



- 1 What is multimedia information retrieval?
 - 1.1 Information retrieval
 - 1.2 Multimedia
 - 1.3 Semantic Gap?
 - 1.4 Challenges of automated multimedia indexing
- 2 Basic multimedia search technologies
 - 2.1 Meta-data driven retrieval
 - 2.2 Piggy-back text retrieval
 - 2.3 Automated annotation
 - 2.4 Fingerprinting
 - 2.5 Content-based retrieval
 - 2.6 Implementation Issues
- 3 Evaluation of MIR Systems
- 4 Added value



- 1 What is multimedia information retrieval?
 - 1.1 Information retrieval
 - 1.2 Multimedia
 - 1.3 Semantic Gap?
 - 1.4 Challenges of automated multimedia indexing
- 2 Basic multimedia search technologies**
 - 2.1 Meta-data driven retrieval
 - 2.2 Piggy-back text retrieval
 - 2.3 Automated annotation
 - 2.4 Fingerprinting
 - 2.5 Content-based retrieval**
 - 2.6 Implementation Issues**
- 3 Evaluation of MIR Systems
- 4 Added value



Why content-based?

Actually, what is content-based search?

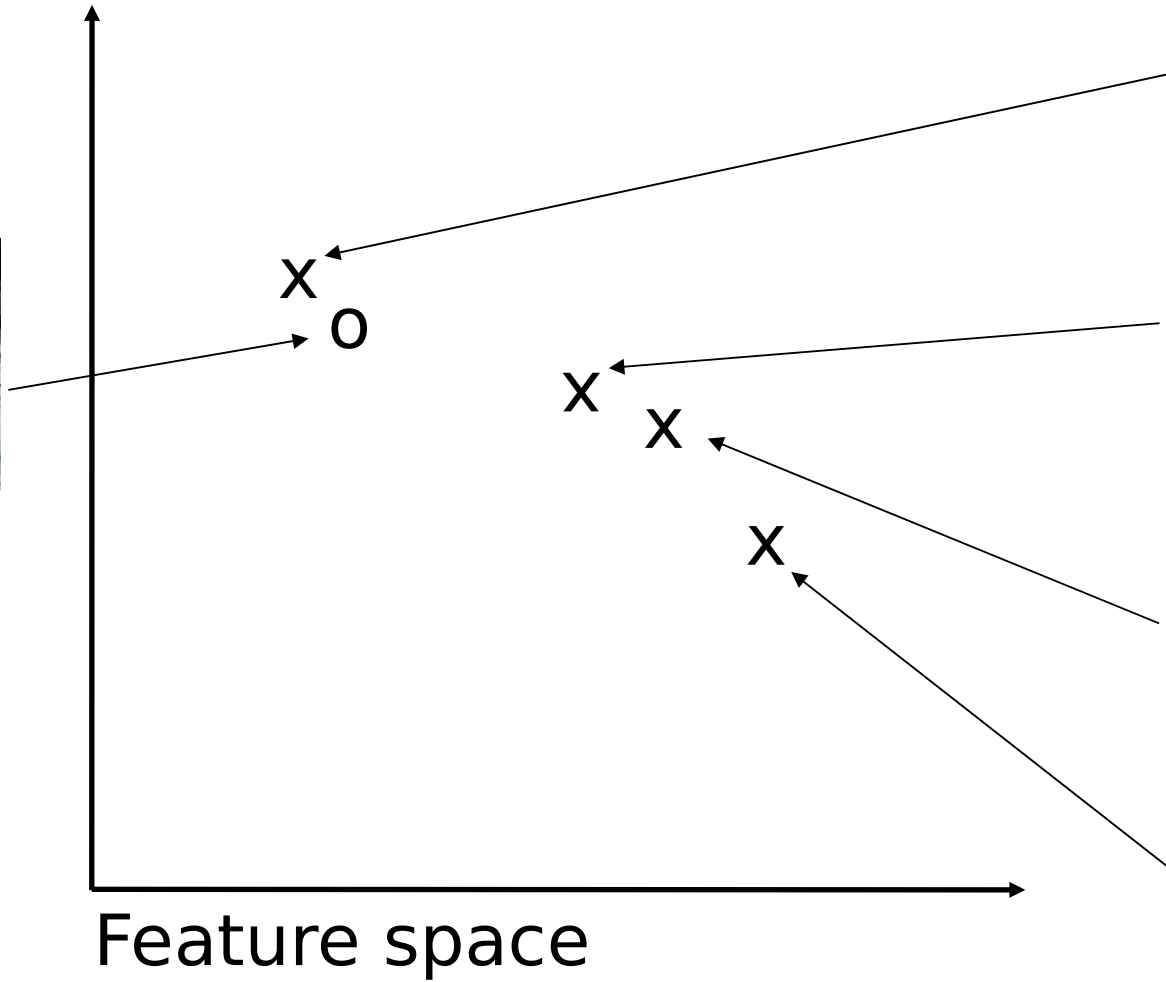
Is human thinking content-based?

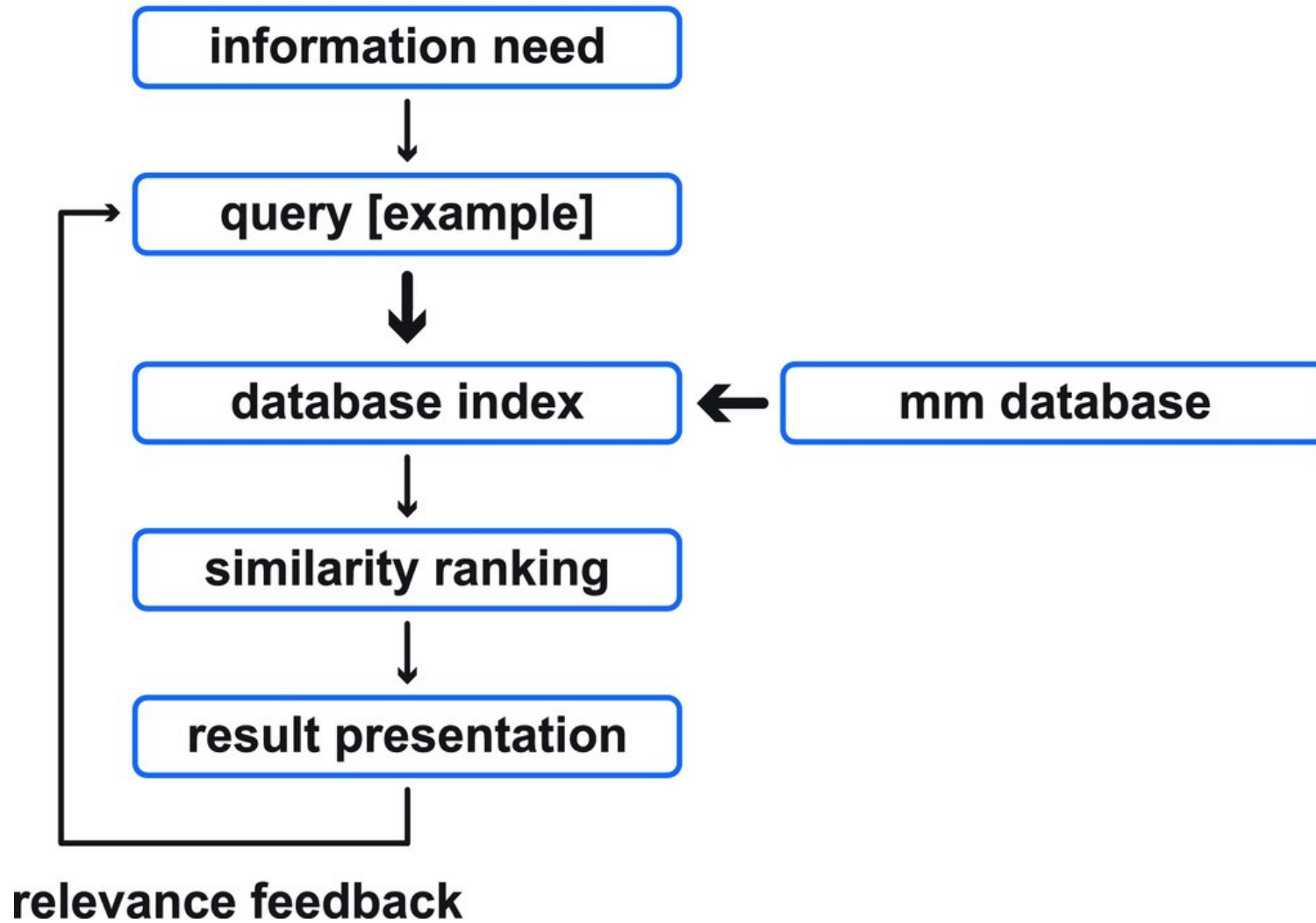
Metadata annotation (text) is good but

-
-
-
-



Features and distances







Visual

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

Audio

Temporal

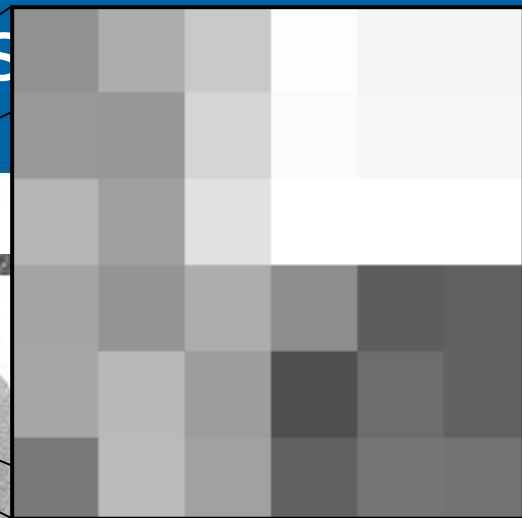
How to describe the features?

For people

For computers

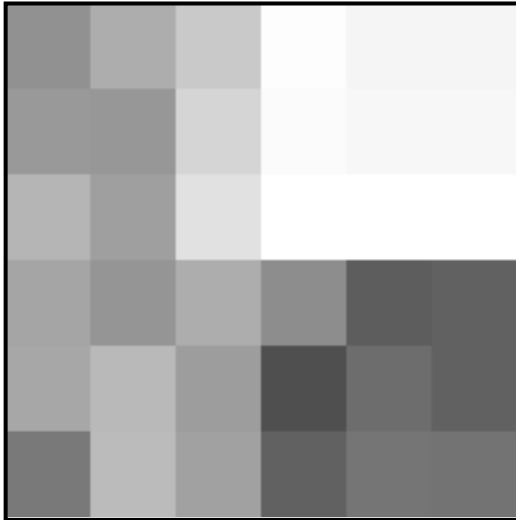


Digital Images





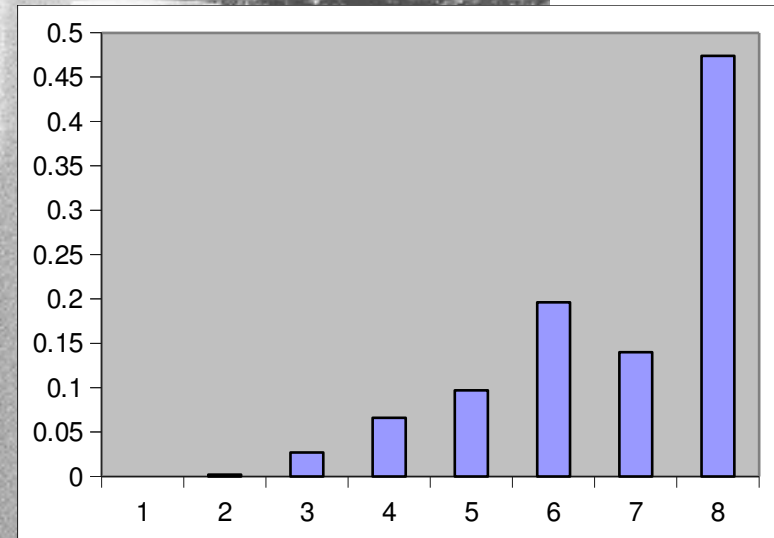
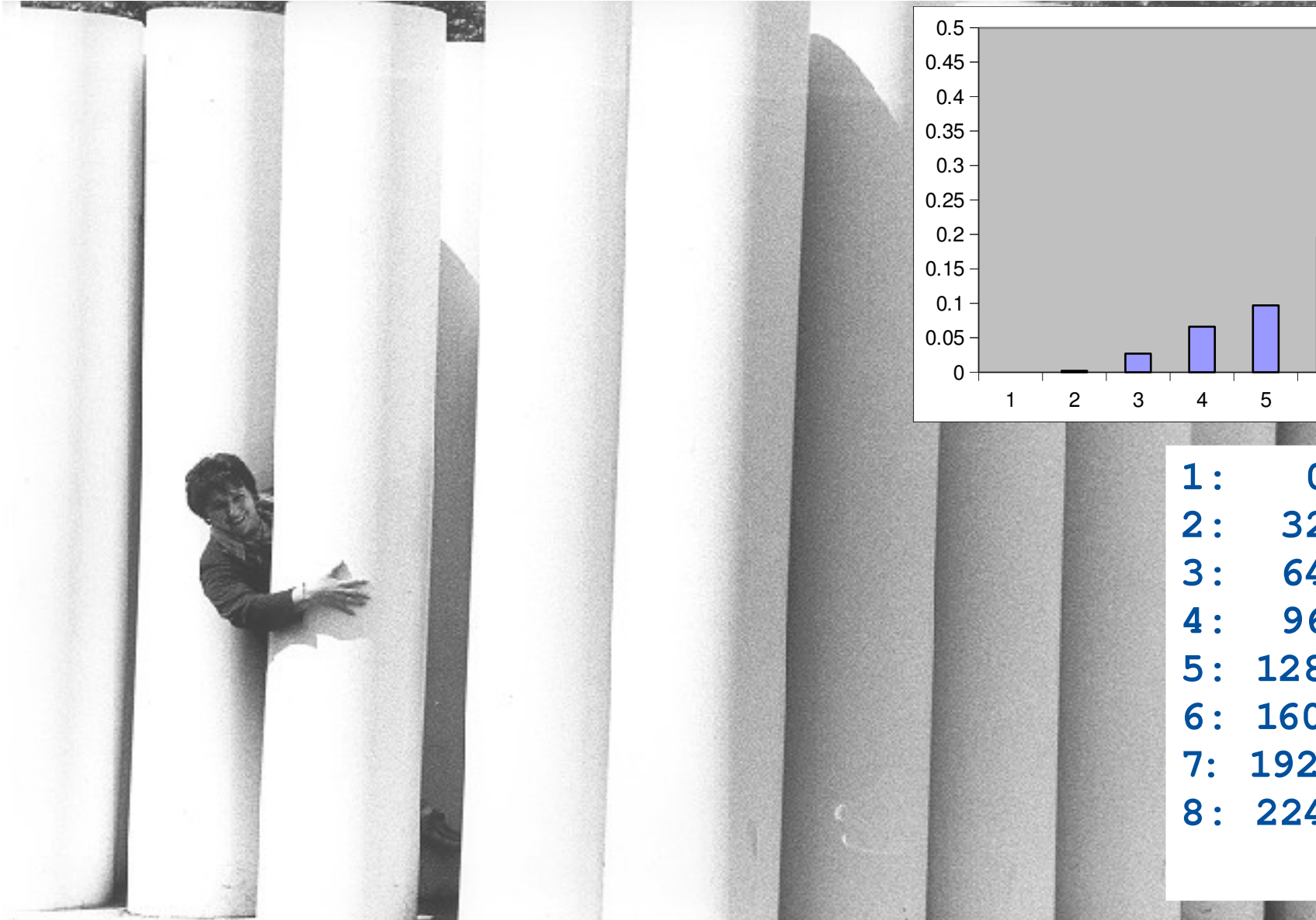
Content of an image



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

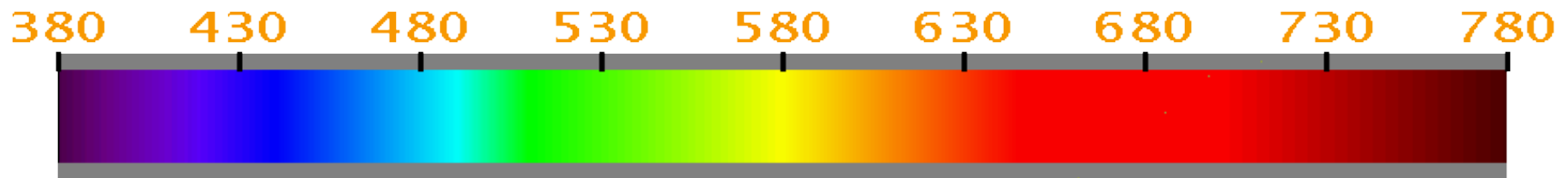


1:	0	-	31
2:	32	-	63
3:	64	-	95
4:	96	-	127
5:	128	-	159
6:	160	-	191
7:	192	-	223
8:	224	-	255



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

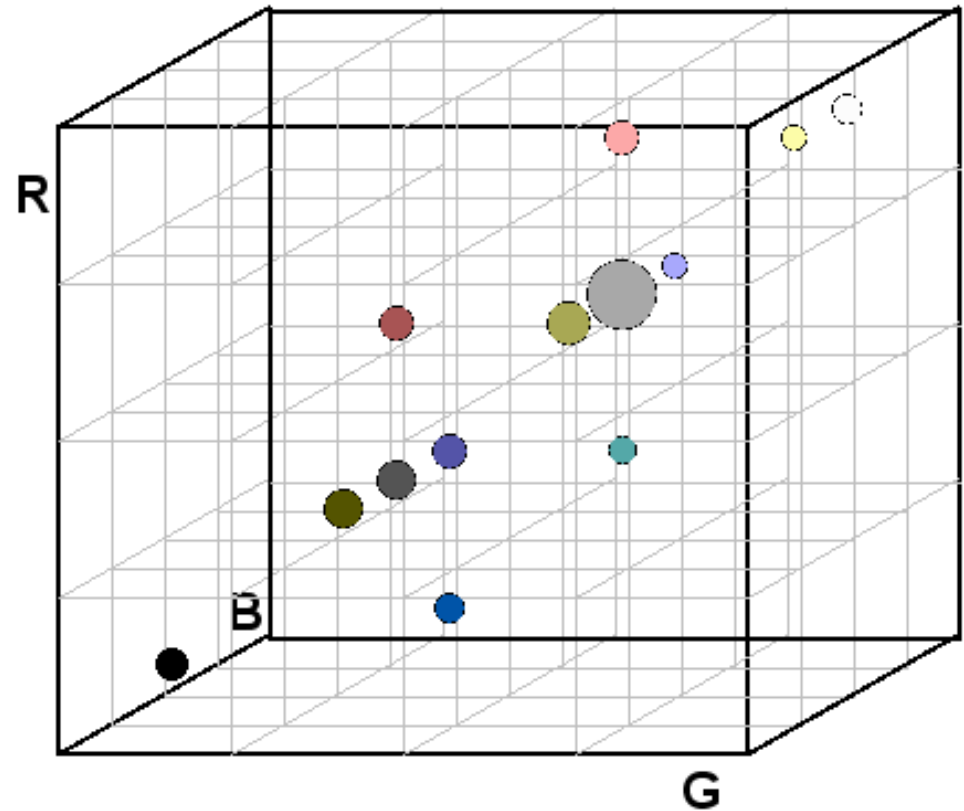
blue cyan green yellow red



spectral colours: wavelength (nm)











Colour histogram





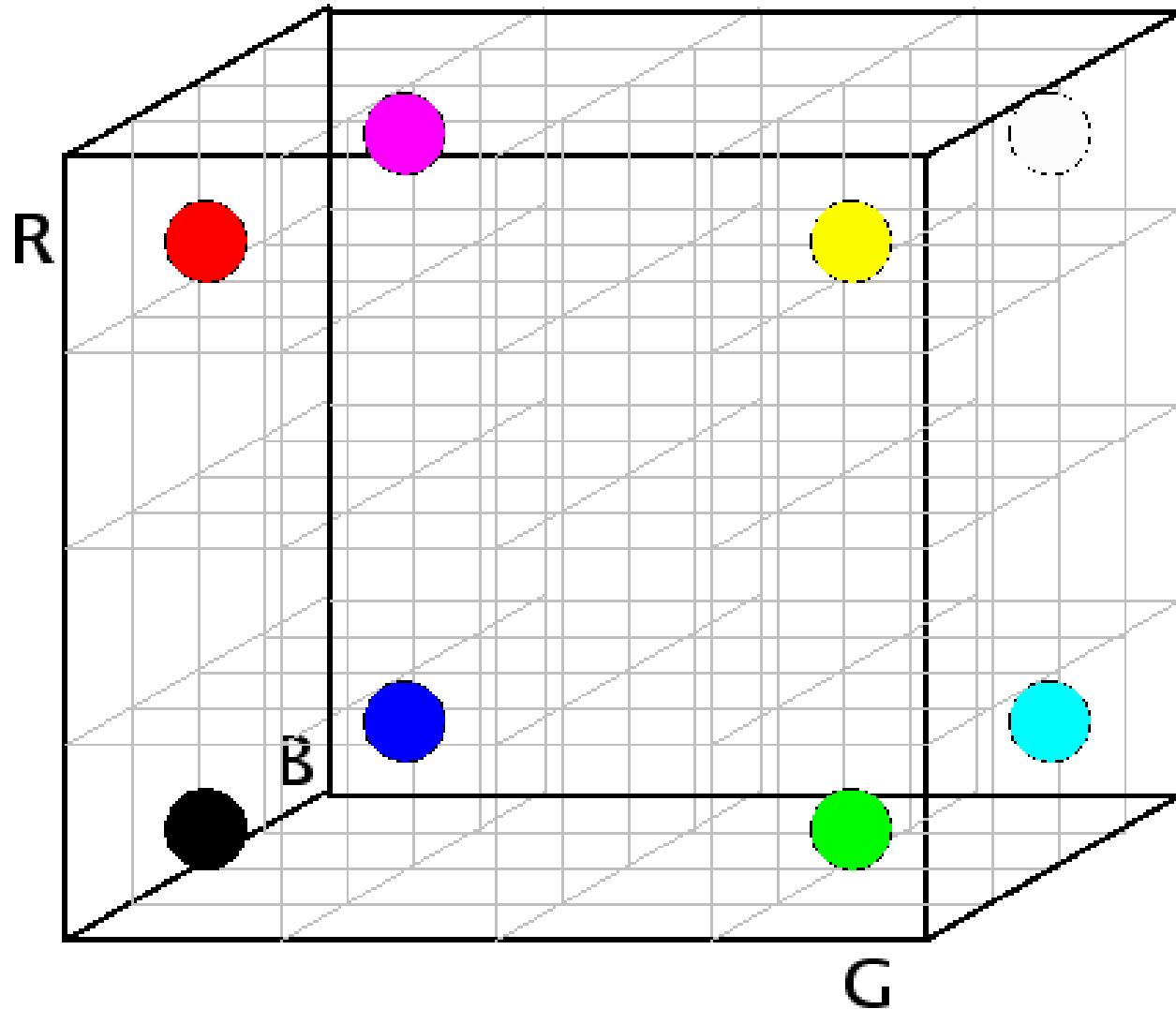
Sketch a 3D colour histogram for



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



Solution

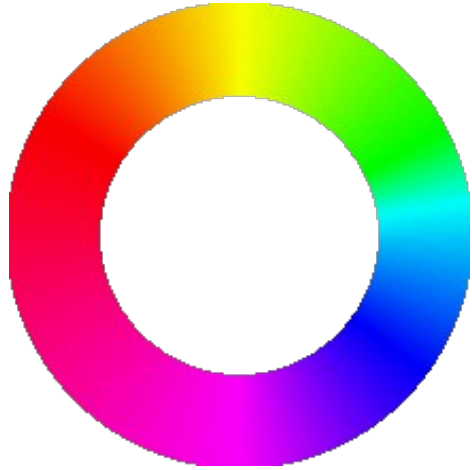




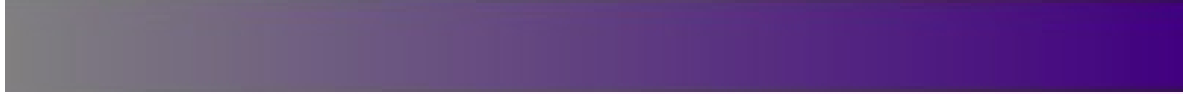
HSV, HSL, CIELAB/CIELUV




HSB colour model



hue (0° - 360°)
spectral colour



saturation (0% - 100%)
= spectral purity

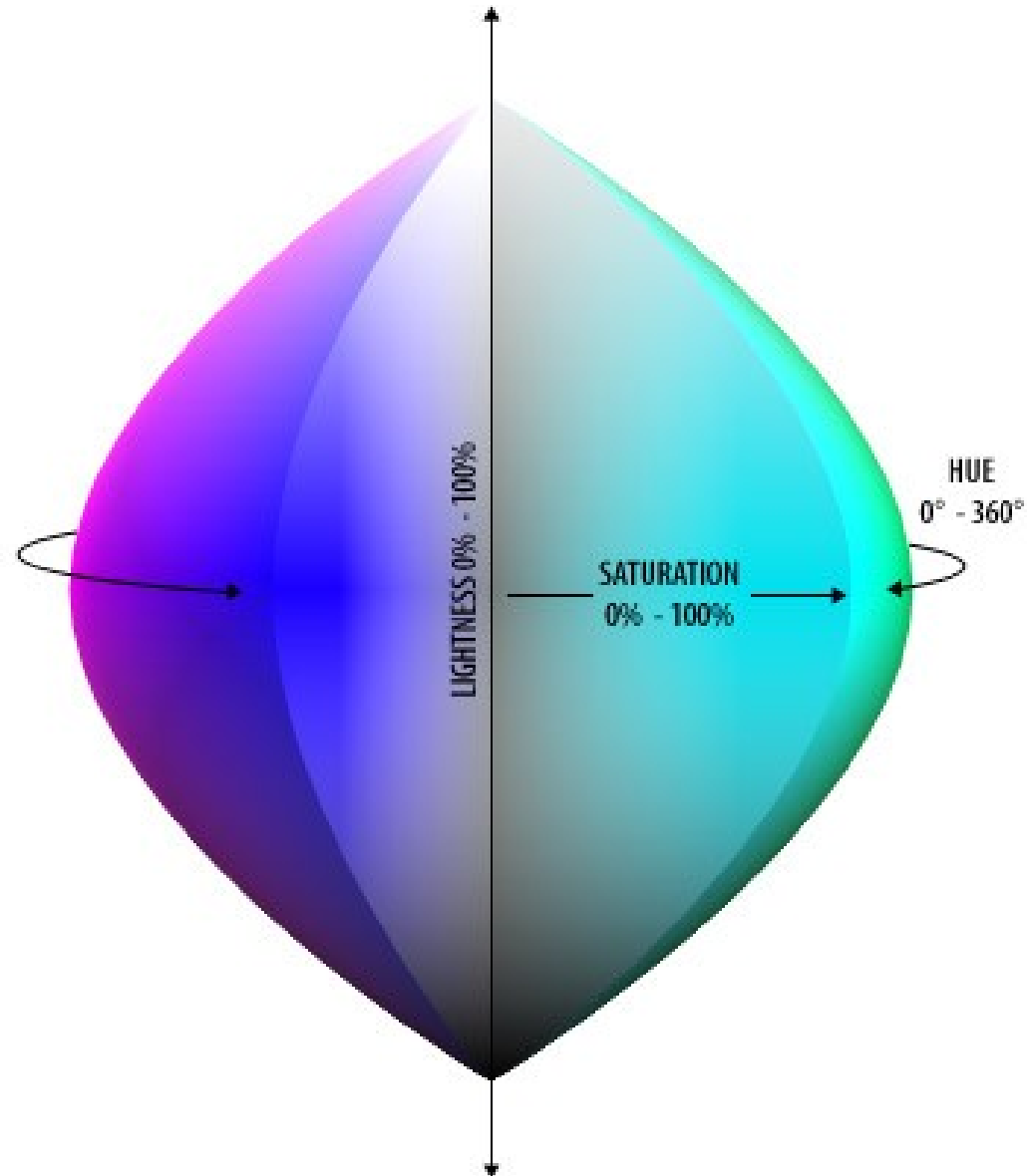


brightness (0% - 100%)
= energy or luminance

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:

red = (0,100%,50%) = (360,100%,50%)

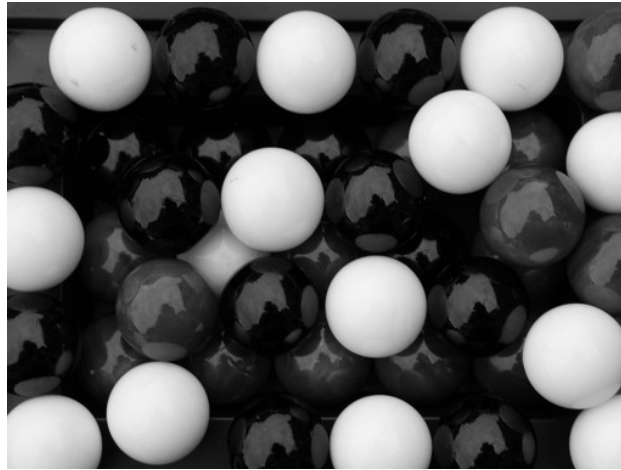
advantage: it is more natural to describe colour changes “brighter blue”, “purer magenta”, etc



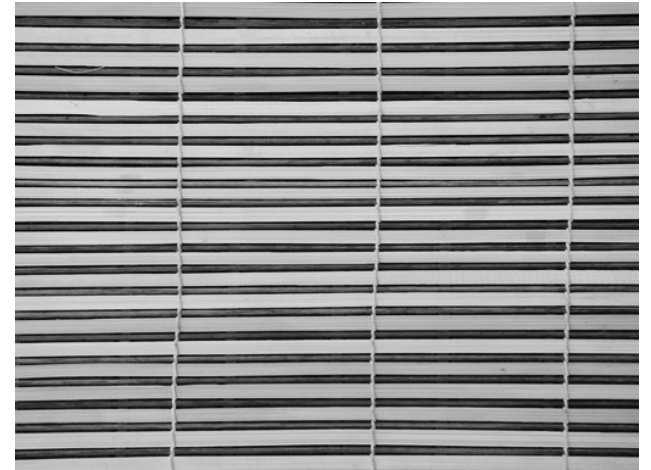
Texture



coarseness



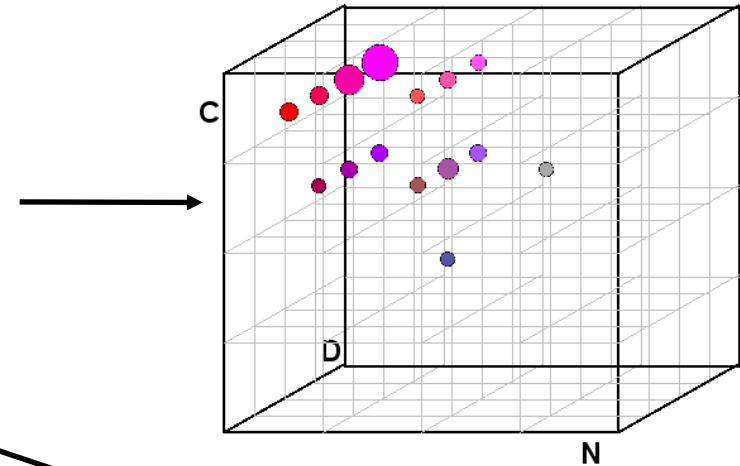
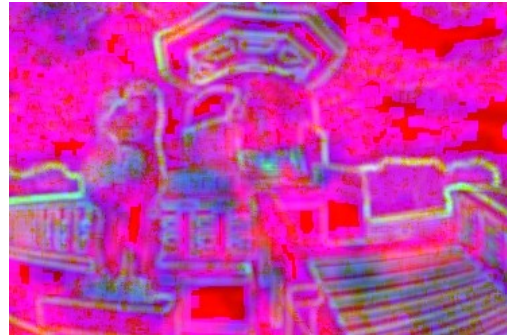
contrast



directionality



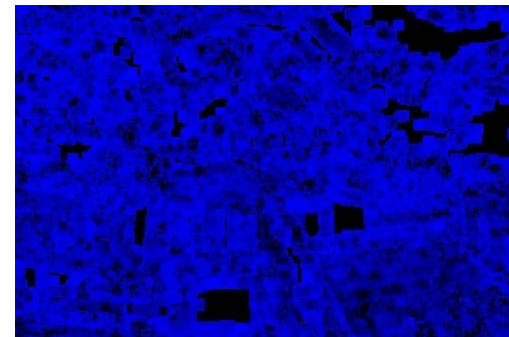
Texture histograms



Coarseness



coNtrast



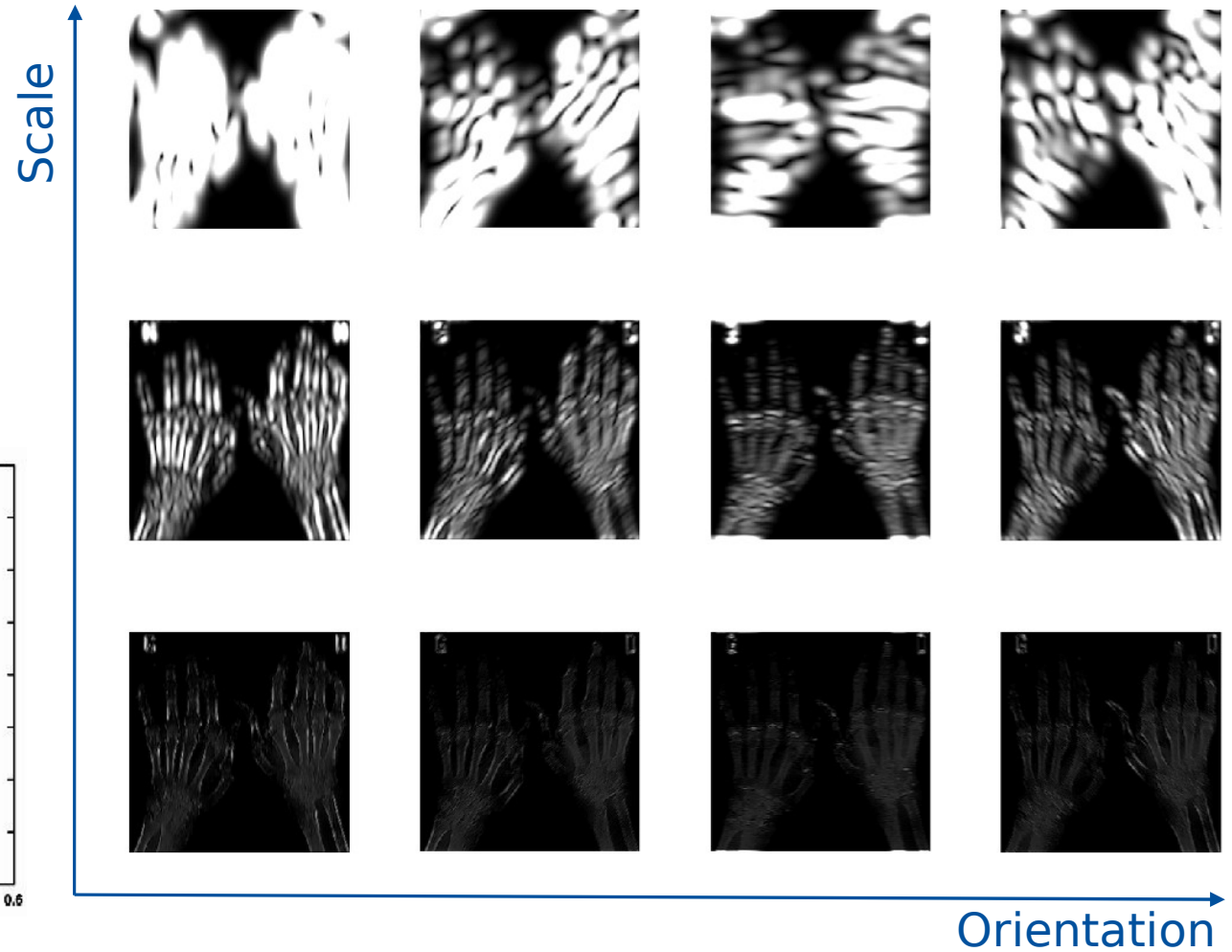
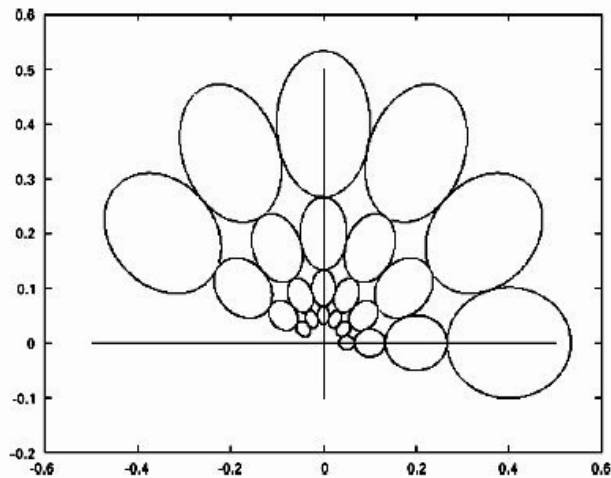
Directionality

[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)

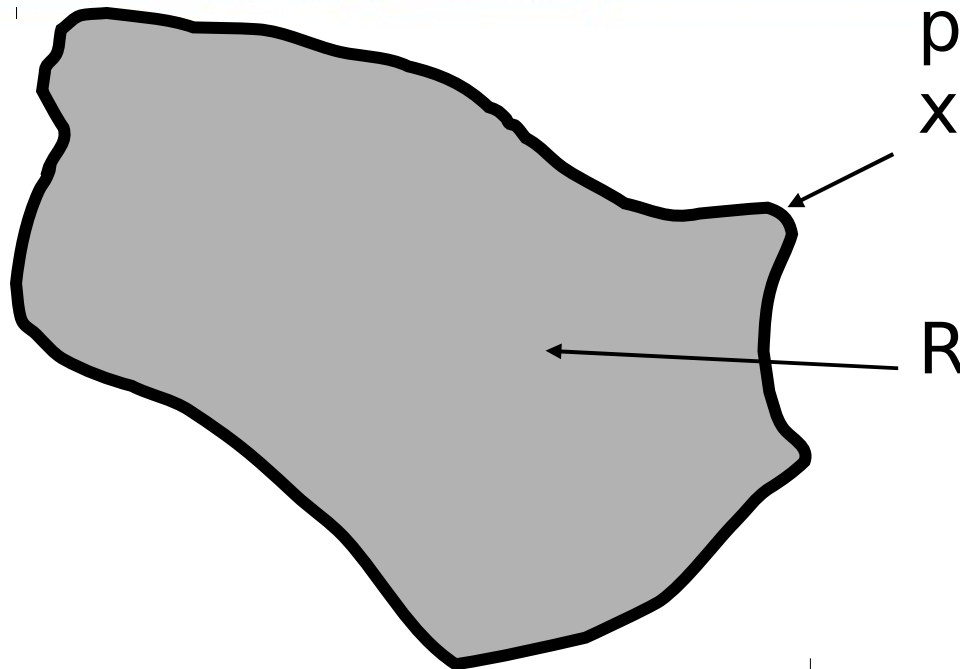
region based (boundary+interior)



- boundary based
 - perimeter & area
 - corner points
 - circularity
 - chain codes
- region based (considering interior and holes, ...)
 - not covered here



Perimeter and area



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

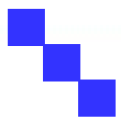
R

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

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

~~boundary pixel count~~

count pixels in area



VS





Circularity

$$T = 4\pi \frac{A}{P^2}$$

A=area, P=perimeter

T is 1 for a circle

T is smaller than 1 for all other shapes

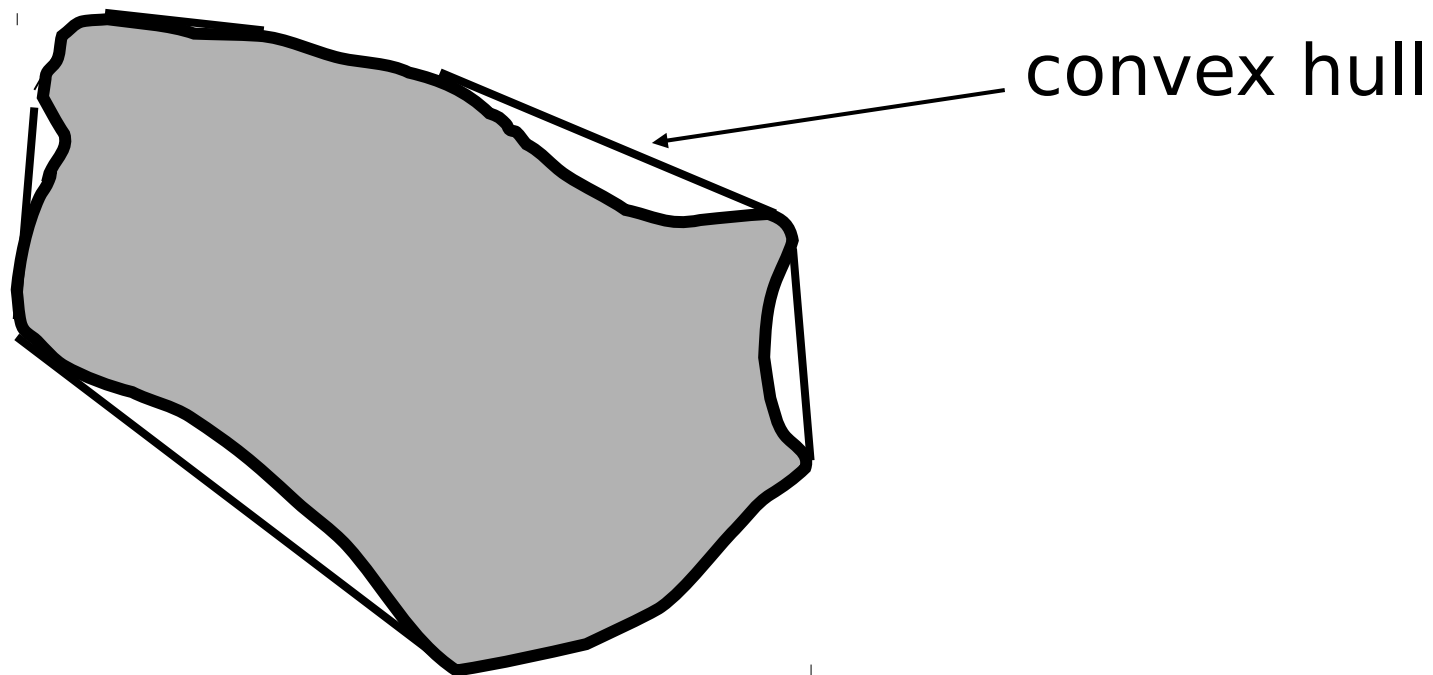
circularity is aka compactness



Convexity

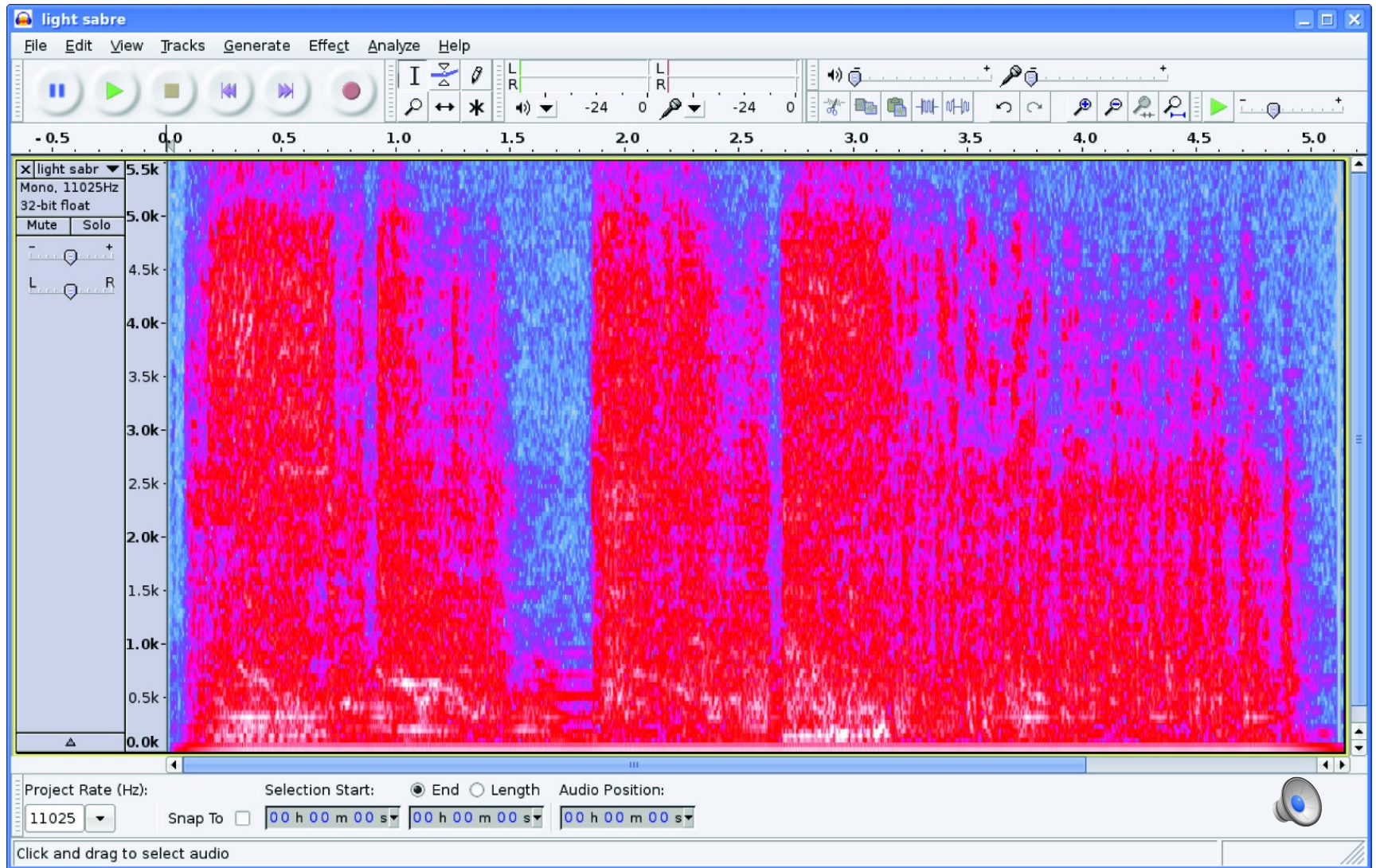
ratio of perimeter of convex hull and the original curve

1 for convex shapes, less than 1 otherwise





Sound





- Spectrogram
 - graph of frequencies/energy/time
- tempo, pitch, mode

- See

Z Liu,

Y Wang and T Chen (1998). Audio feature extraction and analysis for scene segmentation and classification. *VLSI Signal Processing 20*(1-2), 69-79.

Histograms

Condensed

Content-based

Real-valued vector

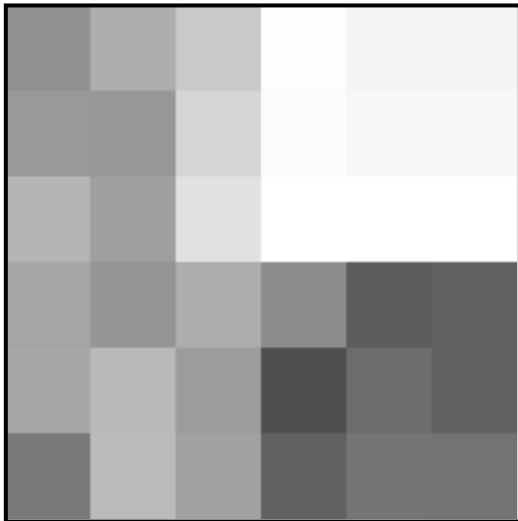
Summarising

Sparseness

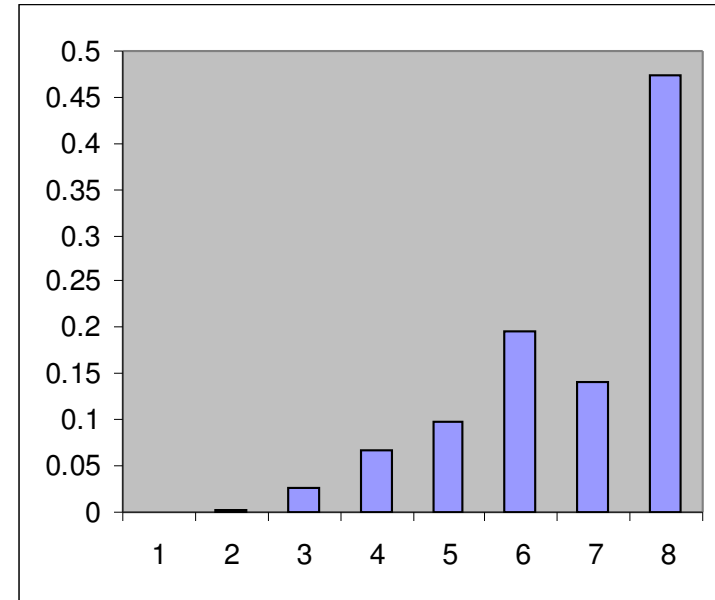
Statistical moments



Feature vectors → histograms



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



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



Simple statistics

μ Mean

\bar{p}_2 Variance (squared standard deviation)

\bar{p}_3 3rd central moment (skewness)

$$\bar{p}_n = \frac{1}{wh} \sum_{i=1}^w \sum_{j=1}^h (p(i, j) - \mu)^n$$

where w is image width and h is image height



Moment features

$$(\mu, \sqrt{\bar{p}_2}, \text{sign}(\bar{p}_3) \sqrt[3]{|\bar{p}_3|})$$



Moment features

$$\left(\begin{array}{l} \mu_r, \sqrt{\bar{r}_2}, \text{sign}(\bar{r}_3) \sqrt[3]{|\bar{r}_3|}, \\ \mu_g, \sqrt{\bar{g}_2}, \text{sign}(\bar{g}_3) \sqrt[3]{|\bar{g}_3|}, \\ \mu_b, \sqrt{\bar{b}_2}, \text{sign}(\bar{b}_3) \sqrt[3]{|\bar{b}_3|} \end{array} \right)$$



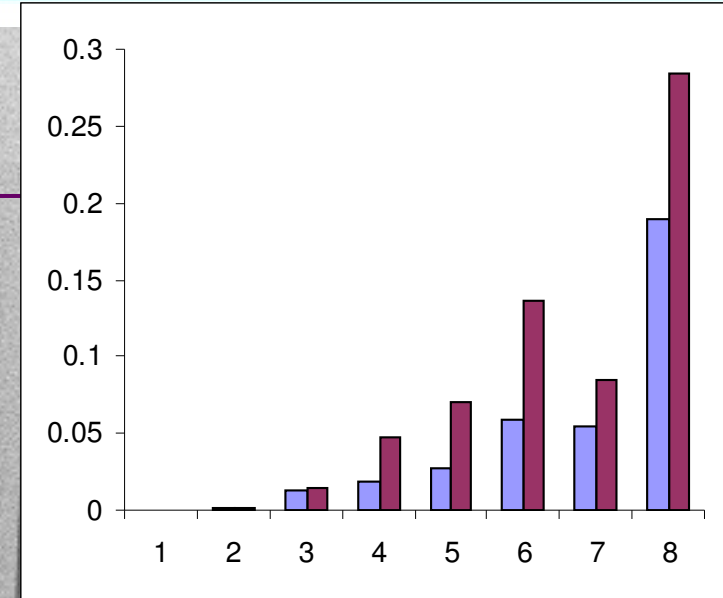
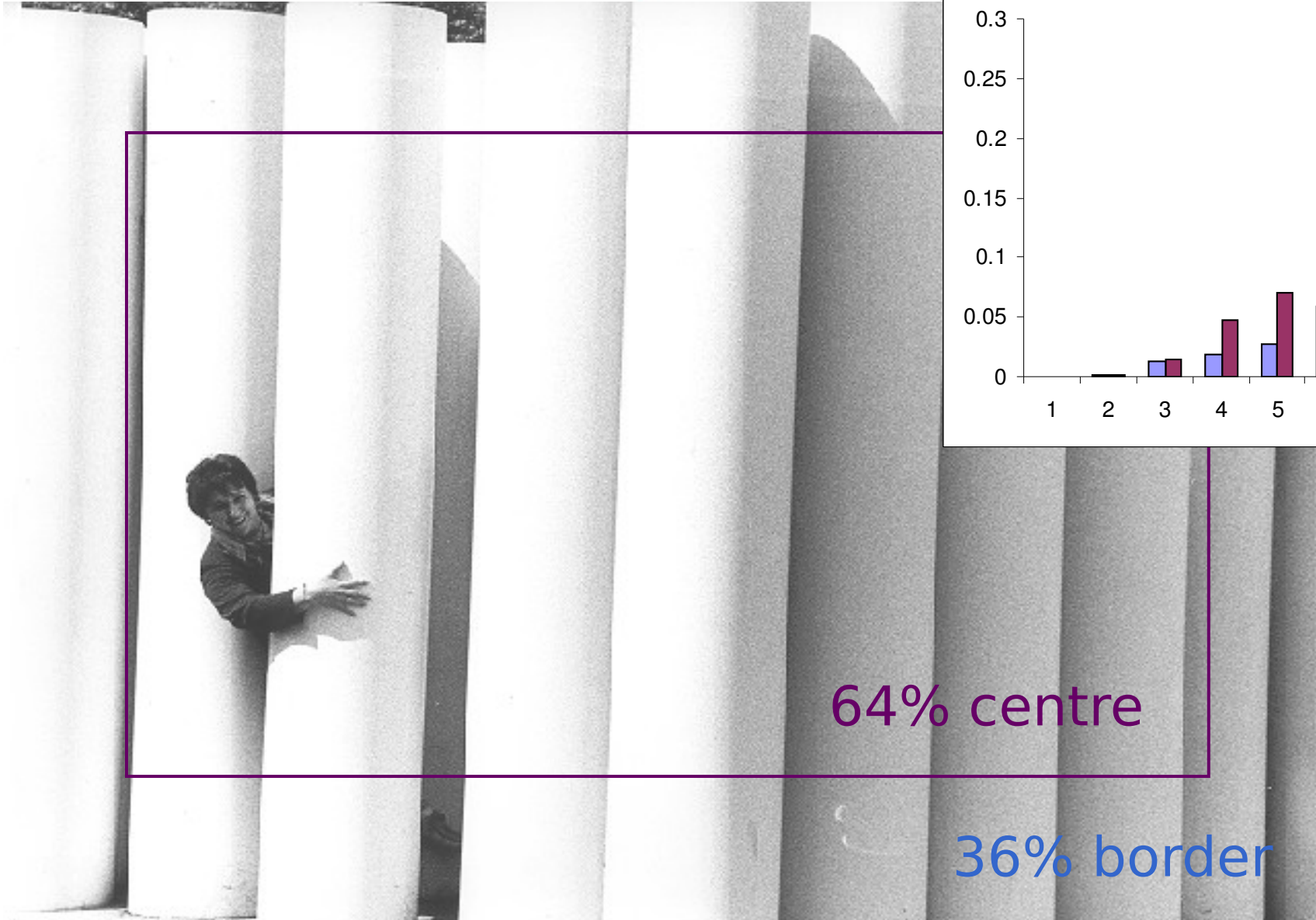
Global vs local



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



Localisation

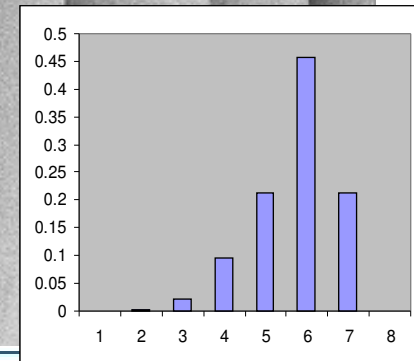
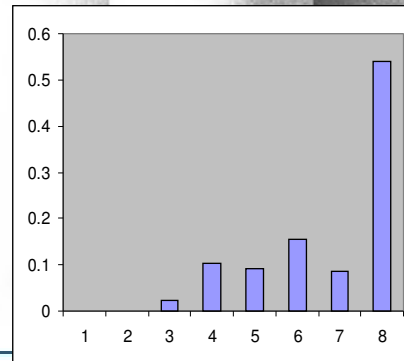
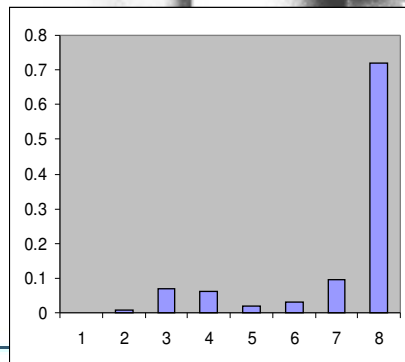
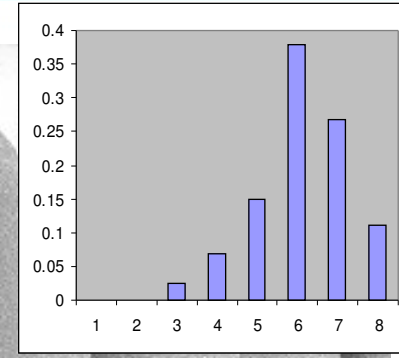
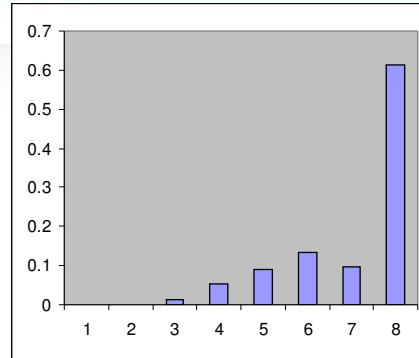
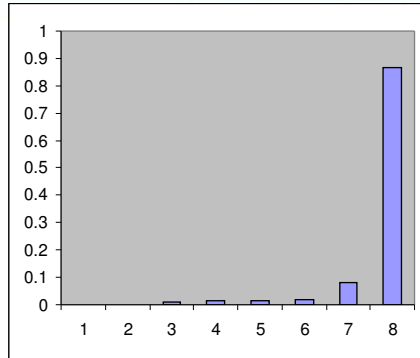


64% centre

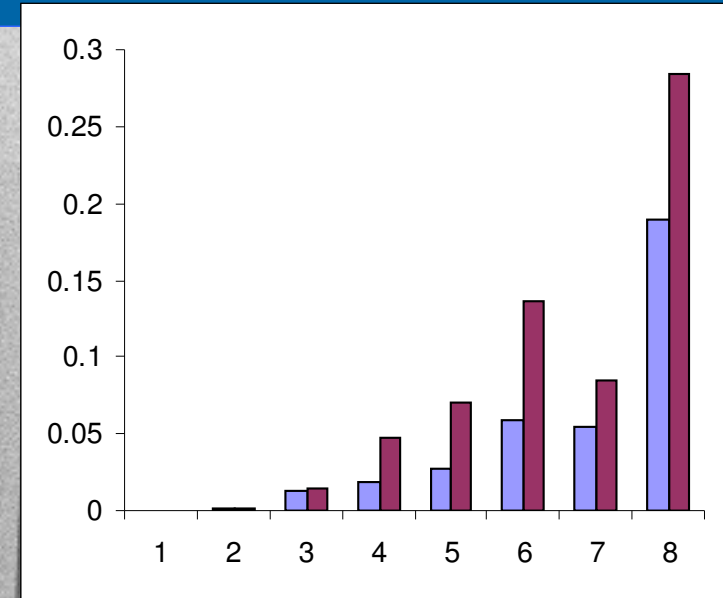
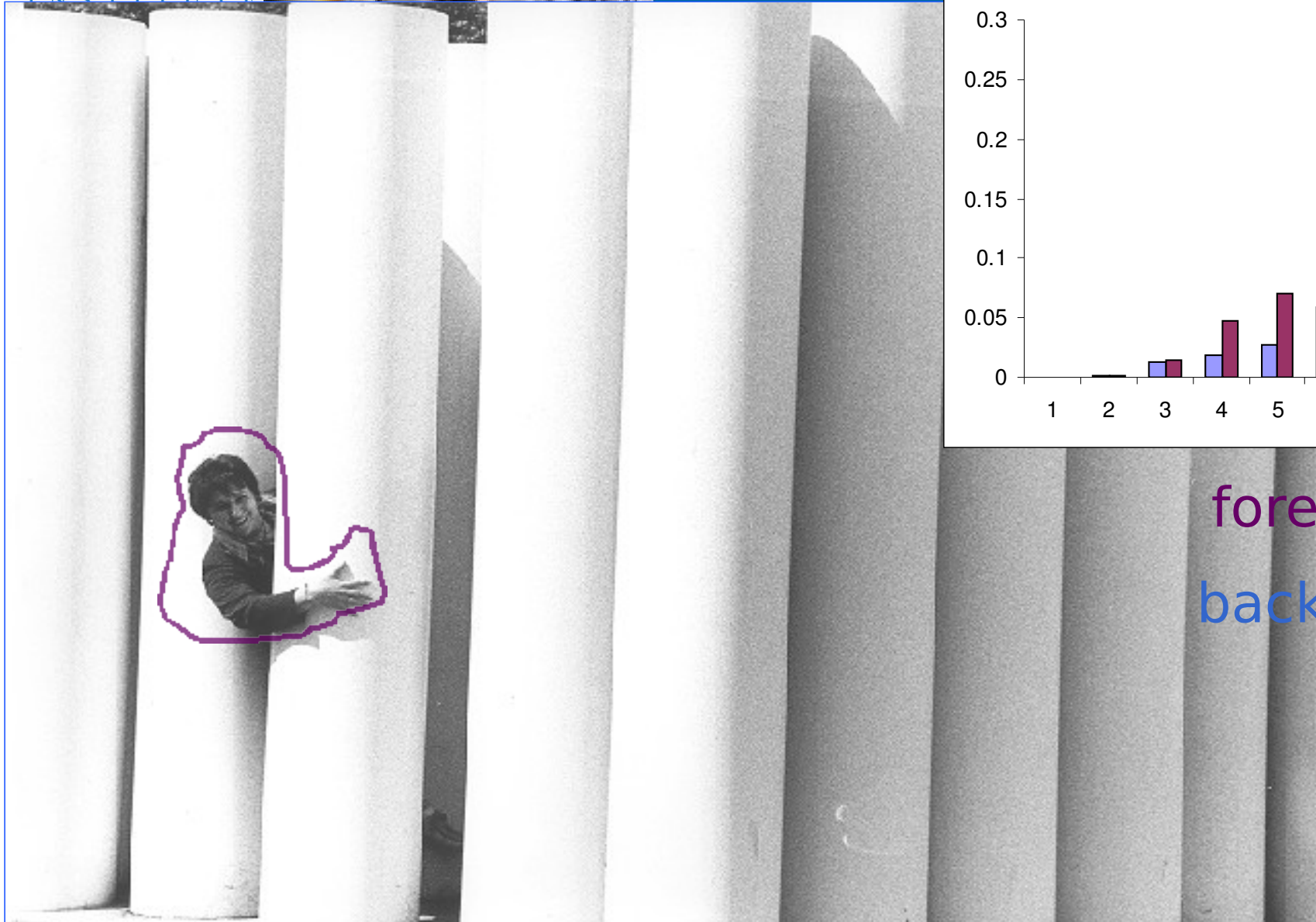
36% border



Tiled Histograms



Segmentation



foreground
background

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



“Bag of Words”



- <http://192.168.1.5:8080/uBase>
- Find an example query image that works well
- Find an example query image that doesn't work
- Try changing the features weights, can you improve the results?



- Anticipation Trailer
- Segmentation Equations



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

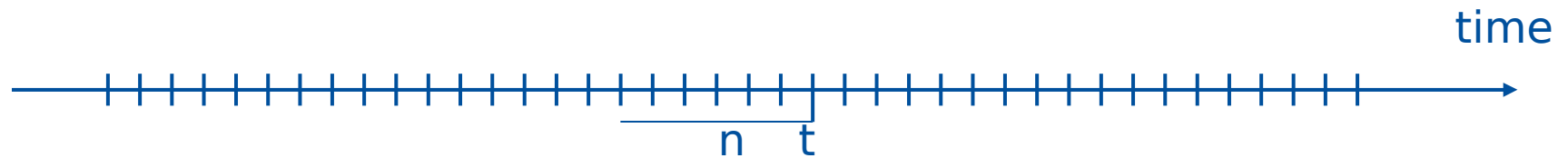
- MPEG provides some motion vectors



[Vlad Tanasescu: Anticipation, SCiFi trailer]



At time t define distance $d_n(t)$



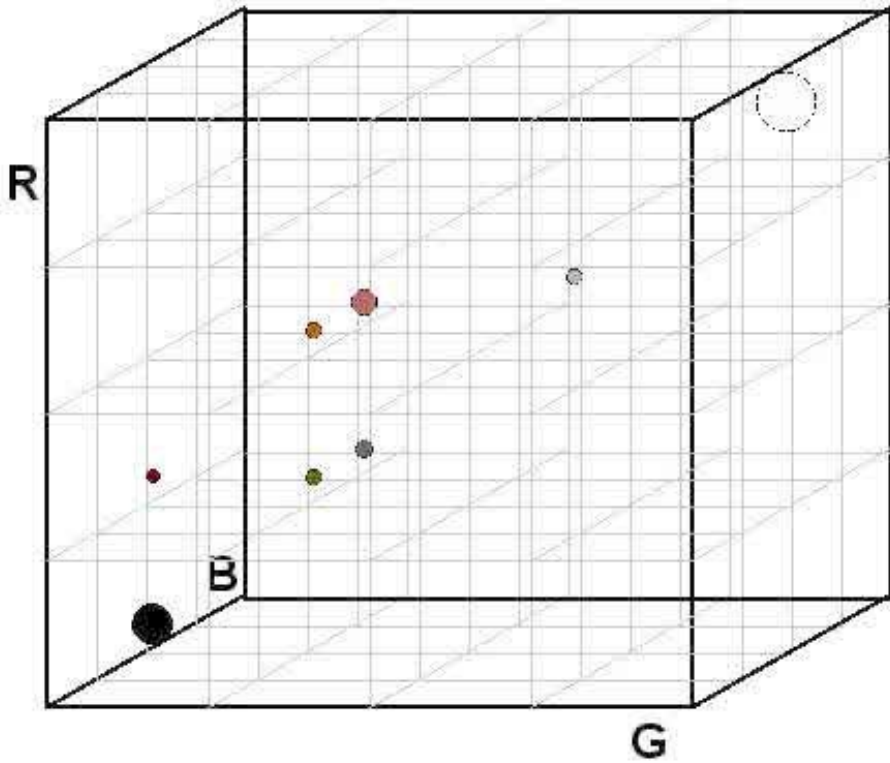
- compare frames $t-n+i$ and $t+i$ ($i=0, \dots, n-1$)
- average their respective distances over i

Peak in $d_n(t)$ detected if

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

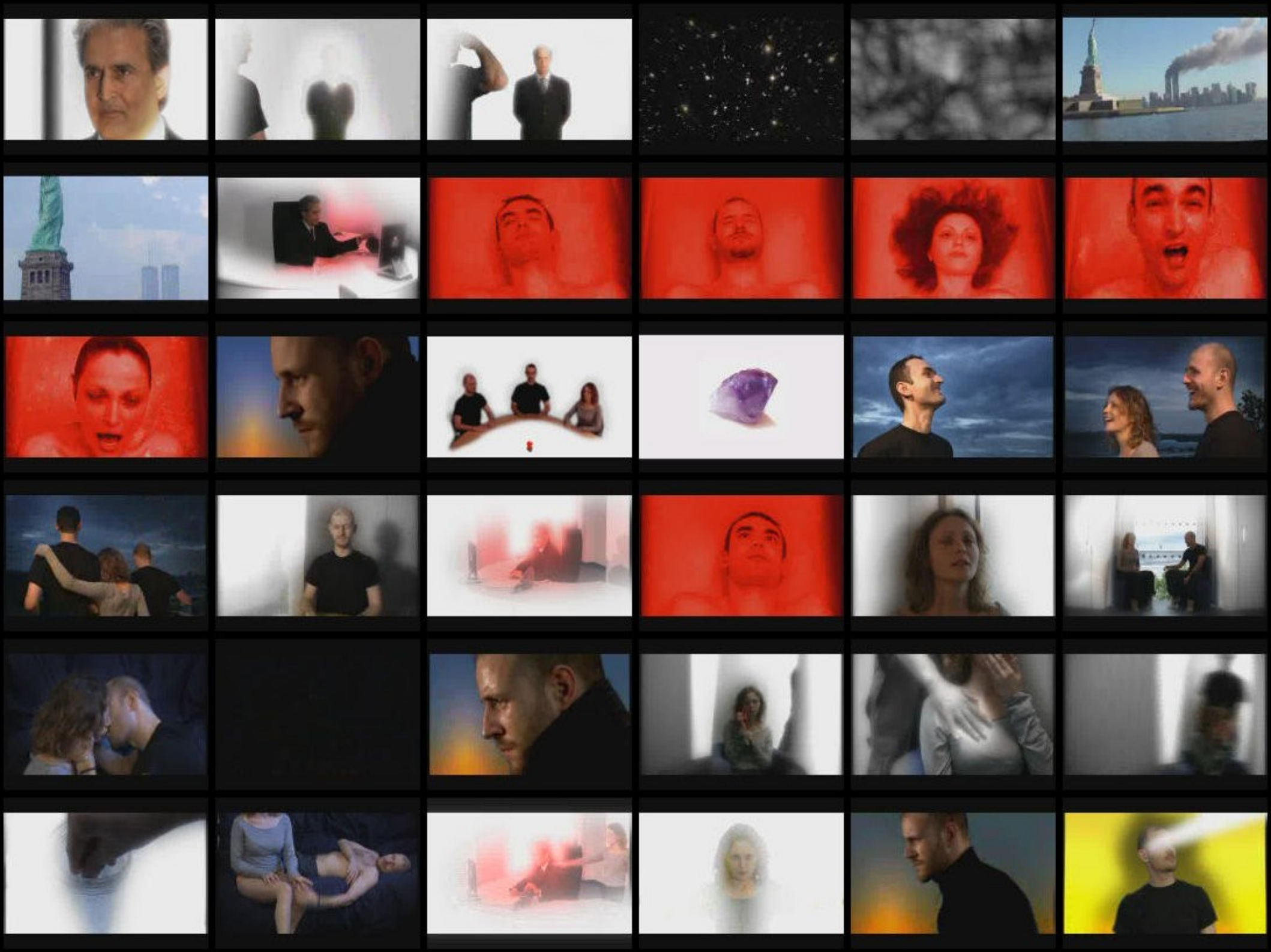
$d_n(t) > d_n(s)$ for all neighbouring s

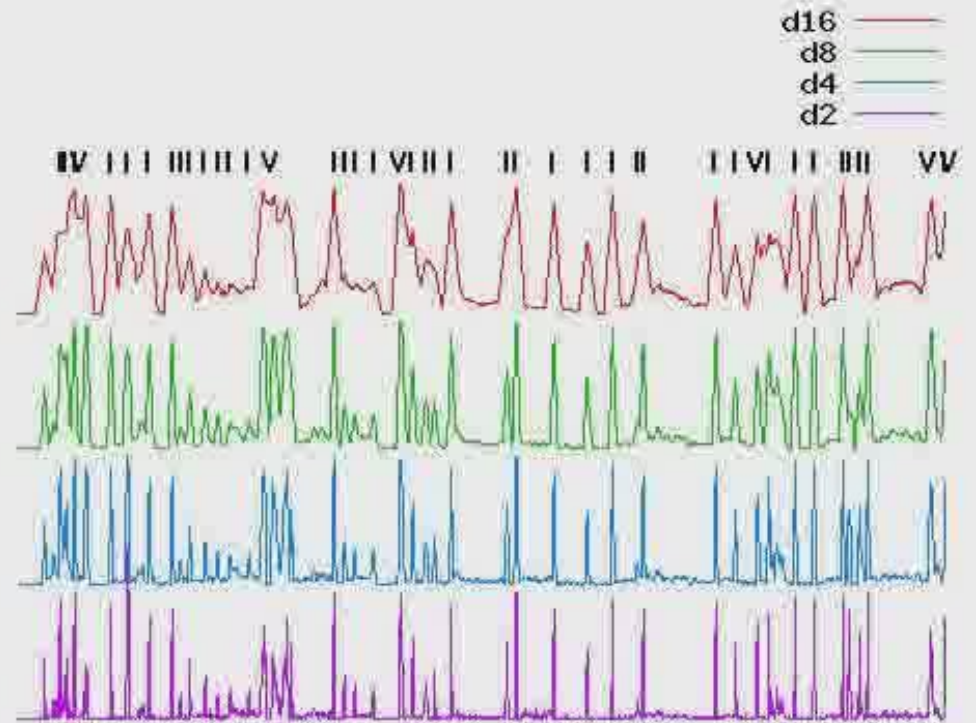
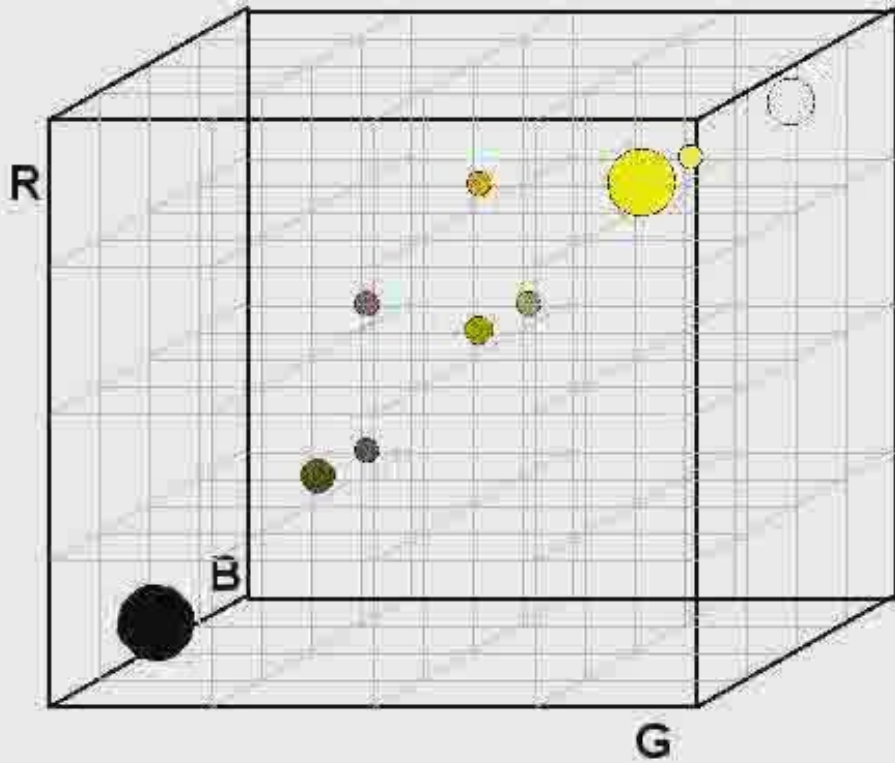
Shot = near-coincident peaks of d_{16} and d_8



d16 —
d8 —
d4 —
d2 —

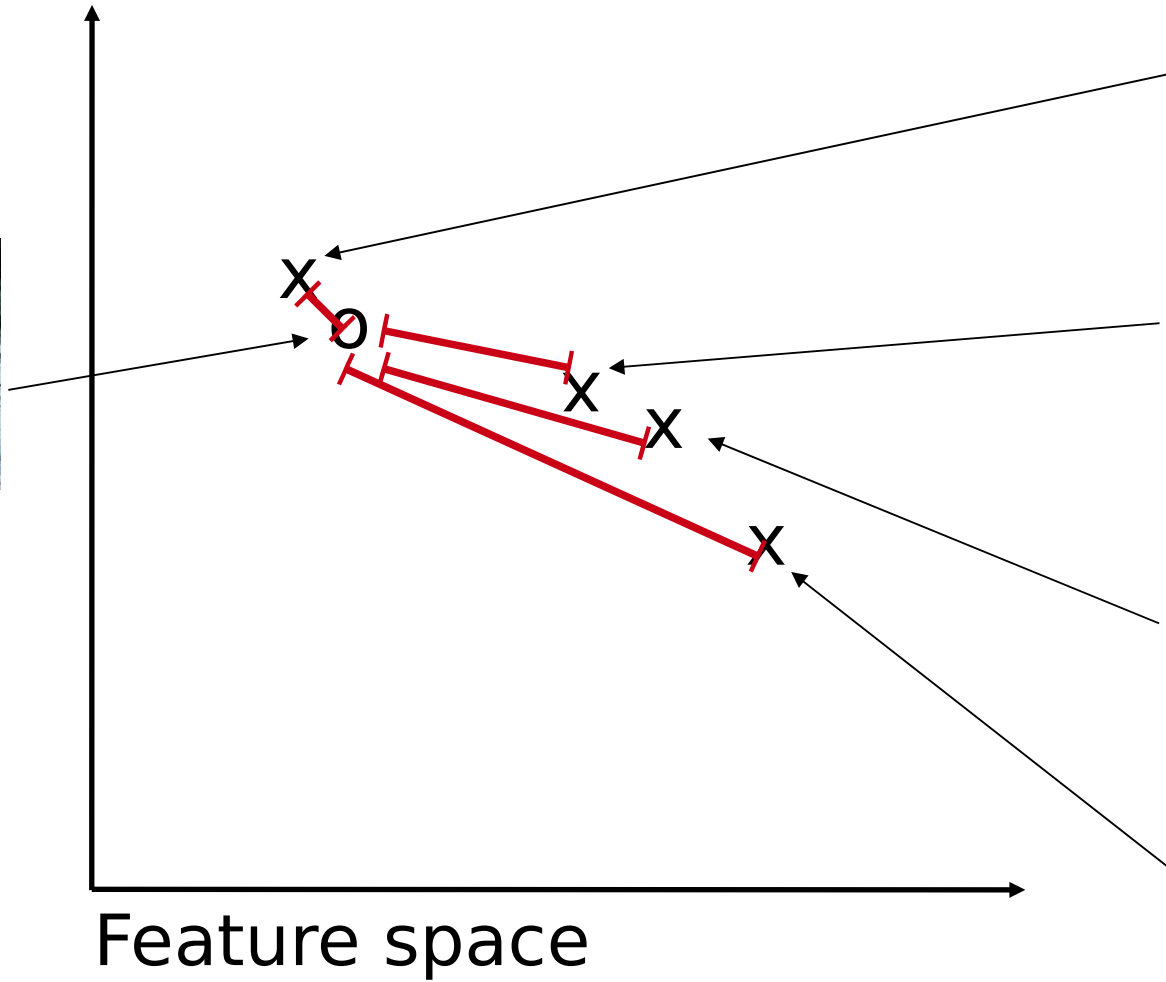
11
12
13
14
15







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

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



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

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

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

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



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

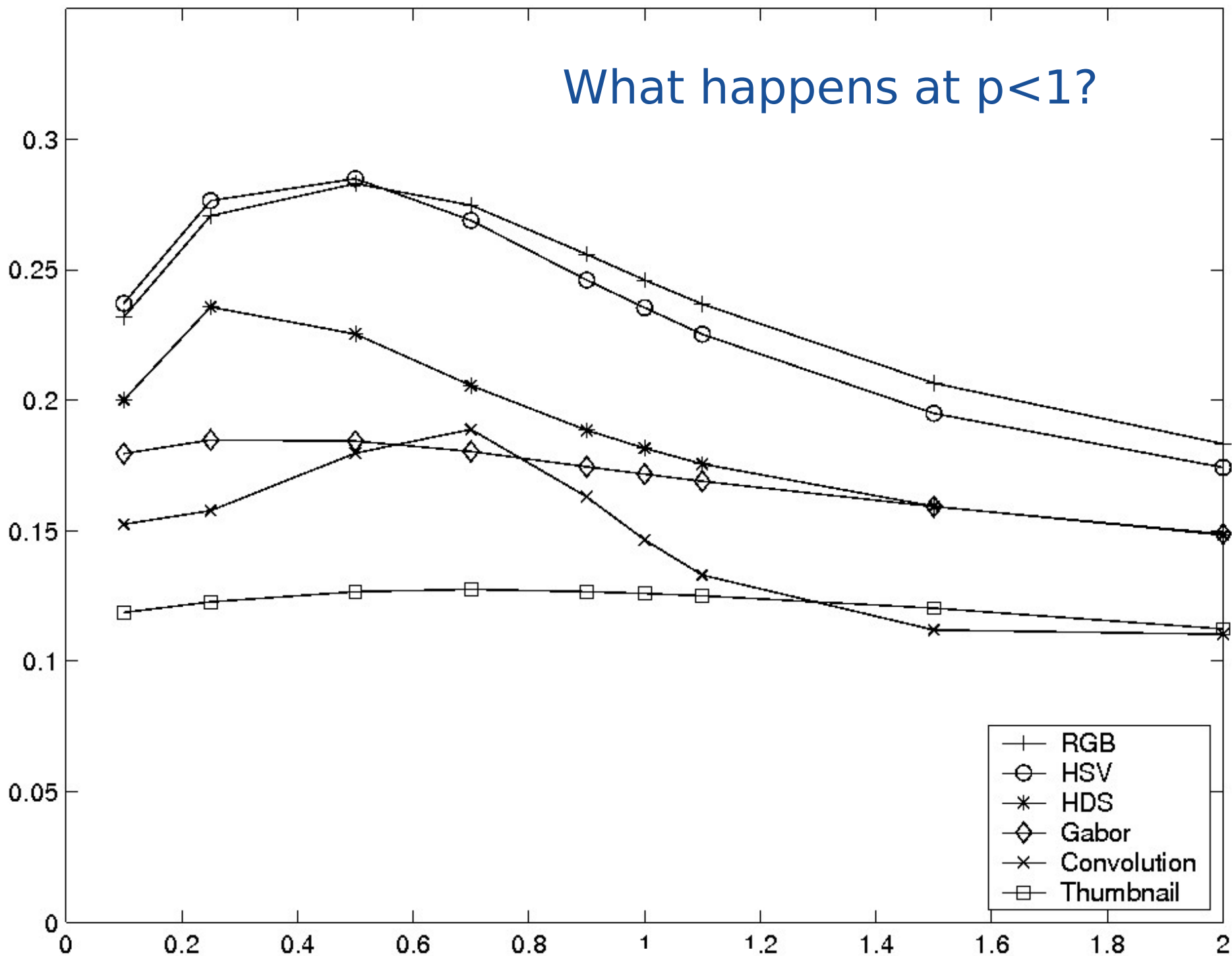
$$d_2(\mathbf{v}, \mathbf{w}) = \sqrt{\sum_i |v_i - w_i|^2}$$

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

$$d_\infty(\mathbf{v}, \mathbf{w}) = \max_i |v_i - w_i|$$

What happens at $p < 1$?

Mean average precision



[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)



speed vs flexibility vs precision

Process:

1. best abstracted representation of your media
2. best method for calculating difference/similarity
3. implement efficiently, considering responsiveness and scalability



Sketch a block diagram showing how you would implement a multimedia information retrieval system for one of these scenarios:

1. Browsing wallpaper patterns in a home decorator store
2. Finding “interesting” photos in a personal collection of holiday snaps
3. Managing industrial design pattern templates for a manufacturing company

Think about:

what types of features you might use
what would the query be
the user interface



For example

“Where is the big pineapple?”



Specific (“known item”)

“Family group photo taken last Christmas”

“The song I heard at the restaurant yesterday”

General

“Family vacation pics at Surfers – like this one”

“Music to go with my vacation photo slide show”



- 1 What is multimedia information retrieval?
 - 1.1 Information retrieval
 - 1.2 Multimedia
 - 1.3 Semantic Gap?
 - 1.4 Challenges of automated multimedia indexing
- 2 Basic multimedia search technologies**
 - 2.1 Meta-data driven retrieval
 - 2.2 Piggy-back text retrieval
 - 2.3 Automated annotation
 - 2.4 Fingerprinting
 - 2.5 Content-based retrieval**
 - 2.6 Implementation Issues**
- 3 Evaluation of MIR Systems
- 4 Added value