

➤ Web Retrieval Multimedia Search

Frank Hopfgartner
Institute for Web Science and Technologies

Intended Learning Outcomes

At the end of this lecture, you will be able to:

- Describe the differences between metadata based and content-based retrieval
- Describe different types of image, audio and video retrieval
- Describe, at a high level, how content-based search of image, audio and video works
- Discuss the advantages and disadvantages of content-based search, for images, audio and video content

Outline

- Motivation and Challenges of Multimedia Search
- Image Retrieval
- Audio Retrieval
- Video Retrieval
- Summary

➤ 1. Motivation and Challenges of Multimedia Search

The “Grand Challenge”

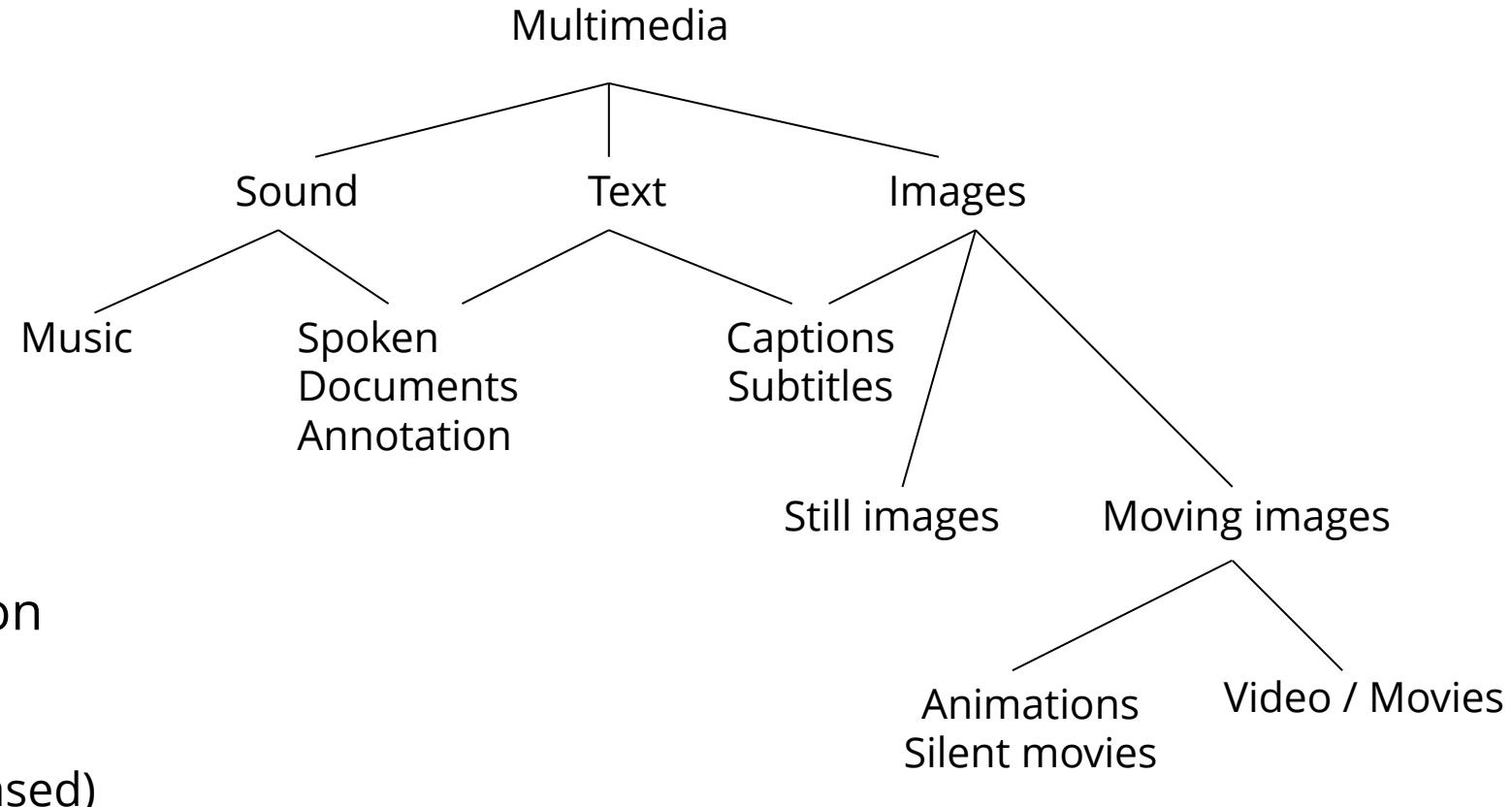
*“Given a query in **any medium** and any language, select relevant items from a multilingual **multimedia collection** which can be in **any medium** and any language, and present them in the style or order most likely to be useful to the querier, with identical or near identical objects in **different media** or languages appropriately identified.”*

What is multimedia?

- Often defined as a combination of two or more media
 - E.g. combination of audio and moving images into a video, many webpages (text + images) can also be considered as “multimedia”
 - See Jaimes et. al. (2005) for some definitions from a number of prominent researchers in the field
- In IR, “multimedia” is often used to refer to the retrieval of non-text items

Why multimedia?

- Digital multimedia content is growing rapidly in both domestic and commercial sectors
 - driving new forms of interaction with images, speech, video, text and other forms of unstructured data
 - Domestic consumer market fuelling the creation, use and sharing of digital multimedia
 - video (camcorders, digital cameras)
 - images (digital cameras, mobile phones)
 - audio (music and speech)
 - Applications include
 - entertainment, e-commerce, digital photos, music downloads, e-learning, mobile media
- How many
multimedia devices do
you own?*



Access methods depend on
information available

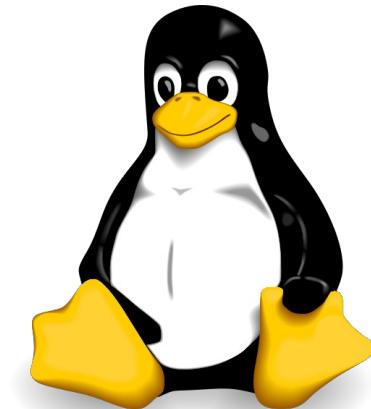
- about an object (metadata)
- within an object (feature-based)

➤ 2. Image Retrieval



Images

- Bitmap or raster
 - Image is a grid of pixels (dots)
 - Produced by digital cameras and software such as Adobe Photoshop or Microsoft paint
- Vector
 - Image is a sequence of drawing instructions
 - Produced by drawing software such as Adobe Illustrator



Have you ever carried out an image search using Google images, Bing images, or similar?



university of koblenz

Search



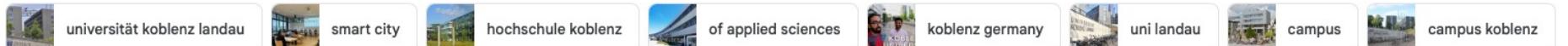
Sign in

iversität
blenz

All Maps Images News Videos : More

Tools

SafeSearch ▾



• Detail | uni-assist
Detail | uni-assist



• Shiksha
University of Koblenz-Landau: Ra...



Q Quora
study at University of Koble...

UK universität
koblenz

weiter:denken
www.uni-koblenz.de
University of Koblenz



• Ukrainian-American Concordia U...
Koblenz – Concordia University



• Universität Koblenz
University of Koblenz ...



w Wikipedia
Universität Koblenz – Wikipedia



w Wikipedia
University of Koblenz and Lan...



• Gerber Architekten
Koblenz-Landau - Gerber Archite...



• www.hs-koblenz.de
Our campuses



• Ukrainian-American Concordia Univ...
University of Koblenz-Landau (Ge...



• Avit Bhowmik
Dr. Avit Bhowmik ...



• Alamy
University koblenz landau h...



• Detail | uni-assist
Detail | uni-assist



• Forschung & Lehre
Verwaltung der Uni Koblenz-Landau wir...



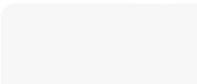
• Universität Koblenz
Onboarding | Universität Koblenz



• SWR.de
Uni Koblenz ist ab Januar eigenständig ...



• YouTube
University of Koblenz- Lan...



Have you ever carried out an image search using Flickr, Getty images, or similar?



Advanced

Any license ▾

SafeSearch on ▾

Share

Relevant ▾



Sponsored images from iStock.

LIMITED DEAL: 20% off with code FLICKR20



Upgrade to Flickr Pro to hide these ads

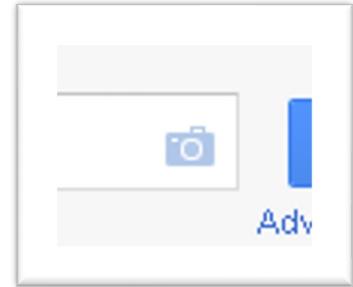
Everyone's photos

View all 262



Have you ever carried out an image search using TinEye or
Google Lens?

Aside: Google Lens: the little blue camera icon on the right of the search bar





Search any image with Google Lens X

 Drag an image here or [upload a file](#)

OR

Ctrl + V

Upload



Sign in

Google

Find image source



Universität Koblenz

4.2 ★★★★☆

University

Search



W Wikipedia

File:UK Logo CMYK.png
- Wikipedia

See exact matches



Universität Koblenz
Fachtag Demokratie | Universität Koblenz



Did you find these results useful?

Related searches



X
Universität Koblenz on X: "Die neue..."



LinkedIn
Aravindhan S P - Graduate Engineering...



Top Universities
Jan Kochanowski University in Kielce :...



LinkedIn
Sadaaf Chowdhury - Software Engineer - Ci...



- universität

Yes No

Different types of image search

1. Search for images embedded in textual or hypertext documents
E.g. Google/Bing image search

2. Search metadata associated with images
E.g. Flickr/Getty images

3. Search using the image content
E.g. TinEye, Google Lens

Representing an image by its content

- If an image is part of a larger document, such as a web page, we can describe it using features such as:
 - The link text which points to the image
 - “Alt” text which is provided in the HTML
 - The filename of the image
 - The text which surrounds the image
 - etc...

Research

WeST members perform research on a large number of topics that evolve around the World Wide Web and how people use it. The web is a virtual space for people to interact and share experiences. A core focus is on the analysis of data-driven applications, their application in the context of digital transformation, and their wide ranging and often disruptive impact on society.

The WeST Institute has a strong record of collaborations with national and international partners from the scientific community and industry. These fruitful collaborations allow us to bring our cutting-edge technologies into the world.

For our research, we can rely on an IaaS cloud infrastructure capable of handling big data. In addition, we are in progress of setting up a state-of-the-art usability lab to enable further research on data-driven applications.

Transparency, reproducibility, and open access are at the core of our work. Most of our research papers can be downloaded from open access repositories. In addition, we have made available various research [prototypes](#) and [datasets](#) under open licenses. Researchers all over the world use resources that got prepared and published by WeST.



➤ [Research Projects](#)

Research

WeST members perform research on a large number of topics that evolve around the World Wide Web and how people use it. The web is a virtual space for people to interact and share experiences. A core focus is on the analysis of data-driven applications, their application in the context of digital transformation, and their wide ranging and often disruptive impact on society.

The WeST Institute has a strong record of collaborations with national and international partners from the scientific community and industry. These fruitful collaborations allow us to bring our cutting-edge technologies into the world.

For our research, we can rely on an IaaS cloud infrastructure capable of handling big data. In addition, we are in progress of setting up a state-of-the-art usability lab to enable further research on data-driven applications.

Transparency, reproducibility, and open access are at the core of our work. Most of our research papers can be downloaded from open access repositories. In addition, we have made available various research [prototypes](#) and [datasets](#) under open licenses. Researchers all over the world use resources that got prepared and published by WeST.



➤ [Research Projects](#)

Research

Note: the content of the image is not used...

WeST members perform research on a large number of topics that evolve around the World Wide Web and how people use it. The web is a virtual space for people to interact and share experiences. A core focus is on the analysis of data-driven applications, their application in the context of digital transformation, and their wide ranging and often disruptive impact on society.

The WeST Institute has a strong record of collaborations with national and international partners from the scientific community and industry. These

```
► <div class="six wide computer twelve wide mobile six wide tablet column">...</div>
▼ <div class="six wide computer twelve wide mobile six wide tablet column">
  <div>
    <p class="block image align center">
       event
    </p>
    <br>
  </div>
</div>
```

resources that got prepared and published by WeST.



➤ Research Projects

Search metadata associated with images

flickr Explore Prints Get Pro

koblenz Log In Sign Up

Photos People Groups

Orientation Minimum size Date taken Content Search in

Advanced All Tags

Any license SafeSearch on Share Relevant

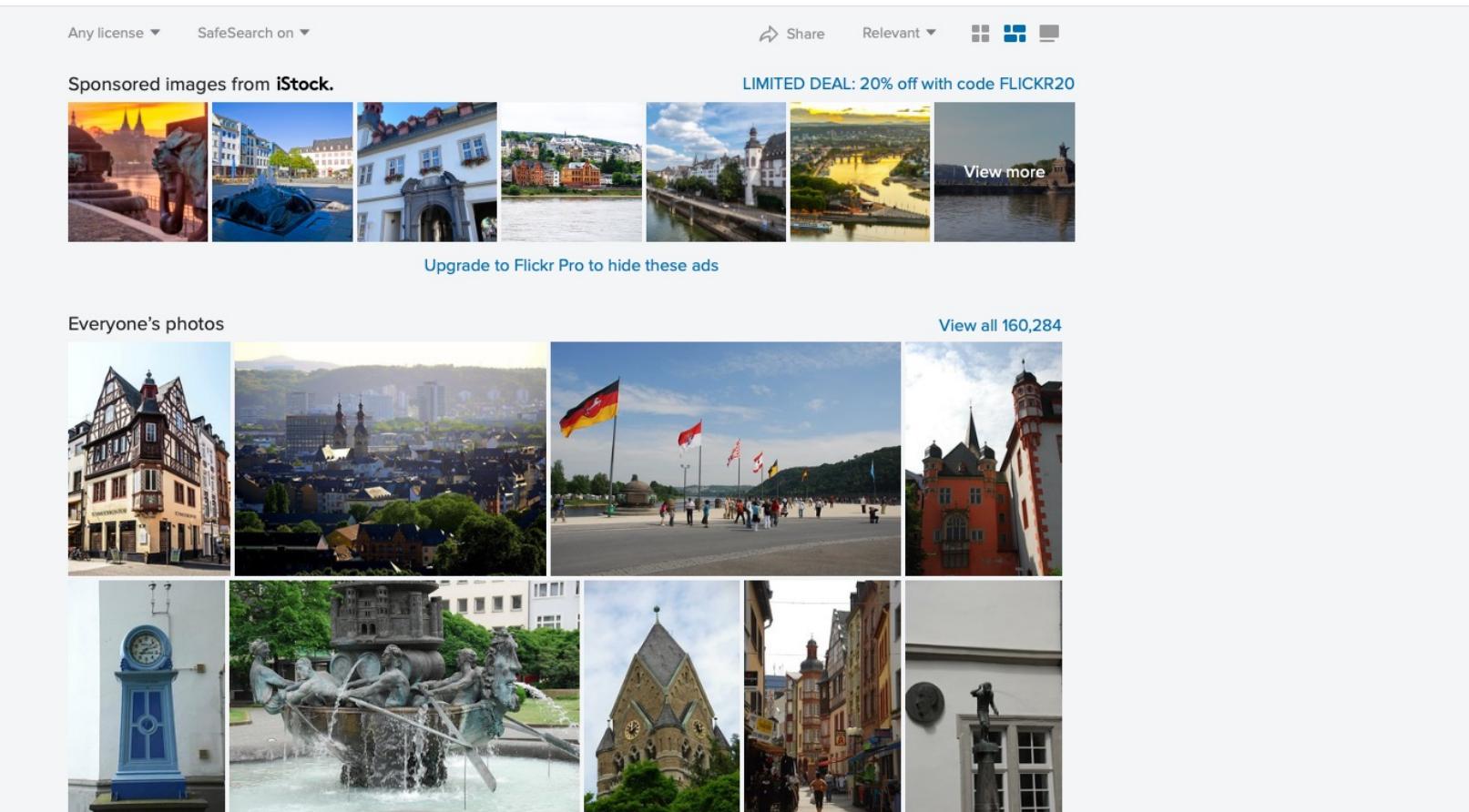
Sponsored images from iStock.

LIMITED DEAL: 20% off with code FLICKR20

View more

Upgrade to Flickr Pro to hide these ads

Everyone's photos View all 160,284



Metadata-based image retrieval

- Images indexed using associated text (metadata)
- Metadata is “data about data”
 - Controlled vocabularies
 - Uncontrolled vocabularies (folksonomies)
 - Embedded data (e.g. Exif)
 - Textual descriptions
 - File (and folder) names
- Typically assigned manually
- Retrieval uses traditional text retrieval or database approaches
- This approach is popular online (e.g. Flickr) where there are many image uploaders who can add metadata to their own images

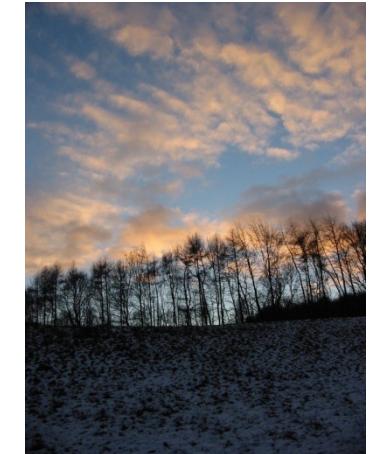


Photo
Trees
Sunset
Dusk

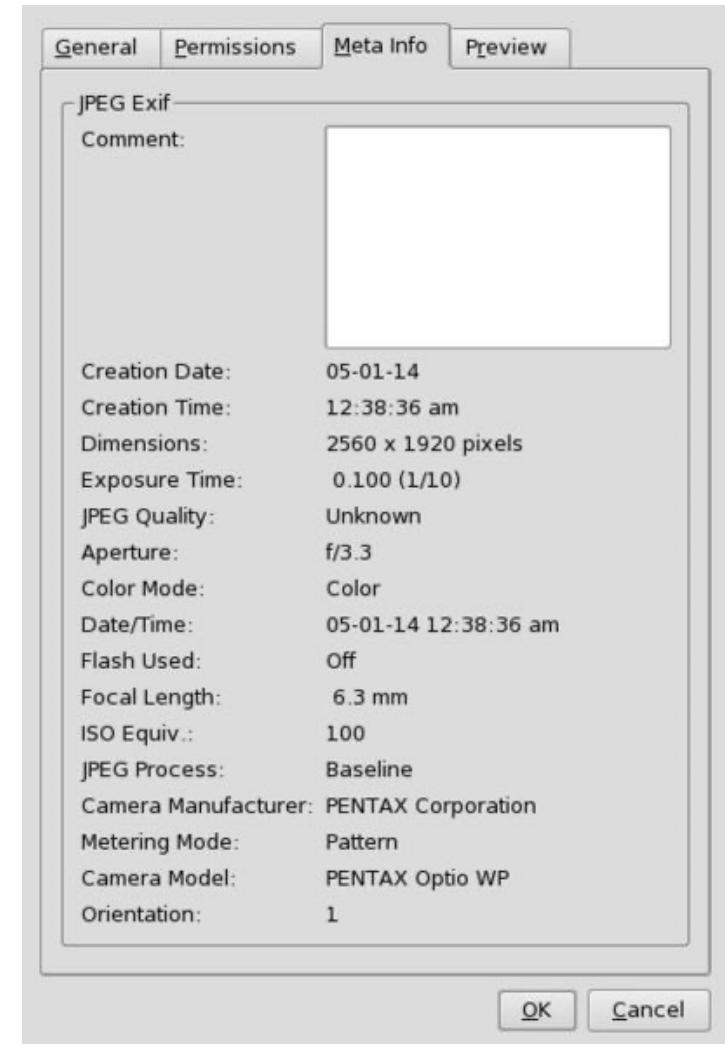
Example metadata

| | |
|----------------------|---|
| Record ID: | JV-.044809 |
| Short title: | The Smeaton Tower, Plymouth. |
| Long title: | Plymouth Hoe. The Smeaton [Lighthouse] Tower. |
| Location: | Devonshire, England |
| Description: | Red and white striped lighthouse on coastal cliff with harbour and town beyond, and substantial building on cliff terrace below. |
| Date: | Registered 1904 |
| Photographer: | J Valentine & Co |
| Categories: | [lighthouses] [beacons & lighthouses] [Devon all views] [Collection - J Valentine & Co] |
| Notes: | JV-44809 pc/mb(or possibly 44810)TECH: Coloured. |



Automatically assigned metadata: EXIF

- Exchangeable Image File Format (EXIF)
 - A *standard* for storing interchange information *in* image files, especially those using JPEG compression
- Most digital cameras now use EXIF format to store
 - Date and time information
 - Camera settings (e.g. orientation, aperture, shutter speed, focal length)
 - Geolocation (based on GPS)



[flickr](#) Explore Prints Get Pro

Photos, people, or groups Log In Sign Up



[Add comment](#)

cavium + Follow
_MG_4461_2
Koblenz, Uni, Bibliothek, Metternich

145 views 0 faves 0 comments Uploaded on January 22, 2008
Taken on January 21, 2008
© All rights reserved

Login to comment

Add comment

Canon EOS 350D Digital

18.0 mm 13
ISO 100 Flash (off, did not fire)

[Show EXIF](#)



Canon EOS
350D Digital

18.0 mm 13

ISO 100 Flash (off,
did not fire)

[Hide EXIF](#)

Compression - JPEG (old-style)

Make - Canon

Orientation - Horizontal (normal)

X-Resolution - 72 dpi

Y-Resolution - 72 dpi

Software - QuickTime 7.4

Date and Time (Modified) -
2008:01:22 00:00:41

Host Computer - Mac OS X 10.5.1

YCbCr Positioning - Centered

ISO Speed - 100

Exif Version - 0220

Date and Time (Original) -
2008:01:21 21:03:13

Date and Time (Digitized) -
2008:01:21 21:03:13

Components Configuration - Y, Cb,
Cr, -

Exposure Bias - +2/3 EV

Metering Mode - Multi-segment

Flashpix Version - 0100

Color Space - Uncalibrated

Camera ID - 31

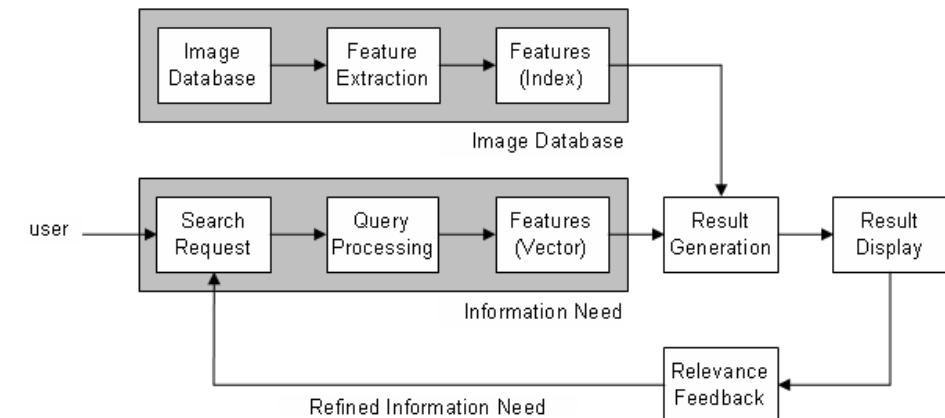
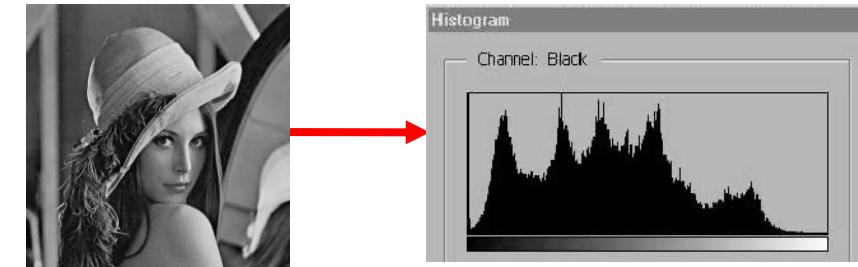
Camera Type - Digital SLR

Problems with metadata-based image retrieval

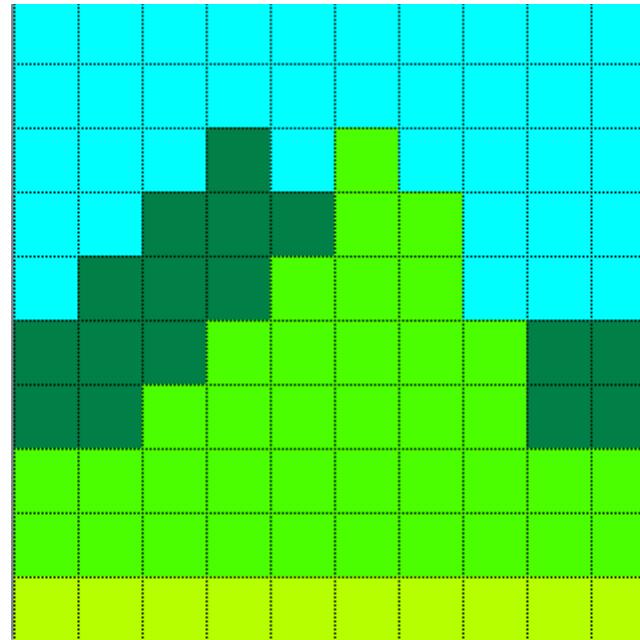
- Manual annotation is expensive (time)
- Meaning of image difficult to interpret/express in written form
- Indexed text might not relate to image content
- Often short texts to index with (e.g. captions)
- Manual annotation is subjective and suffers from low agreement between individuals (and groups)
- Vocabulary mismatch between indexer and user
- Difficult to express more abstract needs (e.g. “visual” or emotive queries)
- Notions of relevance (i.e. based on factors other than topic/theme of image)

Content-based image retrieval

- Images are represented using "low-level" features which are automatically generated
 - E.g. based on colour, shape, texture, etc.
 - Image represented as "feature vector"
- System compares the "feature vector" of the query and to the "feature vector" of all stored images
 - Uses a similarity/distance measure
 - This generates a score (a number)
- A ranking is then produced by sorting this score, e.g. from largest to smallest



Examples of a very simple global “colour” feature

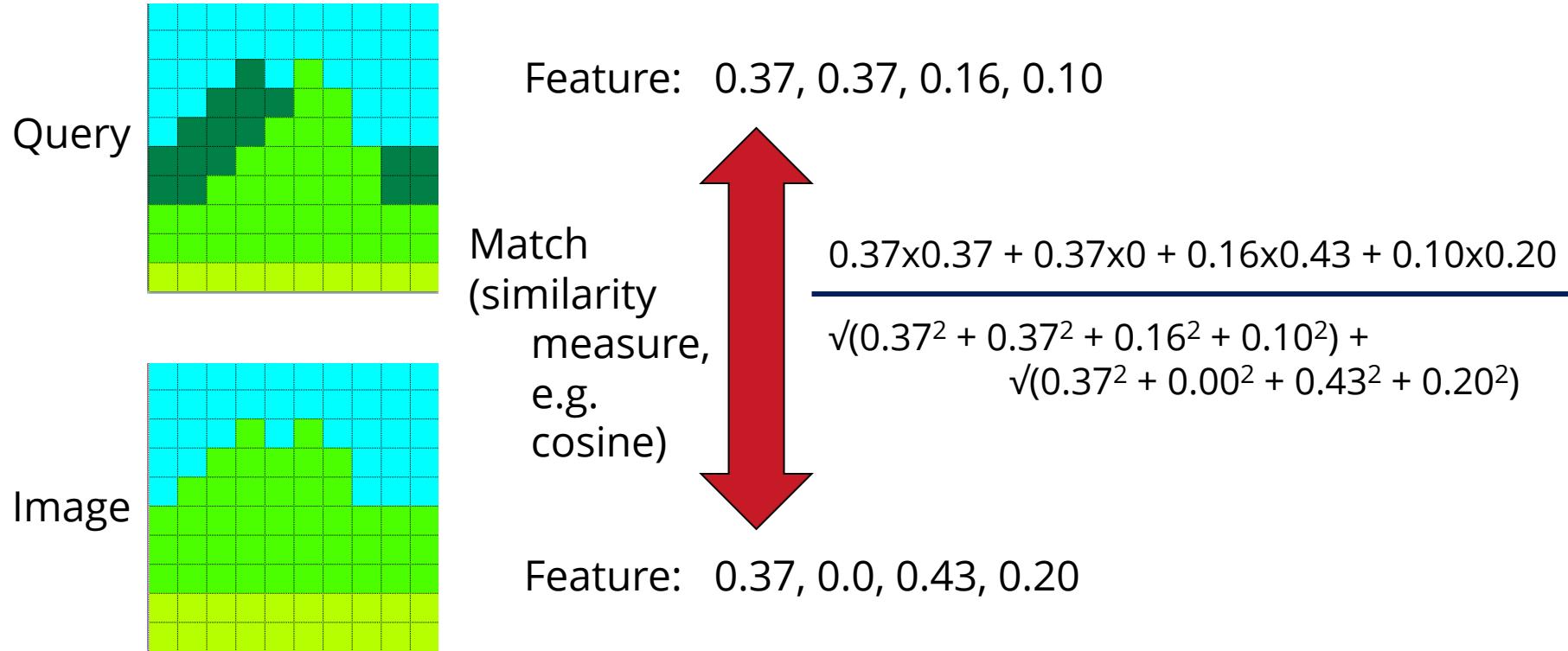


Frequency 37 37 16 10



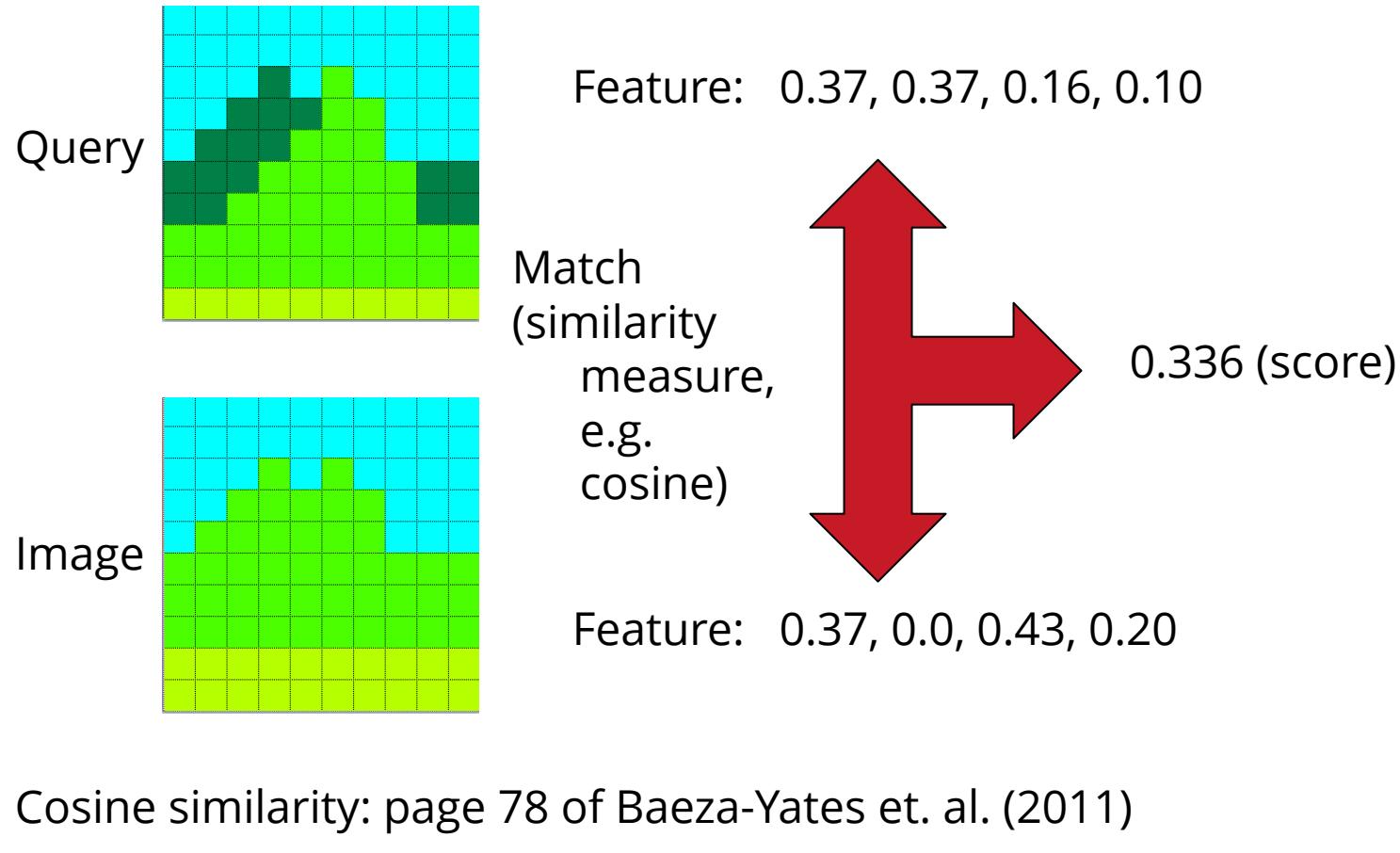
Feature vector: (0.37, 0.37, 0.16, 0.10)

(scaled by the image size of 10 by 10 pixels,
i.e. 37/100, 37/100, 16/100, 10/100)



Cosine similarity: page 78 of Baeza-Yates et. al. (2011)

Simple matching



- Examples of "global features" which can be extracted from images:
- Colour describes pixels (primitive feature)
 - Models include HSV (Hue, Saturation, Value) and RGB (Red, Green, Blue)
 - Define distribution of colour pixels in an image (histogram)
- Texture describes (small-scale) regions
 - Smoothness, contrast, regularity, directionality etc.
 - Often of limited value on its own
- Shape (image segmentation)
 - Represent the size, shape and orientation of objects (blocks & regions)
 - Represent relative position of objects

- **Query by example**
 - One or more example images
- **Query by sketch**
 - Allow the user to “sketch” an image, which is then used for search
- **Relevance feedback**
 - Can browse a collection by selecting one or more existing results, which will then be used as the query

Query by example

Google  JPG [houses of parliament](#)   

All **Images** Maps Shopping More ▾ Search tools SafeSearch on 

About 25,270,000,000 results (0.96 seconds)


Image size:
660 × 371
Find other sizes of this image:
All sizes - Small - Medium - Large

Best guess for this image: **[houses of parliament](#)**

Visit - UK Parliament
www.parliament.uk/visiting/ ▾
Visit the Houses of Parliament, watch committee hearings or take a tour of the Elizabeth Tower to see Big Ben.

Palace of Westminster - Wikipedia
https://en.wikipedia.org/wiki/Palace_of_Westminster ▾
The Palace of Westminster is the meeting place of the House of Commons and the House of Lords, the two houses of the Parliament of the United Kingdom.

Visually similar images 



Pages that include matching images

MPs will demand a heavy price for a clean Brexit – POLITICO
www.politico.eu/article/mps-will-demand-a-heavy-price-for-a-clean-brexit/ ▾
1160 × 772 - 3 Nov 2016 - In other words, a simple motion will not suffice — both houses

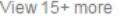
Big Ben 



Big Ben is the nickname for the Great Bell of the clock at the north end of the Palace of Westminster in London, and often extended to refer to the clock and the clock tower. [Wikipedia](#)

Height: 96 m
Opened: 1859
Architectural style: Gothic Revival architecture
Architects: Augustus Pugin, Charles Barry
London borough: City of Westminster

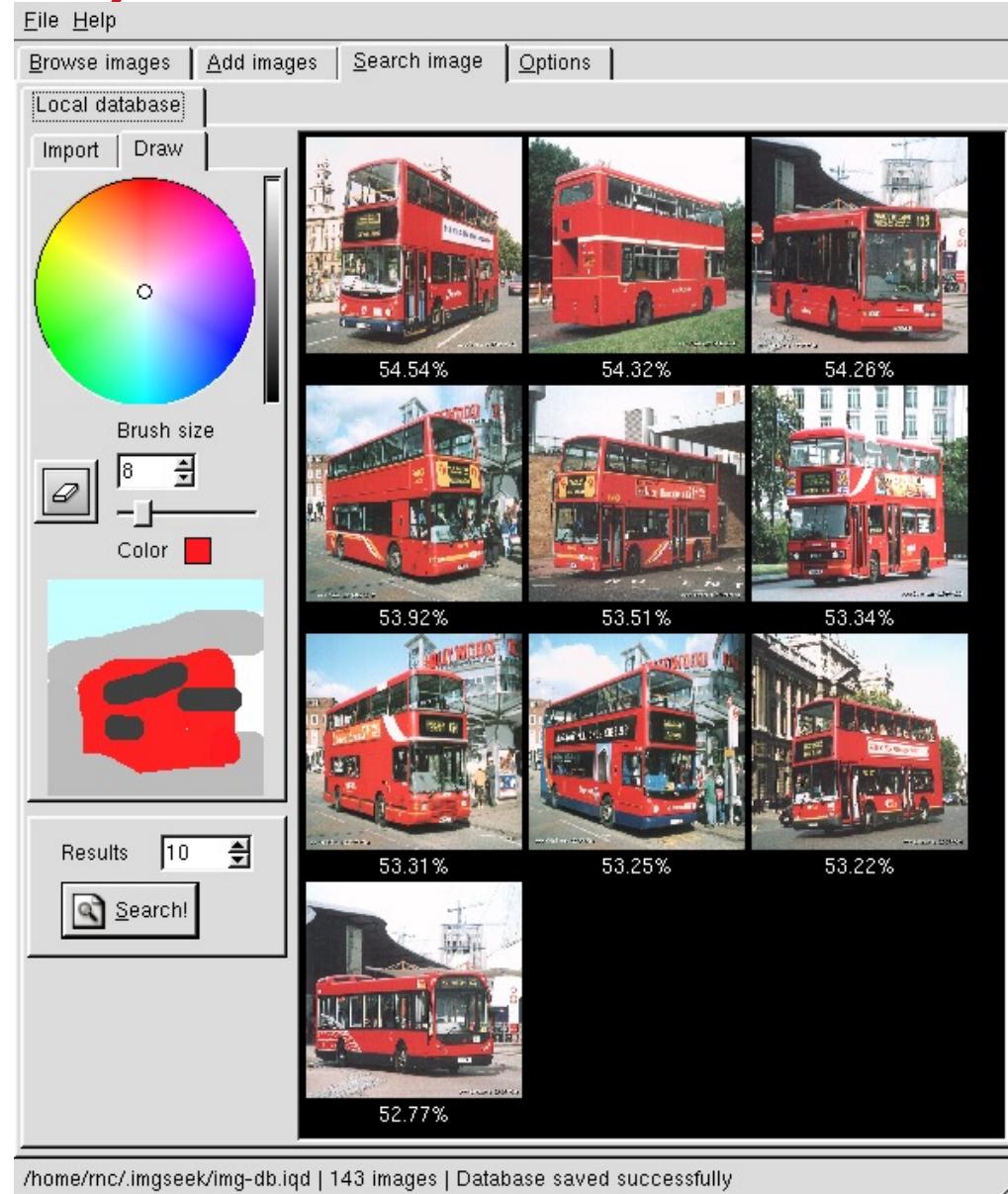
Profiles
 Twitter

People also search for 

    
London Eye Buckingham... Palace Tower Bridge Palace of Westmin... Abbey

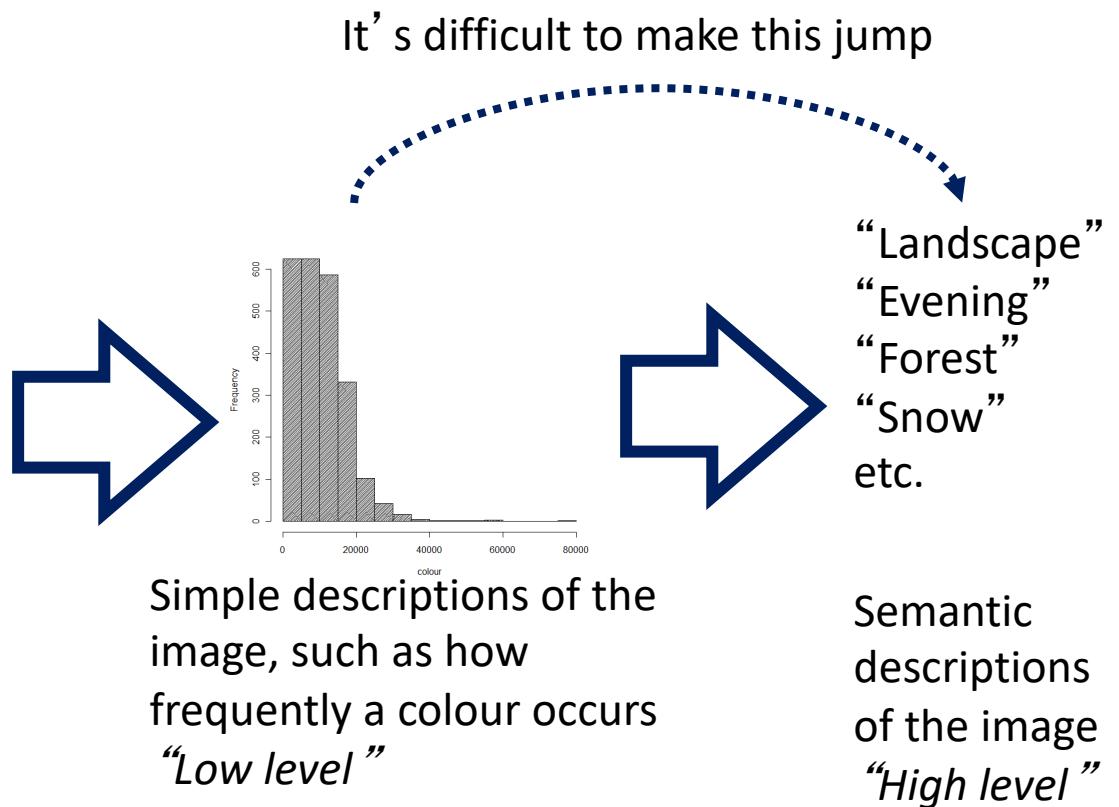
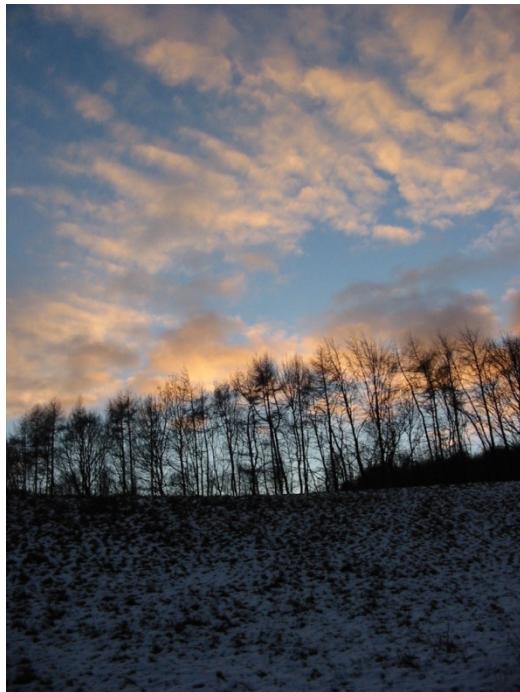
[Feedback](#)

Query by sketch (imgseek)



Big problem: The Semantic Gap

Probably the biggest challenge in multimedia retrieval



Other problems with content-based image retrieval

- Difficult to express information needs visually
 - May not always have an image at hand
 - The same objects may be visually very different
 - Drawing or sketching is difficult for most people
 - Difficult to express more abstract needs
- May require training
 - e.g. to find a picture of "Justin Bieber", need lots of examples of his image to train a classifier
- Generally poorer performance!

➤ 3. Audio retrieval

- Two main types:
 - Techniques which assume speech in the audio
 - Use Automatic Speech Recognition (ASR) to generate text transcripts
 - Techniques which do not assume speech!
 - E.g. Music retrieval

What is audio?

- Analogue audio signal which is sampled and digitised
 - E.g. MP3 music file, most digital audio
- Symbolic music
 - E.g. MIDI (Musical Instrument Digital Interface)
- [concentrate here on the former]

Speech recognition

- Much audio is speech based
 - Interviews, lectures, radio plays, etc.
- Can use Automatic Speech Recognition (ASR) systems to transcribe the speech
 - Feed out audio files to an ASR system
 - Index the resulting transcripts using text retrieval techniques
 - Query the text transcripts, which can include timestamps back into the audio stream

Problems and issues

- ASR systems have improved considerably, but their performance can still be less than ideal
 - The audio may contain a wide range of different accents, in varying audio quality
 - ASR systems work best with *training*
 - Ideally on individuals, but also on types of accents
 - An ASR system trained on American accents is unlikely to work well with Yorkshire accents
 - Audio often involves multiple speakers
 - E.g. consider a TV or radio interview, such as on news programs

Other associated technologies

- Speech discrimination
 - Aims to discriminate speech from non-vocal music or other sounds
 - Important to disable ASR system when there is no speech
- Speaker identification
 - Aims to identify who is speaking
 - Works from a set of known voices ("closed set identification")
 - Related area of speaker verification aims to determine whether a speaker is who he/she claims to be

Non-speech approaches

- Some audio does not consist of speech
 - E.g. Music, special effects, etc.
- Many approaches, including:
 - Audio classification
 - Query by Singing/Humming
 - Music fingerprinting

➤ 3. Video retrieval



Text-based video retrieval

The screenshot shows a YouTube search results page for the query "university of koblenz". The search bar at the top contains the text "university of koblenz". Below the search bar, there are several video thumbnails. The first video thumbnail on the left is for "Latest from Universität Koblenz" and features a collage of three people. The second video thumbnail on the right is for "Pflege studieren an der Uni Koblenz: Was erwartet mich?" and shows a man speaking. The third video thumbnail at the bottom is for "Universität Koblenz 2023 - Der Jahresrückblick" and shows a large red "UK" logo on a wall with the text "2023 Unser erstes Jahr". A "Subscribe" button is visible on the right side of the page.

YouTube DE

university of koblenz

Home

Shorts

Subscriptions

You >

Your channel

History

Playlists

Your videos

Your movies and TV

Watch Later

Liked videos

Subscriptions

Matthias Schwar... •
The Tonight Show... •
Cassie Kozyrkov
Yoyo Chinese
LastWeekTonight •
Browse channels

Explore

Trending

University of Koblenz

@universitaetkoblenz • 939 subscribers

Verantwortlich für den Inhalt und das Angebot: Der Präsident der Universität Koblenz Prof. Dr. Stefan Wehner
Universitätsstraße 1 ...

Subscribe

Pflege studieren an der Uni Koblenz: Was erwartet mich?

81 views • 1 month ago

Universität Koblenz

Schonmal dran gedacht, Pflege zu studieren? An der Universität Koblenz ist das möglich. Das Institut für Pflegewissenschaft ...

4K

Universität Koblenz 2023 - Der Jahresrückblick

1.2K views • 4 months ago

Universität Koblenz

Das erste Jahr als eigenständige Uni in Bildern: Ein Rückblick auf zwölf ereignisreiche Monate an der neuen Universität Koblenz.

+8 More

- Entire video is retrieved using assigned metadata
 - Popular for online video sharing systems (e.g. Youtube)
 - Suffice for short videos, but offer limited interactivity with the video itself
- Problems:
 - What if your videos are very long?
 - Looking for part of a video, buried somewhere in the middle?
 - Who assigns the metadata?
 - On YouTube, it's the user who uploads the video

- Aims to work automatically, without human effort
 - No need for a human to assign metadata
- Indexes and retrieves sections of videos, rather than whole videos
 - Tries to give direct access to the most relevant parts of a video
- Provides ways in which a video can be visualised without watching it
 - Watching video takes time (even if it's enjoyable)

What is a video?

1. A sequence of image frames



2. An associated audio track, which can contain speech, music, effects, etc.

“The quick brown fox jumped over the lazy dog . . . ”

Content-based video retrieval typically uses the visual frames, plus the speech, to index and search videos

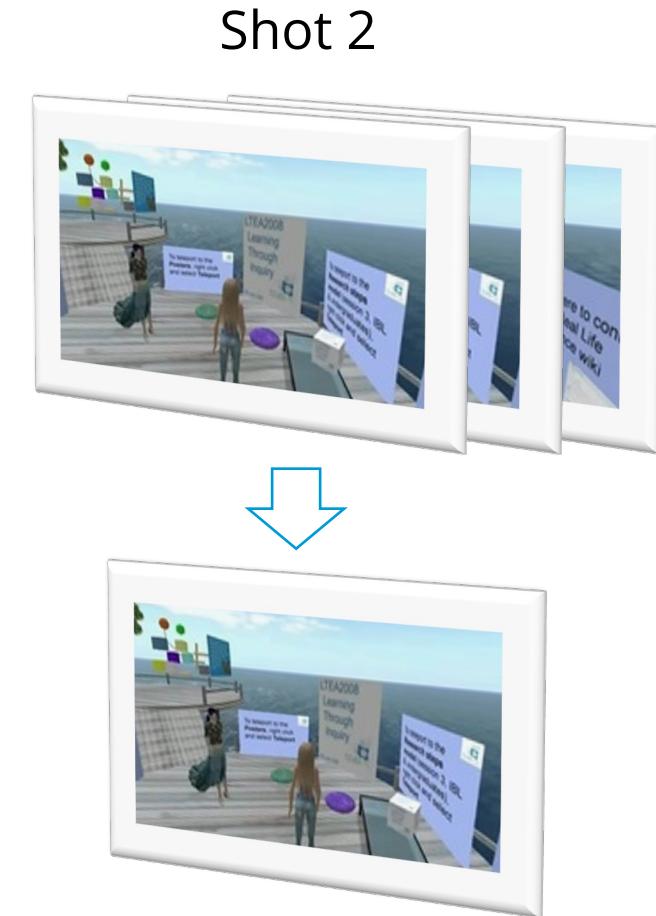
Shot segmentation

- Splits the video into sequences of visually related frames, called *shots*
 - Pairs of frames are compared, large differences between frames indicate shot boundaries
 - Some shot boundaries are more difficult to detect, such as gradual fades from one shot to the next



Keyframes: Representing shots

- To visually represent each shot, one or more frames can be extracted
 - Often take the “iframes” from the compressed mpeg video stream
 - Called “keyframes”
- A single image can then be used to represent the shot when searching or visualising the video



Frame 2 of shot 2

Audio stream: text transcripts

- A text transcript of the audio stream is often an important source of information for retrieval
 - Automatic Speech Recognition (ASR) systems can be used
 - Can be very unreliable without training (as previously mentioned)
 - Closed captions (subtitles)
 - Much more reliable, but not always available
 - e.g. TV rushes, pre-production video
- Not all videos have speech
 - E.g. some cartoons for pre-school children, the films of “Charlie Chaplain”, etc.

Representing shots using a text transcript

A text transcript, such as generated by an ASR system, has timestamps for each word which can be aligned with the shots



Summary: What do we have?

- A video has been segmented into a sequence of shots
- Each shot can be represented by:
 - One or more *keyframe* images
 - Associated text from a transcript

- Summary of a video created to help users make sense of a video without watching it
 - Keyframe
 - Storyboards
 - Collage of image
 - Fast-forward/skimming
 - Keywords/descriptions

What can we do now?

- Search:
 - Can now index this information to allow it to be searched
 - The transcripts can be searched by text
 - The keyframes can be visually searched or browsed using image retrieval methods
 - E.g. "find me another shot visually similar to *this one*"
- Presentation/visualisation
 - Keyframes and transcripts can be used visualise videos and parts of videos

Keyframes, title, description, tags



Mediating Awareness and Communication through Digital but Physical Surrogates (1996)

Video demonstration from the 1999 CHI conference. Digital but physical surrogates are tangible representations of remote people positioned within an office and under digital control. Surrogates selectively collect and...

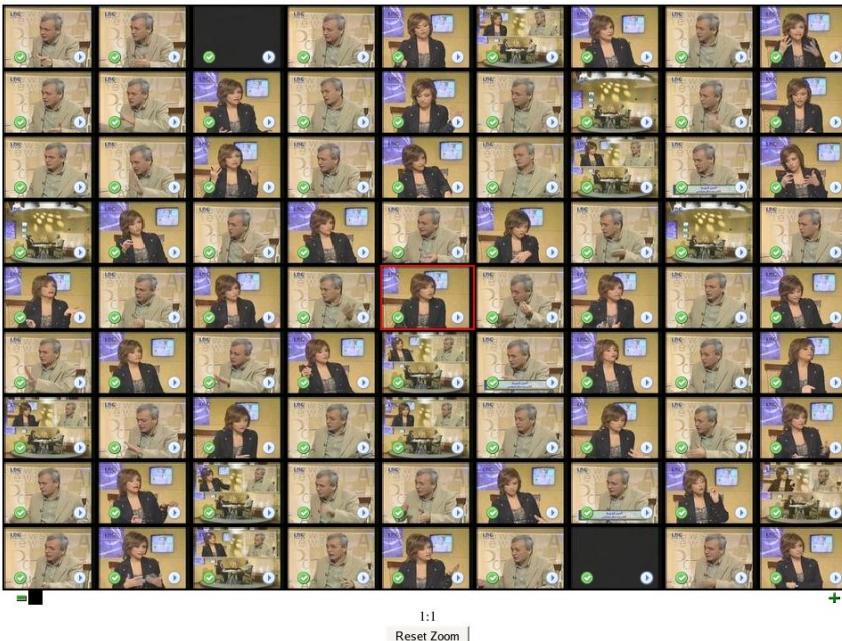
Genre: Educational

Keywords: CHI; UID; Surrogates; Remote Sensing

Duration: 00:06:53

Popularity (downloads): 377

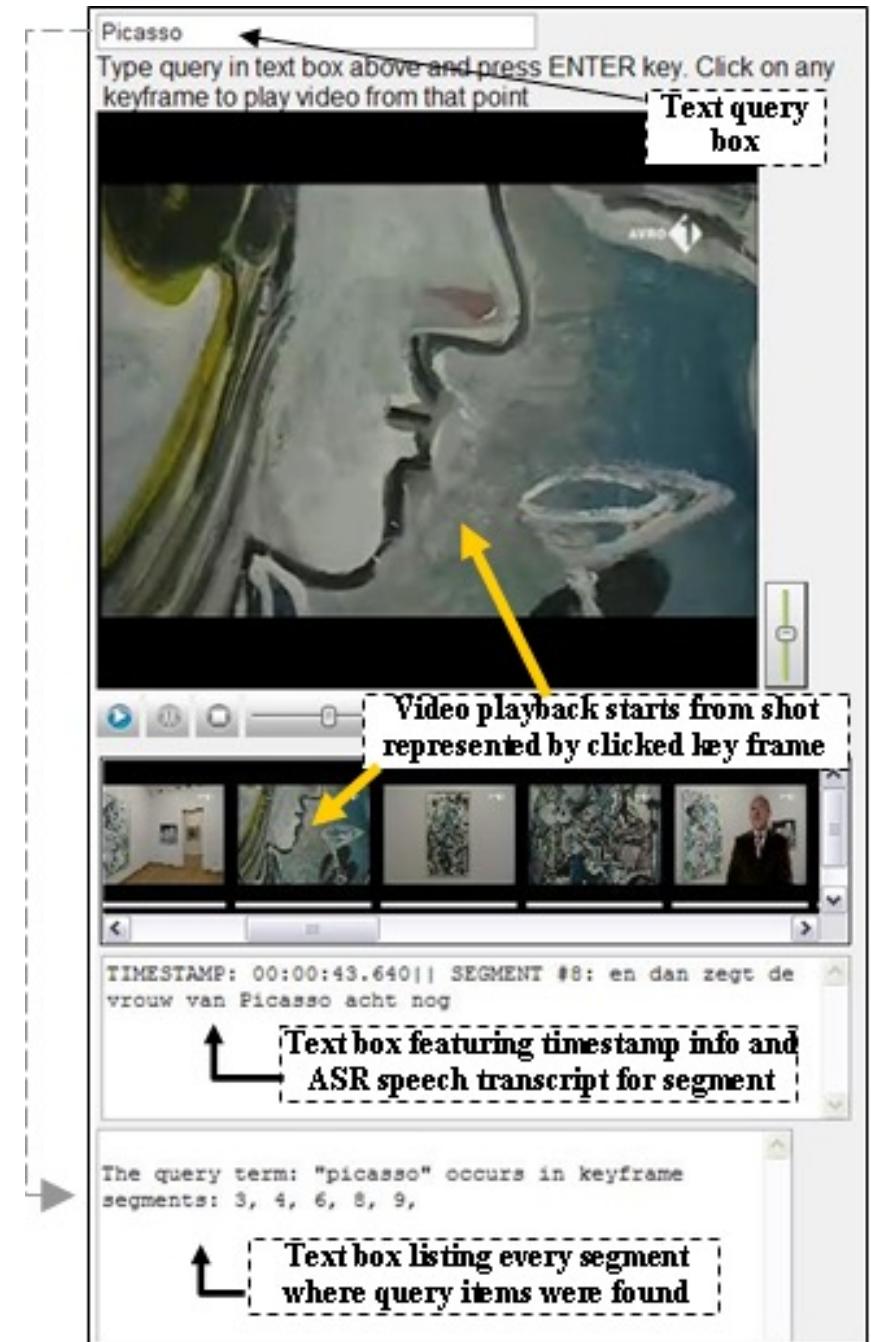
Storyboards

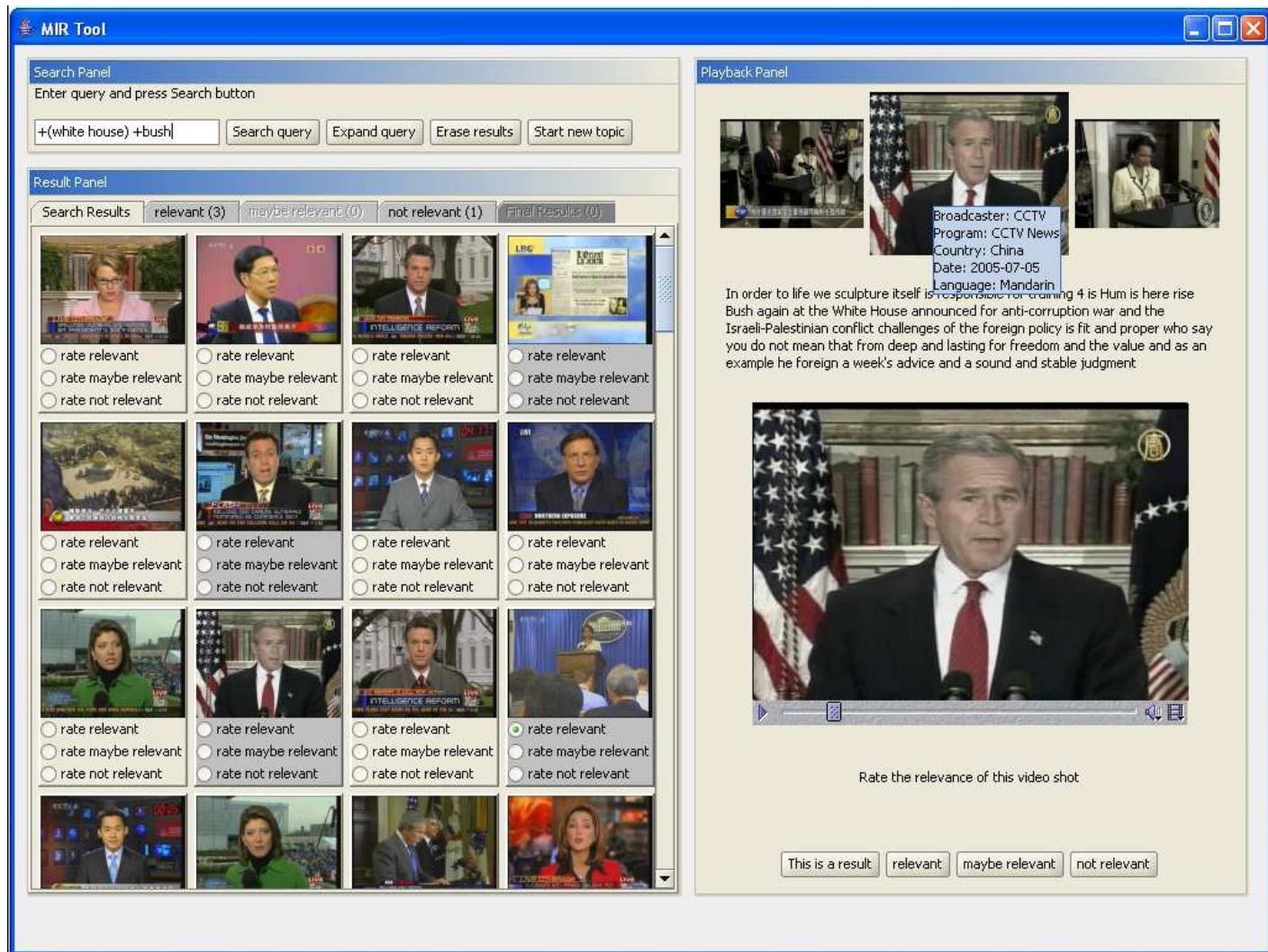


Peter Wilkins, Raphael Troncy, Martin Halvey,
Darragh Byrne, Alia Amin, P. Punitha, Alan Smeaton,
Robert Villa. User Variance and its Impact on Video
Retrieval Benchmarking ACM International
Conference on Image and Video Retrieval
(CIVR 2009), July 2009

Examples search interfaces

- Screenshot of the Multimatch Video Interface (with selected video document already loaded)
 - Carmichael, J., Larson, M., Marlow, J., Newman, E., Clough, P., Oomen, O., and Sav, S. (2008), Multimodal Indexing of Digital Audio-visual Documents: A Case Study for Cultural Heritage Data, In Proceedings of the Sixth International Workshop on Content-Based Multimedia Indexing (CBMI2008), London, UK, 18-20th June, pp. 93-100.





Challenges

- Similar to those for image retrieval
- Robust shot boundary detection
 - Combing audiovisual cues (over low-level features alone)
- Story segmentation
 - Context and application-dependent
- High-level feature and semantic concept extraction
 - e.g. using the 1,000 concepts from LSCOM (Large Scale Concept Ontology for Multimedia)
- Intuitive search and browse interfaces
 - Personalising search

➤ 4. Summary



- Multimedia is everywhere, it's ubiquitous
 - People want access to it
- Multimedia retrieval is typically based on
 - Indexing associated metadata (often simpler)
 - Processing the content of the media (harder)
- There are many challenges in multimedia search
- Online applications provide limited interactivity
- Future work is moving towards:
 - Integrating text and visual features and improving user's experience
 - Improving automatic annotation of concept with high-level concepts