



# Multimedia Information Retrieval



Prof Stefan Ruger  
Multimedia and Information Systems  
Knowledge Media Institute  
The Open University  
<http://kmi.open.ac.uk/mmis>

# MMIS

Multimedia and Information Systems



## Working at KMi

[Job Vacancies \[2\]](#)

[Studentships](#)

[Visitors / Internships](#)

## About KMi

[Contact Us](#)

[Find Us](#)

[Events](#)

[Concept Videos](#)

## KMi Tools



[Compendium](#)



[FM Technology](#)



[Magpie](#)

# About



The Knowledge Media Institute (KMi) was set up in 1995 in recognition of the need for the Open University to be at the forefront of research and development in a convergence of areas that impacted on the OU's very nature: Cognitive and Learning Sciences, Artificial Intelligence and Semantic Technologies, and Multimedia. We chose to call this convergence Knowledge Media.

Knowledge Media is about the processes of generating, understanding and sharing knowledge using several different media, as well as understanding how the use of different media shape these processes.

KMi operates with reference to a number of basic tenets, which define the context in which we formulate and pursue our research objectives:

## Strategic Threads

Our research is aligned with a number of broad strategic threads, currently [Future Internet](#), [Knowledge Management](#), [Multimedia & Information Systems](#), [Narrative Hypermedia](#), [New Media Systems](#), [Semantic Web & Knowledge Services](#) and [Social Software](#).

[Future Internet](#)

[Knowledge Management](#)

[Multimedia & Information Systems](#)

[Narrative Hypermedia](#)

[New Media Systems](#)

[Semantic Web & Knowledge Services](#)

[Social Software](#)





## Projects

All Projects [117]

● [Hot](#) [50]

● [Active](#) [17]

● [Classics](#) [50]

## Research Themes

[Future Internet](#) [4]

[Knowledge Management](#) [9]

[Multimedia & Information Systems](#) [6]

[Narrative Hypermedia](#) [13]

[New Media Systems](#) [9]

[Semantic Web & Knowledge Services](#) [16]

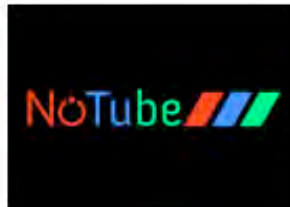
[Social Software](#) [9]

# Projects

ABCDEFGHIJKLMNOPQRSTUVWXYZ

## Projects | Hot

SPOTLIGHTED HOT PROJECTS



### NoTube

[Future Internet](#) [Semantic Web and Knowledge Services](#)

*Networks and Ontologies for the Transformation and Unification of Broadcasting and the intErnet*

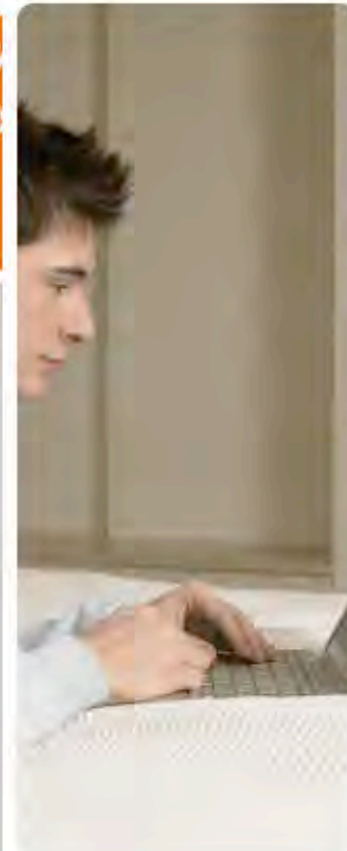


<http://www.mmkm.org>

### The UK Multimedia Knowledge Management Network

[Multimedia and Information Systems](#)

*Enhance communication between the experts in both academia and industry*



[Future Internet](#)

[Knowledge Management](#)

[Multimedia & Information Systems](#)

[Narrative Hypermedia](#)

[New Media Systems](#)

[Semantic Web & Knowledge Services](#)

[Social Software](#)



Since 1995: 117 projects & 67 technologies

Current year

17 live projects

typically £2.5m ext, £1m internal

- 10 EU
- 3 UK
- 1 US
- 3 internal (iTunes U, SocialLearn)



1. What is multimedia information retrieval?
2. Metadata and piggyback retrieval
3. Multimedia fingerprinting
4. Automated annotation
5. Content-based retrieval



# What is Multimedia?

Within this lecture:

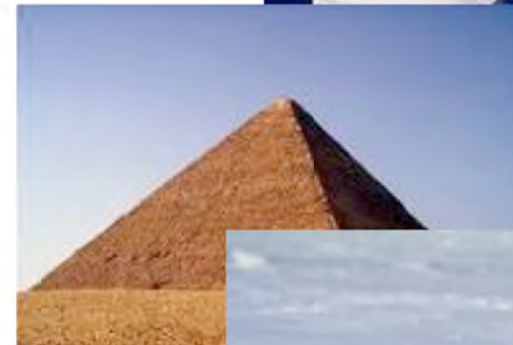
One or more media

Possibly interlinked

Digital

For communication

(not only entertainment)





The Twelve Collegia building on Vasilievsky Island in Saint Petersburg is the university's main building and the seat of the rector and administration (the building was constructed on the orders of Peter the Great)



“twelve collegia building”

Google Images



Bing Images



Flickr



Yahoo Images



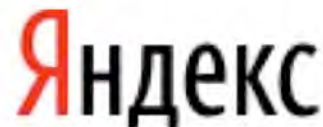
ImageToss



Yandex



Яндекс







Best current practice is a text search:  
Find text in filename, anchor text, caption, ...

Text search works by creating a large index:

368

## GENERAL INDEX

---

Henderson **85**  
Henderson, Louise 30  
Henderson Valley wine 35  
Henley Lake Park (Masterton) 171  
Heritage Expeditions 337  
Heritage trails  
- Buller Coalfields 233  
Hokitika Heritage Trail 235  
Hermitage (Mount Cook) 250, 307

Fiesta (Hamilton) 40, 116  
Hot springs  
Hanmer Springs 231  
Maruia Springs Thermal Resort 231  
Hot Water Beach 123  
Ketetahi Hot Springs 140  
Miranda Hot Springs 120  
Mokoia Island 134  
Morere Hot Springs 131



# New search types

	text	image	location	speech	sound	humming	motion	query / doc
	●							text
					●			video
								images
	●							speech
						●		music
								sketches
								multimedia

## Example

conventional  
 text retrieval  
 you type "bar" and  
 get a wildlife  
 documentary  
 type "radio" and  
 get BBC  
 radio news  
 and get a  
 music piece



Organise yourself in groups

Discuss with neighbours

- *Two* Examples for different query/doc modes?
- How hard is this? Which techniques are involved?
- *One* example combining different modes



# Exercise

text	image	location	speech	sound	humming	motion	query / doc	Discuss
							text	- 2 examples
							video	- How hard is it?
							images	- 1 combination
							speech	
							music	
							sketches	
							multimedia	



# Leaf detection What are the challenges?





# Venation pattern and shape

Shape is key





# The semantic gap



1m pixels with a spatial  
colour distribution

faces & vase-like object



# Polysemy







# Multimedia information retrieval

1. What is multimedia information retrieval?
2. Metadata and piggyback retrieval
3. Multimedia fingerprinting
4. Automated annotation
5. Content-based retrieval



## Dublin Core

simple common denominator: 15 elements such as title, creator, subject, description, ...

## METS

Metadata Encoding and Transmission Standard

## MARC 21

MAchine Readable Cataloguing (harmonised)

## MPEG-7

Multimedia specific metadata standard



- Moving Picture Experts Group “Multimedia Content Description Interface”
- Not an encoding method like MPEG-1, MPEG-2 or MPEG-4!
- Usually represented in XML format
- Full MPEG-7 description is complex and comprehensive
- Detailed Audiovisual Profile (DAVP)

[P Schallauer, W Bailer, G Thallinger, “A description infrastructure for audiovisual media processing systems based on MPEG-7”, Journal of Universal Knowledge Management, 2006]



```
<Mpeg7 xsi:schemaLocation="urn:mpeg:mpeg7:schema:2004 ./davp-2005.xsd" ... >
<Description xsi:type="ContentEntityType">
  <MultimediaContent xsi:type="AudioVisualType">
    <AudioVisual>
      <StructuralUnit href="urn:x-mpeg-7-pharos:cs:AudioVisualSegmentationCS:root"/>
      <MediaSourceDecomposition criteria="kmi image annotation segment">
        <StillRegion>
          <MediaLocator><MediaUri>http://...392099.jpg</MediaUri></MediaLocator>
          <StructuralUnit href="urn:x-mpeg-7-pharos:cs:SegmentationCS:image"/>
          <TextAnnotation type="urn:x-mpeg-7-pharos:cs:TextAnnotationCS:
            image:keyword:kmi:annotation_1" confidence="0.87">
            <FreeTextAnnotation>tree</FreeTextAnnotation>
          </TextAnnotation>
          <TextAnnotation type="urn:x-mpeg-7-pharos:cs:TextAnnotationCS:
            image:keyword:kmi:annotation_2" confidence="0.72">
            <FreeTextAnnotation>field</FreeTextAnnotation>
          </TextAnnotation>
        </StillRegion>
      </MediaSourceDecomposition>
    </AudioVisual>
  </MultimediaContent> </Description> </Mpeg7>
```



```
<Mpeg7 xsi:schemaLocation="urn:mpeg:mpeg7:schema:2004 ./davp-2005.xsd" ... >
<Description xsi:type="ContentEntityType">
  <MultimediaContent xsi:type="AudioVisualType">
    <AudioVisual>
      <StructuralUnit href="urn:x-mpeg-7-pharos:cs:AudioVisualSegmentationCS:root"/>
      <MediaSourceDecomposition criteria="kmi image annotation segment">
        <StillRegion>
          <MediaLocator><MediaUri>http://...392099.jpg</MediaUri></MediaLocator>
          <StructuralUnit href="urn:x-mpeg-7-pharos:cs:SegmentationCS:image"/>
          <TextAnnotation type="urn:x-mpeg-7-pharos:cs:TextAnnotationCS:
            image:keyword:kmi:annotation_1" confidence="0.87">
            <FreeTextAnnotation>tree</FreeTextAnnotation>
          </TextAnnotation>
          <TextAnnotation type="urn:x-mpeg-7-pharos:cs:TextAnnotationCS:
            image:keyword:kmi:annotation_2" confidence="0.72">
            <FreeTextAnnotation>field</FreeTextAnnotation>
          </TextAnnotation>
        </StillRegion>
      </MediaSourceDecomposition>
    </AudioVisual>
  </MultimediaContent> </Description> </Mpeg7>
```



Manage document repositories and their metadata

Greenstone digital library suite

<http://www.greenstone.org/>

interface in 50+ languages (documented in 5)

knows metadata

understands multimedia

XML or text retrieval



# Piggy-back retrieval

text	image	location	speech	sound	humming	motion	query
							doc
							text
							video
							images
							speech
							music
							sketches
							multimedia

text



# Music to text

0 +7 0 +2 0 -2 0 -2 0 -1 0 -2 0 +2 -4

Z G Z B Z b Z b Z a Z b Z B d

ZGZB

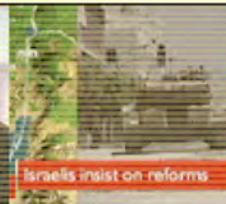
GZBZ

ZBZb

[with Doraisamy, J of Intellig Inf Systems 21(1), 2003; Doraisamy PhD thesis 2004]



Automatic  
News  
Summarization  
Extraction  
System



Search news:

Go

Sort by:  Date  Relevance

From: 1 ▾ Jan ▾ 2003 ▾

To: 3 ▾ Jun ▾ 2008 ▾

[technology licensed by Imperial Innovations]

[patent 2004]

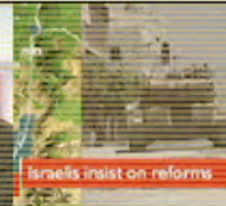
[finished PhD: Pickering]

[with Wong and Pickering, CIVR 2003]

[with Lal, DUC 2002]

[Pickering: best UK CS student project 2000 – **national prize**]

Automatic  
News  
Summarization  
Extraction  
System



Search news:

Sort by:  Date  Relevance

From:

To:

[technology licensed by Imperial Innovations]

[patent 2004]

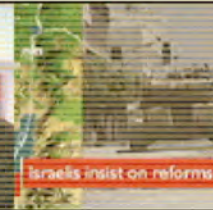
[finished PhD: Pickering]

[with Wong and Pickering, CIVR 2003]

[with Lal, DUC 2002]

[Pickering: best UK CS student project 2000 – **national prize**]

Automatic  
News  
Summarization  
Extraction  
System



Search news:

Go

Sort by:  Date  Relevance

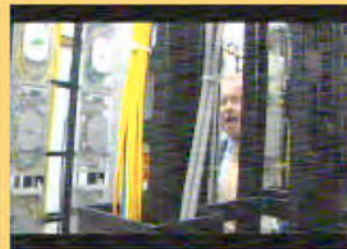
From:

To:

Search results **1** to **10** (out of **23**) for **microsoft**

<1-10> | <11-20> | <21-23>

Organisations  People  Locations  Dates



Play this story



Browse other news on Sun May 4 2008

**Organisations:** AOL,  
**Microsoft**, Police,  
Yahoo

**Date :** Sun May 4 2008

**Length :** 217.65 seconds

**Full Story :** [Link](#)

**People:** Bill, Jay, Leah,  
Paul Ross, Warner, bo,  
ina, olin

**Summary :** **Microsoft** has pulled out of a deal to buy Yahoo, the offer was rejected because it wasn't enough. In trying to buy Yahoo, **Microsoft** wanted to set up a rival to google, which dominates the internet advertising. While some Yahoo executives

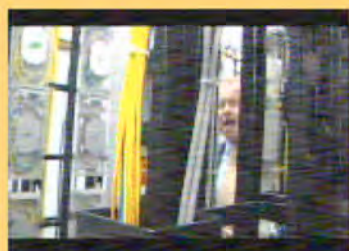
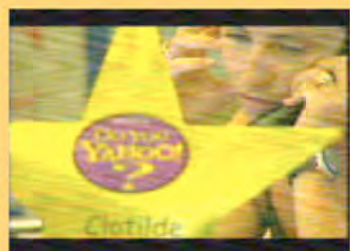
**Locations:** Britain,

might be celebrating their continued independence today, having seen off **Microsoft's**

Search results **1** to **10** (out of **23**) for **microsoft**

&lt;1-10&gt; | &lt;11-20&gt; | &lt;21-23&gt;

Organisations | People | Locations | Dates



[Play this story](#) [Browse other news on Sun May 4 2008](#)

**Organisations:** AOL, **Microsoft**, Police, Yahoo

**People:** Bill, Jay, Leah, Paul Ross, Warner, bo, ina, olin

**Locations:** Britain, Glasgow

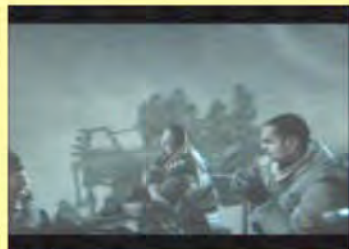
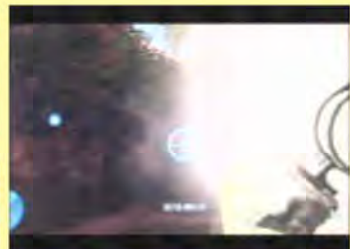
**Dates:** today, tomorrow, yesterday evening

**Date :** Sun May 4 2008

**Length :** 217.65 seconds

**Full Story :** [Link](#)

**Summary :** **Microsoft** has pulled out of a deal to buy Yahoo, the offer was rejected because it wasn't enough. In trying to buy Yahoo, **Microsoft** wanted to set up a rival to google, which dominates the internet advertising. While some Yahoo executives might be celebrating their continued independence today, having seen off **Microsoft's** unwanted attentions, they might already been dreading stock markets pening tomorrow. Both **Microsoft** and Yahoo have come a long way since being ormed in garages, both sets have earned billions along the way. Alternative Leah yahoo may look merge with AOL, owned by Time mre whAO own by ime Warner, but it would have to fast, because AOL might also be under **Microsoft's** radar.



[Play this story](#) [Browse other news on Tue Sep 25 2007](#)



1. What is multimedia information retrieval?
2. Metadata and piggyback retrieval
3. **Multimedia fingerprinting**
4. Automated annotation
5. Content-based retrieval



## Near-duplicate detection: Cool access mode!



MORGAN & CLAYPOOL PUBLISHERS

# Multimedia Information Retrieval

Stefan Rüger

*SYNTHESIS LECTURES ON INFORMATION  
CONCEPTS, RETRIEVAL, AND SERVICES*

Gary Marchionini, *Series Editor*

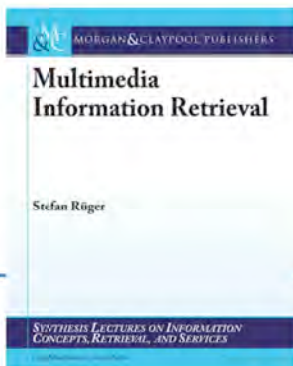
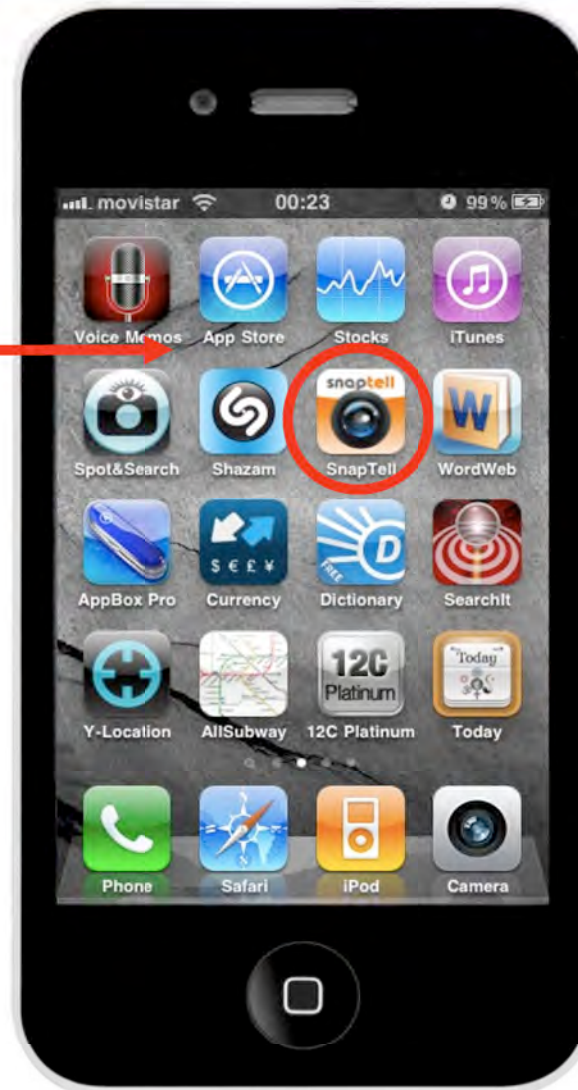
Copyright © 2010 Morgan & Claypool Publishers



The Open  
University

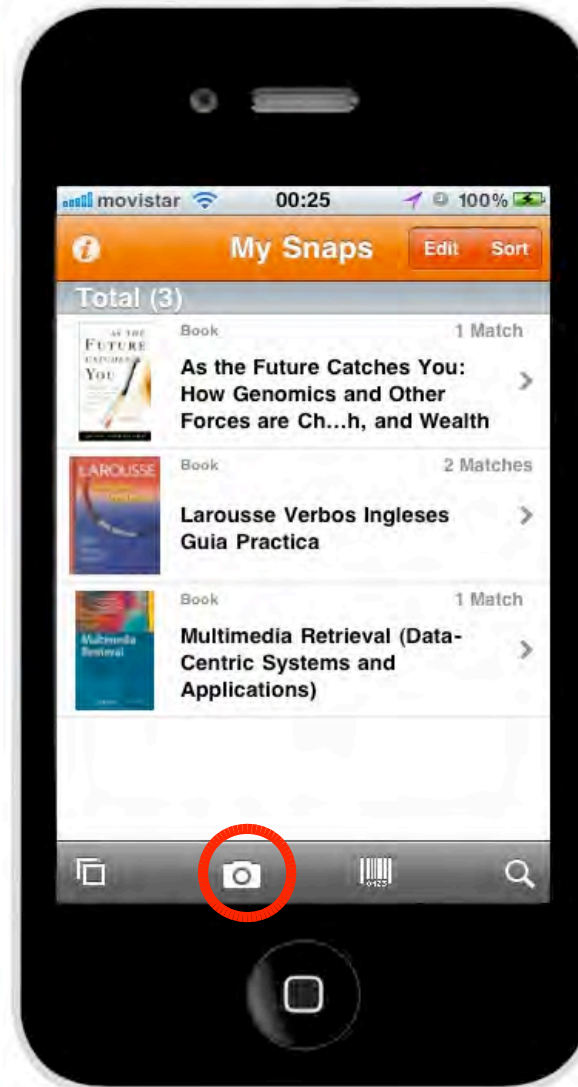


# Snaptell: Book, CD and DVD covers





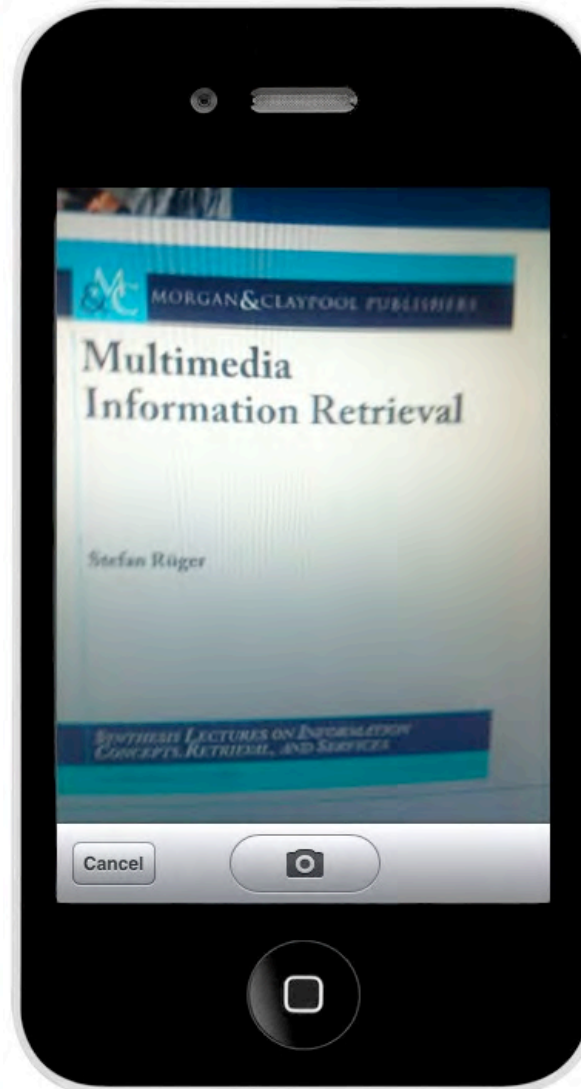
# Snaptell: Book, CD and DVD covers







# Snaptell: Book, CD and DVD covers





# SnapTell: Book, CD and DVD covers



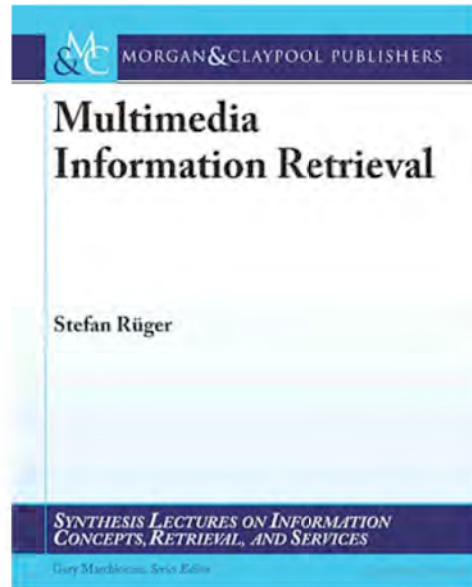


# Snaptell: Book, CD and DVD covers





# Link from real world to databases





# The Open University's Spot & Search



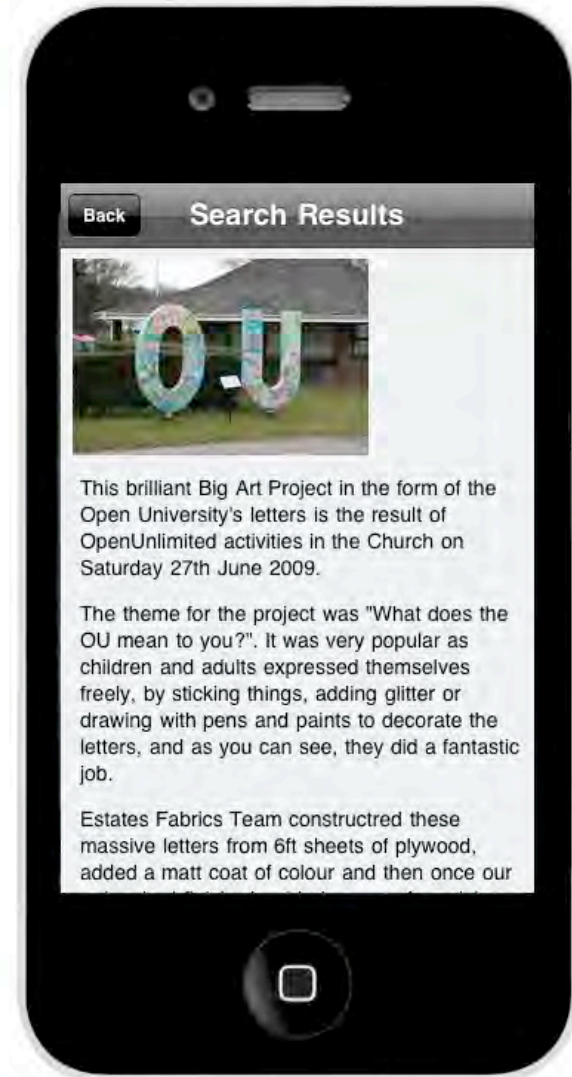
**Scott Forrest:**  $E=MC$  squared

"Between finished surface texture and raw quarried stone. Between hard materials and soft concepts. Between text and context."

[More information](#)



# Spot & Search





## Near duplicate detection

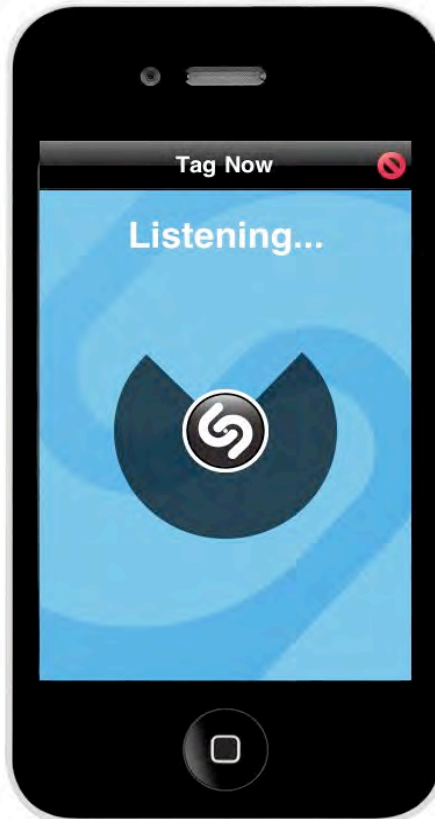
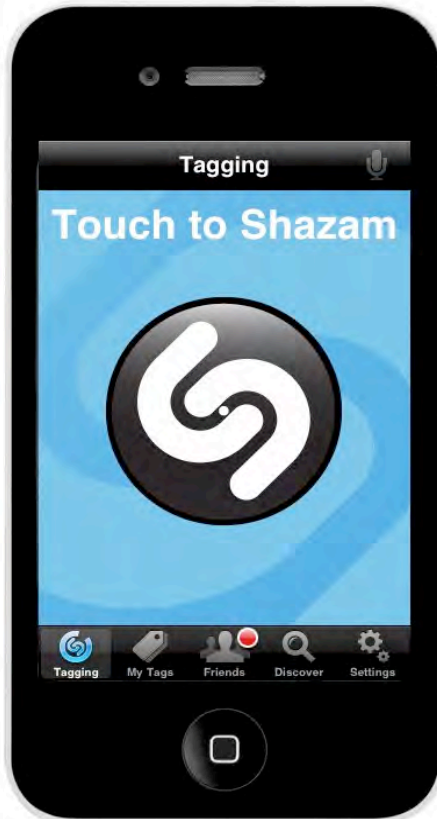
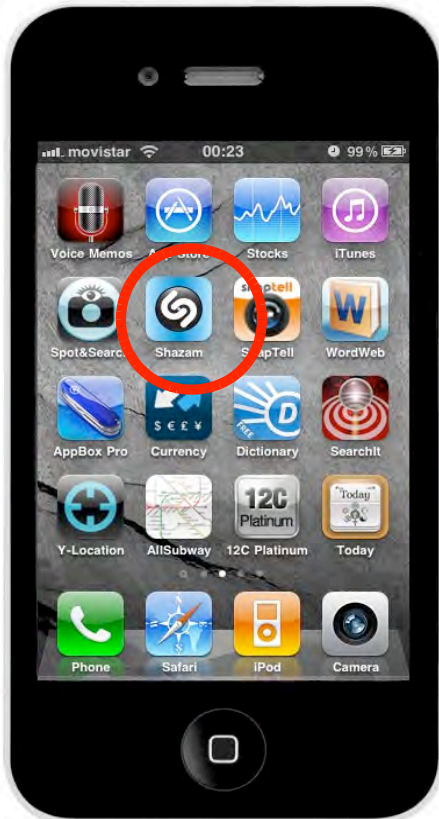
Works well in 2d: CD covers, wine labels, signs, ...

Less so in near 2d: buildings, vases, ...

Not so well in 3d: faces, complex objects, ...



# Shazam







Find applications for near-duplicate detection

- be imaginative: the more “outrageous” the better
- can be other media types (audio, smells, haptic, ...)
- can be hard to do



## Near-duplicate detection Where are the challenges?



[Victoria and Albert museum, London, ceramics collection, 2010]

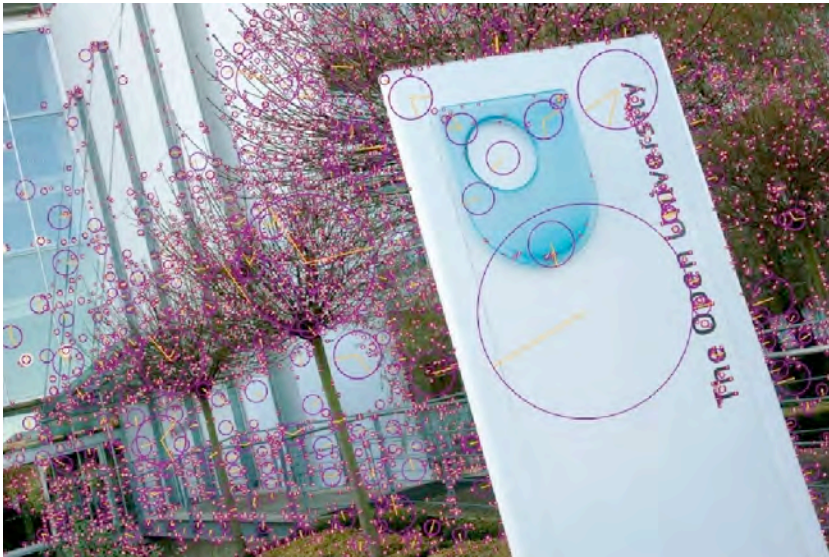


## Fingerprinting technique

- 1 Compute salient points
- 2 Extract “characteristics” from vicinity (feature)
- 3 Make invariant under rotation & scaling
- 4 Quantise: create visterms
- 5 Index as in text search engines
- 6 Check/enforce spatial constraints after retrieval



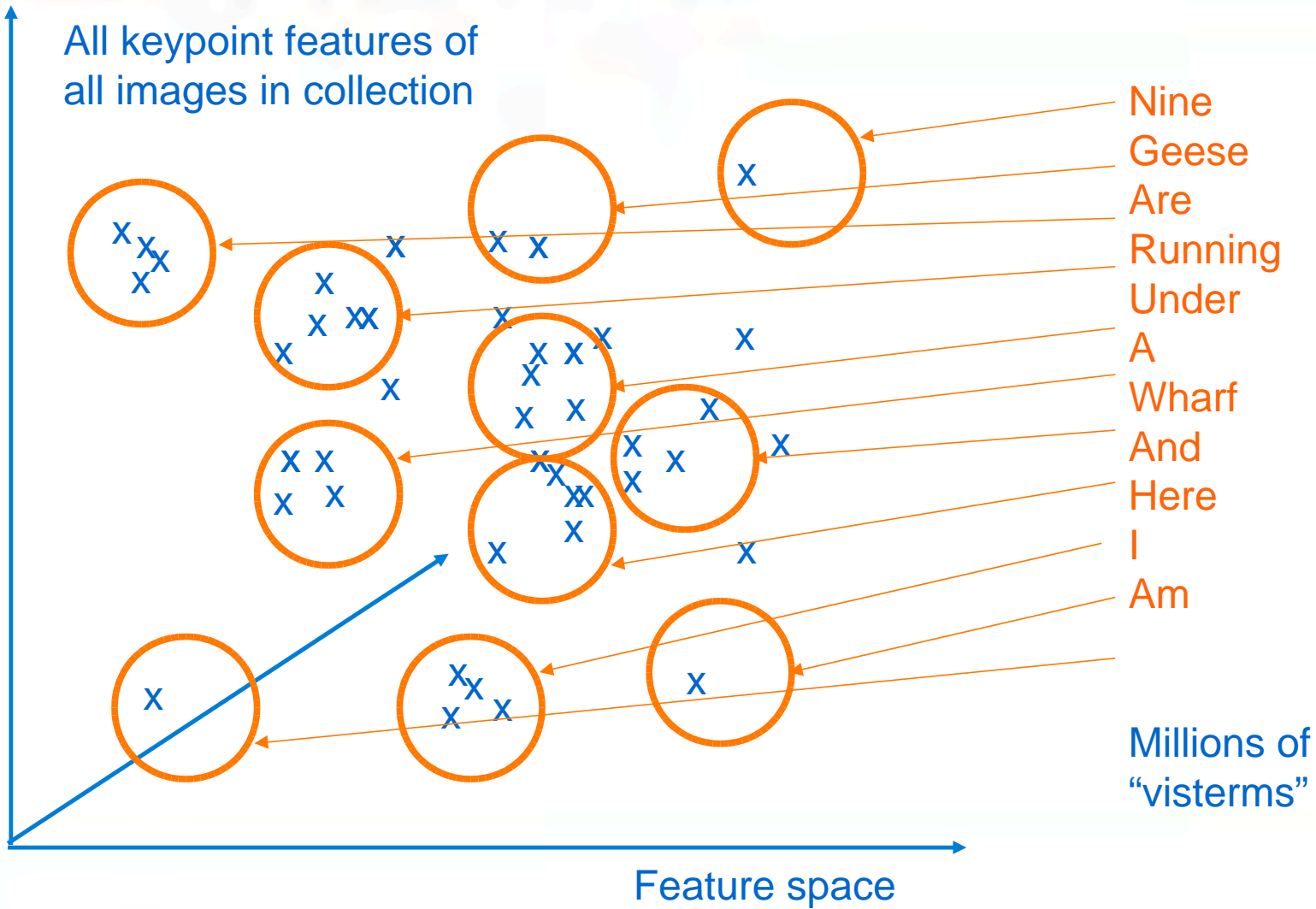
## NDD: Compute salient points and features



Eg, SIFT features: each salient point described by a feature vector of 128 numbers; the vector is invariant to scaling and rotation

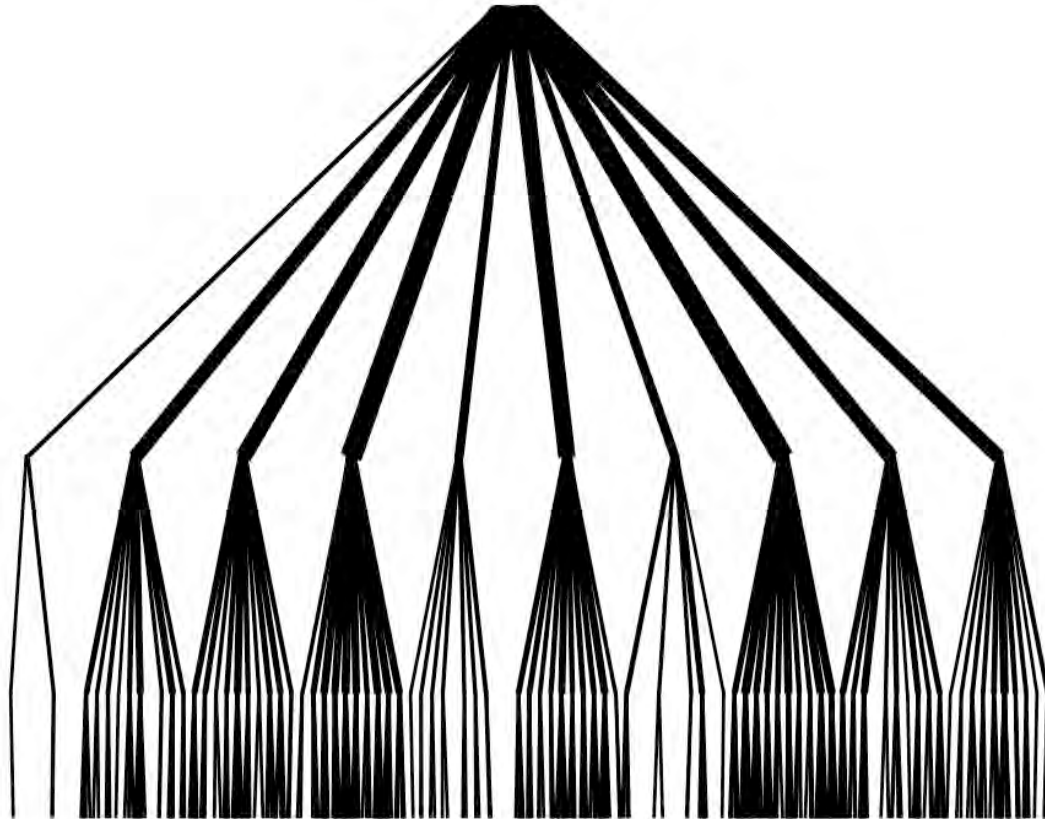


# NDD: Keypoint feature space clustering





# Clustering Hierarchical k-means





## NDD: Encode all images with visterms



Jkjh Geese Bjkj Wharf  
Ojkkjhjhj Kssn Klkekjl Here  
Lkjkll Wjjkll Kkjlk Bnm  
Kllkgjg Lwoe Boerm ...



## NDD: query like text

At query time compute salient points,  
keypoint features and visterms

Query against database of images  
represented as bag of visterms

### Query



Joiu Gddwd Bipoi Wueft  
Oiooiuui Kwwn Kpodoip Hdfd  
Loiopp Wiiopp Koipo Bnm  
Kppoyiy Lsld Bldfm ...





# NDD: Check spatial constraints

social

sociallearn

Learning in an Open World  
Open University: Online Conference, June 2010

CC

Web2.0

SocialLearn Team: Simon Buckingham Shum, Rebecca Ferguson, Mark Glaister, Thanh Le

<http://www.open.ac.uk/olconf2010/>



The Open University





## Fingerprinting technique

- 1 Compute salient points
- 2 Extract “characteristics” from vicinity (feature)
- 3 Make invariant under rotation & scaling
- 4 Quantise: create visterms
- 5 Index as in text search engines
- 6 Check/enforce spatial constraints after retrieval



## How Shazam works - Spectrogram

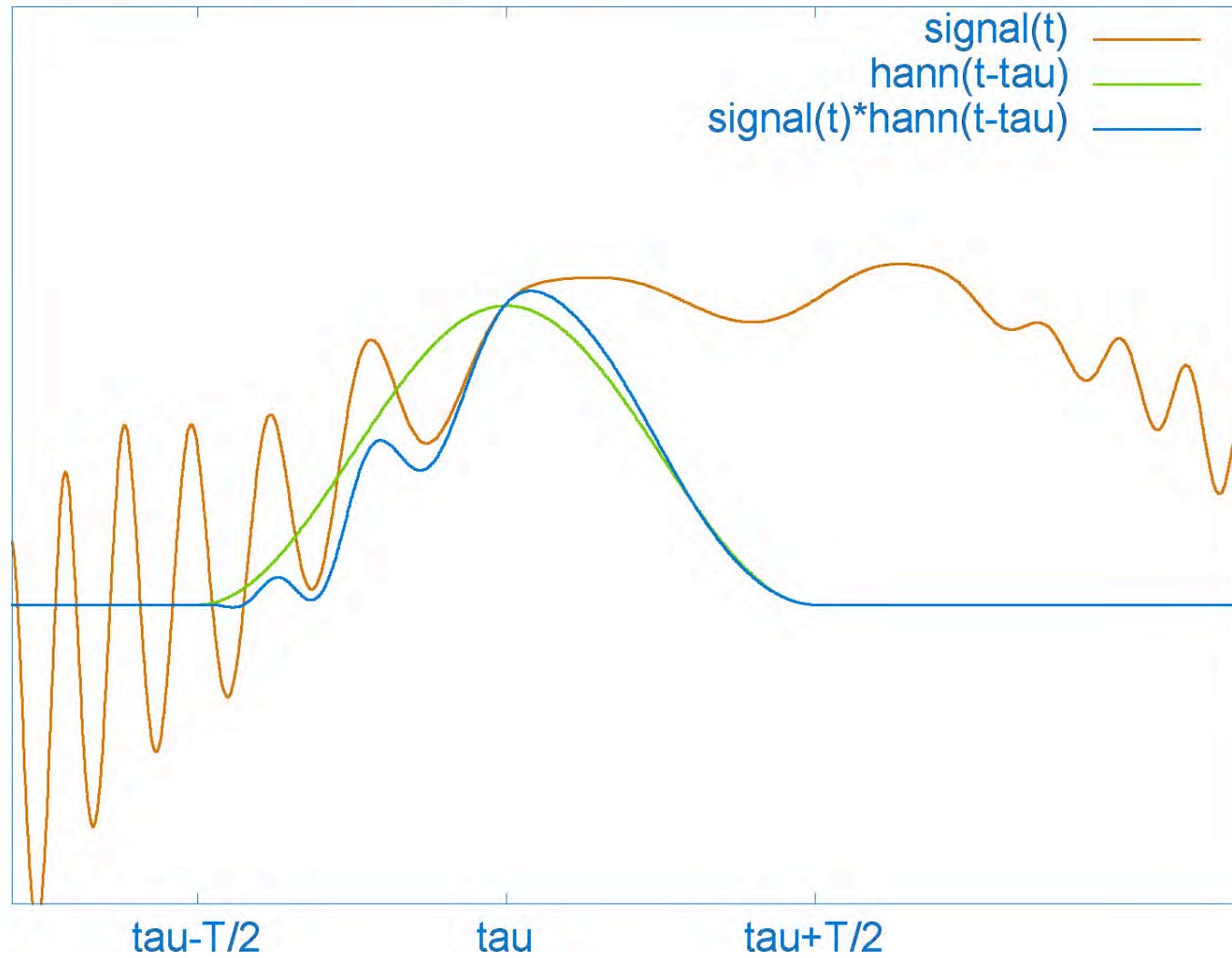
Compute energy for all (frequency,time) pairs  
 using a Fourier transform under a Hann window  $w$

$$\text{spectrogram}(f, \tau) = \left| \int_{-\infty}^{\infty} s(t)w(t - \tau)e^{jft} dt \right|^2$$

$$w(t) = \begin{cases} \frac{1}{2} + \frac{1}{2} \cos\left(\frac{2\pi t}{T}\right) & \text{if } t \in [-T/2, T/2] \\ 0 & \text{otherwise} \end{cases}$$

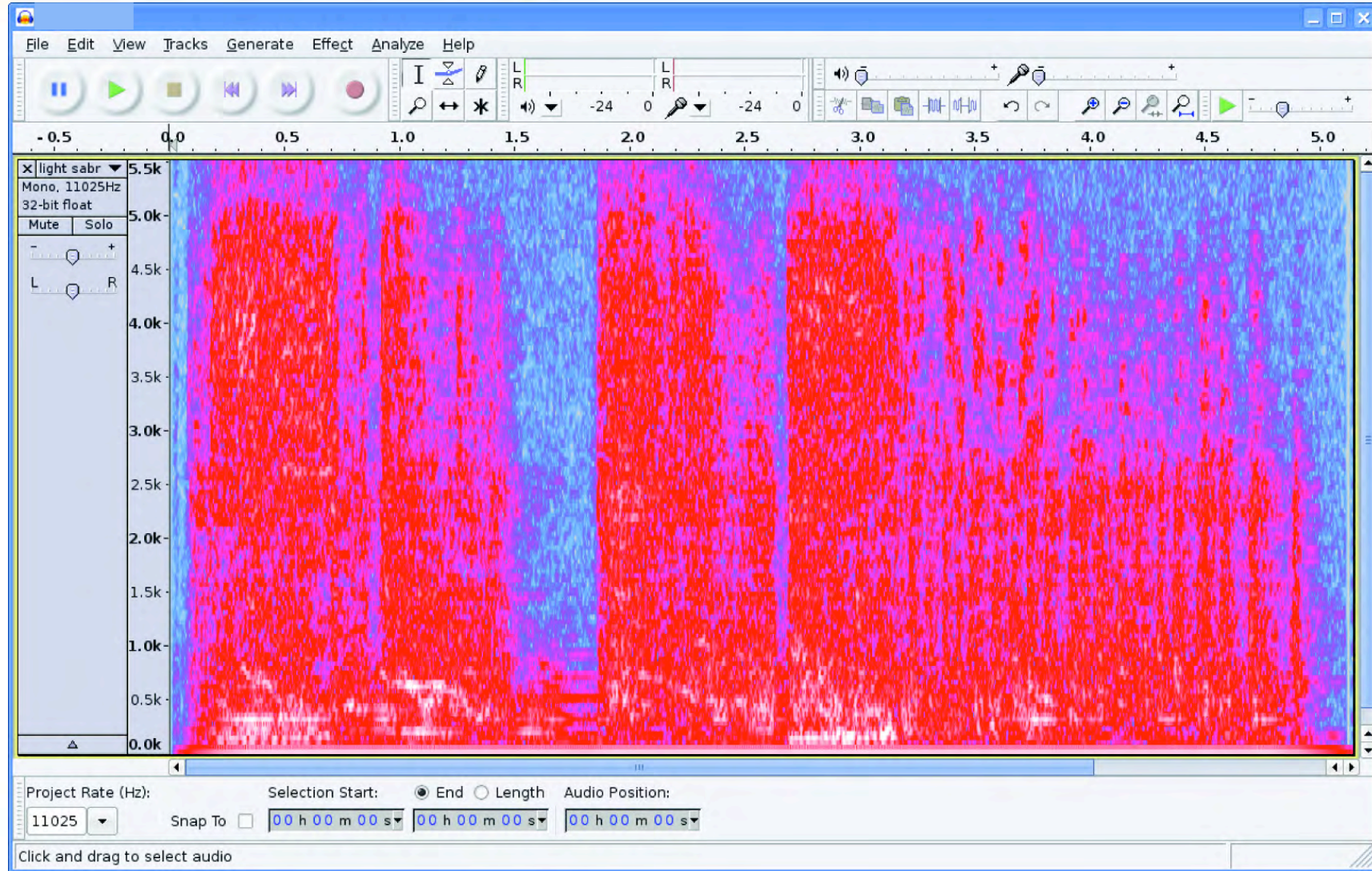


# Hann window application



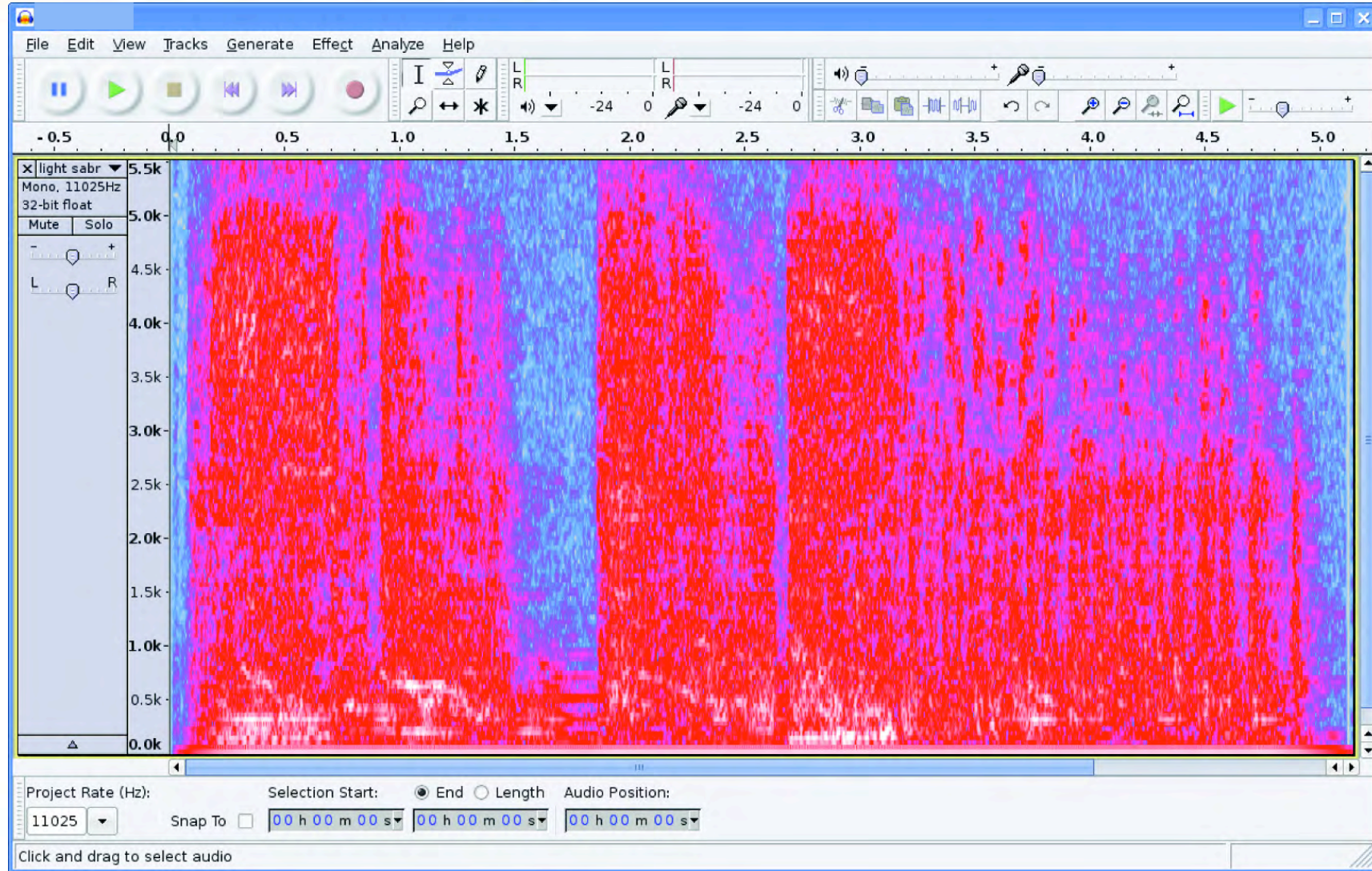


# How Shazam works: audio fingerprinting



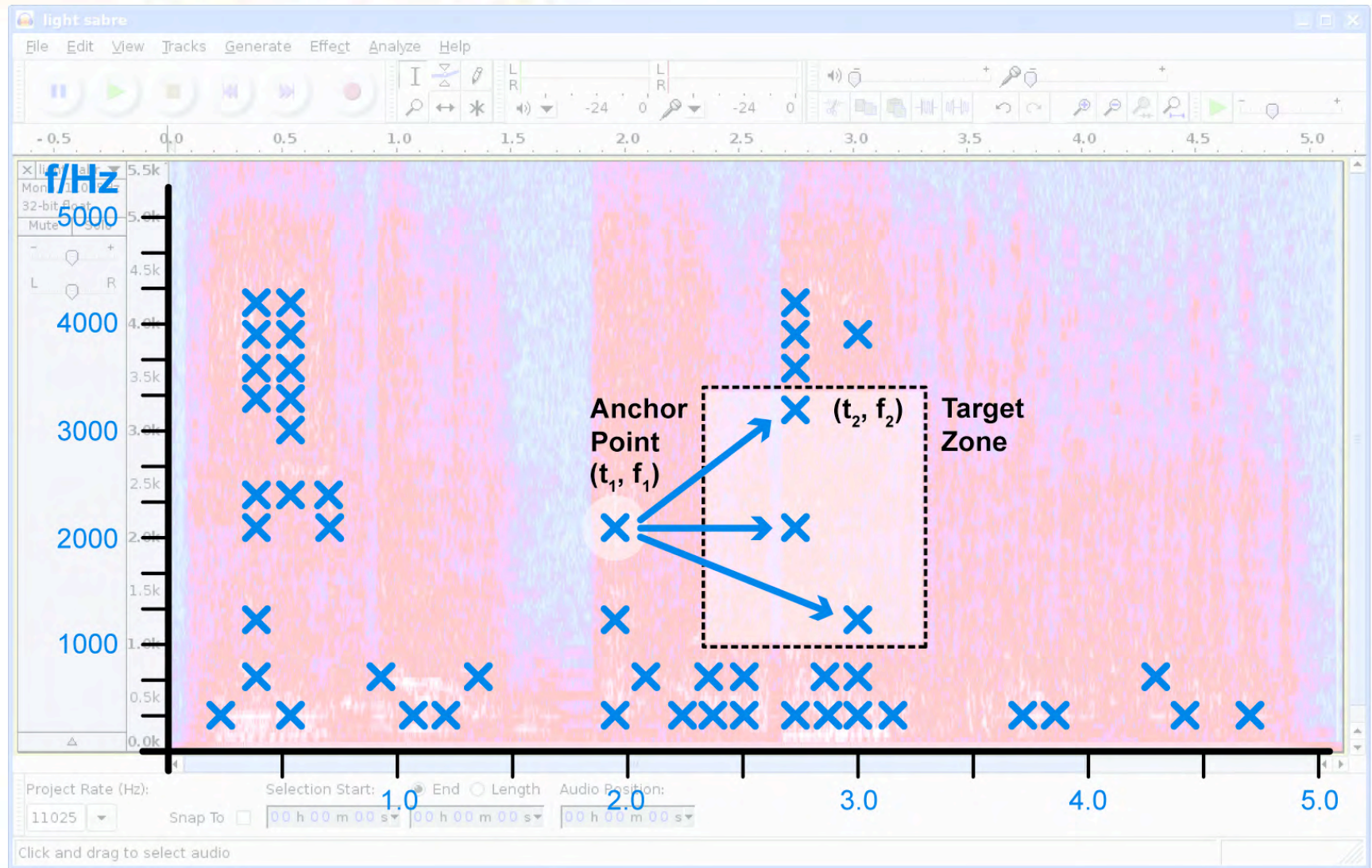


# How Shazam works: audio fingerprinting





# Salient points



Encoding:  $(f_1, f_2, t_2 - t_1)$  hashes to  $(t_1, id)$

[Wang(2003), An industrial-strength search algorithm, ISMIR]



## Temporal consistency check of query

Every query vector  $(f_1, f_2, t_2^q - t_1^q)$  is matched to the database.

You get a list of possible  $(t_1^{id}, id)$  values (some are false positives).

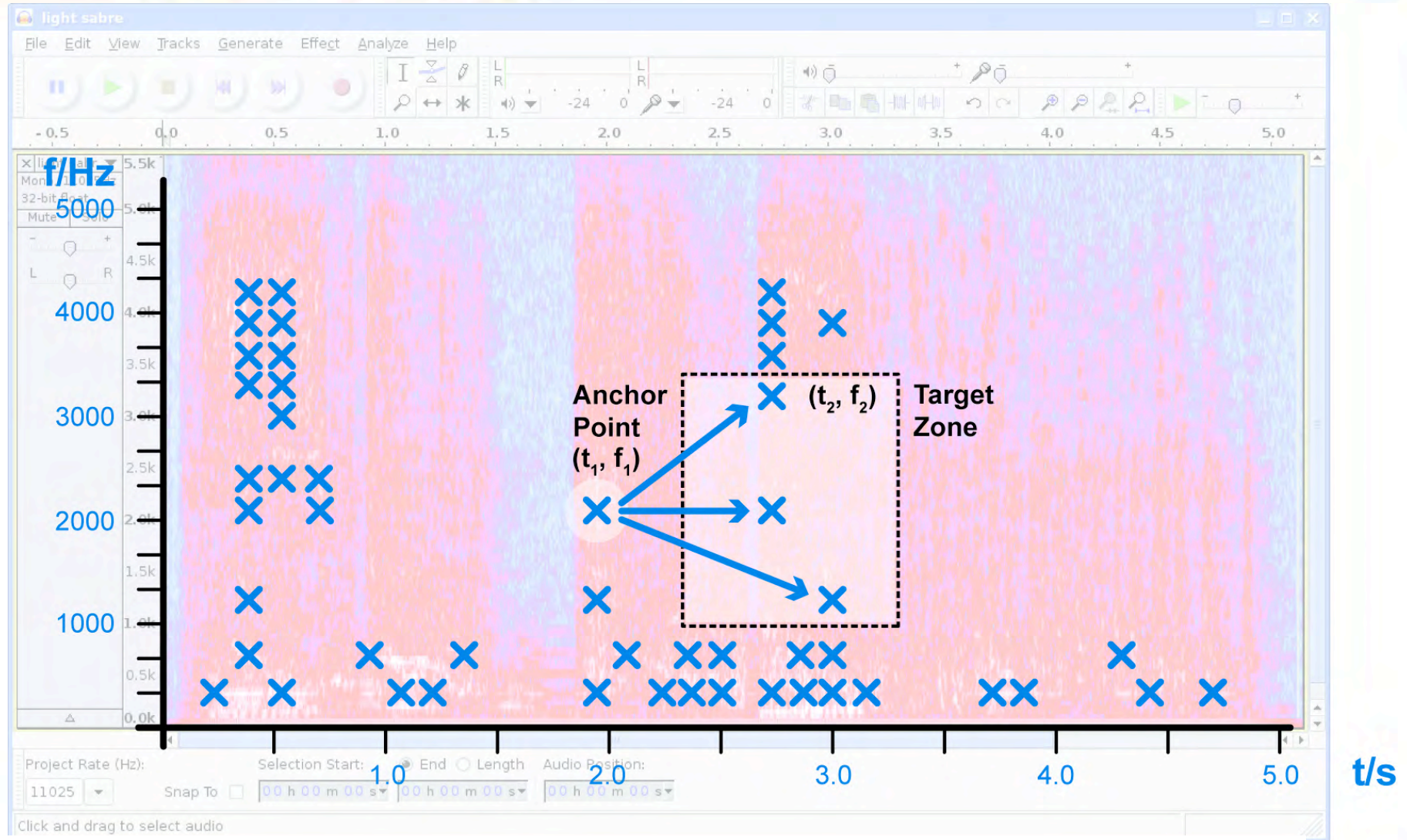
Create a histogram of  $t_1^{id} - t_1^q$  values (temporal consistency check!)

A substantial peak in this histogram means that the query has matched song id at time offset  $t_1^{id} - t_1^q$ .





# Entropy considerations



Specificity: Encoding  $(f_1, f_2, t_2 - t_1)$  to use 30 bit



## Exercise Shazam's constellation pairs

Assume that the typical survival probability of each 30-bit constellation pair after deformations that we still want to recognise is  $p$ , and that this process is independent per pair. Which encoding density, ie, the number of constellation pairs per second, would you need on average so that a typical query of 10 seconds exhibits at least 10 matches in the right song with a probability of at least 99.99%?

Under these assumptions, further assuming that the constellation pair extraction looks like a random independent and identically distributed number, what is the false positive rate for a database of 4 million songs each of which is 5 minutes long on average?



## Exercise Shazam's constellation pairs

Which encoding density would you need on average so that a typical query of 10 seconds exhibits at least 10 matches in the right song with a probability of at least 99.99%?

- approximately 1 match per second needed ( $n = \text{pairs/second}$ ):

$$1 - (1 - p)^n = 0.9999$$

$$n = \log(0.0001) / \log(1 - p)$$



## Exercise Shazam's constellation pairs

Which encoding density would you need on average so that a typical query of 10 seconds exhibits at least 10 matches in the right song with a probability of at least 99.99%?

- Exact solution: binomial distribution

$$n = \min \left\{ m \mid \sum_{i \geq 10} \binom{10m}{i} (1-p)^{10m-i} p^i \geq 0.9999 \right\}$$



## Exercise Shazam's constellation pairs

Which encoding density would you need on average so that a typical query of 10 seconds exhibits at least 10 matches in the right song with a probability of at least 99.99%?

- Large  $n$ : approximate binomial distribution with  $N(np, \sqrt{np(1-p)})$



## Exercise Shazam's constellation pairs

Assuming that the constellation pair extraction looks like a random independent and identically distributed number, what is the false positive rate for a database of 4 million songs each of which is 5 minutes long on average?

Zero:

5min = 30\*10sec

(assume distinctive  $2^{30}$ )

$m = 2^{-30}$

$p(\text{query matches one segment}) \approx m^{10} \approx 2^{-300}$

$1 - (1 - p(\text{qms}))^{(30 \cdot 4e6)} \approx 120e6 \cdot m^{10}$

still near zero



Divide frequency scale into 33 frequency bands between 300 Hz and 2000 Hz  
Logarithmic spread – each frequency step is 1/12 octave, ie, one semitone

Divide time axis into blocks of 256 windows of 11.6 ms (3 seconds)

$E(m,n)$  is the energy of the  $m$ -th frequency at  $n$ -th time in spectrogram

For each block extract 256 sub-fingerprints of 32 bits each

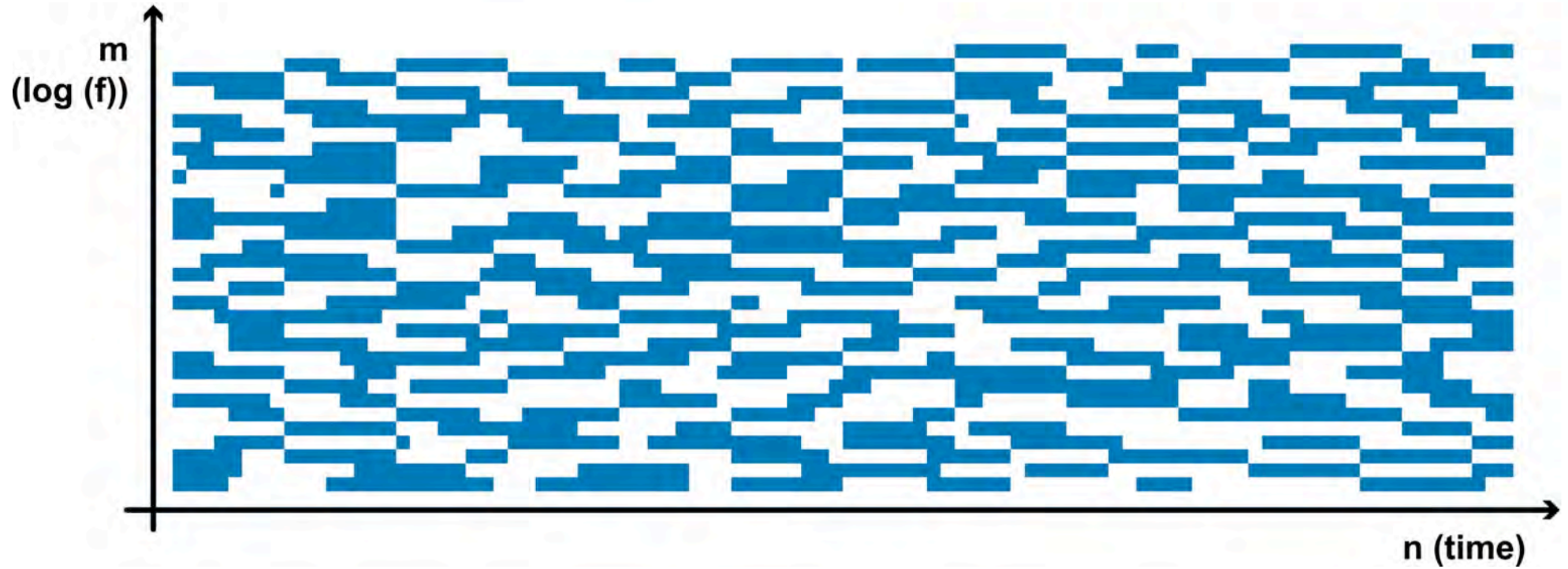
$$b(m, n) = \text{sign} ([E(m, n) - E(m + 1, n)] - [E(m, n + 1) - E(m + 1, n + 1)])$$

$$0 \leq m \leq 31 \text{ (frequency)}$$

$$0 \leq n \leq 255 \text{ (time)}$$



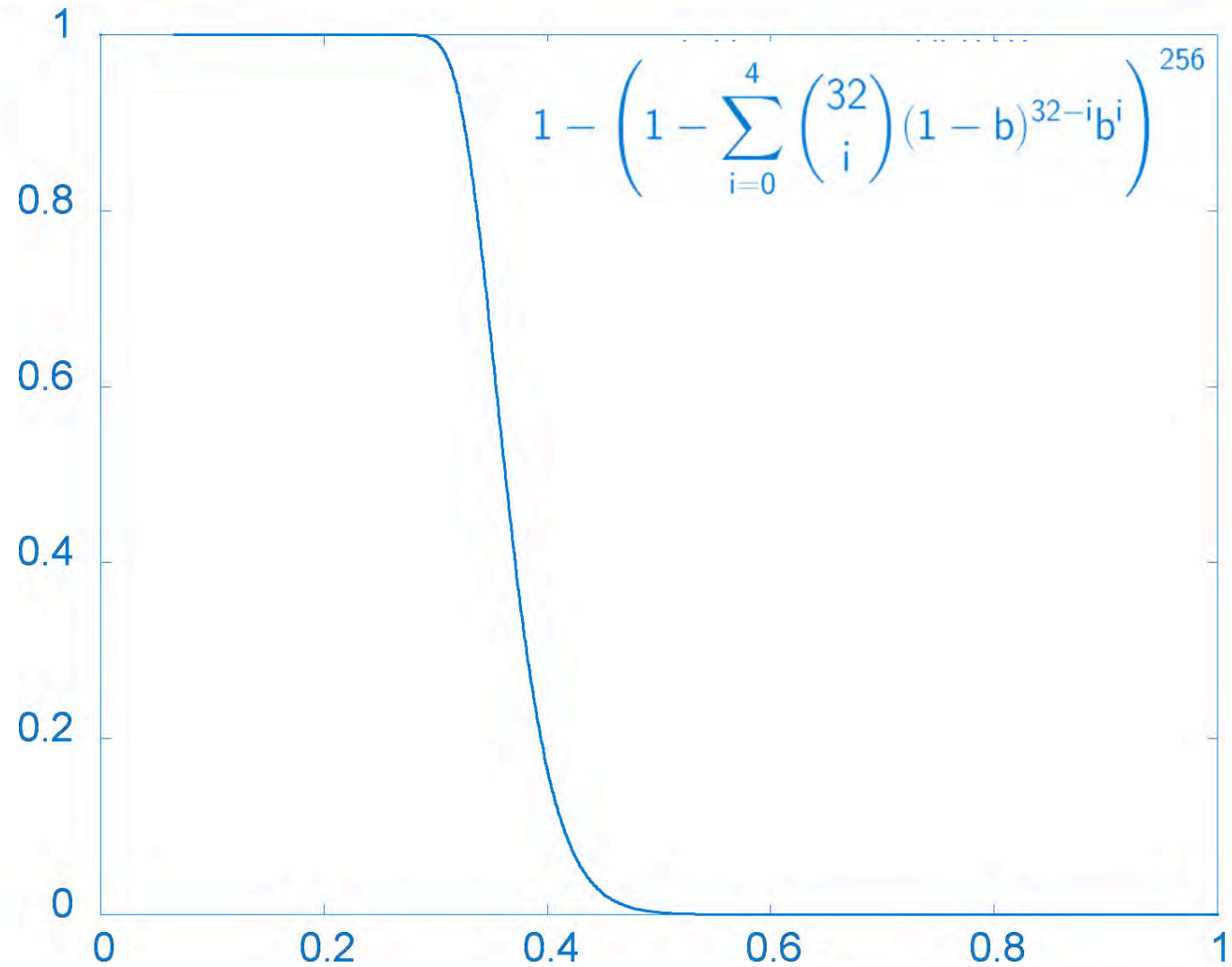
# Partial fingerprint block







# Probability of at least one sub-fingerprint surviving with no more than 4 errors





## Quantisation through locality sensitive hashing (LSH)

$$h^i: \mathbb{R}^d \rightarrow \mathbb{Z}$$

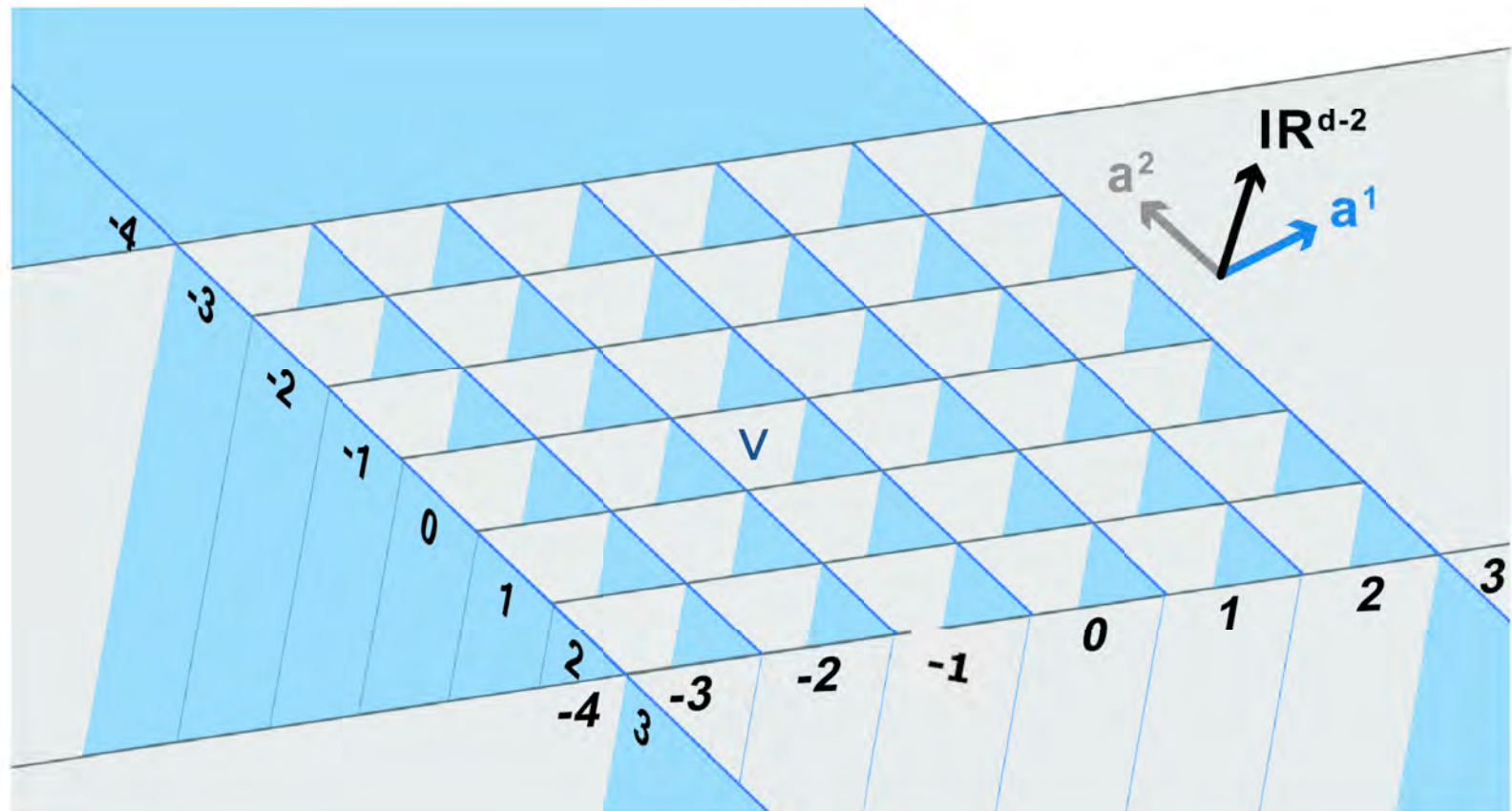
$$\mathbf{v} \mapsto h^i(\mathbf{v}) = \left\lfloor \frac{\mathbf{a}^i \mathbf{v} + b^i}{w} \right\rfloor$$

$\mathbf{a}^i \in \mathbb{R}^d$  is a random Gaussian-distributed vector

$w \in \mathbb{R}^+$  is a constant

$b^i \in [0, w)$  is a random number

$\mathbf{h}(\mathbf{v}) = (h^1(\mathbf{v}), h^2(\mathbf{v}), \dots, h^k(\mathbf{v}))$  is the LSH hash vector.



Vector  $v$ :  $h(v) = (-1, 0)$



Use  $L$  independent hash vectors of  $k$  components each both for the query and for each multimedia object.

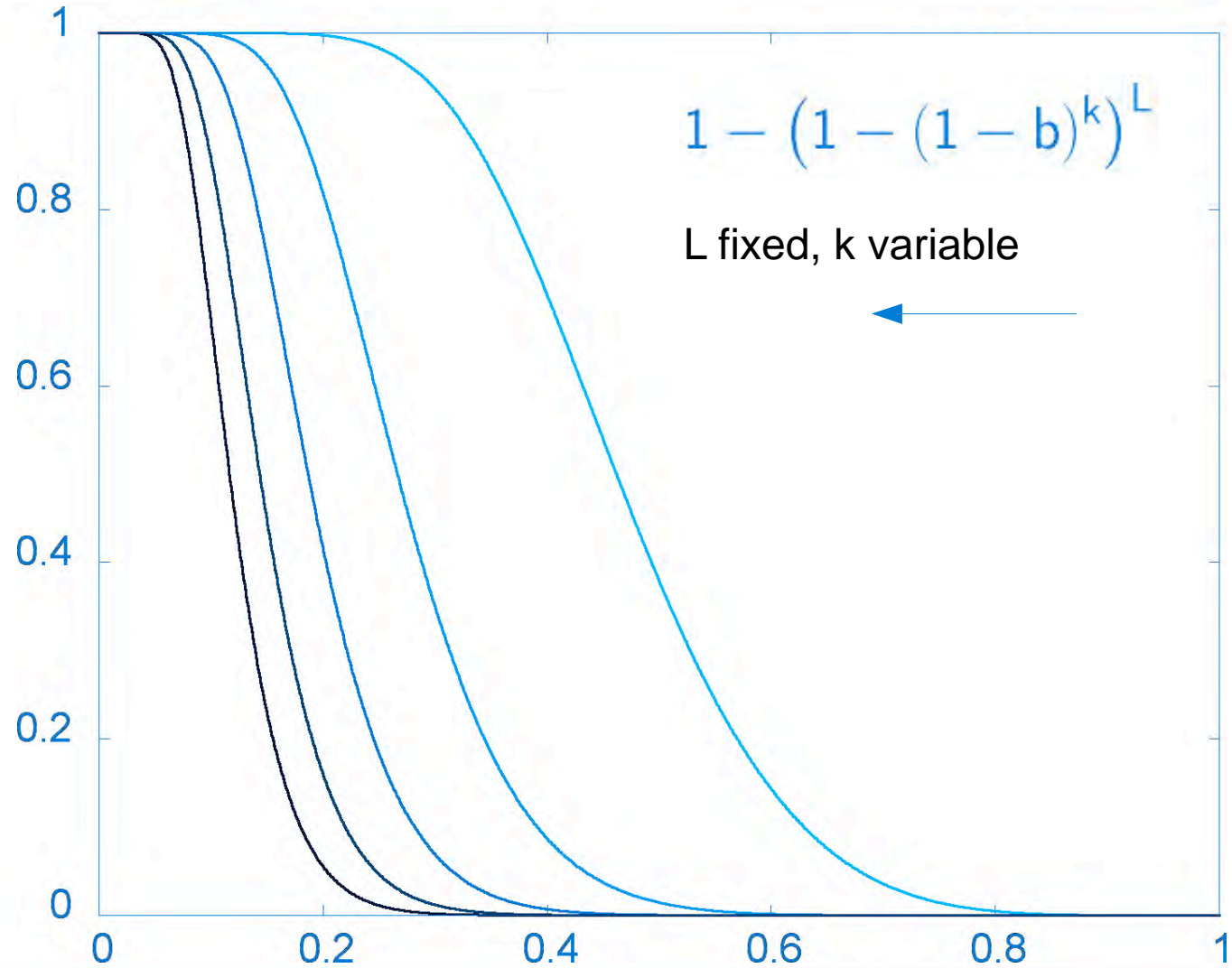
Database elements that match at least  $m$  out of  $L$  times are candidates for nearest neighbours.

Chose  $w$ ,  $k$  and  $L$  (wisely) at runtime

- $w$  determines granularity of bins, ie, # of bits for  $h^i(v)$
- $k$  and  $L$  determine probability of matching

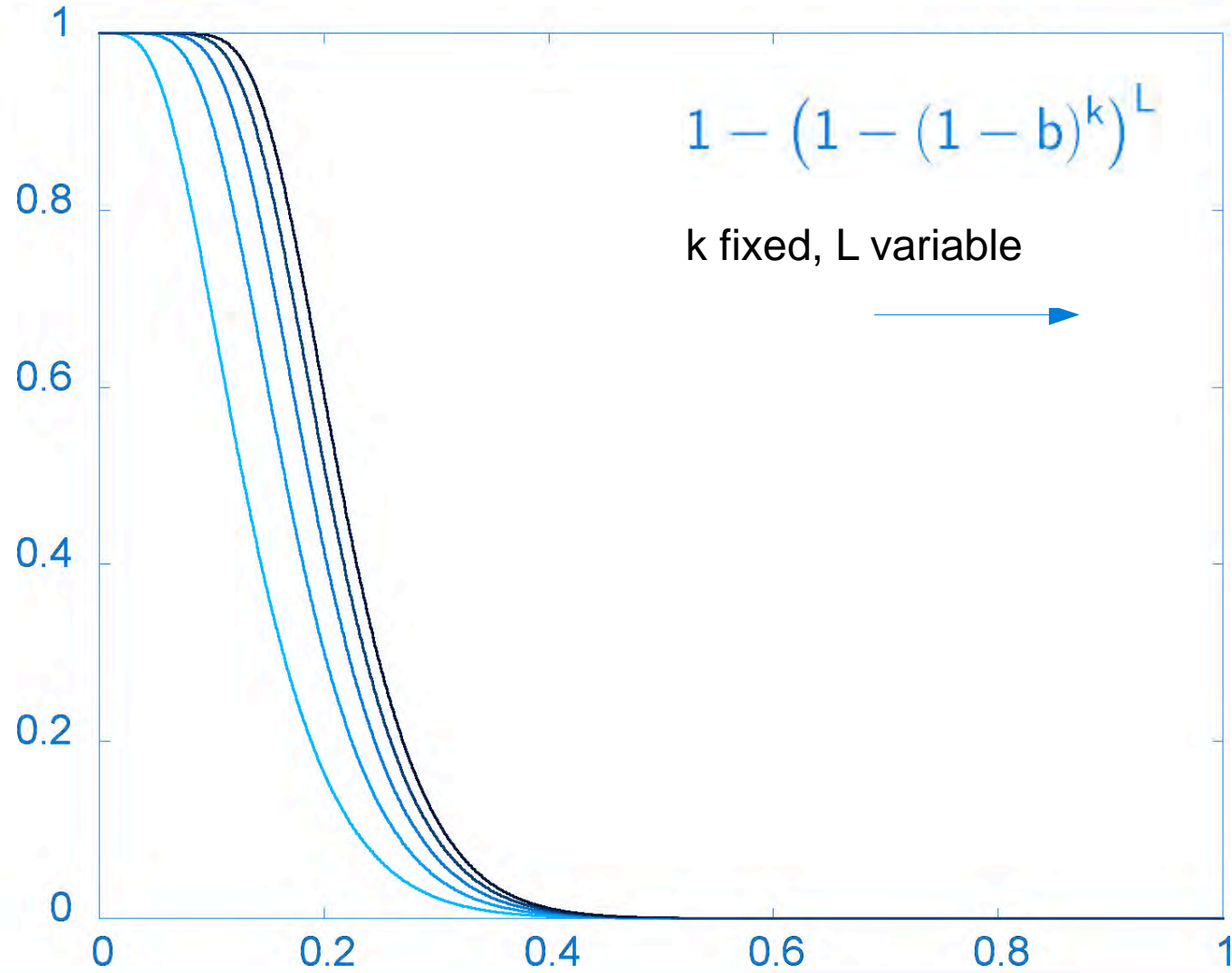


# Prob(min 1 match out of L)



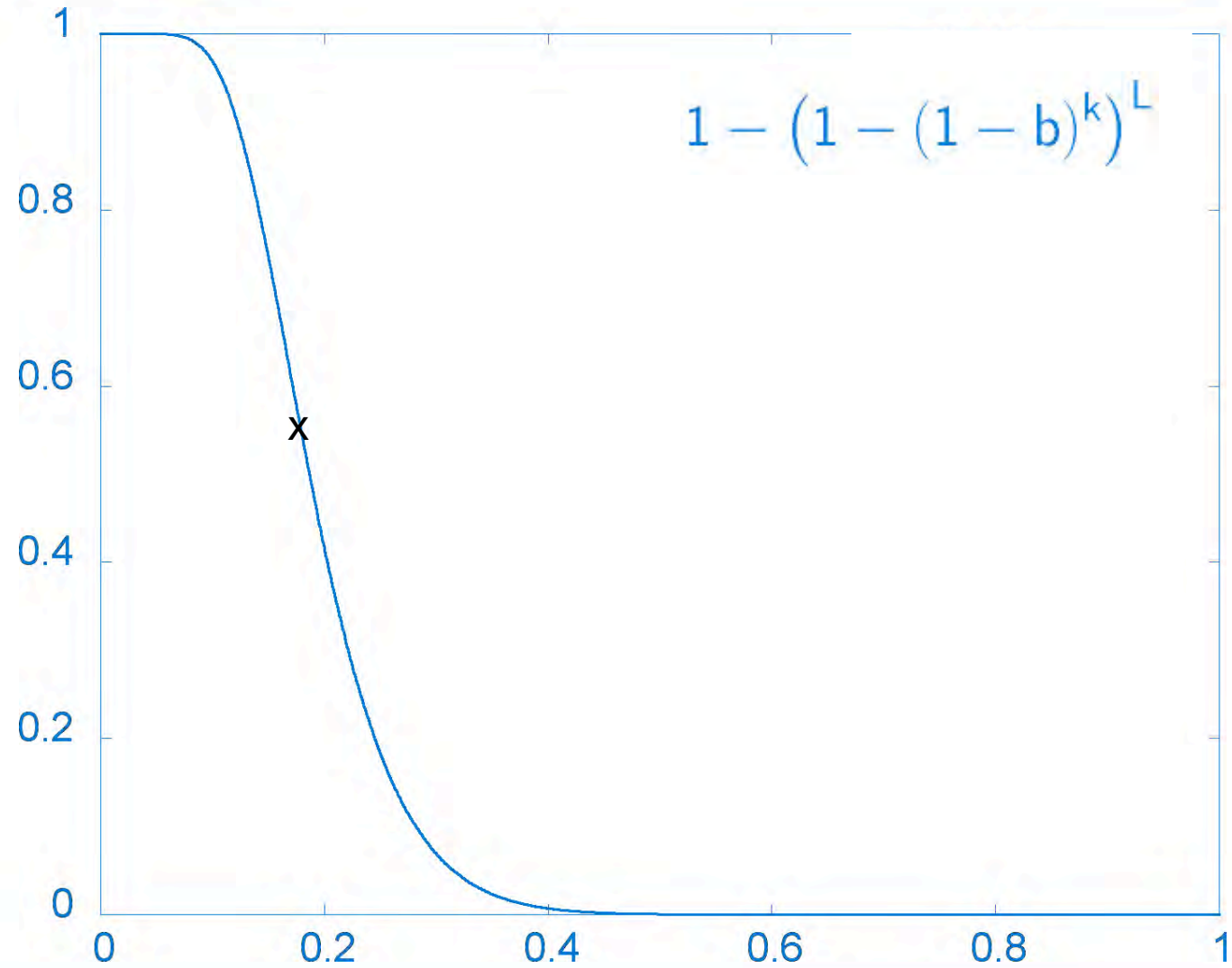


# Prob(min 1 match out of L)





# Exercise: compute inflection point





# Min hash

## Estimate discrete set overlap

$$\text{sim}(A_i, A_j) = \frac{|A_i \cap A_j|}{|A_i \cup A_j|}$$





## An example 4 documents

$D_1$  = Humpty Dumpty sat on a wall,

$D_2$  = Humpty Dumpty had a great fall.

$D_3$  = All the King's horses, And all the King's men

$D_4$  = Couldn't put Humpty together again!



## Surrogate docs after stop word removal and stemming

$A_1 = \{\text{humpty, dumpty, sat, wall}\}$

$A_2 = \{\text{humpty, dumpty, great, fall}\}$

$A_3 = \{\text{all, king, horse, men}\}$

$A_4 = \{\text{put, humpty, together, again}\}$



# Equivalent term-document matrix

	$A_1$	$A_2$	$A_3$	$A_4$
humpty	1	1	0	1
dumpty	1	1	0	0
sat	1	0	0	0
wall	1	0	0	0
great	0	1	0	0
fall	0	1	0	0
all	0	0	1	0
king	0	0	1	0
horse	0	0	1	0
men	0	0	1	0
put	0	0	0	1
together	0	0	0	1
again	0	0	0	1



## Similarity between two docs

$$\text{sim}(A_i, A_j) = \frac{c_{11}}{c_{11} + c_{10} + c_{01}}$$

$c_{xy}$  = number of (x,y) rows

Important observation  
 $c_{00}$  is unused!

	$A_1$	$A_2$	$A_3$	$A_4$
humpty	1	1	0	1
dumpty	1	1	0	0
sat	1	0	0	0
wall	1	0	0	0
great	0	1	0	0
fall	0	1	0	0
all	0	0	1	0
king	0	0	1	0
horse	0	0	1	0
men	0	0	1	0
put	0	0	0	1
together	0	0	0	1
again	0	0	0	1



## Estimation of similarity through random permutations

$\pi_1 = (\text{dumpty, men, again, put, great, humpty, wall, horse, king, sat, fall, together, all})$

$\pi_2 = (\text{fall, put, all, again, dumpty, sat, men, great, wall, king, horse, humpty, together})$

$\pi_3 = (\text{horse, dumpty, wall, humpty, great, again, sat, all, men, together, put, king, fall})$

$\pi_4 = (\text{king, humpty, men, together, great, fall, horse, all, dumpty, wall, sat, again, put})$



# Surrogate documents form random permutations

Keep first occurring word of  $A_i$  in  $\pi_j$  for dense surrogate representation

	$A_1$	$A_2$	$A_3$	$A_4$
$\pi_1$	dumpty	dumpty	men	again
$\pi_2$	dumpty	fall	all	put
$\pi_3$	dumpty	dumpty	horse	humpty
$\pi_4$	humpty	humpty	king	humpty



# Surrogate documents form random permutations

	$A_1$	$A_2$	$A_3$	$A_4$
$\pi_1$	dumpty	dumpty	men	again
$\pi_2$	dumpty	fall	all	put
$\pi_3$	dumpty	dumpty	horse	humpty
$\pi_4$	humpty	humpty	king	humpty

Estimate  $\text{sim}(A_2, A_4) = 1/4$

(proportion of co-occurring words)



## Scale Invariant Feature Transform

“distinctive invariant image features that can be used to perform reliable matching between different views of an object or scene.”

Invariant to image scale and rotation.

Robust to substantial range of affine distortion, changes in 3D viewpoint, addition of noise and change in illumination.

[Lowe, D.G. (2004). Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision*, 60, 2, pp. 91-110.]





For a given image:

Detect scale space extrema

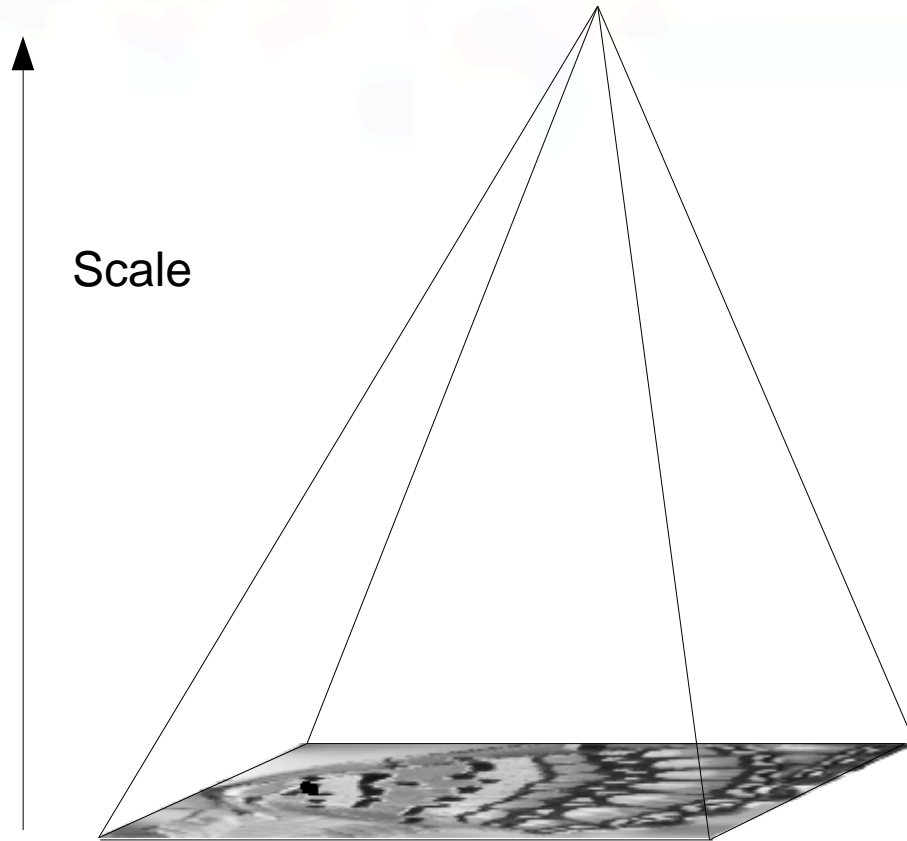
Localise candidate keypoints

Assign an orientation to each keypoint

Produce keypoint descriptor

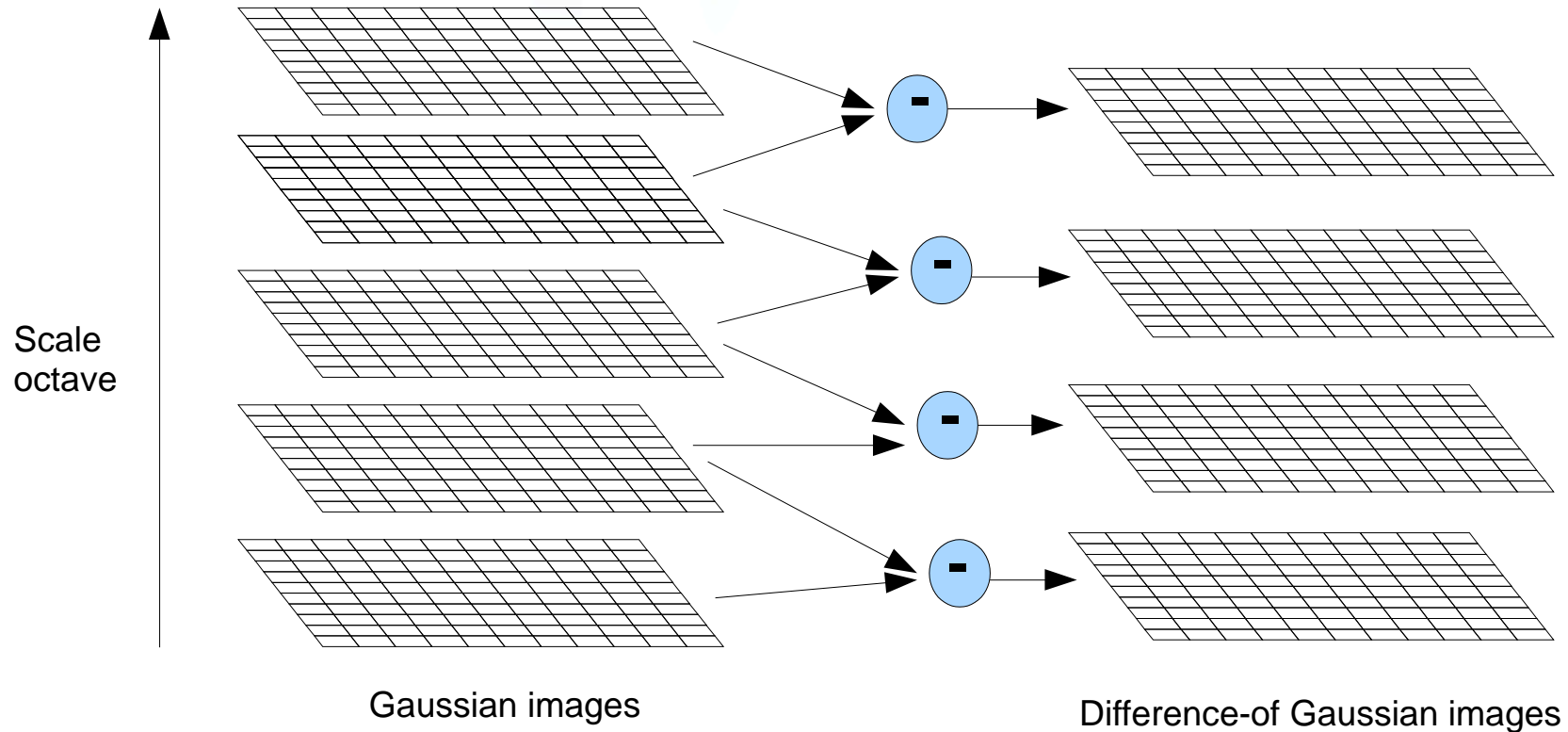


# A scale space visualisation





# Difference of Gaussian image creation





# Gaussian blur illustration





# Difference of Gaussian illustration



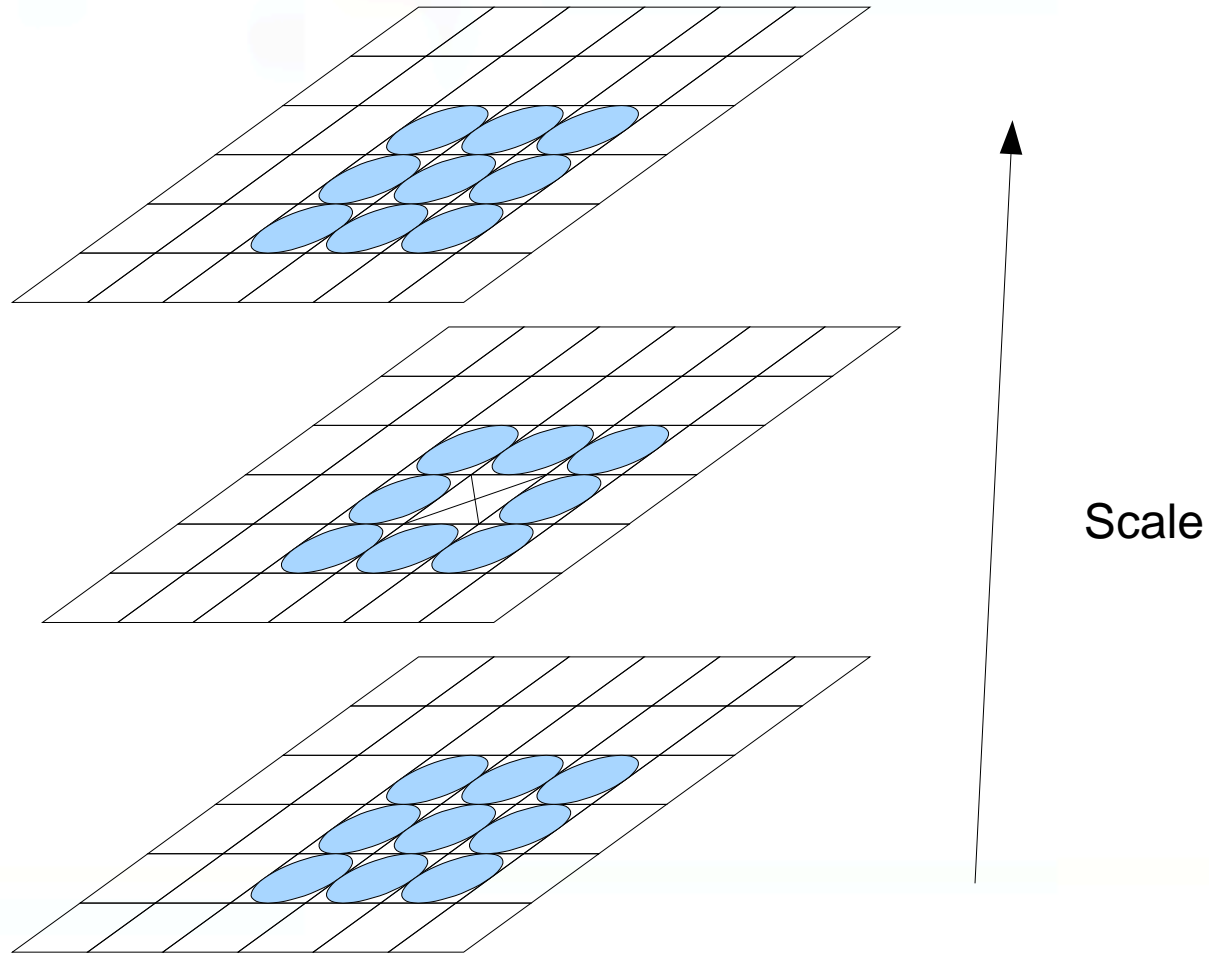


Once the Difference of Gaussian images have been generated:

- Each pixel in the images is compared to 8 neighbours at same scale.
- Also compared to 9 corresponding neighbours in scale above and 9 corresponding neighbours in the scale below.
- Each pixel is compared to 26 neighbouring pixels in 3x3 regions across scales, as it is not compared to itself at the current scale.
- A pixel is selected as a SIFT keypoint only either if its intensity value is extreme.

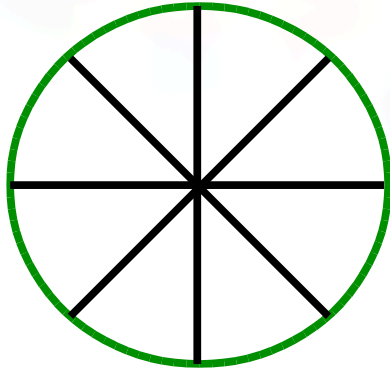


# Pixel neighbourhood comparison



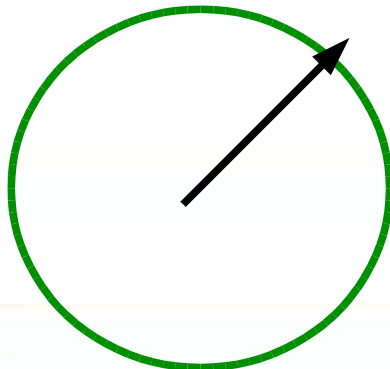


# Orientation assignment

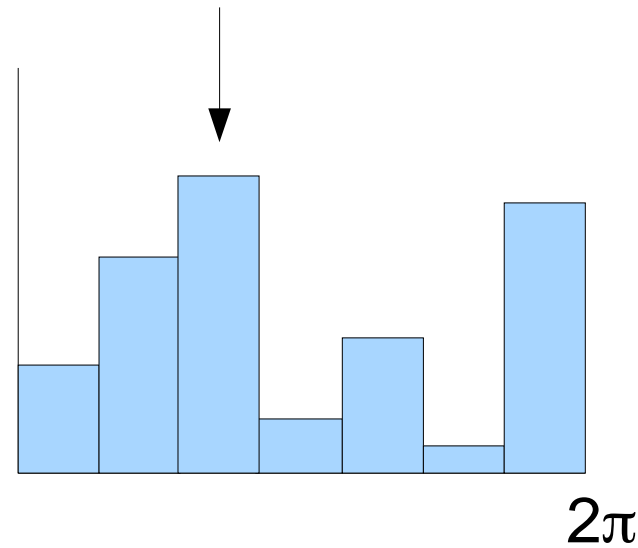


Each sample weighted by gradient magnitude and Gaussian window.

Canonical orientation at peak of Smoothed histogram.



Orientation histogram with 36 bins – one per 10 degrees.



Where two or more orientations are detected, keypoints created for each orientation.





We now have location, scale and orientation for each SIFT keypoint (“keypoint frame”).

→ descriptor for local image region is required.

Must be as invariant as possible to changes in illumination and 3D viewpoint.

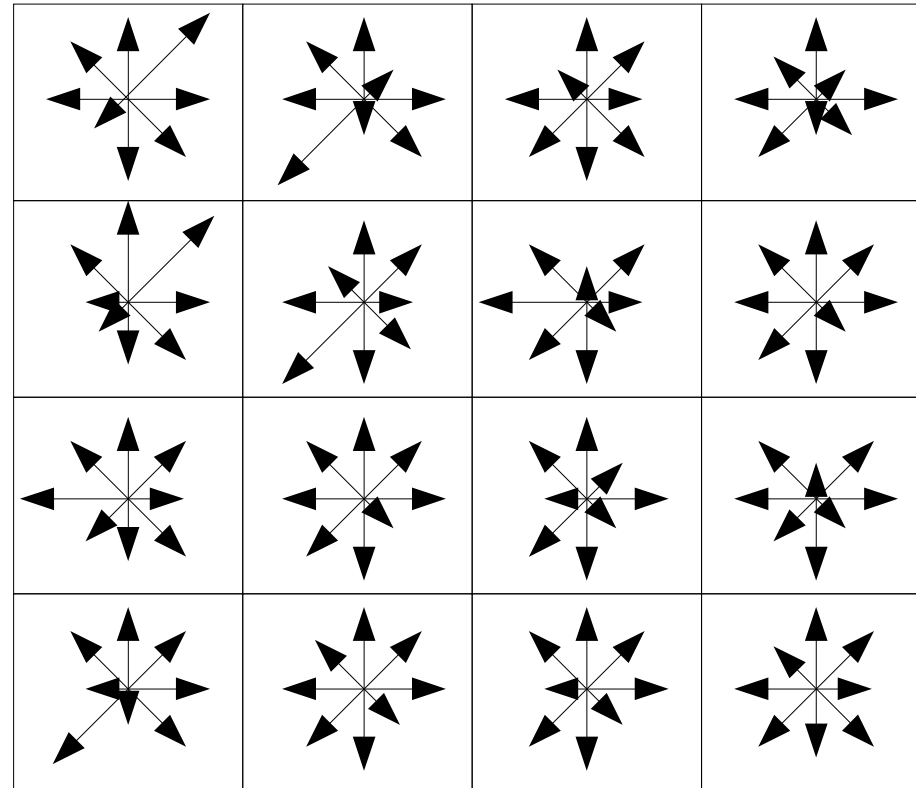
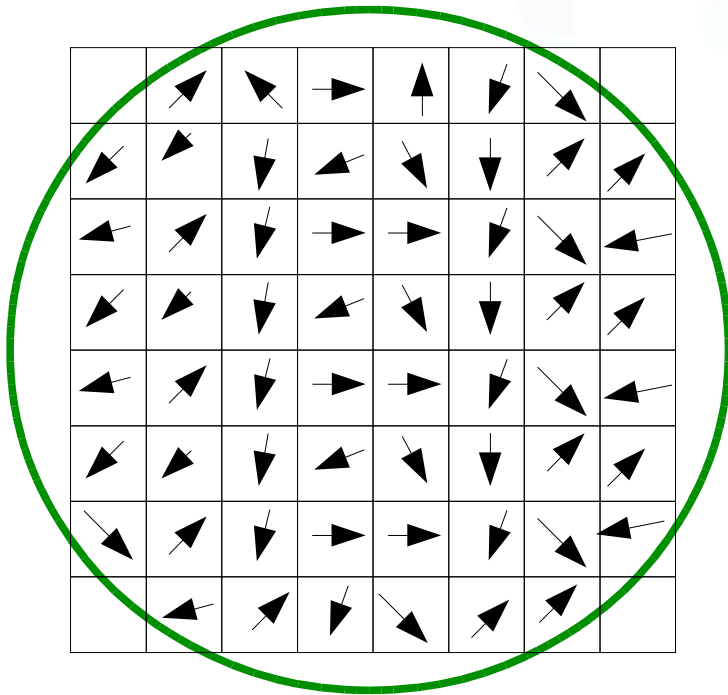
Set of orientation histograms are computed on 4x4 pixel areas.

Each gradient histogram contains 8 bins and each descriptor contains an array of 4 histograms.

→ SIFT descriptor as 128 (4 x 4 x 8) element histogram



# Visualising the keypoint descriptor





# Example SIFT keypoints



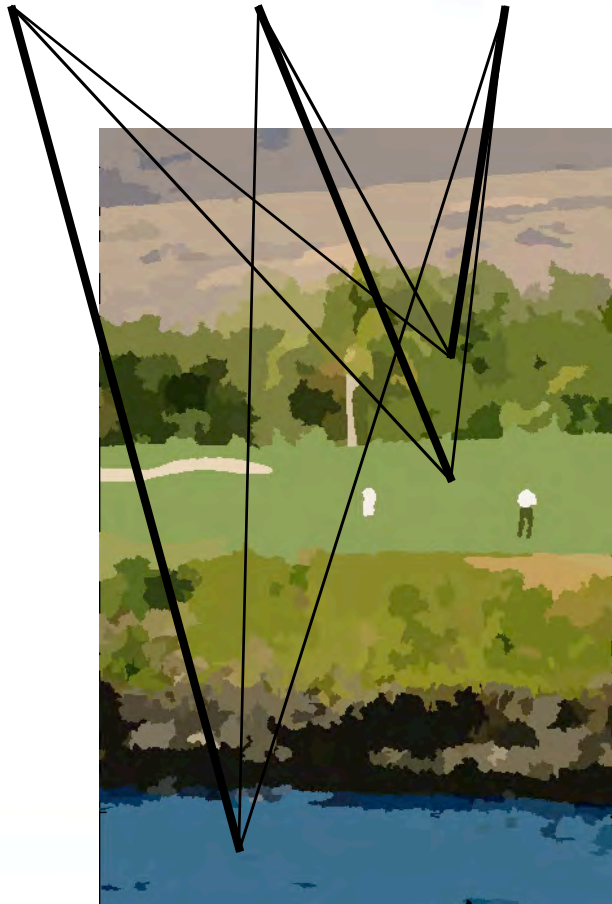


1. What is multimedia information retrieval?
2. Metadata and piggyback retrieval
3. Multimedia fingerprinting
4. **Automated annotation**
5. Content-based retrieval

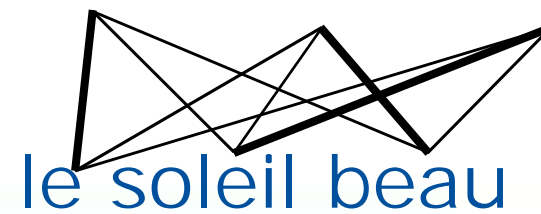


# Automated annotation as machine translation

water grass trees



the beautiful sun





## Probabilistic models:

maximum entropy models

models for joint and conditional probabilities

evidence combination with Support Vector Machines

[with Magalhães, SIGIR 2005]

[with Yavlinsky and Schofield, CIVR 2005]

[with Yavlinsky, Heesch and Pickering: ICASSP May 2004]

[with Yavlinsky et al CIVR 2005]

[with Yavlinsky SPIE 2007]

[with Magalhães CIVR 2007, *best paper*]



$$\begin{aligned}
 P(\mathbf{w}|\mathbf{I}) &= \frac{P(\mathbf{w}, \mathbf{I})}{P(\mathbf{I})} = \frac{\sum_J P(\mathbf{w}, \mathbf{I}|\mathbf{J})P(\mathbf{J})}{\sum_J P(\mathbf{I}|\mathbf{J})P(\mathbf{J})} \\
 &= \frac{\sum_J P(\mathbf{I}|\mathbf{w}, \mathbf{J})P(\mathbf{w}|\mathbf{J})P(\mathbf{J})}{\sum_J \sum_{\mathbf{w}} P(\mathbf{I}|\mathbf{w}, \mathbf{J})P(\mathbf{w}|\mathbf{J})P(\mathbf{J})}
 \end{aligned}$$

Use training data  $\mathbf{J}$  and annotations  $\mathbf{w}$

$P(\mathbf{w}|\mathbf{I})$  is probability of word  $\mathbf{w}$  given unseen image  $\mathbf{I}$

The model is an empirical distribution  $(\mathbf{w}, \mathbf{J})$



# Automated annotation



[with Yavlinsky et al CIVR 2005]  
[with Yavlinsky SPIE 2007]  
[with Magalhaes CIVR 2007, best paper]

Automated: water buildings city sunset aerial

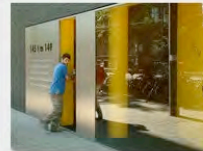
[Corel Gallery 380,000]





# The good

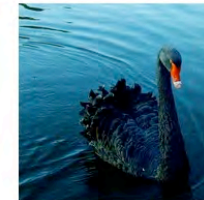
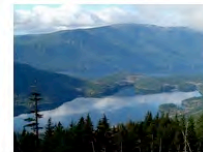
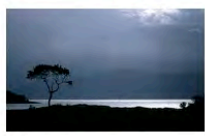
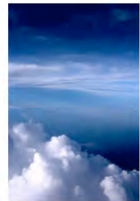
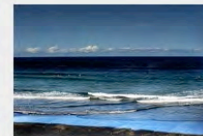
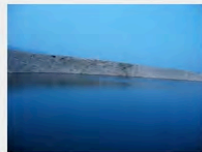
door





# The bad

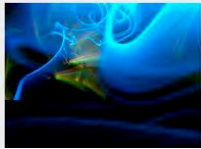
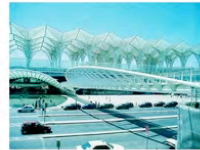
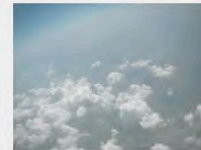
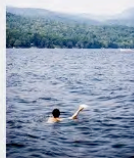
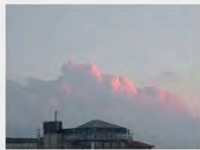
wave





# The ugly

iceberg





# Multimedia information retrieval

1. What is multimedia information retrieval?
2. Metadata and piggyback retrieval
3. Multimedia fingerprinting
4. Automated annotation
5. Content-based retrieval



Give examples where we remember details by

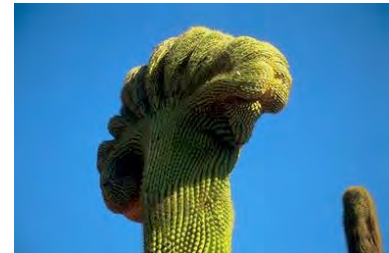
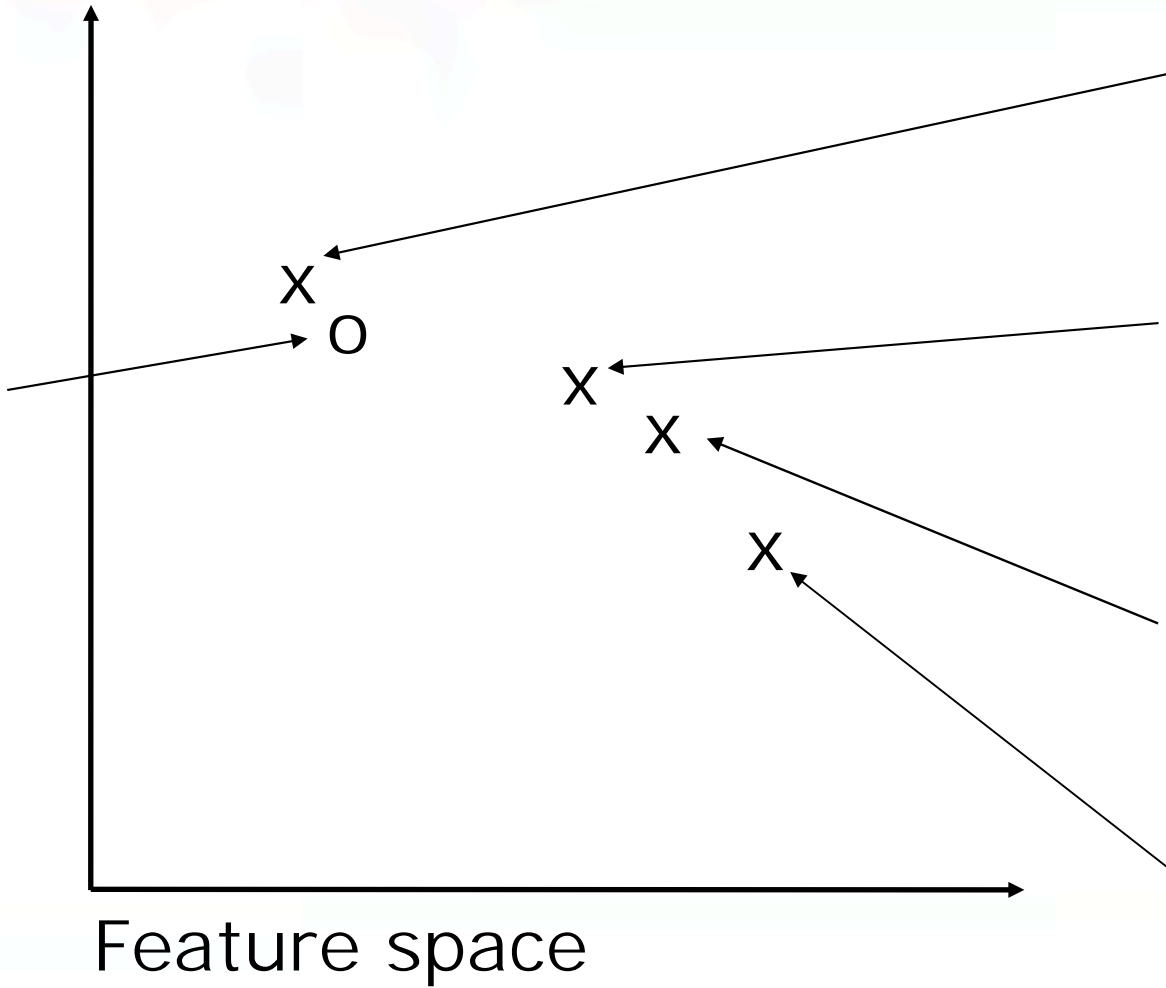
- metadata?
- context?
- content (eg, "x" belongs to "y")?

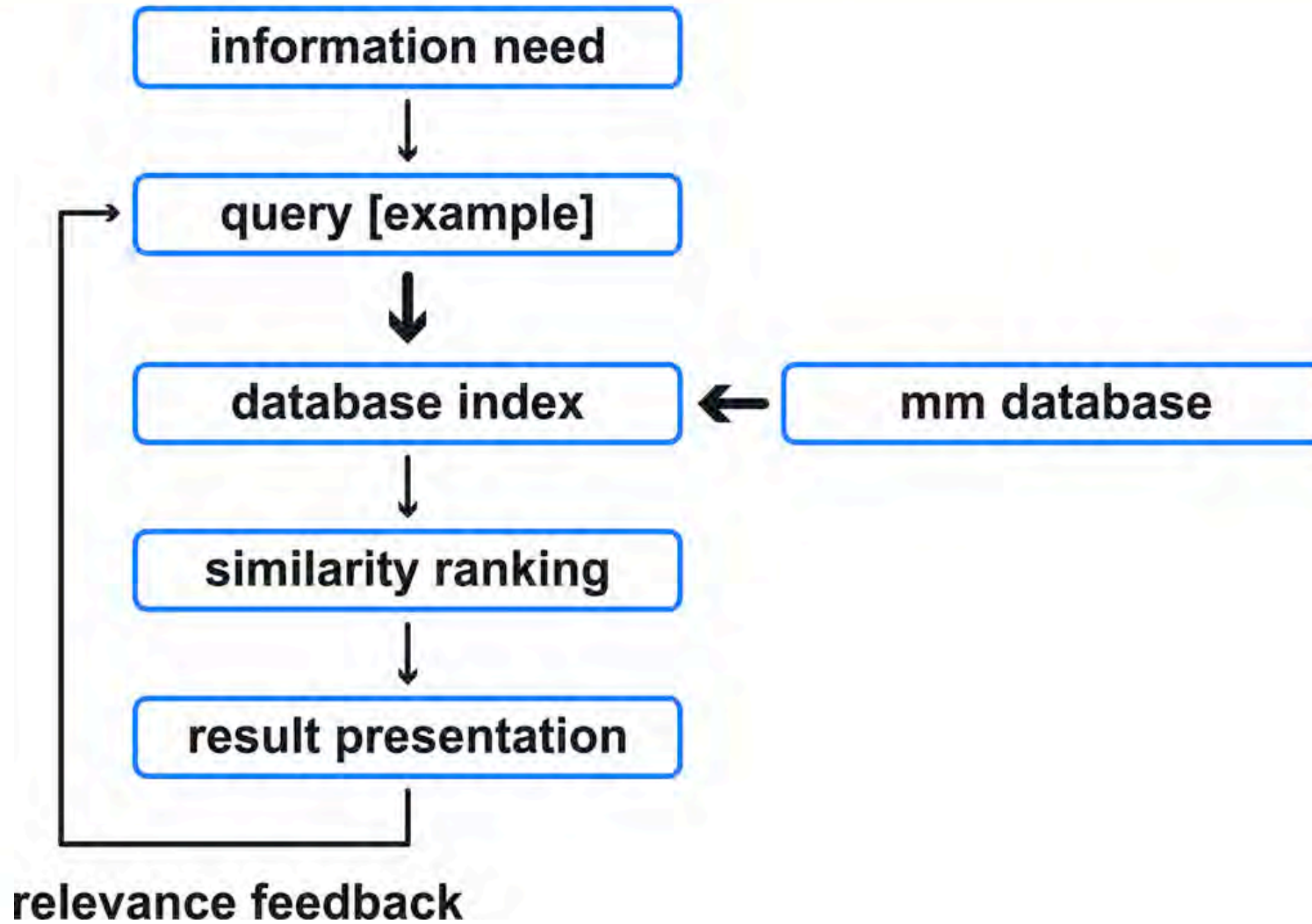
Metadata versus content-based: pro and con

- 
- 
- 
-



# Content-based retrieval: features and distances







Visual

Colour, texture, shape, edge detection, SIFT/SURF

Audio

Temporal

How to describe the features?

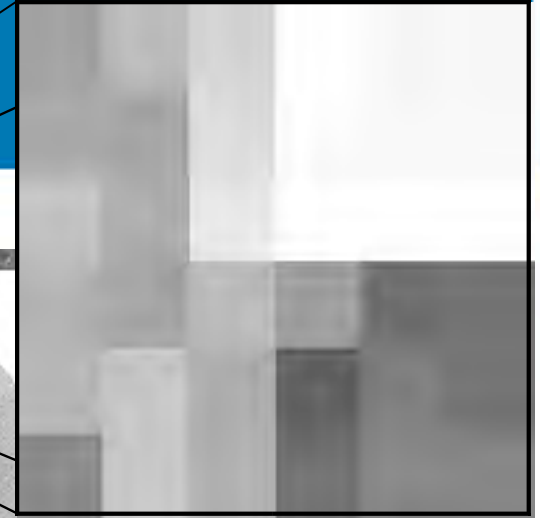
For people

For computers





# Digital Images





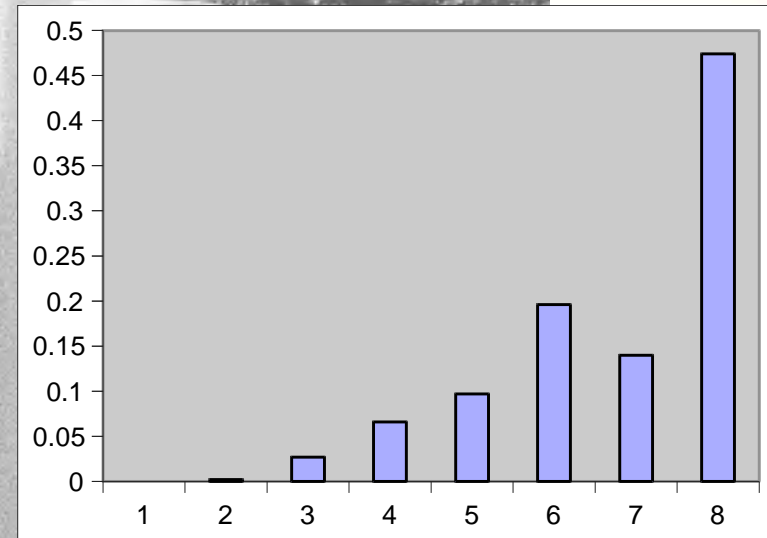
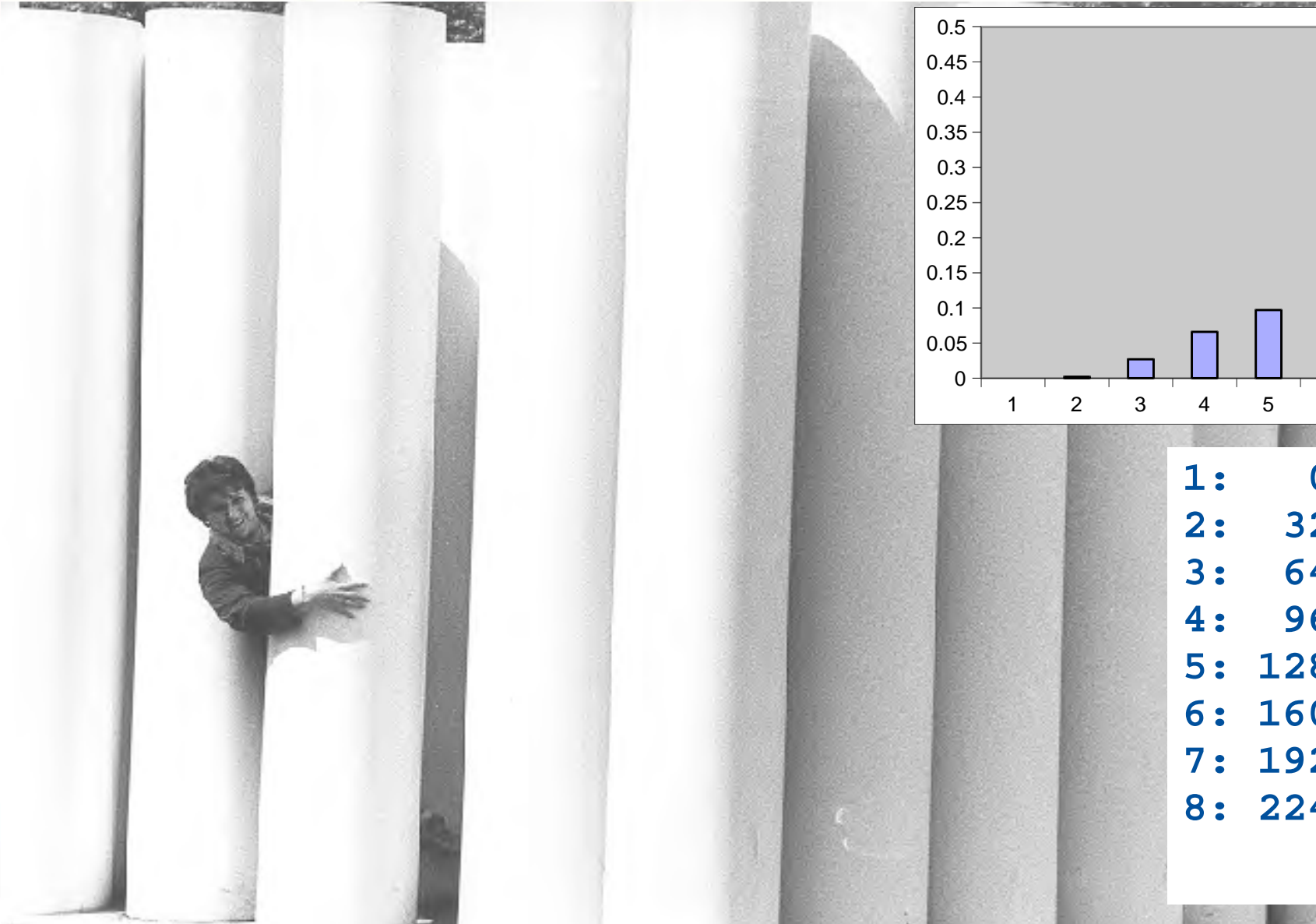
## Content of an image



145	173	201	253	245	245
153	151	213	251	247	247
181	159	225	255	255	255
165	149	173	141	93	97
167	185	157	79	109	97
121	187	161	97	117	115



# Histogram

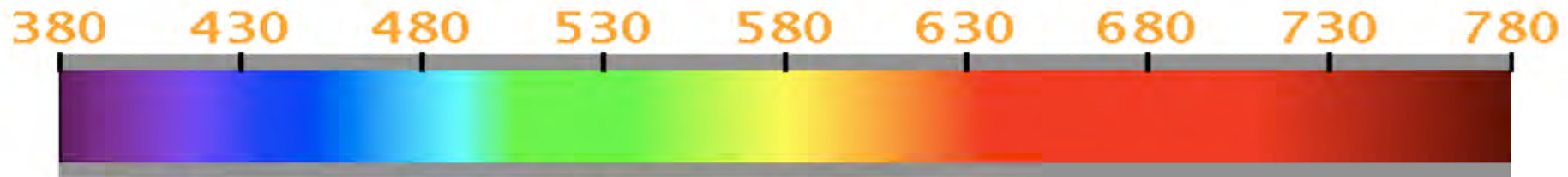


1:	0	-	31
2:	32	-	63
3:	64	-	95
4:	96	-	127
5:	128	-	159
6:	160	-	191
7:	192	-	223
8:	224	-	255



phenomenon of human perception  
three-dimensional (RGB/CMY/HSB)  
spectral colour: pure light of one wavelength

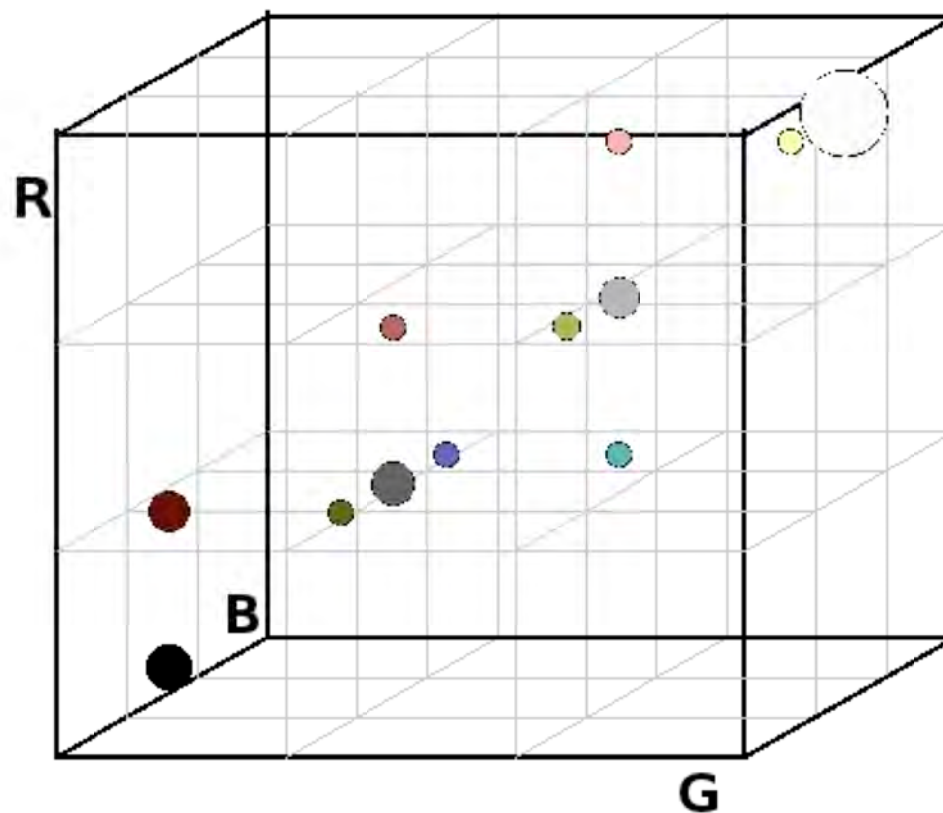
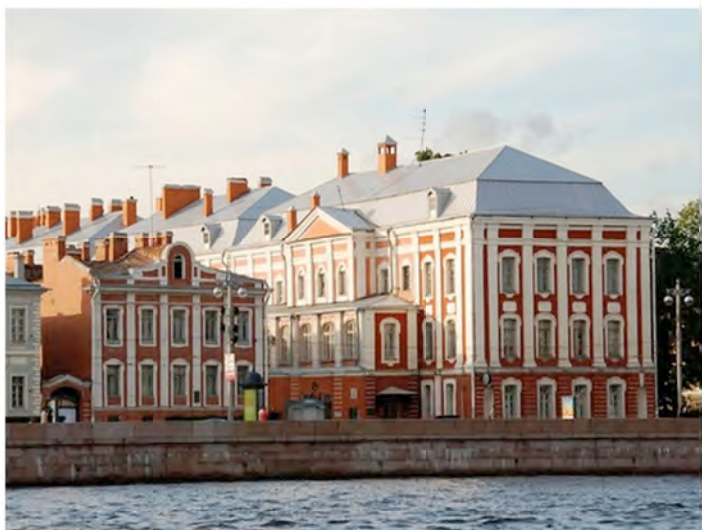
blue cyan green yellow red



spectral colours: wavelength (nm)












# Colour histogram

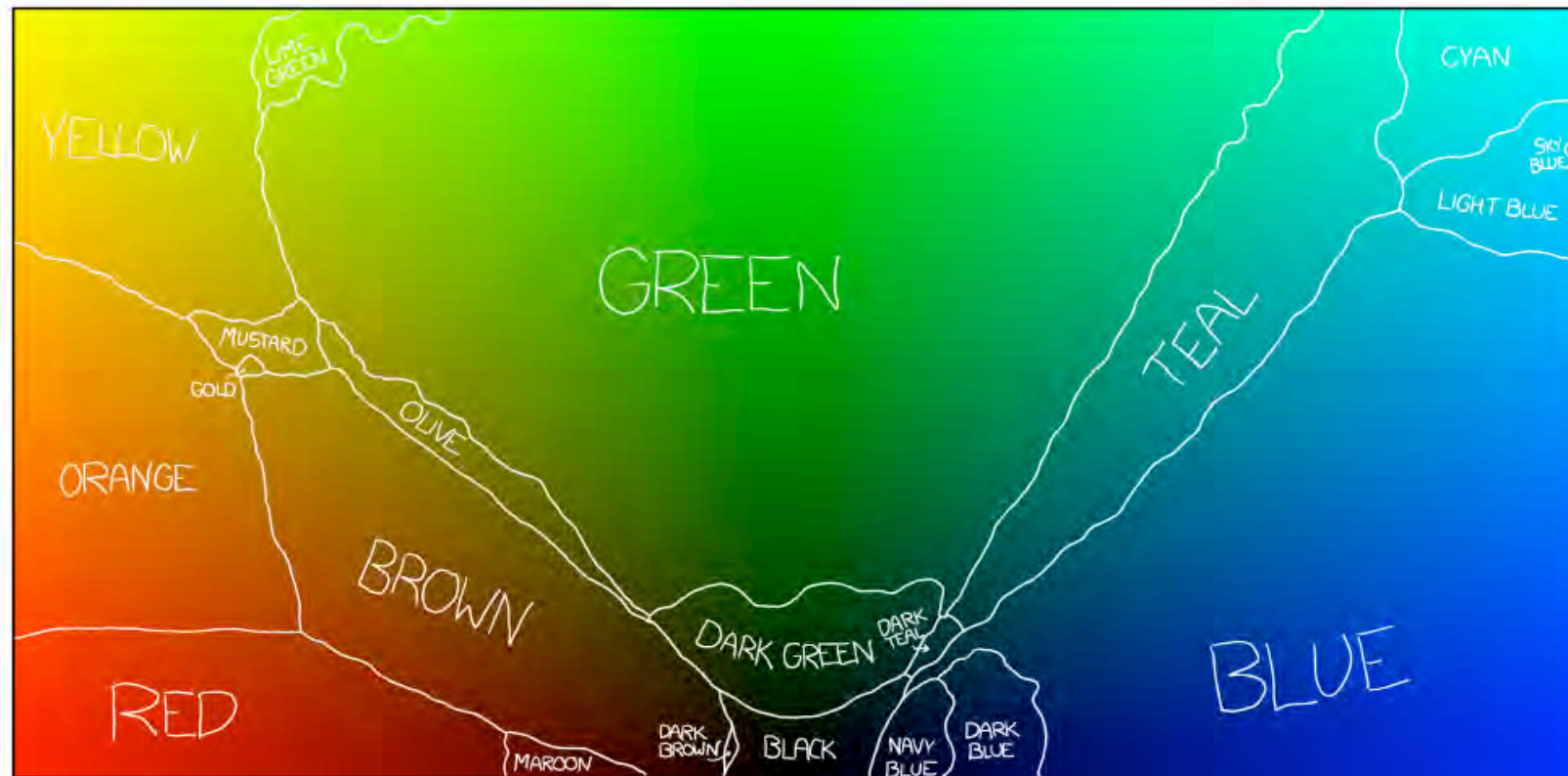




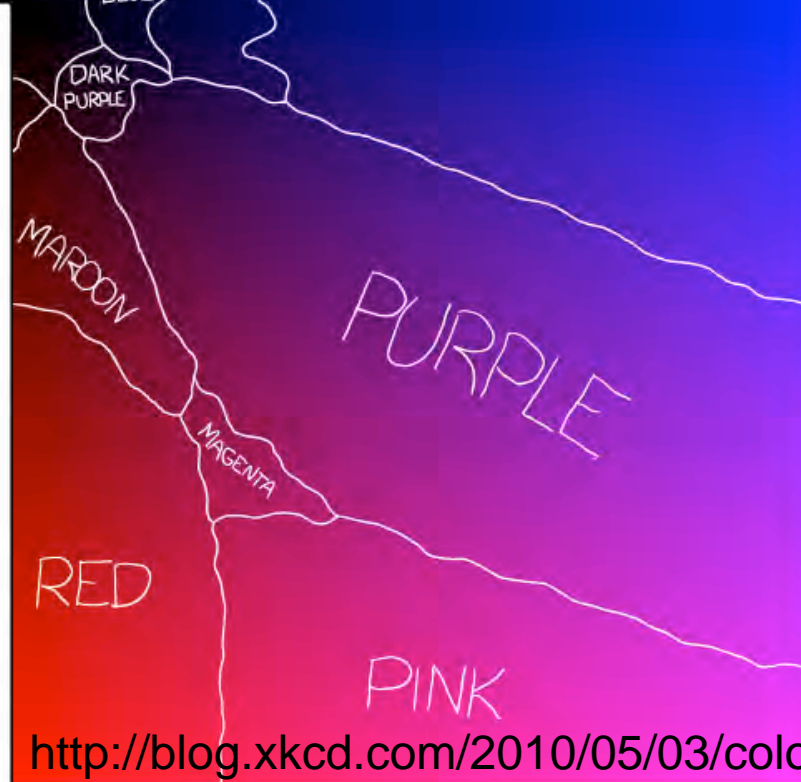
Sketch a 3D colour histogram for



R	G	B		
0	0	0		black
255	0	0		red
0	255	0		green
0	0	255		blue
0	255	255		cyan
255	0	255		magenta
255	255	0		yellow
255	255	255		white
				



THIS CHART SHOWS THE DOMINANT COLOR NAMES OVER THE THREE FULLY-SATURATED FACES OF THE RGB CUBE (COLORS WHERE ONE OF THE RGB VALUES IS ZERO)





hue ( $0^{\circ}$ - $360^{\circ}$ )  
spectral colour



saturation ( $0\%$  -  $100\%$ )  
= spectral purity



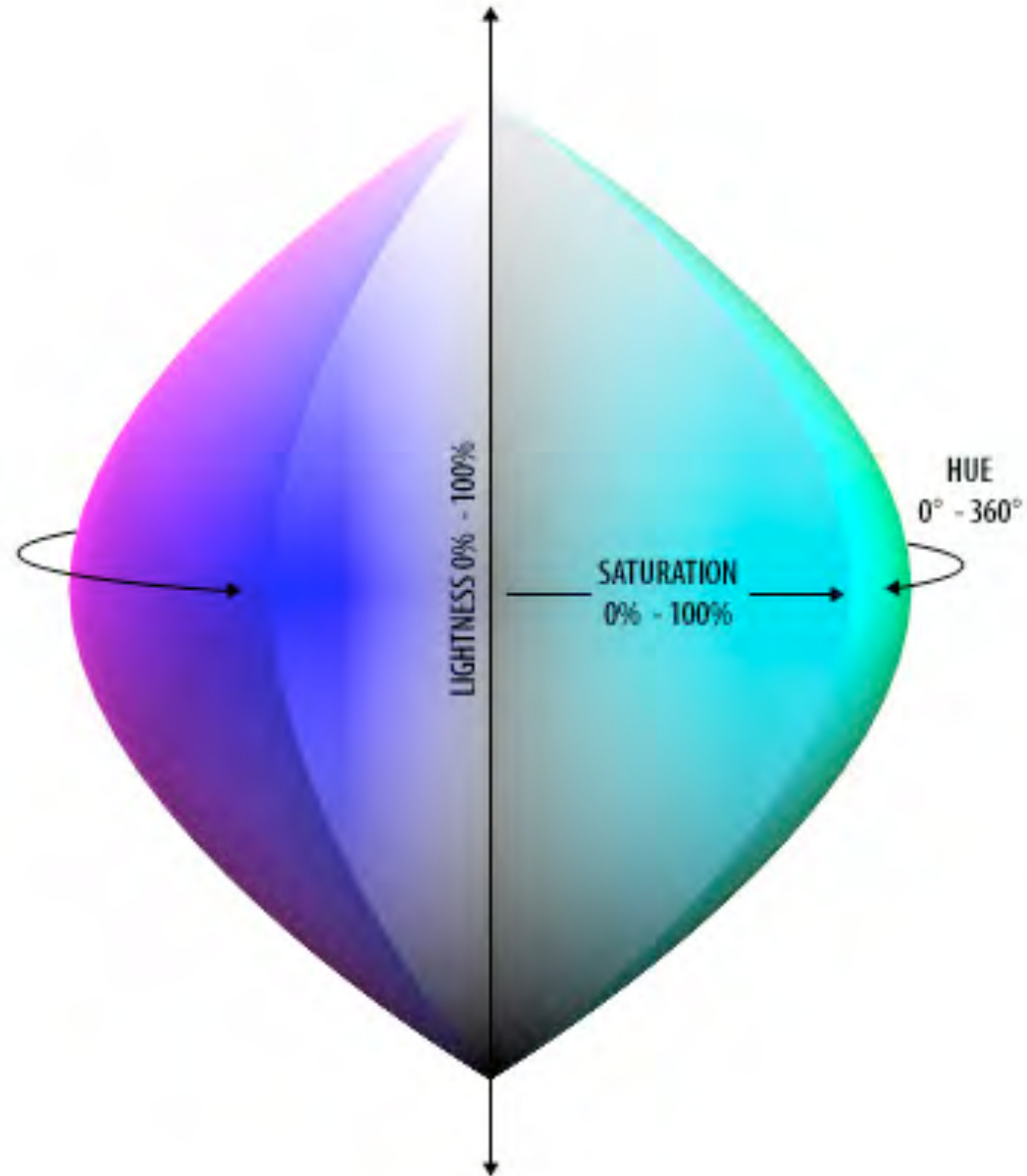
brightness ( $0\%$  -  $100\%$ )  
= energy or luminance

chromaticity = hue+saturation





# HSB colour model





disadvantage: hue coordinate is not continuous

0 and 360 degrees have the same meaning

but there is a huge difference in terms of numeric distance

example:

$$\text{red} = (0^\circ, 100\%, 50\%) = (360^\circ, 100\%, 50\%)$$

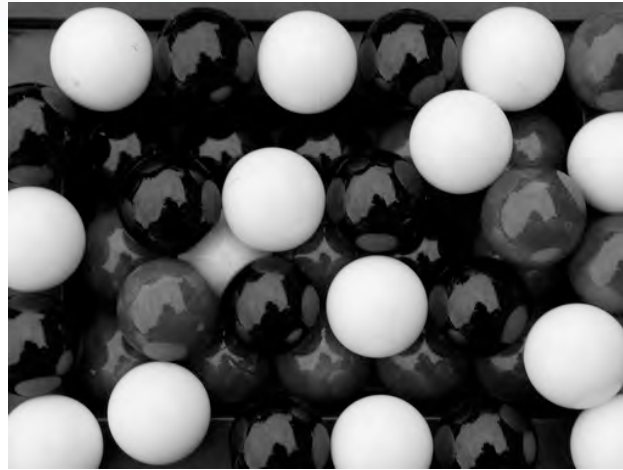
advantage: it is more natural to describe colour changes "brighter blue", "purer magenta", etc



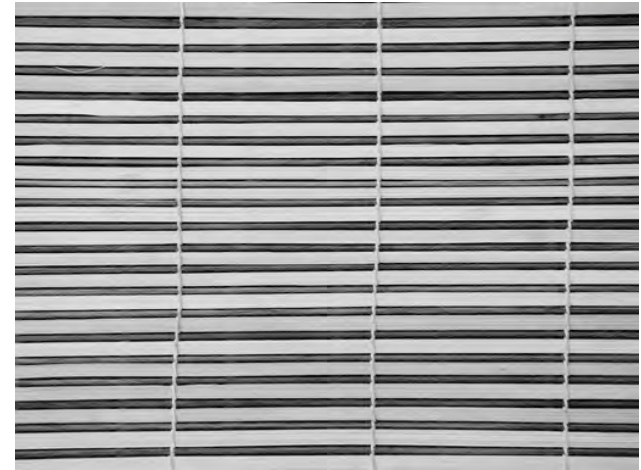
# Texture



coarseness



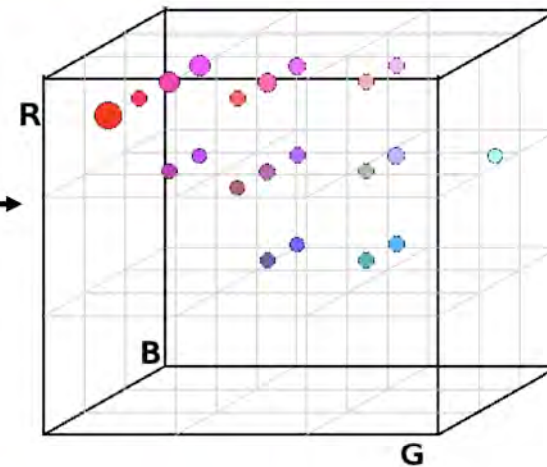
contrast



directionality



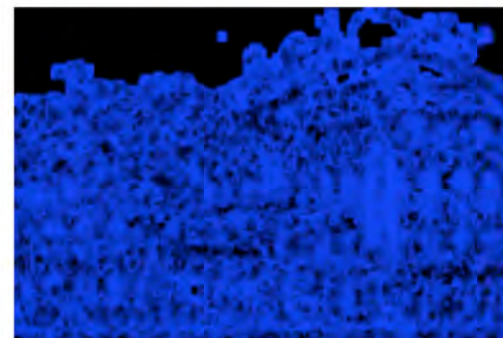
# Texture histograms



Coarseness



coNtrast



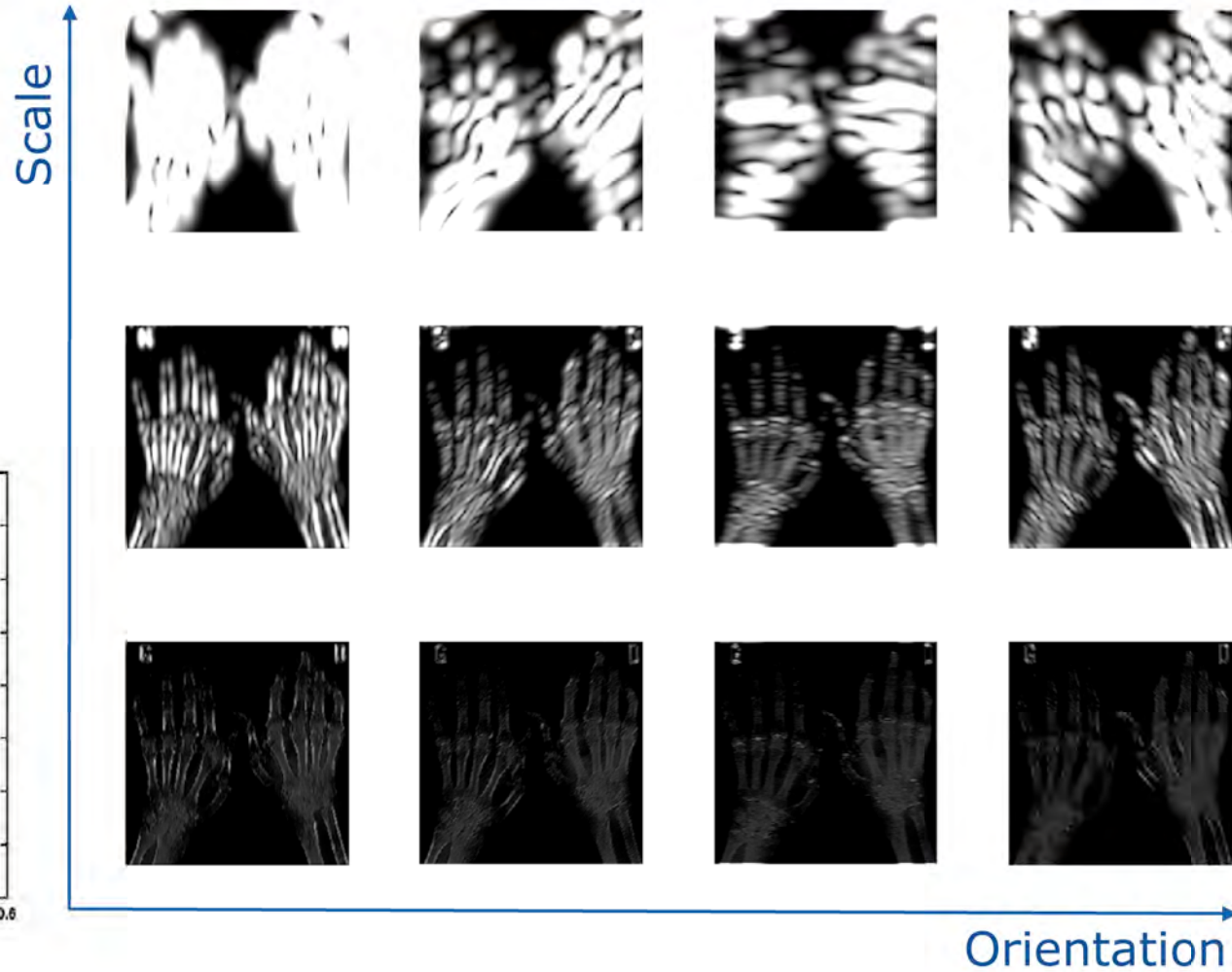
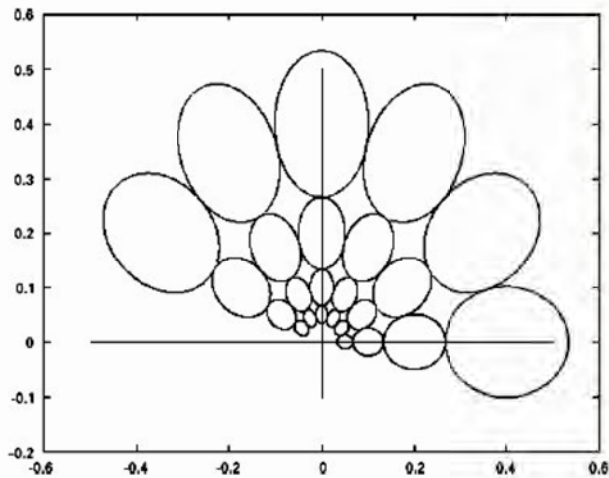
Directionality

[with Howarth, *IEE Vision, Image & Signal Proc* 15(6) 2004; Howarth PhD thesis]



# Gabor filter

## Query



[with Howarth, CLEF 2004]



shape = class of geometric objects invariant under translation

scale (changes keeping the aspect ratio)

rotations

information preserving description  
(for compression)

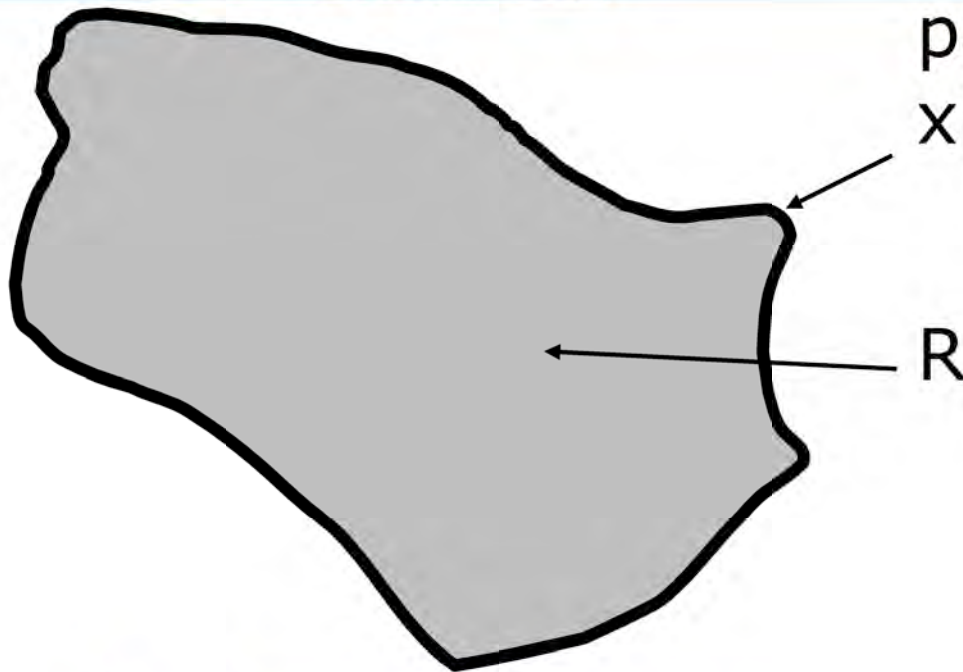
non-information preserving (for retrieval)

boundary based (ignore interior)

region based (boundary+interior)



# Perimeter and area



parameterised curve  
 $x(t), y(t)$

R

$$P = \int \sqrt{x'^2(t) + y'^2(t)} dt$$

~~$$A = \iint_R dx dy$$~~

~~boundary pixel count~~

count pixels in area



VS





## Global vs local

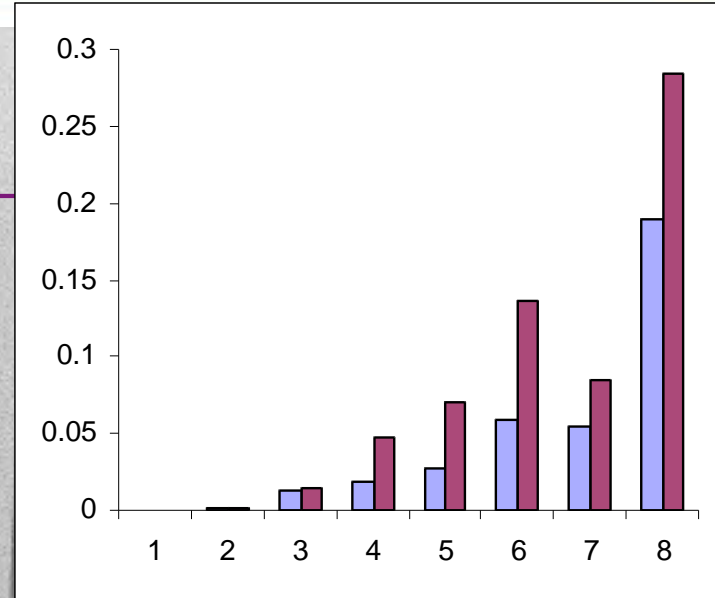
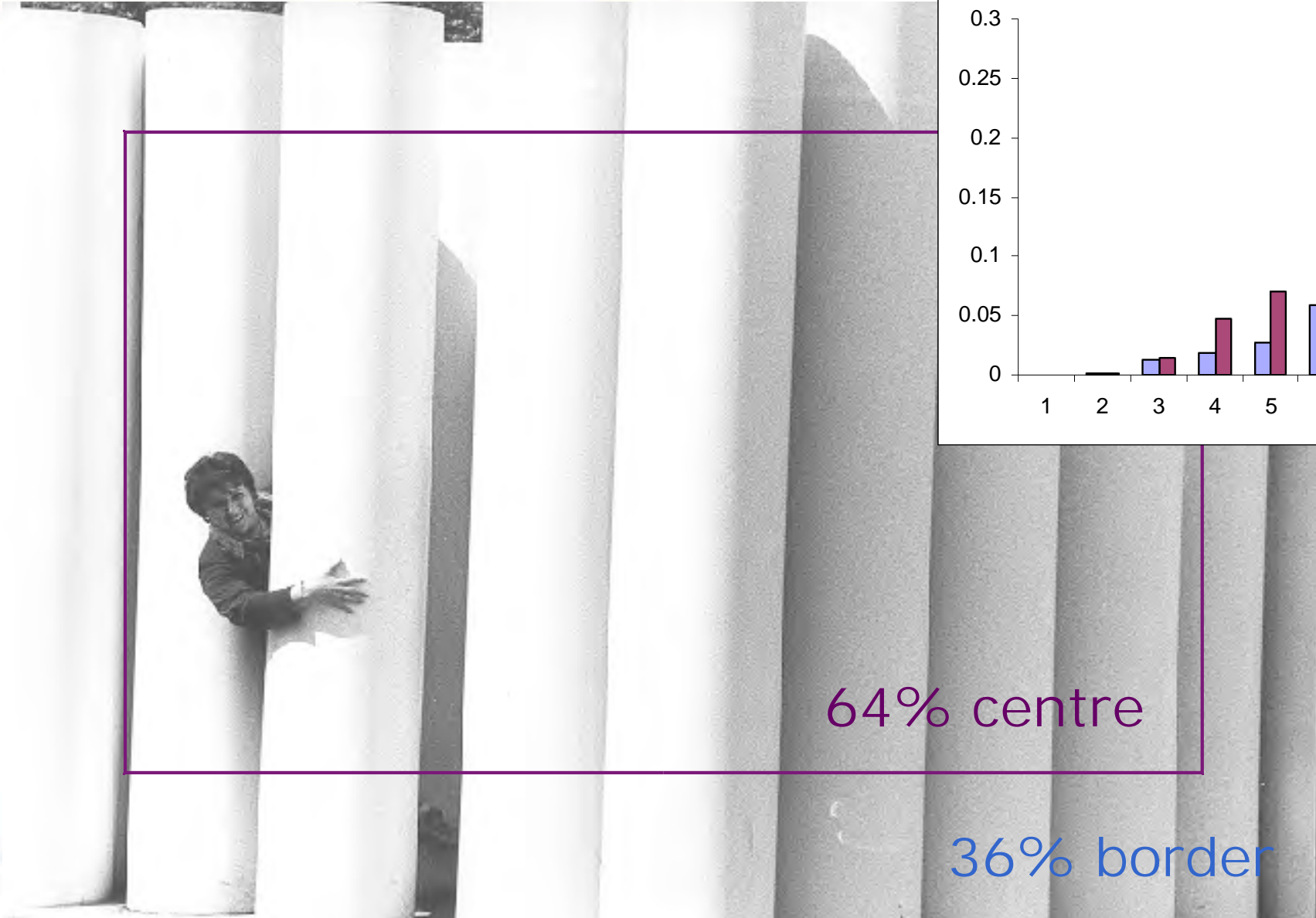


Global histogram also matches polar bears, marble floors, ...





# Localisation

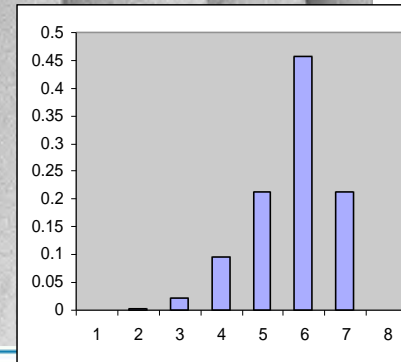
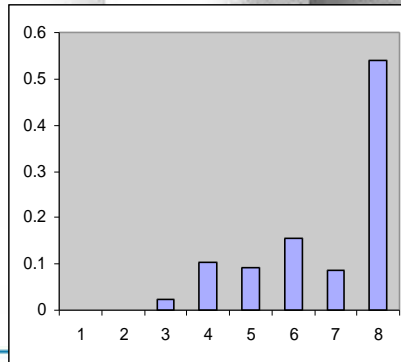
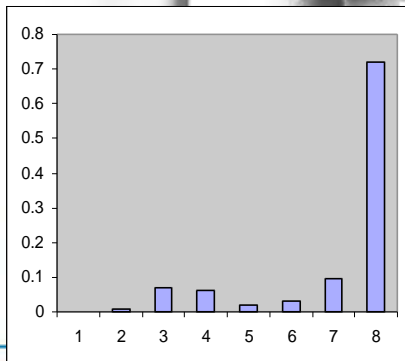
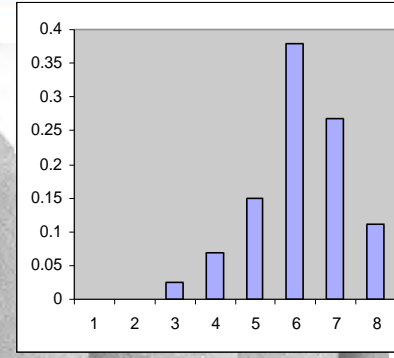
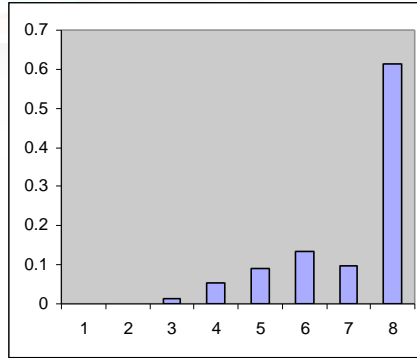
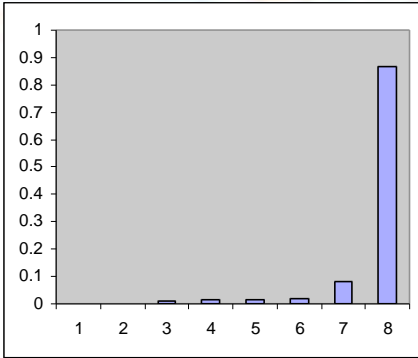


64% centre

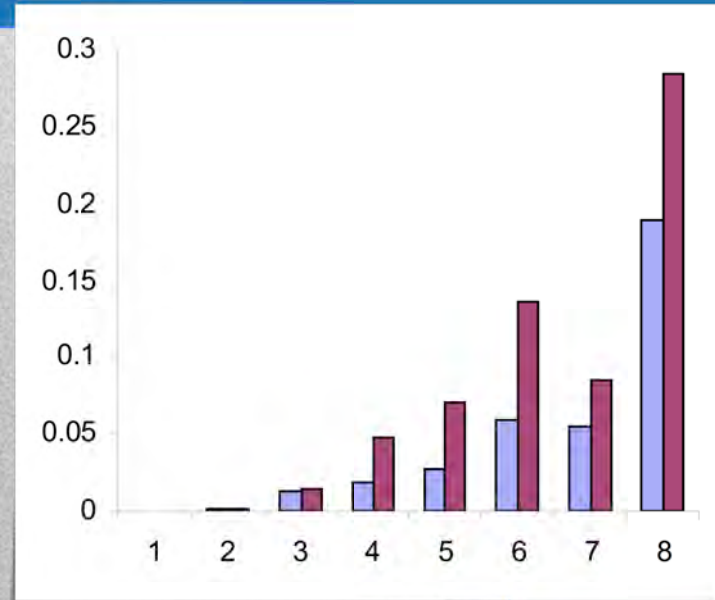
36% border



# Tiled Histograms



# Segmentation



foreground  
background

Many PoI, ie, many feature vectors  
Quantised feature vectors  $\approx$  words  
Bag of word model  $\approx$  text retrieval



gradual transition detection (eg, fade)

- accumulate distances

- long-range comparison

audio cues

- silence and/or speaker change

motion detection and analysis

- camera motion, zoom, object motion

- MPEG provides some motion vectors



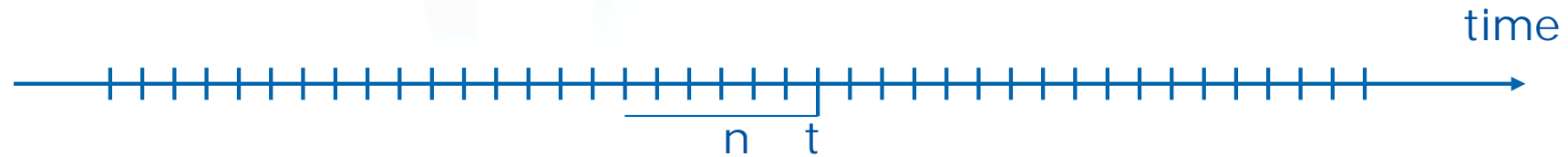
[Vlad Tanasescu: Anticipation, SCiFi trailer]



[Vlad Tanasescu: Anticipation, SCiFi trailer]



At time  $t$  define distance  $d_n(t)$



- compare frames  $t-n+i$  and  $t+i$  ( $i=0, \dots, n-1$ )
- average their respective distances over  $i$

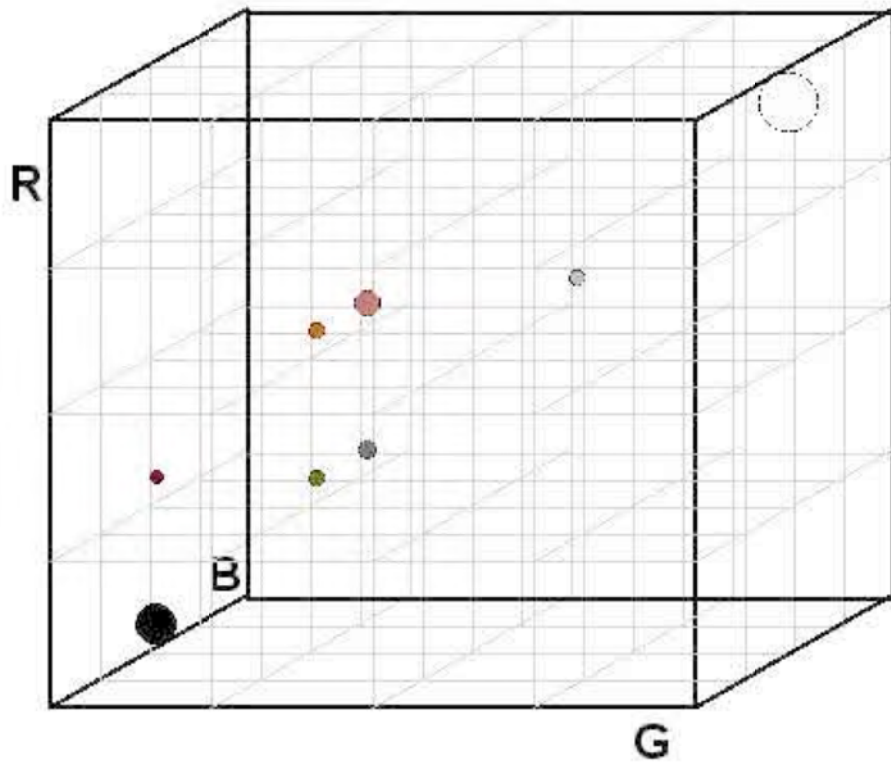
Peak in  $d_n(t)$  detected if

$d_n(t) > \text{threshold}$  and

$d_n(t) > d_n(s)$  for all neighbouring  $s$

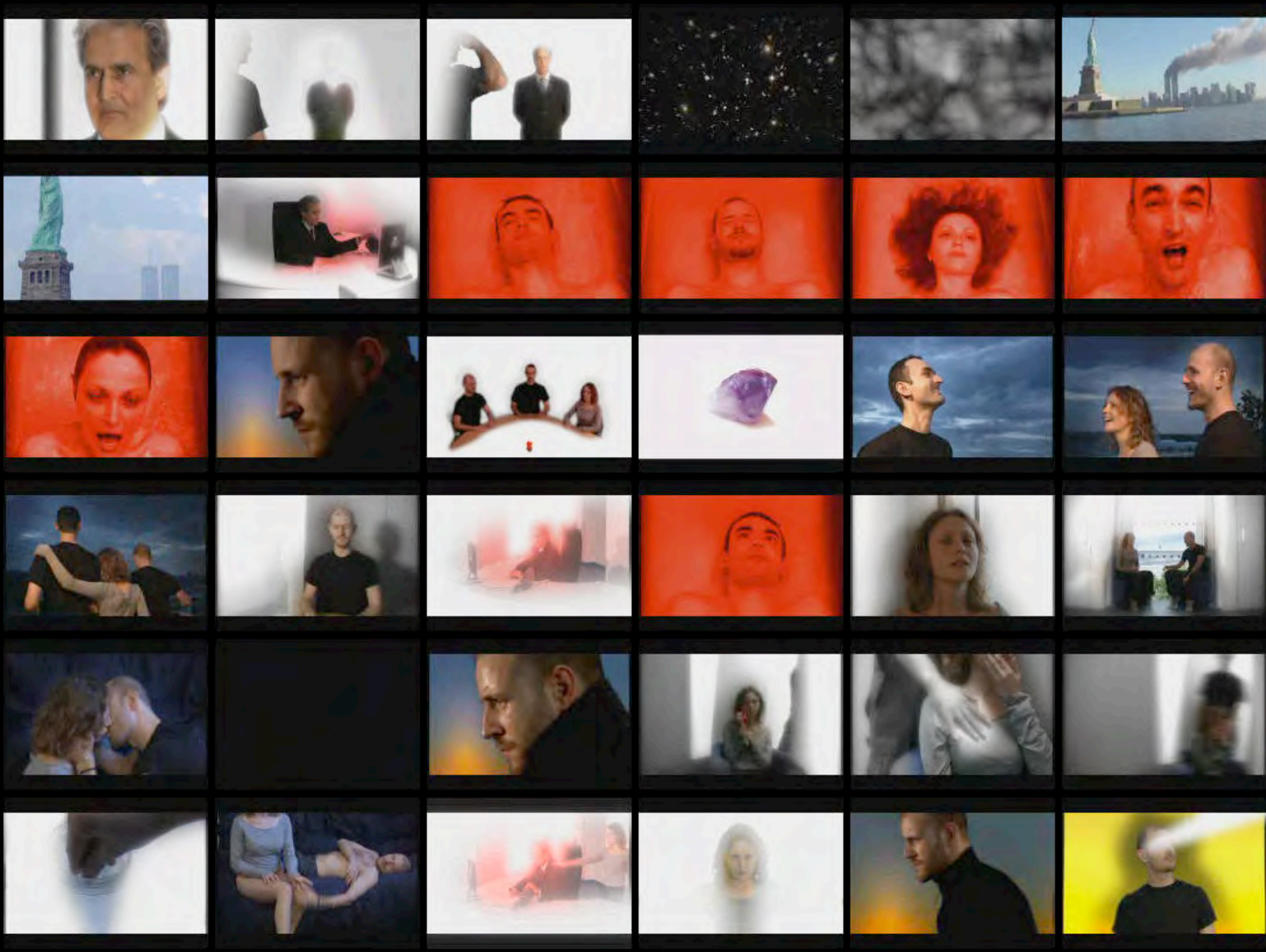
Shot = near-coincident peaks of  $d_{16}$  and  $d_8$

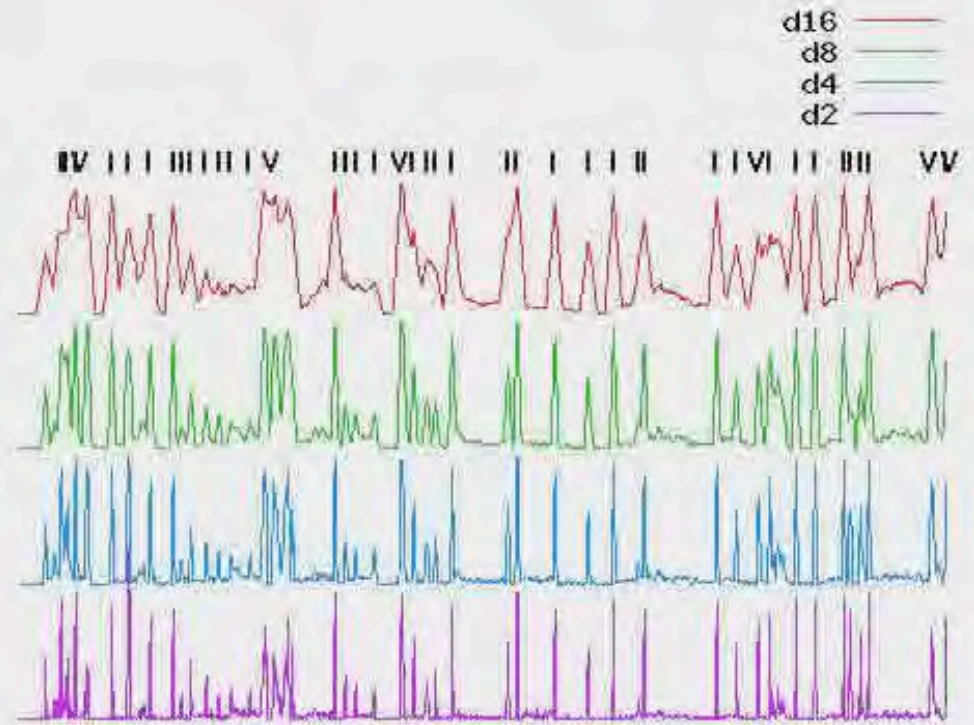
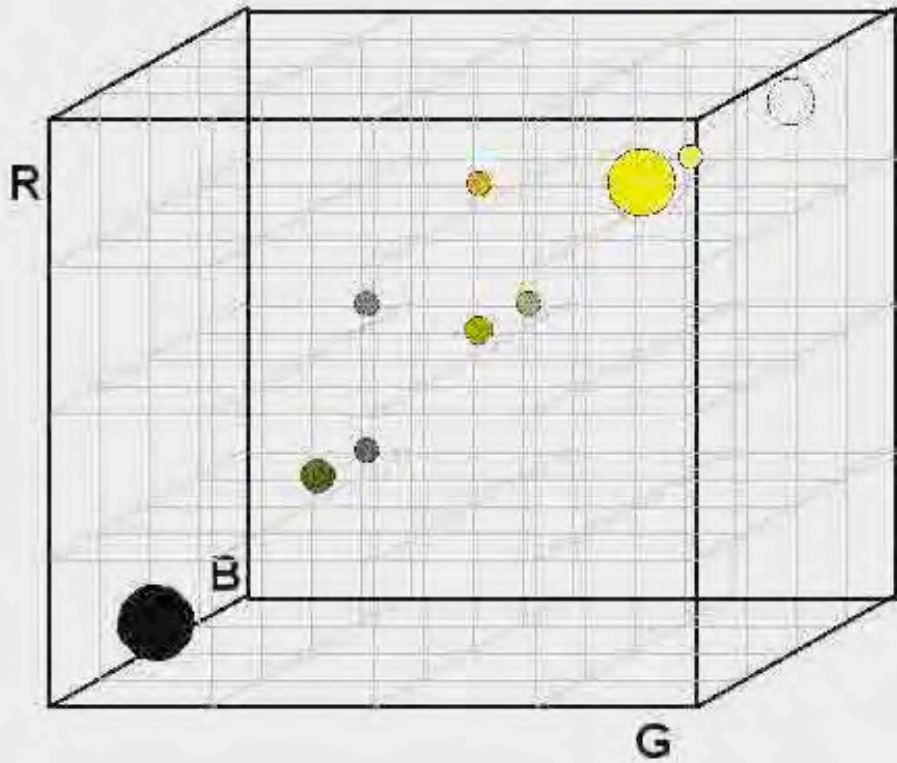
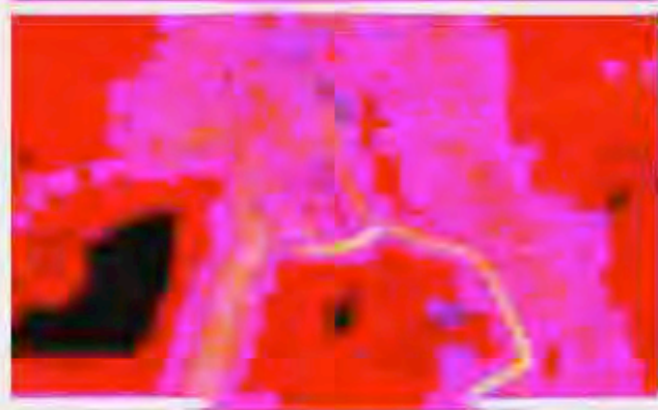




- d16 —
- d8 —
- d4 —
- d2 —

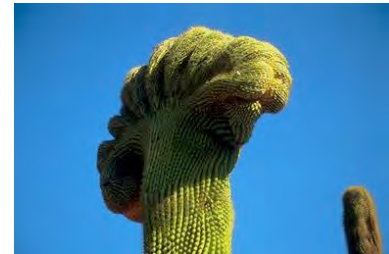
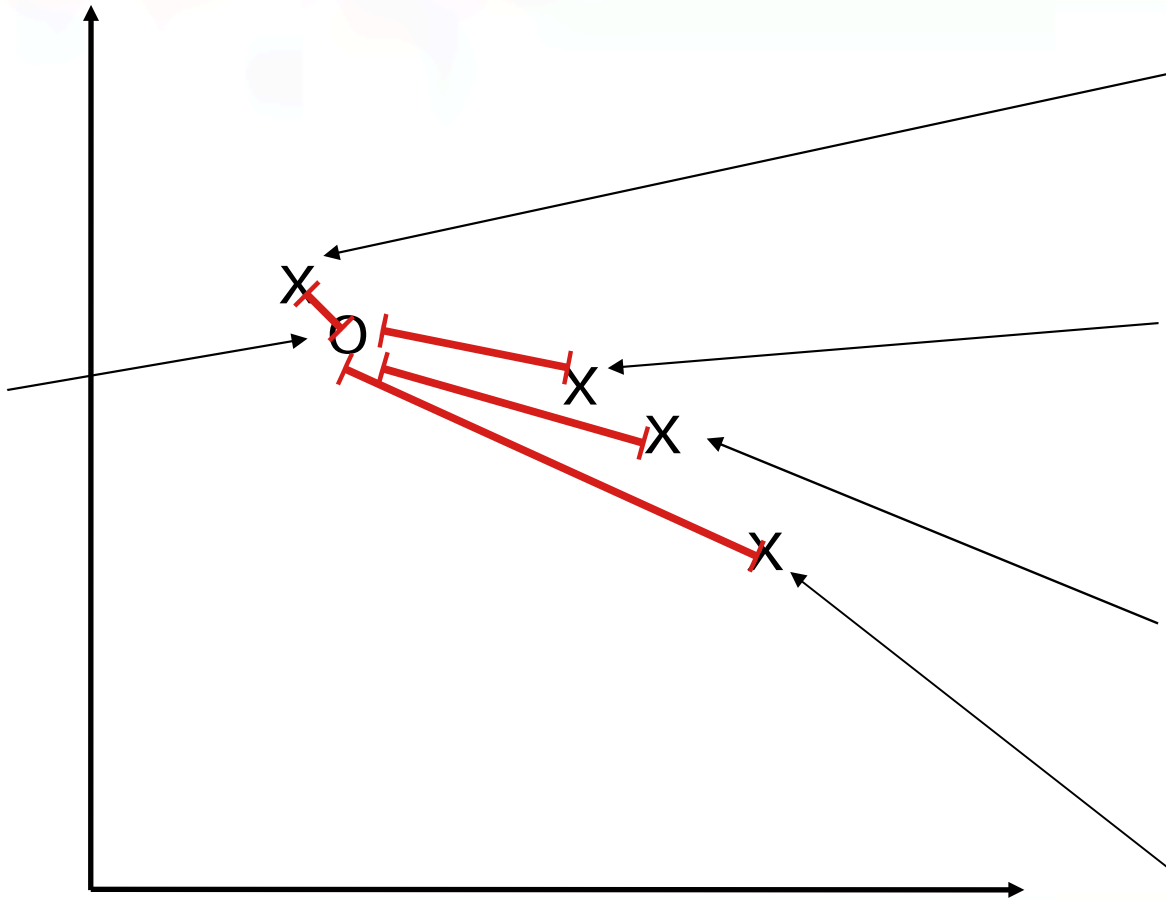








# Features and distances





assumes coding of MM objects as data vectors

distance measures

Euclidean, Manhattan

correlation measures

Cosine similarity measure

histogram intersection for normalised histograms

$$\text{sim}(h, q) = \sum_i \min(h_i, q_i)$$



$$d_p(\mathbf{v}, \mathbf{w}) = \sqrt[p]{\sum_i |\mathbf{v}_i - \mathbf{w}_i|^p}, \quad p \geq 1$$

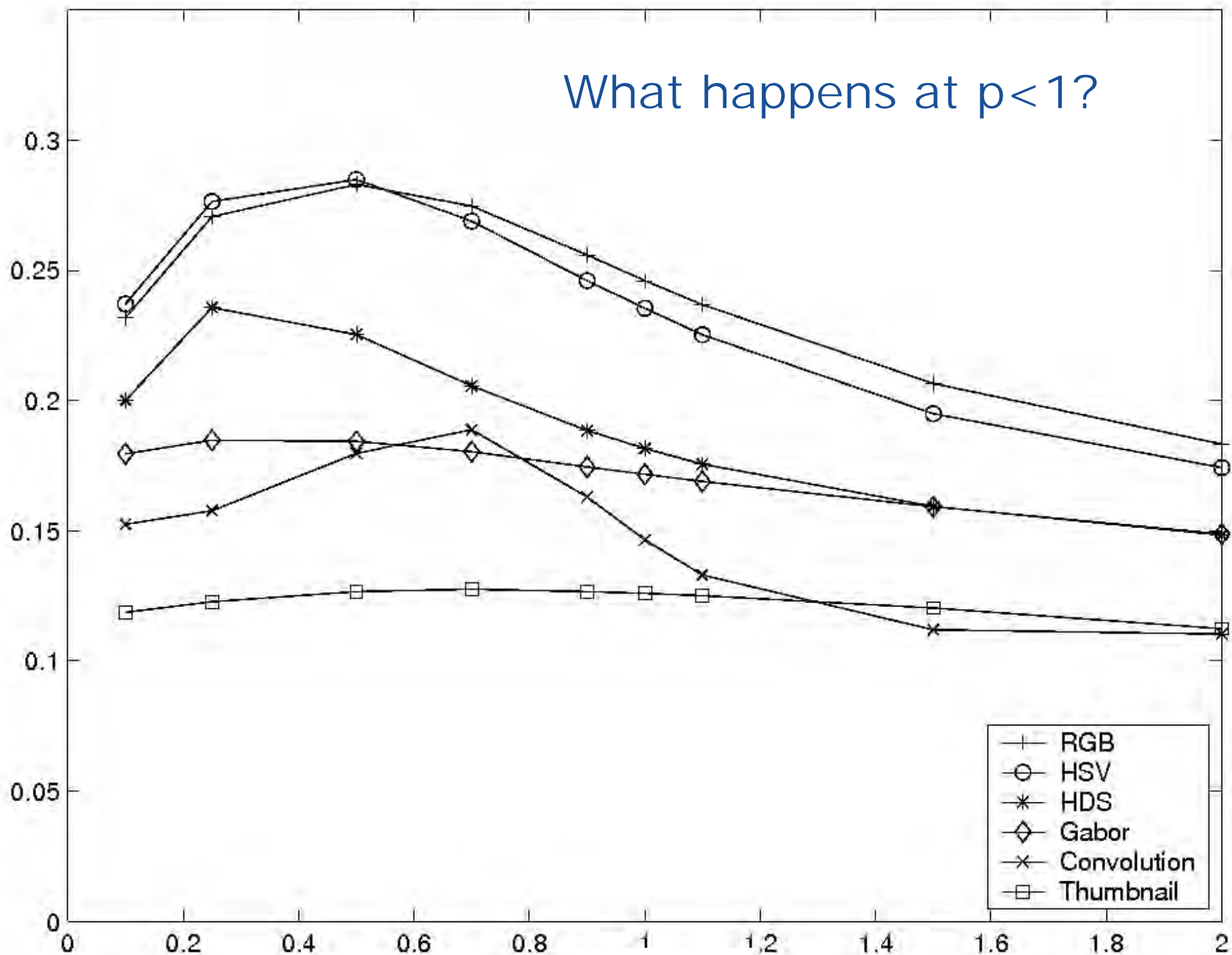
$$d_2(\mathbf{v}, \mathbf{w}) = \sqrt{\sum_i |\mathbf{v}_i - \mathbf{w}_i|^2}$$

$$d_1(\mathbf{v}, \mathbf{w}) = \sqrt[1]{\sum_i |\mathbf{v}_i - \mathbf{w}_i|^1} = \sum_i |\mathbf{v}_i - \mathbf{w}_i|$$

$$d_\infty(\mathbf{v}, \mathbf{w}) =$$

What happens at  $p < 1$ ?

Mean average precision



$p$

[with Howarth, ECIR 2005]





- Squared chord
- Earth Mover's Distance
- Chi squared distance
- Kullback-Leibler divergence (not a true distance)
- Ordinal distances (for string values)



## Squared chord

$$d_{sc}(v, w) = \sum_{i=1}^n (\sqrt{v_i} - \sqrt{w_i})^2$$

[with Liu et al, AIRS 2008; with Hu et al, ICME 2008]



## Recap: Multimedia information retrieval

1. What is multimedia information retrieval?
2. Metadata and piggyback retrieval
3. Multimedia fingerprinting
4. Automated annotation
5. Content-based retrieval



# Multimedia Information Retrieval



Prof Stefan Ruger  
Multimedia and Information Systems  
Knowledge Media Institute  
The Open University  
<http://kmi.open.ac.uk/mmis>

# MMIS

**Multimedia and Information Systems**