



INTERSPEECH 2020

OCTOBER 25-29/ SHANGHAI, CHINA
SHANGHAI INTERNATIONAL CONVENTION CENTER

An Efficient Temporal Modeling Approach for Speech Emotion Recognition by Mapping Varied Duration Sentences into Fixed Number of Chunks

Wei-Cheng Lin and Carlos Busso

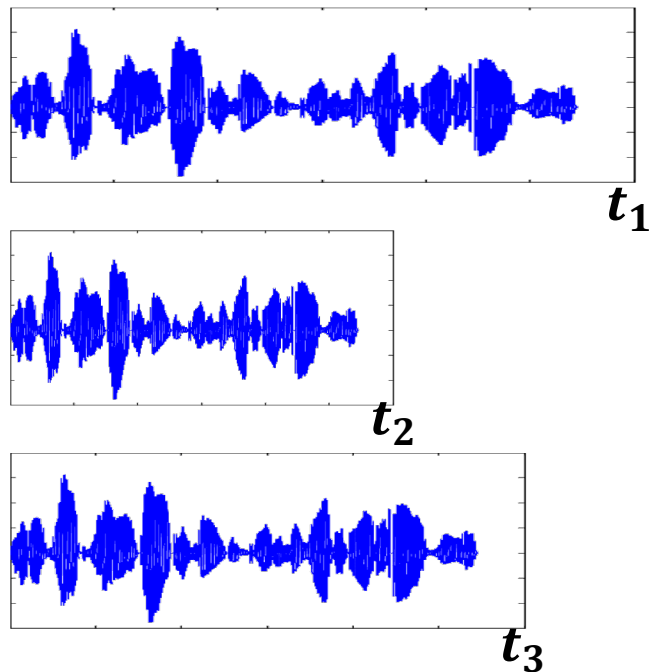


Outline

- 1. Background**
- 2. Proposed Methodology**
- 3. Experimental Results and Analysis**
- 4. Conclusions & Future Works**

■ Sequence-to-One Problem

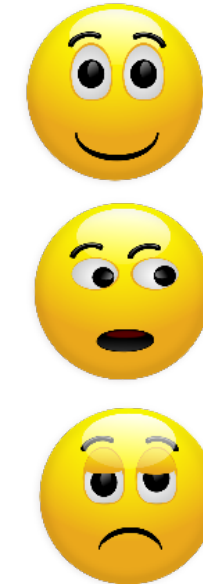
Varied Duration Inputs



Temporal Modeling



Sentence-level Emotions



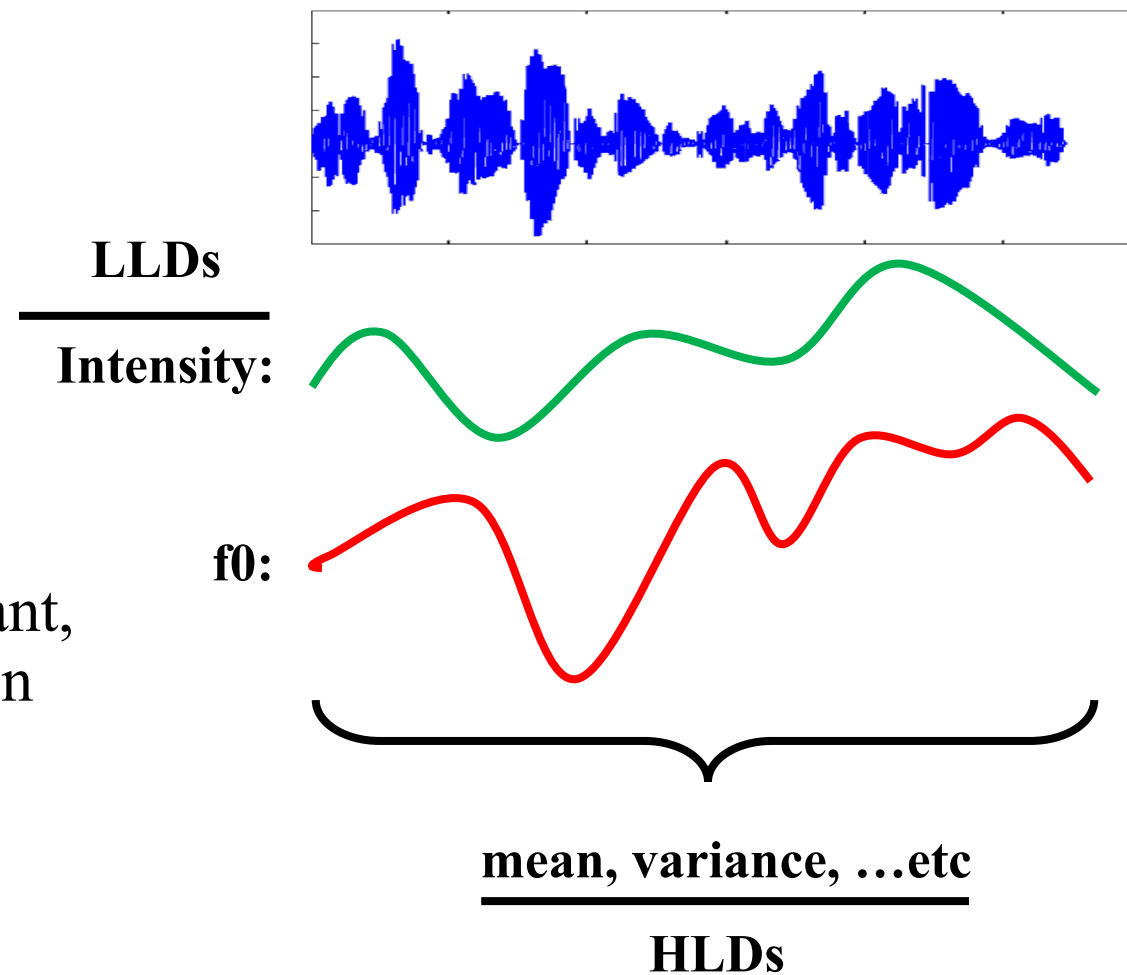
**Sentence-level
Representation**
But How?

Traditional Approach

- Frame-level LLDs (e.g., f_0 , MFCCs, energy)
- Sentence-level HLDs (e.g., mean, variance)
- Learning models (e.g., SVM, FCNN)

Issues

- HLDs assume all the frames are equally important, ignoring non-uniform externalization of emotion
- Static-encoding vector
- It is not able to dynamically reflect emotional changes over time

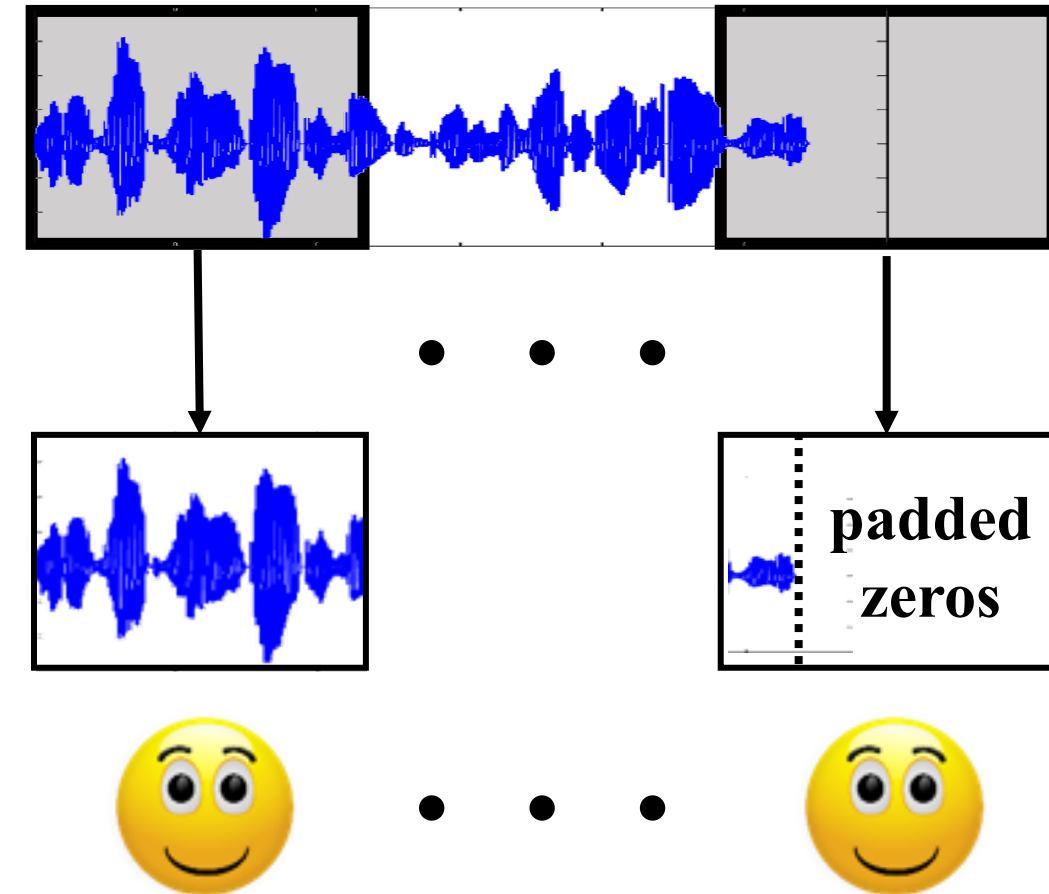


Deep Learning Approach

- Jointly trained feature extractor with discriminator
 - Powerful feature representation
- Cropping and Zero-Padding to deal with sentences with different length
- CNN, LSTM or hybrid CNN-LSTM
 - Temporal mean pooling
 - Majority voting or averaging outputs

Issues

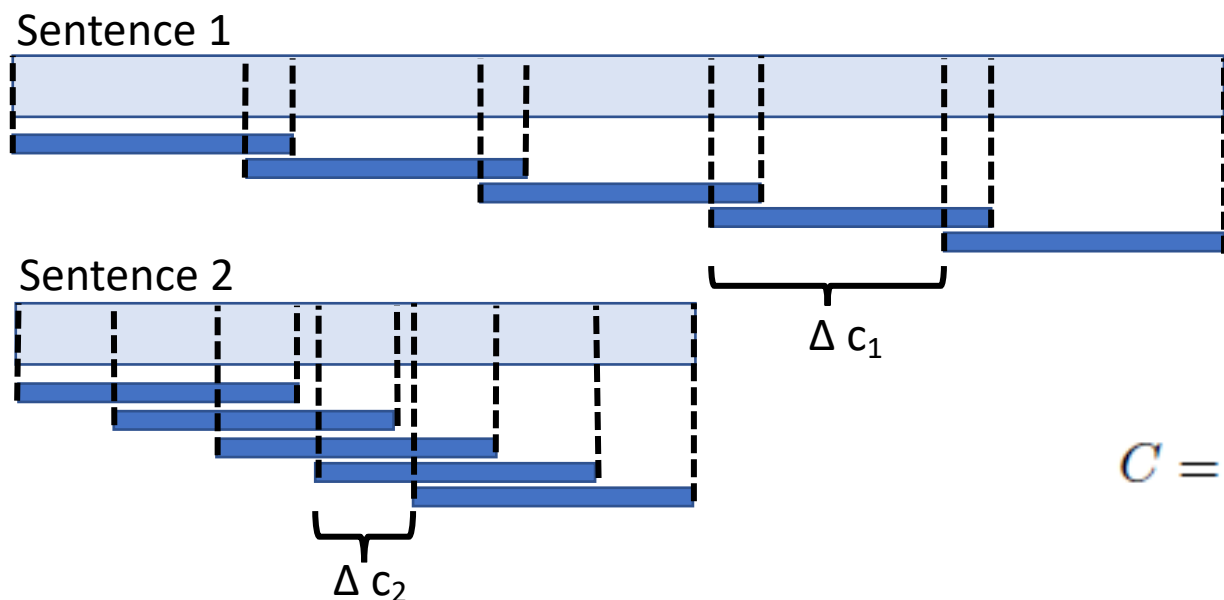
- Truncate original temporal information
- Mean pooling/majority voting/averaging outputs still treat every segment the same (i.e., **non-dynamic temporal modeling**)



- **Goal:** dynamic temporal modeling
- **Key Problem:** sentences have varied number of segments/frames
- **Proposed:** Novel Chunk Segmentation Process

**Fixed size and fixed number of chunks
for different duration sentences**

■ Chunk Segmentation Process



No zero-padding is required!


Pre-defined Parameters:


1. T_{max} (sec): maximum sentence duration in the corpus
2. w_c (sec): desired chunk window length

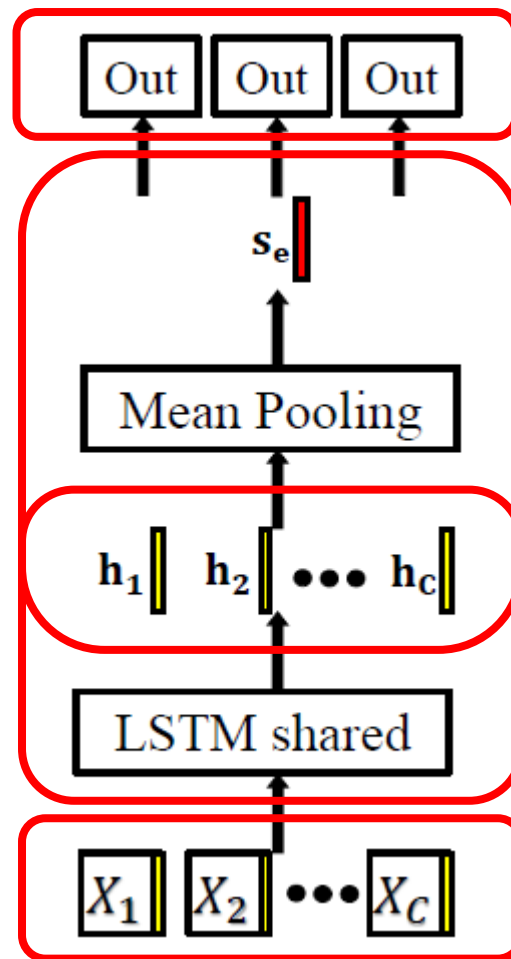
$$C = \left\lceil \frac{T_{max}}{w_c} \right\rceil : \text{number of chunks per sentence}$$

$$\Delta c_i = \frac{T_i - w_c}{C - 1} \text{ (sec): chunk step size depends on sentence duration}$$

Proposed Methodology

 : Chunk-level
feature representation

 : Sentence-level
feature representation



Step4: multi-task learning outputs

**Step3: sentence-level
feature representation**

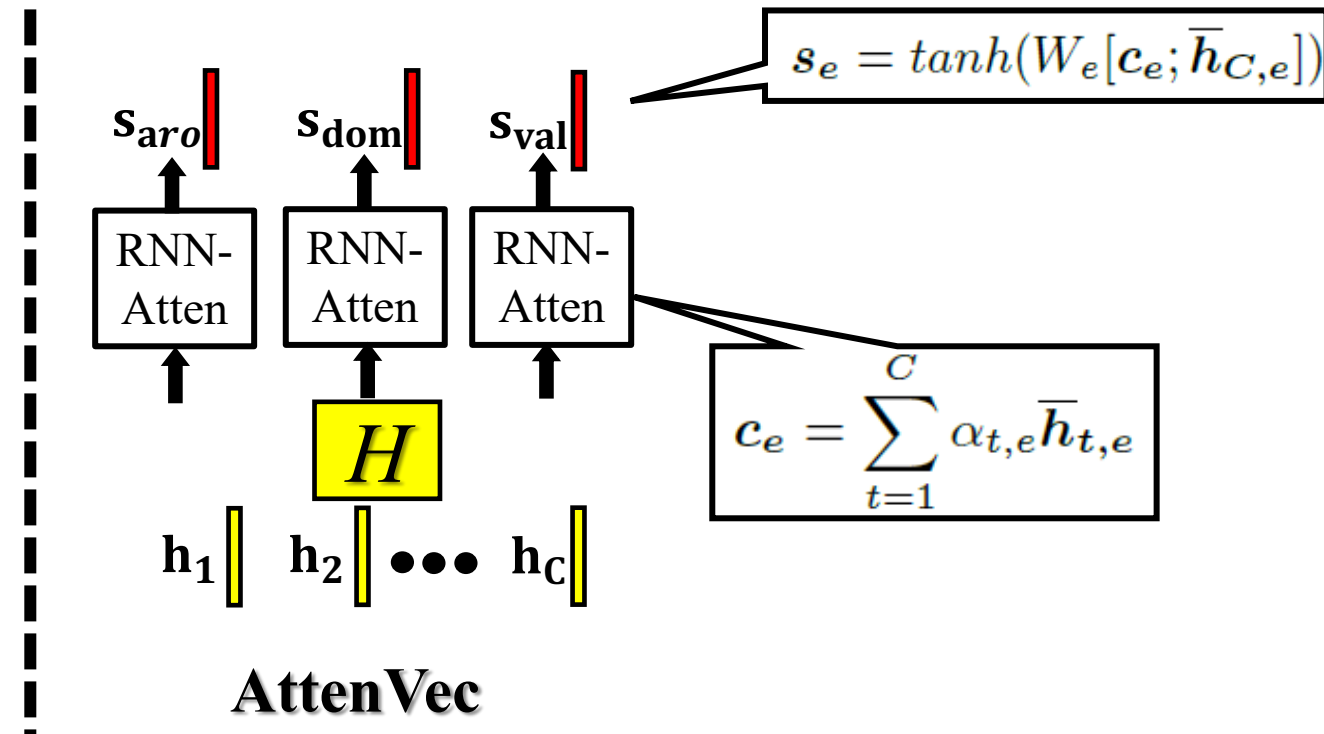
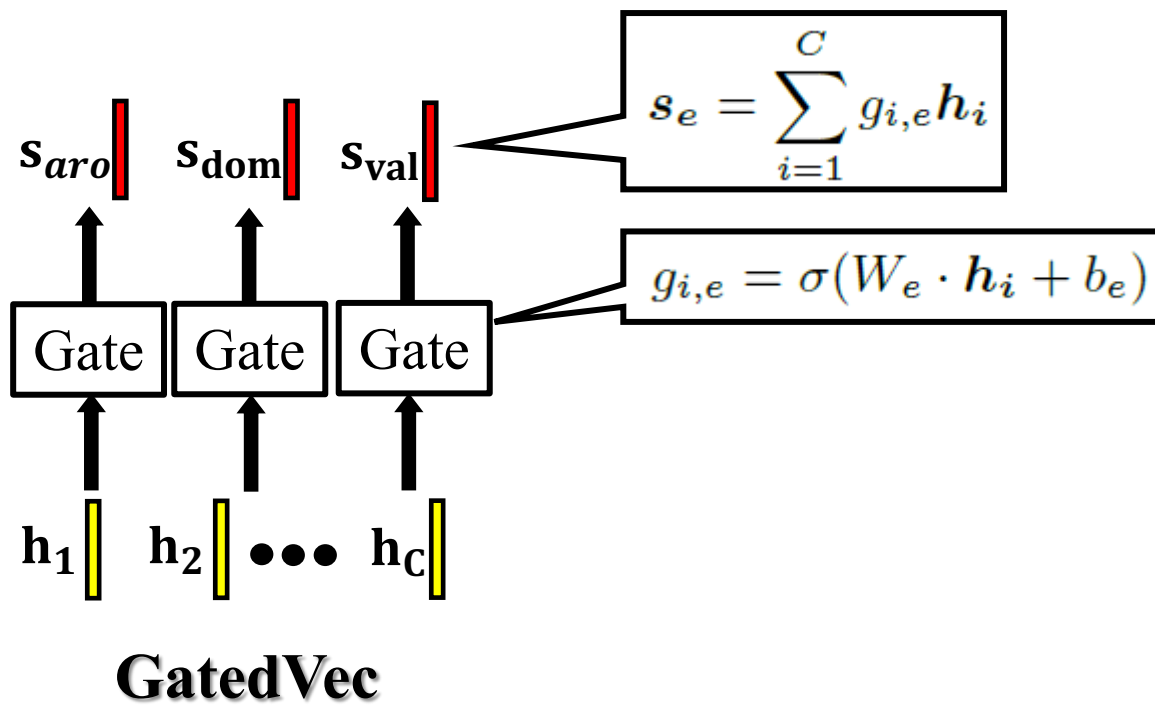
**Step2: chunk-level
feature representation**

Step1: chunk segmentation

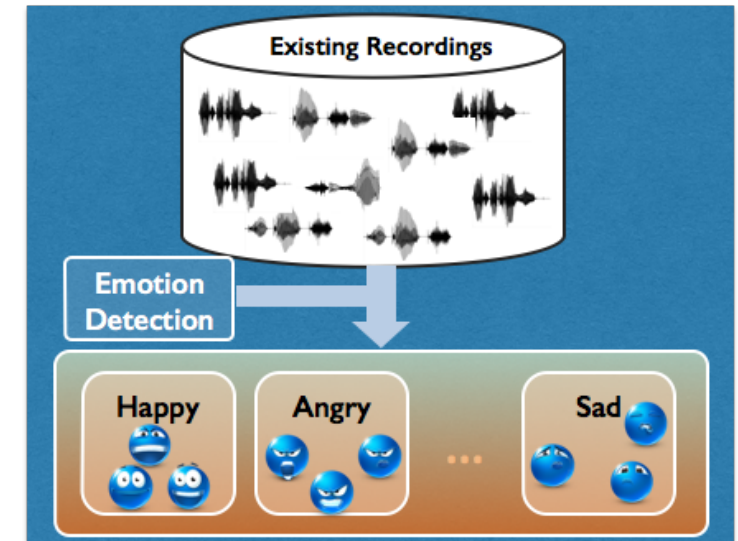
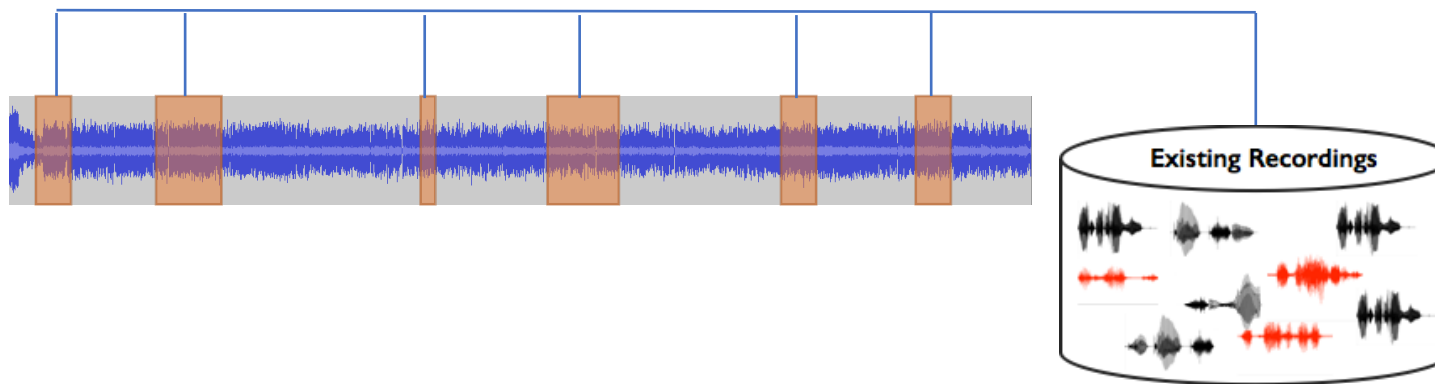
NonAtten

Proposed Methodology

- Dynamically Combining Chunk-Based Feature Representations

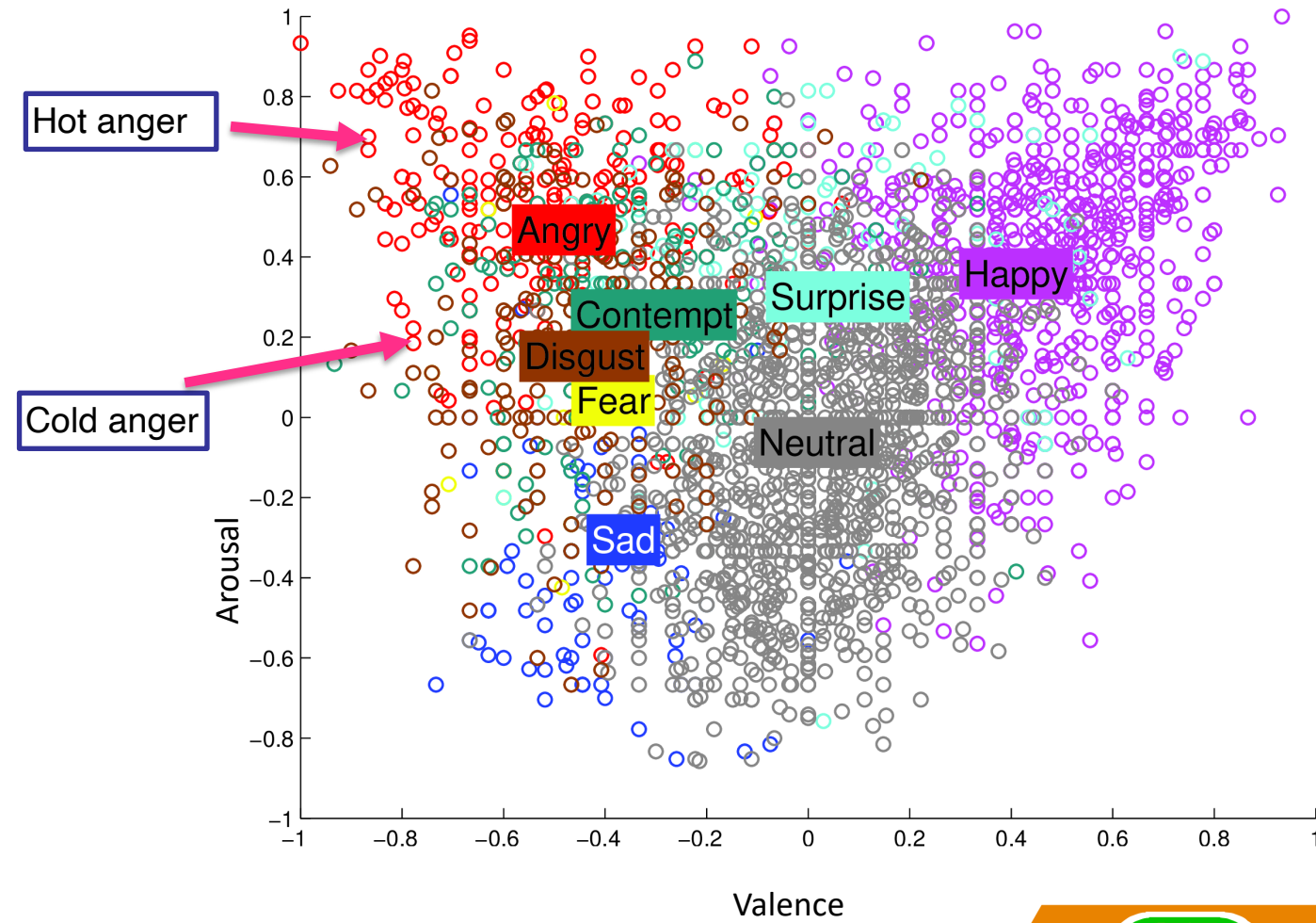


- Corpus: The MSP-Podcast v1.6
 - Use existing podcast recordings
 - Divide into speaker turns
 - Emotion retrieval to balance the emotional content
 - Annotate using crowdsourcing framework



Experimental Settings

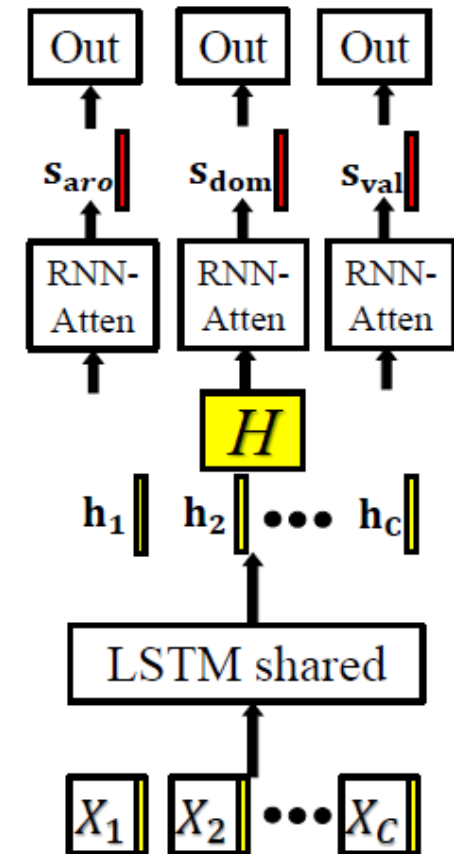
- The MSP-Podcast v1.6
 - 50,362 (83h,29m)
 - duration range: 2.75 ~ 11 secs
- Corpus partition with minimal speaker overlap sets:
 - Test data
 - 10,124 samples from 50 speakers (25 males, 25 females)
 - Validation data
 - 5,958 samples from 40 speakers (20 males, 20 females)
 - Train data
 - Remaining 34,280 samples



- Acoustic Features
 - Opensmile
 - Interspeech 2013 computational paralinguistics challenge (IS13ComparE)
 - 65 low level descriptors (LLDs) + its delta value = 130-dimensions in total
 - mel frequency cepstral coefficients (MFCCs)
 - Fundamental frequency (f0)
 - Intensity (energy)
 - ...

■ Parameters Settings:

- MTL tasks: arousal (Aro.), dominance (Dom.) and valence (Val.)
- w_c : 1 (sec)
- T_{max} : 11 (secs)
- $C = 11$ (chunks/per sentence)
- Network nodes: 130 for all layers (the same as input-dim)
- Adam optimizer
- 128 batch size
- Loss function: concordance correlation coefficient (CCC)
 - Same as the evaluation metric



■ Baseline Models:

- Padding zeros to the max length (i.e., 11 secs) for all sentences
- *LSTM(130)*, number of nodes in the LSTM shared layers is 130 number of nodes
- *LSTM(260)*, number of nodes in the LSTM shared layers is 130 number of nodes

■ *Best performance for all emotional attributes*

Model	Aro [CCC]	Dom [CCC]	Val [CCC]
<i>LSTM (130)</i>	0.6520	0.5711	0.2031
<i>LSTM (260)</i>	0.6875	0.6045	0.2847
<i>NonAtten</i>	0.6781	0.6019	0.2925
<i>GatedVec</i>	0.6747	0.5944	0.3199
<i>AttenVec</i>	0.6947	0.6132	0.3072

- **Improvement of model efficiency**

- *Chunks are parallel processed by GPU*
- *Significant reduction in MFLOPs (i.e., roughly $C=11$ times faster)*

Model	# of Par. [10^6]	MFLOPs [MFLOPS]	Train [sec/epoch]	Online [ms/uttr]
<i>LSTM (130)</i>	0.323	5.67	437.1	547.5
<i>LSTM (260)</i>	1.052	18.49	439.4	598.1
<i>NonAtten</i>	0.323	0.49	74.9	42.2
<i>GatedVec</i>	0.324	0.49	246.6	44.6
<i>AttenVec</i>	0.577	1.50	353.1	45.6

Experimental Results

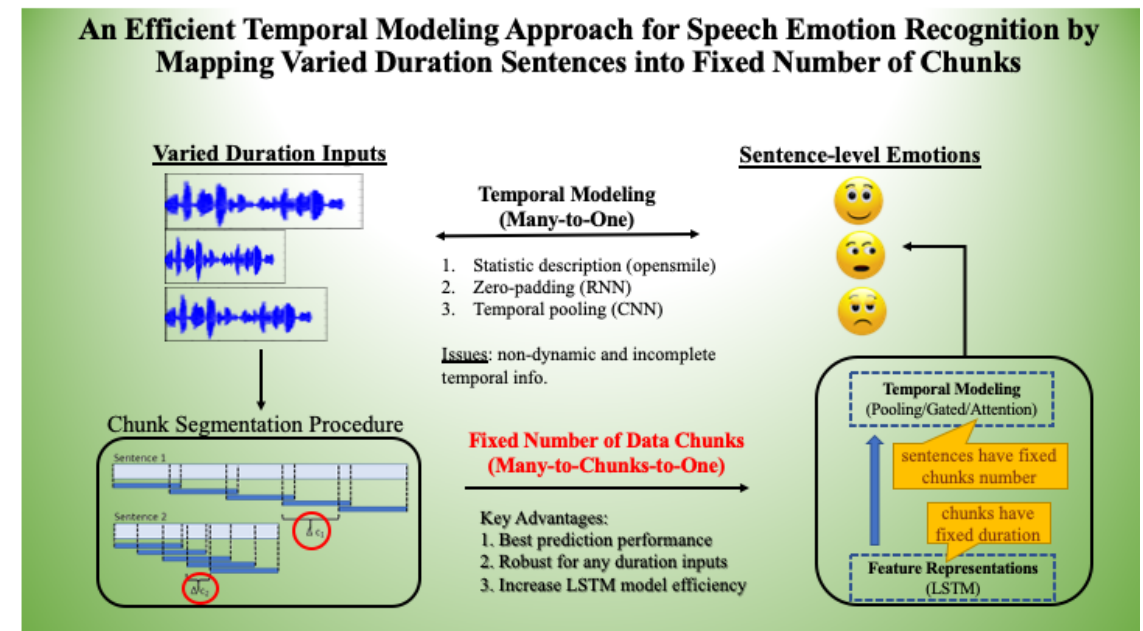
- Robustness for different duration
 - Short (≤ 5 secs): 4,280 sentences
 - Middle (5-8 secs): 3,684 sentences
 - Long (≥ 8 secs): 2,160 sentences
- ***Robust to different duration inputs (especially for long sequences)***

Short (< 5sec)	Aro-CCC	Dom-CCC	Val-CCC
LSTM(130)	0.6636	0.5812	0.2389
NonAtten	0.6761	0.6077	0.3129
GatedVec	0.6621	0.5865	0.3263
AttenVec	0.7003	0.6192	0.3363

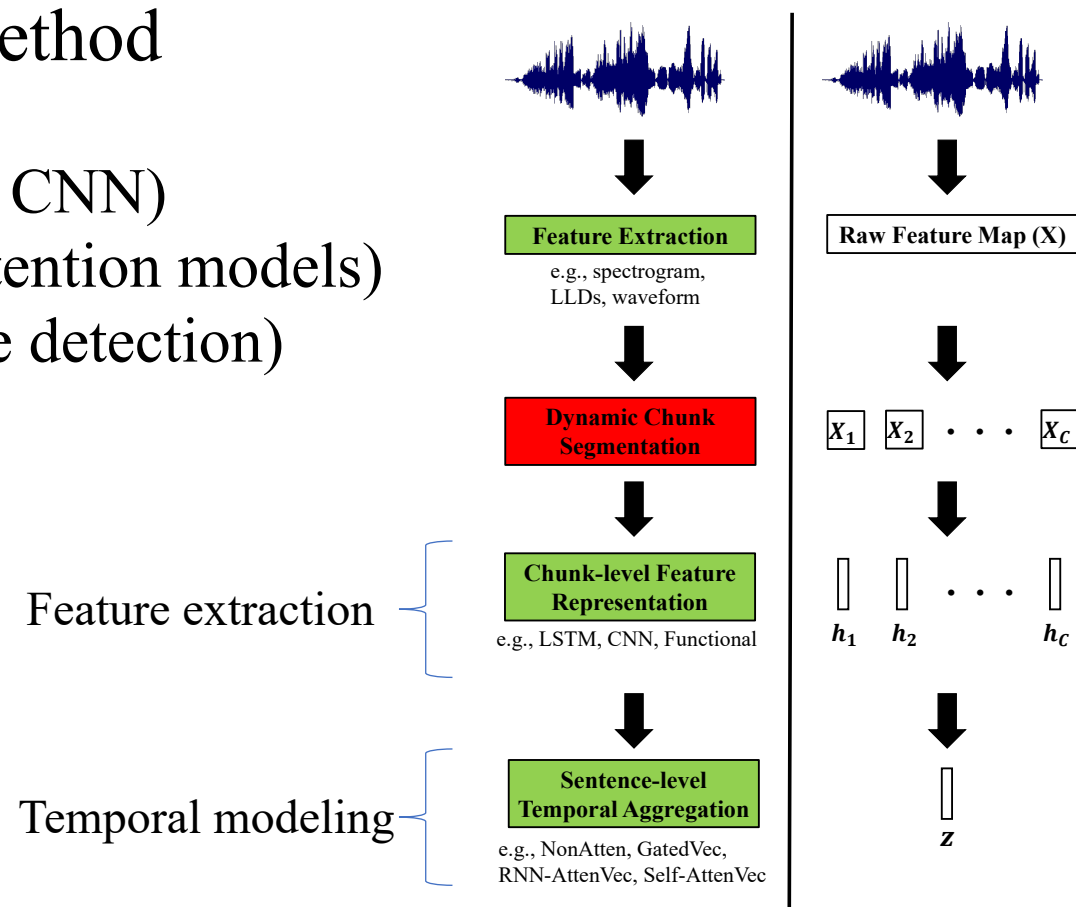
Middle (5~8sec)	Aro-CCC	Dom-CCC	Val-CCC
LSTM(130)	0.6484	0.5642	0.1735
NonAtten	0.6779	0.6071	0.2839
GatedVec	0.6807	0.6042	0.3279
AttenVec	0.6880	0.6129	0.2978

Long (> 8sec)	Aro-CCC	Dom-CCC	Val-CCC
LSTM(130)	0.6314	0.5559	0.1737
NonAtten	0.6811	0.5822	0.2331
GatedVec	0.6933	0.6030	0.2835
AttenVec	0.6912	0.5989	0.2539

- Novel segmentation approach that can split a sentence into a **fixed number** of chunks, which have **the same** duration
- **Flexibly** and **dynamically** combine temporal information
- Best prediction **performance**
- Improve model **efficiency**
- **Robust** for different duration sentences



- General framework of the proposed method
 - Multiple datasets
 - Different feature extraction models (e.g., CNN)
 - Different temporal modeling (LSTM, Attention models)
 - Different sequence-to-one tasks (e.g., age detection)



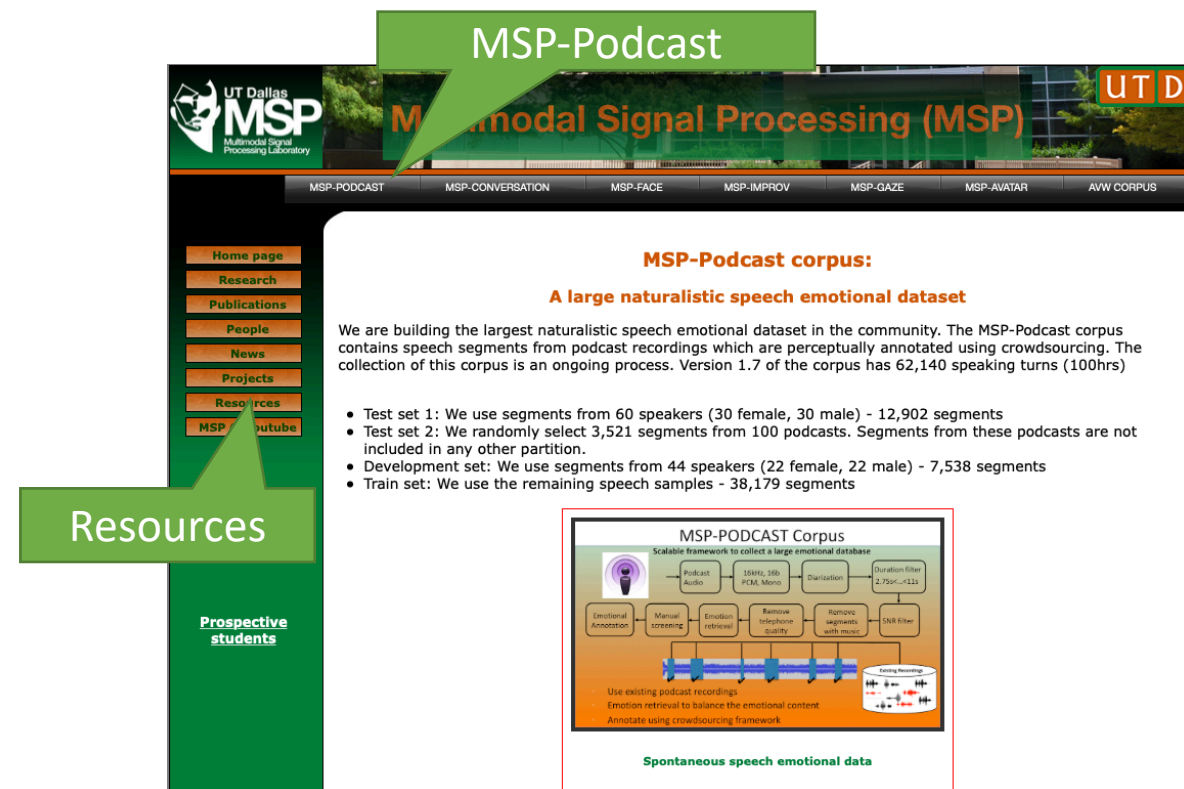
Release of the MSP-Podcast Corpus

Academic license

- Federal Demonstration Partnership (FDP) Data Transfer and Use Agreement
- Free access to the corpus

Commercial license

- Commercial license through UT Dallas



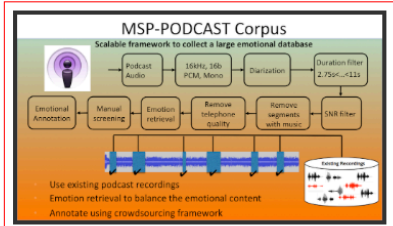
MSP-Podcast

Resources

MSP-Podcast corpus:
A large naturalistic speech emotional dataset

We are building the largest naturalistic speech emotional dataset in the community. The MSP-Podcast corpus contains speech segments from podcast recordings which are perceptually annotated using crowdsourcing. The collection of this corpus is an ongoing process. Version 1.7 of the corpus has 62,140 speaking turns (100hrs)

- Test set 1: We use segments from 60 speakers (30 female, 30 male) - 12,902 segments
- Test set 2: We randomly select 3,521 segments from 100 podcasts. Segments from these podcasts are not included in any other partition.
- Development set: We use segments from 44 speakers (22 female, 22 male) - 7,538 segments
- Train set: We use the remaining speech samples - 38,179 segments



Spontaneous speech emotional data

<https://msp.utdallas.edu>

**Thank you for your
attention !**

This work was funded by NSF
(CNS-1823166; IIS-1453781)



Questions or Contact:
wei-cheng.lin@utdallas.edu