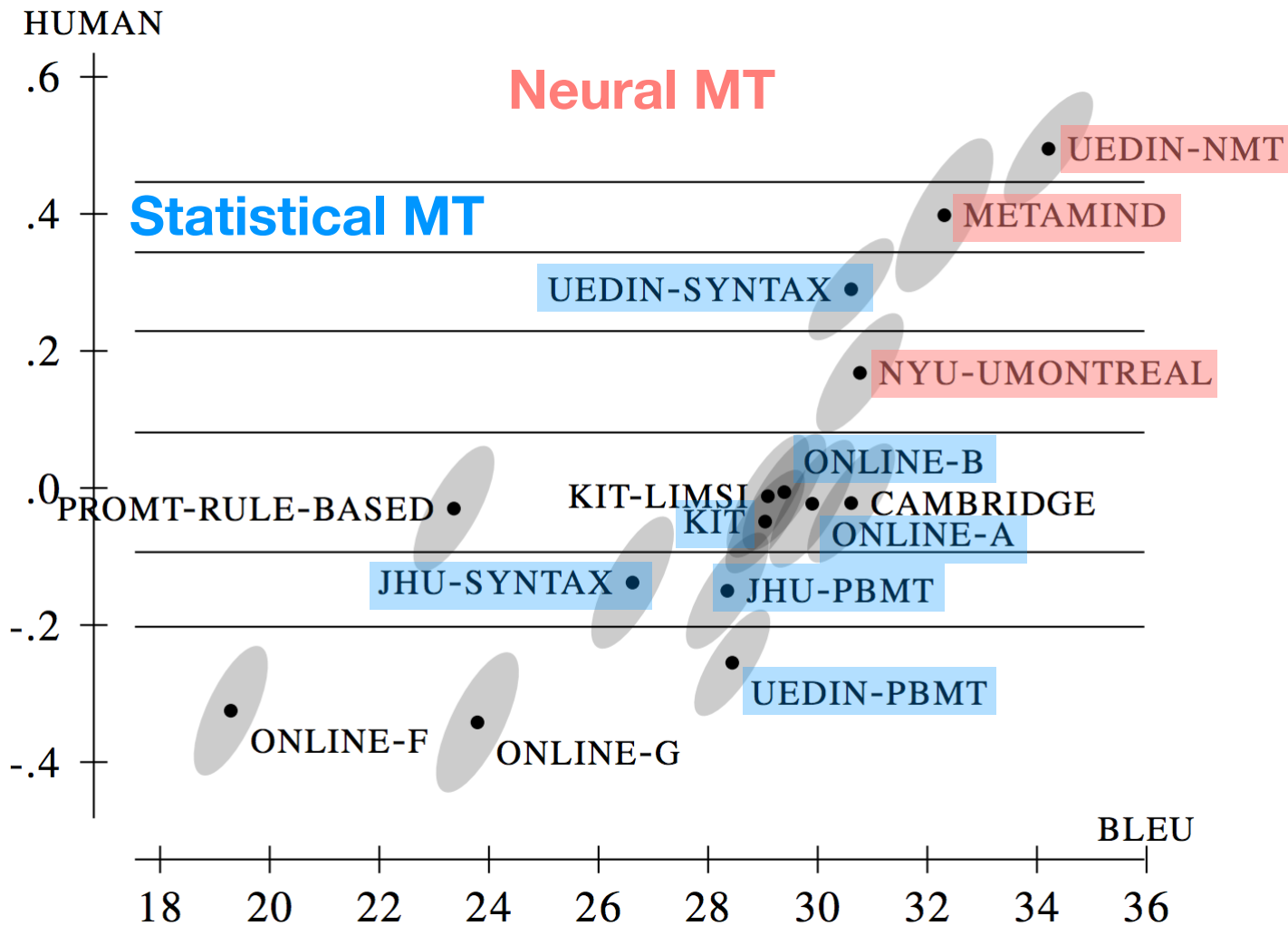

Current Challenges

Philipp Koehn

5 November 2020

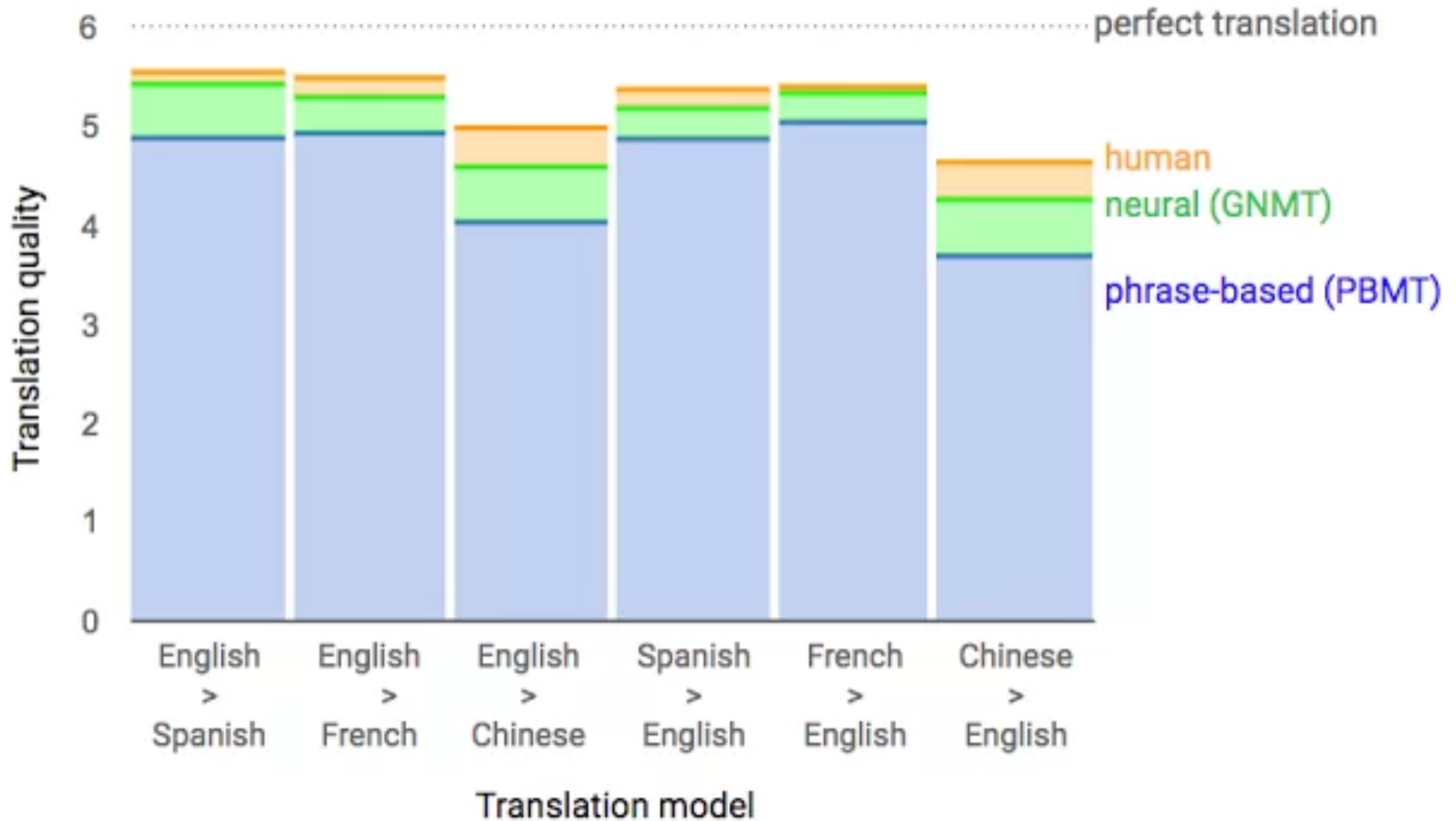


WMT 2016



(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.

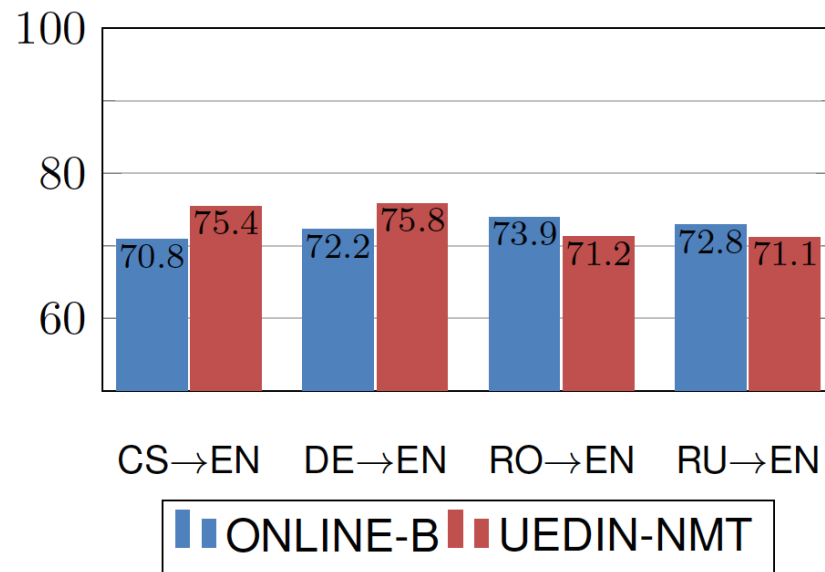
SDL Cracks Russian to English Neural Machine Translation

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

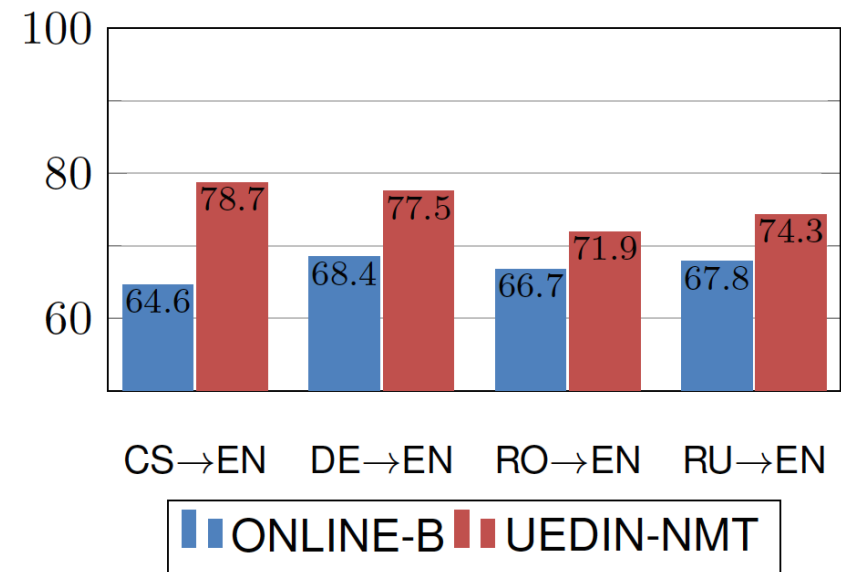
“90% of the system’s output labelled as perfect by professional Russian-English translators”

Just Better Fluency?

Adequacy +1%



Fluency +13%



(from: Sennrich and Haddow, 2017)

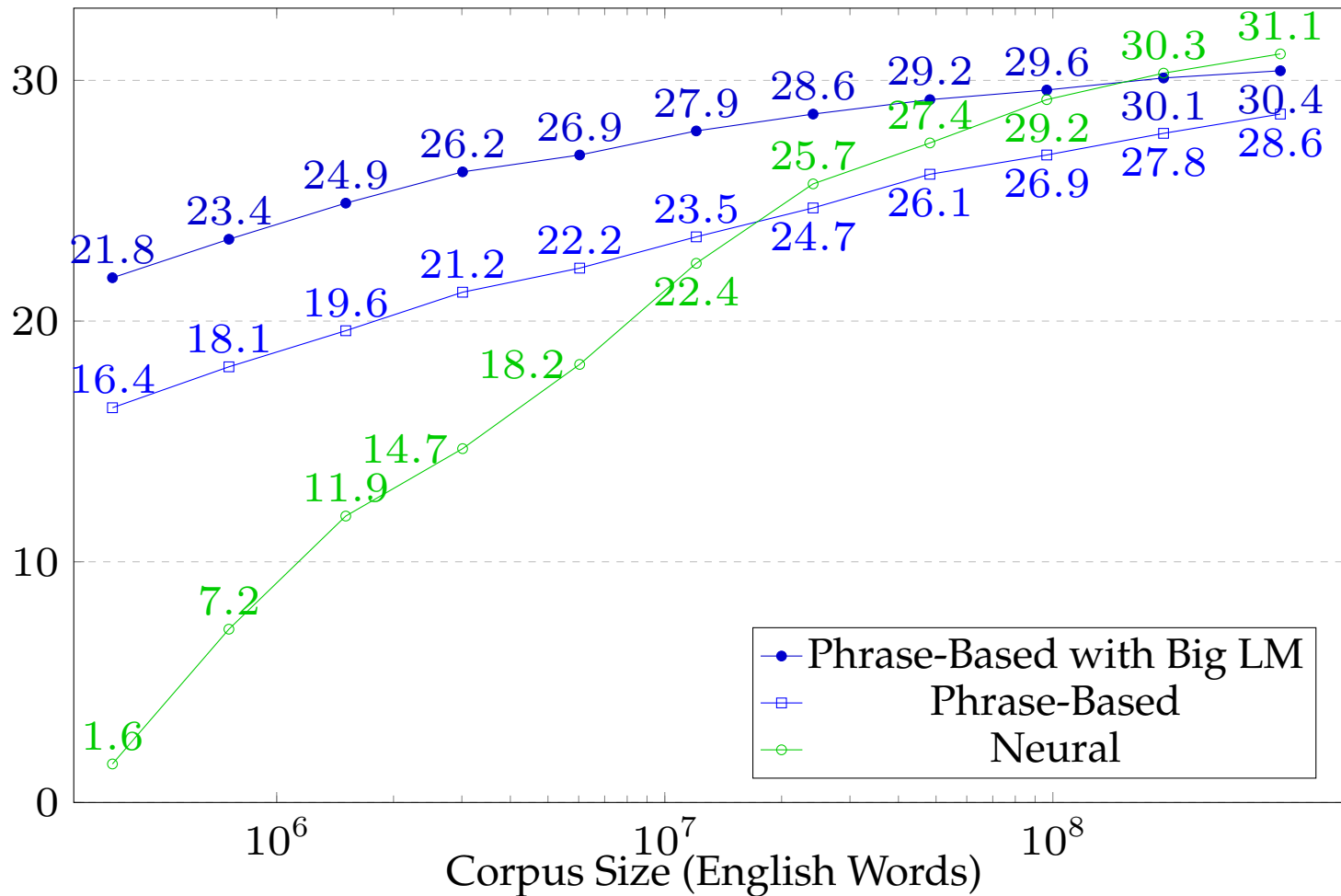
Challenges



- 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

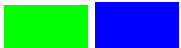
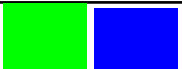
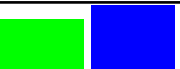



























Translation Examples



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
All Data	 30.5 32.8	 45.1 42.2	 35.3 44.7	 17.9 17.9	 26.4 20.8
Law	 31.1 34.4	 12.1 18.2	 3.5 6.9	 1.3 2.2	 2.8 6.0
Medical	 3.9 10.2	 39.4 43.5	 2.0 8.5	 0.6 2.0	 1.4 5.8
IT	 1.9 3.7	 6.5 5.3	 42.1 39.8	 1.8 1.6	 3.9 4.7
Koran	 0.4 1.8	 0.0 2.1	 0.0 2.3	 15.9 18.8	 1.0 5.5
Subtitles	 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

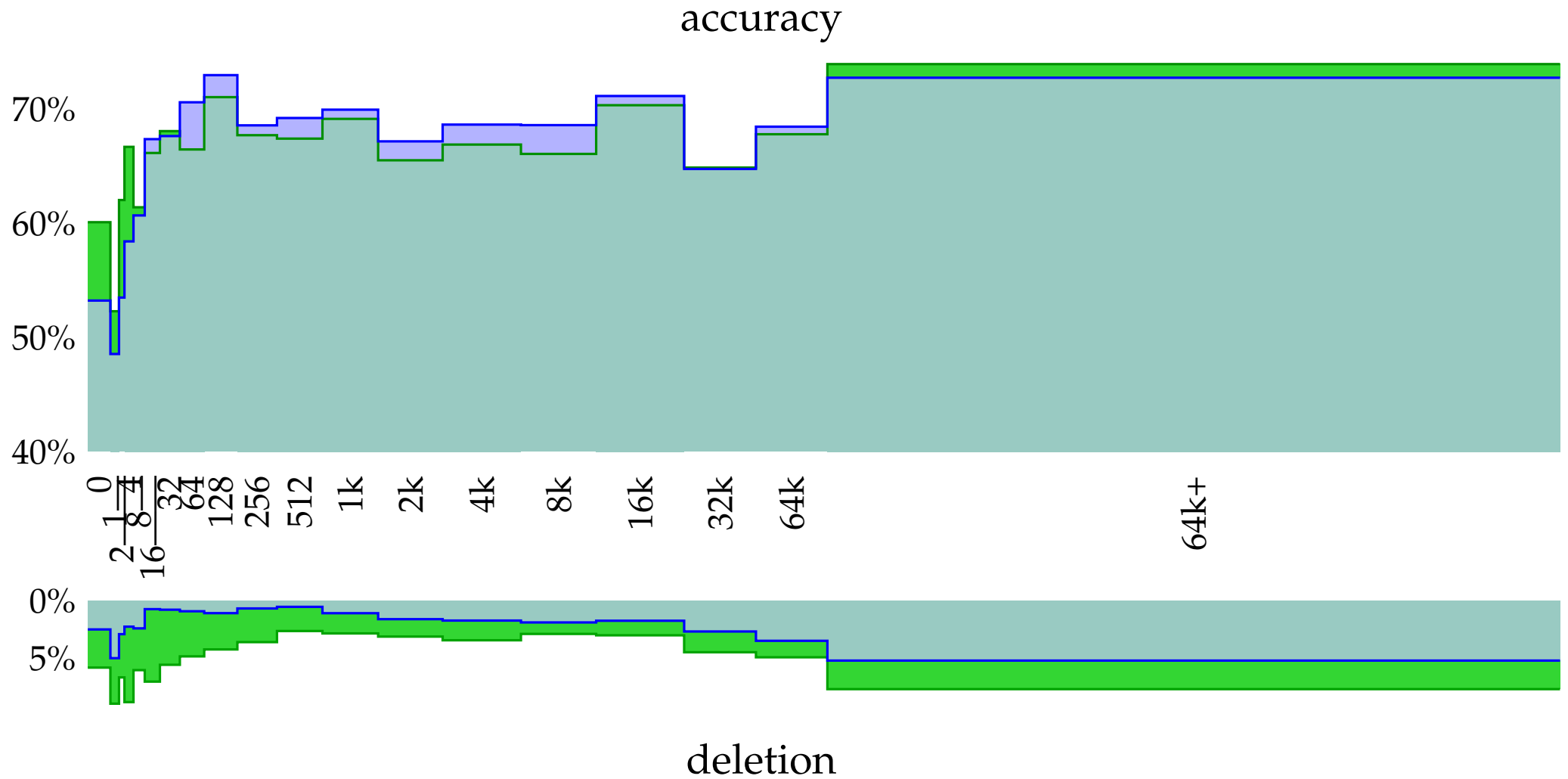
Rare Words

- More frequent in training → 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

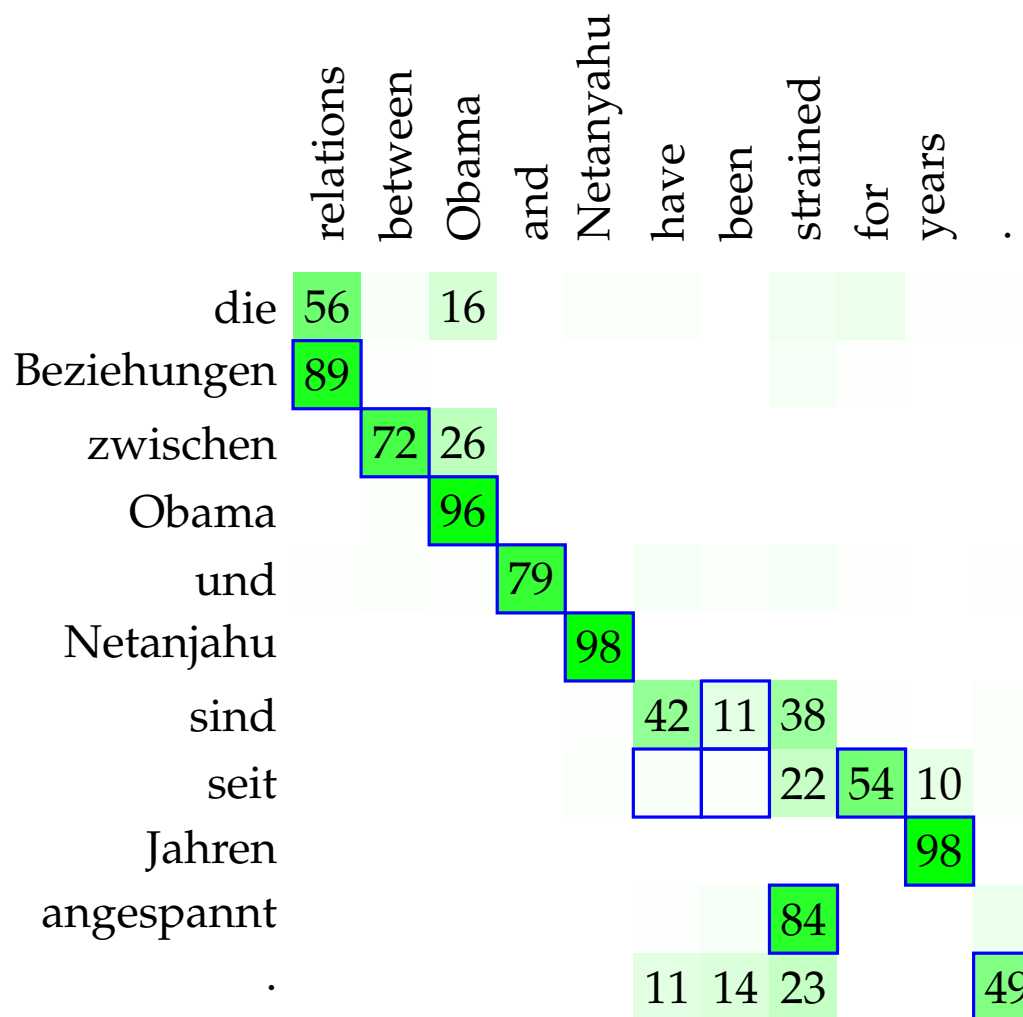
- Generate word alignment between input and output words
- Look up count of input word in training
- Link to output word via word alignment
- Check if it is also in the reference translation■
- A lot of tedious special cases
 - one-to-many alignment, only some output words in reference
 - input word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

Count vs. Accuracy

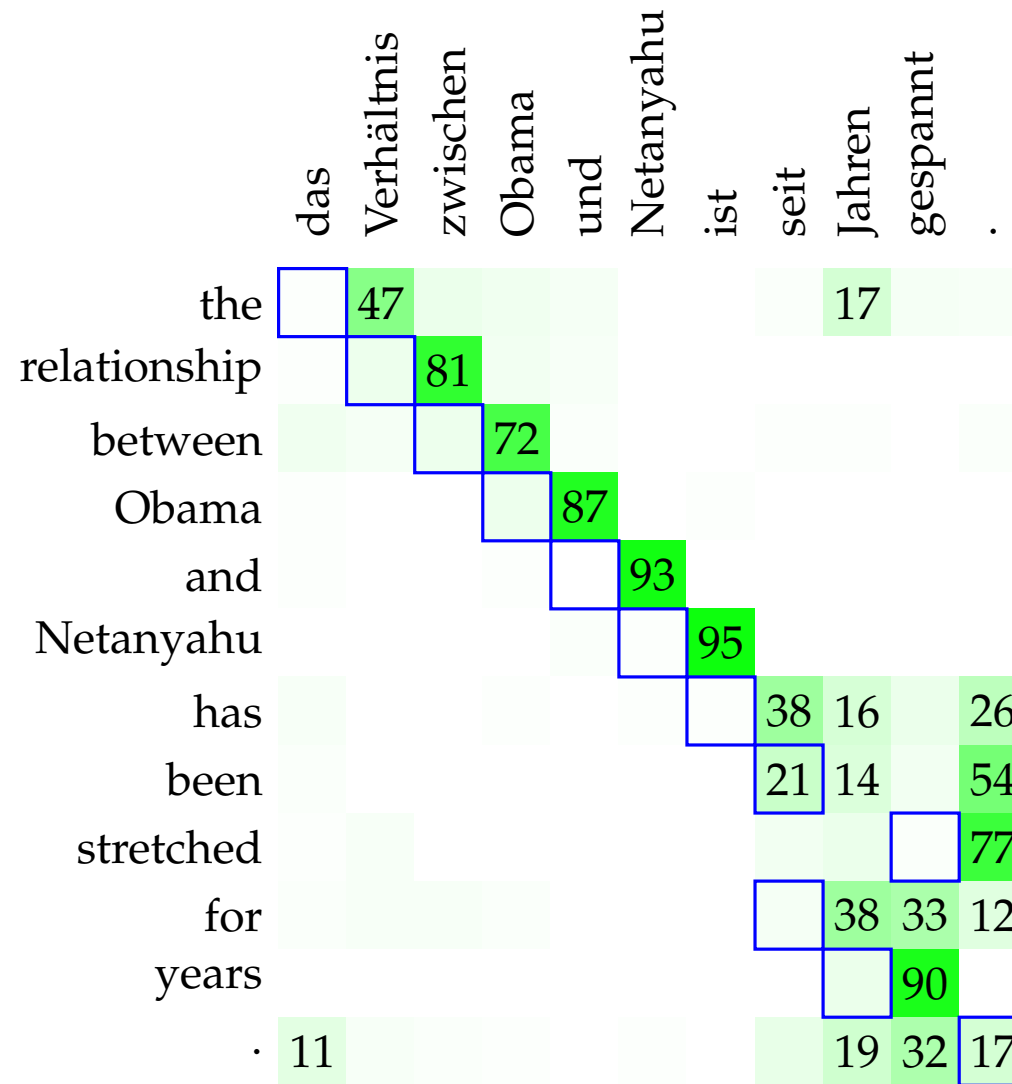


word alignment

Word Alignment

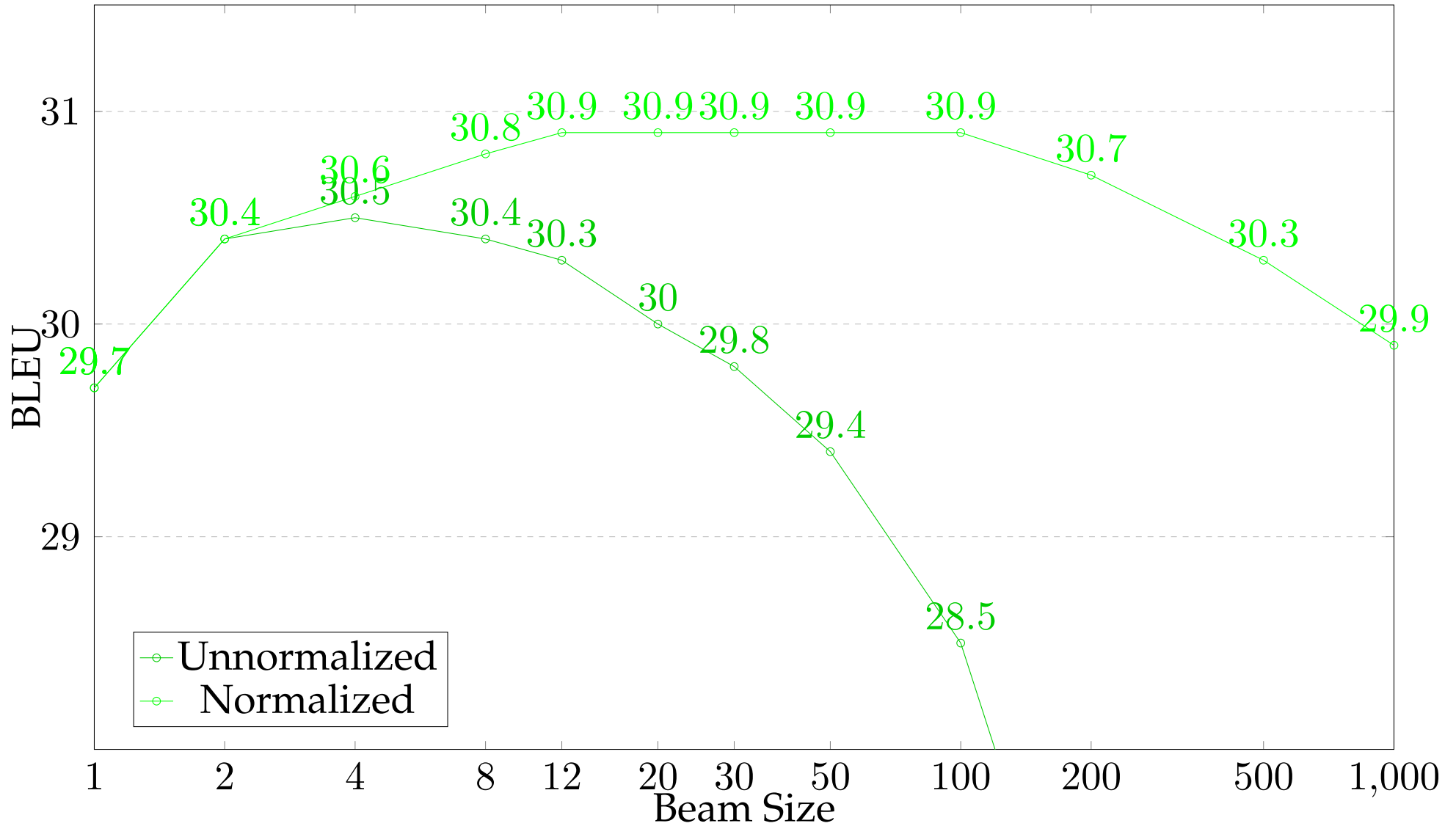


Word Alignment?



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	27.4	24.2	26.6	24.2	24.7	24.4	20.9	24.8	17.3	25.2
	+0.2	+0.2	-0.9	+0.2	-2.5	+0.4	-6.3	+0.8	-9.9	+1.2

- 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%
$\frac{24.0}{-0.0}$	$\frac{24.0}{-0.0}$	$\frac{23.9}{-0.1}$	$\frac{26.1}{-1.1}$ $\frac{23.9}{-0.1}$	$\frac{25.3}{-1.9}$ $\frac{23.4}{-0.6}$

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

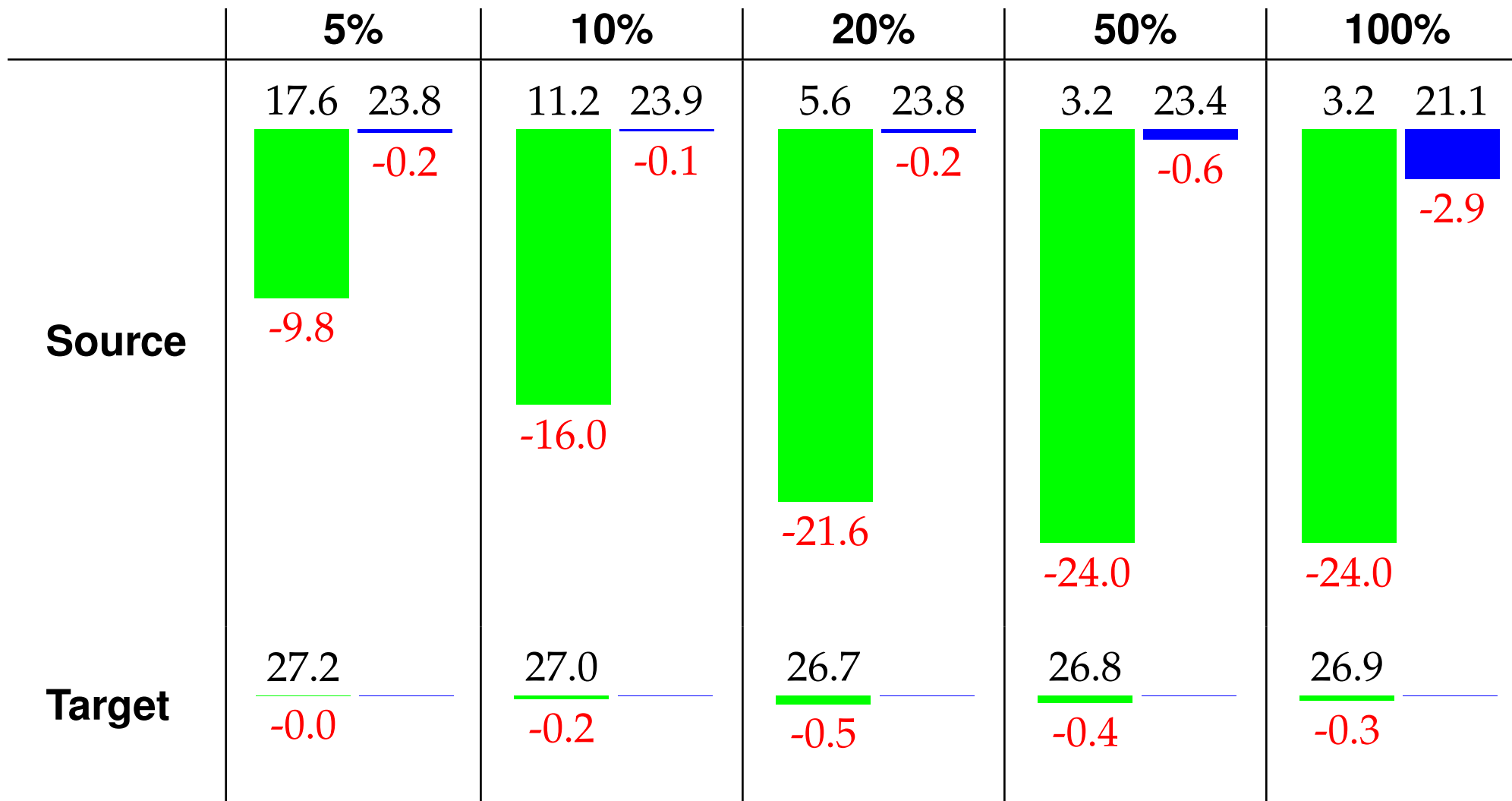
Misordered Words

- Artificial created by randomly shuffling words in each sentence

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

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

Untranslated Sentences



Wrong Language

	5%	10%	20%	50%	100%
fr source	<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
fr target	<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%
1-2 words	$\frac{27.1}{-0.1} \quad \frac{24.1}{+0.1}$	$\frac{26.5}{-0.7} \quad \frac{23.9}{-0.1}$	$\frac{26.7}{-0.5} \quad \frac{23.8}{-0.2}$	
1-5 words	$\frac{27.8}{+0.6} \quad \frac{24.2}{+0.2}$	$\frac{27.6}{+0.4} \quad \frac{24.5}{+0.5}$	$\frac{28.0}{+0.8} \quad \frac{24.5}{+0.5}$	$\frac{26.6}{-0.6} \quad \frac{24.2}{+0.2}$

- 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>`)

- 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



questions?