

Large-Scale Video Retrieval Using Image Queries

André Filgueiras de Araujo

Department of Electrical Engineering
Stanford University

The “Dark Matter” of the Digital Age



85% of data in the form of multimedia



400+ hours of video uploaded per minute



8+ billion video views per day



100+ hours of video uploaded per minute

Key problem: How can we make sense of these data?

Automatic Visual Recognition



Image classification

- *Is this an urban landscape?*

Object detection

- *Does this image contain a bus? Where?*

Instance recognition (a.k.a. “visual search”)

- *Does this image contain the “Wicked” billboard?*

Visual Search

Image query



Retrieval
System

Database of images



Product recognition
[Tsai et al., MM'08, MM'10]



Location recognition
[Chen et al., CVPR'11]



Commercial applications

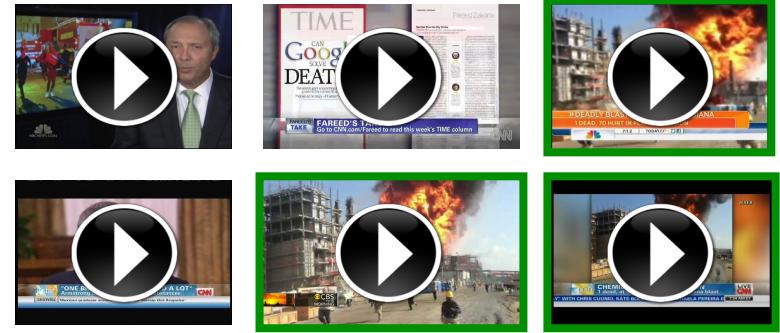
Video Retrieval Using Image Queries

Image query



Retrieval
System

Database of video clips



Applications:

- Brand monitoring: search YouTube using product images
- News videos: search event footage using photos
- Online education: search lectures using slides

Online Prototype

Search news videos with images

Please choose a query image file:

Choose File No file chosen

Search Video Database

Please enter a query image URL:

http://stanford.edu/image.jpg

Search Video Database

Trending Images on TV (click to explore)



Random Query Images (click to explore)

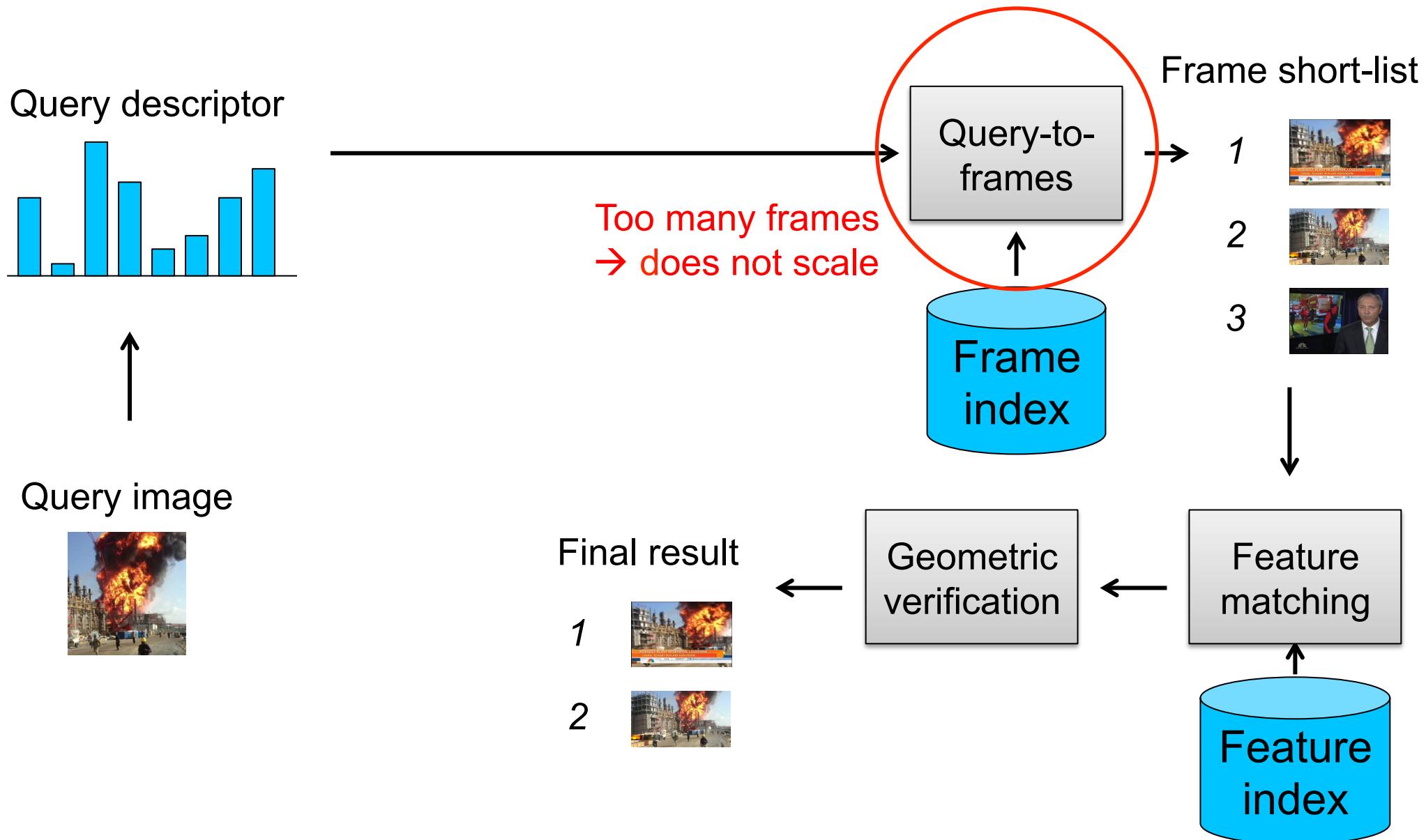


Lost? Watch our [demo video](#) and read our [paper](#).

[Andre Araujo](#), [David Chen](#), [Peter Vajda](#), [Bernd Girod](#)
Image, Video, and Multimedia Systems Group at [Stanford University](#)

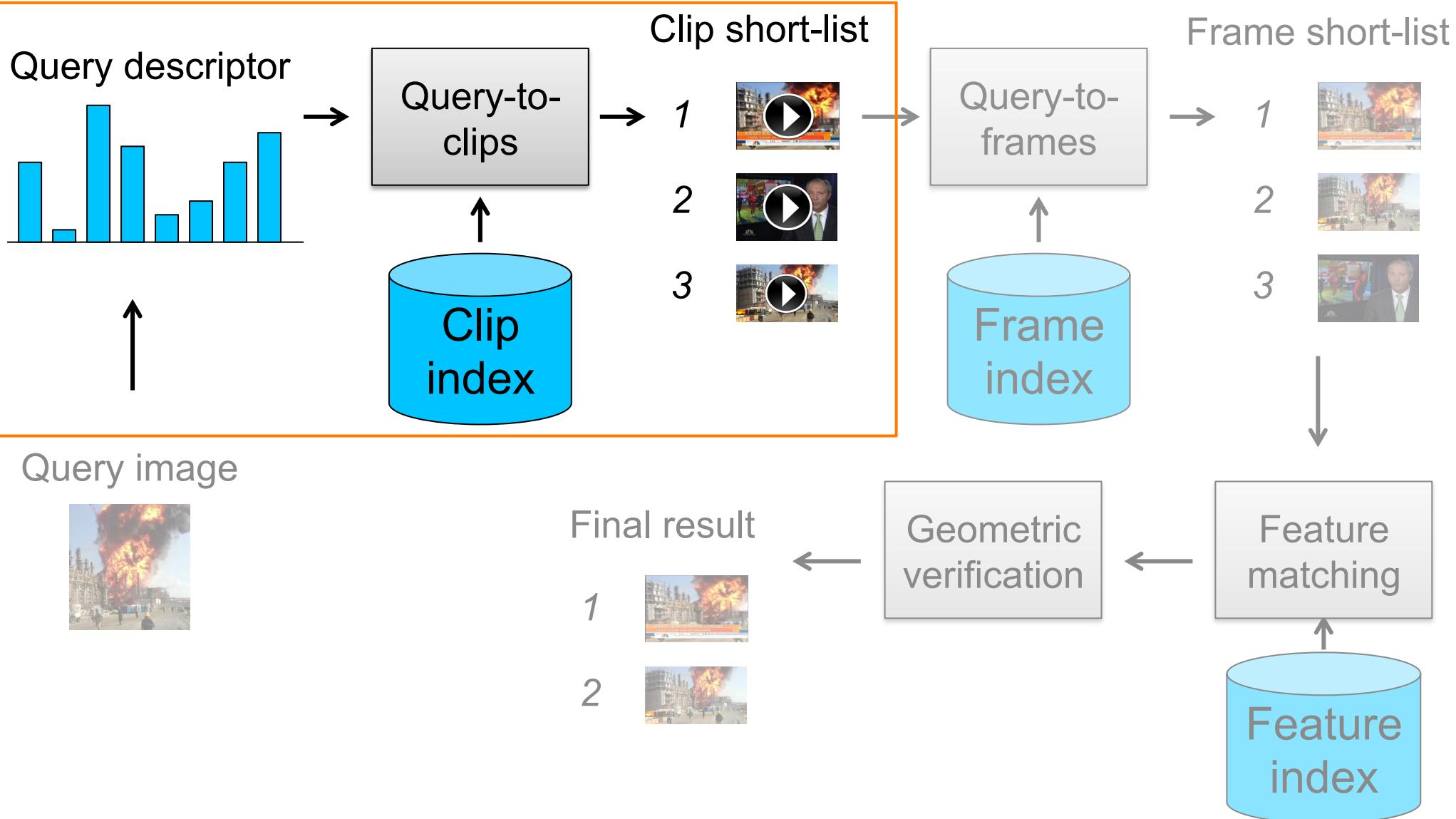
<http://videosearch.stanford.edu>

Simple Architecture

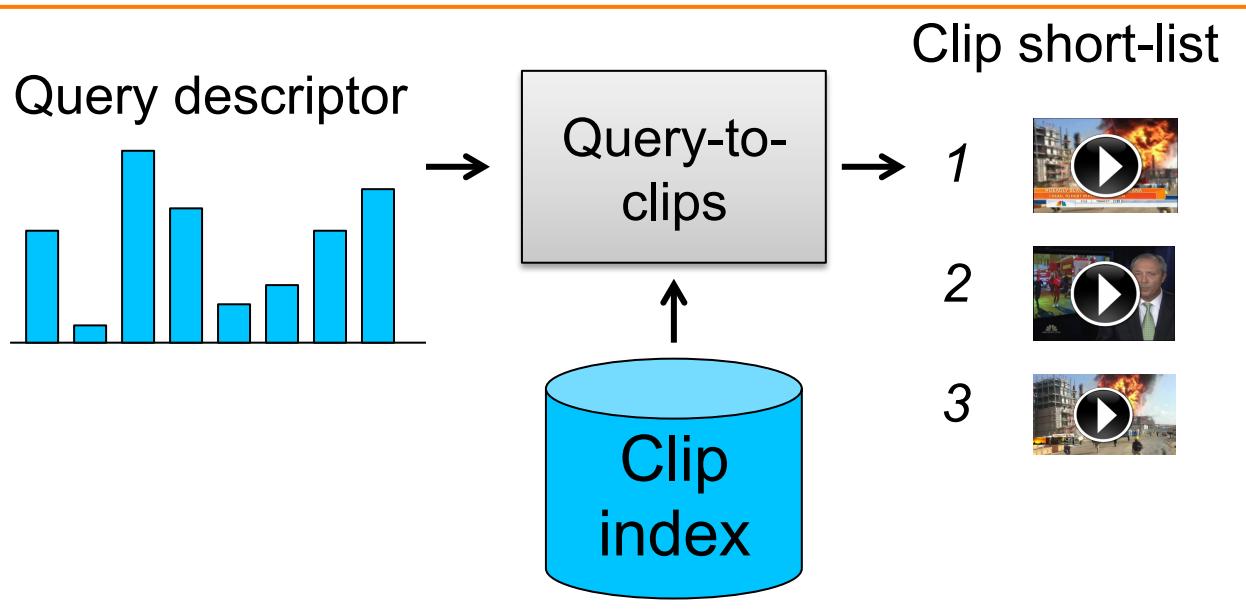


Large-Scale Architecture

Focus of this work



Video Retrieval Using Image Queries



Main challenges:

- Asymmetry: how can we compare images to videos?
- Temporal aggregation: how can we describe a video clip for query-by-image retrieval?

Contributions

Fisher Vector Comparisons

- Asymmetric comparisons for Fisher vectors
- Cluttered query or database images

Fisher Vector Aggregation

- Fisher vector descriptors for video segments
- Compact database for large-scale retrieval

Bloom Filter Aggregation

- Bloom filter descriptors for video segments
- Fast and accurate large-scale retrieval

Related Work: Visual Search

Query

Video	Augmented Reality	Content Tracking
	TCD [Makar et al., '12]	Frame Mat. + ST [Douze et al., '10]
Image	Hybrid Vis. Search [Chen et al., '14]	TRECVID-CCD [Over et al., '12]
	Traditional Visual Search	<u>Video Retrieval by Image</u>
Images	FV [Perronnin et al., '07]	<i>Discussed on next slide</i>
	BoW [Sivic et al., '03]	
	SIFT [Lowe, '04]	

Images

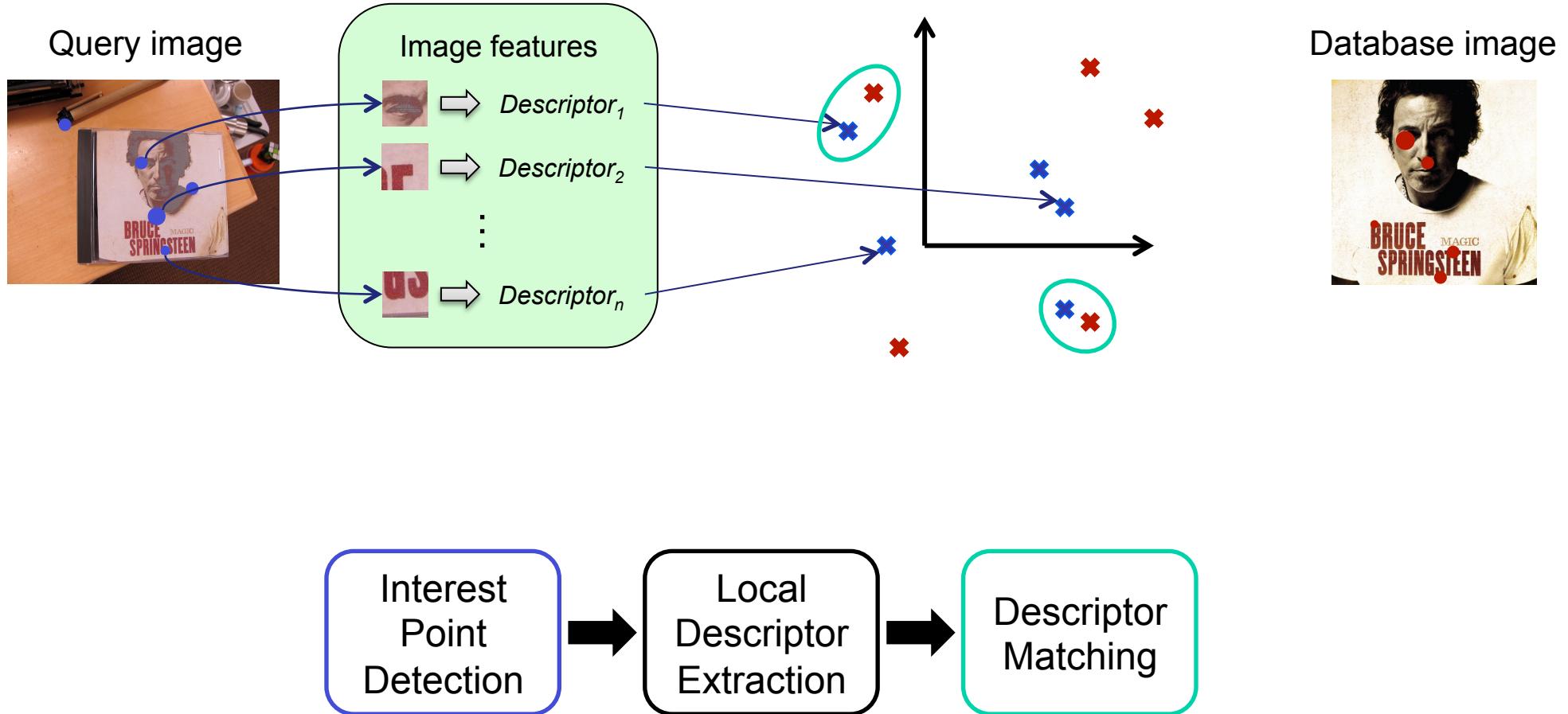
Videos

Database

Related Work: Video Retrieval Using Images

- Early work
 - BoW retrieval of movie frames [*Sivic and Zisserman, ICCV'03*]
 - Object-level retrieval of movie shots [*Sivic et al., ECCV'04*]
- TRECVID Instance Search Challenge [*Over et al., TRECVID'10-15*]
 - Frame-based BoW with Color SIFT [*Le et al., '10-11*]
 - Shot-based aggregation using BoW [*Zhu et al., '13*] [*Ballas et al., '14*]
 - BoW query-adaptive asymmetrical dissimilarities [*Zhu et al., '13*]
- Object localization in videos
 - SURF-based matching per shot [*Apostolidis et al., ICME'13*]
 - Optimal path using dynamic programming [*Meng et al. ICIIP'15*]

Background: Pairwise Image Matching



Background: Fisher Vector (FV)

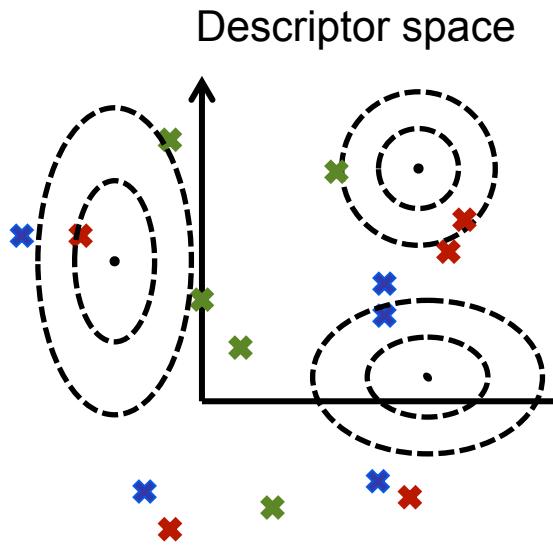
[Perronnin and Dance, CVPR'07]

- State-of-the-art technique for large-scale retrieval
- Key property: represent a set of local descriptors by a compact fixed-length vector
 - Two images can be compared by comparing their Fisher vectors
- Construction: describe an image with aggregated Fisher scores of its local descriptors
 - Local descriptor distribution: Gaussian Mixture Model (GMM)
 - Usually only Gaussian means are taken into account
- Extension of Bag-of-Words technique [Sivic and Zisserman, ICCV'03]

Background: Fisher Vector (FV)

[Perronnin and Dance, CVPR'07]

Query image



Database image 1



Query FV

-0.2	0.2	-0.3	-0.3	-0.3	0.8
------	-----	------	------	------	-----

DB Im. 1 FV

-0.3	0.3	0.3	-0.6	-0.3	0.3
------	-----	-----	------	------	-----

DB Im. 2 FV

0.5	-0.2	-0.7	0.1	-0.6	0
-----	------	------	-----	------	---

Database image 2



⋮

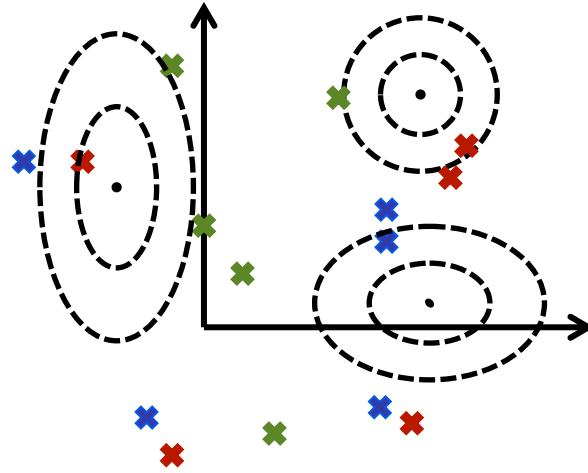
Background: Binarized Fisher Vector (FV*)

[Perronnin et al., CVPR'10]

Query image



Descriptor space



Database image 1



Query FV^*

0	0	0	0	0	1
---	---	---	---	---	---

DB Im. 1 FV^*

0	1	1	0	0	1
---	---	---	---	---	---

DB Im. 2 FV^*

1	0	0	1	0	0
---	---	---	---	---	---

Database image 2



⋮

Contribution 1

Fisher Vector Comparisons

- Asymmetric comparisons for Fisher vectors
- Cluttered query or database images

Fisher Vector Aggregation

- Fisher vector descriptors for video segments
- Compact database for large-scale retrieval

Bloom Filter Aggregation

- Bloom filter descriptors for video segments
- Fast and accurate large-scale retrieval

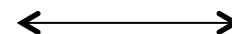
Asymmetric Image Comparison

*Object retrieval
application*

Query image



Database image



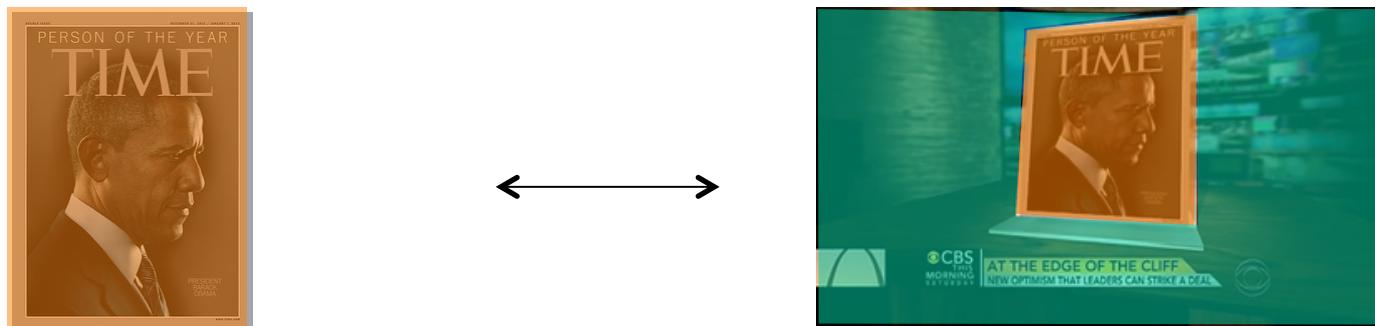
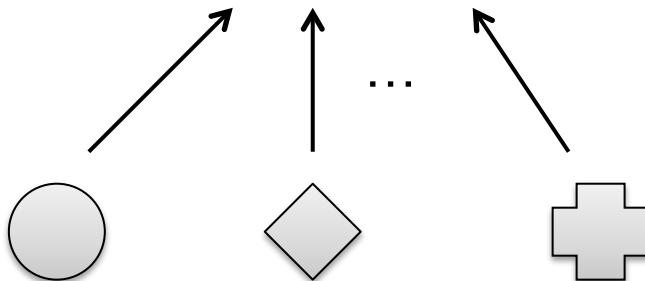
*Video bookmarking
application*



How can we incorporate asymmetry in FV comparisons?

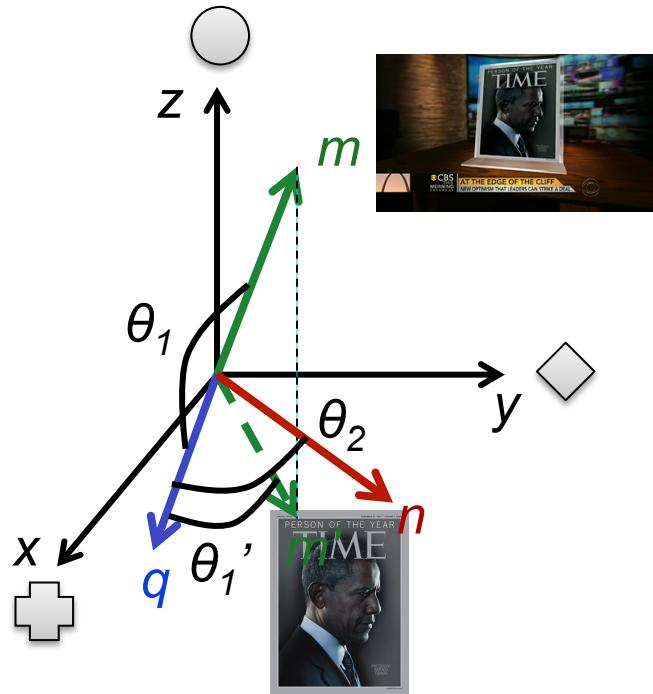
Asymmetric Comparison for FV

Fisher vector = $[v_1, v_2, \dots, v_K]$



Regions and have different statistics
→ features from are usually not present in

Asymmetric Comparison for FV



q query

m correct match in database

n incorrect match in database

θ_1 = $\text{angle}(q, m)$

θ_2 = $\text{angle}(q, n)$

θ_1' = $\text{angle}(q, m')$

- FV comparison metric: cosine similarity

- We want: $\theta_1 < \theta_2$

- Common failure case:

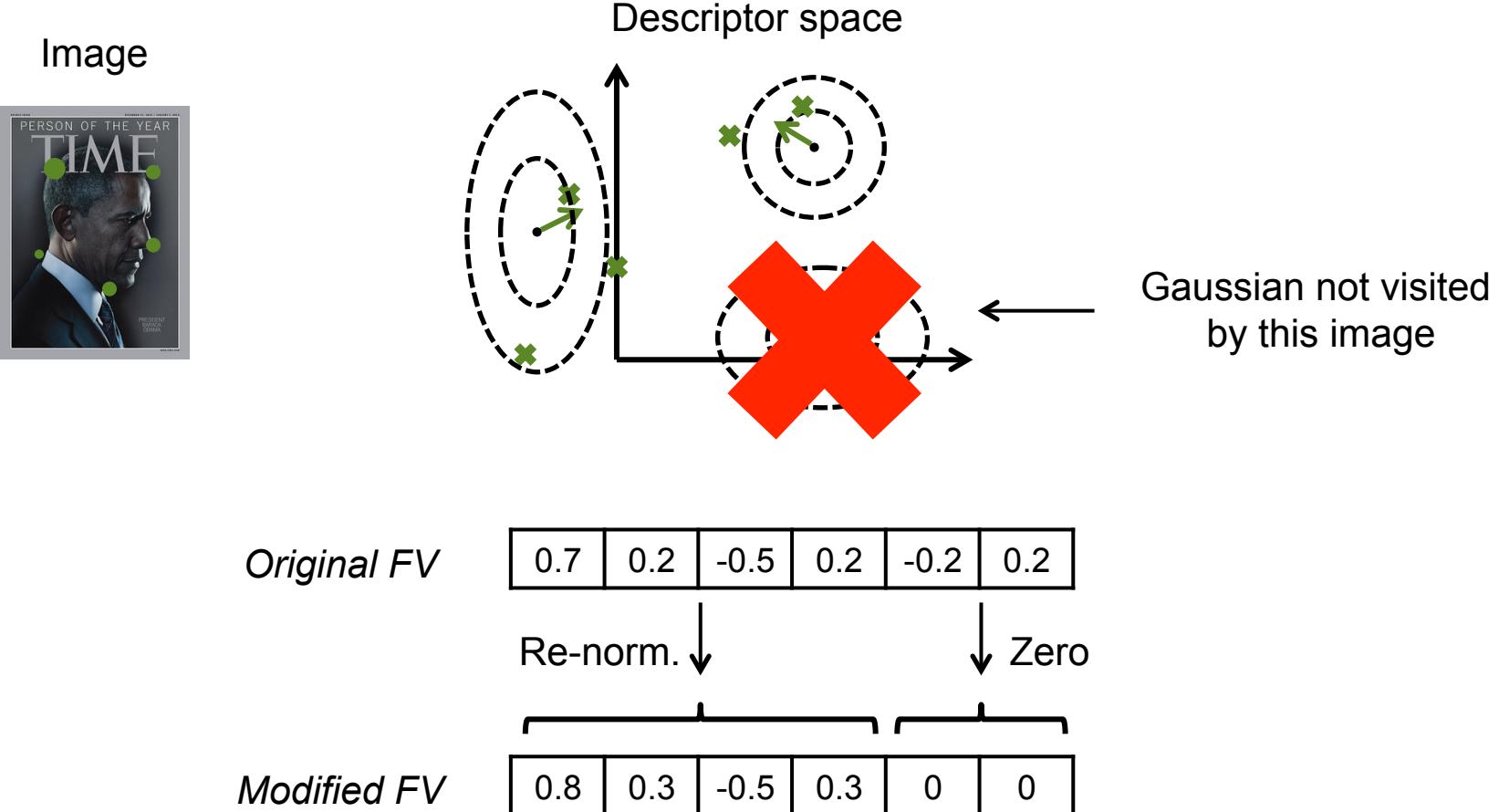
$$\theta_1 > \theta_2 \quad \text{but} \quad \theta_1' < \theta_2$$

- Insight:

Compare query and database based on their projections to the x-y plane

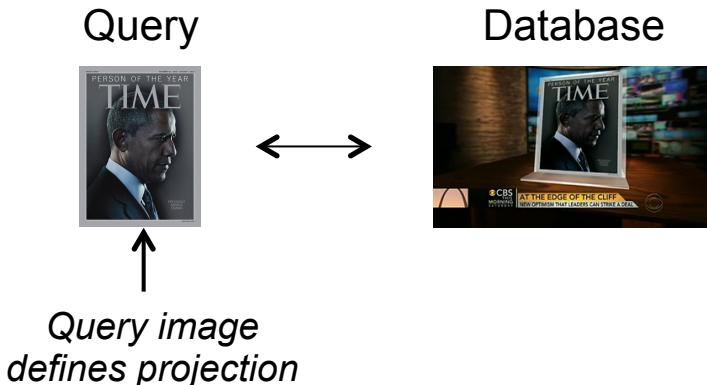
(i.e., using only Gaussians visited by query)

Asymmetric Comparison for FV



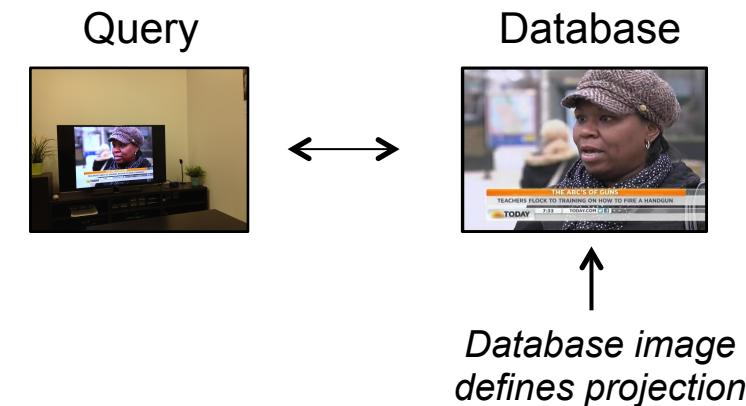
Asymmetric Comparison for FV

- Two retrieval problems
 - Query contained in database
All database images compared to query based on the same subspace

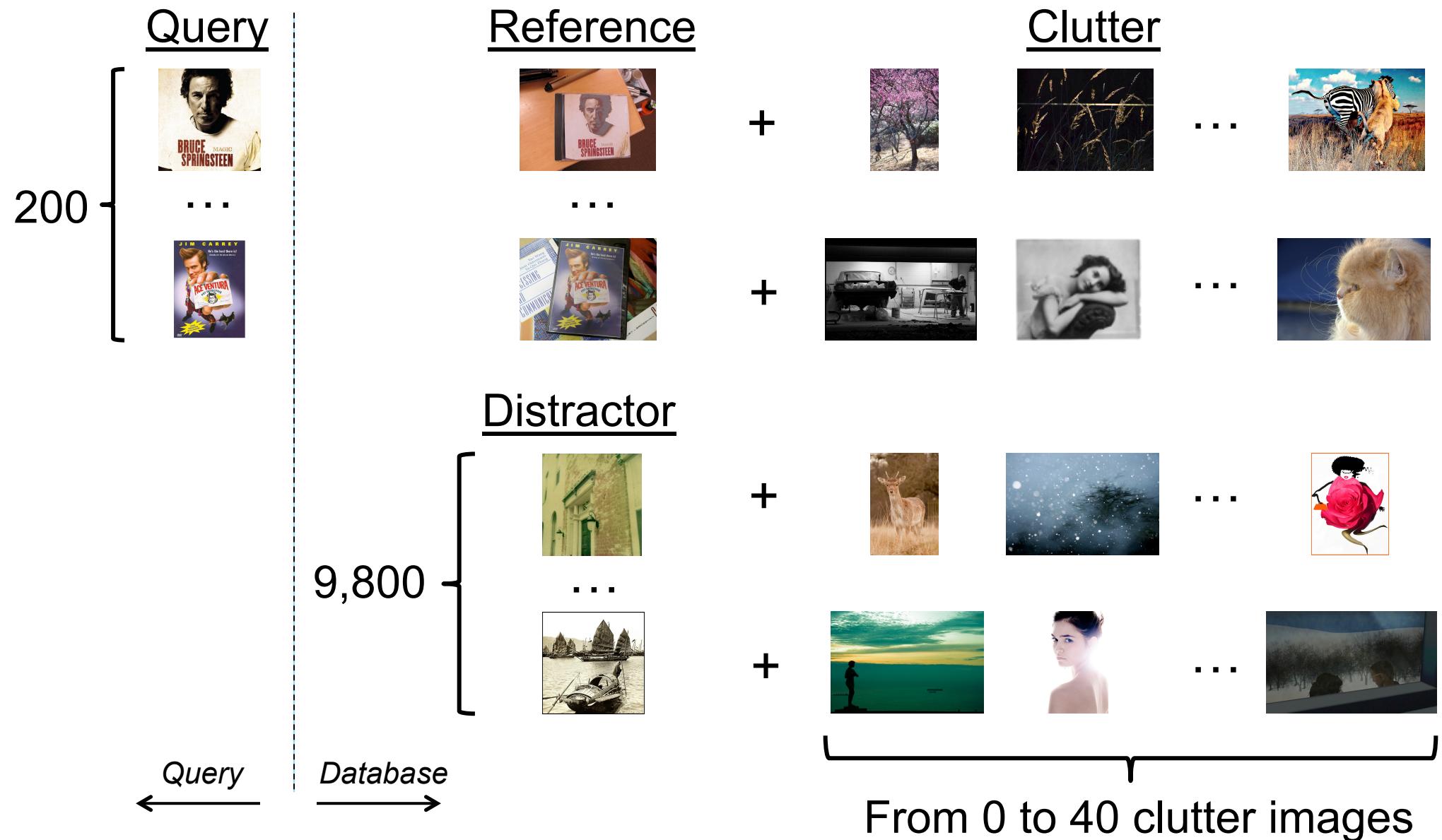


- Database contained in query
Problem: each database image is compared to the query based on different subspaces

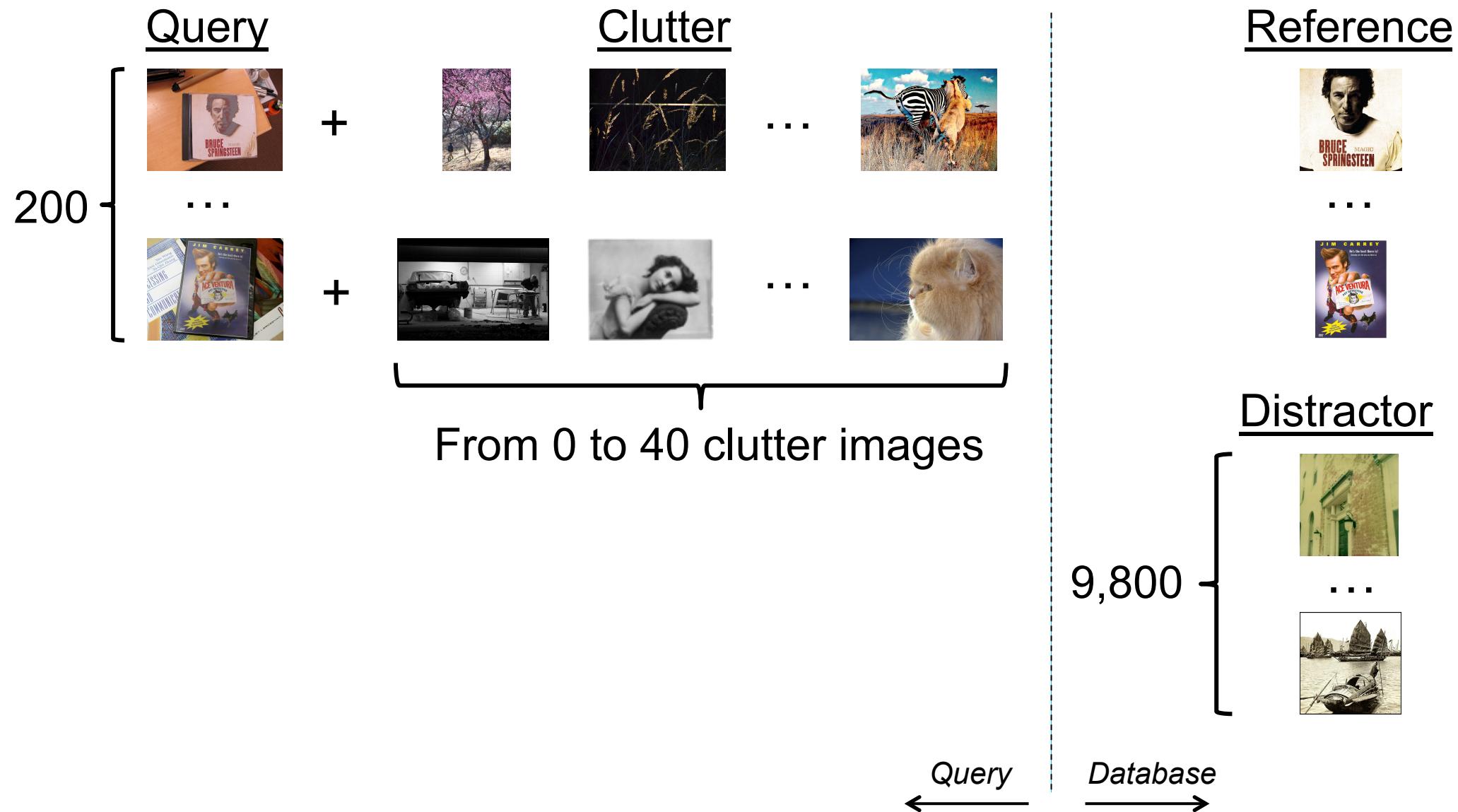
Solution: introduce weight to favor database images with more visited Gaussians



Dataset: Query Contained in Database

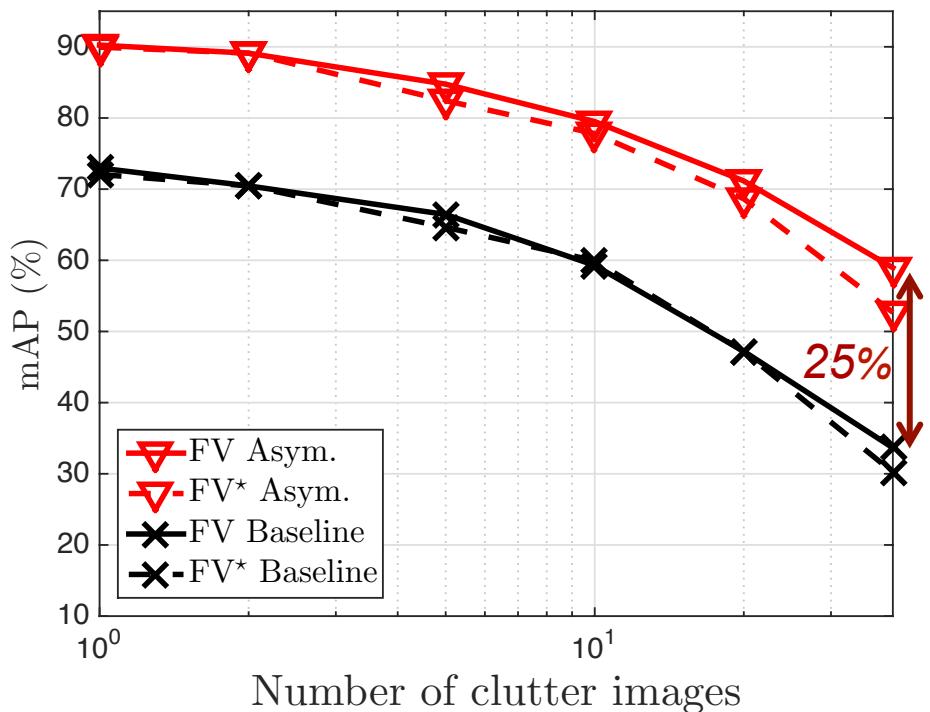


Dataset: Database Contained in Query

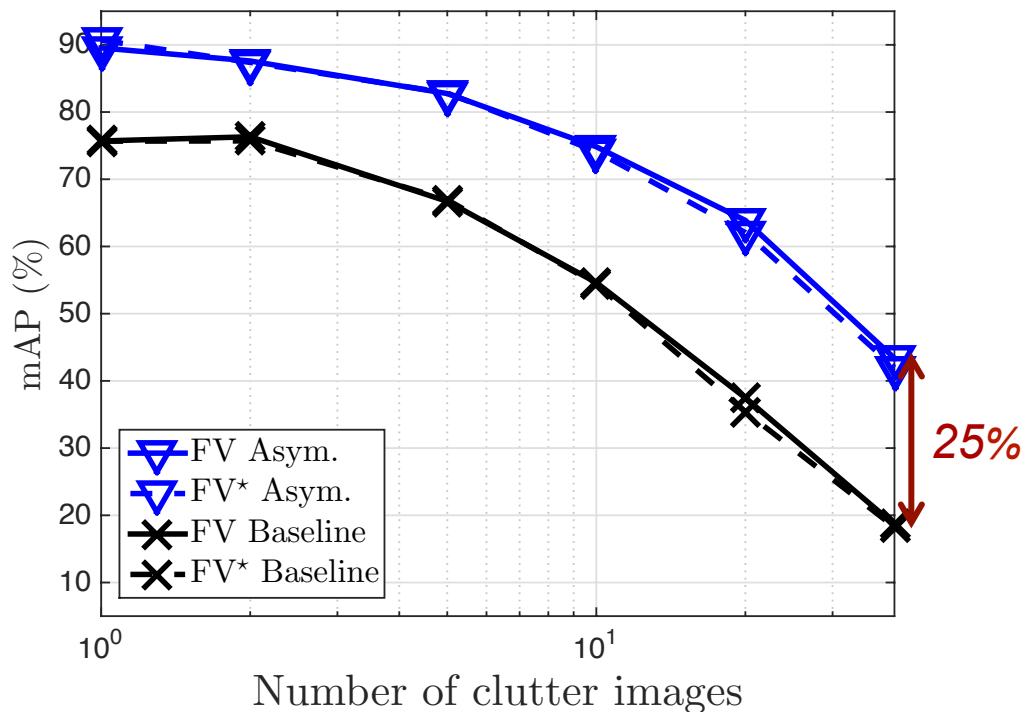


Experiments: Asymmetric FV Comparisons

Query contained in database



Database contained in query



Contribution 2

Fisher Vector Comparisons

- Asymmetric comparisons for Fisher vectors
- Cluttered query or database images

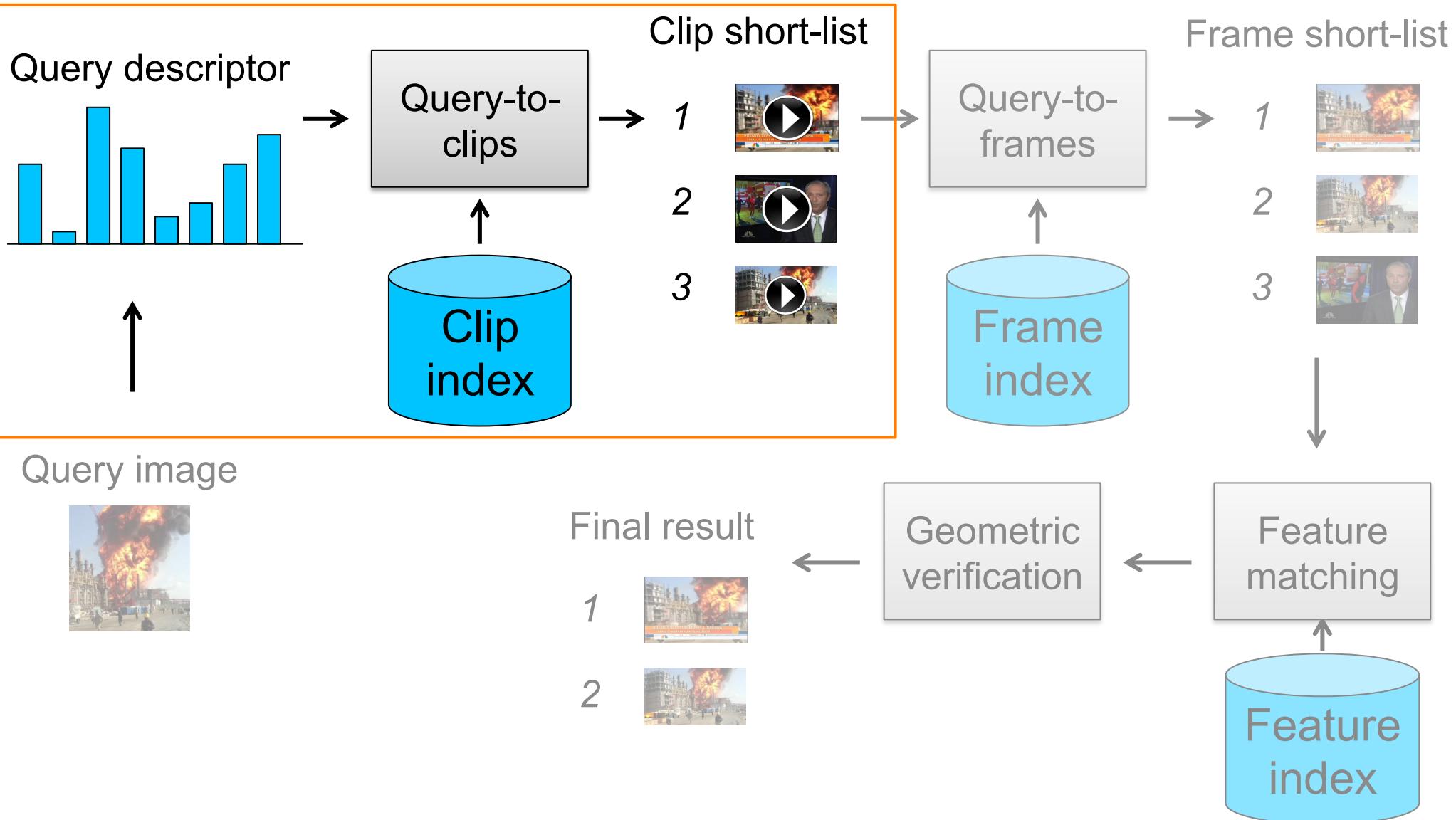
Fisher Vector Aggregation

- Fisher vector descriptors for video segments
- Compact database for large-scale retrieval

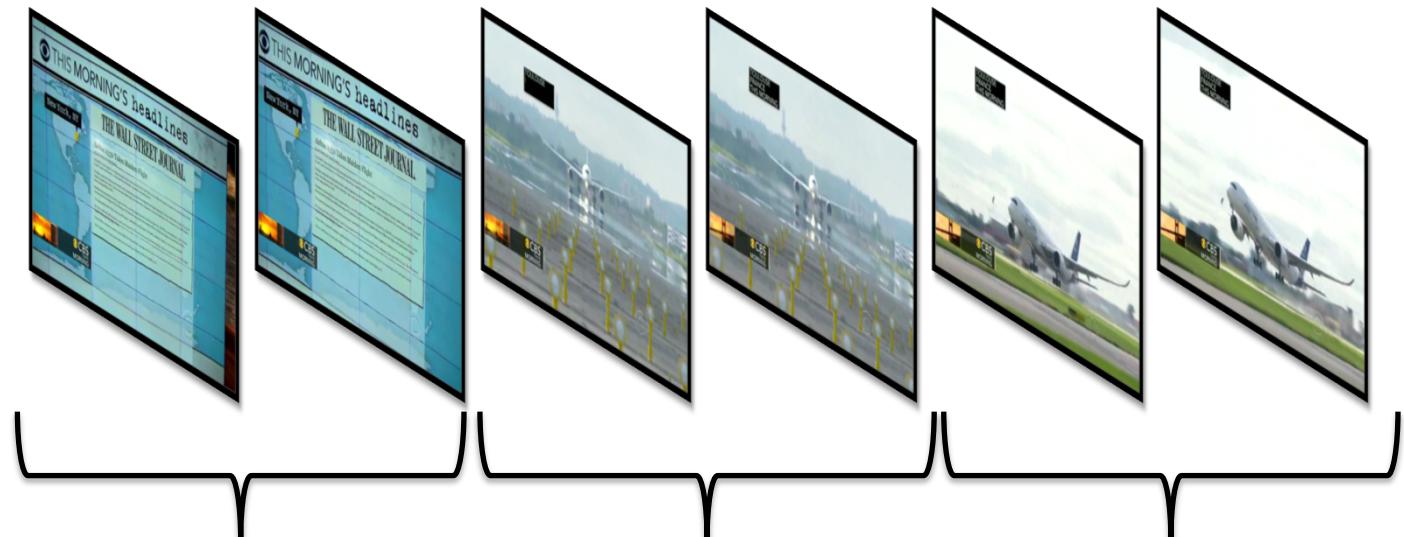
Bloom Filter Aggregation

- Bloom filter descriptors for video segments
- Fast and accurate large-scale retrieval

Large-Scale Architecture



Temporal Structure



Frames
1 fps

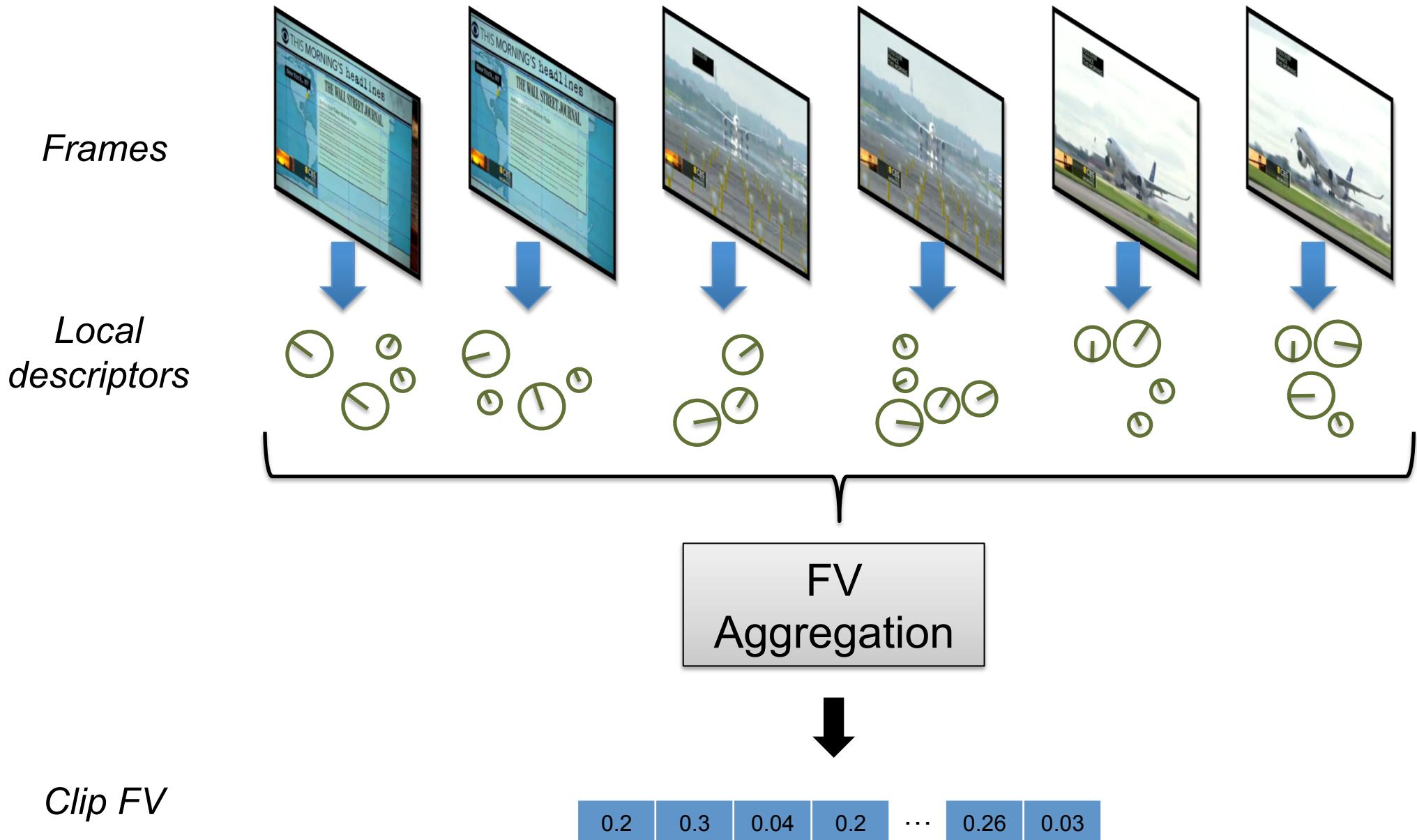


Shots
Contain similar frames
Length of seconds

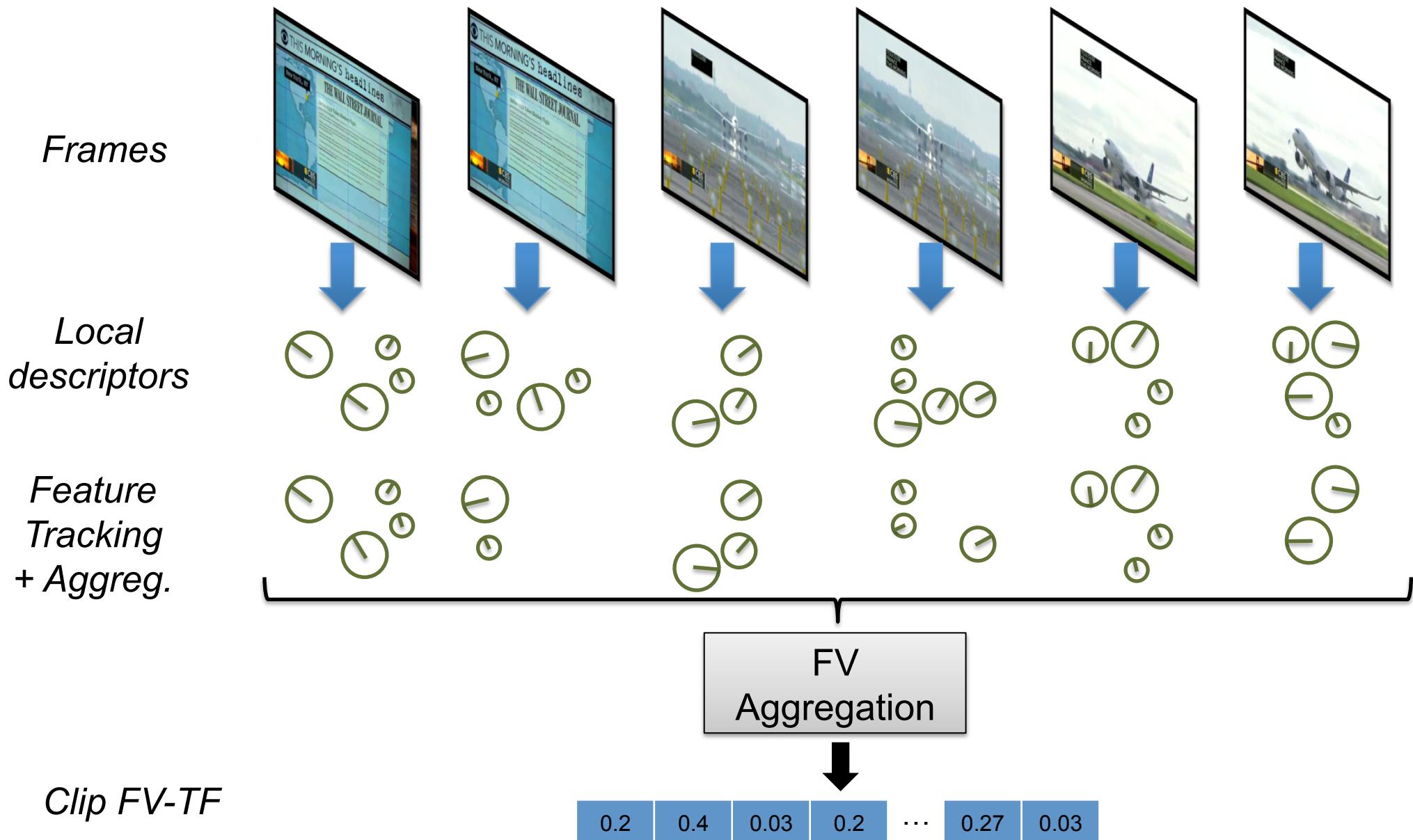


Clips
Contain diverse shots
Length of minutes

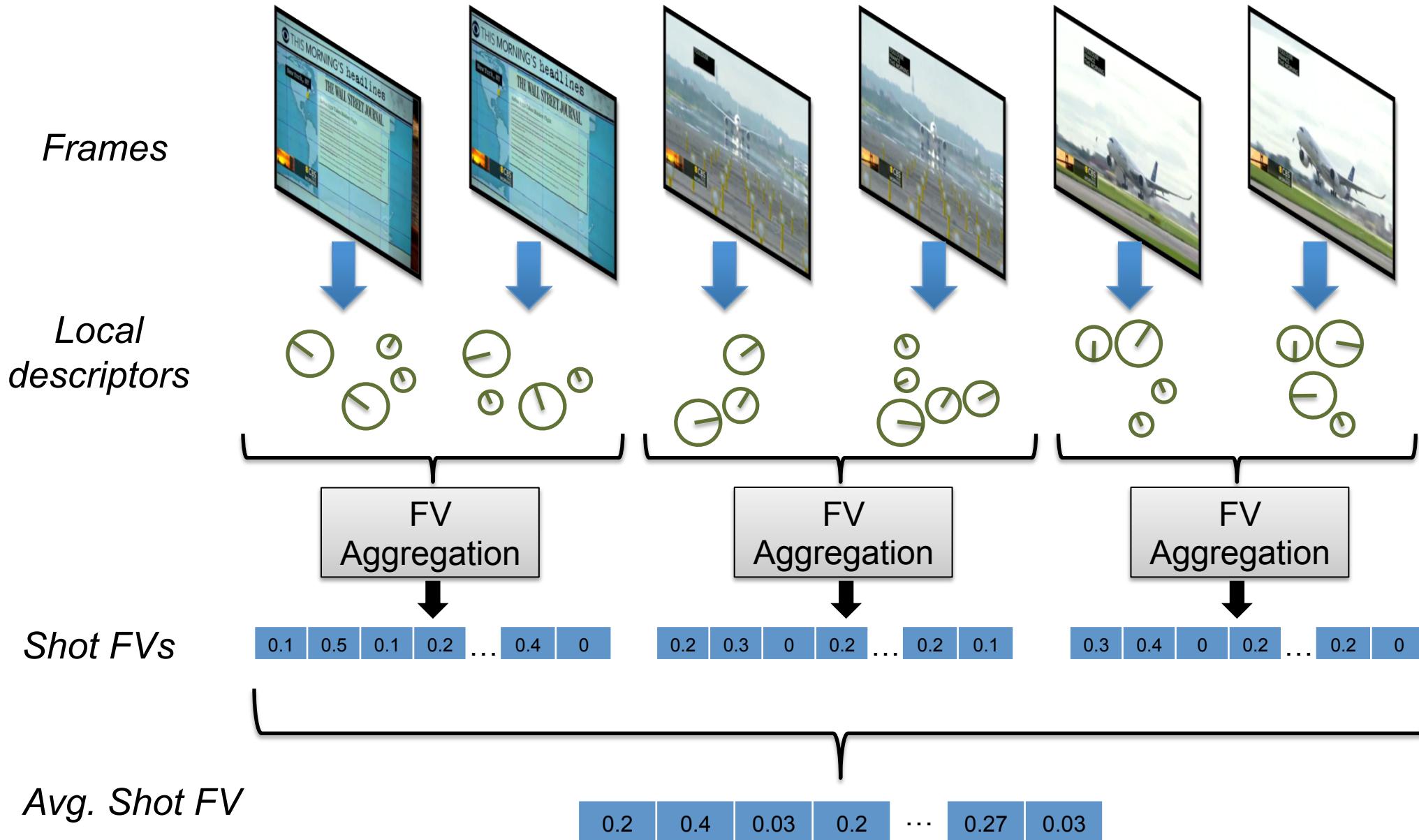
Clip Fisher Vector



Clip Fisher Vector with Tracked Features



Averaged Shot Fisher Vectors



Datasets

News Videos

Query



Database



Video Bookmarking

Query

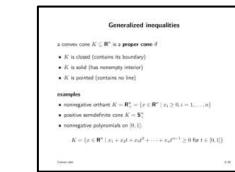
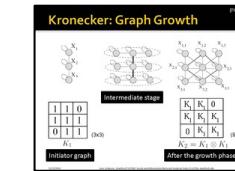
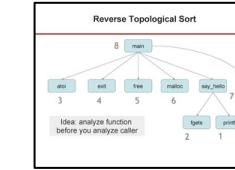


Database

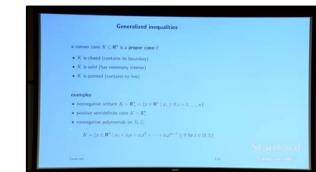
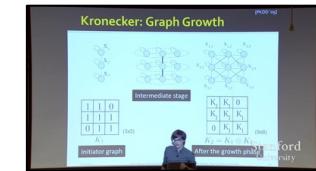
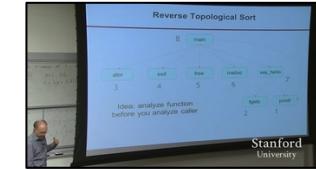


Lecture Videos

Query



Database

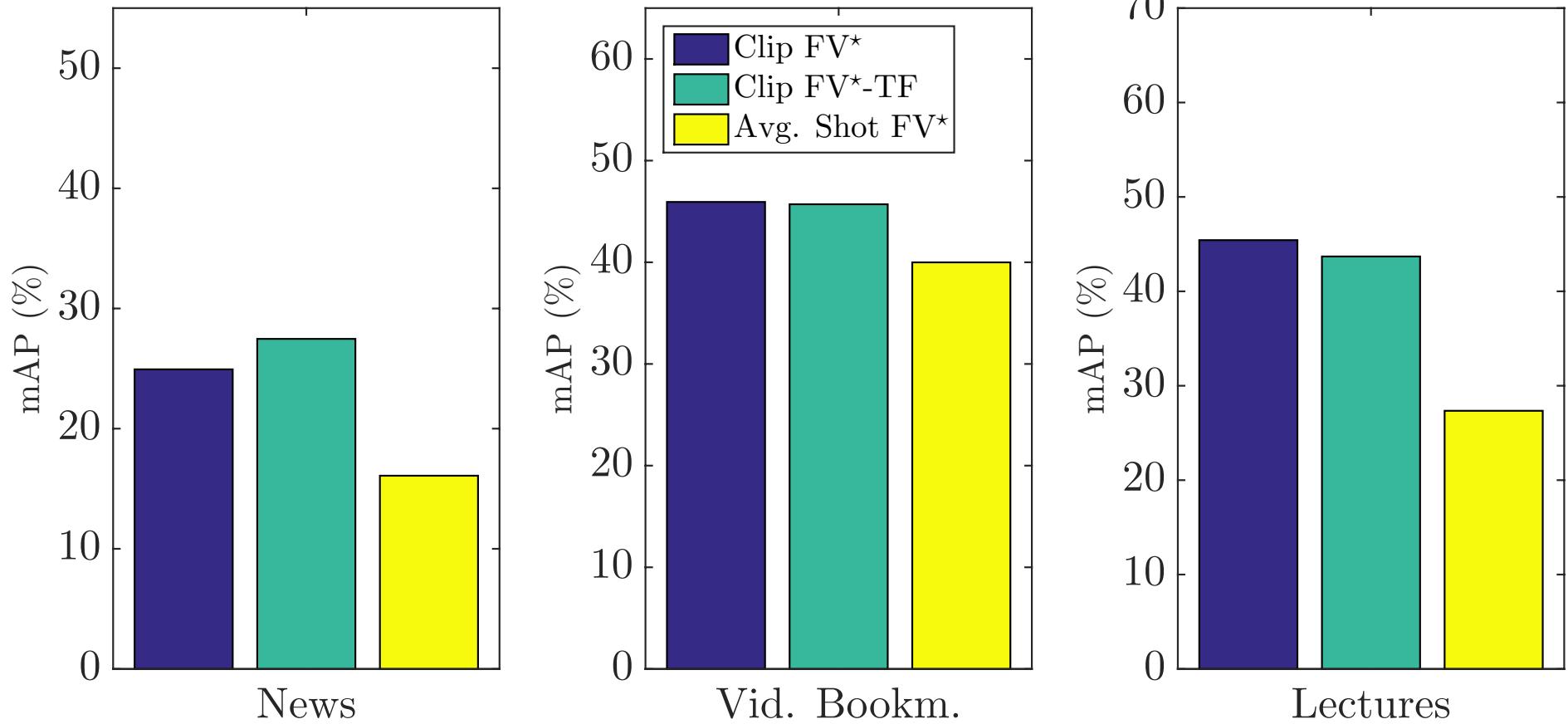


- 229 queries (from web)
- 2.7 minutes/clip
- 50.6 shots/clip
- Versions
 - 600k frames, 164h, 3.4k clips
 - 4M frames, 1,079h, 24.3k clips

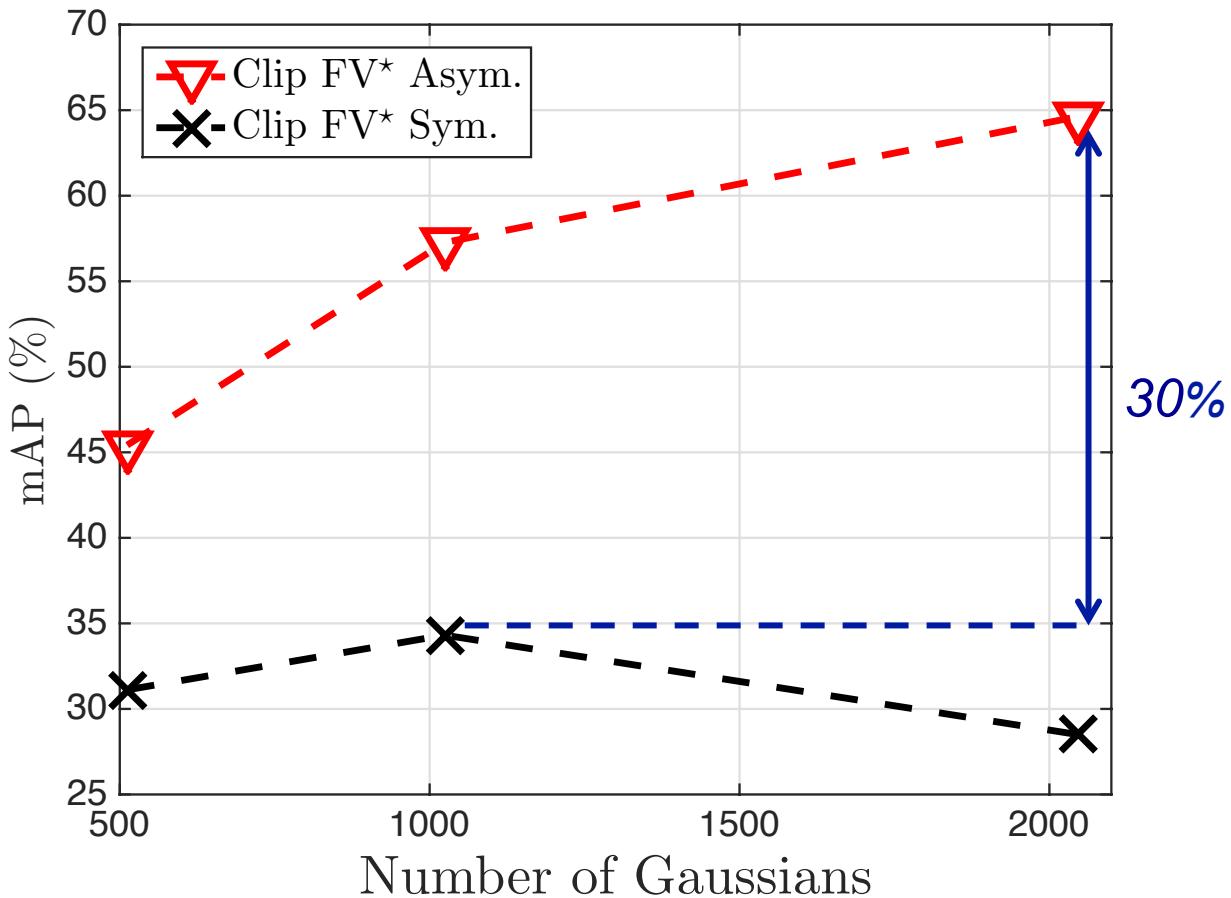
- 282 queries (smartphone pics)
- 2.7 minutes/clip
- 50.6 shots/clip
- Versions
 - 600k frames, 164h, 3.4k clips
 - 4M frames, 1,079h, 24.3k clips

- 258 queries (slides)
- 8.2 minutes/clip
- 58.8 shots/clip
- Versions
 - 600k frames, 169h, 1.1k clips
 - 1.5M frames, 408h, 2.9k clips

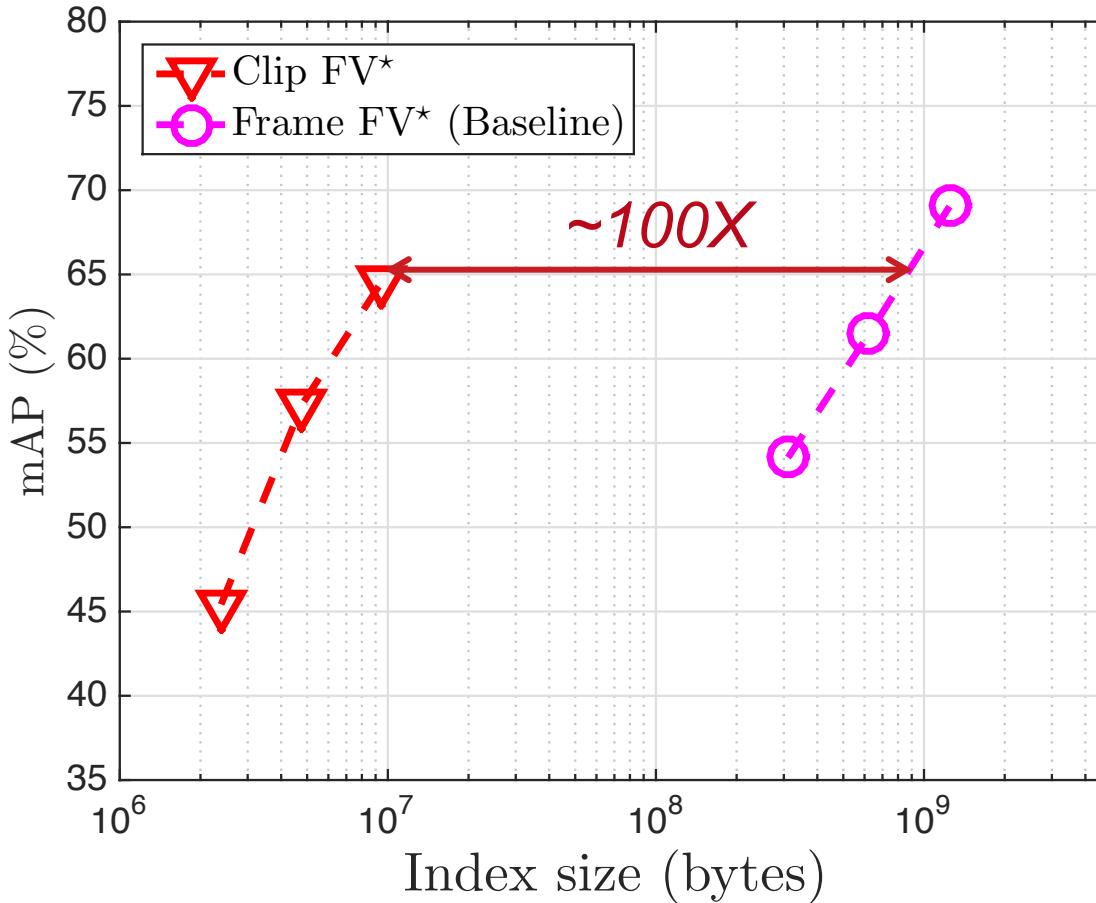
Experiments: Comparison of Techniques



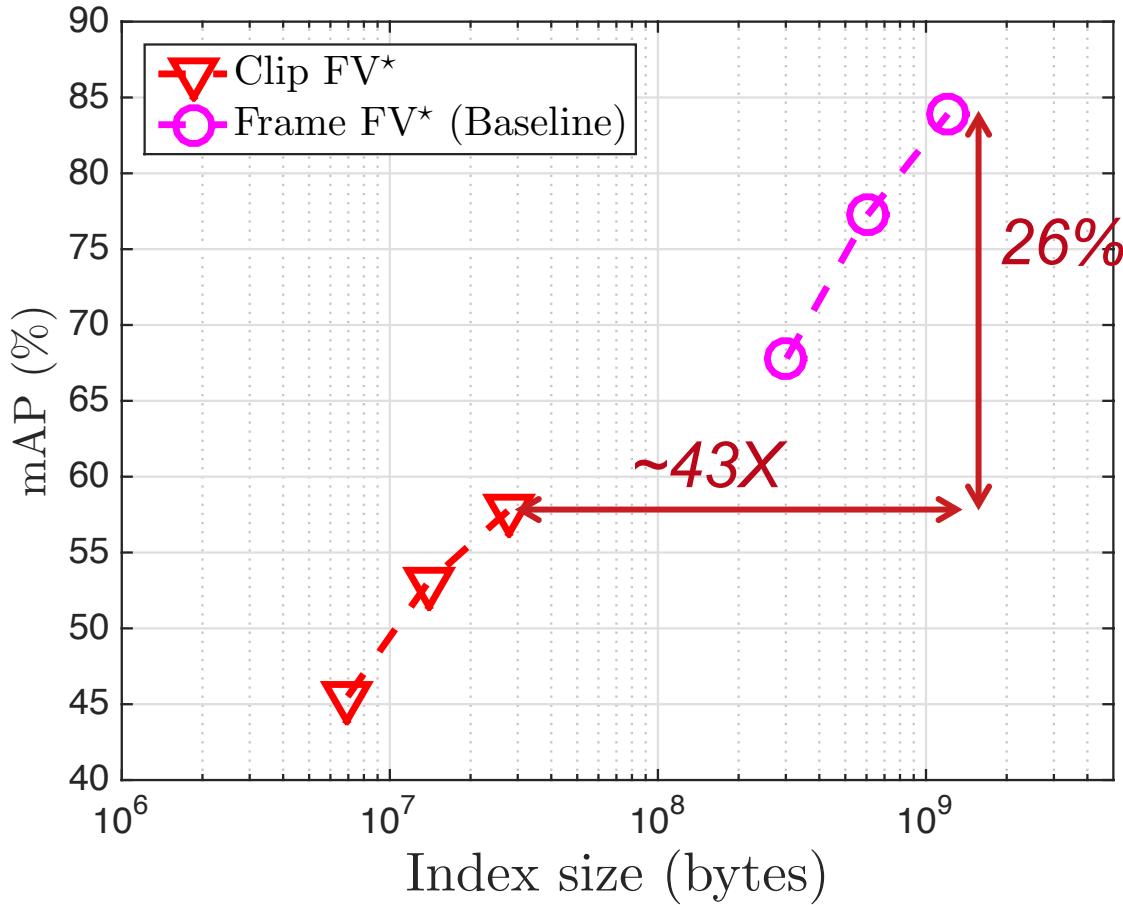
Experiments: Lecture Videos Dataset



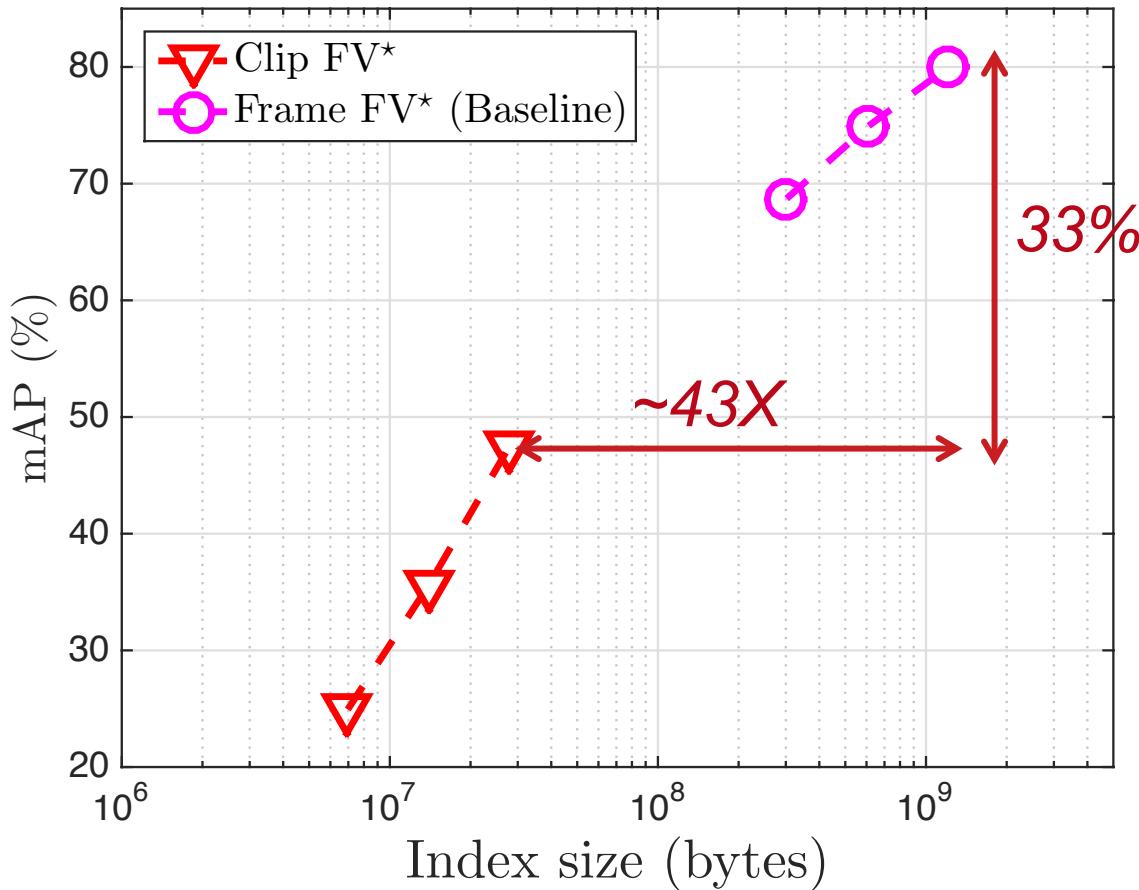
Experiments: Lecture Videos Dataset



Experiments: Video Bookmarking Dataset



Experiments: News Videos Dataset



Contribution 3

Fisher Vector Comparisons

- Asymmetric comparisons for Fisher vectors
- Cluttered query or database images

Fisher Vector Aggregation

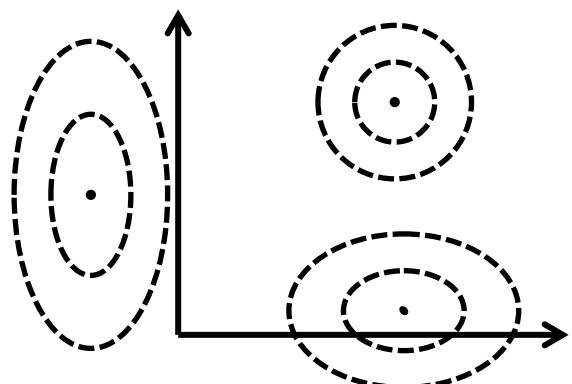
- Fisher vector descriptors for video segments
- Compact database for large-scale retrieval

Bloom Filter Aggregation

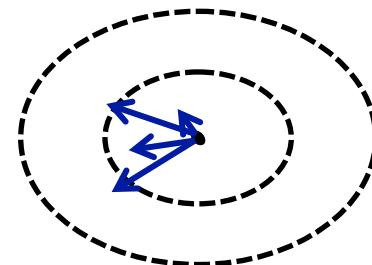
- Bloom filter descriptors for video segments
- Fast and accurate large-scale retrieval

Aggregation Methods

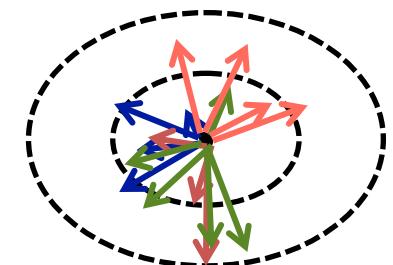
Descriptor space



Frame residuals

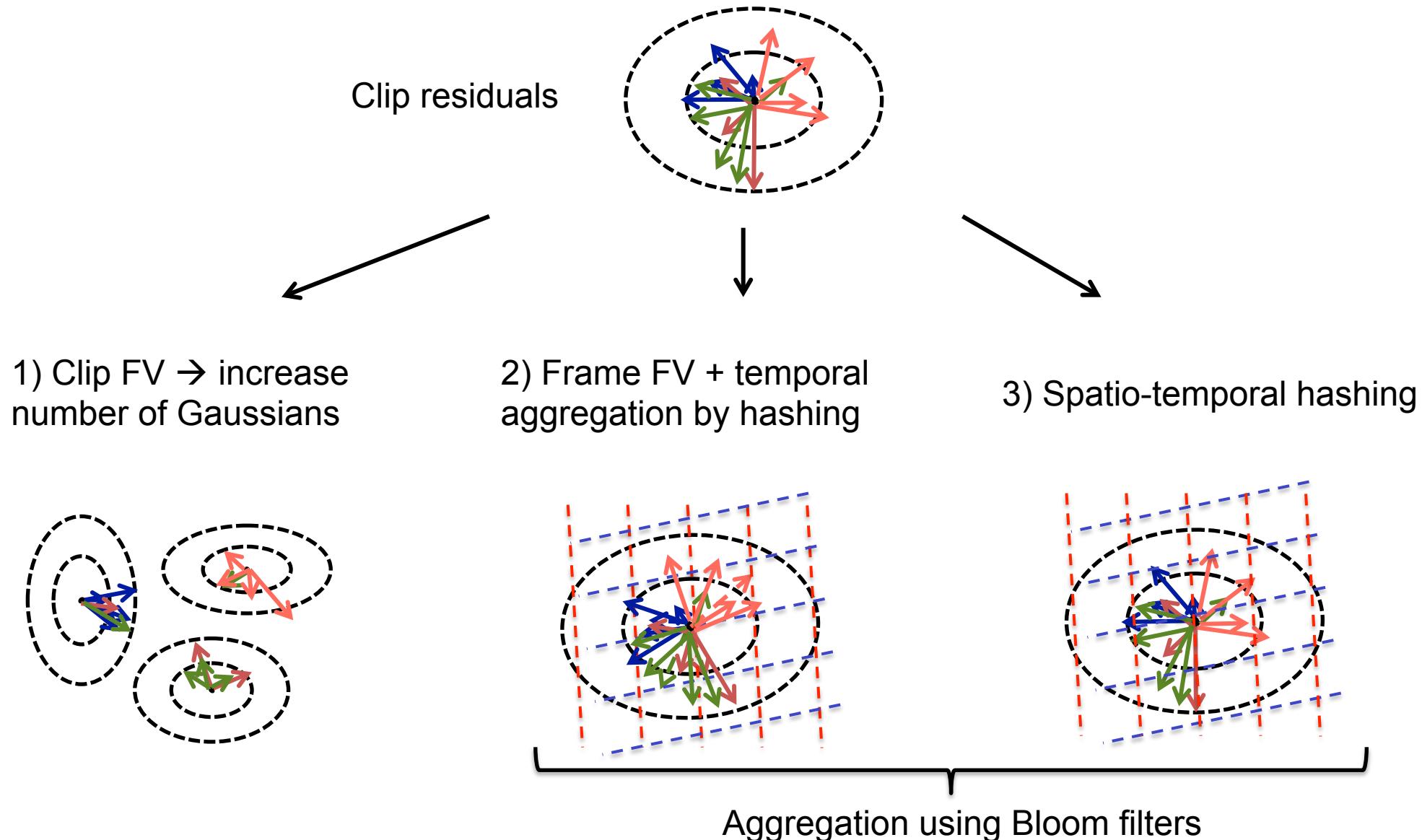


Clip residuals



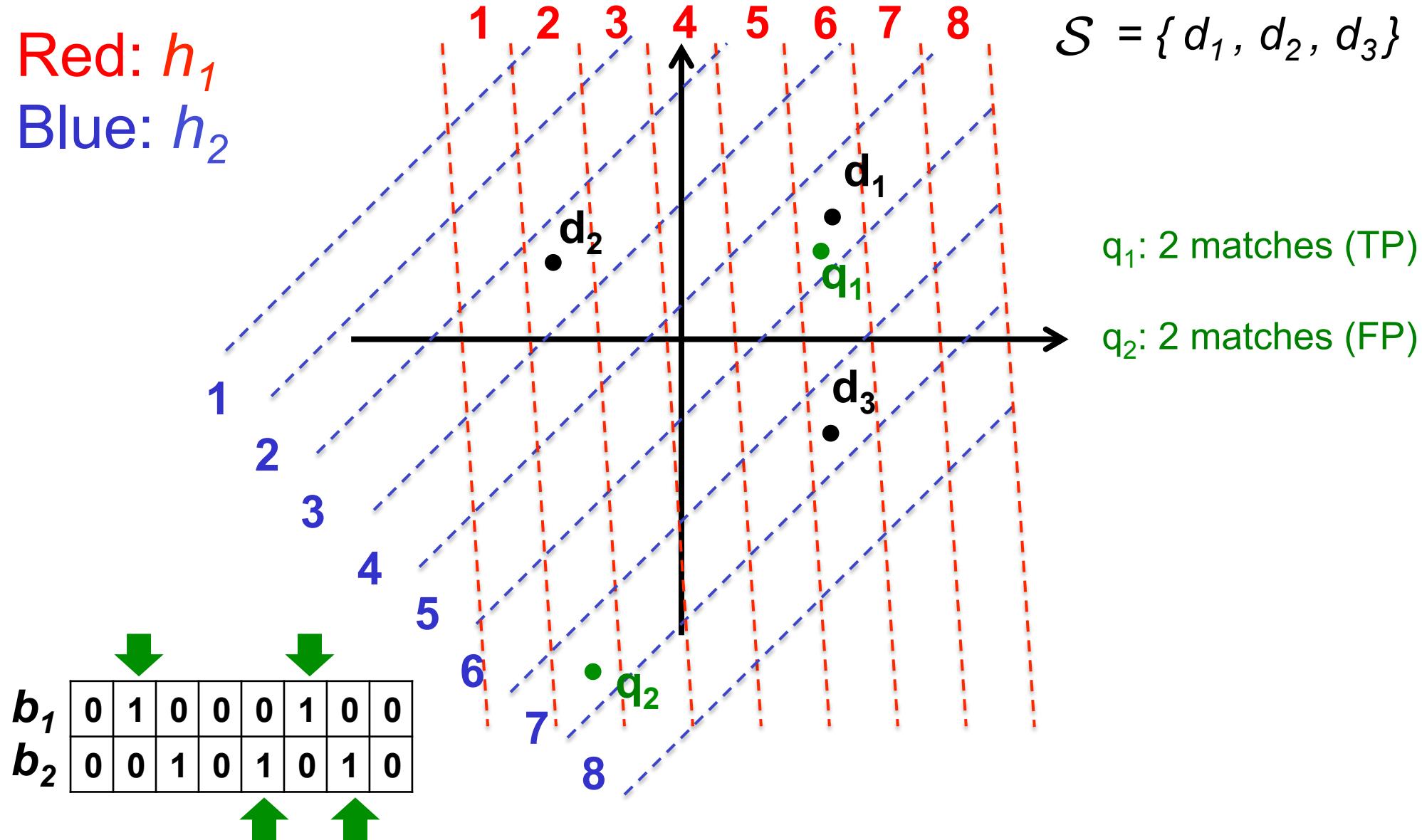
Zoom

Aggregation Methods



Bloom Filter (BF)

Red: h_1
Blue: h_2



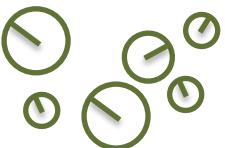
$$\mathcal{S} = \{d_1, d_2, d_3\}$$

q_1 : 2 matches (TP)

q_2 : 2 matches (FP)

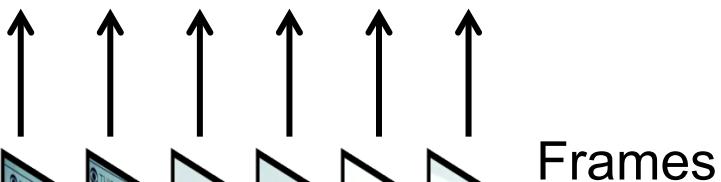
BF using Global Descriptors (BF-GD)

Features per frame



Fisher embedding per feature

0.2	-0.6	0.1	-0.4	...	0	0
0	0	0	0	...	0.4	0.06



Video clip

FV aggregation per frame

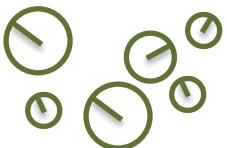
0.1	-0.3	0.05	-0.2	...	0.2	0.03
-----	------	------	------	-----	-----	------

Hash functions
 $h_m(v), m = 1, \dots, M$

Bloom filter

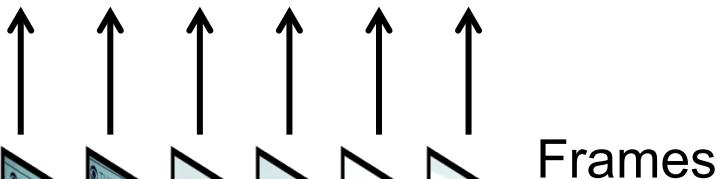
BF using Point-Indexed Descriptors (BF-PI)

Features per frame



Fisher embedding per feature

0.2	-0.6	0.1	-0.4	...	0	0
0	0	0	0	...	0.4	0.06



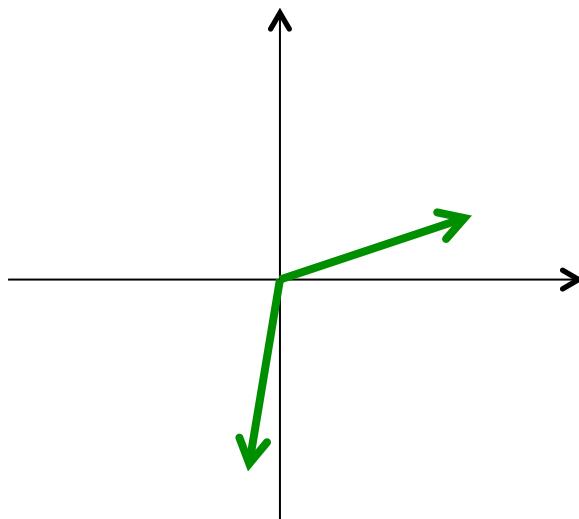
Video clip

Hash functions
 $h_m(v), m = 1, \dots, M$

**Bloom
filter**

Hash Functions

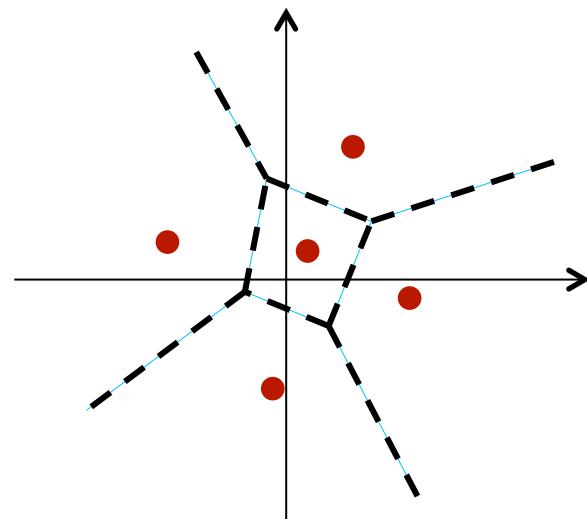
Locality-Sensitive Hashing (LSH)



Random hyperplanes

One hyperplane per bit

Vector Quantization (VQ)

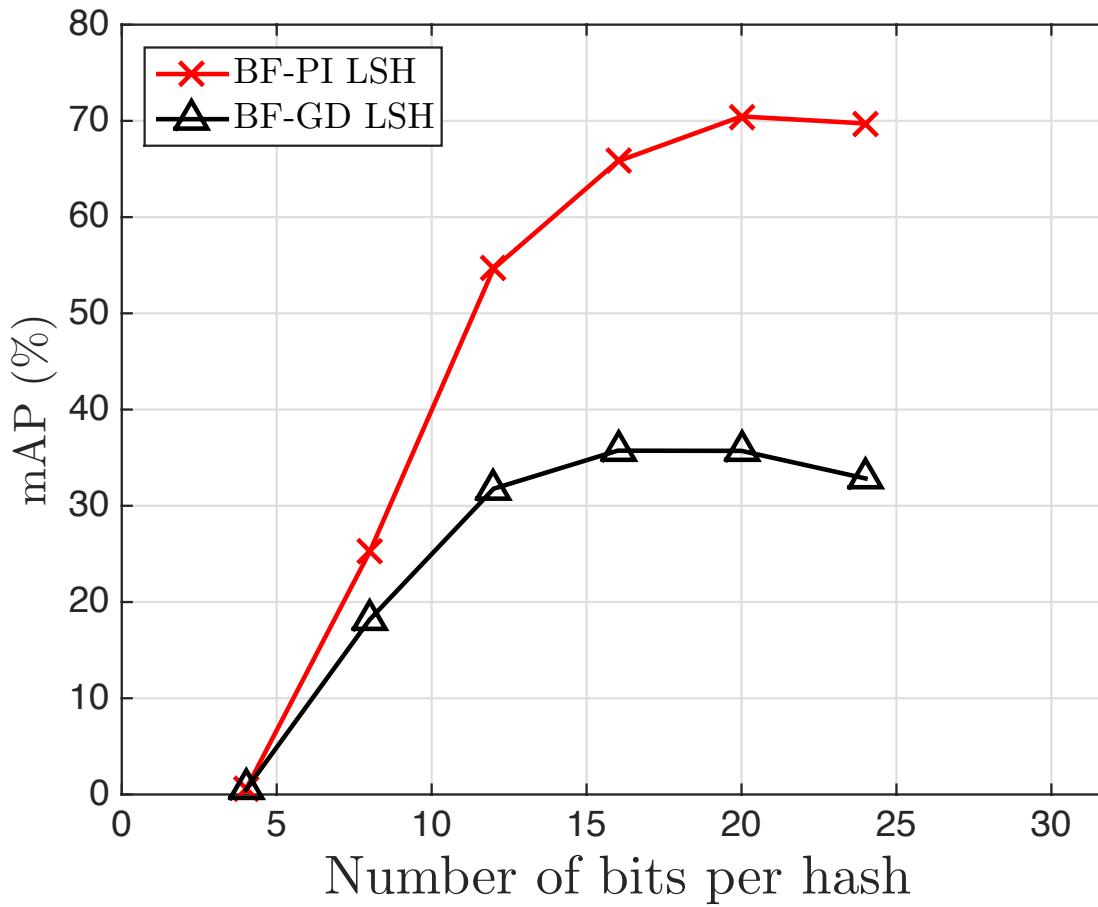


Trained using
Approximate K-Means

bits = $\log_2(\# \text{ centroids})$

Experiments: BF-GD vs BF-PI

Dataset: News Videos – 600k



Visual Invariance

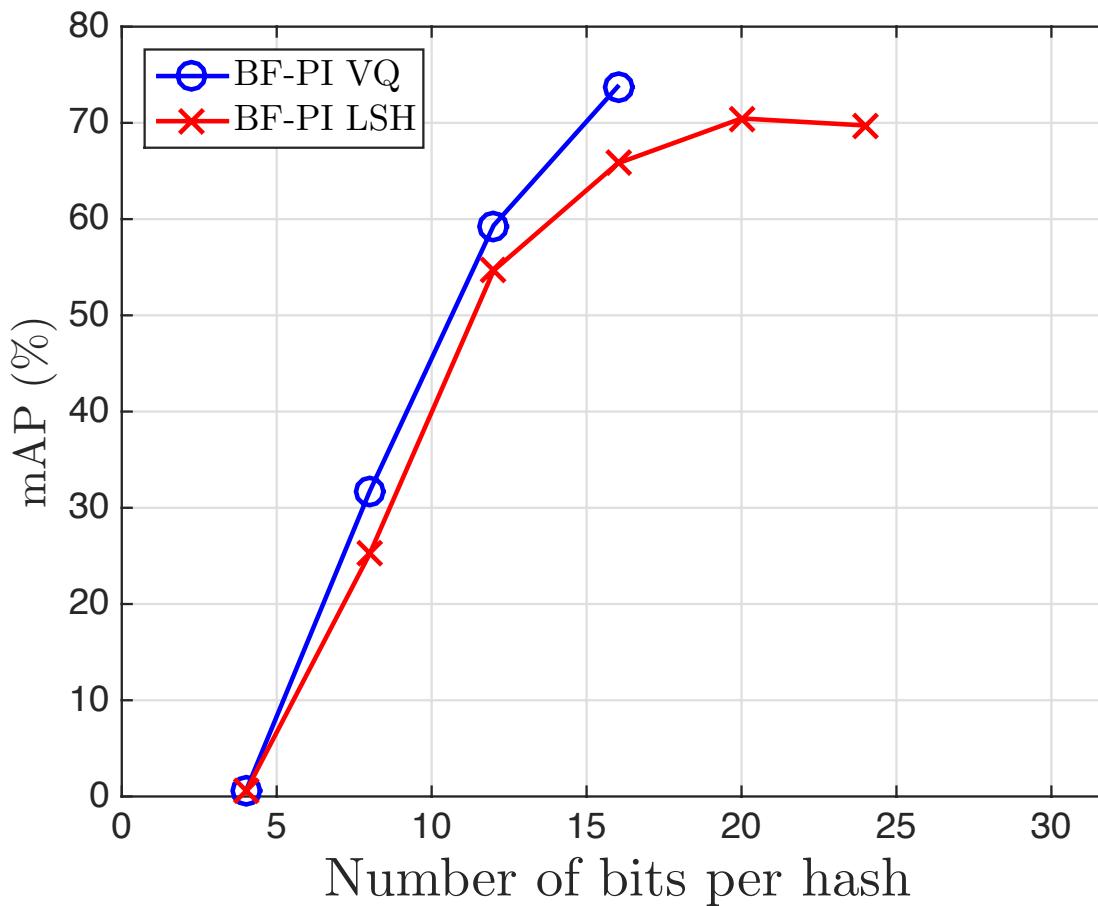


Visual Discriminativeness

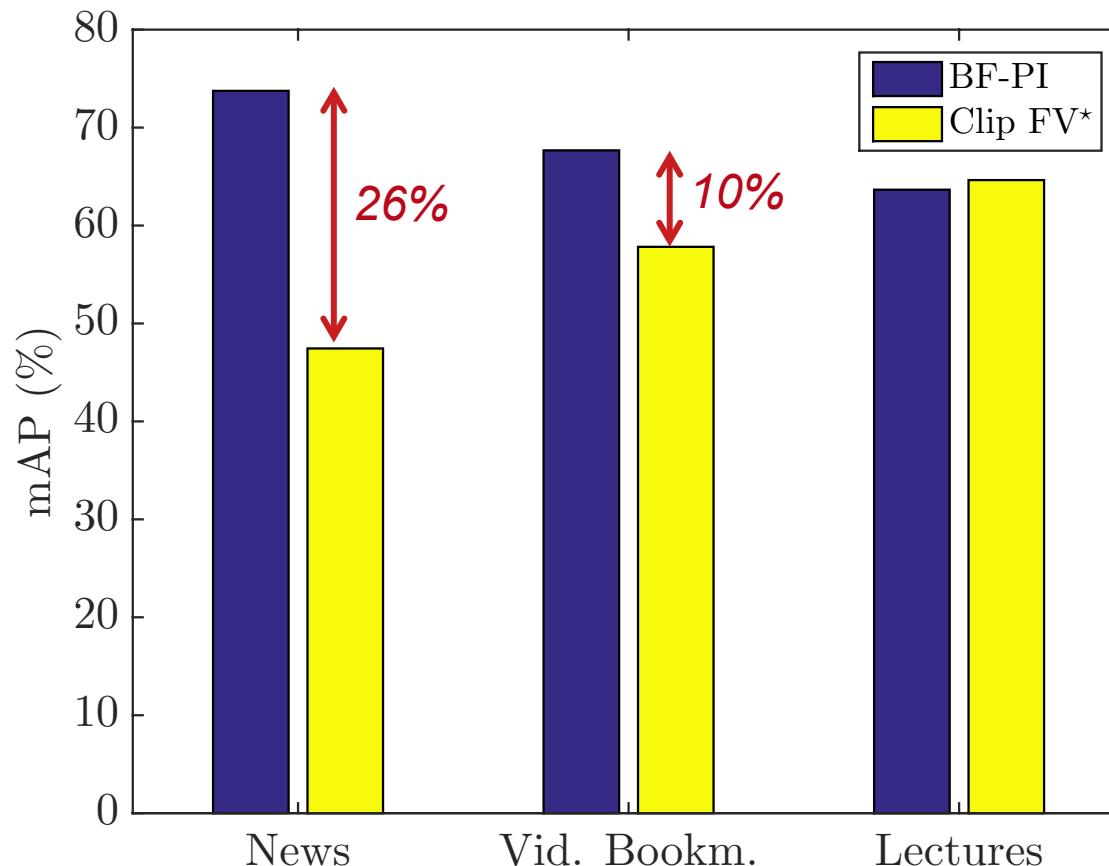


Experiments: BF-PI with Different Hashes

Dataset: News Videos – 600k

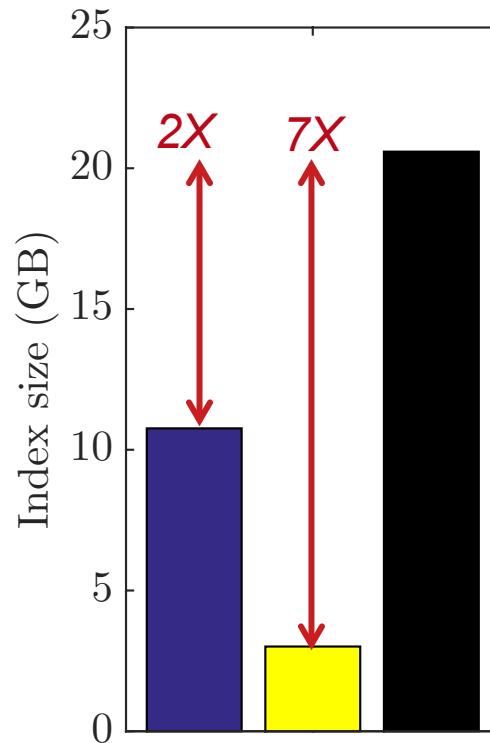
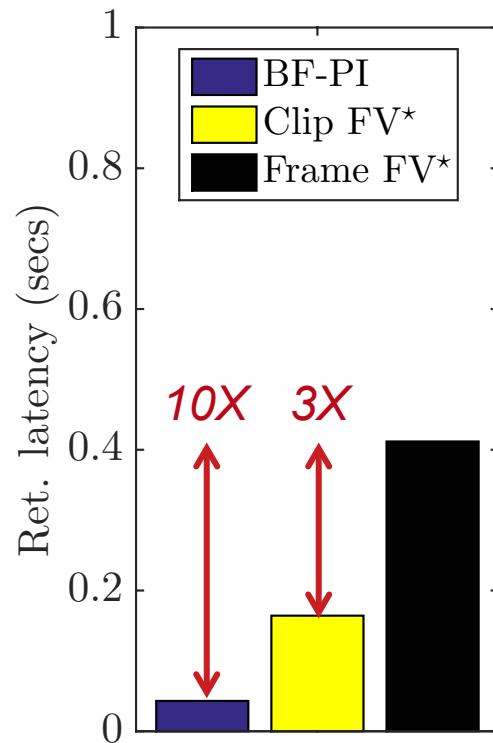
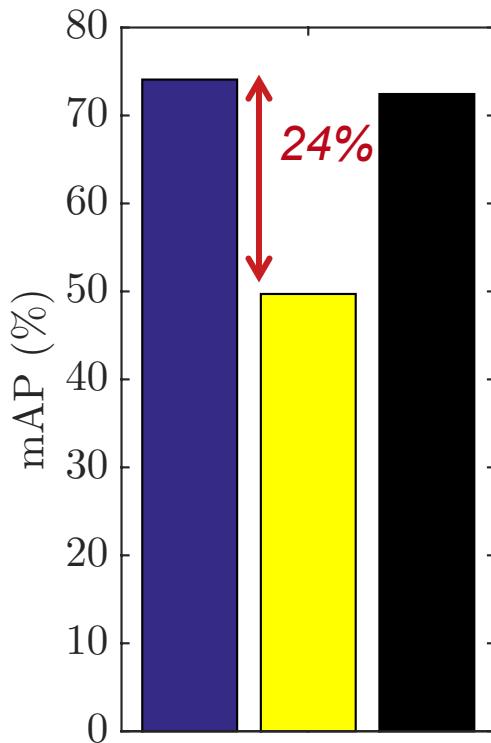


Experiments: Results on 600k Datasets



Experiments: Large-Scale with Re-Ranking

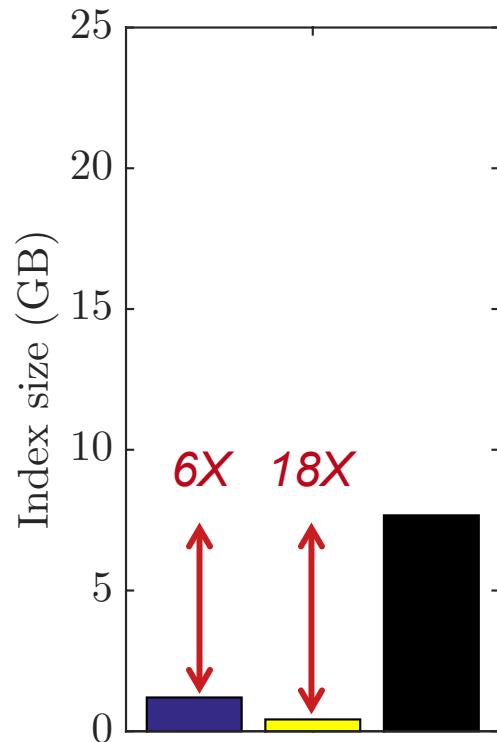
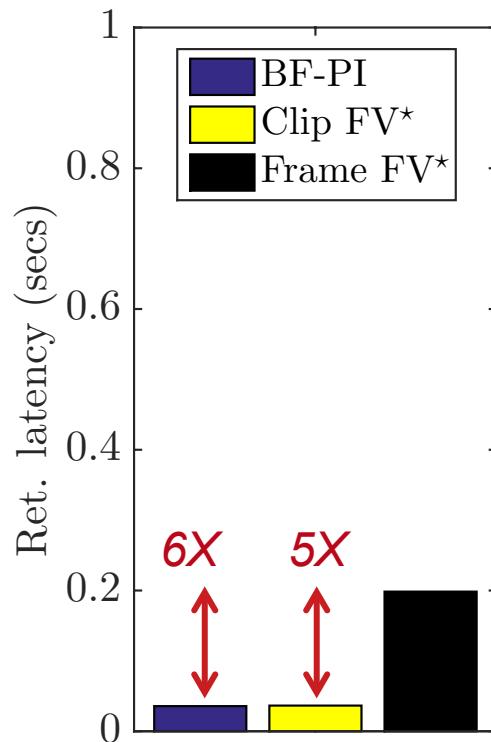
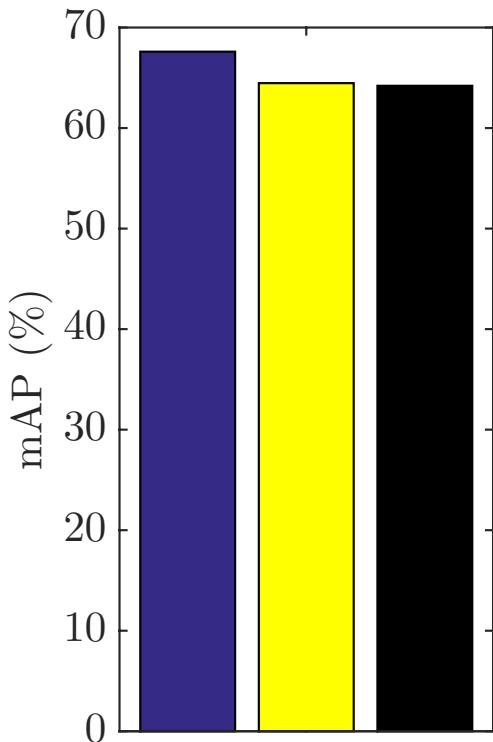
Dataset: News Videos – 4M



BF-PI and *Clip FV** results are re-ranked using *Shot-FV** descriptors

Experiments: Large-Scale with Re-Ranking

Dataset: Lecture Videos – 1.5M



BF-PI and *Clip FV** results are re-ranked using *Shot-FV** descriptors

Conclusions

Fisher Vector Comparisons

- Asymmetric comparisons by projecting cluttered FVs
- Studied two asymmetric retrieval problems
- Large retrieval gains (up to 25% mAP) in both cases

Fisher Vector Aggregation

- Fisher vector aggregation over video segments
- Simple aggregation outperforms other techniques
- Effective retrieval with 100X compression for lectures dataset

Bloom Filter Aggregation

- Bloom filter aggregation over video segments
- Studied hash functions and spatio-temporal aggregation schemes
- Lighter, faster than frame-based schemes, with similar accuracy