

# Visual Landmark Recognition from Internet Photo Collections: A Large-Scale Evaluation

Tobias Weyand, Bastian Leibe

*Computer Vision Group  
RWTH Aachen University  
Germany*

---

## Abstract

The task of a visual landmark recognition system is to identify photographed buildings or objects in query photos and to provide the user with relevant information on them. With their increasing coverage of the world's landmark buildings and objects, Internet photo collections are now being used as a source for building such systems in a fully automatic fashion. This process typically consists of three steps: clustering large amounts of images by the objects they depict; determining object names from user-provided tags; and building a robust, compact, and efficient recognition index. To this date, however, there is little empirical information on how well current approaches for those steps perform in a large-scale open-set mining and recognition task. Furthermore, there is little empirical information on how recognition performance varies for different types of landmark objects and where there is still potential for improvement. With this paper, we intend to fill these gaps. Using a dataset of 500k images from Paris, we analyze each component of the landmark recognition pipeline in order to answer the following questions: How many and what kinds of objects can be discovered automatically? How can we best use the resulting image clusters to recognize the object in a query? How can the object be efficiently represented in memory for recognition? How reliably can semantic information be extracted? And finally: What are the limiting factors in the resulting pipeline from query to semantics? We evaluate how different choices of methods and parameters for the individual pipeline steps affect overall system performance and examine their effects for different query categories such as buildings, paintings or sculptures.

*Keywords:* landmark recognition, image clustering, image retrieval, semantic annotation, compact image retrieval indices

---

## 1. Introduction

Recognizing the object in a photo is one of the fundamental problems of computer vision. One generally distinguishes between object categorization and specific object recognition. Object categorization means recognizing the *class* that an object belongs to, *e.g.* painting or building, while specific object recognition means recognizing a specific object *instance*, such as the Mona Lisa or the Eiffel Tower. In this paper, we consider the latter task, *i.e.*, specific object recognition. In particular, we are interested in two applications, namely photo auto-annotation and mobile visual search. A photo auto-annotation system recognizes objects in a user's photo

albums and labels them automatically, saving the user the effort of manually labeling them. A mobile visual search system provides a user with information on an object that they took a picture of with their smartphone. Because a large part of the photos in these applications are typically tourist photos, many of the objects that such systems need to recognize are landmarks. Therefore, the problem is typically referred to as *landmark recognition*. However, many other types of objects, such as paintings, sculptures or murals, can also be recognized by such systems.

The first step of building a landmark recognition system is to compile a database consisting of one or more photos of each object that shall be recognized. However, since the number of objects that can possibly appear in a user's photos is virtually infinite, it is impossible to con-

---

*Email addresses:* weyand@vision.rwth-aachen.de (Tobias Weyand), leibe@vision.rwth-aachen.de (Bastian Leibe)

struct and maintain such a database by hand. An elegant solution is to build the database from the data it is meant to be applied to, namely public photos from Internet photo collections such as Flickr, Picasa or Panoramio. This approach has several attractive properties: (i) Objects are discovered in an unsupervised, fully automatic way, making it unnecessary to manually create a list of objects and collecting photos for each of them. (ii) The resulting set of objects is likely to be much better adapted to the queries a photo auto-annotation or visual search system might receive than a hand-collected set of objects. (iii) The level of detail of object representation is automatically adapted to the demand. The most popular objects will be represented by the most photos in the database, increasing their chance of successful recognition, while only little memory is used on less popular objects. This approach has gained popularity in the research community [1, 2, 3, 4, 5] and is also being used in applications such as Google Goggles [6].

Constructing such landmark recognition systems on a large scale involves three main research problems: (i) Finding interesting structures in internet image collections, (ii) automatically connecting them to associated semantic content by examining user-provided titles and tags, and (iii) compactly representing them for efficient retrieval. While each of those problems has already been studied in isolation, so far there has not been a systematic evaluation of all three aspects in the context of a fully automatic pipeline. Image retrieval approaches like [7, 8] find matching images for a query, but have no notion of the semantics of the depicted content. Tag mining approaches [3, 9, 4, 10] try to find a description of an image cluster, but so far the large effort to evaluate tag quality has prevented quantitative evaluations in a large-scale setting. Landmark object discovery approaches [1, 11, 4, 5, 6] aim at finding interesting buildings and other objects, but no systematic evaluation has been performed that analyzes what types of objects can be discovered and how recognition performance varies with these types. Furthermore, it is still largely unclear what is the best strategy to determine the identity of the recognized object based on the set of retrieved database images (which becomes a non-trivial problem whenever image clusters overlap and images may contain multiple landmark objects).

In this paper, we evaluate the whole process of constructing landmark recognition engines from Internet photo collections. To do this in a realistic large-scale setting we require a dataset containing thousands of objects. Moreover, in order to create a realistic application scenario, the target database objects should not be specified by hand (as in many other datasets), but

should be mined automatically. Last but not least, the dataset should not be limited to buildings, but also contain smaller objects, such as paintings or statues.

Our evaluation is based on the PARIS 500K dataset [12] containing 500k photos from the inner city of Paris, which was mined from Flickr and Panoramio using a geographic bounding region rather than keyword queries to obtain a distribution unbiased towards specific landmarks. Thus, in contrast to other common datasets, there is no bias on tag annotations or content. In order to evaluate landmark recognition in a realistic setting, we additionally collected a query set of almost 3,000 Flickr images from Paris that is disjoint from the original dataset. Our evaluation thus mimics the task of photo auto-annotation where a photo uploaded to a photo sharing website is automatically annotated with the object it depicts.

To evaluate the performance of landmark recognition, we use a recent landmark discovery algorithm [5] to discover landmarks in the dataset. We created an exhaustive ground truth for the relevance of each of the discovered landmarks with respect to each of the 3,000 queries, which involved significant manual effort. This is the first ground truth for evaluating landmark recognition on an unbiased and realistic dataset. To enable the comparison of other approaches with the ones evaluated in this paper, the ground truth is publicly available.

To give a detailed performance analysis for different types of objects, we introduce a taxonomy for the objects landmark recognition systems are able to recognize. Throughout our evaluation, we report both summary performances over the entire database and detailed findings for different object categories that show how their recognition is affected by the different stages of the system. As our results show, the observed effects vary considerably between query categories, justifying this approach. We give detailed results for each category, and use the four use cases of *Landmark Buildings*, *Paintings*, *Building Details* and *Windows* as representatives for different challenges. The taxonomy is available along with the ground truth.

Note that our goal is not primarily to propose novel methods (although some of the methods evaluated in Sec. 6 and Sec. 7 are indeed novel), but to provide answers to the following questions:

- How many and what kinds of objects are present in Internet photo collections and what is the difficulty of discovering objects of different landmark types (Sec. 5)?
- How to decide which landmark was recognized given a list of retrieved images (Sec. 6)?

- How to efficiently represent the discovered objects in memory for recognition (Sec. 7)?
- Are the user-provided tags reliable enough for determining accurate object names (Sec. 8)?
- Given the entire retrieval, recognition and semantic labeling pipeline, what are the factors effectively limiting the recognition of different object categories (Sec. 9)?

Our analysis provides several interesting insights, for example:

- Semantic annotation is the main bottleneck for system performance. In many cases, the correct object is visually recognized, but the name of the object cannot be determined due to the sparsity and amount of noise of user-provided image titles and tags.
- Different bottlenecks exist for different object categories. For example, *Murals* are easy to recognize using the standard visual words pipeline, but reliable semantic information is often missing for them. For other objects like museum exhibits, the opposite is the case: While semantic information is readily available, they are hard to recognize visually due to their spatial structure and scarce visual examples.
- When the desired application is building recognition, a seeding-based clustering method can bring significant computational savings, since buildings are already discovered when using few seeds, while smaller objects require orders of magnitude more seeds.
- Different techniques for compactly representing object clusters are optimal for different object types.

As a result of this evaluation, we can identify several interesting directions, where progress can still be made.

## 2. Engine Architecture

The architecture of a typical landmark recognition engine such as [1, 13, 4, 6] is shown in Fig. 1. Large amounts of tourist photos are clustered, resulting in a set of *objects*. By *object*, we denote a cluster of images that show the same entity. We will refer to the images in each object cluster as its *representatives*. Since the clusters may overlap, a representative can belong to multiple objects. Each object is then associated with *semantics* (typically its name), *e.g.*, by mining frequently used image tags. The set of representatives for each cluster

is often decimated by eliminating redundant images in order to save memory and computation time. To recognize the object in a query image, a visual search index [14, 8, 15] containing all *representatives* is queried, producing a ranked list of matches. Based on this list, *objects* are ranked w.r.t. their relevance to the query and the corresponding *semantics* are returned.

In this paper, we evaluate different choices for the components of this framework and demonstrate how they affect the system’s overall performance. Sec. 5 considers the stage of determining a set of objects by clustering images from internet photo collections and shows how many objects from which categories can be discovered. Given a ranking of the representatives for a query, Sec. 6 analyzes different schemes for determining the object shown in the query image. In Sec. 7 we consider different ways of speeding up search and reducing memory requirements by removing redundant representatives. Finally, in Sec. 8 we analyze the stage of semantic annotation based on frequent tags and perform an end-to-end analysis of the performance of the whole pipeline from query to semantics.

## 3. Related Work

We now give an overview how the individual parts of the pipeline introduced above have been approached in previous work.

### 3.1. Datasets

We created our own query set and ground truth for this paper, because available benchmarks do not support such an evaluation. Most datasets only cover very few, mostly building-scale, landmarks (*e.g.*, EUROPEAN CITIES 1M [1], STATUE OF LIBERTY, NOTRE DAME and SAN MARCO [16], OXFORD BUILDINGS [8], PARIS BUILDINGS [17]). Another problem is that their ground truths are designed for other tasks. Image retrieval datasets (*e.g.* OXFORD BUILDINGS [8], PARIS BUILDINGS [17], INRIA HOLIDAYS [18]) are not suitable for our evaluation, because we want to evaluate object recognition, *i.e.* recognizing the object(s) in a query image, and not image retrieval, *i.e.* retrieving images similar to a query from a database. Image-based localization datasets (AACHEN [19], VIENNA [20], DUBROVNIK and ROME [21]) evaluate how accurately the camera pose of the query image can be estimated. While this is more related to our problem, our goal differs from pose estimation, because camera pose does not necessarily determine what object the camera is really seeing (See Sec. 3.3 for more details.) The SAN FRANCISCO [22] and LANDMARKS 1K [23] datasets are

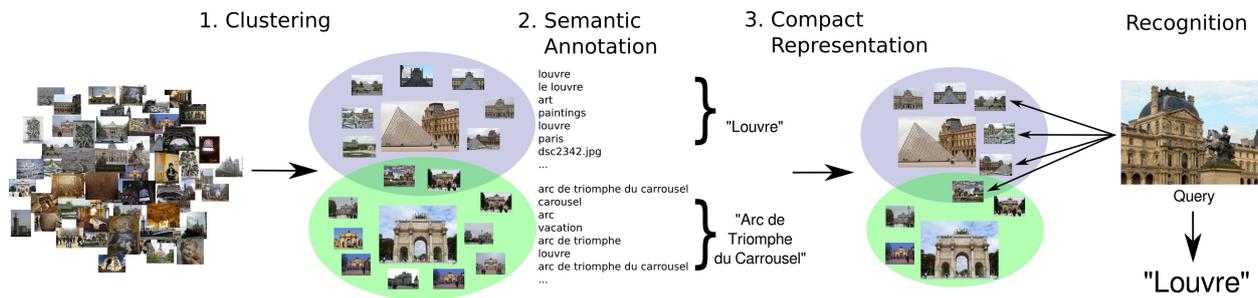


Figure 1: The architecture of a landmark recognition system. 1. *Objects* are discovered by visually clustering touristic photos (Sec. 5). We call the photos in an object cluster its *representatives*. 2. *Semantic* annotations are mined from user-provided tags (Sec. 8). 3. The search index is made more compact by eliminating redundancy (Sec. 7). 4. The object in a query photo is recognized by retrieving similar photos and exploiting the knowledge of their cluster memberships (Sec. 6).

closest to our requirements, but both of them focus on large, building-level landmarks while we are explicitly interested in also evaluating the recognition of smaller, non-building objects.

### 3.2. Landmark Discovery

Landmark Recognition Engines are typically based on a visual search index from objects discovered in Internet photo collections [1, 13, 4, 6]. The underlying *landmark discovery* approaches perform visual [24, 16, 25, 4, 5, 26] or geographical [9, 27] clustering, or a combination of the two [1, 6, 28]. Zheng *et al.* [6] show that online tourist guides can be a valuable additional data source, and Gammeter *et al.* [13] use descriptions determined from user-provided tags to search for additional images on the web. In this work however, we focus on methods based solely on the images from Internet photo collections and their metadata.

Chum *et al.* [24] use *Min-Hash* to find *seed* images and grow landmark clusters by query expansion [7]. Philbin *et al.* [25] over-segment the matching graph using spectral clustering and merge clusters of the same object based on image overlap. Gammeter *et al.* [2] and Quack *et al.* [4] perform hierarchical agglomerative clustering in a local matching graph. Avrithis *et al.* [1] use Kernel Vector Quantization to create a clustering with an upper bound on intra-cluster dissimilarity. *Iconoid Shift* by Weyand *et al.* [5] finds popular objects at different scales using mode search based on a homography overlap distance. We choose *Iconoid Shift* as our analysis tool, because it produces an overlapping clustering and discovers landmarks at varying levels of granularity, thus also discovering, *e.g.*, building details.

### 3.3. Landmark Recognition

Based on the clusters resulting from landmark discovery, it is now possible to recognize the landmark in

a query image. There are three predominant approaches for this in the literature: Image Retrieval, Classification and Pose Estimation.

**Image Retrieval.** Most *Image Retrieval* based approaches use efficient *specific object retrieval* methods [14, 8, 15] that allow searching for images matching a query in a database consisting of potentially millions of images. Several approaches [4, 6, 2, 8, 15] implement a *best match* strategy (Sec. 6.2) where the query is matched against the database of representatives and the object cluster corresponding to the best match is returned. While Quack *et al.* [4] and Zheng *et al.* [6] perform a precise but computationally expensive direct feature matching, Gammeter *et al.* [2] retrieve images using inverted indexing and bags-of-visual-words (BoVWs) [8, 15]. Li *et al.* [16] only want to decide *whether* the query image contains a specific landmark. Given a dataset of photos of one landmark, they perform image retrieval based on both Gist features and BoVWs and apply a threshold to the retrieval score to decide if the query contains the object. Both Avrithis *et al.* [1] and Johns *et al.* [26] compress the images in a cluster into a joint BoVW representation and perform inverted file retrieval to find the best matching scene models for a query image.

**Classification.** An alternative approach is to view the task as a classification problem where each landmark is a class. Gronat *et al.* [29] learn exemplar SVMs based on the BoVWs of the visual features of the database images. Li *et al.* [27] learn a multi-class SVM and additionally use the BoWs (bags-of-words) of the textual tags of the images as features. Bergamo *et al.* [30] use a similar approach, but perform classification using 1-vs-all SVMs. Instead of using approximate k-Means [8] for feature quantization, they reconstruct the landmarks using structure-from-motion and train random forests

on the descriptors of each structure-from-motion feature track. These random forests are then used for quantizing descriptors. While discriminative methods often yield higher accuracy than nearest neighbor matching, they also have disadvantages. For example, they assign every image a landmark label regardless of whether it contains a landmark. Moreover, discriminative models need to be re-trained every time new images and landmarks are added.

**Pose Estimation.** The goal of pose estimation is to determine the camera location and orientation for a given query image. There are several approaches for solving this task by matching the query against street level imagery such as Google Street View panoramas [22, 31, 32, 33, 34, 35] using local feature based image retrieval [14, 8, 15]. Other approaches are based on 3D point clouds created by applying structure-from-motion on Internet photo collections or manually collected photos [36, 37, 21, 23]. Since image retrieval methods cannot be applied here, these approaches directly match the query descriptors against the descriptors of the image features that the 3D points were reconstructed from. After a set of 2D-3D correspondences has been established, the camera pose is determined by solving the perspective-n-point (PnP) problem [38]. Since the descriptor matching problem becomes computationally expensive when matching against very large 3D models, hybrid methods have been proposed that [39, 20, 19] first perform efficient image retrieval using inverted files and then solve the PnP problem based on the relatively small set of 3D points associated with the 2D features of the retrieved images.

It is important to realize that camera pose does not necessarily determine what is visible in the image. Even though the camera is in front of a landmark, the landmark might not be visible due to occlusion or the user might be taking a picture of a non-stationary object or an event near that landmark. Moreover, pose estimation relies on either regularly sampled images (*e.g.* Google Street View panoramas), which are not available everywhere, or structure-from-motion reconstructions, which are not always possible to compute robustly.

Because of the disadvantages of *Classification* and *Pose Estimation* based approaches, we focus on *Image Retrieval* based approaches in this evaluation.

### 3.4. Eliminating Redundancy

Several methods have been proposed to reduce the size of the visual search index. An obvious method is to apply standard compression techniques [40], which

reduces memory consumption at the cost of computational efficiency. Instead, we are interested in *eliminating redundancy* already before index construction.

Several works have addressed this problem at the *image level*, *i.e.* by removing redundant images from the index. Li *et al.* [16] summarize the input image collection in a set of *iconic* images by applying k-means clustering based on Gist descriptors, and use only these images to represent a landmark in retrieval. Gammeter *et al.* [13] identify sets of very similar images using complete-link hierarchical agglomerative clustering and replace them by just one image. This step yields a slight compression of the index without loss in performance. Instead of performing clustering Yang *et al.* [41] only determine a set of canonical views by applying PageRank on the matching graph of the image collection. They then discard all other views and match the query only against the canonical views.

Other works have addressed the problem at the *feature level*. Turcot *et al.* [42] perform a full pairwise matching of the images in the dataset and remove all features that are not at least once inliers w.r.t. a homography. They report a significant reduction of the number of features while maintaining similar retrieval performance. Avrithis *et al.* [1] and Johns *et al.* [26] combine the images in a cluster into a joint BoVW representation. Avrithis *et al.* [1] use Kernel Vector Quantization to cluster redundant features and keep only the cluster centers. While this method only yields a slight compression, the aggregation of features into a Scene Map brings significant improvements in recognition performance. Johns *et al.* [26] performs structure-from-motion and summarize features that are part of the same feature track. Gammeter *et al.* [2] estimate bounding boxes around the landmark in each image in a cluster and remove every visual word from the index that never occurs inside a bounding box. This is reported to yield an index size reduction of about a third with decreasing precision.

There is also work in *pose estimation* that aims to eliminate redundancy in the dataset. In their hybrid 2D-3D pose estimation approach, Irschara *et al.* [20] generate a set of synthetic views by projecting the SfM points onto a set of virtual cameras placed at regular intervals in the scene. They then decimate the set of synthetic views using a greedy set cover approach that finds a minimal subset of views such that each view in the subset has at least 150 3D points in common with an original view. Cao *et al.* [43] use a similar criterion, but instead of views, they decimate the set of points in an SfM point cloud used for localization. Instead of set cover, they use a probabilistic variant of the K-Cover

algorithm.

### 3.5. Semantic Annotation

The most common approach to perform *semantic annotation* of the discovered landmark clusters is by statistical analysis of user-provided image tags, titles and descriptions. In order to remove uninformative tags like “vacation”, Quack *et al.* [4] first apply a stoplist and then perform frequent itemset analysis to generate candidate names. These names are verified by querying Wikipedia and matching images from retrieved articles against the landmark cluster. Zheng *et al.* [6] also apply a stoplist and then simply use the most frequent n-gram in the cluster. Crandall *et al.* [9] deal with uninformative tags in a more general way by dividing the number of occurrences of a tag in a cluster by its total number of occurrences in the dataset. Simon *et al.* [10] additionally account for tags that are only used by individual users by computing a conditional probability for a cluster given a tag, marginalizing out the users.

Unfortunately, a much larger problem, also observed by Simon *et al.* [10], exists for the task of semantic assignment that is much harder to fix: For most clusters accurate tags are simply not available. In our analysis (Sec. 8), we will show for which clusters these methods will still result in accurate descriptions and point out the sources of this problem.

## 4. Evaluation Setup

### 4.1. Dataset

Our evaluation is based on the PARIS 500k dataset [12] consisting of 500k images of the inner city of Paris collected from Flickr and Panoramio. In contrast to many other datasets [16, 8] the images were retrieved using a geographic bounding box query rather than keyword queries to ensure an unbiased distribution of touristic photos.

### 4.2. Query set, Categories and Evaluation

To collect realistic queries for the task of automatic annotation of photos uploaded to a photo sharing website, we downloaded 10k images from the same geographic region as PARIS 500k from Flickr and ensured that they were no (near) duplicates of any image in the original dataset. Since we consider the task of *specific object recognition*, not *object categorization*, we filtered out unsuitable queries like food, pets, plants or cars. To only include objects that have a *chance* of being recognized based on the PARIS 500k dataset, we also exclude

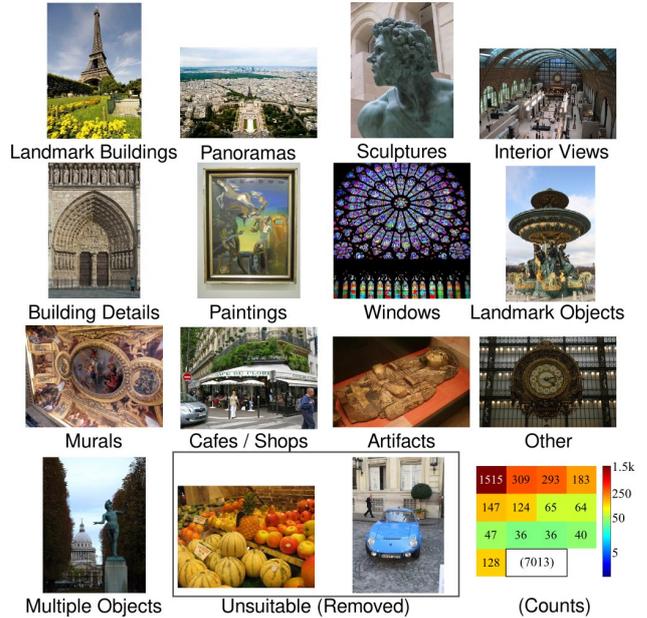


Figure 2: Our set of query categories and the number of query images in each category (bottom right).

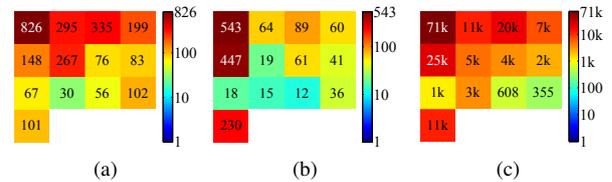


Figure 3: (a) Number of Iconoid clusters in each category. (b) Average cluster size. (c) Total number of images in clusters. Categories are in the same order as in Fig. 2.

queries that do not match *any* image in PARIS 500k, leaving 2,987 queries. We manually grouped the queries into the categories of Fig. 2 in order to enable an analysis of the recognition performance for each query type. We summarize non-building objects such as bridges, fountains or columns under the *Landmark Objects* category. The *Artifacts* category contains historic objects such as sarcophagi or ancient tools. Objects that do not fit into any other category were categorized as *Other*. Note that there is a large variance in the number of query images for each category (given on the bottom right of Fig. 2). The average scores over all query images we provide in this paper therefore have a bias towards the larger categories. This effect is desired, since we want the query distribution to be representative of a real application in a photo auto-annotation system. In addition to this, however, we will also provide a detailed anal-

ysis for all 13 categories, focusing on four categories representative of different use cases, namely *Landmark Buildings*, *Paintings*, *Building Details* and *Windows*.

### 4.3. Image Retrieval

Because it has become the de-facto standard in landmark recognition [1, 2, 13, 3], we use the vector space model for image retrieval [15, 8]. We extract SIFT [44] features from each image in the dataset and jointly cluster them using approximate k-means [8] to generate a dictionary of 1M visual words. We then vector-quantize the SIFT descriptors based on this vocabulary and represent each image as a BoVW histogram. We build an inverted file index from the BoVW histograms of the entire database [15] to enable efficient retrieval. Given a query image, we extract its BoVW as described above, query the index and rank the retrieved images by their  $tf * idf$  scores w.r.t. the query [15]. For the top-300 matches we attempt fitting a homography using SCRAMSAC [45] and consider a match verified if it has 15 or more inliers with the query. We then rank verified images above unverified images [8].

### 4.4. Scoring

We would like to evaluate the performance of a landmark recognition system in a realistic scenario. In a photo auto-annotation application, the system should assign a user’s photos reliable labels without supervision. A mobile visual search app like Google Goggles can also give the user a small selection of objects and let them pick the correct one. Therefore, we consider only the top-3 objects returned by the system.

For annotating recognition results, we showed the query and the iconic image of the recognized object to raters and asked them to rate the object’s relevance to the query as “good” if it is the exact object in the query image, “ok”, if it is somewhat relevant to the query, and “bad” if it is irrelevant. An object should be rated as “ok”, *e.g.* if the query image shows a whole building, but the match only shows a detail of that building, or vice versa. In case the query is a detail of a building and the recognized object is a different detail of the same building, the match should be rated as “bad”. If the query shows multiple landmarks, and the object is one of them, the match should still be rated as “good”.

Based on this rating, we define four scores: *good-1* is the fraction of queries with a “good” top-1 match; *ok-1* is the fraction of queries with an “ok” or “good” top-1 match; *good-3* is the fraction of queries with a “good” match in the top-3; and *ok-3* is the fraction of queries with an “ok” or “good” match in the top-3.

Category	%good-1	%ok-1	%good-3	%ok-3
Landmark Buildings	94.32	98.22	97.62	98.42
Panoramas	87.70	95.15	91.59	95.79
Sculptures	92.15	95.56	94.20	96.25
Interior Views	85.25	89.07	89.07	92.35
Building Details	87.76	91.84	89.80	91.84
Paintings	97.58	98.39	98.39	98.39
Windows	95.38	95.38	95.38	95.38
Landmark Objects	93.75	96.88	96.88	96.88
Murals	100.00	100.00	100.00	100.00
Cafes / Shops	80.56	80.56	83.33	83.33
Artifacts	91.67	91.67	94.44	94.44
Other	92.50	97.50	97.50	97.50
Multiple Objects	98.44	98.44	98.44	98.44
Total	92.74	96.42	95.58	96.92

Table 1: Performance of plain image retrieval using the full dataset.

### 4.5. Baseline Recognition Performance

For an estimate of the difficulty of the different query categories, we perform image retrieval against the full PARIS 500k dataset and manually rate the relevance of the top-3 images for each query according to the above scheme (Tab. 1). Note that these results only show the relevance of *retrieved images*, not *recognized objects*, but can serve as upper bounds for the recognition performance for each category. In total, the top-1 match was “good” for 92.74% and at least “ok” for 96.42% of the queries. Since images that did not have a match in the database are not used in the query set, the remaining 3.58% had only false-positive matches in the top-3.

## 5. Landmark Object Discovery

The first step of building a landmark recognition system is to cluster the image collection into *objects* [24, 16, 25, 4, 5, 26, 9, 27, 1, 6, 28]. A guiding question for our evaluation is: What object types can be discovered by such a clustering? As we motivated in Sec. 3.2, we choose Iconoid Shift [5] as our analysis tool to answer this question, since it produces a set of overlapping clusters, which can represent “overlapping” objects, *e.g.*, both the entire facade of Notre Dame and individual statue groups on it. In addition, it has intuitive parameters for controlling the granularity and number of discovered clusters.

### 5.1. Iconoid Shift

Iconoid Shift [5] is a mode finding algorithm based on Medoid Shift [46], designed for efficient discovery of (potentially overlapping) object clusters in large image collections. Starting from a seed image, the algorithm performs Medoid Shift in image overlap space until it converges onto a discovered object, for which it returns

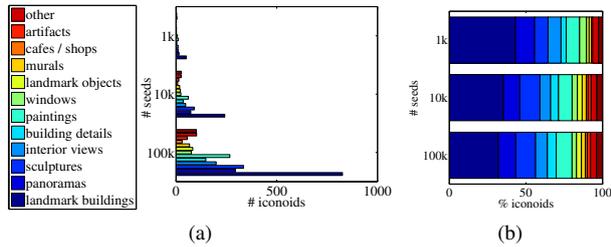


Figure 4: Distribution of categories for different numbers of seeds counting only clusters of size 5 or larger. (a) absolute (b) relative.

an iconic image (or *Iconoid*) as well as a support set containing all images having a certain minimum overlap with the Iconoid. We take this support set as the cluster. Iconoid Shift alternates between two steps: (i) Exploration of all images that overlap with the selected image by recursive image retrieval. (ii) Selection of the image that has the highest overlap with all the explored images. Empirically, the returned clusters often correspond to distinctive views of individual objects or buildings, with the Iconoid picking out the most central view. Like Mean Shift [47, 48] and the object discovery approach of Chum *et al.* [24, 49], Iconoid Shift is typically initialized with a set of *seed* images and is then run until convergence for each seed. A larger number of seeds therefore results in higher computational demands, but also causes more objects to be discovered.

### 5.2. Clusters Discovered per Category

To analyze what objects can be discovered by visual clustering, we run Iconoid Shift on the PARIS 500k dataset. Following [5], we choose a kernel bandwidth of  $\beta = 0.9$ , meaning that an image needs to have at least 10% overlap with an Iconoid to belong to its cluster. We perform several runs of the algorithm using different numbers of seed images selected randomly from PARIS 500k in order to analyze the tradeoff between runtime and the number of objects discovered.

To examine what kinds of objects the algorithm finds, we categorize all resulting Iconoid clusters of at least size 5 using the scheme from Fig. 2. Fig. 3 shows the number of discovered clusters for each category, their average size and the number of images covered by clusters. *Landmark Buildings* are the largest category with 826 clusters covering 71k images. The average cluster size of *Building Detail* is surprisingly large, because some clusters include many photos of the full facades due to our low choice of overlap threshold for Iconoid Shift. *Painting* clusters are small on average, while *Windows* have fewer but larger clusters.

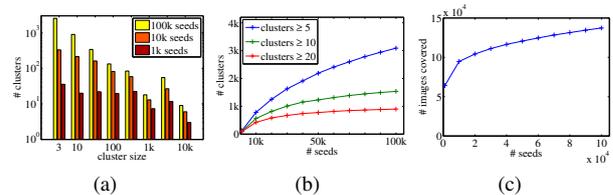


Figure 5: (a) Cluster size distribution for different numbers of seeds (Note that both axes are logarithmic.) (b) Comparison of growth rates for different cluster sizes. (c) Total number of images covered by the clustering.

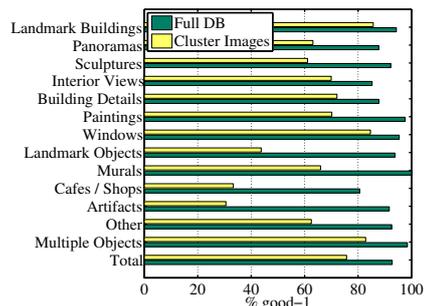


Figure 6: *good-1* retrieval performance for the full database vs. only the images discovered by Iconoid Shift using 100k seeds.

Fig. 4 shows the effect of the number of seeds on the number of clusters discovered per category. When using only 1k or 10k seeds, the category distribution remains relatively constant. The share of *Landmark Building* clusters decreases with increasing number of seeds (Fig. 4b), since more of the smaller objects such as *Paintings* or *Sculptures* are discovered.

### 5.3. Distribution of Cluster Sizes and Performance Gap

The effect of the number of seeds on the distribution of cluster sizes is shown in Fig. 5a. As reported by [13], the cluster sizes are power law distributed. We can observe that the distribution shifts towards smaller clusters when more seeds are used. Fig. 5b shows that the number of large clusters flattens out more quickly than the number of small clusters when increasing the number of seeds. This is because large landmark clusters are found first, but more seeds help find more obscure places and objects. When using 100k seeds, 12,776 Iconoids are returned, but only 3,088 of them contain 5 or more images.

Fig. 5c shows the total number of database images covered for different numbers of seeds. Using 100k seeds, a total of 137,291 images (27.4% of the dataset) are covered by the clustering. The remaining images are either irrelevant or missed by the clustering. To

estimate how many and which objects the clustering missed, we compare the retrieval performance using this reduced set of images to the full set (*cf.* Tab. 1). The result (Fig. 6) shows in which categories performance is lost and gives an upper bound on what a landmark recognition system can achieve based on this clustering. While for some categories, such as *Landmark Buildings* or *Windows*, the loss in performance is small, other categories like *Paintings*, *Landmark Objects* or *Cafes* show a strong decrease, because their clusters are small and thus more likely to be missed by the seeding process. This effect is the main cause of the total performance gap of 17.05% between the full database and the cluster images.

#### 5.4. Discussion

Since often-photographed objects are discovered first, seed-based clustering can be computationally much more efficient than computing the whole matching graph. To sufficiently cover seldom photographed objects such as museum exhibits, a larger number of seeds is necessary. The coverage of such objects could also be increased by seeding strategies that avoid the bias to large clusters. Small object discovery approaches [49, 50] or approaches that crawl tourist guide websites [6] might also help cover these objects better and close the above performance gap further. Whatever strategy is chosen, the results in Fig. 6 show that such additional steps are necessary if *Landmark Objects*, *Cafes/Shops* or *Artifacts* shall be recognized. For most of the following experiments, we choose to use 100k seeds to ensure good coverage of details and small objects.

## 6. Landmark Object Recognition

By clustering a large collection of tourist photos, we have discovered numerous interesting objects and determined a set of representative images for each of them. To now recognize a new object in a query image, a landmark recognition system performs retrieval in the set of discovered object representatives (Fig. 1). The open question here is: Given a ranking of representatives, how to rank the objects they belong to by their relevance to the query? To this end, we compare five object scoring methods and evaluate their respective tradeoffs of performance vs. database size and their suitability for different object categories.

### 6.1. Ground Truth Generation

For this evaluation, we introduce a new ground truth containing relevance ratings (Sec. 4.4) of the Iconoids

discovered with 100k seeds w.r.t. the query set. An exhaustive relevance annotation of the 12,776 Iconoids for each of the 2,987 query images would require about 883 person-days of human work, assuming 2s of annotation effort per query-Iconoid pair. Therefore, we took two measures to reduce the amount of manual labor. (1) We summarized queries showing exactly the same view into 2,042 groups, since the same Iconoids are relevant for them. For this, we computed a pairwise matching of the queries and manually inspected each pair of matching images, discarding all pairs that do not show exactly the same view. We then constructed a matching graph from the verified edges and computed its connected components. During annotation, each group was represented by one image, and annotations for it were transferred to all other members of the group. (2) We automatically rated an Iconoid as irrelevant for a group of queries if *none* of the Iconoid’s representatives were spatially verified *at least once* when querying an image retrieval system with *each* query in the group. To avoid false negatives in image retrieval, we performed an exact spatial verification by establishing correspondences using matching SIFT features that pass the SIFT ratio test [44]. Since the landmark with the largest number of images in the PARIS 500k dataset is the Eiffel Tower with about 20k images, any query can have at most 20k relevant images. To leave some room for ranking errors, we performed spatial verification for the top-30k retrieved images ranked w.r.t. their *tf\*idf* scores. The 26.8k remaining pairs of query groups and Iconoids were manually annotated according to the rating scheme introduced in Sec. 4.4. Annotations were performed by 28 people over a period of 8 weeks. Each pair was shown to 3 people who were asked to rate it as *good*, *ok* or *bad* according to the scoring scheme we introduced in Sec. 4.4, and the final annotation was decided by majority voting. In the 1.8k cases where all three annotations were inconsistent, the image pair was passed to a fourth annotation for a definite annotation.

This ground truth, including the query images, the Iconoid clusters, the query-Iconoid relevance annotations, and the query category annotations is publicly available under <sup>1</sup>. We believe that this ground truth will be useful to the landmark recognition community, since (i) so far, there is no landmark recognition dataset at this scale (3k queries, 13k clusters and 137k representative images) (ii) this is the first dataset where the clusters were produced by an actual landmark clustering algorithm instead of keyword searches on Internet

<sup>1</sup><http://www.vision.rwth-aachen.de/data/paris500k/paris-dataset>

photo collections (iii) the category annotations allow a detailed performance analysis for different use cases, *e.g.*, paintings or statues, while existing datasets only focus on buildings.

## 6.2. Methods

We evaluate five object scoring methods:

**Center** represents objects only by the cluster center (*i.e.*, Iconoid) of each object and discards all other representatives, as done in Yang *et al.* [41]. The object ranking is then simply the same as the representative ranking.

**Size** returns all objects with at least one matching representative, and scores them by their cluster size, *i.e.*, by the number of times they were photographed.

**Voting** lets each matching representative cast a vote for each object it belongs to (note that clusters can overlap). Objects are then ranked by their number of votes.

**Best Match** returns the object with the highest scoring representative, as done in, *e.g.*, [4, 6, 13, 2]. A difference in our case is that we are using a soft clustering, so a representative can belong to multiple objects. In this case, we return the object with the largest cluster size. This method can therefore also be viewed as a variant of the *Size* method that only uses the best matching representative.

**Overlap** uses *Homography Overlap Propagation* [5] to compute the overlap of the query with each Iconoid. The method first computes the overlap region of the query with the matching representative and then propagating this region into the Iconoid via the shortest path in the Iconoid cluster’s matching graph. This is done using the homography overlap propagation (HoP) algorithm from [5]. If a query matches multiple representatives of the same Iconoid, the overlap is computed for each of them and the largest overlap is used. Objects are then ranked in decreasing order of their Iconoids’ overlaps with the query.

Note that, while *Center* and *Best Match* evaluate strategies used in the literature, *Overlap* is a novel strategy. We now first compare the above methods using the objects discovered with 100k Iconoid Shift seeds and then show the effect of the number of seeds on their performance.

## 6.3. Results

We evaluated the five approaches based on the ground truth introduced in Sec. 6.1. To estimate the error introduced by the SIFT matching pre-filter we used to reduce annotation effort (Sec. 6.1), we performed a small control experiment. We manually rated the relevance

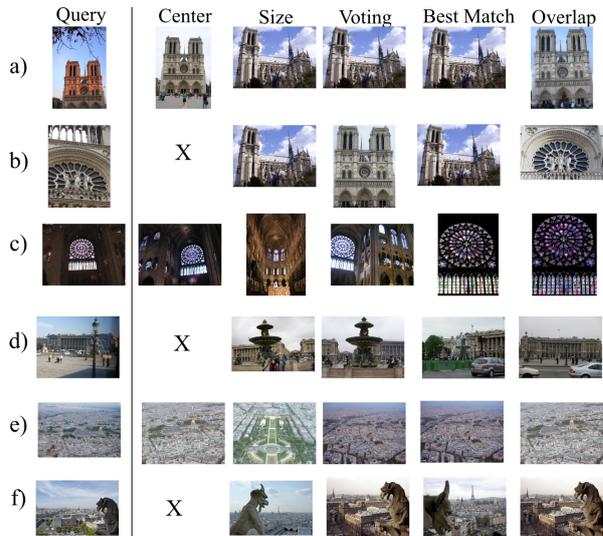


Figure 7: Top-scoring objects for different object scoring methods. “X” means that no objects were recognized.

of the top-3 Iconoids retrieved for each query using the *Voting* method (Sec. 6.2) and found that 0.7% of the Iconoids rated “good” or “ok” were filtered out. We believe this small false-negative rate is still acceptable, since the simplified annotation procedure significantly reduced the amount of manual labor.

The recognition performance of the different methods is compared in Tab. 2. Fig. 7 shows the top scoring objects for typical queries. *Center* finds images closely resembling the query, but often fails to find *any* matching objects, because the cluster centers are not sufficient to recognize all objects under different viewing conditions due to the limited invariance of the matching process. However, since it only requires one image per object, it is by far the fastest and most memory efficient method. Because *Size* chooses the largest cluster, it often finds a viewpoint more popular than the query (popular landmarks are often represented by multiple clusters from different viewpoints). Sometimes this effect is desired since the largest cluster usually depicts the object best, but it can also cause drift (Fig. 7b-f), *i.e.* instead of the query object, a nearby object is recognized. *Voting* finds the clusters with the most matching representatives. This makes it less prone to drift (Fig. 7b,e,f) and causes it to achieve higher performance than *Size*. Despite being simpler than *Size*, *Best Match* outperforms it. The reason is that *Best Match* considers only the closest matching representative to the query, making it less prone to drift than *Size* that also looks at farther away matches (Fig. 7c,e). *Overlap* has the best *good-1* performance, because it computes the actual overlap of the

	% good-1	% ok-1	% good-3	% ok-3
Centers	39.60	45.56	42.99	46.03
Size	57.11	73.32	66.42	76.26
Voting	59.42	<b>76.00</b>	69.80	<b>77.64</b>
Best Match	60.40	75.39	67.26	77.10
Overlap	<b>63.71</b>	75.93	<b>71.78</b>	77.13

Table 2: Performance of different object scoring methods.

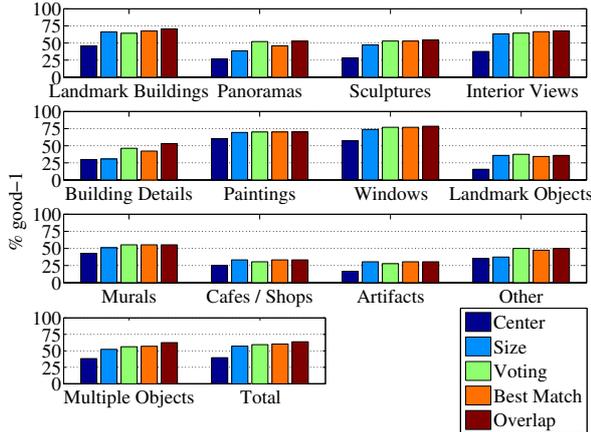


Figure 8: Performance of object scoring methods by query type.

query image with each cluster’s iconic image and selects the iconic whose view is closest to the query (Fig. 7a-f).

The *good-1* performances by query type are shown in Fig. 8. *Size* compares well to the other methods on *Landmark Buildings* and *Cafes / Shops*, where the largest cluster is often the correct one (e.g., the full view of a facade). On *Paintings*, all methods including *Center* have similar performance, because *Paintings* are flat objects and are usually photographed under the same viewing conditions. Since this makes painting retrieval very easy, multiple representatives do not bring an advantage, explaining the relatively good performance of *Center*. *Size* performs worse than *Voting* and *Best Match* on *Panoramas* (Fig. 7e), because it tends to drift to more popular nearby views, moving the query object out of the field of view. The same effect occurs on *Building Details* (Fig. 7b), where *Size* tends to return views of the whole building instead. *Overlap* has a particular advantage on classes where other methods tend to drift (*Panoramas*, *Building Details*, *Multiple Objects*), since it usually finds the iconic image that best matches the photographed part (Fig. 7a-f). *Center* works relatively well for *Windows* and *Murals*, because, like *Paintings*, they are only photographed from a limited range of viewing angles, making them easy to match.

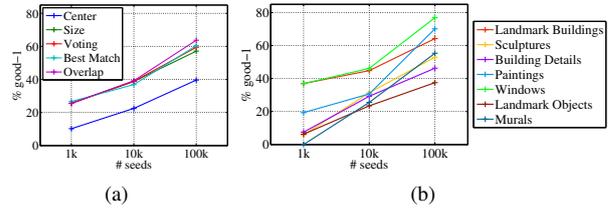


Figure 9: (a) *good-1* performance of object scoring methods for different numbers of Iconoid Shift seeds. (b) *good-1* performance by query type using the *Voting* method.

#### 6.4. Effect of the Number of Seeds

We now analyze the effect of the number of object clusters (which depends on the number of seeds) on the performance of the five methods and compare object retrieval performance by category. As Fig. 9a shows, the performance of the methods does not differ much for 1k and 10k seeds, except that *Center* consistently performs worst for the reasons explained above. The differences become slightly more pronounced at 100k seeds, where the density of objects is very high and methods that are less prone to drift to nearby objects gain an advantage. Conversely, this shows that simpler methods are sufficient if the object density is low. Fig. 9b shows the performance for different query types when using the *Voting* method. *Sculptures*, *Paintings*, *Windows* and *Murals* show the steepest improvement since they require more seeds to be sufficiently covered by the clustering, while *Landmark Buildings* can already be recognized when using a smaller number of seeds since they form large clusters that are discovered early. Surprisingly, *Windows* have the highest recognition rate overall. The reasons for this are that (i) they are easy to recognize since they are flat, highly textured objects, and (ii) they get discovered already with few seeds, since *Window* clusters are almost three times the size of *Painting* clusters (Fig. 3b).

#### 6.5. Discussion

The ideal choice of method varies by application: *Center* provides high performance for flat objects at low computational and memory cost. *Voting* has high accuracy across all categories and its speed makes it applicable, e.g., for mobile visual search. The popular *Best Match* method has similar performance and efficiency, but is outperformed by *Overlap*, which has the highest performance overall but also the highest computational cost. *Overlap* is therefore better suited for offline applications requiring high accuracy, e.g., photo auto-annotation.

Even when looking at the best performing method, *Overlap*, there is still a difference of 11.98 percent points between the *good-1* performance of object retrieval (63.71%) and the *good-1* performance of plain image retrieval (75.69% of plain image retrieval (Fig. 6). The cause for this *clustering gap* could be either a too coarse clustering or imprecise object ranking. If we consider that the difference between *good-1* and *good-3* for *Overlap* (8.07%) is larger than the gap of 3.91 percent points between the *good-3* performance of *Overlap* and the *good-1* performance of plain image retrieval, it becomes apparent that the main part of the clustering gap is due to the object ranking. Hence, there is still room for improved object ranking methods to close this gap.

Fig. 9a shows diminishing returns in performance when the number of clusters increases (note the logarithmic x-axis), since the most popular queries are covered first and the long tail of queries requires exponentially more effort.

The performance gap between *Center* and other methods shows that the representatives are necessary for ensuring invariance to different viewing conditions. However, it is desirable to reduce the set of representatives, since it determines the memory use and speed of the retrieval index. This is examined more closely in the following section.

## 7. Efficient Representations for Retrieval

Since the discovered landmark representatives are highly redundant, subsampling them can save memory and computation time. The goal here is to reduce the set of representatives in a way that still preserves as much visual variability as possible in order to ensure good retrieval performance. The methods we present in the following work by summarizing groups of similar images, as in *e.g.* [1, 13], which can be done efficiently by exploiting the similarity information that is already available from the *matching graph* constructed during clustering [1, 13, 4, 5, 6]. In this graph, every matching pair of images is linked by an edge whose weight is their matching score. We evaluate four approaches and compare them against a random baseline. Following [1, 13, 4] we use the number of homography inliers as edge weights.

### 7.1. Methods

We compare the following five methods for reducing redundancy in the sets of representatives:

**Complete-Link** (Gammeter *et al.* [13]) performs hierarchical agglomerative clustering and replaces each

complete link component containing at least 3 images by its image with the most neighbors.

**Kernel Vector Quantization (KVQ)** [51] is a clustering method that selects a minimum number of points such that each point in the dataset is within radius  $r$  of at least one selected point. It is used by Avrithis *et al.* [1] to reduce the set of features in a cluster after projecting features from all images into a single iconic image, called *Scene Map*. Applying the same idea on the image level, we use it to find a minimal subset of representative images such that each image in a cluster has a given minimum matching score with at least one image in the subset.

**Dominating Set** chooses a subsample such that each representative is adjacent to at least one image in the original cluster. This subset is found by solving the corresponding set cover problem using the greedy set cover algorithm [52].

**Fine Iconoids** performs a second, finer Iconoid Shift [5] clustering at bandwidth  $\beta = 0.7$  that covers the image collection at a very fine granularity. The representatives for each (coarse) cluster are then chosen to be all the fine Iconoids in it. This is similar to the approach of Raguram *et al.* [53] that represents objects by a set of iconic images found by clustering Gist descriptors.

**Random** is the baseline method that simply draws a random subsample of the representative images.

### 7.2. Results

We now compare the tradeoffs between the number of representatives and recognition performance of these methods using the *Voting* method (Sec. 6.2) for object scoring (Fig. 10). The number of representatives that the *Complete-Link*, *KVQ* and *Dominating Set* methods return can be controlled by first deleting edges below a certain edge weight threshold to make the matching graph sparser and then running the algorithm. We generated 3 sets of representatives with each method by applying thresholds of 15, 30 and 50 inliers.

As reported in [13], *Complete Link* only slightly reduces the representative set while maintaining high recognition performance. *Dominating Set* and *KVQ* yield comparable results since they optimize similar criteria. They represent a good tradeoff, allowing for a reduction to about 40% while still achieving a 60.9% *good-1* performance (at threshold 50). An interesting result is that *Fine Iconoids* performs well on *Paintings* and *Building Details*, but lower than *Random* on *Landmark Buildings*. Visual inspection showed that Iconoids do not cover the more obscure viewing conditions necessary for robust recognition, since the algorithm is de-

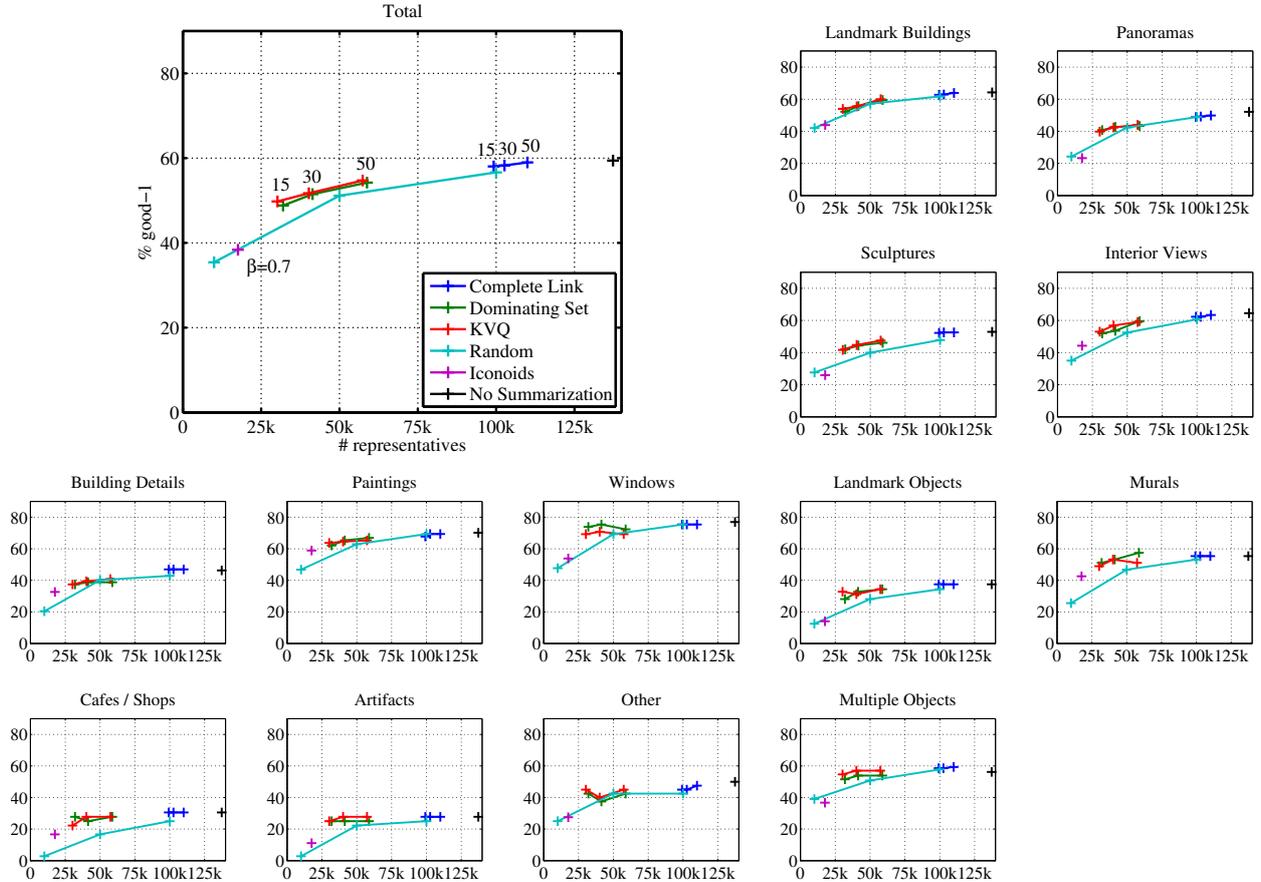


Figure 10: Performance / index size tradeoff of cluster summarization methods using edge weight thresholds of 15, 30 and 50.

signed to converge to popular views. It therefore performs higher on categories with a limited range of possible viewing conditions.

### 7.3. Discussion

In this evaluation, we focused on reduction methods that work on the *image level*. We have shown that in particular *KVQ* and *Dominating Set* methods can achieve high compression at only a small loss in precision. Some recent approaches also perform this reduction on the feature level, usually combined with *offline query expansion*, *i.e.*, projecting features into matching images, which is reported to even improve precision over baseline retrieval [1, 42]. An evaluation of these methods would be an interesting task for future work.

## 8. Interfacing Images with Semantics

To find suitable descriptions for the discovered objects, and thus to enable linking them with informa-

tion on the web, the usual method is to perform statistical analysis of the tags and titles that users provided for the photos in a cluster [9, 4, 10, 6]. Quack *et al.* [4] mine *frequent itemsets* in all tags of a cluster to generate candidate names. The top-15 candidates are then used to query Wikipedia, and the retrieved articles are verified by matching the images occurring in them against the images in the cluster. In a small informal experiment we found that frequent itemsets returns many noisy and non-descriptive names like “vacation”, “photo”, “canon”, or “europe”, which need to be filtered by a comprehensive stoplist. Furthermore, tags that are frequently used by the same user like “summer vacation 2008” can be ranked higher than correct but less frequent terms.

The method of Simon *et al.* [10] is specifically designed to handle both of these problems. It probabilistically computes a score  $score(c, t)$  for each pair of tag  $t$  and cluster  $c$ . This score is based on the conditional

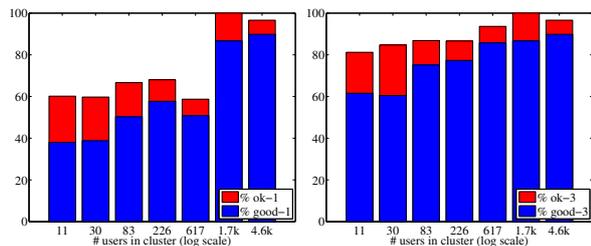


Figure 11: Tag quality as a function of the number of individual users in the cluster. Left: top-1 tag, right: top-3 tags

probability of cluster  $c$  given tag  $t$ , resulting in tags that mainly occur in cluster  $c$ . By marginalizing over the users, tags that are frequent in the cluster, but used by only few users are ranked low. We use the method of Simon *et al.* [10] in our evaluation since we found that it yields more reliable tags than the method of Quack *et al.* [4]. We analyze for which objects we can reliably find semantics and examine the performance gap between object retrieval and semantic annotation.

### 8.1. Data Preparation

We first need to define a set of tags for each image based on the metadata provided by the respective photo sharing website. The PARIS 500k dataset consists of images from Flickr and Panoramio. Since images on Flickr have both tags and a title, we treated the title as an additional tag. Since images on Panoramio have no tags, we used the titles as their sole tag. We preprocessed image tags by applying a very small stoplist containing terms such as “Paris” and “France” and removing filenames like “DSC002342.JPG”.

### 8.2. Tag Quality Annotation

In order to analyze tag quality, we manually rated the quality of the top 3 tags for each object cluster containing 6 or more images. On average, annotation of a single object-tag pair took about 30-60 seconds, since often a web search was necessary to verify the correctness of a tag. This annotation was performed by six people who annotated a total of 2,536 objects. The annotators were asked to rate each image-tag pair as “good” if the tag accurately describes what is visible in the iconic image (*e.g.*, the full name of a building, the title and painter of a painting) and as “ok” if it provides at least some helpful information, such as the creator of a sculpture, but not its name, or if the tag is accurate, but contains noise terms like “me in front of Notre Dame”. This annotation allows us to re-use the evaluation measures we used for object retrieval in Sec. 6.

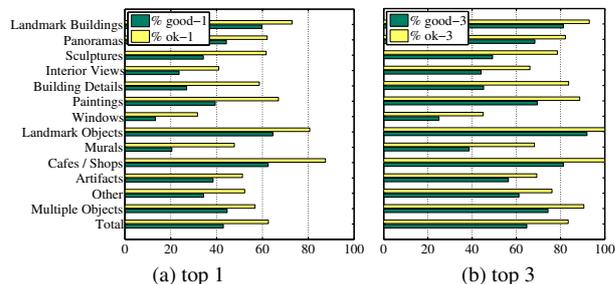


Figure 12: Tag quality of different object categories for Iconoid clusters of size 6 and higher.

museedorস্য musee d'orsay mus e-d'orsay mus e d'orsay museo d'orsay  
muse'e d'orsay mus ed'orsay mus ed'orsay mus e d'orsay museo orsay  
musee d'orsay mus edorস্য mus e d'orsay mus e d'orsay museu de orsay  
musee d'orsay mus e d'orsay mus e d'orsay mus ed'orsay museu d'orsay  
musee d'orsay mus e d'orsay mus e d'orsay mus e orsay museu d'orsay  
musee d'orsay mus e d'orsay mus e d'orsay mus e orsay museu d'orsay  
musee d'orsay mus e d'orsay mus ed'orsay museo orsay museumdorsay  
musee d'orsay mus e d'orsay mus ed'orsay museo d'orsay museum d'orsay

Figure 13: Different spellings of “Mus e d’Orsay” encountered in the dataset.

### 8.3. Tag Mining

We first analyze the influence of the number of users contributing photos to a cluster on the reliability of automatic semantic annotation. Simon *et al.*’s method [10] outputs a ranking of potential names for each cluster, allowing us to examine the accuracy of the top-1 and the top-3 tags. Fig. 11 shows that tag reliability clearly increases with more users. In particular, over 80% of clusters with over 1k users have a *good* top-1 tag. With increasing user count, descriptions also become more precise, since the fraction of *ok* tags decreases. Fig. 12 shows the tag quality for different categories. The tags determined for *Landmark Buildings*, *Landmark Objects* and *Cafes / Shops* are most reliable, since their names are typically well-known. *Cafes / Shops* are particularly easy to tag since their name is usually directly visible. *Murals*, *Windows*, *Sculptures* and *Building Details* are usually lacking proper annotations since photographers often do not know their names and only label them with generic tags. For example, *Building Details* are often tagged with the name of the entire building. This causes the large difference in the *good-1* and *ok-1* scores of these categories.

### 8.4. Discussion

While for most large clusters suitable semantics can be found, for small and medium sized clusters the difference in *ok-1* and *good-1* scores (Fig. 11) suggests that better tag ranking could greatly help the recognition of less popular objects. Significant improvements

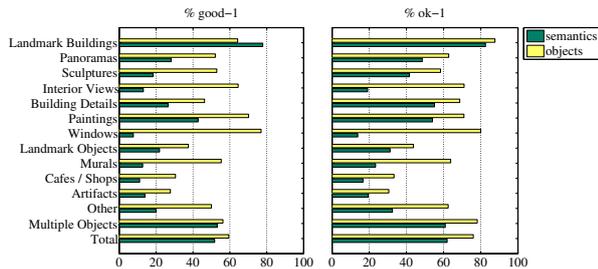


Figure 14: Quality of good (left) and ok (right) semantics assigned to queries, compared to the quality of the objects retrieved.

can likely be made by increasing robustness w.r.t. different languages, spelling errors and tag noise. For example, Fig. 13 shows a selection of different spellings of “Musée d’Orsay” from our dataset that would all be treated as separate tags by current methods. Mining Wikipedia [4] or tourist guide websites [6], or performing specialized per-category metadata mining as done by Arandjelović *et al.* [54] for sculptures might also help in naming the less popular objects.

## 9. End-to-End Analysis

In Sec. 6, we analyzed different methods for assigning *objects* to queries and measured accuracy on a visual level. In this section, we perform this analysis on a *semantic* level, based on the labels assigned in the previous section.

### 9.1. Setup

To evaluate the system from end to end, we cluster the PARIS500k dataset with *Iconoid Shift* using 100k seeds and mine semantics for the clusters using the method of Simon *et al.* [10]. We use the *Voting* object scoring method (Sec. 6) to rank objects w.r.t. a query. Subsampling (Sec. 7) is not used. We then rate the relevance of the top-scoring tag of the top-scoring object for each query as either *good*, *ok* or *bad* (Sec. 8.2).

### 9.2. Results

Fig. 14 shows the results of semantic annotation and compares it to the plain object recognition performance (yellow bars in Fig. 8). *Landmark Buildings* have the highest performance, because their large cluster size enables robust recognition and semantic assignment. The reason semantic assignment has even higher performance than plain object recognition is that semantic assignment sometimes corrects errors of object recognition: It often happens that the query image shows the whole building, but the matching object is a detail

of the building, *e.g.* a door of Notre Dame or a leg of the Eiffel Tower. This occurs, *e.g.*, because the detail is most prominent in the query due to perspective, or because large parts of the building are occluded, but the detail is still visible. However, photos of building details are often labeled with the name of the whole building, since photographers do not know, *e.g.*, the name of a particular door of Notre Dame. Therefore, the query is correctly assigned the name of the whole building even though the recognized object was a detail.

For categories with smaller clusters, there are different bottlenecks: *Artifacts* and *Landmark Objects* are hard to recognize, because they are compact objects and form only small clusters. However, they typically have high quality tags (Fig. 12). In other cases, a relevant object can reliably be retrieved, but low tag quality prevents successful semantic annotation. This problem is most prominent for *Windows* and *Murals* (Fig. 14). They are flat and photographed under a limited range of conditions, making them easy to recognize, but they suffer from low quality semantics, since information on them is not easily available. In total, 51.8% of the queries could be assigned good semantics, while for 59.4% a good object was retrieved; 61.9% of queries were assigned ok tags, while for 76.0% an ok object was retrieved. In the following, we summarize the causes of this gap and discuss possible solutions.

## 10. Discussion and Conclusion

We now sum up the findings of our evaluation by answering the questions posed in the introduction and discuss the areas where improvements can still be made.

*How many and what kinds of objects are present in Internet photo collections and what is the difficulty of discovering objects of different object categories?*

The question how many objects there are cannot be answered based directly on the number of clusters, because there can be multiple clusters of the same object showing different views, and clusters of non-objects, *e.g.* party photos, pictures of animals and food, or photo bursts. We performed an annotation experiment (Sec. 5.2) and labeled the 3,088 clusters containing five or more images discovered in the PARIS500k dataset. 2,585 (83.7%) clusters were labeled as objects and 503 (16.3%) were labeled as non-objects (Note that this ratio will shift more strongly towards non-objects when also considering clusters containing less than 5 images.) Approaches for detecting and removing such non-object clusters could be a direction for future work, since they

unnecessarily increase the size of the retrieval index and increase the chance of recognition errors. The distribution of object categories (Fig. 3a) shows that, not surprisingly, *Landmark Buildings* from the largest category, followed by *Sculptures*, *Panoramas* and *Paintings*. Seed-based object discovery algorithms [5, 24, 49] find the most photographed objects first, because when drawing a random image from an Internet photo collection, the likelihood of drawing an often photographed object is higher. Therefore, much more effort is required to also discover objects in the long tail of the size distribution. This could be addressed in future work, e.g. by seeding methods that avoid the bias to large clusters, or methods that explicitly mine for small objects [49, 50].

*How to decide which landmark was recognized given a list of retrieved images?*

We analyzed five methods for this task (Sec. 6). Our experiments clearly showed that having a set of representative images for each cluster is necessary to recognize it under difficult conditions such as extreme view differences, occlusion, lighting changes, blur, etc. (first two rows of Fig. 15). This can be viewed as a form of offline query expansion. However, like query expansion, this method is also prone to *drift* (Fig. 16, top row), which can cause confusion between nearby objects. While the often used *Best Match* method [4, 6, 13, 2] avoids drift better than some other methods, we found that by using a method that explicitly maximizes the *Overlap* between the query and the object’s iconic image, even higher precision can be achieved. However, the ranking gap (*good-1* vs. *good-3* performance) remains relatively large, suggesting that there is potential for more accurate object ranking methods.

*How to efficiently represent the discovered objects in memory for recognition?*

We analyzed four *image level* techniques for eliminating redundancy in the database (Sec. 7). Our analysis revealed that it is important to keep representatives showing obscure views and extreme lighting conditions. Therefore, representing the object by a set of popular views (as done in our *Fine Iconoids* method or in [53]) did not perform better than random subsampling. We proposed two methods that achieve an acceptable trade-off between database size and recognition performance, but observed that there is still potential for better methods. A comparative analysis of methods that eliminate redundancy at the *feature level* [1, 2, 42] would also be an interesting future direction.

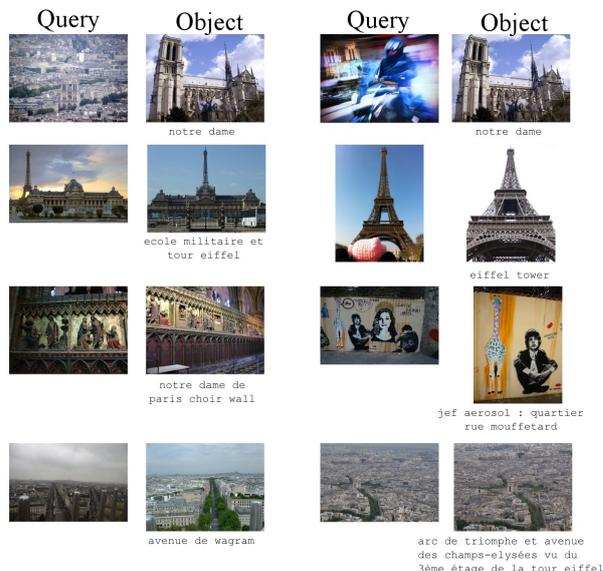


Figure 15: Examples of successful recognition and annotation.

*Are the user-provided tags reliable enough for determining accurate object names?*

The largest room for improvement of the overall performance of landmark recognition is in semantic annotation. We determined two reasons for the performance loss at this step (Sec. 8): (1) User-provided tags come in different languages, have spelling errors and contain noise terms (Fig. 13). Methods robust to these factors could provide more reliable semantics even for small clusters. (2) Often, insufficient information is available to photographers, causing non-descriptive tags (Fig. 16, bottom right). This might be addressed by crawling relevant encyclopedia [4] or tourist guide articles [6] or using image search engines [54]. However, sometimes, the presence of accurate tags in small clusters can also lead to surprisingly accurate results (Fig. 15, rows 3 and 4).

*What are the factors effectively limiting the recognition of different landmark types?*

Our end-to-end analysis (Sec. 9) showed that the factors that limit recognition performance are quite varied and strongly depend on the object category. Some objects like *Windows* or *Murals* are easy to recognize *visually*, but often lack accurate tags, which prevents semantic assignment. This could be addressed using the methods mentioned above. Other objects such as *Artifacts* or *Landmark Objects* do have accurate tags, but are harder to recognize due to their spatial structure and small cluster size. Their recognition could be improved

by mining more photos of them from the web [13]. Finally, improvements to image retrieval and matching, *e.g.* improved feature representation [55], improved feature quantization [18, 56], or ranking methods robust to problems like repeated patterns [57] (Fig. 16, bottom left), will directly benefit both landmark clustering and recognition.

### 10.1. Limitations

While our evaluation has brought to light several opportunities for progress, its scope could still be broadened in future work. Due to our choice of dataset our taxonomy of queries is certainly biased towards the landmarks of Paris. A larger dataset from several cities would increase the generality of the evaluation. Our query set was collected from Internet photo collections and is therefore representative for the task of photo auto-annotation. While this bias only affects the score average and not the per-category scores, a second query set for the task of mobile visual search would make it possible to identify problems specific for that task. The set of methods we analyzed was carefully chosen, but an analysis of other approaches (*e.g.*, for clustering or semantic annotation) may bring further insights into how the component choices affect overall performance.

### 10.2. Conclusion

In this work, we have evaluated the automatic construction of visual landmark recognition engines from Internet image collections. We used a large-scale dataset of 500k photos from Paris, collected a set of 3k typical query images, and created a ground truth for evaluating large-scale photo auto-annotation that we made publicly available. For each component of the pipeline, we evaluated how different methods and parameters affect overall performance as well as the performance for individual query categories. We proposed several novel methods for various sub-tasks, some of which outperform literature approaches. In our analysis, we have identified areas where such a system performs well, as well as areas where improvement is still possible.

### Acknowledgments

This project has been funded, in parts, by a Google Faculty Research Award and by the cluster of excellence UMIC (DFG EXC 89).



Figure 16: Examples of failure modes. Top row: Drift due to dominance of nearby objects. Bottom left: Failed retrieval caused by repeating patterns. Bottom right: Too generic description.

### References

- [1] Y. Avrithis, Y. Kalantidis, G. Toliás, E. Spyrou, Retrieving Landmark and Non-Landmark Images from Community Photo Collections, in: ACM Multimedia, 2010, pp. 153–162.
- [2] S. Gammeter, T. Quack, L. Van Gool, I Know What You Did Last Summer: Object-Level Auto-Annotation of Holiday Snaps, in: International Conference on Computer Vision, 2009, pp. 614–621.
- [3] Y. Kalantidis, G. Toliás, Y. Avrithis, M. Phinikettos, E. Spyrou, P. Mylonas, S. Kollias, Viral: Visual image retrieval and localization, Multimedia Tools and Applications 51 (2011) 555–592.
- [4] T. Quack, B. Leibe, L. Van Gool, World-Scale Mining of Objects and Events from Community Photo Collections, in: Conference on Image and Video Retrieval, 2008, pp. 47–56.
- [5] T. Weyand, B. Leibe, Discovering Favorite Views of Popular Places with Iconoid Shift, in: International Conference on Computer Vision, 2011, pp. 1132–1139.
- [6] Y.-T. Zheng, M. Zhao, Y. Song, H. Adam, U. Buddemeier, A. Bissacco, F. Brucher, T.-S. Chua, H. Neven, Tour the world: Building a Web-Scale Landmark Recognition Engine, in: IEEE Conference on Computer Vision and Pattern Recognition, 2009, pp. 961–962.
- [7] O. Chum, J. Philbin, J. Sivic, A. Zisserman, Total Recall: Automatic Query Expansion with a Generative Feature Model for Object Retrieval, in: International Conference on Computer Vision, 2007.
- [8] J. Philbin, O. Chum, M. Isard, J. Sivic, A. Zisserman, Object Retrieval with Large Vocabularies and Fast Spatial Matching, in: IEEE Conference on Computer Vision and Pattern Recognition, 2007.
- [9] D. Crandall, L. Backstrom, D. Huttenlocher, J. Kleinberg, Mapping the World’s Photos, in: World Wide Web Conference, 2009, pp. 761–770.
- [10] I. Simon, N. Snavely, S. Seitz, Scene Summarization for Online Image Collections, in: International Conference on Computer Vision, 2007.
- [11] J.-M. Frahm, P. Georgel, D. Gallup, T. Johnson, R. Raguram, C. Wu, Y.-H. Jen, E. Dunn, B. Clipp, S. Lazebnik, M. Pollefeys, Building Rome on a Cloudless Day, in: European Conference on Computer Vision, 2010, pp. 368–381.
- [12] T. Weyand, J. Hosang, B. Leibe, An Evaluation of Two Landmark Building Discovery Algorithms, in: ECCV’10 Workshop on Reconstruction and Modeling of Large-Scale Virtual Environments, 2010, pp. 310–323.
- [13] S. Gammeter, T. Quack, D. Tingdahl, L. Van Gool, Size Does Matter: Improving Object Recognition and 3D Reconstruction

- with Cross-Media Analysis of Image Clusters, in: European Conference on Computer Vision, 2010, pp. 734–747.
- [14] D. Nistér, H. Stewénius, Scalable recognition with a vocabulary tree, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2006, pp. 2161–2168.
- [15] J. Sivic, A. Zisserman, Video Google: A Text Retrieval Approach to Object Matching in Videos, in: International Conference on Computer Vision, volume 2, 2003, pp. 1470–1477.
- [16] X. Li, C. Wu, C. Zach, S. Lazebnik, J.-M. Frahm, Modeling and Recognition of Landmark Image Collections Using Iconic Scene Graphs, in: European Conference on Computer Vision, 2008, pp. 427–440.
- [17] J. Philbin, O. Chum, M. Isard, J. Sivic, A. Zisserman, Lost in Quantization: Improving Particular Object Retrieval in Large Scale Image Databases, in: IEEE Conference on Computer Vision and Pattern Recognition, 2008.
- [18] H. Jégou, M. Douze, C. Schmid, Hamming embedding and weak geometric consistency for large scale image search, in: European Conference on Computer Vision, 2008.
- [19] T. Sattler, T. Weyand, B. Leibe, L. Kobbelt, Image Retrieval for Image-Based Localization Revisited, in: British Machine Vision Conference, 2012, pp. 76.1–76.12.
- [20] A. Irschara, C. Zach, J.-M. Frahm, H. Bischof, From structure-from-motion point clouds to fast location recognition, in: IEEE Conference on Computer Vision and Pattern Recognition, 2009.
- [21] Y. Li, N. Snavely, D. Huttenlocher, P. Fua, Location Recognition using Prioritized Feature Matching, in: European Conference on Computer Vision, 2010.
- [22] G. Baatz, K. Koester, D. Chen, R. Grzeszczuk, M. Pollefeys, Handling urban location recognition as a 2D homothetic problem, in: European Conference on Computer Vision, 2010.
- [23] Y. Li, N. Snavely, D. Huttenlocher, P. Fua, Worldwide Pose Estimation Using 3D Point Clouds, in: European Conference on Computer Vision, 2012.
- [24] O. Chum, J. Matas, Large-scale discovery of spatially related images, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32 (2010) 371–377.
- [25] J. Philbin, A. Zisserman, Object Mining using a Matching Graph on Very Large Image Collections, in: Indian Conference on Computer Vision, Graphics and Image Processing, 2008, pp. 738–745.
- [26] E. Johns, G.-Z. Yang, From images to scenes: Compressing an image cluster into a single scene model for place recognition, in: International Conference on Computer Vision, 2011, pp. 874–881.
- [27] Y. Li, D. J. Crandall, D. P. Huttenlocher, Landmark classification in large-scale image collections, in: International Conference on Computer Vision, 2009, pp. 1957–1964.
- [28] R. Ji, L.-Y. Duan, J. Chen, H. Yao, J. Yuan, Y. Rui, W. Gao, Location Discriminative Vocabulary Coding for Mobile Landmark Search, *International Journal of Computer Vision* 96 (2012) 290–314.
- [29] P. Gronat, G. Obozinski, J. Sivic, T. Pajdla, Learning per-location classifiers for visual place recognition, in: IEEE Conference on Computer Vision and Pattern Recognition, 2013.
- [30] A. Bergamo, S. N. Sinha, L. Torresani, Leveraging Structure from Motion to Learn Discriminative Codebooks for Scalable Landmark Classification, in: IEEE Conference on Computer Vision and Pattern Recognition, 2013, pp. 763–770.
- [31] D. Chen, G. Baatz, K. Köser, S. Tsai, R. Vedantham, T. Pylvänäinen, K. Roimela, X. Chen, J. Bach, M. Pollefeys, B. Girod, R. Grzeszczuk, City-scale landmark identification on mobile devices, in: IEEE Conference on Computer Vision and Pattern Recognition, 2011, pp. 737–744.
- [32] J. Knopp, J. Sivic, T. Pajdla, Avoiding confusing features in place recognition, in: European Conference on Computer Vision, 2010, pp. 748–761.
- [33] A. Torii, J. Sivic, T. Pajdla, Visual localization by linear combination of image descriptors, in: ICCV Mobile Vision Workshop, 2011, pp. 102–109.
- [34] G. Schindler, M. Brown, R. Szeliski, City-Scale Location Recognition, in: IEEE Conference on Computer Vision and Pattern Recognition, 2007.
- [35] E. Johns, G.-Z. Yang, Generative Methods for Long-Term Place Recognition in Dynamic Scenes, *International Journal of Computer Vision* 106 (2014) 297–314.
- [36] T. Sattler, B. Leibe, L. Kobbelt, Fast image-based localization using direct 2d-to-3d matching, in: International Conference on Computer Vision, 2011, pp. 667–674.
- [37] T. Sattler, B. Leibe, L. Kobbelt, Improving Image-Based Localization by Active Correspondence Search, in: European Conference on Computer Vision, 2012, pp. 752–765.
- [38] R. Hartley, A. Zisserman, *Multiple View Geometry in Computer Vision*, Cambridge University Press, ISBN: 0521540518, 2004.
- [39] S. Cao, N. Snavely, Graph-Based Discriminative Learning for Location Recognition, in: IEEE Conference on Computer Vision and Pattern Recognition, 2013, pp. 700–707.
- [40] H. Jégou, M. Douze, C. Schmid, Packing Bag-of-features, in: International Conference on Computer Vision, 2009, pp. 2357–2364.
- [41] J. J. Lin Yang, C. Zhang, Efficient place recognition with canonical views, in: International Conference on Multimedia & Expo, 2011.
- [42] P. Turcot, D. Lowe, Better matching with fewer features: The selection of useful features in large database recognition problems, in: ICCV Workshop on Emergent Issues in Large Amounts of Visual Data, 2009, pp. 2109–2116.
- [43] S. Cao, N. Snavely, Minimal Scene Descriptions from Structure from Motion Models, in: IEEE Conference on Computer Vision and Pattern Recognition, 2014.
- [44] D. Lowe, Distinctive Image Features from Scale-Invariant Keypoints, *International Journal of Computer Vision* 60 (2004).
- [45] T. Sattler, B. Leibe, L. Kobbelt, SCRAMSAC: Improving RANSAC’s Efficiency with a Spatial Consistency Filter, in: International Conference on Computer Vision, 2009, pp. 2090–2097.
- [46] Y. Sheikh, E. Khan, T. Kanade, Mode-Seeking by Medoidshifts, in: IEEE Conference on Computer Vision and Pattern Recognition, 2007.
- [47] Y. Cheng, Mean shift mode seeking and clustering, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 17 (1995) 790–799.
- [48] D. Comaniciu, P. Meer, Mean Shift: A Robust Approach Toward Feature Space Analysis, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 24 (2002) 603–619.
- [49] O. Chum, M. Perdoch, J. Matas, Geometric min-Hashing: Finding a (Thick) Needle in a Haystack, in: International Conference on Computer Vision, 2007.
- [50] P. Letessier, O. Buisson, A. Joly, Scalable mining of small visual objects, in: ACM Multimedia, 2012, pp. 599–608.
- [51] M. Tipping, B. Schölkopf, A kernel approach for vector quantization with guaranteed distortion bounds, in: *Artificial Intelligence and Statistics*, 2001, pp. 129–134.
- [52] D. S. Johnson, Approximation algorithms for combinatorial problems, *Journal of Computer and System Sciences* 9 (1974) 256–278.
- [53] R. Raguram, C. Wu, J.-M. Frahm, S. Lazebnik, Modeling and Recognition of Landmark Image Collections Using Iconic Scene Graphs, *International Journal of Computer Vision* 95 (2011) 213–239.

- [54] R. Arandjelović, A. Zisserman, Name that Sculpture, in: ACM International Conference on Multimedia Retrieval, 2012.
- [55] R. Arandjelović, A. Zisserman, Three things everyone should know to improve object retrieval, in: IEEE Conference on Computer Vision and Pattern Recognition, 2012.
- [56] H. Jégou, M. Douze, C. Schmid, Product Quantization for Nearest Neighbor Search, IEEE Transactions on Pattern Analysis and Machine Intelligence 33 (2011) 117–128.
- [57] H. Jégou, M. Douze, C. Schmid, On the burstiness of visual elements, in: IEEE Conference on Computer Vision and Pattern Recognition, 2009.