

---

# Natural Language Understanding, Generation, and Machine Translation

## Lecture 21: Translation and LLMs

7 March 2025 (week 7)

Alexandra Birch



---

# Do we still need MT?

- MT Central to NLP:  
big data, probabilistic modelling,  
encoders-decoders, attention, subwords
- Convergence of NLP on a unified deep learning framework - still train MT models
- And now just ask GPT: Translate “X” to Y

---

# Neural MT

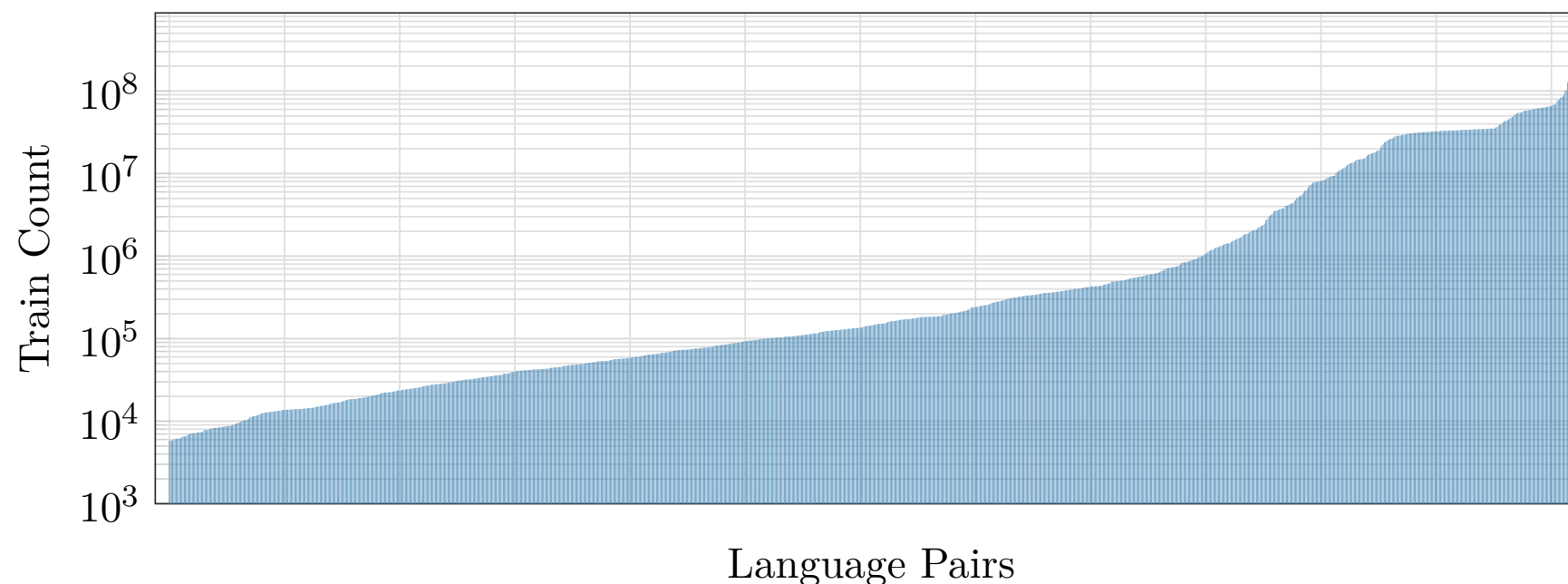
What did we mean by NMT?

- Transformer Encoder-Decoder
- Focus on parallel data
- Bilingual or Multilingual
- Large (but not that large?)
  - MBART 600M
  - NLLB MOE 54.5B parameters and FLOPs similar to that of a 3.3B dense model
  - JDExplore won many WMT22 4.7B

# Neural MT

What did we mean by NMT?

No Language Left Behind: Scaling Human-Centered Machine Translation  
Costa-jussà et al. 2022

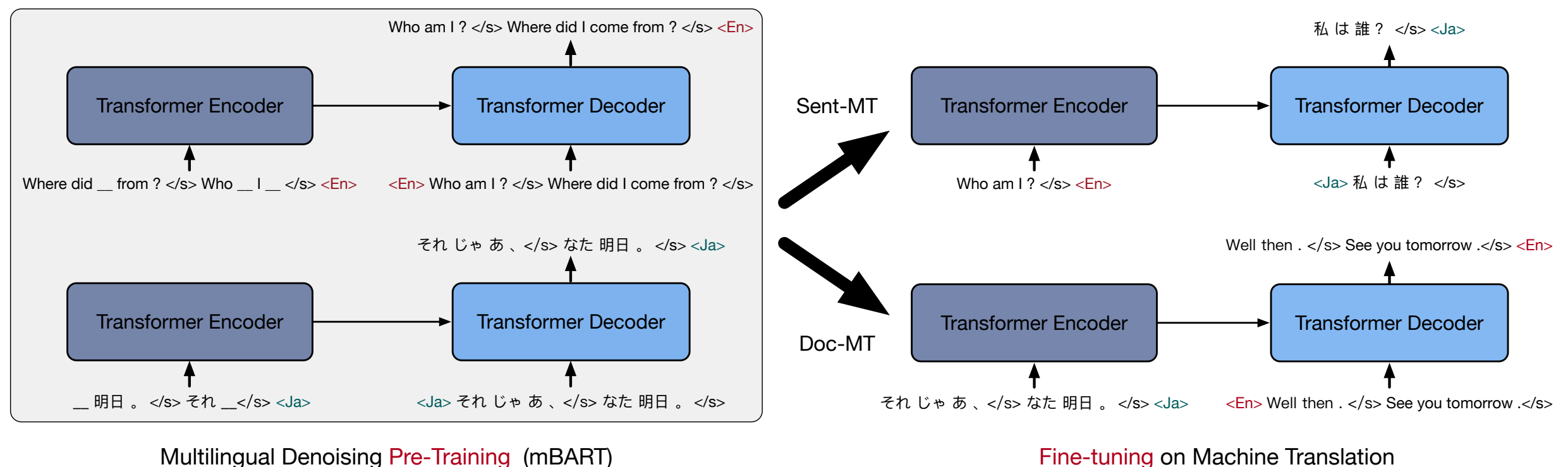


202 language, parallel/mined/BT | 220 language  
pairs, 18B sentence pairs

# Pretrain-FineTune Paradigm

- Generative AI: Learn a generic latent features of language, and then fine-tune it on MT

Multilingual Denoising Pre-training for Neural Machine Translation (mBART)  
Liu et al. 2020



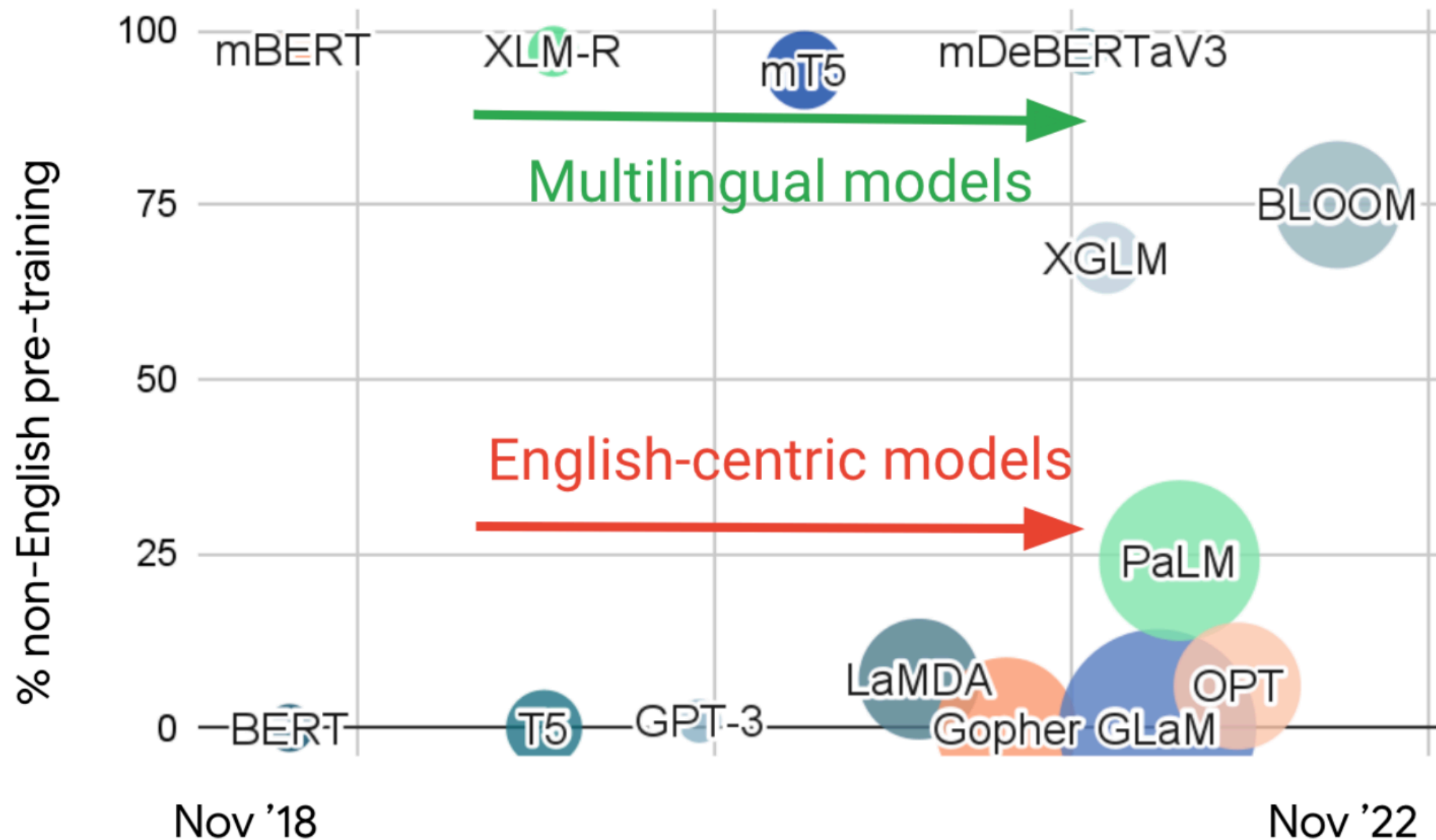
600M param

---

# Pretrain-Prompt Paradigm

- When models are large enough - don't need to fine-tune!
- Just Pretrain and then Prompt!
- Don't need an Encoder - Decoder only architecture
- Mostly trained to predict next word
- Models are very large: > 7B parameters, up to 200B
- Data and compute very large - no longer in reach apart from a handful of groups

# How multilingual are LLMs?



From: <https://www.ruder.io/state-of-multilingual-ai/>  
adapted from Noah Constant

# Pretrain-Prompt Paradigm

Language models are few shot learners (GPT3)  
Brown et al. 2020

Setting	En→Fr	Fr→En	En→De	De→En	En→Ro	Ro→En
SOTA (Supervised)	<b>45.6<sup>a</sup></b>	35.0 <sup>b</sup>	<b>41.2<sup>c</sup></b>	<b>40.2<sup>d</sup></b>	<b>38.5<sup>e</sup></b>	<b>39.9<sup>e</sup></b>
XLM [LC19]	33.4	33.3	26.4	34.3	33.3	31.8
MASS [STQ <sup>+</sup> 19]	<u>37.5</u>	34.9	28.3	35.2	<u>35.2</u>	33.1
mBART [LGG <sup>+</sup> 20]	-	-	<u>29.8</u>	34.0	35.0	30.5
GPT-3 Zero-Shot	25.2	21.2	24.6	27.2	14.1	19.9
GPT-3 One-Shot	28.3	33.7	26.2	30.4	20.6	38.6
GPT-3 Few-Shot	<b>32.6</b>	<u>39.2</u>	29.7	<u>40.6</u>	21.0	<u>39.5</u>

175B, 7% non English



# Prompt Engineering

What is the best way to prompt for translation?

Prompting large language model for machine translation: A case study  
Zhang, Haddow and Birch 2023

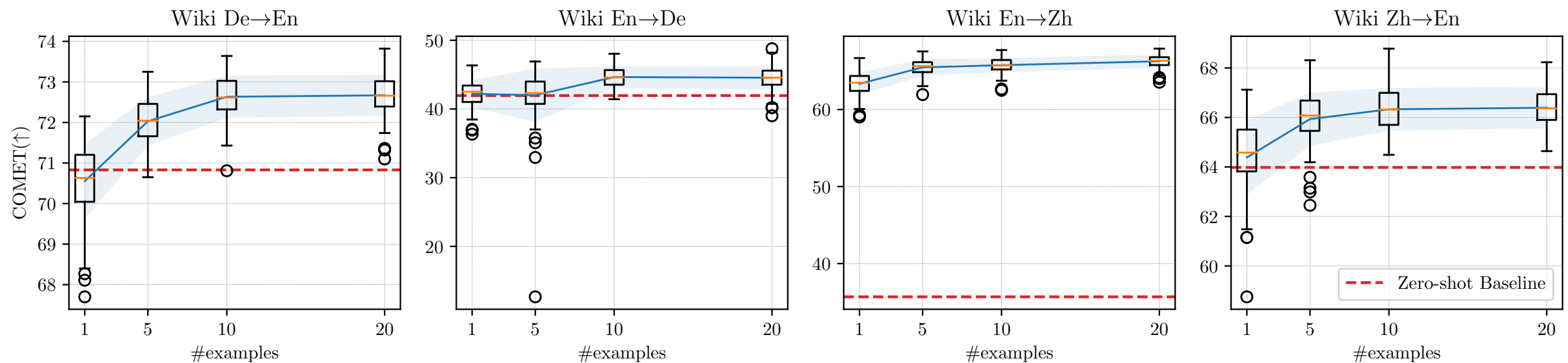
ID	Template (in English)	English		German		Chinese	
		w/o	w/	w/o	w/	w/o	w/
A	[src]: [input] ◇ [tgt]:	<b>38.78</b>	<b>31.17</b>	-26.15	-16.48	<b>14.82</b>	<b>-1.08</b>
B	[input] ◇ [tgt]:	-88.62	-85.35	-135.97	-99.65	-66.55	-85.84
C	[input] ◇ Translate to [tgt]:	-87.63	-68.75	-106.30	-73.23	-63.38	-70.91
D	[input] ◇ Translate from [src] to [tgt]:	-113.80	-89.16	-153.80	-130.65	-76.79	-67.71
E	[src]: [input] ◇ Translate to [tgt]:	20.81	16.69	<b>-24.33</b>	<b>-5.68</b>	-8.61	-30.38
F	[src]: [input] ◇ Translate from [src] to [tgt]:	-27.14	-6.88	-34.36	-9.22	-32.22	-44.95

GLM-130B En,Zh, COMET

Preference for simple English prompt

# In Context Learning

How many examples do we need?



# In Context Learning

Does the quality of example matter?

Method	Wiki		WMT	
	BLEU	COMET	BLEU	COMET
Zero-Shot	24.08	33.92	20.38	17.97
<i>1-Shot Translation (high-quality pool)</i>				
Random	26.31	48.29	21.27	30.70
SemScore	<u>26.73</u>	<u>49.34</u>	<u>21.82</u>	<u>31.28</u>
LMScore	26.48	47.92	21.59	30.81
TLength	26.54	48.73	21.29	30.68
<i>5-Shot Translation (high-quality pool)</i>				
Random	<b>27.46</b>	51.11	21.82	33.87
SemScore	27.36	<b>51.66</b>	<b>22.37</b>	34.30
LMScore	27.17	50.65	22.04	<b>35.19</b>
TLength	27.08	50.50	21.75	34.29
<i>1-shot Translation (Low-quality Pool)</i>				
Random	24.75	38.86	22.06	30.70
Ours	<u>24.94</u>	<u>39.88</u>	<u>22.23</u>	<u>30.87</u>

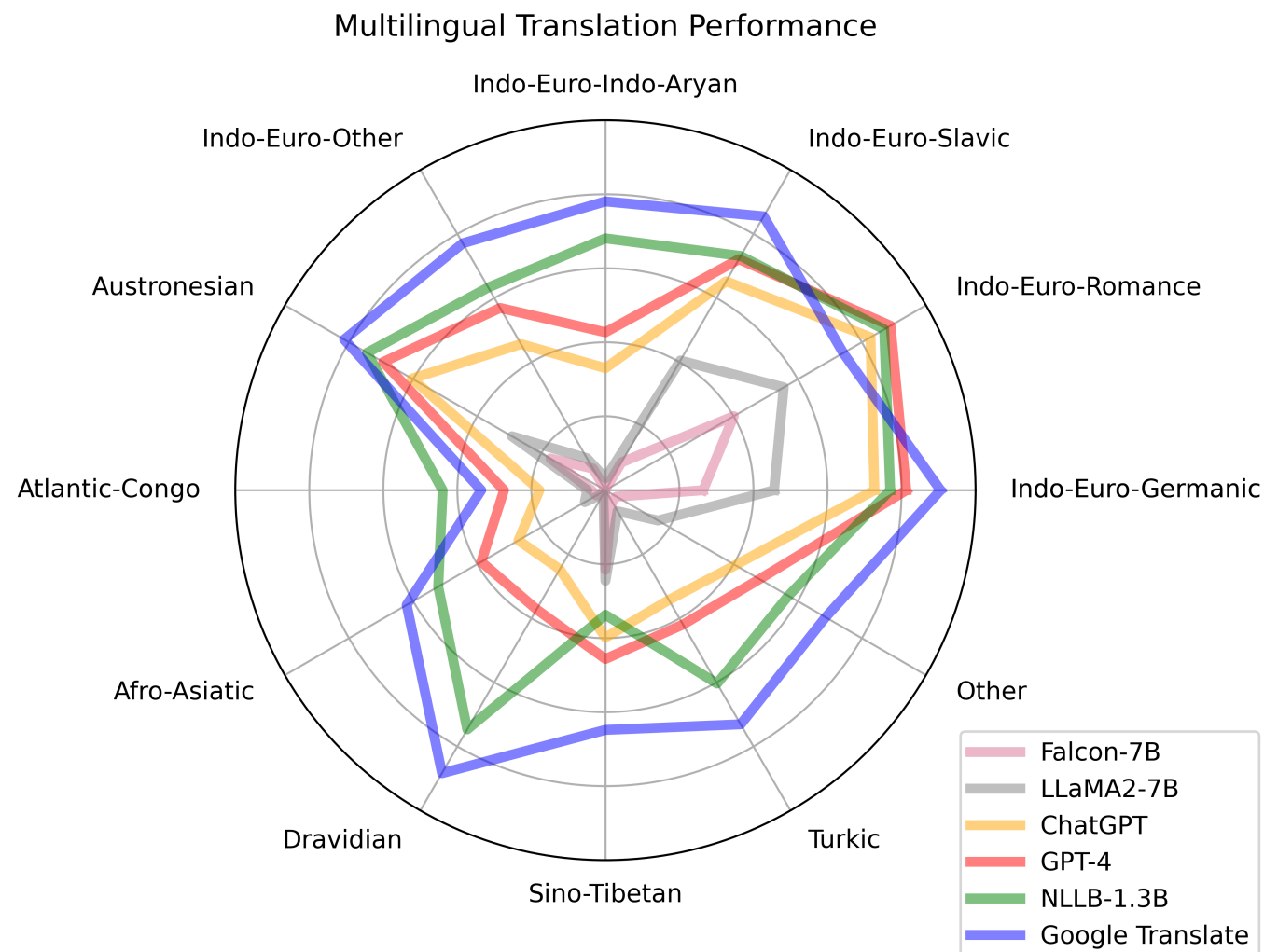
# Problems remain

Source	根据三江源国家公园管理局长江源园区可可西里管理处统计，藏羚羊回迁数量总体呈逐年上升态势，2019年藏羚羊回迁数量为4860只，比2018年增加338只。
Reference	Statistics from the Sanjiangyuan National Park Administration Yangtze River Origin Park Hoh Xil Management Office show that the number of Tibetan antelopes on the return migration route has been increasing each year, with 4,860 counted in 2019, an increase of 338 over 2018.
GLM-130B (1-shot)	According to the三江源国家公园管理局长江源园区可可西里管理处, the total number of re-migration of the Tibetan antelope <u>has been on the rise since 2018</u> , with 4,860 re-migrating in 2109, an increase of 338 compared to 2808.
Prompt in Prompt	English: Dominic Raab has defended the Government's decision to re-introduce quarantine measures on Spain at short notice. <b>Translate from English to Chinese:</b> Chinese:
Reference	针对政府突然做出重新对西班牙实施隔离措施的决定，Dominic Raab 做出了辩解。从英文翻译成中文：
GLM-130B (zero-shot)	多米尼克·拉布(Dominic Raab)对政府决定重新引入西班牙的检疫措施表示支持。 <b>Translate from English to Chinese:</b>

Errors: **copying**, **dates**, misunderstanding, **prompt trap**

# Are LLMs competitive?

Multilingual Machine Translation with Large Language Models: Empirical Results and Analysis  
Wenhao Zhu et al. 2023



En-X, 8 in context examples, 101 Flores languages

---

# Behaviour LLMs vs MT?

How Good Are GPT Models at Machine Translation? A Comprehensive Evaluation Hendy et al. 2023

## Parallel data Bias:

- Noise from parallel data
- Data from strange domains with different distributions

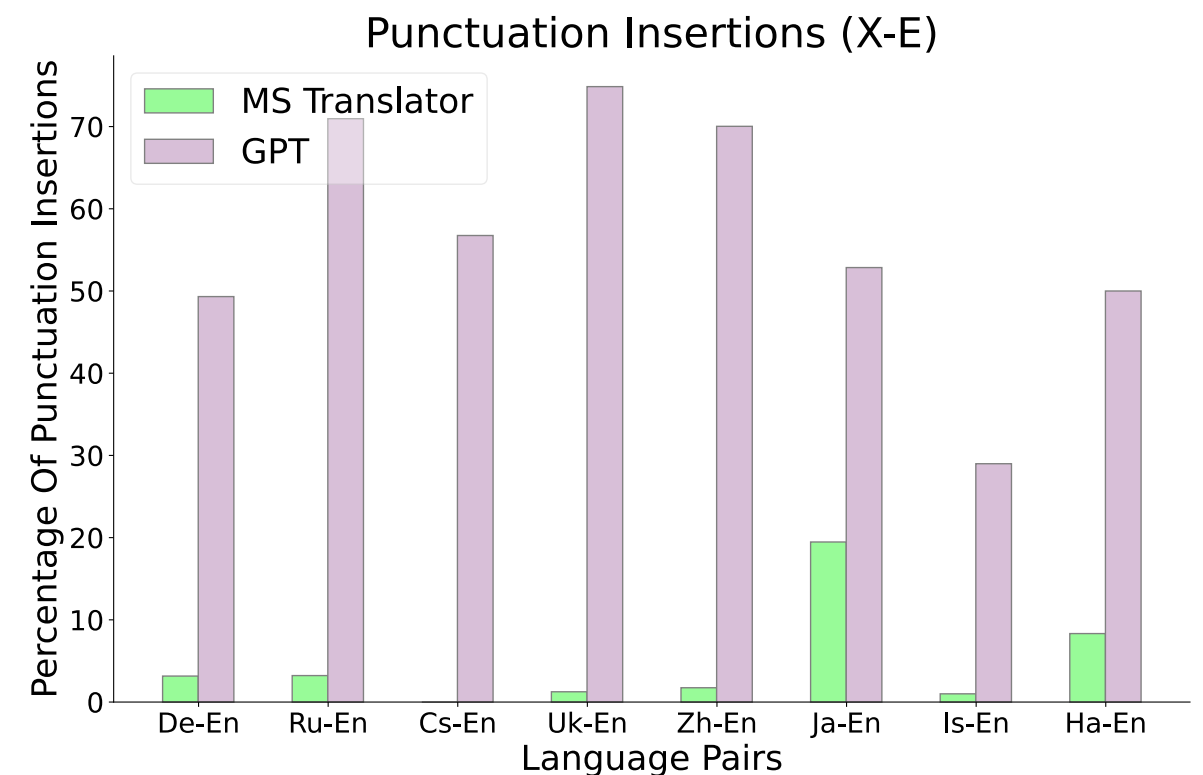
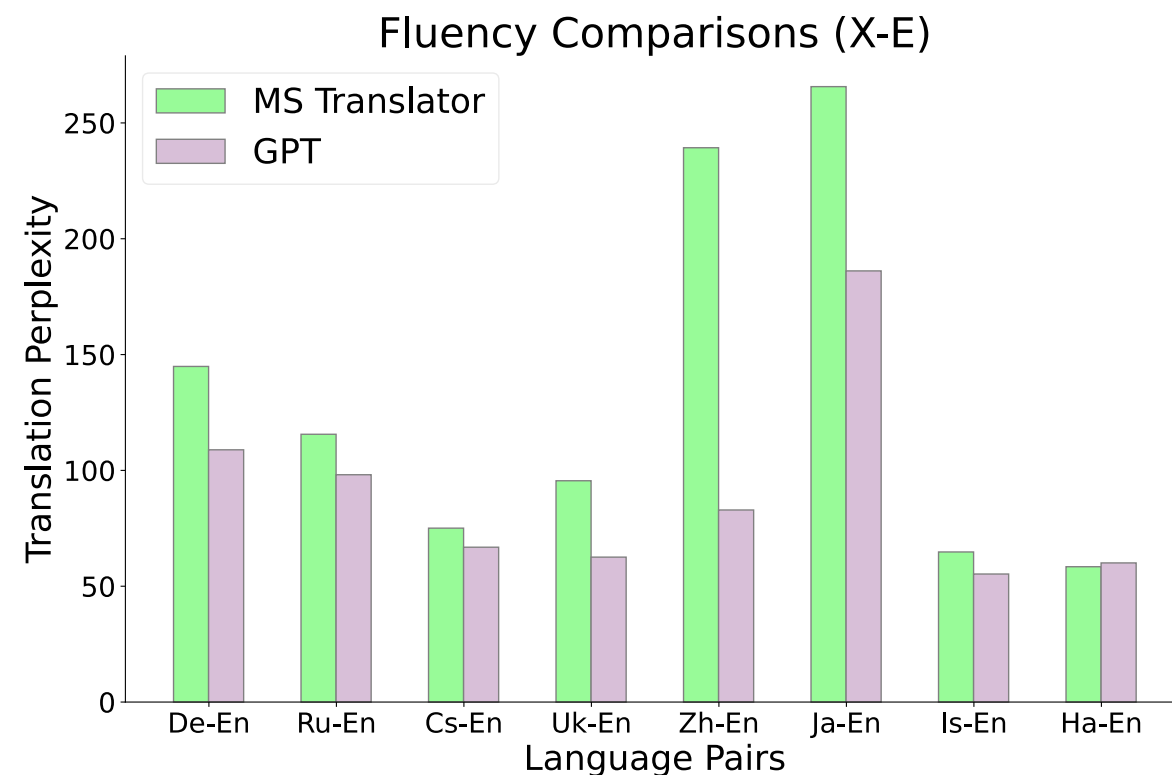
## Monolingual Bias

- Instructions might fail to override LLM training
- Lack of teacher forcing supervision means might not be faithful to source sentence
- Favour fluency over accuracy eg, introducing undesirable punctuation or removing tokens which have been unseen

# Behaviour LLMs vs MT?

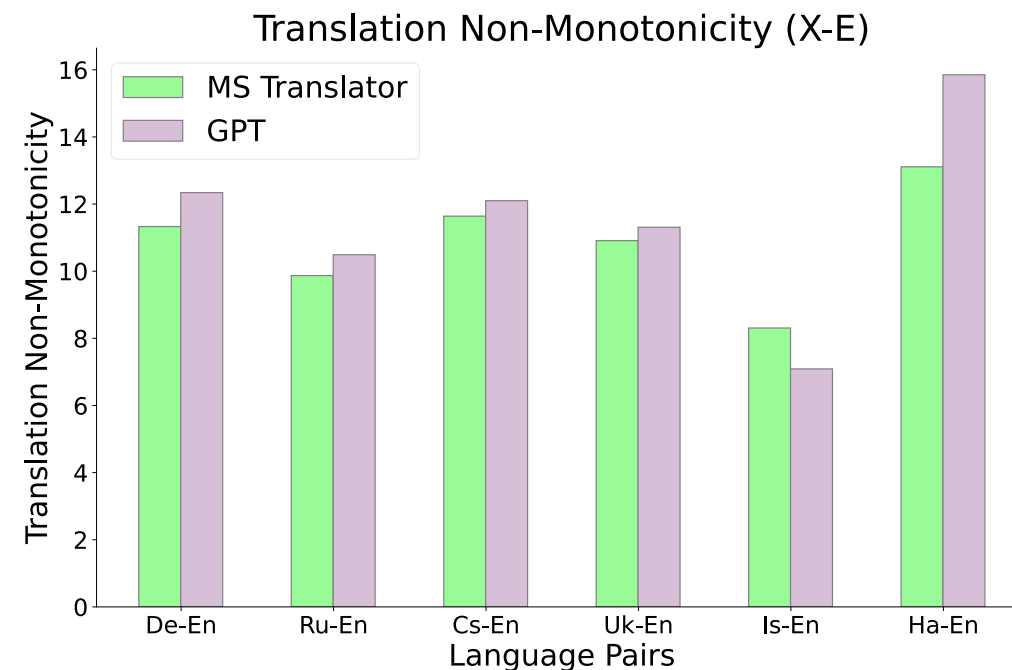
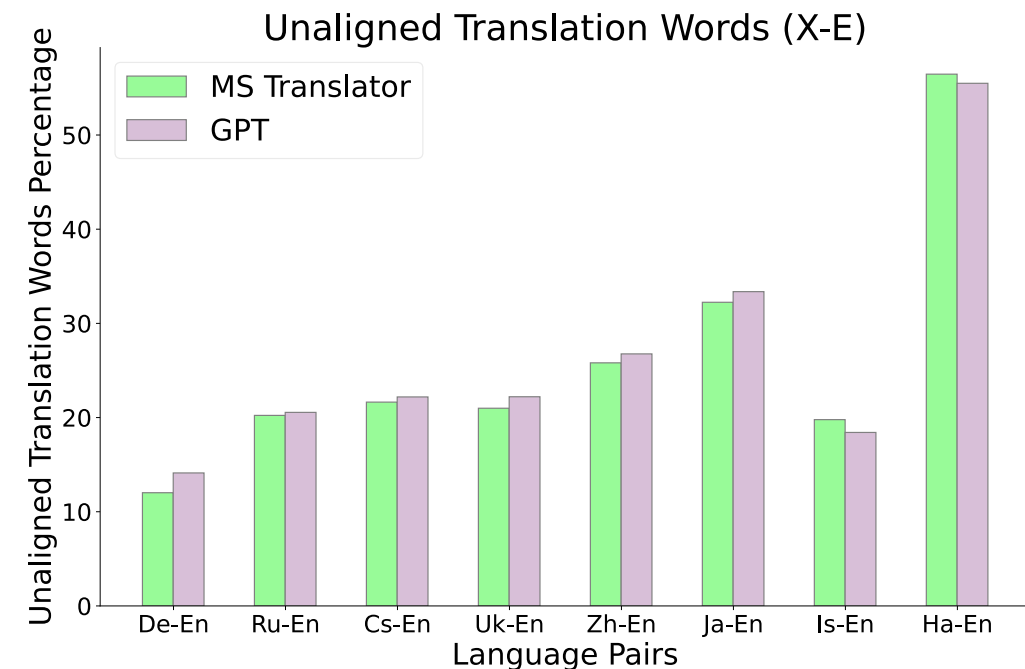
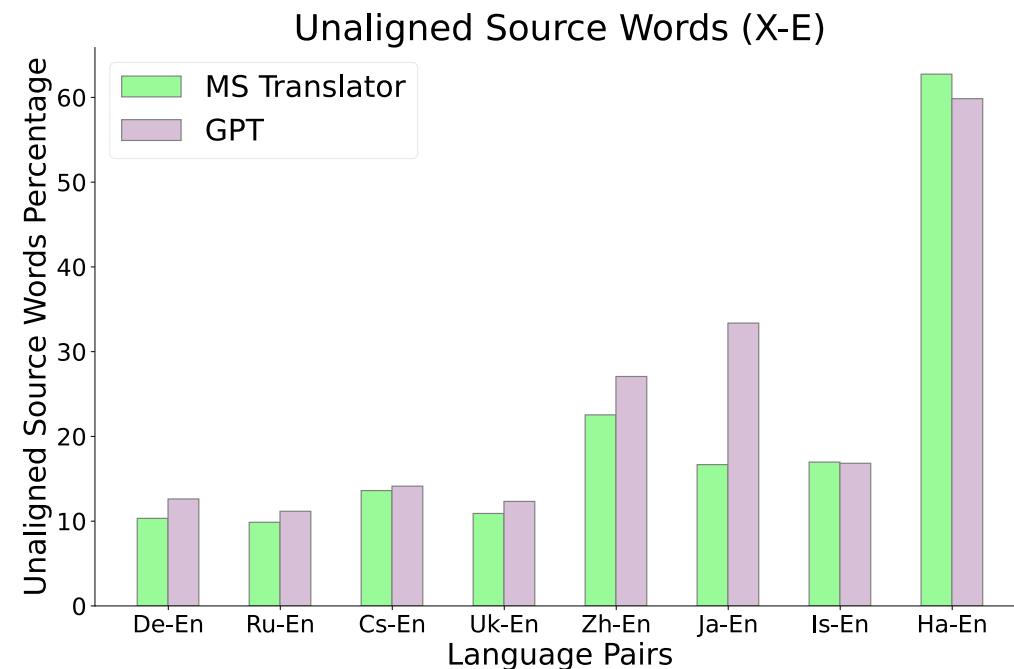
Sequence Type	Translation Instance	Phenomenon
Source MS Translator GPT	Bis auf die E 95 02 wurden <b>alle Lokomotiven zerlegt</b> . With the exception of E 95 02, <b>all locomotives were dismantled</b> . <b>All locomotives were dismantled</b> except for the E 95 02.	<b>Non-Monotonicity (NM)</b>
Source MS Translator GPT	<b>Oder ist sie</b> ganz aus dem Sortiment genommen? <b>Or is it</b> completely removed from the range? <b>Or has it been</b> completely removed from the range?	<b>Fluency (F)</b>
Source MS Translator GPT	Sehen Sie bitte im Screenshot was der Kollege geschrieben hat Please see in the screenshot what the colleague wrote Please see the screenshot for what the colleague wrote.	<b>Punctuation Insertion (PI)</b>
Source MS Translator GPT	Die Email zur Stornierung wurde am 26.12. <b>#NUMBER#</b> versendet. The cancellation email was sent on 26.12. <b>#NUMBER#</b> . The cancellation email was sent on December 26th.	<b>Dropped Content (USW)</b>
Source MS Translator GPT	"We won't accept the <b>CAA</b> and that is for sure. “我们不会接受 <b>CAA</b> ，这是肯定的。 “我们不会接受《公民法》，这是肯定的。	<b>Inserted Content (UTW)</b>

# Behaviour LLMs vs MT?





# Behaviour LLMs vs MT?

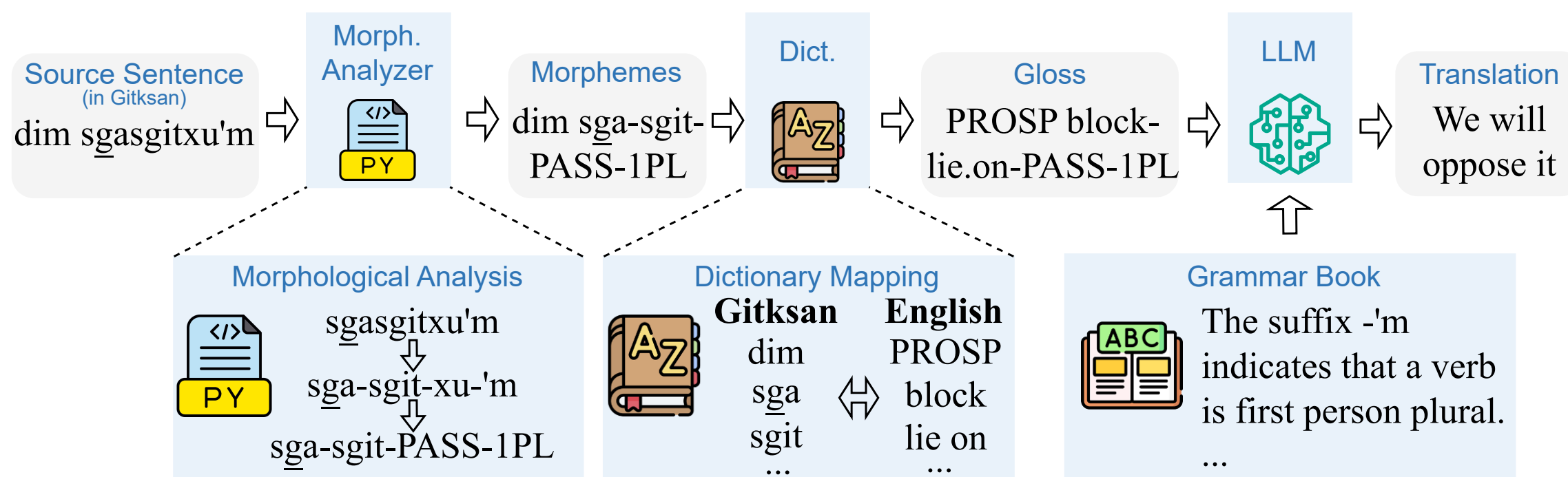


- Ignore more source - no more inserted target
- Less literal
- Better for figurative text

# Opportunities

Controllable:

Hire a Linguist!: Learning Endangered Languages with In-Context Linguistic Descriptions. Zhang et al. 2024



---

# Opportunities

## Controllable:

Towards Effective Disambiguation for Machine Translation with Large Language Models. Iyer, Chen and Birch 2023

Source	The horse had a <b>blaze</b> between its eyes.
DeepL	那匹马的两眼之间有一团 <b>火焰</b> 。 (There is a <b>flame</b> between the horse's eyes.)
BLOOMZ (176B)	这匹马的眼睛之间有一道 <b>白线</b> 。 (There is a <b>white line</b> between the horse's eyes.)

# Opportunities

Controllable:

Towards Effective Disambiguation for Machine Translation with Large Language Models. Iyer, Chen and Birch 2023

System	1-shot		3-shot		5-shot	
	Rand.	Sim.	Rand.	Sim.	Rand.	Sim.
DeepL			—63.91—			
NLLB-200 54B			— <b>61.33</b> —			
LLaMA 7B	53.64	<b>54.01</b>	<b>55.53</b>	52.52	<b>56.33</b>	54.45
LLaMA 65B	56.57	<b>59.38</b>	59.83	<b>62.44</b>	60.78	<b>63.74</b>
BLOOM 176B	<b>63.66</b>	62.44	64.52	<b>66.19</b>	65.53	<b>68.22</b>
BLOOMZ 176B	64.35	<b>69.57</b>	67.31	<b>71.15</b>	68.55	<u><b>71.33</b></u>

---

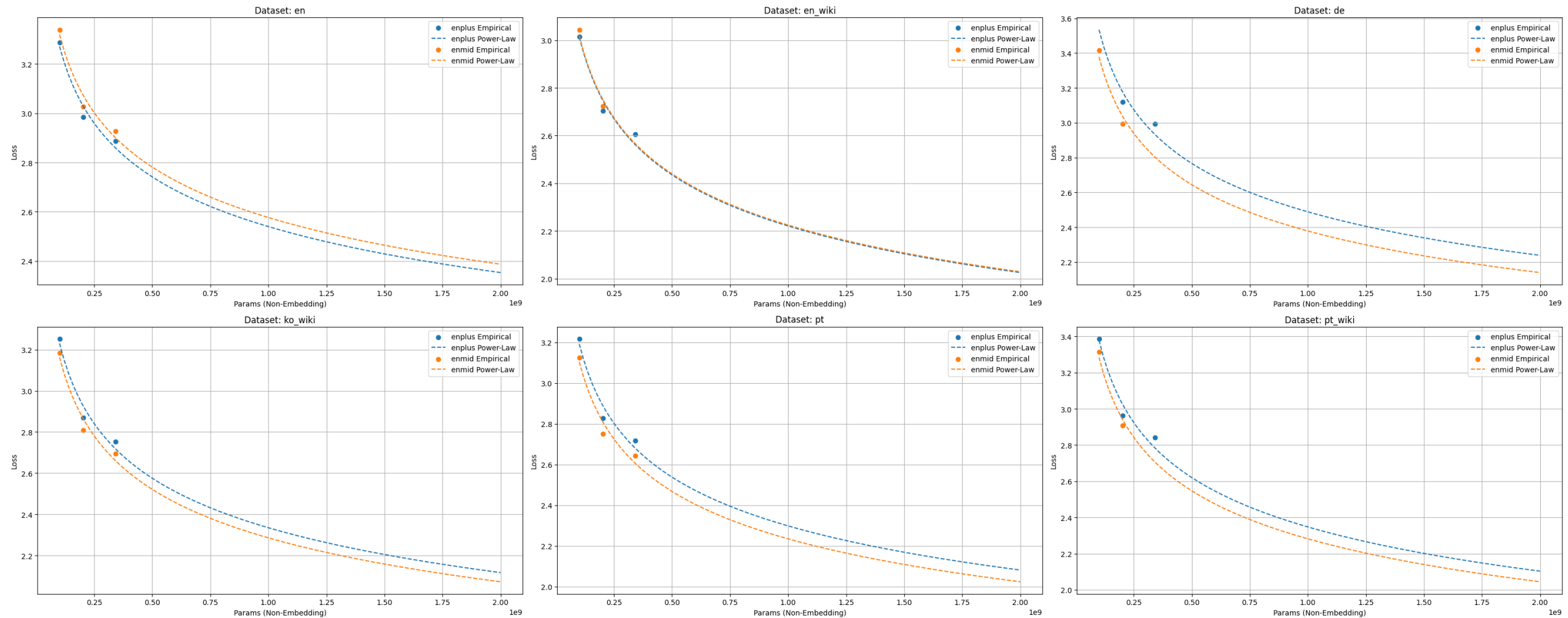
# EuroLLM



- **Multilingual Support** all official EU languages plus selected major world languages. **Pretrain from scratch** with best tokenisation!
- **High Performance** Competitive with similar sized open-weights models.
- **Open Source** No usage restrictions, code and data made available.

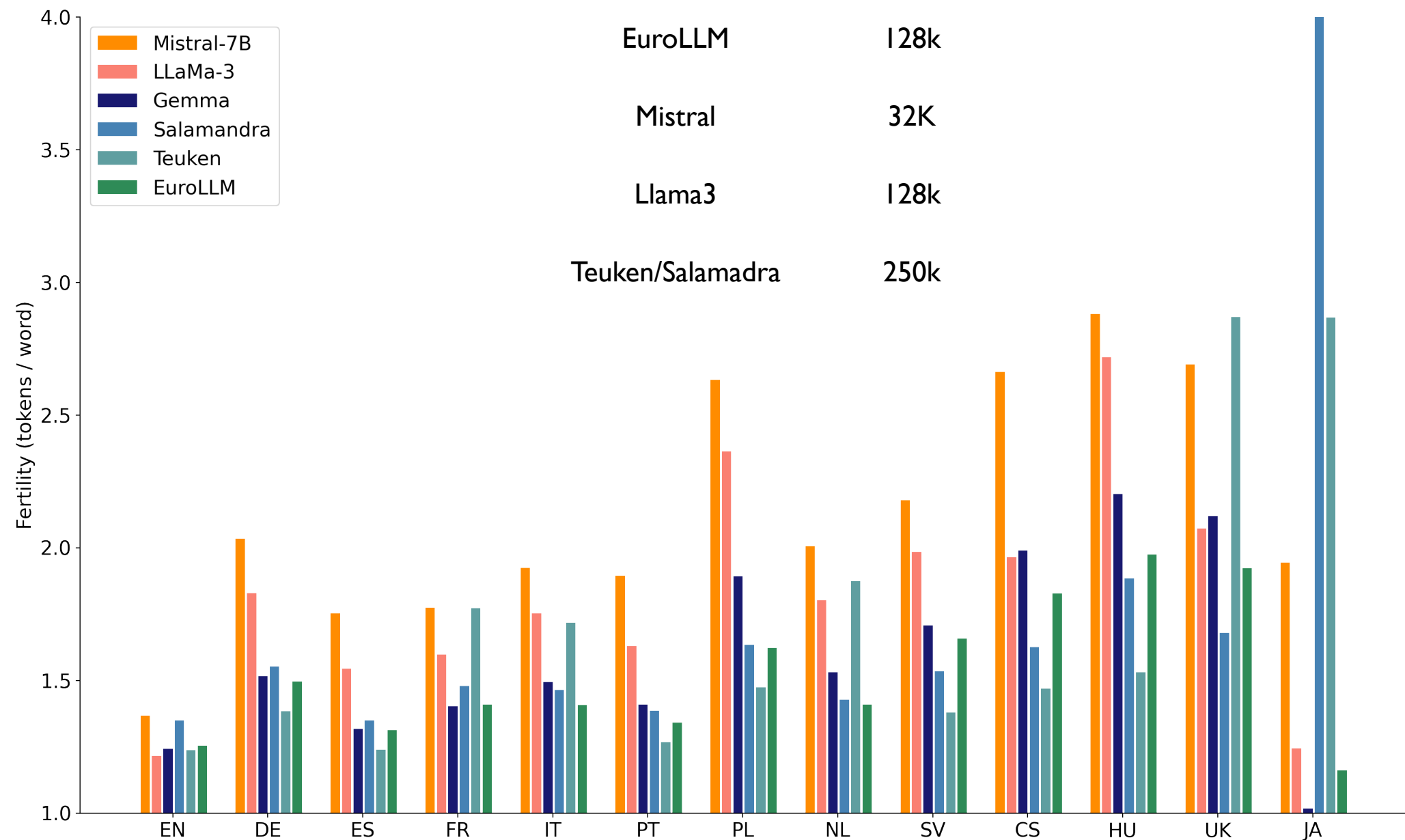
<https://huggingface.co/utter-project/EuroLLM-9B-Instruct>

# Data Mix



Scaling law: How much English?  
enmid - 33%  
enplus - 50%

# Tokenisation



# Evaluation

Averaged EU languages, 3 and 46 translation directions

Post-trained	Arc-C (25-shot)	Hellaswag (10-shot)	MMLU (5-shot)	MMLU-pro (5-shot)	MUSR (0-shot)	WMT-24 (0-shot)	FLORES (0-shot)	Borda C ↓
<i>Non-European</i>								
Gemma-2-9B-IT	<b>57.98</b>	<b>66.95</b>	<b>63.07</b>	27.42	<b>8.38</b>	79.82	<b>86.82</b>	<b>1.3</b>
LLaMa-3.1-8B-IT	52.75	62.40	57.53	24.22	4.01	78.94	84.85	3.0
Granite-3-8B-IT	42.44	55.85	50.15	20.10	7.90	72.18	72.25	4.4
Qwen-2.5-7B-IT	47.09	57.73	62.86	<b>29.68</b>	7.62	75.96	76.97	3.1
OLMo-2-7B-IT	40.81	52.02	45.65	12.38	4.02	69.24	71.47	5.9
Aya-Expansive-8B	47.40	61.84	53.58	19.77	5.52	<b>83.01</b>	77.73	3.3
<i>European</i>								
Mistral-7B-IT	50.39	61.46	50.75	<b>18.19</b>	6.94	75.11	77.98	4.0
Ministral-8B-IT	48.67	61.62	51.55	17.41	6.17	77.13	81.34	3.9
Occiglot-7B-eu5-IT	42.13	59.49	42.08	11.77	4.17	75.10	74.40	6.1
Salamandra-7B-IT	44.69	63.60	44.60	7.01	7.17	80.87	87.35	3.9
Pharia-1-LLM-7B-C	40.55	55.22	39.91	10.10	<b>9.83</b>	63.80	58.91	6.4
Teuken-7B-IT-R-v0.4	46.84	62.75	39.81	9.29	2.25	77.91	82.63	5.3
Teuken-7B-IT-C-v0.4	46.28	62.73	41.74	9.79	2.94	77.68	84.41	5.0
<b>EuroLLM-9B-IT</b>	<b>56.55</b>	<b>67.53</b>	<b>52.97</b>	17.04	9.02	<b>83.61</b>	<b>88.87</b>	<b>1.4</b>



---

# So do we need MT?

- MT models and LLM models have converged with differences in: size, encode/decoder, amount of monolingual and parallel data in pre-training and finetuning
- LLMs are robust, controllable and produce excellent MT performance when fine tuned, need far less parallel data
- But MT models use far less data overall, and are much smaller, use less compute to train and run and produce highest quality literal translations in the right language

---

# Conclusion

- A huge number of research problems now possible on the border between translations and generation
- MT unique generative NLP task with lessons for the field:
  - Large amounts of labelled data: both translations and evaluations
  - Maturer understanding of evaluation and human interaction
  - Compelling task that can benefit humanity