

# Speech Analysis and Linguistics

## Applications: Automatic Speech Recognition & Mispronunciation Detection and Diagnosis (MDD)

**Instructor: Chiranjeevi Yarra**

Speech lab, LTRC



Aug 11, 2022

# Outline

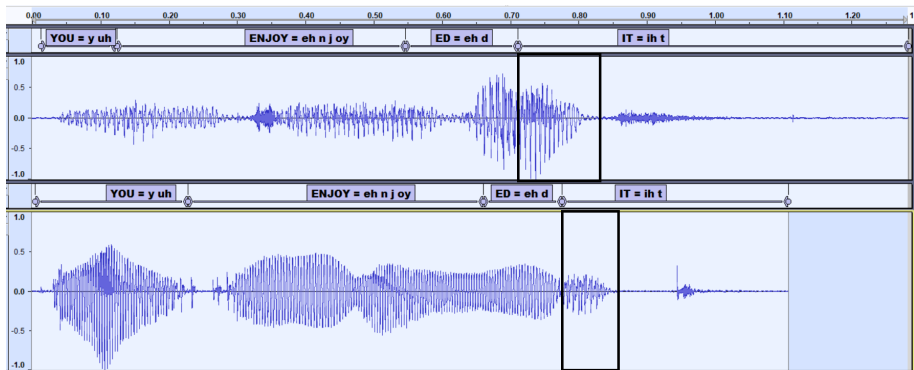
- 1 Phoneme Recognizer
- 2 Word Recognizer
- 3 Mispronunciation Detection and Diagnosis

## 1 Phoneme Recognizer

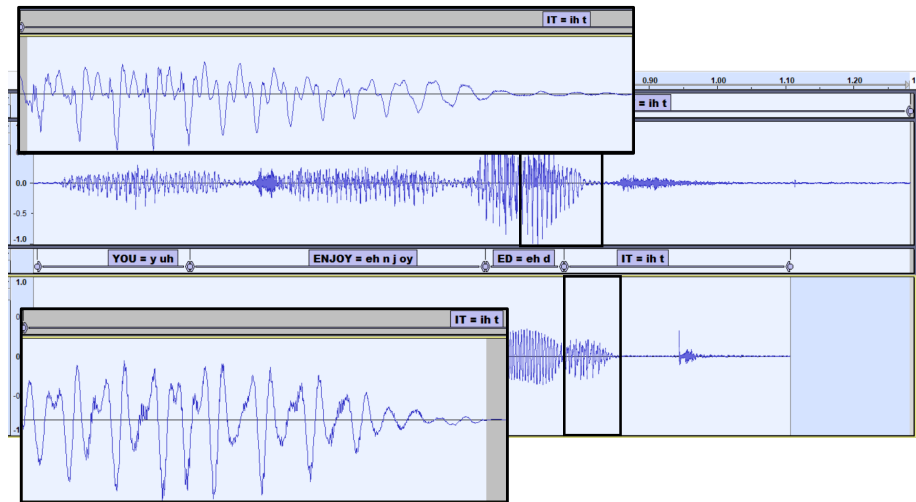
## 2 Word Recognizer

## 3 Mispronunciation Detection and Diagnosis

# What is to be model in ASR?

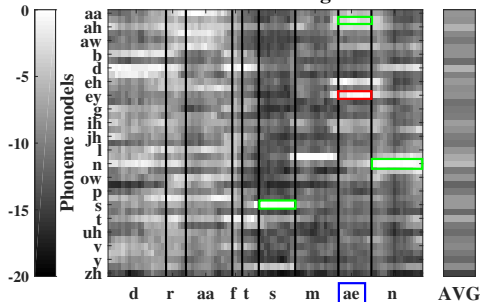


# What is to be model in ASR?

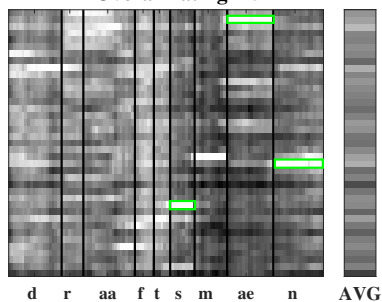


# Errors in the estimation

Overall rating = 1

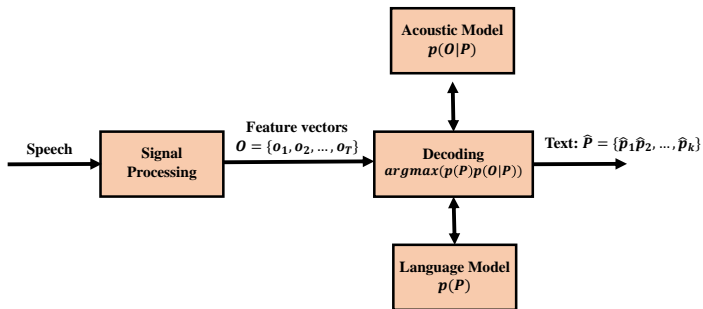


Overall rating = 5

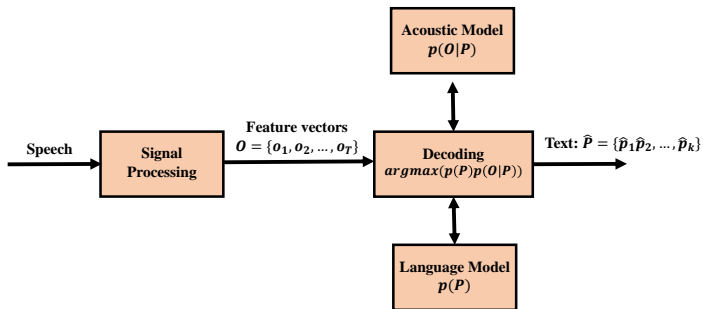


— Phoneme boundaries — Mismatched phoneme model — Matched phoneme model — Mispronounced phoneme

# Building blocks



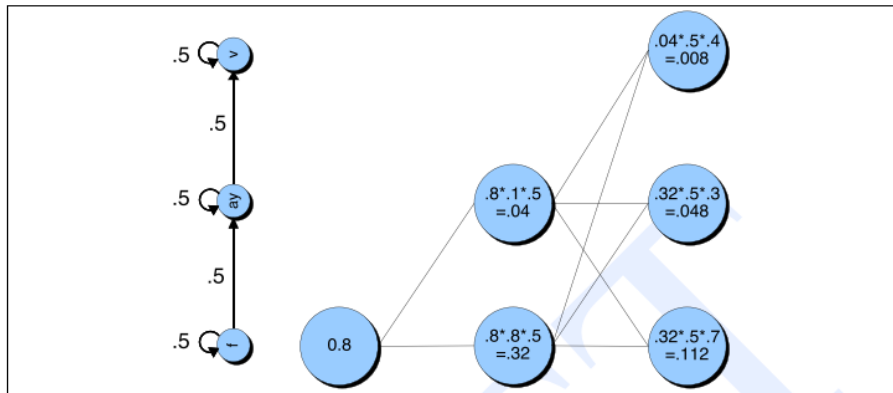
# Building blocks



1 How to model  $p(P)$  and  $p(O|P)$ ?



# Decoding



## Decoding illustration

<b>V</b>	0	0	0.008	0.0072	0.00672	0.00403	0.00188	0.00161	0.000667	0.000493
<b>AY</b>	0	0.04	0.048	0.0448	0.0269	0.0125	0.00538	0.00167	0.000428	8.78e-05
<b>F</b>	0.8	0.32	0.112	0.0224	0.00448	0.000896	0.000179	4.48e-05	1.12e-05	2.8e-06
<b>Time</b>	1	2	3	4	5	6	7	8	9	10
<b>B</b>	<i>f</i> 0.8	<i>f</i> 0.8	<i>f</i> 0.7	<i>f</i> 0.4	<i>f</i> 0.4	<i>f</i> 0.4	<i>f</i> 0.4	<i>f</i> 0.5	<i>f</i> 0.5	<i>f</i> 0.5
	<i>ay</i> 0.1	<i>ay</i> 0.1	<i>ay</i> 0.3	<i>ay</i> 0.8	<i>ay</i> 0.8	<i>ay</i> 0.8	<i>ay</i> 0.8	<i>ay</i> 0.6	<i>ay</i> 0.5	<i>ay</i> 0.4
	<i>v</i> 0.6	<i>v</i> 0.6	<i>v</i> 0.4	<i>v</i> 0.3	<i>v</i> 0.3	<i>v</i> 0.3	<i>v</i> 0.3	<i>v</i> 0.6	<i>v</i> 0.8	<i>v</i> 0.9
	<i>p</i> 0.4	<i>p</i> 0.4	<i>p</i> 0.2	<i>p</i> 0.1	<i>p</i> 0.1	<i>p</i> 0.1	<i>p</i> 0.1	<i>p</i> 0.1	<i>p</i> 0.3	<i>p</i> 0.3
	<i>iy</i> 0.1	<i>iy</i> 0.1	<i>iy</i> 0.3	<i>iy</i> 0.6	<i>iy</i> 0.6	<i>iy</i> 0.6	<i>iy</i> 0.6	<i>iy</i> 0.5	<i>iy</i> 0.5	<i>iy</i> 0.4

1 Phoneme Recognizer

2 Word Recognizer

3 Mispronunciation Detection and Diagnosis

# ASR equation

$$\hat{W} = \operatorname{argmax}\{p(W|O)\}$$

# ASR equation

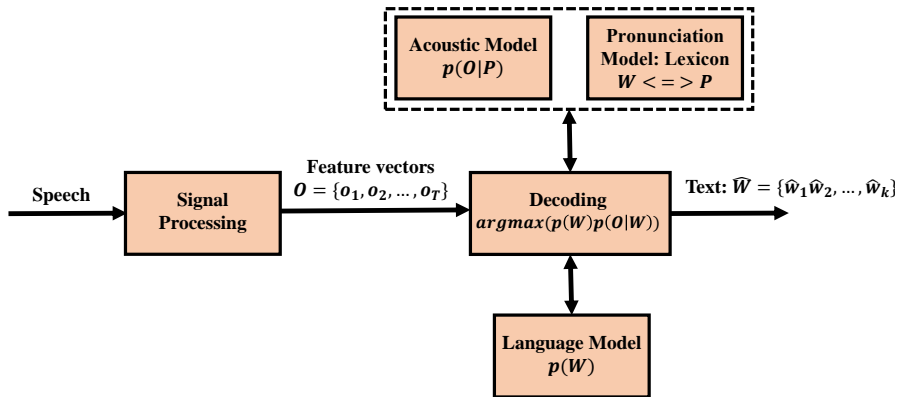
$$\begin{aligned}\hat{W} &= \operatorname{argmax}\{p(W|O)\} \\ &= \operatorname{argmax}\{p(W)p(O|W)\}\end{aligned}$$

# ASR equation

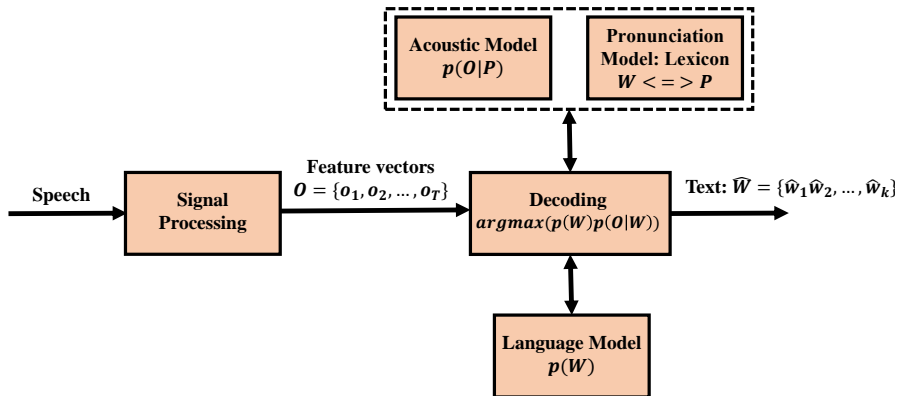
$$\begin{aligned}\hat{W} &= \operatorname{argmax}\{p(W|O)\} \\ &= \operatorname{argmax}\{p(W)p(O|W)\}\end{aligned}$$

**1** What data is needed for modelling  $p(W)$  and  $p(O|W)$ ?

## ASR building blocks



# ASR building blocks



1 Why modelling  $p(O|P)$  is efficient?



# Word and phoneme sequence mapping: Lexicon?

## An exemplary Lexicon

```
BRANER  B R EY1 N ER0
BRANFORD B R AE1 N F ER0 D
BRANHAM B R AE1 N HH AH0 M
BRANI   B R AE1 N IY0
BRANIFF B R AE1 N IH0 F
BRANIFF'S B R AE1 N IH0 F S
BRANIGAN B R AE1 N IH0 G AH0 N
BRANILLO B R AH0 N IH1 L OW0
BRANIN  B R AE1 N IH0 N
BRANISLOV B R AE1 N IH0 S L AA2 V
BRANITZKY B R AH0 N IH1 T S K IY1
BRANK   B R AE1 NG K
BRANK'S B R AE1 NG K S
BRANKI  B R AE1 NG K IY0
BRANKO  B R AE1 NG K OW0
BRANKS  B R AE1 NG K S
BRANN   B R AE1 N
BRANNA  B R AE1 N AH0
BRANNAM B R AE1 N AH0 M
```

# Word and phoneme sequence mapping: Lexicon?

## An exemplary Lexicon

```

BRANER  B R EY1 N ER0
BRANFORD B R AE1 N F ER0 D
BRANHAM B R AE1 N HH AH0 M
BRANI   B R AE1 N IY0
BRANIFF B R AE1 N IH0 F
BRANIFF'S B R AE1 N IH0 F S
BRANIGAN B R AE1 N IH0 G AH0 N
BRANILLO B R AH0 N IH1 L OW0
BRANIN  B R AE1 N IH0 N
BRANISLOV B R AE1 N IH0 S L AA2 V
BRANITZKY B R AH0 N IH1 T S K IY1
BRANK   B R AE1 NG K
BRANK'S B R AE1 NG K S
BRANKI  B R AE1 NG K IY0
BRANKO  B R AE1 NG K OW0
BRANKS  B R AE1 NG K S
BRANN   B R AE1 N
BRANNA  B R AE1 N AH0
BRANNAM B R AE1 N AH0 M

```

- 1 Knowing  $p(O|P)$  sufficient for computing  $p(O|W)$ ?

# Word and phoneme sequence mapping: Lexicon?

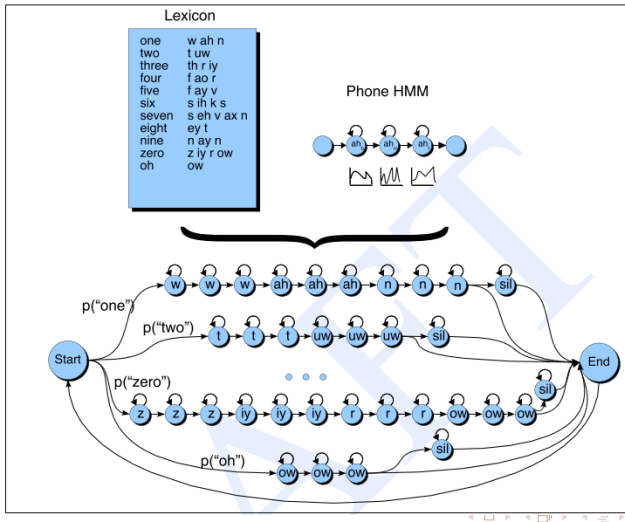
## An exemplary Lexicon

```

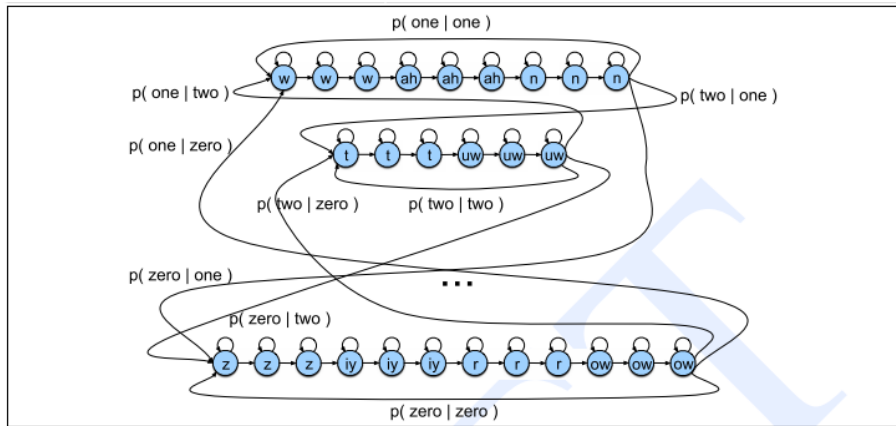
BRANER  B R EY1 N ER0
BRANFORD B R AE1 N F ER0 D
BRANHAM B R AE1 N HH AH0 M
BRANI   B R AE1 N IY0
BRANIFF B R AE1 N IH0 F
BRANIFF'S B R AE1 N IH0 F S
BRANIGAN B R AE1 N IH0 G AH0 N
BRANILLO B R AH0 N IH1 L OW0
BRANIN  B R AE1 N IH0 N
BRANISLOV B R AE1 N IH0 S L AA2 V
BRANITZKY B R AH0 N IH1 T S K IY1
BRANK   B R AE1 NG K
BRANK'S B R AE1 NG K S
BRANKI  B R AE1 NG K IY0
BRANKO  B R AE1 NG K OW0
BRANKS  B R AE1 NG K S
BRANN   B R AE1 N
BRANNA  B R AE1 N AH0
BRANNAM B R AE1 N AH0 M
  
```

- 1 Knowing  $p(O|P)$  sufficient for computing  $p(O|W)$ ?
- 2 Knowing  $W$  sufficient for obtaining  $P$ ?

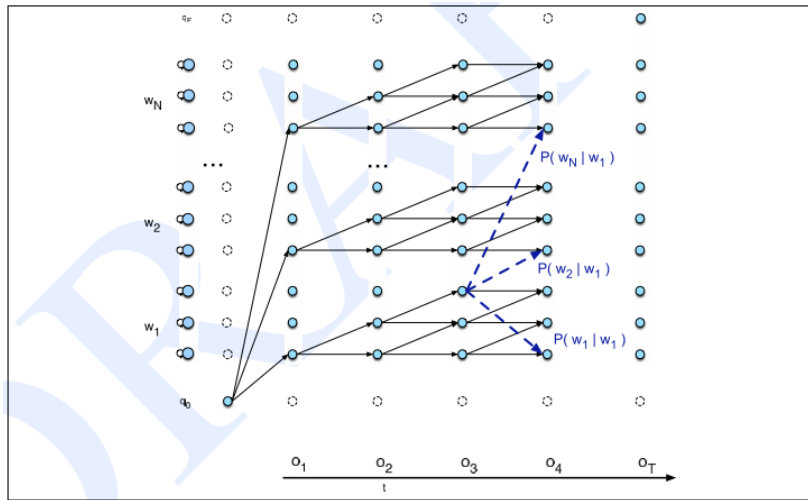
## An example



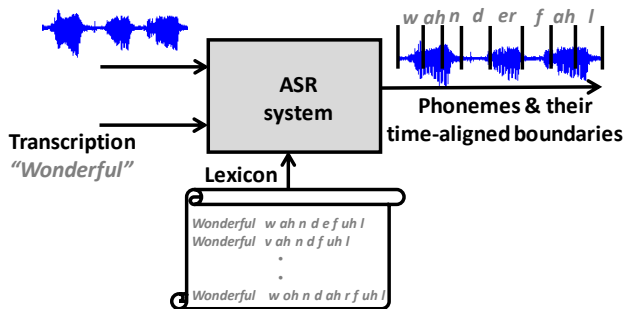
# Language model for Digit Recognition



## Viterbi decoding



# Forced-alignment process



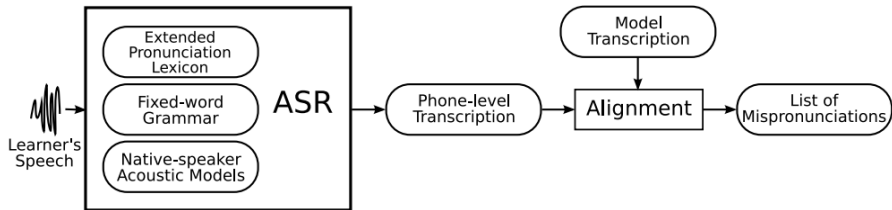
1 Phoneme Recognizer

2 Word Recognizer

3 Mispronunciation Detection and Diagnosis



# Mispronunciation Detection and Diagnosis (MDD)



Thank you