



Speech Recognition and its Applications to Computer-Assisted Language Learning

語音辨識與其在電腦輔助語言學習之應用

Berlin Chen (陳柏琳)

Professor, Department of Computer Science & Information Engineering

National Taiwan Normal University

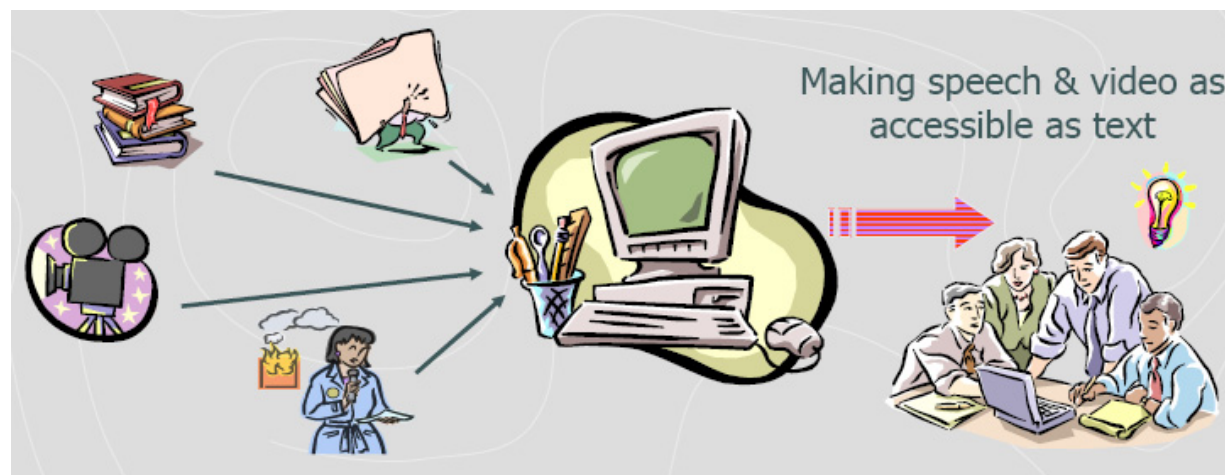
2013/06/08

Introduction (1/3)

- Communication and search are by far the most popular activities in our daily lives
 - Speech is the most natural and convenient means for communications among humans (and between humans and machines in the future)
 - A spoken language interface could be more convenient than a visual interface on a small (or hand-held) device
 - Provide "*anytime*" and "*anywhere*" access to information
 - Already over half of the internet traffic consists of video data
 - Though visual cues are important for search, the associated spoken documents often provide a rich set of semantic cues (e.g., *transcripts, speakers, emotions, and scenes*) for the data

Introduction (2/3)

- Automatic speech recognition (ASR)
 - Transcribe the **linguistic contents** of speech utterances
 - Play a vital role in multimedia information retrieval, summarization, organization, among others
 - Such as the transcription of spoken documents and recognition of spoken queries

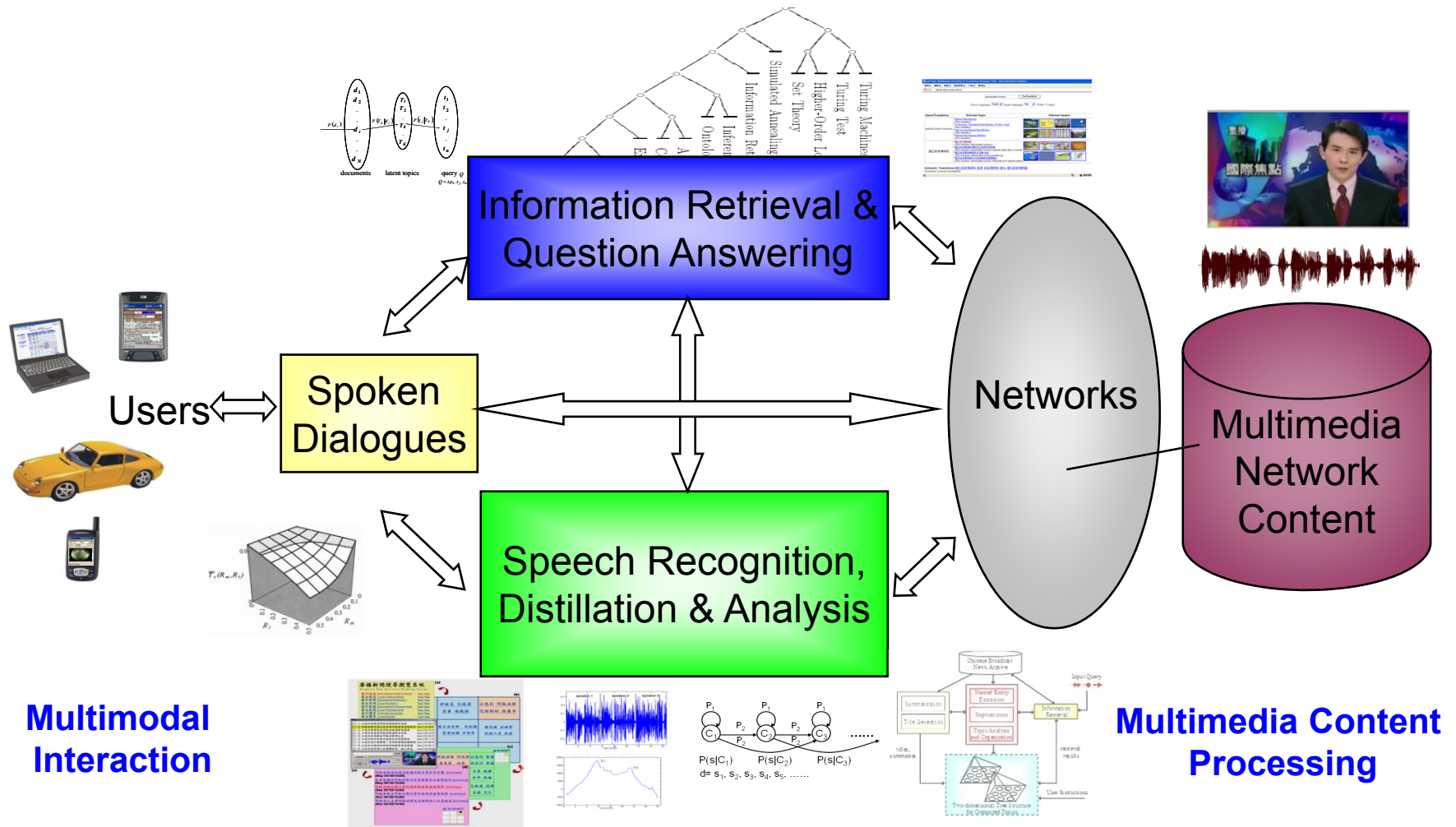


The figure is adapted from the presentation slides of Prof. Ostendorf at *Interspeech 2009*.

Introduction (3/3)

- Text Processing vs. Speech 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 the same 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

Multimodal Access to Multimedia in the Future



Automatic Speech Recognition (ASR)

- Bayes Decision Rule (Risk Minimization)

$$W_{opt} = \arg \min_{W \in \mathbf{W}} Risk (W | O)$$

$$= \arg \min_{W \in \mathbf{W}} \sum_{W' \in \mathbf{W}} Loss (W, W') P (W' | O)$$

$$\approx \arg \max_{W \in \mathbf{W}} P (W | O) \text{ Assumption of Using the "0-1" Loss Function}$$

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

$$= \arg \max_{W \in \mathbf{W}} p(O | W) P(W) \quad \text{Linguistic Decoding}$$

Feature Extraction & Acoustic Modeling

Language Modeling

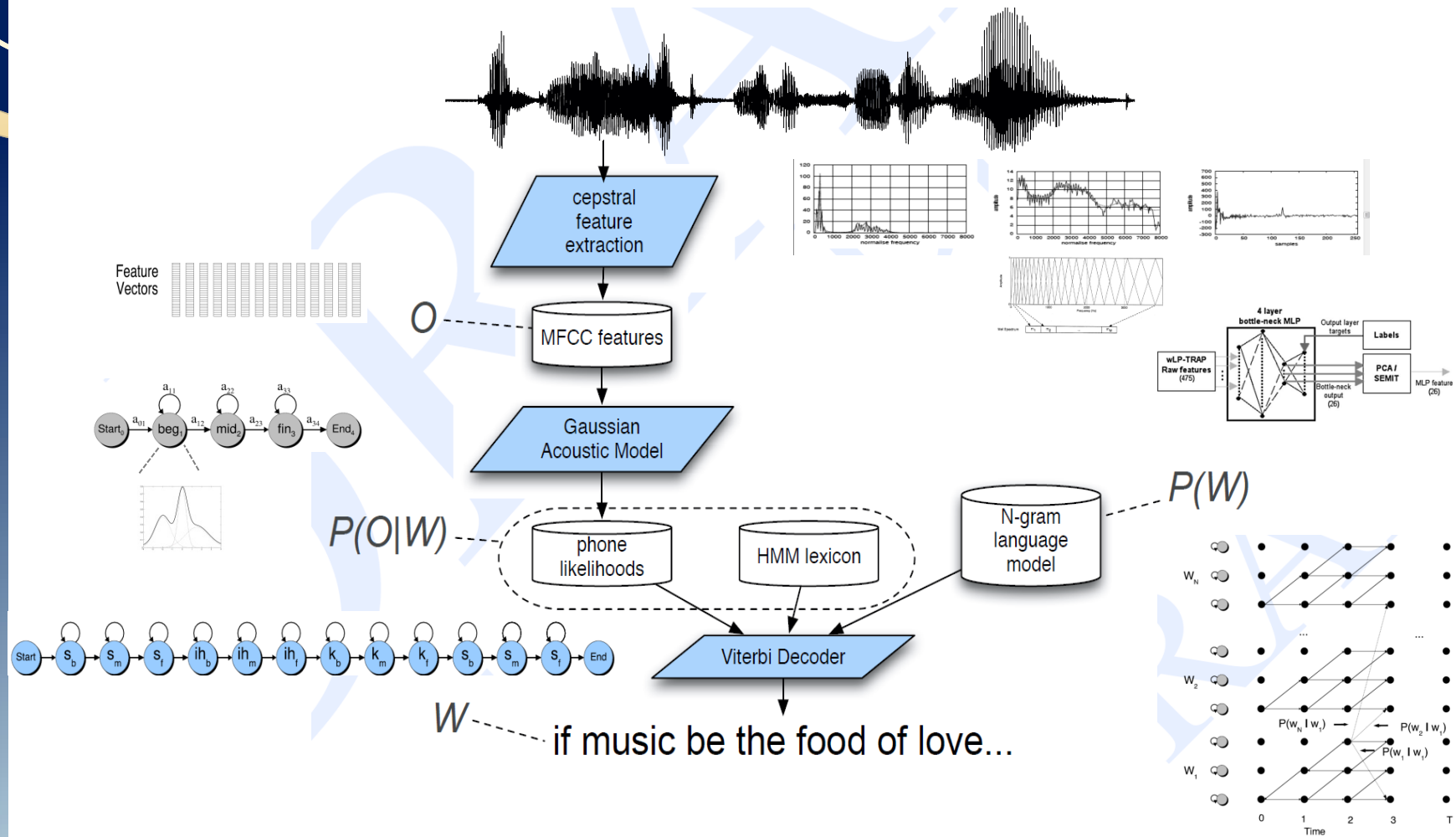
Possible variations

speaker, pronunciation, environment, context, etc.

and

OOV, domain, topic, style, etc.

Schematic Diagram of ASR



F. Valente et al., "Transcribing Mandarin broadcast speech using multi-layer perceptron acoustic features," *IEEE Transactions on Audio, Speech, and Language Processing*, 2011.

Core Components of ASR

- Feature Extraction
 - Convert a speech signal into a sequence of feature vectors describing the inherent acoustic and phonetic properties
- Acoustic modeling
 - Construct a set of statistical models representing various sounds (or phonetic units) of the language
- Language modeling
 - Construct a set of statistical models to constrain the acoustic analysis, guide the search through multiple candidate word strings, and quantify the acceptability of the final output from a speech recognizer
- Robustness
 - Eliminate varying sources of environmental (e.g., channel and background, pronunciation, speaker and context) variations

Applications of ASR

- Multimedia (spoken document) retrieval and organization
 - Speech-driven Interface and multimedia content processing
 - Work in concert with natural language processing (NLP) and information retrieval (IR) techniques
 - A wild variety of potential applications
- Computer-Aided Language Learning (CALL)
 - Speech-driven Interface and multimedia content processing
 - Work in in association with natural language processing techniques
 - Applications
 - Automatic pronunciation assessment/scoring (CAPT)
 - Synchronization of audio/video learning materials
 - Estimation of document (writing) readability
 - Automated reading tutor (with spoken dialogues)
- Others

Speech-driven Multimedia Retrieval & Organization

- Continuous and substantial efforts have been paid to speech-driven multimedia retrieval and organization in the recent past
 - *Informedia* System at Carnegie Mellon Univ.
 - *Rough'n'Ready* System at BBN Technologies
 - IBM Speech Search for Call-Center Conversations & Call-Routing, Voicemails, Monitoring Global Video and Web News Sources (*TALES*)
 - Google Voice Search (*GOOG-411*, *Audio Indexing*, *Translation*)
 - Microsoft Research *Bing Mobile Voice Search*
 - Apple's *Siri* (QA)
 - MIT Lecture Browser

We are witnessing the golden age of ASR!

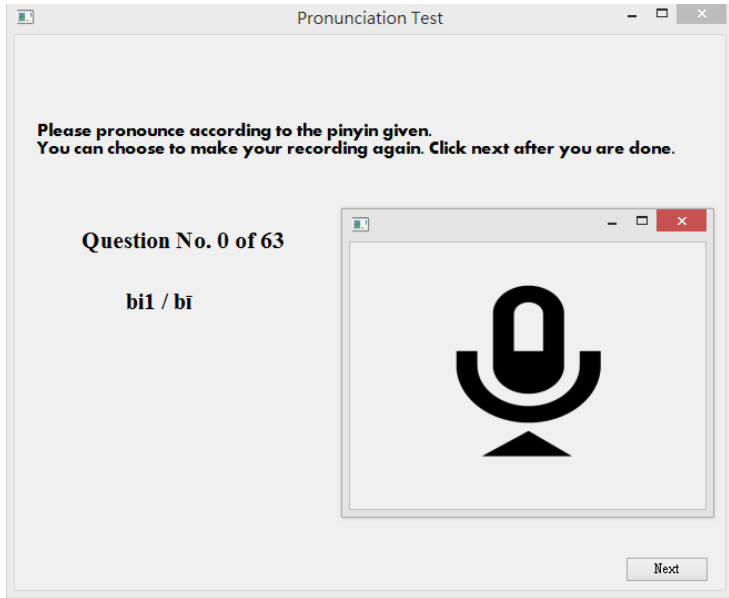


ASR for CALL Applications

- The major task of applying ASR for **CAPT (Computer Assisted Pronunciation Training)** is to automatically detect pronunciation errors and evaluate pronunciation quality
- Facets for ASR-based CAPT in **Mandarin Chinese**
 - **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)

Detection of Tone/Phone Mispronunciations (1/4)

- Detect possible mispronunciations and corresponding error patterns for a Chinese-language learner



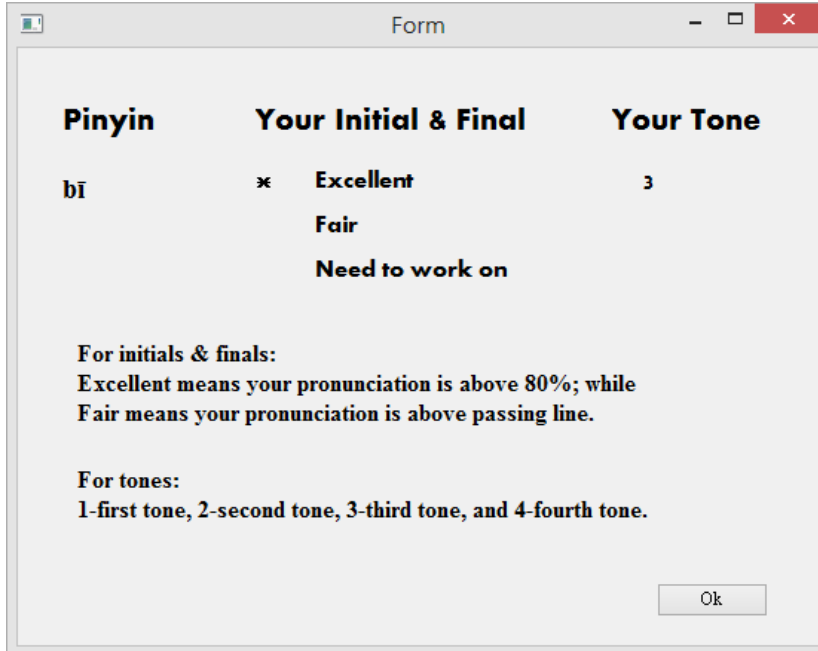

Pronunciation Test

Please pronounce according to the pinyin given.
You can choose to make your recording again. Click next after you are done.

Question No. 0 of 63

bi1 / bi

Next



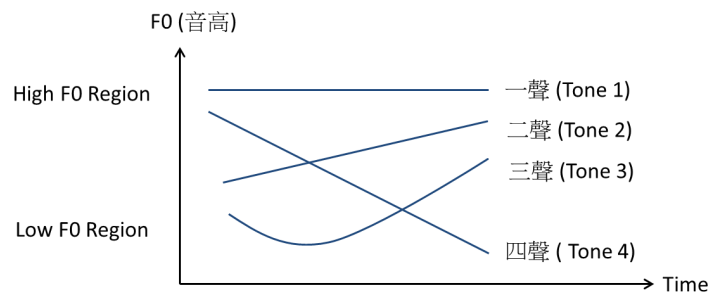
Form

Pinyin	Your Initial & Final	Your Tone
bi	✘ Excellent Fair Need to work on	3

For initials & finals:
Excellent means your pronunciation is above 80%; while
Fair means your pronunciation is above passing line.

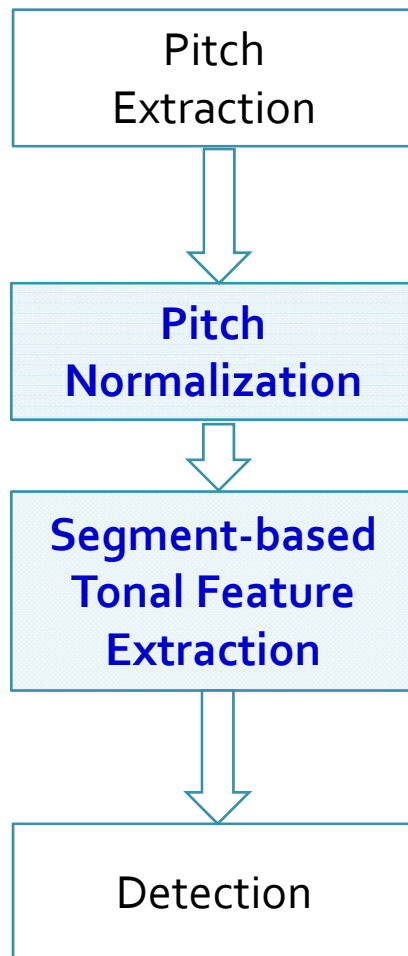
For tones:
1-first tone, 2-second tone, 3-third tone, and 4-fourth tone.

Ok

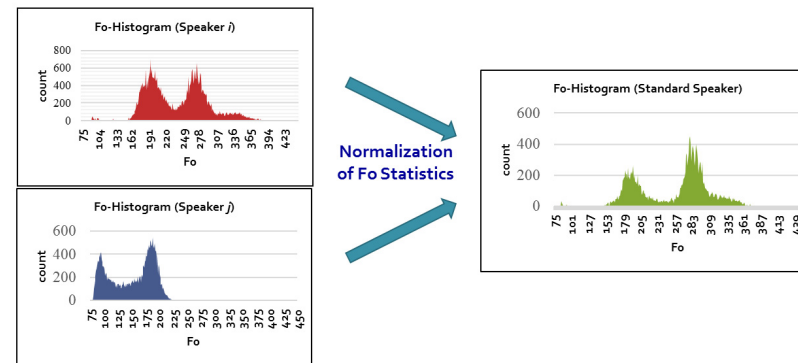


Detection of Tone/Phone Mispronunciations (2/4)

- Typical Steps for Lexical Tone Mispronunciation Detection



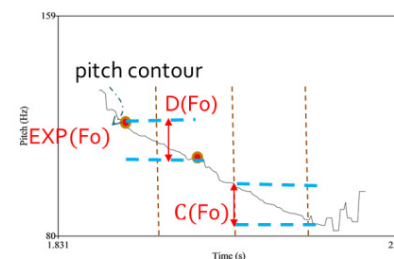
How to mitigate the negative effects caused by speaker and environmental variations?



How does the subtleness (granularity) of tonal features affect mispronunciation detection?

Types of Features

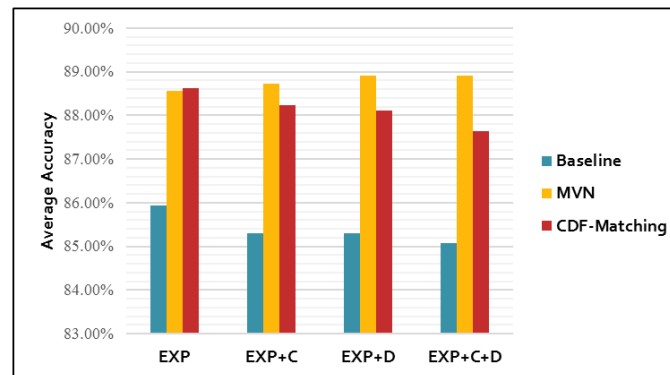
EXP[Fo]	Mean Fo in each segment
C(Fo) (Within-Segment Δ Fo)	Difference of beginning and ending Fo values within each segment
D(Fo) (Between-Segment Δ Fo)	Difference of EXP[Fo] values between any pair of segments



The pitch contour is extracted with the **RAPT** (Robust Algorithm for Pitch Tracking) method.

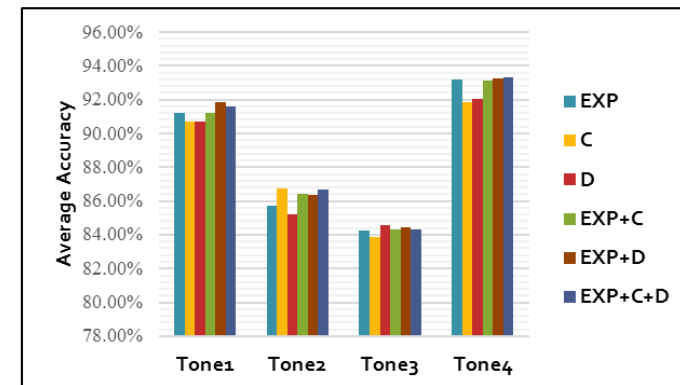
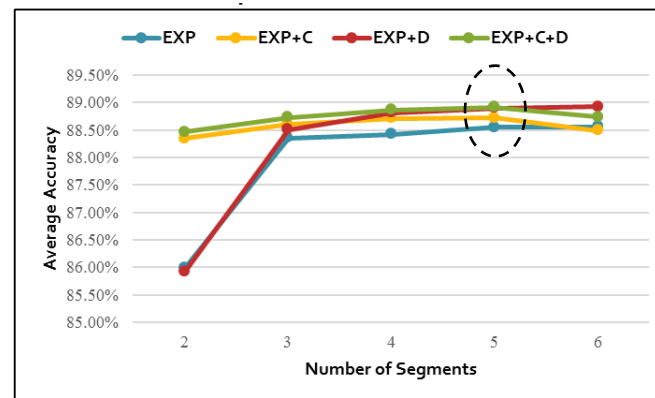
Detection of Tone/Phone Mispronunciations (3/4)

- Results on Automatic Detection of Tone Mispronunciations
 - Comparisons Among Different Normalization Methods



Both methods (MVN and CDF-matching) can offer significant performance boots compared to the baseline (without normalization)

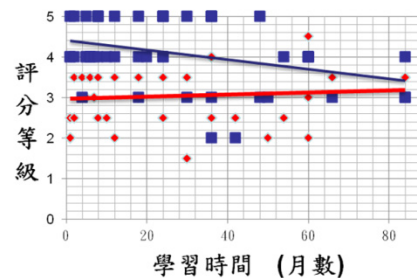
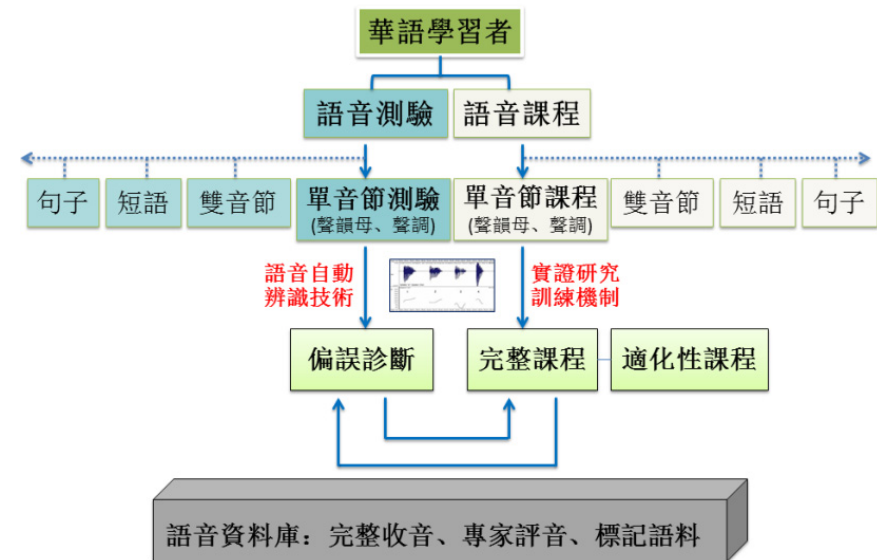
- Different Numbers of Segments for Tonal Feature Extraction



(with 5 segments)

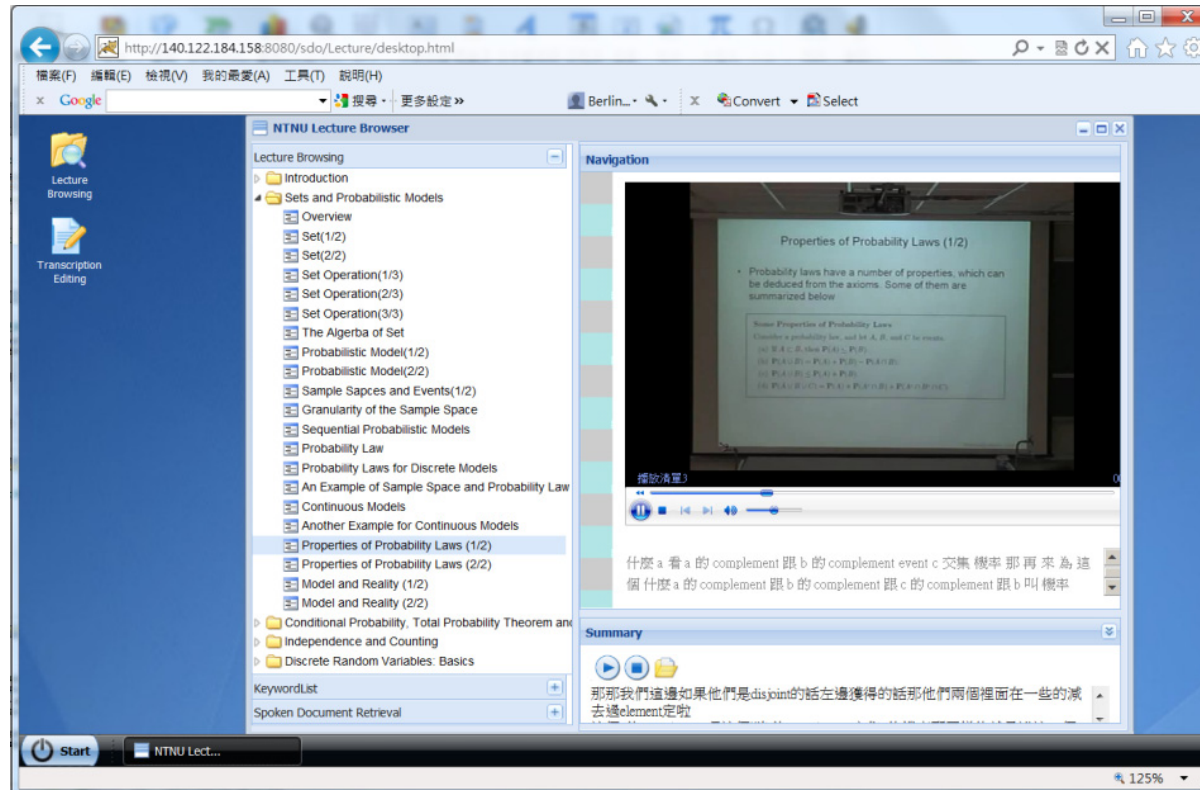
Detection of Tone/Phone Mispronunciations (4/4)

- We are now conducting a project to build a **Chinese Learning and Assessment System**



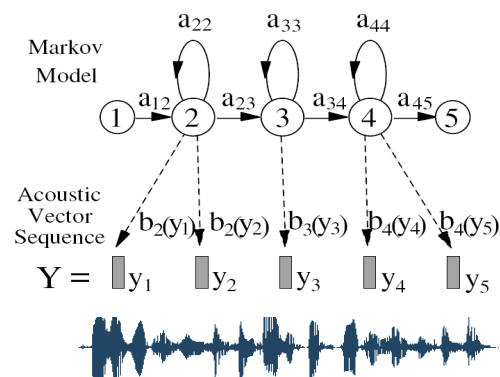
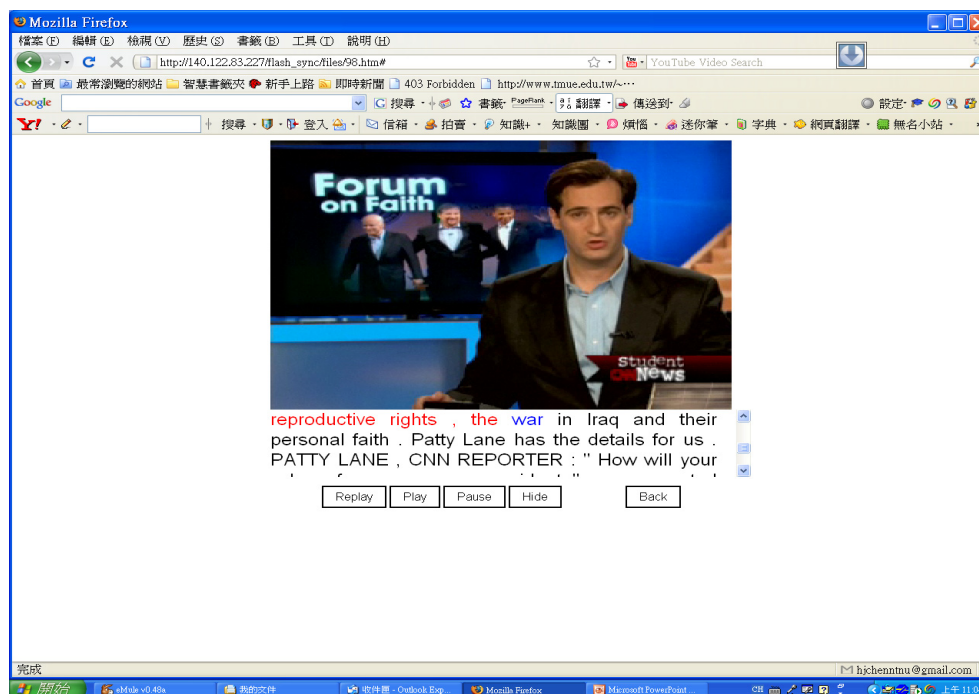
- ◆ 自評程度
- 系統滿意度
- 線性(自評程度)
- 線性(系統滿意度)

NTNU Lecture Browser for Self-Learning



- Featured with ASR-based automatic transcription and speech summarization and retrieval functionality

Video and Script Synchronization for Spoken English Learning



- This figure adapted from:

<http://webho1.ua.ac.be/linguapolis/call2008/CALL%202008%20synctolearn.ppt>

Spoken Document Organization (NTU & NTNU)

廣播新聞搜尋瀏覽系統
Broadcast News Retrieval/Browsing System

[國外政治 \[International Political News\]](#) Topic Map
[國內政治 \[Local Political News\]](#) Topic Map
[國外財經 \[International Business\]](#) Topic Map
[國內財經 \[Local Business\]](#) Topic Map
[國外影劇 \[International Entertainment\]](#) Topic Map
[國內影劇 \[Local Entertainment\]](#) Topic Map
[國外體育 \[International Sports\]](#) Topic Map
[國內體育 \[Local Sports\]](#) Topic Map

[1] 以色列結束對阿拉法特總部的包圍 [sum.] 02.09.21
 [2] 阿拉法特反對以色列保所提結束包圍條件 [sum.] 02.09.21
 [3] 以色列部隊進攻阿拉法特總部後撤軍 [sum.] 02.10.22
 [4] 以色列結束對阿拉法特總部的包圍 [sum.] 02.10.01
 [5] 以色列坦克撤出阿拉法特辦公區 [sum.] 02.09.21
 [6] 以色列與巴勒斯坦展開安全問題會議 [sum.] 02.11.23
 [7] 以色列在加薩擊斃一名回教聖戰組織領袖 [sum.] 02.06.05
 [8] 以色列巴勒斯坦就伯利恆撤軍達成協議 [sum.] 02.02.12
 [9] 以色列坦克闖入加薩難民營 兩人喪生 [sum.] 02.04.20

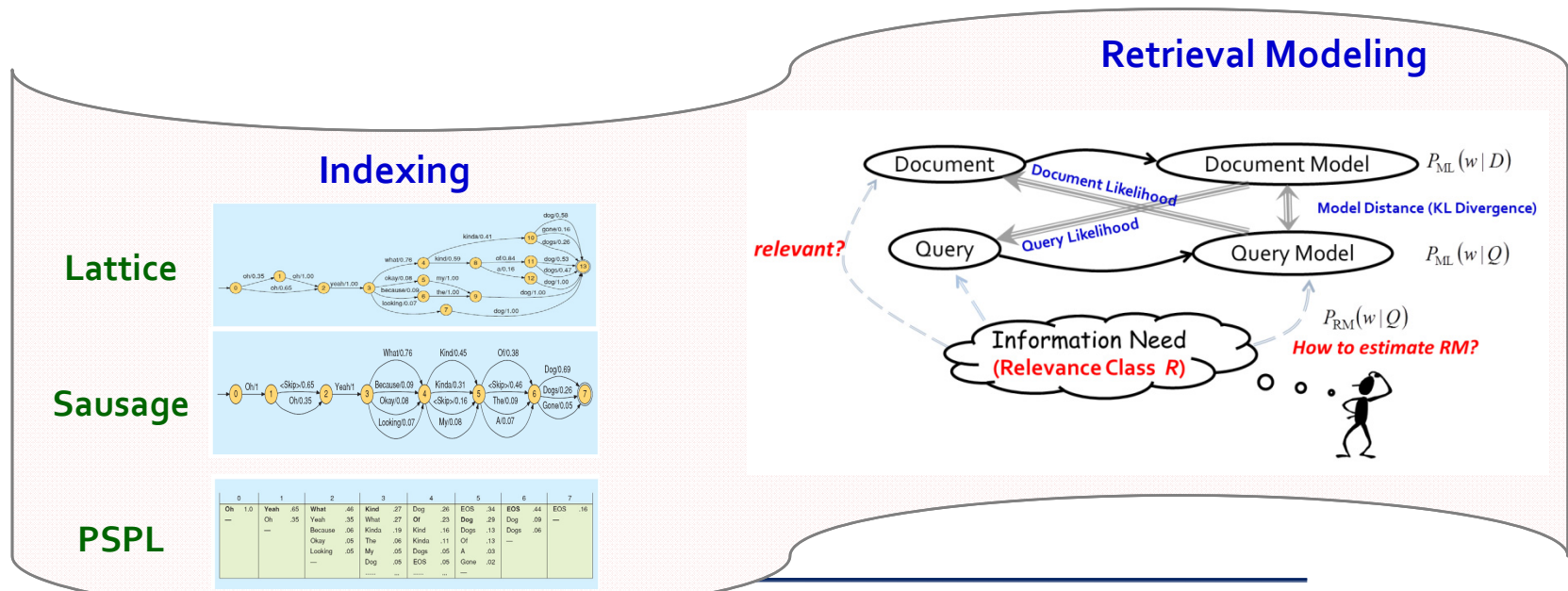
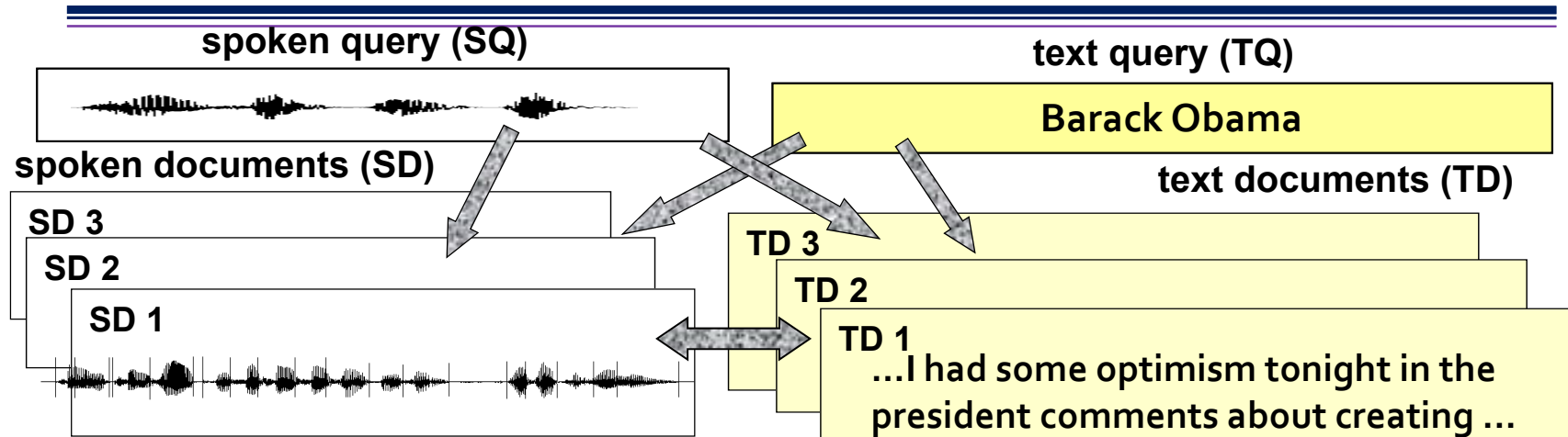
阿拉法特 阿巴斯 以色列 夏隆
 雷馬拉 任命 約旦河 英國
 中東 鮑爾
 和平 路線
 巴格達 炸彈
 自殺 巴士

伊拉克 巴格達 以色列 阿拉法特
 美軍 陸戰隊 巴勒斯坦 迦薩市
 國土安全部 民航機 聯合國 安理會
 蓋達組織 中情局 武檢人員 武器

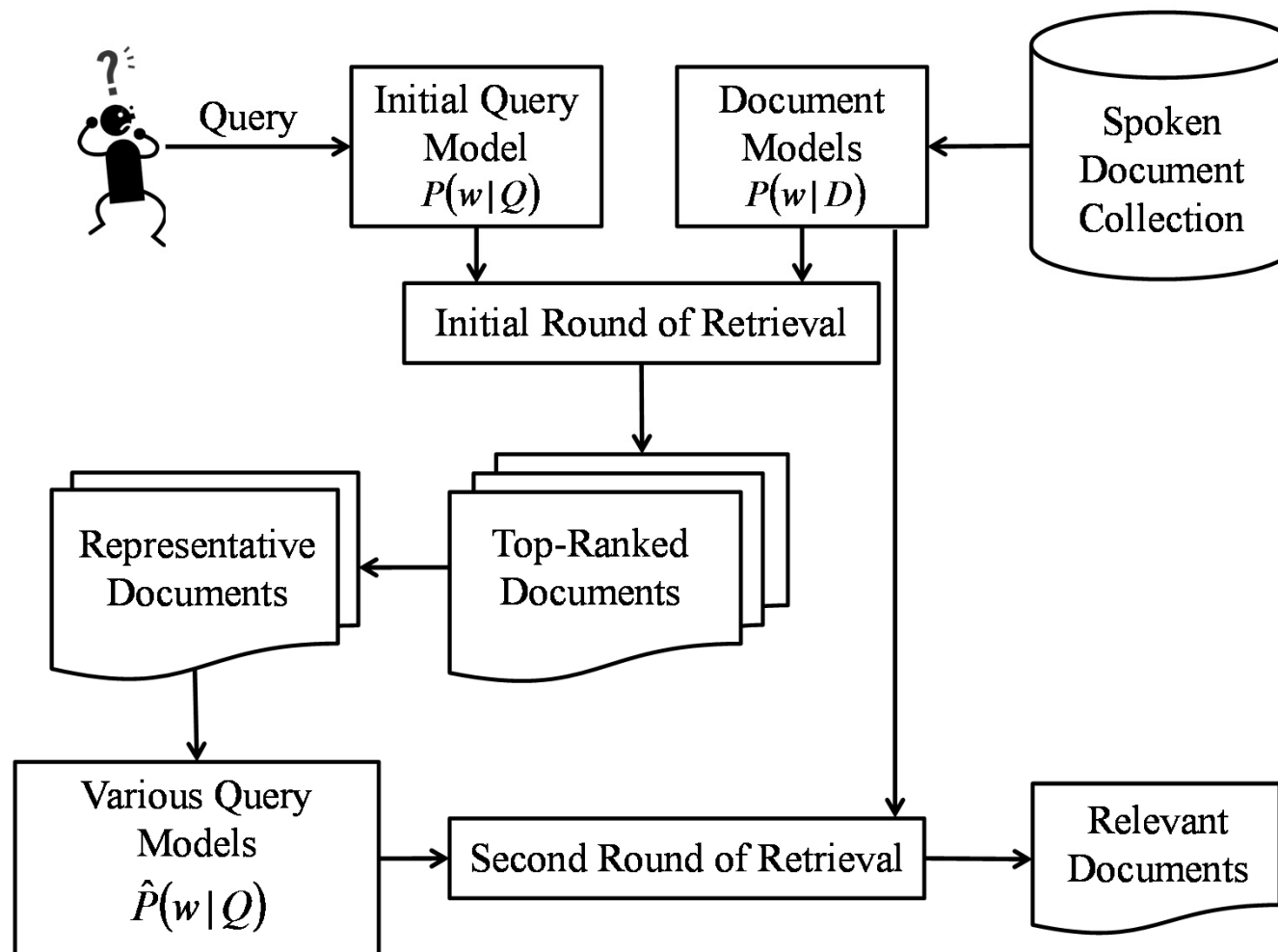
阿拉法特原則接受歐盟所提中東和平計畫 [summary]
 (May 03/02/12:00)
 英美就解決阿拉法特所受包圍與巴方展開談判 [summary]
 (May 06/02/12:00)
 阿拉法特反對以色列保所提結束包圍條件 [summary]
 (Sep 20/02/12:00)
 阿拉法特宣布新內閣引發巴勒斯坦國會激辯 [summary]
 (Oct 30/02/12:00)
 阿拉伯人支持阿拉法特及巴勒斯坦人正當抵抗 [summary]
 (Nov 02/02/12:00)

go to Level-1
 go to Level-2

Speech Retrieval: Scenarios and Methodologies



Speech Retrieval: Pseudo-relevance Feedback



Speech Summarization

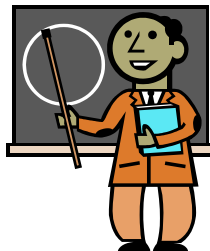
conversations



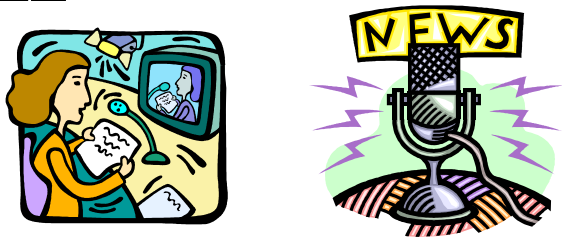
meetings



lectures



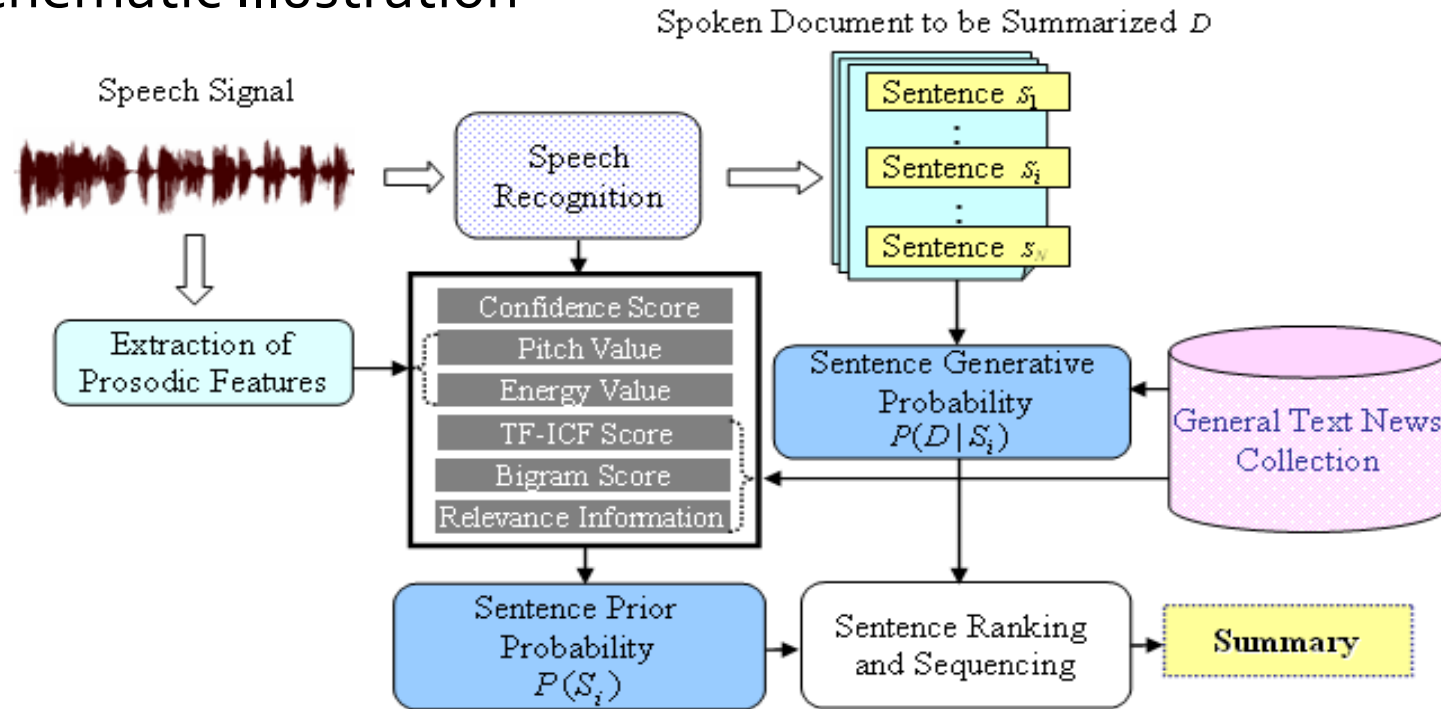
broadcast
and TV news



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

Risk Minimization-based Speech Summarization

Schematic Illustration



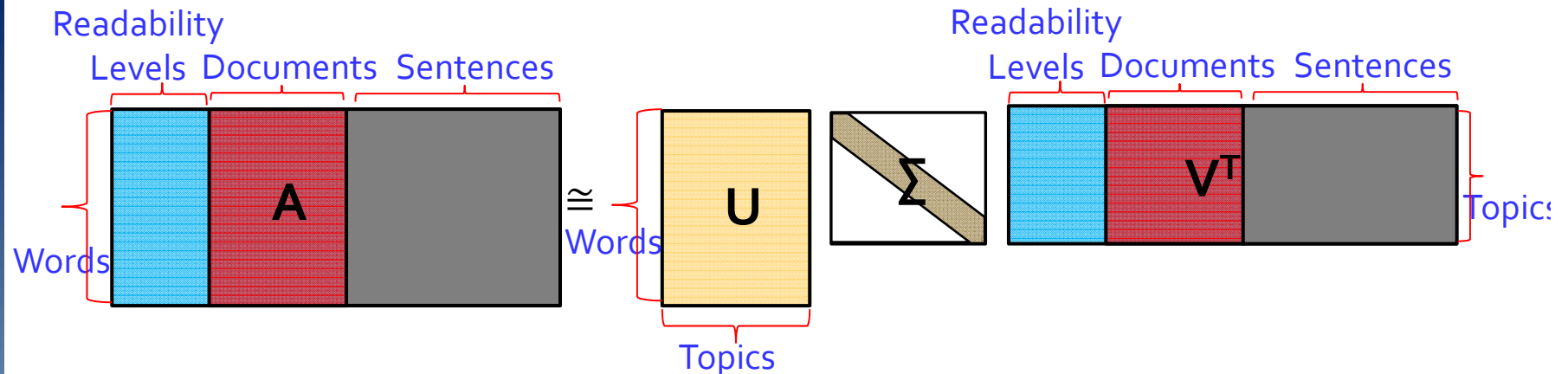
$$S^* = \arg \min_{S_i \in \tilde{D}} \sum_{S_j \in \tilde{D}} \text{Loss}(S_i, S_j) \cdot P(S_j | \tilde{D})$$

$$= \arg \min_{S_i \in \tilde{D}} \sum_{S_j \in \tilde{D}} \text{Loss}(S_i, S_j) \cdot \frac{P(\tilde{D} | S_j) P(S_j)}{\sum_{S_m \in \tilde{D}} P(\tilde{D} | S_m) P(S_m)}$$

Types	Description
Structural feature	1. Duration of the current sentence (S1)
Lexical features	1. Number of named entities (L1) 2. Number of stop words (L2) 3. Bigram language model scores (L3) 4. Normalized bigram scores (L4)
Acoustic features	1. The 1st formant (F1-1 to F1-5) 2. The 2nd formant (F2-1 to F2-5) 3. The pitch value (P-1 to P-5) 4. The peak normalized cross-correlation of pitch (C-1 to C-5)
Relevance features	1. Relevance score obtained by WTM 2. Relevance score obtained by VSM 3. Relevance score obtained by LSA 4. Relevance score obtained by MRW

Readability Classification

- Leverage the **LSA (Latent Semantic Analysis)** based language modeling technique to extract “word-readability level”, “word-document” and “word sentence” co-occurrence relationships



- Very Preliminary Results (10-fold tests; w.r.t. classification accuracy (%))

	NHK98 (410 documents)	國編版 (265 documents)
“word-readability level” relationship (dimensionality=6)	0.329	0.260
“word-readability level” & “word-document” relationships (dimensionality=20)	0.346	0.426

Conclusion and Outlook

- Multimedia information (knowledge) access using speech will be very promising in the near future
 - Speech is the key for multimedia understanding and organization
 - Several task domains still remain challenging and warrant further investigation
- ASR technologies are expected to play an essential role in computer-aided (language) learning
 - As to future work, we would like to apply and extend our methods to automatic pronunciation scoring for sub-word (syllable, INITIAL/FINAL) units, and overall pronunciation quality evaluation
 - In addition, we are planning to leverage more state-of-the-art machine learning techniques for CAPT in Mandarin Chinese