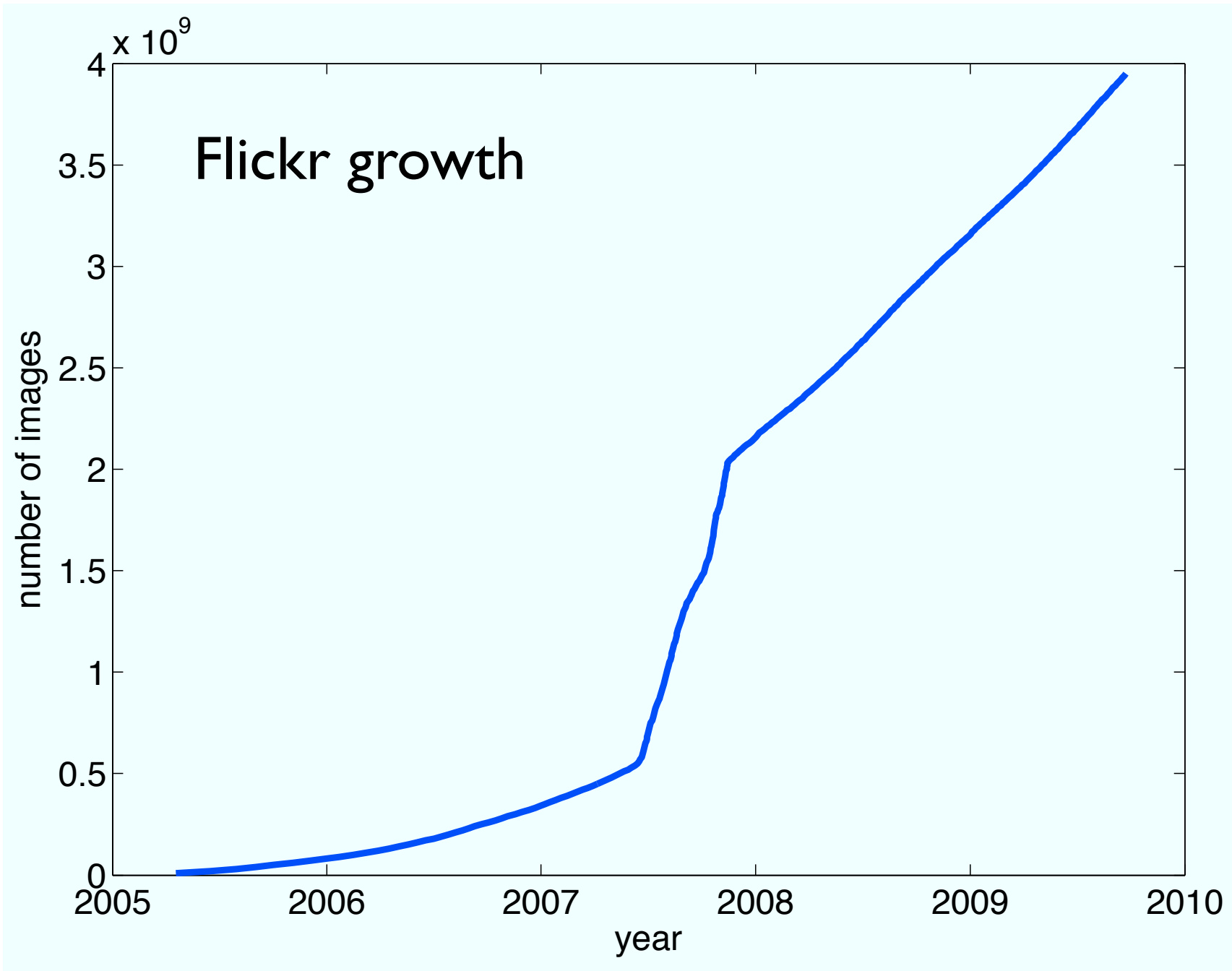# Scaling Object Recognition: Benchmark of Current State of the Art Techniques

Mohamed Aly[1], Peter Welinder[1],
Mario Munich[2], Pietro Perona[1]

(1) California Institute of Technology
(2) Evolution Robotics

Sunday, October 4, 2009

# Motivation



Flickr growth

y-axis: number of images ($\times 10^9$), from 0 to 4
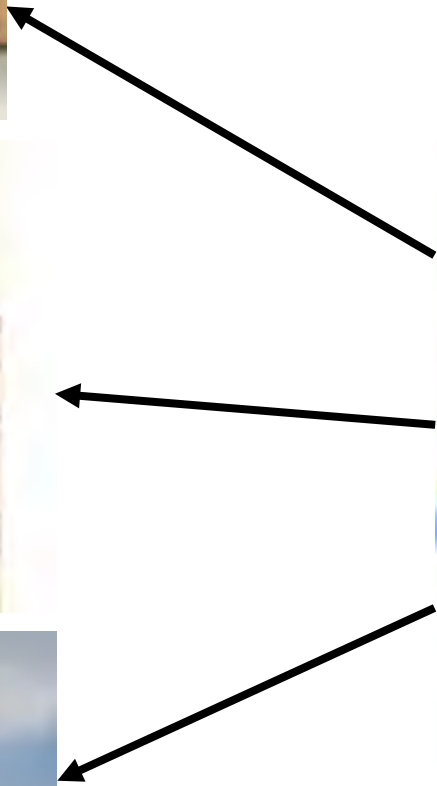
x-axis: year, from 2005 to 2010

# Motivation

Large image collections:

- 20 billions on ImageShack

- 15 billions on Facebook

- 7 billions on Photobucket

- 4 billions on Flickr

[http://www.techcrunch.com April 2009]

# Motivation

# Individual object recognition:

# How do current methods scale?

- CPU cycles

- RAM

- Precision / recall
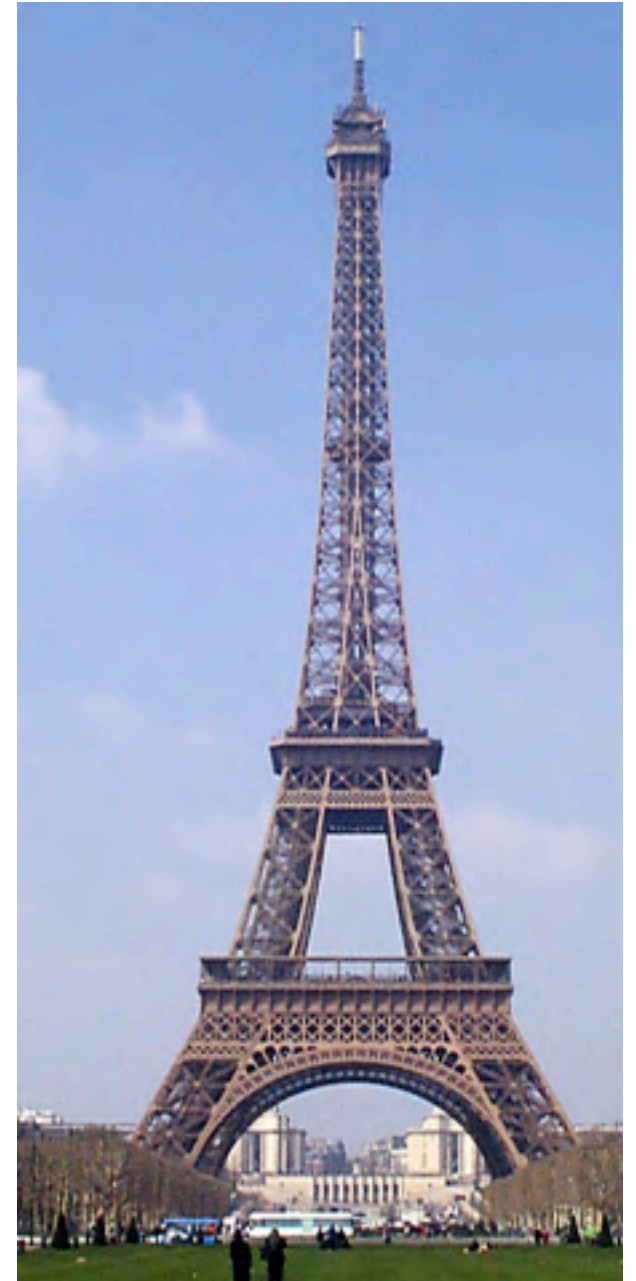
# Outline

- Datasets

- Recognition Methods

- Experimental Setup

- Results

- Conclusions
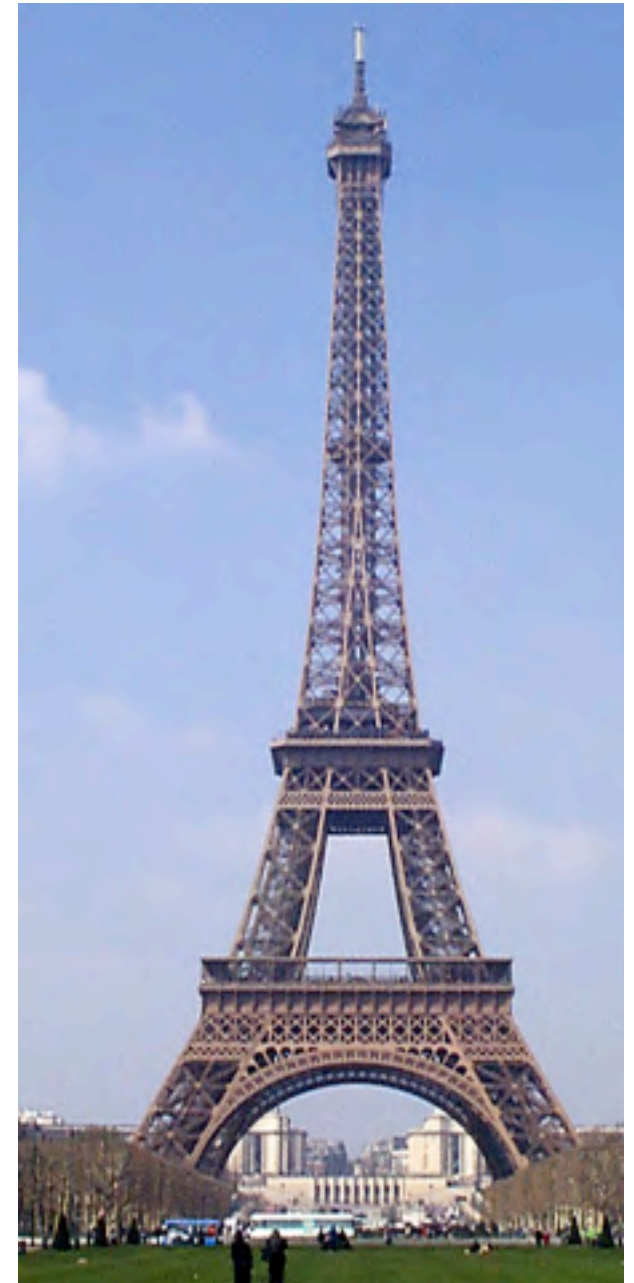
# Three flavors:

# Item to item



Probe

Model

# Scene to item



Probe

Model

# Item to scene



Probe

Model

# Three flavors:

- Item to item
- Item to scene
- Scene to item

# Dataset 1: CD Covers
## Model Set: ~ 132,000 unique images



downloaded from freecovers.net (available on vision.caltech.edu)
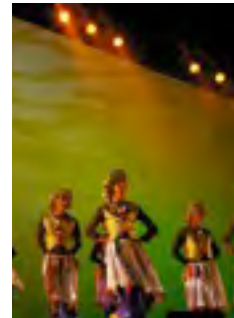
# Dataset 1: CD Covers

## Probe sets



500 Synthetic transformations

388 Photographs [Nister 06]

# Dataset 2: Pasadena Houses

Model set: ~$10^5$ photographs



125 pictures of
LA houses

$10^5$ flickr
photographs

# Dataset 2: Pasadena Houses

Probe set: 625 pictures of Pasadena houses



Query 1

Model

Query 5

Different:
- viewpoint
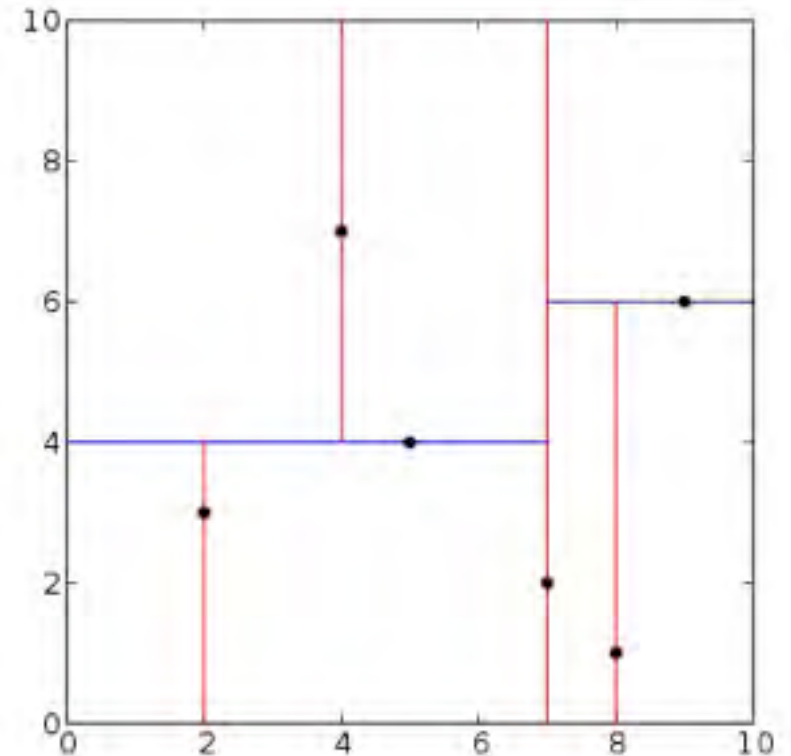- time of day
- camera

# Dataset 2: Pasadena Houses

# Recognition Approaches

- Sift/NN/Hough/RANSAC  [Lowe '04]

- Sift/Quantize/Rank  [VideoGoogle '03]

# Nearest-Neighbor 1: Kd-tree

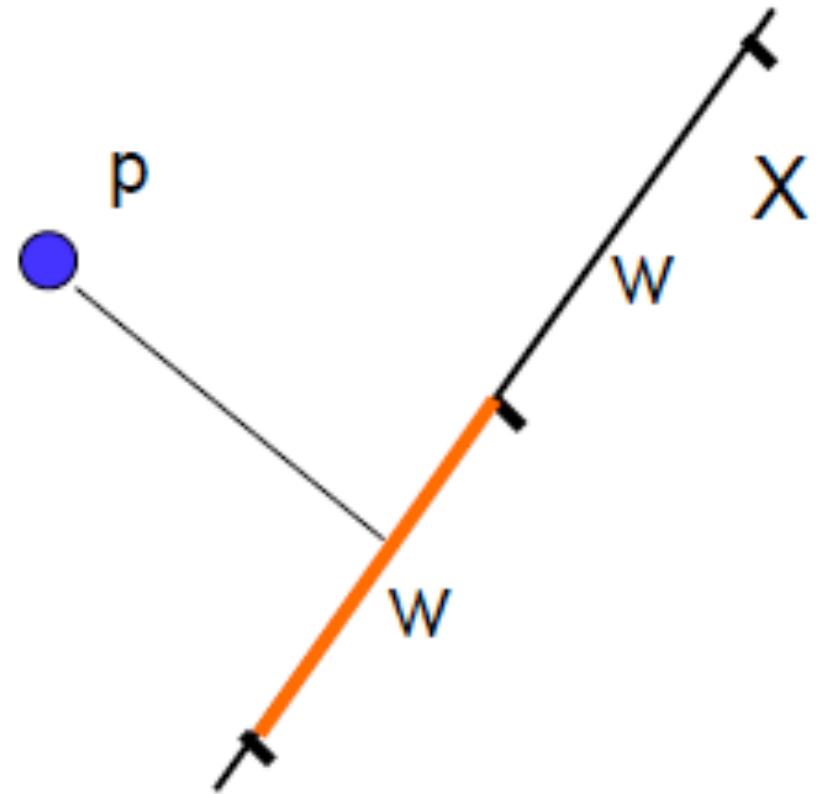- Approximate

- Build *O(d N)*

- Search *O(log(N))*

- Multiple trees: Kd-forest



N ~ number of prototypes

# Nearest-Neighbor 2: LSH

- E$^2$LSH package [Andoni 2004]

- Build $O(N)$

- Search $O(b) \sim O(N)$

# Method 3: Bag-of-Words



| query | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | word 1 | doc 1 | weight | doc 3 | weight | ... | doc M-3 | weight |
| word 2 | word 2 | doc 3 | weight | doc 4 | weight | ... | doc M | weight |
| word 3 | word 3 | doc 2 | weight | doc 4 | weight | ... | doc M-2 | weight |
| | word 4 | doc 1 | weight | doc 2 | weight | ... | doc M | weight |
| ⋮ | ⋮ | | | | | | | |
| | word N-1 | doc 4 | weight | doc 6 | weight | ... | doc M | weight |
| word N | word N | doc 1 | weight | doc 2 | weight | ... | doc M-1 | weight |

[Sivic et al., Video Goggle '03]

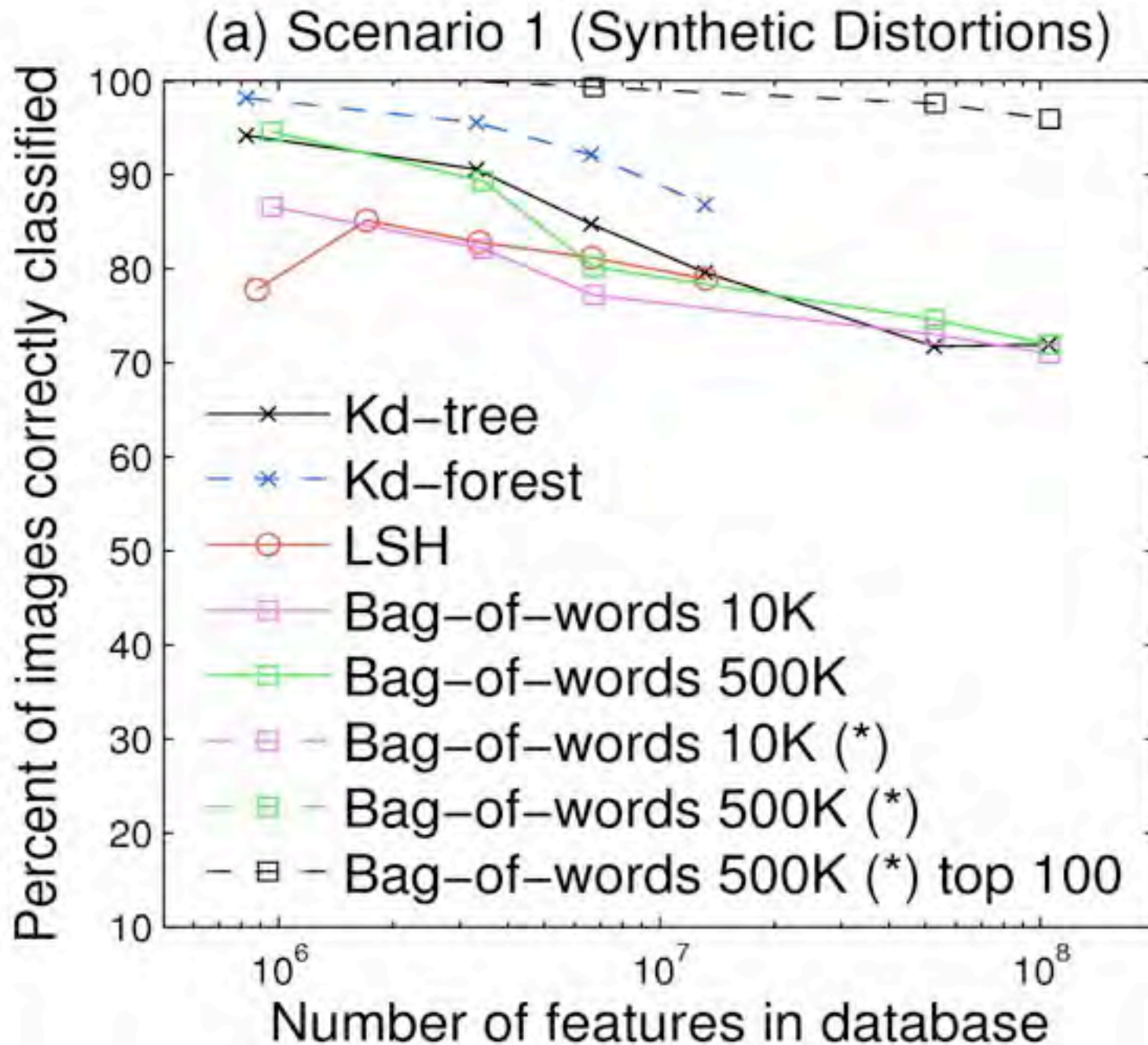# Method 3: Bag-of-Words

- Extract Sift Features

- Quantize using Approximate K-means with Kd-forest  *[Philbin et al. '07]*

- Compute word histograms *[Dorko-Schmid '03]*

- Search *O(N)*

- Fast search using Inverted File *[Sivic et al. '03]*

# Experimental Setup

- Datasets:

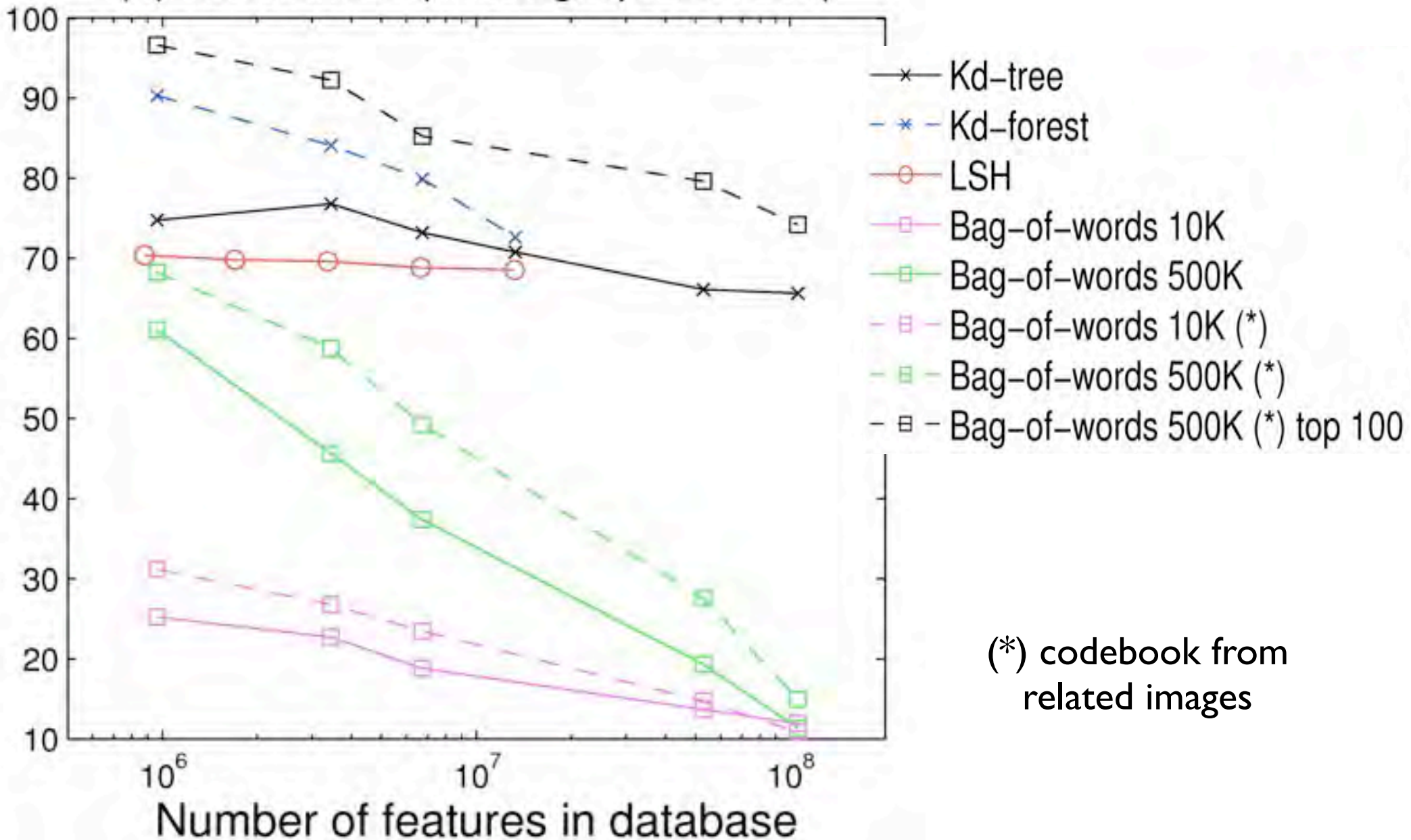| | Model Set | Probe Set | # images |
|---|---|---|---|
| Scenario 1 | Covers | Synthetic | 500 |
| Scenario 2 | Covers | Photographed | 388 |
| Scenario 3 | Flickr | Houses | 625 |

- One image/object in Model Set

- Rest in Probe Set

- Increase model set size: 1k, 4k, 8k, 16k, 32k, 64k, 128k images

# Results: Recognition



(a) Scenario 1 (Synthetic Distortions)

Legend:
- Kd–tree
- Kd–forest
- LSH
- Bag–of–words 10K
- Bag–of–words 500K
- Bag–of–words 10K (*)
- Bag–of–words 500K (*)
- Bag–of–words 500K (*) top 100

Y-axis: Percent of images correctly classified
X-axis: Number of features in database

# Results: Recognition



(b) Scenario 2 (Photographed CDs)

Legend:
- Kd-tree
- Kd-forest
- LSH
- Bag-of-words 10K
- Bag-of-words 500K
- Bag-of-words 10K (*)
- Bag-of-words 500K (*)
- Bag-of-words 500K (*) top 100

Number of features in database

(*) codebook from related images

# Results: Recognition



(c) Scenario 3 (Pasadena Buildings)

Legend:
- Kd–tree
- Kd–forest
- LSH
- Bag-of-words 10K
- Bag-of-words 500K
- Bag-of-words 10K (*)
- Bag-of-words 500K (*)
- Bag-of-words 500K (*) top 100
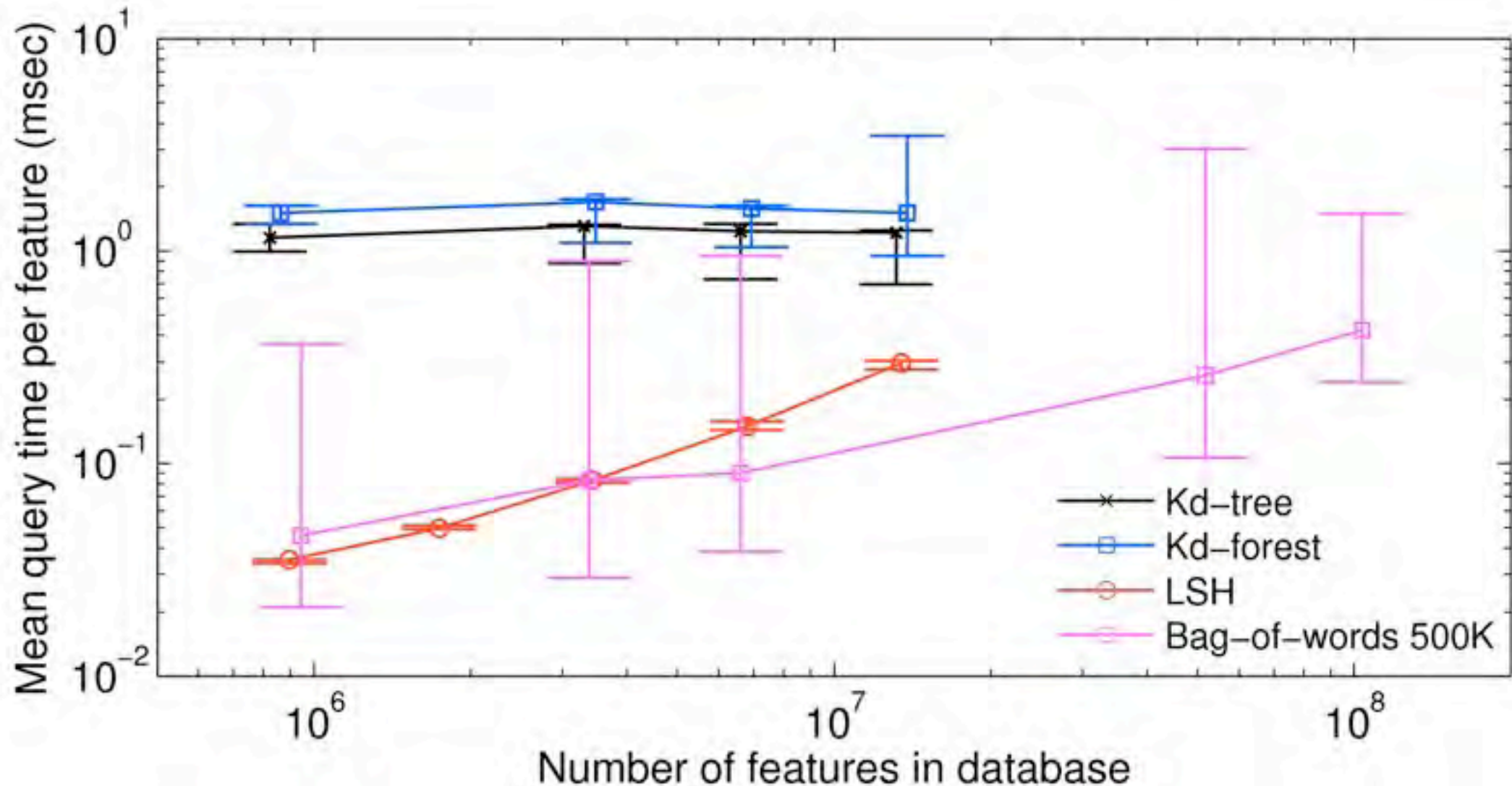
x-axis: Number of features in database

(*) codebook from related images

# Conclusions I

- Synthetic distortions are useless

- Performance drops w.r. to n. of features

- LSH scales best

- KD forest scales OK

- Bag-of-words scales poorly

# Results: Query Time



(Database: synthetic CD covers)

# Conclusions 2

- LHS scales linearly (ouch)

- Bag-of-words scales like sqrt(N)

- KD forests have constant cost O(1)

# Conclusions

- Importance of diverse datasets, natural probes

- Recall overall disappointing

- Nowhere close to $10^{10}$ images

- Bag-of-words recall does not scale well

- Kd-trees cost constant w.r. to N, unlike LSH

- Only bag-of-words fits in RAM beyond $10^5$ images

- Much work still ahead of us!