

Task Introduction

Task: Multiclass Classification

Framewise phoneme prediction from speech.



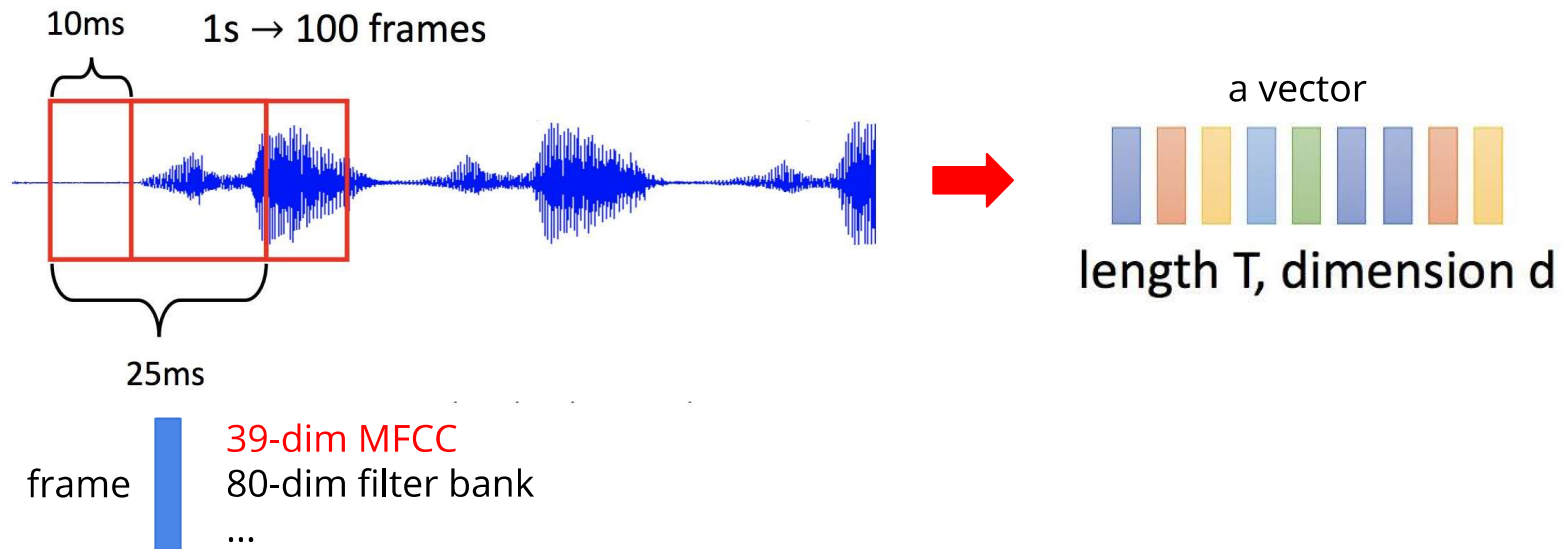
What is a phoneme?

A unit of speech sound in a language that can serve to distinguish one word from the other.

- bat / pat , bad / bed
- Machine Learning → M AH SH IH N L ER N IH NG

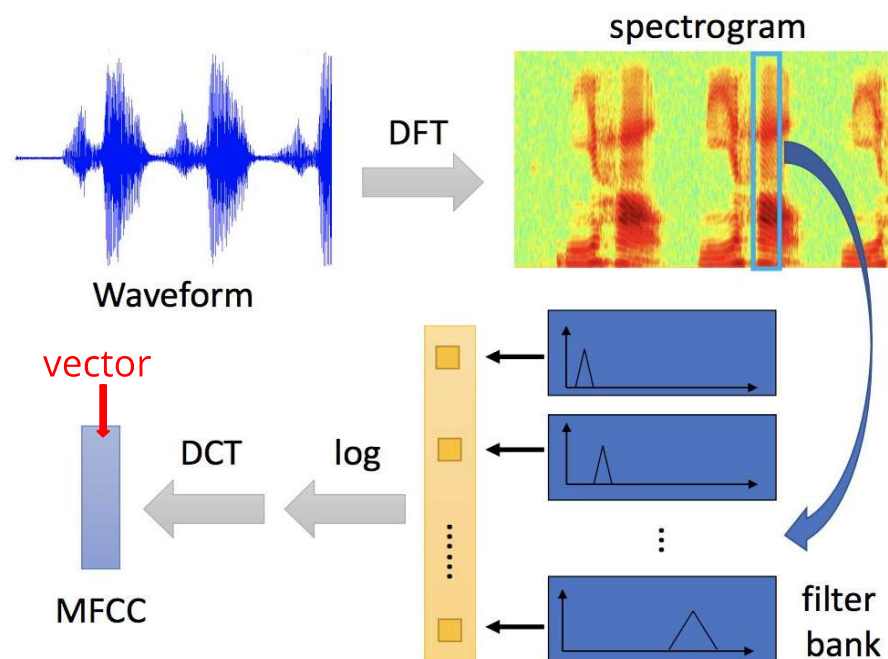
Task Introduction

Data Preprocessing



Task Introduction

Acoustic Features - MFCCs (Mel Frequency Cepstral Coefficients)



For more details,
please refer to Prof. Lin-Shan Lee's
[\[Introduction to Digital Speech Processing\]
Chap.7](#)

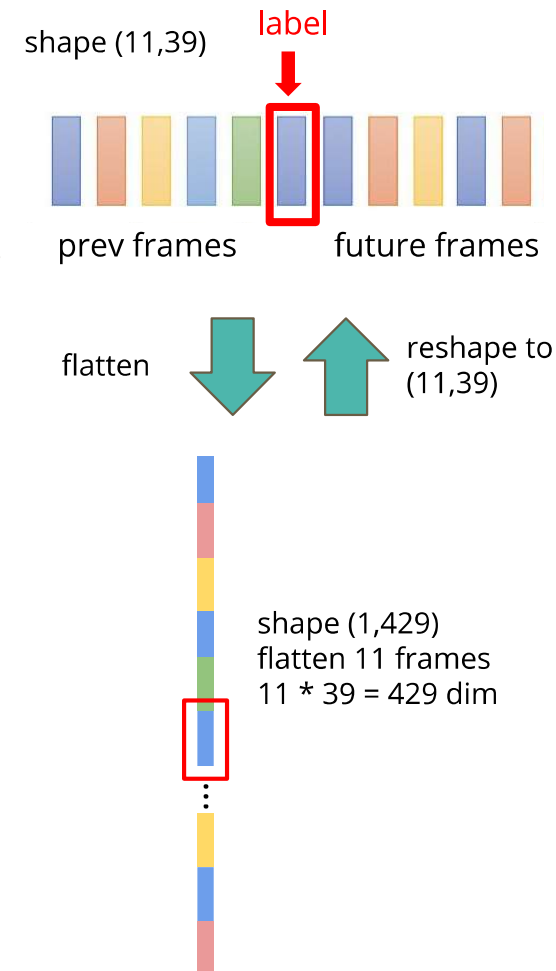
Image ref.
Prof. Hung-Yi Lee
[\[2020Spring DLHLP\] Speech Recognition](#)

More Information About the Data

Since each frame only contains 25 ms of speech, a single frame is unlikely to represent a complete phoneme

- Usually, a phoneme will span several frames
 - Hint: post-processing may help
- Concatenate the neighboring phonemes for training
 - In this HW, we concatenate the past and the future five frames for training (total 11 frames)
 - You may reshape the input (1,429) back to (11,39) to get separated 11 frames
 - Just remember that the label corresponds to the center frame
- Finding testing labels or doing human labeling are strictly prohibited!

[Introduction to Digital Speech Processing](#)



Dataset & Data Format

- Dataset: TIMIT Acoustic-Phonetic Continuous Speech Corpus
 - Phonetically balanced for English
 - Data Format (The TAs have already preprocessed the data)
timit_11/
 - **train_11.npy** → **training data (# of training frames, 11 x feature dim)**
 - **train_label_11.npy** → **framewise phoneme label (0-38)**
 - **test_11.npy** → **testing data (# of testing frames, 11 x feature dim)**
 - Acoustic features (39-dim MFCC)
 - Concatenate the past and the future five frames (feature dim = 11 x 39)
 - The phoneme label of each input corresponds to the center frame
 - **Using additional data is prohibited.** Your final grade will be multiplied by 0.9!
-

Class	Phoneme	Example	Class	Phoneme	Example	Class	Phoneme	Example
0	iy	<i>beet</i>	13	l	<i>lay</i>	26	dx	<i>mudd^y</i>
1	ih	<i>bit</i>	14	r	<i>ray</i>	27	g	<i>gay</i>
2	eh	<i>bet</i>	15	y	<i>yacht</i>	28	p	<i>pea</i>
3	ae	<i>bat</i>	16	w	<i>way</i>	29	t	<i>tea</i>
4	ah	<i>but</i>	17	er	<i>bird</i>	30	k	<i>key</i>
5	uw	<i>boot</i>	18	m	<i>mom</i>	31	z	<i>zone</i>
6	uh	<i>book</i>	19	n	<i>noon</i>	32	v	<i>van</i>
7	aa	<i>bob</i>	20	ng	<i>sing</i>	33	f	<i>fin</i>
8	ey	<i>bait</i>	21	ch	<i>choke</i>	34	th	<i>thin</i>
9	ay	<i>bite</i>	22	jh	<i>joke</i>	35	s	<i>sea</i>
10	oy	<i>boy</i>	23	dh	<i>then</i>	36	sh	<i>she</i>
11	aw	<i>bout</i>	24	b	<i>bee</i>	37	hh	<i>hay</i>
12	ow	<i>boat</i>	25	d	<i>day</i>	38	sil	silence/closure sounds