# Enhancing Resilience to Missing Data in Audio-Text Emotion Recognition with Multi-Scale Chunk Regularization

### UT Dallas MSP — Multimodal Signal Processing Laboratory
### THE UNIVERSITY OF TEXAS AT DALLAS

Wei-Cheng Lin, Lucas Goncalves, Carlos Busso

Erik Jonsson School of Engineering & Computer Science at the University of Texas at Dallas, Richardson, Texas 75080, USA

## Motivation

**Background:**

- Audio-Text Multimodal Emotion Recognition
  - Unclear role of temporal synchronization (i.e., alignment) between the input audio and text sequences
  - Current model-level and feature-level fusion techniques cannot investigate this research question
  - Multimodal modeling can effectively improve recognition performance but also reduces the model robustness against missing data scenario
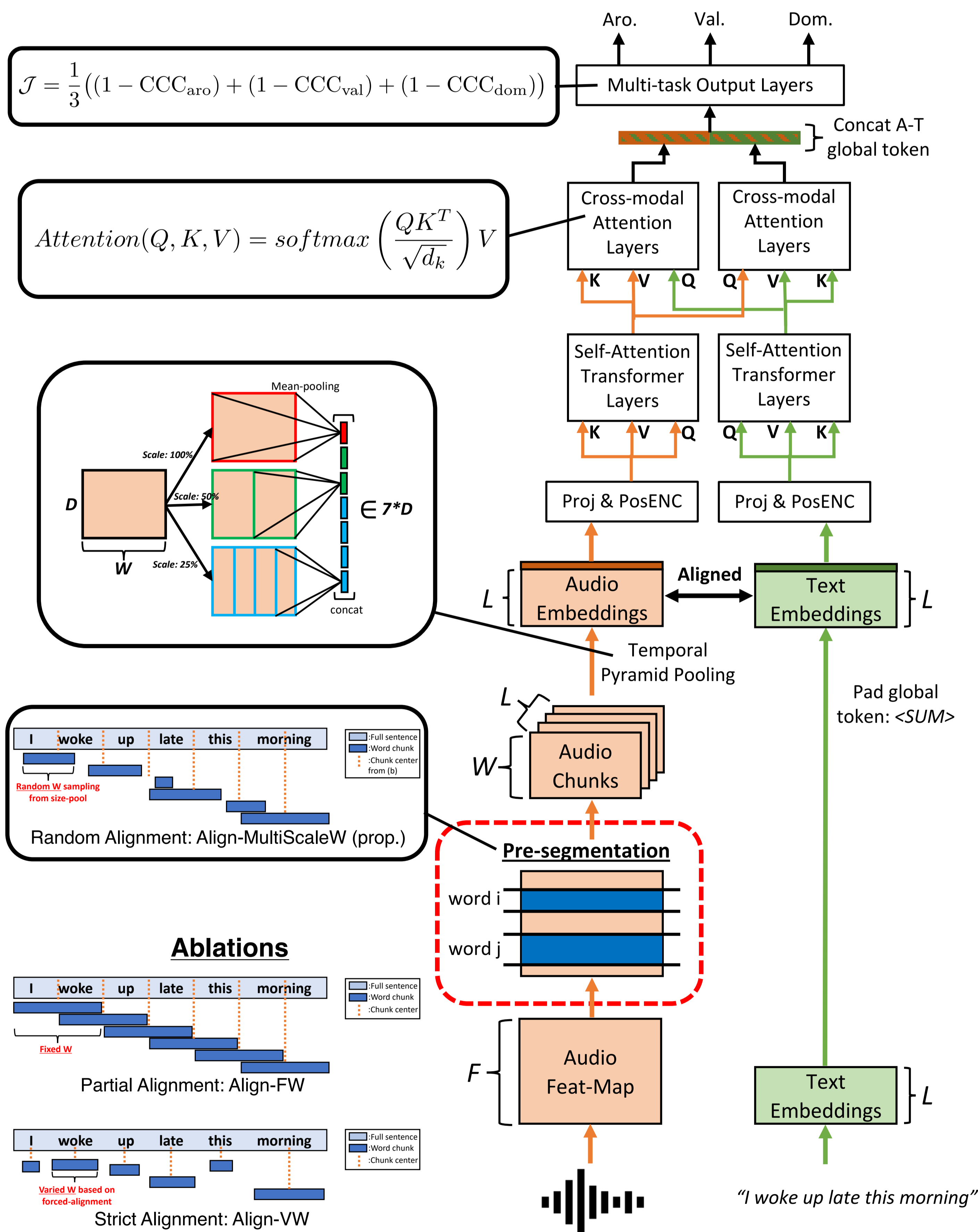
**Our Work:**

- Proposes a novel word-chunk modeling concept for audio-text emotion recognition
- Confirms that the self-attention mechanism is powerful enough to capture temporal alignment across audio-text
- Propose to leverage multi-scale chunk regularization to improve model's robustness against missing data
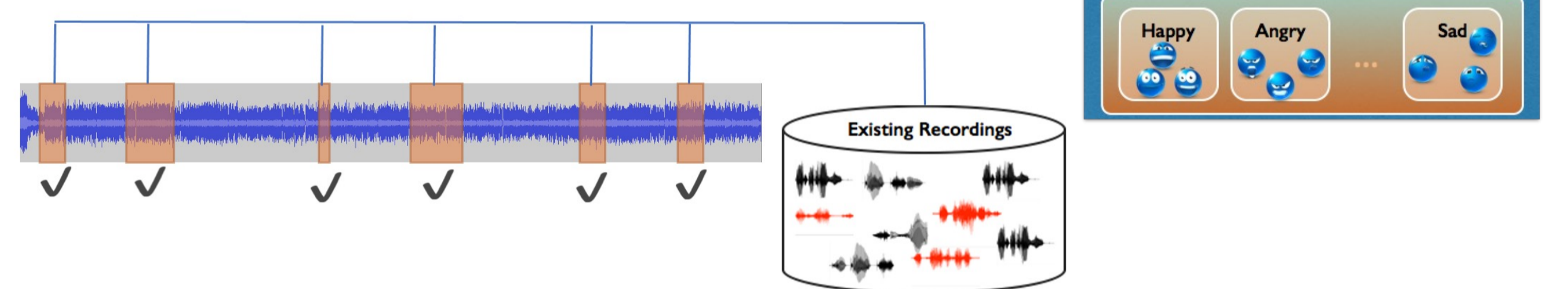
## Proposed Framework

**Setups:**

- Input Features
  - Wav2vec2-large-robust (Audio: 1,024D)
  - RoBERTa-base (Text: 768D)
- Model Architecture
  - [3 X Self-Attention Blocks] + [3 X Cross-Attention Blocks]
  - Hidden dim: 256D

$$\mathcal{J} = \frac{1}{3}\left((1-\text{CCC}_{\text{aro}}) + (1-\text{CCC}_{\text{val}}) + (1-\text{CCC}_{\text{dom}})\right)$$

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Aro.   Val.   Dom.

Multi-task Output Layers

Concat A-T global token

Cross-modal Attention Layers — Cross-modal Attention Layers
K V Q — Q V K

Self-Attention Transformer Layers — Self-Attention Transformer Layers
K V Q — Q V K

Proj & PosENC — Proj & PosENC

Audio Embeddings — **Aligned** — Text Embeddings

$L$ — $L$

Mean-pooling
Scale: 100%
Scale: 50%
Scale: 25%
$\in 7*D$
concat
$D$ / $W$

Temporal Pyramid Pooling

Pad global token: *<SUM>*

Audio Chunks
$W$ / $L$

Audio Feat-Map
$F$

Text Embeddings $L$

Pre-segmentation
word i
word j

*"I woke up late this morning"*

**Random Alignment: Align-MultiScaleW (prop.)**

I woke up late this morning
Full sentence / Word chunk / Chunk center from (b)
Random W sampling from size-pool

**Ablations**

I woke up late this morning
Full sentence / Word chunk / Chunk center
Fixed W

**Partial Alignment: Align-FW**

I woke up late this morning
Full sentence / Word chunk / Chunk center
Varied W based on forced-alignment

**Strict Alignment: Align-VW**

## Resources

**Dataset:**

- The MSP-PODCAST v1.10 corpus
  - Largest spontaneous speech emotion corpus collected from existing podcast recordings
  - Annotated using crowdsourcing framework (AMT)
  - Includes train/dev/test: 63,076/10,999/16,903 clips (~166hrs)
  - Regression problem: arousal, dominance, and valence
  - Montreal-forced-aligner (MFA)
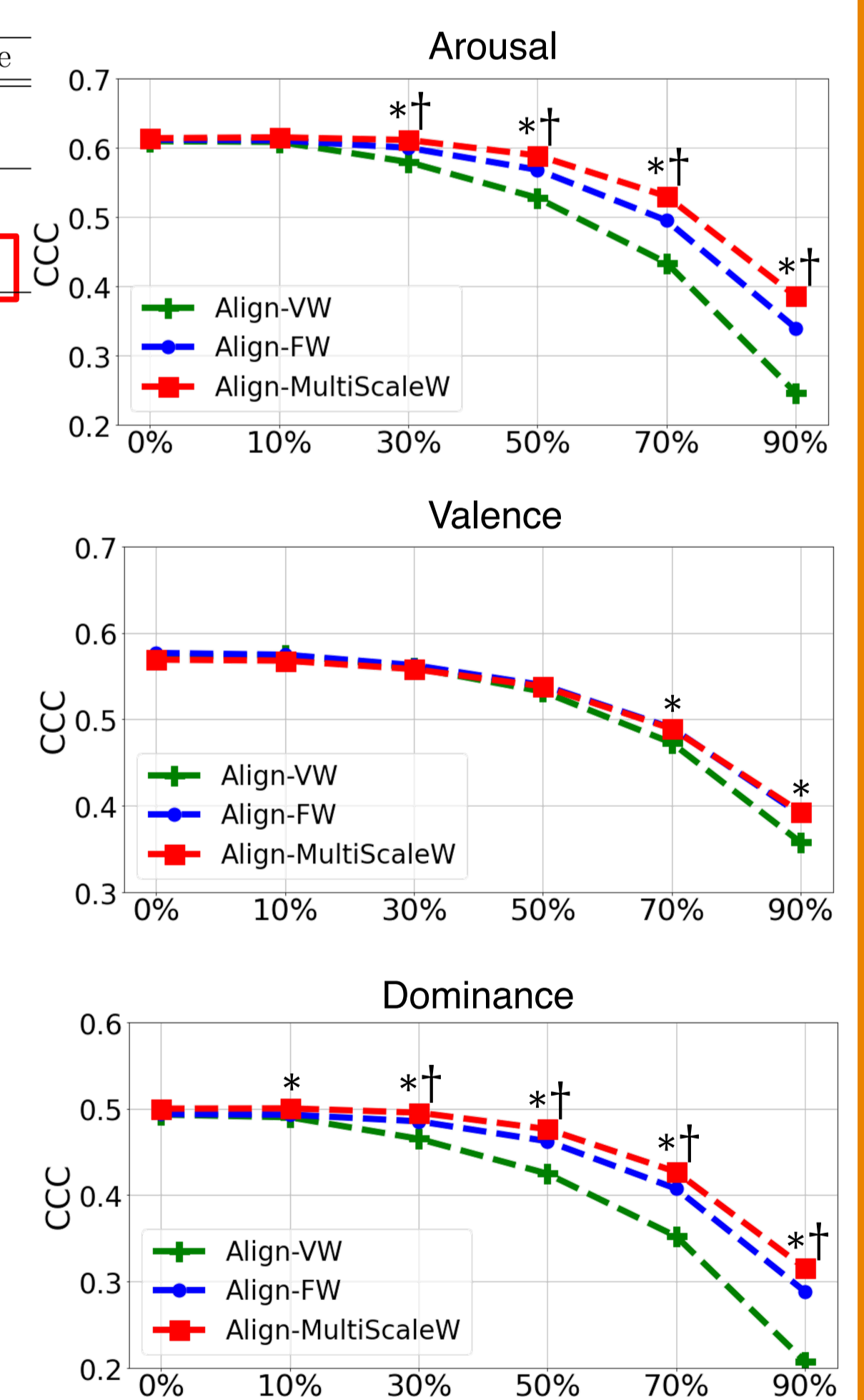  - Divide into speaker turns (2.75 - 11 secs)

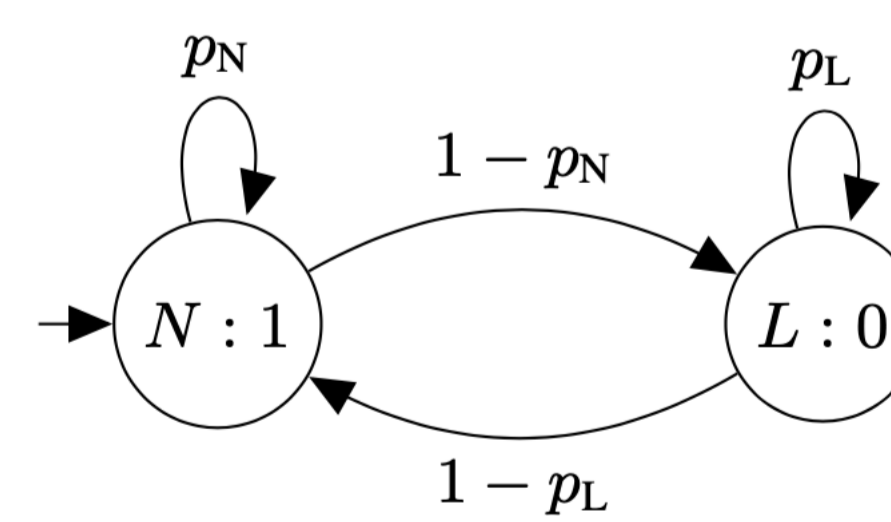## Experimental Results

### SER Performance Comparison

| Approach | Arousal | Valence | Dominance |
|---|---|---|---|
| *MDRE-GRU* [52] | 0.6090 | 0.5592 | 0.4842 |
| *MDREA-GRU* [52] | 0.6029 | 0.5603 | 0.4792 |
| *Align-VW* | 0.6094 | 0.5723*† | 0.4932*† |
| *Align-FW* | 0.6117† | **0.5767**\*† | 0.4930*† |
| *Align-MultiScaleW* | **0.6137**† | 0.5691*† | **0.5000**\*† |

**Improved and robust performance for all three emotions under different missing data testing conditions**

### Missing Data Robustness for Audio-Text: Random Drop Words

Arousal — CCC vs 0%–90%
Align-VW / Align-FW / Align-MultiScaleW

Valence — CCC vs 0%–90%

Dominance — CCC vs 0%–90%

### Missing Data Robustness for Audio-Only: Packet-Loss

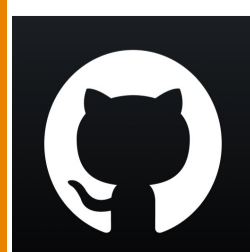$p_N$ — $N:1$ — $1 - p_N$ — $L:0$ — $p_L$
$1 - p_L$

**Obtain the same robust performance trend as well as for Audio-Only SER**

## Conclusions

- Word-chunk concept can explicitly model and control the alignment level between audio-text sequences
- Multi-scale chunk regularization can effectively improve model robustness against missing data conditions, which is valid for both the audio-text and audio-only scenarios

**Future Work**

- Extending word-chunk concept to more modalities (e.g., audio-text-video) for temporal synchronization
- Apply multi-scale chunk regularization based on different modalities to improve model performance

**Github link:** https://github.com/winston-lin-wei-cheng/MultiScale-Chunk-Regularization