# Information Retrieval and Extraction
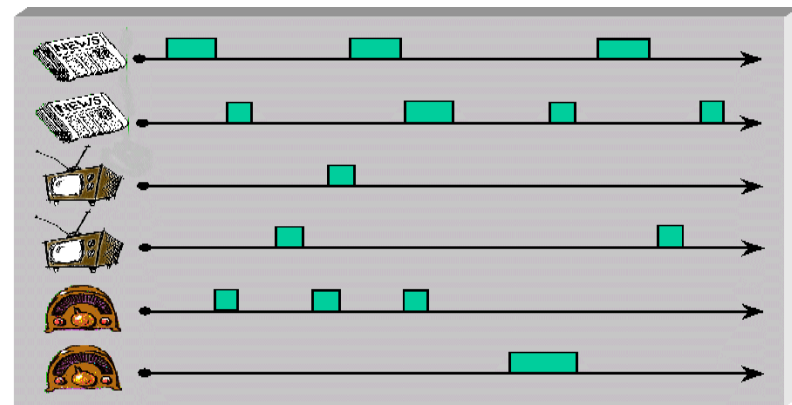
Berlin Chen 2007
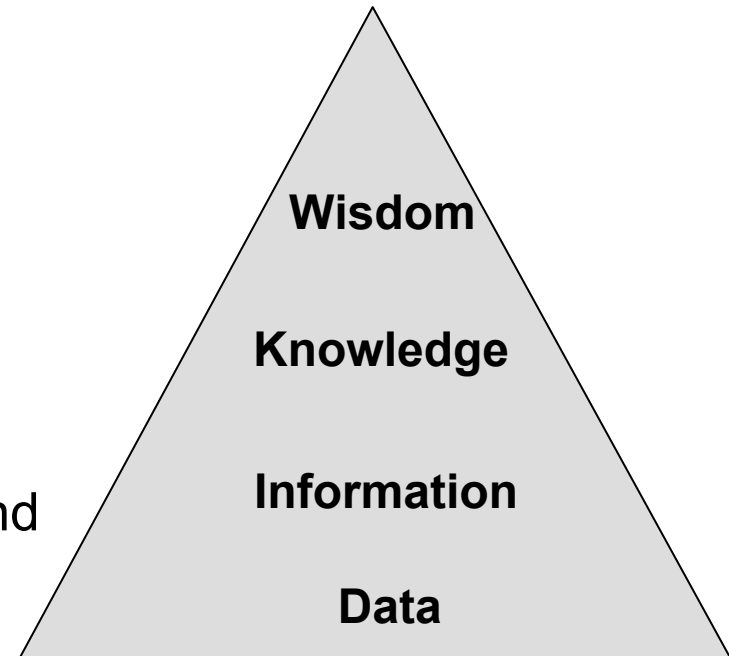
(Picture from the TREC web site)

# Textbook and References

- Textbook
  - R. Baeza-Yates and B. Ribeiro-Neto. *Modern Information Retrieval.* Addison Wesley Longman, 1999

- References
  - D. A. Grossman, O. Frieder, *Information Retrieval: Algorithms and Heuristics*, Springer. 2004
  - Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, *Introduction to Information Retrieval*, Cambridge University Press. 2008
  - W. B. Croft and J. Lafferty (Editors). *Language Modeling for Information Retrieval*. Kluwer-Academic Publishers, July 2003
  - I. H. Witten, A. Moffat, and T. C. Bell. *Managing Gigabytes: Compressing and Indexing Documents and Images*. Morgan Kaufmann Publishing, 1999
  - C. Manning and H. Schutze. *Foundations of Statistical Natural Language Processing*. MIT Press, 1999

# Motivation (1/2)

- Information Hierarchy
  - Data
    - The raw material of information
  - Information
    - Data organized and presented by someone
  - Knowledge
    - Information read, heard or seen and understood
  - Wisdom
    - Distilled and integrated knowledge and understanding

**Wisdom**

**Knowledge**

**Information**

**Data**

# Motivation (2/2)

- User information need
  - Find all docs containing information on college tennis teams which:

    (1) are maintained by a USA university and

    (2) participate in the NCAA tournament

    (3) National ranking in last three years and
        contact information

Query

Emphasis is on the retrieval of
information (not data)

Search engine/IR system

# Information Retrieval

- Deal with the representation, storage, organization of, and access to information items

- Focus is on the user information need
  - Information about a subject or topic
  - Semantics is frequently loose
  - Small errors are tolerated

- Handle natural language text which is not always well structured and could be semantically ambiguous

# Data Retrieval

- Determine which document of a collection contain the *keywords* in the user query

- Retrieve all objects (attributes) which satisfy clearly defined conditions in a regular expression or a relational algebra expression
    - Which documents contain a set of keywords?
    - Well defined semantics
    - A single erroneous object implies failure!

# IR system

- Interpret contents of information items (docs)

- Generate a ranking which reflects relevance

- Notion of *relevance* is most important

# IR at the Center of the Stage

- IR in the last 20 years:
  - Modelng, classification, clustering, filtering
  - User interfaces and visualization
  - Systems and languages

- WWW environment (90~)
  - Universal repository of knowledge and culture
  - Without frontiers: free universal access
  - Lack of well-defined data model

# IR Main Issues

- The effective retrieval of relevant information affected by

  - The user task

  - Logical view of the documents

# The User Task

- Translate the information need into a query in the language provided by the system
  - A set of words conveying the semantics of the information need

- Browse the retrieved documents

Retrieval

Browsing

1. Doc *i*
2. Doc *j*
3. Doc *k*

F1 racing

Directions to
Le Mans
Tourism in France

Information Records

# Logical View of the Documents (1/2)

- A full text view (representation)
    - Represent document by its whole set of words
        - Complete but higher computational cost

- A set of index terms by a human subject
    - Derived automatically or generated by a specialist
        - Concise but may poor

- An intermediate representation with feasible *text operations*

# Logical View of the Documents (2/2)

- ## Text operations
  - Elimination of stop-words (e.g. articles, connectives, …)
  - The use of stemming (e.g. tense, …)
  - The identification of noun groups
  - Compression ….

- ## Text structure (chapters, sections, …)

# Different Views of the IR Problem

- Computer-centered (commercial perspective)
  - Efficient indexing approaches
  - High performance matching ranking algorithms

- Human-centered (academic perceptive)
  - Studies of user behaviors
  - Understanding of user needs

Library science
psychology
….

# IR for Web and Digital Libraries

- Questions should be addressed
  - Still difficult to retrieve information relevant to user needs
  - Quick response is becoming more and more a pressing factor (Precision vs. Recall)
  - The user interaction with the system (HCI, Human Computer Interaction)

- Other concerns
  - Security and privacy
  - Copyright and patent

# The Retrieval Process (1/2)

# The Retrieval Process (2/2)

- In current retrieval systems

  - Users almost never declare his information need

    - Only a short queries composed few words
      (typically fewer than 4 words)

  - Users have no knowledge of the text or query operations

  Poor formulated queries lead to poor retrieval !

# Major Topics (1/2)

- Four Main Topics



**Figure 1.4** Topics which compose the book and their relationships.

# Major Topics (2/2)

- ## Text IR
  - Retrieval models, evaluation methods, indexing

- ## Human-Computer Interaction (HCI)
  - Improved user interfaces and better data visualization tools

- ## Multimedia IR
  - Text, speech, audio and video contents
  - Multidisciplinary approaches
  - Can multimedia be treated in a unified manner?

- ## Applications
  - Web, bibliographic systems, digital libraries

# Textbook Topics

# Text Information Retrieval (1/4)

- Internet searching engine



**Web**

**Spider**

**Mirrored Web Page Repository**

**Indexer**

Queries

Ranked Docs

**Search Engine**

# Text Information Retrieval (2/4)

- http://www.google.com

# Text Information Retrieval (3/4)

- http://www.openfind.com.tw (Service is No Longer Available)

# Text Information Retrieval (4/4)

- http://www.baidu.com

# Speech Information Retrieval (1/4)



speech information

Text-to-Speech Synthesis

text information

speech

Spoken Dialogue

Information Retrieval

Internet

Public Services/ Information/ Knowledge

Private Services/ Databases/ Applications

text, image, video, speech, …

speech 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 ....

# Speech Information Retrieval (2/4)

- ## HP Research Group – Speechbot System
  (Service is No Longer Available)

  – Broadcast news speech recognition, Information retrieval, and topic segmentation (SIGIR2001)

  – Currently indexes **14,791 hours of content** (2004/09/22,

# Speech Information Retrieval (3/4)

・輸入聲音問句：“請幫我查總統府升旗典禮”



中文影音多媒體資訊檢索雛形展示系統。

# Speech Information Retrieval (4/4)

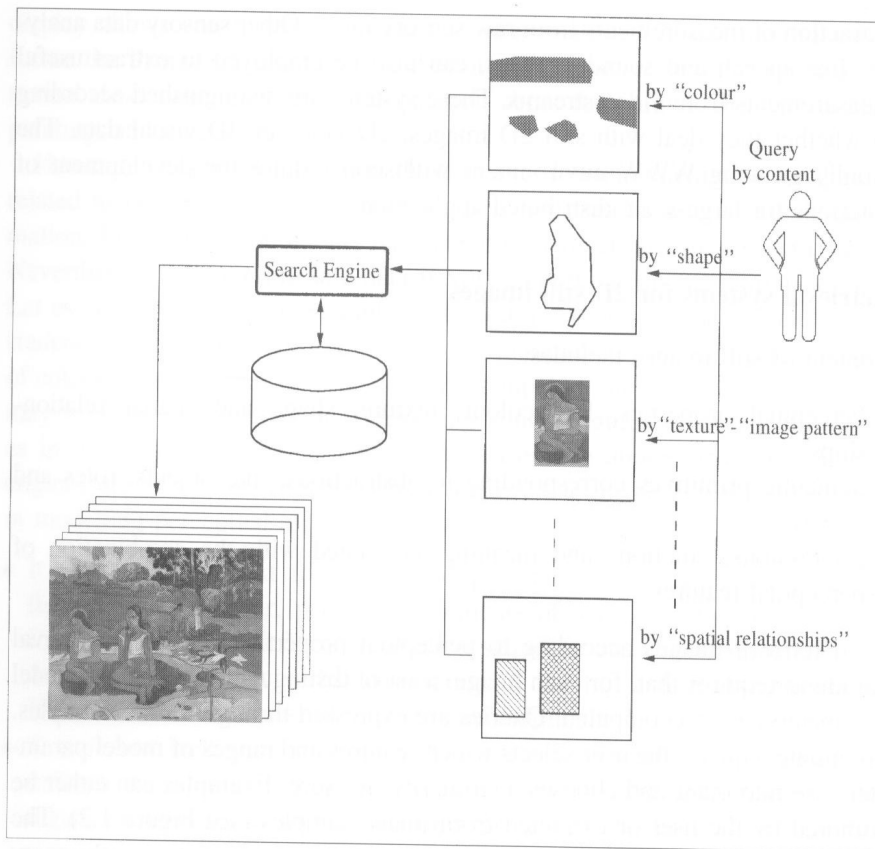# Visual Information Retrieval (1/4)

- Content-based approach



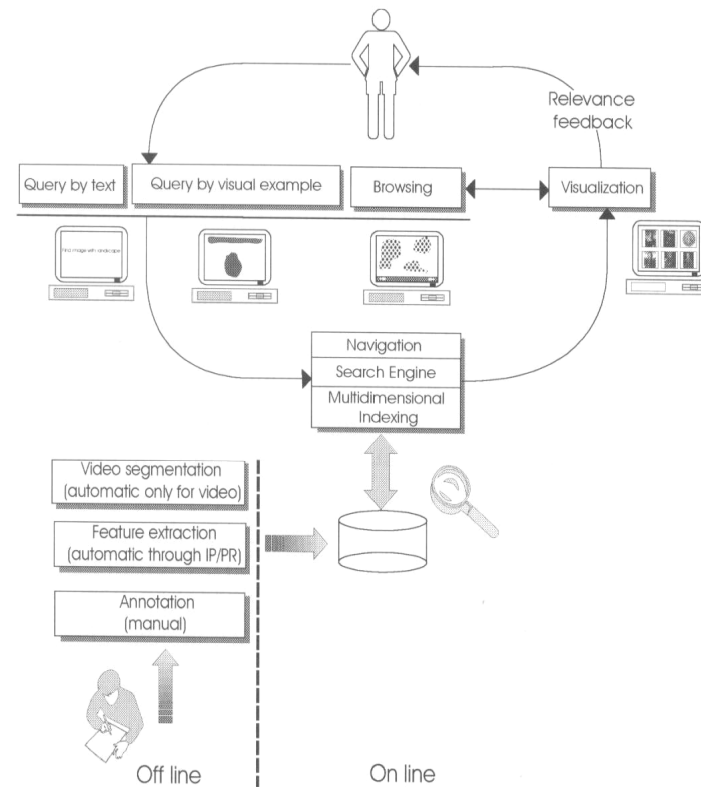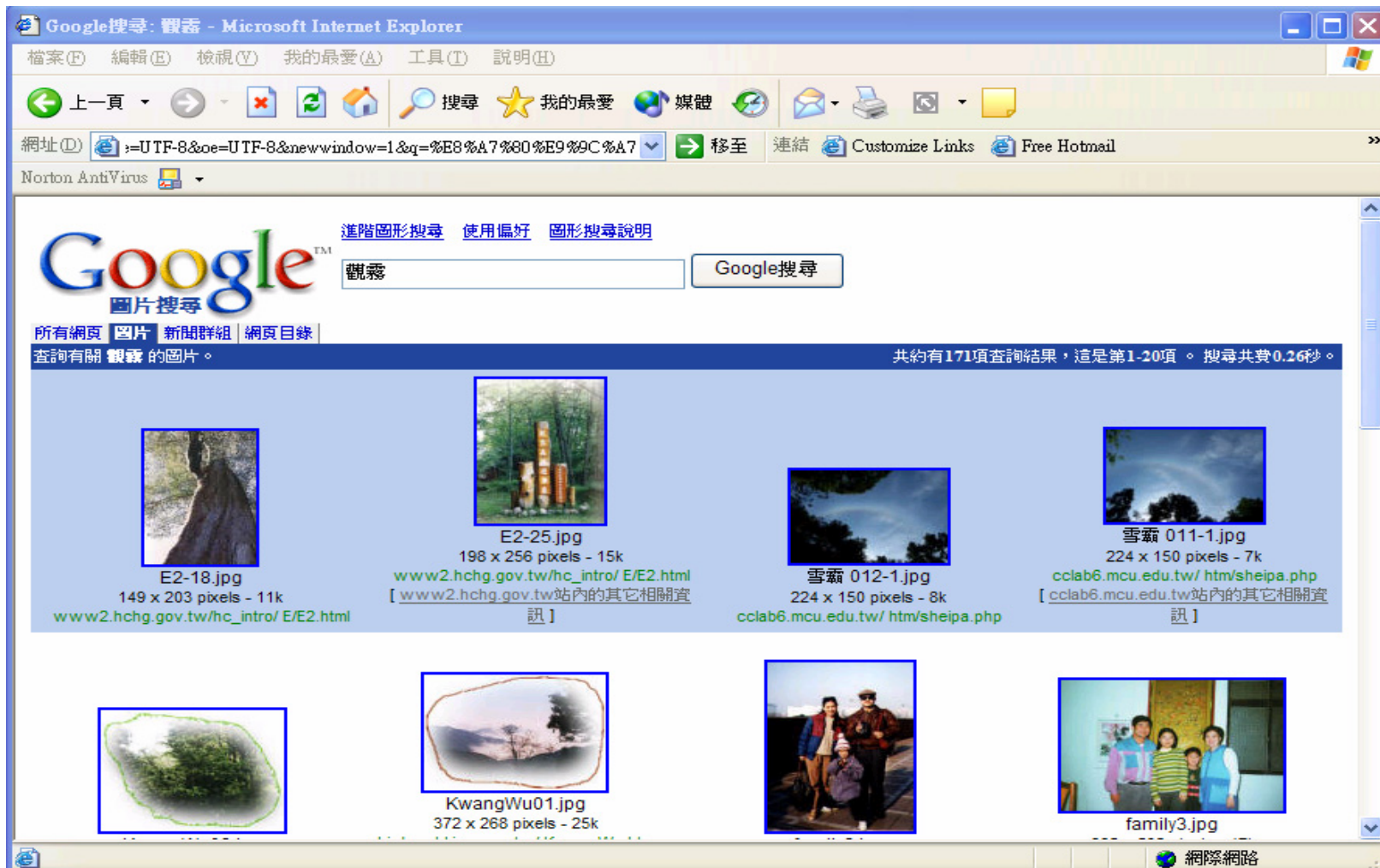Figure 1.2 Different types of query by example.



Figure 1.5 Sketch of a new-generation visual information retrieval system for video.

# Visual Information Retrieval (2/4)
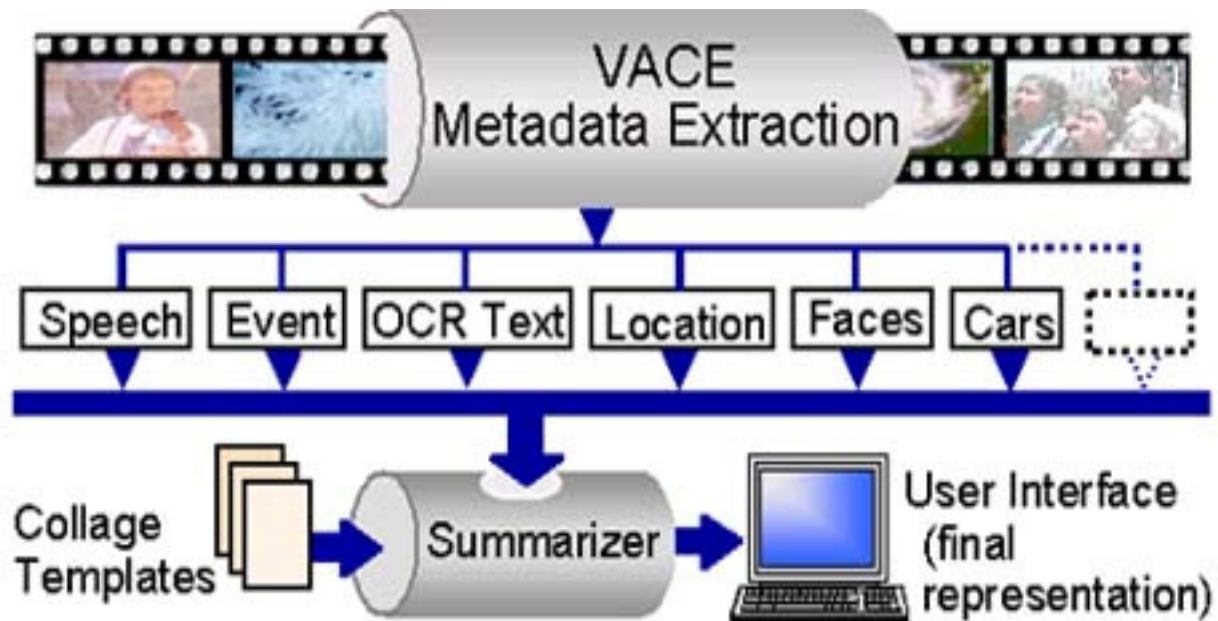
- Images with Texts

# Visual Information Retrieval (3/4)

- Content-based Image Retrieval

# Visual Information Retrieval (4/4)

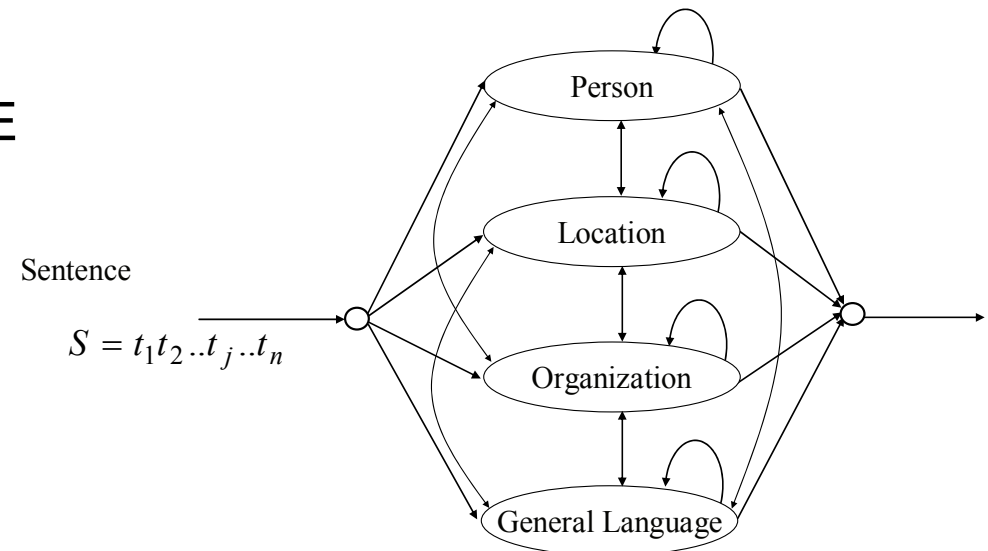**Video Analysis and Content Extraction**

# Other IR-Related Tasks

- Information filtering and routing
- **Term/Document categorization**
- **Term/Document clustering**
- **Document summarization**
- **Information extraction**
- Question answering
- Crosslingual information retrieval
- …..

# Document Summarization

- Audience
  - Generic summarization
  - User-focused summarization
    - Query-focused summarization
    - Topic-focused summarization

- Function
  - Indicative summarization
  - Informative summarization

- Extracts vs. abstracts
  - Extract: consists wholly of portions from the source
  - Abstract: contains material which is not present in the source

- Output modality
  - Speech-to-text summarization
  - Speech-to-speech summarization

- Single vs. multiple documents

# Information Extraction

- E.g., Named-Entity Extraction
  - NE has it origin from the Message Understanding Conferences (MUC) sponsored by U.S. DARPA program
    - Began in the 1990's
    - Aimed at extraction of information from text documents
    - Extended to many other languages and spoken documents (mainly broadcast news)

  - Common approaches to NE
    - Rule-based approach
    - Model-based approach
    - Combined approach

Sentence

$$S = t_1 t_2 .. t_j .. t_n$$

Person

Location

Organization

General Language

# Cross-lingual Information Retrieval

- E.g., Automatic Term Translation
  - Discovering translations of unknown query terms in different languages
  - E.g., The Live Query Term Translation System (LiveTrans) developed at Academia Sinica/by Dr. Chien Lee-Feng



Machine-Extracted Translation

# Multidisciplinary Approaches



**Natural Language Processing**

**Multimedia Processing**

**Networking**

**IR**

**Machine Learning**

**Artificial Intelligence**

# Resources

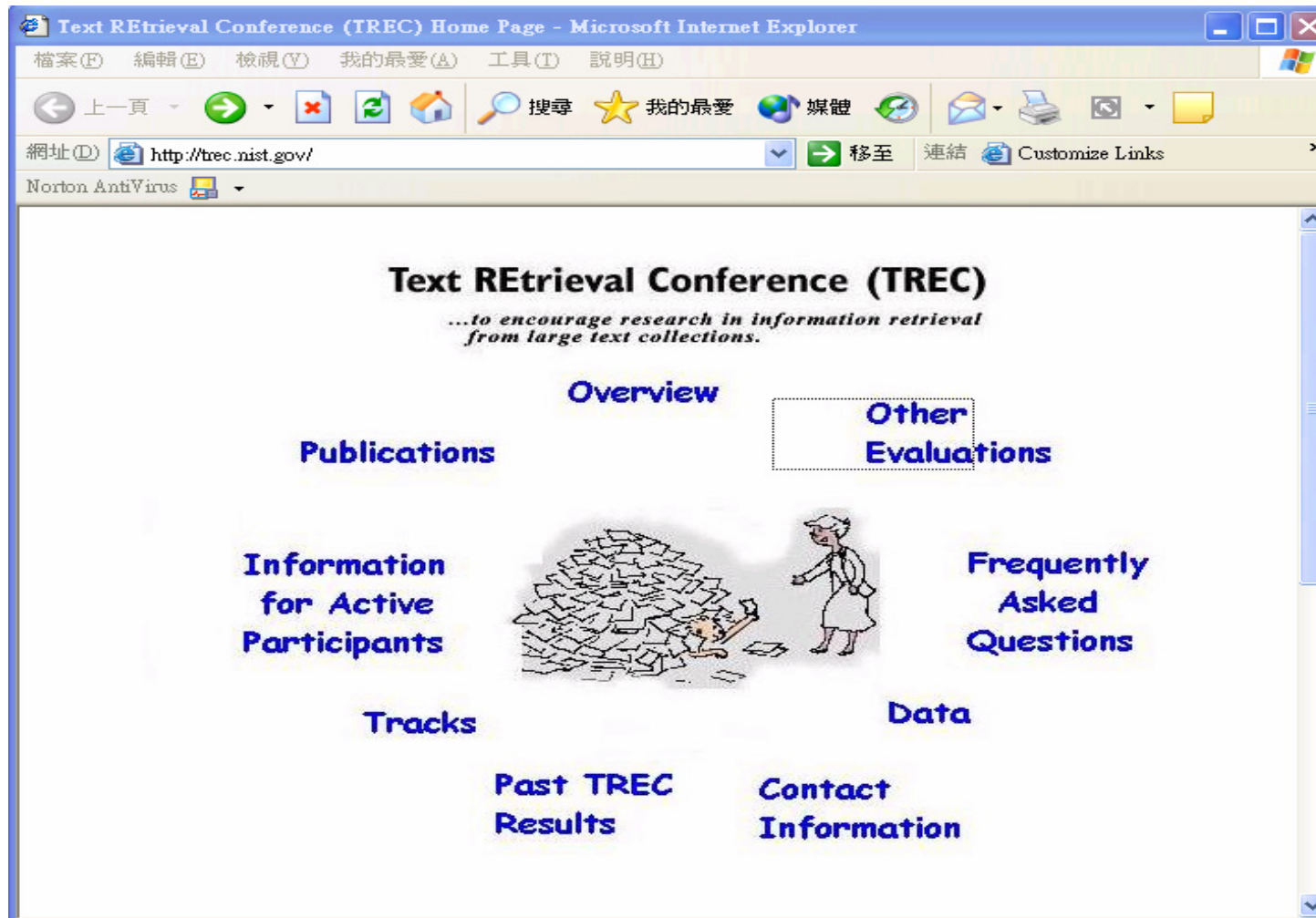- Corpora (Speech/Language resources)
  - Refer speech waveforms, machine-readable text, dictionaries, thesauri as well as tools for processing them
    - LDC - Linguistic Data Consortium

# Contests (1/2)

- [Text REtrieval Conference](#) (TREC)

# Contests (2/2)

- <u>US National Institute of Standards and Technology</u>

# Conferences/Journals

- ## Conferences
  - ACM Annual International Conference on Research and Development in Information Retrieval (SIGIR )
  - ACM Conference on Information Knowledge Management (CIKM)
  - …

- ## Journals
  - ACM Transactions on Information Systems (TOIS)
  - ACM Transactions on Asian Language Information Processing (TALIP)
  - Information Processing and Management (IP&M)
  - Journal of the American Society for Information Science (JASIS)
  - …

# Tentative Topic List

| 1. | Course Overview & Introduction |
|----|-------------------------------|
| 2. | Retrieval Performance Evaluation - Measures & Reference Collections |
| 3. | Retrieval Models (I) - Classic Retrieval Models (Boolean, Vector Space, Probabilistic Model, Fuzzy Set, Extended Boolean, Generalized Vector Space Model, Latent Semantic Analysis, etc) |
| 4. | Query Operations - Query Expansion and Term Re-weighting |
| 5. | Retrieval Models (II)- Language Model Approaches (HMM/N-gram, Probabilistic Latent Semantic Analysis, Latent Dirichlet Allocation, etc.) |
| 6. | Text Statistics and Operations |
| 7. | Text Classification & Clustering (SVM, Naïve Bayes, etc.; HAC, k-means, EM, etc.) |
| 8. | Word Sense Disambiguation |
| 9. | Indexing and Searching |
| 10. | Web Search and Link Analysis |
| 11. | Text and Spoken Document Summarization |

# Grading (Tentative)

- Midterm (or Final): 20%

- Homework/Projects: 50%

- Presentation: 20%

- Attendance/Other: 10%

- TA: 朱芳輝同學
  - E-mail: g94470144@mail.csie.ntnu.edu.tw
  - Tel: 29322411ext 208 (資工系208室)