

# Modeling Context With Linear Attention for Scalable Document-Level Translation

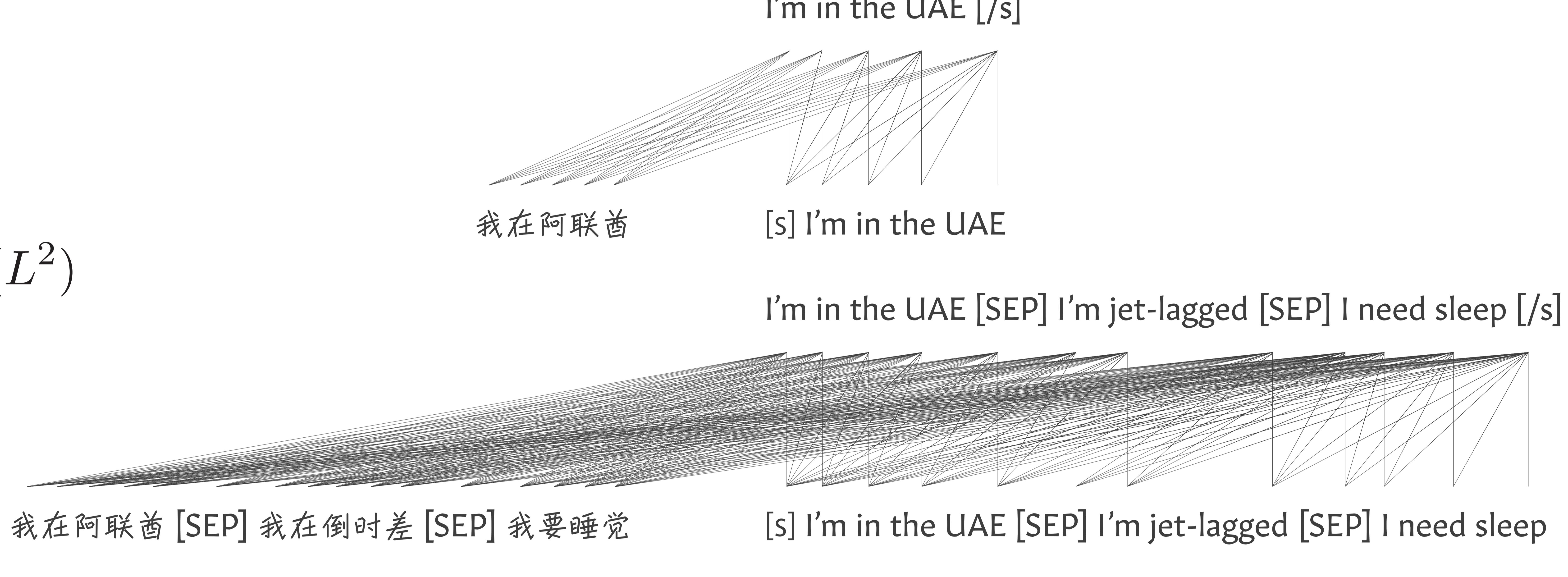
Zhaofeng Wu, Hao Peng, Nikolaos Pappas, Noah Smith

zfw@csail.mit.edu

Replacing transformer attention to a linear approximation leads to substantial speedup on doc-level translation with similar/better quality

## DOC-LEVEL TRANSLATION

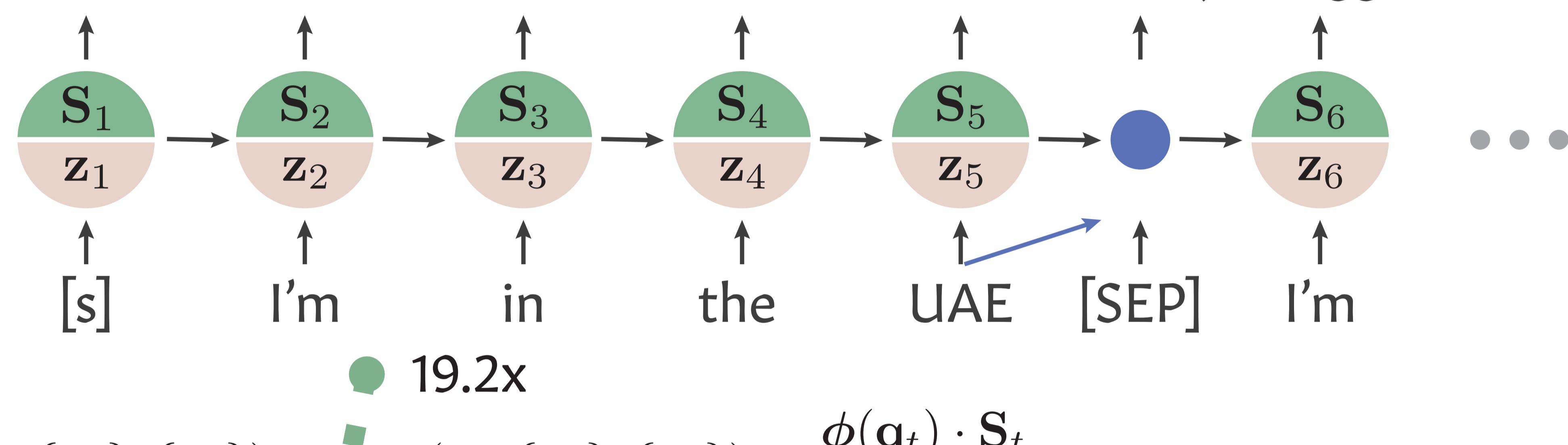
$\mathcal{O}(L^2)$



## RANDOM FEATURE ATTENTION

$\mathcal{O}(L)$

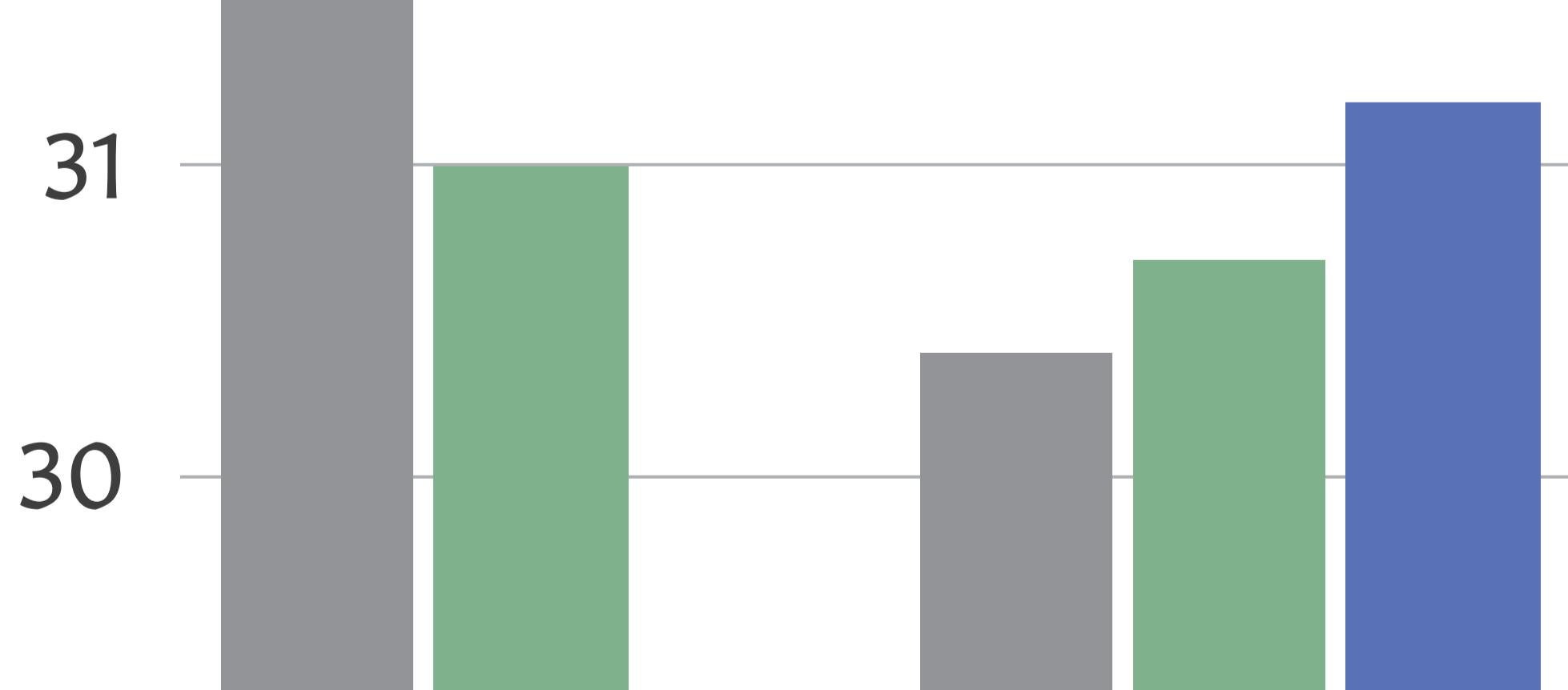
我在阿联酋 [SEP] 我在倒时差 [SEP] 我要睡觉



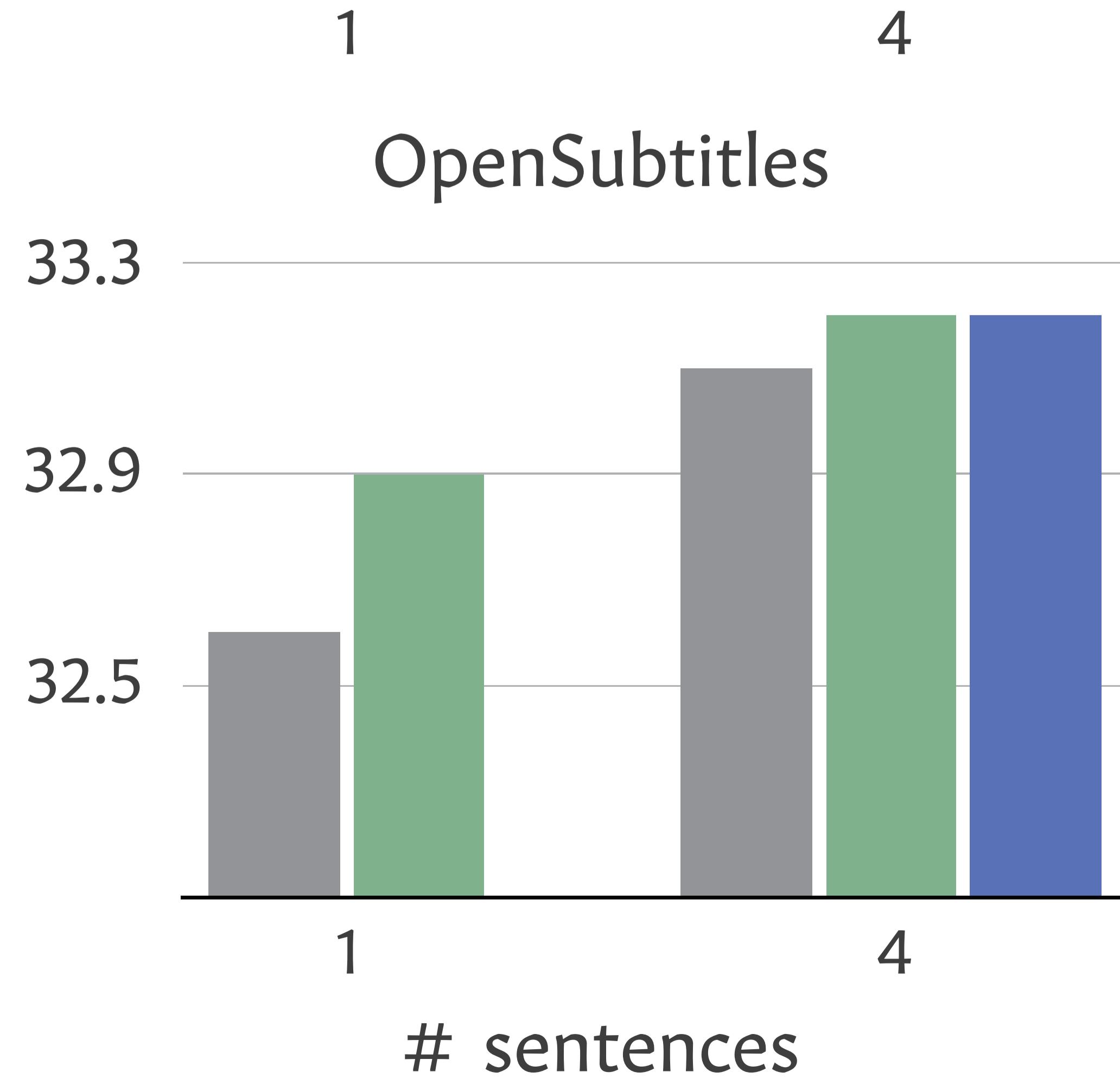
## RESULTS

### BLEU

IWSLT



OpenSubtitles

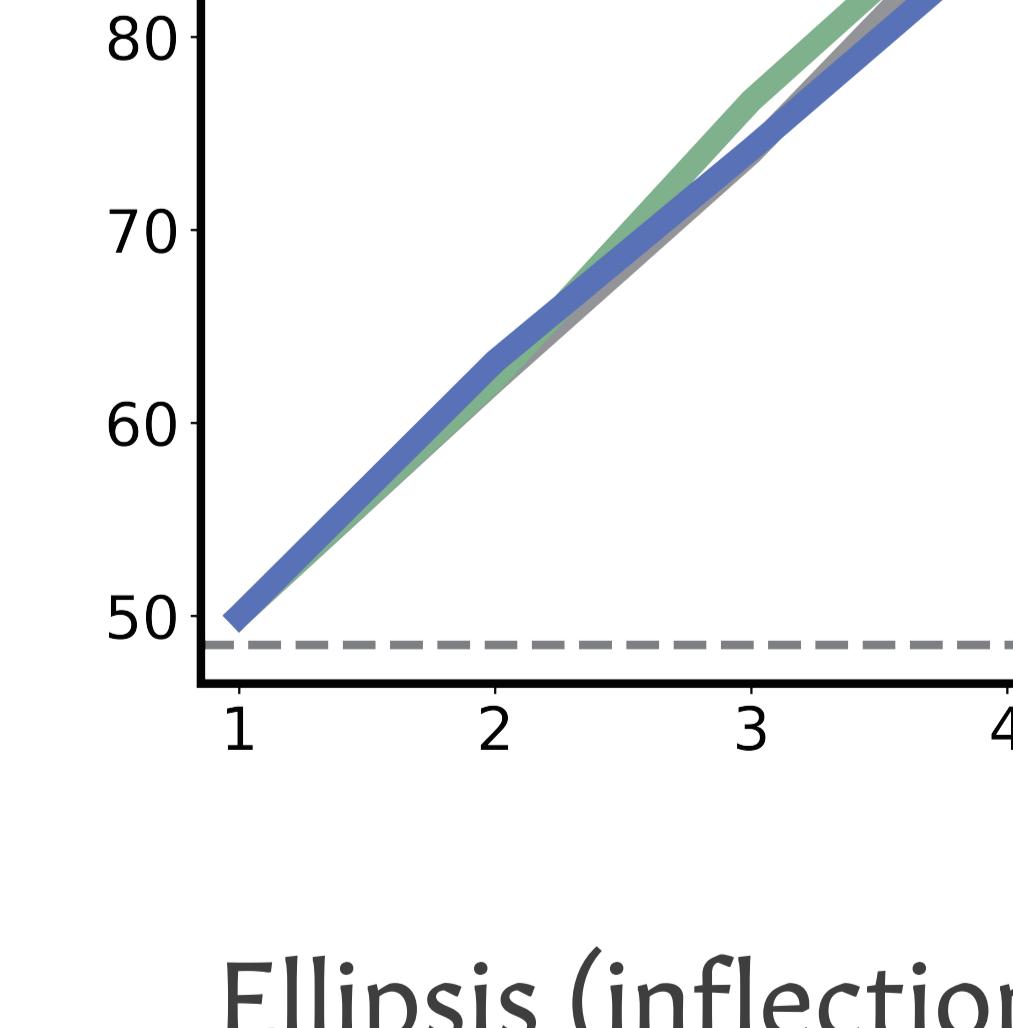


Transformer RFA RFA+gate

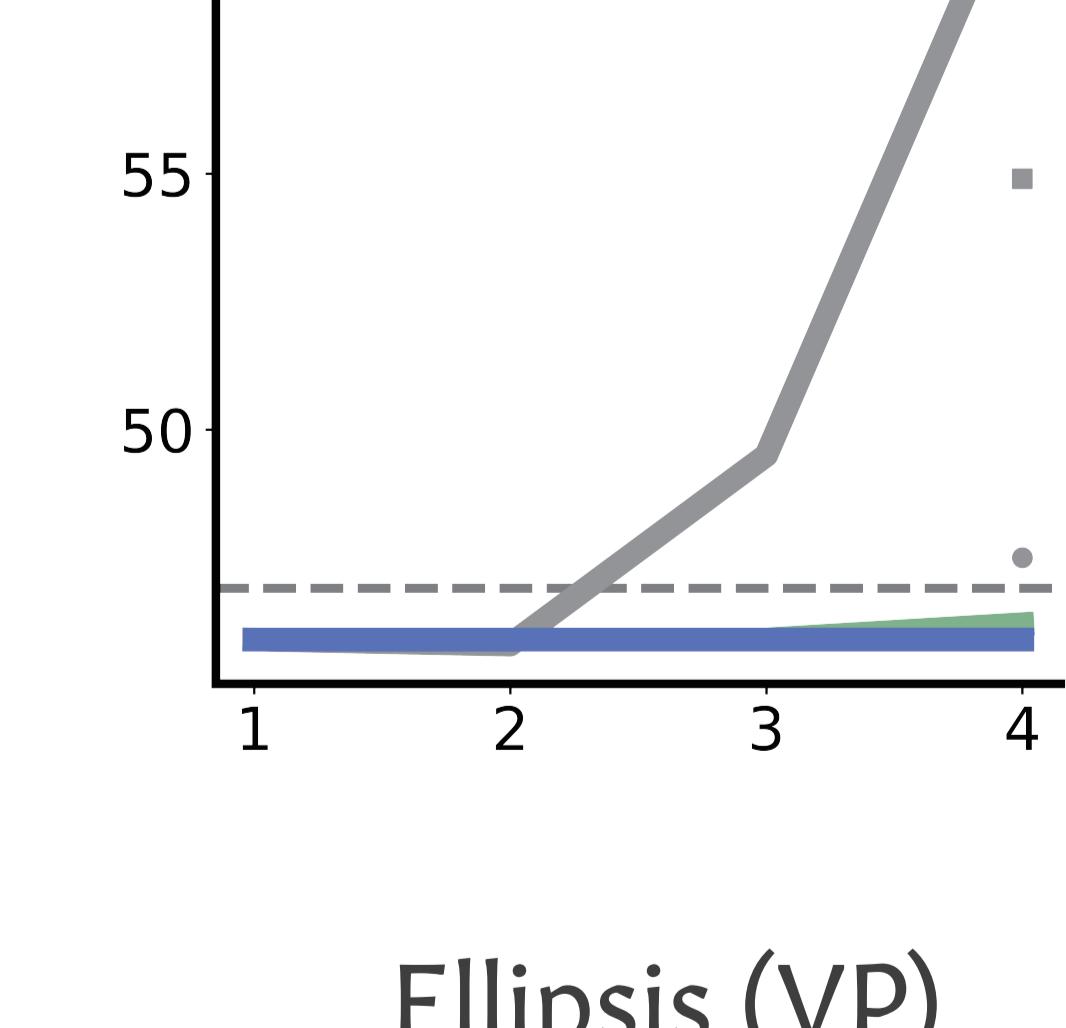
### Consistency

I have received your gift [SEP] Thank you  
我 收到 你 的 礼物 了 [SEP] 谢谢 您  
I receive you-GEN gift-PFV Thank-you-FOR

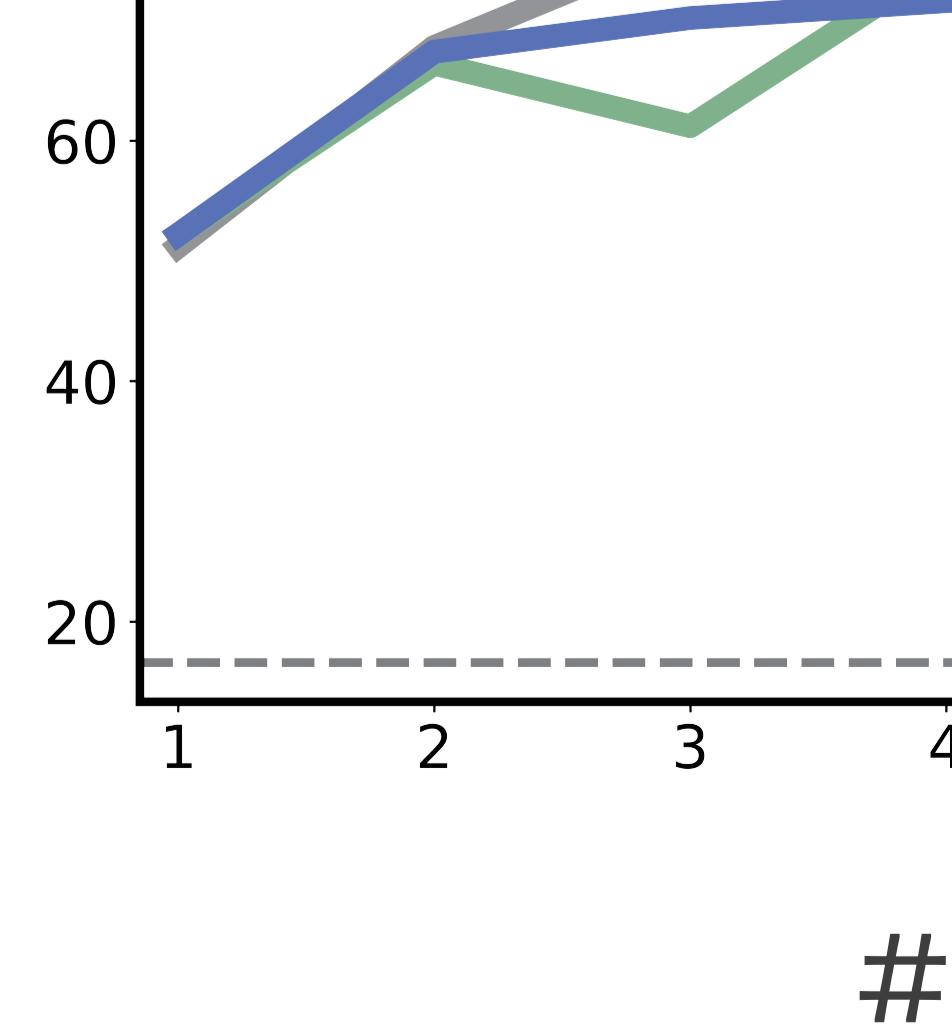
Deixis



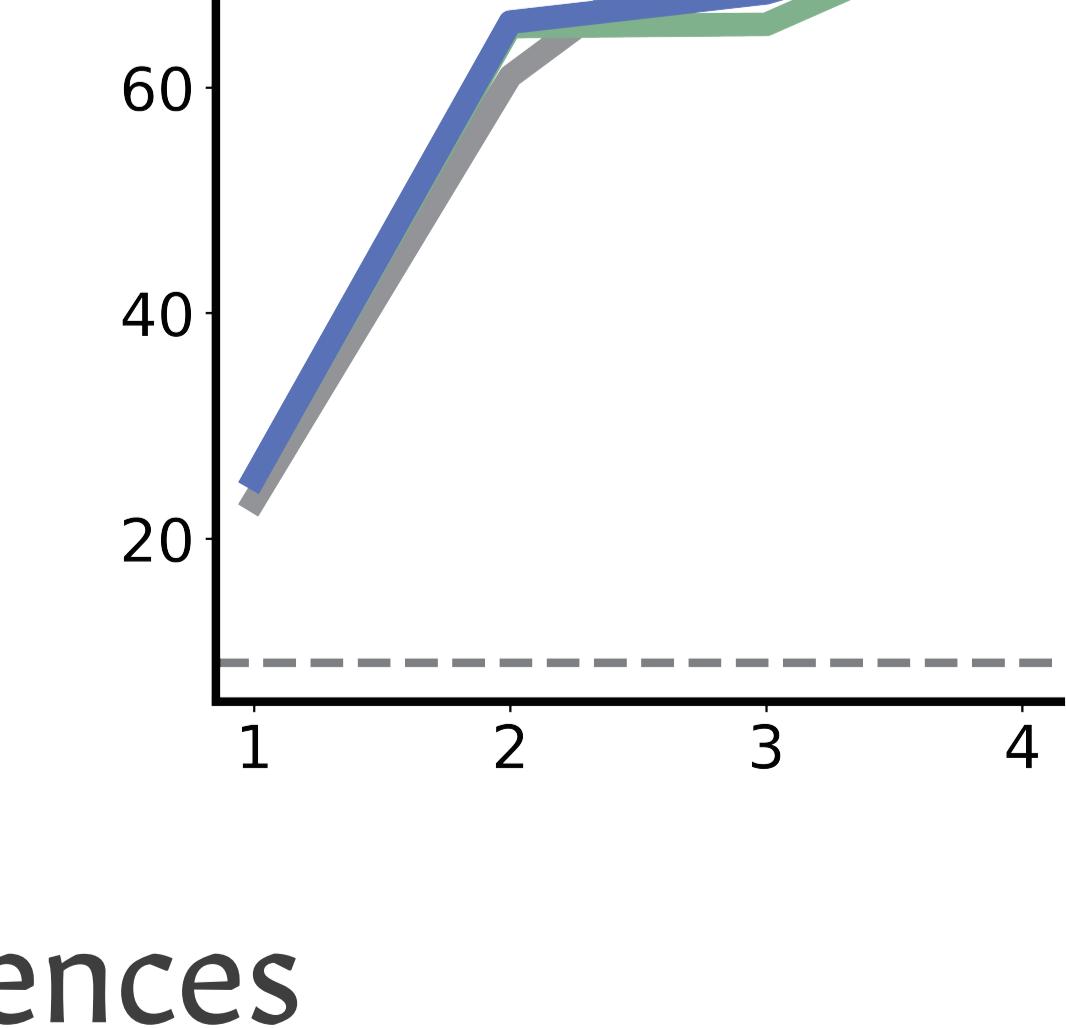
Lexical Cohesion



Ellipsis (inflection)



Ellipsis (VP)



# sentences

Transformer RFA RFA+gate

paper

