

# Speech Processing

Berlin Chen

Department of Computer Science & Information Engineering  
National Taiwan Normal University

# Course Contents

- Both the theoretical and practical issues for spoken language processing will be considered
- Technologies for Automatic Speech Recognition (ASR) and associated applications will be further emphasized
- Topics to be covered
  - Fundamentals and Statistical Modeling Paradigms
    - Spoken Language Structure
    - Hidden Markov Models
    - Speech Signal Analysis and Feature Extraction
    - Acoustic and Language Modeling
    - Search/Decoding Algorithms
  - Systems and Applications
    - Keyword Spotting, Dictation, Speaker Recognition, Spoken Dialogue, Speech-based Information Retrieval, etc.

# Some Textbooks and References (1/3)

- References books
  - X. Huang, A. Acero, H. Hon. Spoken Language Processing, Prentice Hall, 2001
  - L. Rabiner, R. Schafer, Theory and Applications of Digital Speech Processing, Pearson, 2011
  - D. Yu and L. Deng, Automatic Speech Recognition: A Deep Learning Approach, Springer, 2015
  - Jacob Benesty (ed.), M. Mohan Sondhi (ed.), Yiteng Huang (ed.), Springer Handbook of Speech Processing, Springer, 2007
  - M.J.F. Gales and S.J. Young. The Application of Hidden Markov Models in Speech Recognition. Foundations and Trends in Signal Processing, 2008
  - C. Manning and H. Schutze. Foundations of Statistical Natural Language Processing. MIT Press, 1999
  - J. R. Deller, J. H. L. Hansen, J. G. Proakis. Discrete-Time Processing of Speech Signals. IEEE Press, 2000
  - F. Jelinek. Statistical Methods for Speech Recognition. MIT Press, 1999
  - L. Rabiner, B.H. Juang. Fundamentals of Speech Recognition. Prentice Hall, 1993
  - 王小川教授，語音訊號處理，全華圖書 2004

# Some Textbooks and References (2/3)

- Reference papers


1. L. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," Proceedings of the IEEE, vol. 77, No. 2, February 1989
2. A. Dempster, N. Laird, and D. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," J. Royal Stat. Soc., Series B, vol. 39, pp. 1-38, 1977
3. Jeff A. Bilmes "A Gentle Tutorial of the EM Algorithm and its Application to Parameter Estimation for Gaussian Mixture and Hidden Markov Models," U.C. Berkeley TR-97-021
4. J. W. Picone, "Signal modeling techniques in speech recognition," proceedings of the IEEE, September 1993, pp. 1215-1247
5. R. Rosenfeld, "Two Decades of Statistical Language Modeling: Where Do We Go from Here?," Proceedings of IEEE, August, 2000
6. H. Ney, "Progress in Dynamic Programming Search for LVCSR," Proceedings of the IEEE, August 2000
7. H. Hermansky, "Should Recognizers Have Ears?", Speech Communication, 25(1-3), 1998

## Some Textbooks and References(3/3)

8. Frederick Jelinek, "[The Dawn of Statistical ASR and MT](#)," Computational Linguistics, Vol. 35, No. 4. (1 December 2009), pp. 483-494
9. L.S. Lee and B. Chen, "Spoken document understanding and organization," *IEEE Signal Processing Magazine*, vol. 22, no. 5, pp. 42-60, Sept. 2005
10. M. Gilbert and J. Feng, "Speech and Language Processing over the Web," *IEEE Signal Processing Magazine* 25 (3), May 2008
11. C. Chelba, T.J. Hazen, and M. Saraclar. Retrieval and Browsing of Spoken Content. *IEEE Signal Processing Magazine* 25 (3), May 2008
12. S. Young et al., The HTK Book. Version 3.4: <http://htk.eng.cam.ac.uk>
13. J. Schalkwyk et al., "[Google Search by Voice: A case study](#)," 2010

# Website for This Course

- Visit <http://berlin.csie.ntnu.edu.tw/> and then click the link “Spring 2016: Speech Processing”



The screenshot shows a web browser window with the address bar containing the URL [berlin.csie.ntnu.edu.tw/Courses/Speech%20Processing/Speech%20Processing\\_Main\\_2016S.htm](http://berlin.csie.ntnu.edu.tw/Courses/Speech%20Processing/Speech%20Processing_Main_2016S.htm). The page content includes the course title "Speech Processing" in blue, followed by "Spring 2016" and the schedule "9:10 ~12:10 am, Mondays". The instructor is listed as "Dr. Berlin Chen (陳柏琳)".

**Topic List and Schedule:**

02/22	<a href="#">Course Overview &amp; Introduction</a>	<b>Readings:</b> 1. F. Jelinek, The Speech Recognition Problem, Chapter 1 of the book "Statistical Methods for Speech Recognition" 2. L. Rabiner. <a href="#">The Power of Speech</a> . Science, Vol. 301, pp. 1494-1495, Sep. 2003. 3. S. Young. " <a href="#">Talking to Machines</a> ," Royal Academy of Engineering Ingenia, 54, pp. 40-46, 2013. 4. Frederick Jelinek, " <a href="#">The Dawn of Statistical ASR and MT</a> ," Computational Linguistics, Vol. 35, No. 4. (1 December 1979)

**Reference Books:**

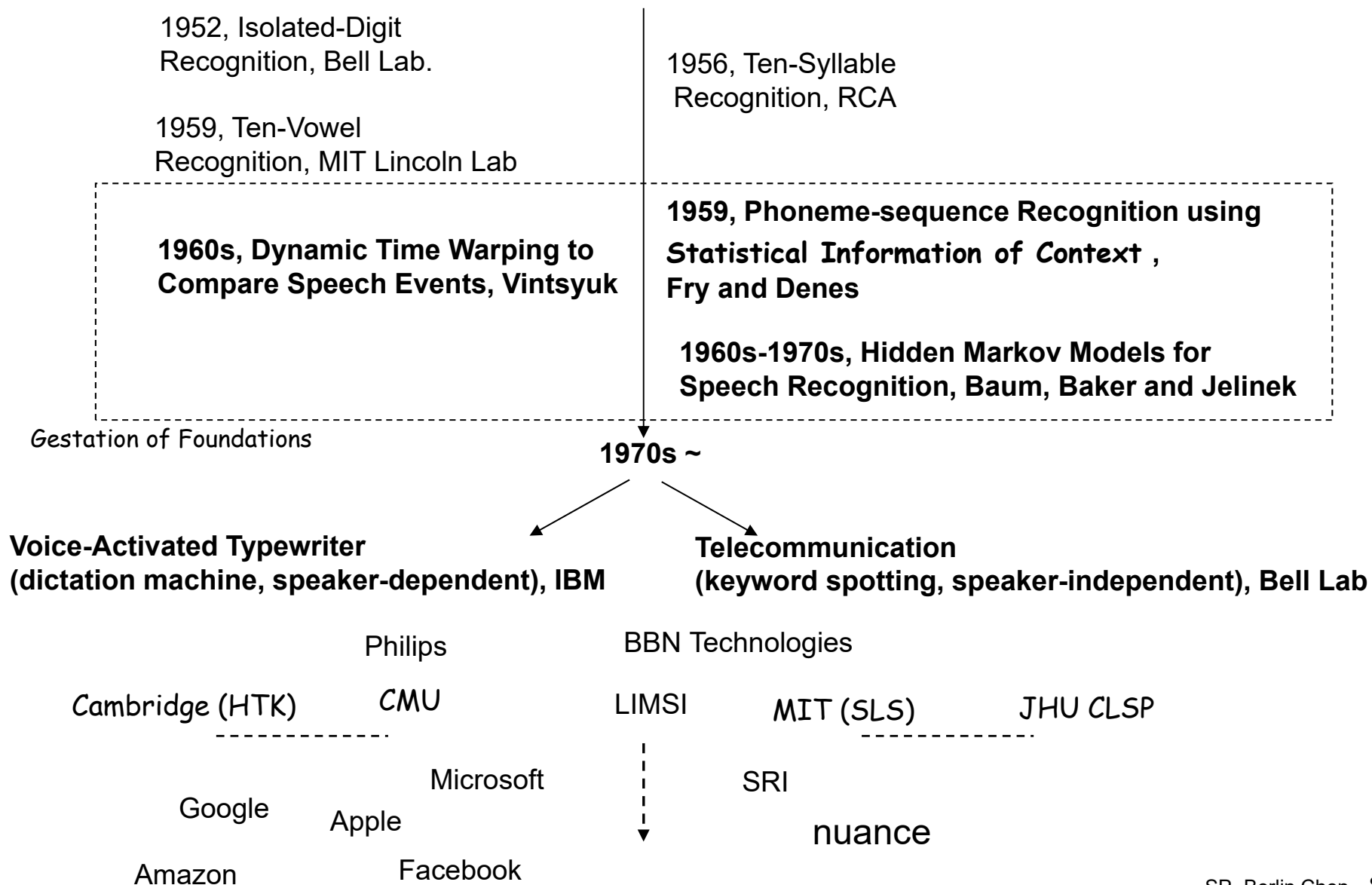
- L. Rabiner, R. Schafer, Theory and Applications of Digital Speech Processing, Pearson, 2011
- X. Huang, A. Acero, H. Hon, [Spoken Language Processing: A Guide to Theory, Algorithm and System Development](#), Prentice Hall, 2001
- Jacob Benesty, M. Mohan Sondhi, Yiteng Huang (ed.), [Springer Handbook of Speech Processing](#), Springer, 2007
- Tuomas Virtanen, Rita Singh, Bhiksha Raj (ed.), [Techniques for Noise Robustness in Automatic Speech Recognition](#), John Wiley & Sons, 2013
- L. Rabiner, B.H. Juang, "Fundamentals of Speech Recognition", Prentice Hall, 1993
- [M.J.F. Gales](#) and [S.J. Young](#). [The Application of Hidden Markov Models in Speech Recognition](#). Foundations and Trends in Signal Processing, 2008
- L. Rabiner and R.W. Schafer. [Introduction to Digital Speech Processing](#). Foundations and Trends in Signal Processing, 2007
- W. Chou., [B.H. Juang](#). [Pattern Recognition in Speech and Language Processing](#). CRC Press, 2003
- S. Young et al., "The HTK Book", Version 3.2, 2002. "<http://htk.eng.cam.ac.uk>"
- T. F. Quatieri, "Discrete-Time Speech Signal Processing - Principles and Practice," Prentice Hall, 2002
- F. Jelinek, "Statistical Methods for Speech Recognition," The MIT Press, 1999

# Introduction

## References:

1. B. H. Juang and S. Furui, "Automatic Recognition and Understanding of Spoken Language - A First Step Toward Natural Human-Machine Communication," *Proceedings of IEEE*, August, 2000
2. I. Marsic, A. Medl, and J. Flanagan, "Natural Communication with Informatio Systems," *Proceedings of IEEE*, August, 2000

# Historical Review





# Speech Processing vs. Text Processing

## – Recognition, Analysis and Understanding

- **Text**: analyze and understand text
- **Speech**: recognize speech (i.e., ASR), and subsequently analyze and understand the recognized text (propagations of ASR errors)

## – Variability

- **Text**: different synonyms to refer to a specific semantic object or meaning, such as 台灣師範大學, 師大, 教育界龍頭, etc.
- **Speech**: an infinite number of utterances with respect to the same word (e.g., 台灣師範大學)
  - Manifested by a wide variety of oral phenomena such as disfluences (hesitations), repetitions, restarts, and corrections
  - Gender, age, emotional and environmental variations further complicate ASR
  - No punctuation marks (delimiters) or/and structural information cues exist in speech

# Areas for Speech Processing

- Production, Perception, and Modeling of Speech (phonetics and phonology)
- Signal Processing for Speech
- Speech Coding
- Speech Synthesis (Text-to-Speech)
- Speech Recognition (Speech-to-Text) and Understanding
- Speaker Recognition
- Language Recognition
- Speech Enhancement
- ....

C.f. Jacob Benesty (ed.), M. Mohan Sondhi (ed.), Yiteng Huang (ed.), Springer Handbook of Speech Processing, Springer, 2007



# Progress of Technology (1/6)

- US. National Institute of Standards and Technology (NIST)



The screenshot shows the NIST Information Technology Laboratory website. At the top, there is a navigation bar with the NIST logo and links for "NIST Time", "NIST Home", "About NIST", "Contact Us", and "A-Z Site Index". Below this is a blue banner with the text "Information Technology Laboratory". Underneath the banner is a secondary navigation bar with links for "About ITL", "Publications", "Topic/Subject Areas", "Products/Services", "News/Multimedia", and "Programs/Projects". A breadcrumb trail reads: "NIST Home > ITL > Information Access Division > Multimodal Information Group > Benchmark Tests". The main content area is divided into two sections: "Ongoing Benchmark Tests" and "Past Benchmark Tests".

**Ongoing Benchmark Tests**

- GALE Translation (2006 - present)
- Language Recognition (1996 - present)
- Machine Translation (2001 - present)
- Metrics for Machine Translation (2008 - present)
- Rich Transcription (2003 - present)
- Speaker Recognition (1996 - present)
- TRECvid Event Detection (2008-present)
- MADCAT (2008-present)
- Multiple Camera Single Person Tracking (2009-present)

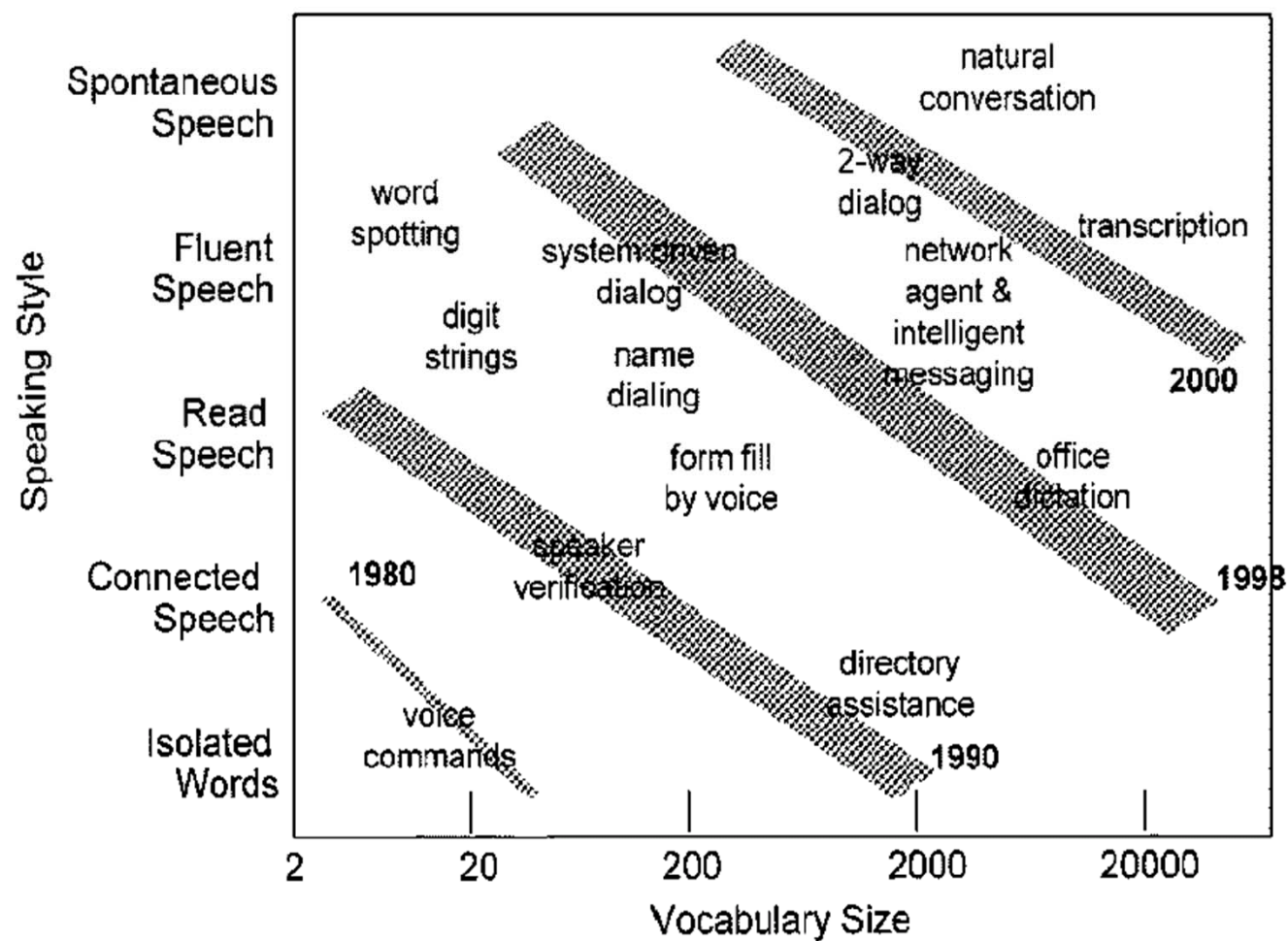
**Past Benchmark Tests**

- CLEAR (2006 - 2007)
- Spoken Term Detection (2006)
- Broadcast News Recognition (1996 - 1999)
- Conversational Telephone Recognition (1997 - 2001)
- Spoken Document Retrieval (1997 - 2000)
- Topic Detection and Tracking (1998 - 2004)
- Automatic Content Extraction (1999 - 2008)

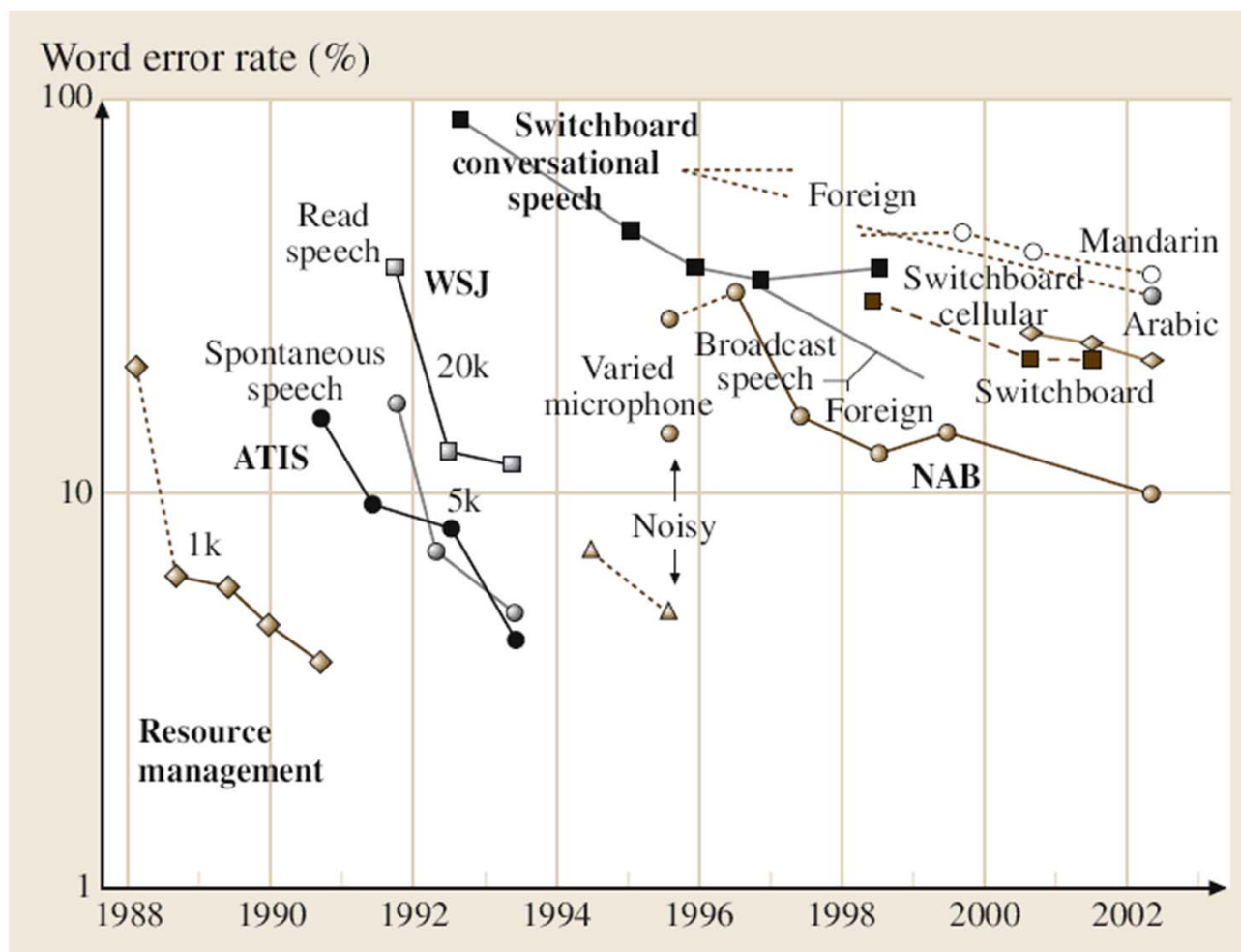
<http://www.nist.gov/itl/iad/mig/bmt.cfm>

# Progress of Technology (2/6)

- Generic Application Areas (vocabulary vs. speaking style)



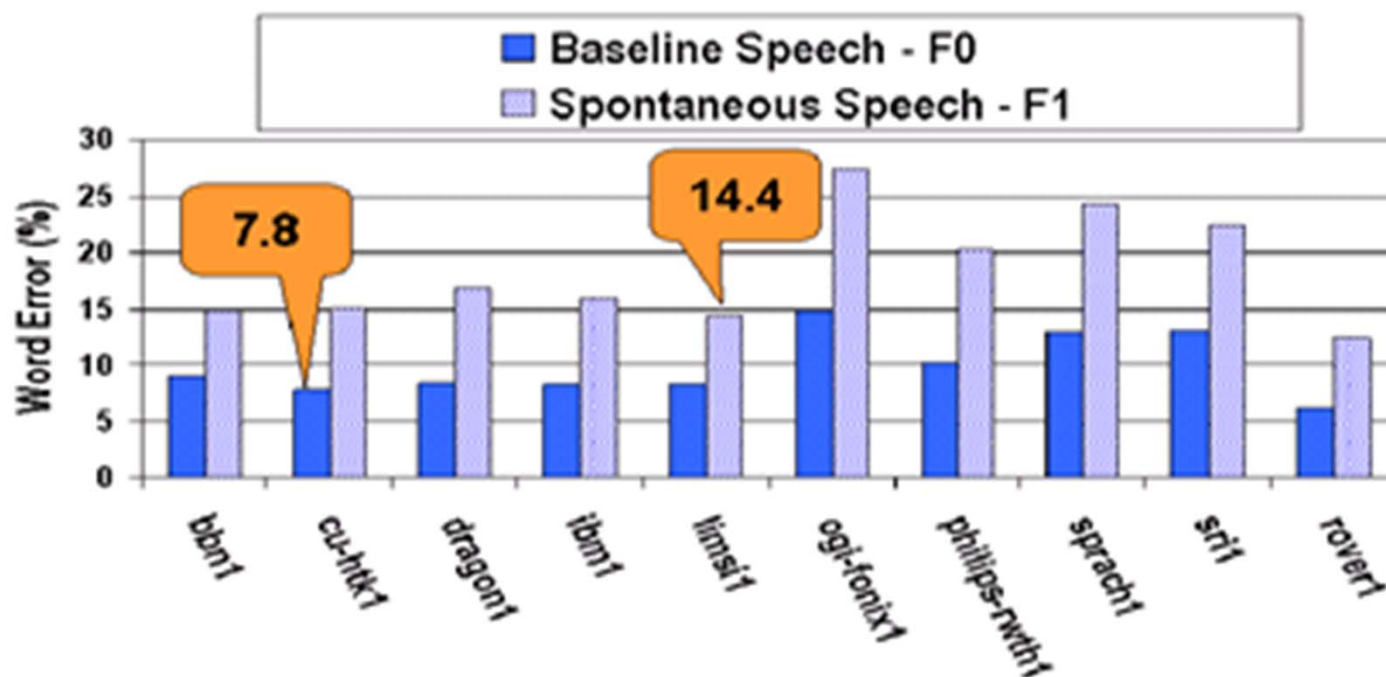
# Progress of Technology (3/6)



L. Rabiner, B.-H. Juang, "Historical Perspective of the Field of ASR/NLU" Chapter 26 in the book "Springer Handbook of Speech Processing"

# Progress of Technology (4/6)

- Benchmarks of ASR performance: Broadcast News Speech



FO: anchor speakers

F1: field reports and interviewees

## Progress of Technology (5/6)

- Benchmarks of ASR performance: Conversational Speech

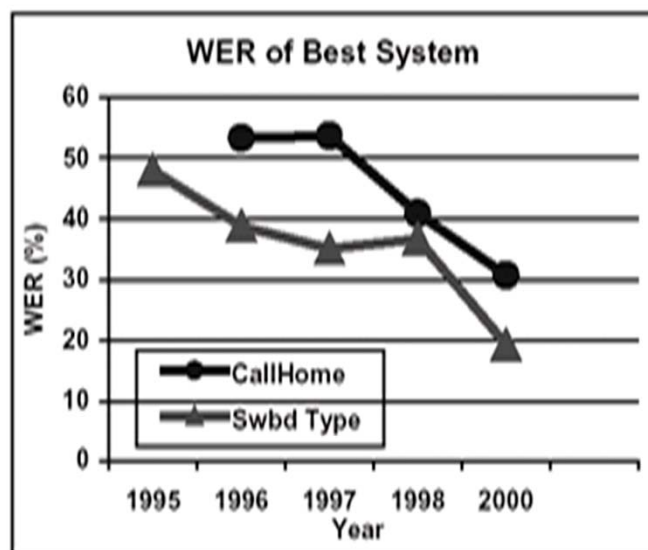


Figure 4 History of lowest word error rates (WER) obtained in NIST conversational speech evaluations on Switchboard and CallHome type conversations in English [26].

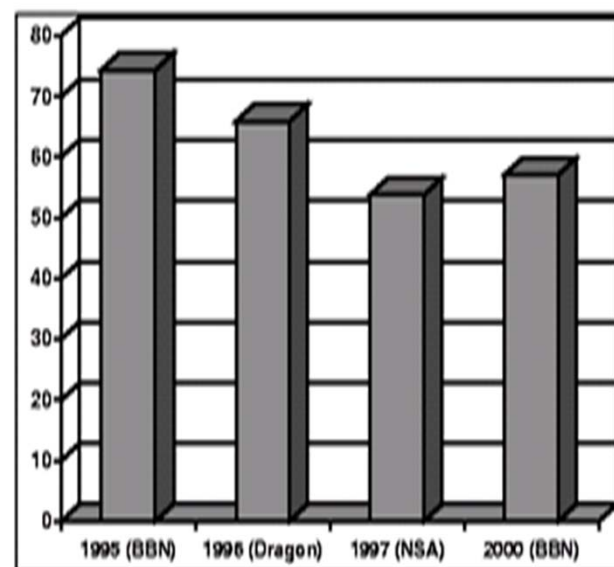


Figure 5 Chinese Character error rates of the best performing evaluation system in NIST Mandarin conversational speech evaluations 1995-2000 [26].



# Progress of Technology (6/6)

- Mandarin Conversational Speech (2003 Evaluation)
  - Acoustic/Training Test Data:
    - training data: 34.9 hours, 379 sides, from LDC CallHome (22.4hrs) and CallFriend (12.5hrs), 451K Words (+7K English word), 628K Characters
    - development data: dev02 1.94 hours from CallFriend

		CER (%)	
		dev02	eval03
P1	trans for VTLN	55.1	54.7
P2	trans for MLLR	50.8	51.3
P3	lat gen (bg)	49.3	50.5
	tgintcat rescore	48.9	49.8
P4	lat MLLR	48.6	49.5
CN	P4	47.9	48.6

%CER on dev02 and eval03 for all stages of 2003 system

– Adopted from

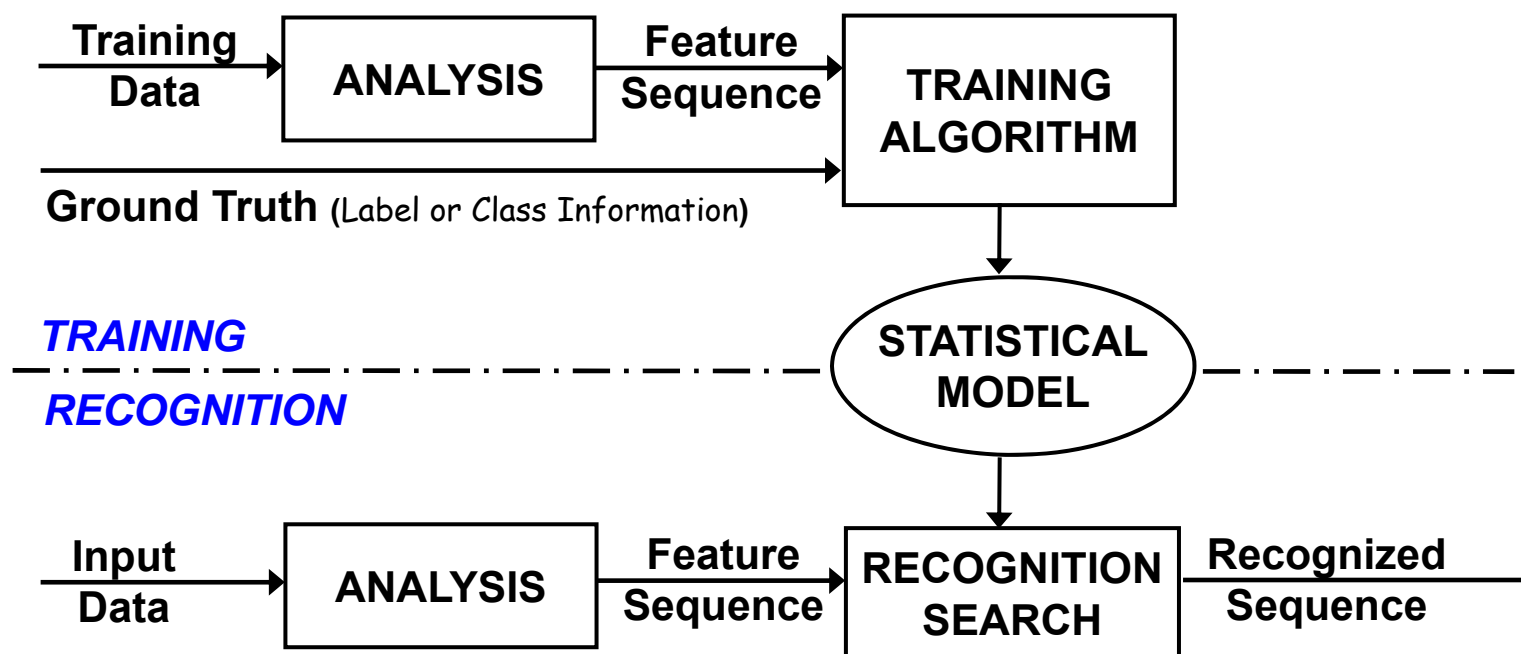


Cambridge University  
Engineering Department

Rich Transcription Workshop 2003

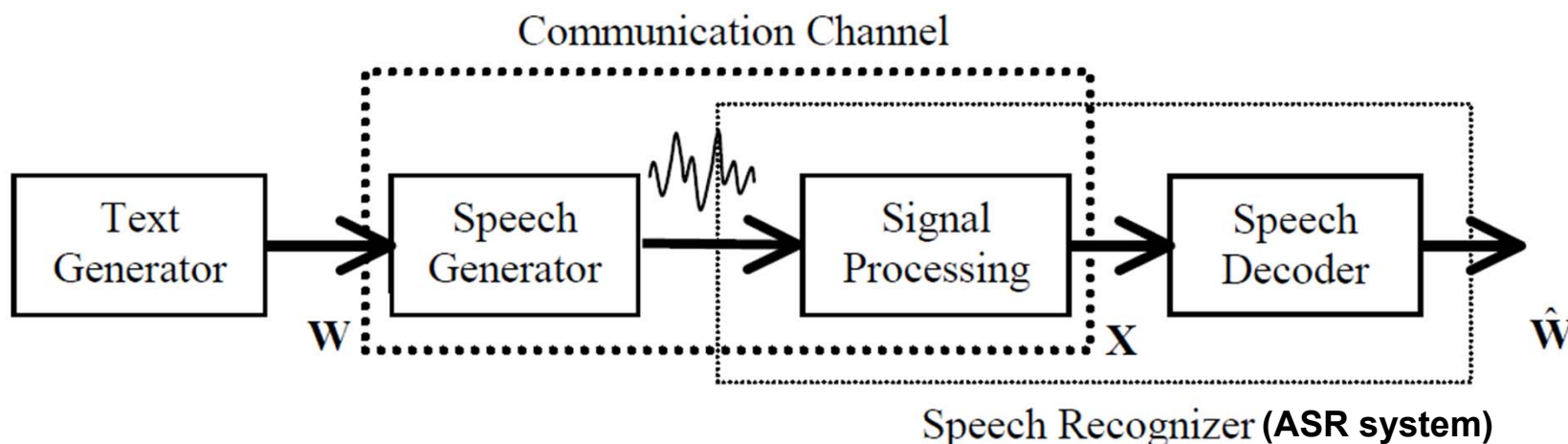
# Statistical Modeling Paradigm

- Most approaches to speech and language processing generally follow the statistical modeling paradigm



- Data-driven approaches: automatically extract “knowledge” from the data
- It would be better to pair data-driven approaches with rule-based ones

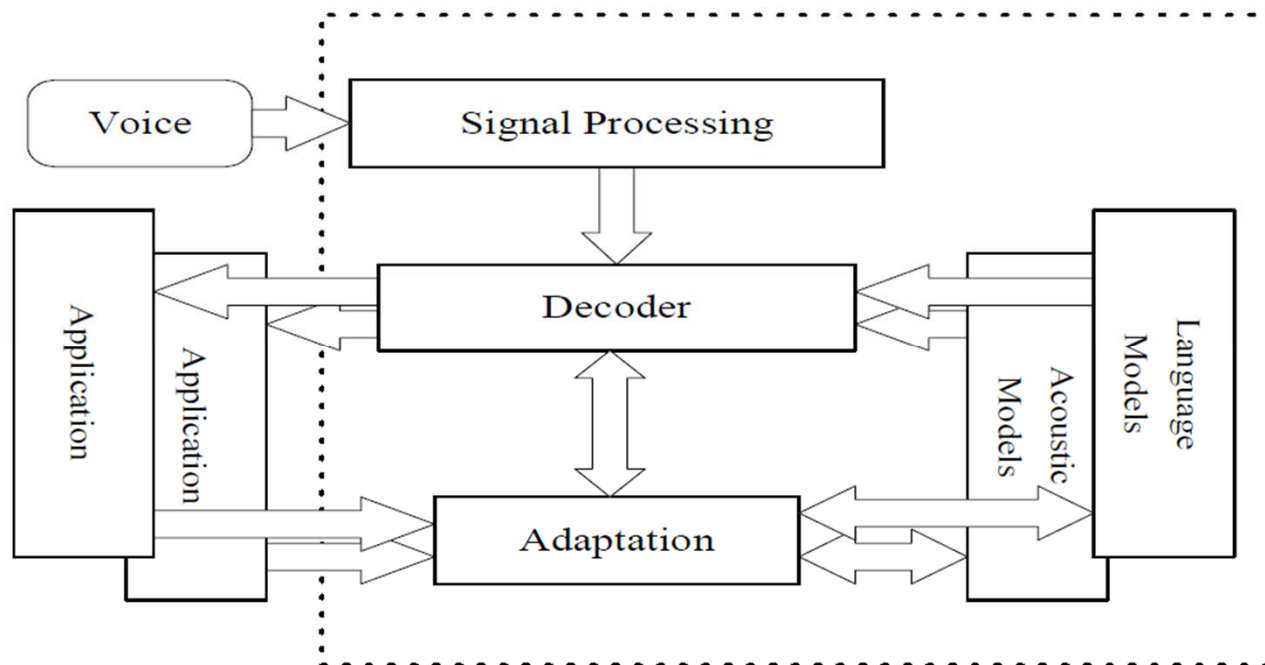
# A Source-Channel Model for ASR



- Communication channel consists of speaker's vocal apparatus to produce speech (the waveform) and the signal processing component of the speech recognizer
- The speech decoder aims to decode the acoustic signal  $X$  into a word sequence  $\hat{W}$  (Hopefully,  $\hat{W} \approx W$ .)

Uncertainties to be contended with: unknown words, grammatical variation, noise interference, acoustic variation, to name a few

# Basic Architecture of ASR System



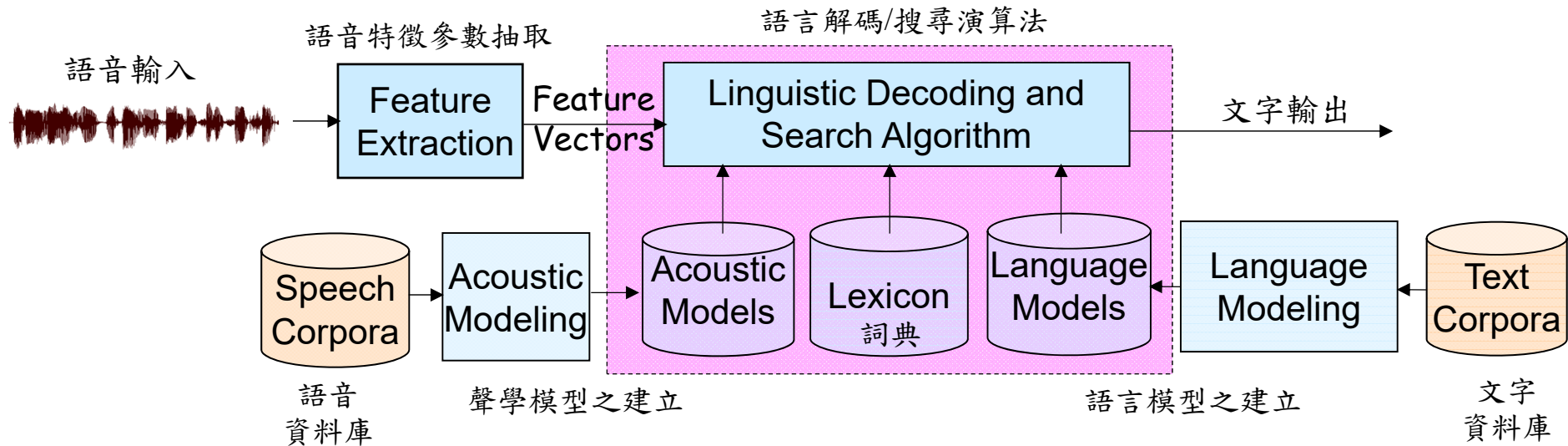
- **Signal processing:** extract salient features for the decoder
- **Decoder:** use both acoustic and language models to generate the “best” word sequence in response to the input voice
- **Adaptation:** modify either acoustic or language models so that improved performance can be obtained

# ASR: Applications

- E.g., Transcription of Broadcast News Speech



# ASR: A Bit of Terminology



可能詞句      語音輸入

$$\hat{\mathbf{W}} = \arg \max_{\mathbf{W}} P(\mathbf{W} | \mathbf{X}) \quad \text{Bayes Decision Theory}$$

$$= \arg \max_{\mathbf{W}} \frac{p(\mathbf{X} | \mathbf{W}) P(\mathbf{W})}{P(\mathbf{X})} \quad \text{Bayes Rule}$$

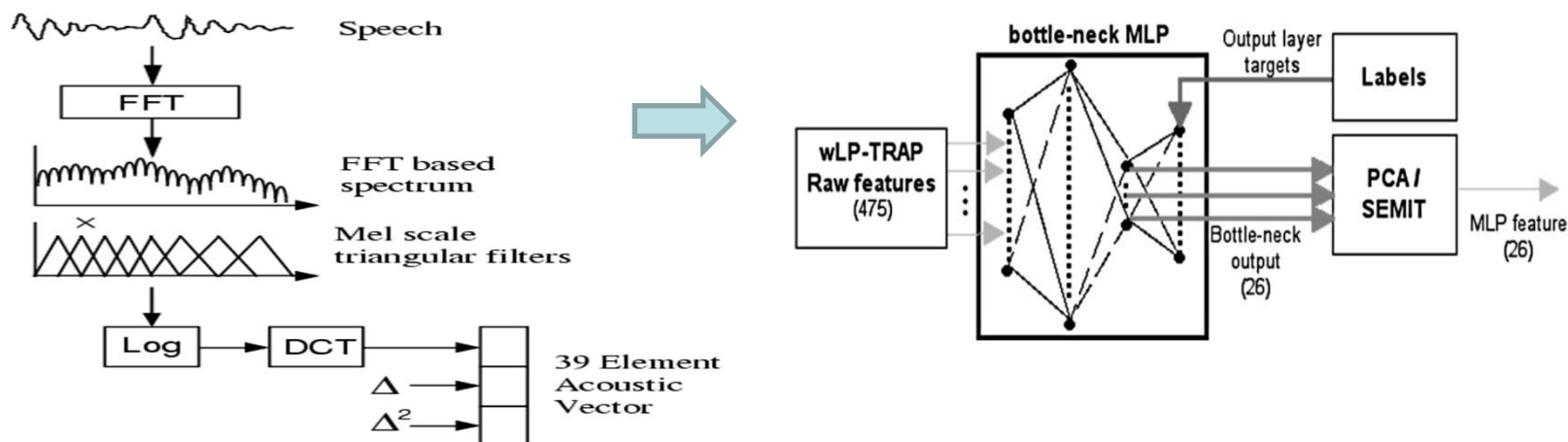
$$= \arg \max_{\mathbf{W}} p(\mathbf{X} | \mathbf{W}) P(\mathbf{W})$$

Decoding

Acoustic Modeling      Language Modeling

# Speech Feature Extraction

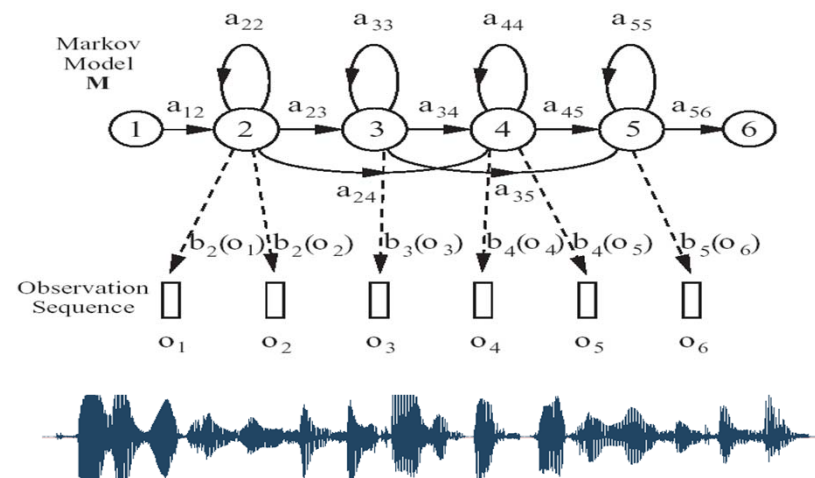
- The raw speech waveform is passed through feature extraction to generate relatively compact feature vectors at a frame rate of around 100 Hz
  - Parameterization: an acoustic speech feature is a simple compact representation of speech and can be modeled by cepstral features such as the Mel-frequency cepstral coefficient (MFCC)



raw (perception-driven) features vs. discriminant (posterior) features

# ASR: Acoustic Modeling

- Construct **a set of statistical models** representing various sounds (or phonetic units) of the language
  - Approaches based on Hidden Markov Models (HMMs) dominate the area of speech recognition
  - HMMs are based on rigorous mathematical theory built on several decades of mathematical results developed in other fields
  - HMMs are constructed by the process of training on a large corpus of real speech data






# ASR: Language Modeling

- Constrain the acoustic analysis, guide the search through multiple candidate word strings, and quantify the acceptability of the final word string output from a speech recognizer

$$W = w_1 w_2 \dots w_L \implies P(W) = ?$$

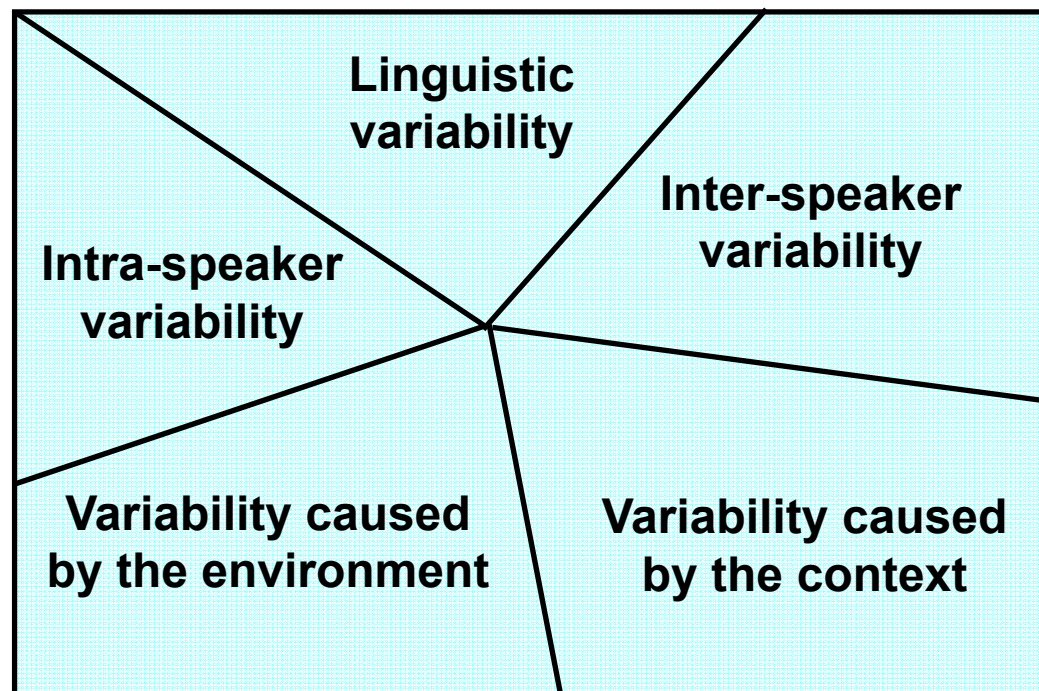
- The  $n$ -gram language model that follows a statistical modeling paradigm is the most prominently-used in ASR

**bigram modeling**


$$P(w_1 w_2 \dots w_L) = P(w_1) P(w_2 | w_1) P(w_3 | w_1 w_2) \dots P(w_L | w_1 w_2 \dots w_{L-1})$$
$$P(w_1 w_2 \dots w_L) = P(w_1) P(w_2 | w_1) P(w_3 | w_2) \dots P(w_L | w_{L-1})$$

# Difficulties: Speech Variability

**Pronunciation  
Variation**



**Speaker-independency  
Speaker-adaptation  
Speaker-dependency**

**Robustness  
Enhancement**

**Context-Dependent  
Acoustic Modeling**

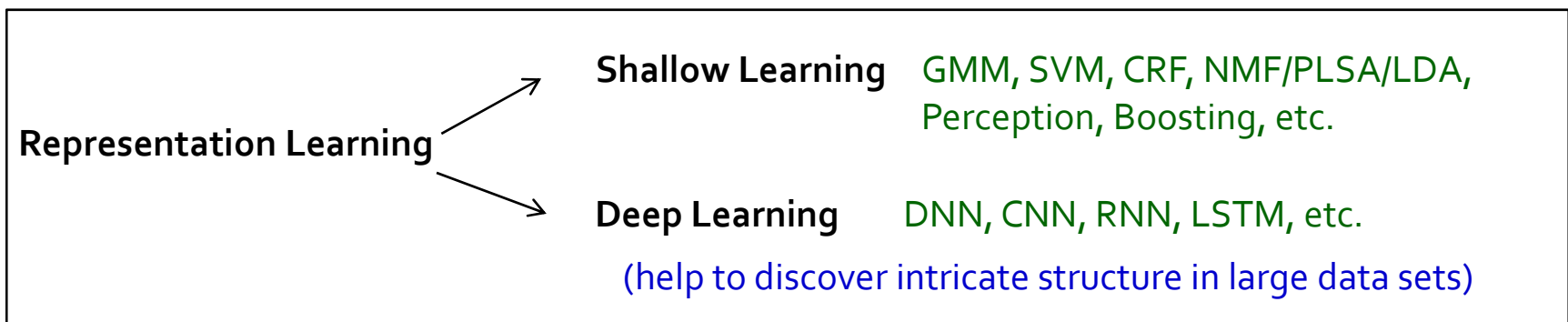
# What is Deep Learning?



## Deep learning

From Wikipedia, the free encyclopedia

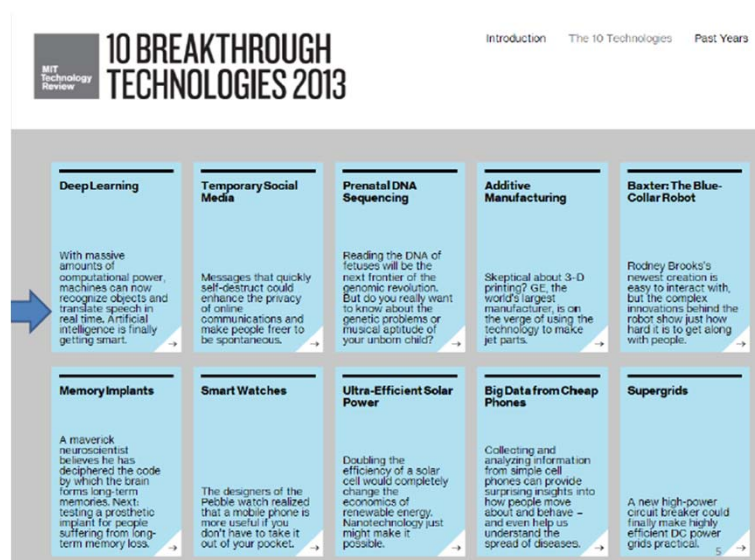
**Deep learning** (*deep machine learning*, or *deep structured learning*, or *hierarchical learning*, or sometimes *DL*) is a branch of [machine learning](#) based on a set of [algorithms](#) that attempt to model high-level abstractions in data by using multiple processing layers with complex structures or otherwise, composed of multiple non-[linear transformations](#).<sup>[1](p198)[2][3][4][5]</sup>



Deeper is better? vs. Simple is elegant?

# Deep Learning and its Applications to ASR (1/5)

- **Deep Learning** is concerned with learning multiple levels of representation and abstraction that help to make sense of data such as images, sound, and text
- By virtue of **Deep Learning**
  - Our computers can learn and grow on their own
  - Our computers are able to understand complex, massive amount of data (**deep learning is the holy grail of big data?**)



# Deep Learning and its Applications to ASR (2/5)

**MIT  
Technology  
Review**

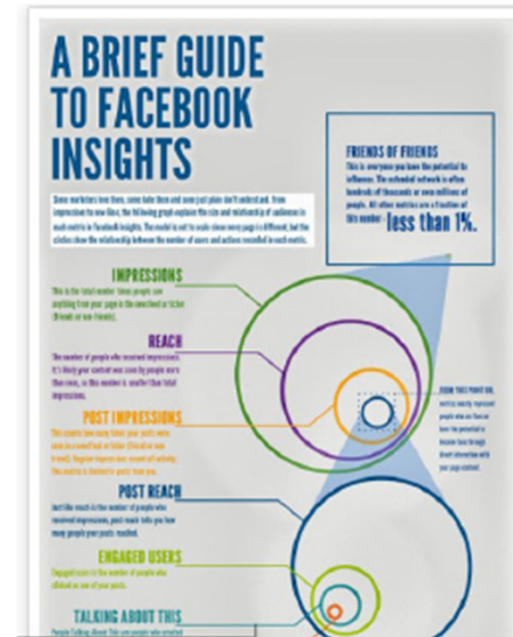
## Facebook Launches Advanced AI Effort to Find Meaning in Your Posts

September 20, 2013

A technique called deep learning could help Facebook understand its users and their data better.

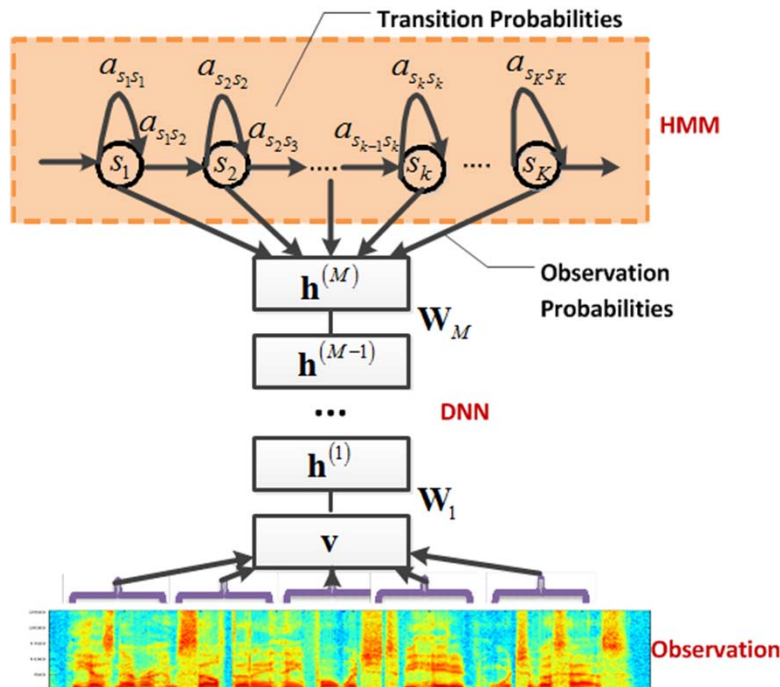
By Tom Simonite on September 20, 2013

..... Facebook's foray into deep learning sees it following its competitors **Google and Microsoft**, which have used the approach to impressive effect in the past year. Google has hired and acquired leading talent in the field (see "[10 Breakthrough Technologies 2013: Deep Learning](#)"), and last year created software that taught itself to recognize cats and other objects by reviewing stills from YouTube videos. The underlying deep learning technology was later used to slash the error rate of Google's voice recognition services (see "[Google's Virtual Brain Goes to Work](#)").... **Researchers at Microsoft have used deep learning** to build a system that translates speech from English to Mandarin Chinese in real time (see "[Microsoft Brings Star Trek's Voice Translator to Life](#)"). Chinese Web giant Baidu also recently established a Silicon Valley research lab to work on deep learning.



# Deep Learning and its Applications to ASR (3/5)

- **Deep Learning** is the cutting edge!
  - Use deep neural network hidden Markov model (DNN-HMM) hybrid architecture to train DNN to produce a distribution over senones (tied triphone states) as its output



deeper layers,  
longer features &  
wider temporal contexts

$$b_{s_i}(\mathbf{o}) = p(\mathbf{o} | s_i) = \frac{P_{\text{DNN}}(s_i | \mathbf{o}) p(\mathbf{o})}{P_{\text{ML}}(s_i)} \propto \frac{P_{\text{DNN}}(s_i | \mathbf{o})}{P_{\text{ML}}(s_i)}$$

$$P_{\text{DNN}}(s_i | \mathbf{o}) = v_i^L = \text{softmax}_i(\mathbf{z}^L) = \frac{e^{z_i^L}}{\sum_j e^{z_j^L}}$$

$$\mathbf{v}^\ell = f(\mathbf{z}^\ell) = f(\mathbf{W}^\ell \mathbf{v}^{\ell-1} + \mathbf{b}^\ell), \text{ for } 0 < \ell < L$$

$f(\cdot)$ : sigmoid, hyperbolic, or rectified linear unit (ReLU) functions

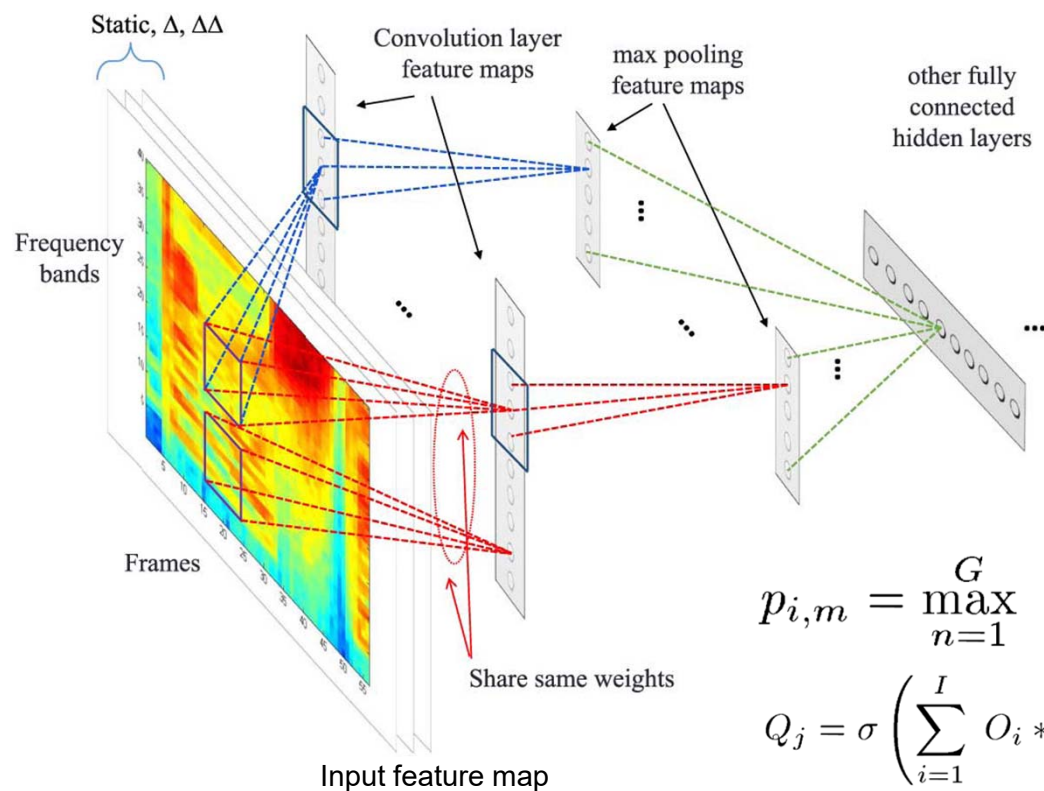
Model parameters of DNN can be estimated with **the error back-propagation algorithm and stochastic gradient descent (SGD)**.

G. Dahl, D. Yu, L. Deng, and A. Acero, "Context-Dependent Pre-trained Deep Neural Networks for Large Vocabulary Speech Recognition," IEEE Transactions on Audio, Speech, and Language Processing, Vol. 20, No. 1. pp. 30-42, 2012

# Deep Learning and its Applications to ASR (4/5)

- **CNN-HMM**

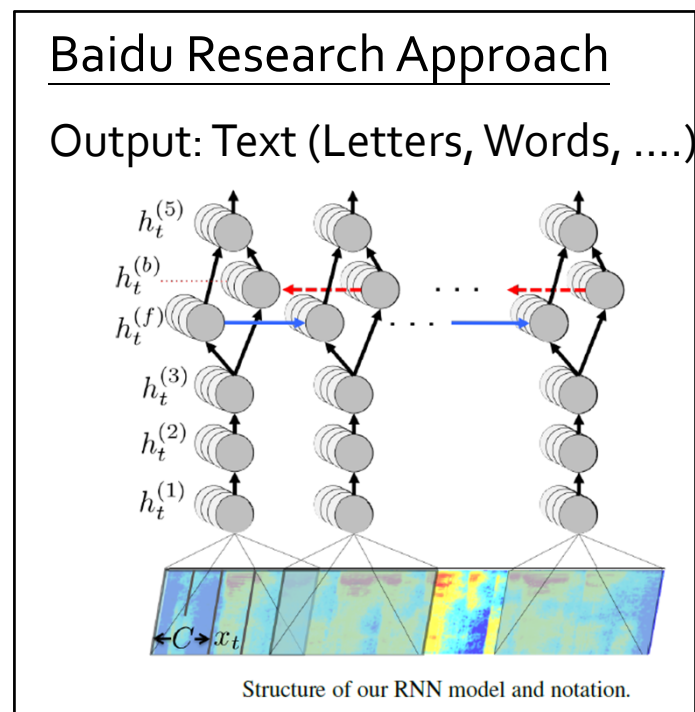
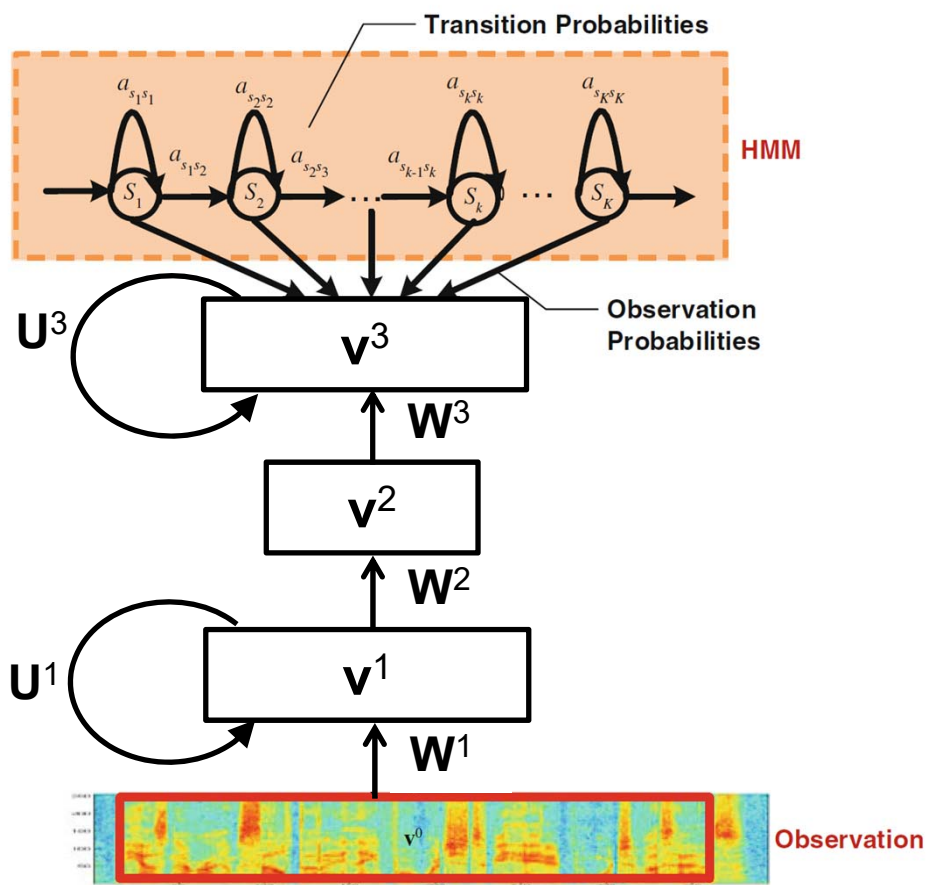
- CNN: Convolutional Neural Networks



$$p_{i,m} = \max_{n=1}^G q_{i,(m-1) \times s + n}$$
$$Q_j = \sigma \left( \sum_{i=1}^I O_i * \mathbf{w}_{i,j} \right) \quad (j = 1, \dots, J)$$

# Deep Learning and its Applications to ASR (5/5)

- Recurrent Neural Networks (**RNN-HMM**)



A. Hannun et al. (Lead by Andrew Ng), "Deep Speech: Scaling up end-to-end speech recognition," arXiv:1412.5567v2, December 2014.



# Example: Automatic Meeting Transcription

## Manual Transcripts

A: 那會在二 a 那個那叫什麼二 b 啊二 a  
 A: vip vip room  
 B: 欸  
 A: 就是大家開 all hands meeting 那裡  
 C: 錄音的話就只能用八爪魚喔  
 A: 錄音就對啊那場就反正錄下來就好了對  
 A: 好一開始  
 D: 請問一下  
 D: 上次二 a. 的時候那個圓方不是有來教我們  
 怎麼用八爪魚錄音所以那個測試設定都  
 沒有動  
 D: 就直接麥克風可以把聲音收進來  
 A: 圓圓形會議對啊圓形會議是這樣  
 D: 好好  
 A: 可是我們這一次不是在圓形我們這次是  
 在呃 vip  
 A: 就是董事長開會的地方

## Automatic Transcripts

A: 那會在二 a. h 那個資料怎麼二的啊把二  
 a.  
 A: 七 vip 喔 vip vip room  
 B: 嘿  
 A: 可是打開過 hand meeting 那裡  
 C: 錄音的話是怎麼用滑動語料  
 A: 錄音就對啊那一場就反正錄下就好了  
 A: 好一開始了  
 D: 請問一下  
 D: 上是二月的時候那個員工不是來教我們  
 怎麼跟八爪魚錄音最那個測試設定檔秒  
 鐘  
 D: 就支麥克風可以把聲音投進來  
 A: 每圓形會議對啊圓形會議室這樣  
 D: 好  
 A: 可是我們這次不是在圓形我  
 edge vip  
 A: 是董事會開會的地方



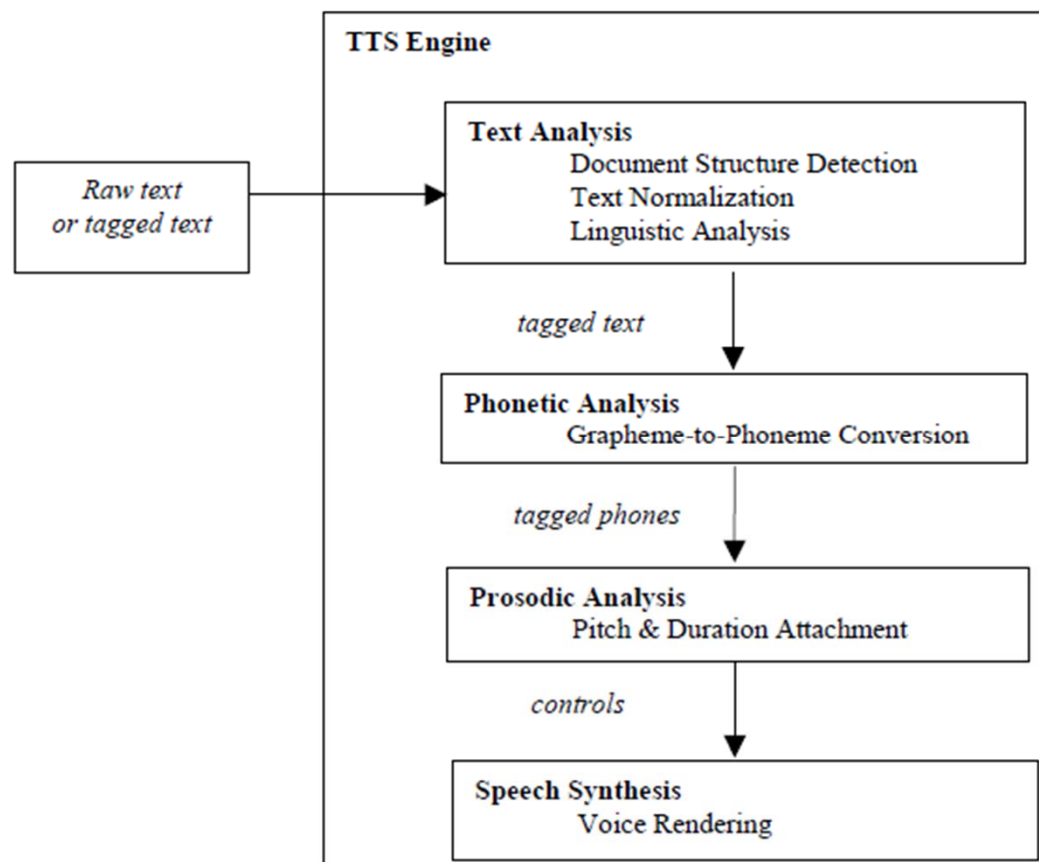
Chinese Character Error Rate (CER%)



GMM-HMM	DNN-HMM	CNN-HMM	RNN(LSTM)-HMM	RNN(BLSTM)-HMM
49.74	40.95	35.41	51.63	41.87

# Text to Speech (1/2)

- Text to speech (TTS) can be viewed as ASR in reverse



- We are now able to general high-quality TTS systems, although the quality is inferior to human speech for general-purpose applications

# Text to Speech (2/2)

- Example 1

- 青少年在成長的過程中，非常需要角色模範的引導、族群的認同及自我的肯定，所以我一直在找這方面的好書來幫助孩子。

- Original Speech: 

- Synthesized Speech: 

- Example 2

- 新北市市長朱立倫昨天邀台北市市長柯文哲參加新北市天燈節第三場活動，兩人在廿呎高的剪紙天燈上寫下「雙北合作」「神采飛羊」，柯則寫下「天佑台灣」，大小天燈齊放升空，照亮平溪夜空。

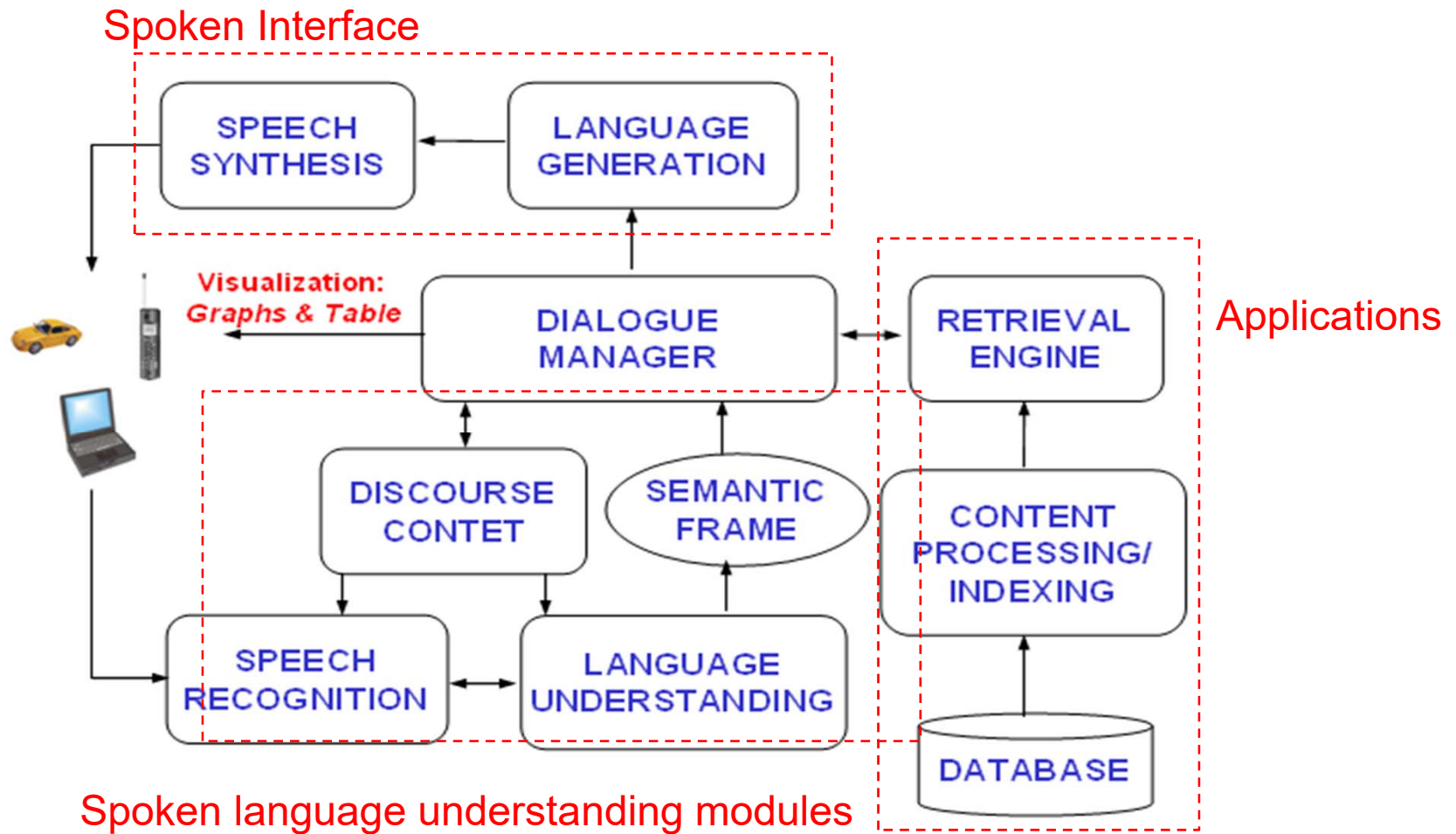
- Synthesized Speech: 

# Spoken Dialogue: CMU's Systems

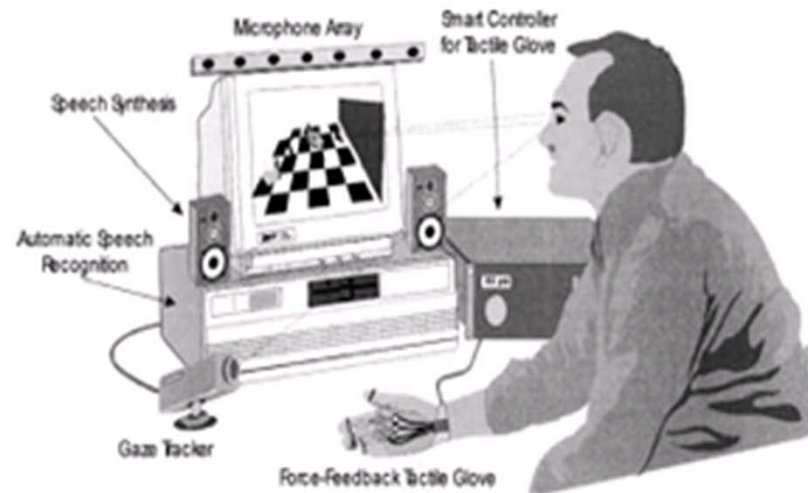
- Spoken language is attractive because it is the most natural, convenient and inexpensive means of exchanging information for humans
- In mobilizing situations, using keystrokes and mouse clicks could be impractical for rapid information access through small handheld devices like PDAs, cellular phones, etc.



# Spoken Dialogue: Basic System Architecture



# Spoken Dialogue: Multimodality of Input and Output



Experimental client workstation incorporating sight, sound, and touch modalities for human/machine communication. The eye tracker provides a gaze-controlled cursor for indicating objects in the display. The tactile force-feedback glove allows displayed objects to be grasped, “felt,” and moved. Hands-free speech recognition and synthesis provides natural conversational interaction [7].

I. Marsic, A. Medl, and J. Flanagan, Natural Communication with Information Systems. Proceedings of the IEEE, Vol. 88, No. 8, August 2000

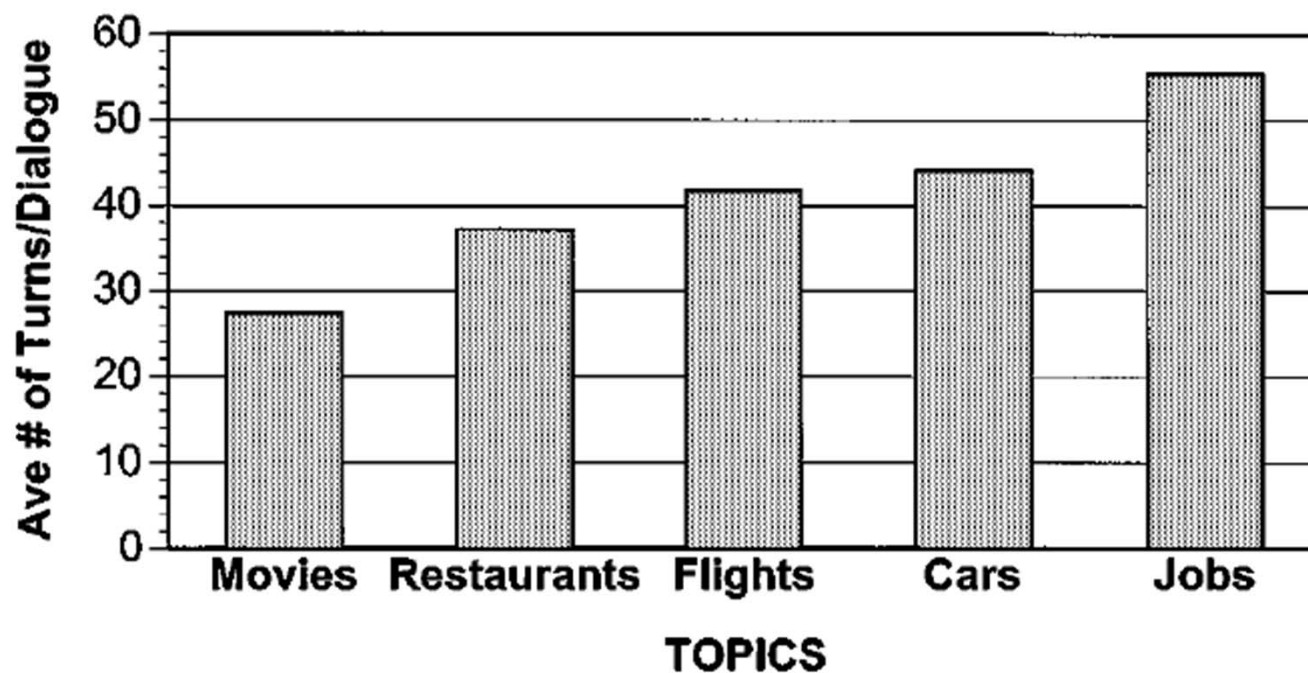
# Spoken Dialogue: Some Deployed Systems

- Complexity Analysis

Domain	Language	Vocabulary Size	Average	
			Words/Utt	Utts/Dialogue
CSELT Train Timetable Info	Italian	760	1.6	6.6
SpeechWorks Air Travel Reservation	English	1000	1.9	10.6
Philips Train Timetable Info	German	1850	2.7	7.0
CMU Movie Information	English	757	3.5	9.2
CMU Air Travel Reservation	English	2851	3.6	12.0
LIMSI Train Timetable Info	French	1800	4.4	14.6
MIT Weather Information	English	1963	5.2	5.6
MIT Air Travel Reservation	English	1100	5.3	14.1
AT&T Operator Assistance	English	4000	7.0	3.0
Air Travel Reservations (human)	English	?	8.0	27.5

# Spoken Dialogue: Some Statistics

- Topics vs. Dialogue Terms

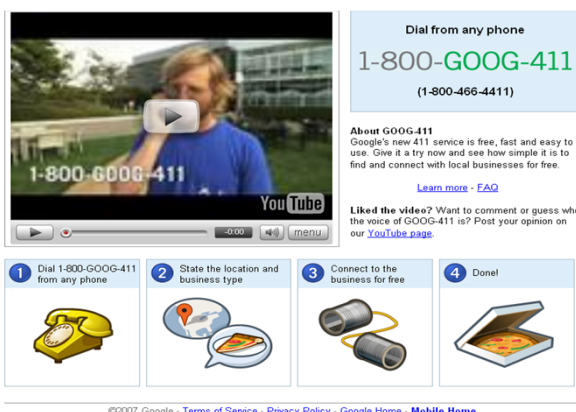




# Current Deployed Speech Retrieval and Spoken Dialogue Systems

- Google, Apple, Microsoft and Amazon's Deployed Services

Google-411:  
Finding and connecting to  
local business



Google Voice Search

<http://www.google.com/mobile/voice-search/>



Apple Siri

<http://www.apple.com/iphone/features/siri.html>



Microsoft Cortana

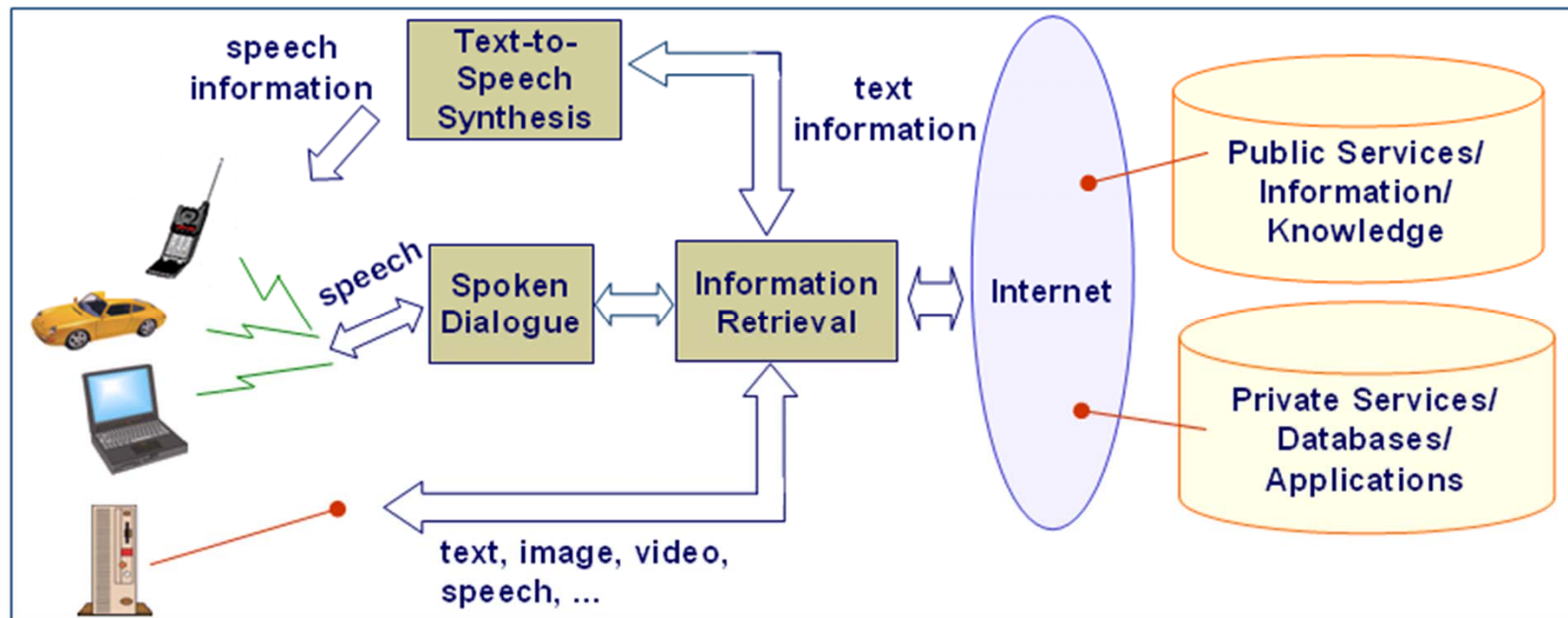
[http://zh.wikipedia.org/wiki/Microsoft\\_Cortana](http://zh.wikipedia.org/wiki/Microsoft_Cortana)



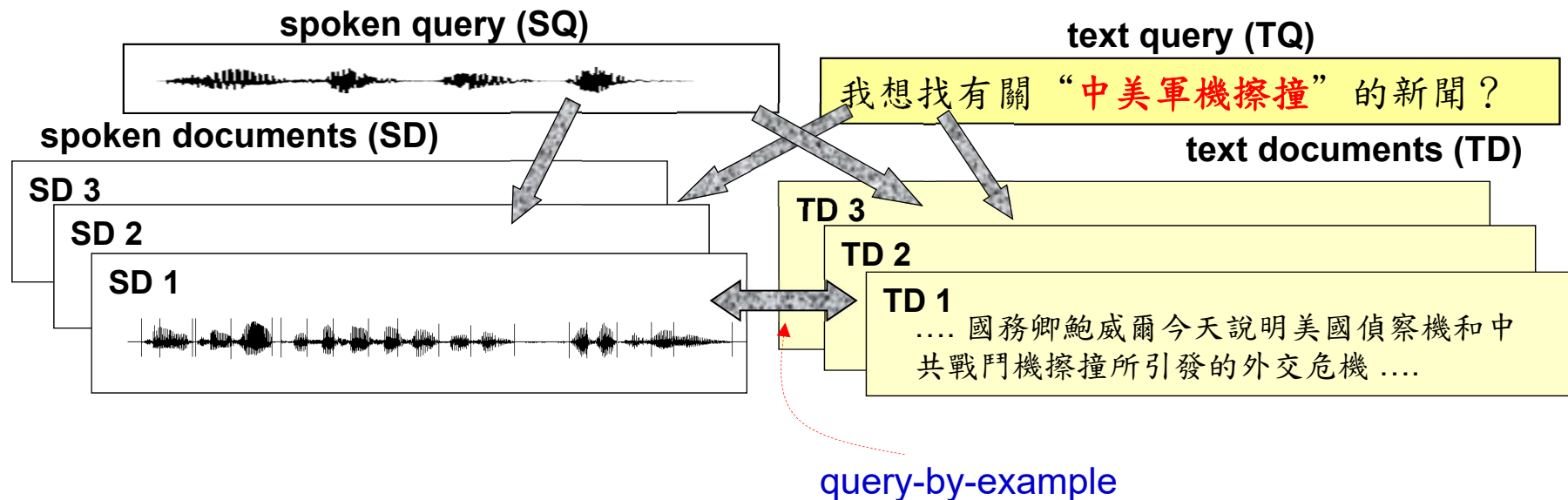
*We are witnessing the golden age of ASR!*

# Speech-based Information Retrieval (1/5)

- Task :
  - Automatically indexing a collection of spoken documents with speech recognition techniques
  - Retrieving relevant documents in response to a text/speech query



# Speech-based Information Retrieval (2/5)



- SQ/SD is the most difficult
- TQ/SD is studied most of the time
- Query-by-example
  - Attempt to retrieve relevant documents when users provide some specific query exemplars describing their information needs
  - Useful for news monitoring and tracking

# Speech-based Information Retrieval (3/5)

輸入聲音問句：“請幫我查總統府升旗典禮”

中文電視廣播新聞檢索系統 2002v1 - Berlin Chen & Lin-shan Lee

辨識1 等待輸入指令...

測靜音 收音 解開 載入新聞

3.70秒

語音辨識結果

總統府升旗典禮 ← 聲音問句的語音辨識結果

Viterbi->End Time= 100  
TotalFrame=362 1. (接受) 幫我找 8340.57 (時間) 28 100

文字檢索

語音辨識結果

FILE (Erroneous Transcription): FTV2002-004.txt

中華民國就是明年元旦總統府升旗典禮即將在下而星期二登場  
而今年首度社教有民間工商團體來舉辦  
新科立委金素梅將帶著貴為原住民亦同高唱國歌  
展現多元文化的特性有以今年的元旦升旗典禮將打破傳統方式  
經紀人龍門一千人到新竹美勞他擔任市為原住民

檢索到新聞的語音辨識結果

可以選擇同時使用音節、字、詞等三種索引特徵

QueryByExpr	檢索結果之排名
[ 1 ] FTV2002-004 3.59164e-001	
[ 2 ] N200201211200-01 1.11802e-001	
[ 3 ] N200201091200-12 1.91467e-001	
[ 4 ] N200110051200-09 1.89940e-001	
[ 5 ] N200109061200-07 1.66562e-001	
[ 6 ] T20020111200-06 6.60336e-001	
[ 7 ] N200111071200-04 1.60011e-001	
[ 8 ] N200111131200-04 1.57109e-001	
[ 9 ] T200201211200-04 1.51319e-001	
[ 10 ] N200110031200-03 1.47177e-001	
[ 11 ] N200201171200-11 1.44006e-001	
[ 12 ] N200105071400-02 1.41382e-001	
[ 13 ] T200106191000-02 1.39268e-001	
[ 14 ] N200110291200-01 1.38799e-001	
[ 15 ] N200104301230-05 1.36488e-001	
[ 16 ] N200109051200-05 1.33595e-001	
[ 17 ] N200109141200-18 1.33158e-001	
[ 18 ] N200105142000-05 1.32321e-001	
[ 19 ] FTV2002-064 1.32147e-001	
[ 20 ] N200201181200-11 1.31223e-001	
[ 21 ]	

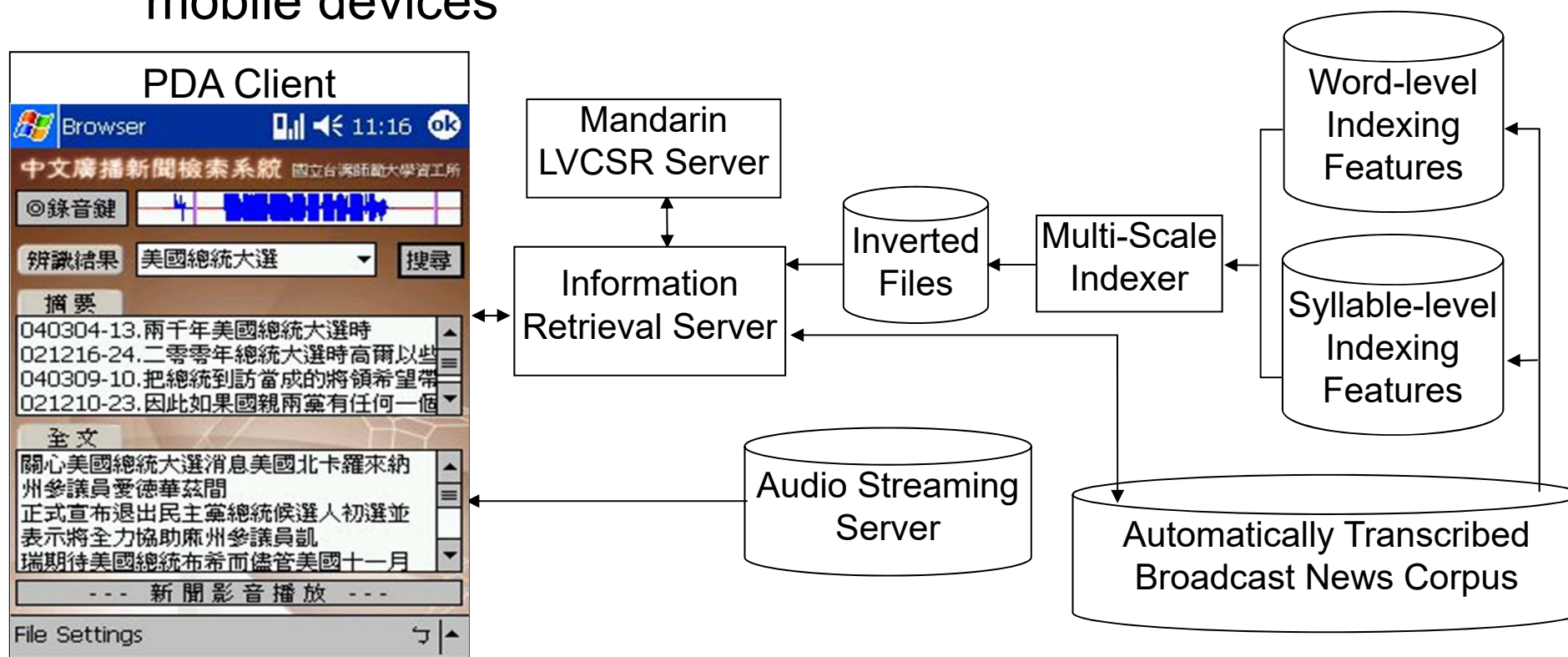
檢索到新聞的影音

中文語音資訊檢索雛形展示系統。

C.f. B. Chen, H.M. Wang, Lin-shan Lee, "Discriminating capabilities of syllable-based features and approaches of utilizing them for voice retrieval of speech information in Mandarin Chinese", IEEE Transactions on Speech and Audio Processing, Vol. 10, No. 5, pp. 303-314, July 2002.

# Speech-based Information Retrieval (4/5)

- Spoken queries retrieving text news documents via mobile devices

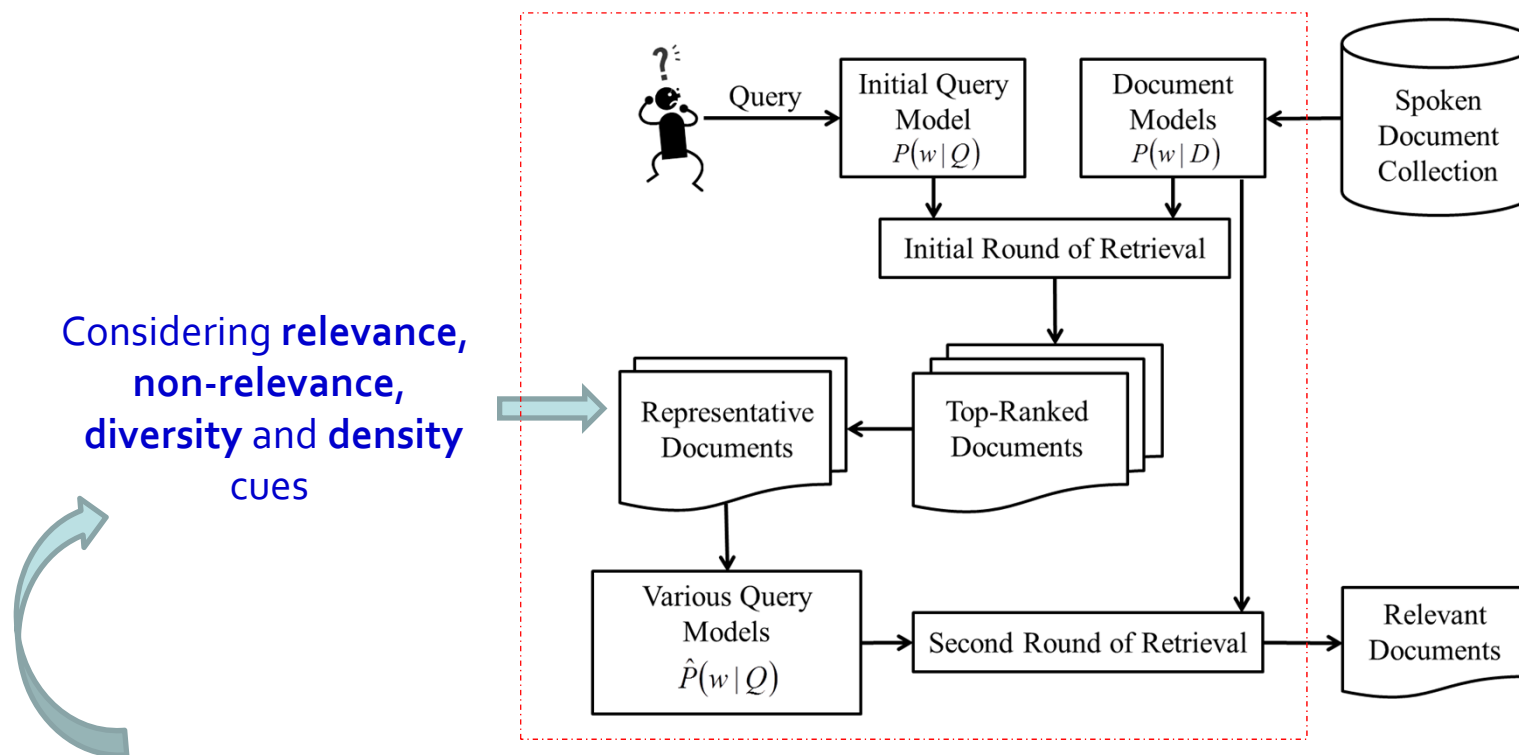


C.f. B. Chen, Y..T. Chen, C.H. Chang, H.B. Chen, "Speech Retrieval of Mandarin Broadcast News via Mobile Devices," Interspeech2005

Chang, E., Seide, F., Meng, H., Chen, Z., Shi, Y., And Li, Y. C. 2002. A system for spoken query information retrieval on mobile devices. IEEE Trans. on Speech and Audio Processing 10, 8 (2002), 531-541.

# Speech-based Information Retrieval (5/5)

- Query modeling for information retrieval

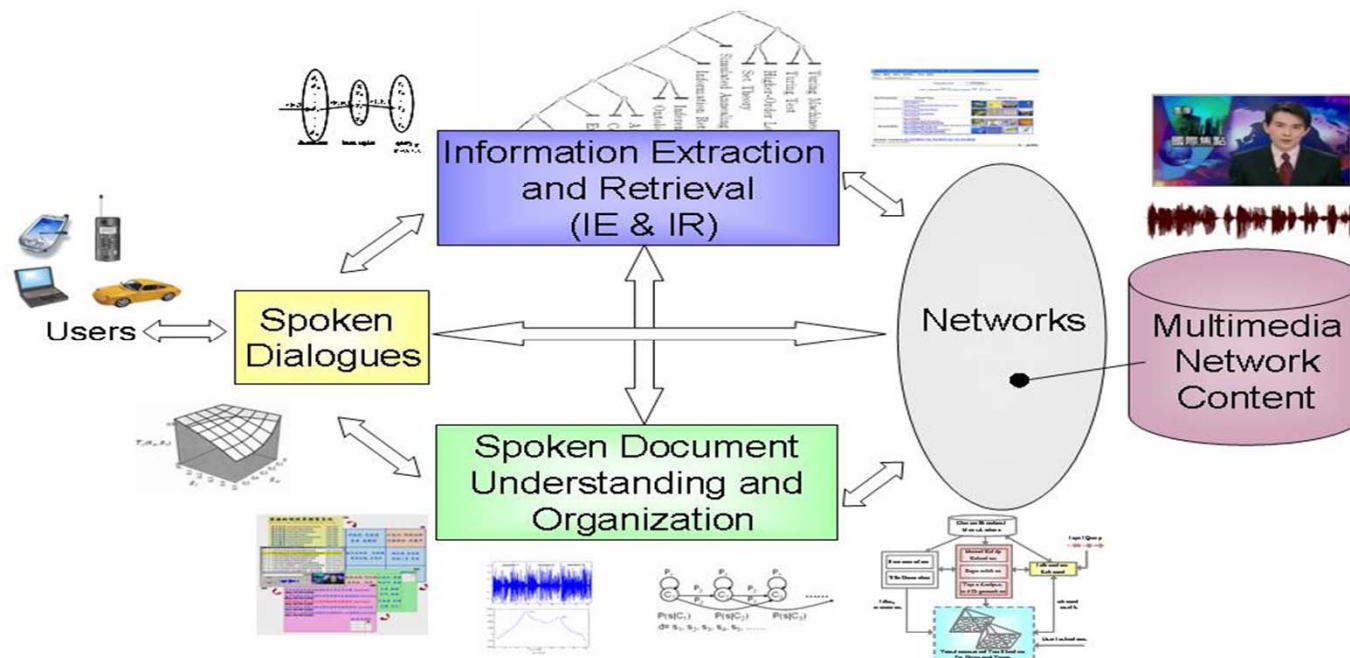


$$D^* = \arg \max_{D \in \mathbf{D}_{\text{Top}} - \mathbf{D}_P} \left[ (1 - \alpha - \beta - \gamma) \cdot M_{\text{Rel}}(Q, D) + \alpha \cdot M_{\text{NR}}(Q, D) + \beta \cdot M_{\text{Diversity}}(D) + \gamma \cdot M_{\text{Density}}(D) \right]$$

C.f. B. Chen, K.-Y. Chen, P.-N. Chen, Y.-W. Chen, "Spoken document retrieval with unsupervised query modeling techniques," IEEE Transactions on Audio, Speech and Language Processing, Vol. 20, No. 9, pp. 2602-2612, 2012

# Spoken Document Organization and Understanding (1/2)

- Problems
  - The content of multimedia documents very often described by the associated speech information
  - Unlike text documents with paragraphs/titles easy to look through at a glance, multimedia/spoken documents are unstructured and difficult to retrieve/browse



C.f. L.S. Lee and B. Chen, "Spoken document understanding and organization," IEEE Signal Processing Magazine, vol. 22, no. 5, pp. 42-60, Sept. 2005

# Spoken Document Organization and Understanding (2/2)

- Speech Summarization

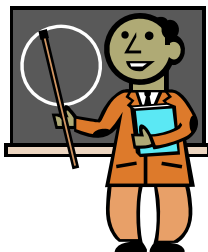
conversations



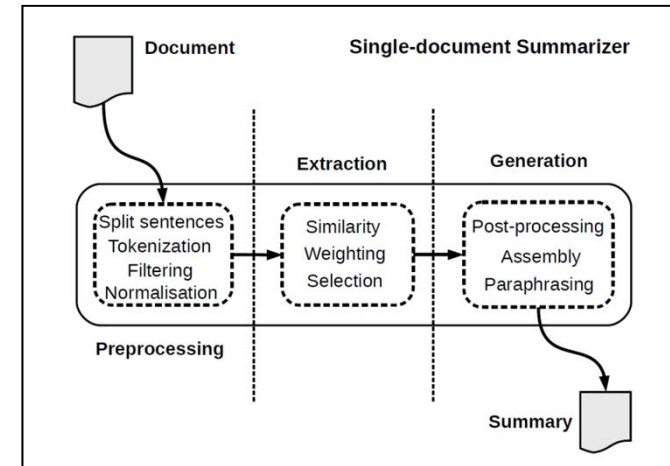
meetings



lectures



broadcast  
and TV news



distilling

important information  
*abstractive vs. extractive*  
*generic vs. query-oriented*  
*single- vs. multi-documents*

C.f. Y. Liu and D. Hakkani-Tür, "Speech summarization," Chapter 13 in Spoken Language Understanding: Systems for Extracting Semantic Information from Speech, G. Tur and Renato D. Mori (eds.), Wiley, 2011.



# Speech-to-Speech Translation (1/2)

- Multilingual interactive speech translation
  - Aim at the achievement of a communication system for precise recognition and translation of spoken utterances for several conversational topics and environments by using human language knowledge synthetically (adopted form ATR-SLT )



ATR-SLT



IBM Mastor Project

# Speech-to-Speech Translation (2/2)

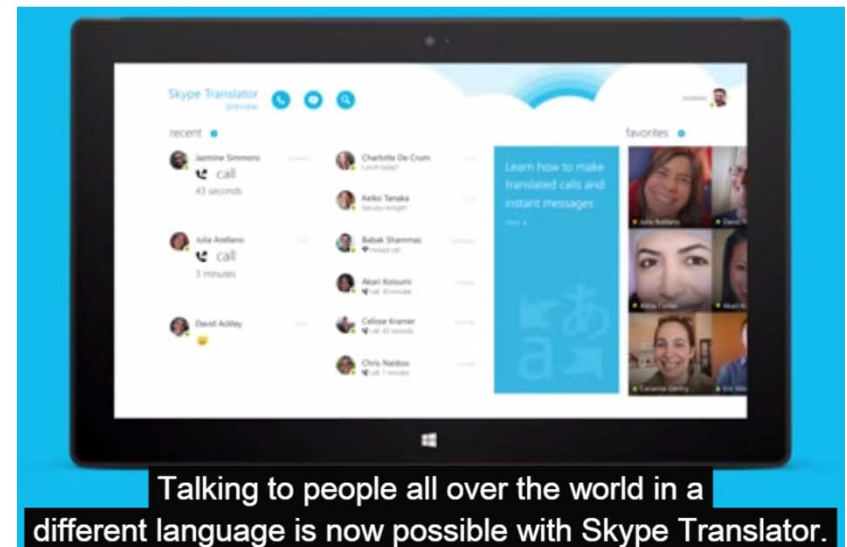
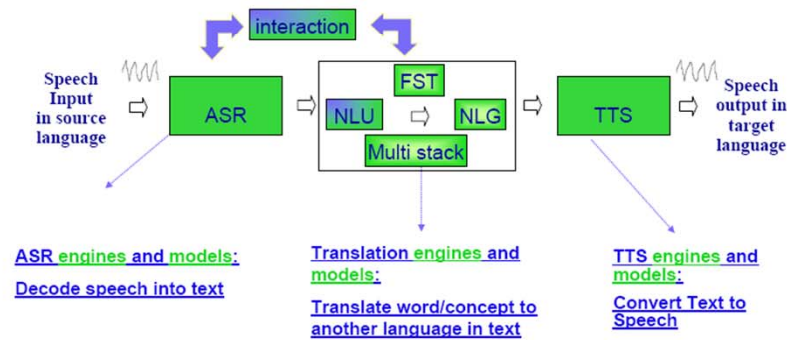
## Handheld System



## Laptop systems - hands-free, eyes-free function



## IBM Advanced Speech-to-Speech Translation Techniques

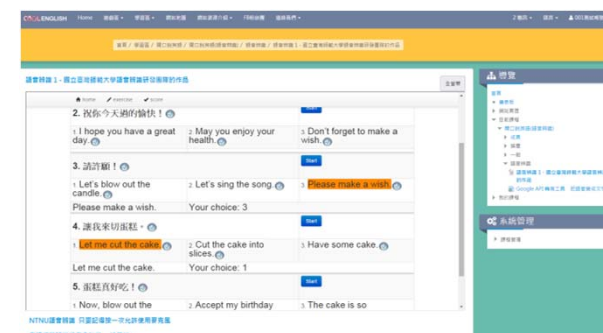


# Car Play Systems

- Aim to s provide a smarter, safer way to use your communication devices in vehicle
- E.g. Apple Car Play



# Computer-Assisted Pronunciation Training (CAPT)



網站介紹 | 網站統計分析

- **Pronunciation of Lexical Tones:** Detection and Assessment
- **Pronunciation of Sub-word (Syllable, INITIAL/FINAL) Units:** Detection and Assessment
- **Duration/ Speaking Rate (Fluency/Proficiency):** Detection and Assessment
- **Overall Scoring (word-, phrase-, sentence-levels)**

專為國內國中、小學生所設計的  
**英語線上學習平臺**

Cool English主要分為「學習區」與「遊戲區」，提供多元的學習單元與內容，主要的特色與目標成效如下：

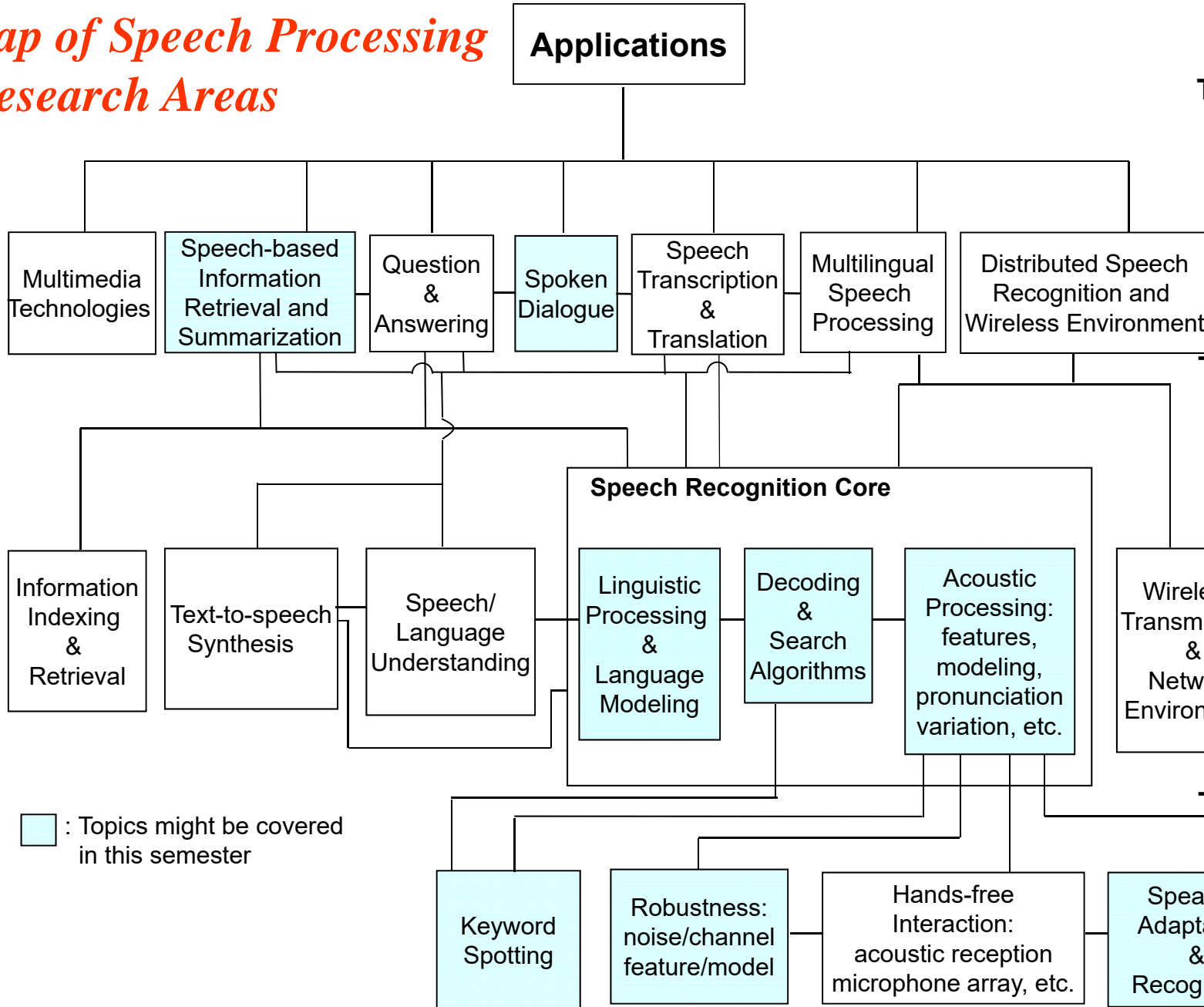
- 基本學習內容，動畫版，提供基本學習教材讓學生進行課後練習。
- 情境動畫聽講，提升學習者的聽講能力。
- 開口說英語，透過語音辨識系統的研發讓學習者開口說英語。
- 圖解式閱讀，輔助學習者理解閱讀內容。
- 主題式字彙練習，配合基本1200單字發展課程遊戲，讓學習字彙更有趣。
- 文法關卡挑戰，提供互動式文法練習題，可讓學生自我反覆練習。
- 國中課本中英雙語字彙檢索，幫助學生了解英文詞彙的用法。
- 國內外英語學習相關資源，讓學生獲得更多種豐富的英語學習資源。
- 遊戲嘉年华，精選Flash遊戲，讓英語學習變得更有趣。
- 冒險式RPG角色扮演遊戲，讓學習者沈浸於使用英語的學習環境。
- APP字彙遊戲下載，讓學生可以享受英語行動學習的樂趣。



1. Mandarin Chinese CAPT: <http://140.122.96.191/ALS/assessment.aspx>
2. English CAPT: <http://www.coolenglish.edu.tw/>

# Map of Speech Processing Research Areas

Emerging Technologies



Applied Technologies

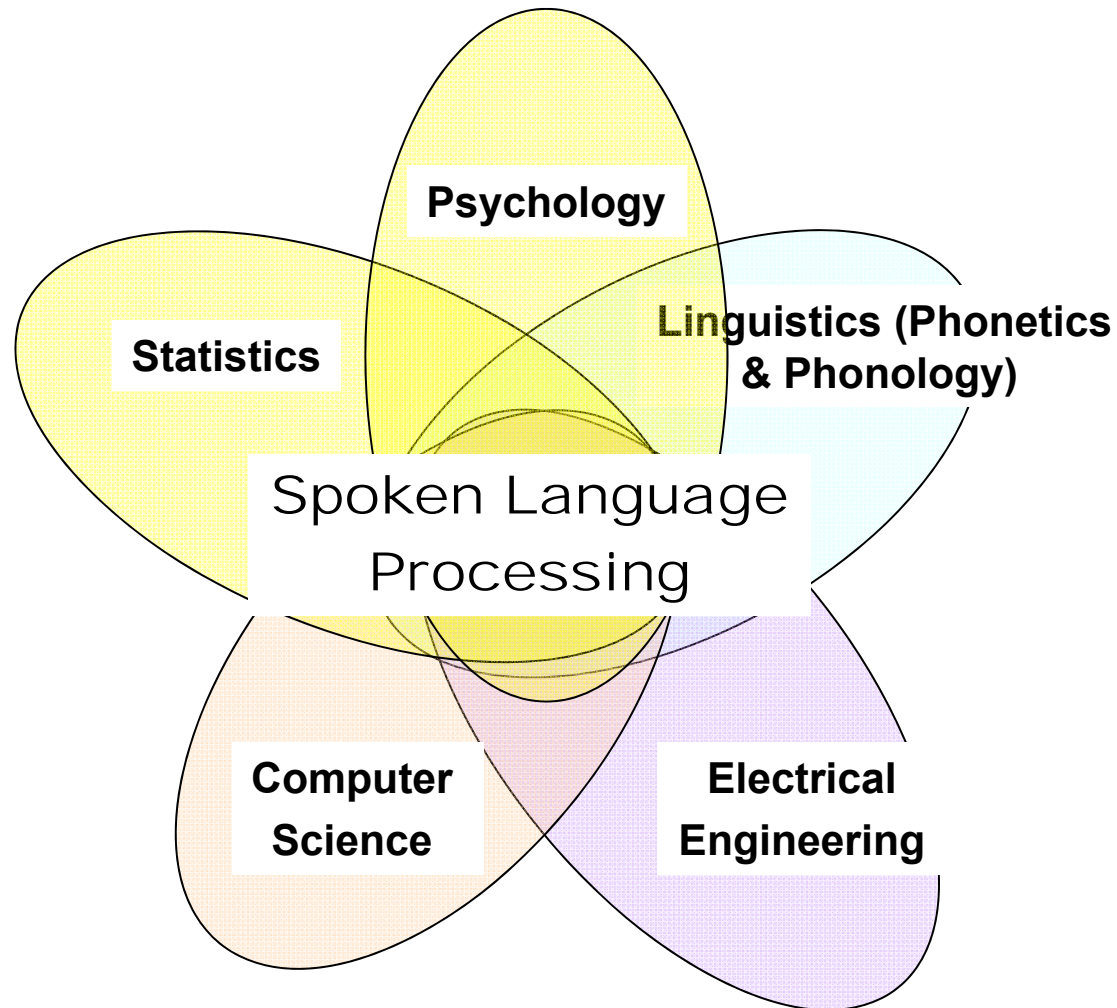
Integrated Technologies

Basic Technologies

□ : Topics might be covered in this semester

# Different Academic Disciplines

- The foundations of spoken language processing lies in

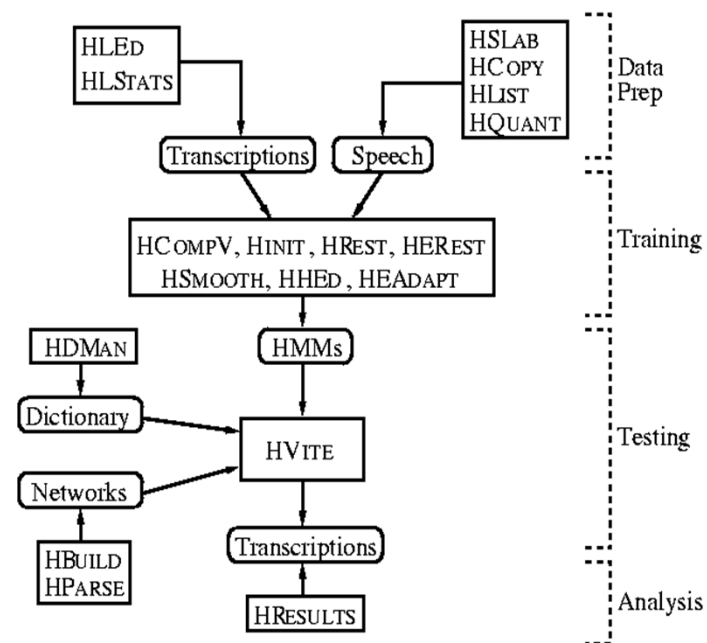


# Speech Processing Toolkit (1/2)

- **HTK (Hidden Markov Model ToolKit)**
  - A toolkit for building Hidden Markov Models (HMMs)
  - The HMM can be used to model any time series and the core of HTK is similarly general-purpose
  - In particular, for the acoustic feature extraction, HMM-based acoustic model training and HMM network decoding

# Speech Processing Toolkit (2/2)

- HTK (**H**idden **M**arkov **M**odel **T**ool**K**it)



- Nowadays, **Kaldi** emerges as a cutting-edge toolkit for developing speech recognition tasks

<http://kaldi.sourceforge.net/>



# Journals & Conferences

- Journals

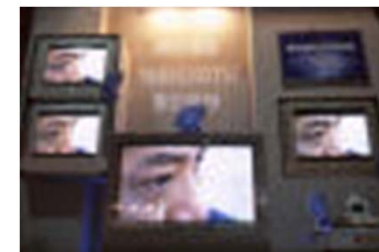
- IEEE Transactions on Audio, Speech and Language Processing
- Computer Speech & Language
- Speech Communication
- Proceedings of the IEEE
- IEEE Signal Processing Magazine
- ACM Transactions on Speech and Language Processing
- ACM Transactions on Asian and Low-Resource Language Information Processing
- ...

- Conferences

- IEEE International Conference on Acoustics, Speech, Signal processing (ICASSP)
- Annual Conference of the International Speech Communication Association (Interspeech)
- IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)
- IEEE Workshop on Spoken Language Technology (SLT)
- International Symposium on Chinese Spoken Language Processing (ISCSLP)
- ROCLING Conference on Computational Linguistics and Speech Processing
- ...

# Speech Industry (1/3)

- Telecommunication
- Information Appliance
- Interactive Voice Response
- Voice Portal
- Multimedia Database
- Education
- .....



# Tentative Schedule

Topics to be Covered
Overview & Introduction
Hidden Markov Models
Spoken Language Structure
Acoustic Modeling & HTK Toolkit & Kaldi Toolkit
Statistical Language Modeling & SRI LM Toolkit
Speech Signal Representations
Digit Recognition, Word Recognition and Keyword Spotting
Large Vocabulary Continuous Speech Recognition (LVCSR)
Speech Enhancement and Environment Robustness
Model Training and Adaptation Techniques
Utterance Verification and Confidence Measures