
Speech Generation & Recognition

Reid Simmons
Illah Nourbakhsh

Speech Generation

Desirable Speech Characteristics

- Naturalness
 - Sounds human-like
- Intelligibility
 - Easily understandable



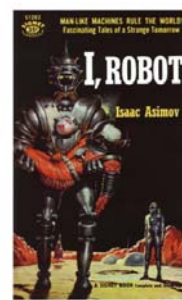
HRI: Speech

3

Simmons, Nourbakhsh : Spring 2015

Speech Synthesis (Text-to-Speech)

- Why is this difficult, technically?



HRI: Speech

4

Simmons, Nourbakhsh : Spring 2015

Speech Synthesis

- Normalization
 - Pre-process text to contain only words
- Text Analysis
 - Syntactic parsing
 - Semantic parsing
- Text-to-Phoneme
 - Pronunciation
- Prosody
 - Pitch, loudness, duration => stress
- **Affect**

HRI: Speech

5

Simmons, Nourbakhsh : Spring 2015

Issues with Normalization

- Numbers
- Abbreviations
- Acronyms

HRI: Speech

6

Simmons, Nourbakhsh : Spring 2015

Issues with Pronunciation

- Ambiguous words
 - Often depends on part-of-speech
 - May depend on semantics
 - May depend on tense!
- Many exceptions to “sounding out” rules
 - though, through, bough, cough, tough
 - comb, tomb, bomb
 - dose, hose, lose

HRI: Speech

7

Simmons, Nourbakhsh : Spring 2015

Issues with Prosody

- Punctuation
 - Pause, after comma
 - Rising tone for questions?
- Syllabic stress
- Word stress

HRI: Speech

8

Simmons, Nourbakhsh : Spring 2015

Issues with Affect

- Speed of speech
- Emotional content of speech

Exercise

- Get into your teams
- Create a four word sentence
 - noun-verb-noun-adverb or noun-verb-adjective-noun
 - Read and record using four different emotions (happy, angry, sad, disgusted, fearful, surprised)
 - listen and analyze how prosody changes
- Take notes: We will then discuss

Techniques for Speech Generation

- **Formant** (rule-based)
 - Use acoustic models
 - Compact program
 - Tends to be quite intelligible, but limited prosody
- **Concatenation** (unit selection)
 - Use human speech, “sliced and diced”
 - phones, diphones, triphones, ...
 - Layer on prosody using signal processing
 - Domain-specific synthesis

HRI: Speech

11

Simmons, Nourbakhsh : Spring 2015

SSML

- Semi-standard markup language for specifying pronunciation and prosody
 - <emphasis level=“strong”>
 - <break time=“4500ms”>
 - <prosody rate=“fast”>
 - <prosody pitch=“+25Hz”>
 - <prosody volume=“33%”>
 - _{Dr.}
 - <phoneme ph=“t ah0 m ey1 t ow0”>tomato</phoneme>
 - <say-as interpret-as=“digits”>123</say-as>
 - <say-as interpret-as=“number:ordinal”>VIII</say-as>

HRI: Speech

12

Simmons, Nourbakhsh : Spring 2015

Speech Recognition

Speech Recognition

- Why is this difficult, technically?

Diverse Sources of Ambiguity

- Acoustic/Phonetic
 - Let us pray
 - Lettuce spray
- Syntactic
 - Meet her at the end of Main Street
 - Meter at the end of Main Street
- Semantic
 - Is the baby crying
 - Is the bay bee crying
- Discourse Context
 - It is hard to recognize speech
 - It is hard to wreck a nice beach

HRI: Speech

15

Simmons, Nourbakhsh : Spring 2015

Phonemes

- ~40-45 phonemes in English
 - Variance depends mostly on dialect
- Voiced vs. unvoiced
 - vowels vs. consonants
- Phonetic sounds may differ based on preceding and succeeding phonemes

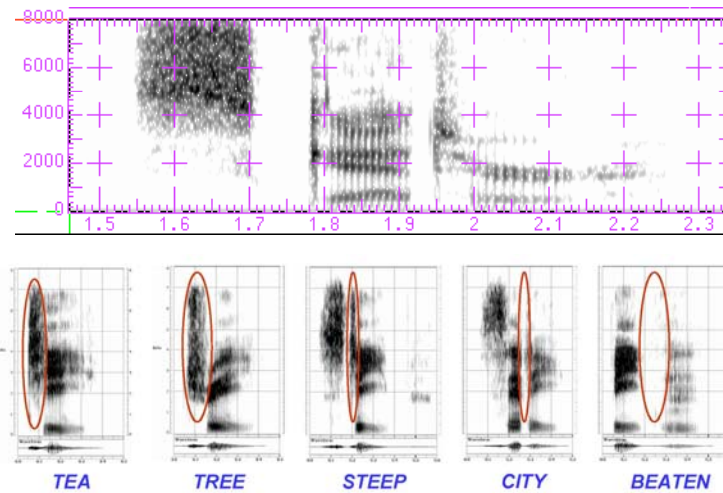
HRI: Speech

16

Simmons, Nourbakhsh : Spring 2015

The Acoustic Signal

Sad:



HRI: Speech

17

Simmons, Nourbakhsh : Spring 2015

Techniques for Speech Recognition

- Almost all current approaches use statistical modeling and massive amounts of data
- Maximize probability of word sequence
 - $P(W^* | A) \approx P(A | W^*)P(W^*)$
- Typically, language model uses **trigrams**
 - Probable sequence of phonemes
 - $P(W_n | W_{n-1}, W_{n-2})$
- Other constraints
 - Syntax
 - Semantics
 - Domain / context

HRI: Speech

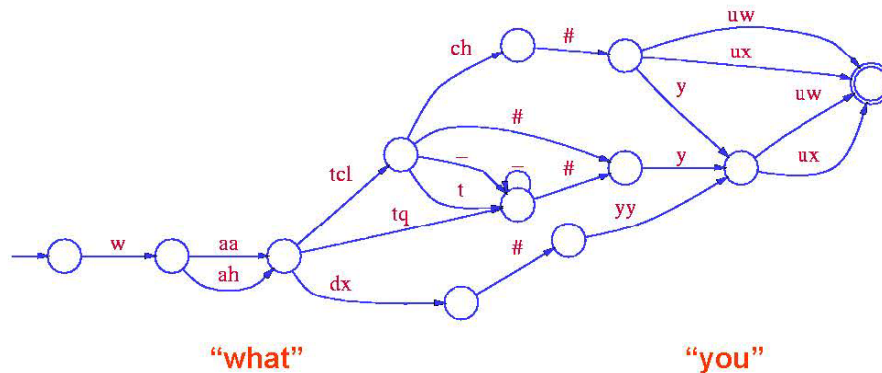
18

Simmons, Nourbakhsh : Spring 2015

Representing Language Model

- Hierarchical **Hidden Markov Model** (HMM)
 - Atomic units (sub-phoneme)
 - ~20ms slices, characterized by power in bands of frequencies
 - HMM of phonemes
 - HMM of diphones or triphones
 - HMM of words
 - HMM of phrases (trigrams)
 - Put together into single HMM
- Use Viterbi algorithm to find best path
 - May use backwards search to refine path

Simple Word-Level HMM



Grammatical Issues

- Incomplete Sentences
- Non-Grammatical Sentences
- Fillers
 - er, um, ...
- Disfluencies
 - cutting off mid-word
 - corrections
 - hesitations

Further Issues

- Recognizing Prosody
 - Stress is important in interpreting pragmatics
- Recognizing Emotion/Affect
- Current Status
 - Siri and Google Voice have been able to use millions of training examples to create fairly good continuous, speaker-independent speech recognition