

A Brief Introduction to Automatic Speech Recognition

Jim Glass (glass@mit.edu)

MIT Computer Science and Artificial Intelligence Laboratory

March 20, 2007

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 1

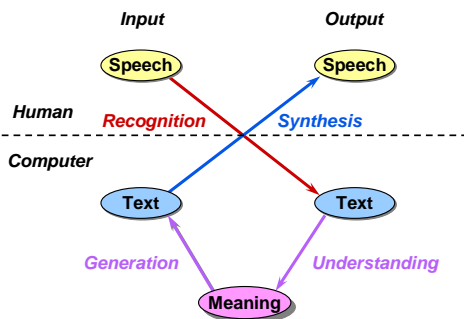
Overview

- Introduction
- Speech
- Models
- Search
- Representations

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 2

Communication via Spoken Language



Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 3

Virtues of Spoken Language

- | | |
|--------------------|---|
| Natural: | Requires no special training |
| Flexible: | Leaves hands and eyes free |
| Efficient: | Has high data rate |
| Economical: | Can be transmitted/received inexpensively |

Speech interfaces are ideal for information access and management when:

- The information space is broad and complex,
- The users are technically naive, or
- Only telephones are available.

video

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 4

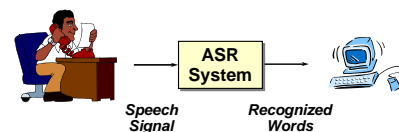
Diverse Sources of Knowledge for Spoken Language Communication

- | | |
|---------------------------|--|
| 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 easy to recognize speech
It is easy to wreck a nice beach |
| Others: | I'm <i>flying</i> to Chicago tomorrow
I'm flying to <i>Chicago</i> tomorrow |

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 5

Automatic Speech Recognition



- An ASR system converts the speech signal into words
- The recognized words can be
 - The final output, or
 - The input to natural language processing

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 6

MIT Application Areas for Speech Interfaces

- **Mostly input (recognition only)**
 - Simple command and control
 - Simple data entry (over the phone)
 - Dictation
- **Interactive conversation (understanding needed)**
 - Information kiosks
 - Transactional processing
 - Intelligent agents

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 7

MIT Parameters that Characterize the Capabilities of ASR Systems

Parameters	Range
Speaking Mode:	Isolated word to continuous speech
Speaking Style:	Read speech to spontaneous speech
Enrollment:	Speaker-dependent to speaker-independent
Vocabulary:	Small (<20 words) to large (>50,000 words)
Language Model:	Finite-state to context-sensitive
Perplexity:	Low (<10) to high (>200)
SNR:	High (>30dB) to low (<10dB)
Transducer:	Noise-canceling microphone to cell phone

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 8

MIT Read versus Spontaneous Speech

Filled and unfilled pauses: read, spontaneous
Lengthened words: read, spontaneous
False starts: read, spontaneous

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 9

MIT Speech Recognition: Where Are We Now?

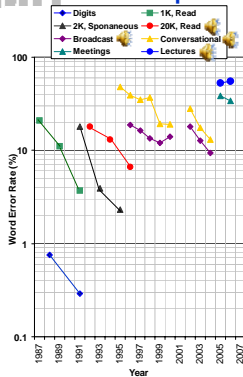
- **High performance, speaker-independent speech recognition is now possible**
 - Large vocabulary (for cooperative speakers in benign environments)
 - Moderate vocabulary (for spontaneous speech over the phone)
- **Commercial recognition systems are now available**
 - Dictation (e.g., IBM, Microsoft, Nuance, etc.)
 - Telephone transactions (e.g., AT&T, Nuance, VST, etc.)
- **When well-matched to applications, technology is able to help perform real work**
- **Demos:**
 - Speaker-independent, medium-vocabulary, small footprint ASR
 - Dynamic vocabulary speech recognition with constrained grammar (<http://web.sls.csail.mit.edu/city>)
 - Academic spoken lecture transcription and retrieval (<http://web.sls.csail.mit.edu/lectures>)

video
video

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 10

MIT Examples of ASR Performance



- Telephone digit recognition has word error rates of 0.3%
- Error rate for spontaneous speech twice that of read speech
- Error rate cut in half every two years for moderate vocabularies
- Corpora range in size from tens to thousands of hours
- Conversational speech from many speakers with noise remains a research challenge
 - Current focus on meetings & lectures

Intelligent Multimodal Interfaces (6.870)

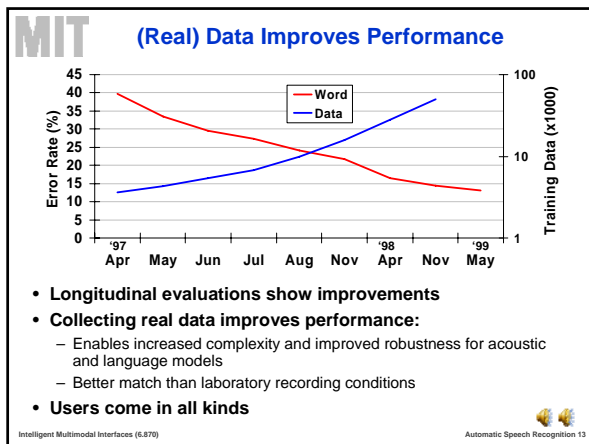
Automatic Speech Recognition 11

MIT The Importance of Data

- **We need data for analysis, modeling, training, and evaluation**
 - "There is no data like more data"
- **However, we need to have the right kind of data**
 - From real users
 - Solving real problems
- **Conduct research within the context of real application domains**
 - Forces us to confront critical technical issues (e.g., rejection, new word problem)
 - Provides a rich and continuing source of useful data
 - Demonstrates the usefulness of the technology
 - Facilitates technology transfer

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 12



MIT Real Data will Dictate Technology Needs

TECHNOLOGY REQUIRED	EXAMPLE
Simple word spotting	Um, Braintree
Complex word spotting	Eh yes, Avis rent-a-car in Boston
	Hello, please Brighton, uh, can I have the number of Earthscape, in, uh, on Nonantum Street
Speech understanding	Woburn, uh, Somerville. I'm sorry

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 14

MIT Important Lessons Learned

- Statistical modeling and data-driven approaches have proved to be powerful
- Research infrastructure is crucial:
 - Large amounts of linguistic data
 - Evaluation methodologies
- Availability and affordability of computing power lead to shorter technology development cycles and real-time systems
- Performance-driven paradigm accelerates technology development
- Interdisciplinary collaboration produces enhanced capabilities (e.g., spoken language understanding)

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 15

MIT ASR Trends*: Then and Now

	before mid 70's	mid 70's - mid 80's	after mid 80's
Recognition Units:	whole-word and sub-word units	sub-word units	sub-word units
Modeling Approaches:	heuristic and ad hoc	template matching	mathematical and formal
	rule-based and declarative	deterministic and data-driven	probabilistic and data-driven
Knowledge Representation:	heterogeneous and complex	homogeneous and simple	homogeneous and simple
Knowledge Acquisition:	intense knowledge engineering	embedded in simple structure	automatic learning

* There are, of course, many exceptions.

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 16

MIT But We Are Far from Done!

Corpus	Speech Type	Lexicon Size	Word Error Rate (%)	Human Error Rate (%) *
Digit Strings (phone)	spontaneous	10	0.3	0.009
Resource Management	read	1000	3.6	0.1
ATIS	spontaneous	2000	2	--
Wall Street Journal	read	~20K	6.6	1
Broadcast News	mixed	~64K	9.4	--
Switchboard (phone)	conversation	~25K	13.1	4
Meetings	conversation	~25K	30	--

* Lippmann, 1997

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 17

MIT What Makes Speech Recognition Hard?

- Phonological variations**
 - Local and global contexts, ...
- Individual differences**
 - Anatomy, socio-linguistic factors, ...
- Environmental factors**
 - Transducers, noise, ...
- Diversity of language use**
 - Syntax, semantics, discourse, ...
- Real-world issues**
 - Disfluencies, new words, ...
- ...

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 18

MIT ASR Is All About Utilizing Constraints

- **Acoustic**
 - Speech signal is generated by the human vocal apparatus
- **Phonetic**
 - /s/ in word initial /st/ cluster is unaspirated (e.g. “stay”)
- **Phonological**
 - /s/-/S/ sequence can turn into a long /S/ (e.g., “gas shortage”)
- **Lexical**
 - Words in a language are limited (e.g., “blit” and “vnuk” are not English words)
- **Language**
 - Probability of a word depends on its predecessors (e.g., “you” is the most likely word to follow “thank”)
 - A sentence must be syntactically and semantically well formed (e.g., subject-verb agreement)
- ...

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 19

MIT Major Components in a Speech Recognizer

```

graph TD
    TD[Training Data] --> AM[Acoustic Models]
    TD --> LM[Lexical Models]
    TD --> LangM[Language Models]
    AM --> S[Search]
    LM --> S
    LangM --> S
    C[Applying Constraints] --> S
    SS[Speech Signal] --> R[Representation]
    R --> S
    S --> RW[Recognized Words]
  
```

- **Speech recognition is the problem of deciding on**
 - How to *represent* the signal
 - How to *model* the constraints
 - How to *search* for the most optimal answer

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 20

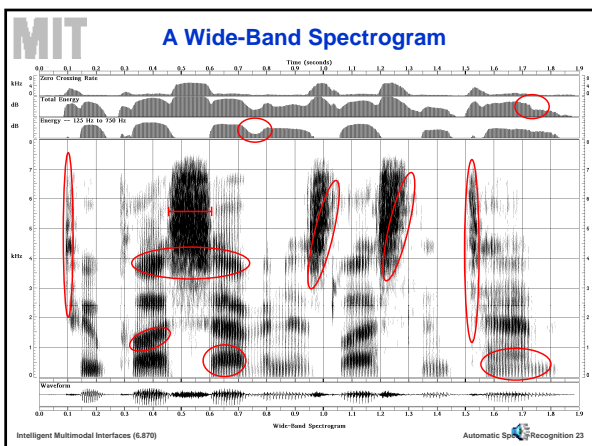
MIT Speech

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 21

MIT Speech Production

- Speech produced via coordinated movement of articulators
- Spectral characteristics of speech influenced by source, vocal tract shape, and radiation characteristics
- Speech articulation characterized by manner and place
 - Vowels: No significant constriction in the vocal tract; usually voiced
 - Fricatives: turbulence produced at a narrow constriction
 - Stops: complete closure in the vocal tract; pressure build up
 - Nasals: velum lowering results in airflow through nasal cavity
 - Semivowels: constriction in vocal tract, no turbulence

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 22



MIT Phonological Variation

- The acoustic realization of a phoneme depends strongly on the context in which it occurs

TEA TREE STEEP CITY BEATEN

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 24

MIT **Signal Processing**

Waveform

- Frame-based spectral feature vectors (typically every 10 milliseconds)
- Efficiently represented with Mel-frequency scale cepstral coefficients
 - Typically ~13 MFCCs used per frame

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 25

MIT

Models

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 26

MIT **Statistical Approach to ASR**

- Given acoustic observations, A , choose word sequence, W^* , which maximizes *a posteriori* probability, $P(W|A)$

$$W^* = \underset{W}{\operatorname{argmax}} P(W|A)$$
- Bayes rule is typically used to decompose $P(W|A)$ into acoustic and linguistic terms

$$P(W|A) = \frac{P(A|W)P(W)}{P(A)}$$

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 27

MIT **Probabilistic Framework**

- Words are typically represented as sequence of phonetic units
- Using phonetic units, U , expression expands to:

- Search must efficiently find most likely U and W
- Pronunciation and language models encoded in a graph

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 28

MIT **Language Modeling**

- ASR systems constrain possible word combinations by way of simple, but powerful, language models:
 - Finite-state network
 - Deterministic, sequential constraints (e.g., word-pair)
 - Probabilistic, sequential constraints (e.g., bigram, trigram)
- Trigram is the dominant language model for ASR:

$$P(w_n | w_{n-2}, w_{n-1})$$
- Much effort has gone into smoothing techniques for sparse data
- Task difficulty is measured by perplexity

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 29

MIT **Acoustic Modeling**

Waveform

- Feature vector scoring:

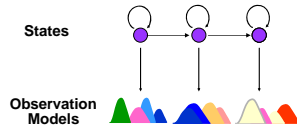
$$P(A|U) = \prod_{i=0}^N P(\bar{x}_i | u_i)$$
- Each phonetic unit modeled w/ a mixture of Gaussians:

$$P(\bar{x} | u) = \sum_{j=0}^M w_j N(\bar{x} | \mu_j, \Sigma_j)$$

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 30

Hidden Markov Models

- Dominant modeling framework used for speech recognition
- Generative model that predicts likelihood of observation sequence O being generated by state sequence Q
 - Either discrete or continuous observation models can be used



- HMMs can model words or sub-words (e.g., phones)
 - Sub-word HMMs concatenated to create larger word-based HMMs

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 31

Phonological Modeling

- Words described by phonemic baseforms
- Phonological rules expand baseforms into graph, e.g.,
 - Deletion of stop bursts in syllable coda (e.g., *laptop*)
 - Deletion of /t/ in various environments (e.g., *intersection*, *crafts*)
 - Gemination of fricatives and nasals (e.g., *this_side*, *in_nome*)
 - Place assimilation (e.g., *did_you* (/d ih jh uw/))

batter : b æ t f er

This can be realized phonetically as:

bcl b æ tcl t er Standard /t/

or as:

bcl b æ dx er Flapped /t/

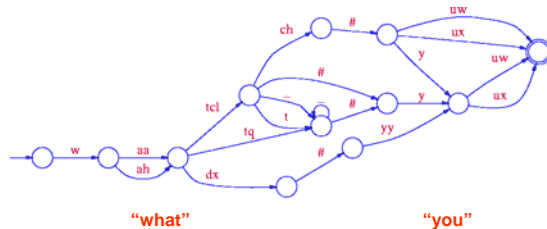
- Arc probabilities can be trained (i.e., $P(U|W)$)

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 32

Phonological Example

- Example of “what you” expanded with phonological rules
 - Final /t/ in “what” can be realized as released, unreleased, palatalized, or glottal stop, or flap



Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 33

Search

Intelligent Multimodal Interfaces (6.870)

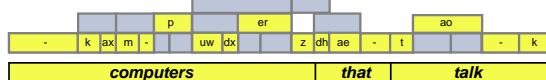
Automatic Speech Recognition 34

A Simple View of Speech Recognition

Waveform

Frame-based measurements

Acoustic models generate phonetic likelihoods



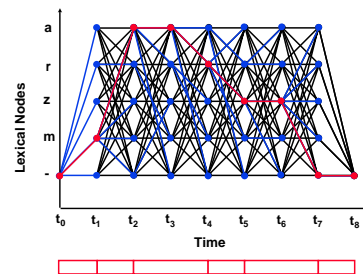
Probabilistic search finds most likely phone & word strings

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 35

Viterbi Search Example

- Viterbi search typically used in first-pass to find best path



- Relative and absolute thresholds used to speed-up search

Automatic Speech Recognition

36

MIT **A* Search Example**

- Second pass uses **backwards** A* search to find *N*-best paths
- Viterbi backtrace is used as future estimate for path scores

Automatic Speech Recognition 37

MIT **Search Issues**

- Search often uses forward and backward passes, e.g.,
 - Forward Viterbi search using bigram
 - Backwards A* search using bigram to create a word graph
 - Rescore word graph with trigram (i.e., subtract bigram scores)
 - Backwards A* search using trigram to create *N*-best outputs
- Search relies on two types of pruning:
 - Pruning based on relative likelihood score
 - Pruning based maximum number of hypotheses
 - Pruning provides tradeoff between speed and accuracy
- Multiple searches is a form of successive refinement
 - More sophisticated models can be used in each iteration

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 38

MIT

Representations

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 39

MIT **Finite-State Transducers**

- Most speech recognition constraints and results can be represented as finite-state automata:
 - Language models (e.g., n-grams and word networks)
 - Lexicons
 - Phonological rules
 - N-best lists
 - Word graphs
 - Recognition paths
- Common representation and algorithms desirable
 - Consistency
 - Powerful algorithms can be employed throughout system
 - Flexibility to combine or factor in unforeseen ways
- Finite-state transducers (FSTs) are effective for defining weighted relationships between regular languages
 - Extend FSAs by enabling transduction between input and output strings
 - Pioneered by researchers at AT&T for use in speech recognition

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 40

MIT **Example FST Operations**

- Construction (produce new functionality)
 - Union: $A \cup B$
 - Composition: $A \circ B$
- Optimization (retain original functionality)
 - Determinization
 - Minimization

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 41

MIT **Speech Recognition as Cascade of FSTs**

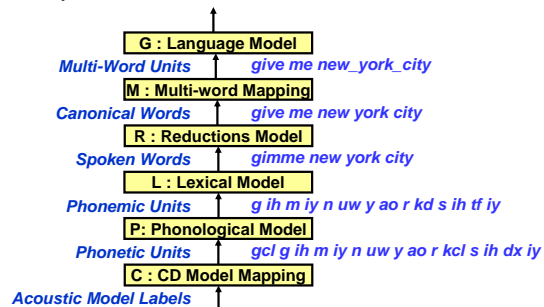
- Cascade of FSTs

$$O \circ (M \circ P \circ L \circ G)$$
 - G: language model (weighted words \leftarrow words)
 - L: lexicon (phonemes \leftarrow words)
 - P: phonological rule application (phones \leftarrow phonemes)
 - M: model topology (e.g., HMM) (states \leftarrow phones)
 - O: observations with acoustic model scores
- $(M \circ P \circ L \circ G)$ is single FST seen by search
- Search performs composition of O with $(M \circ P \circ L \circ G)$
- Gives great flexibility in how components are combined

Intelligent Multimodal Interfaces (6.870) Automatic Speech Recognition 42

Expanded FST Representation

- FST representation can be expanded for more efficient representation of lexical variation



Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 43

Related Areas of Research

- Speech understanding and spoken dialogue
- Multimodal interaction
- Audio-visual analysis (e.g., AVSR)
- Spoken document retrieval
- Speaker identification and verification
- Paralinguistic analysis (e.g., emotion)
- Acoustic scene analysis (e.g., CASA)
- ...

Intelligent Multimodal Interfaces (6.870)

Automatic Speech Recognition 44