



# Multimedia Information Retrieval

- 1 What is multimedia information retrieval?**
- 2 Basic Multimedia Search Technologies**
- 3 Evaluation of MIR Systems**
- 4 Added Value - user interaction, visualisation and the MIR research landscape**



# Multimedia Information Retrieval

- 1 What is multimedia information retrieval?
- 2 Basic Multimedia Search Technologies
- 3 Evaluation of MIR Systems**
  - 3.1 Metrics**
  - 3.2 Calculating and Comparing**
  - 3.3 Evaluation Campaigns**
- 4 Added Value - user interaction, visualisation and the MIR research landscape



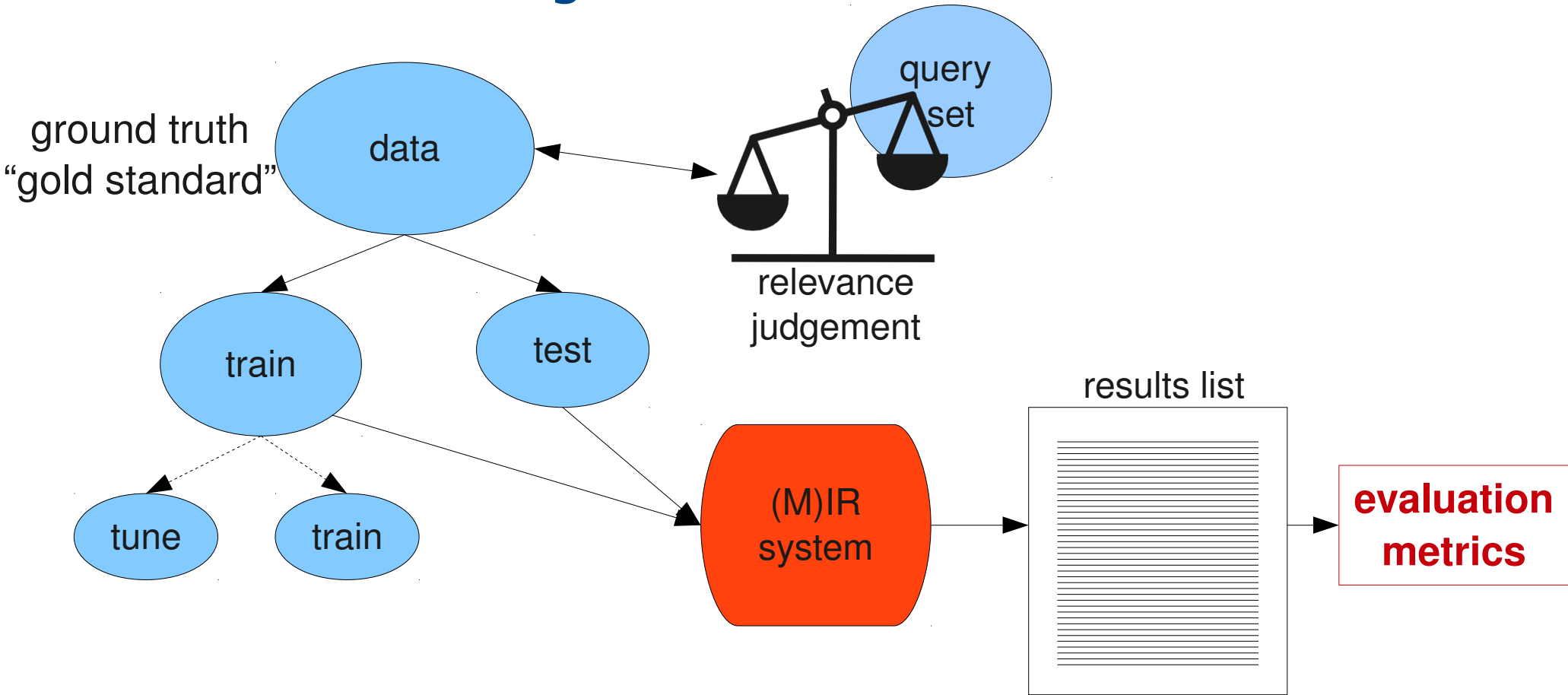
How do we know if our MIR system is effective?

Why do we care about quantifying the performance?

*“If you can not measure it, you can not improve it.” –  
Lord Kelvin*

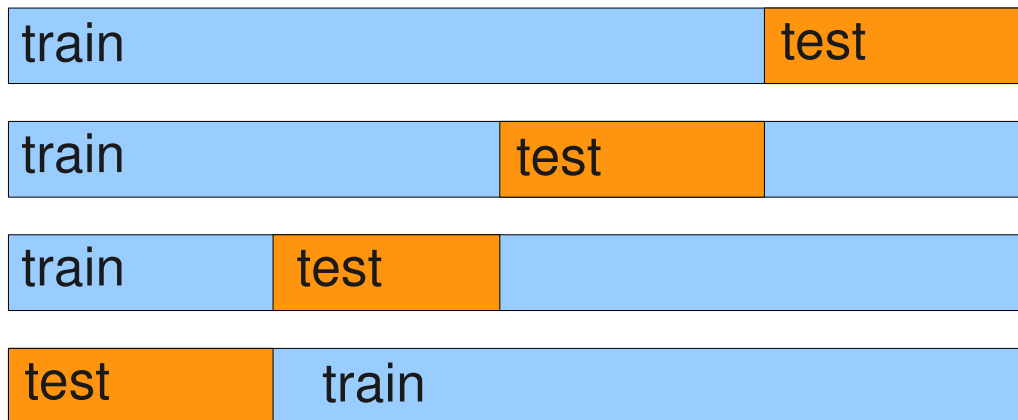


## “Cranfield Paradigm”





## Cross-validation



Randomise data and divide  
Train-test 4 times  
Average all metrics  
4-fold cross-validation

Extreme is Leave-One-Out : test size = 1



# Relevance?



## Find me pictures of triumph

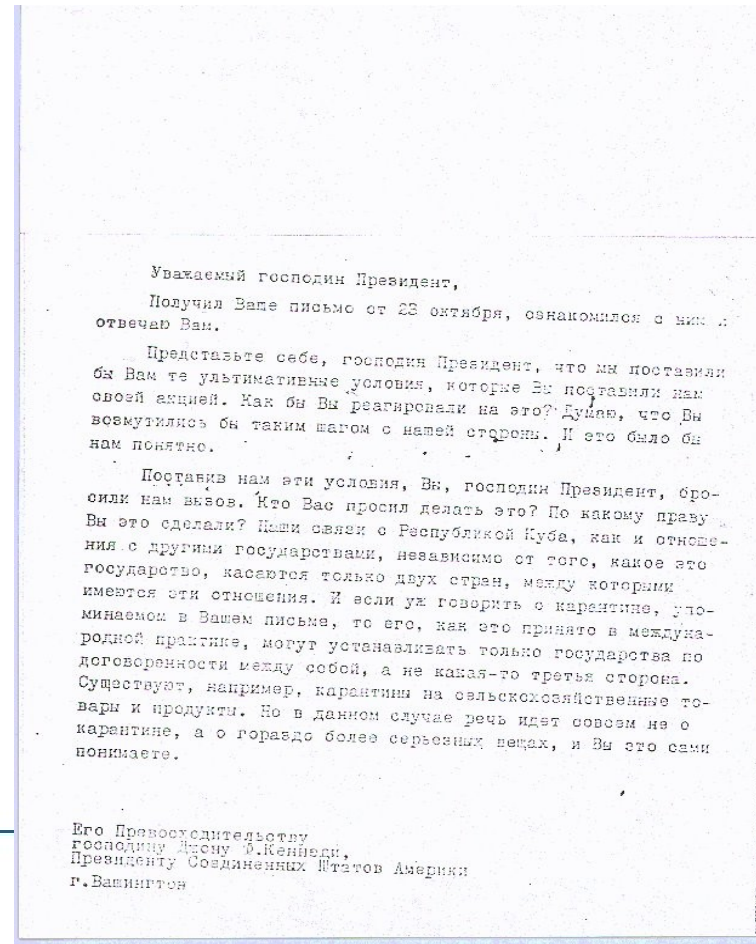
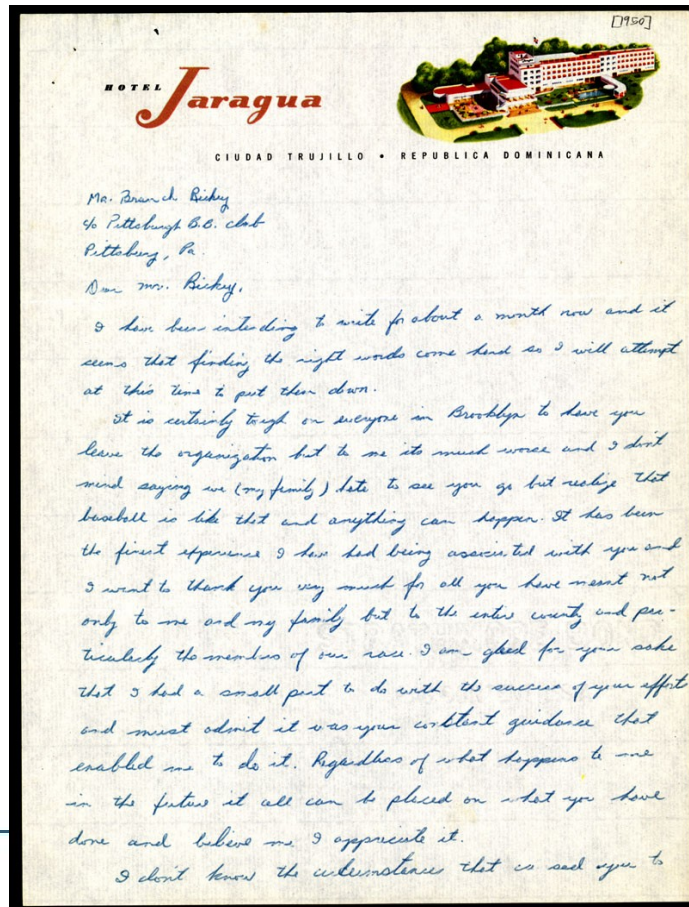


\*<http://www.flickr.com/photos/ricardodiaz/>



# Exercise

Find shots of printed, typed, or handwritten text, filling more than half of the frame area





	<b>Relevant</b>	<b>Irrelevant</b>
<b>Retrieved</b>	True Positive (tp)	False Positive (fp)
<b>Not Retrieved</b>	False Negative (fn)	True Negative (tn)

Precision (P) = fraction retrieved that are relevant

$$P = tp / (tp + fp)$$

Recall (R) = fraction relevant that are retrieved

$$R = tp / (tp + fn)$$





# Precision or Recall?

What about accuracy?

$$\text{Accuracy} = (tp+tn)/(tp+fp+fn+tn)$$

Is precision or recall more useful/important

if I'm doing a web search on Gold Coast accommodation?

if I'm a paralegal researching case precedence?

How could I make a system with 100% recall?

$F_1$ -measure (weighted harmonic mean of P & R)

$$F_1 = \frac{2 \cdot \textit{precision} \cdot \textit{recall}}{(\textit{precision} + \textit{recall})}$$



# Exercise

An IR system returns 8 relevant documents and 10 irrelevant documents. There are a total of 20 relevant documents in the collection. Calculate the precision and recall.



An IR system returns 8 relevant documents and 10 irrelevant documents. There are a total of 20 relevant documents in the collection. Calculate the precision and recall.

$tp = 8$	$fp = 10$
$fn = 12$	$tn = (\text{unknown})$

$$P = tp/(tp+fp) = 8/(8+10) = 8/18 = 0.44$$

$$R = tp/(tp+fn) = 8/(8+12) = 8/20 = 0.40$$

$$F_1\text{-measure would be } 2 \times 0.44 \times 0.40 / (0.44 + 0.40) = 0.42$$



Which is better? There are 5 relevant documents to be found.

## System A

1. Relevant
2. Relevant
3. Irrelevant
4. Irrelevant
5. Relevant
6. Relevant

$$\text{Precision} = 4/6 = 0.66$$

$$\text{Recall} = 4/5 = 0.80$$

## System B

1. Relevant
2. Irrelevant
3. Relevant
4. Relevant
5. Relevant
6. Irrelevant

$$\text{Precision} = 4/6 = 0.66$$

$$\text{Recall} = 4/5 = 0.80$$



Precision @ N

Precision/Recall graphs

Mean Average Precision



Which is better? There are 5 relevant documents to be found.

## System A

1. Relevant
2. Relevant
3. Irrelevant
4. Irrelevant
5. Relevant
6. Relevant

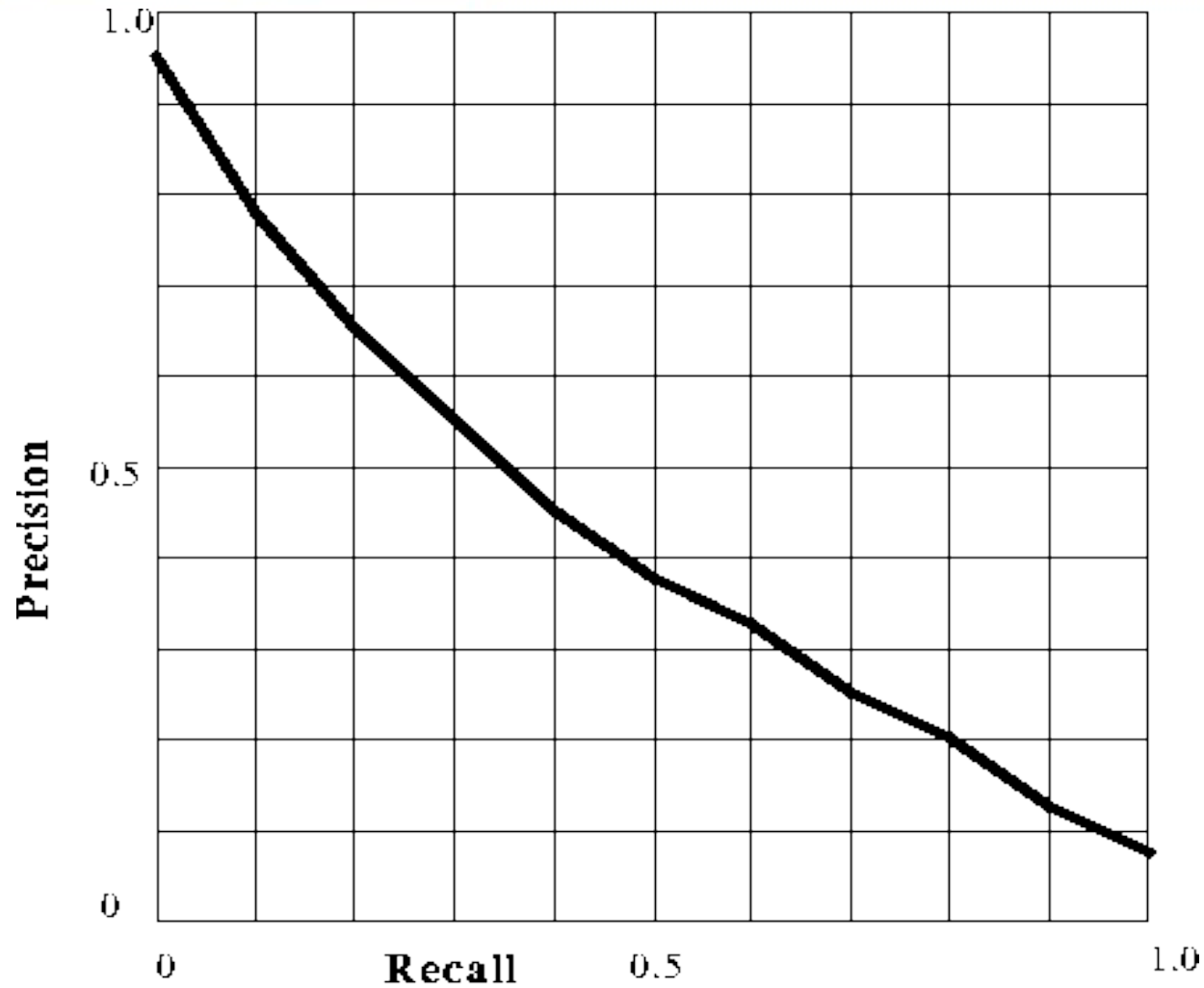
## System B

1. Relevant
2. Irrelevant
3. Relevant
4. Relevant
5. Relevant
6. Irrelevant

P@1  
P@2  
P@3  
P@4  
P@5



# Precision/Recall Curve



*Fig. 1. A typical precision-recall graph*



## System A

- |               |          |
|---------------|----------|
| 1. Relevant   | P = 1    |
| 2. Relevant   | P = 1    |
| 3. Irrelevant | -        |
| 4. Irrelevant | -        |
| 5. Relevant   | P = 0.6  |
| 6. Relevant   | P = 0.67 |

$$(1+1+0.6+0.67)/4 = 0.82$$

## System B

- |               |          |
|---------------|----------|
| 1. Relevant   | P = 1    |
| 2. Irrelevant | -        |
| 3. Relevant   | P = 0.67 |
| 4. Relevant   | P = 0.75 |
| 5. Relevant   | P = 0.8  |
| 6. Irrelevant | -        |

$$(1+0.67+0.75+0.8)/4 = 0.69$$





Which is better? There are 5 relevant documents to be found.

## System A

1. Relevant
2. Relevant
3. Irrelevant
4. Irrelevant
5. Relevant
6. Relevant

AP = 0.82

## System B

1. Relevant
2. Irrelevant
3. Relevant
4. Relevant
5. Relevant
6. Irrelevant

AP = 0.69



Use the results (exercises/evaluation/) from 2 image search engines and calculate the performance. Which is better?

Spreadsheet



Overfitting to limited training data → unbalanced, fragile system

Unrealistic training data

Difficulty in finding training data

Comparison and competition

Numbers not users



TRECVID

ImageCLEF

MediaEval

MIREX





# TREC Video retrieval conferences

Organised by NIST with support from other U.S. government agencies - <http://www-nlpir.nist.gov/projects/trecvid/>

Objective is to encourage research in information retrieval by:  
Providing a **large** test collection.

Uniform **scoring** procedures.

Forum for organizations interested in **comparing** their results.

Tasks:

Shot boundary detections (retired)

High-level feature extraction (semantic annotation)

Search (interactive, manually-assisted or fully automated)

Rushes summarisation



In the first few years of TRECVID video retrieval was best done with “text only”

Image analysis did not help in early years

**BUT** situation has changed!

Combination of weak classifiers to corroborate evidence

The number of visual concepts has increased; see, eg, LSCOM



## TRECvid example queries

“Find shots of a road taken from a moving vehicle through the front window”



“Find shots of a person talking behind a microphone”

“Find shots of a street scene at night”





CLEF = Cross Language Evaluation Forum

Process is modelled from TREC

ImageCLEF started in 2003

Tasks:

Image retrieval (queries in different languages)

Medical Image Annotation

Annotation of photographs

Geographic retrieval (GeoCLEF)

Video retrieval (VideoCLEF/MediaEval)





## System issues

Indexing speed

Scalability

Robustness

Query expressiveness

## User issues

Diversity, Responsiveness

“happiness” ?

The interface vs IR performance



# Multimedia Information Retrieval

- 1 What is multimedia information retrieval?
- 2 Basic Multimedia Search Technologies
- 3 Evaluation of MIR Systems**
  - 3.1 Metrics**
  - 3.2 Calculating and Comparing**
  - 3.3 Evaluation Campaigns**
- 4 Added Value - user interaction, visualisation and the MIR research landscape