



Role of Lexical Boundary Information in Chunk-Level Segmentation for Speech Emotion Recognition

Wei-Cheng Lin and Carlos Busso

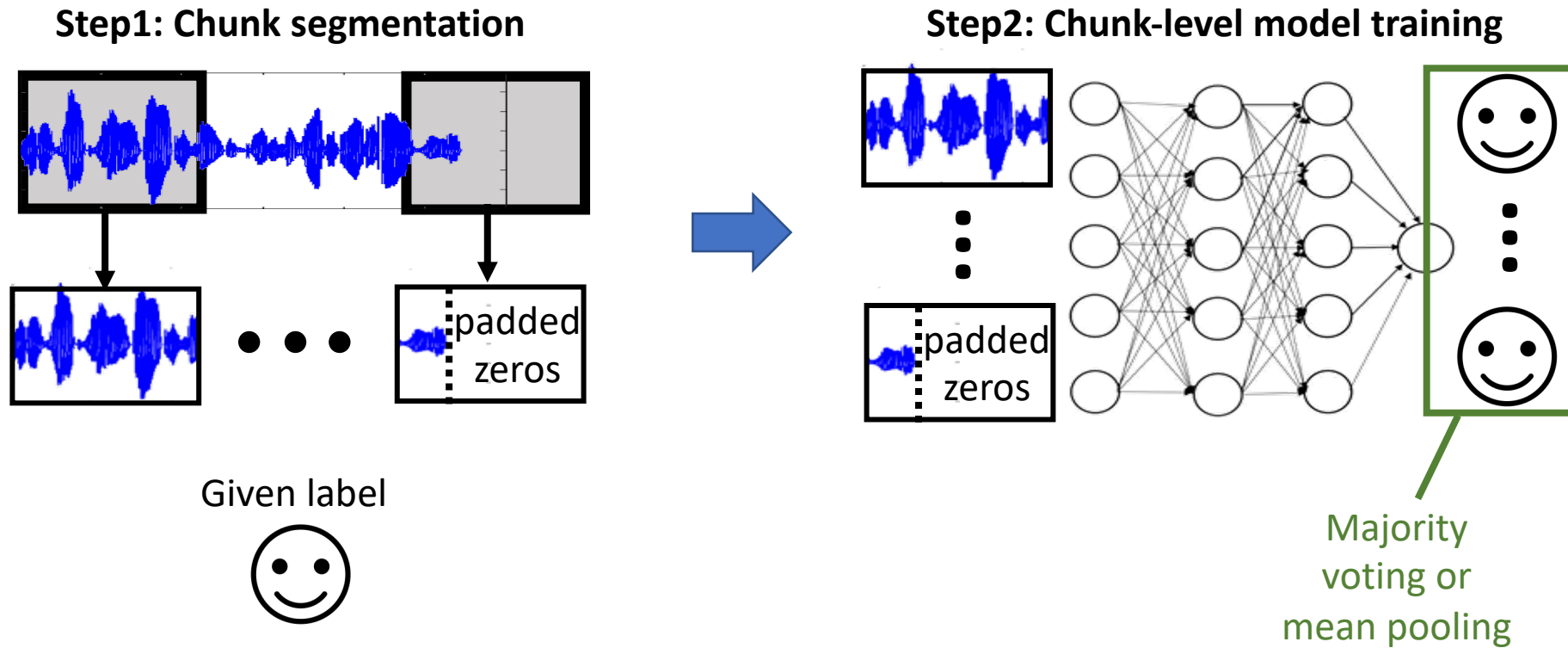


Outline

- 1. Motivation**
- 2. Lexical-based Segmentation**
- 3. Experimental Results and Analysis**
- 4. Conclusions & Future Works**

Motivation

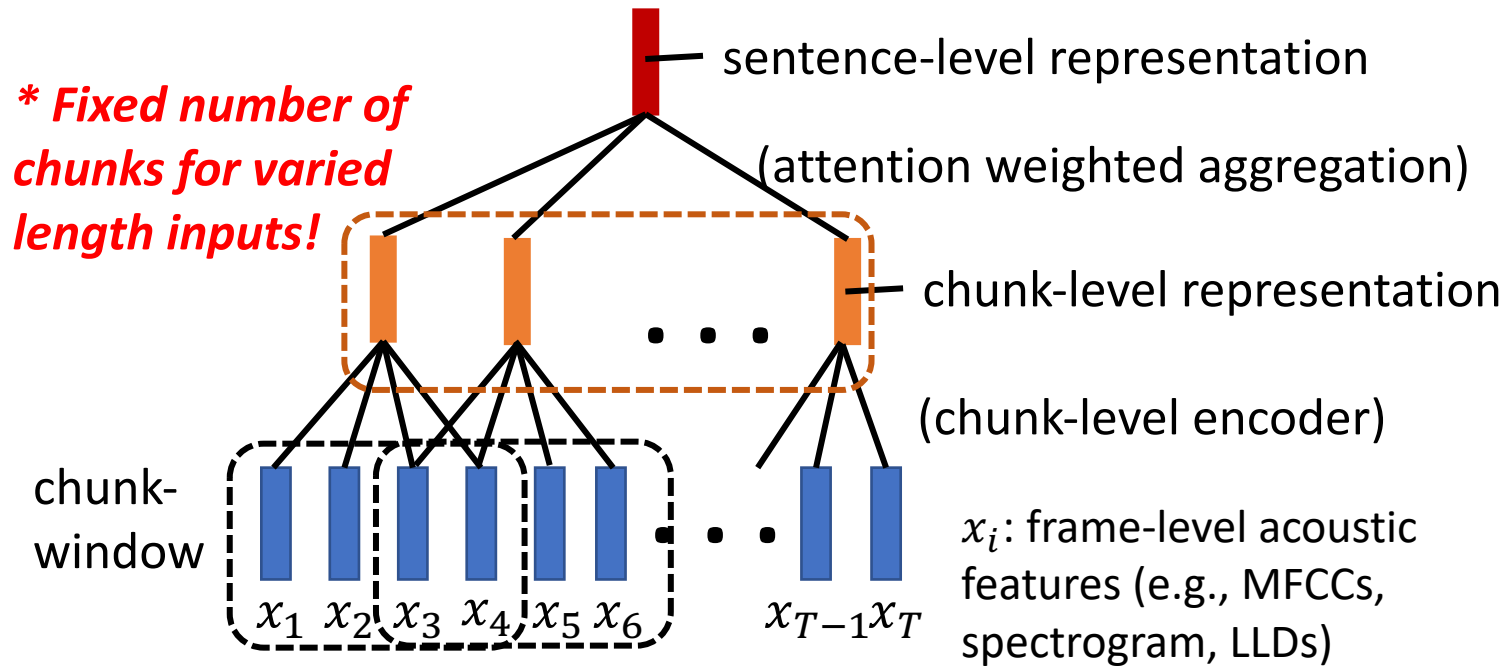
- **Chunk-level modeling for *speech emotion recognition (SER)***



Motivation

- Dynamic chunk segmentation [1]

Hierarchical temporal-info summarization

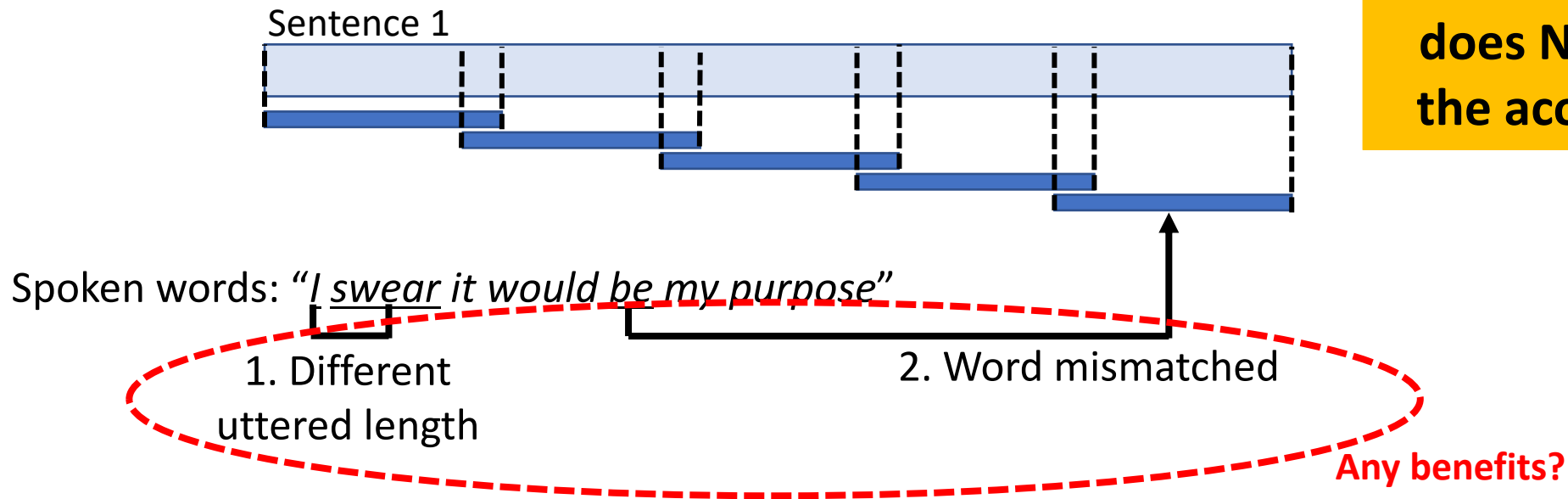


Motivation

- **Time-based chunk segmentation**

- Split data chunks without considering their actual lexical boundaries
- E.g., dynamic chunk segmentation

Lexical information does NOT align with the acoustic frames!



- **Lexical chunk segmentation**

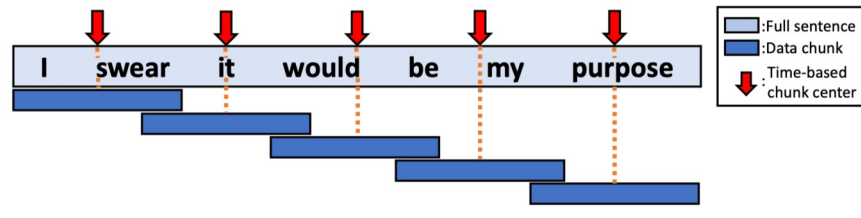
- Factor C : number of chunks to split
=> fixed or number of spoken words
- Factor W : chunk window size
=> fixed or word-level alignment boundary
- Result in different combinations, e.g., Fixed W -Varied C (FW-VC)

The model does NOT consider the semantic meaning of words, still an SER task!

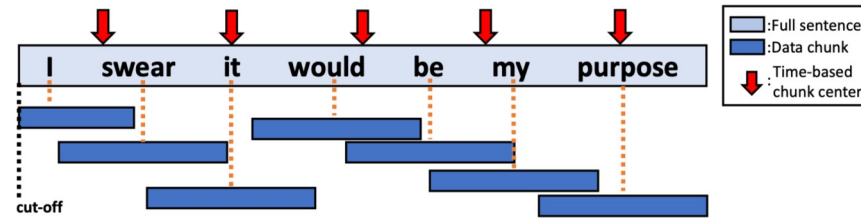
Lexical-based Segmentation

Visualization of different options

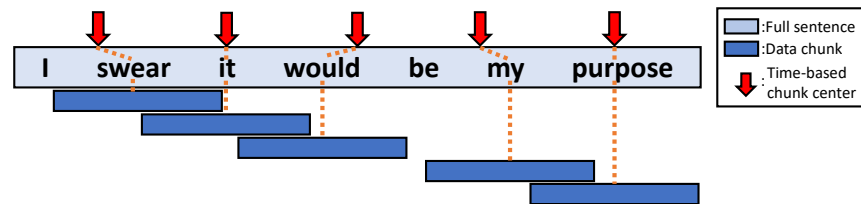
Time-based FixedW-FixedC (tFW-FC)



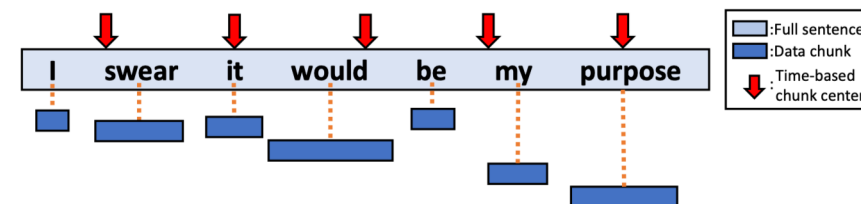
Lexical FixedW-VariedC (FW-VC)



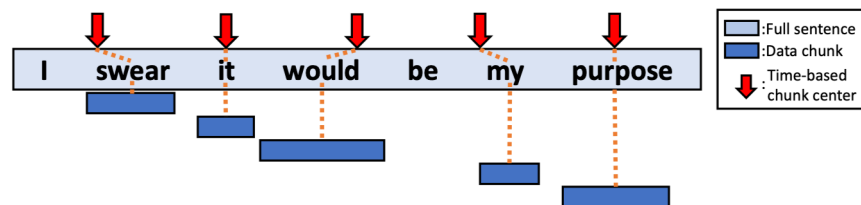
Lexical FixedW-FixedC (FW-FC)



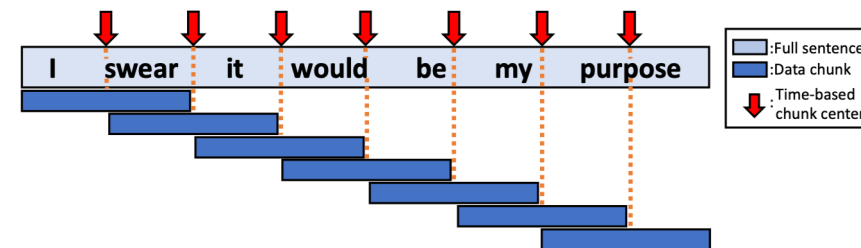
Lexical VariedW-VariedC (VW-VC)



Lexical VariedW-FixedC (VW-FC)



Combine FixedW-VariedC (cFW-VC)



■ Experimental Setting

- Datasets: IEMOCAP and MSP-Podcast v1.10 [2]
- Acoustic features: LLDs and wav2vec2 [3]
- Model arch: LSTM chunk-level encoder + Multi-Heads Self-Attention aggregation
- Emotion: arousal, dominance, and valence
- Evaluation metric: concordance correlation coefficient (CCC)

The only difference is how we segment the data chunks!

[2] R. Lotfian and C. Busso, "Building naturalistic emotionally balanced speech corpus by retrieving emotional speech from existing podcast recordings," IEEE Transactions on Affective Computing, vol. 10, no. 4, pp. 471–483, October-December 2019.

[3] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," in Advances in Neural Information Processing Systems (NeurIPS 2020), vol. 33, Virtual, December 2020, pp. 12 449–12 460.

Experimental Results and Analysis

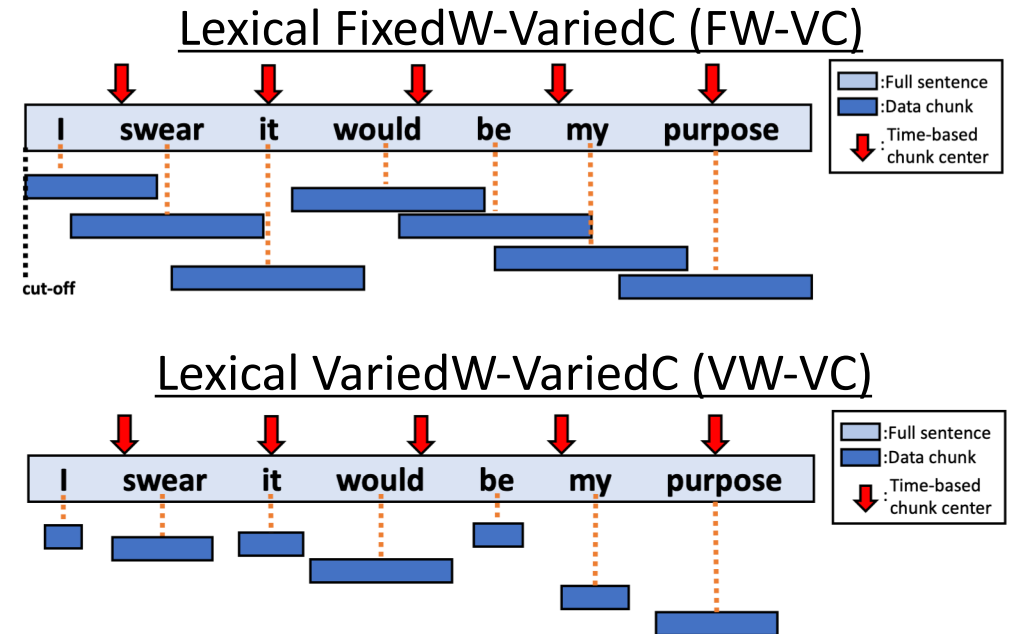
Performance comparison summary

MSP-Podcast v1.10								
Method	CovrR/OvrR [%]	LLDs (CCC)			Wav2Vec2 (CCC)			
		Aro.	Val.	Dom.	Aro.	Val.	Dom.	
<i>tFW-FC</i>	99 / 50	0.528	0.216*	0.430	0.604	0.352	0.478	
<i>FW-FC</i>	90 / 56	0.529	0.191	0.423	0.598	0.344	0.475	
<i>VW-FC</i>	57 / 35	0.534	0.170	0.427	0.595	0.352	0.460	
<i>FW-VC</i>	94 / 68	0.544*	0.141	0.455*	0.620*	0.349	0.497*	
<i>VW-VC</i>	82 / 0	0.546*	0.118	0.459*	0.613*	0.336	0.492*	
<i>cFW-VC</i>	99 / 65	0.562†	0.207*	0.468†	0.616*	0.343	0.499*	

IEMOCAP								
Method	CovrR/OvrR [%]	LLDs (CCC)			Wav2Vec2 (CCC)			
		Aro.	Val.	Dom.	Aro.	Val.	Dom.	
<i>tFW-FC</i>	99 / 79	0.614	0.353	0.406	0.709*	0.554	0.531	
<i>FW-FC</i>	82 / 83	0.593	0.257	0.411	0.700	0.537	0.539	
<i>VW-FC</i>	47 / 71	0.595	0.279	0.409	0.688	0.532	0.526	
<i>FW-VC</i>	81 / 67	0.633*	0.378	0.451*	0.713*	0.582*	0.538	
<i>VW-VC</i>	61 / 0	0.626*	0.395*	0.433	0.719*	0.577*	0.558*	
<i>cFW-VC</i>	99 / 60	0.636*	0.404†	0.463†	0.718*	0.584*	0.549*	

* means statistically significant better performance over other approaches without a marker

† means the results are statistically significant better than all other approaches



Knowing the precise word boundary (i.e., VariedW) does NOT bring significant performance benefits!

Experimental Results and Analysis

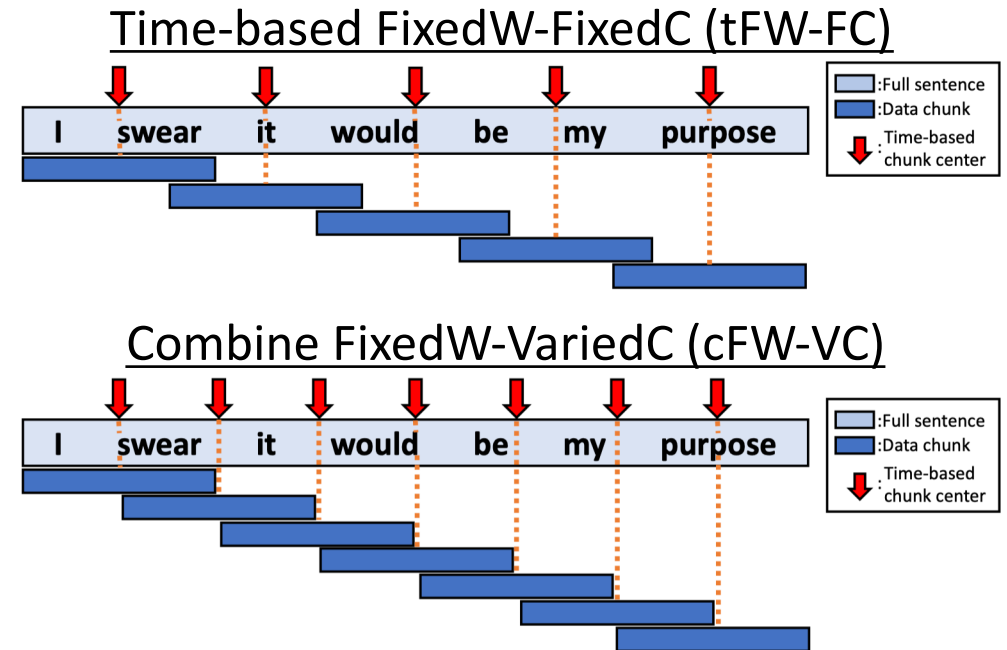
Performance comparison summary

MSP-Podcast v1.10							
Method	CovrR/OvrR [%]	LLDs (CCC)			Wav2Vec2 (CCC)		
		Aro.	Val.	Dom.	Aro.	Val.	Dom.
<i>tFW-FC</i>	99 / 50	0.528	0.216*	0.430	0.604	0.352	0.478
<i>FW-FC</i>	90 / 56	0.529	0.191	0.423	0.598	0.344	0.475
<i>VW-FC</i>	57 / 35	0.534	0.170	0.427	0.595	0.352	0.460
<i>FW-VC</i>	94 / 68	0.544*	0.141	0.455*	0.620*	0.349	0.497*
<i>VW-VC</i>	82 / 0	0.546*	0.118	0.459*	0.613*	0.336	0.492*
<i>cFW-VC</i>	99 / 65	0.562†	0.207*	0.468†	0.616*	0.343	0.499*

IEMOCAP							
Method	CovrR/OvrR [%]	LLDs (CCC)			Wav2Vec2 (CCC)		
		Aro.	Val.	Dom.	Aro.	Val.	Dom.
<i>tFW-FC</i>	99 / 79	0.614	0.353	0.406	0.709*	0.554	0.531
<i>FW-FC</i>	82 / 83	0.593	0.257	0.411	0.700	0.537	0.539
<i>VW-FC</i>	47 / 71	0.595	0.279	0.409	0.688	0.532	0.526
<i>FW-VC</i>	81 / 67	0.633*	0.378	0.451*	0.713*	0.582*	0.538
<i>VW-VC</i>	61 / 0	0.626*	0.395*	0.433	0.719*	0.577*	0.558*
<i>cFW-VC</i>	99 / 60	0.636*	0.404†	0.463†	0.718*	0.584*	0.549*

* means statistically significant better performance over other approaches without a marker

† means the results are statistically significant better than all other approaches



Knowing how many chunks use to split (i.e., depending on the number of words) is crucial!

■ Conclusion

- We found a minor role of word-level timing boundaries for chunk-level SER
- We found that splitting data chunks according to number of words is the key leads to better chunk-level SER

■ Future Work

- Multimodal segmentation and benefits for modeling (e.g., video-speech-text)

Welcome to join our poster session in ICASSP 2023

Paper ID: 4354

Session Date/Time: 6/8/2023 14:00:00 (EEST)

Session Name: Speech Emotion Recognition: Multimodality



**Thank you for your
attention !**

This work was funded by NSF
(CNS-2016719; IIS-1453781)



Questions or Contact:

wei-cheng.lin@utdallas.edu

busso@utdallas.edu