

Jan. 16th.

CSC2518, LEC02

## Phonemes.

phone  $\rightarrow$  triphone  $\rightarrow$  State  
pentaphone

① replaced by DNNs.

$P(x | "t-d+uw") \rightarrow \sim$  Mixture of Gaussians (GMM)  
 $P(x | "d-uw-y") \rightarrow$  pattern classification / recognition.  
 $\rightarrow$  Tied-state  $\leftarrow$  eg. (frame classification)

## Lexical Model.

sequence acoustic model  $\rightarrow$  [ ]  $\rightarrow$  language model (HMM)

② replaced with end-to-end model

Decision-Tree-Based state clustering.  $\rightarrow$  large space.

$\rightarrow$  Fully connected Layers  $\rightarrow 2 \times 5, 1D, 11$  frames.

$\rightarrow$  RNN (LSTM)  $\rightarrow$  # of frame (context) can be dynamically adjusted.

$\rightarrow$  Connectionist Temporal Classification (CTC)  
Bi-directional LSTM is essential for CTC.

## Attention

## DNN

Speaker Adaptation eg. VTLN

$\hookrightarrow$  { Conservative Training.  
Transformation methods.

Speaker-aware Training (noise/device-aware)  
only need acoustic info.