# Recent Developments in Speech Retrieval and Related Applications

Berlin Chen (陳柏琳)

Department of Computer Science & Information Engineering

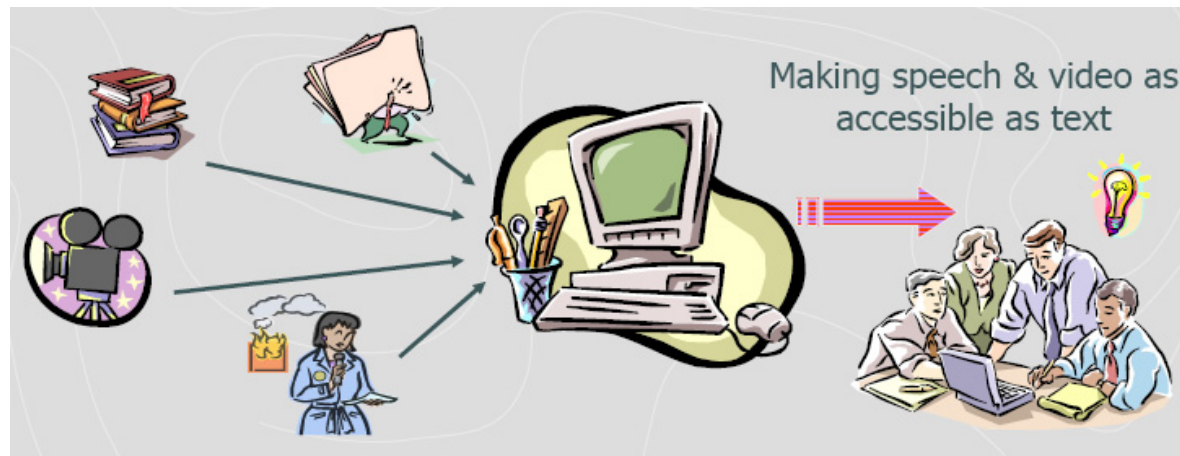National Taiwan Normal University

2012/05/29

# Introduction (1/3)

- Communication and search are by far the most popular activities in our daily lives

  ◦ Speech is the most nature and convenient means of communication between humans, and between humans and machines

    • A spoken language interface could be more convenient than a visual interface on a small 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 linguistic 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 and mining
    - Such as the transcription of spoken documents and recognition of spoken queries
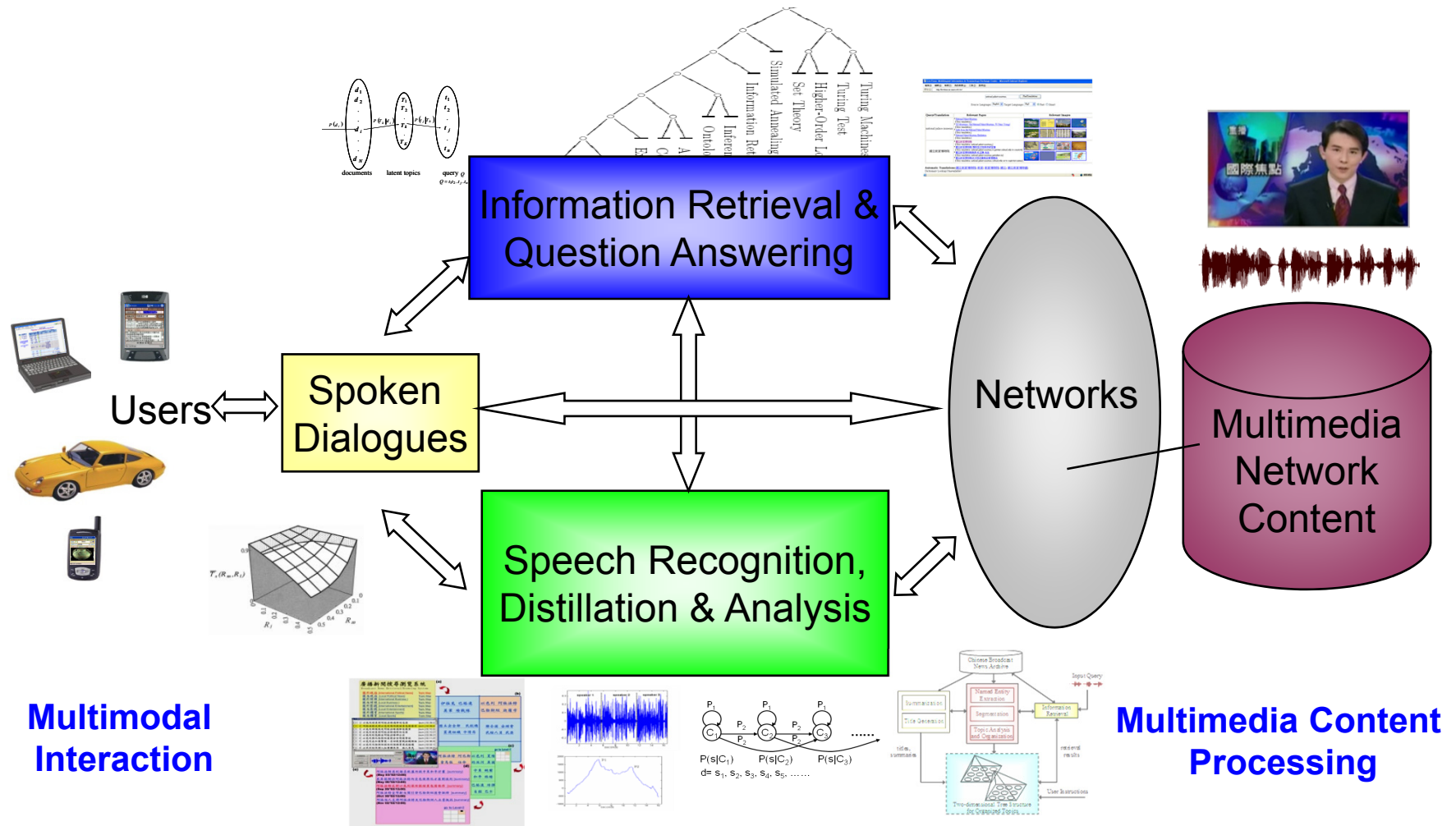


The figure is adapted from the presentation slides of Prof. Ostendorf at *Interspeech 2009*.

# 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 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 pertain 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)

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

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

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

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

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

Feature Extraction & Acoustic Modeling   Language Modeling

- There is an emerging trend of "direct modeling" the discriminant function $P\left(W \mid O\right)$

# 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) variations

M.J.F. Gales and S.J. Young. *The Application of Hidden Markov Models in Speech Recognition. Foundations and Trends in Signal Processing*, 2008

# 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 (to be introduced later)

- Computer-Aided Language Learning (CALL)
  - Speech-driven Interface and multimedia content processing
  - Work in in association with natural language processing techniques
  - Applications
    - Synchronization of audio/video learning materials
    - Automatic pronunciation assessment/scoring
    - Read student essays and grade them
    - Automated reading tutor

- 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
  - MIT Lecture Browser
  - 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*, Audio-Video Indexing System (*MAVIS*)
  - Apple's *Siri* (*QA*)

  *We are witnessing the golden age of ASR!*

IEEE SLTC eNewsletter - Spring 2010 : *Following Global Events with IBM Translingual Automatic Language Exploration System (TALES)*

9

# World-wide Speech Research Projects

- There also are several research projects conducted on a wide array of spoken document processing tasks, e.g.,

  - Rich Transcription Project[1] in the United States (2002-)
    - Creation of recognition technologies that will produce transcriptions which are more readable by humans and more useful for machines

  - TC-STAR Project[2] (Technology and Corpora for Speech to Speech Translation) in Europe (2004-2007)
    - Translation of speeches recorded at European Parliament, between Spanish and English, and of broadcast news by Voice of America, from Mandarin to English

  - "Spontaneous Speech: Corpus and Processing Technology" Project in Japan (1999-2004)
    - 700 hours of lectures, presentations, and news commentaries
    - Automatic transcription, analysis (tagging), retrieval and summarization of spoken documents

# Key Technologies (1/2)

- ## Automatic Speech Recognition (ASR)
  - Automatically convert speech signals into sequences of words or other suitable units for further processing

- ## Spoken Document Segmentation
  - Automatically segment speech signals (or automatically transcribed word sequences) into a set of documents (or short paragraphs) each of which has a central topic

- ## Named Entity Extraction from Spoken Documents
  - Personal names, organization names, location names, event names
  - Very often out-of-vocabulary (OOV) words, difficult for recognition
    - E.g., "蔡煌郎", "九二共識", "烏普薩拉(Uppsala)" etc.

- ## Speech Retrieval
  - Robust representations of spoken queries and spoken documents
  - Matching between (spoken) queries and spoken documents

# Key Technologies (2/2)

- Topic Analysis and Organization for Spoken Documents
  - Analyze the subject topics for (retrieved) documents
  - Organize the subject topics of documents into graphic structures for efficient browsing
- Title Generation for Multi-media/Spoken Documents
  - Automatically generate a title (in text/speech form) for each short document; i.e., a very concise summary indicating the themes of the documents
- Speech Summarization
  - Automatically generate a summary (in text or speech form) for each spoken document or a set of topic-coherent documents
- Information Extraction for Spoken Documents
  - Extraction of key information such as who, when, where, what, why and how for the information described by spoken documents
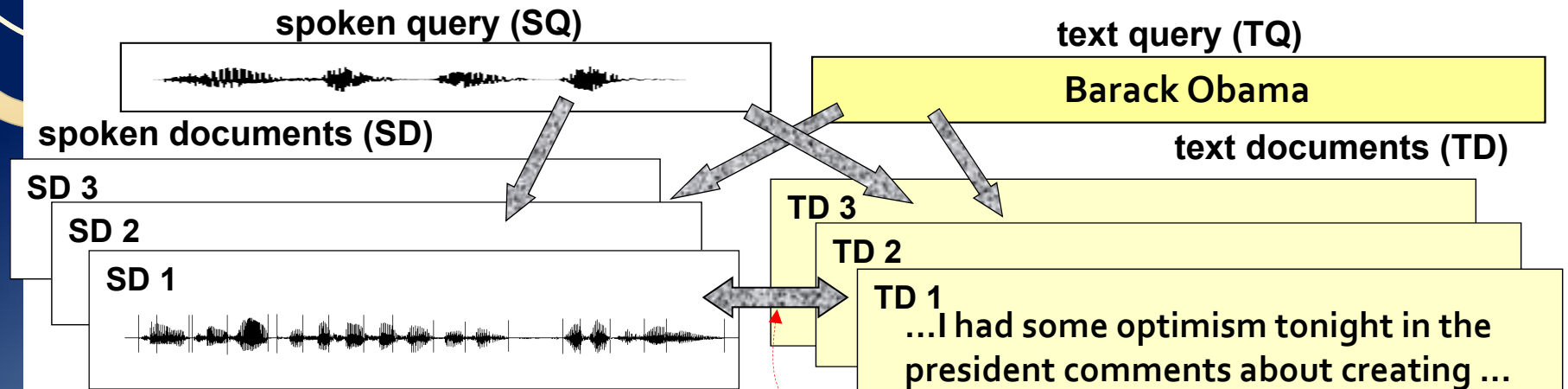
# I. Speech Retrieval

Tur and Mori, *Spoken language understanding – systems for extracting Semantic Information from speech*, Wiley 2011.

# Task Definition of Speech Retrieval

- Robustly Index spoken documents with speech recognition techniques
  - Explore better ways to represent multiple recognition hypotheses of spoken documents beyond the top scoring ones
  - Hybrid of words and subwords (phone/ syllable/ character/ morpheme $n$-grams) for indexing
- Retrieve relevant spoken documents in response to a user query
  - Spoken Document Retrieval (SDR)
    - Find spoken documents/utteramces that are "topically related" to a given query
  - Spoken Term Detection (STD)
    - Find "literally matched" spoken documents where all/most query terms should be present (much like Web search)

# Scenarios of Spoken Document Retrieval

- Scenarios

spoken query (SQ)

text query (TQ)

**Barack Obama**

spoken documents (SD)

text documents (TD)

**SD 3**

**SD 2**

**SD 1**

**TD 3**

**TD 2**

**TD 1**

**...I had some optimism tonight in the president comments about creating ...**

query-by-example



- ◦ SQ/SD is the most difficult
- ◦ SQ/TD and TQ/SD are studied most of the time
  - SQ/TD: viz. voice search (like directory assistance)
  - TQ/SD: e.g., "query-by-exemplar," using text news documents to retrieve relevant broadcast news documents
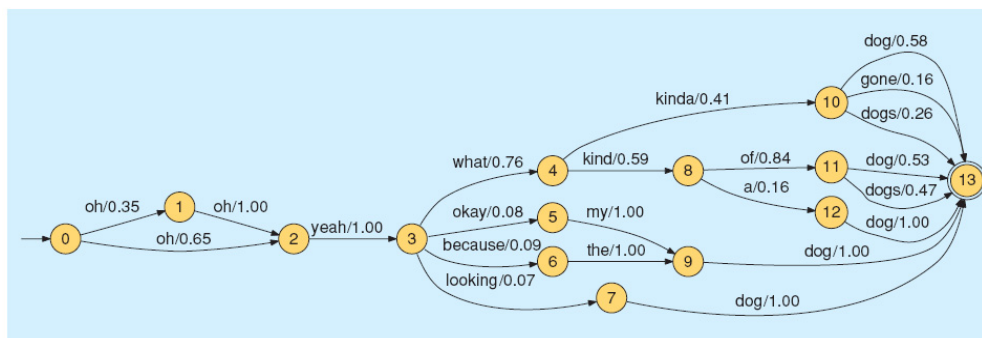    - Useful for news monitoring and tracking

# Representations of Spoken Queries and Documents

- Lattice/confusion network structures for retaining multiple recognition hypotheses

Lattice



Confusion Network



Position-Specific Posterior Probability Lattices

| 0 | | 1 | | 2 | | 3 | | 4 | | 5 | | 6 | | 7 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Oh** | 1.0 | **Yeah** | .65 | **What** | .46 | **Kind** | .27 | Dog | .26 | EOS | .34 | **EOS** | .44 | EOS | .16 |
| — | | Oh | .35 | Yeah | .35 | What | .27 | **Of** | .23 | **Dog** | .29 | Dog | .09 | — | |
| | | — | | Because | .06 | Kinda | .19 | Kind | .16 | Dogs | .13 | Dogs | .06 | | |
| | | | | Okay | .05 | The | .06 | Kinda | .11 | Of | .13 | — | | | |
| | | | | Looking | .05 | My | .05 | Dogs | .05 | A | .03 | | | | |
| | | | | — | | Dog | .05 | EOS | .05 | Gone | .02 | | | | |
| | | | | | | ...... | ... | ...... | ... | — | | | | | |

# Retrieval Models

- Information retrieval (IR) models can be characterized by two different matching strategies

  ○ Literal term matching

    • Match queries and documents in an index term space

  ○ Concept matching

    • Match queries and documents in a latent semantic space

對岸新一代
空軍戰力

*relevant ?*

香港星島日報篇報導引述軍事觀察家的話表示，到二零零五年台灣將完全喪失空中優勢，原因是中國大陸戰機不論是數量或是性能上都將超越台灣，報導指出中國在大量引進俄羅斯先進武器的同時也得加快研發自製武器系統，目前西安飛機製造廠任職的改進型飛豹戰機即將部署尚未與蘇愷三十通道地對地攻擊住宅飛機，以督促遇到挫折的監控其戰機目前也已經取得了重大階段性的認知成果。根據日本媒體報導在台海戰爭隨時可能爆發情況之下北京方面的基本方針，使用高科技答應局部戰爭。因此，解放軍打算在二零零四年前又有包括蘇愷三十二期在內的兩百架蘇霍伊戰鬥機。

# Retrieval Models: Literal Term Matching (1/2)

- Vector Space Model (VSM)
  - Vector representations are used for queries and documents
  - Each dimension is associated with a index term (TF-IDF weighting), describing the intra-/inter-document statistics between all terms and all documents (bag-of-words modeling)
  - Cosine measure for query-document relevance

$$sim\,(D_j, Q)$$

$$= \text{cosine}\,(\Theta) = \frac{\vec{D}_j \bullet \vec{Q}}{|\vec{D}_j| \times |\vec{Q}|}$$

$$= \frac{\sum_{i=1}^{n} w_{i,j} \times w_{i,q}}{\sqrt{\sum_{i=1}^{n} w_{i,j}^2} \times \sqrt{\sum_{i=1}^{n} w_{i,Q}^2}}$$



  - VSM can be implemented with an inverted file structure for efficient document search (instead of exhaustive search)

# Retrieval Models: Literal Term Matching (2/2)

- ## Hidden Markov Model (HMM)
  - Also known as the Language Model (LM)
  - A language model consists of a set of *n*-gram distributions is established for each document for predicting the query

A (unigram) document model

Query

$Q = w_1 w_2 .... w_L$

$\lambda$ → $P(w_i | \mathrm{M}_D)$

$1 - \lambda$

$P(w_i | \mathrm{M}_C)$

$$P_{\mathrm{HMM}}(Q | \mathrm{M}_D) = \prod_{i=1}^{L} [\lambda \cdot P(w_i | \mathrm{M}_D) + (1 - \lambda) \cdot P(w_i | \mathrm{M}_C)]$$

  - Such models can be optimized with different criteria
  - Provide a potentially effective and theoretically attractive probabilistic framework for studying IR problems

Zhai, *Statistical Language Models for Information Retrieval* (Synthesis Lectures Series on Human Language Technologies), Morgan & Claypool Publishers, 2008

# Retrieval Models: Concept Matching (1/2)

- Latent Semantic Analysis (LSA)
  - Start with a matrix ($A$) describing the intra-/inter-document statistics between all terms and all documents
  - Singular value decomposition (SVD) is then performed on the matrix to project all term and document vectors onto a reduced latent topical space $A \approx U \Sigma V^T$
  - Matching between words/queries and documents can be carried out in this latent topical space



$$sim\left(\hat{q}, \hat{d}\right) = coine\ (\hat{q}\Sigma, \hat{d}\Sigma) = \frac{\hat{q}\Sigma^2 \hat{d}^T}{\left|\hat{q}\Sigma\right|\left|\hat{d}\Sigma\right|}$$

Landauer, McNamara, Dennis, and Kintsch (eds.) , *Handbook of Latent Semantic Analysis*, Lawrence Erlbaum, 2007

# Retrieval Models: Concept Matching (2/2)

- Recently, several probabilistic counterparts of LSA have proposed and demonstrated with good success

- Each document as a whole consists of a set of shared latent topics with different weights -- A document topic model (DTM)

  - Each topic offers a unigram (multinomial) distribution for observing a query word

  $$P_{\text{PLSA/LDA}}\left(Q \mid \mathbf{M}_D\right) = \prod_{i=1}^{L}\left[\sum_{k=1}^{K} P\left(w_i \mid T_k\right) P\left(T_k \mid \mathbf{M}_D\right)\right]$$

  - PLSA (Probabilistic Latent Semantic Analysis) and its extension, LDA (Latent Dirichlet Allocation), are the two good examples
    - Mainly differ in inference of model parameters: fixed & unknown vs. Dirichlet distributed

# Word Topic Models (WTM)

- Each word of language is treated as a word topic model (WTM) for predicting the occurrences of other words

$$P_{\text{WTM}}\left(w_i \mid M_{w_j}\right) = \sum_{k=1}^{K} P\left(w_i \mid T_k\right) P\left(T_k \mid M_{w_j}\right)$$

- The relevance measure between a query and a document can be expressed by

$$P_{\text{WTM}}\left(Q \mid D\right) = \prod_{i=1}^{L}\left[\sum_{w_j \in D} P_{\text{WTM}}\left(w_i \mid M_{w_j}\right) P\left(w_j \mid D\right)\right]$$

  - A spoken document can be viewed as a composite WTM
  - WTM is a kind of LM for translating words in the document to words in the query
  - $P\left(w_j \mid D\right)$ is estimated according to the frequency of $w_j$ in $D$

# Unsupervised Training of WTM

- The WTM $P_{\mathrm{WTM}}\left(w_i \mid \mathrm{M}_{w_j}\right)$ of each word can be trained with maximum likelihood estimation (MLE)

  ◦ By concatenating those words occurring within a context window around each occurrence of the word, which are assumed to be relevant to the word, to form the training observation

$$Q_{w_j,1} \qquad Q_{w_j,2} \qquad\qquad Q_{w_j,N} \qquad Q_{w_j} = Q_{w_j,1}, Q_{w_j,2}, \cdots, Q_{w_j,N}$$

$$\boxed{w_j} \quad \text{-----} \quad \boxed{w_j} \quad \text{--------------} \quad \boxed{w_j}$$

$$\log L_{\mathbf{w}} = \sum_{w_j \in \mathbf{w}} \log P_{\mathrm{WTM}}\left(O_{w_j} \middle| \mathrm{M}_{w_j}\right) = \sum_{w_j \in \mathbf{w}} \sum_{w_i \in Q_{w_j}} c\left(w_i, O_{w_j}\right) \log P_{\mathrm{WTM}}\left(w_i \middle| \mathrm{M}_{w_j}\right)$$

  - $\mathbf{W}$ : the set of distinct words in the language

  ◦ WTM was trained to optimize its prediction power over the observation

# Comparison Between WTM and DTM -- Probabilistic Matrix Decompositions

PLSA/LDA

documents

words $\mathbf{A}$ $\approx$ words $\mathbf{G}$

topics

topics $\mathbf{H}^{\mathrm{T}}$

documents

mixture weights

normalized "word-document" co-occurrence matrix

mixture components

$$P_{\mathrm{PLSA/LDA}}\left(w_i \mid \mathrm{M}_D\right) = \sum_{k=1}^{K} P\left(w_i \mid T_k\right) P\left(T_k \mid \mathrm{M}_D\right)$$

WTM

vicinities of words

words $\mathbf{B}$ $\approx$ words $\mathbf{Q}$

topics

topics $\mathbf{Q'}^{\mathrm{T}}$

vicinities of words

mixture weights

normalized "word-word" co-occurrence matrix

mixture components

$$P_{\mathrm{WTM}}\left(w_i \mid \mathrm{M}_{w_j}\right) = \sum_{k=1}^{K} P\left(w_i \mid T_k\right) P\left(T_k \mid \mathrm{M}_{w_j}\right)$$

Chen, "Word topic models for spoken document retrieval and transcription," *ACM TALIP*, 8(1), March 2009

# Example Topic Distributions of WTM

| Topic 13 | | Topic 14 | | Topic 23 | |
|---|---|---|---|---|---|
| **word** | **weight** | **word** | **weight** | **word** | **weight** |
| Vena (靜脈) | 1.202 | Land tax (土地稅) | 0.704 | Cholera (霍亂) | 0.752 |
| Resection (切除) | 0.674 | Tobacco and alcohol tax law (菸酒稅法) | 0.489 | Colorectal cancer (大腸直腸癌) | 0.681 |
| Myoma (肌瘤) | 0.668 | Tax (財稅) | 0.457 | Salmonella enterica (沙門氏菌) | 0.471 |
| Cephalitis (腦炎) | 0.618 | Amend drafts (修正草案) | 0.446 | Aphtae epizooticae (口蹄疫) | 0.337 |
| Uterus (子宮) | 0.501 | Acquisition (購併) | 0.396 | Thyroid (甲狀腺) | 0.303 |
| Bronchus (支氣管) | 0.500 | Insurance law (保險法) | 0.373 | Gastric cancer (胃癌) | 0.298 |

# Some Extensions of DTM and WTM

- ## Hybrid of Word- and Syllable-level Features by using DTM/WTM

documents

topics

documents

words

words

topics

DTM

syllable
pairs

$\mathbf{A}$ $\approx$ syllable
pairs $\mathbf{G}$ $\mathbf{H}^{\mathrm{T}}$

mixture weights

"word-document" &
"syllable pair-document"
co-occurrence matrix

mixture
components

- ## Hybrid of DTM and WTM by Sharing the Same Latent Topics

documents

vicinity
documents

topics

documents

vicinity
documents

words

PLSA WTM $\approx$ words $P(w|T)$ Topics $P(T|D)$ $P(T|\mathrm{M}_W)$

mixture weights

normalized
"word-document" & "word-word"
co-occurrence matrix

mixture
components

Lin and Chen, "Topic modeling for spoken document retrieval using word- and syllable-level information," *SSCS 2009*
Chen et al. "Latent topic modeling of word vicinity information for speech recognition," *ICASSP 2010*

# Supervised Training of WTM

- ## Maximum Likelihood Estimation (MLE)
  - Maximize the log-likelihood of an outside training set of (~800) query exemplars generated by their relevant documents

$$\log L_{\mathbf{Q}_{TrainSet}} = \sum_{Q \in \mathbf{Q}_{TrainSet}} \sum_{D_r \in \mathbf{D}_{R \text{ to } Q}} \log P_{WTM}\left(Q \middle| M_{D_r}\right)$$

- ## Minimum Classification Error Training (MCE)
  - Given a training query exemplar, we can instead minimize the following error function

relevant document        irrelevant document

$$E(Q, D_r, D_{irr}) = \frac{1}{|Q|}\left[ -\log P_{WTM}\left(Q \middle| M_{D_r}\right) + \max_{D_{irr}} \log P_{WTM}\left(Q \middle| M_{D_{irr}}\right) \right]$$

Other irrelevant documents for the training query can be taken into consideration
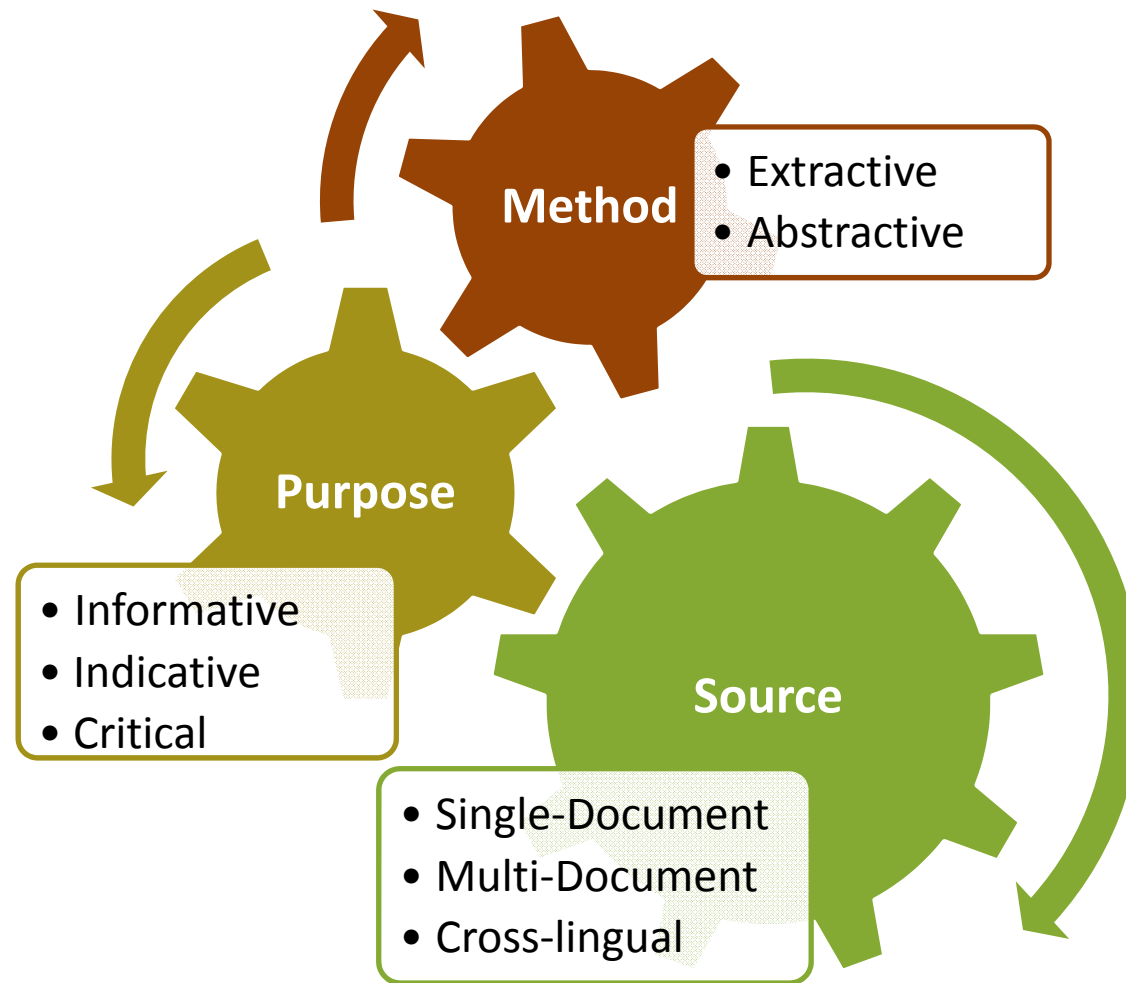
- Further converted to a loss function with a Sigmoid operator
- Corresponding parameters of WTM then are updated with a generalized probabilistic descent (GPD) procedure
- *Learning to rank* !

Associate documents with queries even if they do not share any of the query words!

# II. Speech Summarization

Tur and Mori, *Spoken language understanding – systems for extracting Semantic Information from speech*, Wiley 2011.

# Spectrum of Summarization Research



**Method**
- Extractive
- Abstractive

**Purpose**
- Informative
- Indicative
- Critical

**Source**
- Single-Document
- Multi-Document
- Cross-lingual

# Flowchart of Extractive Speech Summarization

**Pre-processing**

- Speech Detection
- Speaker Identification
- Speech Recognition
- Spontaneous Effect Removal
- Sentence Boundary Detection
- Summarization Unit Selection

**Feature Extraction**

- Structural Info. Extraction
- Prosodic Info. Extraction
- Lexical Info. Extraction
- Acoustic Info. Extraction
- Discourse Info. Extraction

**Summarization**

- Summarization Algorithms

**Post-processing**

- Compaction
- Representation
- Evaluation

# Generic, Extractive Speech Summarization

- A summary is formed by selecting some salient sentences from the original spoken document
- A sentence to be selected as part of the summary is usually being considered by its
  - Significance
    - How importance of the sentence itself
  - Relevance
    - The degree of the similarity between the sentence and other sentences in the document
  - Redundancy
    - The information carried by the candidate summary sentence and the already selected summary sentences should be as dissimilar as possible

# Related Work (1/3)

- **Supervised Machine-Learning Methods** (Significance)
  - The summarization task is usually cast as a two-class sentence-classification problem
  - A sentence is characterized by a set of indicative features
    - Acoustic cues, lexical cues, structural cues or discourse cues
  - Bayesian classifier (BC), support vector machine (SVM), conditional random fields (CRF)
  - Drawbacks
    - "*Bag-of-sentences*" assumption
    - Require manually labeled training data
    - Less generalization capability

# Related Work (2/3)

- **Unsupervised Machine-Learning Methods** (Relevance)
  - Based on the concept of sentence *centrality*
    - Sentences more similar to others are deemed more salient to the main theme of the document
  - Get around the need for manually labeled training data
  - Vector Space Model (VSM), Language Modeling (LM), Graph-based Algorithm (e.g., PageRank)
  - Drawbacks
    - The performance is usually worse than that of supervised summarizers
    - Most of methods constructed solely on the basis of the lexical information
      - Would be adversely affected by imperfect speech recognition

# Related Work (3/3)

- Maximum Marginal Relevance (MMR) (Relevance + Redundancy)
  - Perform sentence selection iteratively with the criteria of topic relevance and coverage
  - A summary sentence is selected according to
    - Whether it is more similar to the whole document than the other sentences (Relevance)
    - Whether it is less similar to the set of sentences selected so far than the other sentences (Redundancy)

$$S_{MMR} = \arg\max_{S_i} \left[ \beta \cdot Sim(S_i, D) - (1 - \beta) \cdot \max_{S' \in \mathbf{Summ}} Sim(S_i, S') \right]$$
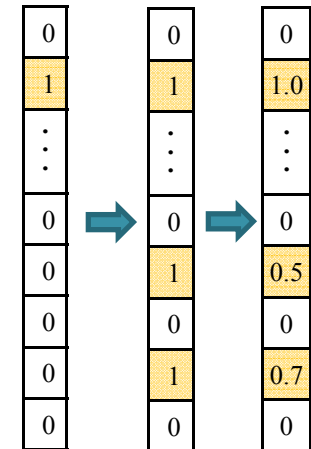
- None of the abovementioned methods fully address these three (Significance, Relevance, Redundancy) factors

# A Risk Minimization Framework (1/4)

- Extractive summarization can be alternatively viewed as **a decision making process**
  - Select a representative subset of sentences or paragraphs from the original documents ➔ action

- **Bayes decision theory** can be employed to guide the summarizer in choosing a course of action
  - It quantifies the tradeoff between
    - Various decisions and the potential cost that accompanies each decision
  - The optimum decision can be made by contemplating each action
    - Choose the action that has the minimum expected risk

Lin and Chen, "A risk minimization framework for extractive speech summarization," *ACL 2010*
Chen and Lin, "A risk-aware modeling framework for speech summarization," IEEE Transactions on Audio, Speech and Language Processing, 2012.

# A Risk Minimization Framework (2/4)

- Without loss of generality, let us denote $\pi \in \mathbf{\Pi}$ as a selection strategy

  ◦ It comprises a set of indicators to address the importance of each sentence $S_i$ in a document $D$ to be summarized

  ◦ The feasible selection strategy can be fairly arbitrary according to the underlying principle

    • E.g., sentence-wise selection vs. list-wise selection

| 0 |
| 1 |
| ⋮ |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |

| 0 |
| 1 |
| ⋮ |
| 0 |
| 1 |
| 0 |
| 1 |
| 0 |

| 0 |
| 1.0 |
| ⋮ |
| 0 |
| 0.5 |
| 0 |
| 0.7 |
| 0 |

- Moreover, we refer to the $k$-th action $a_k$ as choosing the $k$-th selection strategy $\pi_k$, and the observation $O$ as the document $D$

# A Risk Minimization Framework (3/4)

- The expected risk of a certain selection strategy $\pi_k$

$$R(\pi_k \mid D) = \int_\pi L(\pi_k, \pi) p(\pi \mid D) d\pi$$

- Therefore, the ultimate goal of extractive summarization could be stated as

  ◦ The search of the best selection strategy $\pi_{opt}$ from the space of all possible selection strategies that minimizes the expected risk

$$\pi_{opt} = \arg\min_{\pi_k} R(\pi_k \mid D)$$

$$= \arg\min_{\pi_k} \int_\pi L(\pi_k, \pi) p(\pi \mid D) d\pi$$

# A Risk Minimization Framework (4/4)

- Sentence-wise (iterative) selection

$$S^* = \arg\min_{S_i \in \widetilde{D}} R\left(S_i \mid \widetilde{D}\right)$$

$$= \arg\min_{S_i \in \widetilde{D}} \sum_{S_j \in \widetilde{D}} L\left(S_i, S_j\right) P\left(S_j \mid \widetilde{D}\right)$$

- ∘ $\widetilde{D}$ denotes the "residual" document

- By applying the Bayes' rule, the final selection strategy for extractive summarization is stated as

**Relevance/Redundancy**   **Relevance**   **Significance**

$$S^* = \arg\min_{S_i \in \widetilde{D}} \sum_{S_j \in \widetilde{D}} L\left(S_i, S_j\right) \frac{P\left(\widetilde{D} \mid S_j\right) P\left(S_j\right)}{\sum_{S_m \in \widetilde{D}} P\left(\widetilde{D} \mid S_m\right) P\left(S_m\right)}$$

$\pi_2$

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

# Relation to Other Summarization Models

- The use of "0-1" loss function

$$S^* = \underset{S_i \in \widetilde{D}}{\arg\max} \frac{P\left(\widetilde{D} \mid S_i\right)P(S_i)}{\sum\limits_{S_m \in \widetilde{D}} P\left(\widetilde{D} \mid S_m\right)P(S_m)} = \underset{S_i \in \widetilde{D}}{\arg\max} \, P\left(\widetilde{D} \mid S_i\right)P(S_i)$$

  ◦ A natural integration of the supervised and unsupervised summarizers

- Uniform prior distribution

  ◦ Estimate the relevance between the document and sentence using $P\left(\widetilde{D} \mid S_i\right)$

- Equal document-likelihood

  ◦ Sentences are selected solely based on the prior probability $P(S_i)$

# Implementation Details (1/4)

- ## Sentence Generative Model $P\left(\widetilde{D} \mid S_i\right)$

  - We explore the language modeling (LM) approach

    - Each sentence is simply regarded as a probabilistic generative model consisting of a unigram distribution for generating the document

    $$P\left(\widetilde{D}\middle|S_i\right) = \prod_{w \in \widetilde{D}} P\left(w\middle|S_i\right)^{c\left(w,\widetilde{D}\right)}$$

    - Maximum Likelihood Estimation (MLE) of $P\left(w\middle|S_i\right)$

      - It may suffer from the problem of unreliable model estimation
      - Enhanced via topic modeling (PLSA, LDA, WTM, etc.)
      - Enhanced by incorporating relevance information

Chen et al., "A probabilistic generative framework for extractive broadcast news speech summarization," *IEEE TASLP, 17(1), 2009*

# Implementation Details (2/4)

- **Sentence Prior Model** $P(S_i)$

  ◦ We assume the sentence prior probability is in proportion to the posterior probability of a sentence being included in the summary class

$$P(S_i) \approx \frac{p(X_i \mid \mathbf{S})P(\mathbf{S})}{P(X_i \mid \mathbf{S})P(\mathbf{S}) + P(X_i \mid \overline{\mathbf{S}})P(\overline{\mathbf{S}})}$$

  - $\mathbf{S}$ and $\overline{\mathbf{S}}$ : summary and non-summary classes
  - $X_i$ : a set of indicative (prosodic/lexical/structural) features used for representing sentence $S_i$

  ◦ Several popular supervised classifiers can be leveraged for this purpose
  - Bayesian Classifier (BC), Support Vector Machine (SVM), etc.

# Implementation Details (3/4)

- **Loss Function**
  - VSM-based loss function $L(S_i, S_j)$
    - We use the "*TF-IDF*" weighting to calculate the cosine similarity
    - If a sentence is more dissimilar from most of the other sentences, it may incur a higher loss

$$L(S_i, S_j) = 1 - Sim(S_i, S_j)$$

  - MMR-based loss function
    - Additionally address the "redundancy" issue

$$L(S_i, S_j) = 1 - \left[ \beta \cdot Sim(S_i, S_j) - (1 - \beta) \cdot \max_{S' \in \mathbf{Summ}} Sim(S_i, S') \right]$$

    - **Summ** the set of already selected summary sentences

# Summarization Experiments (1/4)

- ## MATBN corpus
  - A subset of 205 broadcast news documents was reserved for the summarization experiments
  - The average Chinese character error rate (CER) is about 35%
  - Three subjects were asked to create summaries of the 205 spoken documents
  - The assessment of summarization performance is based on the widely-used ROUGE measure

| | ROGUE-1 | ROUGE-2 | ROUGE-L |
|---|---|---|---|
| Agreement | 0.600 | 0.532 | 0.527 |

*The agreement among the subjects for important sentence ranking for the evaluation set.

# Summarization Experiments (2/4)

$$S^* = \arg\min_{S_i \in \tilde{D}} \sum_{S_j \in \tilde{D}} L(S_i, S_j) \frac{P(\tilde{D} \mid S_j) P(S_j)}{\sum_{S_m \in \tilde{D}} P(\tilde{D} \mid S_m) P(S_m)}$$

- **Baseline experiments**
  - Supervised summarizer – A Bayesian classifier (BC) with 28 indicative features determines the sentence prior probability $P(S_i)$
  - Unsupervised summarizer –A (unigram) language modeling approach determines the document-likelihood $P(D \mid S_i)$

| | Text Document (TD) | | | Spoken Document (SD) | | |
|---|---|---|---|---|---|---|
| | ROGUE-1 | ROUGE-2 | ROUGE-L | ROGUE-1 | ROUGE-2 | ROUGE-L |
| BC | 0.445 (0.390 - 0.504) | 0.346 (0.201 - 0.415) | 0.404 (0.348 - 0.468) | 0.369 (0.316 - 0.426) | 0.241 (0.183 - 0.302) | 0.321 (0.268 - 0.378) |
| LM | 0.387 (0.302 - 0.474) | 0.264 (0.168 - 0.366) | 0.334 (0.251 - 0.415) | 0.319 (0.274 - 0.367) | 0.164 (0.115 - 0.224) | 0.253 (0.215 - 0.301) |

  - Erroneous transcripts cause significant performance degradation
  - BC outperforms LM
    - BC is trained with the handcrafted document-summary data
    - BC utilizes a rich set of features

# Summarization Experiments (3/4)

$$S^* = \arg\min_{S_i \in \tilde{D}} \sum_{S_j \in \tilde{D}} L(S_i, S_j) \frac{P(\tilde{D} \mid S_j) P(S_j)}{\sum_{S_m \in \tilde{D}} P(\tilde{D} \mid S_m) P(S_m)}$$

- Experiments on proposed methods

| Prior | Loss | Text Document (TD) | | | Spoken Document (SD) | | |
|---|---|---|---|---|---|---|---|
| | | ROGUE-1 | ROUGE-2 | ROUGE-L | ROGUE-1 | ROUGE-2 | ROUGE-L |
| | 0-1 | 0.501 | 0.401 | 0.459 | 0.417 | 0.281 | 0.356 |
| BC | SIM | 0.524 | 0.425 | 0.473 | 0.475 | 0.351 | 0.420 |
| | MMR | 0.529 | 0.426 | 0.479 | 0.475 | 0.351 | 0.420 |

- Simple "0-1 Loss" gives about 4-5% absolute improvements as compared to the results of BC
- "SIM/MMR Loss" results in higher performance
  - MMR is slightly better than SIM
- The performance gaps between the TD and SD cases are reduced to a good extent

# Summarization Experiments (4/4)

$$S^* = \arg\min_{S_i \in \tilde{D}} \sum_{S_j \in \tilde{D}} L(S_i, S_j) \frac{P(\tilde{D}|S_j)P(S_j)}{\sum_{S_m \in \tilde{D}} P(\tilde{D}|S_m)P(S_m)}$$

- Experiments on proposed methods

| Prior | Loss | Text Document (TD) | | | Spoken Document (SD) | | |
|---|---|---|---|---|---|---|---|
| | | ROGUE-1 | ROUGE-2 | ROUGE-L | ROGUE-1 | ROUGE-2 | ROUGE-L |
| Uniform | SIM | 0.405 | 0.281 | 0.348 | 0.365 | 0.209 | 0.305 |
| | MMR | 0.417 | 0.282 | 0.359 | 0.391 | 0.236 | 0.338 |

- Assume the sentence prior probability $P(S_i)$ is uniformly distributed
  - The importance of a sentence is considered from two angles
    - Relationship between a sentence and the whole document
    - Relationship between the sentence and the other individual sentences
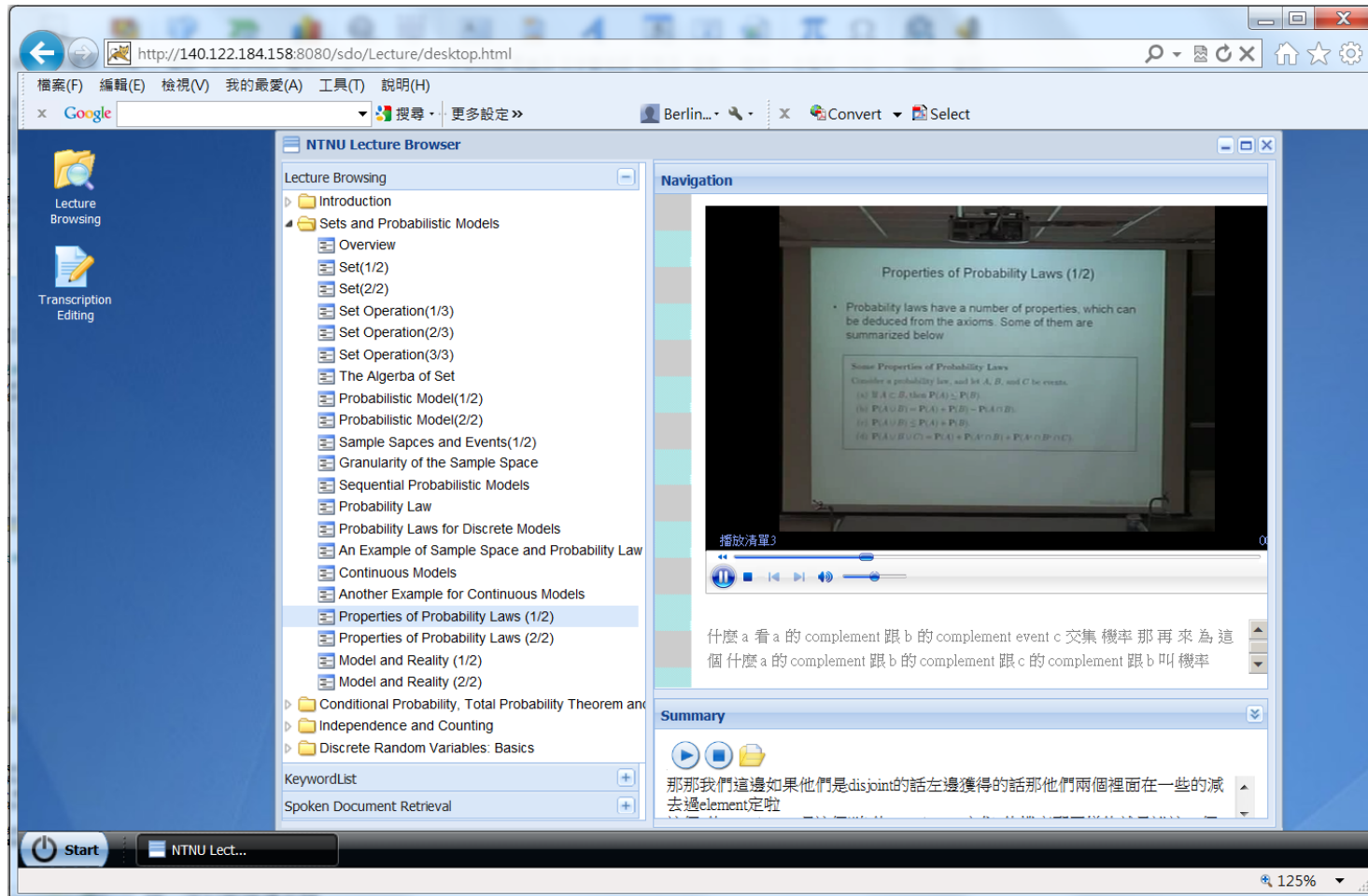- Additional consideration of the "sentence-sentence" relationship appears to be beneficial

# Future Work on Speech Summarization

- Look for different selection strategies
  - e.g., the listwise strategy

$$Summary = \arg\min_{\psi_i \in \Psi_D} \sum_{\psi_j \in \Psi_D} L(\psi_i, \psi_j) \frac{P(D \mid \psi_j)P(\psi_j)}{\sum_{\psi_m \in \Psi_D} P(D \mid \psi_m)P(\psi_m)}$$

- Explore different modeling approaches and indicative features for the component models

- Investigate discriminative training criteria for training the component models

- Robustly represent the recognition hypotheses of spoken documents beyond the top scoring ones

- Extend and apply the proposed framework to multi-document summarization tasks

- ...

Chen et al., "Extractive speech summarization using evaluation metric-related training criteria," *Information Processing & Management*, available online 16 January 2012.

# NTNU Lecture/News Browsing System

# Conclusions

- Multimedia information access (over the Web) 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

- Voice search provides good assistance for companies, for instance, in
  - Contact (Call)-center conservations: monitor agent conduct and customer satisfaction, increase service efficiency
  - Content-providing services: such as MOD (Multimedia on Demand): provide a better way to retrieve and browse descried program contents

- Speech processing technologies are expected to play an essential role in computer-aided (language) learning

*Thank You!*

# More on Language Modeling for IR/SDR

- LM approaches have been introduced to IR (and SDR), and demonstrated with good success

$$P_{\mathrm{LM}}(D|Q) = \frac{P(Q|D)P(D)}{P(Q)} \propto P(Q|D)$$

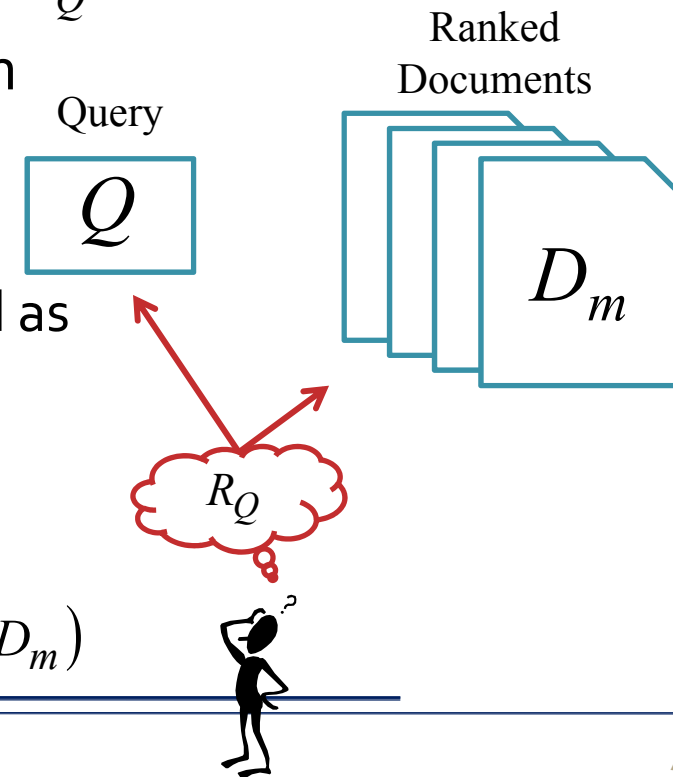- The Kullback-Leibler (KL)-Divergence measure is another basic formulation of LM for IR

$$\mathrm{KL}(Q\|D) = \sum_{w \in V} P(w|Q)\log\frac{P(w|Q)}{P(w|D)} \propto -\sum_{w \in V} P(w|Q)\log P(w|D)$$

  ◦ A query is treated as a probabilistic model rather than simply an observation
  ◦ KL-divergence supports us to improve not only the document model but also the **query model** for better document ranking

# Relevance Modeling (RM)

- In the conventional relevance modeling

  ◦ Each query $Q$ is assumed to be associated with an unknown relevance class $R_Q$

  ◦ Documents that are relevant to the information need expressed in the query are samples drawn from $R_Q$

- The document ranking problem can be reduced to determine the probability $P_{\mathrm{RM}}(w\,|\,Q)$

  ◦ The relevance model can be defined as the probability of the word selected from relevance documents

$$P_{\mathrm{RM}}(w\,|\,Q) \propto \sum_{m=1}^{M} P(D_m)P(q_1,\ldots,q_L,w\,|\,D_m)$$

$$= \sum_{m=1}^{M} P(D_m)P(w\,|\,D_m)\prod_{l=1}^{L} P(q_l\,|\,D_m)$$

Query

$$Q$$

Ranked Documents

$$D_m$$

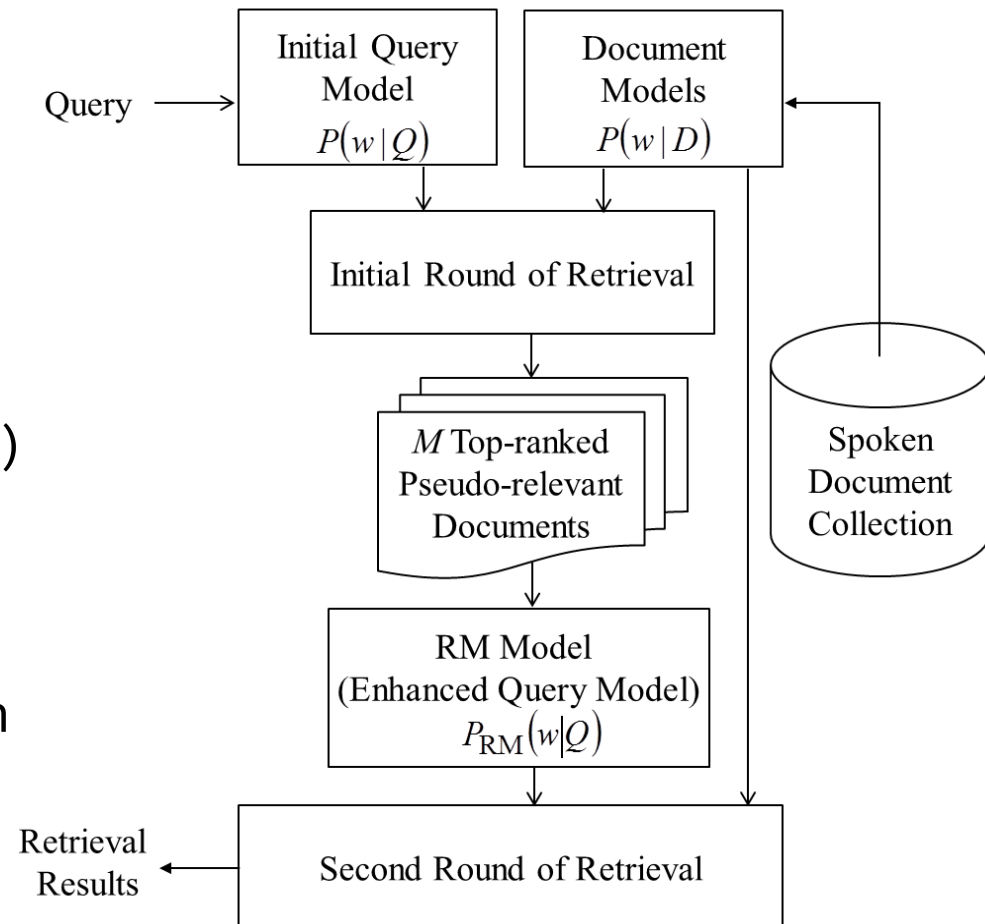$$R_Q$$

# Incorporating Topical Information in RM

- Topic-based relevance model (TRM) makes a step forward by incorporating latent topic information into RM
  - As conventional topic models, the probability that a word occurs is estimated from a set of latent topics

$$P_{\text{TRM}}(w \mid Q) \propto \Sigma_{m=1}^{M} \Sigma_{k=1}^{K} P(D_m) P(T_k \mid D_m) P(w \mid T_k) \Pi_{l=1}^{L} P(q_l \mid T_k)$$

- TRM has some assumptions and properties:
  - Relevant documents are assumed to share a set of pre-defined latent topic variables $\{T_1, \ldots, T_K\}$
  - When given a latent topic, words and documents are independent of each other
  - TRM assumes that the additional cues of how words are distributed across a set of latent topics can carry useful global topic structure for relevance modeling

# Inference of the Various Relevance Models

- In practice, the relevant documents are unknown in advance
  - First-round retrieval with the "query-likelihood" LM approach is applied to obtain a set of top-ranked (pseudo-relevant) documents to approximate the relevance class
  - Second-run retrieval with the "KL-divergence" measure is used to re-rank the spoken documents

Query $\longrightarrow$

| Initial Query Model $P(w\,|\,Q)$ | Document Models $P(w\,|\,D)$ |
|---|---|

Initial Round of Retrieval

$M$ Top-ranked Pseudo-relevant Documents

Spoken Document Collection

RM Model (Enhanced Query Model) $P_{\mathrm{RM}}(w|Q)$

Retrieval Results $\longleftarrow$ Second Round of Retrieval

# Incorporating Non-Relevance Information (1/2)

- Further, in addition to using the relevance information, we also hypothesize that the non-relevant (low-ranked) documents can provide extra useful cues

- For this idea to work, we attempt to estimate a non-relevance model $P(w \mid \mathrm{NR}_Q)$ for each test query $Q$
  - The non-relevance model can be estimated simply based on the ML criterion or be further optimized with the EM algorithm
  - E-step:

$$P(\mathrm{NR}_Q \mid w) = \frac{\lambda \cdot P(w \mid \mathrm{NR}_Q)}{\lambda \cdot P(w \mid \mathrm{NR}_Q) + (1 - \lambda) \cdot P(w \mid \mathrm{BG})}$$

  - M-step

$$P(w \mid \mathrm{NR}_Q) = \frac{\sum_{D' \in \mathbf{D}_{\mathrm{Low}}} c(w, D') \cdot P(\mathrm{NR}_Q \mid w)}{\sum_w \sum_{D' \in \mathbf{D}_{\mathrm{Low}}} c(w, D') \cdot P(\mathrm{NR}_Q \mid w)}$$

# Incorporating Non-Relevance Information (2/2)

- The similarity measure between query $Q$ and any document $D$ thus can be computed as follows:

$$\mathrm{SIM}(Q,D) = -\mathrm{KL}(Q \| D) + \alpha \cdot \mathrm{KL}(\mathrm{NR}_Q \| D)$$

**Relevance Information**     **Penalty Factor**    **Non-Relevance Information**

- Note also that
  - Here we adopt an unsupervised way to estimate the non-relevance model
  - We intend to explore whether the relevance and non-relevance cues of a test query can conspire to enhance the SDR performance

# Experimental Results on the TDT2 Collection

- we investigate the joint exploration of relevance and non-relevance cues for query modeling

- We also consider using different levels of index features (viz. word-level features, syllable-level features and their combination) for SDR

| SD | Word | Syllable | Combination |
|----|------|----------|-------------|
| ULM | 0.323 | 0.330 | - |
| RM | 0.364 | 0.378 | 0.396 |
| TRM | 0.394 | 0.383 | 0.412 |
| RM+NR | 0.392 | 0.405 | 0.426 |
| TRM+NR | 0.402 | 0.415 | 0.441 |

Chen et al., "Query modeling for spoken document retrieval," *ASRU 2011*