# Speech Recognition

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

- Technology for **Automatic Speech Recognition** (ASR) 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
  - 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
  - T. F. Quatieri. Discrete-Time Speech Signal Processing - Principles and Practice. Prentice Hall, 2002
  - 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 Star. 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 "Fall 2012: Speech Recognition"

# 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

1952, Isolated-Digit
Recognition, Bell Lab.

1956, Ten-Syllable
Recognition, RCA

1959, Ten-Vowel
Recognition, MIT Lincoln Lab

**1959, Phoneme-sequence Recognition using Statistical Information of Context, Fry and Denes**

**1960s, Dynamic Time Warping to Compare Speech Events, Vintsyuk**

**1960s-1970s, Hidden Markov Models for Speech Recognition, Baum, Baker and Jelinek**

Gestation of Foundations

**1970s ~**

**Voice-Activated Typewriter (dictation machine, speaker-dependent), IBM**

**Telecommunication (keyword spotting, speaker-independent), Bell Lab**

Philips

BBN Technologies

Cambridge (HTK)

CMU

LIMSI

MIT (SLS)

JHU CLSP

Microsoft

SRI

Google

Apple

nuance

# 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)



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 Switchboad 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 $\mathbf{X}$ into a word sequence $\hat{\mathbf{W}}$ (Hopefully, $\hat{\mathbf{W}} \approx \mathbf{W}$.)

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



語音特徵參數抽取

語言解碼/搜尋演算法

語音輸入

Feature Extraction

Feature Vectors

Linguistic Decoding and Search Algorithm

文字輸出

Speech Corpora

Acoustic Modeling

Acoustic Models

Lexicon 詞典

Language Models

Language Modeling

Text Corpora

語音 資料庫

聲學模型之建立

語言模型之建立

文字 資料庫

可能詞句

語音輸入

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

Bayes Rule

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

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

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
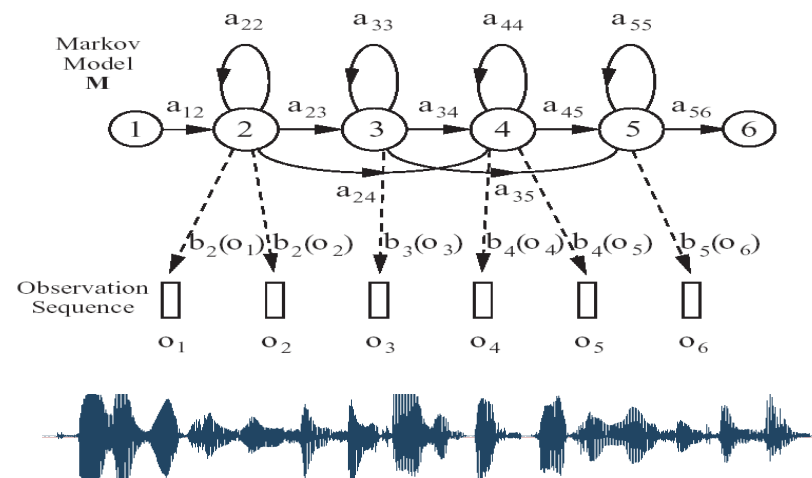  - HMMs are based on rigorous mathematical theory built on several decades of mathematical results developed in other fields
  - HMMs are constructed by the process of training on a large corpus of real speech data

# 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

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

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

**bigram modeling**

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

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

# Difficulties: Speech Variability



Pronunciation Variation

Speaker-independency
Speaker-adaptation
Speaker-dependency

Linguistic variability

Inter-speaker variability

Intra-speaker variability

Variability caused by the environment

Variability caused by the context

Robustness Enhancement

Context-Dependent Acoustic Modeling

# Deep Learning (Neural Networks) for ASR

- 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



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

# Text to Speech (TTS)

- 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

# 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 Interface

Applications

Visualization: Graphs & Table

Spoken language understanding modules

V. Zue, J.R. Glass, Conversational Interfaces: Advances and Challenges. Proceedings of the IEEE, Vol. 88, No. 8, August 2000

# 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
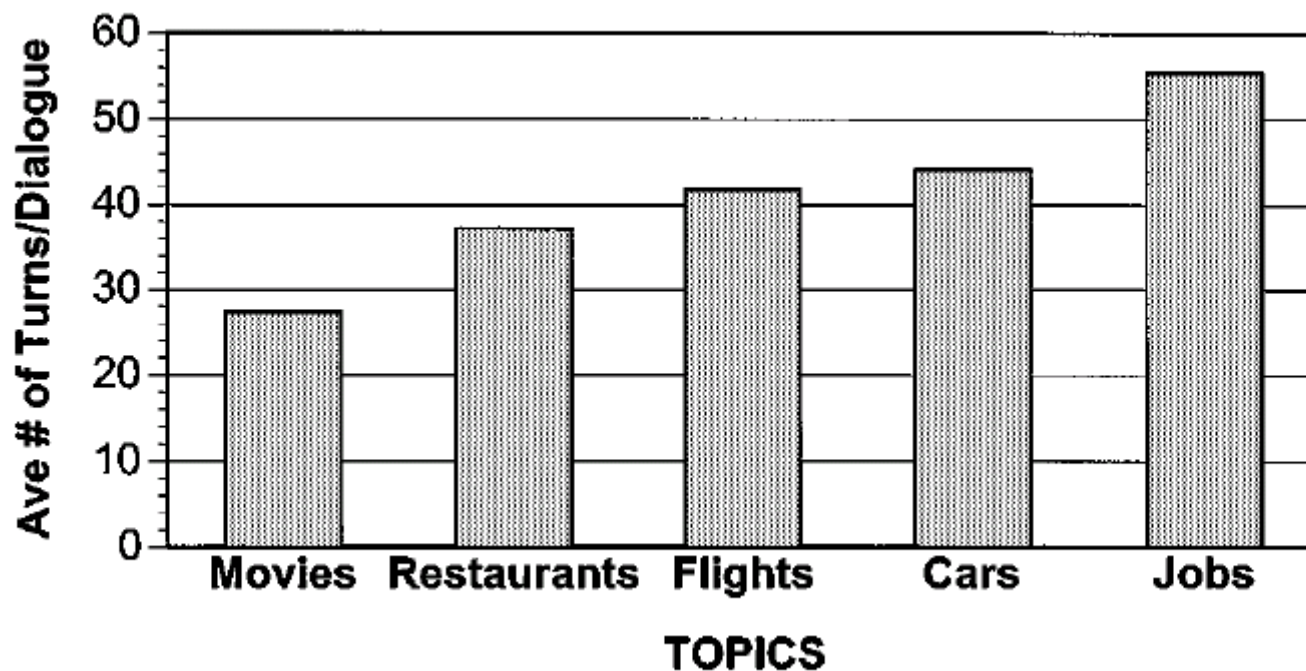
# Spoken Dialogue: Some Deployed Systems

- Complexity Analysis

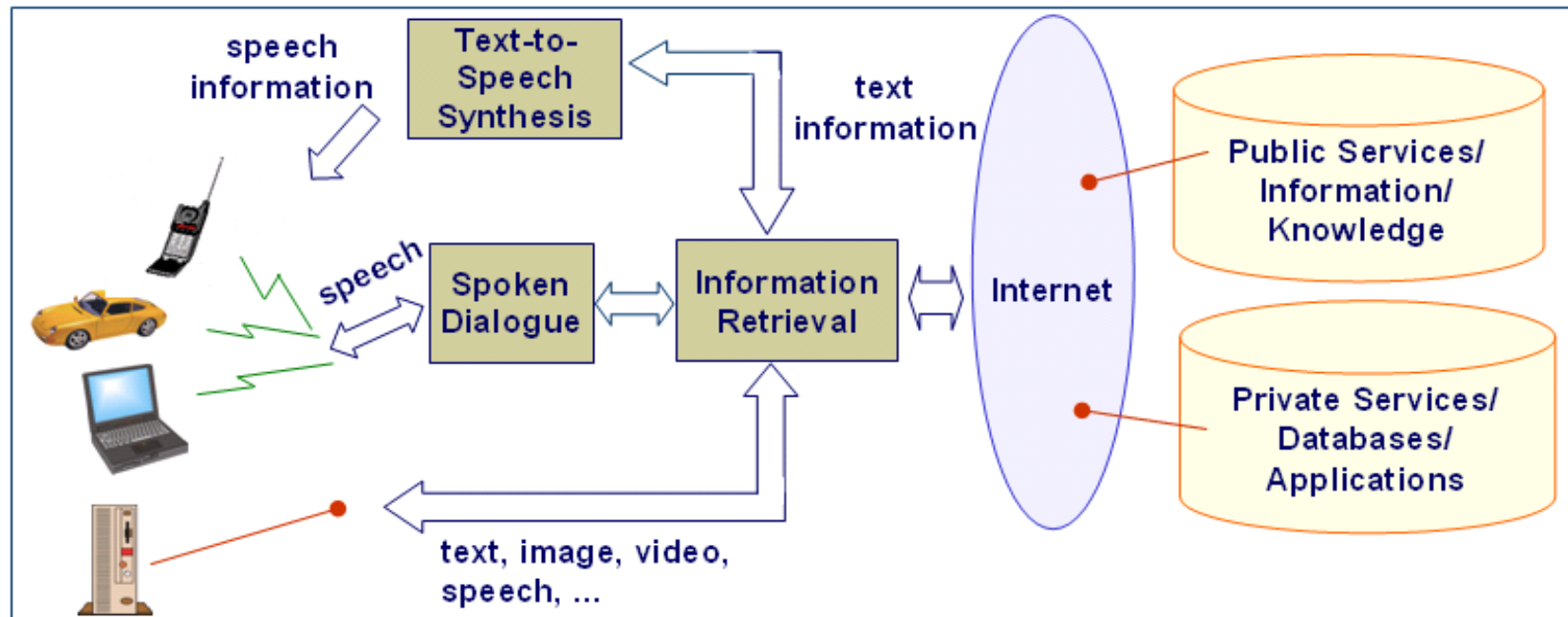| Domain | Language | Vocabulary Size | Average | |
|---|---|---|---|---|
| | | | Words/Utt | Utts/Dialogue |
| CSELT Train Timetable Info | Italian | 760 | 1.6 | 6.6 |
| SpeechWorks Air Travel Reservation | English | 1000 | 1.9 | 10.6 |
| Philips Train Timetable Info | German | 1850 | 2.7 | 7.0 |
| 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
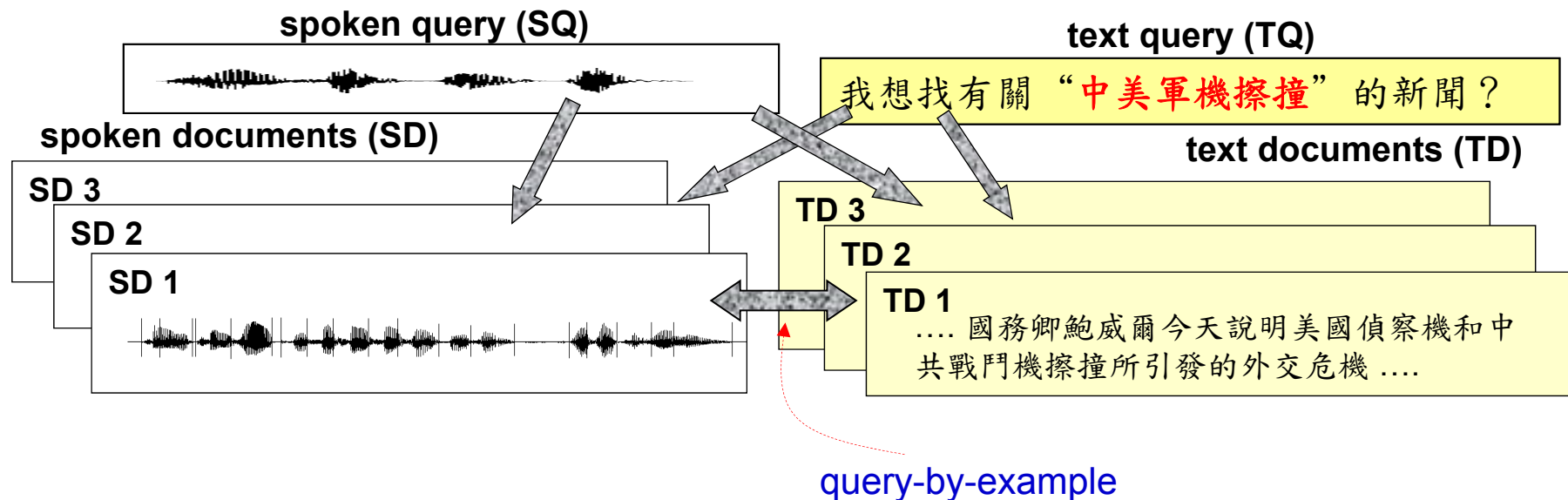
- Topics vs. Dialogue Terms

# 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)

**spoken query (SQ)**

**text query (TQ)**

我想找有關"中美軍機擦撞"的新聞？

**spoken documents (SD)**

**text documents (TD)**

SD 3

SD 2

SD 1

TD 3

TD 2

TD 1

.... 國務卿鮑威爾今天說明美國偵察機和中共戰鬥機擦撞所引發的外交危機 ....

query-by-example

- – 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

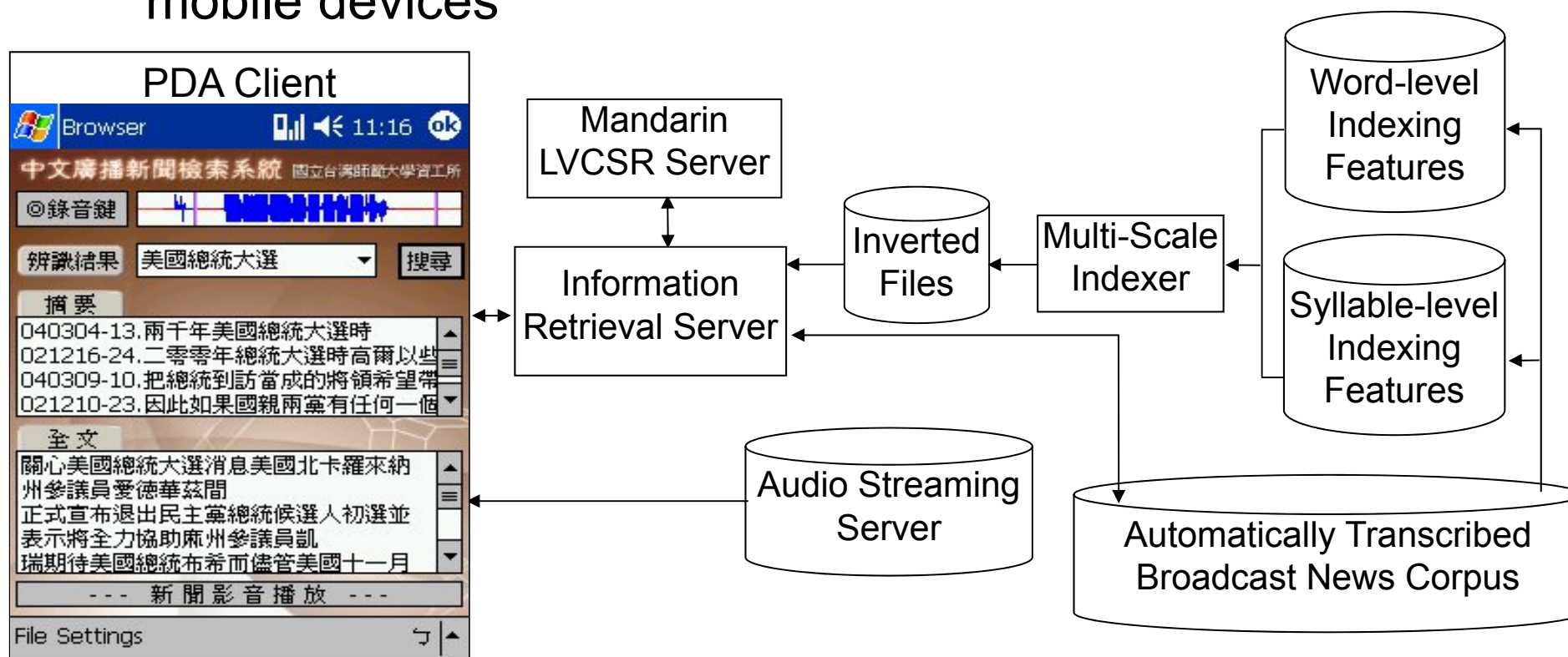# Speech-based Information Retrieval (3/5)

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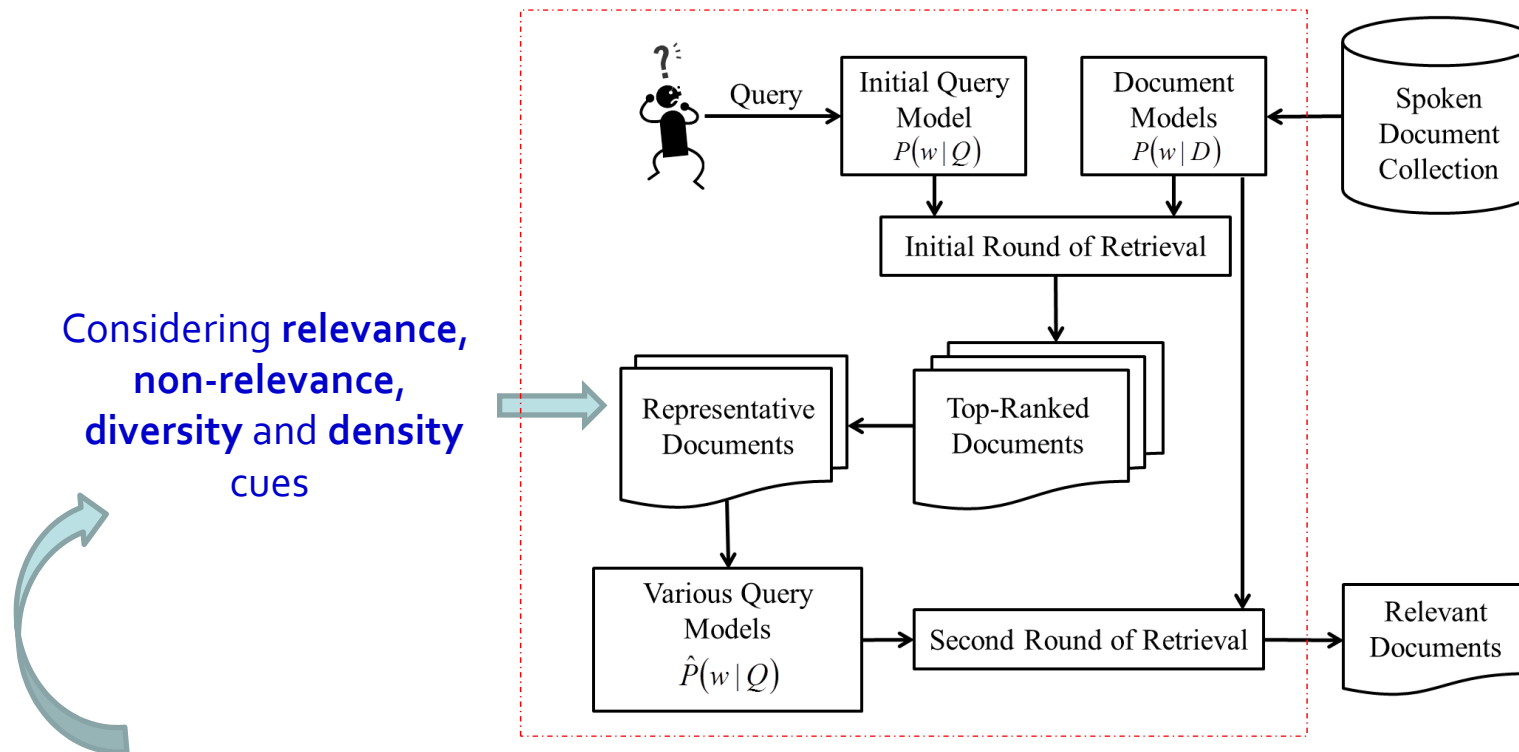
C.f. B. Chen, Y..T. Chen, C.H. Chang, H.B. Chen, "Speech Retrieval of Mandarin Broadcast News via Mobile Devices," Interspeech2005

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



Considering **relevance**, **non-relevance**, **diversity** and **density** cues

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

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 Dialogue: Google Voice Search



**Google Audio Indexing:**
**Searching what people are saying inside YouTube videos (currently only for what the politicians are saying)**

**Google-411:**
**Finding and connecting to local business**

# 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

- 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

**Map of Speech Processing Research Areas**

Applications

Applied Technologies

Integrated Technologies

Basic Technologies

Emerging Technologies

Multimedia Technologies

Speech-based Information Retrieval and Summarization

Question & Answering

Spoken Dialogue

Speech Transcription & Translation

Multilingual Speech Processing

Distributed Speech Recognition and Wireless Environment

**Speech Recognition Core**

Information Indexing & Retrieval

Text-to-speech Synthesis

Speech/ Language Understanding

Linguistic Processing & Language Modeling

Decoding & Search Algorithms

Acoustic Processing: features, modeling, pronunciation variation, etc.

Wireless Transmission & Network Environment

☐ : Topics might be covered in this semester

Keyword Spotting

Robustness: noise/channel feature/model

Hands-free Interaction: acoustic reception microphone array, etc.

Speaker Adaptation & Recognition

Adapted from Prof. Lin-shan Lee

SP- Berlin Chen   41

# Different Academic Disciplines

- The foundations of spoken language processing lies in

# Speech Processing Toolkit (1/2)

- HTK (**H**idden Markov Model **T**ool**K**it)
  - 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 Markov Model **T**ool**K**it)
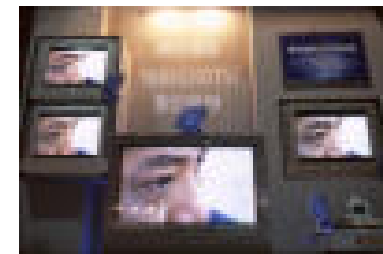
# Journals & Conferences

- ## 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 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
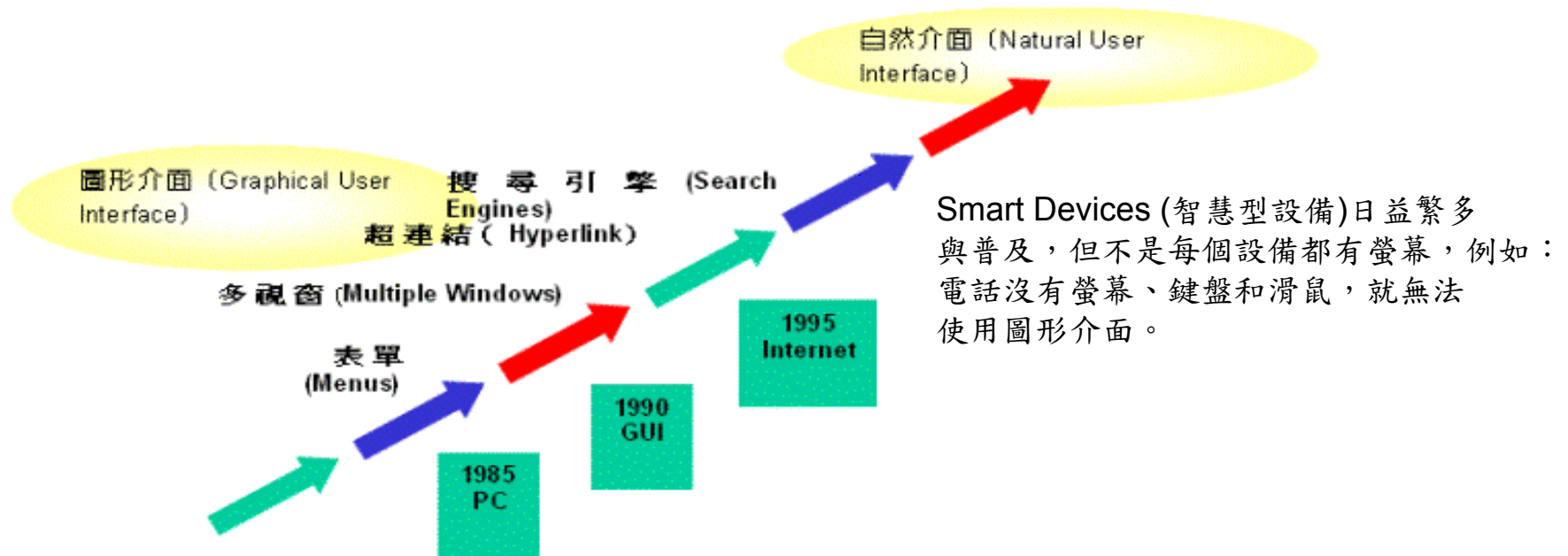
- …..

# Tentative Schedule

| Topics to be Covered |
|---|
| Overview & Introduction |
| Hidden Markov Models |
| Spoken Language Structure |
| Acoustic Modeling & HTK Toolkit |
| 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 |

# Speech Industry (2/3)

- Microsoft: Smart Device/Natural UI

使用介面的發展



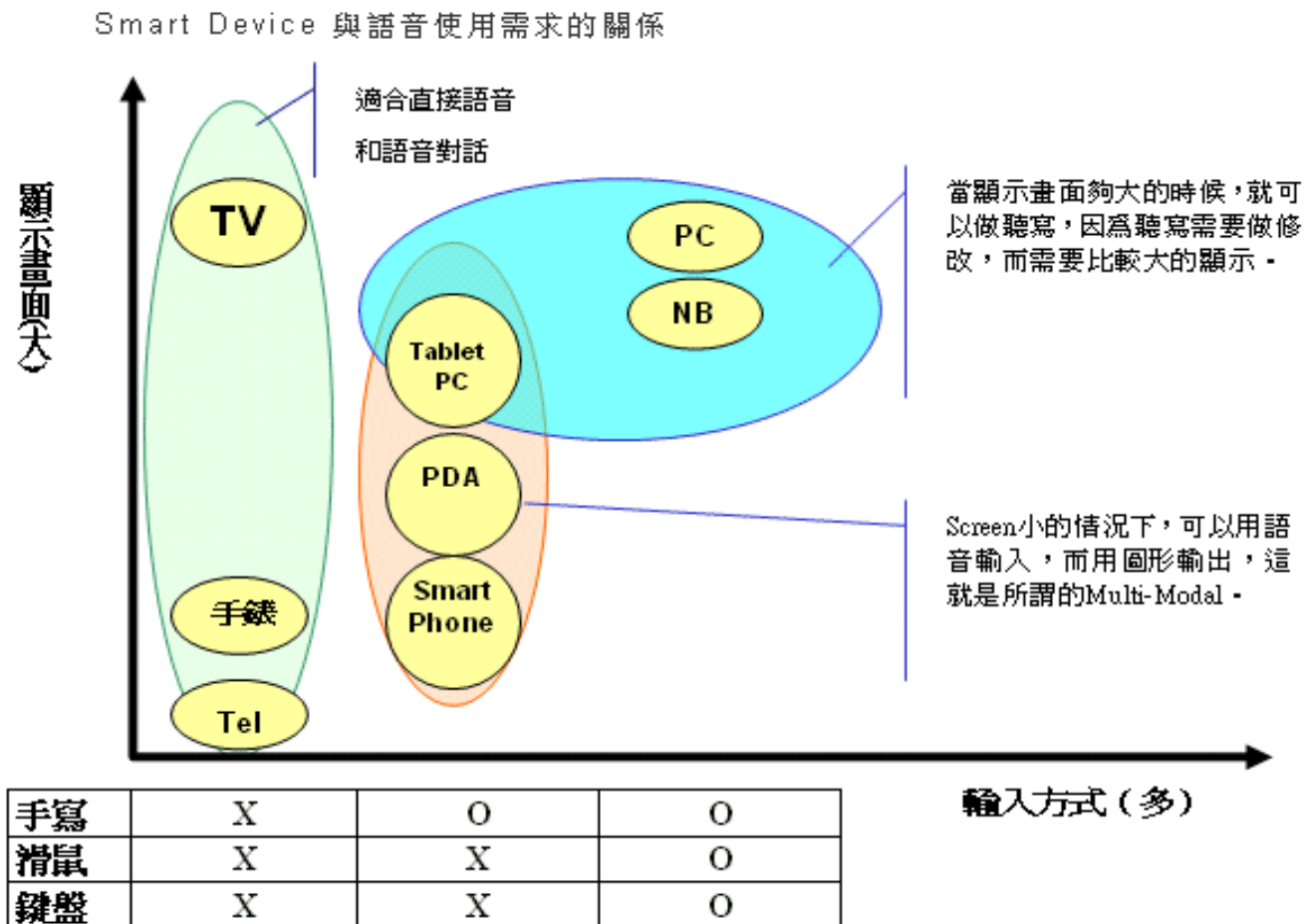Smart Devices (智慧型設備)日益繁多與普及，但不是每個設備都有螢幕，例如：電話沒有螢幕、鍵盤和滑鼠，就無法使用圖形介面。

Source：微軟自然互動服務產品部門 (NISD)副總裁李開複博士講稿， 2003/04

.NET 的最初構想，以符合人類需求的自然介面，其包括 –
- 語音合成
- 語音辨識技術
- 結合XML為基礎的網路服務

# Speech Industry (3/3)

- Microsoft: Smart Device/Natural UI



Smart Device 與語音使用需求的關係

適合直接語音和語音對話

當顯示畫面夠大的時候，就可以做聽寫，因為聽寫需要做修改，而需要比較大的顯示。

Screen小的情況下，可以用語音輸入，而用圖形輸出，這就是所謂的Multi-Modal。

顯示畫面(大)

輸入方式（多）

| 手寫 | X | O | O |
|---|---|---|---|
| 滑鼠 | X | X | O |
| 鍵盤 | X | X | O |

# Good Words

- Your attitude determines your altitude.

- Stay Hungry; Stay Foolish

- Every job is a self-portrait of those who did it. Autograph your work with quality.

- ….