



# Pronunciation Error Detection using DNN Articulatory Model based on Multi-lingual and Multi-task Learning

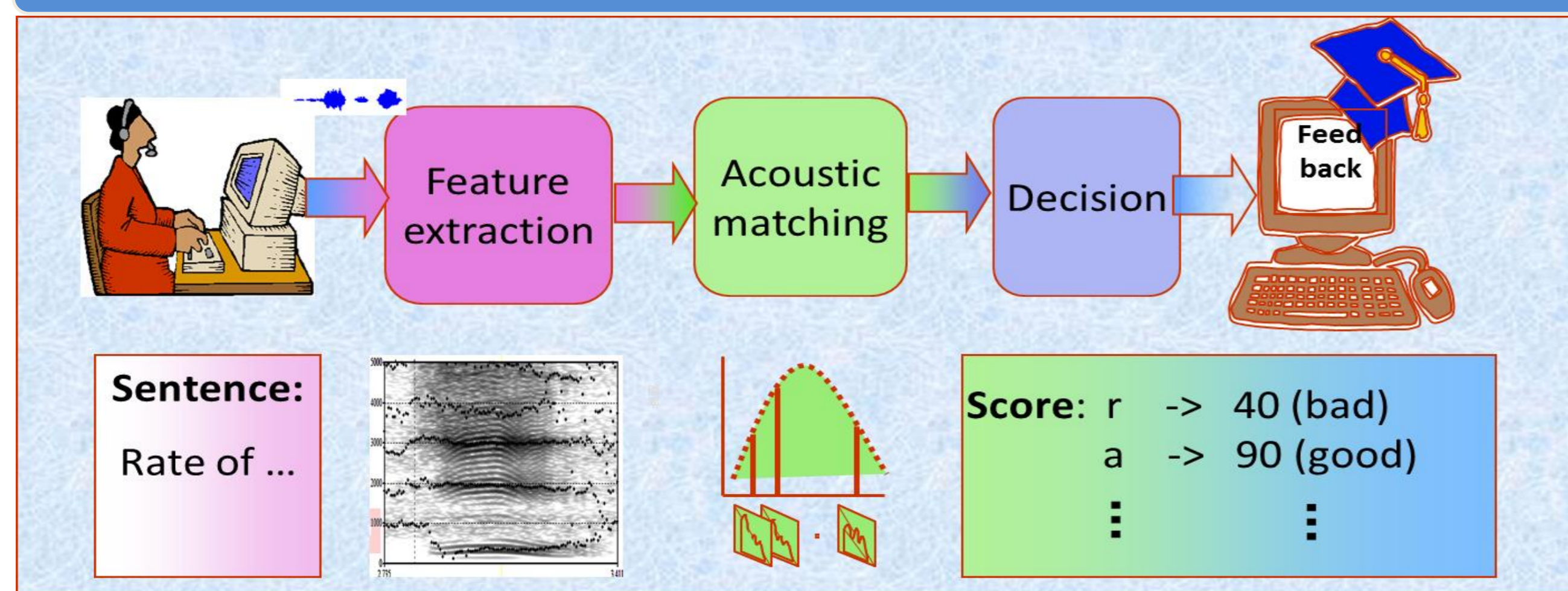
Richeng Duan<sup>1</sup> Tatsuya Kawahara<sup>1</sup> Masatake Dantsuji<sup>2</sup> Jinsong Zhang<sup>3</sup>

<sup>1</sup>School of Informatics, Kyoto University

<sup>2</sup>Academic Center for Computing and Media Studies, Kyoto University

<sup>3</sup>School of Information Science, Beijing Language and Culture University

## Introduction



■ Providing feedbacks directly related with articulation.

## Challenge & Proposed method

- **Challenge**
  - Non-native corpus collection in a large scale is not easy.
  - Precisely annotating non-native speech is difficult.
- **Proposed method**
  - Modeling articulatory attributes **without** non-native training data.
  - Enhancing articulatory models with multi-task learning.
  - Learning better feature representation using Multi-lingual learning.

## Definition of articulatory attribute

CH (Chinese) learning by JP (Japanese) students

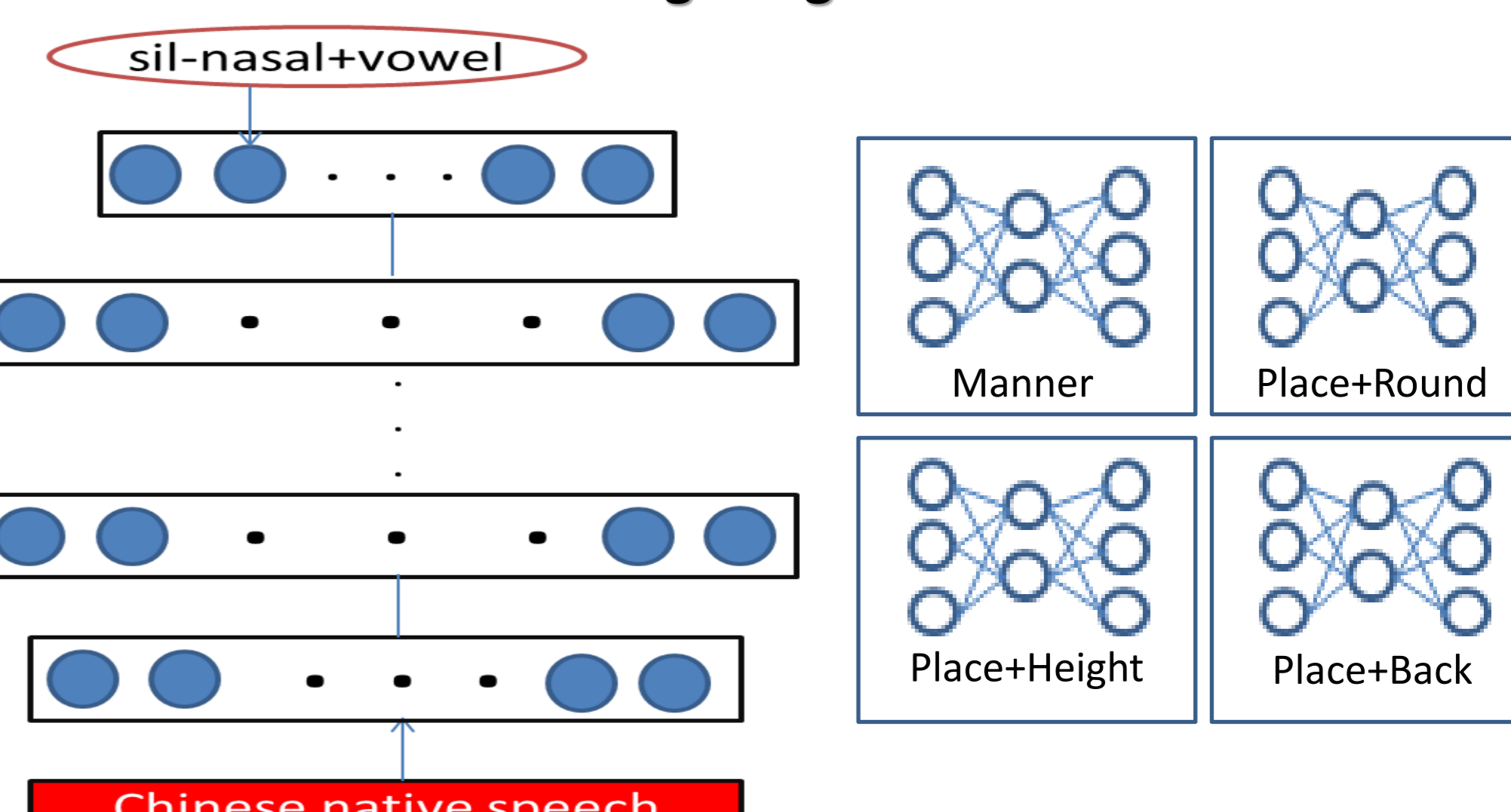
Manner	Phone set	Attribute	Phone set
Aspirated-stop	CH: p t k	Backless	CH: i ü JP: i e
Unaspirated-stop	CH: b d g		CH: a JP: a u
nasal	CH: m n		CH: e u o JP: o
Unvoiced-fricative	JP: m n N	Height	CH: i u ü JP: i u
	CH: f s sh		CH: o e JP: e o
Unvoiced-stop	JP: p t k	Low	CH: a JP: a
Voiced-stop	JP: b d g	Roundedness	CH: a i e JP: a i e
Place	Phone set		CH: o u ü JP: o
Retroflex	CH: zh ch sh r	Unroundedness	
Bilabial	CH: b p m	Roundedness	
	JP: b p m		
Velar	CH: g k h		
	JP: g k		
glottal	JP: h		

## Context-dependent attribute modeling

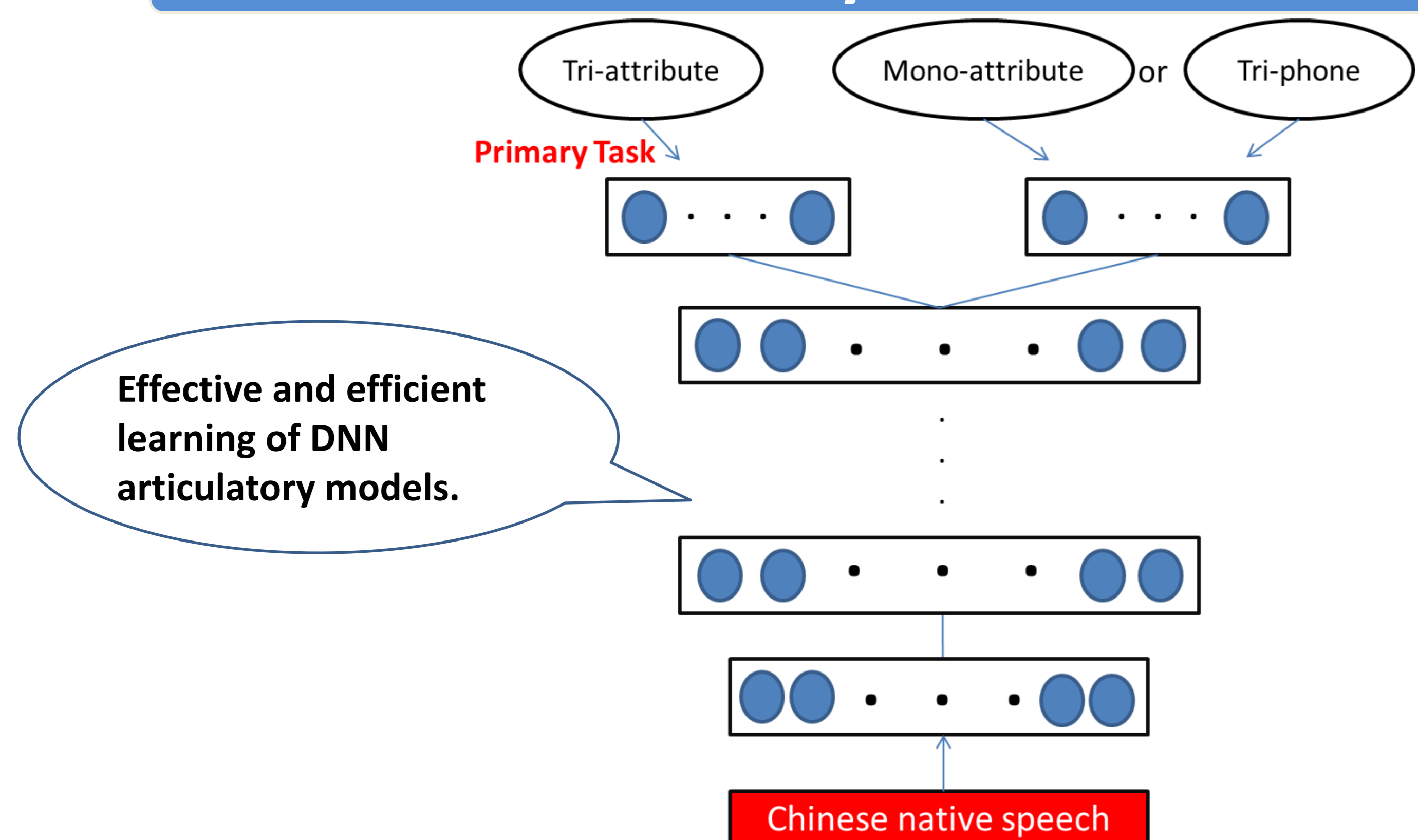
Articulatory attribute transcription

Sentence	你好 (HELLO)					
Phone	sil	n	i	h	ao	sil
Manner	sil	nasal	vowel	unvoiced-fricative	vowel	sil
Place & Backness	sil	alveolar	anterior	velar	central back	sil
Place & Height	sil	alveolar	high	velar	low middle	sil
Place & Roundedness	sil	alveolar	unroundedness	velar	unroundedness roundedness	sil

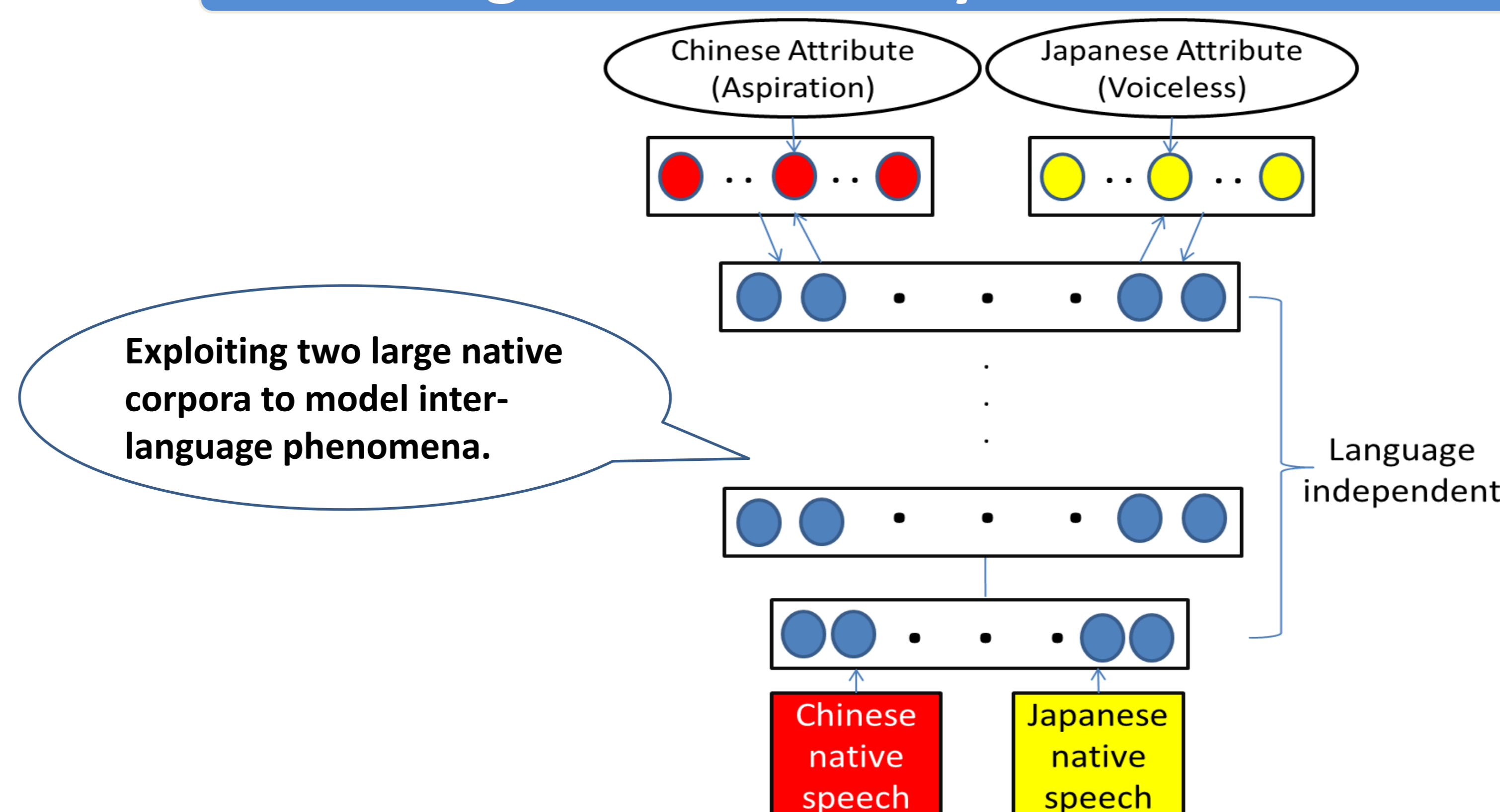
Modeling the golden attribute



## Multi-task articulatory attribute modeling

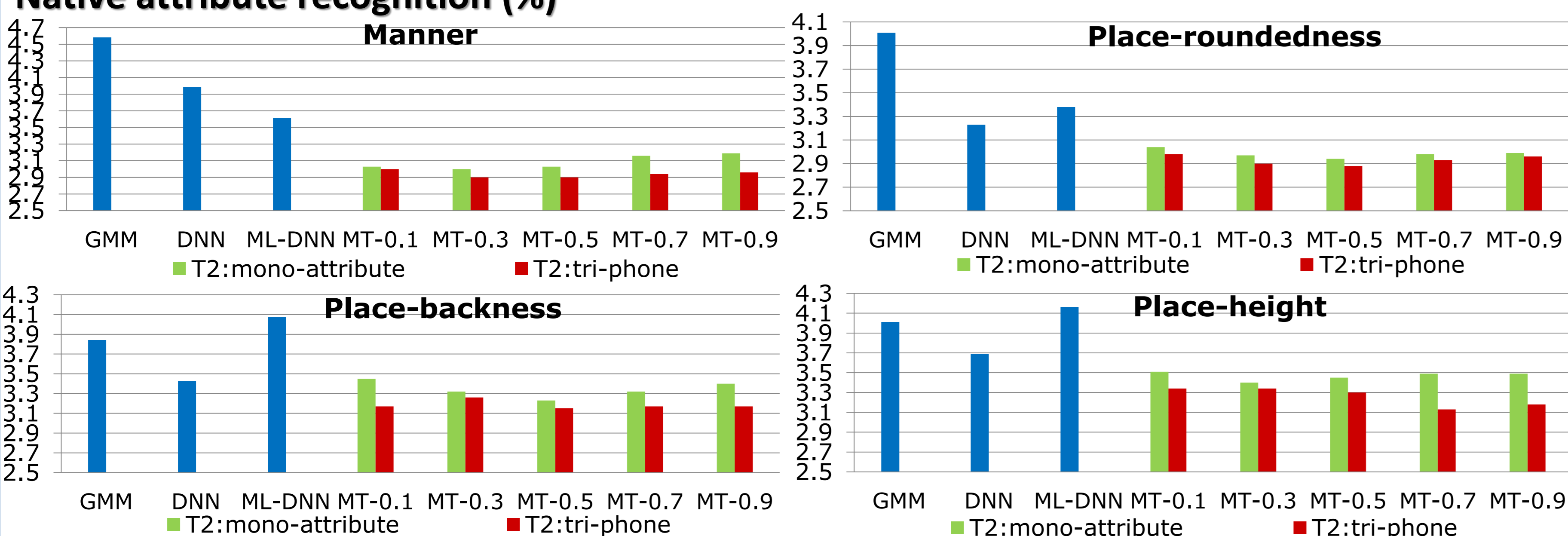


## Multi-lingual articulatory attribute modeling

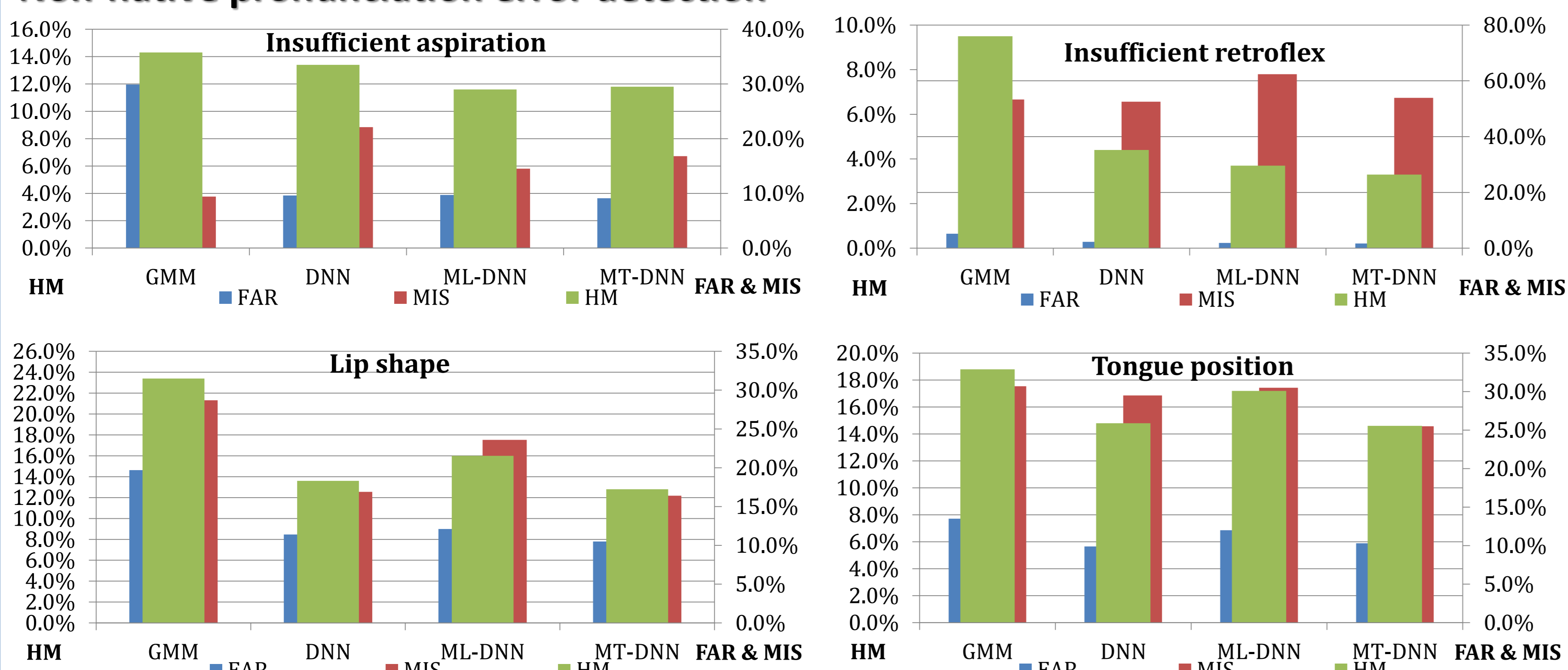


## Experiments

### Native attribute recognition (%)



### Non-native pronunciation error detection



### Experimental data

	Data set	Amount of data
Train	Chinese native training data	42h (28males, 36females)
	Japanese native training data	42h(80males, 73females)
Test	Chinese native testing data	5.3h (5males,3females)
	Chinese non-native testing data	1896 utterances (7 female Japanese students)

### Error type focused

- Insufficient aspiration:* Insufficient aspiration when producing aspirated constants (e.g. p)
- Insufficient retroflex:* Insufficient retroflex when producing retroflex constants (e.g. r)
- Lip roundedness:* Vowels with spread lips have problems of rounded sound (e.g. ü)
- Backness:* Inappropriate tongue position with a little back (e.g. an)

### Evaluation metrics

- *False Alarm Rate (FAR)*: rate of correct pronunciation that are detected as pronunciation errors by the system.
- *Miss Rate (MIS)*: rate of true pronunciation errors that are missed detected by system.
- *Harmonic Mean (HM)*: harmonic mean of FAR and MIS.

## Conclusions

We address effective articulatory models without non-native training data.

Multi-task learning method can enhance DNN articulatory modeling.

Multi-lingual learning method is effective for modeling non-native speech.