

# Whisper

Robust Speech Recognition via Large-Scale Weak  
Supervision

Peter Sullivan  
DL-NLP 10-13-22

# Motivation + High Level Overview

- DL models are 'brittle' → train on a massively diverse data
- How to get massive data? → internet video with existing captions
  - Assume OK quality on human transcriptions, detect and remove ASR captions
- *Weakly supervised* multitask training on
  - ASR (English audio + English captions)
  - Speech Translation (English audio + English captions)
  - Non-English ASR (English audio + *same language* captions)
  - Voice activity detection
- Use prefix tokens to indicate which task

## English transcription



"Ask not what your country can do for ..."



Ask not what your country can do for ...

## Any-to-English speech translation



"El rápido zorro marrón salta sobre ..."



The quick brown fox jumps over ...

## Non-English transcription



"언덕 위에 올라 내려다보면 너무나 넓고 넓은 ..."



언덕 위에 올라 내려다보면 너무나 넓고 넓은 ...

## No speech



(background music playing)



∅

# Data

*“audio that is paired with transcripts on the Internet”*

680,000 Hours Total

117,000 Non-English (covering 96 different languages)

125,000 X  $\rightarrow$  English translated data

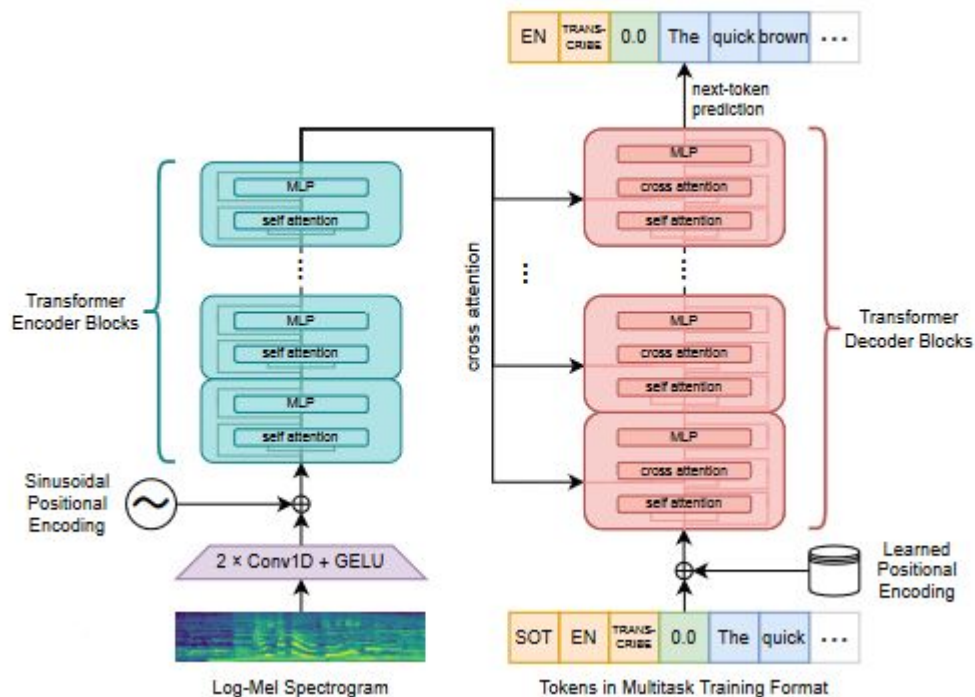
# Model

Observations:

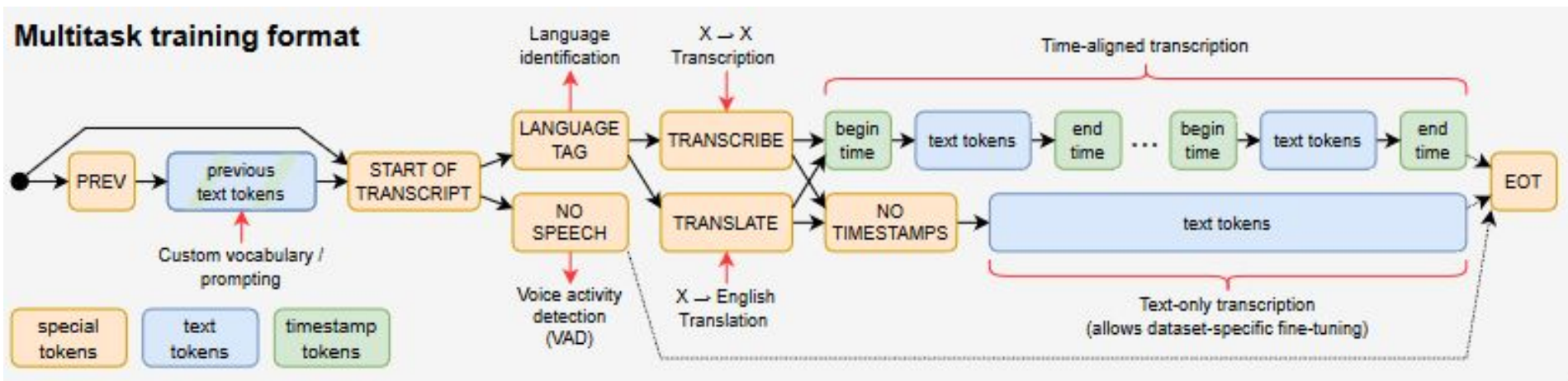
“It’s just a seq2seq transformer”

Novelty: prefix tokens + timestep

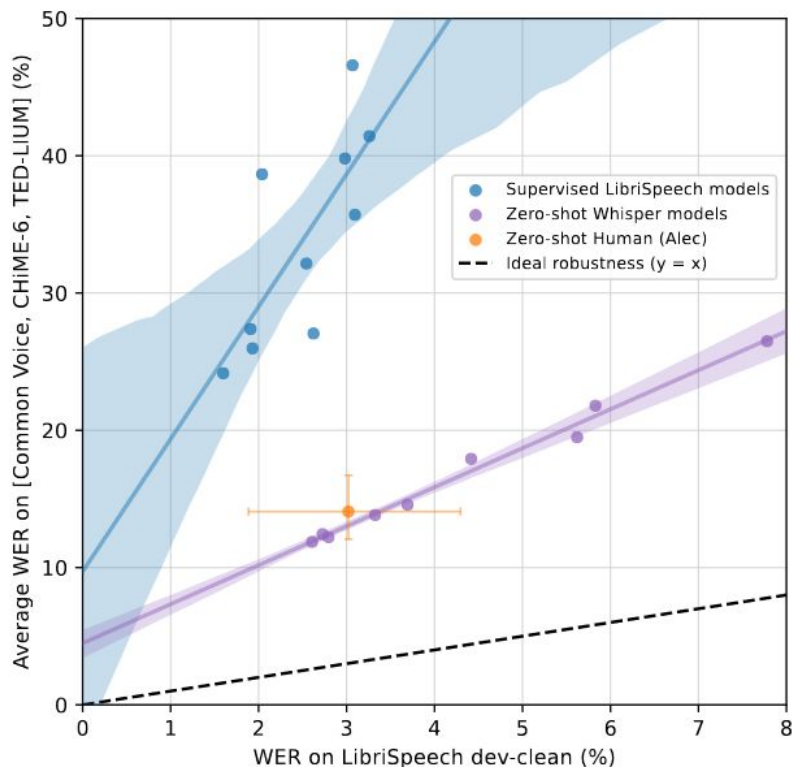
encoder sinusoidal encoding?



# Prefix tokens

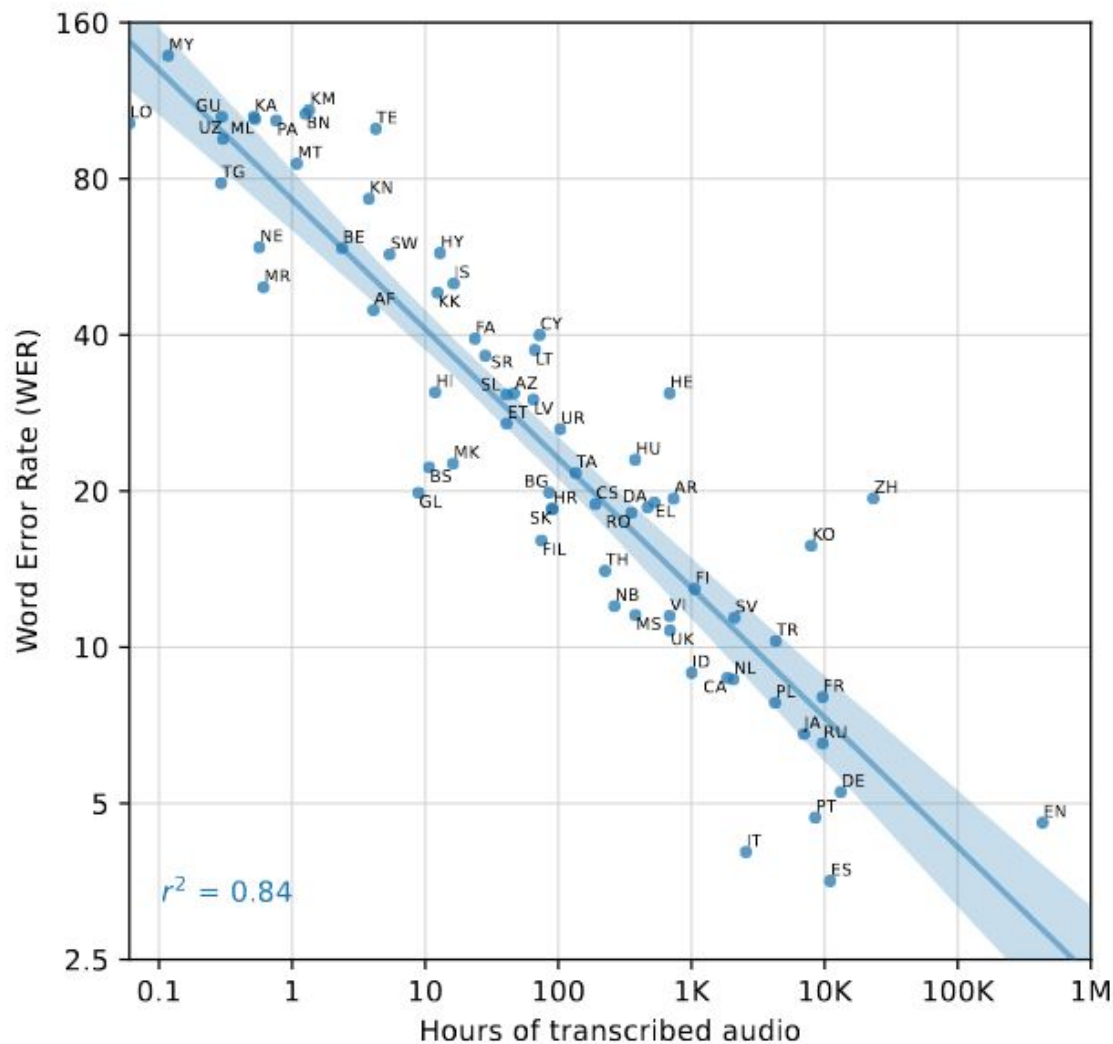


# Robustness?

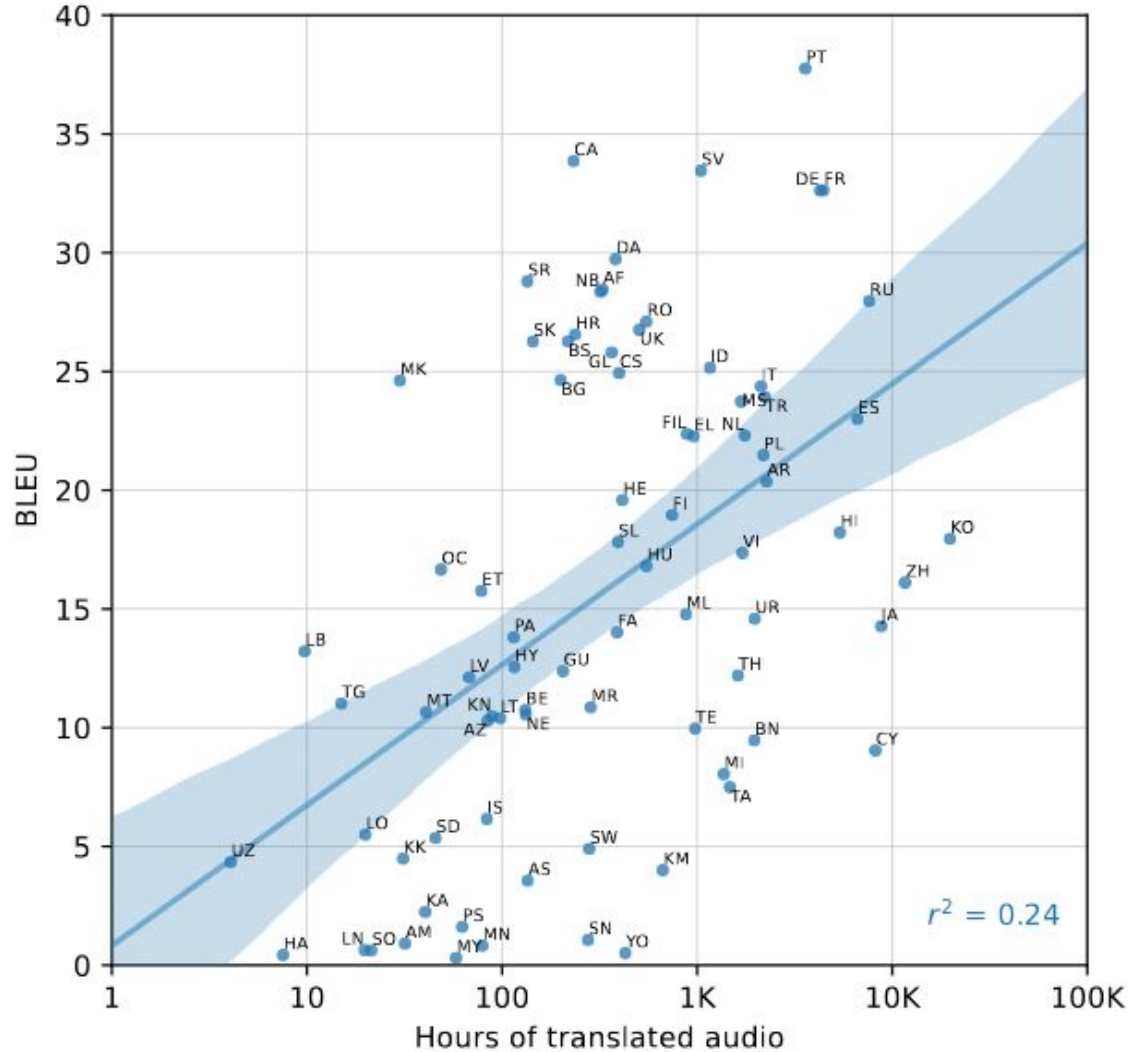


Dataset	wav2vec 2.0 Large 960h	Whisper Large	RER (%)
LibriSpeech test-clean	2.7	2.7	0.0
Artie	24.5	6.7	72.7
Fleurs (English)	14.6	4.6	68.5
Common Voice	29.9	9.5	68.2
Tedlium	10.5	4.0	61.9
CHiME6	65.8	25.6	61.1
WSJ	7.7	3.1	59.7
VoxPopuli (English)	17.9	7.3	59.2
AMI-IHM	37.0	16.4	55.7
CallHome	34.8	15.8	54.6
Switchboard	28.3	13.1	53.7
CORAAL	38.3	19.4	49.3
AMI-SDM1	67.6	36.9	45.4
LibriSpeech test-other	6.2	5.6	9.7
Average	29.5	12.9	55.4

# ASR (Fleurs)



# SLT (Fleurs)





# Aggregate Multilingual ASR

Model	MLS	VoxPopuli
Supervised Baseline	-	37.5
VP-10K + FT	-	15.3
XLS-R (1B)	10.9	10.6
mSLAM-CTC (2B)	9.7	<b>9.1</b>
Zero-Shot Whisper	<b>8.1</b>	15.2

MLS = Multilingual Librispeech

## Average SLT (Covost 2)

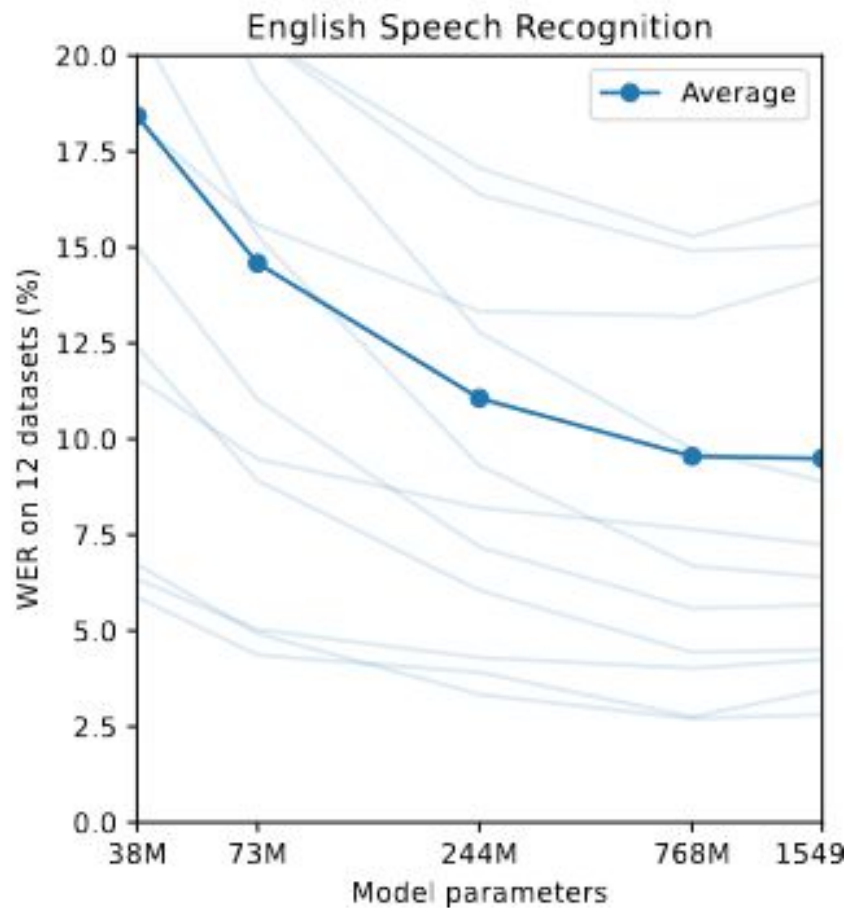
Average BLEU score performance, broken down by High / Mid / Low resource test settings

X → English	High	Mid	Low	All
XMEF-X	34.2	20.2	5.9	14.7
XLS-R (2B)	36.1	27.7	15.1	22.1
mSLAM-CTC (2B)	<b>37.8</b>	29.6	18.5	24.8
Zero-Shot Whisper	35.0	<b>31.1</b>	<b>23.1</b>	<b>27.3</b>

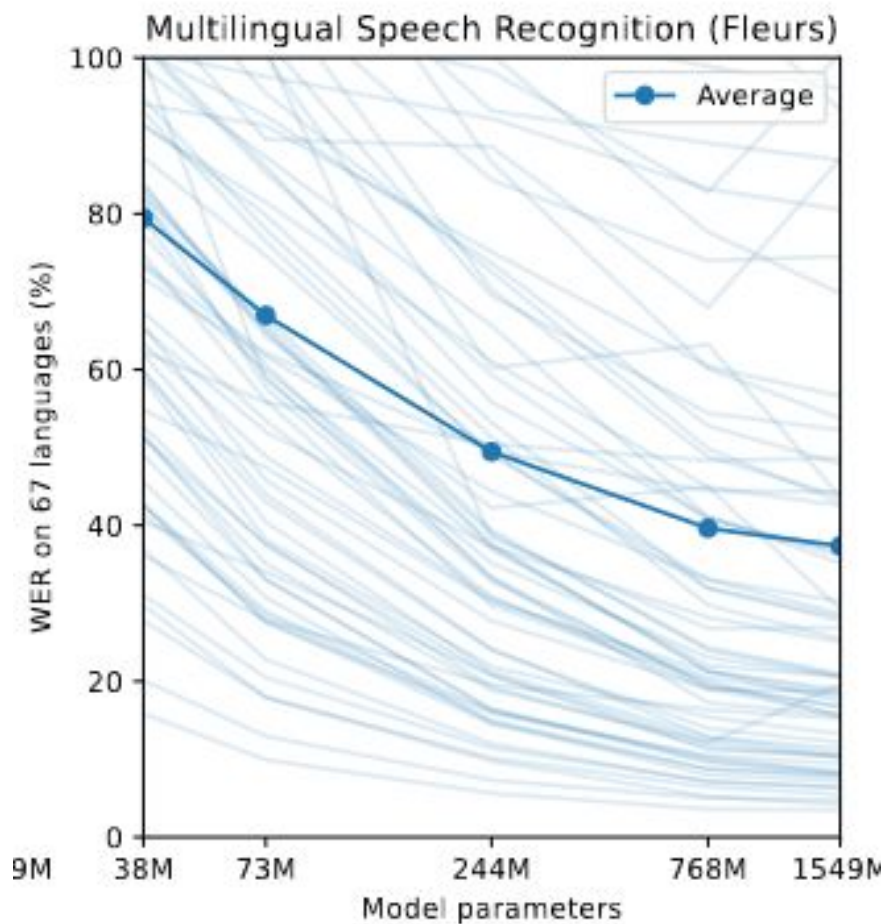
## Scaling (w/ dataset training hours)

Dataset size	English WER (↓)	Multilingual WER (↓)	X→En BLEU (↑)
3405	30.5	92.4	0.2
6811	19.6	72.7	1.7
13621	14.4	56.6	7.9
27243	12.3	45.0	13.9
54486	10.9	36.4	19.2
681070	<b>9.9</b>	<b>29.2</b>	<b>24.8</b>

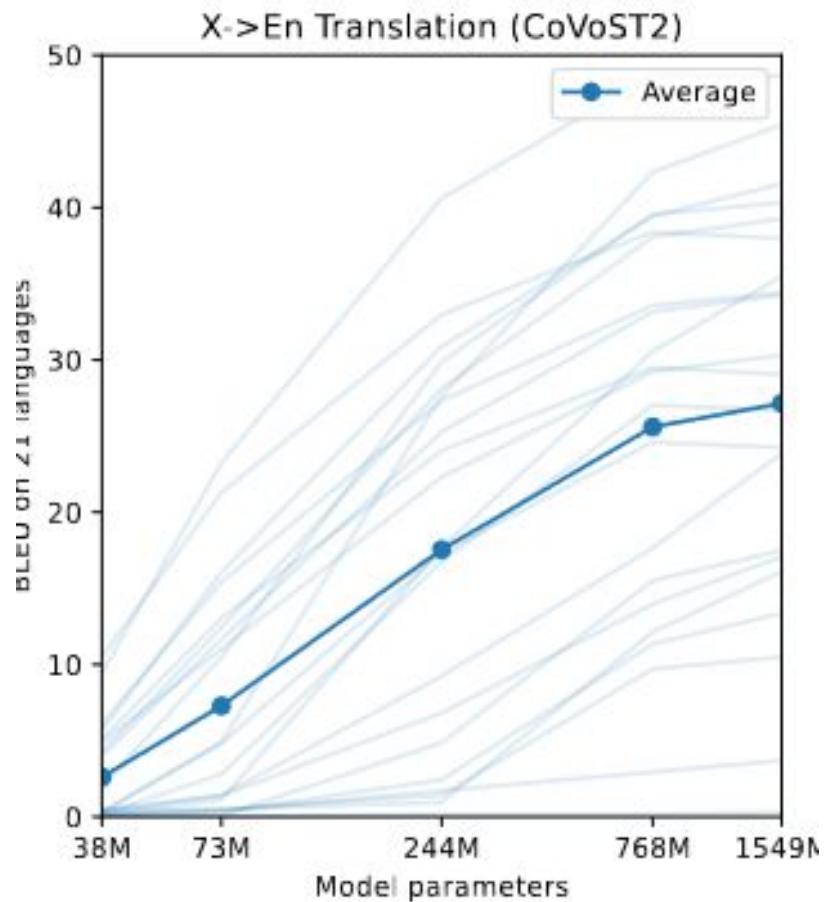
## Scaling cont.



Scaling cont.



## Scaling cont.



# Scaling cont.

Weird

