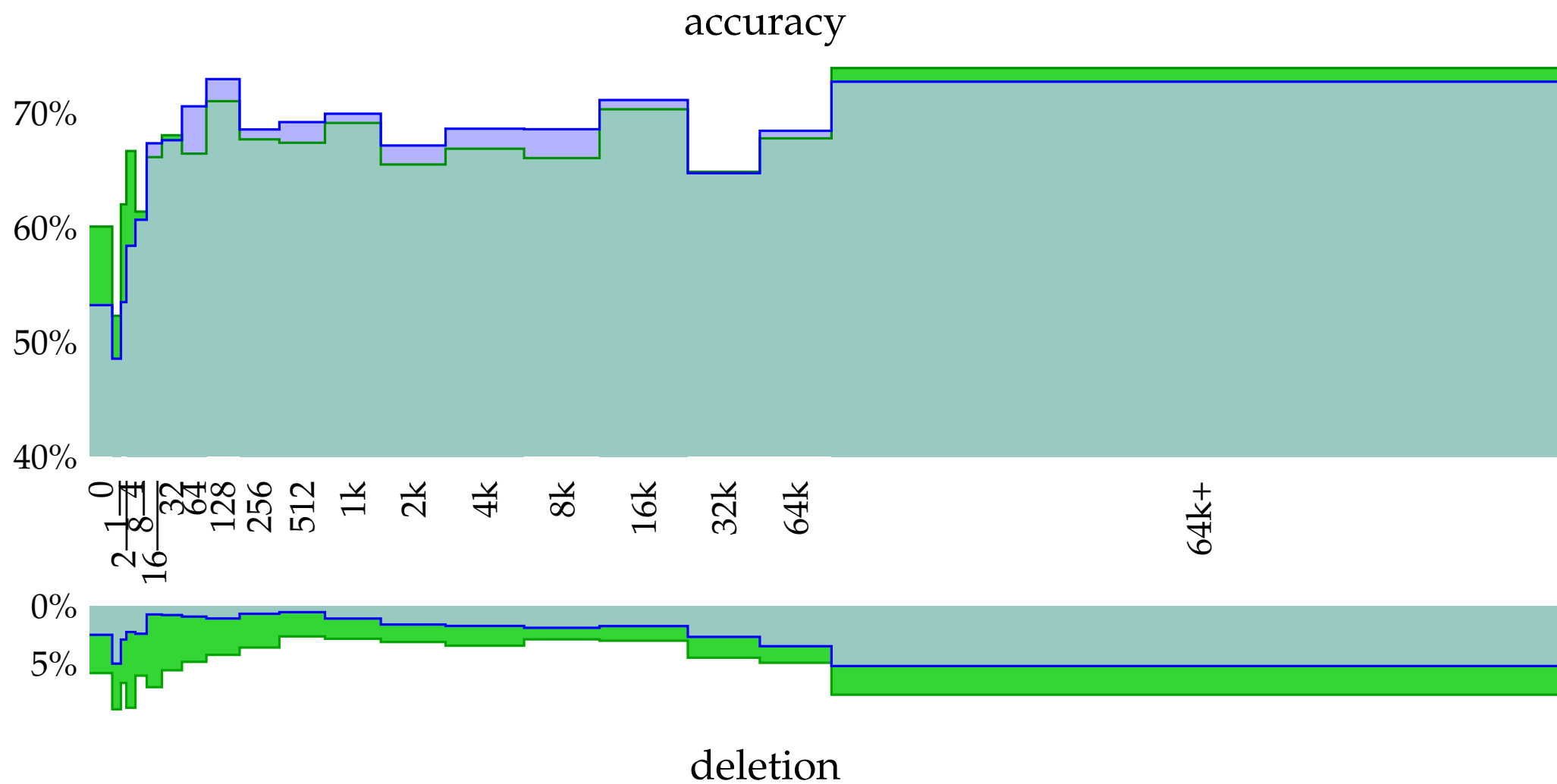
- Generate word alignment between input and output words
- Look up count of input word in training
- Link to output word via word alignment
- Check if it is also in the reference translation



# Translation Accuracy for Input Words

- Generate word alignment between input and output words
- Look up count of input word in training
- Link to output word via word alignment
- Check if it is also in the reference translation
- A lot of tedious special cases
  - one-to-many alignment, only some output words in reference
  - input word not aligned to any target word
  - many-to-one alignment
  - output word occurs multiple time in output or reference sentence

# Count vs. Accuracy





# word alignment

# Word Alignment

	relations	between	Obama	and	Netanyahu	have	been	strained	for	years	.
die	56		16								
Beziehungen	89										
zwischen		72	26								
Obama			96								
und				79							
Netanjahu					98						
sind						42	11	38			
seit								22	54	10	
Jahren										98	
angespannt								84			
.						11	14	23			49

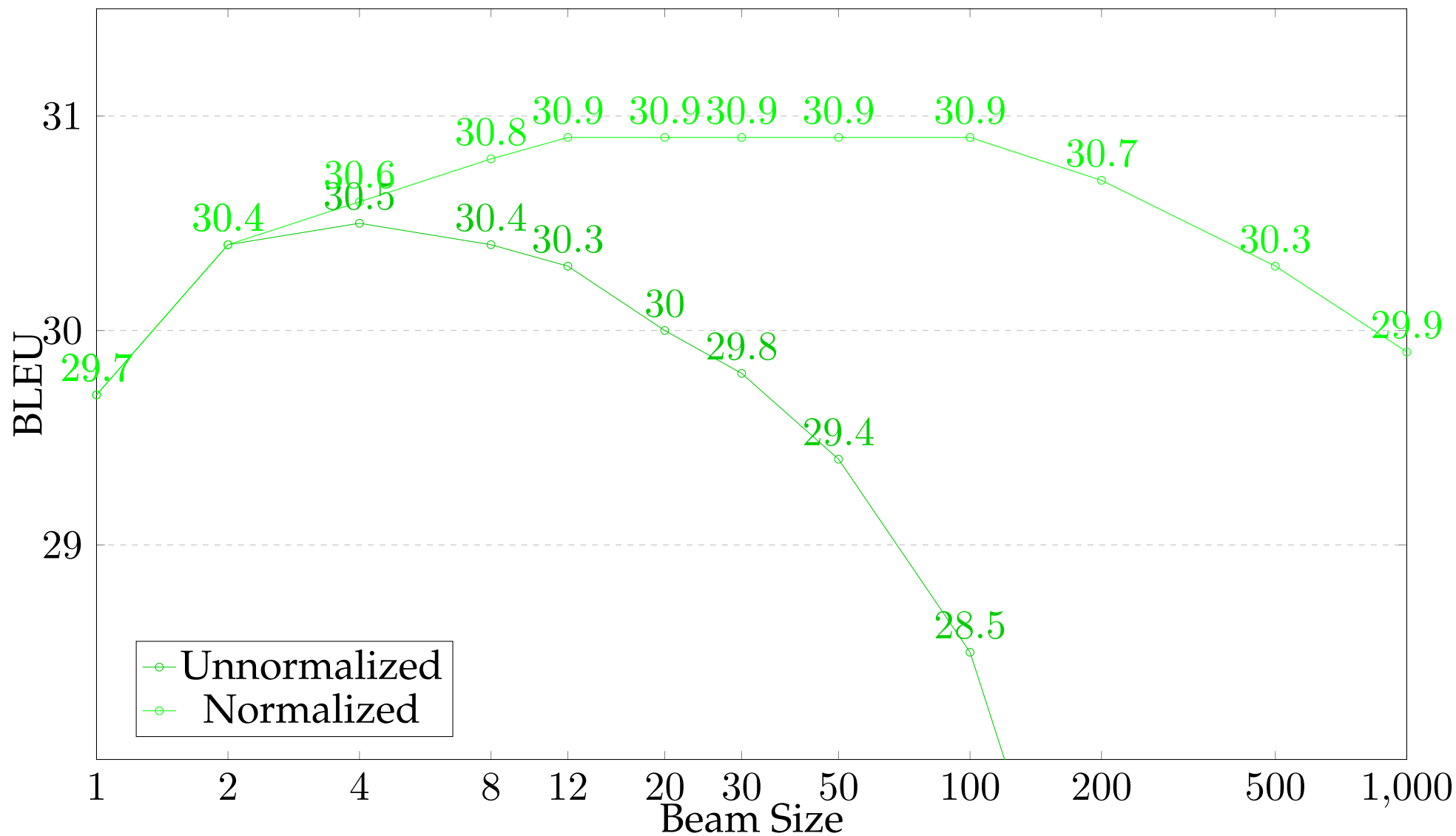
# Word Alignment?

	das	Verhältnis	zwischen	Obama	und	Netanyahu	ist	seit	Jahren	gespannt	.
the		47							17		
relationship			81								
between				72							
Obama					87						
and						93					
Netanyahu							95				
has								38	16		26
been								21	14		54
stretched											77
for									38	33	12
years										90	
.	11								19	32	17



# beam search

# Beam Search





# noisy data



# Noise in Training Data

- Crawled parallel data from the web (very noisy)

	SMT	NMT
WMT17	24.0	27.2
+ Paracrawl	25.2 (+1.2)	17.3 (-9.9)

(German-English, 90m words each of WMT17 and Crawl data)

	5%	10%	20%	50%	100%
Raw crawl data	<div><div>27.4</div><div>24.2</div><div>+0.2 +0.2</div></div>	<div><div>26.6</div><div>24.2</div><div>-0.9 +0.2</div></div>	<div><div>24.7</div><div>24.4</div><div>-2.5 +0.4</div></div>	<div><div>20.9</div><div>24.8</div><div>-6.3 +0.8</div></div>	<div><div>17.3</div><div>25.2</div><div>-9.9 +1.2</div></div>

- Corpus cleaning methods [Xu and Koehn, EMNLP 2017] give improvements

# Types of Noise

- Misaligned sentences
- Disfluent language (from MT, bad translations)
- Wrong language data (e.g., French in German–English corpus)
- Untranslated sentences
- Short segments (e.g., dictionaries)
- Mismatched domain

# Mismatched Sentences

- Artificial created by randomly shuffling sentence order
- Added to existing parallel corpus in different amounts

5%	10%	20%	50%	100%
<div><div>24.0</div><div>-0.0</div></div>	<div><div>24.0</div><div>-0.0</div></div>	<div><div>23.9</div><div>-0.1</div></div>	<div><div>26.1</div><div>-1.1</div></div> <div><div>23.9</div><div>-0.1</div></div>	<div><div>25.3</div><div>-1.9</div></div> <div><div>23.4</div><div>-0.6</div></div>

- Bigger impact on NMT (green, left) than SMT (blue, right)

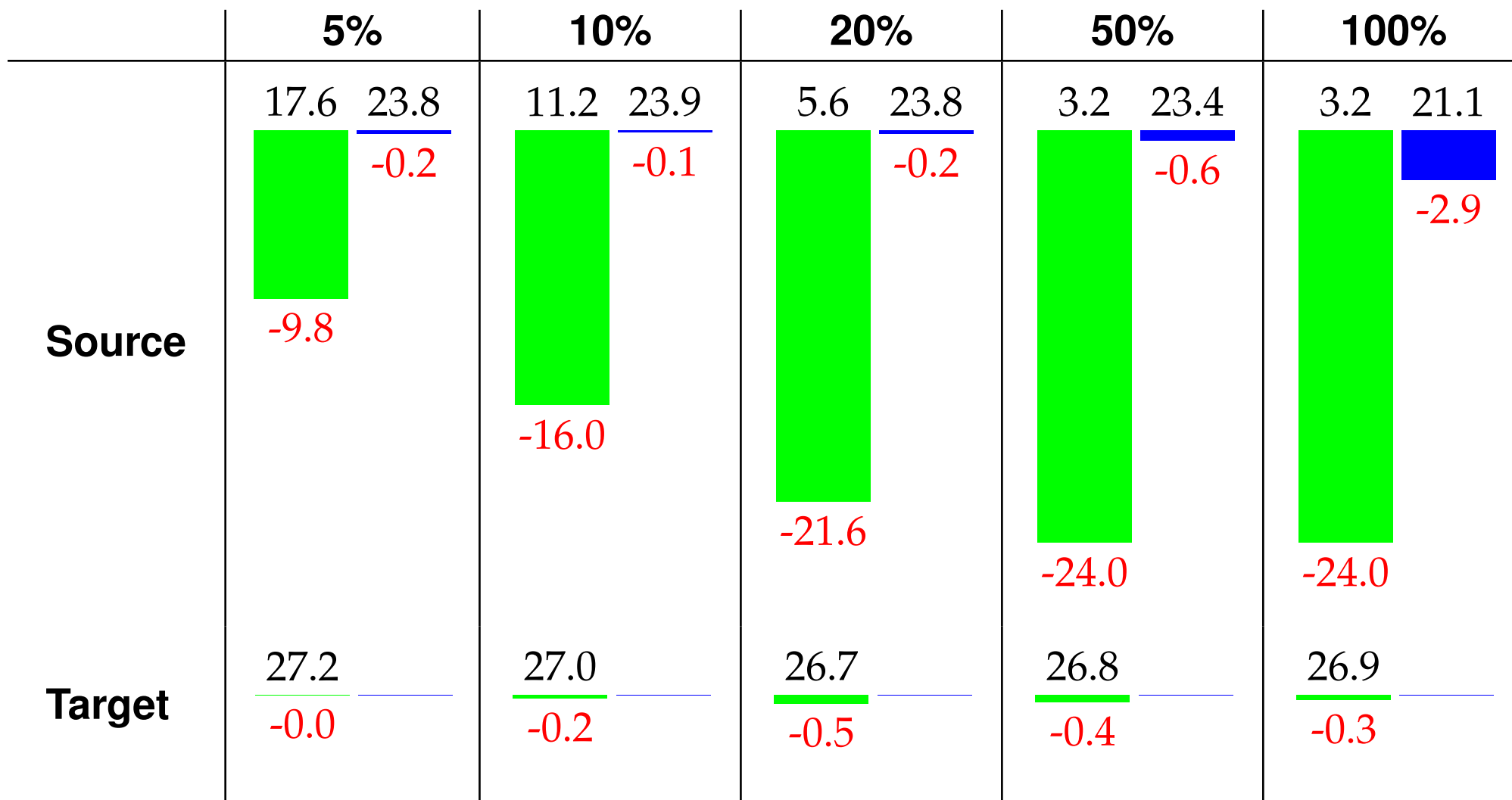
# Misordered Words

- Artificial created by randomly shuffling words in each sentence

	5%	10%	20%	50%		100%	
Source	24.0	23.6	23.9	26.6	23.6	25.5	23.7
	-0.0	-0.4	-0.1	-0.6	-0.4	-1.7	-0.3
Target	24.0	24.0	23.4	26.7	23.2	26.1	22.9
	-0.0	-0.0	-0.6	-0.5	-0.8	-1.1	-1.1

- Similar impact on NMT than SMT, worse for source reshuffle

# Untranslated Sentences



# Wrong Language

	5%	10%	20%	50%	100%
<b>fr source</b>	<u>26.9</u> <u>24.0</u> -0.3 -0.0	<u>26.8</u> <u>23.9</u> -0.4 -0.1	<u>26.8</u> <u>23.9</u> -0.4 -0.1	<u>26.8</u> <u>23.9</u> -0.4 -0.1	<u>26.8</u> <u>23.8</u> -0.4 -0.2
<b>fr target</b>	<u>26.7</u> <u>24.0</u> -0.5 -0.0	<u>26.6</u> <u>23.9</u> -0.6 -0.1	<u>26.7</u> <u>23.8</u> -0.5 -0.2	<u>26.2</u> <u>23.5</u> -1.0 -0.5	<u>25.0</u> <u>23.4</u> -2.2 -0.6

- Surprisingly robust, maybe due to domain mismatch of French data

# Short Sentences

	5%	10%	20%	50%
<b>1-2 words</b>	$\frac{27.1}{-0.1} \frac{24.1}{+0.1}$	$\frac{26.5}{-0.7} \frac{23.9}{-0.1}$	$\frac{26.7}{-0.5} \frac{23.8}{-0.2}$	
<b>1-5 words</b>	$\frac{27.8}{+0.6} \frac{24.2}{+0.2}$	$\frac{27.6}{+0.4} \frac{24.5}{+0.5}$	$\frac{28.0}{+0.8} \frac{24.5}{+0.5}$	$\frac{26.6}{-0.6} \frac{24.2}{+0.2}$

- No harm done

# control over output



# Specifying Decoding Constraints

- Overriding the decisions of the decoder
- Why?
  - ⇒ translations have followed strict terminology
  - ⇒ rule-based translation of dates, quantities, etc.

The `<x translation="Router"> router </x> is <wall/>`  
a model `<zone> Psy X500 Pro </zone> .`

- The XML tags specify to the decoder that
  - the word `router` to be translated as `Router`
  - `The router is,` to be translated before the rest (`<wall/>`)
  - brand name `Psy X500 Pro` to be translated as a unit (`<zone>`, `</zone>`)

# Formal Constraints

- Subtitles
  - translation has to fit into space on screen  
(may have to be shortened)
  - input and output broken up into lines

# Formal Constraints

- Subtitles
  - translation has to fit into space on screen (may have to be shortened)
  - input and output broken up into lines
- Speech translation
  - input often not well-formed
  - real time translation: start while sentence is spoken
  - subtitles: have to be readable in limited time
  - dubbing: sync up with video of speaker's mouth movement

- Subtitles
  - translation has to fit into space on screen (may have to be shortened)
  - input and output broken up into lines
- Speech translation
  - input often not well-formed
  - real time translation: start while sentence is spoken
  - subtitles: have to be readable in limited time
  - dubbing: sync up with video of speaker's mouth movement
- Poetry
  - meter
  - rhyme

# catastrophic errors

News | Science and Technology

## Facebook apologises for rude mistranslation of Xi Jinping's name

*Company blames technical glitch that 'caused incorrect translations' of Chinese leader's name from Burmese to English.*

## Facebook's auto translation AI fail leads to a nightmare for a Palestinian man

The AI feature had "Good morning" in Arabic wrongly translated as "attack them" in Hebrew.

By [Gianluca Mezzofiore](#) on October 24, 2017



Industry News • By [Marion Marking](#) On 3 Aug 2020

## Thai Mistranslation Shows Risk of Auto-Translating Social Media Content



After a machine translation of a post from English into Thai about the King's birthday proved offensive to the Thai monarchy, Facebook Thailand said it was deactivating auto-translate on Facebook and Instagram, revamping machine translation (MT) quality, and offering the Thai people its "profound apology."

# What are Catastrophic Errors?

- Generation of profanity
  - first step: maintain list of offensive words for each language
  - only eliminate these words, if the input did not include such words
  - but: offensive language is not limited to specific words
- Generation of violent / inciting content
- Opposite meaning
- Mistranslation of names

⇒ All this is hard to detect



# robustness

# Robustness to User Generated Content

English ▼

↔

German ▼

daily content of  
#scaramouche  
from genshin  
impact #原神 ★  
mute  
#mouchecc for  
no cc tweets !  
not leak free ★  
<http://dailymouche.carrd.co>

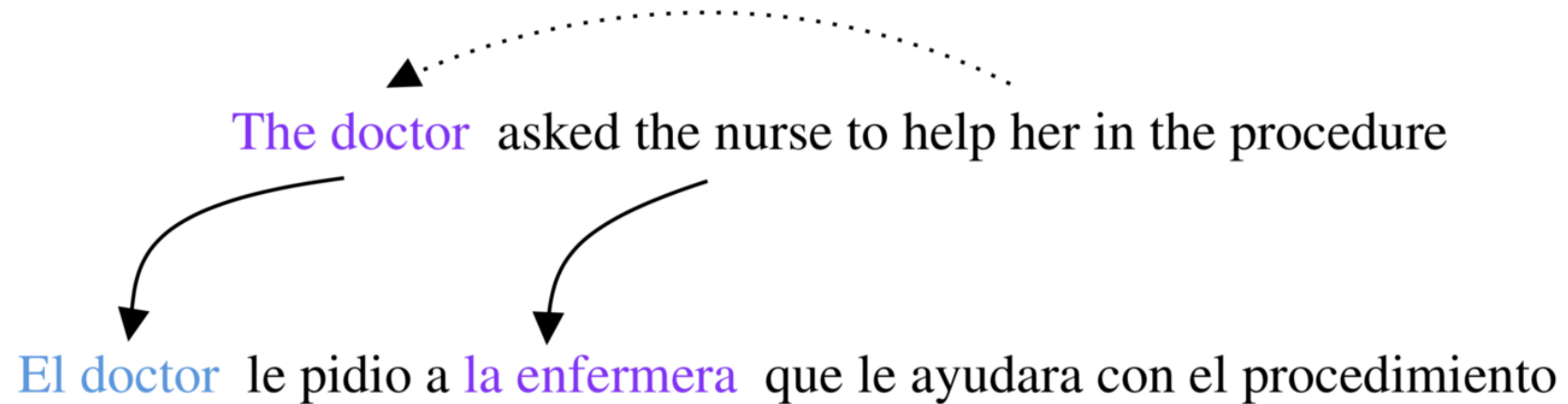
×

täglicher Inhalt von  
#scaramouche von  
genshin impact #原神  
★ stumm  
#mouchecc für keine  
CC-Tweets! nicht  
auslaufsicher ★  
<http://dailymouche.carrd.co>

- Jargon and acronyms
- Misspellings (sometimes intended for effect)
- Mangled grammar
- Special symbols (emojis, etc.)
- Hashtags, URLs, ...
- Use of dialectical languages
- Use of non-standard writing systems (e.g., Latin script due to lack of keyboard)

- Special handling of non-words like emojis, hashtags, URLs
- Creating synthetic noisy training data
- Adversarial training
- Resources
  - Machine translation of noisy text data set (MTNT)
  - WMT 2020 Shared Task on Machine Translation Robustness

bias



# Gender Bias

42



English ▼

↔

Spanish ▼

the doctor said: take the pill.	×	La doctora dijo: toma la píldora. <i>(feminine)</i>
		El doctor dijo: toma la píldora. <i>(masculine)</i>

[Open in Google Translate](#)

[Feedback](#)

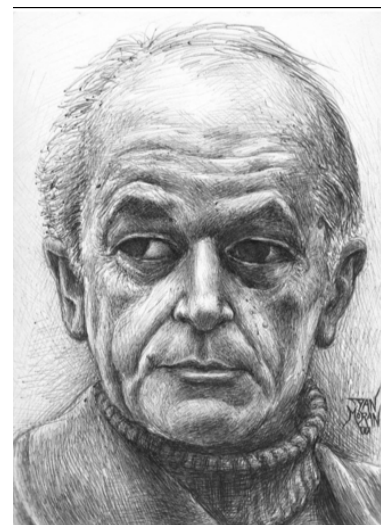
## **“You Sound Just Like Your Father” Commercial Machine Translation Systems Include Stylistic Biases**

**Dirk Hovy**

**Federico Bianchi**  
Bocconi University  
Via Sarfatti 25, 20136  
Milan, Italy

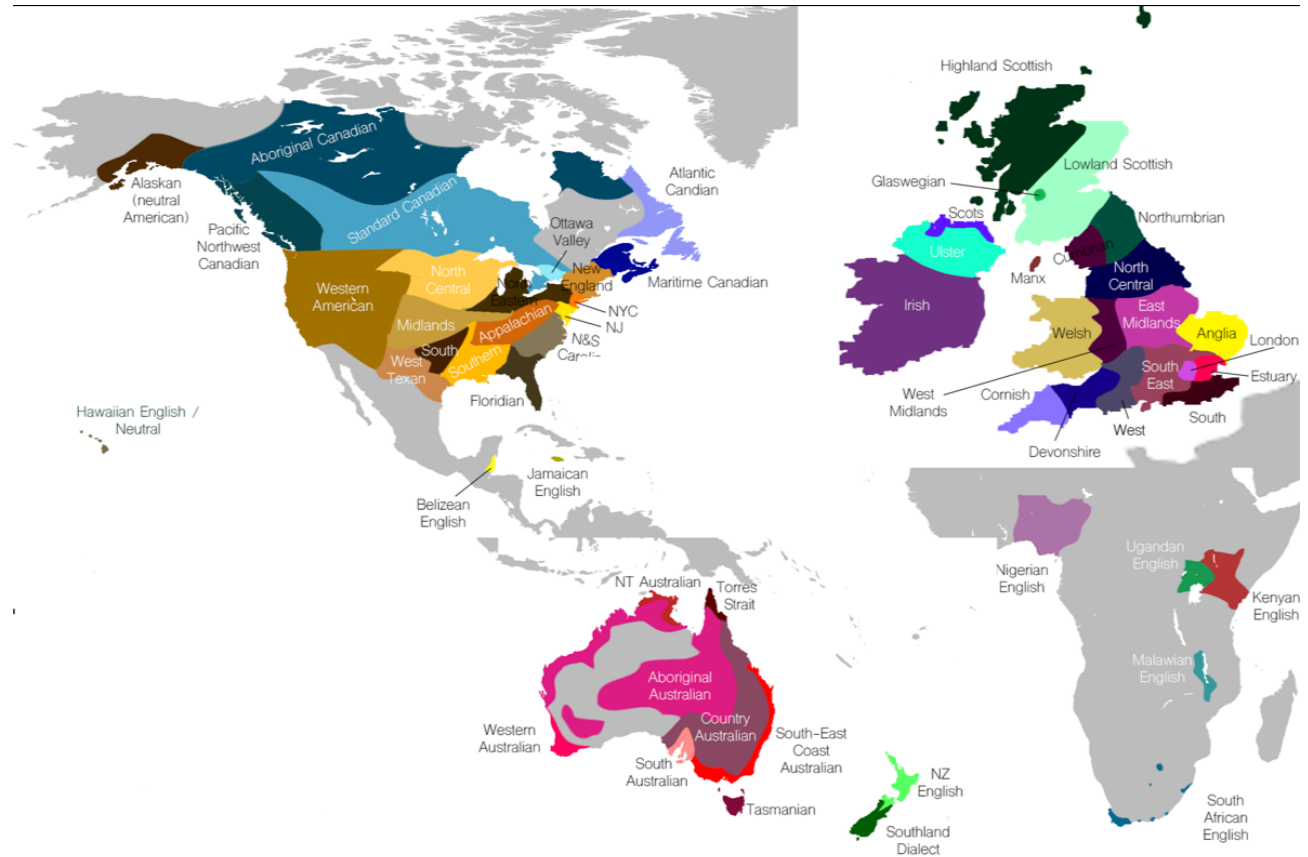
**Tommaso Fornaciari**

{dirk.hovy, f.bianchi, fornaciari.tommaso}@unibocconi.it



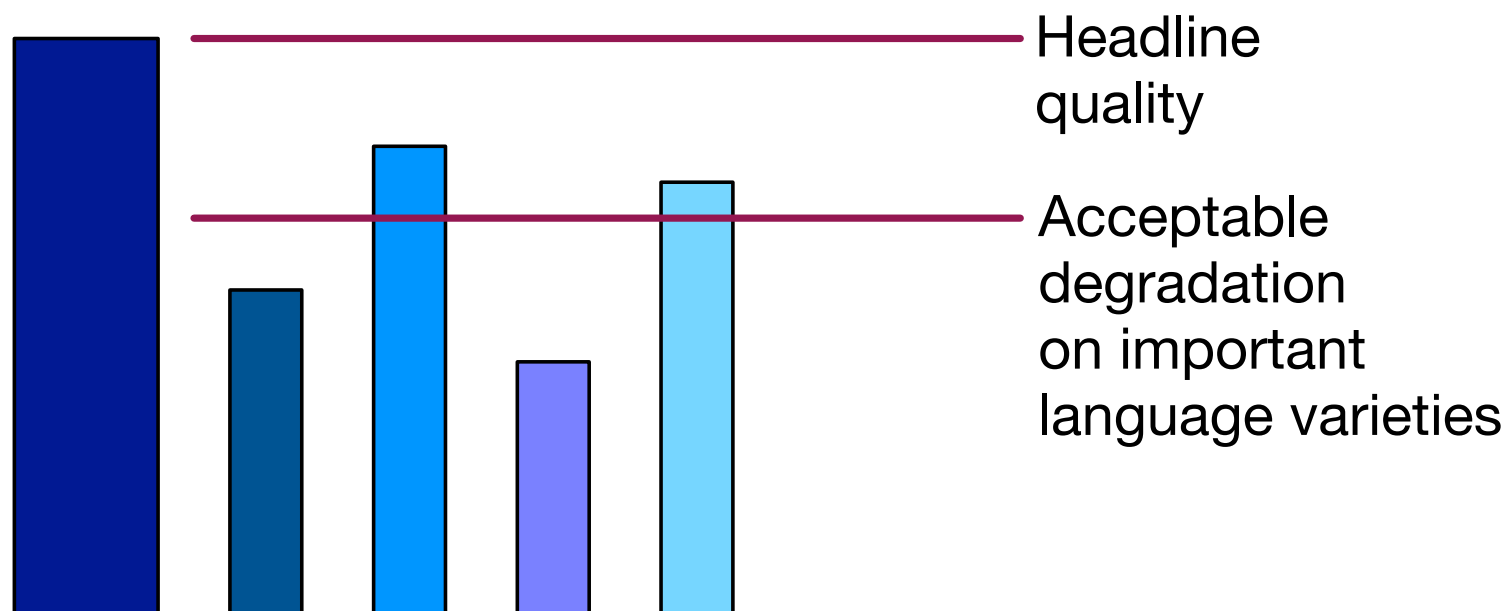


- Models often trained only on standard languages (British, American)
- Work less well on other dialects
- Bigger problem for automatic speech recognition



# Evaluate Across Language Varieties

- BLEU score on standard language is not enough
- Also need test sets for each language variety



# document-level translation

# Document-Level Translation

*The shop is selling a nice table.  
Jane is quite taken by it.  
The table would match the chairs in her living room.*

- Machine translation translates one sentence at a time
- But: surrounding context may help

# Document-Level Translation

*The shop is selling a nice table.  
Jane is quite taken by it.  
The table would match the chairs in her living room.*

- Machine translation translates one sentence at a time
- But: surrounding context may help
  - translation of pronouns may require co-reference

# Document-Level Translation

*The shop is selling a nice table.  
Jane is quite taken by it.  
The table would match the chairs in her living room.*

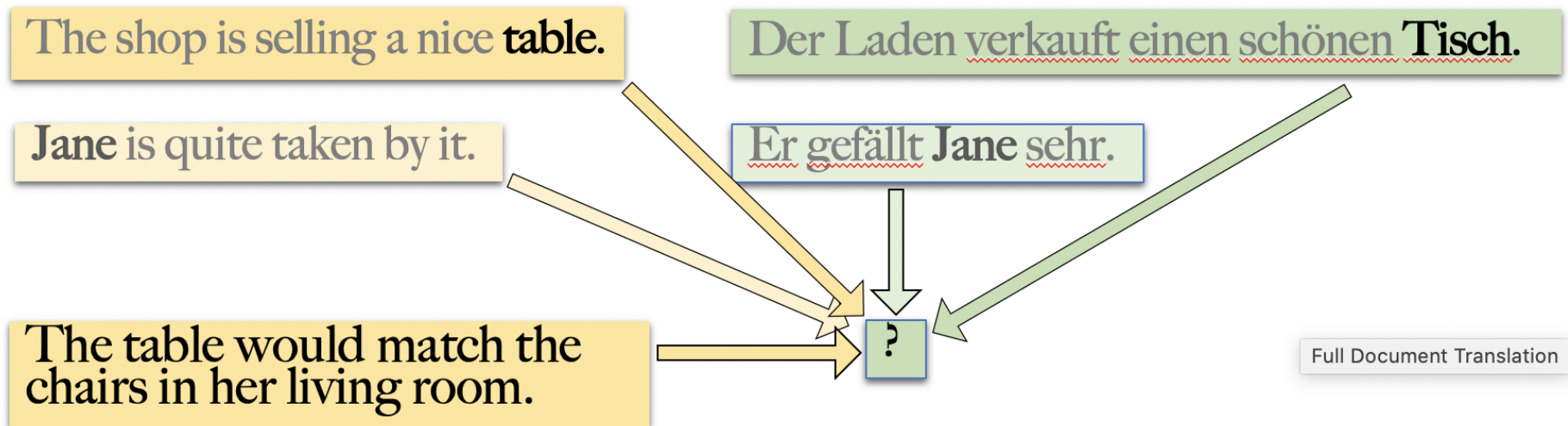
- Machine translation translates one sentence at a time
- But: surrounding context may help
  - translation of pronouns may require co-reference
  - ambiguous words may be informed by broader context

# Document-Level Translation

*The shop is selling a nice **table**.  
Jane is quite taken by it.  
The **table** would match the chairs in her living room.*

- Machine translation translates one sentence at a time
- But: surrounding context may help
  - translation of pronouns may require co-reference
  - ambiguous words may be informed by broader context
  - **consistent translation of repeated words**

# Conditioning on Broader Context

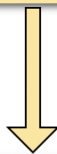


- Hierarchical attention
  - compute which previous sentences matter most
  - compute which words in these sentences matter most



# Conditioning on Broader Context

The shop is selling a nice table. <s> Jane is quite taken by it. <s> The table would match the chairs in her living room.



Der Laden verkauft einen schönen Tisch. <s> Er gefällt Jane sehr. <s> ...

- Concatenate all sentences together
  - document = very long sentence
  - special treatment for sentence boundaries
  - requires scaling of neural decoding implementation

questions?