

# Unified self-supervised learning for speech, vision and NLP

Meta AI



Arun Babu



Alexis  
Conneau



Steffen  
Schneider



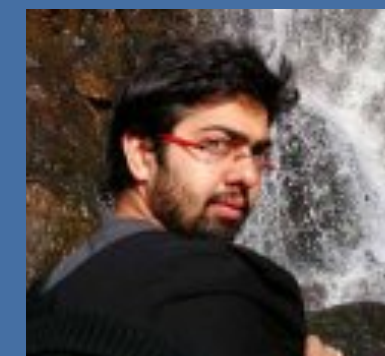
Henry Zhou



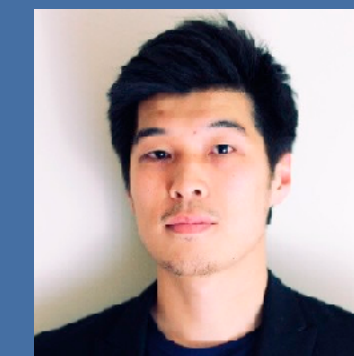
Abdelrahman  
Mohamed



Jiatao Gu



Naman  
Goyal



Wei-Ning Hsu



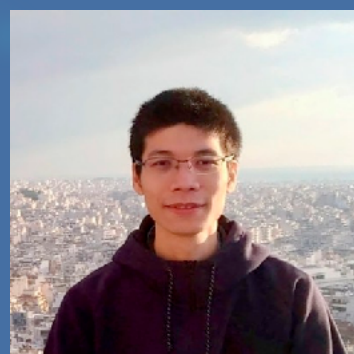
Alexei Baevski



Michael Auli



Kushal  
Lakhotia



Andros Tjandra



Kritika Singh



Yatharth Saraf



Geoffrey Zweig



Qiantong Xu



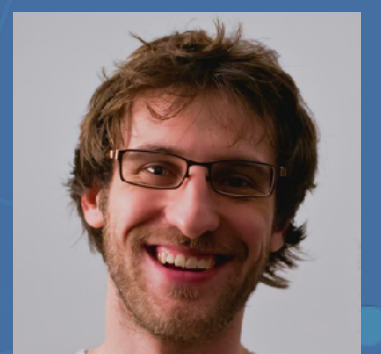
Tatiana  
Likhomanenko



Paden  
Tomasello



Ronan  
Collobert



Gabriel  
Synnaeve



# Why Self-supervised Learning?

# Supervised Machine Learning



potential train/test mismatch



Need to annotate lots of data!



# Supervised Machine Learning

(  , cat )

Not how humans learn!

potential train/test mismatch



Need to annotate lots of data!



# Supervised Machine Learning?

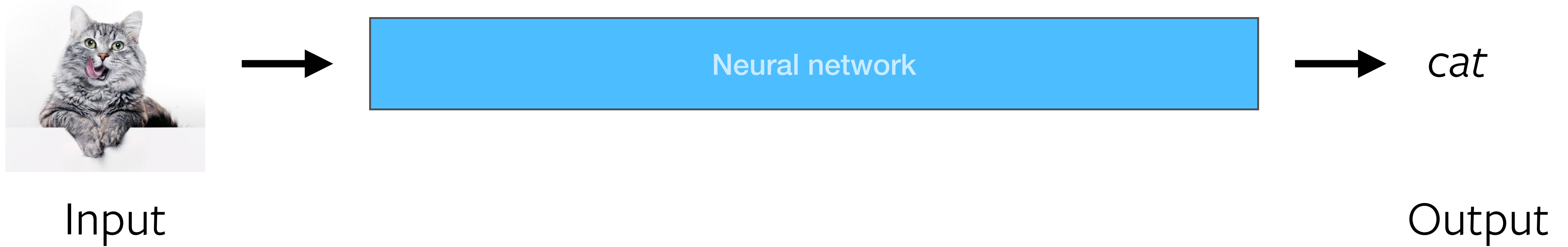




# Self-supervised Learning

- Learn good data representations (structure, features etc.) **without labels**
- $|\text{Unlabeled data}| \gg |\text{Labeled data}|$
- Use representations to solve the task





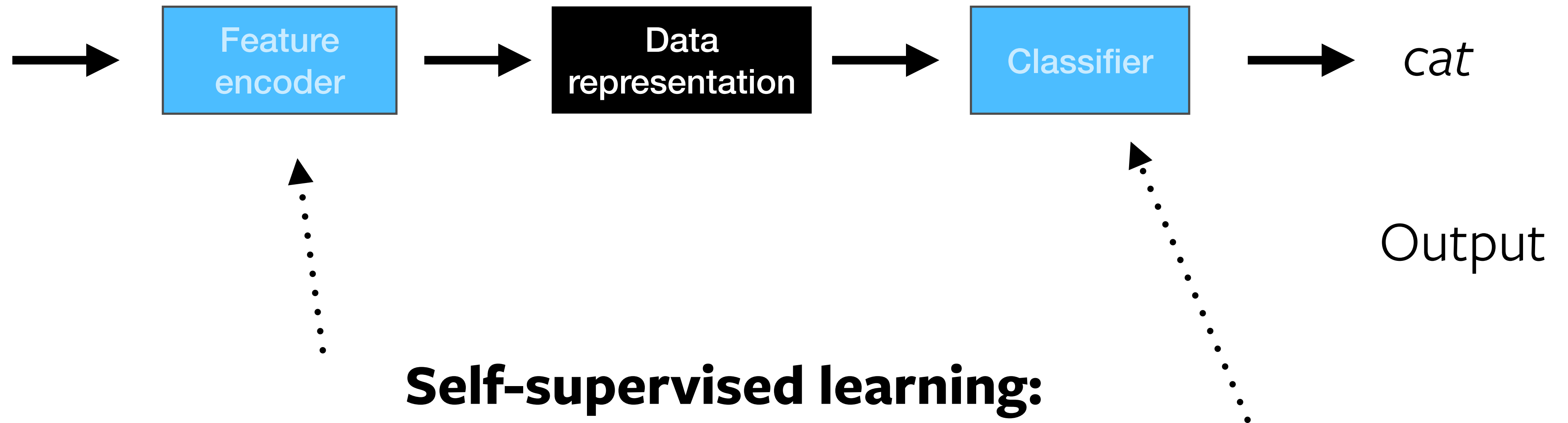
**Supervised learning** simultaneously performs representation learning of the data and associating these features with labels

**Limitation:** relies on labeled data to learn feature encoding





Input

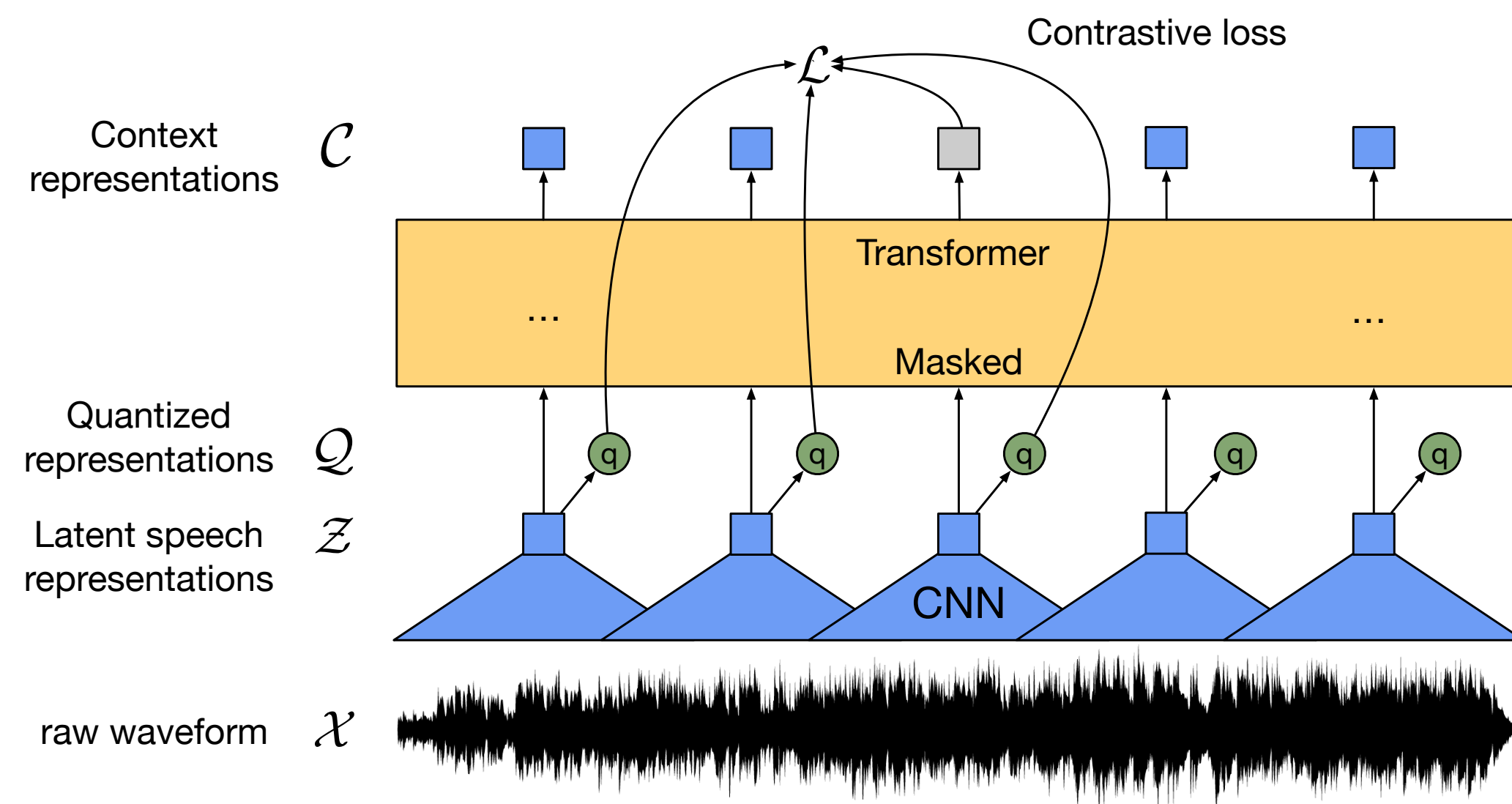


- 1/ representation learning of the data
- 2/ learn to associate labels with the representations

Reduces reliance on labeled data!

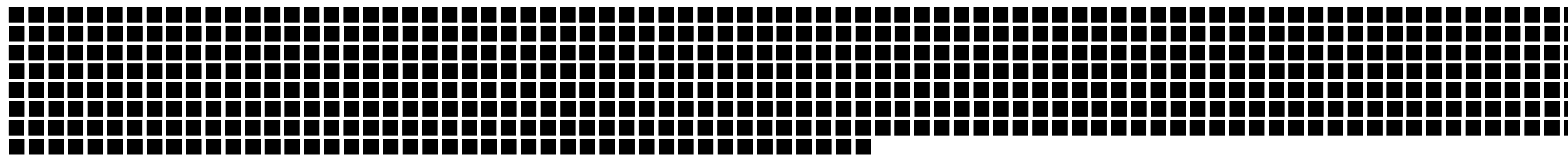


# wav2vec 2.0



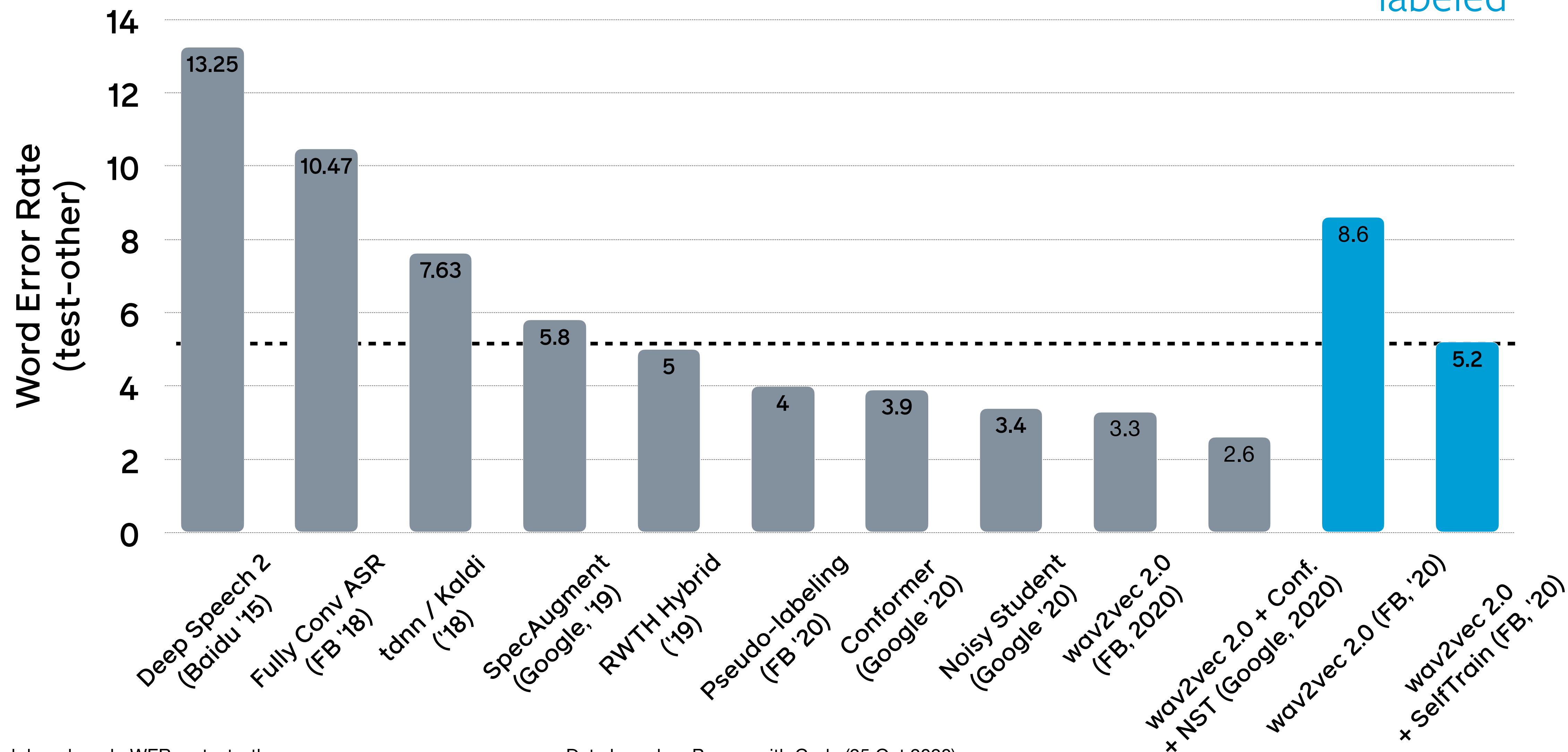
- Masked prediction with Transformer (similar to BERT).
- But predict what? Learned inventory of speech units with vector quantization!
- Learning task: Joint VQ & masked prediction

Amount of  
labeled  
data used



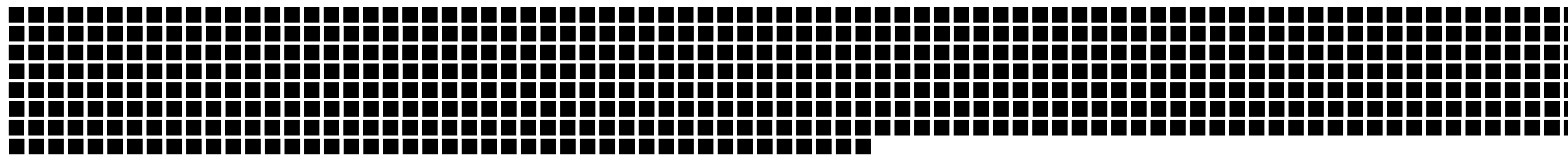
960h labeled

10min  
labeled



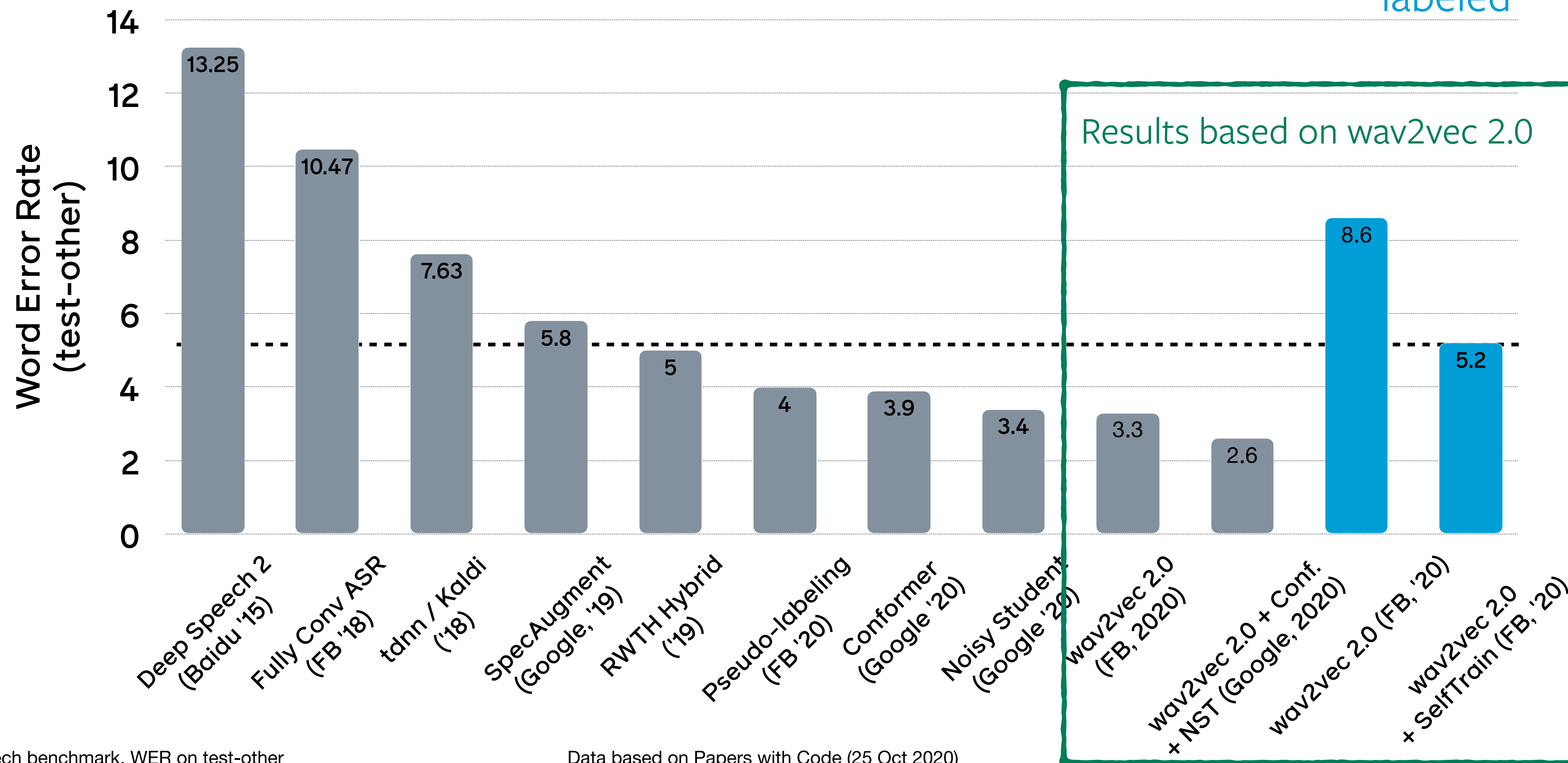


Amount of  
labeled  
data used



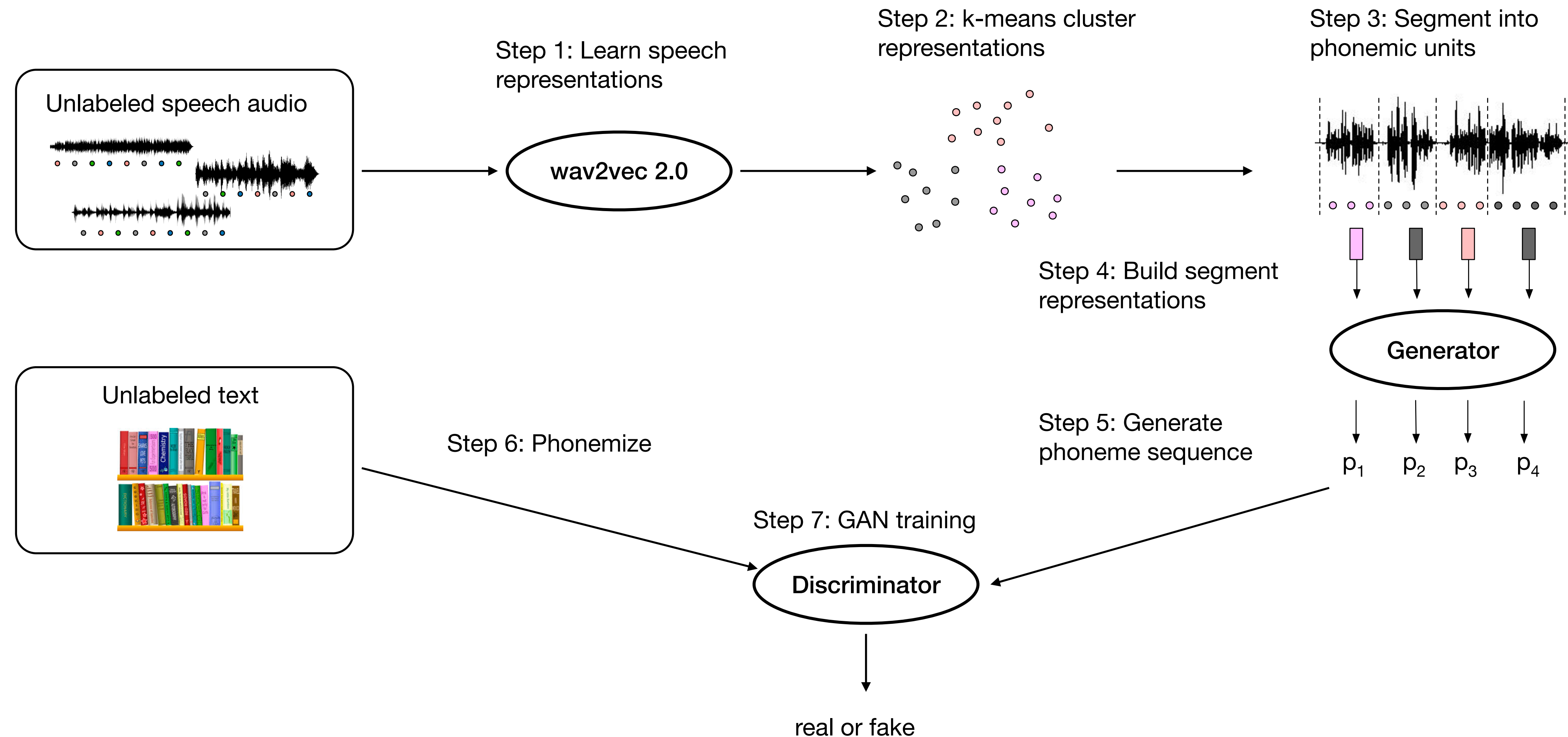
960h labeled

10min  
labeled



# Unsupervised Speech Recognition

- Speech-to-Text with no labels.

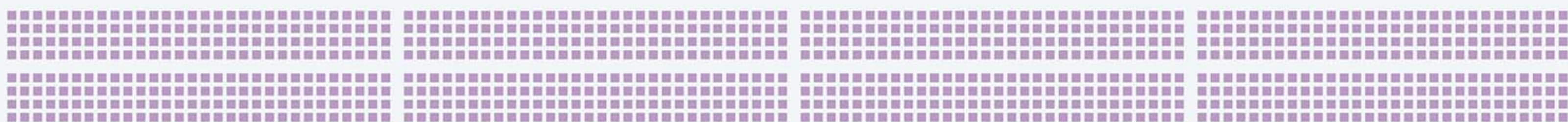




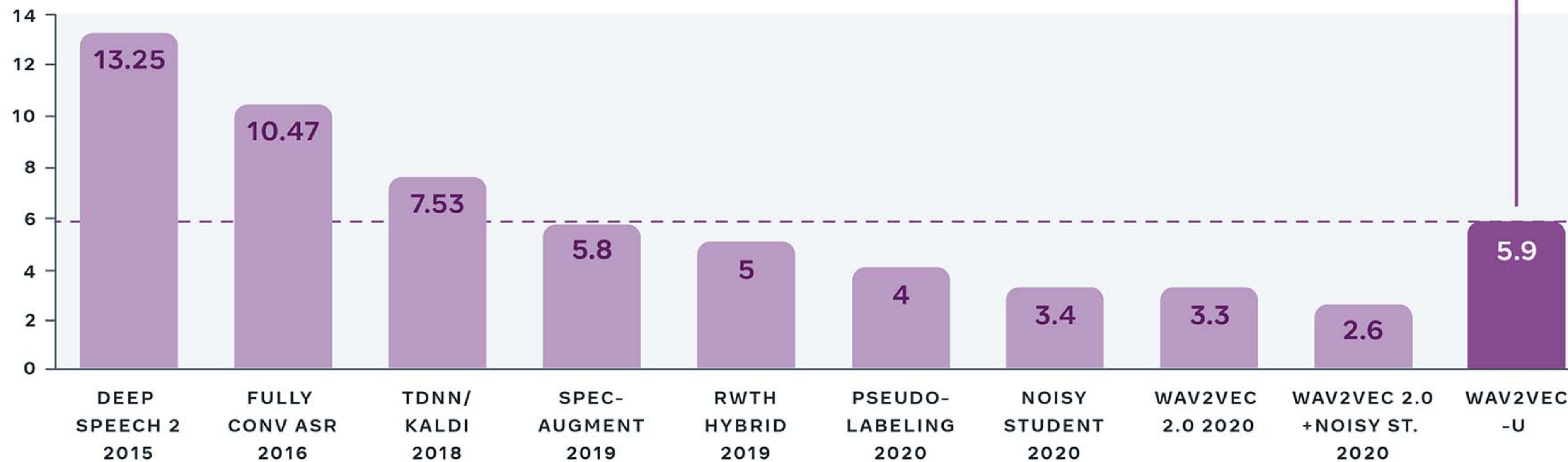
# Comparison to Best Supervised Systems

Amount of labeled data used

960 hrs.+ ■=1hr.



Word error rate



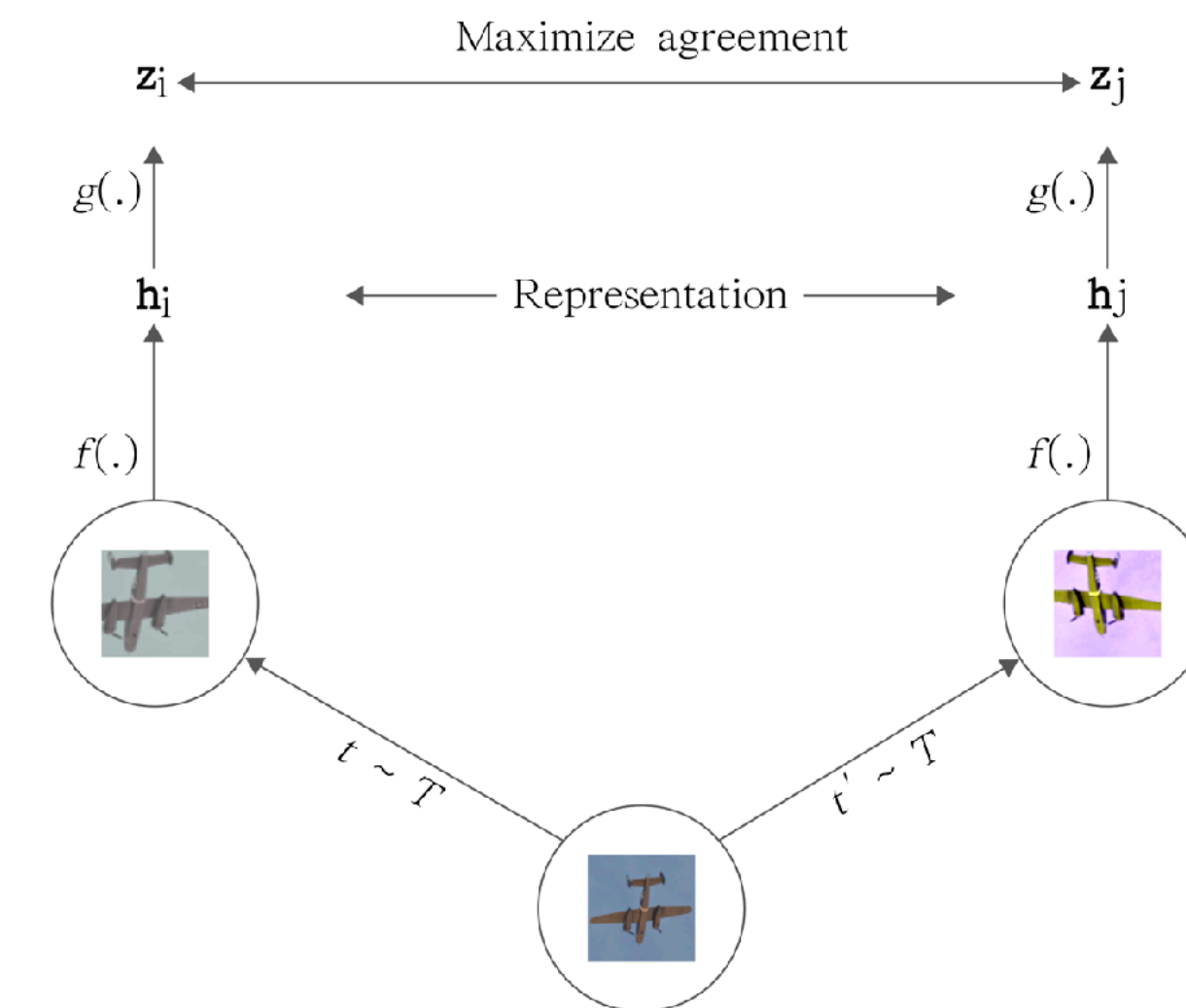
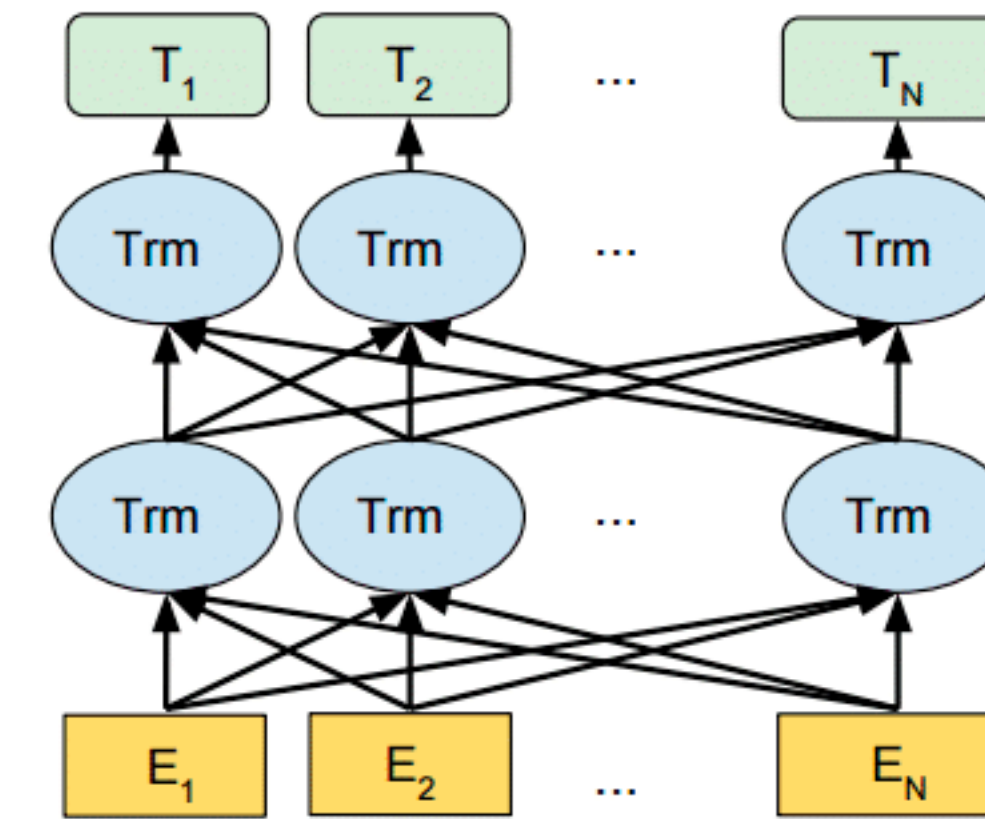
The background is a solid blue color with several faint, light-blue geometric patterns. These patterns consist of interconnected lines and dots, resembling a network or a series of overlapping triangles and polygons, scattered across the slide.

# data2vec: A Unified Objective for Self-supervised Learning



# Self-supervised Learning

- NLP: BERT, GPT, ...
- Vision: MoCo, SimCLR, BYOL, DINO, MAE, ....
- Speech: wav2vec, CPC, APC, HuBERT, ...



# Current State of Self-supervised Learning

- Many different algorithms
- Most algorithms developed for particular modality
- Little focus on algorithms that generalize across modalities



# Underlying Learning Mechanisms

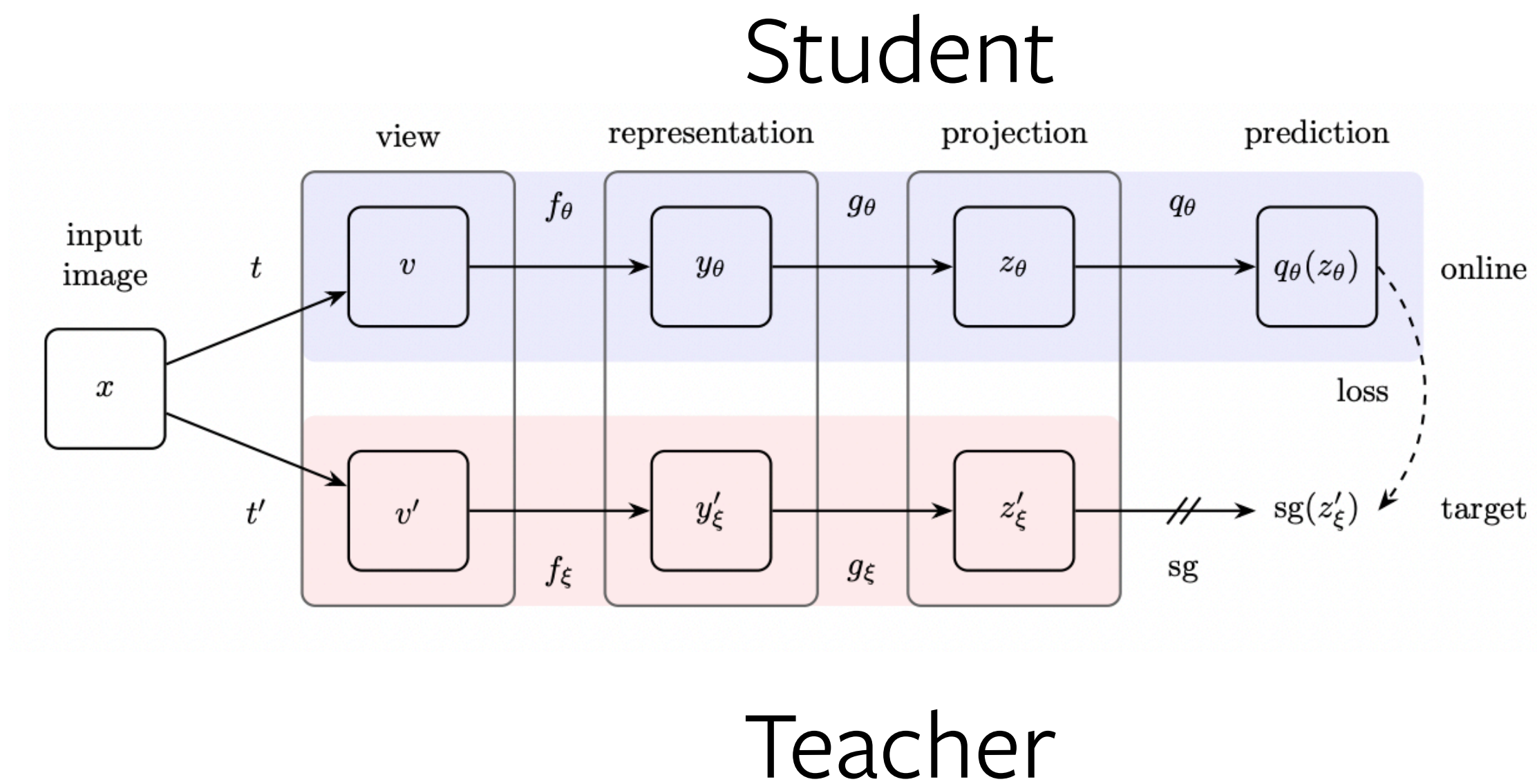


# data2vec

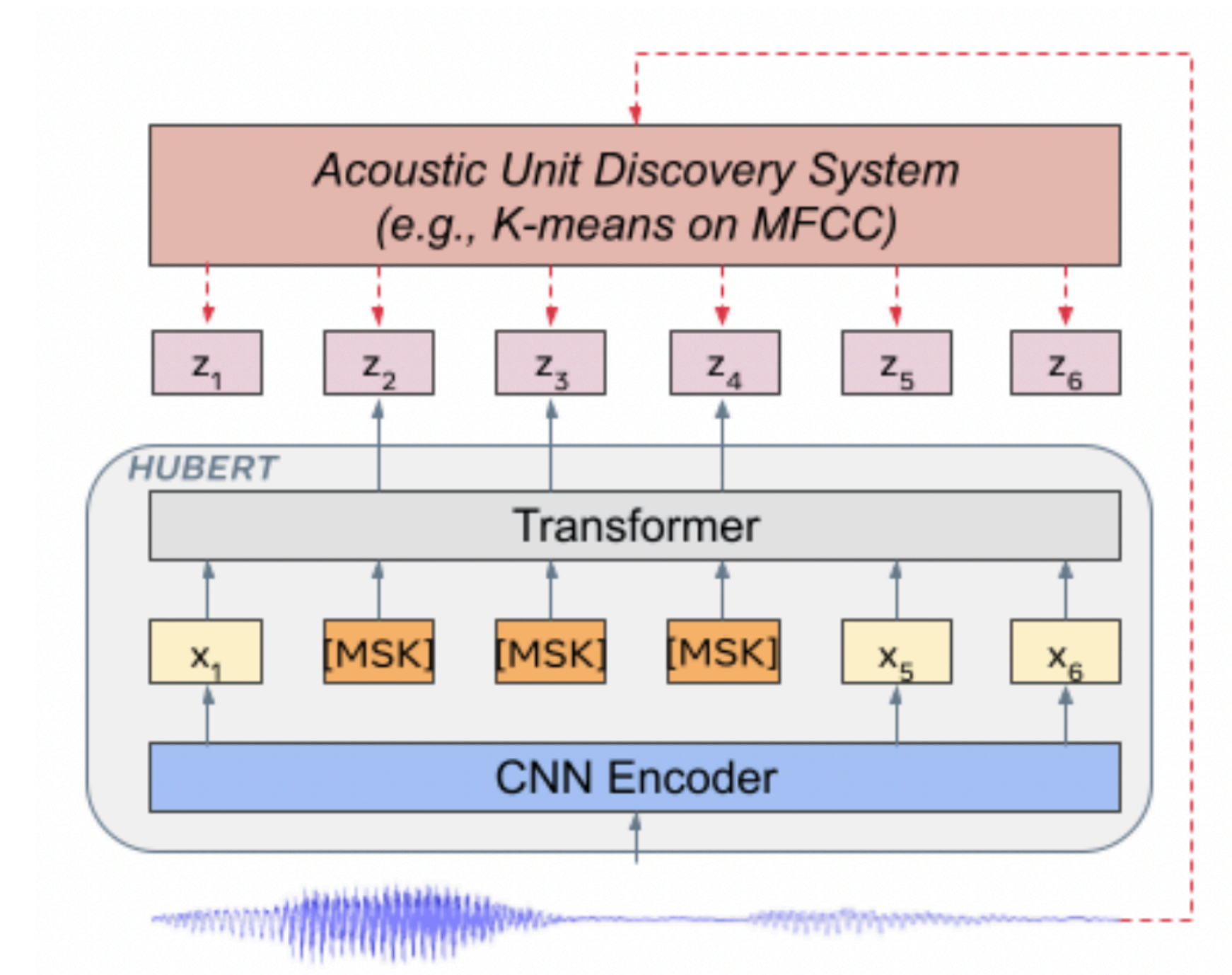
- General algorithm that works very well across modalities  
(Outperforms best algorithms in speech/vision and competitive in NLP)
- Same learning objective for each modality
- Idea: self-distillation of contextualized representations in a masked prediction setup

# Related Work

- Momentum teacher  
(Grill et al., '20, Caron et al., '21)

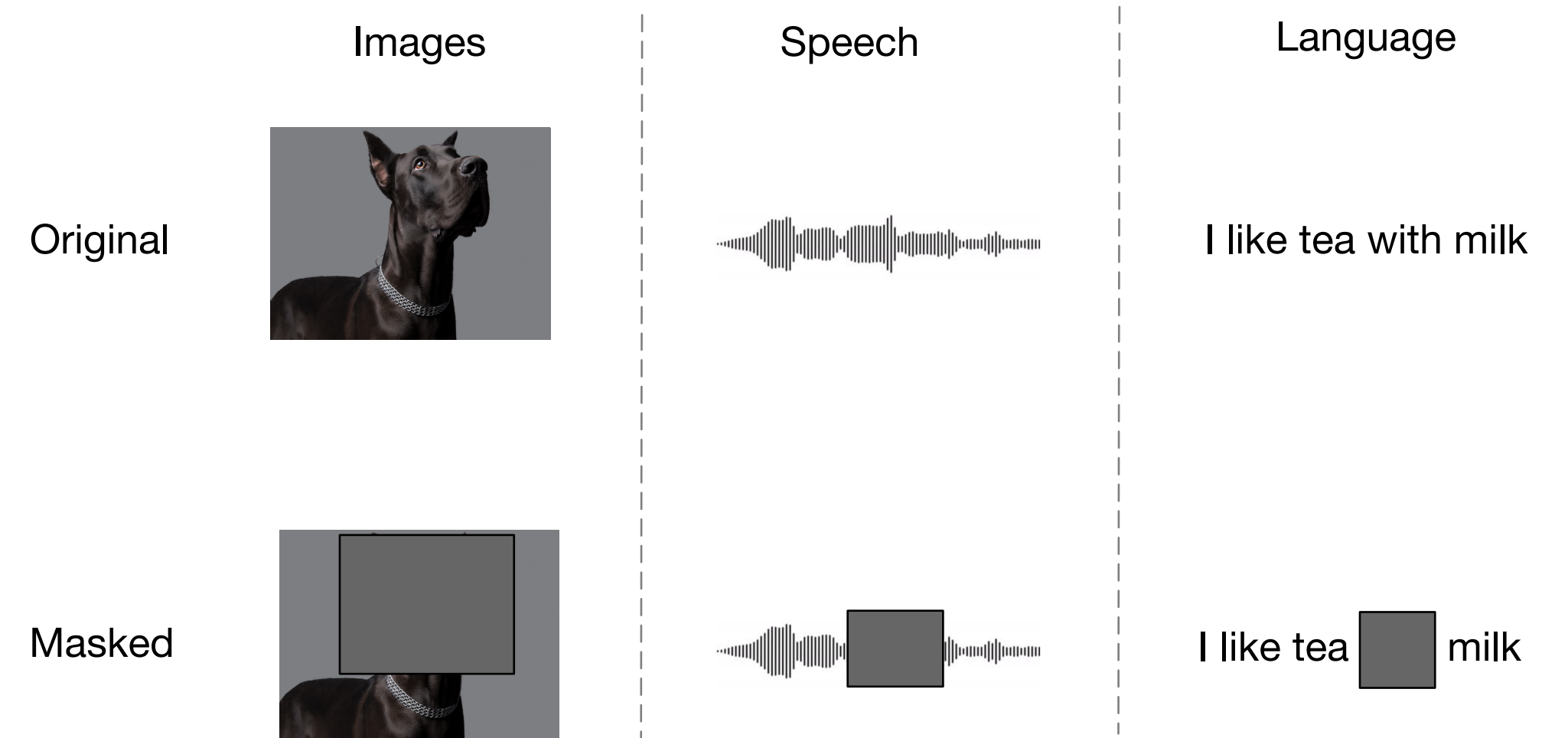


- Contextualized targets  
(Hsu et al., '21)

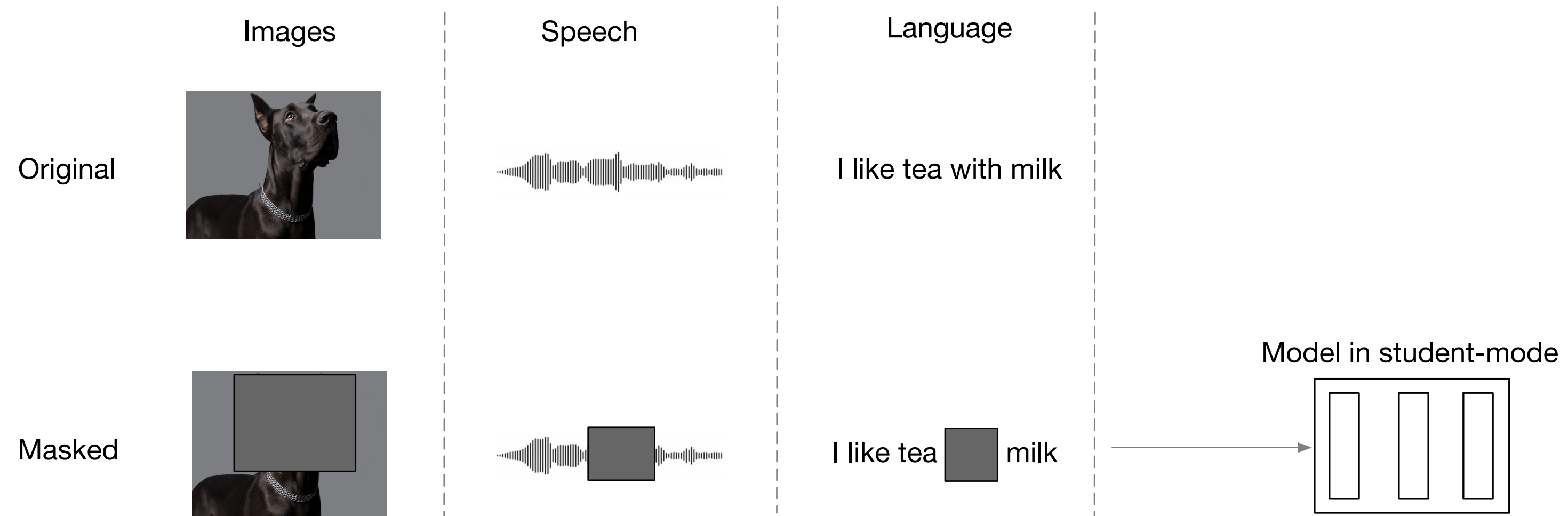




# data2vec

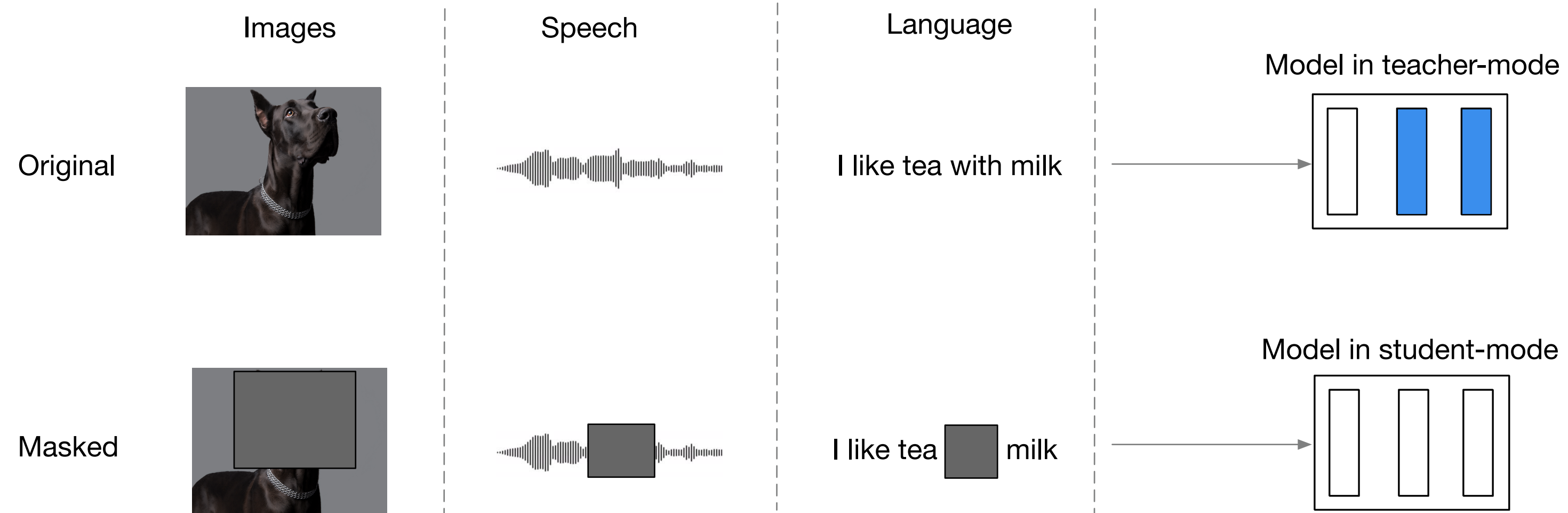


# data2vec



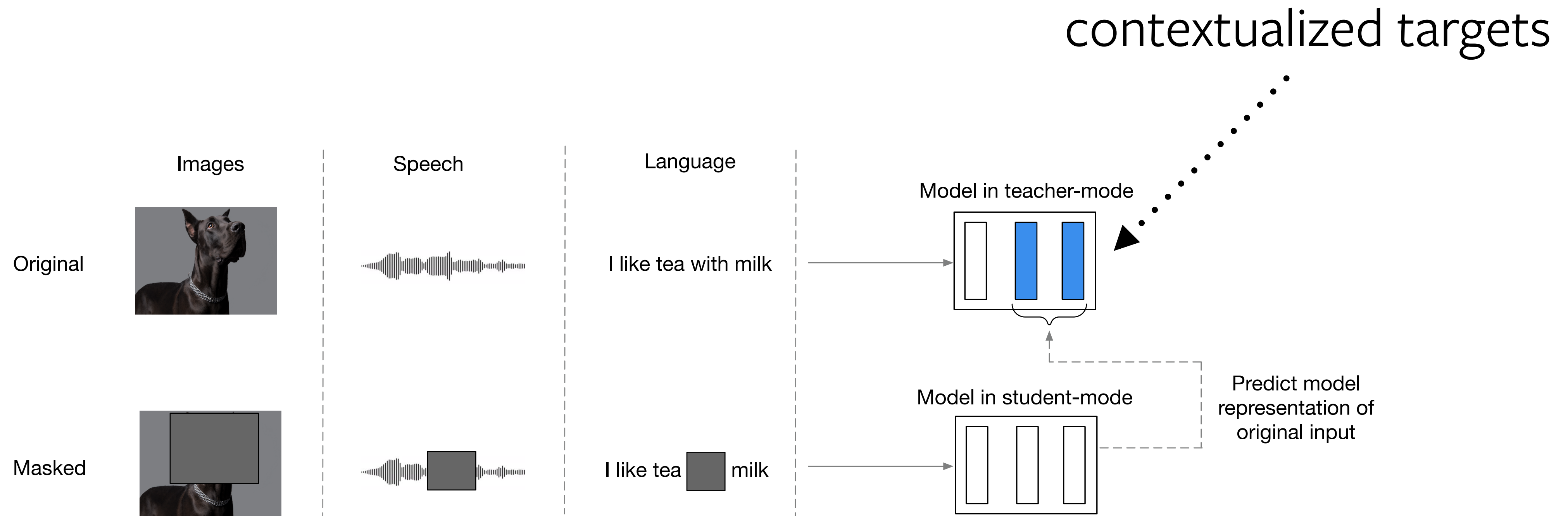
masked prediction

# data2vec

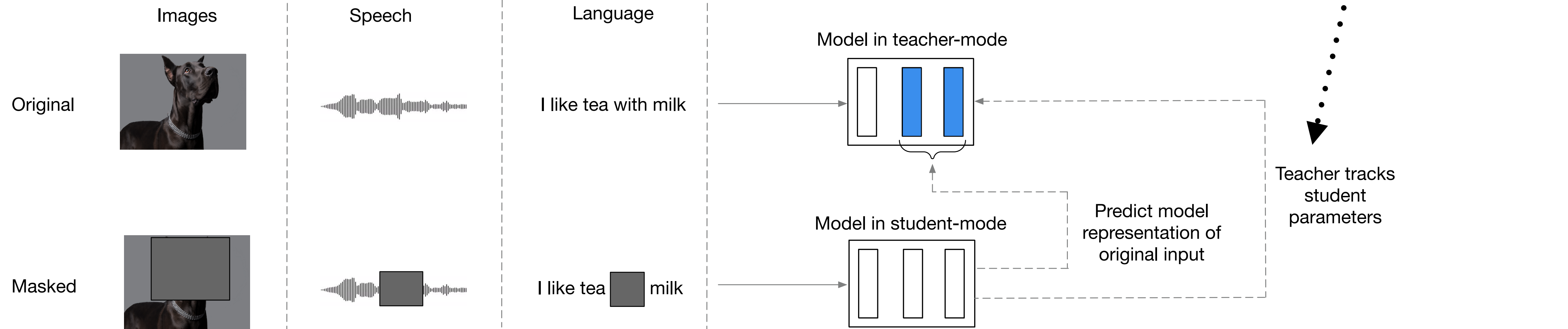




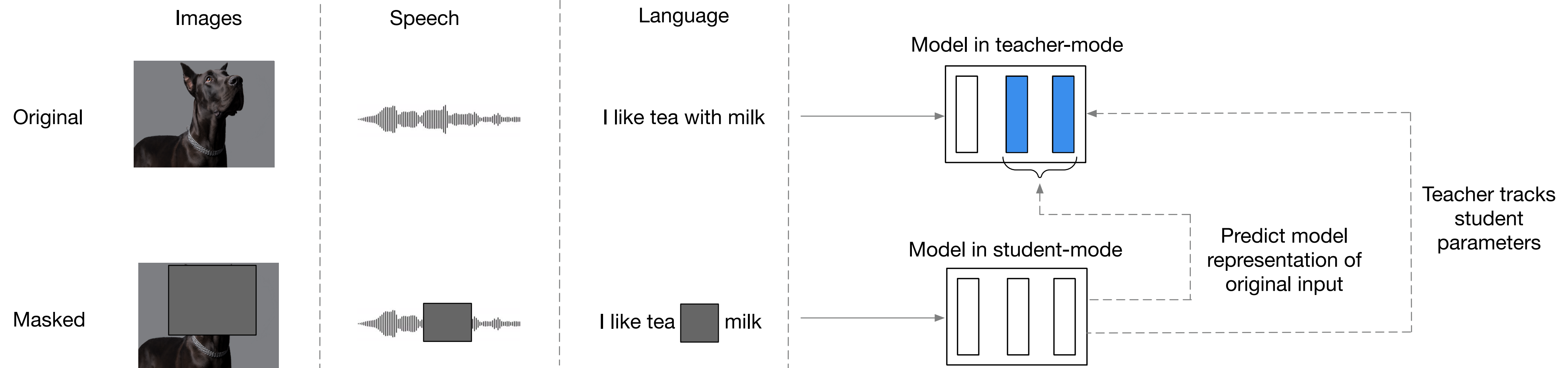
# data2vec



# data2vec



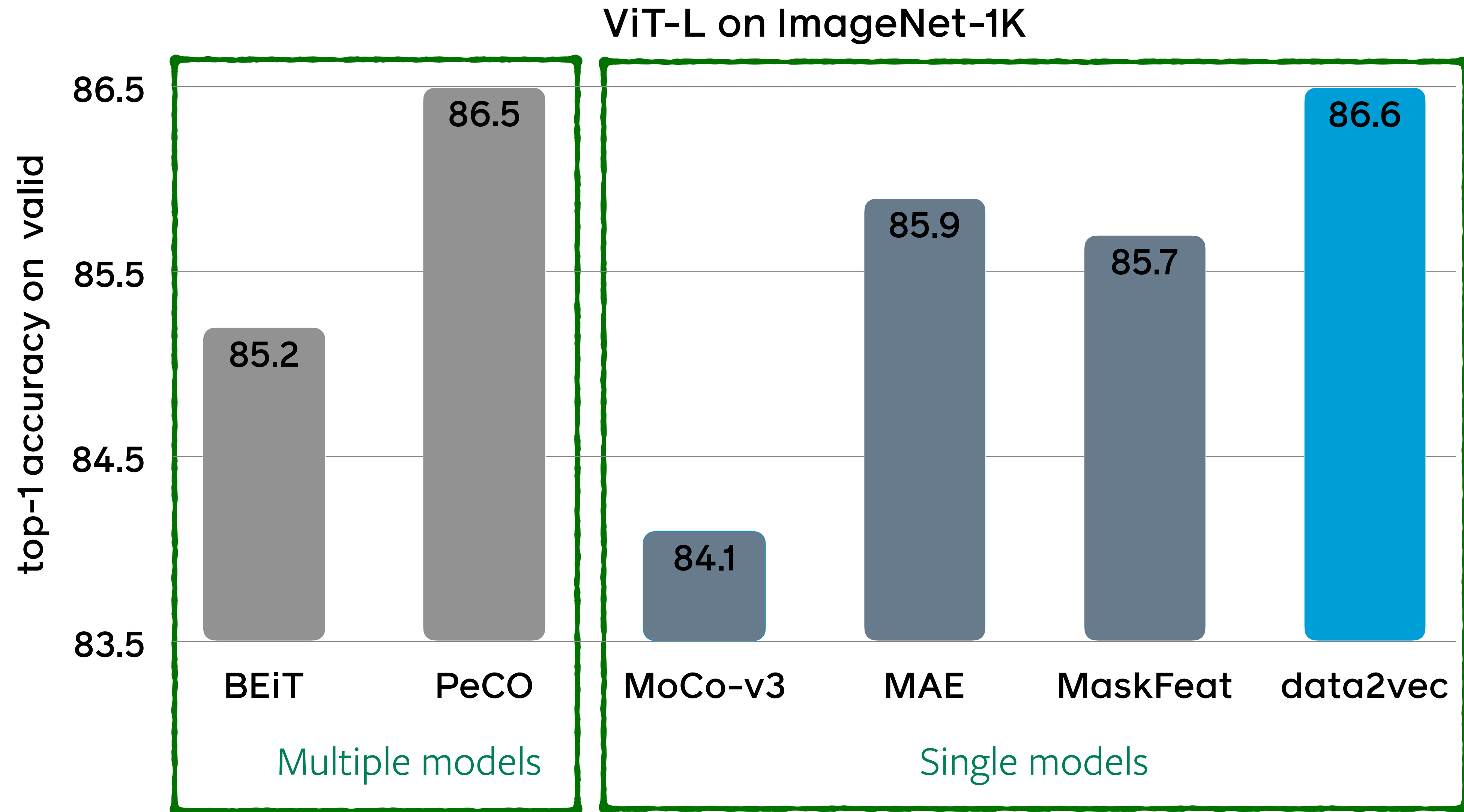
# data2vec



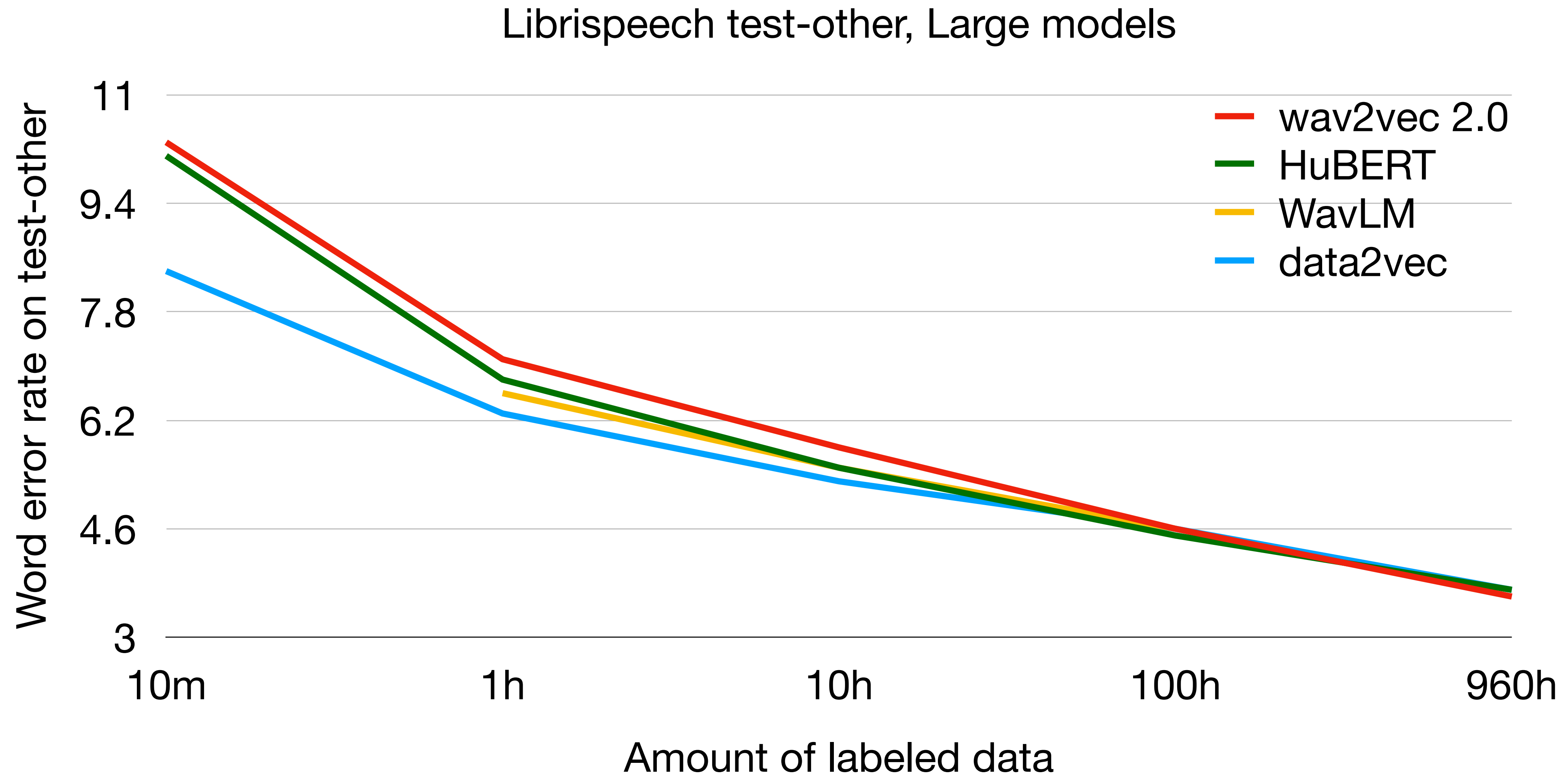
- Modality specific feature encoder (CNN, embedding table, patch mapping)
- Common masking policy, but modality/dataset specific parameterization
- Identical context encoder (Transformer)
- Identical learning task



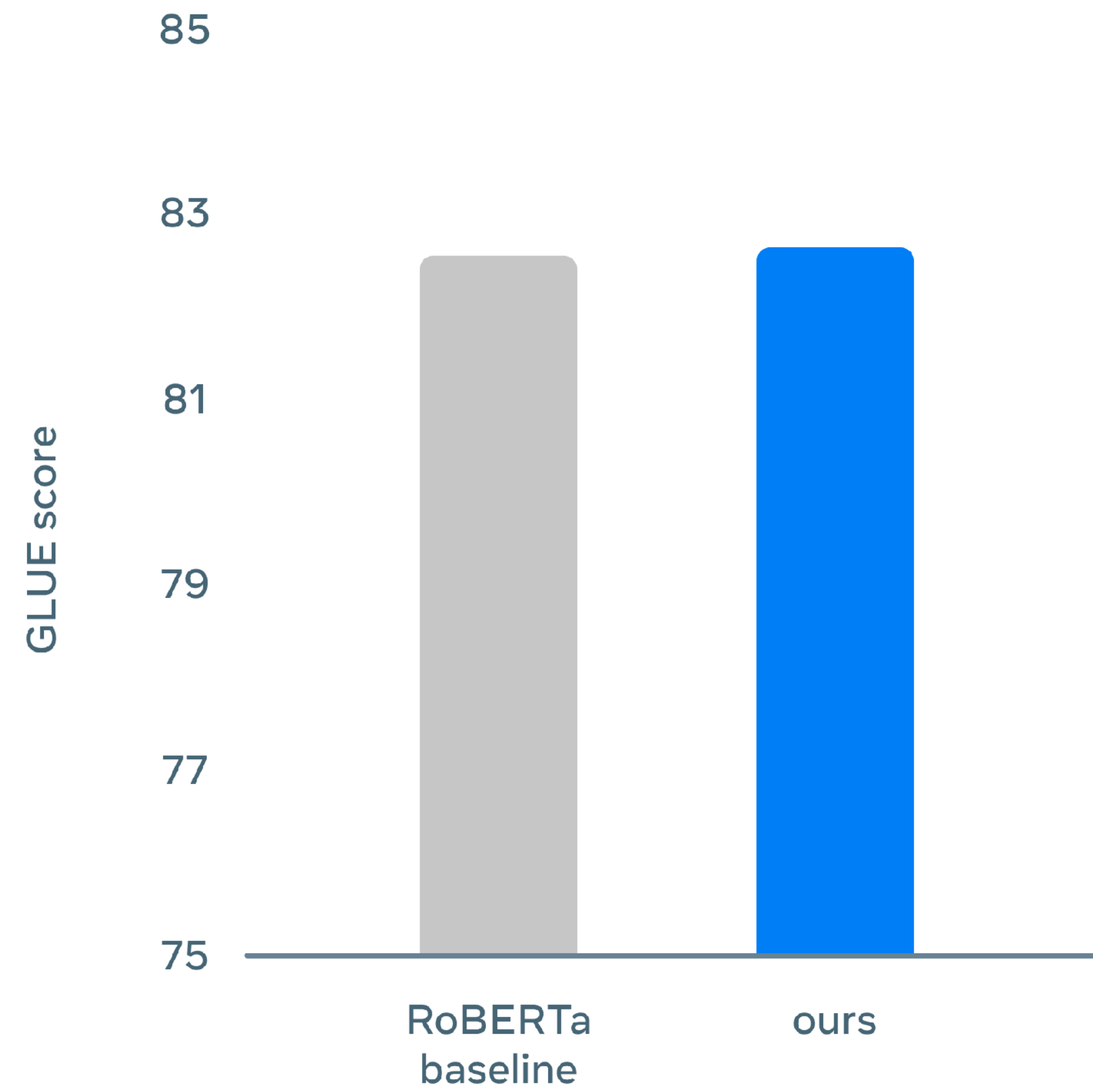
# Vision Results



# Speech & NLP Results

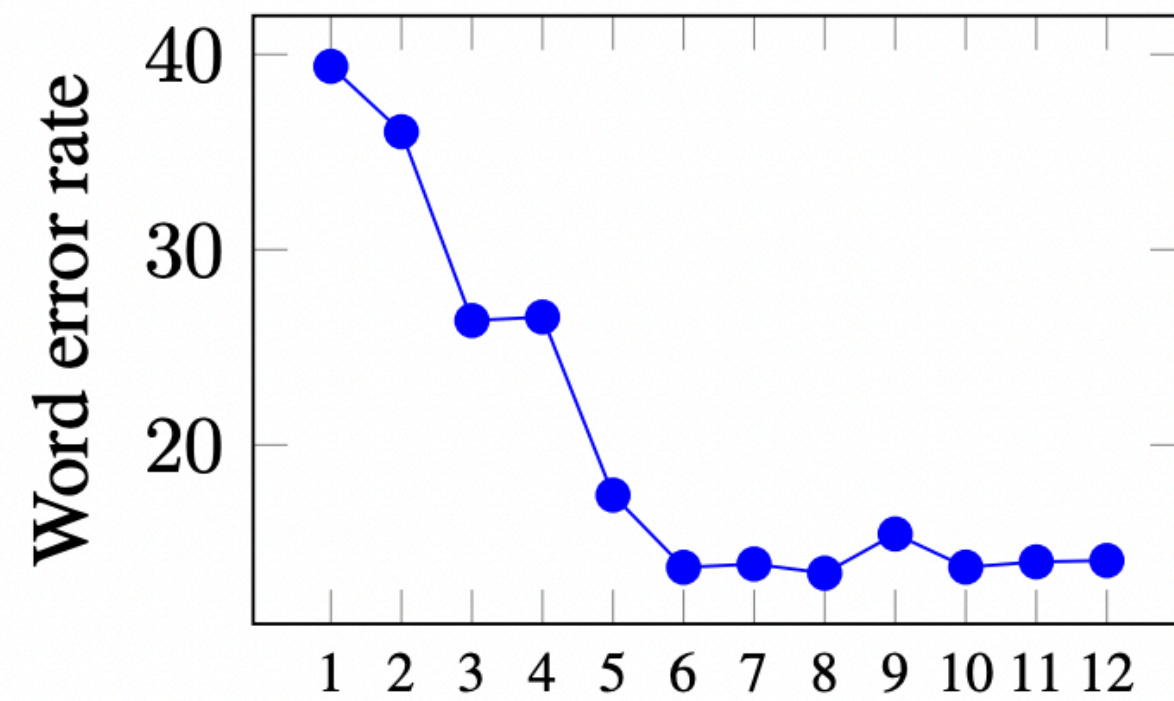
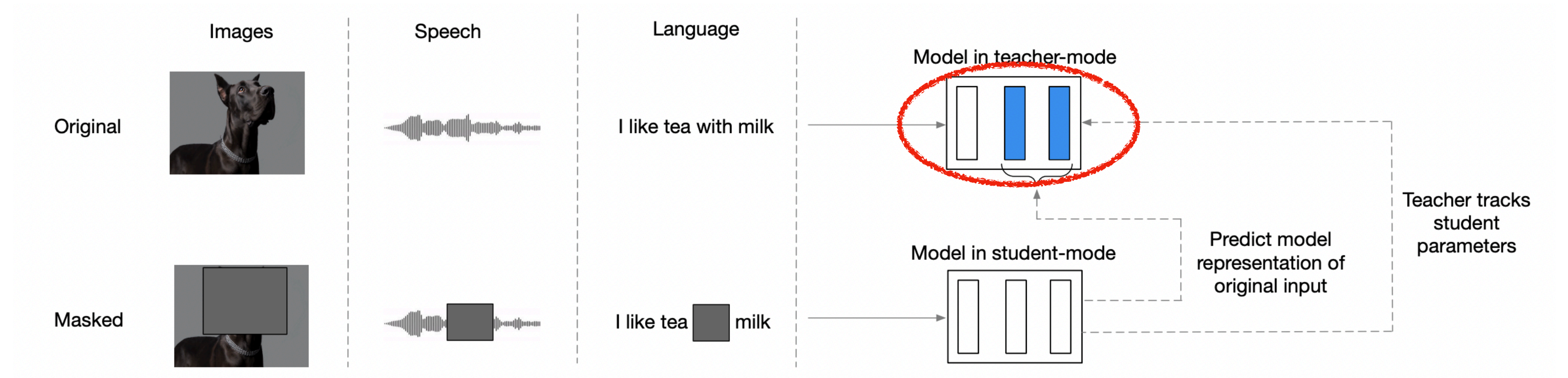


# NLP Results

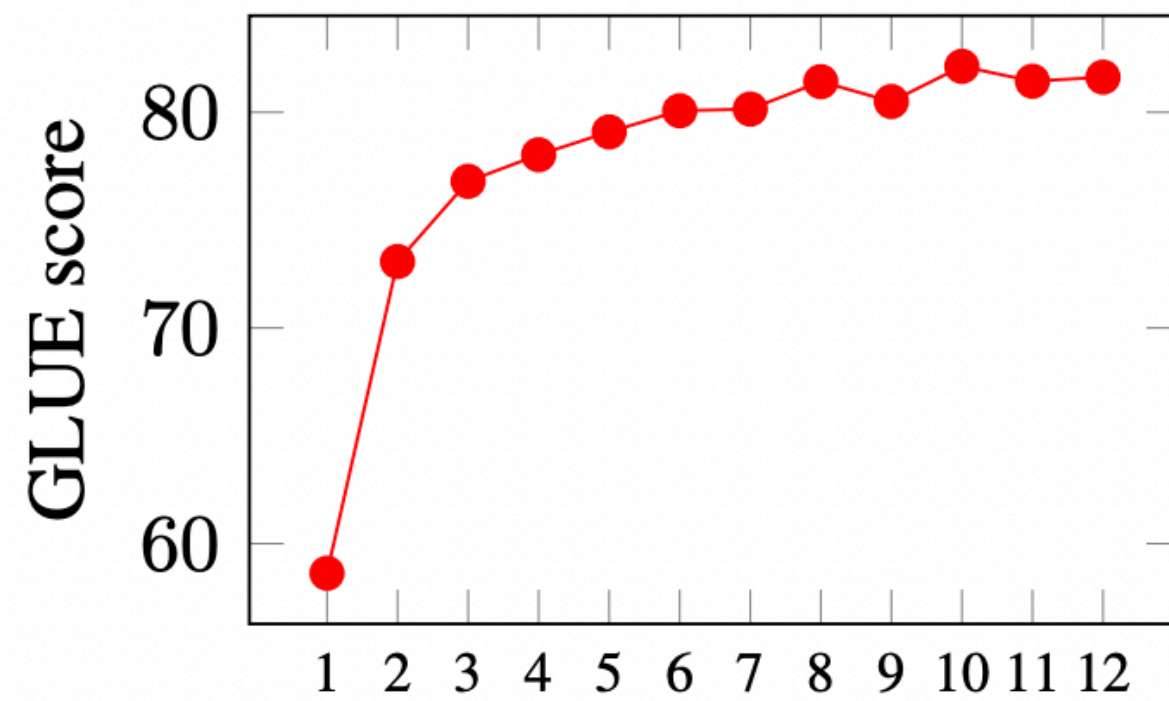




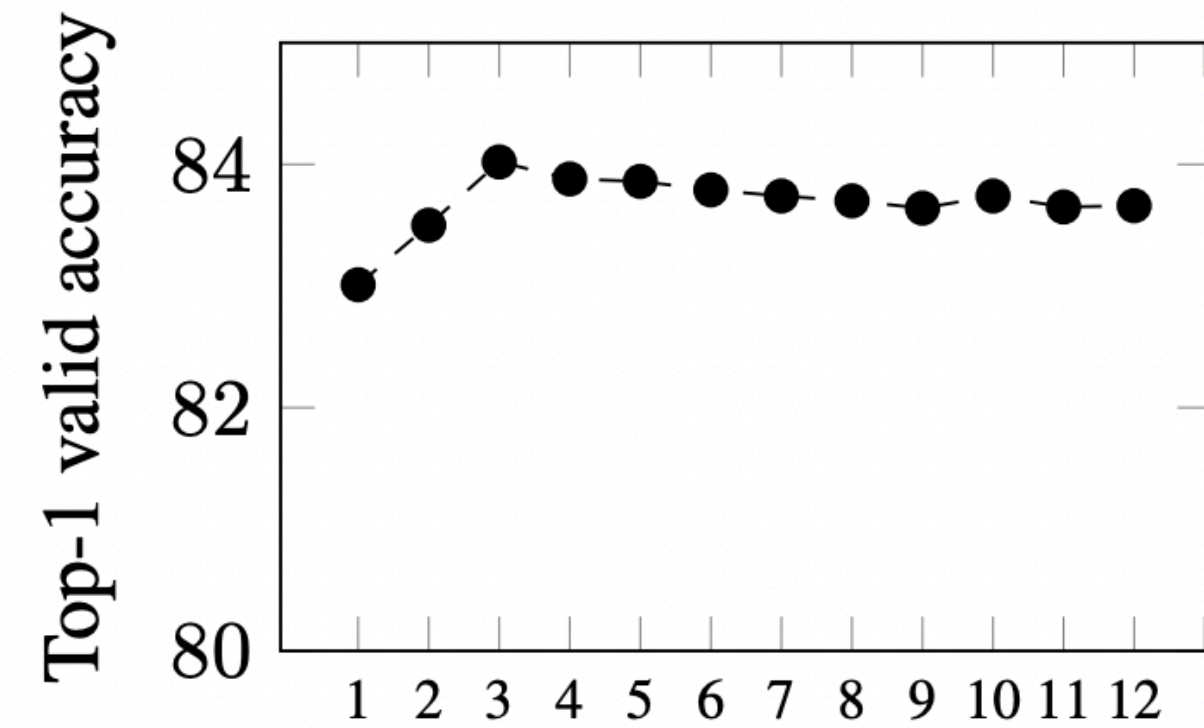
# Teacher Representation Construction



(a) Speech



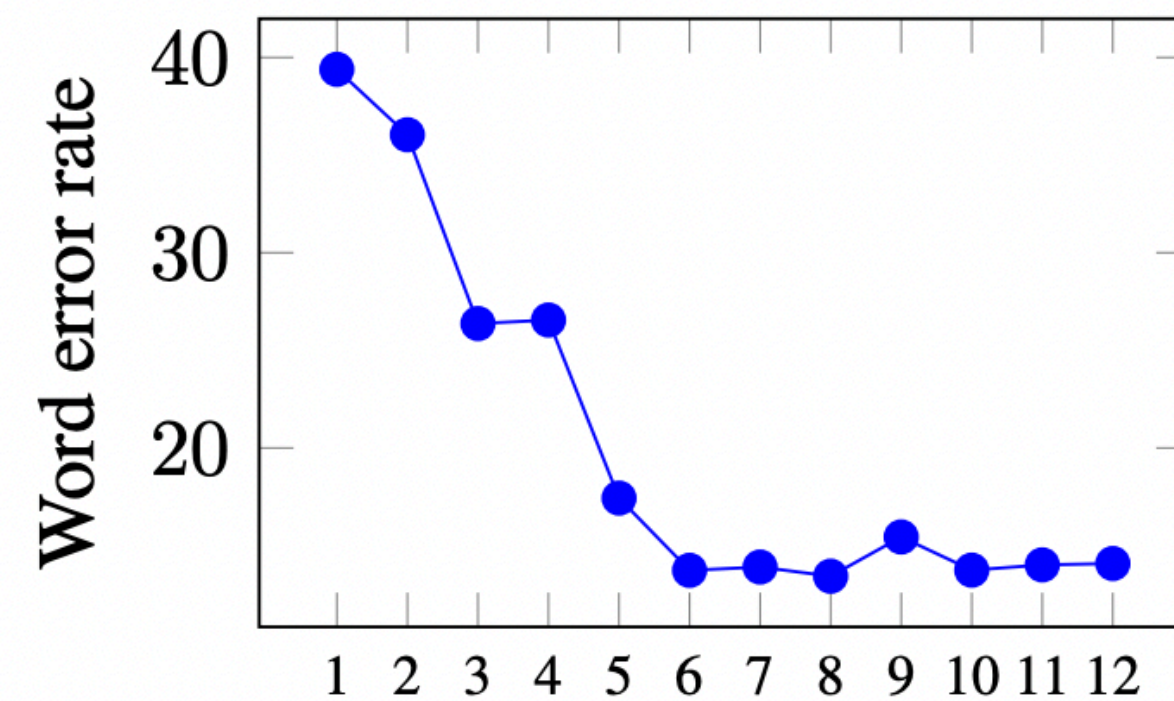
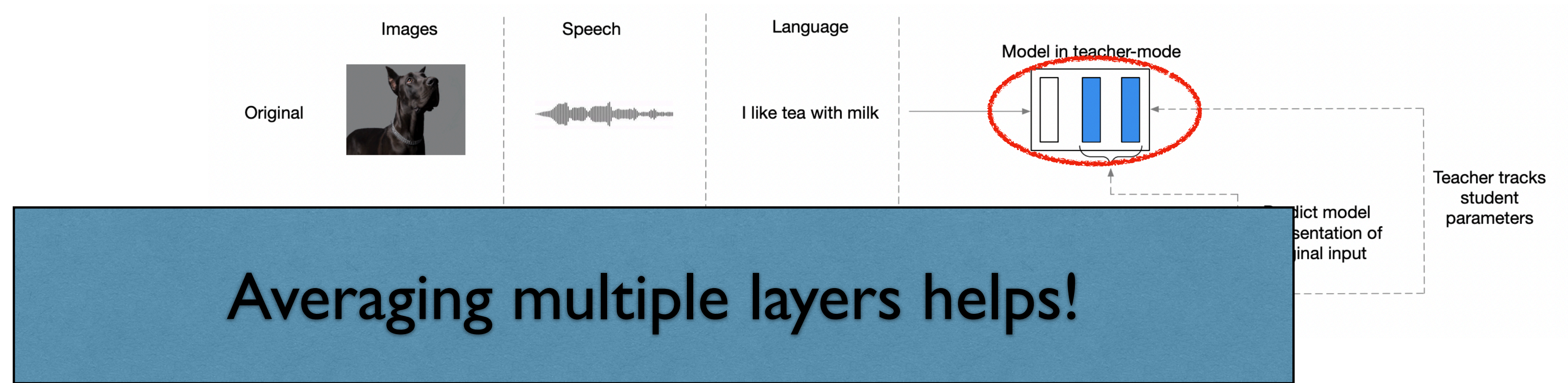
(b) NLP



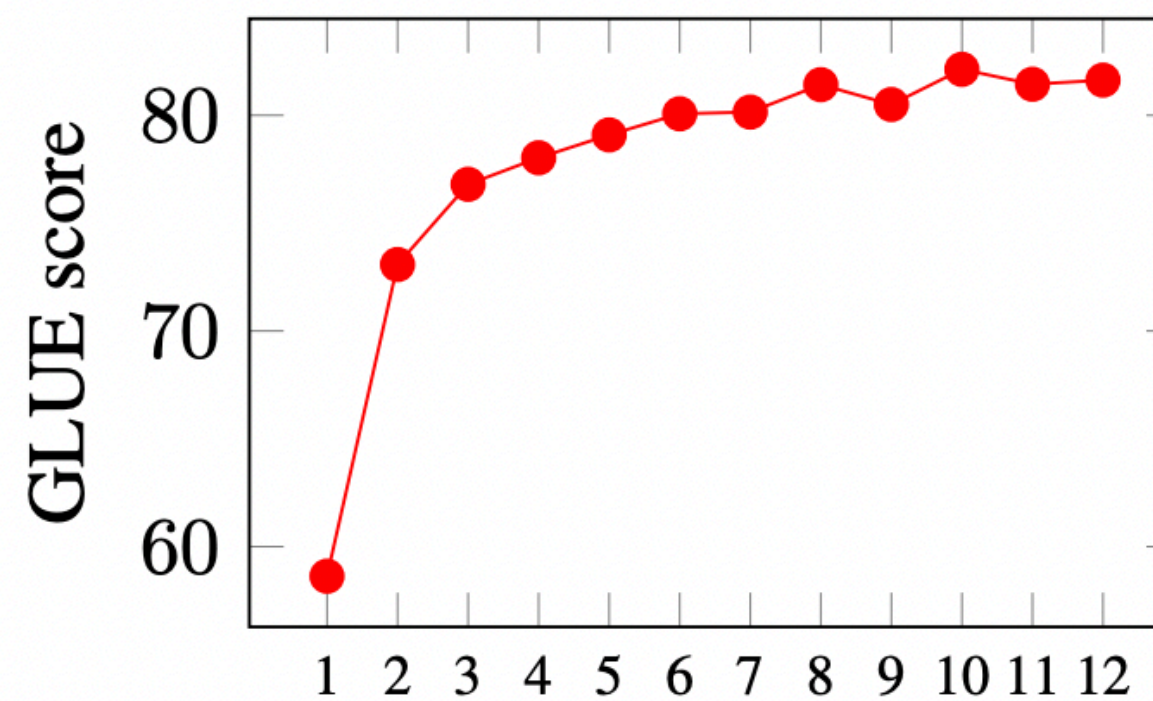
(c) Vision



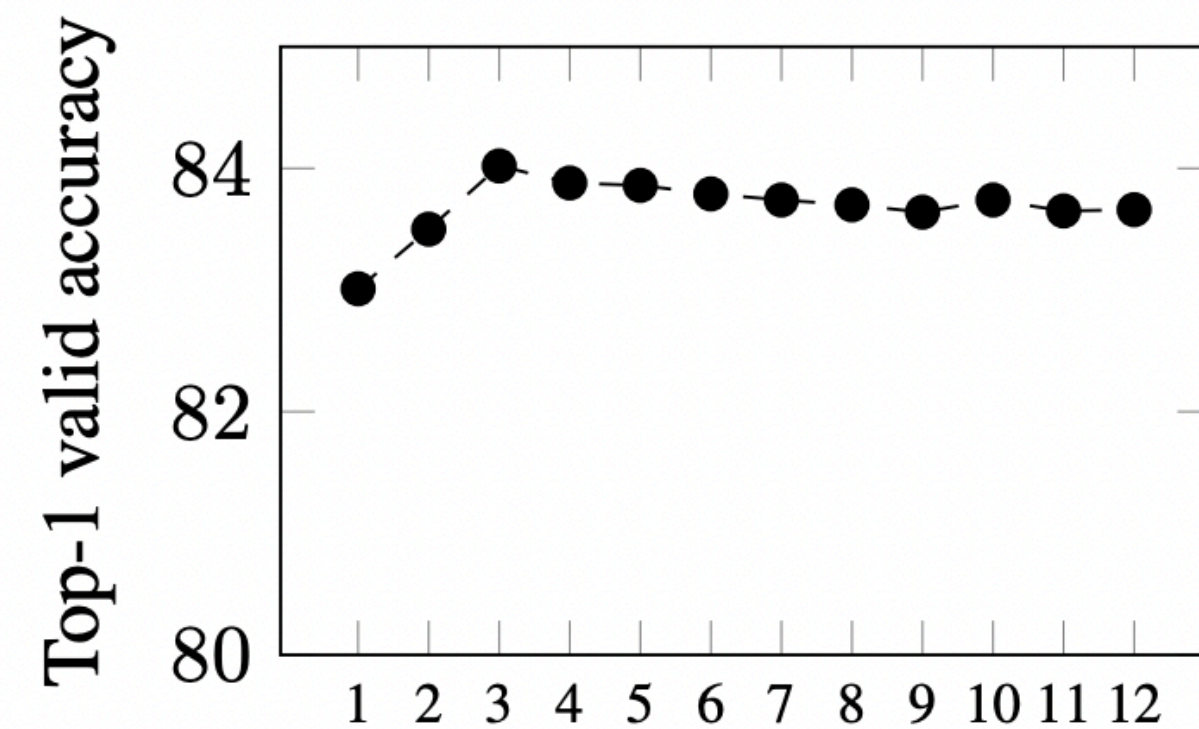
# Teacher Representation Construction



(a) Speech



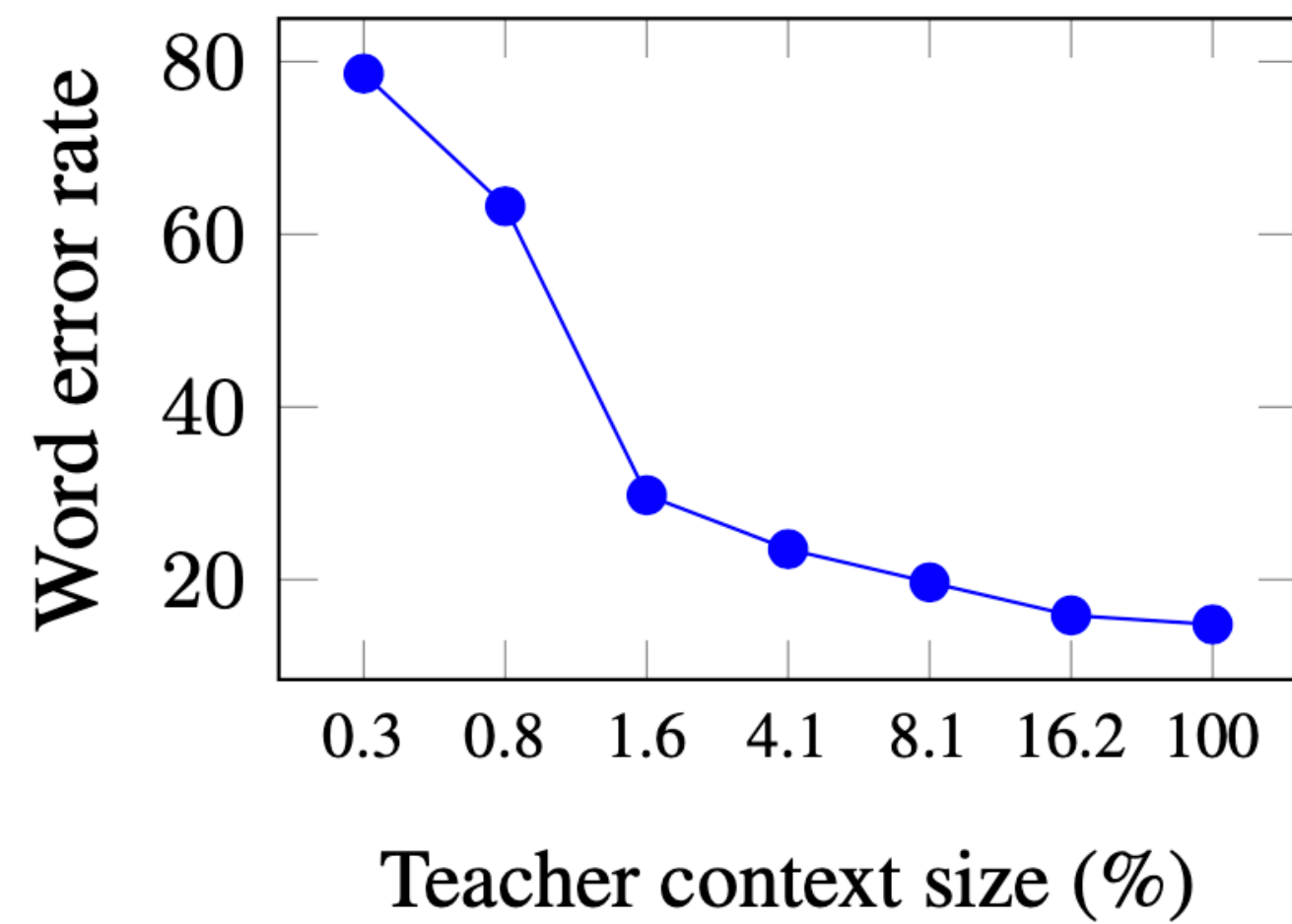
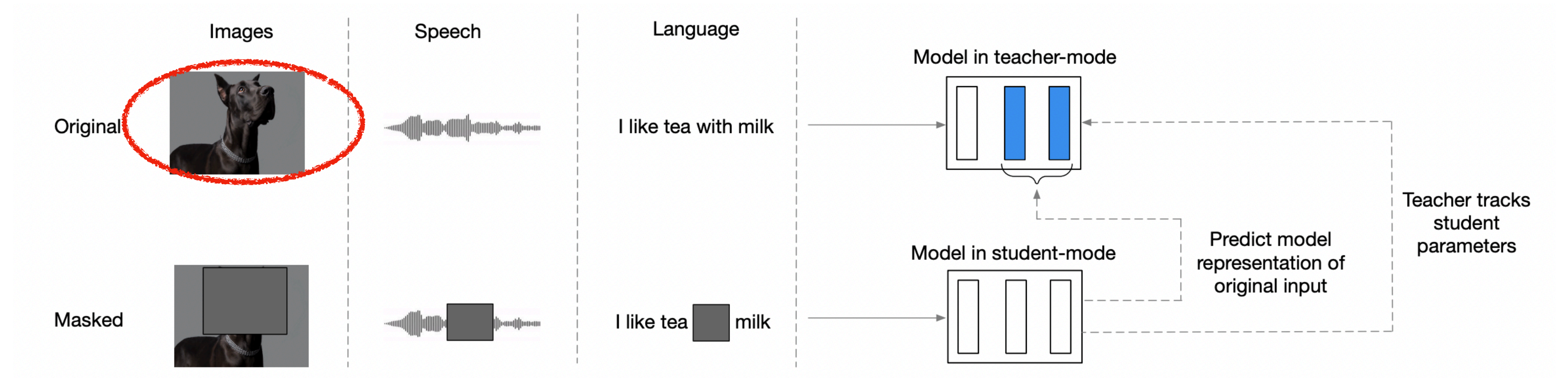
(b) NLP



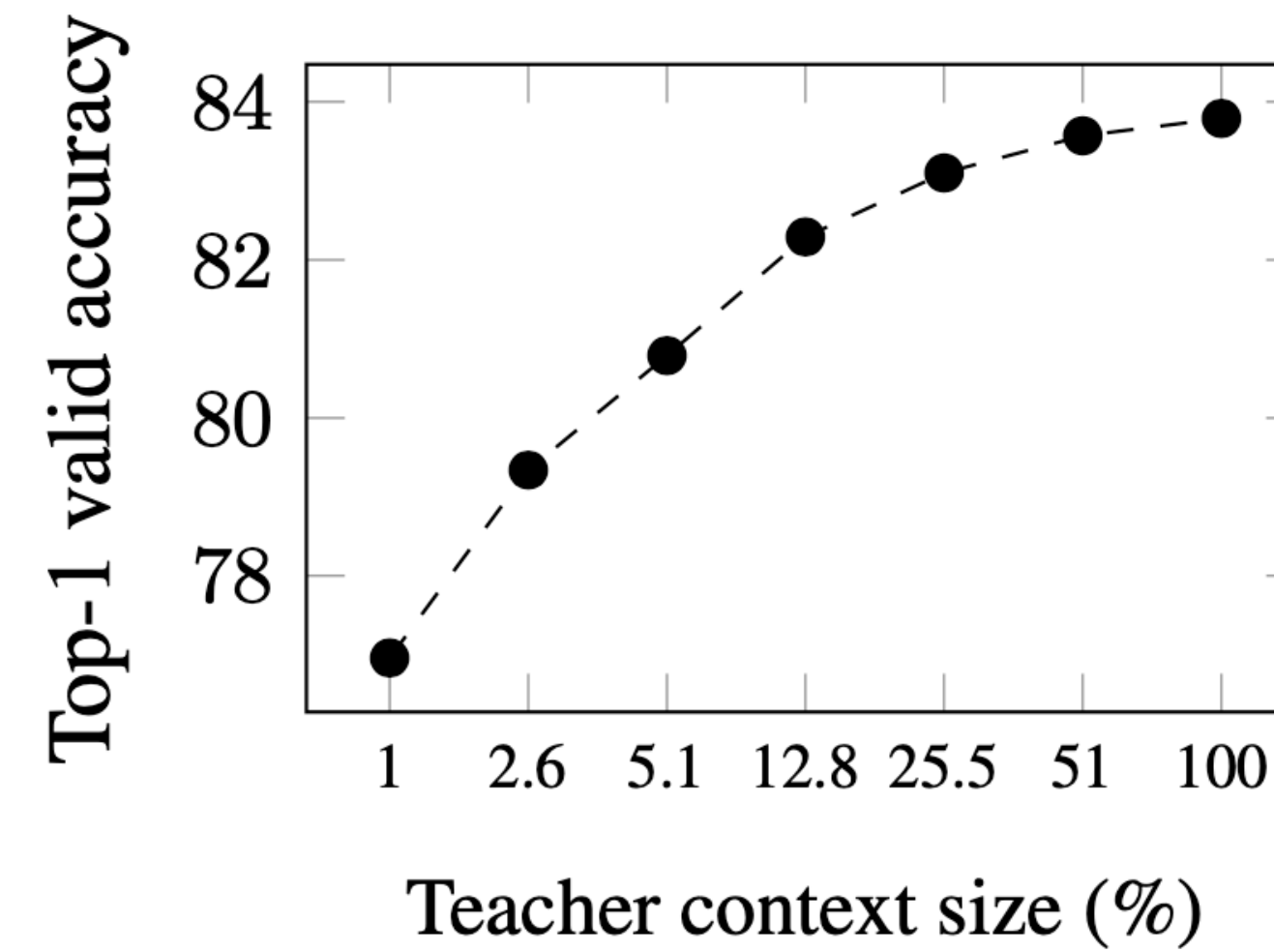
(c) Vision



# Target Context Size



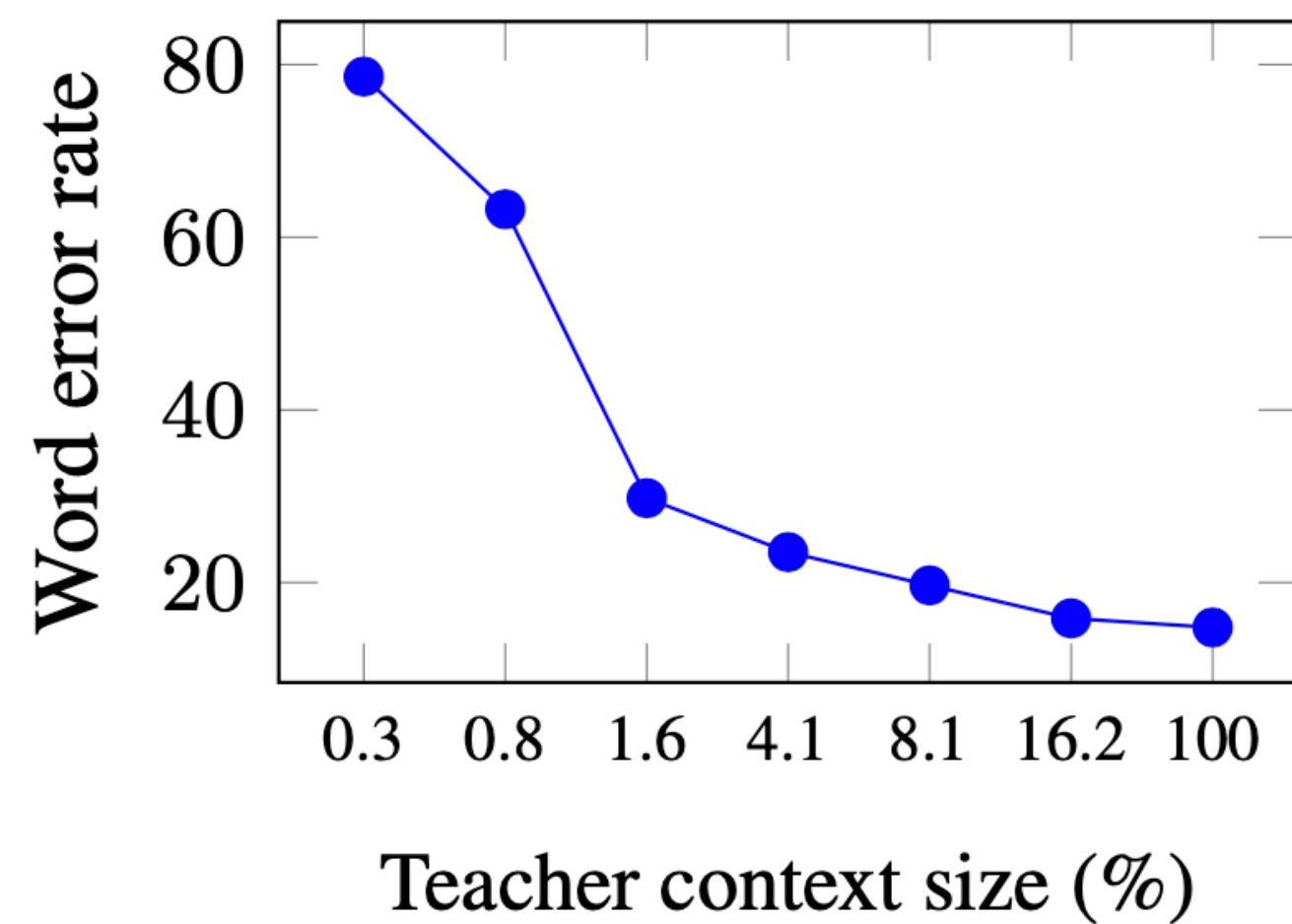
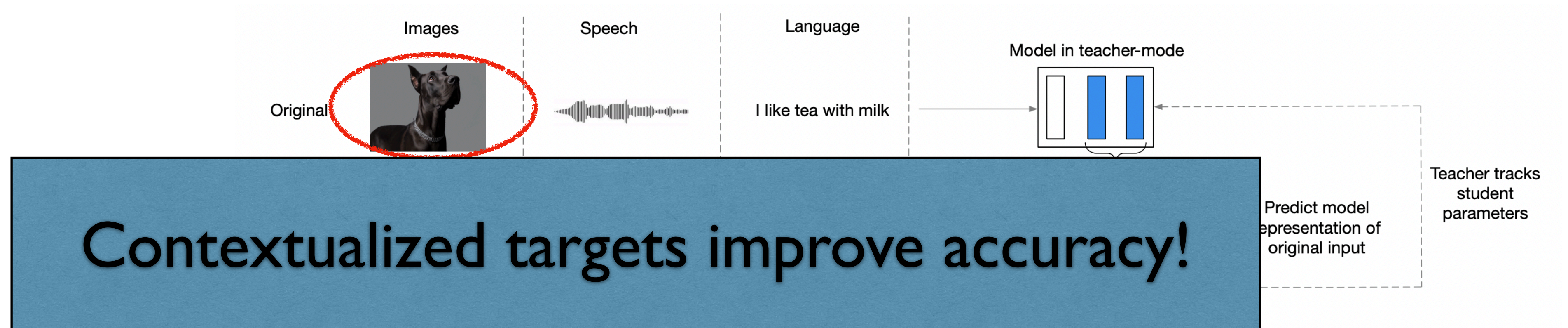
(a) Speech



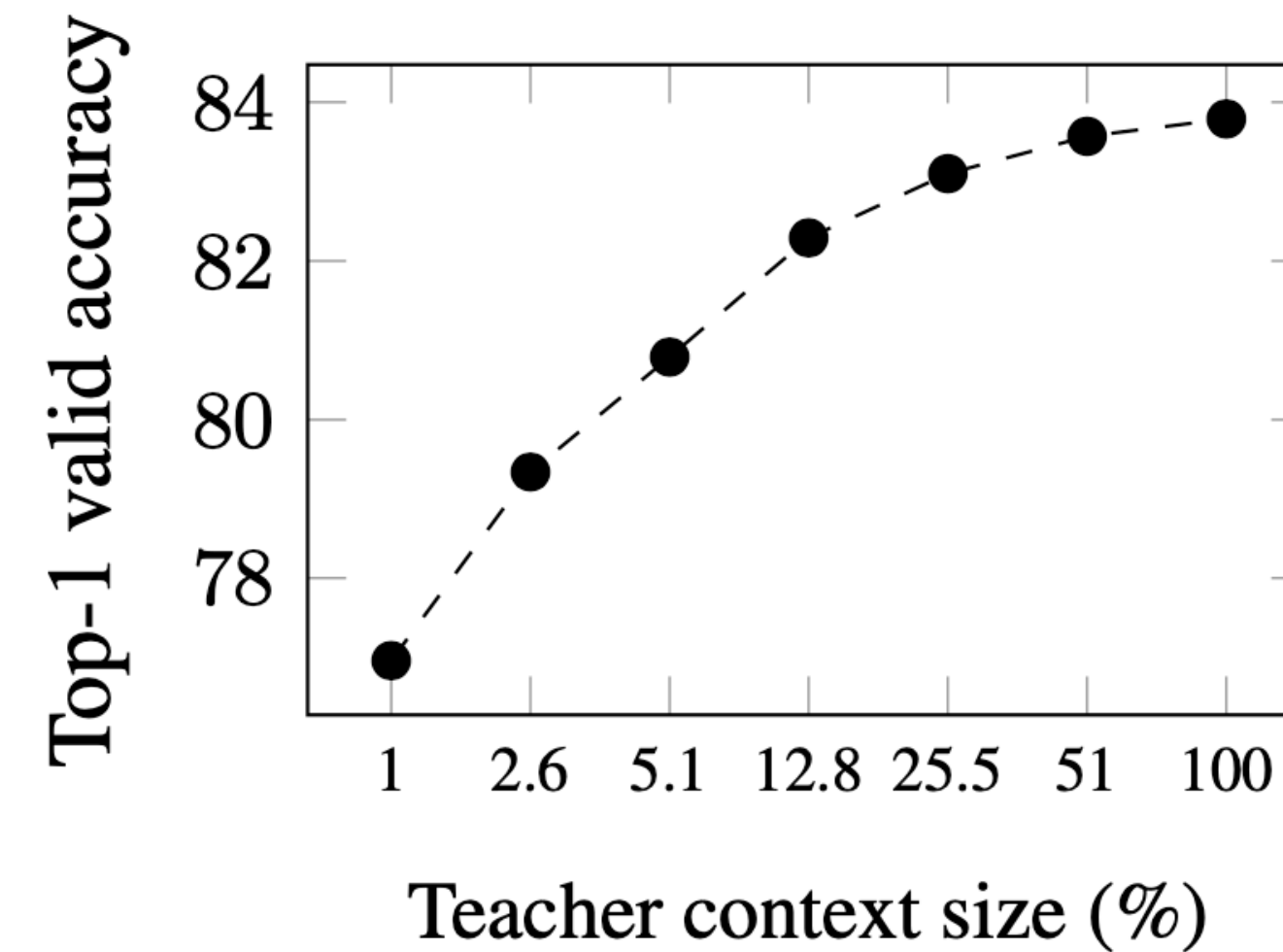
(b) Vision



# Target Context Size



(a) Speech



(b) Vision

# Limitations

- Modality specific feature encoder and masking parameters
- Requires two forward-passes

# Conclusion

- A single learning objective can outperform the best modality-specific algorithms for vision/speech while being competitive on NLP.
- Contextualized targets lead to a rich SSL task and improve performance.
- Future work:
  - Thinks about multiple modalities from the outset
  - unified architectures / objectives (Perceiver IO etc.)



# Thank you



Arun Babu



Alexis  
Conneau



Steffen  
Schneider



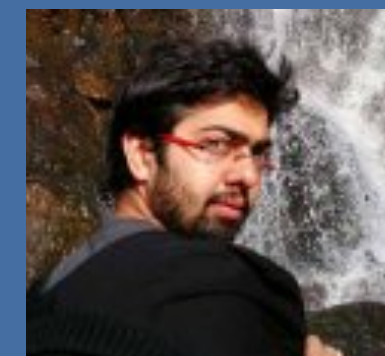
Henry Zhou



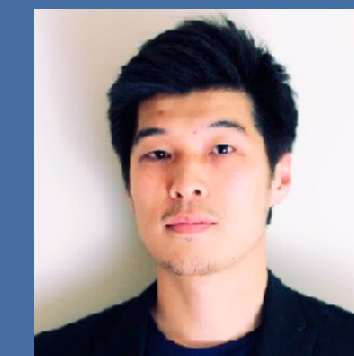
Abdelrahman  
Mohamed



Jiatao Gu



Naman  
Goyal



Wei-Ning Hsu



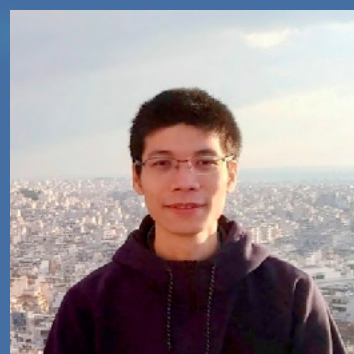
Alexei Baevski



Michael Auli



Kushal  
Lakhotia



Andros Tjandra



Kritika Singh



Yatharth Saraf



Geoffrey Zweig



Qiantong Xu



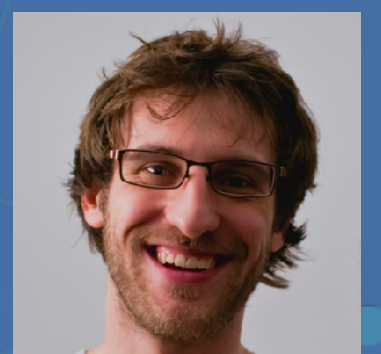
Tatiana  
Likhomanenko



Paden  
Tomasello



Ronan  
Collobert



Gabriel  
Synnaeve