

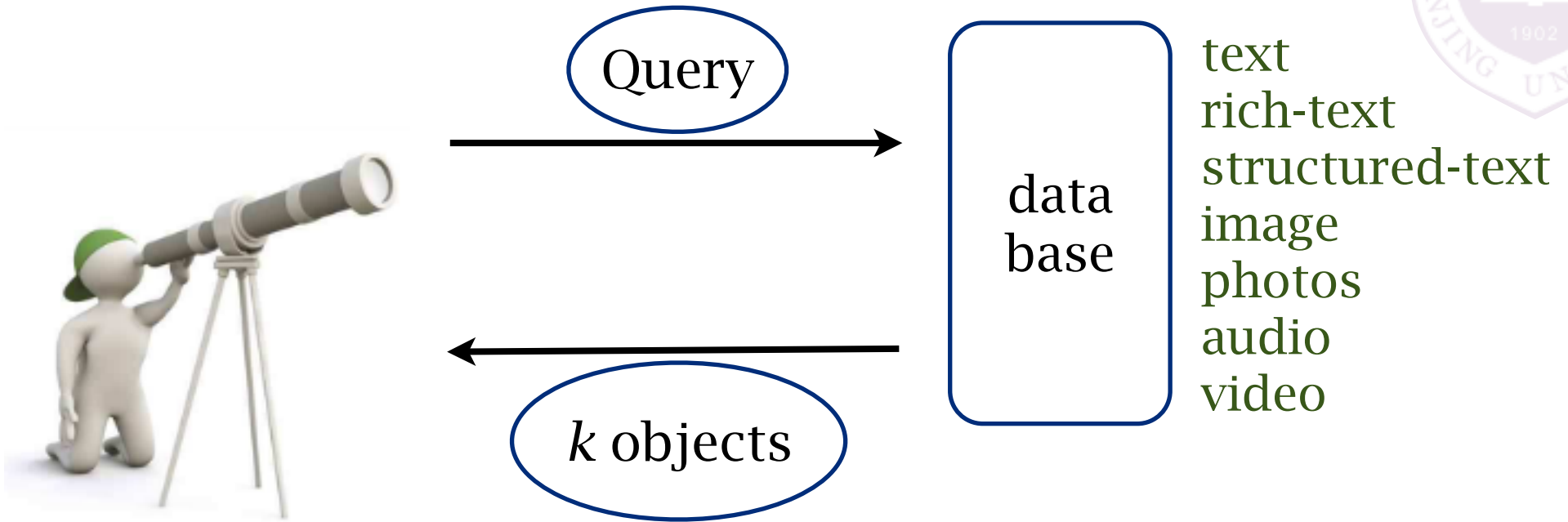


Lecture 12: Some Applications

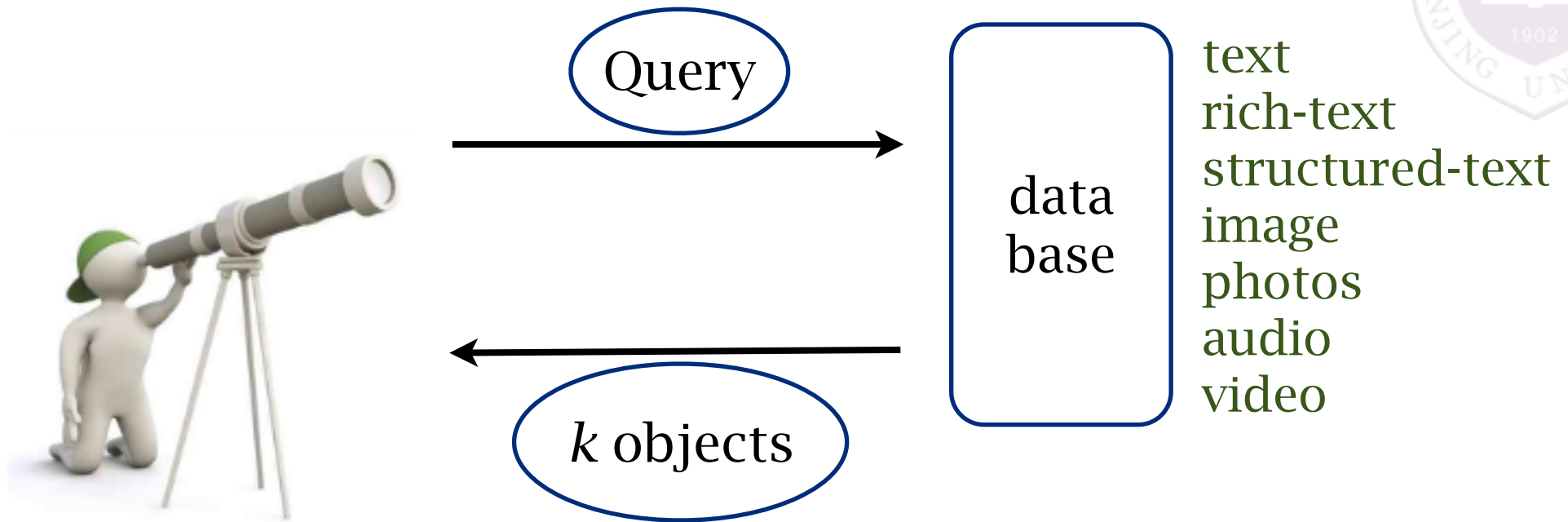
http://cs.nju.edu.cn/yuy/course_dm12.ashx



Information retrieval systems

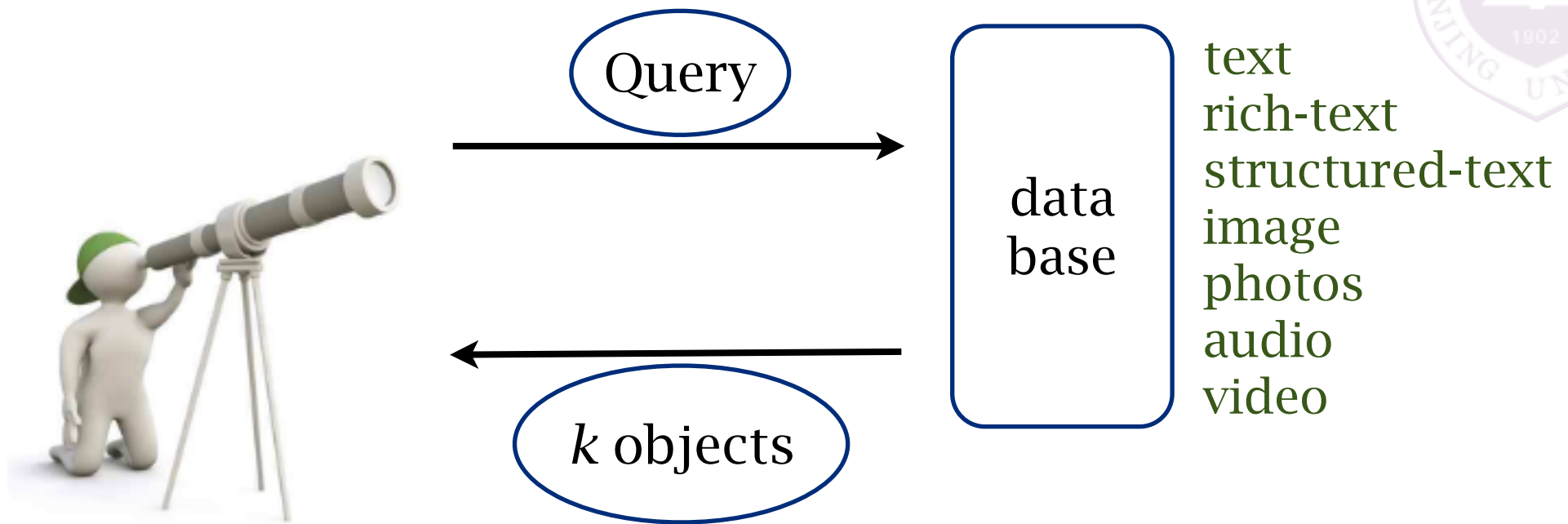


Information retrieval systems



Content-based information retrieval:
for objects with rich semantics
find top k objects most similar to the query

Information retrieval systems



Content-based information retrieval:
for objects with rich semantics
find top k objects most similar to the query

- ▶ searching historical records of the Dow Jones index for past occurrences of a particular time series pattern
- ▶ searching a database of satellite images for any images which contain evidence of recent volcano eruptions in Central America
- ▶ searching the Internet for online documents that provide reviews of restaurants in Helsinki

Evaluation

how good is an retrieval system?



unlike classification where labels are given

Evaluation

how good is an retrieval system?



QUERY



N/R

R

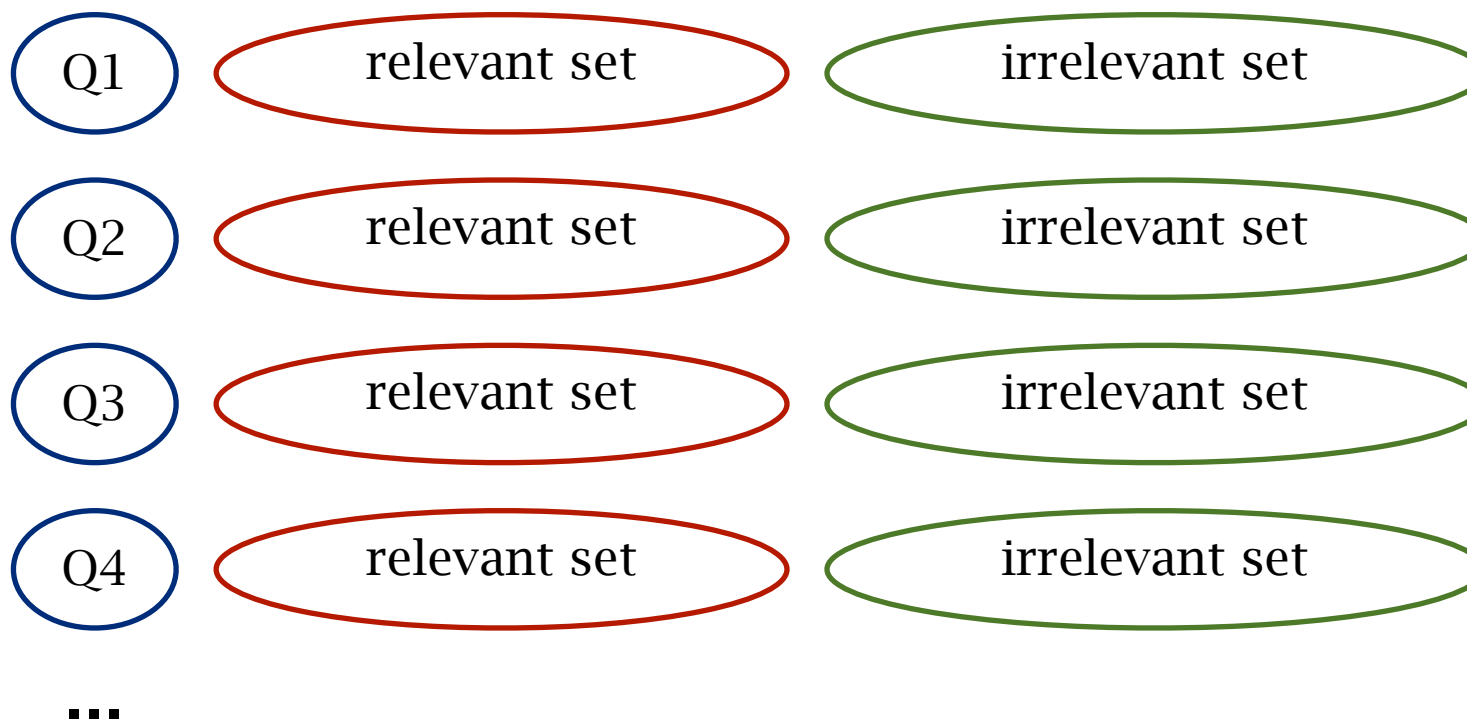


for a particular query, objects can be categorized into “relevant” and “irrelevant”

Evaluation



a set queries and pre-labeled relevant/
irrelevant objects



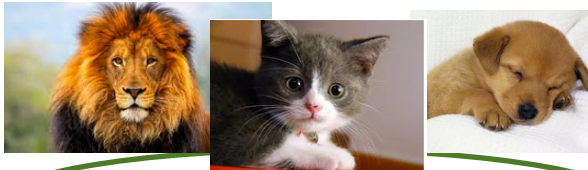
Configurable output



Q



relevant set



irrelevant set



output:



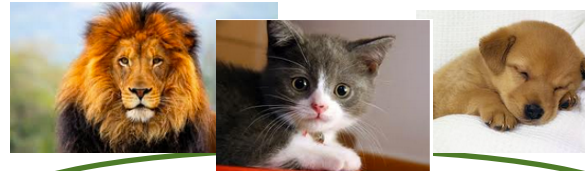
Configurable output



Q



relevant set



irrelevant set



output:



$k=1$

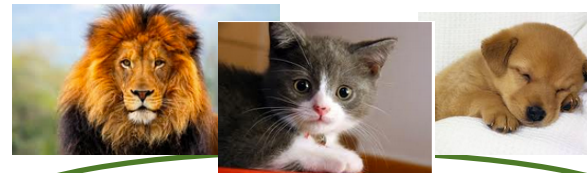
Configurable output



Q



relevant set



irrelevant set



output:



$k=1$

$k=2$

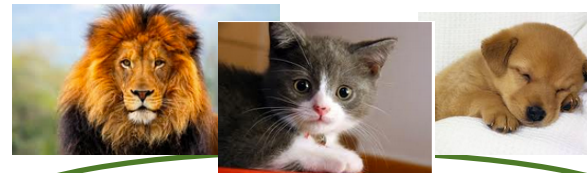
Configurable output



Q



relevant set



irrelevant set



output:



$k=1$



$k=2$



$k=\max$

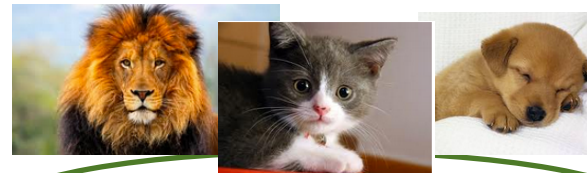
Configurable output



Q



relevant set



irrelevant set



output:



$k=1$



$k=2$



$k=\max$

usually a retrieval system evaluates all objects and rank them according to the similarity

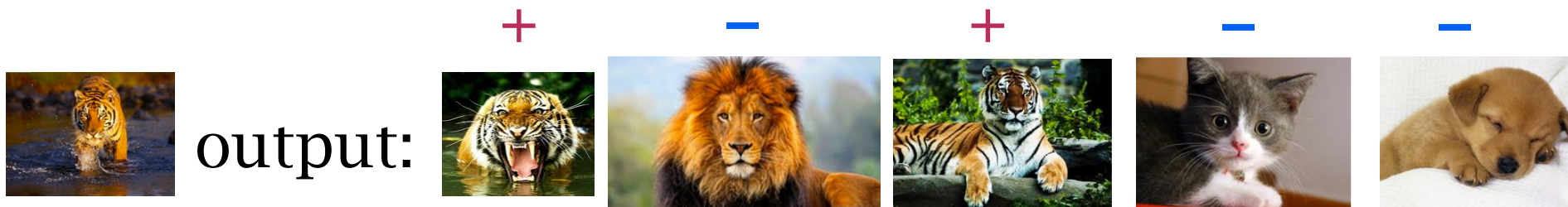
classification error?

Precision and recall



Precision: relevant outputs / all outputs

Recall: relevant outputs / all relevant objects



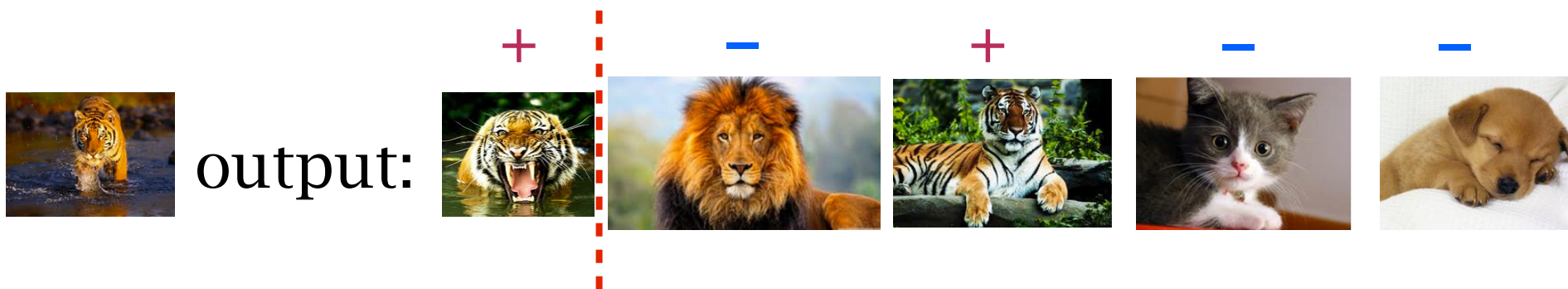
should be averaged over all test queries

Precision and recall



Precision: relevant outputs / all outputs

Recall: relevant outputs / all relevant objects



$k=1$

P: 1

R: 0.5

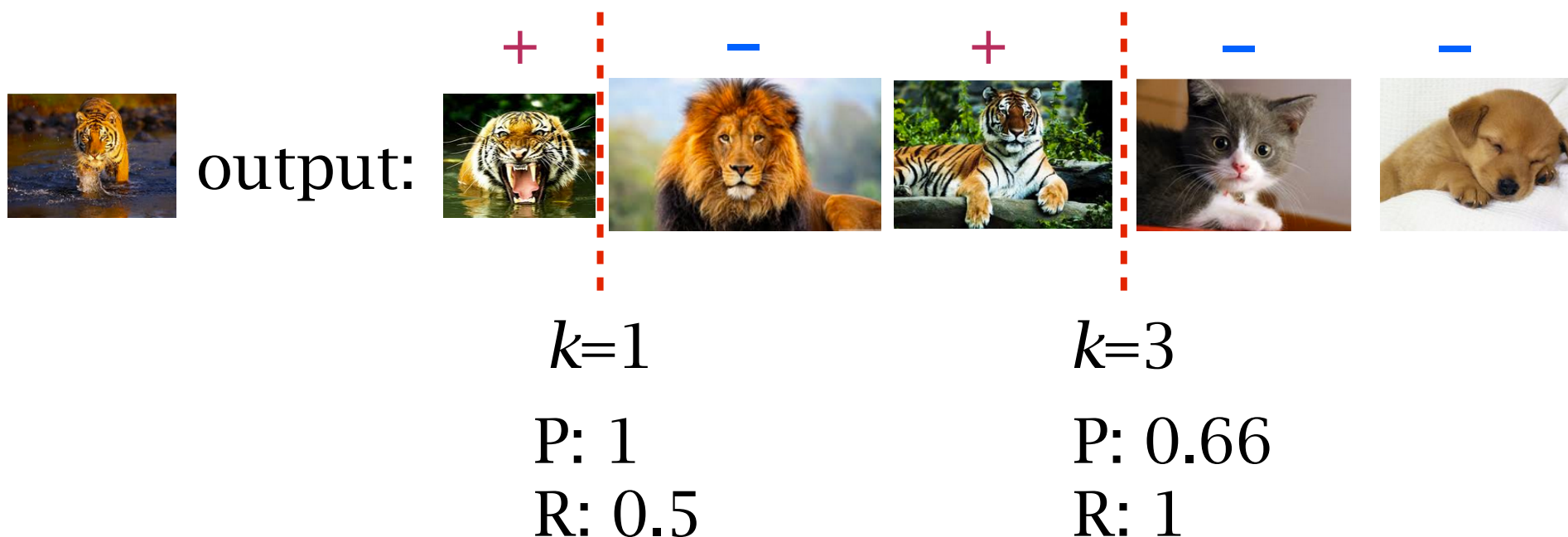
should be averaged over all test queries

Precision and recall



Precision: relevant outputs / all outputs

Recall: relevant outputs / all relevant objects



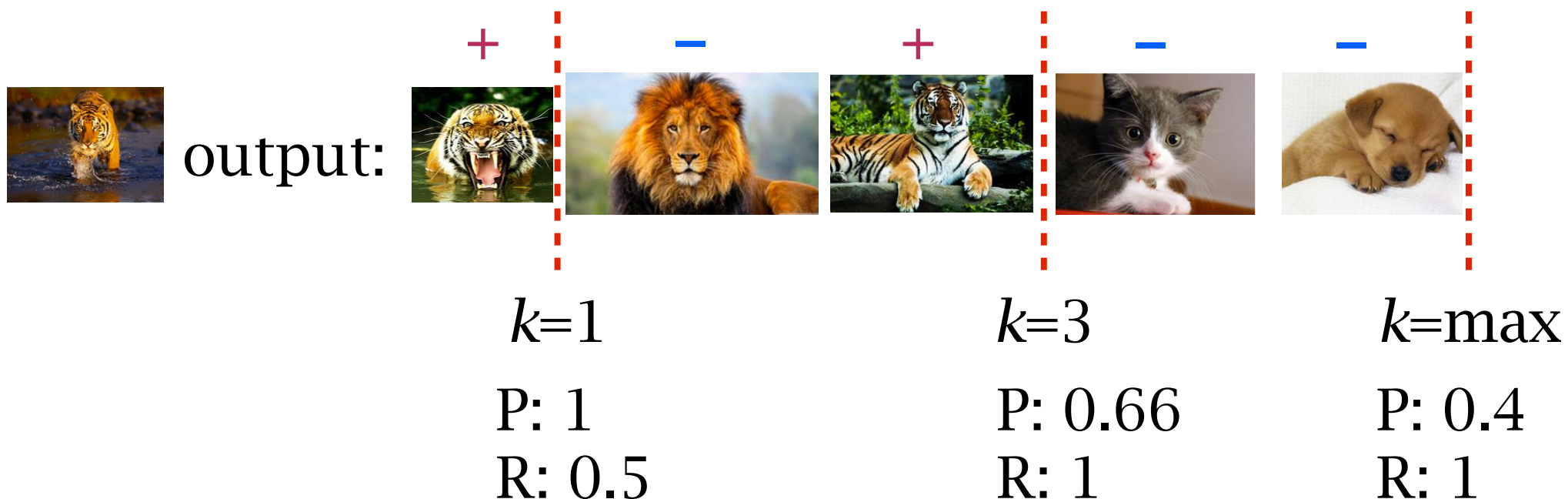
should be averaged over all test queries

Precision and recall



Precision: relevant outputs / all outputs

Recall: relevant outputs / all relevant objects

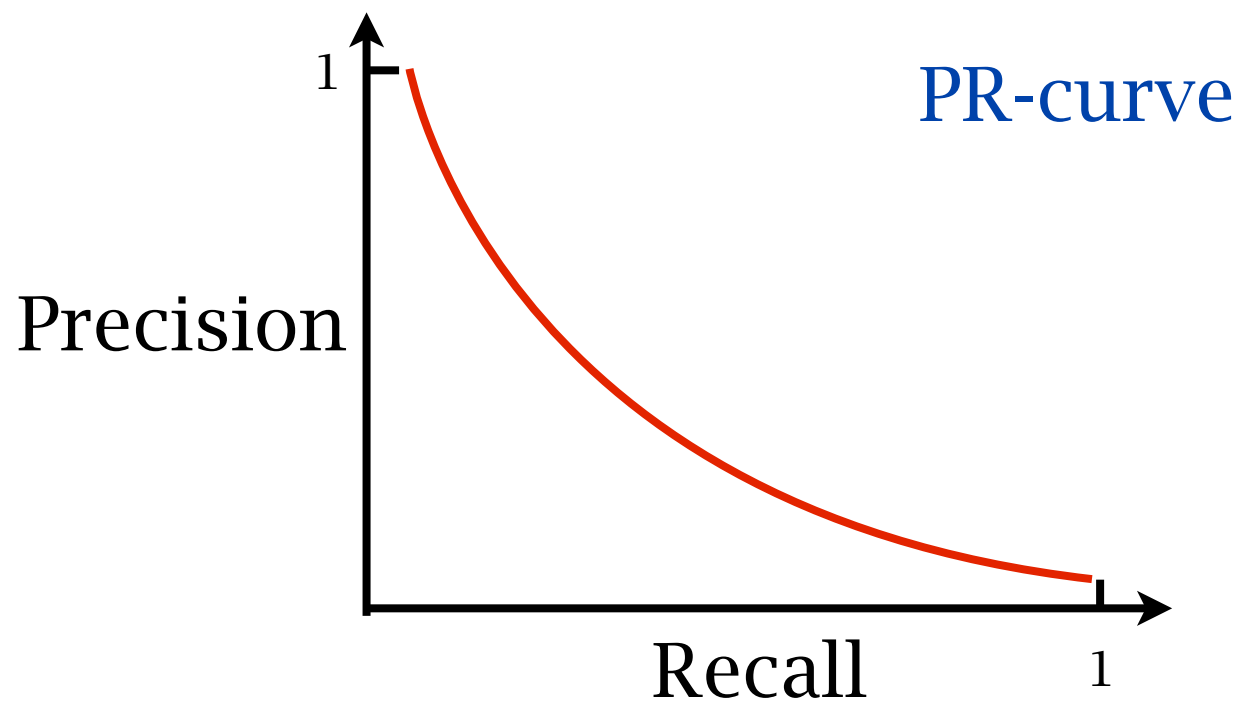


should be averaged over all test queries

Precision and recall



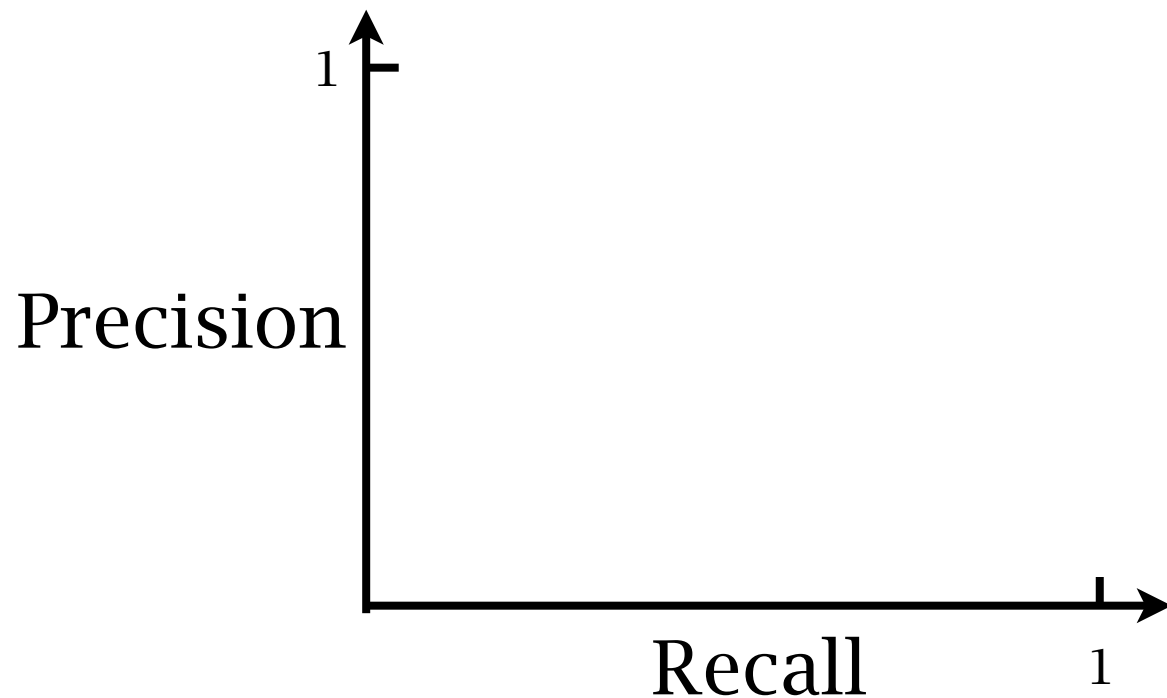
Enumerate all k to produce a set of (P,R) pairs



Precision and recall



Compare retrieval systems

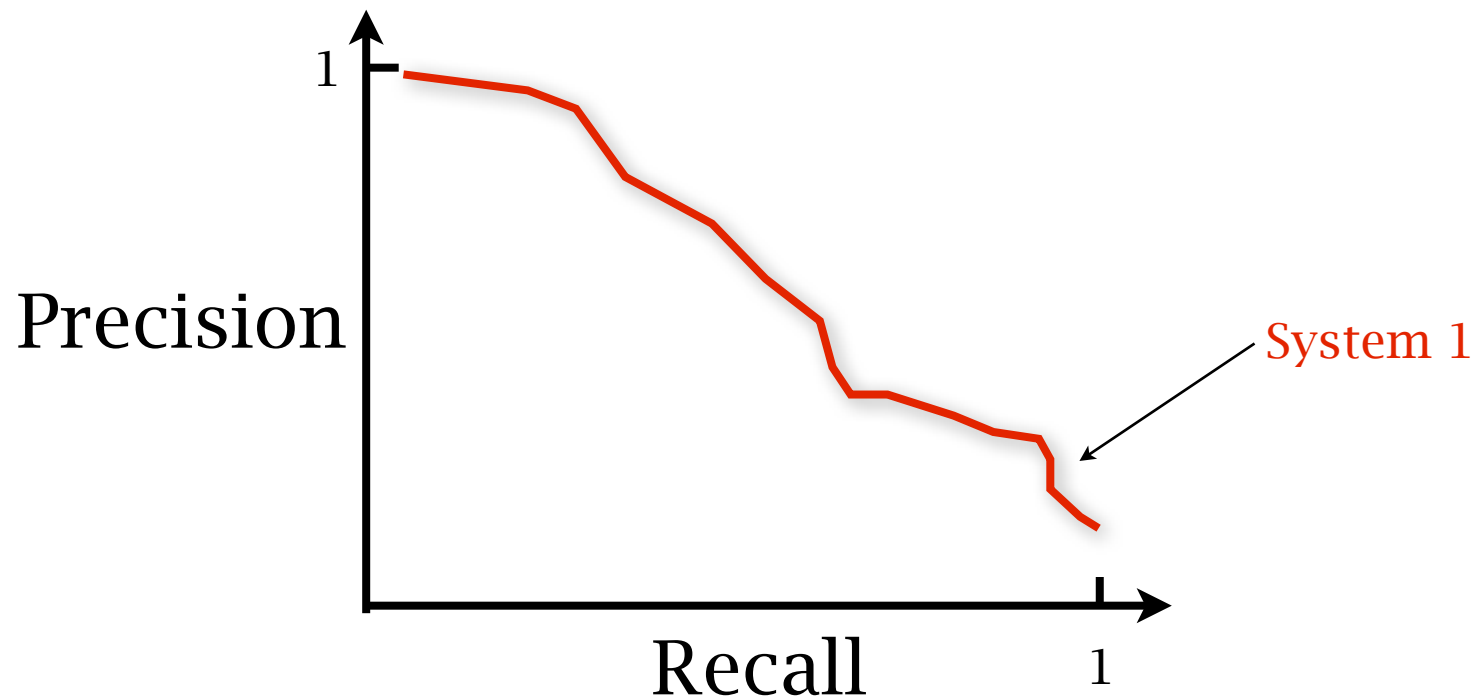


System 1 is better than System 2
System 1 v.s. System 3?

Precision and recall



Compare retrieval systems

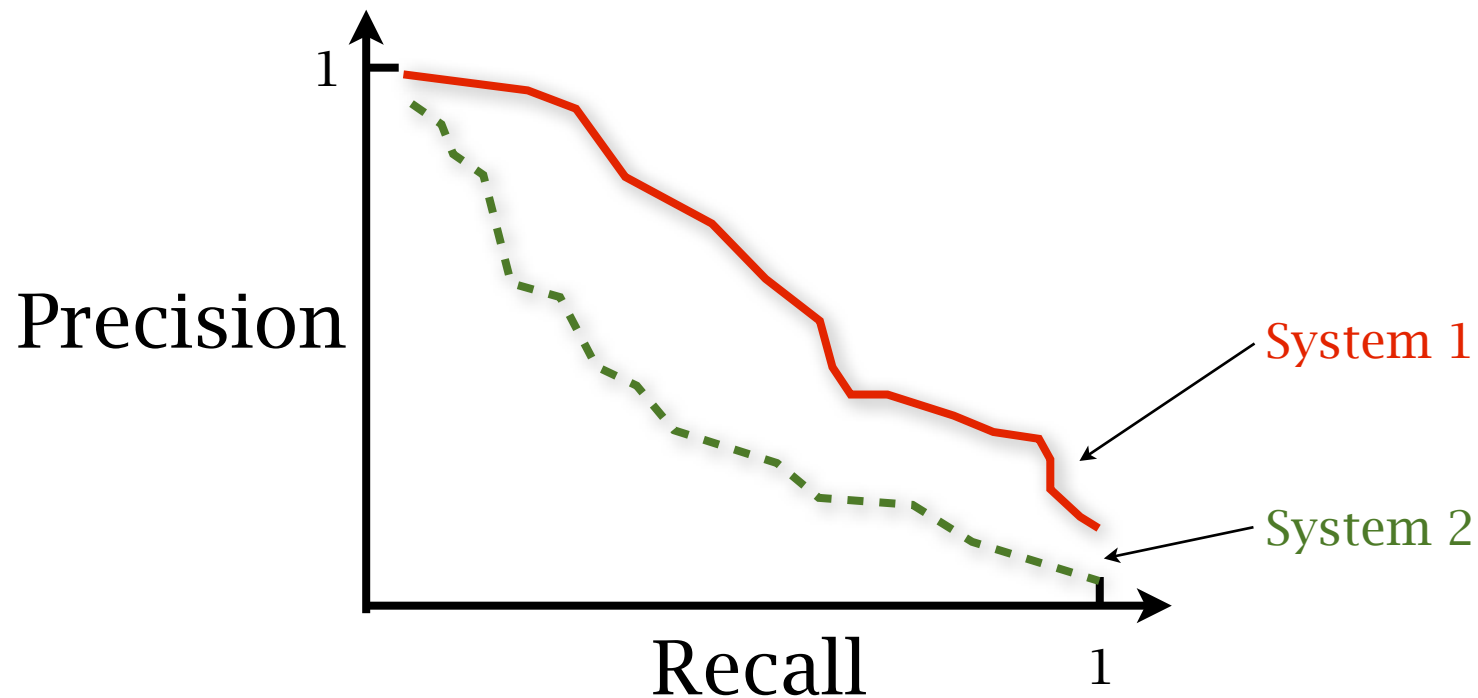


System 1 is better than System 2
System 1 v.s. System 3?

Precision and recall



Compare retrieval systems

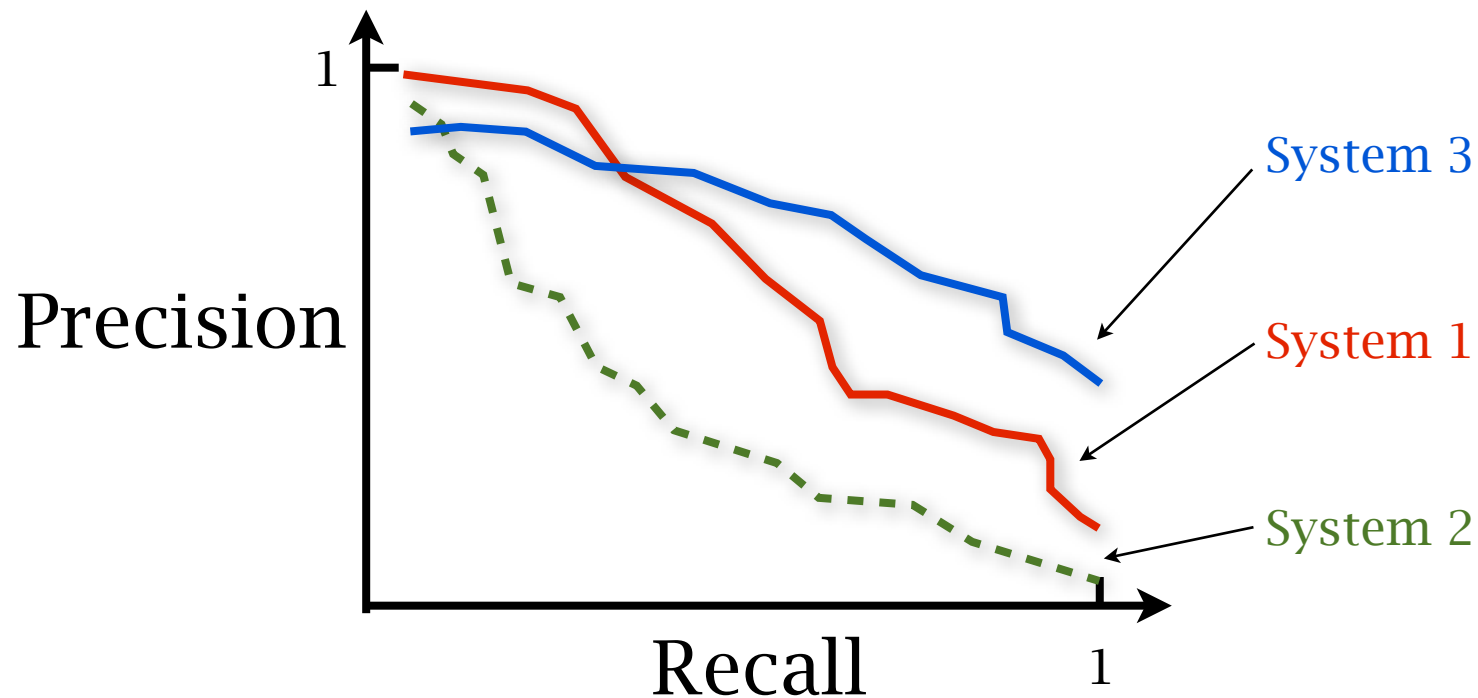


System 1 is better than System 2
System 1 v.s. System 3?

Precision and recall



Compare retrieval systems



System 1 is better than System 2
System 1 v.s. System 3?

Precision and recall



Compare retrieval systems

Precision/recall at a fixed k

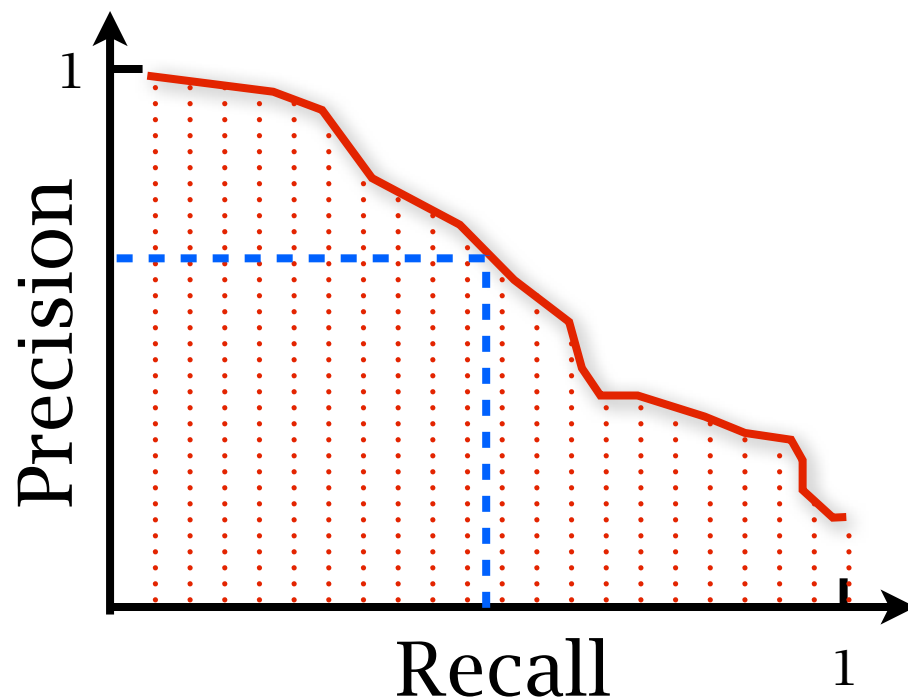
Area under PR-Curve:

Position where $P=R$

F-measure:

for arbitrary cut-point

$$F = \frac{1}{\frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right)}$$



Harmonic mean: the probability of the binary random variable whose expectation equals the average expectation of two binary random variables

Precision v.s. recall



application dependent

Criminal face retrieval: high recall



output:



Recommendation in social network: high precision



output:

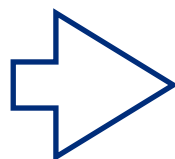


Text retrieval system



Retrieval from a text database

苏州



hit of words:
too many candidates

汪洋：穿中山装离开带走广东文化中国特色

12月18日，汪洋回顾在广东五年的工作时深情地说：五年前我是穿着西装来与大家首次见面的，今天我将穿着中山装离开，带走的是广东文化，中国特色！我在这里受到的熏陶，将使我终身受益。

胡春华：期盼汪洋到中央工作以后继续关注广东

中新网广州12月18日电 (索有为 奚婉婷 岳宗) 18日下午，中共广东省委召开全省领导干部会议，中央政治局委员、中央书记处书记、中央组织部部长赵乐际受中央委派，在会上宣布了中央关于广东省委书记调整的决定：汪洋不再兼任广东省委书记、常委、

苏州老子雕塑卖萌 背对裤衩楼“吐舌扮鬼脸” (图)

老子雕像继裸女座椅雕塑之后，苏州金鸡湖畔的一尊老子雕塑再度引发争议。道家创始人老子以朴素辩证法思想和无为而治的政治主张，润泽千年，成为中华文化不可或缺的瑰宝。然而，就是这样一个万民敬仰的圣贤，在这尊雕塑上却眼睛紧闭，舌头伸出，

我国五省区党委书记职务调整

日前，中共中央决定：汪洋同志不再兼任广东省委书记、常委、委员职务；胡春华同志兼任广东省委常委、常委、书记，不再兼任内蒙古自治区党委书记、常委、委员职务；王君同志任内蒙古自治区党委书记、常委、书记；王儒林同志任吉林省委书记；赵正永同志任陕西省委书记；夏宝龙同志任浙江省委书记。

Vector representation



Dictionary: (1,苏州) (2. 南京) (3. 孔子) (4. 老子)...

苏州老子雕塑卖萌 背对裤衩楼 “吐舌扮鬼脸” (图)

老子雕像继裸女座椅雕塑之后，苏州金鸡湖畔的一尊老子雕塑再度引发争议。道家创始人老子以朴素辩证法思想和无为而治的政治主张，润泽千年，成为中华文化不可或缺的瑰宝。然而，就是这样一个万民敬仰的圣贤，在这尊雕塑上却眼睛紧闭，舌头伸出，露出嘴中一个大门牙，作出一副“龇牙吐舌”的怪状，雷倒了许多路过的市民和游客。昨日，这尊老子“龇牙吐舌”的雕塑在微博上被众多网友转发，一度引起广泛关注。



Vector representation



Document-term frequency matrix

	t1	t2	t3	t4	t5
D1	24	21	9	0	0
D2	32	10	5	0	3
D3	12	16	5	0	0
D4	6	7	2	0	0
D5	43	31	20	0	3
D6	2	0	0	18	7
D7	0	0	1	32	12
D8	3	0	0	22	4
D9	1	0	0	34	27

cosine similarity:

$$\cos(q, x) = \frac{q^T x}{\|q\| \cdot \|x\|}$$

Query:

(0,0,1,1,0)

features are important to the performance of a retrieval system

Vector representation



Inverse Document Frequency

$$IDF(t) = \log \left(\frac{\text{Number of total documents}}{\text{Number of documents containing } t} \right)$$

Document-term frequency (TF)

	t1	t2	t3	t4	t5	t6
D1	24	21	9	0	0	3
D2	32	10	5	0	3	0
D3	12	16	5	0	0	0
D4	6	7	2	0	0	0
D5	43	31	20	0	3	0
D6	2	0	0	18	7	16
D7	0	0	1	32	12	0
D8	3	0	0	22	4	2
D9	1	0	0	34	27	25
D10	6	0	0	17	4	23

$$IDF(t1) = \log \frac{10}{9} = 0.1520$$

$$IDF(t6) = \log_2 \frac{10}{5} = 1$$

multiply

Document-term TF-IDF matrix

	t1	t2	t3	t4	t5	t6
D1	3.7	21	6.6	0	0	3
D2	4.9	10	3.7	0	1.5	0
D3	1.8	16	3.7	0	0	0
D4	0.9	7	1.5	0	0	0
D5	6.5	31	15	0	1.5	0
D6	0.3	0	0	18	3.6	16
D7	0	0	0.7	32	6.2	0
D8	0.5	0	0	22	2.1	2
D9	0.2	0	0	34	14	25
D10	0.9	0	0	17	2.1	23

Vector representation



Many ways to form features

Table 4. Performance results for eight term-weighting methods averaged over 5 collections

Term-weighting methods	Rank of method and ave. precision	CACM 3204 docs 64 queries	CISI 1460 docs 112 queries	CRAN 1397 docs 225 queries	INSPEC 12,684 docs 84 queries	MED 1033 docs 30 queries	Averages for 5 collections
1. Best fully weighted ($tf \cdot nfx$)	Rank P	1 0.3630	14 0.2189	19 0.3841	3 0.2626	19 0.5628	11.2
2. Weighted with inverse frequency f not used for docs ($txc \cdot nfx$)	Rank P	25 0.3252	14 0.2189	7 0.3950	4 0.2626	32 0.5542	16.4
3. Classical $tf \times idf$ No normalization ($tfx \cdot tfx$)	Rank P	29 0.3248	22 0.2166	219 0.2991	45 0.2365	132 0.5177	84.4
4. Best weighted probabilistic ($nxx \cdot bpx$)	Rank P	55 0.3090	208 0.1441	11 0.3899	97 0.2093	60 0.5449	86.2
5. Classical idf without normalization ($bfx \cdot bfx$)	Rank P	143 0.2535	247 0.1410	183 0.3184	160 0.1781	178 0.5062	182
6. Binary independence probabilistic ($bxx \cdot bpx$)	Rank P	166 0.2376	262 0.1233	154 0.3266	195 0.1563	147 0.5116	159
7. Standard weights cosine normalization (original Smart) ($txc \cdot txx$)	Rank P	178 0.2102	173 0.1539	137 0.3408	187 0.1620	246 0.4641	184
8. Coordination level binary vectors ($bxx \cdot bxx$)	Rank P	196 0.1848	284 0.1033	280 0.2414	258 0.0944	281 0.4132	260

[Salton and Buckley, 88]

Image retrieval system

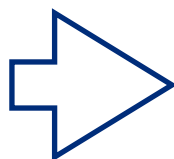


Retrieval from an image database

description: tiger



example:



sketch:



similar system structure as text retrieval system
the hardest part: features

Vector representation



common ingredient:

colors

RGB, HSV, LIB...

texture

Fourier transformation, wavelets

gradients

edges, descriptors

Vector representation



Global features

1. 3-D color feature vector
 - Spatially averaged over the whole image
 - Euclidean distance
2. k-dimensional color histogram
 - bins selected by partition based-based clustering algorithm such as k means
 - k is application dependent
 - Mahanalobis distance using inverse variances
3. 3-D Texture Vector
 - coarseness/scale, directionality, contrast
4. 20-dimensional shape feature based on area, circularity, eccentricity, axis orientation, moments

Vector representation



Local features

bag-of-words

split the images into small pieces
extract a feature vector per piece
clustering to find centers of feature vectors
each image by a vector of frequency of centers



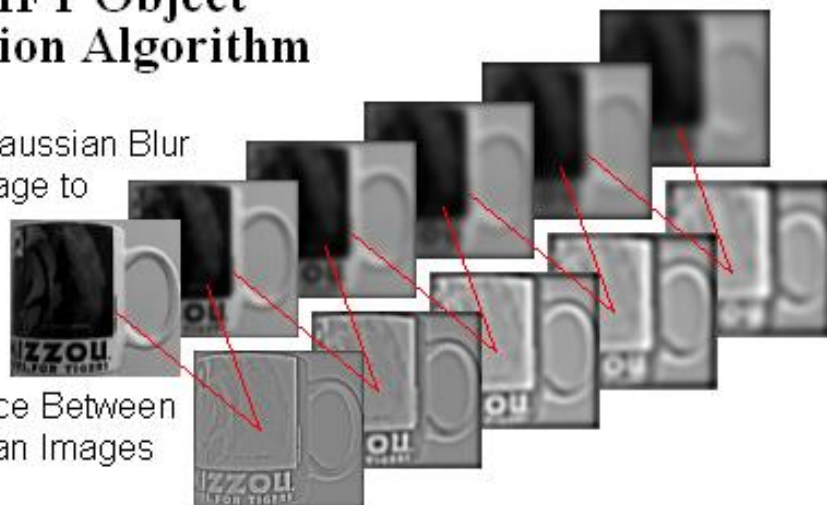
Vector representation



Local features

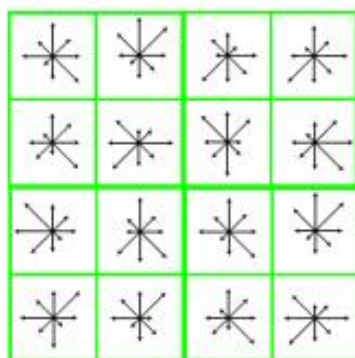
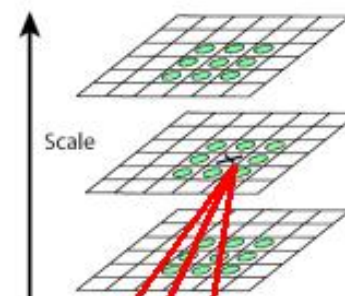
The SIFT Object Recognition Algorithm

Incrementally Gaussian Blur
The Original Image to
Create a Scale
Space

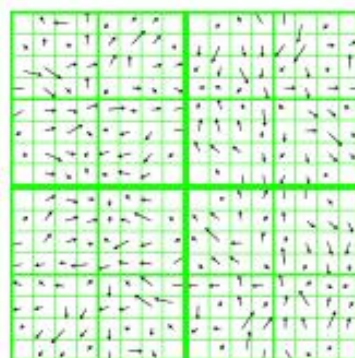


Find the Difference Between
Adjacent Gaussian Images
in Scale Space

Keypoints are Pixels
in Difference Images
That are Larger Than
or Smaller Than all 26
Neighbors



Sixteen Histograms are
Created Using The Gradients.
Using 8 Orientations, This
Makes 128-D Feature Vectors.



The Gradient of Pixels Around
Each Keypoint is Determined
At the Gaussian Scale at Which
It Was Found



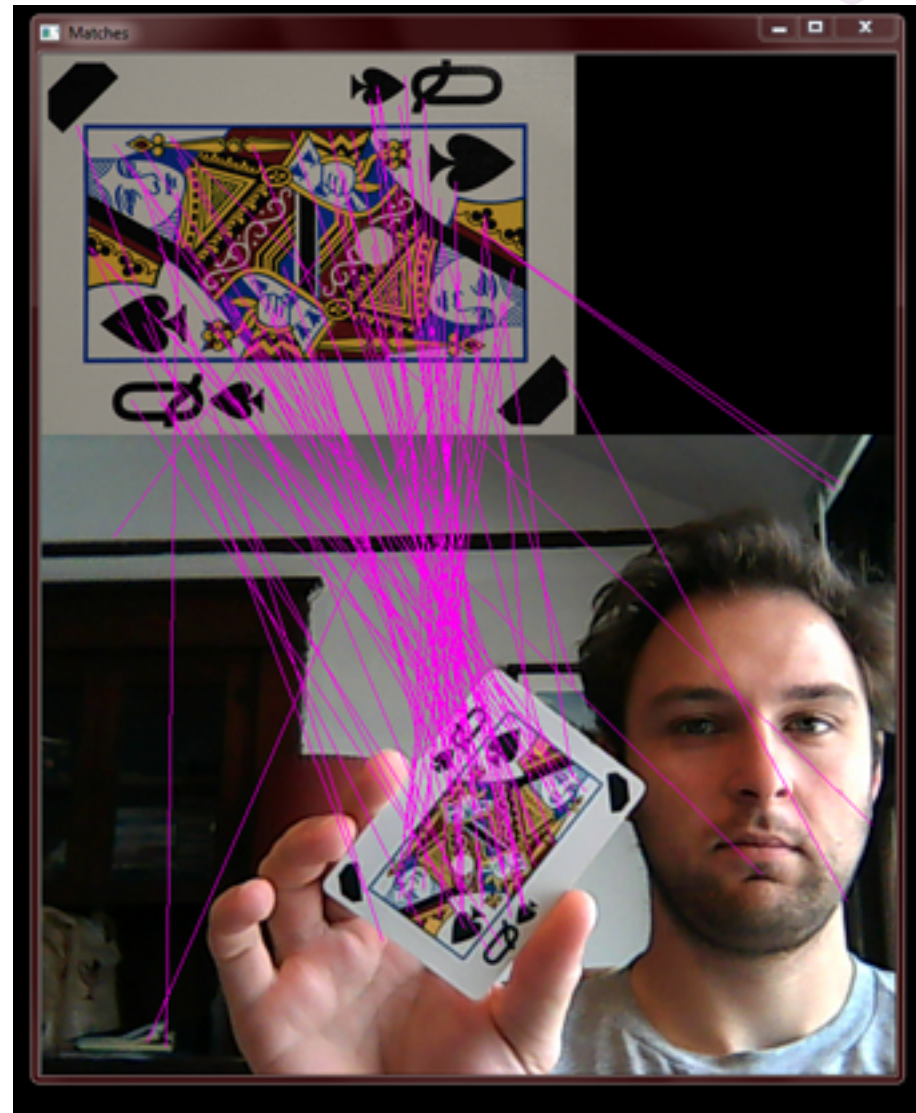
Hundreds of
Keypoints are Found

Vector representation



Local features

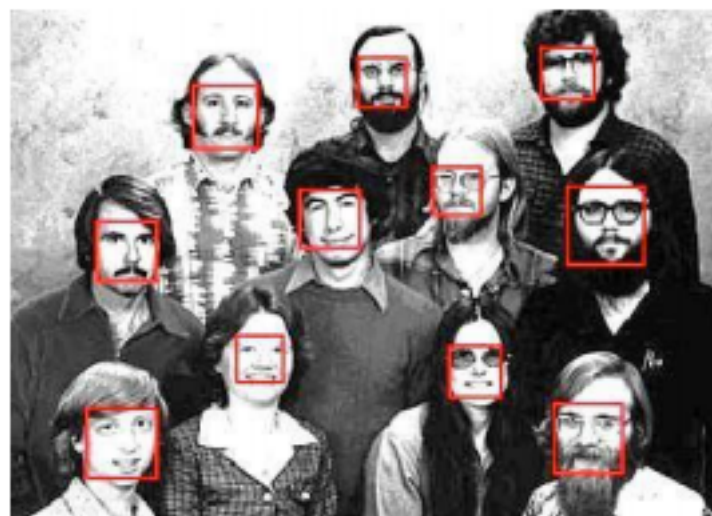
Bag of words of SIFT
vectors



Face detection



find faces in a given photo



sliding window



What is a face?



What is a face?



What is a face?



What is a face?

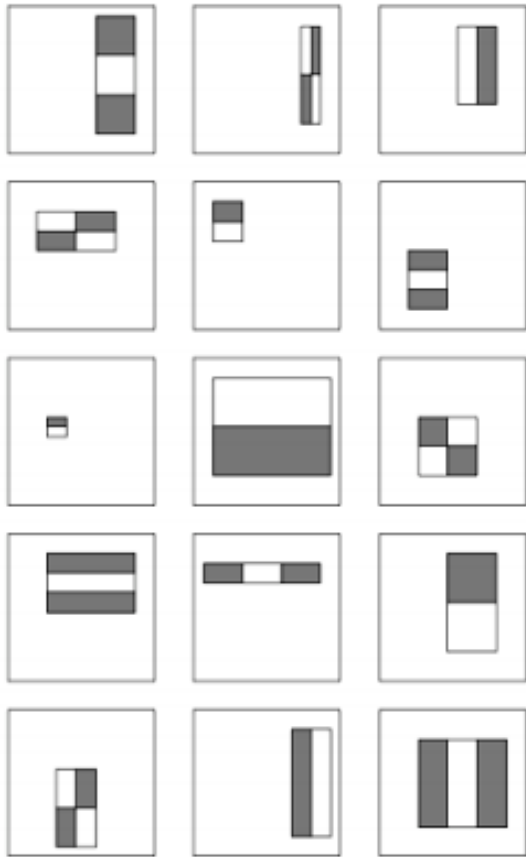


similar to positive face
rather than negative face

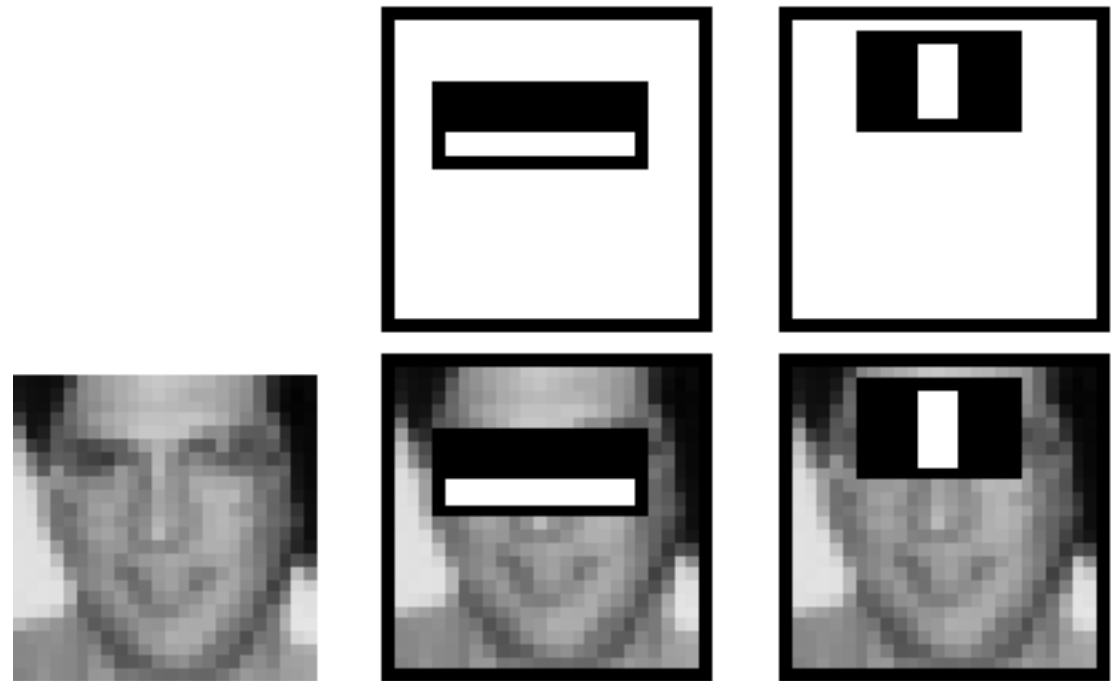
Viola&Jones face features



features: simple templates



for each sliding window apply templates to calculate features



conceptually forms a vector:

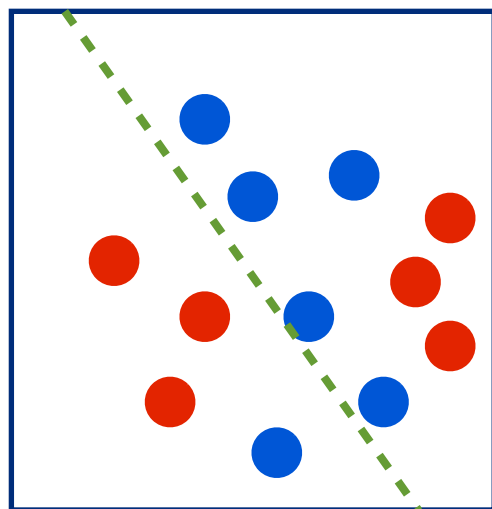
(200, 50, 90,)



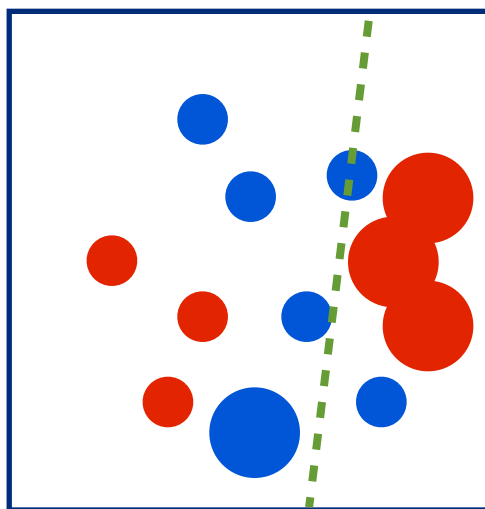
AdaBoost



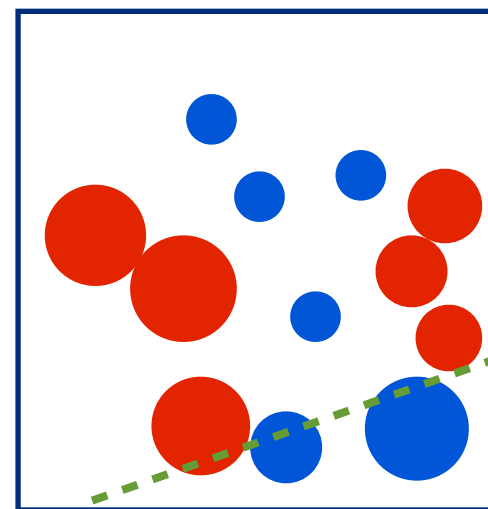
classifier 1



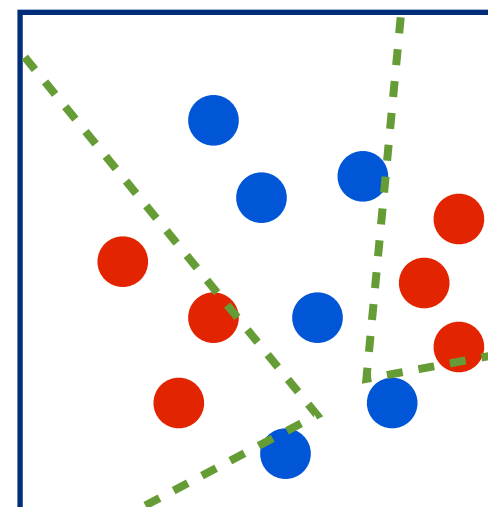
classifier 2



classifier 3



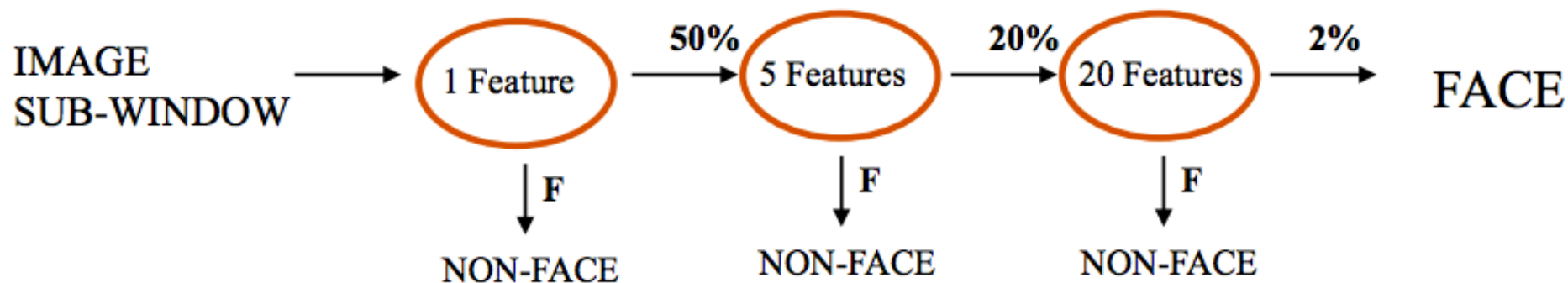
final classifier



In V&J's system, each classifier is one feature

AdaBoost selects a small subset of features

Viola&Jonse face detector



“15 times faster” than a state-of-the-art while keeping the accuracy”



THANKS

习题



对于用户的一条查询，数据库中总共有100个相关对象，系统返回了10个对象，其中不相关的有3个，请问对于这一条查询，系统的查准率 (Precision) 和查全率 (Recall) 各是多少？