

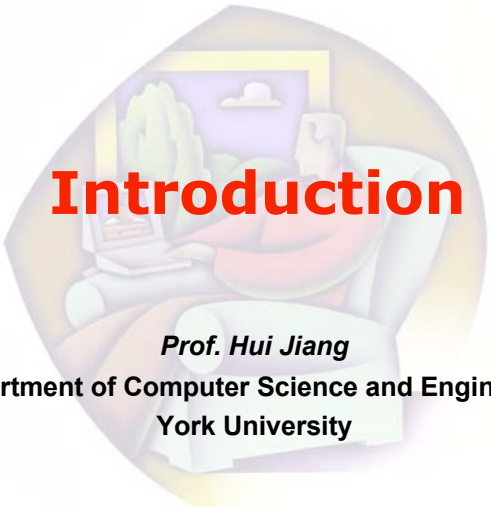
CSE6328.3
Speech & Language Processing

No.1

Introduction

Prof. Hui Jiang
Department of Computer Science and Engineering
York University

YORK UNIVERSITY redefine THE POSSIBLE.



CSE6328 Course Outline:
"Speech & Language Processing"

- Part I: Introduction (2 weeks)
 - Overview of speech and language technologies
 - Basic Knowledge of speech and spoken language
 - Math foundation: review
- Part II: Basic theory of pattern classification/verification (4—5 weeks)
 - Bayesian decision rule; Model estimation methods
 - Generative models: Gaussian, GMM, Markov Chain, HMM, Graphical models
 - Discriminative models: SVM, Neural networks (NN) and beyond
- Part III: Case studies (3—4 weeks)
 - Automatic speech recognition
 - Spoken language processing
- Part IV: Advanced topics – YOUR PARTICIPATION !! (1—2 weeks)
 - Choose an advanced topic in machine learning
 - Choose a journal article in speech/language or other applications
 - Self-study and oral presentation in class

Course Info

- Course Web site: <http://www.cs.yorku.ca/course/6328/>
- Course Format:
 - Lectures (10 weeks):
 - Covers basic data modeling, pattern classification theory;
 - Introduces some selected applications in speech recognition and spoken language processing.
 - Students' in-class presentations (2 weeks): based on basic theories in class, self-study recently published technical article and orally present it in class.
 - Choose an advanced topic in machine learning
 - Choose an application in speech/language or other domains
- Evaluation:
 - One assignment (10%) (roughly first 1/3 of the course)
 - Two lab projects (50%): report + oral presentation(?)
 - Advanced topic self-study and in-class presentation (30%)
 - Class Participation (10%)

Reference Materials

- Lecture notes
- Assigned reading materials through the course
- Reference books:
 - [1] Pattern Recognition and Machine Learning by C. M. Bishop. (Springer, ISBN 0-387-31073-8)
 - [2] *Pattern Classification* (2nd Edition) by R. O. Duda, P. Hart and D. Stork. (John Wiley & Sons, Inc., ISBN 0-471-05669-3)
 - [3] *Spoken Language Processing: a guide to theory, algorithm, and system development* by X.D. Huang, A. Acero, H.W. Hon. (Prentice Hall PTR, ISBN 0-13-022616-5)
 - [4] *Foundations of Statistical Natural Language Processing* by C. D. Manning and H. Schutze. (The MIT Press, ISBN 0-262-13360-1)
- Prerequisite:
 - First course in probability or statistics
 - First course in linear algebra or matrix theory
 - C/C++/Java and perl/shell programming skill (for project)

Speech Research and Technology

- **Speech Communication**
- **Speech Production and Perception**
- **Speech Analysis and Synthesis**
- **Speech and Audio Coding & Compression**
- **Speech Recognition and Understanding**
- **Speaker Identification and Verification**
- **Speech Enhancement**
- **Language Identification**
- **Dialogue Processing**

Language Research and Technology

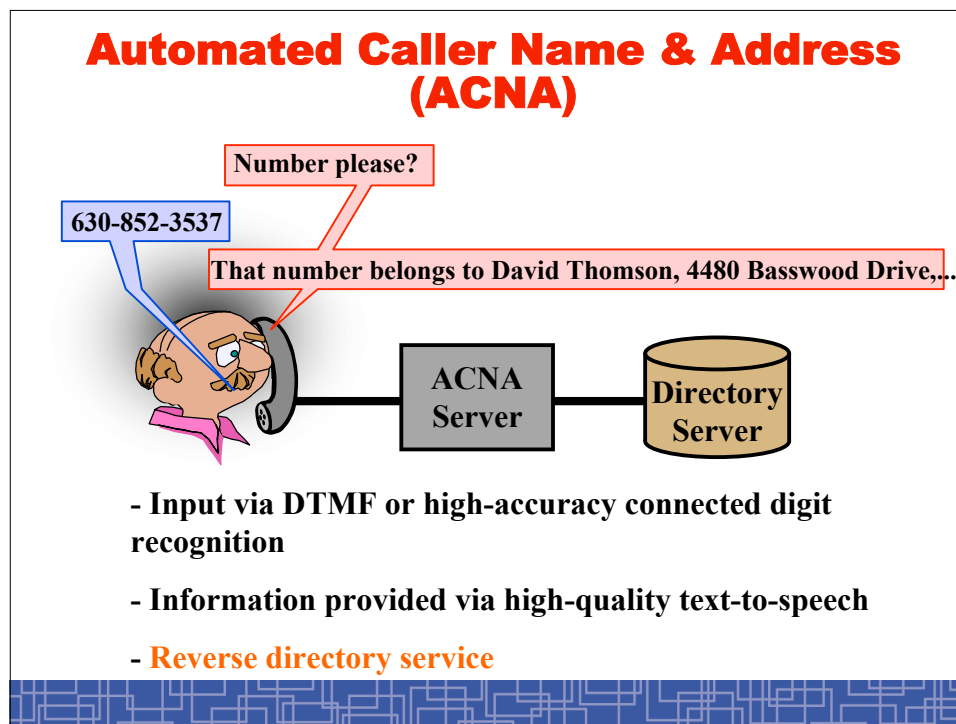
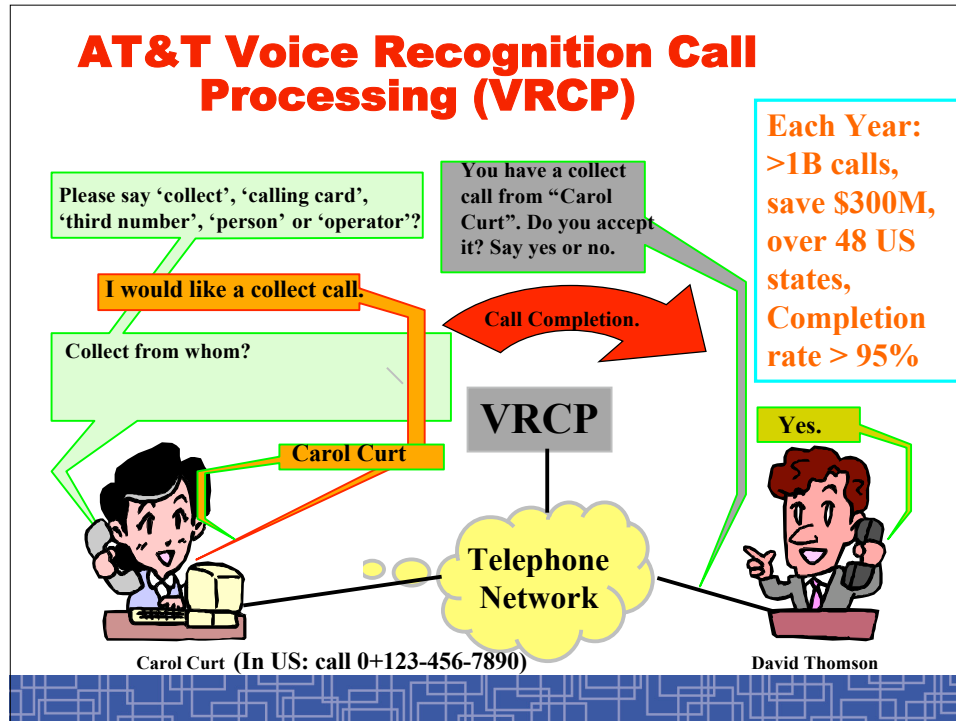
- **Written vs. Spoken Languages**
- **Computational Linguistics**
- **Corpus-Based Language Technologies**
- **Statistical Language Modeling**
- **Language Analysis and Generation**
- **Statistical Part-of-Speech Tagging**
- **Modeling Syntax and Semantics**
- **Statistical Text Understanding / Text Mining**
- **Probabilistic Parsing**
- **Text Categorization**
- **Statistical Machine Translation**
- **Information Retrieval**

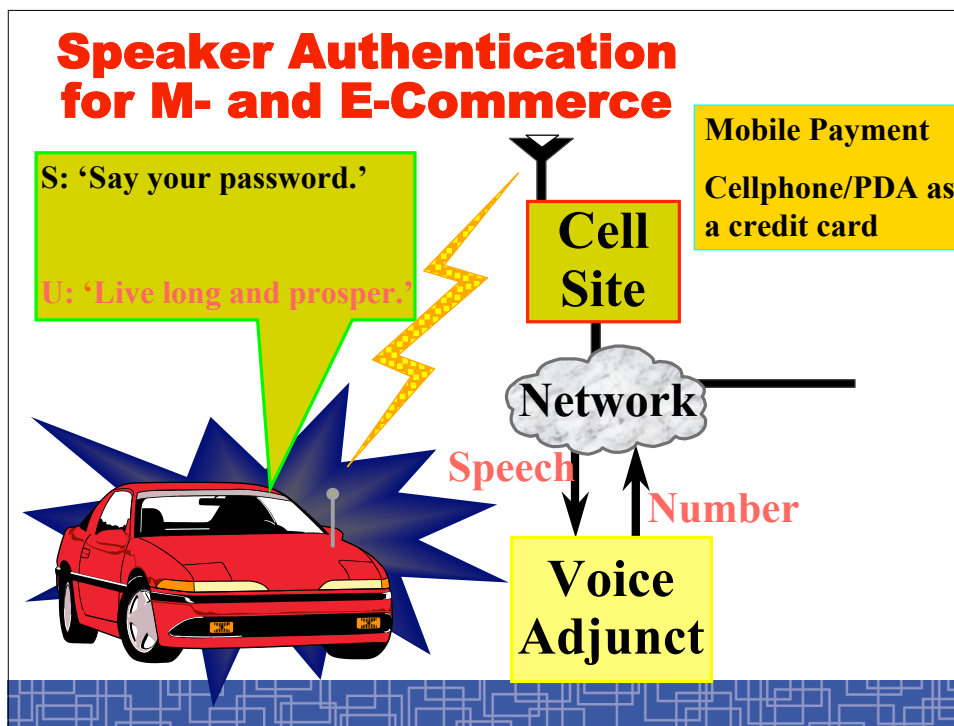
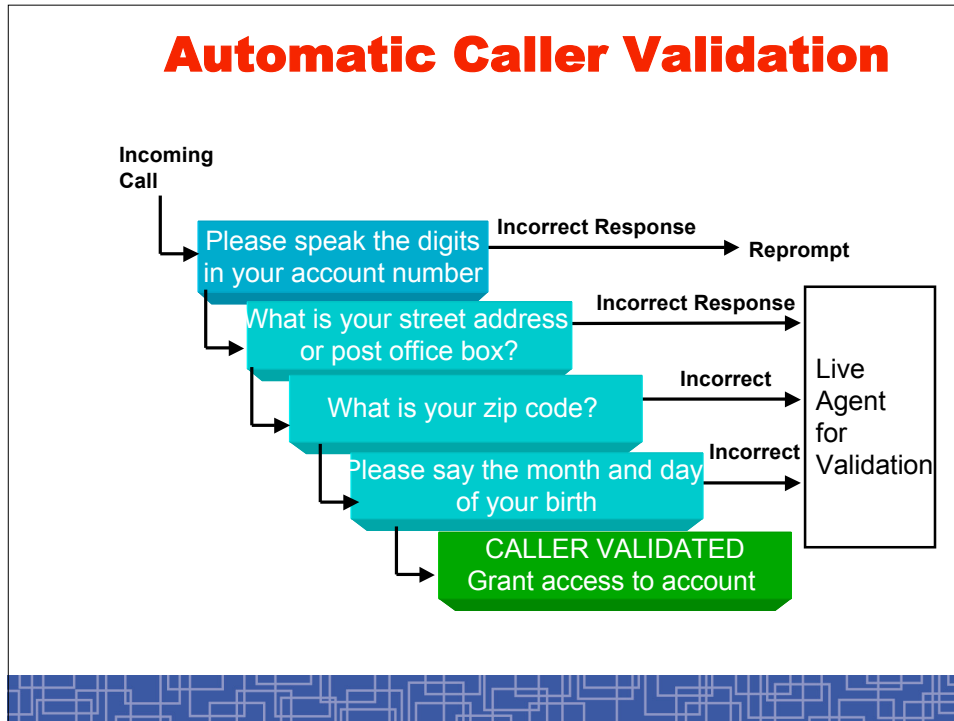
Applications of Speech and Language Technologies

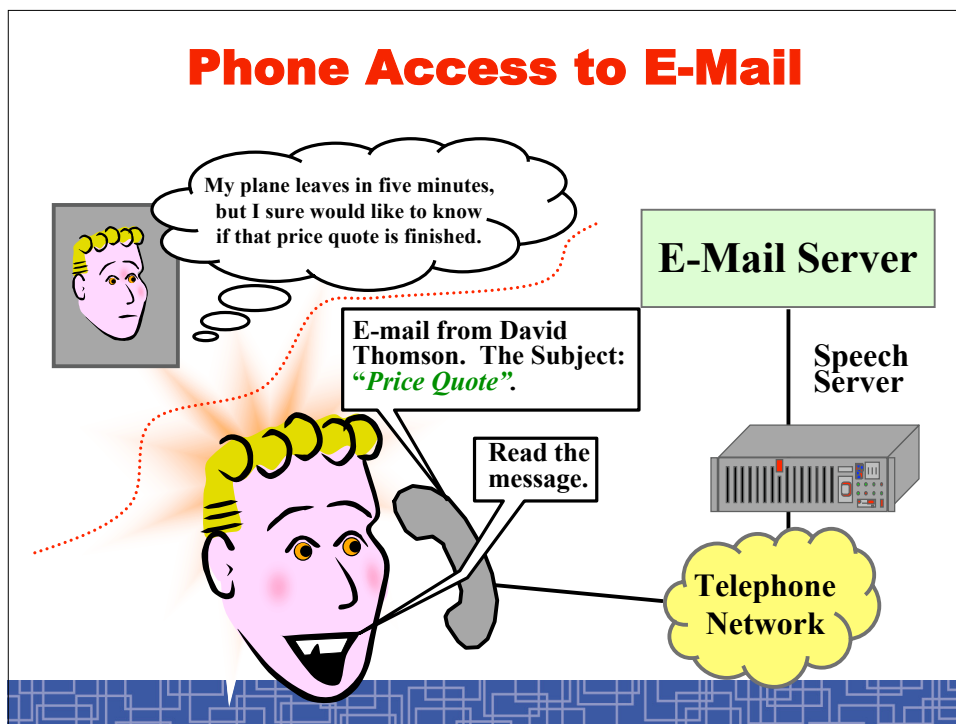
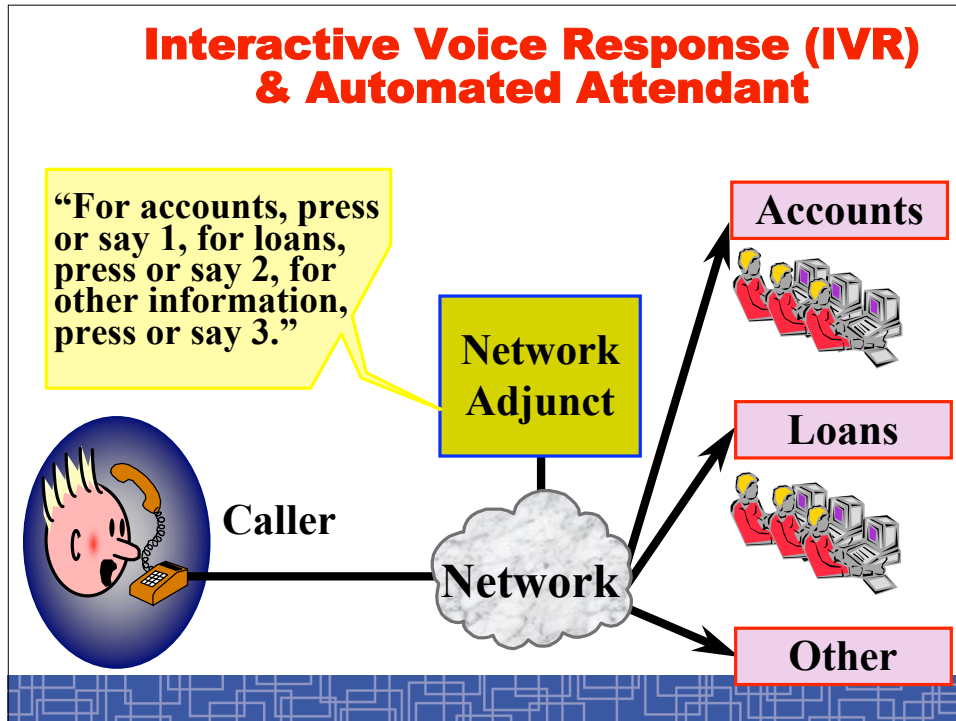
- Voice typewriter – dictation systems:
 - IBM, Microsoft, Nuance, etc.
- Applications in telecommunications:
 - AT&T, Lucent Bell Labs, Nuance, Philips, Motorola, etc.
 - Automatic Call Centers.
 - Google 411, Microsoft Tellme, Free 411.
- Applications related to the Internet:
 - Google’s voice search
 - Apple’s Siri (voice assistant)
 - More and more to emerge

Voice Typewriter (Dictation system)

- IBM ViaVoice® System
<http://www.ibm.com/software/voice/viavoice/>
- Microsoft Speech SDK and Whisper® System
<http://www.microsoft.com/speech/>
- Nuance’s Naturally Speaking System
<http://www.dragontalk.com/NATURAL.htm>
<http://www.scansoft.com/naturallspeaking/>







Natural Language Example: Movie Locator

What movies are playing at
the Rice Lake Square
theater in Wheaton?



Moviefone (777-film)



Other applications:

- “This is the operator, how can I help you?”
- “Hi, I’d like a large pizza with pepperoni with mushrooms toppings.”
- “Play all messages from Tom Smith.”
- Business (restaurant) locator, yellow pages
- Travel information systems (train/flight)
- TellMe, UA FlightInfo, Google-411
- L&H, Nuance, SpeechWorks, Philips, etc.

The Bell Labs “Natural Language Call Router”

- **Input:** user request (in speech or text)
- **Output:** desired destination related to the request (in a call center)
- **Data Preparation:** user (request, destination) pairs are grouped to train routing matrix using a data-driven approach
- **Technologies:** speech recognition, language modeling, call routing, dialogue generation

The United Airlines Flight Information System

- **Input:** user request (in speech)
- **Output:** desired flight information
- **Data Preparation:** flight schedule info is converted into groups of finite state grammars, prompts are used to guide users
- **Technologies:** speech recognition, language modeling, user modeling

The Stock Information System

- **Input:** user request (in speech)
- **Output:** desired update stock info
- **Data Preparation:** convert stock names into pronunciation entries
- **Technologies:** speech recognition, pronunciation modeling, database, text-to-speech synthesis

The Google 411 Service

- <http://www.google.com/goog411/>

Google GOOG-411 1-800-GOOG-411

Dial from any phone
1-800-GOOG-411
(1-800-466-4411)

411
411 FREE FROM Google

YouTube

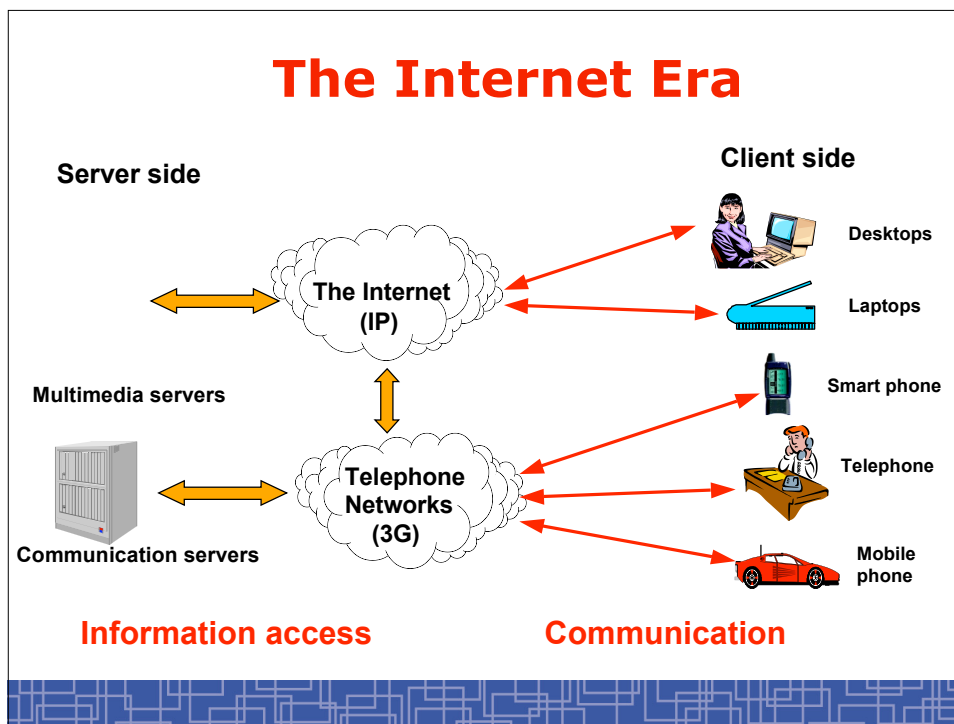
0:02 menu

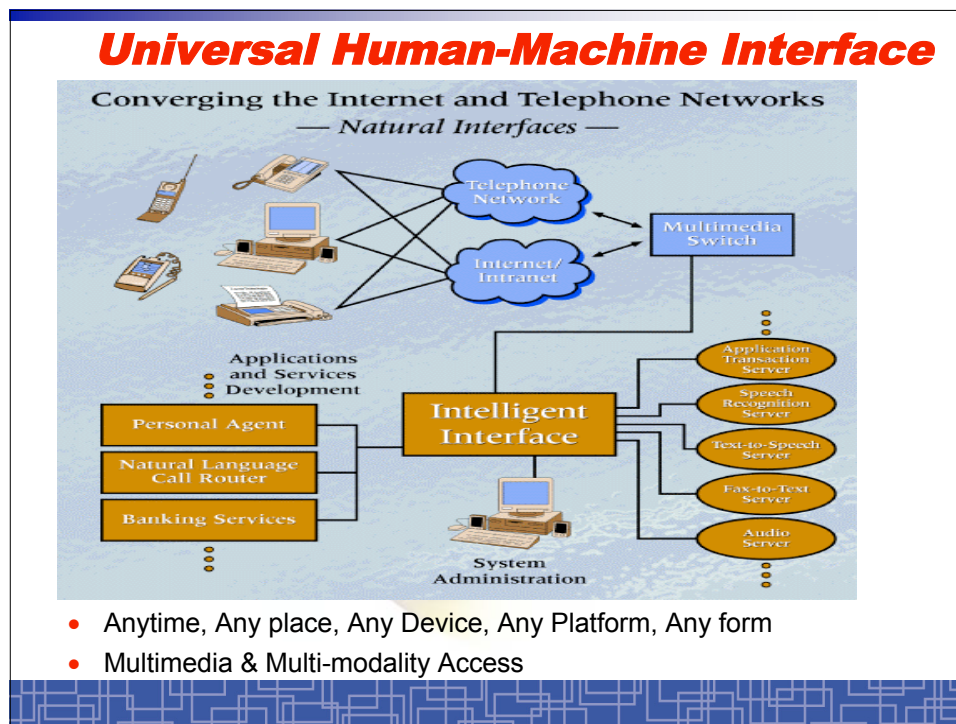
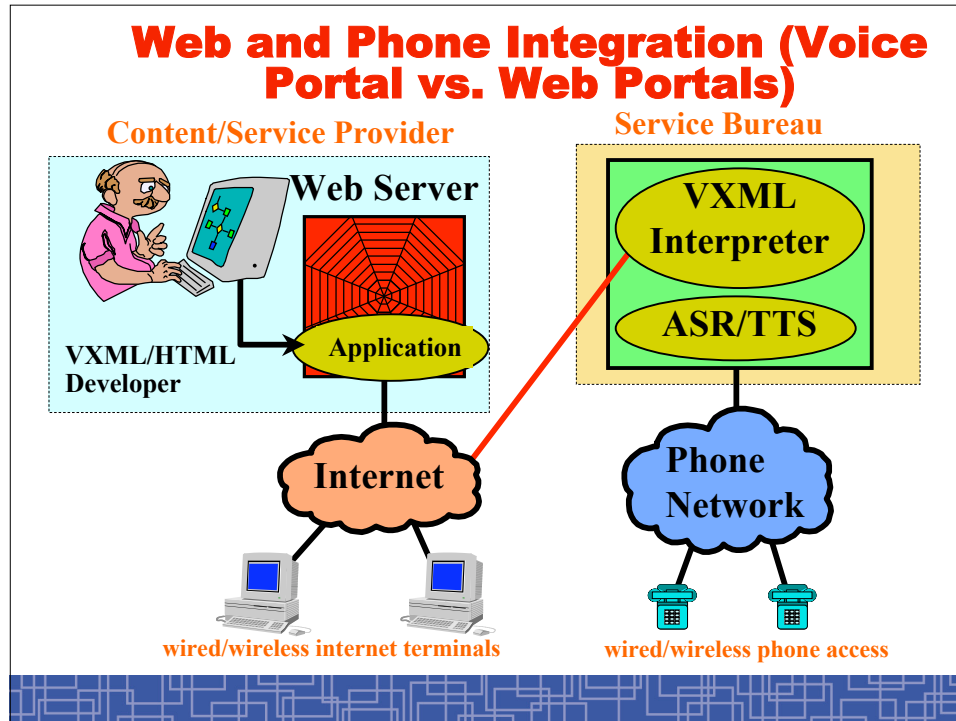
- 1 Dial 1-800-GOOG-411 from any phone
- 2 State the location and business type
- 3 Connect to the business for free
- 4 Done!

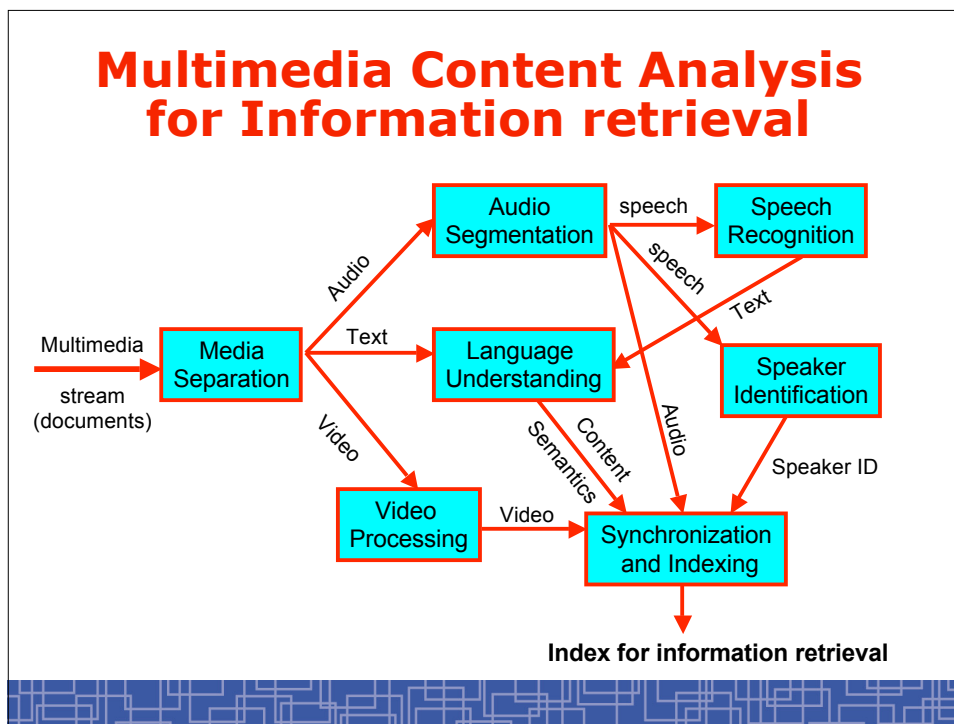
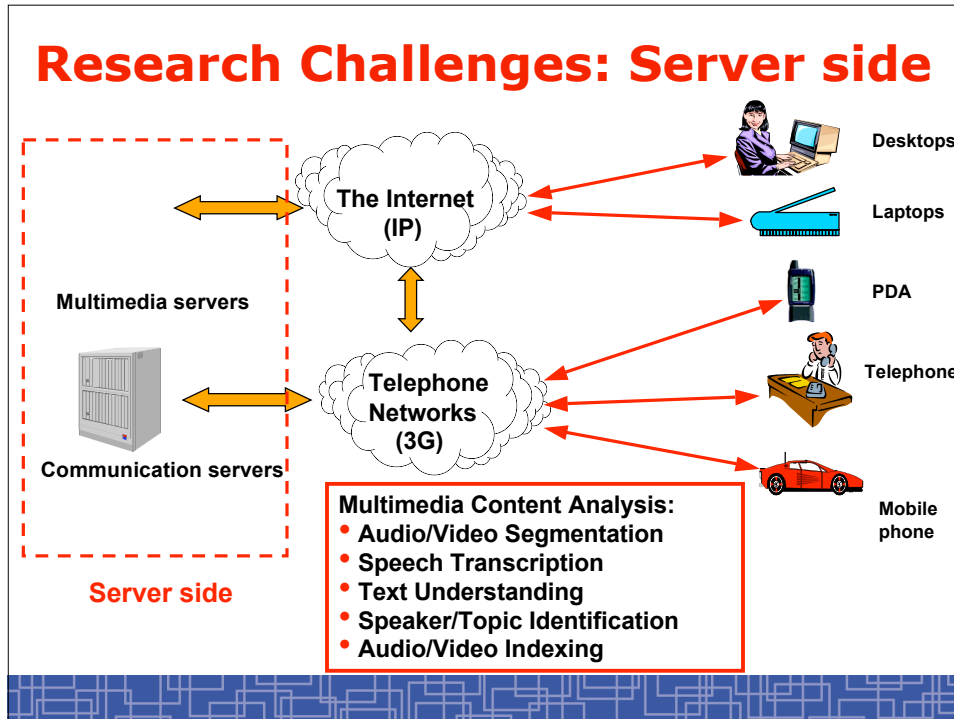
About GOOG-411
Google's new 411 service is free, fast and easy to use. Give it a try now and see how simple it is to find and connect with local businesses for free.

[Learn more - FAQ](#)

Liked the video? Want to comment or guess who the voice of GOOG-411 is? Post your opinion on our [YouTube page](#).







Video and Audio Segmentation

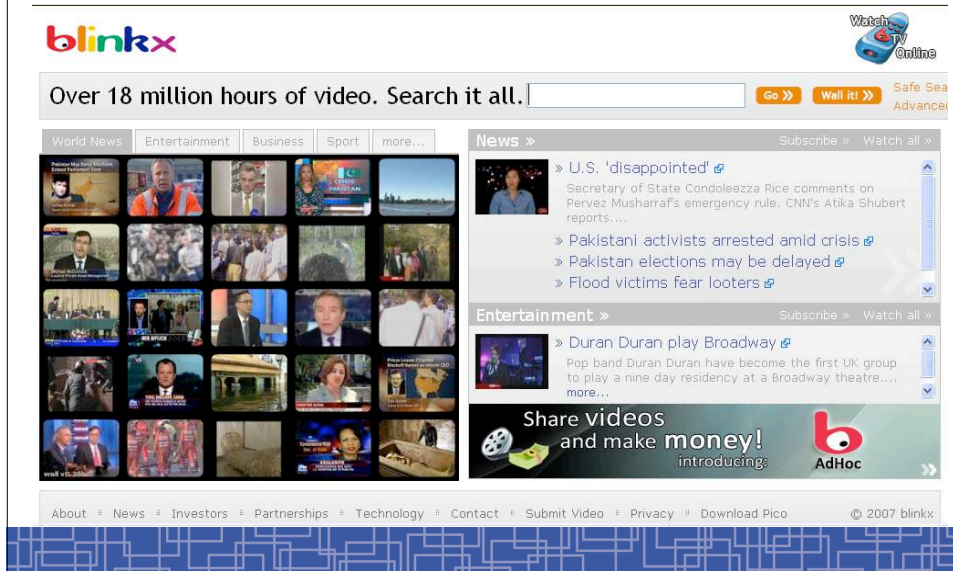


Archiving & Browsing Multimedia Data

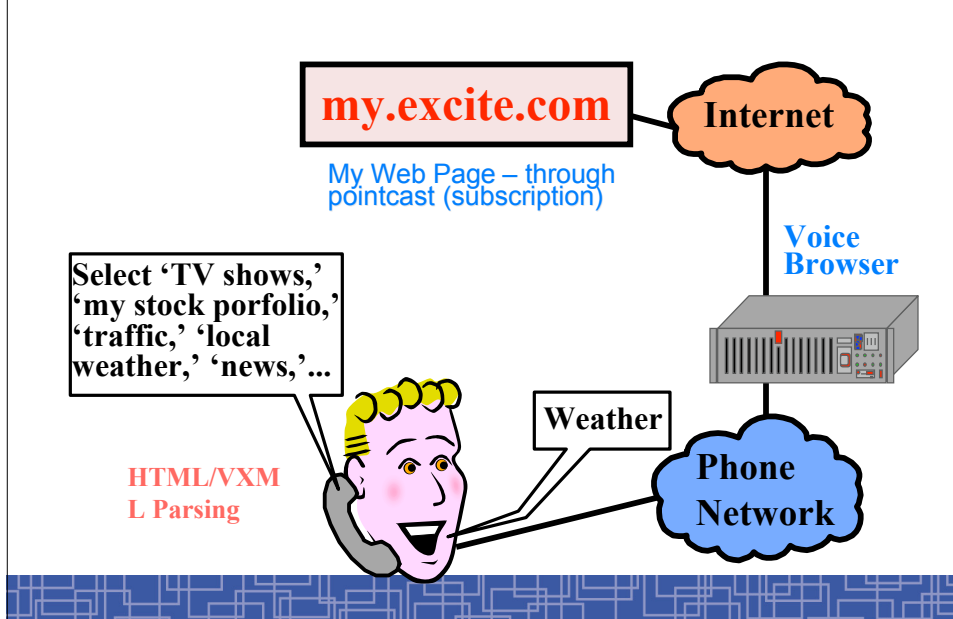
- **Input:** user request (in speech or text)
- **Output:** desired audio/video segments
- **Data Preparation:** video/audio segments with semantic description and (recognized) text for easy browsing, like MPEG7 descriptions (creation of indexing info for access is key)
- **Technologies:** speech recognition, video processing, multimedia segmentation and data mining, fusion of audio/video/caption information and presentation, etc.

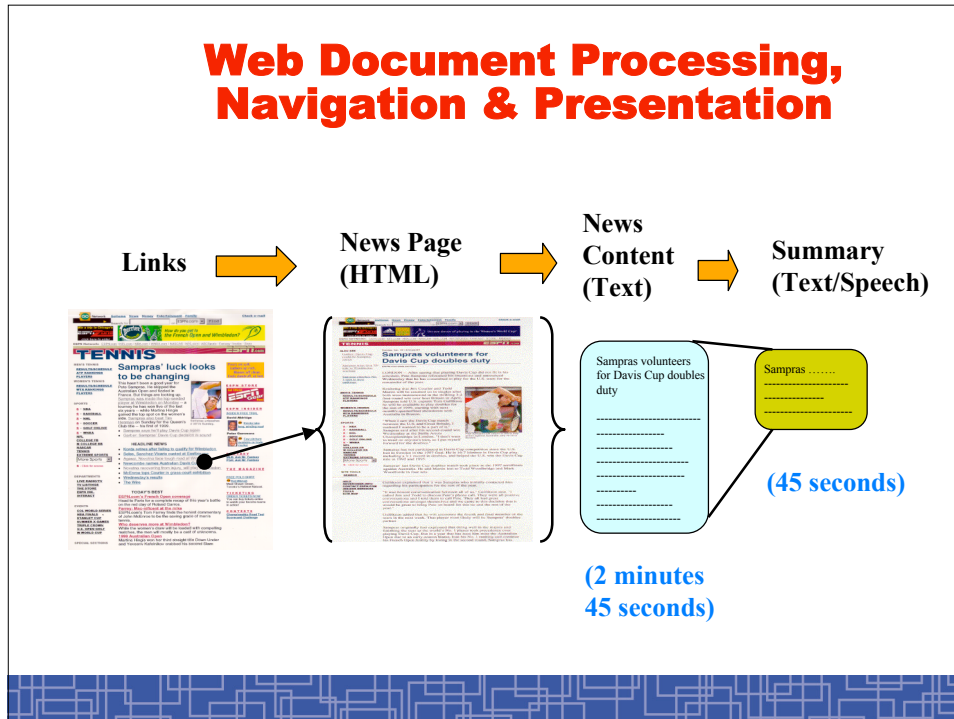
Video Search -- Blinkx

- WWW: <http://www.blinkx.com/>



Voice Browsing Applications





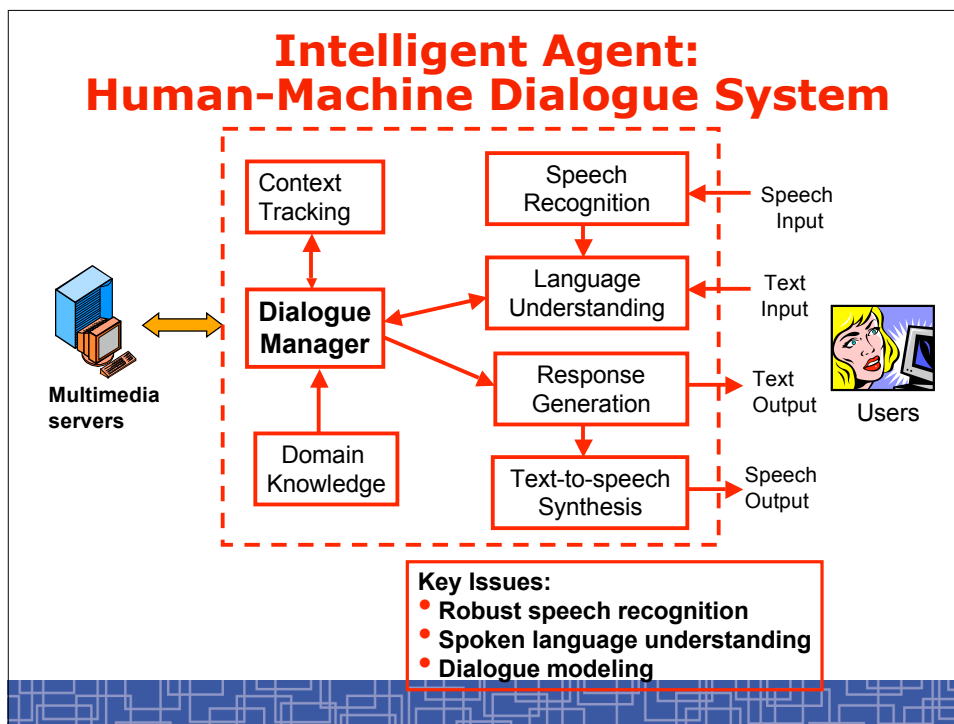
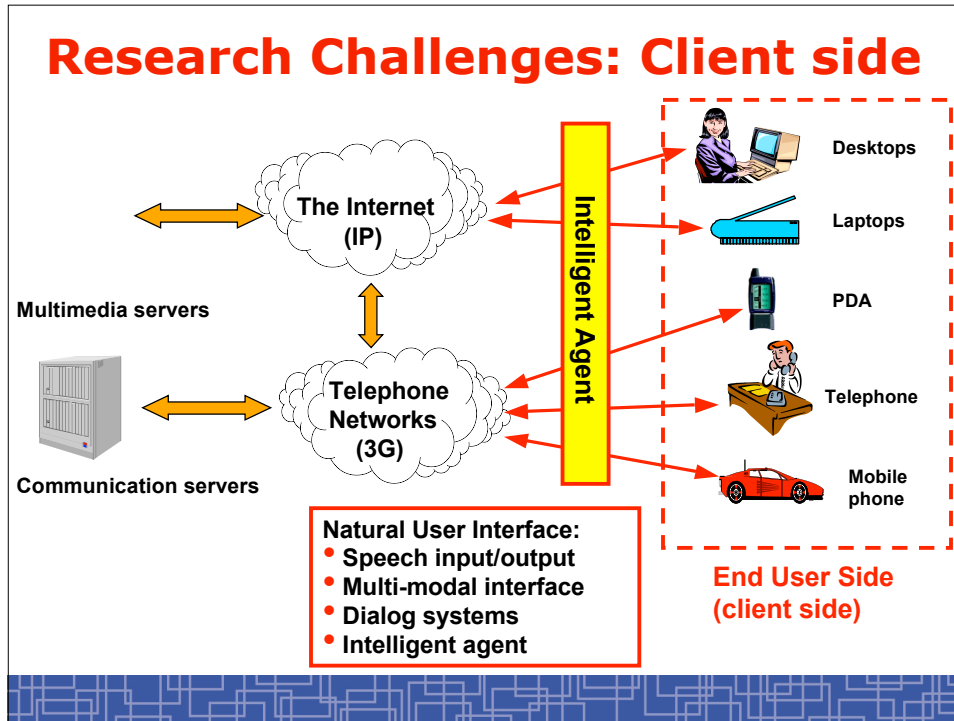
Google's Voice Search

- Google's voice search

Speak now



Search by voice
Speak your queries instead of typing.

[Watch video](#) [Learn more](#)



Apple's Siri

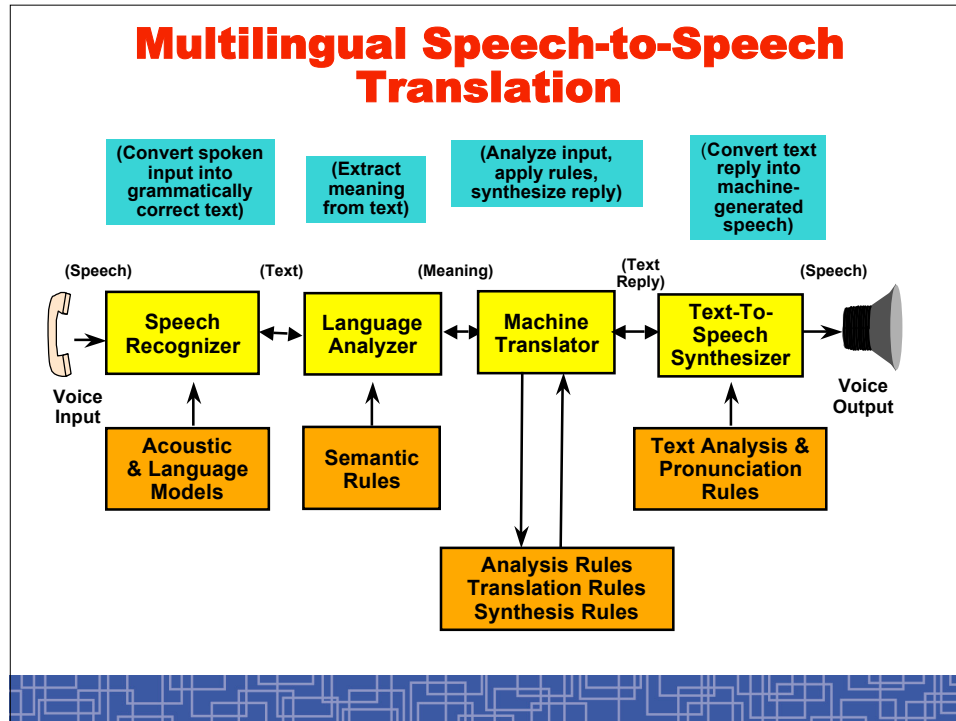
- Voice Assistant: *Siri*



Speech-To-Speech Translation

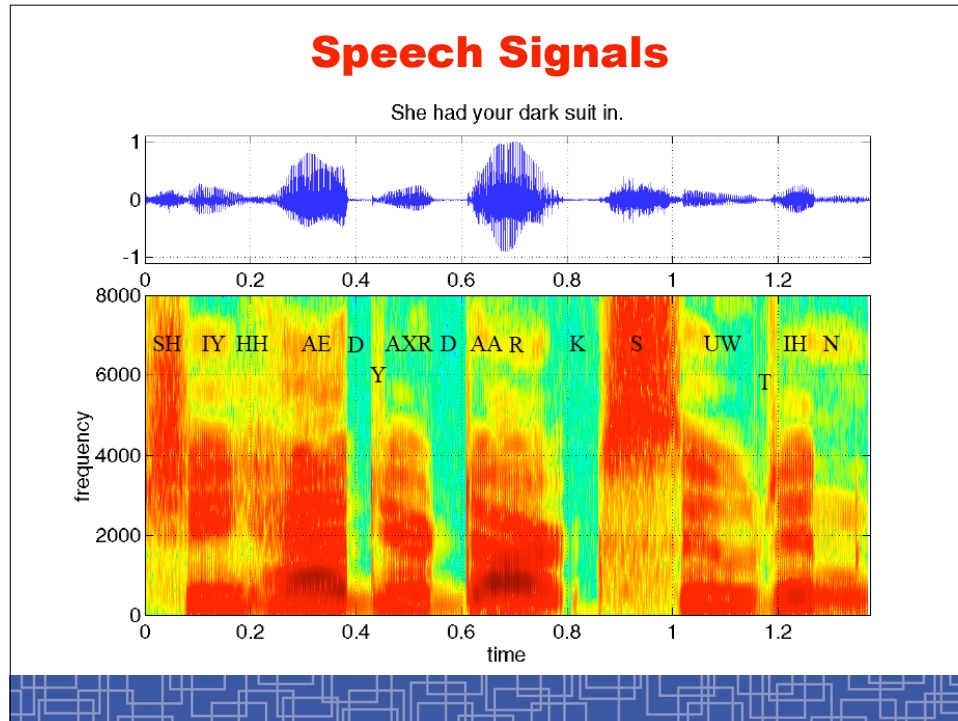
Bi-lingual Conversation





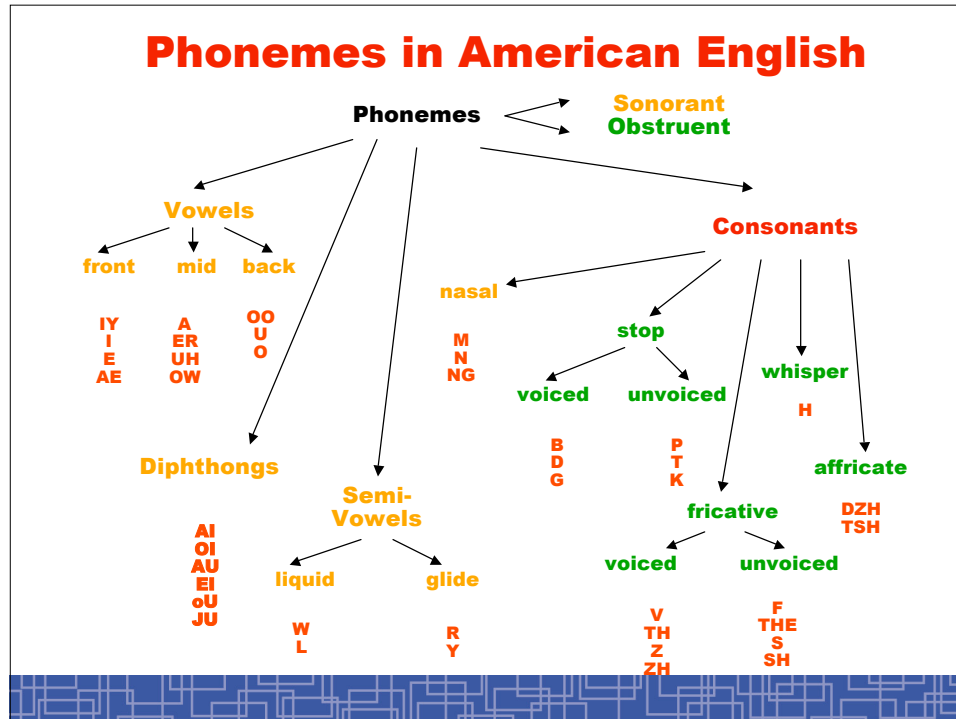
Statistical Pattern Classification

- Feature extraction:
 - Need to know objects to extract good features
 - Varies a lot among different applications (speech, audio, text, image, video, biological sequences, etc)
 - Statistical model training/learning
 - Inference, matching, decision
- } The basic theories common to various applications



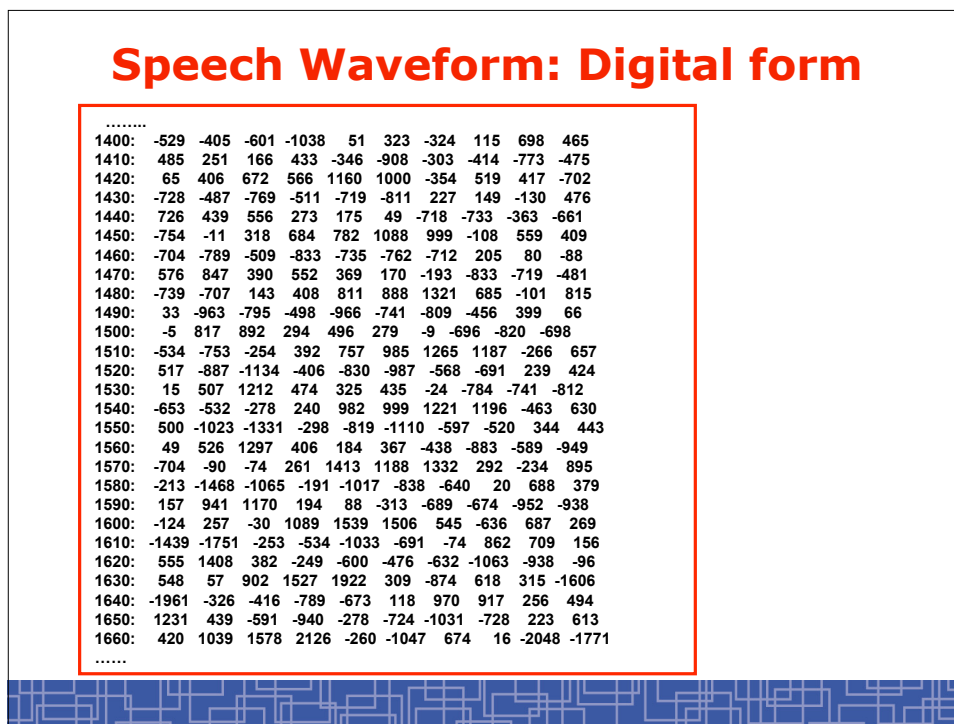
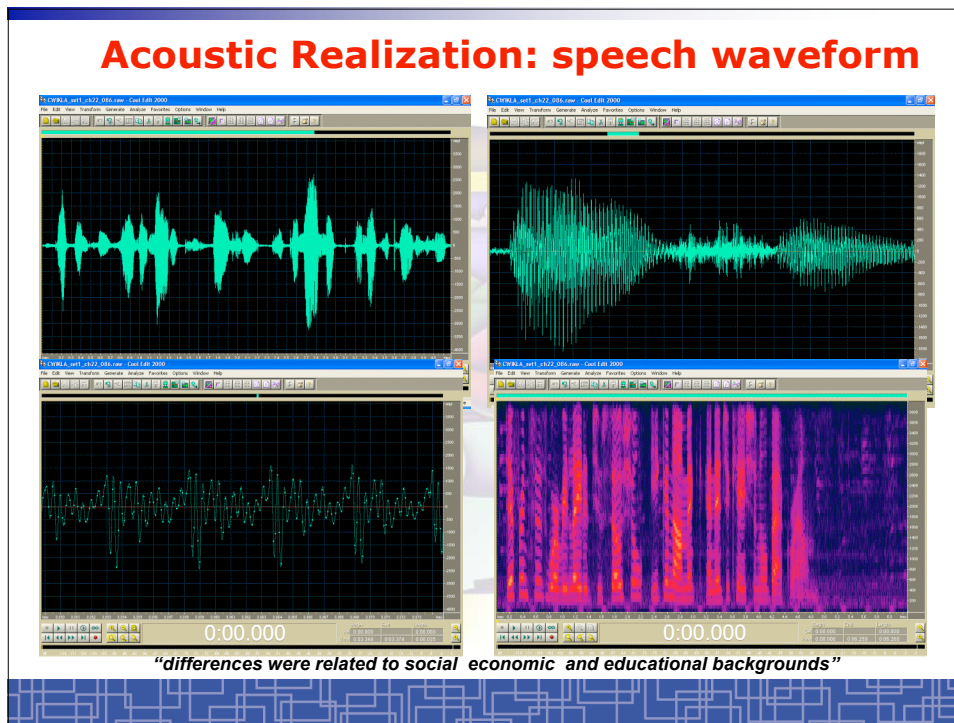
Fundamental Speech Units

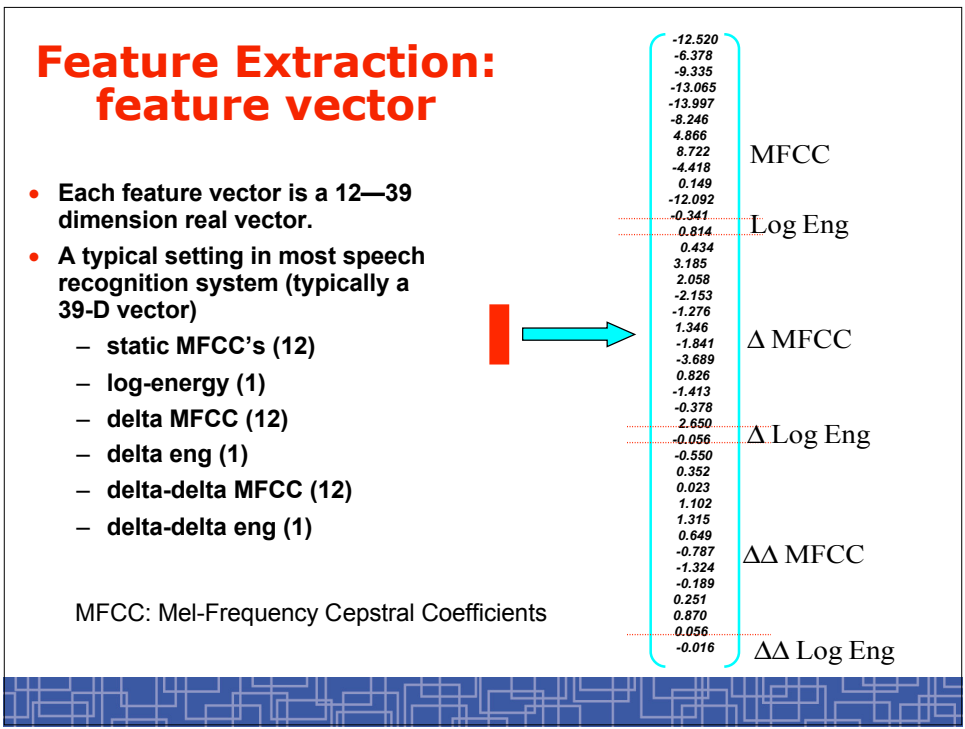
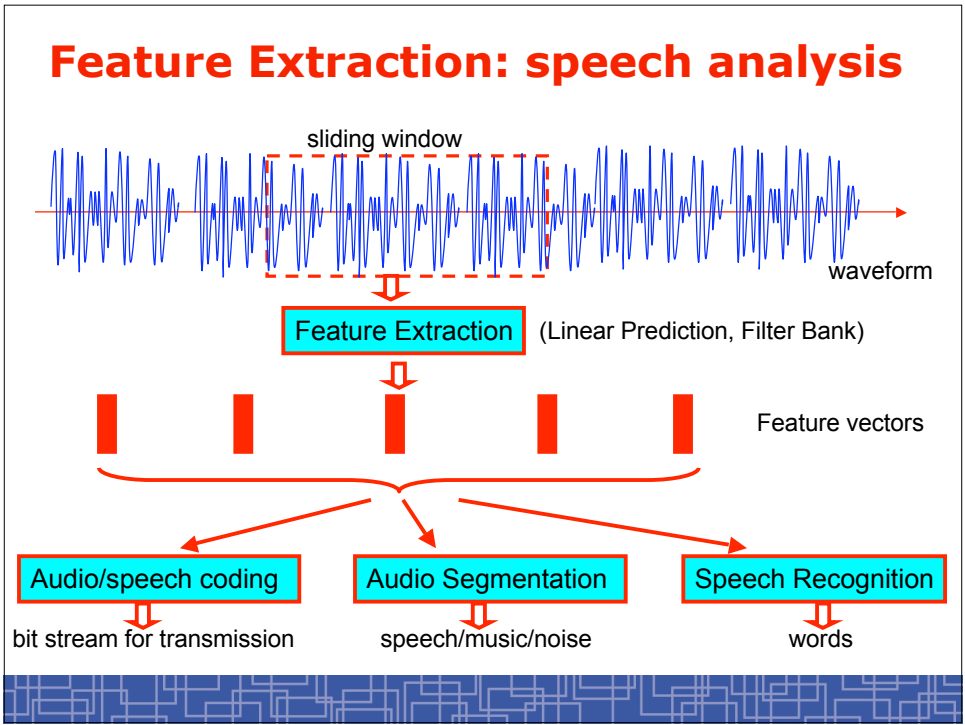
- Sentence/utterance → Phrase → Word → Syllable → Phone
- Phone
 - abstract name is called “phoneme”
 - infinite number of acoustic realization
 - monophone: context-independent phone
 - allophone: context-dependent phone
- Other considerations:
 - Language dependency
 - Task or vocabulary dependency
 - digit (small size but critically important)
- Example:
 - Sentence: How do they turn out later ?*
 - Syllables: How do they turn out la-ter?*
 - Phones: h aw d uh dh eh t er n aw t l ai t er*



Coarticulation

- Phones exhibit consistent acoustic characteristics if pronounced in isolated; but large acoustic variations may appear if uttered in different contexts.
- Coarticulation: acoustic realization of a phone is largely affected by its neighboring contexts.
- Reason: in speech production, articulatory gestures follow dynamics constrained by mechanical time constants associated with the articulator to keep the effort of muscles to a minimum.
- In speech recognition, how to model a phone:
 - Context-independent phone modeling – monophone: treat each phoneme equally no matter where it appears.
(in American English → 42 distinct phone units to model)
 - Context-dependent phone modeling:
 - Left (or right) biphone: a phone unit varies based on left(right) neighboring phone. (American English → 42X42 distinct units)
 - Triphone: a phone varies based on both left and right adjacent phones. (American English → 42X42X42 distinct units)





What's MFCC?

Step 1:

Fig. 5.3 Mel-Scale Filter Bank

Step 2: DCT (Discrete Cosine Transform) to de-correlate

$$c_i = \sqrt{\frac{2}{N}} \sum_{j=1}^N m_j \cos\left(\frac{\pi i}{N}(j - 0.5)\right)$$

MFCC: Mel-Frequency Cepstral Coefficients

Energy Measure

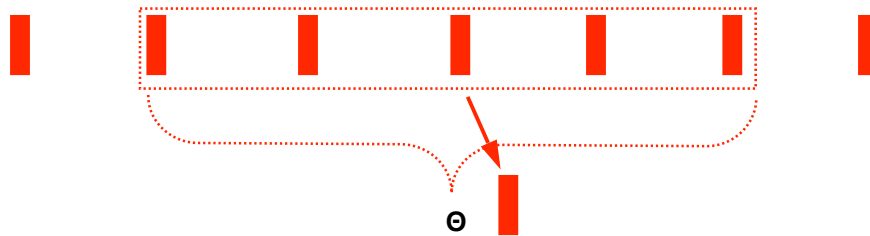
- For each frame, log-energy is calculated as:

$$E = \log \sum_{n=1}^N s_n^2$$

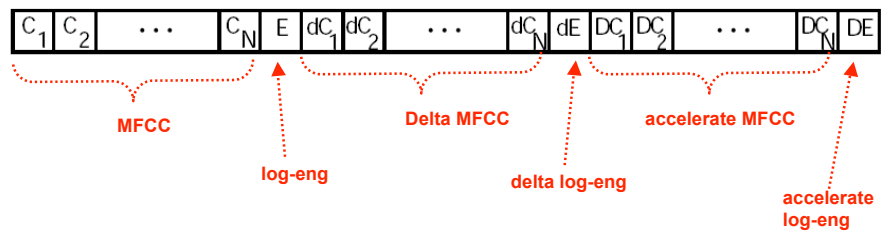
What Delta (Δ) and Acceleration ($\Delta\Delta$) Coefficients?

1. Delta coefficients: difference of MFCC among consecutive frames.
2. Acceleration coefficients: difference of delta among consecutive frames

$$d_t = \frac{\sum_{\theta=1}^{\Theta} \theta (c_{t+\theta} - c_{t-\theta})}{2 \sum_{\theta=1}^{\Theta} \theta^2}$$



Feature Vector Layout for each frame of speech



In most cases, $N=12 \rightarrow$ 39-dimension feature vector for each frame

Spoken Language Processing

- Style: written vs. spoken language
 - Written → formal; spoken → casual
- Disfluencies in spoken language:
 - Filled pauses: *um*
 - Repetitions: ... *the—the* ...
 - Repairs: ... *on Thursday – on Friday* ...
 - False starts: *I like ... – what I always get is ...*
- Lots of ungrammatical sentences exist in spoken language.
- In spoken language system:
 - speech recognition errors
- Obviously, spoken language processing is much harder.
- Our goal: build spoken language systems in some very constrained domains to perform shallow understanding:
 - Topic identification
 - Key-word spotting → to obtain gist and/or key message.
 - etc.

Feature Extraction to represent Text Document

- Text document ==> bag of words
- Key words vs. stop words
- "Russia cancels an **exhibition** of Russian **art** in London over fears the **art** could be seized to settle legal claims."
- => "exhibition(1) art(2)" <=> "art(2) exhibition(1)"
- Raw Feature vector:
- May need some normalization:
 - e.g. *TF-IDF*,...

| |
|---|
| 0 |
| ⋮ |
| 2 |
| 0 |
| ⋮ |
| 1 |
| 0 |
| ⋮ |

Statistical Machine Translation

- Score and rank all possible combinations...

