

Articulatory Model & Synthesis

Vocal Tract Analysis with Linear Prediction





Articulatory Model Parameters

- lip opening W
- lip protrusion length L
- tongue body height and length – Y, X
- velar closure K
- tongue tip height and length
 B, R
- jaw raising (dependent parameter)
- velum opening N

Georgialnstitute of Technology

1.5 15.0Cross-sectional area function [cm²] AVELUM [cm²] FRACT 0.0 0.0 GLOTTIS i→ LIPS 40 [dB] HTRACT Speech Model Spectrum -40 2 n f [kHz]→ Georgialnstitute Tech

Continuing Evolution (1959-1987)

- Haskins, 1959 🏾 🍕
- KTH Stockholm, 1962 🍕
- Bell Labs, 1973 🏻 🍕
- MIT, 1976 🍕
- MIT-talk, 1979 🛛 🐠
- Speak 'N spell, 1980 🛛 🐠
- Bell Labs, 1985 🛛 🍕
- Dec talk, 1987 🍕 🍕 🍕

Paradigm Shift in Speech Synthesis

- Parametric TTS has become highly intelligible but still extremely unnatural;
- <u>Parameterization amounts to data</u> reduction, causing loss of subtle cues on naturalness;
- Rapid growth in computation and memory diminishes the need in parameterization; an 80G HD can store 11,100 hours of good quality speech!



2.5" hard disc, 64.8 GB/in², 0.95cm (h) 80G < \$200

• New concatenative systems use segments of **recorded speech** to ensure better quality sound than any synthesized one with remaining challenge in **control**.



GeorgiaInstitute of Technology



GeorgiaInstitute of Technology



An early 20th century transcribing pool at Sears, Roebuck and Co. courtesy www.recording-history.org



Regularity and Variability in Human Speech

- Regularity (Structure in Speech)
 - Linguistic hierarchy exists for any language: phonemes, lexicon, regular expressions, ...
 - Speakers are taught to speak according to the convention of a language
- Variability (Randomness in Speech)
 - Casual use of words, improper pronunciation & word grouping, completeness in sentential structure or usual disfluency (um, repair) – to err is human
 - Unregulated speaking pace & breath grouping relaxed pronunciation – to rest is human

Speech Recognition Techniques



Hidden Markov Models



$$P(X | W, \lambda_x) = \sum_{\mathbf{q}} P(X, \mathbf{q} | W, \lambda_x)$$

$$P(X, \mathbf{q} | W, \lambda_x) = a_0 \prod_{t=1}^T a_{q_{t-1}q_t} b_{q_t}(x_t)$$

$$X = (x_1, x_2, \dots, x_T) \qquad \mathbf{q} = (q_0, q_1, q_2, \dots, q_T)$$



- Each state represents a process of measurable observations.
- Inter-process transition is governed by a finite state Markov chain.
- Processes are stochastic and individual observations do not immediately identify the state.

Embedding Acoustics in Finite-State Graphs



Composite Finite State Phone Network

GeorgiaInstitute

Language Models

- $P(X, W \mid \Lambda) = P(W \mid \lambda_w) P(X \mid W, \lambda_x)$
- N-gram: $P(W | \lambda_w) = \prod_{i=1}^{L} P(w_i | w_{i-1}, \dots, w_{i-N+1})$
- General finite state: $P(W | \lambda_w) = \prod_{i=1}^{L} P(w_i | \lambda_w, H_i)$

where H_i is the state of history for word w_i

Domain dependency:



DARPA Speech-to-Text Benchmark Tests



Statistical Pattern Recognition?

Bayes decision needs knowledge of *a posteriori* probabilities → conventional problem of distribution estimation.



Paradigm shift – look at the "error function" as an objective for optimization.

• Minimum error discriminative training leads to much improved recognition accuracy.

Georgialnstitute

One Success Story – Call Automation

- Telephone companies used Operator Support Position System to handle operator calls, a precursor to call center
- Automatic speech recognition w/ hidden Markov models reached maturity by late 80s for some telecom applications
- Voice Recognition Call Processing (VRCP) by spotting five essential keywords handles billions of 0+ calls annually.



GeorgiaInstitute of Technology

Current/Overall Technology Assessment



Human Still Outperforms Machines



Moving Forward

- Learn from human why human does better?
 - Biological system models
 - Detection-based paradigm
- Address naturalness in human-machine interactions
 - Speech act and hierarchy; semantics
 - Dialog management
- Broaden "circle of impact" or usefulness
 - Ability to handle any kind of speech

Learning from Human – The Ear?





Learning from Human – The Cortex?



"Cortical" Regions of Interest (ROI)

Towards Understanding & Conversation



Speech Everywhere

- Speech is being uttered everywhere, most likely mixed with noise, "competing voice" & distortion
 - Speech enhancement, noise suppression, dereverberation?
- Many sounds exist at any given time how do we make sense from them?
 - Source separation? Information fusion?
- How does the machine know humans are trying to communicate with it?
 - Modeling of speech act?
- "Living room" at NTT Communication Science Research

A Tour for <u>Demo</u> on Use of Reference & Context

Login

Self Awareness I What are you doing? Calendar Arithmetic What's the date? ... tomorrow's? How many days next month? What day is the 1st of next month? Self Awareness II What have you done this session? What dates have we discussed? Simple Sequence Create <tomorrow eve; 9:00; 30min> <Repeat> What are you doing? Stop doing that. Create a 2-hr appt tomorrow @ noon

Self Awareness III What did you just do? What day are we talking about? Selective Access What time does that appt end? Selective Modification Extend the duration by one hour. Change it back to what it was. Resch'ed it for 1st Sat in April. Resch'ed it for next Monday. Conclusion Create a 2-hour appt at 2:00 on the last Tuesday of the month after next. Goodbye, Daisy.

> Georgialnstitute of Technology

Summary

- Speech technologies are of human communication and for human communication; naturalness is important.
- Speech research has expanded beyond simple conversion between sound and text, to wit,

From Machines That Convert to Machines That Converse

• Speech is an essential transport of intelligence and knowledge; an intelligent machine needs to know how to talk and listen.

