

# Multi-modal query expansion for video object instances retrieval

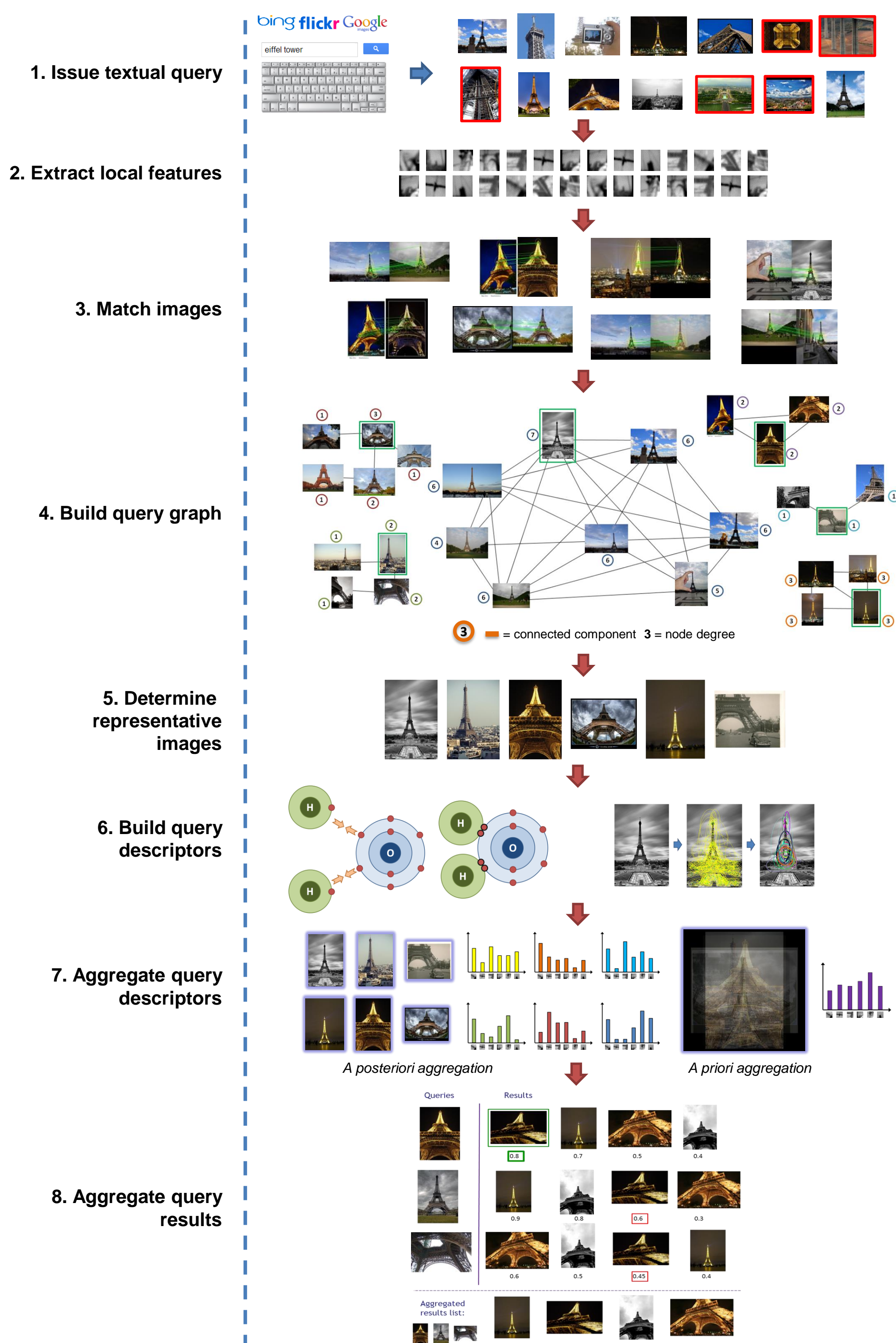
## Authors

Andrei BURSUC  
Titus ZAHARIA

## Objectives

- Retrieve object instances from a large video repository starting from minimum, user-provided textual information
- Leverage on users' affinity for textual queries and crawl images from the Internet
- Remove outliers from retrieved data and identify representative instances for the topic given by the user
- Build visual descriptors from filtered representative instances and use them for querying the video repository

## Approach overview



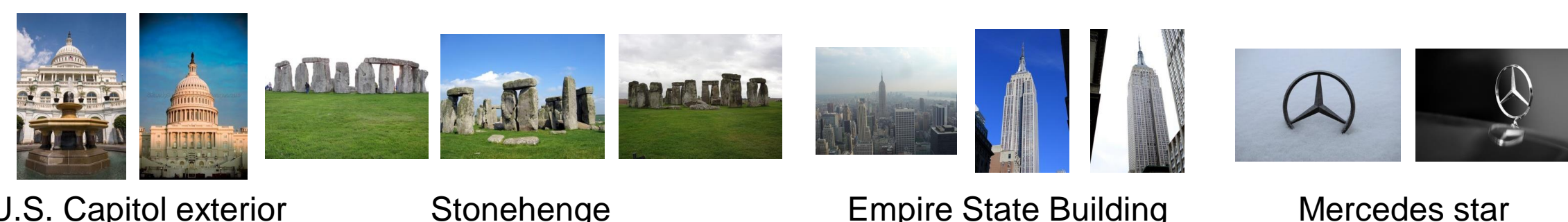
## Results

### Retrieval performance (mean Average Precision)

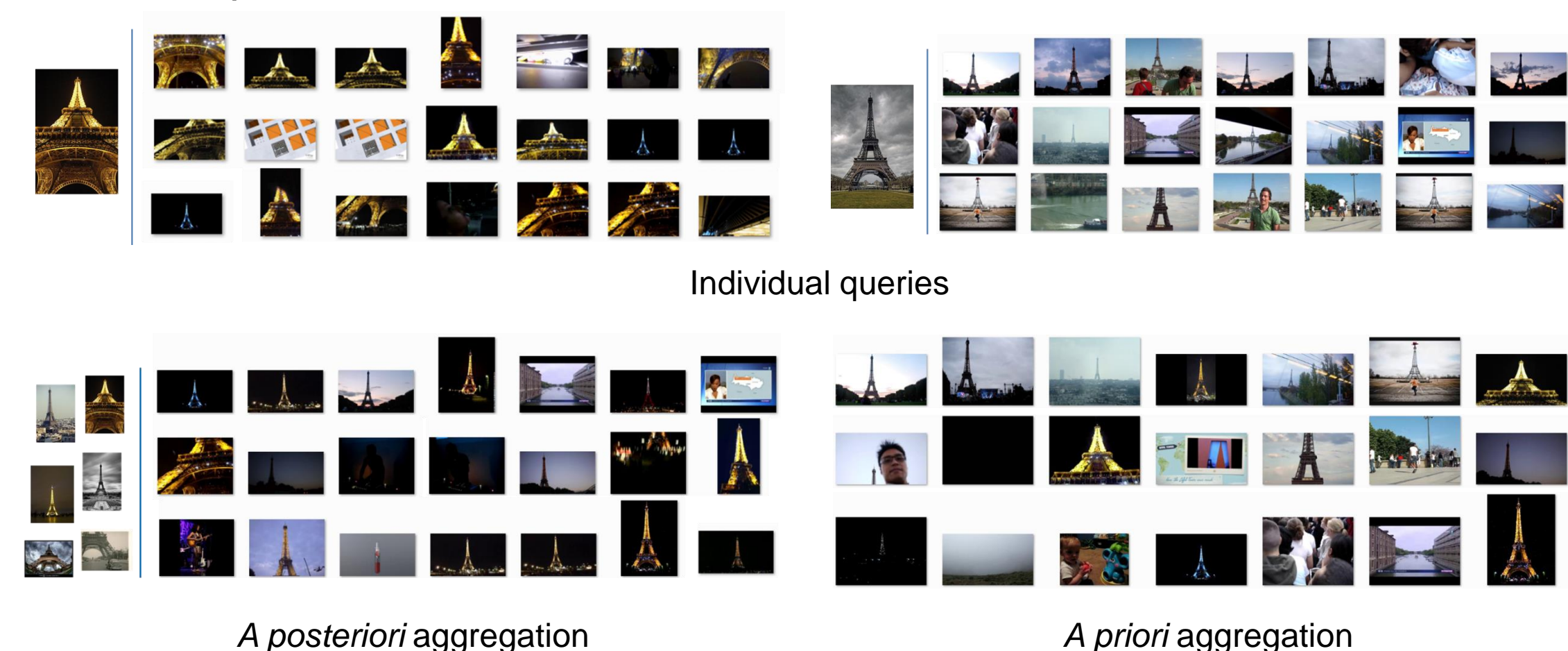
- Centered representative query:** Consists of the union of all points from the representative image that have been matched/shared with at least one neighboring image from the query graph.
- Distributed representative query:** Consists of the union of all points from every neighboring image that have been matched with points from the representative image

Expansion method	Number of mined images	Aggregation strategy	mean Average Precision
Centered representative query	25	<i>A posteriori</i>	0.0455
		<i>A priori</i>	0.0476
	50	<i>A posteriori</i>	0.0585
		<i>A priori</i>	0.0583
	100	<i>A posteriori</i>	0.0689
		<i>A priori</i>	0.0688
Distributed representative query	25	<i>A posteriori</i>	0.0540
		<i>A priori</i>	0.0558
	50	<i>A posteriori</i>	0.0756
		<i>A priori</i>	0.0787
	100	<i>A posteriori</i>	0.0871
		<i>A priori</i>	<b>0.0967</b>
TRECVID 2012 Median mean Average Precision			0.0795
Baseline Bag-of-Words			0.095

### Representative images

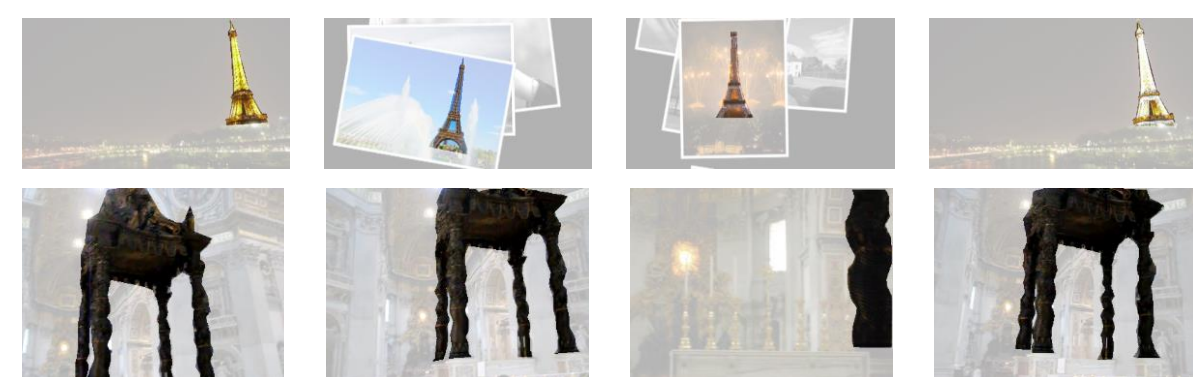


### Examples



## Evaluation

- Trecvid 2012 Instance Search Task Flickr dataset
- 74,958 videos mined from Flickr
- 22 query topics with up to 9 example images with precise object annotation and basic textual description: 102 query images
- Hessian Affine regions + RootSIFT descriptors from 683,433 keyframes
- Bag-of-Words with vocabulary of 1M visual words



Eiffel Tower

Baldachin in Saint  
Peter's Basilica

## Conclusion and perspectives

- Novel multi-modal query definition and expansion method:  
text  $\Rightarrow$  image  $\Rightarrow$  video
- Good object retrieval performance even when using only textual data
- Distributed query descriptors with *a priori* aggregation provide better results while reducing the number of query operations
- Extend method for multiple Internet sources
- Use an ad-hoc SVM classifier on representative images
- Integrate other image metadata for validating positive instances (geotags, image popularity, uploader reputation)