# Information Retrieval and Extraction
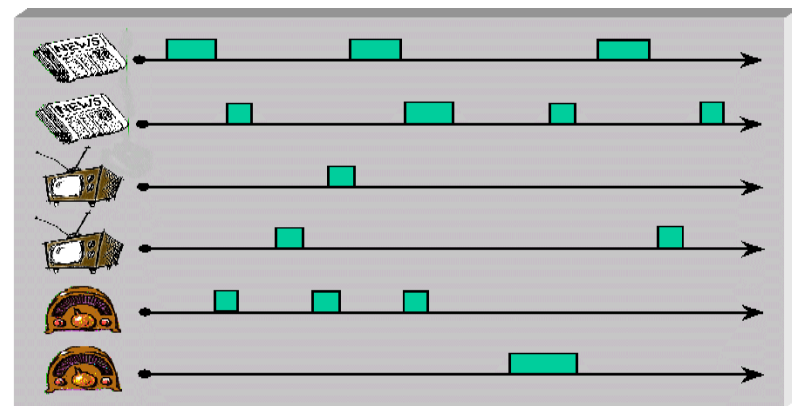
Berlin Chen 2004
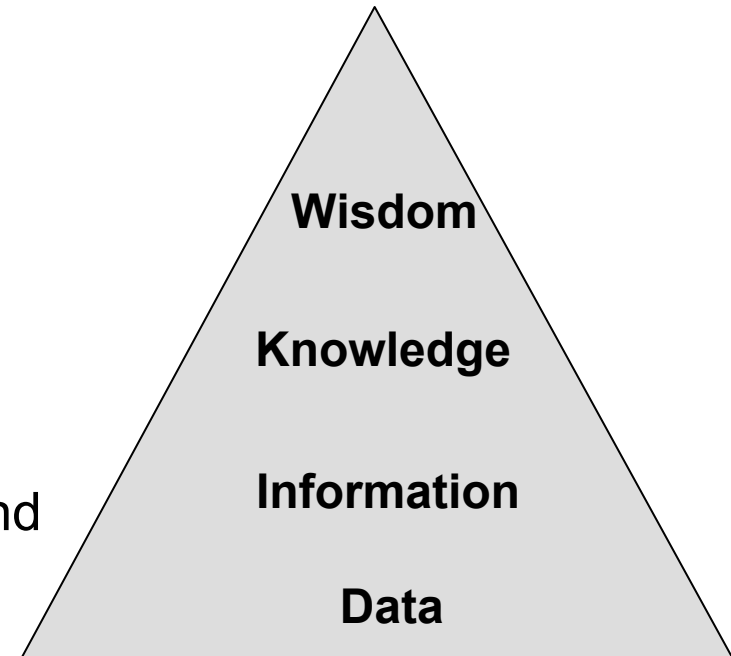
(Picture from the TREC web site)

# Textbook and References

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

- References
  - W. B. Croft and J. Lafferty (Editors). *Language Modeling for Information Retrieval*. Kluwer-Academic Publishers, July 2003
  - W. B. Frakes and R. Baeza-Yates. *Information Retrieval: Data Structures & Algorithms*. Prentice-Hall, 1992
  - 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
  - A. D. Bimbo. *Visual Information Retrieval*. Morgan Kaufmann, 1999

# Motivation

- 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 (cont.)

- 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**

**Search engine/IR system**

Emphasis is on the retrieval of information (not data)

# 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!

# Motivation (cont.)

- 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

- 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 (cont.)

- ## 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, …)

```
Docs → accents, spacing, etc. → stopwords → Noun groups → stemming → Manual indexing
```

text + structure

structure → text

**structure**        **Full text** ---------------------------------→ **Index terms**

# 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



User
Interface

Text

4, 10

user need

Text

Text   Operations

6, 7

logical view

logical view

Query
Operations

user feedback

5

Indexing

DB Manager
Module

8

query

inverted file

Searching

8

Index

retrieved docs

Text
Database

Ranking

2

ranked docs

# The Retrieval Process (cont.)

- 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

- Four Main Topics



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

# Major Topics (cont.)

- 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

- Applications
  - Web, bibliographic systems, digital libraries

# Textbook Topics

# Text Information Retrieval

- Internet searching engine



Web

Spider

Mirrored Web Page Repository

Indexer

Queries

Ranked Docs

Search Engine

# Text Information Retrieval (cont.)

# Speech Information Retrieval



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 (cont.)

- Compaq Research Group – Speechbot System
  - Broadcast news speech recognition, Information retrieval, and topic segmentation (SIGIR2001)
  - Currently indexes **14,791 hours of content** (2004/09/22, http://speechbot.research.compaq.com/)

# Speech Information Retrieval (cont.)

・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 (cont.)

# Visual Information Retrieval
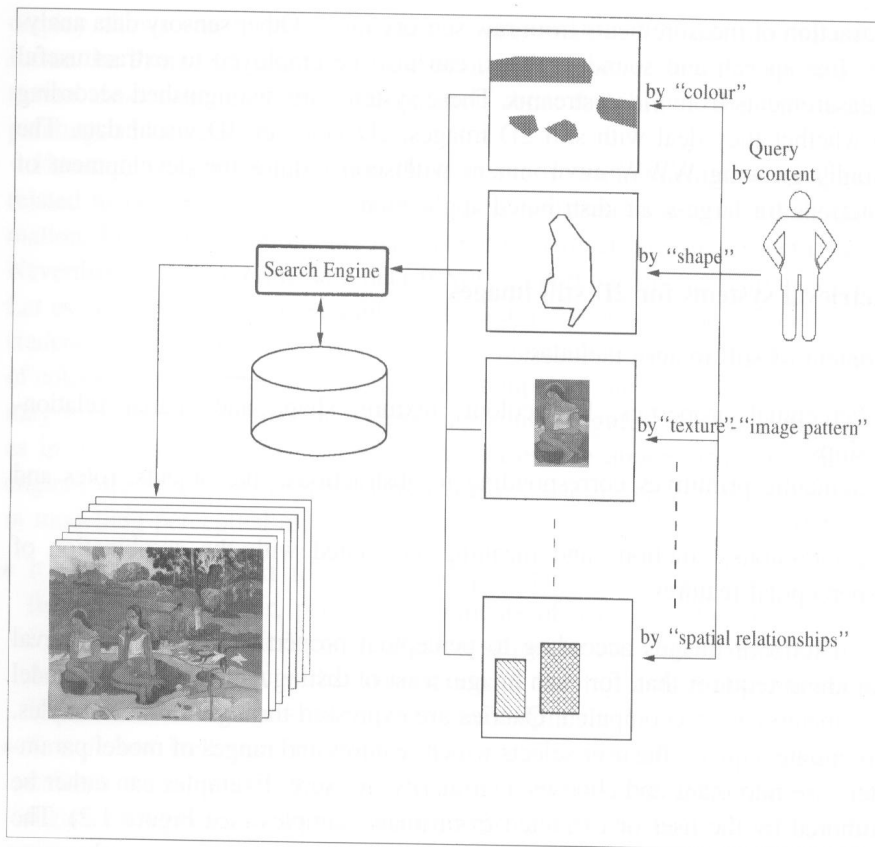
- Content-based approach



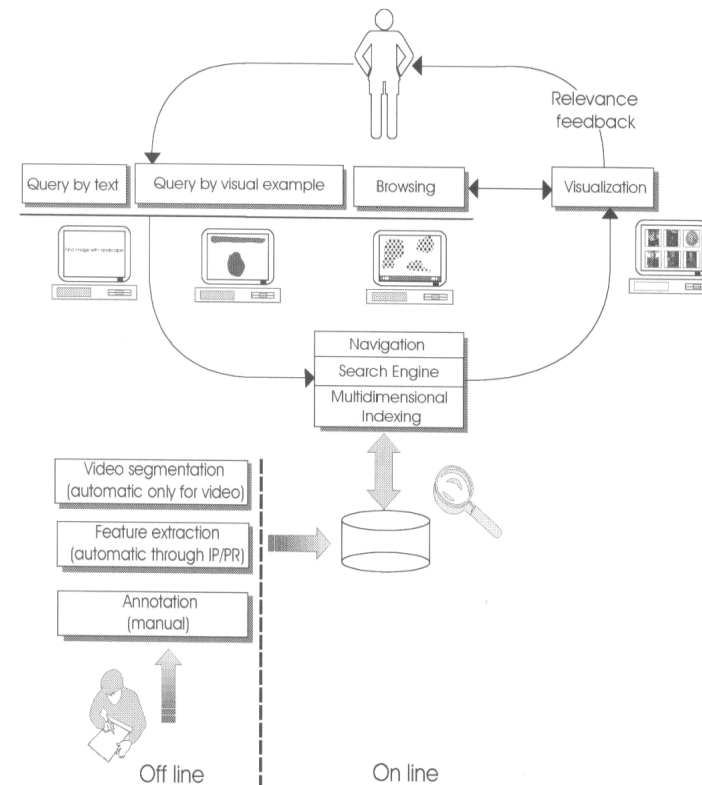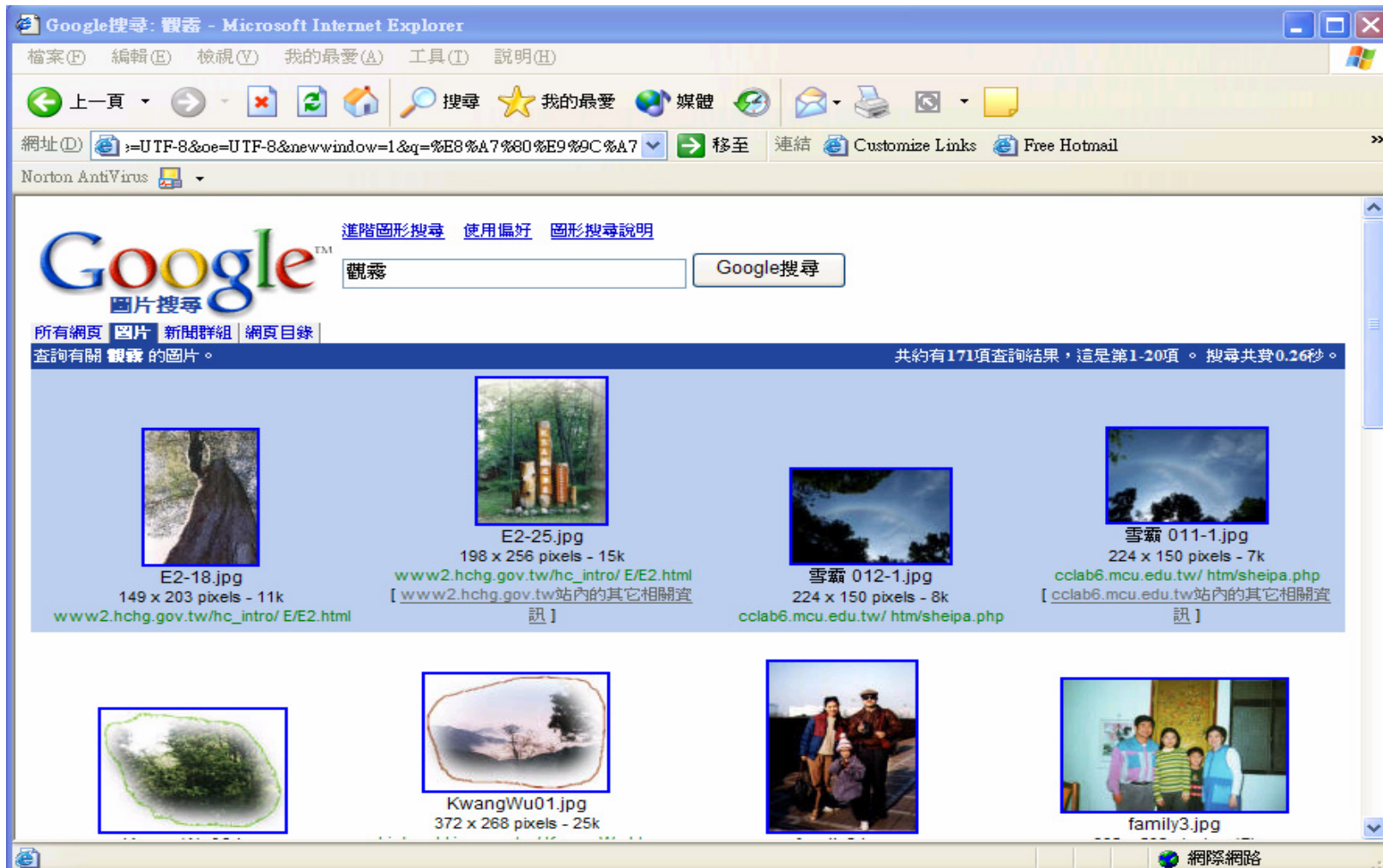Figure 1.2 Different types of query by example.



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

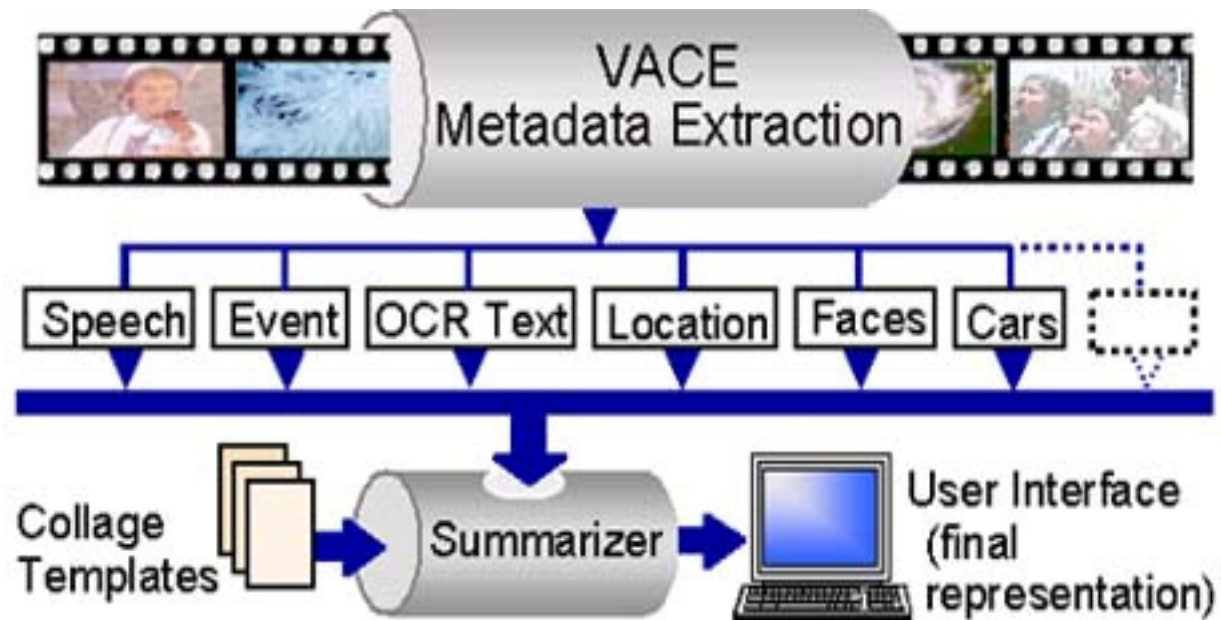# Visual Information Retrieval (cont.)

- Images with Texts

# Visual Information Retrieval (cont.)

- Content-based Image Retrieval

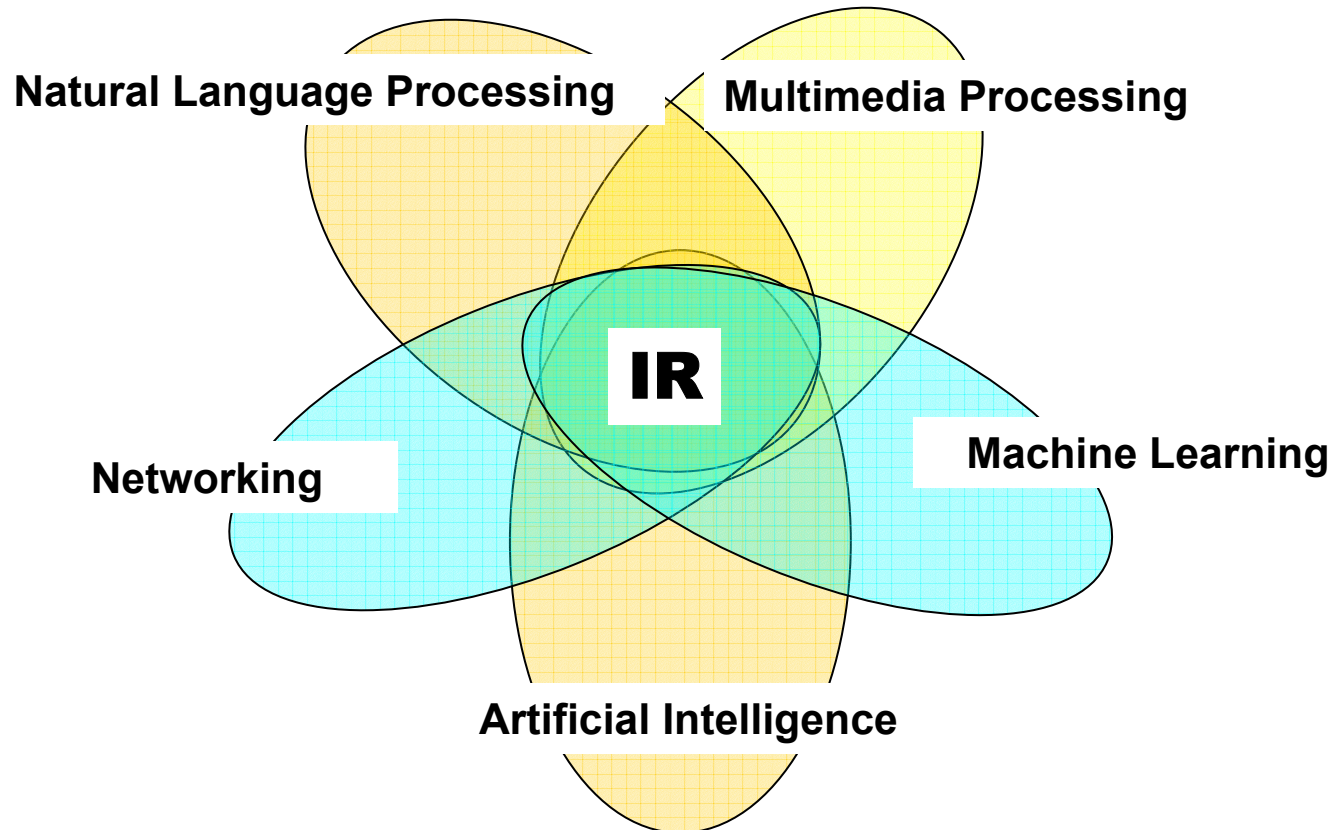# Visual Information Retrieval (cont.)

**Video Analysis and Content Extraction**
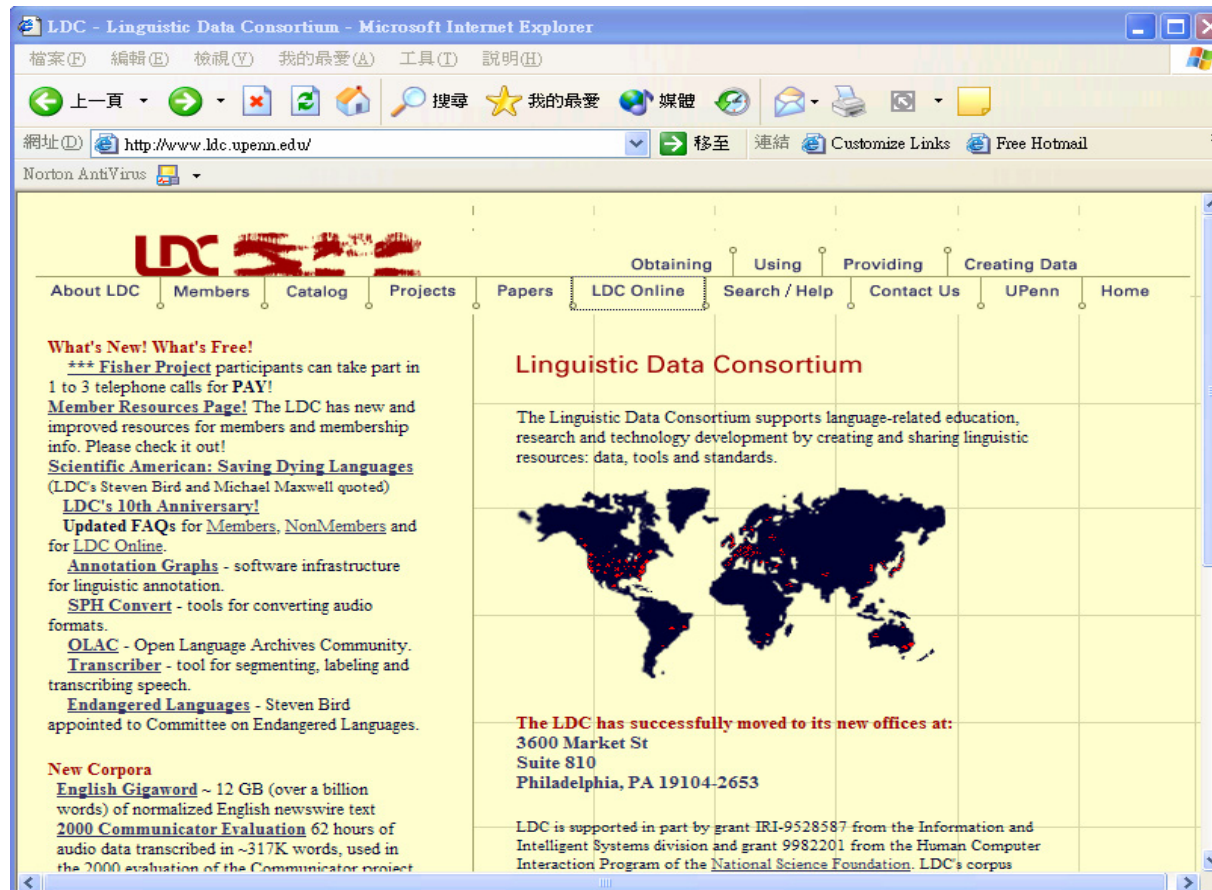
# Other IR-Related Tasks

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

# Multidisciplinary Approaches



**Natural Language Processing**    **Multimedia Processing**

**IR**

**Networking**                        **Machine Learning**
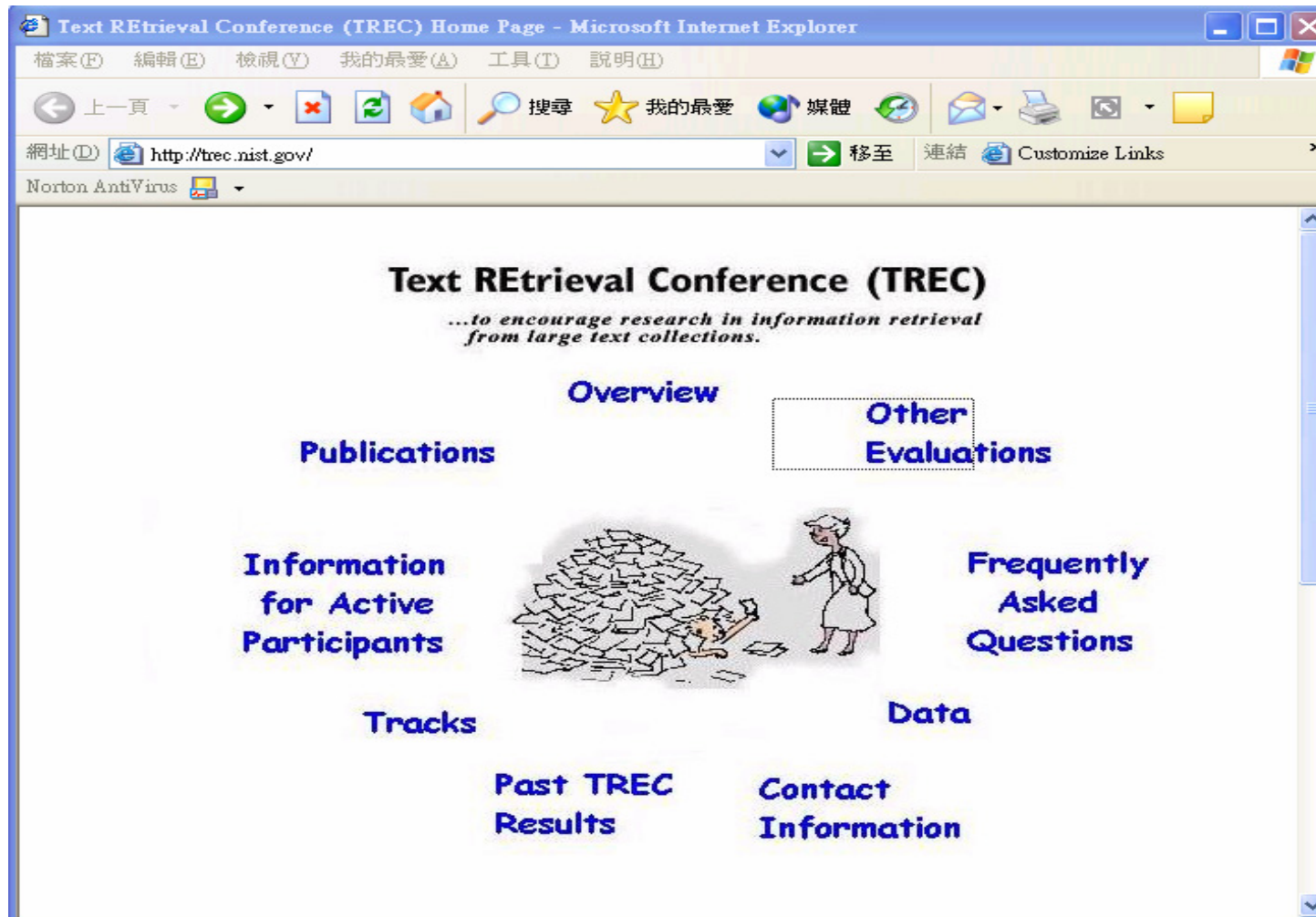
**Artificial Intelligence**

# 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

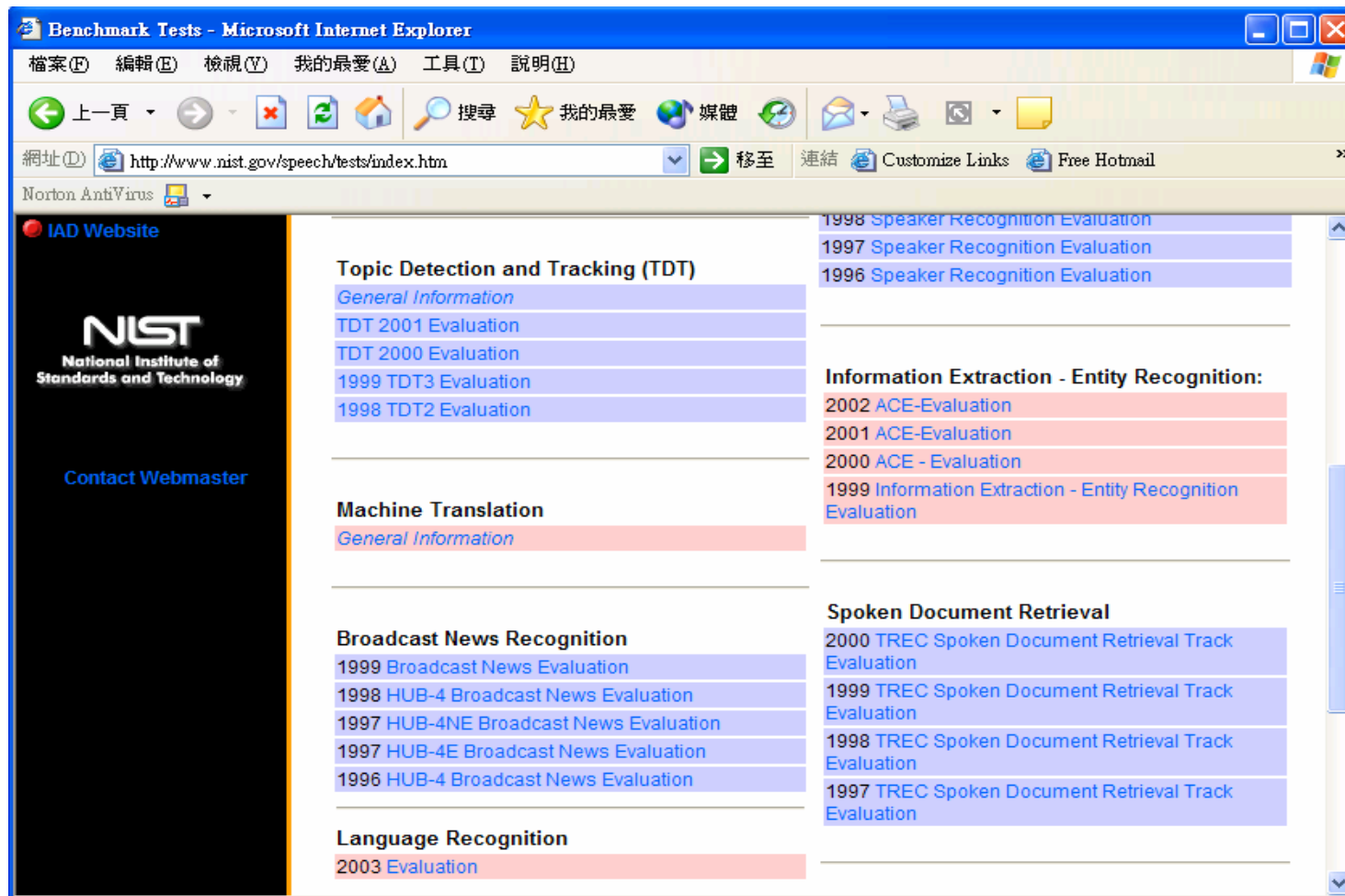- [Text REtrieval Conference](#) (TREC)

# Contests

- **US National Institute of Standards and Technology**

# 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)
  - …