

Evaluating Off-the-Shelf Machine Listening and Natural Language Models for Automated Audio Captioning

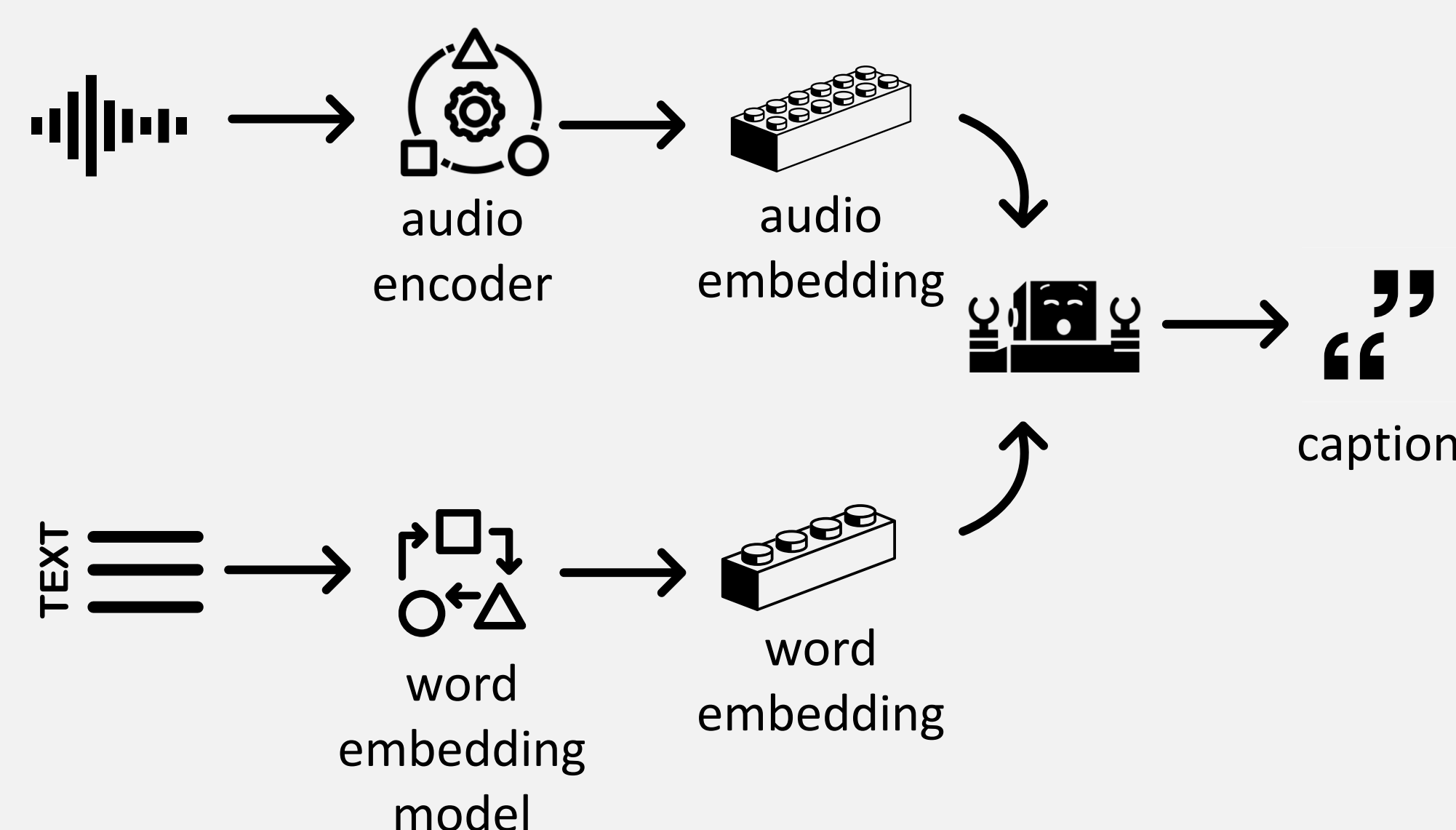
PRESENTER:
Benno Weck

BACKGROUND:

Find resources useful for Audio Captioning by comparing pre-trained models from Natural Language Processing and Machine Listening.

METHODS

1. Transformer-decoder as caption generator
2. Audio embeddings are fixed.
3. Word embeddings: fine-tuned & fixed.
4. Try different adapters on top of the audio embeddings and different hop sizes.
5. Test all combinations and compare by SPIDER score.



RESULTS

- YAMNet is the best audio encoder in our comparison
- BERT is the best word embedding model in our comparison.
- Pre-trained word embeddings work better than randomly initialized.
- They work even better when fine-tuned!
- Extracting overlapping audio embeddings gives a small boost in performance.

YAMNet + BERT are the best pre-trained audio and language models for audio captioning.*

**conditions apply*

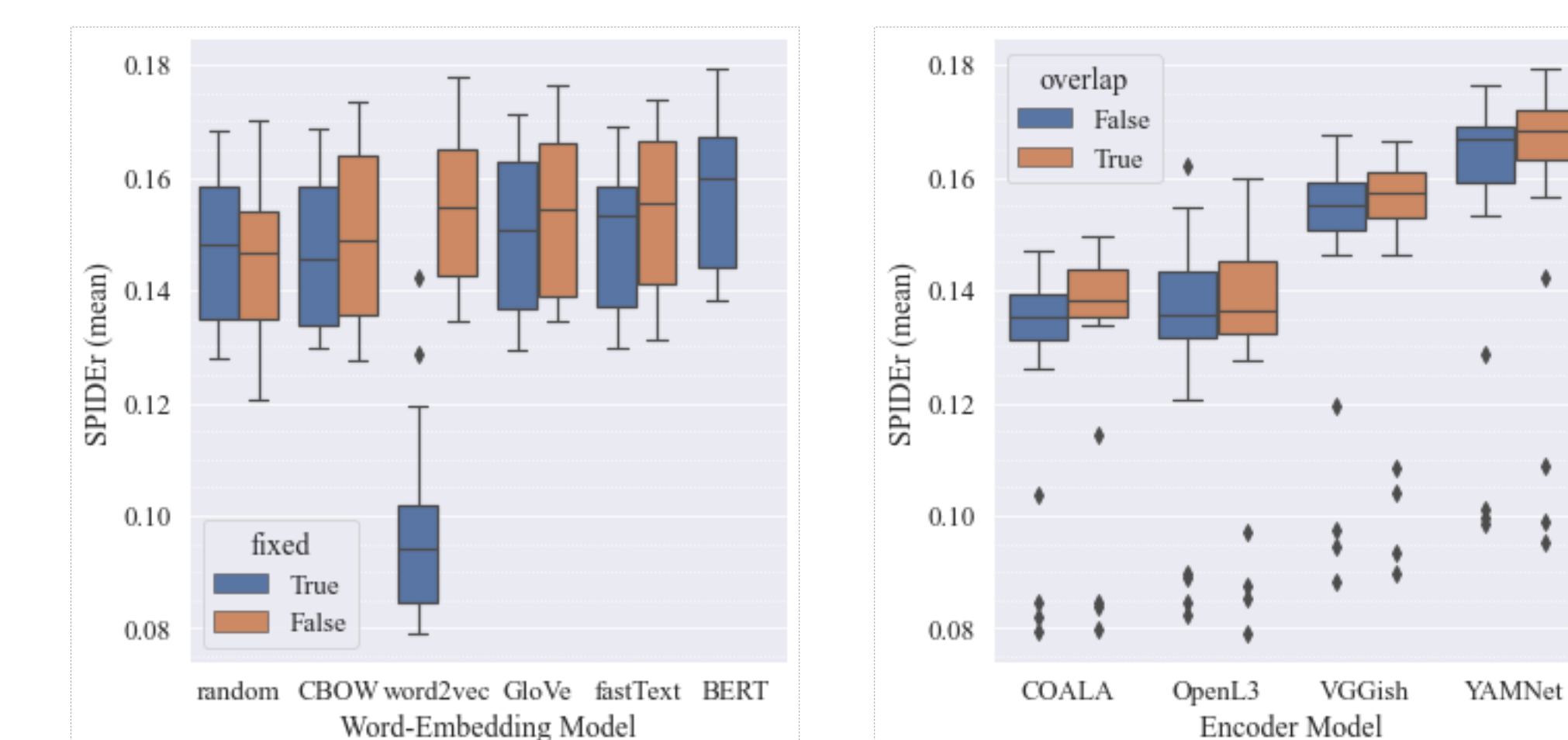
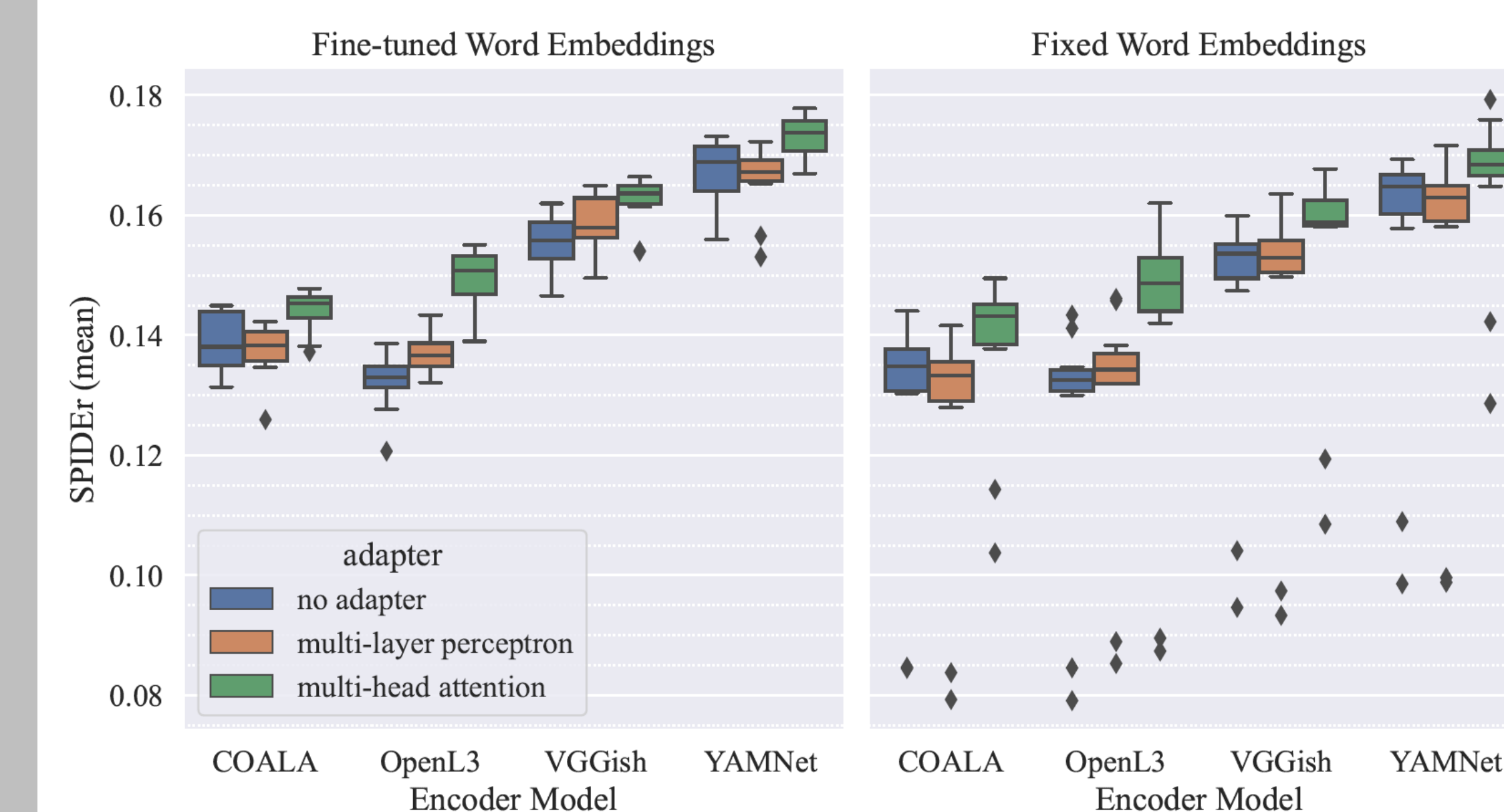
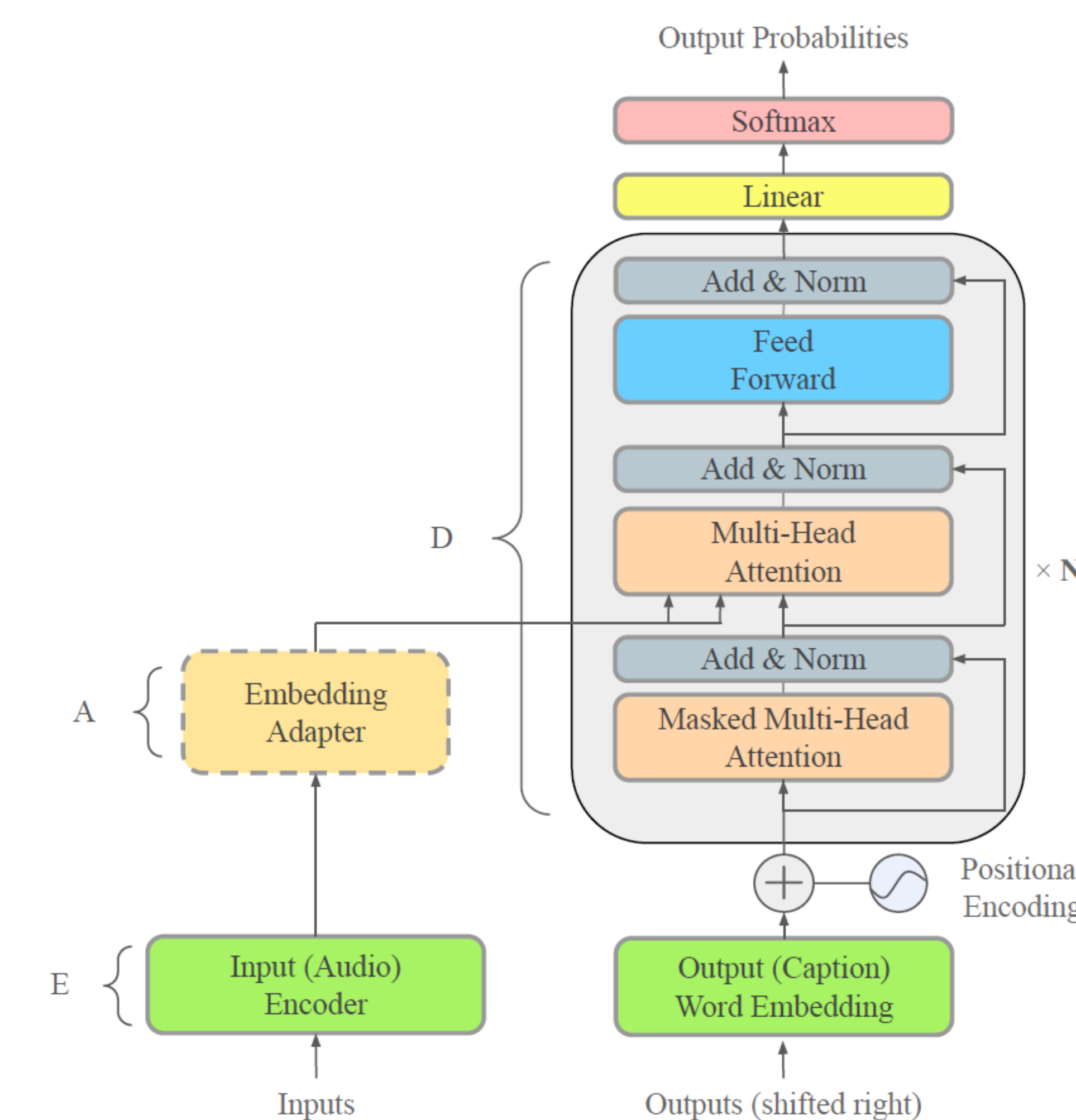
Pre-trained embeddings + Transformer =
 YAMNet + BERT =

- 1 YAMNet
- 2 VGGish
- 3 OpenL3
- COALA

adapted
fine-tuned

- 1 BERT
- 2 fastText
- 3 GloVe
- word2vec
- randomly initialized

fine-tune!



Encoder	Word embedding	Adapter	SPIDER	
			Mean	SD
COALA [†]	BERT	MHA-based	0.1495	0.0044
OpenL3 [*]	BERT	MHA-based	0.1620	0.0051
VGGish [*]	BERT	MHA-based	0.1677	0.0052
YAMNet [†]	BERT	MHA-based	0.1793	0.0066

Benno Weck, Xavier Favory, Konstantinos Drossos, Xavier Serra

