

CSC 2518

Topics:

- Robustness / Learned Representations.
- Reductive Architecture / Under-resourced Language.
- Intelligibility / Enhancement / Adaptation.
- RL
- Time-domain / Phase Representation.

Automatic Speech Recognition.

Dynamic Time Warping (DTW)

Hidden Markov Model (HMM)

Large Vocabulary continuous dictation., Neural acoustic models outperformed by Gaussian mixtures.

~~Discriminative~~ Training, Weighted finite-state transducers (WFSTs)

Neural architectures

Sarkis et al. 2011: 30% RER on Switchboard.

Aim of ASR.

$$\hat{W} = \underset{W}{\operatorname{argmax}} P(W|A) = \underset{W}{\operatorname{argmax}} P(A|W) P(W).$$

W : most likely sentence.

A : speech audio.

Acoustic model: $P(A|W)$

Language model: $P(W)$

Evaluation : $F(x, w) = P(w | x)$.

Inference

$$w = \arg$$

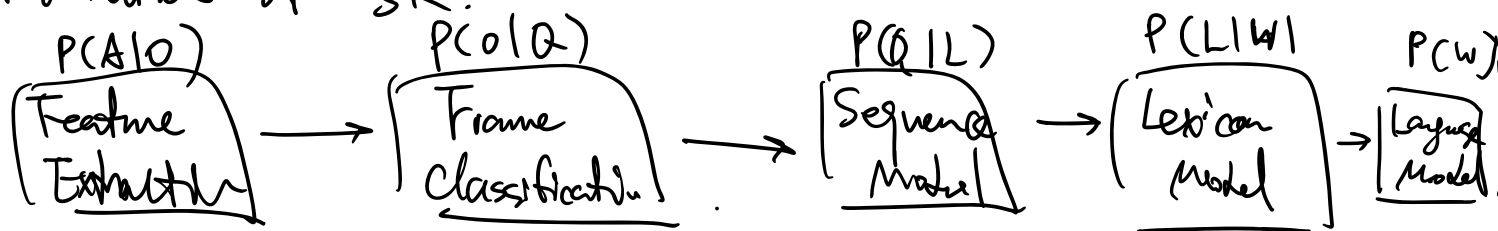
Language Modeling .

- Guide the search algorithm.
- Disambiguate between phrases which are acoustically similar

N-gram.

RNN.

Architecture of SR.



A: Audio

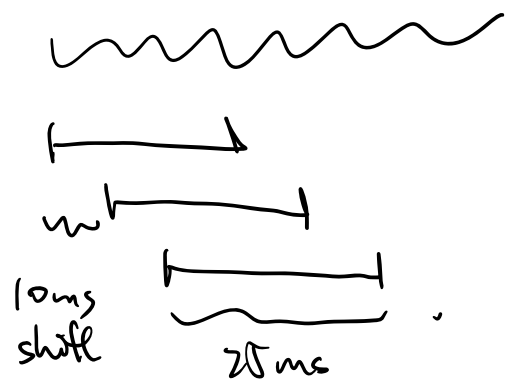
O: Feature Frames

Q: Sequence States

L: ~~Q~~ Phrases

W: Words.

Feature Extraction.



Phoneme
basic unit.