

Contrastive Unsupervised Learning for Speech Emotion Recognition

Mao Li², Bo Yang¹, Joshua Levy¹, Andreas Stolcke¹, Viktor Rozgic¹,
Spyros Matsoukas¹, Constantinos Papayiannis¹, Daniel Bone¹, Chao Wang¹

Amazon Alexa¹ and University of Illinois at Chicago²



Motivation

- Modern speech systems are mainly designed for **speech content understanding**, while **speech emotion recognition (SER)** becomes a key technology to enable natural human-machine communication.
- Application scenarios of SER system:
 - voice assistant
 - human health assistant
 - chat-bots & social robot



[Image credit: Patrick J. Kiger's blog]

Motivation

- Modern speech systems are mainly designed for **speech content understanding**, while **speech emotion recognition (SER)** becomes a key technology to enable natural human-machine communication.
- **The deficiency of emotion annotated data** is the bottleneck for development of SER system.
 - labeling is expensive
 - labeling emotion data is challenging due to annotator disagreements

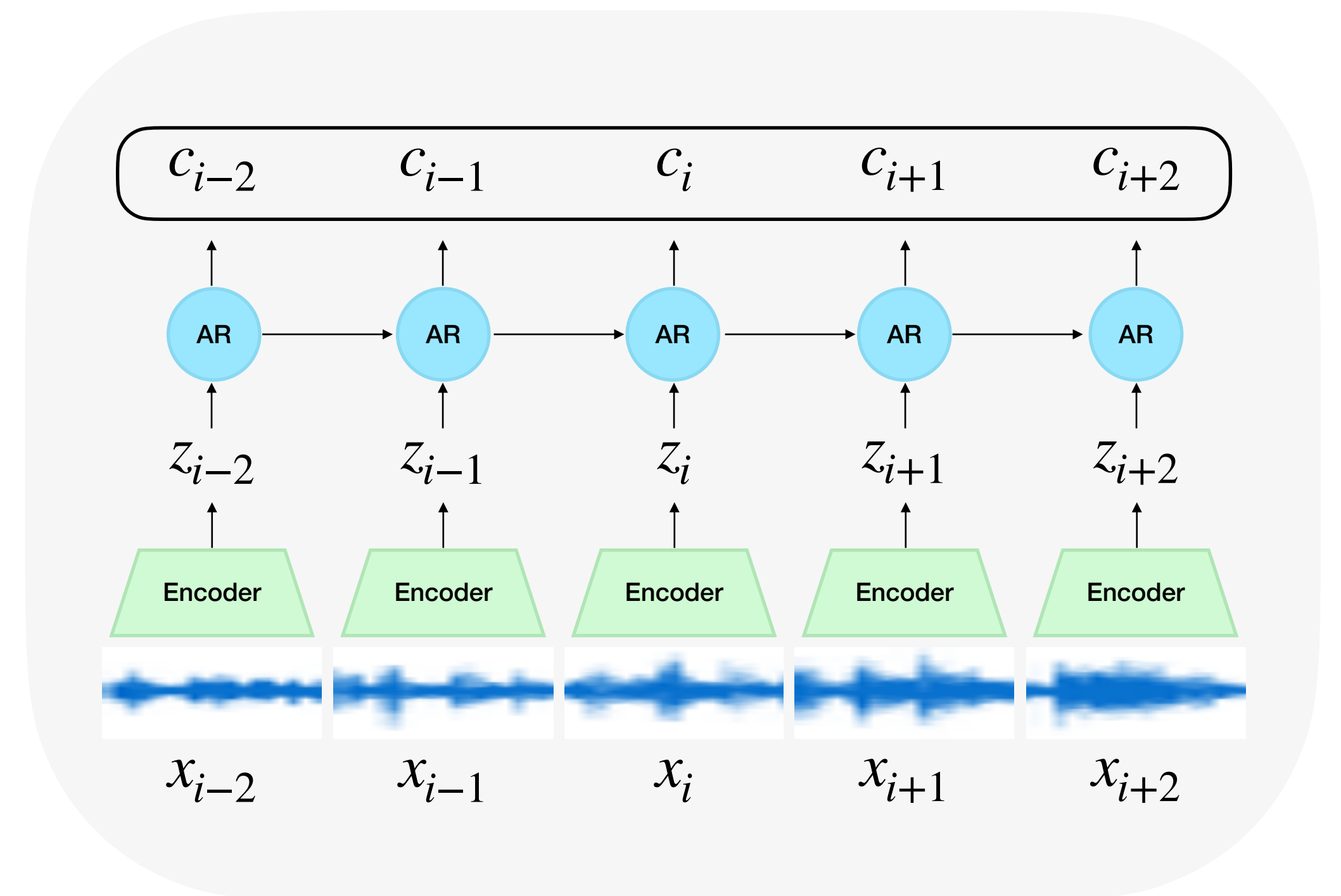
Motivation

- Modern speech systems are mainly designed for **speech content understanding**, while **speech emotion recognition (SER)** becomes a key technology to enable natural human-machine communication.
- The **deficiency of emotion annotated data** is the bottleneck for development of SER system.
- Two solutions:
 - Transfer learning from a related speech task
 - Unsupervised representation learning

Proposed Method

Contrastive Predictive Coding (CPC)

- sequence of audio frames: $\{x_1, x_2, \dots, x_n\}$
- non-linear encoder f
- frame-level representation: $z_i = f(x_i)$
- autoregressive model g
- contextual representation: $c_i = g(z_{\leq i})$

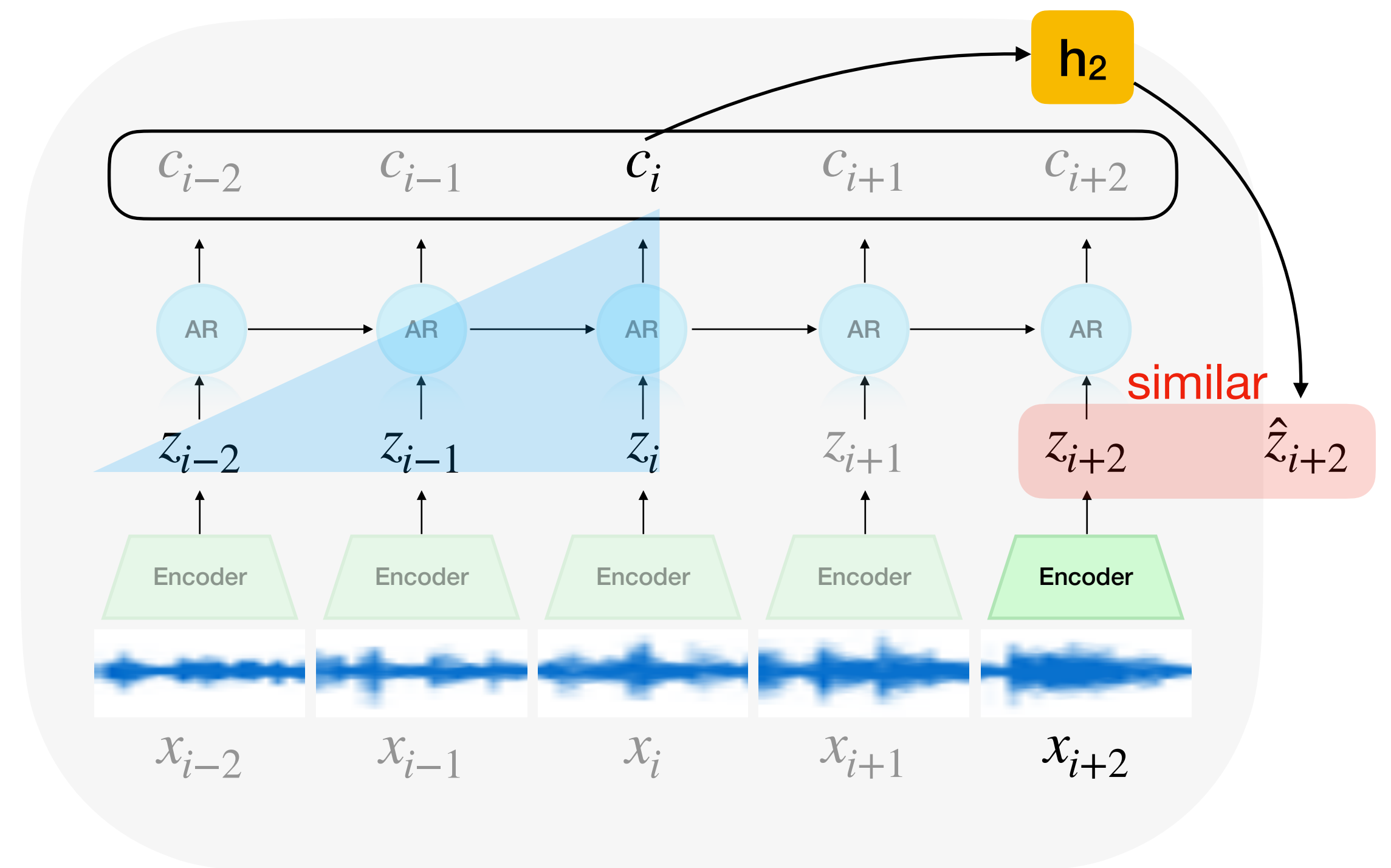


Proposed Method

Contrastive Predictive Coding (CPC)

- frame-level representation: $z_i = f(x_i)$
- contextual representation: $c_i = g(z_{\leq i})$
- prediction function for a specific k: h_k
- predict future : $\hat{z}_{i+k} = h_k(c_i) = h_k(g(z_{\leq i}))$
- InfoNCE Loss:

$$\mathcal{L} = - \sum_{m=1}^k \left[\log \frac{\exp(\hat{z}_{i+m}^\top z_{i+m})}{\exp(\hat{z}_{i+m}^\top z_{i+m}) + \sum_{i=1}^{N-1} \exp(\hat{z}_{i+m}^\top z_i)} \right]$$



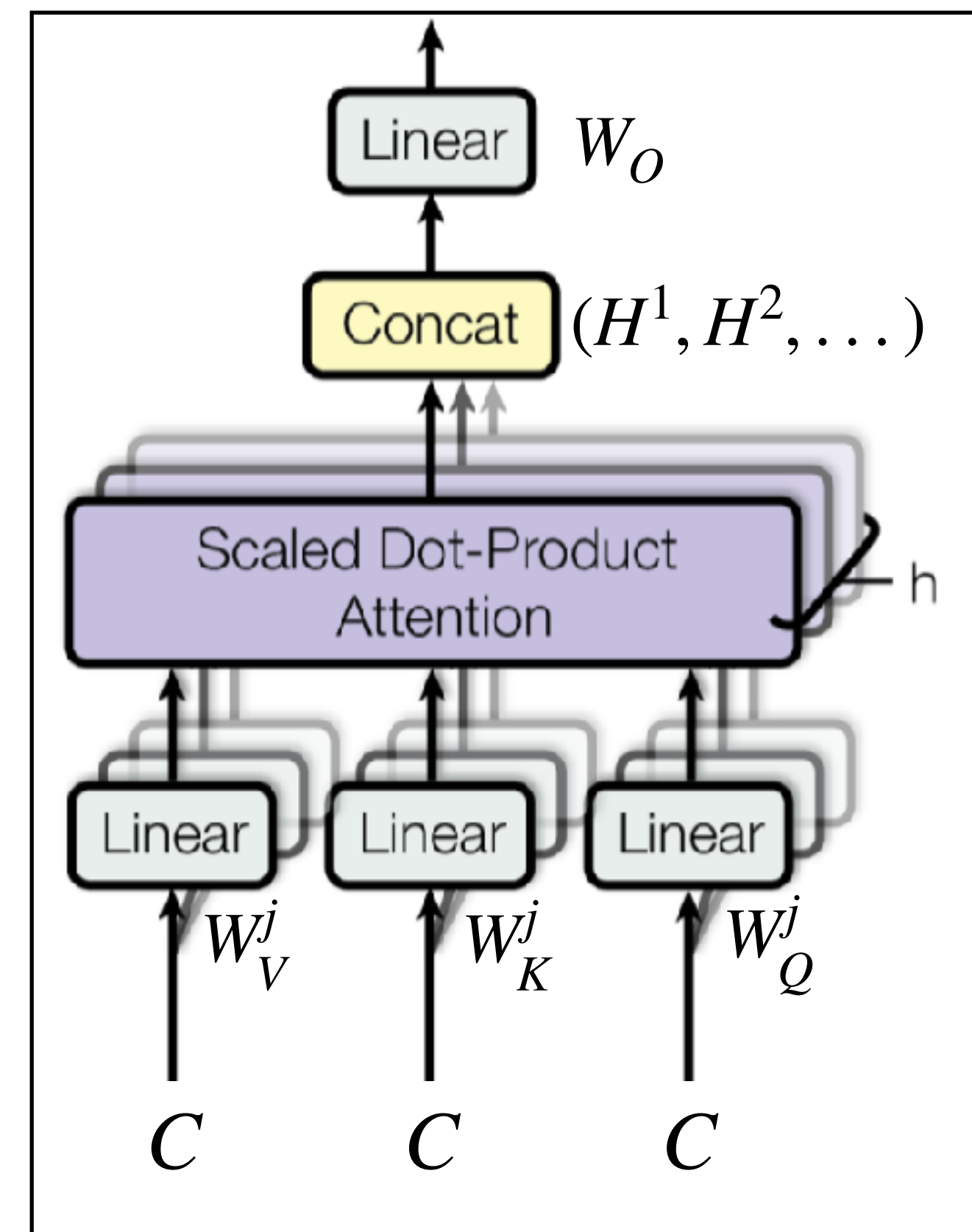
Proposed Method

Attention-based Emotion Recognizer

- Multi-head attention
 - C is the output of CPC
 - W_*^j, W_O are trainable weights
 - H^j is attention score of a single head
 - h is number of heads, d_K is the dimension

$$H^j = \text{softmax} \left(W_Q^j C \left(W_K^j C \right)^\top / \sqrt{d_K} \right) W_V^j C$$
$$U = \text{Concat} (H^1, H^2, \dots, H^n) W_O$$

Multi-head Attention



Proposed Method

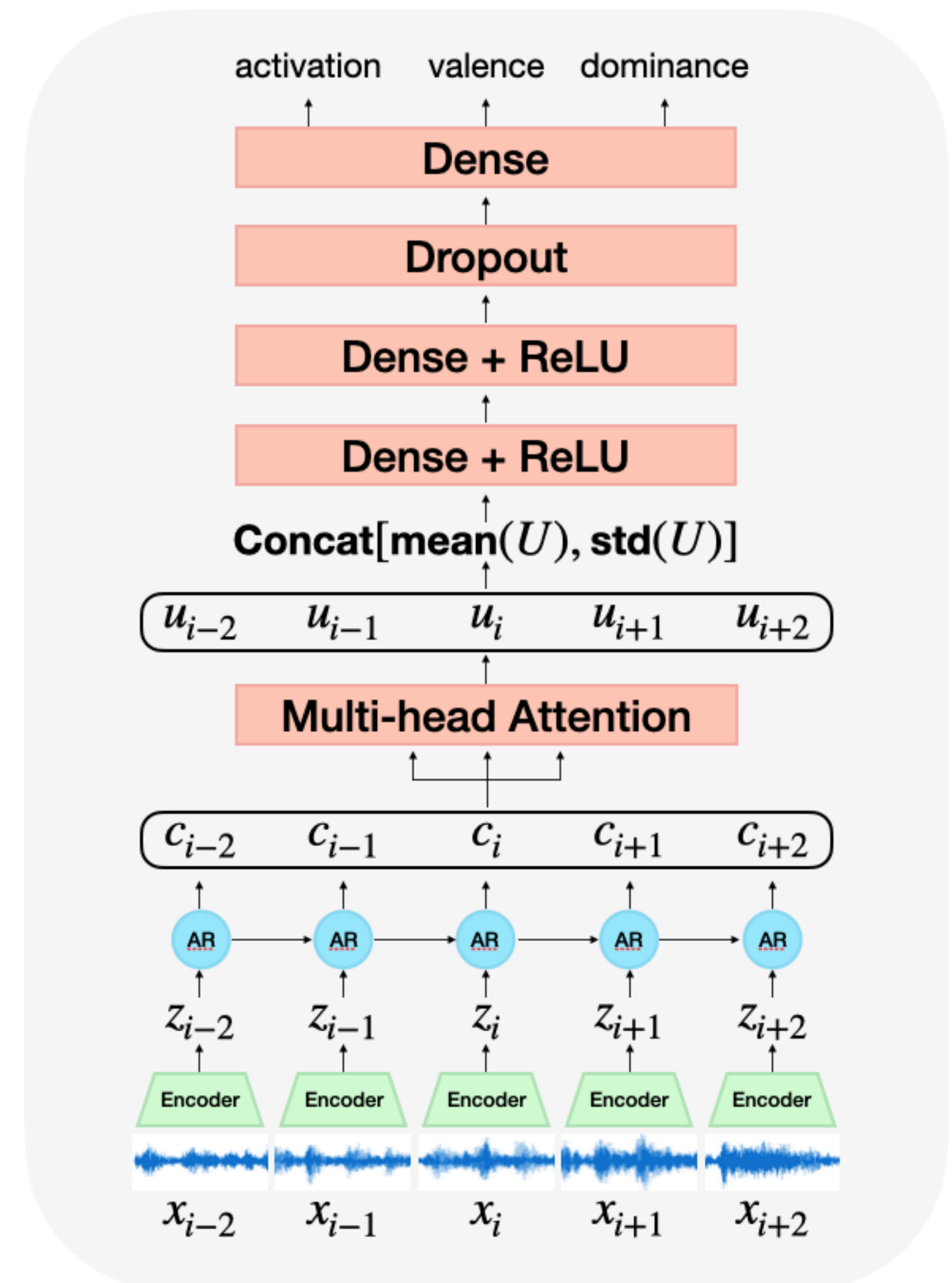
Attention-based Emotion Recognizer

- Utterance embedding $u = [mean(U); std(U)]$
- Concordance Correlation Coefficient (CCC) measures alignment of two random variables:

$$CCC(X, Y) = \rho \frac{2\sigma_X\sigma_Y}{\sigma_X^2 + \sigma_Y^2 + (\mu_X - \mu_Y)^2}, \rho = \frac{\sigma_{XY}}{\sigma_X\sigma_Y}$$

- X: ground truth score; Y: predicted score
- Loss function:

$$\mathcal{L} = 1 - \alpha CCC_{act} - \beta CCC_{val} - (1 - \alpha - \beta) CCC_{dom}$$



Experiments

Datasets

Dataset	Label	Duration	Usage
LibriSpeech	N/A	100 hours of audio	Unsupervised pre-training
IEMOCAP	Primitive labels; Categorical labels	12 hours of audio	Model evaluation; Visualization
MSP-Podcast	Primitive labels; Categorical labels	84 hours of audio	Model evaluation

Experiments

Experimental Setups

- preCPC: pre-trained CPC (LibriSpeech) + supervised (IEMOCAP/MSP-Podcast)



- Sup: supervised only (IEMOCAP/MSP-Podcast)



Experiments

Experimental Setups

- preCPC: pre-trained CPC + supervised
- Sup: supervised only

Hypothesis:

- Representations learned by CPC are superior to handcrafted features for speech emotion recognition task

Table 1: CCC scores (mean/std) on the IEMOCAP dataset

Methods	CCC _{avg}	CCC _{act}	CCC _{val}	CCC _{dom}
Sup	.664 ± .007	.638 ± .017	.718 ± .004	.635 ± .009
jointCPC	.562 ± .012	.549 ± .032	.642 ± .013	.491 ± .016
miniCPC	.660 ± .005	.673 ± .028	.702 ± .009	.606 ± .019
preCPC	.731 ± .003	.752 ± .014	.752 ± .009	.691 ± .009

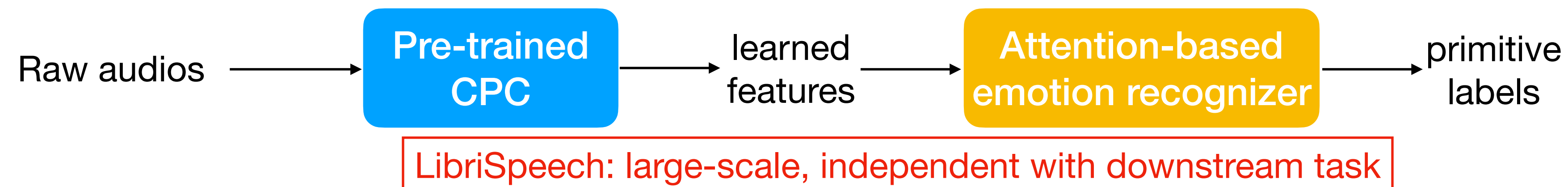
Table 2: CCC scores (mean/std) on the MSP-Podcast dataset

Methods	CCC _{avg}	CCC _{act}	CCC _{val}	CCC _{dom}
Sup	.458 ± .005	.596 ± .007	.266 ± .004	.501 ± .013
jointCPC	.491 ± .008	.628 ± .006	.280 ± .006	.568 ± .007
miniCPC	.549 ± .006	.688 ± .009	.345 ± .005	.615 ± .011
preCPC	.571 ± .004	.706 ± .006	.377 ± .008	.639 ± .012

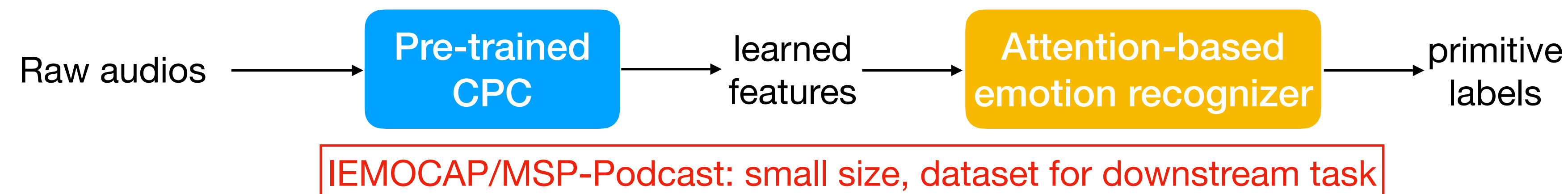
Experiments

Experimental Setups

- preCPC: pre-trained CPC (LibriSpeech) + supervised (IEMOCAP/MSP-Podcast)



- Sup: supervised only (IEMOCAP/MSP-Podcast)
- miniCPC: pre-trained CPC (IEMOCAP/MSP-Podcast)+ supervised (IEMOCAP/MSP-Podcast)



Experiments

Experimental Setups

- preCPC: pre-trained CPC + supervised
- Sup: supervised only
- miniCPC: pre-trained CPC + supervised

Hypothesis:

- Exposing the model to more diverse acoustic conditions and speaker variations (LibriSpeech) is beneficial for learning robust features.

Table 1: CCC scores (mean/std) on the IEMOCAP dataset

Methods	CCC _{avg}	CCC _{act}	CCC _{val}	CCC _{dom}
Sup	.664 ± .007	.638 ± .017	.718 ± .004	.635 ± .009
jointCPC	.562 ± .012	.549 ± .032	.642 ± .013	.491 ± .016
miniCPC	.660 ± .005	.673 ± .028	.702 ± .009	.606 ± .019
preCPC	.731 ± .003	.752 ± .014	.752 ± .009	.691 ± .009

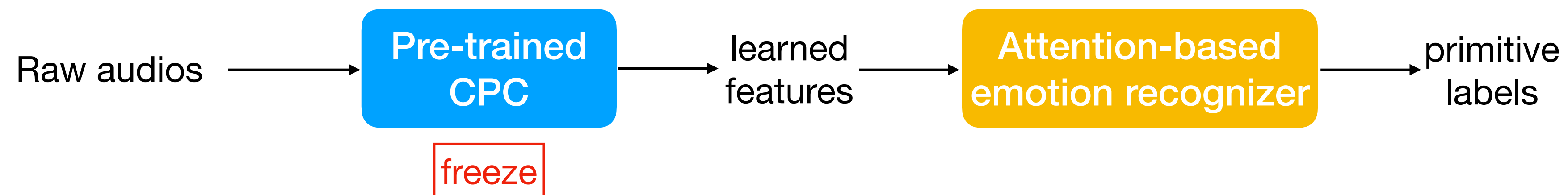
Table 2: CCC scores (mean/std) on the MSP-Podcast dataset

Methods	CCC _{avg}	CCC _{act}	CCC _{val}	CCC _{dom}
Sup	.458 ± .005	.596 ± .007	.266 ± .004	.501 ± .013
jointCPC	.491 ± .008	.628 ± .006	.280 ± .006	.568 ± .007
miniCPC	.549 ± .006	.688 ± .009	.345 ± .005	.615 ± .011
preCPC	.571 ± .004	.706 ± .006	.377 ± .008	.639 ± .012

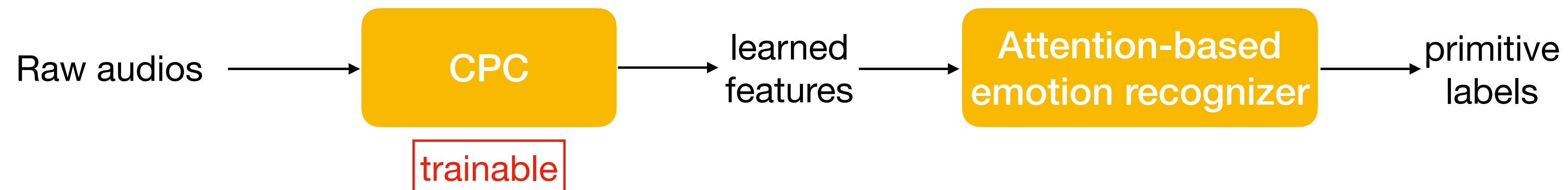
Experiments

Experimental Setups

- preCPC: pre-trained CPC (LibriSpeech) + supervised (IEMOCAP/MSP-Podcast)
- Sup: supervised only (IEMOCAP/MSP-Podcast)
- miniCPC: pre-trained CPC (IEMOCAP/MSP-Podcast)+ supervised (IEMOCAP/MSP-Podcast)



- jointCPC: pre-trained CPC (IEMOCAP/MSP-Podcast)+ supervised (IEMOCAP/MSP-Podcast)



Experiments

Experimental Setups

- preCPC: pre-trained CPC + supervised
- Sup: supervised only
- miniCPC: pre-trained CPC + supervised
- jointCPC: pre-trained CPC + supervised

Hypothesis:

- Unsupervised pre-training produces representations with better generalization which facilitate various downstream tasks

Table 1: CCC scores (mean/std) on the IEMOCAP dataset

Methods	CCC _{avg}	CCC _{act}	CCC _{val}	CCC _{dom}
Sup	.664 ± .007	.638 ± .017	.718 ± .004	.635 ± .009
jointCPC	.562 ± .012	.549 ± .032	.642 ± .013	.491 ± .016
miniCPC	.660 ± .005	.673 ± .028	.702 ± .009	.606 ± .019
preCPC	.731 ± .003	.752 ± .014	.752 ± .009	.691 ± .009

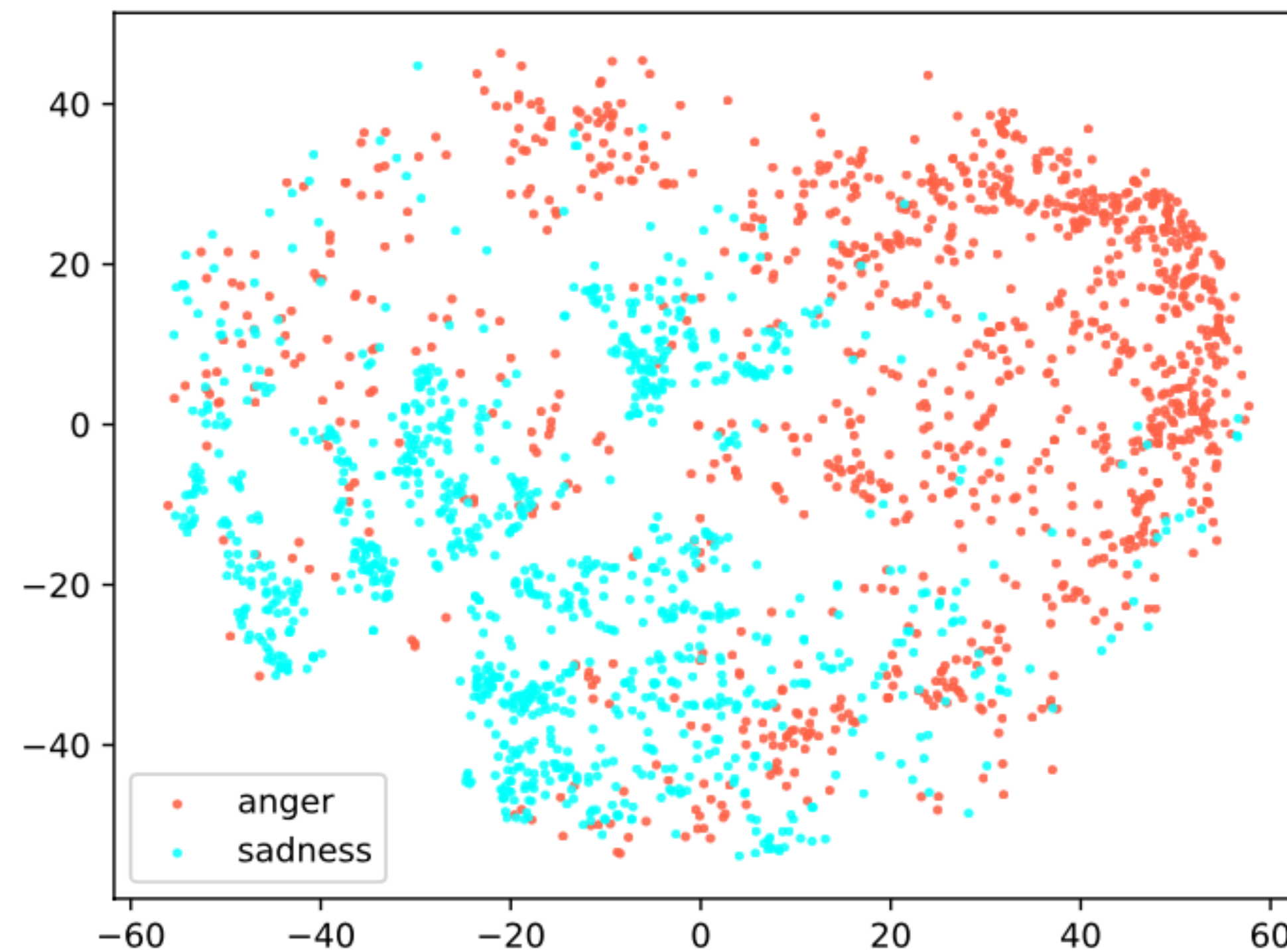
Table 2: CCC scores (mean/std) on the MSP-Podcast dataset

Methods	CCC _{avg}	CCC _{act}	CCC _{val}	CCC _{dom}
Sup	.458 ± .005	.596 ± .007	.266 ± .004	.501 ± .013
jointCPC	.491 ± .008	.628 ± .006	.280 ± .006	.568 ± .007
miniCPC	.549 ± .006	.688 ± .009	.345 ± .005	.615 ± .011
preCPC	.571 ± .004	.706 ± .006	.377 ± .008	.639 ± .012

Experiments

Representation Visualization:

- t-SNE plot of representations that learned from CPC



Data points are well separated, even though trained without emotion labels

Conclusion & Future Work

Conclusion

- CPC can learn salient features from unlabeled speech corpora that benefits emotion recognition task
- Obtained competitive performance on public benchmarks

Future work

- Investigate the impact unsupervised representation learning data on emotion recognition performance (e.g., replace Libri speech with TED data)
- Still to do: end-to-end modeling

You are very welcome to our poster session!

Speech Emotion 3: Emotion Recognition-Representations, Data Augmentation

Wednesday, 9 June, 15:30 - 16:15

Thank you!

