

## Introduction

### Motivation:

- ✗ Implicitly learning prosodic information from audio is often less than optimal due to the discretization of audio signals during training of leading speech models (e.g., HuBERT)
- ✗ Direct fine-tuning of existing speech models originally trained for ASR doesn't perform well on SER.
- ✗ Direct use of transcripts at *run-time* can lead to low performances due to transcription errors.
- ✗ Using both audio and linguistic information at *run-time* requires a multimodal system which can increase computational overhead.

### Our contributions:

- ✓ We introduce EmoDistill, a novel cross-modal Knowledge Distillation (KD) framework for learning unimodal representations from speech that explicitly capture both the linguistic and prosodic aspects of emotions.
- ✓ EmoDistill outperforms previous state-of-the-art methods on IEMOCAP and achieves 77.49% UA and 78.91% WA.

## Experiment Details

### Dataset:

- IEMOCAP benchmark
- 4 emotions (neutral, angry, sad, happy)
- 10-Fold cross-validation
- Subject-independent

### Implementation:

- Prosodic Teacher: 4-layer ResNet2D trained on eGeMAPs LLDs.
- Linguistic Teacher: BERT-base (pre-trained)
- AdamW, base LR of 1e-4
- 4 × NVIDIA A100 GPUs, Batch size = 128

## Method

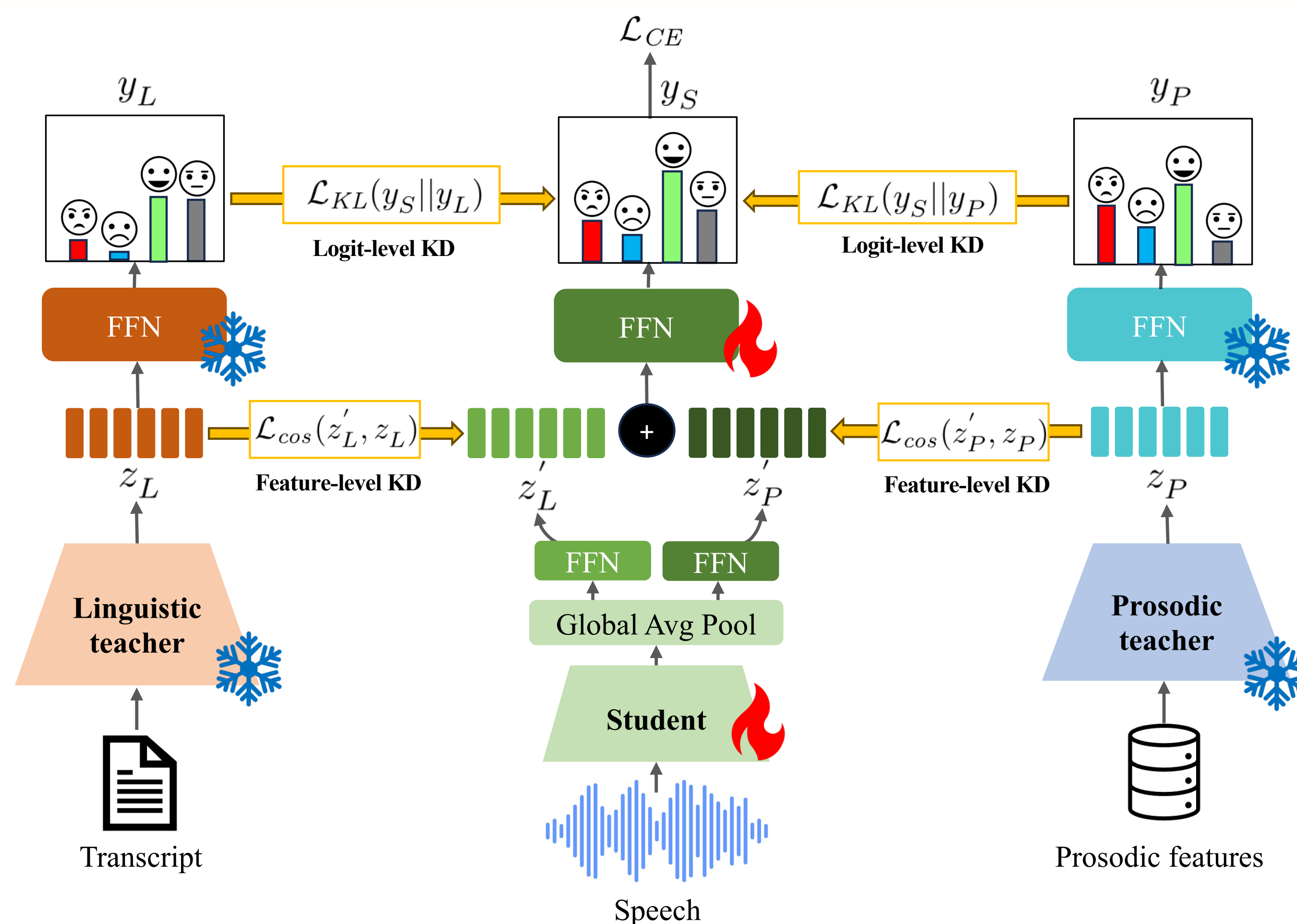


Fig. 1: EmoDistill Framework. Our student network is trained using a distillation of logit-level and embedding-level knowledge from frozen linguistic and prosodic teacher networks, along with standard cross-entropy loss. **During inference, we only use the student network in an unimodal setup, avoiding computational overhead as well as transcription and prosodic feature extraction errors.**

**Logit-level KD.** First, we transfer the logit-level knowledge using traditional KD with temperature-scaled labels [1]. Specifically, we minimize the KL-Divergence  $L_{KL}$  between the predicted logit distributions of teacher and student models, where the objective becomes:

$$L_{logits} = L_{KL}(y_S || y_L) + L_{KL}(y_S || y_P). \quad (1)$$

Here,  $y_S$  refers to the predictions of the student, while  $y_L$  and  $y_P$  represent the predictions of Linguistic and Prosodic teacher models, respectively. In all cases, the predicted logits  $y$  are obtained using temperature parameter  $\tau$  in the output softmax activation function. In practice, we use different values of  $\tau$  for KD from  $f_T^L$  and  $f_T^P$ . Let  $z_c$  be the output logits for class  $c$ , among a total of  $N$  classes. The temperature-scaled logits  $y_c$  are obtained as:

$$y_c = \frac{e^{z_c/\tau}}{\sum_{k=1}^N e^{z_k/\tau}}. \quad (2)$$

**Feature-level KD.** Next, we use embedding-level KD to transfer knowledge to the student model from the latent space of Linguistic and Prosodic teacher models. Let  $z_L$  and  $z_P$  denote the embeddings of Linguistic and Prosodic teachers, while  $z'_L$  and  $z'_P$  denote the embeddings of the student model from linguistic and prosodic projection layers respectively. We minimize the negative cosine similarity  $L_{cos}$  among the teacher and student embeddings as follows:

$$L_{embeddings} = L_{cos}(z'_L, z_L) + L_{cos}(z'_P, z_P). \quad (3)$$

Given two embeddings  $a$  and  $b$ ,  $L_{cos}$  can be defined as:

$$L_{cos}(a, b) = \frac{a \cdot b}{\|a\|_2 \cdot \|b\|_2}, \quad (4)$$

where  $\|\cdot\|_2$  represents  $\ell_2$ -norm.

**Loss objective.**  $L_{EmoDistill} = \alpha L_{logits} + \beta L_{embeddings} + \gamma L_{CE}$  (5)

## Performance Evaluation

Tab. 1: SER results on IEMOCAP. **Bold** denotes the best results while underline denotes the second-best.

Method	Inf. Backbone	Modality	WA	UA
Sun <i>et al.</i> (2021)	CNN+LSTM	Multimodal	61.2	56.01
Heusser <i>et al.</i> (2019)	BiLSTM+XLNet	Multimodal	71.40	68.60
Triantafyllopoulos <i>et al.</i> (2023)	MFCNN+BERT	Multimodal	-	72.60
Ho <i>et al.</i> (2020)	RNN+BERT	Multimodal	73.23	74.33
Aftab <i>et al.</i> (2022)	FCNN	Unimodal	70.23	70.76
Liu <i>et al.</i> (2020)	TFCNN+DenseCap+ELM	Unimodal	70.34	70.78
Cao <i>et al.</i> (2021)	LSTM+Attention	Unimodal	70.50	72.50
Lu <i>et al.</i> (2020)	RNN-T	Unimodal	71.72	72.56
Wu <i>et al.</i> (2021)	CNN-GRU+SeqCap	Unimodal	72.73	59.71
Zou <i>et al.</i> (2022)	Wav2Vec2+CNN+LSTM	Unimodal	71.64	72.70
Ye <i>et al.</i> (2023)	TIM-Net	Unimodal	72.50	71.65
Ours	HuBERT-base	Unimodal	<b>75.16</b>	<b>76.12</b>
Ours	HuBERT-large	Unimodal	<b>77.49</b>	<b>78.91</b>

## Acknowledgement

This work was supported by Mitacs, Vector Institute, and Ingenuity Labs Research Institute.

## Ablation study

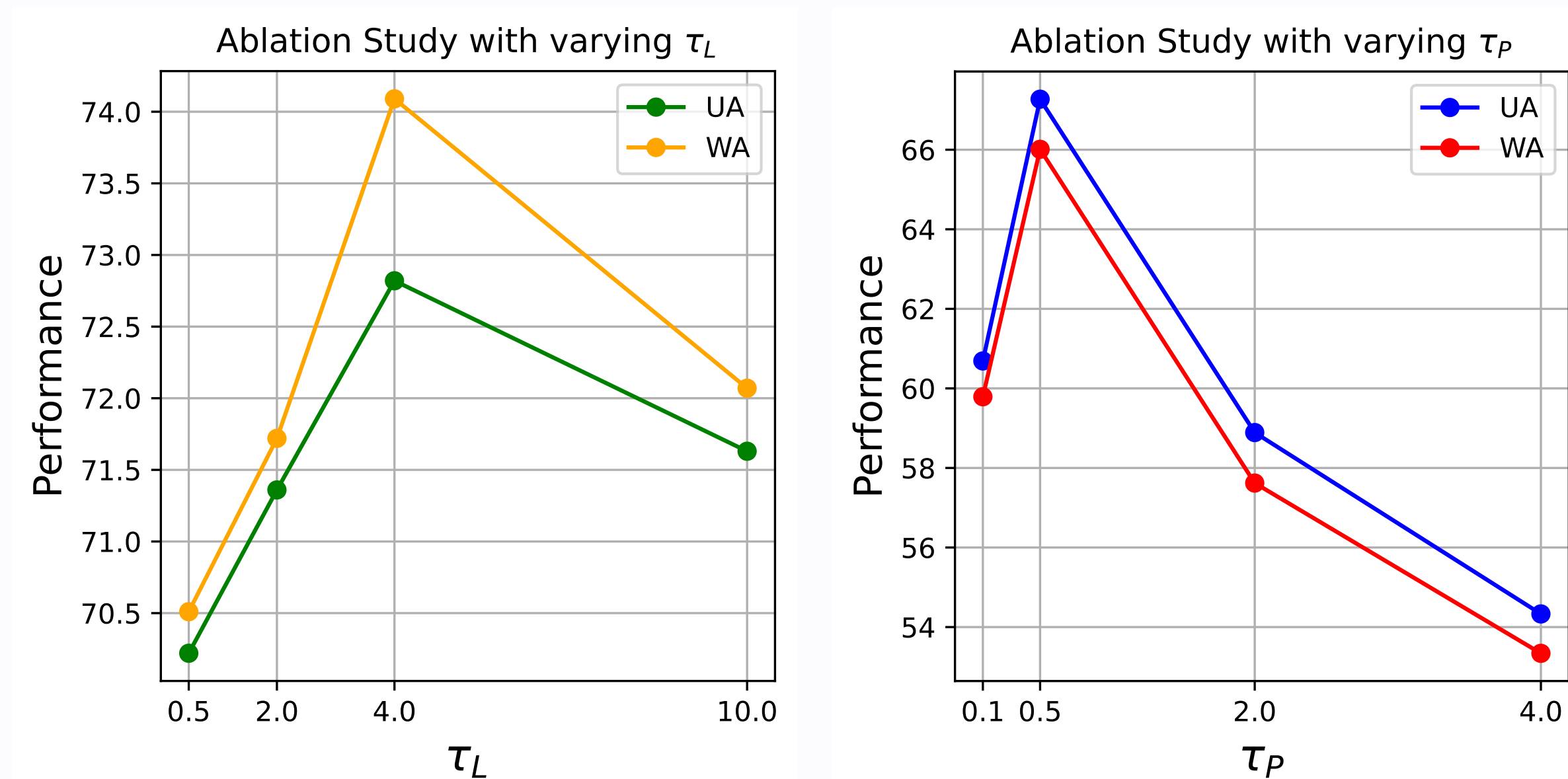


Fig. 2: **Left:** We remove  $f_T^L$  and vary  $\tau_P$ . **Right:** We remove  $f_T^P$  and vary  $\tau_L$ .

Variants	WA	UA
Ours	<b>75.16</b>	<b>76.12</b>
w/o $L_{logits}$	73.94 (↓ 1.22)	74.02 (↓ 2.10)
w/o $L_{embedding}$	73.88 (↓ 1.28)	74.01 (↓ 2.11)
w/o $f_T^P$	74.09 (↓ 1.07)	72.82 (↓ 3.30)
w/o $f_T^L$	66.01 (↓ 9.15)	67.27 (↓ 8.85)
w/o $f_T^P$ and $f_T^L$	69.92 (↓ 5.24)	70.17 (↓ 5.95)
w/o $f_S$ and $f_T^P$	49.42 (↓ 25.74)	50.08 (↓ 26.04)
w/o $f_S$ and $f_T^L$	71.09 (↓ 4.07)	71.83 (↓ 4.29)

- ✓ Linguistic understanding is crucial for SER.
- ✓ Prosodic understanding is complimentary, but leads to a boost in SER performance.
- ✓ Hard logits are better for KD from the prosodic teacher, as it is a weak teacher.
- ✓ Soft logits are better for KD from the linguistic teacher, as it is a strong teacher.
- ✓ Using Logit and Embedding-level KD together improves performance.