



Innovative R&D by NTT

# **Soft-Target Training with Ambiguous Emotional Utterances for DNN-based Speech Emotion Classification**

**NTT Media Intelligence Laboratories, NTT Corporation**

**Atsushi Ando, Satoshi Kobashikawa, Hosana Kamiyama,  
Ryo Masumura, Yusuke Ijima, Yushi Aono**

# Summary



## Purpose

- ✓ Speech emotion classification from acoustic features
  - Task: 4-class classification (*Neutral, Happy, Sad, Angry*)

## Novelty

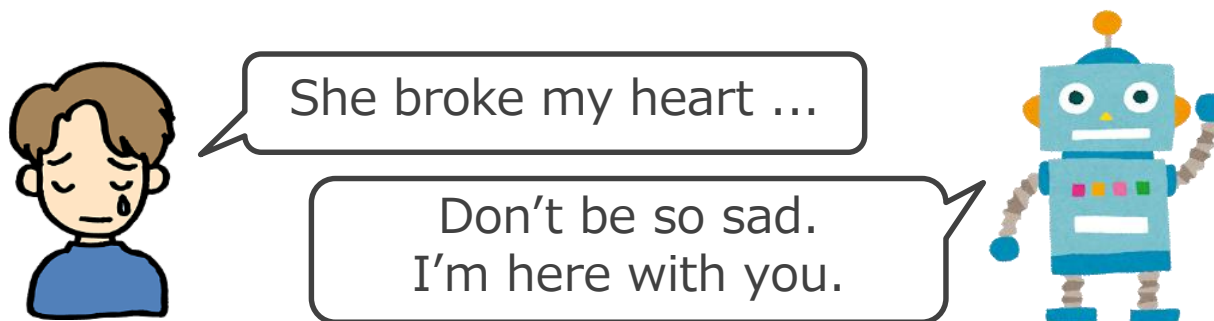
- ✓ To mitigate training data limitation problem, utilizing **ambiguous emotional utterances** (no target emotions are dominant) **which are ignored in the conventional methods**
  - Employ two types of soft-target training

## Results

- ✓ Performance improved
  - Overall Accuracy: 58.6% → 62.6%, Average Recall: 53.7% → 63.7%

## Speech emotion recognition is important technology to understand natural speech

- ✓ Application : “**sympathetic**” spoken dialog system

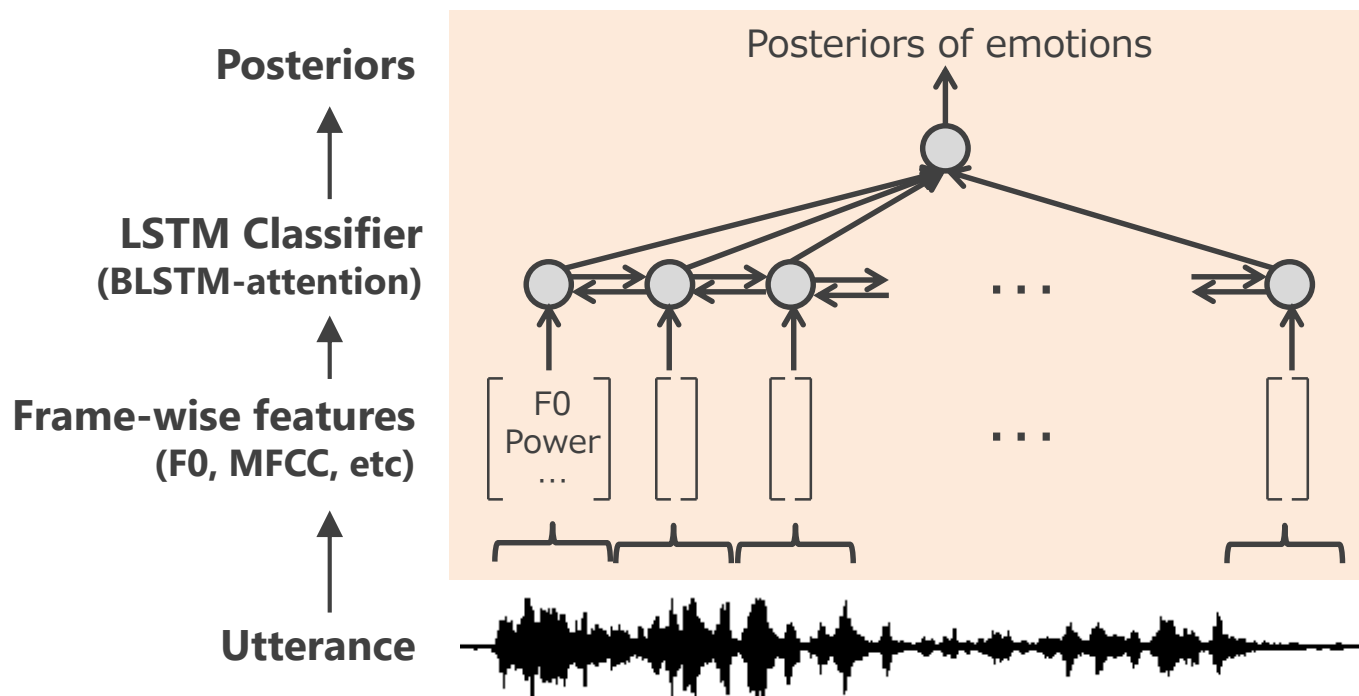


- ✓ Task description

- Input : short utterance (1~10 sec.)
- Target : 4-class speech emotion (*Neutral, Happy, Sad, Angry*)

## Frame-wise acoustic features + BLSTM-RNNs

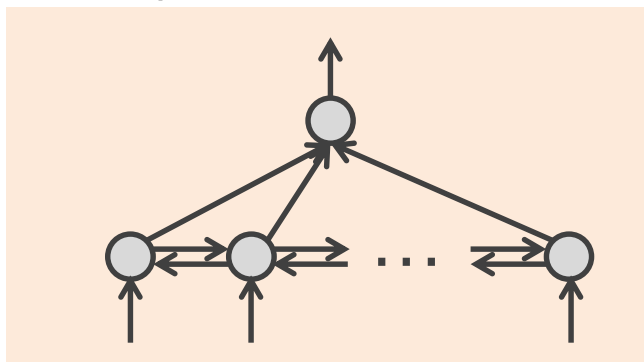
- ✓ Emotion classification by BLSTM w/ attention [Mirsamadi+, 17]
  - Utilizing **local characteristics** of emotions



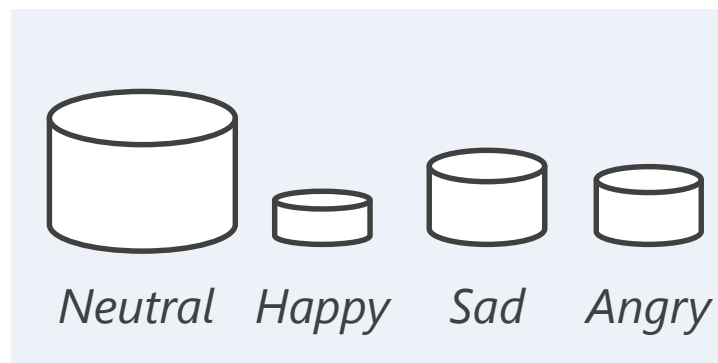
## Training data is usually limited

- ✓ Emotion classification by BLSTM w/ attention [Mirsamadi+, 17]

# of parameters: **100k~**



# of train data: **~5k**



→ Classifier is overfitted / less generalized

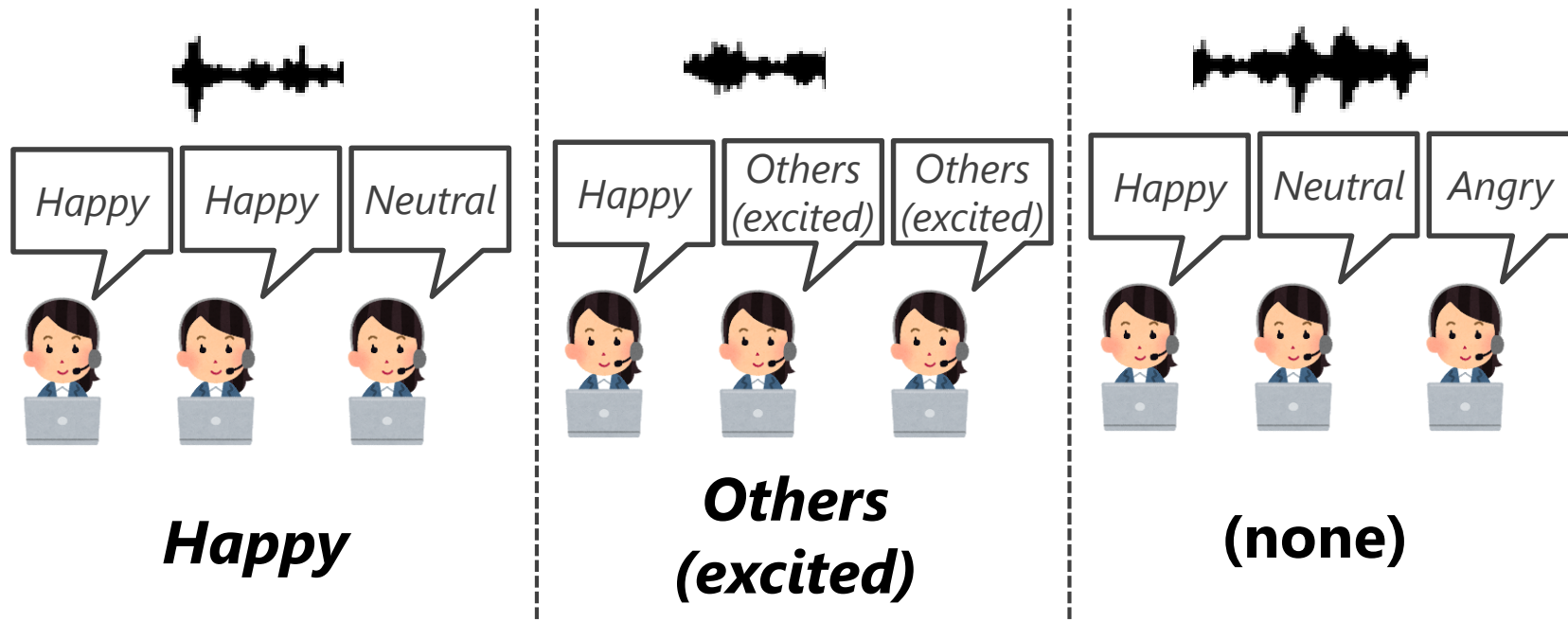
**Issue How to train complex classifier from limited data ?**

# Problem - Why limited?



Ground truths are decided by several annotators.  
**Some utterances are ignored for training**

- ✓ Ground truth = **Dominant emotion** of annotations



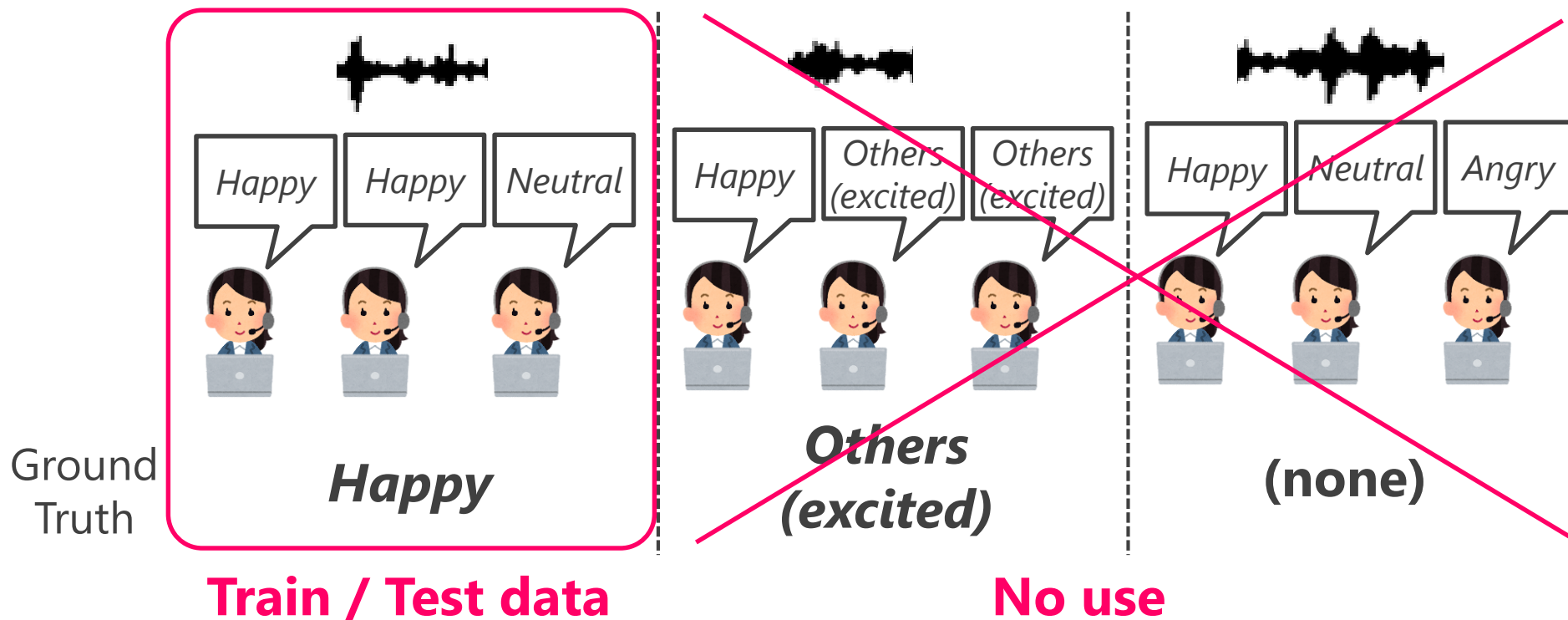
Ground Truth

# Problem - Why limited?



Ground truths are decided by several annotators.  
**Some utterances are ignored for training**

- ✓ Ground truth = **Dominant emotion** of annotations



# Approach (1/2)

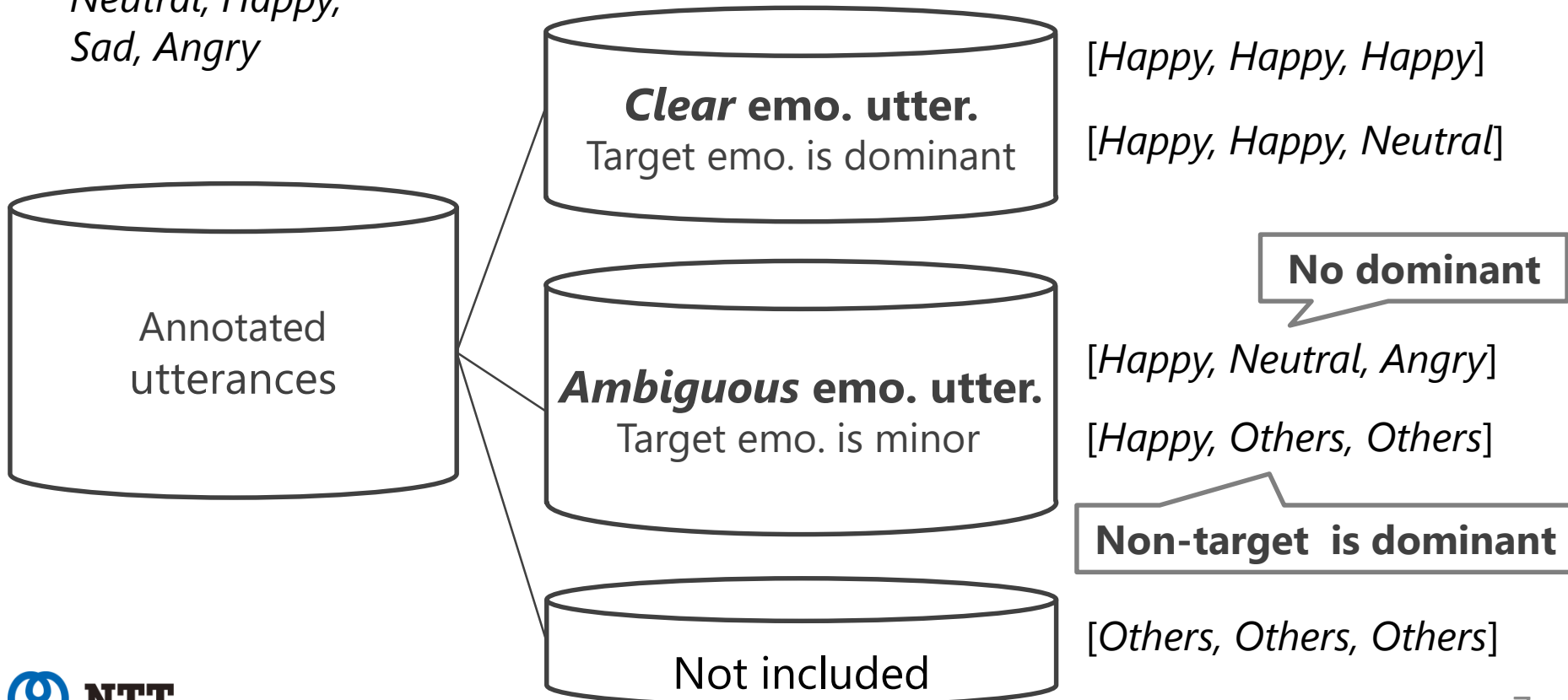


Utilize **ambiguous emotional utterances** (target emo. are minor) to mitigate training data limitation

## Target emotions

*Neutral, Happy,  
Sad, Angry*

## Annotation example





# Approach (1/2)

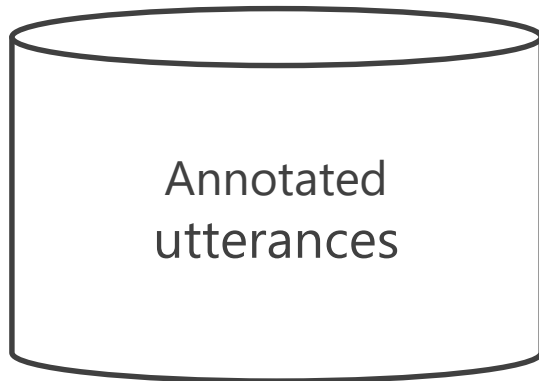


Utilize **ambiguous emotional utterances** (target emo. are minor) to mitigate training data limitation

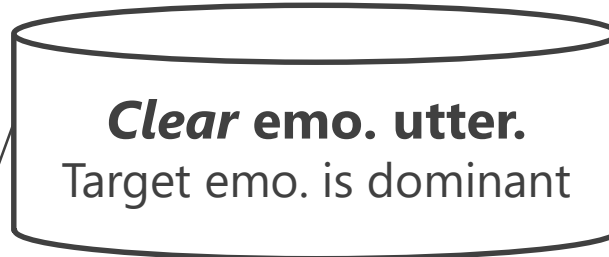
Target emotions

*Neutral, Happy,  
Sad, Angry*

Annotation example

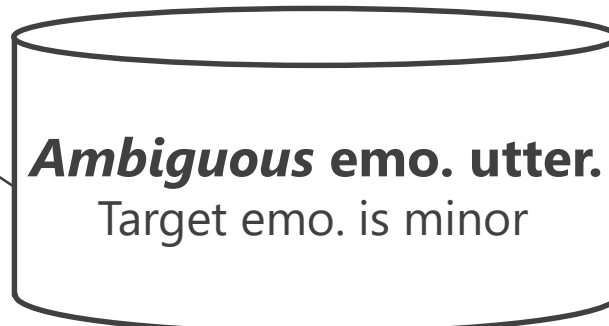


**Conventional training**



*[Happy, Happy, Happy]*

*[Happy, Happy, Neutral]*



*[Happy, Neutral, Angry]*

*[Happy, Others, Others]*



*[Others, Others, Others]*

# Approach (1/2)

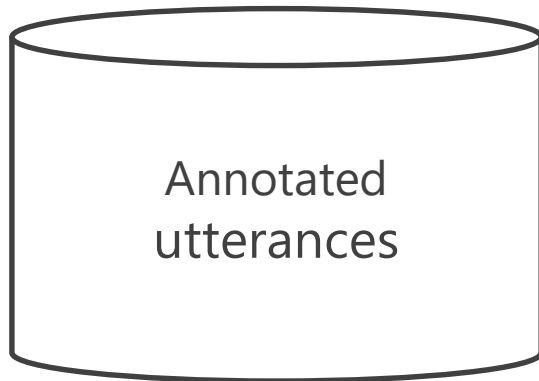


Utilize **ambiguous emotional utterances** (target emo. are minor) to mitigate training data limitation

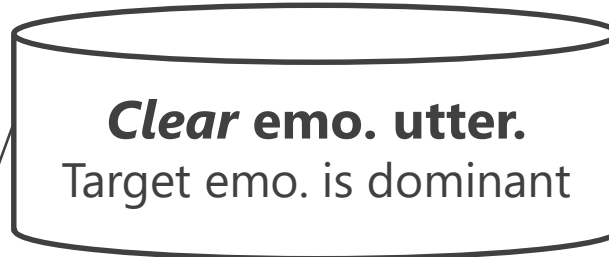
Target emotions

*Neutral, Happy,  
Sad, Angry*

Annotation example

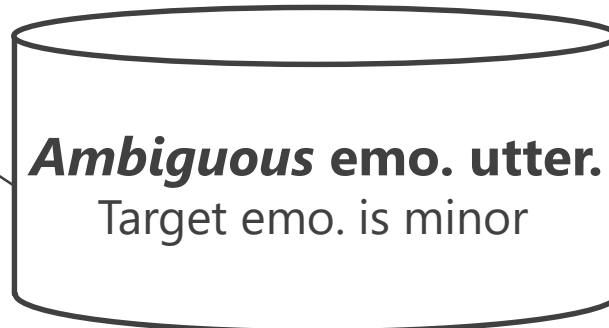


**Conventional training**



*[Happy, Happy, Happy]*

*[Happy, Happy, Neutral]*



*[Happy, Neutral, Angry]*

*[Happy, Others, Others]*



Are there no **Happy** characteristics ?

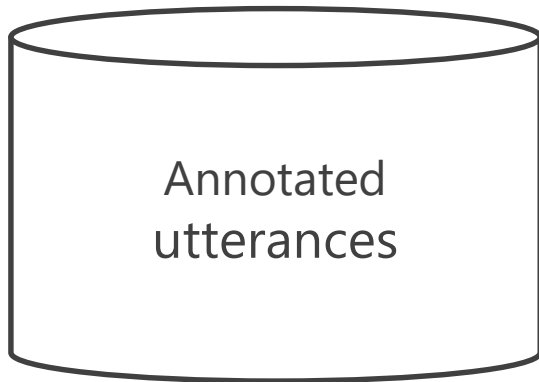
# Approach (1/2)



Utilize **ambiguous emotional utterances** (target emo. are minor) to mitigate training data limitation

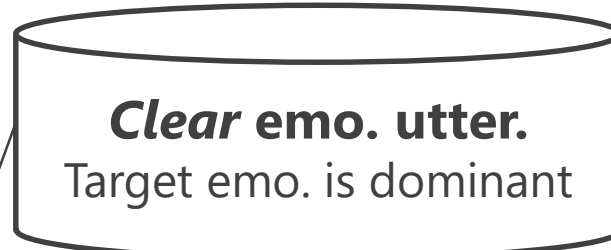
Target emotions

*Neutral, Happy,  
Sad, Angry*



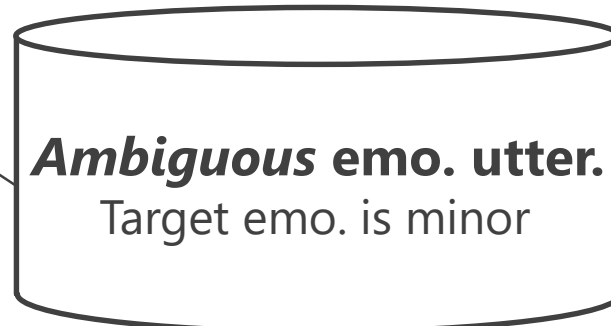
**Proposed training**

**Conventional training**



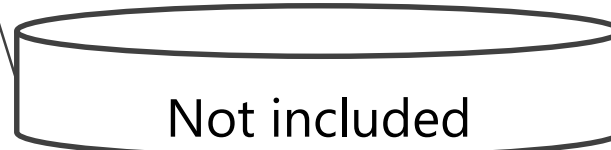
*[Happy, Happy, Happy]*

*[Happy, Happy, Neutral]*



*[Happy, Neutral, Angry]*

*[Happy, Others, Others]*



*[Others, Others, Others]*

# Approach (2/2)



**Control discriminativity** to handle both *clear* and *ambiguous* emotional utterances

**High discriminativity**

Train as **definitely *Happy***

***Clear* emo. utter.**

Target emo. is dominant

[*Happy, Happy, Happy*]

[*Happy, Happy, Neutral*]

**Low discriminativity**

Train as **maybe *Happy***

***Ambiguous* emo. utter.**

Target emo. is minor

[*Happy, Neutral, Angry*]

[*Happy, Others, Others*]

Not included

[*Others, Others, Others*]

## Soft-target training is employed to deal *clear/ambiguous* emotional utterances

### ✓ Two types of soft-target

#### 1. Soft-target [Fayek+, 16]

$$\underline{q(c_k)} = \frac{\sum_n h_k^{(n)}}{\sum_k \sum_n h_k^{(n)}}$$

Annotation frequency (sum=1)

$h_k^{(n)}$  : Binary label-existence (0/1)  
 $n$ -th annotator,  $k$ -th emotion class

$K$  : Total emotion classes

#### 2. Modified soft-target

$$\underline{q(c_k)} = \frac{\alpha + \sum_n h_k^{(n)}}{\alpha K + \sum_k \sum_n h_k^{(n)}}$$

**Additive smoothed form** of conventional soft-target

$\alpha$  : Smoothing coefficient

### ✓ Model parameters are updated by cross-entropy loss

$$L = - \sum_{k=1}^K \underline{q(c_k)} \log p(c_k | \mathbf{X}, \theta)$$

# Proposed: modified soft-target



Modified soft-target is suitable to represent *ambiguous* emotional utterances

✓ Examples of teachers  $q(c_k)$

	Hard-target	Soft-target [Fayek+,16]	Modified Soft-target
[Happy, Happy, Happy]			
[Happy, Happy, Neutral]			
[Happy, <u>Others</u> , <u>Others</u> ]	(no use)		

Non-target

(Smoothing coeff.  $\alpha = 1$ )

# Proposed: modified soft-target



Modified soft-target is suitable to represent *ambiguous* emotional utterances

✓ Examples of teachers  $q(c_k)$

	Hard-target	Soft-target [Fayek+,16]	Modified Soft-target
[Happy, Happy, Happy]			
[Happy, Happy, Neutral]			
[Happy, <u>Others</u> , <u>Others</u> ]	(no use)		

Non-target

**Ambiguous utterances are discarded**

(Smoothing coeff.  $\alpha = 1$ )

# Proposed: modified soft-target



Modified soft-target is suitable to represent *ambiguous* emotional utterances

✓ Examples of teachers  $q(c_k)$

	Hard-target	Soft-target [Fayek+,16]	Modified Soft-target
[Happy, Happy, Happy]			
[Happy, Happy, Neutral]			
[Happy, <u>Others</u> , <u>Others</u> ]	(no use)		

Non-target

Allocate same teacher labels to clear/ambiguous (Smoothing coeff.  $\alpha = 1$ )



# Proposed: modified soft-target



Modified soft-target is suitable to represent *ambiguous* emotional utterances

✓ Examples of teachers  $q(c_k)$

	Hard-target	Soft-target [Fayek+,16]	Modified Soft-target
[Happy, Happy, Happy]			
[Happy, Happy, Neutral]			
[Happy, <u>Others</u> , <u>Others</u> ]	(no use)		

Non-target

Lower discriminativity in ambiguous emo. uttr.

Modified soft-target is regarded as **Maximum a posteriori (MAP) estimation** from annotations

Utterance



"true" distribution  
of target emo.



Annotations

Sampling

(N=# of annotations)

[*Happy, Happy, Sad*]

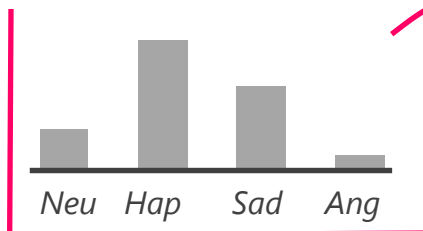
Objective function of the model

Modified soft-target is regarded as **Maximum a posteriori (MAP) estimation** from annotations

## Utterance



"true" distribution of target emo.



## Annotations

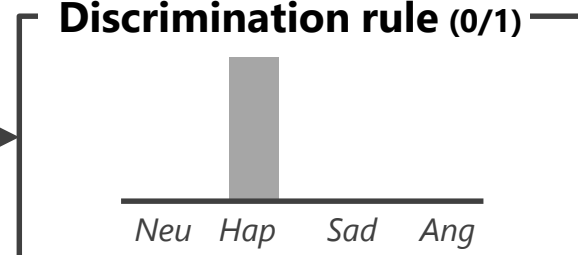
Sampling  
(N=# of annotations)

[Happy, Happy, Sad]

## Objective function of the model

Discrimination rule (0/1)

Hard-target



# Interpretation



Innovative R&D by NTT

Modified soft-target is regarded as **Maximum a posteriori (MAP) estimation** from annotations

## Utterance



"true" distribution of target emo.



## Annotations

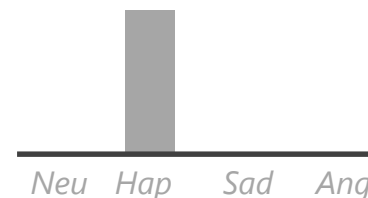
Sampling  
( $N = \#$  of annotations)

[Happy, Happy, Sad]

## Objective function of the model

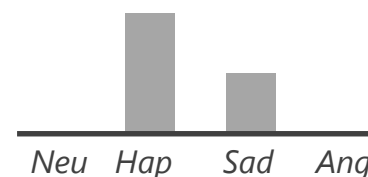
Hard-target

Discrimination rule (0/1)



Soft-target

ML-based distribution



Modified soft-target

MAP-based distribution



# Interpretation



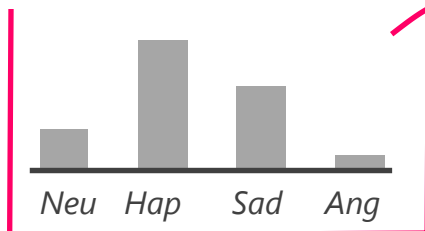
Innovative R&D by NTT

Modified soft-target is regarded as **Maximum a posteriori (MAP) estimation** from annotations

## Utterance



"true" distribution of target emo.



## Annotations

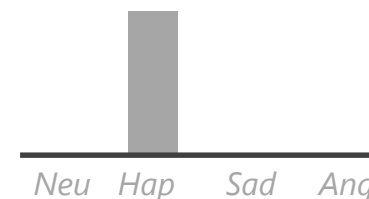
Sampling  
( $N = \#$  of annotations)

[Happy, Happy, Sad]

## Objective function of the model

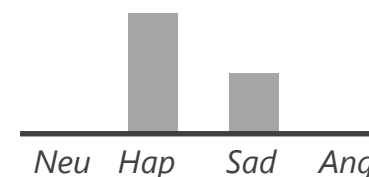
Hard-target

Discrimination rule (0/1)



Soft-target

ML-based distribution



Modified soft-target

MAP-based distribution



Uniform prior

## ✓ Purpose

1. Evaluate effectiveness of *ambiguous* emotional utterances for train
2. Compare teacher labels (hard / soft / modified soft)

## ✓ Dataset: IEMOCAP [Busso+, 08]

- **Task:** 2-speaker dialogue (1 male, 1 female)
- **# of speakers:** 10 (train: 8, test: 2)
- **# of annotators:** 3

*frustrated, excited,  
surprised, fear,  
disgust, no-dominant*

		Total	# of utterances (dominant emotion)				
			<i>Neutral</i>	<i>Happy</i>	<i>Sad</i>	<i>Angry</i>	<i>Others</i>
Train	<i>clear</i>	3548	1324	460	890	874	-
	<i>ambiguous</i>	3693	0	0	0	0	3693
Test		942	384	135	194	229	-

## ✓ **Classifier:** BLSTM + attention [Mirsamadi+,17]

### – **Structure**

➤ Full256-BLSTM128-attention-Full256

### – **Input:** frame-wise acoustic features, 47 dims.

➤ MFCC12,  $\Delta$ MFCC12,  $\Delta\Delta$ MFCC12,  
Loudness,  $\Delta$ Loudness,  $\Delta\Delta$ Loudness,  
F0, VoiceProb, ZCR, HNR,  $\Delta$ F0,  $\Delta$ VoiceProb,  $\Delta$ ZCR,  $\Delta$ HNR

- ### – **Teacher:**
- ① Hard-target
  - ② Soft-target [Fayek+, 16]
  - ③ Modified soft-target
- } **baseline**

### – **Train data:** *clear / ambiguous / clear + ambiguous*

## ✓ **Evaluation measures**

– Weighted Accuracy (WA): overall accuracy

– Unweighted Accuracy (UA): average recall of emotion classes

➤ Average results of 5 trials of training

Moderate performance with *ambiguous* data alone,  
and best with *clear + ambiguous* data

	Teacher	Train set		Accuracy [%]	
		<i>clear</i>	<i>ambig.</i>	WA	UA
MajorityClass (All Neutral)				40.8	25.0
Baseline	Hard-target	✓		58.6	53.7
	Soft-target	✓		58.1	54.9
Proposed	Modified soft-target	✓		58.5	57.4
			✓	53.6	54.0
		✓	✓	<b>62.6</b>	<b>63.7</b>

Overall Acc.

Avg. Recall



Moderate performance with *ambiguous* data alone,  
and best with *clear + ambiguous* data

	Teacher	Train set		Accuracy [%]	
		<i>clear</i>	<i>ambig.</i>	WA	UA
MajorityClass (All Neutral)				40.8	25.0
Baseline	Hard-target	✓		58.6	53.7
	Soft-target	✓		58.1	54.9
Proposed	Modified soft-target	✓		58.5	57.4
			✓	53.6	54.0
		✓	✓	<b>62.6</b>	<b>63.7</b>

**Moderate performance**  
even they have been ignored for training!

Moderate performance with *ambiguous* data alone,  
and best with *clear + ambiguous* data

	Teacher	Train set		Accuracy [%]	
		<i>clear</i>	<i>ambig.</i>	WA	UA
MajorityClass (All Neutral)				40.8	25.0
Baseline	Hard-target	✓		58.6	53.7
	Soft-target	✓		58.1	54.9
Proposed	Modified soft-target	✓		58.5	57.4
			✓	53.6	54.0
		✓	✓	<b>62.6</b>	<b>63.7</b>

Best performance

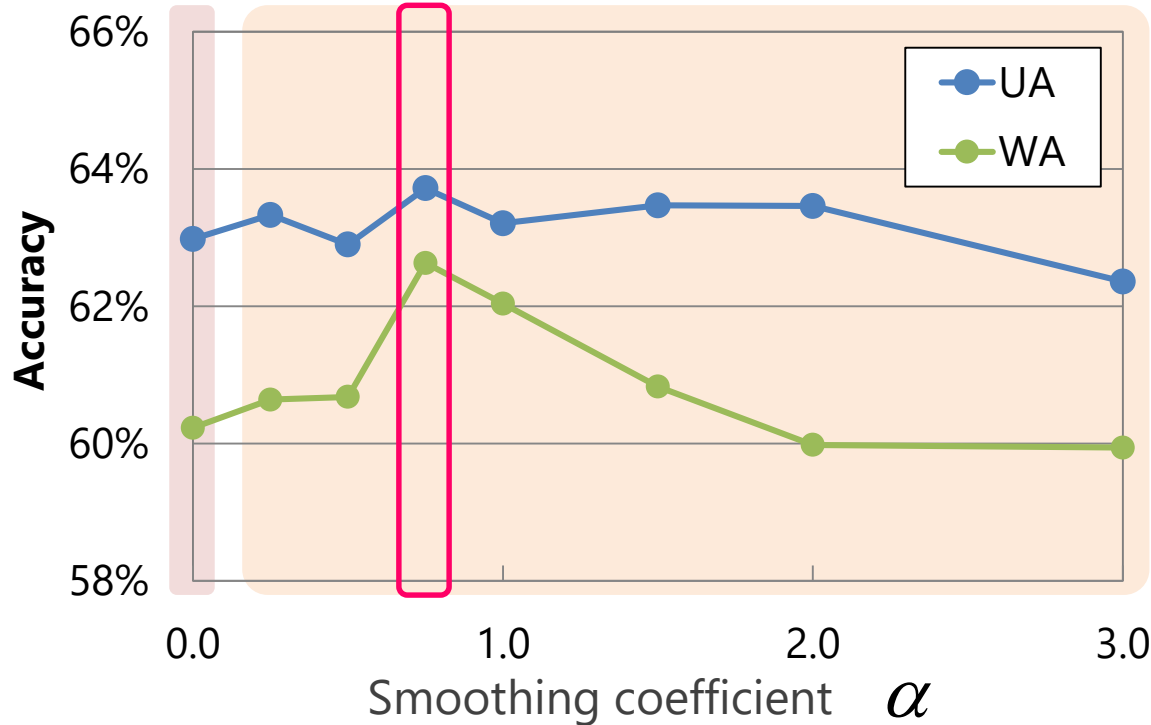
# Comparisons of teacher labels



**Modified soft-target with smoothing coeff. = 0.75 is better than (conventional) soft-target**

Soft-target

Modified soft-target



## Setup

Train: *clear + ambig.*  
Model: BLSTM-att

## ✓ Summary

- **Purpose:** emotion classification from acoustic features
- **Approach:** Utilizing *ambiguous* emotional utterances to mitigate training data limitation problem
- **Method:** Soft-target training which deals both *clear* and *ambiguous* emotional utterances in same criteria
  - Equal to ML/MAP estimation of true emotion distributions
- **Results:** Performances were improved (WA 58.6→62.6%)  
Show the effectiveness of *ambiguous* data for training

## ✓ Future works

- Evaluations by other corpus / emotion set
- Improve modified soft-target (prior distribution of MAP estimation)