



**Politecnico
di Torino**

POLITECNICO DI TORINO

Computer engineering
Academic year 2022/2023
Graduation period December 2022

Master thesis project

Optimizations and efficient retrieval solutions for large-scale visual Geo-localization problems

Supervisors:

Candidate:

Masone Carlo

Gambino Matteo

Contents

I	Abstract	V
II	Acknowledgments	VI
1	Introduction	1
2	State of the art	3
2.1	Datasets	3
2.2	Visual Geo-localization	4
2.2.1	Fine-tuning CNN Image Retrieval with No Human Annotation [8]	5
2.2.2	NetVLAD: CNN architecture for weakly supervised place recognition [12]	5
2.2.3	Learned contextual feature reweighting for image geolocation [13]	6
2.2.4	Self-supervising fine-grained region similarities for large-scale image localization [14]	6
2.2.5	Rethinking Visual Geo-localization for Large-Scale Applications [1]	7
2.2.6	Results comparisons	8
3	Indexing techniques	10
3.1	Exhaustive search	10
3.2	Hierarchical Navigable Small Worlds	11
3.3	Inverted File Index	13
3.4	Product Quantization	15
3.5	Inverted Multi-Index	17
3.6	Product Quantization and re-ranking with source codes	19
3.7	Optimized Product Quantization	21
4	Experiments and results	23
4.1	Testing methodology	23
4.2	Trend of memory usage and time per query with descriptor database size	26
4.3	Inverted File Index with Product Quantization	28
4.4	Inverted File Index with Optimized Product Quantization	30
4.5	Inverted File Index with Optimized Product Quantization and Re-Ranking	32
4.6	Inverted Multi-Index with Optimized Product Quantization	36
4.7	Hierarchical Navigable Small Worlds with Optimized Product Quantization	38
4.8	Results comparison	41
5	Descriptor database reduction	46
6	Conclusions and further developments	52

A	Appendix A - All experimental results with the indexing techniques	54
A.1	Inverted file index with product quantization	54
A.2	Inverted File Index with Optimized Product Quantization	57
A.3	Inverted File Index with Optimized Product Quantization and Re-Ranking	62
A.4	Inverted Multi-Index with Optimized Product Quantization	67
A.5	Hierarchical Navigable Small Worlds with Optimized Product Quantization	71
B	Appendix B - Results with the database reduction techniques	78

List of Figures

1	Example search with the HNSW algorithm	12
2	Voronoi diagram obtained by applying k-means clustering to random points.	14
3	Example representation of the application of product quantization to an 8-dimension vector with m set to 4 and 3 clusters per group	16
4	Trend of the RSS memory usage at the variation of the database size for different numbers of dimensions. The axes are in logarithmic scale for graphical reasons.	26
5	Trend of the time required per each query at the variation of the database size for different numbers of dimensions. The database size axis is in logarithmic scale.	28
6	RSS memory usage and R@1 metric for the IVFPQ method with various configurations.	29
7	RSS memory usage and R@1 metric for the OPQ IVFPQ method with various configurations.	31
8	Comparison between the RSS memory usage and R@1 metric for the OPQ IVFPQ method and IVFPQ with various configurations.	33
9	Comparison between the RSS memory usage and R@1 metric for the IVFPQR method and the OPQ IVFPQ REFINE method with various configurations.	34
10	This figure reports the RSS memory usage and R@1 metric obtained for the IMI method combined with optimized product quantization with various parameters configurations.	36
11	The RSS memory usage and R@1 metric measured for the Hierarchical Navigable Small World (HNSW) method combined with Optimized Product Quantization (OPQ) with various parameters configurations.	39
12	All of the results obtained with the different number of descriptors number of dimensions regardless of the indexing technique used.	41
13	This graph reports all of the results measured for each of the indexing techniques analysed with 1024 dimensions.	42
14	Trend of the recall values with different number of IVF cells probed with the various indexing techniques.	44
15	Trend of the recall values with different number of sub-vectors in the product quantization used on 1024 dimensional descriptors. The configuration used are the one reported in table 8 with only the PQ sub-vecs. parameter changed. The number of bits for each dimensions of the sub-vectors is always 8. The radius of the points is proportional to the approximated RSS memory usage of each configuration.	45
16	An example of panorama used in SF-XL dataset	46
17	On the left the recall measured with the techniques presented for various database sizes, on the right an extract of the graph on the left which focuses on the higher recall values measured.	49

List of Tables

1	Recall values obtained with the presented methods on the various dataset described in section 2.1. * Only the photos taken after 2010 are used for the queries, and the other ones for the database. ** Only the photos taken after 2015 are used for the queries, and the remaining images are used for the database.	9
2	Results of the exhaustive search on different number of descriptors dimensions	26
3	Results obtained with IVFPQ index for some of the best configurations for each descriptor dimensionality.	30
4	Results obtained with OPQ IVFPQ method for some of the best configurations for each descriptor dimensionality.	32
5	Results obtained with IVFPQR and OPQ IVFPQ REFINE method for some of the best configurations for each descriptor dimensionality.	35
6	The results that have been measured for some of the best configuration tested using the Inverted Multi Index with the different number of descriptors dimensions. For the product quantization are always used 8 bits per sub-vector.	37
7	Results obtained with OPQ HNSW method for some of the best configurations for each descriptor dimensionality. The number of bits for product quantization is always 8 and m represent the number of links for each node in the HNSW graph.	40
8	The best results obtained for different configurations of the various indexing techniques. The descriptors used have 1024 dimensions. In the headers of the table the R is indicating the number of sub-vectors used for the re-ranking with compressed error (when this is used the number of bits for product quantization is always 8). The parameter m is representing the number of links per each node in HSNW technique.	43
9	Some of the best results obtained with the database reduction techniques presented. The n parameter is the number of panoramas per zone kept. The time and the RSS memory measured are obtained with the exhaustive search.	50

List of Algorithms

1	HNSW index creation	13
2	IMI index construction pseudo-code	18
3	A pseudo-code version of OPQ algorithm	22

I Abstract

Visual Geo-localization is the task of determining the location where a photo was taken, exploiting only visual information. This task plays an important role in numerous applications, such as in the categorization of images for photo collections, augmented reality, and for the localization of mobile robots, therefore it is an active area of research. The task is commonly approached as an image retrieval problem: given a query, its location is inferred by performing a similarity search, via k-nearest neighbour (k-NN) over a database of Geo-tagged images. While this solution allows to achieve remarkable results in moderately sized problems, it does not scale well to large maps. This is due to two problems:

- the execution of the k-NN requires to keep in memory the embeddings extracted from all the images in the database; as the database increases, the required memory can quickly become infeasible
- the time required to perform the k-NN grows linearly with the dimension of the database; as the database grows, the latency for processing a single query may become not sustainable

This thesis addresses these problems and it investigates various solutions to improve the retrieval pipeline of visual Geo-localization methods. In particular:

- it investigates the impact of using advanced indexing techniques in the similarity search, such as inverted file indexes and product quantization. The effect of these techniques is evaluated in terms of the accuracy of the final results, memory footprint and time required to perform the retrieval. The results of this investigation not only demonstrate that using appropriate indexing techniques can enable VG problems to scale to large database, with minimal loss in accuracy, but they can also serve as a guideline for developers to choose the right solution based on their required trade-off between resources and performance.
- It presents a novel solution for filtering the database images based on their semantic content, in order to reduce the search space for the similarity search, thus ultimately decreasing its memory and time complexity.

All these analyses and solutions have been developed on realistic urban datasets of various scales.

II Acknowledgments

First of all, I want to thank my supervisor prof. Carlo Masone for his help and support during this work as well as Gabriele Berton and Gabriele Trivigno that have helped me with always good suggestions to solve the problems that have raised during the experimental phase of this thesis. Thank also for your infinite availability and promptness to answer all of my requests. I want to thank also all of the members of the Vandal group at Politecnico di Torino that have always included me in their activities and for the good time that we have spent together.

I would like to thank also the HPC@POLITO centre and staff that have given me all of the computational resources fundamental for this thesis project. Thank because without your support and resources this thesis would have not been possible.

After that I want to thank all of my friends that have always been with me. Thank for all of the good moments that we have spent together. Thank for the happy time that we have shared. Part of this result is also yours.

Thank to Marianna that has always supported, helped and made me smile during all of the time of this project. This achievement is also yours.

I want to thank also my family and my sister Sara that have encouraged and pushed me to do my best. Thank for having always believed in me. Without you nothing of this would have been possible. Thank for being constantly with me.

1 Introduction

In the field of computer vision, visual Geo-localization is the task of estimating where a given photo has been taken by only analysing its visual characteristics. This very challenging task has been usually tackled by finding the most similar image inside a huge database of photos where the exact position of where they have been taken is known. As an example, to just apply this technique to the San Francisco city as reported in [1], a database of about 2.8 millions of images is used. A trivial solution to this problem would be to make a similarity comparison with every image inside the database and return only the most similar ones but, given the size of the database involved and the number of images to compare, this is clearly unfeasible.

Most classical approaches to this problem divide the search in two main steps. The first one is the construction of the database used for the retrieval of the most similar image. This part usually consists in extracting visual descriptors from the photos of the dataset that have also the information of where they have been taken. Those descriptors are vectors that encode all of the relevant visual characteristics of the images. This is done to reduce the final size of the database to be used during the following steps and eliminating the features of the photos that do not carry information on the geographical position of where the image was taken. For that reason, it's fundamental to reduce as much as possible the size of the vectors containing the image descriptors limiting the number of dimensions of each of them but, at the same time, to maintain, as much as possible, all of the features that describes the original image. Those descriptors can be extracted using machine learning algorithms, and, more specifically, deep convolutional neural networks (CNN) are actually providing the best results as will be described in section 2.

The second part of the algorithm instead is focused on finding the descriptors, inside the database, which are the most similar to the one of the images to be localized (called queries). Since the number of records present in those databases is usually high, indexes are used to speed up the search inside them. Those indexes must be kept in memory in order to maintain all of the advantages that they provide. Unfortunately, this is not feasible in a large-scale application since the size of those indexes is proportional to the database number of elements and it would be too large to fit in most servers main memory.

For that reason, the aim of this work is to find a solution to drastically reduce the size of the indexes produced by maintaining the same accuracy as the one would have been achieved if a full-size index would have been used and, at the same time, maintain the time needed to retrieve results from the database acceptable.

This thesis is organized by following the chronological development of the project and analyses the techniques known in literature to optimize the memory footprint of an application that uses indexed search to scan the content of a huge database. So, starting from analysing the actual state of the art in visual Geo-localization, by presenting some interesting publications and their results, in order to better compare the one obtained in this work, this project analyses the most relevant indexing techniques present in literature in order to solve the problems presented before in the visual Geo-localization algorithm. After that the project continues by illustrating the experiments conducted in order to assess the performances of each of the techniques analysed by comparing the memory

used and the average time needed to extract the location of a single image. At the end of this thesis project is also presented a novel technique that is useful to reduce the database number of elements by discarding some of the images based on their visual content.

So, summing up, the main objective of this thesis is to find a way to actually keep the efficiency of an indexed search on a large-scale database while reducing the resource usage.

2 State of the art

In order to tackle the problem of visual Geo-localization presented before, many approaches have been proposed in literature. In this section, starting from reporting some examples of datasets widely used, is reported a brief overview of some of the most relevant publications that analyses this topic by presenting the strategy adopted to correctly extract the position of an image by studying its visual characteristics. A common approach of all the methods to predict the location of an image by its only visual characteristics is to split the algorithm in two main steps:

1. **Descriptors database building.** During this step the main objective is to extract vectors of features that describe what is illustrated in the input images (those vectors are commonly called descriptors). This usually require a convolutional neural network in order to analyse the photo and understand what information is important to maintain in order to locate it.
2. **Location retrieval.** After the database of descriptors has been built, it's possible to predict the location where some unseen before photos (commonly called query images) have been taken. In order to do so, the same network or algorithm used in the dataset building phase is used to produce the descriptors from the query image. After that, a K-Nearest Neighbour (K-NN) is used to extract from the descriptors database the K most similar one to the query images. This operation is usually performed with small values of the K parameter like, for example, 1, 5 or 10 and the predicted location of the image can be approximated by using those results. Searching in the database performing a K-NN is computational expensive, since it requires to maintain in main memory the whole database of descriptors and to perform a large number of comparisons between the vectors. For that reason, a standard solution is to adopt optimized nearest neighbours retrieval algorithms that approximate the search procedure of the standard K-NN. Those algorithms use many strategies that include reducing the search space to limit the time needed for the retrieval and quantize the descriptors vectors to reduce memory usage. Those techniques will be widely discussed and analysed in section 3.

Most of the works that will be reported do not provide information on the retrieval part of the algorithm since the standard K-NN search is used. Despite this, it's important to notice, while reporting the various state of the art techniques, the number of dimensions of descriptors used and the results obtained to obtain a wider view on the experiments that will be reported in section 4.

2.1 Datasets

Some of the most relevant datasets built for the visual Geo-localization task are:

- **Pitts250k:** This dataset, presented in [2], contains 250 thousand images taken in the city of Pittsburgh in Pennsylvania. The images contained are taken only in good light conditions and the scenes reported are mostly urban. It exists also a reduced

version of this dataset (called Pitts30k) that contains only a subset of 30 thousand images taken from the full size one.

- **Tokyo 24/7**: this dataset, presented in [3], contains about one thousand images taken by a smartphone camera in the city of Tokyo in Japan. Those images are taken in different illumination conditions and during different hours of the day. This is a very challenging dataset since the variations of light may make harder for the CNN used to extract the relevant features to build the descriptors.
- **MSLS**: this dataset, presented in [4], contains more than 1.6 million images taken from 30 major cities located in six different continents. Due to this sparsity, this dataset contains many different locations and light conditions
- **St. Lucia**: this dataset, presented in [5], contains about 13 thousand images taken in the e suburb of St. Lucia located in Brisbane, Australia. The photos contained are taken from a moving car and represent mostly rural and suburban scenes.
- **SF-XL**: this dataset, presented in [1], contains about 41.2 million images in its training set and 2.8 million in the dataset used at test time. All of these photos are taken in the city of San Francisco in California. The images present in it are taken from Google Street View imagery along different years. This dataset, together with the database images, contains also two different sets of test images.

Those datasets are the one that are most commonly used for this task and all of the state of the art methodologies that will be presented evaluate their performances on those datasets. The SF-XL dataset, due to its dimension, has been selected for this thesis project and it is used, in part or in the full version, in the section 4 in order to evaluate the different indexing techniques analysed.

2.2 Visual Geo-localization

In this section some of the most relevant publications on the task of visual Geo-localization will be reported with the objective to give a strong background and better explain the reasons of the decisions taken during the next sections of this thesis project. Starting from some shallow learning approaches and the usage of hand-crafted descriptors, like the one proposed in [6] or [7], the research has been pushed to find a way to extract better descriptors from the images. In fact, the methods mentioned above use simple feature representations that require human interaction and, for that reason, usually leads to sub-optimal results. More recent works, instead, are concentrated on how to extract representative features from a given image using a convolutional neural network and avoiding the human interaction. In the following sections some of the publications on the usage of CNN in the extraction of descriptors and the novelties and improvements of each of the methodologies presented will be discussed.

2.2.1 Fine-tuning CNN Image Retrieval with No Human Annotation [8]

The method proposed by this paper uses contrastive learning in order to learn which features are really important and should be kept inside descriptors and which are instead not relevant because not representative of the picture location like, for example, the sky or the road. This is done by using a specific loss function to train a fully convolutional network followed by a novel pooling layer. The method proposed uses well known network architectures like Alexnet [9], VGG [10] or ResNet [11] replacing the fully connected layers present in their specifications with a generalized mean pooling layer that produces the final output. This last layer, that is one of the main novelties of this publication, applies the following equation to the feature maps produced by previous layers

$$f = [f_1 \dots f_k \dots f_K], f_k = \left(\frac{1}{X_k} \sum_{x \in X_k} x^{p_k} \right)^{\frac{1}{p_k}} \quad (1)$$

where f represents the final output feature vector, X_k represent the k-th feature map produced by the k filters of the last convolutional layer and p_k that represent the pooling factor. This pooling factor indicates how large are the areas on which the network focuses most on and can be fixed or automatically learned during the training process. Another important characteristic of this publication is the usage of a Siamese learning together with a contrastive loss during the training phase. In practice, the network doesn't directly learn how to produce the correct vectors for each image but tends to produce similar descriptors for images with similar visual characteristics and different vectors for not similar ones. In order to apply that procedure, the training set is composed by couples of images and the network is trained on computing similar vectors for the photos in the pair that have been taken close to a query one and different ones for the further images. Finding the correct couples to be fed to the network is fundamental for a correct training. In fact, it is important that images that are taken on the same place are grouped together especially if their visual characteristics are very different (those images are usually called soft positives) and, instead, photos that do not belong to the same place should be placed together if they are visually similar (called hard negatives). Results obtained by this method are discussed in section 2.2.6

2.2.2 NetVLAD: CNN architecture for weakly supervised place recognition [12]

This publication proposes the introduction of a VLAD layer at the end of a fully convolutional network in order to produce Vector of Locally Aggregated Descriptors (VLAD). Those vectors are commonly used in the field of instance level retrieval and image classification and codify, inside their elements, the sum of differences between the produced local descriptor and the closest cluster centroid (corresponding to a visual word) to which the descriptor is assigned. Those centroids are initialized with random values and, after each training step, are re-computed. It is important to notice that this function is not differentiable for the operation that assigns the local descriptors to the centroids. Since the production of the descriptors is introduced inside the training process of a convolutional neural network, all parts of the function that computes the descriptors must be

differentiable. The authors of the publication solve this problem of not differentiability of the assignment of each local descriptor to a visual word by introducing the concept of soft assignment. This is done by partially assigning the local descriptor to each cluster center by using the soft-max operator following

$$V(j, k) = \sum_{i=1}^N \frac{e^{W_k x_i + b_k}}{\sum_{k'} e^{W_{k'} x_i + b_{k'}}} (x_i(j) - c_k(j)) \quad (2)$$

where $V(i, j)$ represent the element at position j, k of the matrix containing the VLAD descriptors, W_k, b_k, c_k are trainable parameters and x_i is the i -th of the N local descriptor. The resulting layer can be used on top of known architectures after removing all of the fully connected layers. As the method reported in section 2.2.1, the approach illustrated in this publication uses a contrastive approach by refining the method of building the training dataset. In fact, instead of using only couples of images, it computes a set of samples composed by some hard-negative examples (images not corresponding to the location where the target photo has been taken but visually similar to the target one) and some positives examples. The results that can be obtained by using this approach are reported in section 2.2.6.

2.2.3 Learned contextual feature reweighting for image geolocation [13]

The work reported in the section 2.2.2 with the introduction of the VLAD layers has inspired many different publications like the one reported here. In fact, starting from the architecture of a NetVLAD network, this paper has introduced some novelties in order to make the descriptors produced of a higher quality like the usage of attention modules. Those modules are widely used in recent architectures since they help the network to focus more easily on the relevant part of an image and for that reason are suitable for the adoption in a visual Geo-localization application. The authors of the paper [13], propose to use a so-called Contextual Re-weighting Network (CRN) to, as the name of the module says, re-weight the features map by considering the full context of the image. This module can be easily integrated with the NetVLAD implementation since its output can be used as a multiplication factor for the values produced by the soft assignment part of the network (described in 2.2.2). This CRN module is composed by various convolutional filters with different kernel sizes that have the duty to analyse the input features (produced by a standard CNN like ResNet [11] without its fully connected layers) at various scales in order to get the full context of the image and at the same time not to ignore the small details that may be present. As for the approaches analysed in the previous sections, this network is trained by using a contrastive approach. The results of this approach will be discussed in section 2.2.6.

2.2.4 Self-supervising fine-grained region similarities for large-scale image localization [14]

In this work is proposed a technique to localize the query image by progressively refining the image-to-region similarity. In fact, how the authors of the publication notice, GPS

position data introduce noise into the labels produced since two images taken in the same place can be completely different if, for example, are taken in opposite directions. Their method proposes so to use a self-supervised approach (starting from unlabelled data the model learns how to get supervision from the data itself) in order to make the network learn how to produce similar vectors for similar images. In fact, starting from a VLAD network like NetVLAD [12], the authors of this work propose to compute similarity vectors by following

$$S_{\theta_1}(q, p_1, \dots, p_k : \tau_1) = \textit{softmax} \left(\left[\frac{f_{\theta_1}^q \cdot f_{\theta_1}^{p_1}}{\tau_1}, \dots, \frac{f_{\theta_1}^q \cdot f_{\theta_1}^{p_k}}{\tau_1} \right]^T \right) \quad (3)$$

where θ_1 is a set of weights on which the features depends on, q is the query image, p_i is one the N positives per each query, $f_{\theta_1}^{p_i}$ is the feature map produced by using the weights θ_1 on the image p_i , τ_1 is a parameter called temperature that is characterizing the sharpness of the produced similarity vector. The usage of those vectors permits to use also the information coming from the parts of the images that are not similar. The parameter τ is reduced during the training process since a large value of it produces more equally distributed vectors that are not suitable for the task. In addition to that, to make the network more sensible also to the small overlapping regions of the query and positive image, also crops at some fixed positions are used as positives. The results that can be obtained by the usage of this network (called SFRS in the rest of this work) will be presented in the section 2.2.6.

2.2.5 Rethinking Visual Geo-localization for Large-Scale Applications [1]

This publication focuses on two of the main problems of the other methodologies described above: the memory used during the training of the networks and the number of dimensions of descriptors produced. In fact, in all of the works mentioned before are not suitable for world scale applications since for example at training time the whole database of descriptors has to be recomputed periodically in order to update the descriptors with the new learned features and those descriptors must be kept in a cache. This procedure is not applicable for big datasets since the amount of memory required to store that cache grows linearly with the number of training images. In addition to that, networks like NetVLAD produce descriptors with a very high number of dimensions (e.g. if VGG-16 [10] is used as backbone a single descriptor has 32k dimensions without any further elaboration or 4096 dimensions if reduction techniques like PCA are adopted).

By taking inspiration on the face recognition task, the authors of the paper have casted the problem of producing the descriptors vectors for the retrieval step as a classification task. This not only solves the problem of recomputing during training the cache of descriptors but permits also to reduce the dimensionality of descriptors. In order to give a class to each image of the dataset, the geographical UTM coordinates and the orientation of the camera that has taken the photos are quantized, by truncating the division of the actual value of those metrics to a parameter, and the final class is composed by the union of those values. During the training phase, only not adjacent classes are fed to the network, in order to avoid the problem of having very similar images to belong to different

classes (for quantization errors). Also the loss function adopted is different from the one used by the publications mentioned before since it does not use a contrastive approach and do not require caching. This function is known as Large Margin Cosine Loss (LMCL) loss, as described in [15], and it is able to make the network learn how to separate the various classes with a large margin. After the training procedure, the network can be used not to produce a class for each image but only to produce descriptors. The descriptors produced can be characterized with a dimensionality that is much smaller than the one produced by the NetVLAD network for example but at the same time they guarantee performances that are higher as will be discussed in section 2.2.6. This publication is the most suitable for a large-scale application of a visual Geo-localization network among the one analysed since it already produces acceptable descriptors sizes (with 512 dimensions) and the training procedure has been proved to be much faster and less memory intensive of the one of a NetVLAD network. Those characteristics permit to use bigger datasets and for that reason this method has been used in the rest of this work in order to compute the descriptors needed to test the various indexing methodologies.

2.2.6 Results comparisons

In this section is reported the comparison among the various networks reported. The metric that is commonly used in the field of visual Geo-localization is the recall. Given a set of N location prediction extracted from a K-NN algorithm, the value of the recall metric is computed as:

$$R@N = \frac{1}{N} \sum_{i=1}^N c(p_i) \quad (4)$$

where p_i is the i -th of the N predictions and $c(x)$ is described as

$$c(x) = \begin{cases} 1, & \text{if } x \in P \\ 0, & \text{Otherwise} \end{cases} \quad (5)$$

where P is the set of images that are taken in a given range (called positive threshold) from the query image. This metric is computing the ratio between the correctly localized images and the number of predictions done. This metric is also the one that will be used in the rest of this thesis project to report the results obtained. The results that are presented are, when possible, taken from the original publication or, when the paper does not compute them, have been calculated from the authors of [1]. The recall values obtained with various dataset are reported in table 1.

How it's possible to notice from table 1, CosPlace network [1] is the one that is producing best results on most of dataset although it's descriptors number of dimensions is only 512. It's important to notice that the comparison reported in table 1 it's not completely fair since the training dataset used are not the same for all networks. This is not avoidable since some methods can't be trained on some datasets. In fact, many methodologies presented require a high number for the descriptor's number of dimensions and it's not possible to store the whole cache of descriptors in main memory. As an example, the NetVLAD network can't be trained on huge datasets like SF-XL since to store all

Method	Desc.	Train set	Pitts250k		Tokyo 24/7		MSLS		St. Lucia		Average	
			R@1	R@5	R@1	R@5	R@1	R@5	R@1	R@5	R@1	R@5
NetVLAD [12]	32768	Pitts30k	85.9±0.3	93.6±0.2	62.2±0.7	75.4±0.5	54.8±1.1	66.6±1.0	70.8±0.5	81.8±0.7	72.0	82.3
NetVLAD [12]	32768	MSLS	79.7±0.8	90.2±0.9	63.6±1.1	77.5±1.2	75.4±0.4	84.2±0.2	92.8±0.4	97.6±0.3	78.5	88.0
NetVLAD [12]	32768	SF-XL*	81.5±0.1	90.8±0.3	66.3±0.4	78.2±0.5	59.3±0.8	68.7±0.7	80.6±0.5	90.7±0.8	74.0	84.0
NetVLAD [12]	32768	SF-XL**	76.2±0.3	88.8±0.4	56.2±0.8	67.7±0.9	51.4±0.9	60.7±0.7	78.5±0.6	88.3±0.8	68.2	79.1
CRN [13]	32768	Pitts30k	87.0±0.4	94.5±0.1	62.8±0.5	77.4±0.8	57.6±0.3	70.4±0.8	70.9±0.7	82.8±0.7	72.9	83.9
GeM [8]	512	Pitts30k	75.3±0.2	88.4±0.3	46.4±0.9	65.3±0.7	51.8±0.9	64.4±0.9	59.9±1.6	76.3±2.0	62.3	77.0
GeM [8]	512	MSLS	65.3±1.2	81.0±1.6	44.9±1.7	62.6±1.2	66.7±0.7	78.9±0.5	84.6±1.1	93.3±0.7	66.6	80.2
GeM [8]	512	SF-XL*	64.7±0.8	81.4±0.8	37.9±2.3	51.0±2.1	46.8±2.1	58.1±1.2	68.5±2.4	82.7±1.8	57.1	71.4
SFRS [14]	4096	Pitts30k	90.1±0.3	95.8±0.2	78.5±0.8	87.3±0.5	62.8±0.8	73.0±0.8	72.5±3.6	85.4±3.2	78.5	87.1
CosPlace [1]	512	SF-XL	89.7±0.0	94.5±0.0	82.8±0.9	90.0±0.6	79.5±0.1	87.2±0.2	94.3±0.4	97.4±0.3	87.0	93.1

Table 1: Recall values obtained with the presented methods on the various dataset described in section 2.1. * Only the photos taken after 2010 are used for the queries, and the other ones for the database. ** Only the photos taken after 2015 are used for the queries, and the remaining images are used for the database.

descriptors, supposing a number of dimensions of 4096 after the application of PCA to the raw output of the network, would require more than 600GB of main memory. For the analysis that will follow it’s fundamental to notice that all of the methods reported in the table 1 use an exhaustive search, based on a standard K-NN, for the retrieval part of the algorithm and, for that reason, the results reported are the best obtainable with each method.

In the rest of this work the method that will be used to extract descriptors for the comparison between the indexing techniques is CosPlace [1] since the results obtained are representing the state of the art on most of the datasets analysed and this method permit to obtain those results with a very low descriptors dimensionality (compared for example with the one obtained with the SFRS [14] method that require descriptors 8 times larger).

3 Indexing techniques

In this section, the retrieval part of the visual Geo-localization algorithm will be widely analysed especially focusing on the indexing techniques that can be used to reduce the resources needed to extract the N most similar descriptors from the one contained in the pre-computed database of Geo-tagged photos and, in that way, predict the location where the query image fed has been taken. This is one of the most critical parts of the algorithm in a real world application since the dimension of the descriptors dataset (that has to be kept in main memory while performing the nearest neighbour search) would be extremely large and, in order to serve multiple queries by the different users, the efficiency of the retrieval part is crucial.

For that reason, by starting from analysing the simplest approach possible, this section wants to analyse some of the possible different ways to construct an index that is able to reduce firstly the memory footprint of the naive approach and at the same time maintain the execution time acceptable by both reducing the number of dimensions of the descriptors, applying techniques like product quantization, and implement a partial search that reduces the resources used.

3.1 Exhaustive search

This simple approach is the one used by every method presented in section 2 since it does an exhaustive search on the database of descriptors. This method can use an index that simply stores the vectors in an uncompressed way in order to produce the best results obtainable with the descriptors stored and use the standard K-NN search algorithm to retrieve the requested nearest neighbours.

This approach is not memory efficient since the index dimension is exactly the same as the size of the whole database and, for that reason, is not suitable for a real-world application since the amount of memory required would not be acceptable. When a descriptor vector is fed to this indexing technique, a sequential scan is performed on each element of the database in order to find the N most similar vectors. This search strategy also is not efficient since all vectors have to be compared with the query one in order to compute a similarity function like the Euclidean distance between the two vectors under analysis. For that reason, the complexity of retrieving the nearest neighbours from the database is $\mathcal{O}(DN)$ where D is representing the number of dimensions of each vector and N is the size of the database to be analysed. It is also possible to notice that the time required to retrieve the nearest neighbours does not depends on the K parameter of the K-NN since it is possible to scan the descriptors database just one time with any value of K .

The main advantage of using this approach is the fact that produces always the best result possible and, for that reason, it has been used as baseline in the rest of this thesis project in order to compare the different indexing techniques. The algorithm for this approach has been integrated in the test environment built by using the Faiss library [16] developed by the Facebook AI Research group helped by some external contributors. This open source library it's highly optimized and permits to apply the indexing techniques on

a billion scale databases and for that reason perfectly suits the use case proposed in this work.

3.2 Hierarchical Navigable Small Worlds

The first technique that has been analysed to solve the problems reported in the previous sections is the one called Hierarchical Navigable Small Worlds (HNSW). This approach uses graph theory to produce a graph of proximity between the various vectors (usually the proximity is computed as the Euclidean distance between the two vectors and nodes are connected together based on the proximity value) and probabilistic skip lists to speed up the search algorithm avoiding the visit of some nodes in the graph. For that reason, this approach can be classified as an Approximate Nearest Neighbour search (ANN) and the results provided may be lower than the baseline.

The skip lists mentioned before consist in a layered structure of linked lists, one for each vector in the database, connected together by some links where to longer connections correspond to a higher number of nodes not analysed during the search. Only the nodes of the list at the same index can be connected together forming in that way a layer. Each of those layers is a graph that is called Navigable Small Worlds (NSW).

The resulting graph is, by construction, composed by a set of NSWs organized in a way that higher index layers correspond to longer links and, for that reason, to a higher number of nodes skipped. By using this approach, it is possible to guarantee a logarithmic search time complexity. In fact, by starting from a fixed start node in the top layer, the search algorithm proceeds by selecting, among the nodes directly connected to the current one, the closest to the query vector and the search continues by looking at the neighbours of the chosen node. This step is repeated until all nodes, connected to the one under analysis, are further than the current one finding in that way a local minimum of the distance function for the query at the top layer. At this point the search can be refined by moving to the same instance of the node on the layer that is just below the current one and restarting the search on the new NSW. An example of this search approach can be found in figure 1.

The construction of such graphs is expensive since each vector is analysed sequentially and inserted in the graphs. So, given an empty graph composed by L layers, the insertion proceeds as follow:

1. Firstly, the layer where the descriptor under analysis has to be inserted is selected. This is done by randomly selecting the index of the layer among the L possible or, only for the first vector to be inserted, the top layer is used (the one at index $L - 1$).
2. After the layer selection, starting from index $L - 1$, the insertion process begins by searching in this layer the closest node to the one to be added to the graph and, when it is found, the search process is moved to the layer just below the current one. This step is repeated until the insertion layer is reached (the one selected at point 1).
3. When the layer where the node has to be inserted is reached, a set of C search operations is performed on the insertion layer in order to find the M closest nodes

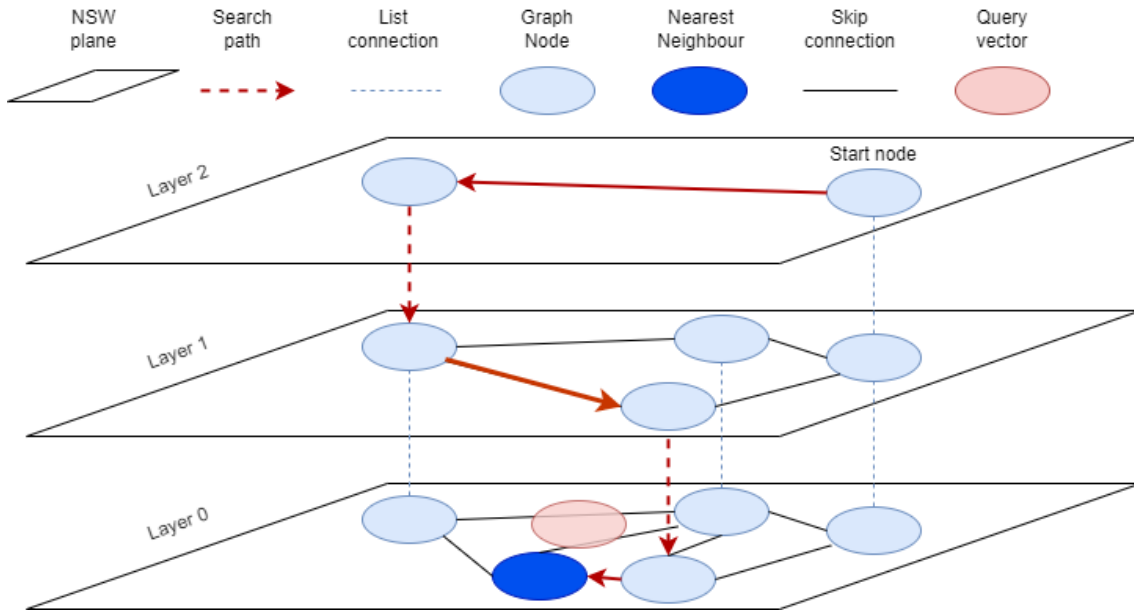


Figure 1: Example search with the HNSW algorithm

to the one to be inserted. The C parameter is usually set higher than the M one in order to permit to the algorithm to choose which node, among the C candidates extracted, has to be connected to the one to be inserted, following a criterion. In the used implementation, only the closest nodes are connected together. This process is actually inserting the descriptor under analysis in the NSW at the insertion layer. It is possible that some nodes have less than M links at the end of the construction of the graph because of two main reasons: the graph when the node has been inserted was containing less than M nodes or some of the nodes among the C extracted does not satisfy the conditions included in the criteria used for the candidate selection.

4. Later, the step 3 is repeated on all of the layers below the insertion one terminating in that way the insertion process of the descriptor under analysis.

This process is then repeated for each of the descriptor in the database producing, in that way, a hierarchical structure of navigable small world. A possible pseudo-code implementation of the described algorithm can be found at algorithm 1.

As it's possible to notice from the graph construction details, this method is not actually reducing the memory footprint of the index compared to the baseline method since for each vector the memory usage in byte can be computed following

$$Mem = 4d + 8M \quad (6)$$

where d represent the number of dimensions of the input vectors and M is the number of connections per node of the resulting graph and it is assumed that each dimension of the vector and each link of the graph is requiring 4 bytes of memory. For that reason, the final

Algorithm 1 HNSW index creation

```

 $X \leftarrow$  set of all database descriptors
 $L \leftarrow$  number of desired layers
for  $x \in X$  do
  if it's not the first vector then
     $l \leftarrow \text{random}(1, L)$ 
  else
     $l \leftarrow 0$ 
  end if
   $N_{start} \leftarrow \text{HNSW}_{search}(x, \text{until layer } l)$ 
  while  $l \geq 0$  do
     $N_{connect} \leftarrow M$  nearest neighbours at layer  $l$ 
    Link all  $N_{connect}$  to  $x$ 
     $l \leftarrow l - 1$ 
  end while
end for

```

memory usage will be higher than the descriptors database size which it is not acceptable for large databases of descriptors.

This technique is still representing a valid approach since this method can be combined with other strategies to reduce the vector dimensionality, and in that way the final size of each descriptors, like product quantization as will be described in section 3.4.

Like in the exhaustive search case, the version of the algorithm used during the experiment phase is the one provided by the Faiss library. At search time it's possible to configure the number of neighbours proposal among the results are extracted. It is fundamental to properly tune this parameter since a too high number of proposals will result in slower search times and a low value for this parameter will produce very low values for the recall metric.

3.3 Inverted File Index

This technique is optimizing the efficiency of the retrieval of the nearest neighbours by reducing the search space that has to be analysed for each query vector. This is performed by introducing in the index building a clustering algorithm that has the duty to group together some of the descriptors present in the database and, in that way, permitting to not analyse some of those groups basing the decision only on the comparison of the centroids of the clusters.

In fact, during the computation of the index, a k-means algorithm is used in order to produce k clusters starting from the database vectors. After this computation, each vector in the descriptors database is assigned to the closest centroid among the calculated one building in that way a Voronoi diagram like the one reported in figure 2. The resulting data structure is used as an index for the descriptor database.

At search time, some of the clusters, represented by the cells of the Voronoi diagram,

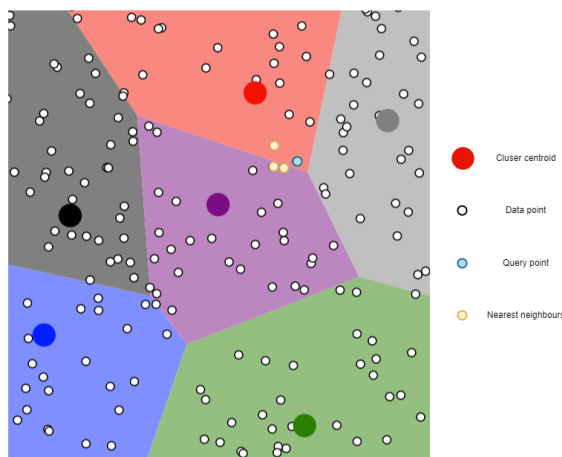


Figure 2: Voronoi diagram obtained by applying k-means clustering to random points.

are selected among the produced one in order to choose the descriptor’s groups that have to analysed to extract the nearest neighbours. This choice is performed by taking the n_{probe} cluster centroids that are most similar to the query descriptor. The query vector is then compared to each point belonging to the selected clusters extracting, in that way, the nearest neighbours desired.

This simple approach can guarantee a high speed in the nearest neighbours computation since the number of database vectors to be analysed can be extremely reduced. The number of clusters analysed while choosing the closest points to the query is usually set greater than 1 since a major problem may arise when the query vector lies very close or on the margin that separates two different clusters among the produced ones. In fact, in this region of space, the assignment to a specific cluster may exclude from the similarity analysis some points that belong to other clusters but have a low value for the Euclidean distance to the query vector as it’s shown if figure 2 where two of the nearest neighbours of the query point would not be reported because belonging to a different cluster. The parameter n_{probe} can vary from 1 to k (and in that case the inverted file index algorithm degrades to an exhaustive search) and it’s fundamental, for obtaining good results, to select a correct value of this parameter in order to analyse the minimum possible number of clusters and at the same time all of the nearest points to the query vector. The Faiss library provide an implementation also for this algorithm that will be used in the various experiment reported in section 4.

It’s important to notice that also this algorithm does not reduce the memory required since no operation is performed on the database vectors that has to be kept in memory for efficient distance computation. This algorithm is instead reducing the time needed to retrieve the nearest neighbours since the number of comparisons done can be much smaller than the ones that would be performed if the standard exhaustive search would have been applied. In fact, the operation of retrieving the most similar vectors from the database has a complexity that is $\mathcal{O}(DN_{dst})$ where D represent the number of dimensions of the input vectors and N_{dst} is the number of elements of the database to be analysed which

can be much smaller if the inverted file index is used. For that reason, it's fundamental to tune the N_{probe} parameter in order to avoid not useful comparisons.

As already discussed for the HNSW index, also this technique can be enriched by the usage of product quantization that would actually reduce the memory footprint of the algorithm. Also this technique is implemented in the Faiss library and it used as the base for many different indexing techniques. In fact, it is possible to apply it, partitioning the search space, and use a HNSW index for each partition combining in that way the benefits of both the techniques.

3.4 Product Quantization

Quantization is a standard technique to reduce the scope of possible values of a continuous value to some fixed symbols. Although this operation reduces the amount of information contained in the output vector, this technique is suitable for the task of reducing the memory footprint of the indexes used in the retrieval part of the visual Geo-localization algorithm since it permits to reduce drastically the size of each descriptor vector.

Starting from a high dimensional vector, like the descriptors, the product quantization algorithm divides each vector in m sub-vector of equal size (in order to permit this operation, the number of dimensions of each vector must be a multiple of the parameter m). Each sub-vector is later associated to a group by considering the position of the partition in the original vector. For all of the resulting m groups, a separate instance of a standard clustering algorithm like k-means is applied. So, for each subspace represented by the sub-vectors belonging to a given position in the original one, a set of clusters centroids is computed. After this operation, each sub-vector is assigned to the closest centroid among the ones calculated for the position in the original descriptors it belongs to. In order to do so, it is sufficient to memorize only the index of the centroid in the list associated to the specific position of the sub-vector in the descriptor thus further reducing the memory required to store each compressed vector. A graphical example of the application of product quantization is reported in figure 3.

In that way, by following the procedure described, it is possible to reduce the final memory usage required to store one of the vectors in the database by a factor

$$R = \frac{nm}{32D} \quad (7)$$

where R is the ratio between the size in bytes of the quantized vector and the original one, 2^n represent the number of centroids used for the clustering algorithm, m is the number of sub-vectors for each descriptor.

This technique can be used before any indexing strategy to reduce the memory used by the index and thus making possible to use large databases of descriptors. It's important to notice that, although the vectors are compressed, it's still possible to approximate the Euclidean distance between two raw vectors in the quantized space. In fact, the output vectors of product quantization contain the indexes of the cluster centroid, for memory efficiency reasons, but it's still possible to approximate the original vector by substituting the index of the centroid with its actual representation in the sub-vectors space. This

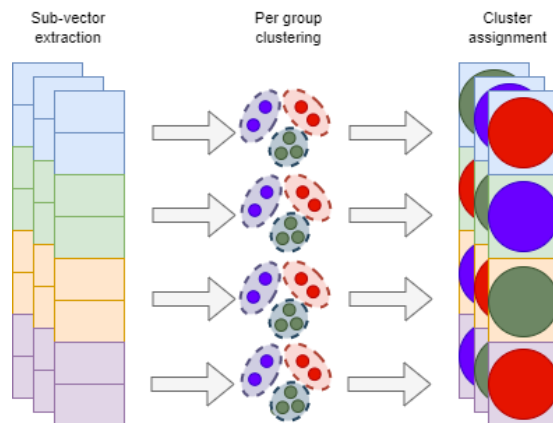


Figure 3: Example representation of the application of product quantization to an 8-dimension vector with m set to 4 and 3 clusters per group

operation can be efficiently performed by using a set of m look-up tables that can be easily stored in main memory since their size is in general small. For example, considering a database of raw vectors composed each by 1024 dimensions and assuming that the product quantization is applied with 128 sub-vectors and for each group of them 256 cluster centroids are computed, the final size of the look-up table to be used for the approximation of the original vector from a compressed one is only 1 MB. This operation of approximating the original uncompressed vector is usually called reconstruction and the main objective of product quantization is to reduce as much as possible the reconstruction error defined as the distance between the original vector and the reconstructed one. This error is influenced by both the number of sub-vectors and the number of clusters computed for each group of sub-vectors. Those parameters do not influence equally the reconstruction error since a higher number of descriptors partitions with a lower number of centroids per group is usually not preferable to a configuration that is requiring the same amount of memory for each descriptor but uses a higher number of centroids and a lower number of sub-vectors. This is due to the fact that a higher number of centroids is associated to a more accurate approximation of the original sub-vectors. Despite it is possible to think that a very large number of centroids should be used to obtain good results, in general, the product quantization technique is applied with a number of clusters for each group that is 2^4 or 2^8 since it can be respective represented in half or in one byte and thus making the storage of the compressed vector memory aligned. For that reason, it is fundamental to properly tune the number of sub-vectors balancing it with the number of centroids.

The distance computation described can be further optimized by avoiding the computation in the original vector space and using a particular distance look-up table to calculate it directly in the compressed space. In fact, during the retrieval part of the algorithm, the query vector has to be compared with many database entries resulting in a cost for the distance computation that is $\mathcal{O}(DN)$ where D is the number of dimensions of each vector and N is the database size. It is possible to pre-compute the distances between the m sub-vectors of the query with the whole set of cluster centroid and storing the results in

a distance look-up table. By using this table, during distance computation it is no more necessary to reconstruct each database vector and the approximated distance can be simply calculated by summing up all the entries in the distance lookup table that correspond to the index of the centroids of the compressed database vector under analysis. In that way, the complexity of the vector distance computation is $\mathcal{O}(2^n m D + m N)$ where 2^n is the number of different centroids per each of the m sub-vectors. By means of this optimized distance computation, it's also possible to reduce the time needed to extract the nearest neighbours of the query from the database with respect to the case of exhaustive search since the m parameter is usually much smaller than the number of dimensions of each descriptor.

The main drawback of the application of product quantization is that some of the information contained in the original descriptors is irremediably lost and so the recall metric value can be lower than baseline. For that reason, it's fundamental to correctly tune the m and n parameters in order to find a good trade-off between the memory reduction and the recall loss produced by the approximations introduced. The product quantization algorithm is integrated in the Faiss library and will be widely used in the section 4 for the extreme memory reduction that can be obtained by adopting by this approach.

3.5 Inverted Multi-Index

The Inverted Multi-Index [17] (IMI) is an indexing technique that tends to fuse the approach of an inverted file index with the main idea behind the product quantization. In fact, inverted indexes tend to produce fine space partitioning in order to reduce the number of points in each partition and so reduce the number of distance computations to be performed. This approach shows some limitations when the database to be indexed is containing a large number of entries since the number of lists that contain the various points belonging to each partition tends to become large. For that reason, the index build and search time tends to become not acceptable because having more lists results in increasing the number of centroids to be calculated by the k-means algorithm during the index creation phase and, at search time, more lists have to be visited to obtain optimal recall values. Given a database of N vectors each composed by D dimensions, an inverted multi-index is computed as follows:

- Firstly, like in the product quantization case, each vector is divided in sub-vectors. Unlike the product quantization case where optimal number of splits to be performed is usually a number between 8 and 32, in this case it is sufficient to split in half the vectors. So, after this step is performed, each descriptor is partitioned in two sub-vectors.
- After that, a standard clustering algorithm like k-means is applied independently to both sets of sub-vectors producing two sets of K clusters centroids (that can be called codebooks). It is now possible to identify a new $K \times K$ space generated by the two sets of centroids and it is possible to visualize it as a grid where each cell corresponds to a different combination of the centroids belonging to the two sets.

- Each descriptor present in the database is then assigned to the closest cell grid by separately considering the closest centroid in both the sets previously built. This procedure is producing a new set of K^2 lists of descriptors.

This result constitutes the inverted multi-index. A possible pseudo-code implementation of the algorithm needed to generate an inverted multi index can be found in algorithm 2.

Algorithm 2 IMI index construction pseudo-code

```

 $X \leftarrow$  set of  $N$  database vectors
 $K \leftarrow$  number of centroids for each group
 $V1 \leftarrow []$ 
 $V2 \leftarrow []$ 
for  $x \in X$  do
   $V1$  append  $x[start : len(x)/2]$ 
   $V2$  append  $x[len(x)/2 : end]$ 
end for
 $C_1 \leftarrow$  k-means( $V1, k$ )
 $C_2 \leftarrow$  k-means( $V2, k$ )
for  $x \in X$  do
   $s1 = x[start : len(x)/2]$ 
   $s2 = x[len(x)/2 : end]$ 
   $i \leftarrow i | c_i = \min(d(s1, c_i)), \forall c_i \in C1$ 
   $j \leftarrow j | c_j = \min(d(s2, c_j)), \forall c_j \in C2$ 
   $l[i, j] \leftarrow x$ 
end for

```

In order to query such an index, a two-step approach is suggested. In fact, as first the query vector is divided in two sub-vectors like was done during the index building phase and, for both of them, the L nearest neighbour are extracted from the cluster centroids previously computed at the position to which the sub-vector belongs to. It is now possible to traverse the lists associated to all the L^2 possible combinations of the selected nearest clusters centroids of the two sub-vectors. This traversal part is done in order. In fact, it is possible to sort the pairs of centroids by considering the sum of distances of the sub-vectors from the analysed centroid. This traversal is, then, stopped when the desired number of nearest neighbours for the query have been chosen. This search part of the algorithm can be efficiently implemented by using priority queues.

How it's possible to notice from the building details of this index, this algorithm produces K^2 different list that, compared with the K that would be obtained with the inverted file index approach, produce a finer partitioning of the descriptors space. This has the major drawback that more memory is required to store the lists. It is possible to overcome this issue by using product quantization on the vectors and store in memory only the compressed versions reducing in that way the memory footprint of the algorithm and it is possible to take advantage of the optimized distance computation permitted by product quantization in order to speed up the nearest neighbours retrieval part. Also this algorithm it's implemented in the Faiss library.

3.6 Product Quantization and re-ranking with source codes

In this section, will be reported how to integrate the product quantization inside an index and how it's possible to refine the approximation of the distance between the query vector and the database ones. Given a set of N database descriptors, it is possible to compress them by using the product quantization algorithm as described in section 3.4. It is possible to substitute the descriptors to be fed to the index to be computed with the compressed version, obtained by applying product quantization, reducing in that way the memory footprint of the whole index. This approach, as discussed before, is introducing an error during the reconstruction phase. For that reason, it is possible to obtain the original descriptor from the quantized version by

$$y = pq(y) + r(y) \quad (8)$$

where y is representing the compressed version of the input vector y , pq is the function that starting from a not quantized version of vector y returns the compressed version of it and $r(y)$ is the residual vector that measures the error related to the loss of information introduced by the product quantization technique. In a simple approach this residual vector is usually ignored accepting that the error introduced by it can be considered negligible, since the product quantization technique tends to minimize it. In addition to that, by increasing the number of sub-vectors and the number of clusters produced per each group of sub-vectors, the reconstruction error asymptotically goes to 0.

In reality, this error value may be considered negligible only if the number of sub-vectors and the number of centroids produced is sufficiently high but, in case of an extreme compression of the database vectors like the one desired in this application, this residual vector should be taken into consideration in order to still obtain high values for the recall metric.

Directly including this vector in the index computation, although possible, is not feasible because it would remove all benefits that can be obtained by applying product quantization since the storage of the vectors that encode the residual error require as much memory as the full-sized descriptor. This approach, although not useful for the analysis that will follow because does not reduce in any way the memory consumption of the indexes, it's available in the Faiss library and the results obtainable with this method will be reported in the section 4 in order to illustrate the results that can be obtained by using also the residual error vector and to compare the results with the method that will follow.

A possible approach to solve the issue of memory usage of the residual error vectors is to use a method called re-ranking with source codes as described in [18]. This approach is based on the usage of a second independent product quantization algorithm instance that has the duty to reduce the memory usage of the storage of the residual vectors. Also for this task, product quantization can be used for its ability to reduce drastically the amount of space needed for the compressed vectors.

The proposed algorithm can be implemented as follows: given the set of residual vectors obtained from the vectors in the database and a parameter m' that indicates the number of dimensions that the compressed vector will have, it is possible to partition the vectors in m' sub-vectors and apply product quantization like described in 3.4. After this step,

the result is a compressed version of the residual vector and m' sets of centroids. It's important to notice that in this case the number of centroids has selected to be always 256 since, in that way, each of the m' dimensions of the output vector can be stored in a single byte. The final dimension on memory of the vectors produced with this approach can be expressed by

$$S = \frac{nm}{8} + m' \quad (9)$$

where the first term in the addition $\frac{nm}{8}$ represent the number of bytes required for the storage of the quantized version of the descriptor and m' is the number of sub-vectors of the residual vector (each of them requires only one byte as mentioned before). The value of the m' parameter has to be selected in order to obtain a good trade-off between the amount of additional memory required and the performance gain that can be obtained by including the residual vector.

At search time, when the k nearest neighbours of a given query image has to be extracted the algorithm works as follows:

1. A set of k' candidates nearest neighbours vectors are extracted for the query by considering only the quantized version of the descriptors and, for the moment, ignoring the compressed residual. This is the standard approach that would have been followed also if the residual would have been ignored. During this step is possible to use the optimized distance computation for product quantization as described in section 3.4.
2. For each of the k' candidates extracted, the distance from the query one is computed taking into consideration also the compressed residual vectors by following

$$d = \|q - (r(x) + r'(x))\| \quad (10)$$

where q is representing the query raw vector, r is the function that reconstruct the original vector from the compressed one, r' represent the function that given the compressed reconstruction error of the database vector x returns the reconstructed one. Also in this case, it is possible to reduce the time needed by the distance computation by computing it directly in the quantized space (although the advantages of this approach are reduced since longer vectors have to be used during distance computation).

3. Among the k' candidates the k requested nearest neighbours are selected by choosing the ones with smallest distance to the query one considering also the compressed residual vector.

This search algorithm is usually requiring a larger amount of time to retrieve the nearest neighbours since the number of operations needed are higher than the one that would be necessary if the standard product quantization technique is used since, before computing the distance, it has to be reconstructed and summed the residual vector to the reconstructed nearest neighbours. The candidates re-ranking and the product quantization can be integrated together in all of the combinations of techniques that will be presented by using the Faiss library.

3.7 Optimized Product Quantization

As discussed in the previous sections, product quantization is providing a solution to the problem of reducing the memory required to store the indexes during the retrieval step of the visual Geo-localization pipeline. It is possible to further optimize the product quantization algorithm in order to make it produce an optimal decomposition of the space by finding the best possible centroids that describe each sub-vector group by using a technique called Optimized Product Quantization [19]. This method is born from the consideration that the product quantization algorithm can be casted as an optimization problem where the objective function is represented by the minimization of the distortion between the raw vector and the corresponding compressed one. This distortion introduced by the quantization approach can be expressed by

$$\mathcal{D} = \frac{1}{N} \sum_{i=1}^N \|x_i - r(c(x_i))\|^2 \quad (11)$$

where \mathcal{D} represent the mean distortion value among all the N vectors of the descriptor database, x_i is the i -th uncompressed vector in the database, $c(x)$ is the function that, given the raw vector, returns the compressed version of it by applying product quantization, $r(x)$ represent the function that reconstruct the original vector from the compressed one and $\|x\|$ represent the Euclidean distance of point x with respect to the origin.

As described in section 3.4, the product quantization algorithm is assigning each sub-vector s_i to the closest centroid c_i in the set of the one computed by the k-means algorithm for the position where the given sub-vector was taken. Given an orthogonal matrix R that represent the operation of partitioning the space described by the descriptor vectors to a set of m sub-spaces, it's possible to notice how if the R matrix is multiplied to both s_i and c_i the assignment does not change. This R matrix is relaxing the constraint inserted when the sub-vectors are created from the descriptors in the original formulation of product quantization by allowing the dimensions of the descriptor to rotate. For that reason, it is possible to minimize the distortion error by finding the orthogonal matrix R that minimizes the equation

$$\min_R \sum_{i=1}^N \|Rx_i - Rr(c(Rx_i))\|^2 \quad (12)$$

It is possible to approximate the solution of the equation 12 in an iterative way by repeating the following steps:

1. Given a fixed R matrix (that may be computed during the previous iteration or, for the first one, it can be the identity matrix), the k cluster centroids are computed in the rotated space obtained by multiplying the R matrix to the input vectors. This is done by using the standard k-means algorithm independently in each of the m sub-spaces generated by the sets of sub-vectors extracted.
2. Given the fixed centroids computed in the previous step, it is now possible to find the matrix that minimizes the equation 12. The solution to that problem (that

is referring to one of the formulations of the Orthogonal Procrustes problem [20]) can be found by applying the Singular Value Decomposition (SVD) to the matrix obtained by multiplying the matrix obtained by stacking all raw vectors X and the one containing the reconstructed one Y by following $XY^T = USV^T$. In that way, the rotation matrix that has to be calculated can be expressed by $R = VU^T$.

These two steps are repeated until a maximum number of iteration (this value usually is set to 100) has been performed finding in that way a good approximation of the matrix that produces the minimum distortion error. A possible pseudo-code version of the reported algorithm can be found in algorithm 3.

Algorithm 3 A pseudo-code version of OPQ algorithm

```

 $X \leftarrow$  matrix of stacked descriptors
 $R \leftarrow I$ 
for  $it \in 1, \dots, 100$  do
     $X' \leftarrow RX$  ▷ Step 1
    for  $i \in 1, \dots, m$  do
        k-means on group of sub-vectors  $i$ 
    end for
     $Y \leftarrow r(c(X))$  ▷ Compute the reconstructed vectors
     $U, S, V^T \leftarrow SVD(XY^T)$  ▷ Step 2
     $R \leftarrow VU^T$ 
end for

```

By using the proposed approach, it is possible to obtain an optimal space partitioning that may help product quantization in finding the cluster's centroids that minimizes the distortion error introduced by the quantization. It is possible to integrate this OPQ algorithm each time product quantization is used.

The main disadvantage of adopting this approach is that it may require more time during the building phase of the index. In fact, especially if the database number of elements is large, the time required for finding this optimal rotation matrix can be high since the SVD algorithm has a complexity that is $\mathcal{O}(D^3)$ where D is the dimensionality of the input vectors and the k-means algorithm, that is run many times during the optimization procedure, has a complexity that grows linearly with the number of vectors to be clustered.

For that reason, the decision of applying an algorithm like this has to be taken with care correctly analysing the benefits that it would really provide to the final results. As for the other approaches reported, the OPQ algorithm is implemented in the Faiss library and will be widely used in the section 4.

4 Experiments and results

In this section will be presented the experiments conducted and presented the results obtained in order to find the most suitable indexing strategy for the retrieval part of the visual Geo-localization algorithm. So, starting from presenting the methodology used to produce the results presented, the analysis will follow by analysing the most relevant indexing techniques present in literature specifically focusing on the amount of memory they require and the time that is needed in order to retrieve the nearest neighbours. At the end of this section will be also presented a possible technique that can be used to reduce the database size before applying the indexing techniques. This method has in fact the main objective to discard some of the images inside the database that can be considered redundant since do not include any visual characteristic that is not already present in other photos.

Before presenting the different indexing techniques, it has been decided that the time required to retrieve the nearest neighbours for each query image should be below 200 ms in order to be considered acceptable. This threshold has been selected since it permits to produce the results with a tolerable latency and, as will be explained later, leaves the possibility to tolerate additional delays due to the fact that in a real-world application only a query is analysed at a time.

Regarding the memory usage, it has been decided that an indexing technique can be considered good if it is able to reduce the amount of memory used below 10% of the size of the raw database used. Selecting this value, permits the usage of very large descriptors databases that are similar to the one that will be used in a real-world application.

4.1 Testing methodology

In order to measure the performances of the various indexing techniques illustrated before, it has been built a testing environment capable to monitor the memory usage and the execution time for each algorithm. All of the code written to produce the results that will be presented in the rest of this work has been written using the Python programming language enriched by the usage of standard libraries. Among the libraries used, the most relevant are the Faiss library already presented before and the PyTorch library [21] that has been used to extract the descriptors from the input images.

In fact, before testing each indexing technique that will be presented it is fundamental to have a pre-computed database containing all of the descriptors extracted by the images inside the dataset. In order to do so, the code used has been taken from the work reported in [1] that, as already mentioned, is providing state of the art results in the visual Geo-localization task and, together with the code, provides also the weights necessary to the neural network used by the mentioned approach to correctly produce the descriptors required. This has been inserted as a starting step of the testing environment produced.

In order to permit also possible future extensions to the mentioned environment, the built framework is providing a common interface that provides a standard implementation of the function needed to measure the memory usage and each of the algorithm to be tested has only to implement, in its specific class, the functions needed to initialize, build

and configure the index. By means of this interface, the memory usage measures are centralized and each measurement is taken in the same way for each method reducing the risk of possible measurement errors. In order to assess how much memory is the current index using, the framework is simulating a real application of the index. In fact, like would be done in a real environment, the measures are divided in two phases:

1. The index under analysis is initialized with the configurable parameters specific for its type. After that the index is trained by providing it the full database of descriptors pre-computed. The resulting data structure it's later saved on disk exactly like would happen in a real application. At this point it is also measured the space required on disk by the serialized version of the data structures that composes the index. This measurement imposes a lower bound on the memory that will be used when the index is loaded since some additional data structures can be used by the Faiss library at run time in order to reduce the query time and manage a set of concurrent thread that is allocated to serve the queries performed.
2. In a different Python process, the descriptors of the queries, that will be used to measure the recall metrics, are loaded on main memory and the initial memory usage is measured by taking the Resident Set Size (RSS) of the current process. This measurement is used to eliminate from the final results the memory used by the process itself and the queries descriptors that do not influence the memory usage of the index. After that, the previously saved index it's loaded from disk and a first measurement of the RSS memory it's performed. Later, for each of the recall values to be tested, the nearest neighbour search is done and a new measurement of the memory usage is recorded. Those further measurements take into consideration the index in a simulated working condition like the one that would be experimented on a real application and so are very useful to correctly estimating the true value of the memory used. The final value for the Resident Set Size used by the index under analysis is computed by averaging out all the previously recorded values.

Also to measure the time needed for each query a centralized strategy has been used. In fact, the piece of code that is extracting the nearest neighbours of the provided query from the database is inserted in a timed section and, for that reason, at each extraction the mean time required by each query vector it's calculated.

It's important to notice that the queries are analysed in batches by the Faiss library. This has the consequence that the times that will be reported are the lower obtainable for the index under analysis since the library it's performing some optimization, that include also the usage of multiple threads to run in parallel the queries analysis. This is not easily applicable in a real usage of an application that is using the indexes analysed since it would be difficult or not possible to collect together the queries before searching the nearest neighbours in the database. The results provided can be still considered valid since they are taken in a situation that model an optimal condition of usage of the indexing techniques and, since each method is evaluated at the same way, they provide a clear indication of how much time the search operation would require making possible to compare the various techniques. In addition to that, the simulation of the time needed

to execute the nearest neighbours search in a real environment is out of the scope of this work since it depends too much on the machines on which the code would run and to the actual load of work produced by the requests of various users.

All of the experiments that will be presented have been run on a node of the Legion cluster of the HPC@POLITO. Each node of the cluster is equipped with 2 Intel Xeon Scalable Processors Gold 6130 and for the experiments it has been used 24 of the 32 cores of those processors. The nodes in the clusters are connected with Infiniband-EDR and they have a distributed shared memory with a total size of 22TB. During the tests performed, part of the shared main memory has been reserved for the processes to be run.

The database used to obtain the results that will be presented is the test set provided by the SF-XL [1] dataset that contains 2.8 million photos, if not differently specified. This set of the dataset has been preferred to the full version composed by 41.2 millions of photos during the experiments reported in order to reduce the memory usage during experiments and permitting in that way a more complete analysis of the indexing techniques. The results provided can be still considered a good approximation of the one that would have been obtained on the full SF-XL since the database size used is sufficiently high to appreciate the differences in both memory usage and time required per query of each indexing technique. The queries used for the various experiments are taken from the first set provided with the dataset that contains one thousand images.

The indexing techniques that will be used in the experiments that will follow are:

- Inverted File Index with Product Quantization (IVFPQ)
- Inverted File Index with Optimized Product Quantization (OPQ IVFPQ)
- Inverted File Index with Optimized Product Quantization and Re-Ranking (with compressed and raw residual error vectors)
- Inverted Multi Index with Optimized Product Quantization (IMI)
- Hierarchical Navigable Small Worlds with Optimized Product Quantization (OPQ HNSW)

How it's possible to notice from the list of techniques combinations analysed, only methods that can potentially reduce the memory footprint of the index compared to the one produced by the exhaustive search are investigated (with the only exception of inverted file index with optimized product quantization and re-ranking with raw residual error since this method is used only to compare the results with the one obtained with the same version of the algorithm but with compressed error). This choice has been taken since the indexes that do not reduce the memory usage are not applicable in a real-world scenario.

In order to evaluate how the various methods behave with different number of dimensions for the descriptor vectors, all the tests performed have been repeated with four different descriptors number of dimensions that have been set to 256, 512, 1024 and 2048 dimensions. This permits also to analyse which descriptor number of dimensions is the most suitable for the task of visual Geo-localization with the CosPlace network [1]. The

N. Dimensions	R@1	RSS memory usage (GB)	Mean time per query (ms)
256	73.6	5.58 ± 0.01	1.01 ± 0.08
512	76.7	10.95 ± 0.01	2.64 ± 1.66
1024	78.1	21.62 ± 0.01	5.96 ± 4.26
2048	76.5	43.02 ± 0.01	194.55 ± 19.23

Table 2: Results of the exhaustive search on different number of descriptors dimensions

different databases obtained with the four descriptors number of dimensions provides different values for the recall metric, as expected the lower the number of dimensions the lower is the recall value but, as reported in table 2, also a too high number of descriptors dimensions can cause sub-optimal value for the recall metric.

It is important to notice that in the graphs and tables that will follow only the value for the recall with one nearest neighbour retrieved are reported. This choice has been done since this metric is the one that provides the major variations for each indexing technique and, for that reason, it is more suitable to evaluate the results of the various indexing techniques since present more clearly the differences presented.

4.2 Trend of memory usage and time per query with descriptor database size

In this section, it has been analysed the trend of memory usage and the time required to retrieve the nearest neighbours of various combinations of techniques, described in section 3, using various descriptor’s database sizes in order to study the relationship that is present between them.

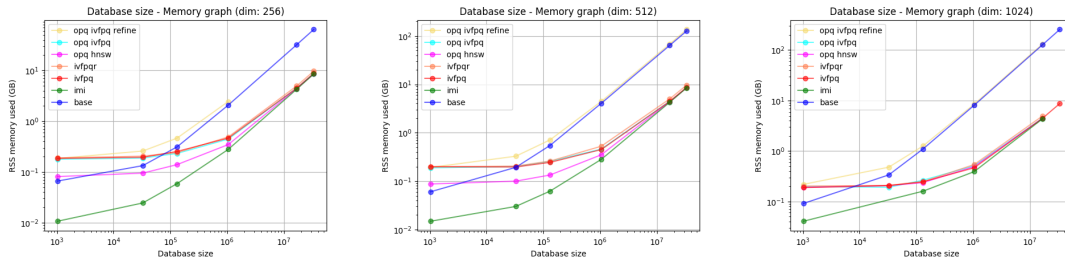


Figure 4: Trend of the RSS memory usage at the variation of the database size for different numbers of dimensions. The axes are in logarithmic scale for graphical reasons.

The graphs contained in figure 4 are reporting the results measured for the trend of RSS memory usage with the database size. As expected, the memory footprint of the various algorithm grows with the size of the database independently from the dimensionality of descriptors. It is possible to notice that all the curves reported, obtained by connecting with straight segments all of the measurements done, have a shape that can be approximated with the union of two half-lines with, respectively, a low and a high slope. It is in fact possible to notice from the figure 4 that the amount of RSS memory used by

the various indexing technique is increasing slowly as the database number of elements is growing until it reaches 10^5 descriptors and, after that point, the main memory required is rising with a higher constant rate. This difference in the slope of the lines is due to the Faiss library that is hiding the real memory usage of the index since it uses threads and additional data structures that requires much more memory than the index itself. For example, if the database is containing a set of 1024 descriptors each of 256 dimensions, the storage of the baseline index would require 1 MB of memory while it is possible to notice, from the leftmost graph in figure 4, that the measured memory is of more than 60 MB. The effect of this memory usage of the library is then becoming less relevant when the number of elements of the database grows.

In addition to that, in order to optimize the results, the Faiss library is using some additional data structures that influences the memory usage. This strategy is used when product quantization is applied by the library to reduce the time needed to compute the distances between the compressed vectors and the query one building a look-up table of distances of the quantized query descriptor from the centroids, as explained in section 3.4. Those tables are automatically built by the Faiss library if their memory usage is estimated to be below 2 GB and, for that reason, are used in any test performed in this thesis project, when product quantization is applied, since the configuration adopted always satisfy this condition.

For all of the reasons mentioned, the trend of the memory usage with the database size is linear with some exceptions due to problems in correctly estimate the memory usage with few elements in the database. This was expected since the index size depends directly on the number of elements present in the database. This finding is fundamental since it permits to approximate the behaviour of the various indexing techniques in an application that uses a different number of images from any of the one tested and it confirms that the usage of the test version of the SF-XL database is sufficient to report the behaviour of the indexes. The indexes reported in figure 4 have been produced by using different configurations so confirming the linear trend found and the different slope of the curves is due to the various indexing techniques that have different results in reducing the memory usage.

It is important to notice how the memory required by the baseline and optimized product quantization with inverted file index and re-ranking with raw residual error is becoming easily unacceptable when the database size increases. This is also more evident with the higher number of dimensions for each descriptor since, how it is shown by the graph in figure 4, those techniques may require more than 200 GB of RSS memory to store the index. Instead, when the indexing technique that is analysed is using product quantization to compress the vectors, the amount of memory needed remains limited and no more dependent on the number of dimensions of the raw vectors.

Regarding the time needed to extract the nearest neighbours per each query, the graphs contained in figure 5 reports the measurements done for three of the different descriptors number of dimensions. As in the case of the graphs in figure 4, these graphs have been taken by connecting each of the data points associated to the various measurements through straight segments. It is possible to notice that, in this case, the points representing the various times measured are aligned thus representing a direct proportionality

between the time needed for each query and the number of descriptors used. As in the case of memory usage, it has been noticed that with a limited number of elements inside the database, the trend identified is not fully explaining the data points reported. This is due to the fact that the Faiss library is using some optimization techniques to reduce the time needed to retrieve the nearest neighbours that seems to be less effective on small database. For that reason, it is possible to approximate the behaviour of the indexing techniques at any database size by following the trend reported.

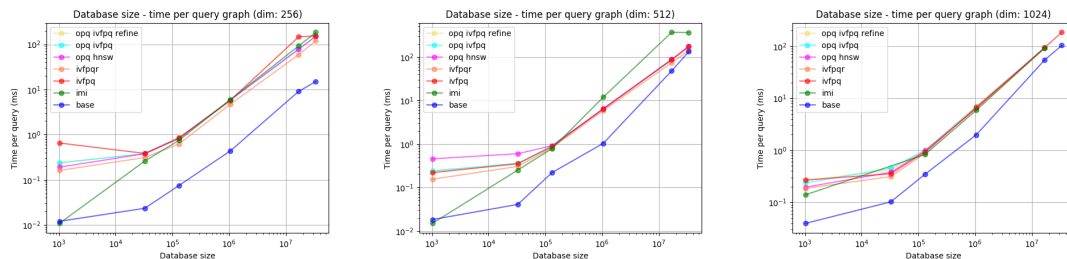


Figure 5: Trend of the time required per each query at the variation of the database size for different numbers of dimensions. The database size axis is in logarithmic scale.

Also in that case, the results that have been reported are in accord to the expectations since each query has to be compared, in the worst case, with all of the database descriptors and also in case of an approximate nearest neighbours search the subset of vectors to be analysed depends on the number of elements inside the database. This is in fact more evident with techniques like the exhaustive search where the number of distance computations to be done is the same as the database size.

From graphs in figure 5, it is also evident that the time required for each query remains acceptable, since below 200 ms, in most of the tests done. This makes the problem of reducing the time needed to retrieve the nearest neighbours less important in the analysis that will follow since this amount of time remains acceptable while the memory usage not. Despite this, it is fundamental to choose a technique that maintains this time low since, how it's possible to notice from figure 5, the indexing methods are, in general, requiring more time than the baseline.

4.3 Inverted File Index with Product Quantization

In this section, the results obtained with the Inverted File Index combined with product quantization (IVFPQ) are reported and discussed. It is possible to notice from the graphs contained in figure 6 that this method is able to reduce the amount of memory required by the index of about 10 times without any loss in the recall value with all descriptors number of dimensions. This is certainly due to the ability of product quantization to compress the input vectors and, at the same time, limiting the amount of information lost. It is possible to reduce even more the amount of memory used at the cost of some loss in the recall value. In fact, if a decrease of 2 percent of recall is considered acceptable, it is possible to reduce the memory footprint of the index by 30 times permitting the usage of a larger

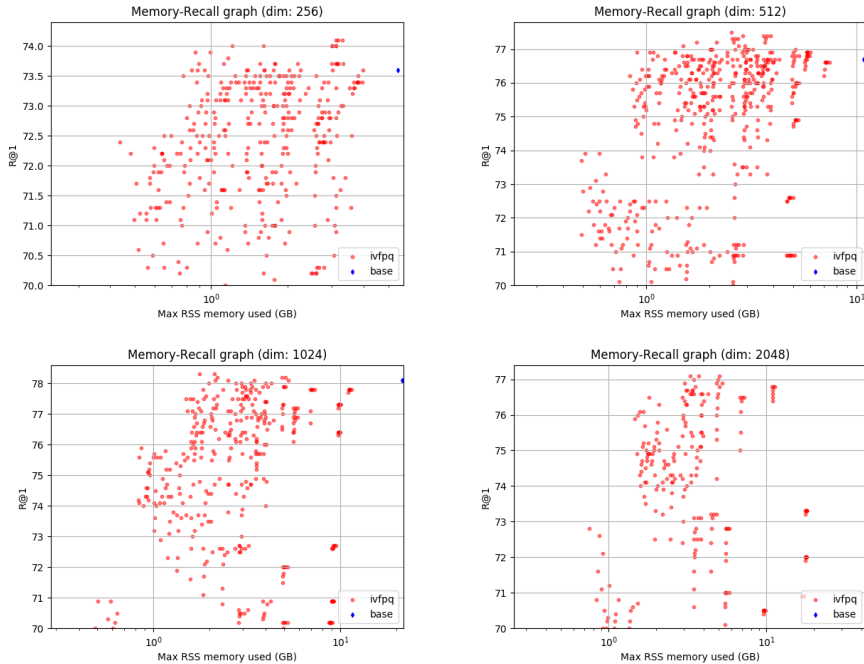


Figure 6: RSS memory usage and R@1 metric for the IVFPQ method with various configurations.

database of descriptors.

It is possible to notice from figure 6 that some configurations provide results that, in terms of recall computed, are higher than baseline. Those results are due to the fact that the compression algorithm has distorted the descriptor vectors and, by chance, this operation has produced vectors that are more similar to the queries provided. This characteristic of the mentioned configurations, although positive for the specific instance of the algorithm, has not a statistic relevance and may not be true with other queries vectors or other database data and, for that reason, can be ignored considering those results as they would have produced results similar to the baseline.

As it's possible to notice from table 3 that reports the detailed results for some of the best configurations tested, the memory usage is extremely reduced only for the higher dimensionalities while for the other descriptors databases the advantages of the usage of product quantization are lower. In fact, it's not possible to reduce the memory footprint of the lower dimensionality vectors and, at the same time, maintain the recall metric high and the mean time per query as low as the baseline since the loss of information produced by the product quantization is more significant in those cases. All of the results obtained with the others parameters configuration tested can be found in appendix A.1.

In table 3, are also reported the mean time per query measured on the 1000 queries. Those times are much higher than the one reported in table 2 for the baseline since many of the optimizations done by the Faiss library are not effective in the indexing technique

Dim.	IVF		PQ		RSS (GB)	Time (ms)	R@1
	clusters	probes	sub-vec.	bits			
256	16384	2048	128	4	0.73 ± 0.14	5.23 ± 0.01	72.2
	8192	1024	64	8	1.09 ± 0.14	1.29 ± 0.01	73.1
	2048	1024	128	8	1.17 ± 0.03	9.49 ± 0.80	73.8
512	32768	1024	128	8	0.98 ± 0.14	6.84 ± 0.73	75.4
	1024	512	128	8	1.02 ± 0.01	9.37 ± 0.92	76.2
	8192	1024	256	4	1.02 ± 0.16	11.43 ± 1.17	76.4
1024	4096	512	256	8	2.57 ± 0.09	6.34 ± 0.84	76.7
	131072	4096	512	4	2.73 ± 0.17	13.57 ± 1.54	77.0
	16384	16384	256	8	1.62 ± 0.18	192.09 ± 1.02	77.1
2048	512	512	512	4	1.61 ± 0.02	139.99 ± 1.57	76.1
	8192	2048	512	8	3.14 ± 0.02	48.56 ± 0.40	76.4
	32768	1024	512	8	3.44 ± 0.16	31.63 ± 0.55	76.6

Table 3: Results obtained with IVFPQ index for some of the best configurations for each descriptor dimensionality.

reported. For that reason, it is fundamental to finely tune the parameter that sets the number of cells of the Voronoi diagram visited during the search phase of the inverted file index in order to avoid spending too much time on the search operation. It is in fact possible to notice from the last row of section of the table 3 that reports the results for 1024 dimensions, that a full visit of the IVF index is not affordable taking the execution time very close to not acceptable values especially when the number of IVF cells is high. Those results show also that, in the application of product quantization, a lower number of sub-vectors with a higher number of bits per vectors has to be preferred to the usage of a higher number of sub-vectors each with a lower number of bits since the recall values obtained with the first approach are usually higher. This comes with the fact that the time needed for each query is usually lower since, like described in section 3.4, the cost of the operation of computing the similarity between the query and the database vector is proportional to the number of dimensions of the compressed vector. By considering instead the recall values reported, it seems that the number of dimensions that provide the best results is 512 since the recall values obtained are similar, although a bit lower, to the one with 1024 dimensions and the memory usage is between the lower obtainable.

4.4 Inverted File Index with Optimized Product Quantization

The effects of optimized product quantization on the inverted file index (this technique is called OPQ IVFPQ in the rest of this thesis project) are analysed in this section. How it's possible to notice from figure 7, the results obtained are, as expected, similar to the one obtained without the optimization of the product quantization with some minor differences. In fact, it's possible to notice that the recall values are slightly higher in the mean case with respect to the one measured without optimization. This is surely due to the lower distortion of the compressed descriptors obtained.

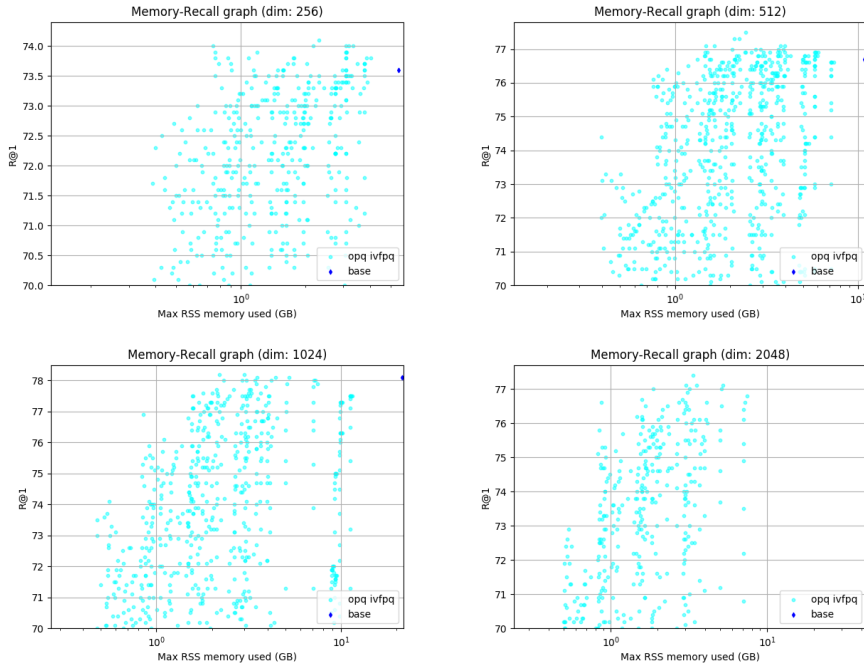


Figure 7: RSS memory usage and R@1 metric for the OPQ IVFPQ method with various configurations.

As already noticed for the inverted file index with product quantization, the memory reduction that can be obtained with this method is of about 90% of the baseline. In fact, since the product quantization, that is the element of the indexing technique that reduces the memory consumption, is common to those methods, it was not possible to further reduce the memory footprint.

As it's possible to notice from table 4, that reports some of the results obtained for the best configurations tested (all of the measurements done are reported in appendix A.2), the mean execution time of the nearest neighbours search is usually higher than the one obtained both in the baseline methods and the IVFPQ index. This is probably due to the matrix multiplication that has to be performed before analysing each query that has to be performed in order to apply the optimization of the product quantization. Those time increments, although tolerable, should be taken into consideration in a real application of this algorithm since they may limit the number of queries that can be served in a given amount of time. A possible approach that can be used in order to solve that issue, is to reduce the number of probes done by the inverted file index search algorithm. In fact, by following this suggestion, the time needed for each query is drastically reduces how clearly shows the fourth and seventh row of table 4 where the only variation of this parameter is reducing the execution time of 94%. Unfortunately, this reduction comes with a cost. In fact, the recall value decreases of about one percent point that, although acceptable, has to be taken into consideration in the trade-off between the time needed and the recall

Dim.	IVF		PQ		RSS (GB)	Time (ms)	R@1
	clusters	probes	sub-vec.	bits			
256	512	256	64	8	0.53 ± 0.01	3.96 ± 0.43	72.5
	512	256	128	8	0.90 ± 0.01	9.67 ± 0.52	73.0
	2048	1024	128	8	1.16 ± 0.03	10.32 ± 0.79	73.6
512	32768	512	128	8	1.10 ± 0.13	2.06 ± 0.15	75.7
	512	512	128	8	0.90 ± 0.01	18.05 ± 2.32	76.4
	32768	8192	128	8	1.10 ± 0.13	32.83 ± 2.60	76.7
1024	32	32	128	8	0.90 ± 0.02	16.68 ± 0.91	76.9
	16384	1024	256	8	1.95 ± 0.12	9.72 ± 0.58	77.0
	8192	1024	512	4	2.06 ± 0.14	22.28 ± 1.90	77.3
2048	64	64	128	8	0.87 ± 0.03	14.31 ± 1.84	75.8
	32	32	128	8	0.91 ± 0.02	19.24 ± 1.38	76.6
	32	32	512	4	1.59 ± 0.01	141.82 ± 1.61	76.7

Table 4: Results obtained with OPQ IVFPQ method for some of the best configurations for each descriptor dimensionality.

obtained.

In contrast with what was obtained in the IVFPQ case, the descriptors dimensionality that is providing the best results is 1024 since it is providing slightly better results with similar memory usage and also the time needed for the nearest neighbours extraction is generally lower than the one obtained with 512 dimensions. With the considerations done, it is now possible to compare the results presented with one obtained with inverted file index with product quantization. In fact, how it’s possible to notice from figure 8, that shows a comparison between the two indexes on the full set of configurations tested, the optimization of product quantization is only in some cases producing results that have a slightly higher recall metric value with lower memory usage with respect to the inverted file index without optimizations. Since the higher time needed to retrieve the nearest neighbours is acceptable, the optimized version of the product quantization has still to be preferred in the index selection since there exists the possibility to obtain better results.

4.5 Inverted File Index with Optimized Product Quantization and Re-Ranking

Re-ranking applied to an inverted file index is a technique that, as already mentioned in section 3.6, by considering also the residual error obtained by product quantization, is used to refine the results produced by an inverted file index. Two possible approaches are possible:

- Consider the raw residual error obtained by applying product quantization without any compression by re-ranking the set of candidates extracted with the inverted file index on the quantized vector using the exact calculation of the similarity between the query and the database descriptors since the full error is available. As it’s possible

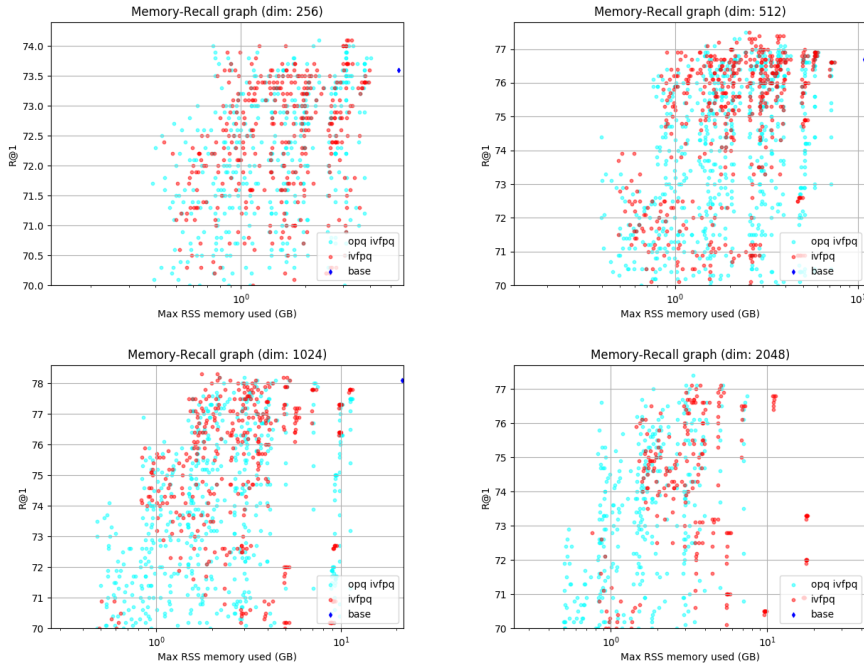


Figure 8: Comparison between the RSS memory usage and R@1 metric for the OPQ IVFPQ method and IVFPQ with various configurations.

to notice, this technique is not reducing the memory usage but it's providing the best results obtainable with the used inverted file index. In the rest of this work this technique will be identified with the word "refine" posed at the end of the indexing technique name.

- Compress and store the residual error on the database vectors in order to estimate the similarity metric. This approach is reducing the memory usage but may also introduce some loss in the recall value as will be discussed in this section. This technique will be identified in the rest of this work by appending as a suffix a "r" to the name of the method analysed.

In fact, how it's possible to notice from the graphs contained in figure 9, that reports the results obtained with various configuration for both the indexing techniques, the usage of the exact distance computation requires the storage of the full-sized descriptors and, consequentially, require the same or higher memory usage than the baseline. Instead, compressing the distortion error, has demonstrated to provide a significant reduction of the memory usage by maintaining the value of the recall metric very similar to the baseline. In fact, with the usage of compressed residual error vectors, is possible to use, for example, less than 1 GB of main memory to store the index and obtain the same result as the baseline in the case of 512 dimensions in descriptor vectors.

From the graphs of figure 9, it is also clear that, although the compression of the residual error vectors done, the recall values obtained are in many cases exactly the same

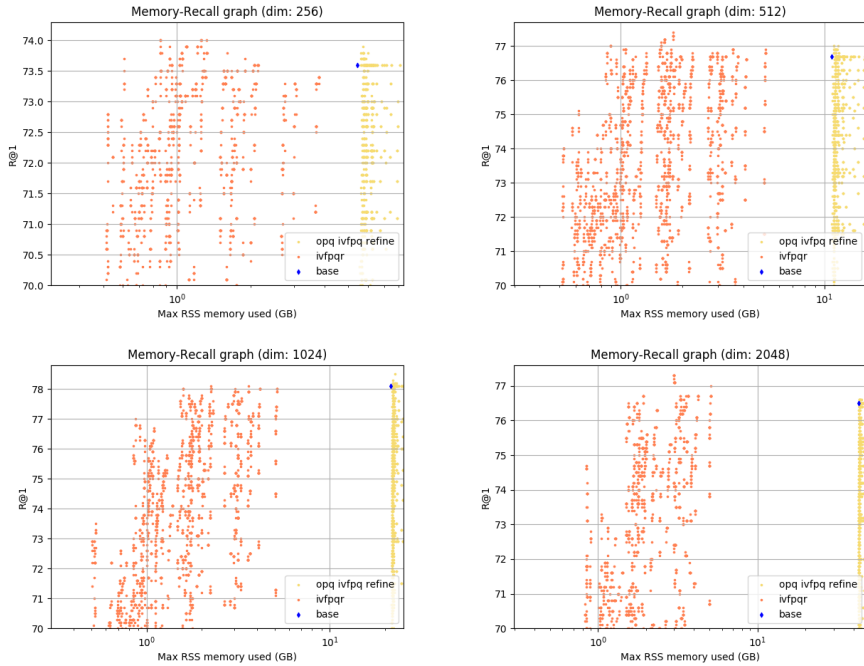


Figure 9: Comparison between the RSS memory usage and R@1 metric for the IVFPQR method and the OPQ IVFPQ REFINE method with various configurations.

as the one that can be measured with the usage of raw error. This last result confirms the fact that using a full-size error is not helpful in order to solve the problem of retrieval the nearest neighbours for visual Geo-localization since the same recall values showed can be obtained using much less resources. For that reason, this technique has not been used in any other part of this work.

Looking more in details to the results obtained with the quantized version of the residual vectors, it is possible to notice how those results are naturally organized in vertical lines. This fact is due to the parameter that is specifying the number of cells visited in the inverted file index Voronoi graph. This is strongly influencing only the recall metrics while the memory usage is not affected at all by this parameter since it is used only at search time without influencing the index construction. In fact, higher values of this parameter usually correspond to higher recall values, as already mentioned in section 4.3 and 4.4, but this is also correlated to a higher time for retrieving the nearest neighbours from the database.

Table 5, that reports some detailed results obtained for some of the best configurations tested for the two techniques under analysis for each descriptor dimensionality, clearly shows that the time needed to extract the nearest neighbours is, in general, higher than the one that can be obtained with the standard inverted file index. In fact, the re-ranking procedure used by those techniques requires a considerable amount of time to be applied since it has to compute the similarities for all the candidates two times.

Dim.	IVF		PQ		PQ re-rank		RSS (GB)	Time (ms)	R@1
	clusters	probes	sub-vec.	bits	M	bits			
IVFPQR									
256	128	128	64	8	16	8	0.53 ± 0.01	7.62 ± 0.59	72.8
	2048	1024	64	8	16	8	0.82 ± 0.01	7.55 ± 0.14	73.4
	256	256	64	8	16	8	0.61 ± 0.01	6.73 ± 0.13	73.7
512	128	128	128	8	16	8	0.86 ± 0.01	20.39 ± 0.73	76.8
	512	512	128	8	16	8	1.03 ± 0.01	70.26 ± 2.82	76.9
	64	64	128	8	16	8	0.89 ± 0.01	16.73 ± 1.66	77.0
1024	128	128	128	8	16	8	0.88 ± 0.01	19.19 ± 1.35	77.0
	128	128	256	8	16	8	1.55 ± 0.01	36.18 ± 0.16	77.2
	64	64	256	8	16	8	1.59 ± 0.01	46.68 ± 1.51	78.0
2048	64	64	512	4	16	8	1.52 ± 0.01	138.74 ± 1.33	75.9
	128	128	256	8	16	8	1.55 ± 0.01	36.62 ± 1.01	76.5
	512	512	512	4	16	8	1.76 ± 0.01	140.39 ± 1.97	76.5
OPQ IVFPQ REFINE									
256	64	64	64	4			5.65 ± 0.01	19.12 ± 1.28	73.6
	32	32	64	4			5.68 ± 0.01	19.01 ± 1.37	73.7
	512	512	32	8			5.79 ± 0.01	3.43 ± 0.29	73.7
512	128	128	64	4			10.98 ± 0.01	18.47 ± 0.90	76.4
	64	64	128	4			11.14 ± 0.01	36.47 ± 1.31	76.9
	128	128	128	4			11.14 ± 0.01	38.31 ± 1.77	77.0
1024	512	512	32	8			21.9 ± 0.1	3.42 ± 0.09	76.7
	32	32	64	4			21.9 ± 0.1	18.44 ± 0.83	77.4
	32	32	128	4			22.0 ± 0.1	36.58 ± 3.57	78.1
2048	32	32	128	4			43.31 ± 0.1	41.93 ± 0.54	76.1
	512	512	128	4			43.32 ± 0.1	36.11 ± 0.22	76.4
	512	512	64	8			43.37 ± 0.1	6.63 ± 0.04	76.5

Table 5: Results obtained with IVFPQR and OPQ IVFPQ REFINE method for some of the best configurations for each descriptor dimensionality.

It is important to notice how, in contrast with the results obtained with the methodologies analysed before, in order to obtain the best results an exhaustive visit of the Voronoi cells of the inverted file index has to be performed. When the search is not exhaustive the recall value measured is in many cases very low as reported in the table that can be found in appendix A.3 that reports all of the measurements done. This, together with the re-ranking procedure, is surely one of the causes of the high times reported in table 5. The times measured are still considerable acceptable since below 200 ms but the adoption of such an index should be considered with care since those times may be even slower if the queries are not analysed in batches, like described in the introduction of section 4.1.

Despite this, the recall values obtained are usually very close to the baseline also when product quantization drastically compress the descriptors vectors. This is probably due

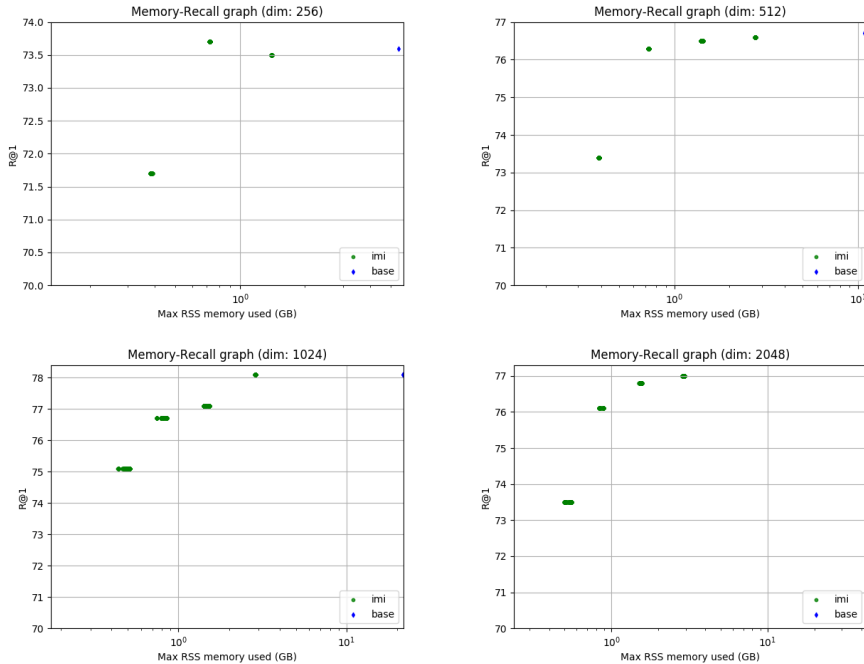


Figure 10: This figure reports the RSS memory usage and R@1 metric obtained for the IMI method combined with optimized product quantization with various parameters configurations.

to the usage of the residual vectors, also if in the compressed form, that permits a better estimate of the first nearest neighbour that mostly influences the R@1 metric reported.

Considering the memory usage of the IVFPQR index, it’s possible to notice that, if compared to the standard IVFPQ index, better recall is obtainable with the same amount of main memory. This actually demonstrate that using the residual error in a quantized form is adding a negligible amount of memory (in the test reported only 16 bytes per vector) and it is providing better results in term of recall metric. For that reason, this indexing technique has to be preferred to the IVFPQ or OPQ IVFPQ techniques especially if the amount of time required for the nearest neighbour retrieval is not constituting a major problem in the particular application that is applying visual Geo-localization.

4.6 Inverted Multi-Index with Optimized Product Quantization

In this section are reported the results obtained with the inverted multi-index with optimized product quantization (IMI). As it’s possible to notice from figure 10 that reports the results obtained, this technique is able to reduce the memory footprint by maintaining the recall values similar to the one obtained with the baseline method. It also evident that the points in the various graphs are less spread than the one obtained with previously analysed techniques like IVFPQ. This behaviour of the index technique itself that is not as sensible to the input parameters like the other approaches tested. Another possible

Dim.	IVF		PQ sub-vec.	RSS (GB)	Time (ms)	R@1
	clusters	probes				
256	32	1	64	0.38 ± 0.01	6.35 ± 0.03	71.7
	512	4	128	0.72 ± 0.01	15.46 ± 0.02	73.7
	131072	1	128	0.72 ± 0.01	15.72 ± 0.09	73.7
512	4096	1	128	0.73 ± 0.01	15.41 ± 0.15	76.3
	32	1	256	1.40 ± 0.01	34.84 ± 0.25	76.5
	512	1	512	2.75 ± 0.01	71.91 ± 0.18	76.6
1024	4096	4	128	0.74 ± 0.01	46.38 ± 1.05	76.7
	8192	4	256	1.42 ± 0.01	71.47 ± 0.95	77.1
	256	4	512	2.84 ± 0.01	72.63 ± 0.10	78.1
2048	64	1	128	0.84 ± 0.01	15.51 ± 0.08	76.1
	32	1	256	1.54 ± 0.01	35.07 ± 0.72	76.8
	512	2	512	2.86 ± 0.01	71.89 ± 0.26	77.0

Table 6: The results that have been measured for some of the best configuration tested using the Inverted Multi Index with the different number of descriptors dimensions. For the product quantization are always used 8 bits per sub-vector.

cause of this particular distribution of points in the graph is that there are only two parameters that can be specified in the construction of this index. Actually only the number of sub-vectors of product quantization and the number of centroids used by the clustering algorithm used by this technique, as described in section 3.5, are the input parameters that can be specified and the effects of the number of centroids on the final memory usage of the index are negligible since it may happen that, when this number is increased, some of the partitions built by the index can remain empty so they are not using any space.

It is possible to notice from the graphs contained in figure 10 that the performances that can be obtained by using this indexing technique are in general better with high dimensional descriptors since the memory reduction that can be obtained by maintaining the recall similar to the baseline is, in proportion, higher with respect to the one that is measured with a lower number of dimensions. That is probably due to the particular construction of this indexing technique that, in contrast with what is happening with standard inverted file indexes, is splitting in a half the descriptors before applying k-means clustering. This result, in general, in better performances of the clustering algorithm since, due to the lower number of dimensions of the vectors, the k-means algorithm is applied in optimal conditions.

How it is possible to notice from table 6, that shows the results measured for some of the best configurations, the results are strongly influenced by the product quantization number of sub-vectors. It is evident that the memory usage tends to be, approximately, two times larger if the number of PQ sub-vectors is doubled. This is due to the fact that the inverted multi index does require only a negligible amount of memory for the storage of the lists inside the particular kind of inverted file index it uses. This is confirmed by the second and third row of table 6 where the number of sub-vectors for the product

quantization is the same and there are three order of magnitude of difference between the number of partitions of the inverted file index with the same memory usage. Those two rows report also an example of the reason why this indexing technique is producing overlapping points in the graphs of figure 10. Actually, this indexing technique is confirmed to be not sensible to the number of partitions in the grid it produces internally. Appendix A.4 is reporting the full set of measurements done with this indexing technique.

The amount of time required by this technique to retrieve the nearest neighbours is usually higher, although still acceptable, than the one that can be obtained with other techniques like inverted file index with product quantization and re-ranking. This is due to search strategy of this index, described in section 3.5, that requires two computationally expensive parts: the independent retrieval of the nearest neighbours of the two sub-vectors of the query and the final comparison with the candidates extracted. This time grows linearly with the number of sub-vectors of product quantization since, as explained before, the distance computation has a cost that is proportional to the number of dimensions of the vectors that are used in the computation.

This indexing technique, although the good results that have been measured for the recall metric, it is not preferable for the visual Geo-localization task since the memory reduction that is obtainable is lower than the one that have been reported for other indexes like IVFPR or OPQ IFPQ when approximately the same recall value is measured. In addition to that, also the time required for each query is not as low as the one that can be obtained with the up-mentioned techniques. For all of those reason, the adoption of this kind of index is not suggested for an application like the one under analysis.

4.7 Hierarchical Navigable Small Worlds with Optimized Product Quantization

In this section, the results obtained with the hierarchical navigable small world method combined with optimized product quantization (OPQ HNSW) are analysed. It has been chosen to use the optimized version of product quantization seen the results obtained with the inverted file index in section 4.4. Since the number of nodes in the graphs used by this technique is large, it has been used an inverted file index in order to partition the space in cells before applying the HNSW index and to reduce the number of links that would be used by the HNSW technique. In fact, in that way the number of nodes in the HNSW index is limited to the number of cells in the Voronoi diagram built by the IVF index and it avoided the storage of the links for all descriptors vectors. This is a standard approach adopted with the hierarchical navigable small world when applied to huge datasets since the number of graphs links is proportional to m times the number of descriptors in the database in the worst case, that corresponds to the connection of each node to the m nearest neighbours. In fact, the storage of those links would require $2s_{link}mN$ bytes where s_{link} is the size of each link in the architecture used. This means that, for example, in case of a database of one billion elements with 32 connections for each node (this is a standard value for the number of links for each database element) would require more than 238GB of memory only to store the connections between the various nodes which is not acceptable. With the application of the inverted file index, instead,

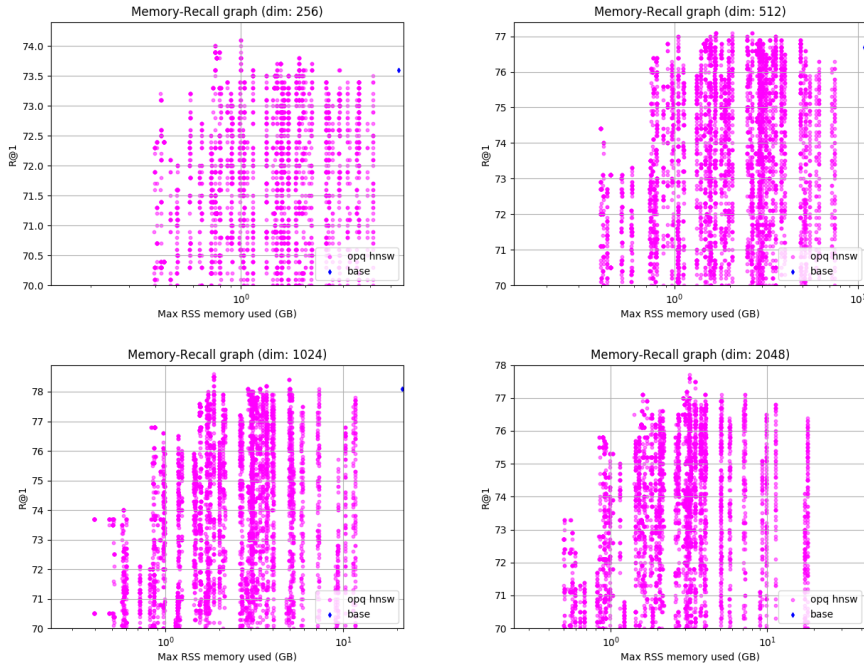


Figure 11: The RSS memory usage and R@1 metric measured for the Hierarchical Navigable Small World (HNSW) method combined with Optimized Product Quantization (OPQ) with various parameters configurations.

the number of links it's no more proportional to the number of elements of the database but to the partitions of the IVF index, which can be set as parameter, and it is usually much smaller than the number of descriptors.

Figure 11 is reporting graphs that shows the results obtained with this technique by using different configurations of the parameters required. It is possible to notice from figure 11 that this indexing technique is able to reduce the memory footprint with respect to the baseline by maintaining a similar recall value but, if compared with other methods, the memory reduction seems, at a first analysis, not as effective as the one that can be obtained with other indexing techniques like IMI or IVFPQ. In fact, it is clear from the graphs reported that most of the results that provide recall similar to the baseline are located in the area that indicates more than 1GB of memory required while, for example, in the IMI case those results are located well below this zone. This behaviour of this indexing technique may be due to the intrinsic characteristics of the navigable small world graph that has to store possibly a large number of links in addition to the memory required by the inverted file index inserted.

Table 7 reports more in detail some of the best results measured by using the hierarchical navigable small worlds combined with the optimized product quantization (all of the results obtained are reported in appendix A.5). How it's possible to notice from the results reported in this table, this approach is providing worse recall metric values

Dim.	IVF		PQ sub-vec.	m	RSS (GB)	Time (ms)	R@1
	clusters	probes					
256	512	256	64	32	0.42 ± 0.01	3.98 ± 0.42	72.1
	2048	1024	64	32	0.58 ± 0.01	4.17 ± 0.35	72.9
	512	256	128	32	0.79 ± 0.01	8.92 ± 0.84	73.3
512	512	256	128	32	0.80 ± 0.01	8.60 ± 0.39	75.4
	1024	512	256	32	1.67 ± 0.01	19.24 ± 0.66	76.2
	131072	4096	256	16	2.06 ± 0.01	34.12 ± 0.80	76.6
1024	16384	2048	256	16	1.71 ± 0.02	26.27 ± 0.88	77.1
	16384	2048	256	32	1.74 ± 0.01	19.06 ± 1.15	77.8
	32768	1024	256	16	1.86 ± 0.01	13.53 ± 1.14	78.0
2048	32768	1024	128	32	1.44 ± 0.01	17.52 ± 0.68	75.4
	1024	512	256	32	1.84 ± 0.01	26.81 ± 1.06	75.9
	2048	1024	256	32	2.06 ± 0.01	26.54 ± 0.56	76.3

Table 7: Results obtained with OPQ HNSW method for some of the best configurations for each descriptor dimensionality. The number of bits for product quantization is always 8 and m represent the number of links for each node in the HNSW graph.

and a higher memory usage with respect to the other indexing techniques. This is surely done to the combined application of the IVF and HNSW techniques which requires the storage of two different indexes as already mentioned. Also the time needed to retrieve the nearest neighbours for each query are higher than the one obtained with other techniques probably for the double index search to be performed. This methodology is usually applied to much larger dataset than the one used for those tests since the advantages of applying it are evident only with partitions with a number of elements larger than the one obtained since, in such a situation, the HNSW index can really help to find the IVF cells to be analysed avoiding not useful vectors comparisons while, in the testing environment used, it is possible to analyse all partitions and, in that way, obtain better recall values. An example of this can be found on the results reported for 2048 dimensions in table 7 where, although only 1024 cells of the 32768 of the IVF index are analysed, the recall value obtained is only one percent lower than the baseline.

A possible solution in order to reduce the amount of memory needed for this approach is to reduce the number of sub-vectors or the number of clusters in product quantization but this approach has not been tested since the recall values are already low and the only possible effect of using a stronger compression is to reduce the recall metric values. For all of those reasons, the usage of this particular kind of index should not be preferred to other indexing techniques like IVFPQR especially if the descriptor database does not contain a huge number of elements and the recall value is considered as the most critical metric in the particular application chosen.

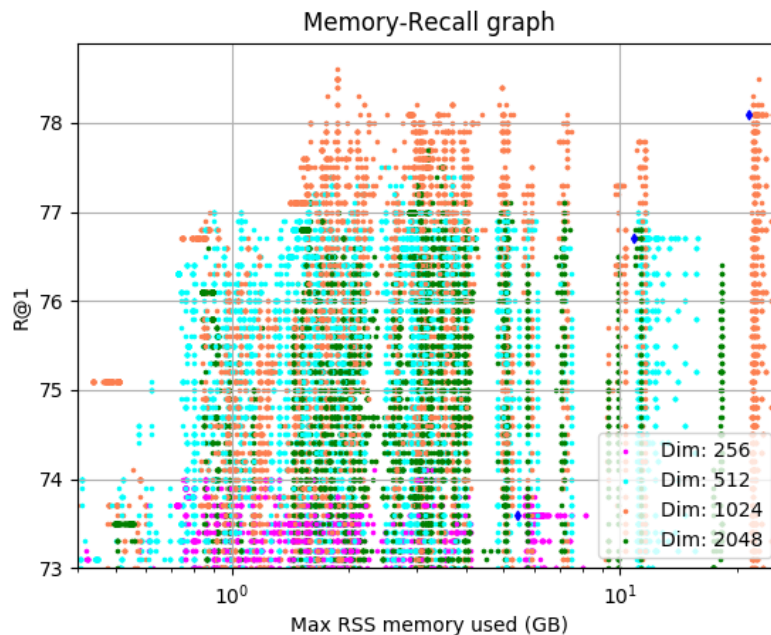


Figure 12: All of the results obtained with the different number of descriptors number of dimensions regardless of the indexing technique used.

4.8 Results comparison

After having analysed in details all of the results obtained with the various techniques adopted, it is possible to compare the measurements presented and, in that way, have a wider view on the chosen methodologies. In fact, in order to study a solution that would be really usable in an application that is providing visual Geo-localization services, it has to be discussed which number of dimensions the descriptors should have and which indexing technique should be adopted in order to obtain a good trade-off between the memory used and the recall metric values obtained. The graph contained in figure 12 is reporting together all of the measurements performed for the various indexing techniques and number of descriptors dimensions presented during this thesis project analysing the recall values measured with the various number of dimensions for each descriptor.

By analysing this graph, it is possible to notice that the recall values that have been obtained with 1024 dimensions are in general higher than the ones obtained with other dimensionalities. This, although due to the result obtained with the exhaustive search which is higher than the ones obtained with the other descriptors number of dimensions, indicates that 1024 represents the most suitable number of dimensions for the visual Geo-localization task if the CosPlace network [1] is used to extract descriptors and an indexing technique that is able to reduce the memory usage of the database is applied.

It is important to notice that, if 512 dimensions for each descriptor is used, results similar to the one obtained with 1024 dimensions when the memory footprint of the index

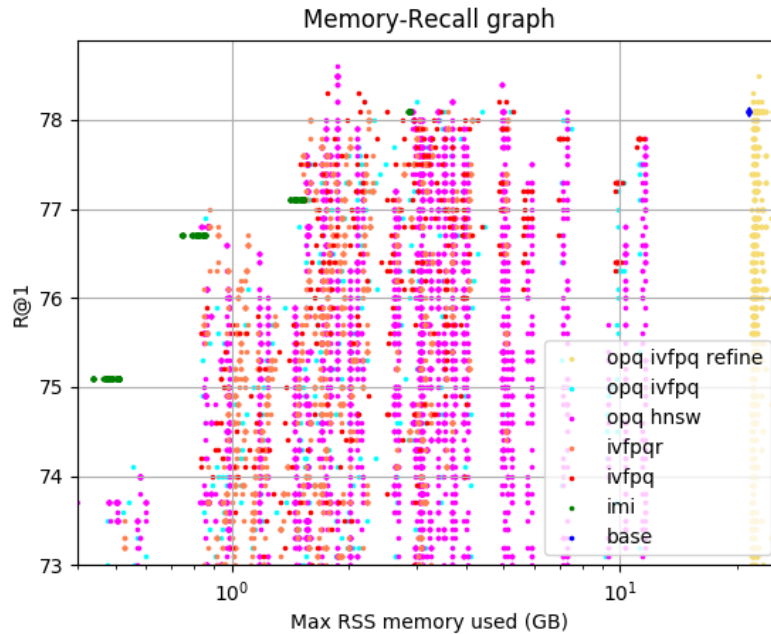


Figure 13: This graph reports all of the results measured for each of the indexing techniques analysed with 1024 dimensions.

goes below $1GB$. Despite this last finding, it is still possible to notice that the recall obtained with 512 dimensions is usually slightly lower than the one measured with 1024 and, for that reason, it is confirmed that the most suitable number of dimensions for the descriptor vectors is 1024. For that reason, this dimensionality is the one that is suggested in a real environment since it can guarantee high recall values and, at the same time, a limited memory usage.

In order to compare the indexing techniques that have been analysed before, it has been decided to compare them by using descriptors with 1024 dimensions that, as discussed before, is guaranteeing the best results obtainable.

How it's possible to notice from the graph contained in figure 13 that shows the results for each of the indexing technique adopted, it is not evident which method is actually providing the best results. In fact, by looking at the zone of the graph below 1 GB of RSS memory used and between a recall value of 76 and 77, that is depicting the best results obtained for this descriptor number of dimensions, there are different indexing techniques that are providing similar results in term of memory reduction and recall metrics. More specifically, those indexing techniques are:

- Inverted file index with product quantization and re-ranking based on the compressed residual error
- Hierarchical navigable small worlds with optimized product quantization

- Inverted file index with optimized product quantization
- Inverted multi-index with optimized product quantization

Those techniques, although using different approaches, are effectively producing recall values that are similar (they have a maximum difference of 0.3 percent that corresponds to a wrong location assignment for just three images of the one thousand queries) with approximately the same memory requirements. Among those techniques the one that is producing slightly better results in term of recall value is the inverted file index with product quantization and re-ranking with compressed residual error.

Technique	IVF		PQ		m	R	R@1	RSS (GB)	Time (ms)
	clusters	probes	sub-vec.	bits					
IMI	8192	1	128	8	-	-	76.7	0.78 ± 0.01	15.86 ± 0.03
IMI	64	1	128	8	-	-	76.7	0.81 ± 0.01	15.53 ± 0.05
IMI	16384	1	128	8	-	-	76.7	0.81 ± 0.01	17.88 ± 0.75
OPQ HNSW	32	32	128	8	16	-	76.8	0.83 ± 0.02	16.60 ± 0.92
OPQ HNSW	32	32	128	8	32	-	76.8	0.84 ± 0.01	20.39 ± 0.80
OPQ IVFPQ	32	32	128	8	-	-	76.9	0.85 ± 0.01	16.68 ± 0.91
IVFPQR	128	128	128	8	-	16	77.0	0.87 ± 0.01	19.19 ± 1.35

Table 8: The best results obtained for different configurations of the various indexing techniques. The descriptors used have 1024 dimensions. In the headers of the table the R is indicating the number of sub-vectors used for the re-ranking with compressed error (when this is used the number of bits for product quantization is always 8). The parameter m is representing the number of links per each node in HSNW technique.

Table 8 is reporting the results obtained for the most relevant configuration for the various indexing techniques (the ones that produces results in the zone described before of the graph contained in figure 13) in order to better explain this last finding. In fact, this table clearly shows that the IVFPQR technique is the one that is producing the higher recall values among the other indexing methods. Despite this, inverted file index with product quantization and re-ranking is also the technique that is requiring the larger amount of memory and the larger amount of time to retrieve the nearest neighbours for each query. For that reason, this technique should be chosen only if the recall value is considered the most critical metric in the target application and other indexes should be preferred if the memory or time per query are also important.

In fact, the most memory efficient technique among the one tested is inverted multi index with optimized product quantization which is producing slightly lower recall values with 10% less memory used than the IVFPQR index. Also the time needed to extract the nearest neighbours is slightly smaller with the IMI technique.

Hierarchical navigable small worlds with optimized product quantization and inverted file index with optimized product quantization are producing results that are in between the IVFPQR and IMI techniques. In fact, table 8 is showing that the OPQ HNSW indexing method is producing recall values that are in general slightly higher than the

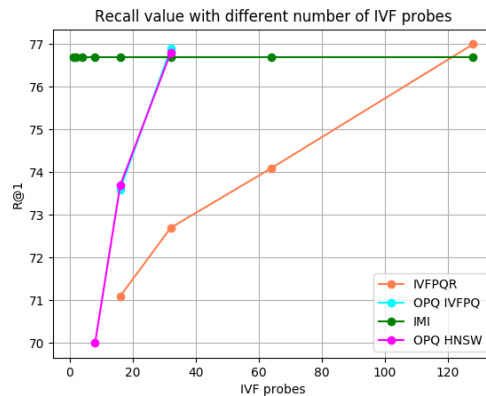


Figure 14: Trend of the recall values with different number of IVF cells probed with the various indexing techniques.

IMI technique but with a memory usage that is smaller than the one measured with the IVFPQR method. OPQ IVFPQ indexing technique is instead producing even higher recall values than the one obtained with the OPQ HNSW index at the cost of a slightly higher amount of memory required. For those reasons, the OPQ HNSW or the OPQ IVFPQ techniques have to be preferred to IMI and IVFPQR without any notable difference among the two if the memory and the time required for each query are considered more critical than the recall value measured.

It is still important to notice that all of the techniques reported in table 8 are analysing all of the cells in the Voronoi diagram produced by the inverted file index that is used with the only exception of the IMI technique. In that case, this approach is possible since product quantization is helping the application reducing the time required to compute the distances between the query and the database vectors limiting in that way the resources needed and making this search strategy feasible. If the number of cells visited is reduced, the recall value is usually rapidly decreasing, as it's possible to notice from figure 14. This makes not affordable to use a partial search when the descriptors are highly compressed, since the recall value becomes easily not acceptable and, in addition, the time required to retrieve the nearest neighbours is higher.

Regarding the product quantization technique, it is possible to notice from table 8 that all of the best configuration tested have a number of sub-vectors equal to 128 with 8 bits per dimension (resulting in 256 clusters per each of the 128 dimensions of the compressed vectors). This configuration permits to obtain a good trade-off between the compression factor and the amount of information kept inside the descriptors needed to extract the correct database vectors for a given query. It is in fact possible to notice from figure 15, that reports the measurements done with different number of sub-vectors in product quantization, how the advantages of using more sub-vectors (maintain more of the original vector information and producing a higher recall value) are always smaller as the number of them increases. In fact, the recall values obtained with 512 sub-vectors for each descriptor

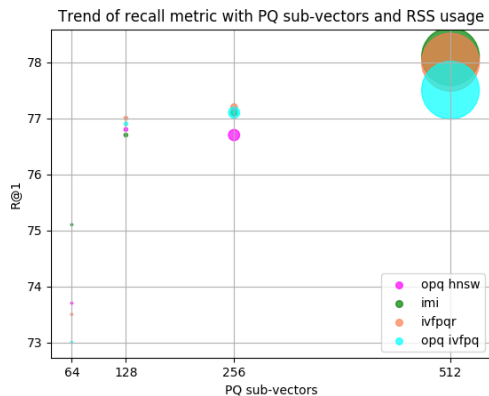


Figure 15: Trend of the recall values with different number of sub-vectors in the product quantization used on 1024 dimensional descriptors. The configuration used are the one reported in table 8 with only the PQ sub-vecs. parameter changed. The number of bits for each dimensions of the sub-vectors is always 8. The radius of the points is proportional to the approximated RSS memory usage of each configuration.

is just about one percent higher than the one obtained with only 128 partitions and with a memory usage that is approximately 25% of the one measured with 512 independently from the indexing technique used.

For that reason, in an application that is providing visual Geo-localization services it is strongly suggested the usage of 128 sub-vectors for each descriptor since this value is providing a recall value that is close to the baseline with a lower memory usage than using more sub-vectors despite, with the latter approach, slightly better performances can be obtained. It is possible to notice from figure 15 that the performances of the various indexing techniques are reduced drastically when the number of sub-vectors for product quantization is set to 64 since the compression is, in that case, introducing a significant error when the vectors are reconstructed to recover the original descriptors making more difficult for the system to extract the correct nearest neighbours from the database for each query.

Summing up all of the results discussed in this section, in order to obtain the best performances in the retrieval part of the visual Geo-localization algorithm, it is suggested to use the optimized version of the product quantization with a number of sub-vectors equal to 128 and coupled to a hierarchical navigable small world index or to an inverted file index. Before applying the indexing techniques mentioned, the descriptors contained in the database should have 1024 dimensions, as discussed before, given that the CosPlace network [1] is used to compute them.

5 Descriptor database reduction

In this section will be analysed a possible approach to be used in order to reduce the size of the descriptor database used for retrieving the nearest neighbours. This technique is born from the observation that many of the images contained in the SF-XL dataset are similar one to other because taken at positions that are very close one to each other and, since they are not adding relevant information because their content is practically the same, can be eliminated from the database and still obtaining similar recall values.

Before diving into this technique, it is important to mention how the SF-XL database is built from Google Street View images. In fact, this procedure is fundamental in order to fully understand the reasons behind the choices that will be presented while reporting the approach followed. The images taken from Google Street View that are used to build the dataset are 360°spherical photos usually taken from a car moving on the street at fixed distances one from the other.

Those images can't be directly used for visual Geo-localization since their size is too big for the usage inside a neural network since the convolutional operations on this kind of photos would require an amount of resources that is not acceptable for a standard GPU and, in addition, the time needed to process a single image would be too high (since the number of operations to produce the descriptors from the photo is proportional to the size of the image). For that reason, the 360°photos (usually called panoramas) must be cropped in different sub-images that are of an acceptable dimension for the next steps in the algorithm.

In order to do so, the angle from the camera and the ground is fixed, like it's possible to notice in the figure 16 that reports an example of the panoramas used, and the 360°photo is cropped in twelve non-overlapping regions of size 512×512 . The resulting images are then ready to use for the next steps of the algorithm. It is possible to reduce the number of images in the database by removing some panoramas from the set used to build the crops. This selection should keep only relevant panoramas, the one representing different places for example, while discarding the images that are too similar. In fact, those images can be considered as duplicates and, by not inserting them into the database, it is possible to reduce the memory usage of the pipeline and to speed up the retrieval part of the algorithm (since less images have to be analysed). It has been chosen to work on panoramas instead of using directly the cropped images since an analysis performed on both sets would have produced practically the same results but, since the number of panoramas to be analysed is lower than the number of images, the first approach is more efficient.

In order to estimate how much similar are two different panoramas three possible



Figure 16: An example of panorama used in SF-XL dataset

approaches have been investigated:

1. It has been considered only the location where the panoramas have been taken without considering the visual features of the images. This strategy is applied by assigning the panoramas to a specific squared zone of space by rounding up the UTM coordinates of the ground truth position of the image. After the zone assignment has been performed, only one of the panoramas belonging to the zone under analysis is selected to be kept in the final database. In order to choose which image has to be taken, the list of panoramas for each zone is sorted by the time when the picture has been taken and only the most recent one is taken. This choice has been done since, during time, the area of where the photos are taken may change (for example if new buildings are built) and the most recent images may be able to correctly insert those changes into the database while not updated photos may introduce noise. This simple approach has been adopted for its efficiency and, since it does not require any similarity computation, can be considered as a reference method on which the other approaches can be compared to. In the following graph and tables this approach is identified as "zone" since, as described while explaining the approach adopted in the selection of panoramas, one of the fundamental parts of this technique consists in dividing the space in zones.
2. The cosine similarity between the descriptors vectors taken from the panoramas has been calculated. In fact, in order to compare the panoramas, it is more important to consider how the network that actually computes the descriptors sees the features present in the image rather than applying a similarity function directly on the raw images since the network may extract different features, and compute different descriptors, for images that are similar to a human eye. Since, as mentioned before, computing the descriptors directly on the whole panorama is requiring many resources, the descriptors of the crops belonging to the same panorama have been fused together in order to produce the descriptor of the panorama. In order to do so, the descriptors of the various crops have been firstly sorted by the crop position in the panorama, to avoid the problem of different possible ordering of the crop descriptors, and then joined together in a $12d$ dimensions vector. After this preliminary step, the panoramas are then divided in zones like described at point 1 and, for each vector that lays in one of the zones computed, the cosine similarity between a selected vector and all the others in the same zone is computed by:

$$s_{i,z} = \sum_{j \in 1, \dots, N | j \neq i} \frac{x_{i,z} \cdot x_{j,z}}{\|x_{i,z}\| \|x_{j,z}\|} \quad (13)$$

where $x_{k,z}$ represent the k -th panorama of the N contained in the same zone z . This equation is computing the sum of similarities between a selected vector and all the others in the same zone in order to analyse how similar is the vector to the others. The panoramas to be kept are the n that have the lower values of the computed similarity. In the following figures and graphs of this work, this technique will be called "cosine" for the similarity function adopted.

-
3. The correlation between each descriptor vector extracted from the panoramas has been used to determine which image should be kept. In order to do so, as already described at point 2, firstly the panoramas descriptors are reconstructed from the one computed with the crops. Then the descriptors are divided in zones based on the geographical coordinates of where the corresponding panoramas have been taken and, after this, only the descriptors with the n lowest sum of correlations with the other vectors belonging to the same zone are kept. To compute the correlation between two indexes, it has been used the Pearson correlation coefficient that can be defined as follows:

$$\rho = \frac{Cov(X, Y)}{\sigma_X \sigma_Y} \quad (14)$$

where x, Y represent the two vectors to be compared and σ_z is the standard deviation of the z vector. This technique has been used since if the descriptor vectors that are kept are not correlated one with the other they are more probably representing different features. In the next figures and graphs, this technique will be called "correlation" for the similarity function adopted.

To measure the performances of those techniques, a subset of the panoramas contained in the train set of the SF-XL dataset have been adopted. This subset contains about 750 thousand panoramas (that correspond to approximately 9 million crops) taken in different zones of the San Francisco city. In order to measure the recall value that can be obtained with the usage of less images on the database, the original query set provided by the SF-XL dataset have been adapted by removing queries belonging to the areas of the city that do not contain any panorama among the selected ones. After this operation, the resulting query set contains about 350 queries of the one thousand in the original set and, for that reason, the recall values reported in the next figures and tables are in general much higher than the ones showed in the previous parts of this work. From those images, before proceeding with the analysis of the performances obtained, the descriptors have been extracted from the crops obtained for each panorama.

As described in the section 4.1, the memory is measured by the testing framework built in terms of resident set memory usage and the time required to extract the nearest neighbours is computed on the whole query set analysed in a single batch. The experiments reported in this section are performed only with descriptors consisting of 512 dimensions since they provide a good compromise between the memory usage and the recall values that can be obtained as shown before. Only the exhaustive search has been used to obtain results that are not influenced by the specific indexing technique used. It is possible to calculate also the theoretical size that the index will use in the server main memory by following:

$$m = n_{dim} s_{float} N_{DB} \quad (15)$$

where m represents the number of bytes required, n_{dim} is the number of dimensions of the descriptors, s_{float} represent the size in bytes that a floating point number is requiring in main memory for the architecture used (assuming that the descriptors are stored as arrays of floating point numbers) and N_{DB} is the number of elements in the database of images. It is important to notice that the equation 15 is not considering the memory required by

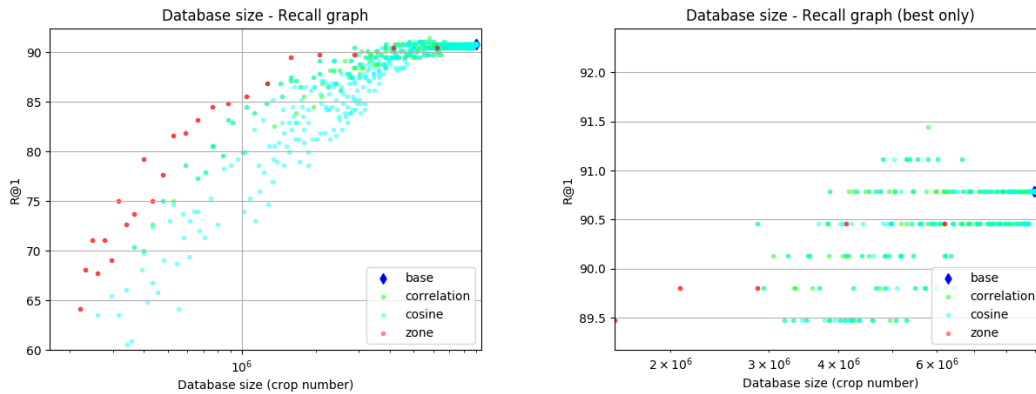


Figure 17: On the left the recall measured with the techniques presented for various database sizes, on the right an extract of the graph on the left which focuses on the higher recall values measured.

the Faiss library for the additional structures used to optimize the search process and, for that reason, the value that can be calculated with it is in most cases different from the one actually measured.

Figure 17 is reporting the values of the recall with a single nearest neighbour extracted for each query for the database reduction techniques presented. As it's possible to notice from the graph on the right of the figure 17, the method proposed to reduce the database can reduce the number of elements of about a half of the original size and, at the same time, obtain exactly the same recall as the baseline. Among the three techniques presented, the one that is considering only the distance and the time when the picture have been taken is the best one when the database size is very low while the other approaches are providing better results when the final number of images is not smaller than 40 percent of the initial size. Analysing the graphs in figure 17 it is possible to notice also that both the cosine and correlation techniques are producing results that are similar one to the other. In fact, many of the points in the graphs of the figure 17 are overlapped. Regarding the recall values measured, it is fundamental not to reduce the recall value from the baseline of more than one percent when the panoramas are filtered out since, this step, is used just before using an indexing technique to search efficiently the nearest neighbours and to further reduce the final memory footprint. The errors that the approximations done during the indexing phase is causing, that have been largely analysed in section 3, has to be summed to the one introduced by the database reduction technique and, since the advantages taken by the indexing techniques are more relevant for reaching the final objective of reducing the memory usage of the retrieval algorithm while maintaining a high recall, it is fundamental that this technique do not reduces too much the recall value in order to leave some margin for the indexing step and, in that way, find a good balance between the memory used and the recall metric obtained.

Table 9 is reporting a more detailed view of the results measured with the various techniques presented. The measurements reported clearly confirms how the correlation

Method	Side (m)	n	R@1	$N_{DB} (\times 10^6)$	RSS (GB)	Time (ms)
base	-	-	90.8	8.98	17.20 ± 0.01	31.41 ± 14.62
zone	2	-	90.5	4.13	7.90 ± 0.01	2.79 ± 0.21
zone	3	-	89.8	2.86	5.48 ± 0.01	1.98 ± 0.09
zone	5	-	89.5	1.59	3.04 ± 0.01	1.03 ± 0.02
correlation	3	1	90.5	2.86	5.48 ± 0.01	1.11 ± 0.06
correlation	15	9	88.8	3.01	5.77 ± 0.01	2.09 ± 0.44
correlation	7	2	87.5	1.95	3.74 ± 0.01	0.86 ± 0.03
cosine	3	1	90.5	2.86	5.48 ± 0.01	3.33 ± 1.11
cosine	4	1	89.1	2.08	3.97 ± 0.01	2.51 ± 0.44
cosine	7	2	87.5	1.95	3.74 ± 0.01	2.55 ± 0.34

Table 9: Some of the best results obtained with the database reduction techniques presented. The n parameter is the number of panoramas per zone kept. The time and the RSS memory measured are obtained with the exhaustive search.

and cosine techniques guarantee results that are very similar. For example, the recall value obtained with squared zones of side $3m$ and only one panorama kept per each zone are exactly the same for both the approaches. Those two techniques are guaranteeing results that are slightly better than the one obtained with the zone approach when the database size is still high. In fact, if compared the various methodologies with the same database size, it is possible to notice from table 9 that, for example, when the length of the side of the zones is set to 3 meters the recall value can be 0.7% higher with the cosine and correlation technique. This shows how effective are those two approaches in choosing the best panorama to be kept in the database. Instead, when the size of the final database is reduced below 2 million elements, the zone technique is producing results that are closer to the baseline than the ones that can be obtained with the cosine or correlation approaches. This behaviour is due to the fact that both of those techniques are not able to correctly choose the panoramas when the zones are too large since the number of images per zone is high and the differences between them may be significant making difficult to choose which one to keep while a simpler approach based on the time when the photo has been taken is guaranteeing always new images that surely contain most of the more recent visual characteristic of the place depicted.

Regarding the memory required, it is possible to notice from table 9 that the RSS used by the exhaustive search approach (that is the one used for the measurements reported in this table) can be reduced to approximately 32% of the memory required by the baseline and at the same time maintaining the loss of the recall limited to 0.3%. This is possible for the ability of the cosine and correlation techniques to choose the best panorama to be kept inside the database. Also the time required to extract the nearest neighbours for each query is taking advantage of this database reduction. In fact, this time it is dependent on the number of elements in the descriptors database and, as discussed in the previous sections, the lower the size of the database the lower the number of distance computations to be done during the nearest neighbours retrieval and also the time required. This is clearly shown in

table 9 that reports how it is easily possible to reduce up to 97% the amount of time needed for each query. All of the measurements performed with the various configuration tested of the three techniques used to reduce the size of the descriptors database are reported in appendix B.

For all of those reasons, it is strongly suggested to use techniques to reduce the database size like the one reported before applying any indexing technique in order to obtain a drastic reduction of the RSS memory required to store the indexes and the time for each query by maintaining at the same time high recall values. Among the approaches presented in this thesis project for this task, the division of the space in squared zones and the selection of the images to be kept based on the cosine similarity between them or the Pearson correlation coefficient are producing the best results when the database size is reduced up to the 70% of the original database. If instead a lower number of elements in the database is desired, the most promising technique is the one that is choosing just the most recent image for each of the zones produced.

6 Conclusions and further developments

In this thesis project has been presented various possible techniques that can be used to reduce the resources required by the retrieval part of the visual Geo-localization algorithm. In fact, it has been proposed some possible solutions that can be used in order to make the visual Geo-localization task is applicable to large-scale datasets obtaining results that are similar to the baseline but with the usage of just a small portion of the memory that would be required by a standard k-Nearest Neighbours approach. Among all of the techniques analysed a particular mention has to be done to the inverted file index with optimized product quantization and hierarchical navigable small world combined with an inverted file index and the optimized product quantization which are respectively producing the best results obtained in terms of recall metric and memory usage, as shown in section 4.8. In fact, those approaches, thanks to the memory reduction obtainable with the product quantization, are the ones that are guaranteeing the lowest memory footprints together with the higher recall values measured during this thesis project. Unfortunately, those techniques have shown also some limitations during the experiments done. In fact, both of those techniques are the ones that are, in the mean case, requiring the longer times to retrieve the nearest neighbours from the database of descriptors. This main limitation can be considered acceptable since the execution times remain in the order of milliseconds.

In the last part of this work, it has been also shown a technique that tries to reduce the size of the database of descriptors that are used by eliminating images that can be considered duplicates in the up-mentioned database. This approach is useful in the whole application since eliminates images that are not relevant for the retrieval of the nearest neighbours and, in that way, it can help the whole process reducing the resources needed. In fact, it has been shown how it is possible to filter the input database of panoramas discarding a significant number of them and obtain recall values that are comparable with the one that would be possible to measure if the full database was used. This, as reported before, helps the retrieval part also reducing the time required to retrieve from the database the required images since the retrieval time is proportional to the number of images to be analysed.

Both those techniques, the usage of indexes that perform an approximate search in the descriptors database and the removal from the database of the panoramas that are not actually useful to the process, can be combined together in order to fuse the advantages deriving from the application of them. Despite all of the methodologies reported in this thesis project, the problem of large-scale visual Geo-localization remains still an unsolved problem for a world scale application since the techniques proposed are able to just mitigate the unacceptable amount of resources required by this process. In fact, it has still to be studied a technique able to split and distribute the retrieval part of the algorithm among many different servers, each in charge of the retrieval of the images in a specific zone for example, being, in that way, capable to scale horizontally the application and at the same time being able to serve many more requests in parallel. In fact, the main limitation of all of the methodologies proposed is that they require that the resources needed have to be located in a single machine that must be capable of load the index needed in the main memory, requiring in that way servers with very high performances.

Despite the limitations reported, this thesis project has shown that it is possible to apply the visual Geo-localization algorithm to large-scale problems by keeping the memory usage of the whole application and at the same time also the time required to retrieve the nearest neighbours for each query acceptable.

A.2 Inverted File Index with Optimized Product Quantization

Dim.	IVF clusters	PQ probes	sub-vec.	bits	RSS (GB)	Time (ms)	R@1	Dim.	IVF clusters	PQ probes	sub-vec.	bits	RSS (GB)	Time (ms)	R@1				
512	1048576	1024	256	4	5.14	0.01	2.42	0.31	74.2	512	1048576	2048	256	4	5.15	0.01	3.25	0.32	74.2
512	1048576	4096	256	4	5.22	0.01	3.45	0.46	74.1	512	1048576	8192	256	4	5.37	0.01	6.58	0.31	74.2
512	1048576	16	512	4	5.76	0.01	1.52	0.14	74.4	512	1048576	32	512	4	5.76	0.01	1.54	0.16	75.3
512	1048576	64	512	4	5.76	0.01	1.54	0.07	75.6	512	1048576	128	512	4	5.76	0.01	1.68	0.19	76.2
512	1048576	256	512	4	5.76	0.01	1.74	0.01	76.4	512	1048576	512	512	4	5.76	0.01	2.39	0.27	76.5
512	1048576	1024	512	4	5.77	0.01	3.03	0.33	76.7	512	1048576	2048	512	4	5.78	0.01	4.69	0.47	74.8
512	1048576	4096	512	4	5.86	0.01	6.81	0.64	76.8	512	1048576	8192	512	4	6.00	0.01	11.57	0.31	76.8
512	1048576	16	128	8	5.09	0.01	1.62	0.18	74.4	512	1048576	32	128	8	5.09	0.01	1.81	0.27	75.1
512	1048576	64	128	8	5.10	0.01	1.89	0.19	75.7	512	1048576	128	128	8	5.10	0.01	2.20	0.14	76.2
512	1048576	256	128	8	5.12	0.01	2.78	0.32	76.3	512	1048576	512	128	8	5.13	0.01	4.13	0.26	76.4
512	1048576	1024	128	8	5.15	0.01	6.67	0.21	76.6	512	1048576	2048	128	8	5.17	0.01	11.63	0.38	76.6
512	1048576	4096	128	8	5.26	0.01	21.46	0.30	76.6	512	1048576	8192	128	8	5.41	0.01	41.61	0.93	76.4
512	1048576	8	256	8	5.76	0.01	1.57	0.18	73.0	512	1048576	16	256	8	5.76	0.01	1.55	0.04	74.6
512	1048576	64	256	8	5.76	0.01	1.65	0.02	75.3	512	1048576	64	256	8	5.77	0.01	2.01	0.23	75.6
512	1048576	256	256	8	5.76	0.01	2.35	0.16	76.5	512	1048576	256	256	8	5.76	0.01	3.15	0.07	76.5
512	1048576	1024	256	8	5.77	0.01	5.08	0.25	76.7	512	1048576	4096	256	8	5.79	0.01	8.32	0.73	76.9
512	1048576	2048	256	8	5.82	0.01	15.77	1.58	76.9	512	1048576	4096	256	8	5.92	0.01	26.91	0.18	76.9
512	1048576	8192	256	8	6.05	0.01	72.01	0.48	76.9	512	1048576	8192	256	8	7.10	0.01	16.89	0.01	75.0
512	1048576	16	512	8	7.10	0.01	1.77	0.17	74.9	512	1048576	32	512	8	7.10	0.01	1.89	0.01	75.0
512	1048576	64	512	8	7.10	0.01	2.35	0.03	75.4	512	1048576	128	512	8	7.10	0.01	3.52	0.29	76.2
512	1048576	256	512	8	7.10	0.01	5.20	0.36	76.2	512	1048576	512	512	8	7.10	0.01	9.63	0.86	76.4
512	1048576	1024	512	8	7.10	0.01	16.81	0.54	76.6	512	1048576	2048	512	8	7.10	0.01	32.15	0.49	76.6
512	1048576	4096	512	8	7.16	0.01	63.50	2.25	76.6	512	1048576	8192	512	8	7.36	0.01	121.54	3.41	76.6

A.3 Inverted File Index with Optimized Product Quantization and Re-Ranking

IVF						PQ						IVF						PQ																																																																																																																																																																																																																																																																																																																																																																																			
Dim.	clusters	probes	sub-vec.	bits	RSS (GB)	Time (ms)	R@1	Dim.	clusters	probes	sub-vec.	bits	RSS (GB)	Time (ms)	R@1	Dim.	clusters	probes	sub-vec.	bits	RSS (GB)	Time (ms)	R@1	Dim.	clusters	probes	sub-vec.	bits	RSS (GB)	Time (ms)	R@1																																																																																																																																																																																																																																																																																																																																																																						
2048	4096	4096	256	8	2.86	0.01	42.10	± 0.50	76.7	2048	4096	64	512	8	4.98	0.01	1.61	± 0.03	73.7	2048	4096	256	512	8	4.99	0.01	6.73	± 0.74	74.6	2048	4096	256	512	8	5.08	0.01	30.52	± 2.69	76.5	2048	4096	256	512	8	5.23	0.01	90.25	± 1.98	77.1	2048	8192	512	256	4	1.13	0.01	10.77	± 0.07	73.4	2048	8192	2048	256	4	1.28	0.01	40.90	± 0.09	73.8	2048	8192	8192	256	4	1.93	0.01	2.69	± 0.04	73.7	2048	8192	256	512	4	1.95	0.01	10.36	± 0.06	75.3	2048	8192	1024	512	4	2.02	0.01	40.98	± 0.07	76.1	2048	8192	4096	512	4	2.31	0.01	140.84	± 0.20	76.5	2048	8192	128	128	8	1.98	0.01	0.91	± 0.04	74.3	2048	8192	1024	128	8	2.01	0.01	3.39	± 0.03	75.0	2048	8192	4096	128	8	2.17	0.01	13.16	± 0.08	75.4	2048	8192	256	256	8	3.70	0.01	1.85	± 0.02	74.2	2048	8192	1024	256	8	3.74	0.01	8.19	± 0.99	75.1	2048	8192	4096	256	8	3.90	0.01	31.10	± 2.29	75.6	2048	8192	128	512	8	3.00	0.01	3.09	± 0.07	74.2	2048	8192	256	512	8	3.01	0.01	12.22	± 0.11	76.0	2048	8192	1024	512	8	3.10	0.01	48.85	± 0.43	76.9	2048	8192	8192	512	8	3.36	0.01	182.21	± 0.48	77.4	2048	16384	1024	256	4	1.43	0.01	7.19	± 0.66	73.8	2048	16384	4096	256	4	1.61	0.01	23.61	± 1.43	73.8	2048	16384	256	512	4	2.28	0.01	2.95	± 0.02	73.2	2048	16384	1024	512	4	2.31	0.01	11.42	± 0.12	74.6	2048	16384	4096	512	4	2.46	0.01	44.35	± 0.36	74.7	2048	16384	8192	512	4	2.60	0.01	14.29	± 0.45	75.4	2048	16384	1024	256	8	1.87	0.01	16.49	± 1.12	75.4	2048	16384	4096	256	8	2.07	0.01	59.02	± 0.87	75.9	2048	16384	128	512	8	3.18	0.01	3.04	± 0.28	73.0	2048	16384	512	512	8	3.19	0.01	74.37	± 0.81	75.4	2048	16384	2048	512	8	3.30	0.01	41.25	± 1.53	76.0	2048	16384	8192	512	8	3.67	0.01	155.39	± 0.71	76.1	2048	131072	1024	128	8	2.99	0.01	13.81	± 0.42	73.8	2048	131072	4096	128	8	3.24	0.01	50.96	± 0.16	74.3	2048	131072	8192	128	8	3.04	0.01	1.01	± 0.01	73.5	2048	262144	256	512	8	7.04	0.01	8.73	± 0.05	75.4	2048	262144	512	512	8	7.04	0.01	8.73	± 0.05	75.4	2048	262144	2048	512	8	7.09	0.01	29.93	± 0.10	76.4

A.3 Inverted File Index with Optimized Product Quantization and Re-Ranking

The following table is reporting all of the results obtained with the IVFPQR and OPQ IVFPQ REFINE techniques. For the inverted file index with product quantization and re-ranking with compressed residual error the number of sub-vectors used for the error has been fixed to 16 and each resulting dimension is using 8 bits to be represented.

Dim.	IVF		PQ		RSS (GB)	Time (ms)	R@1	Dim.	IVF		PQ		RSS (GB)	Time (ms)	R@1																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																										
	clusters	probes	sub-vec.	bits					clusters	probes	sub-vec.	bits																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																													
256	64	64	128	8	0.92	0.01	16.82	± 1.38	73.5	256	64	64	256	4	0.90	0.01	71.62	± 1.49	73.4	256	64	64	256	8	0.85	0.01	16.28	± 0.11	74.0	256	128	128	256	4	1.60	0.01	36.30	± 1.86	73.8	256	256	256	128	8	0.93	0.01	17.66	± 1.00	73.9	256	256	256	256	8	1.67	0.01	36.85	± 1.18	73.5	256	512	512	128	8	1.03	0.01	17.09	± 1.26	73.7	256	512	512	256	4	0.92	0.01	72.30	± 1.27	73.0	256	512	512	256	8	1.05	0.01	16.99	± 0.25	73.7	256	1024	1024	128	8	1.05	0.01	16.99	± 0.25	73.7	256	1024	1024	256	8	1.88	0.01	26.53	± 0.75	73.0	256	1024	1024	256	8	0.80	0.01	4.06	± 0.16	73.0	2048	1024	1024	64	8	1.60	0.01	39.97	± 0.19	73.2	2048	1024	1024	128	8	1.26	0.01	5.20	± 0.50	73.2	2048	1024	1024	128	8	1.32	0.01	18.42	± 1.12	74.0	2048	1024	1024	256	4	1.05	0.01	49.43	± 1.27	73.8	2048	1024	1024	256	8	2.10	0.01	2.56	± 0.17	73.3	2048	1024	1024	64	8	0.99	0.01	4.98	± 0.31	73.2	2048	1024	1024	128	8	1.19	0.01	0.67	± 0.03	73.2	2048	1024	1024	64	8	1.25	0.01	3.00	± 0.26	74.0	2048	1024	1024	128	8	2.04	0.01	1.97	± 0.10	73.2	2048	1024	1024	256	4	1.16	0.01	19.65	± 0.29	73.2	2048	1024	1024	256	8	3.71	0.01	7.34	± 0.73	73.3	2048	16384	1024	128	8	3.03	0.01	2.52	± 0.03	73.3	2048	16384	1024	256	8	1.27	0.01	6.17	± 0.65	73.5	2048	16384	1024	256	8	1.69	0.01	10.80	± 0.86	73.2	2048	16384	1024	64	8	2.72	0.01	1.05	± 0.07	73.1	2048	16384	1024	128	4	1.00	0.01	3.24	± 0.40	73.0	2048	16384	1024	256	8	1.08	0.01	9.09	± 0.69	73.1	2048	16384	1024	256	4	1.58	0.01	6.71	± 0.59	73.4	2048	16384	1024	256	8	1.80	0.01	18.43	± 1.47	73.2	512	64	64	64	8	0.58	0.01	7.57	± 0.73	74.6	512	64	32	128	8	0.89	0.01	9.08	± 1.01	74.9	512	64	32	256	4	0.89	0.01	39.04	± 1.50	74.0	512	64	32	256	8	0.58	0.01	19.22	± 0.69	75.0	512	64	32	512	4	1.55	0.01	88.41	± 2.92	74.8	512	64	32	512	8	2.97	0.01	38.68	± 0.30	74.8	512	64	64	64	8	0.52	0.01	7.87	± 0.57	74.2	512	128	128	128	8	0.85	0.01	10.92	± 0.30	73.3	512	128	128	256	4	0.84	0.01	88.94	± 3.86	76.2	512	128	128	256	8	1.60	0.01	44.73	± 1.50	76.7	512	128	128	256	8	2.98	0.01	39.96	± 2.50	73.2	512	128	128	64	8	0.62	0.01	4.52	± 0.15	73.3	512	256	256	64	8	0.92	0.01	35.75	± 0.46	73.7	512	256	256	128	8	0.67	0.01	20.71	± 0.68	75.8	512	256	256	256	4	0.95	0.01	87.67	± 3.35	76.6	512	256	64	256	8	1.68	0.01	23.51	± 0.68	74.9	512	256	64	512	4	1.62	0.01	36.92	± 0.98	73.7	512	256	256	512	4	1.64	0.01	139.66	± 0.90	77.2	512	256	128	512	8	3.14	0.01	40.47	± 0.99	75.4	512	256	256	512	8	0.64	0.01	8.71	± 0.33	73.4	512	512	64	128	8	0.99	0.01	2.48	± 0.28	73.6	512	512	256	128	8	1.01	0.01	9.48	± 0.95	74.9	512	512	128	256	4	1.01	0.01	11.32	± 0.47	73.5	512	512	128	256	4	1.02	0.01	37.88	± 1.50	75.3	512	512	128	256	8	1.74	0.01	5.33	± 0.53	73.4	512	512	128	256	8	1.75	0.01	19.17	± 0.98	75.4	512	512	128	256	8	1.63	0.01	19.31	± 0.76	75.3	512	512	128	256	4	1.64	0.01	73.27	± 1.36	75.3	512	512	128	256	8	3.26	0.01	11.27	± 1.03	73.3	512	512	128	256	8	3.27	0.01	39.58	± 1.31	75.4	512	512	128	256	8	3.27	0.01	2.11	± 0.23	73.3	512	1024	1024	64	8	0.77	0.01	8.04	± 0.69	74.4	512	1024	1024	128	4	0.69	0.01	37.61	± 1.43	73.9	512	1024	1024	128	8	1.02	0.01	4.60	± 0.19	74.2

A.4 Inverted Multi-Index with Optimized Product Quantization

Dim.	IVF clusters	IVF probes	PQ sub-vec.	RSS (GB)	Time (ms)	R@1	Dim.	IVF clusters	IVF probes	PQ sub-vec.	RSS (GB)	Time (ms)	R@1
1024	1048576	4	512	2.87 ± 0.01	71.86 ± 0.16	78.1	1024	1048576	8	512	2.87 ± 0.01	72.32 ± 0.88	78.1
2048	32	4	64	0.50 ± 0.01	6.39 ± 0.03	73.5	2048	32	8	64	0.51 ± 0.01	6.40 ± 0.04	73.5
1048576	4	512	2.87 ± 0.01	71.86 ± 0.16	78.1	1048576	8	512	2.93 ± 0.01	72.99 ± 1.94	77.0		

A.5 Hierarchical Navigable Small Worlds with Optimized Product Quantization

Dim.	IVF		PQ		m	RSS (GB)		Time (ms)		R@1	Dim.	IVF		PQ		m	RSS (GB)		Time (ms)		R@1
	clusters	probes	sub-vec.											clusters	probes		sub-vec.				
2048	1048576	8192	256	16	18.04	± 0.01	102.39	± 0.64	75.9		2048	1048576	32	256	32	18.20	± 0.01	16.48	± 0.20	73.5	
2048	1048576	64	256	32	18.20	± 0.01	14.29	± 0.17	74.3		2048	1048576	128	256	32	18.20	± 0.01	13.27	± 0.11	75.3	
2048	1048576	256	256	32	18.20	± 0.01	15.18	± 0.13	75.2		2048	1048576	512	256	32	18.20	± 0.01	18.77	± 0.15	75.6	
2048	1048576	1024	256	32	18.20	± 0.01	25.67	± 0.16	76.0		2048	1048576	2048	256	32	18.20	± 0.01	38.75	± 0.30	76.2	
2048	1048576	4096	256	32	18.20	± 0.01	63.57	± 0.18	76.3		2048	1048576	8192	256	32	18.20	± 0.01	112.44	± 0.35	76.4	

Table with 5 columns: Method, s (m), n, R@1, NDB (x10^6), RSS (GB), Time (ms). It contains two main sections of data for 'cosine' and 'zone' methods, comparing various parameters across different configurations. Each row includes a method name, a parameter value in meters (s), a count (n), accuracy (R@1), data rate (NDB), resource usage (RSS), and execution time (Time).

Method	s (m)	n	R@1	NDB ($\times 10^6$)	RSS (GB)	Time (ms)	Method	s (m)	n	R@1	NDB ($\times 10^6$)	RSS (GB)	Time (ms)
zone	6	-	86.8	1.27	2.42 \pm 0.01	1.08 \pm 0.05	zone	7	-	85.5	1.04	1.99 \pm 0.01	0.96 \pm 0.04
zone	8	-	84.9	0.88	1.68 \pm 0.01	0.28 \pm 0.03	zone	9	-	84.5	0.76	1.45 \pm 0.01	0.80 \pm 0.03
zone	10	-	83.2	0.66	1.27 \pm 0.01	0.53 \pm 0.01	zone	11	-	81.9	0.59	1.13 \pm 0.01	0.50 \pm 0.04
zone	12	-	81.6	0.53	1.01 \pm 0.01	0.37 \pm 0.03	zone	13	-	77.6	0.48	0.91 \pm 0.01	0.32 \pm 0.01
zone	14	-	75.0	0.43	0.83 \pm 0.01	0.34 \pm 0.04	zone	15	-	79.3	0.40	0.76 \pm 0.01	0.51 \pm 0.15
zone	16	-	73.7	0.37	0.70 \pm 0.01	0.24 \pm 0.01	zone	17	-	72.7	0.34	0.65 \pm 0.01	0.25 \pm 0.01
zone	18	-	75.0	0.32	0.60 \pm 0.01	0.29 \pm 0.01	zone	19	-	69.1	0.30	0.57 \pm 0.01	0.28 \pm 0.02
zone	20	-	71.1	0.28	0.53 \pm 0.01	0.20 \pm 0.01	zone	21	-	67.8	0.26	0.50 \pm 0.01	0.17 \pm 0.02
zone	22	-	71.1	0.25	0.47 \pm 0.01	0.19 \pm 0.00	zone	23	-	68.1	0.23	0.44 \pm 0.01	0.20 \pm 0.01
zone	24	-	64.1	0.22	0.42 \pm 0.01	0.07 \pm 0.01							

References

- [1] G. Berton, C. Masone, and B. Caputo, “Rethinking visual geo-localization for large-scale applications,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 4878–4888.
- [2] A. Torii, J. Sivic, T. Pajdla, and M. Okutomi, “Visual place recognition with repetitive structures,” in *2013 IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 883–890.
- [3] A. Torii, R. Arandjelović, J. Sivic, M. Okutomi, and T. Pajdla, “24/7 place recognition by view synthesis,” in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1808–1817.
- [4] F. Warburg, S. Hauberg, M. Lopez-Antequera, P. Gargallo, Y. Kuang, and J. Civera, “Mapillary street-level sequences: A dataset for lifelong place recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [5] M. J. Milford and G. F. Wyeth, “Mapping a suburb with a single camera using a biologically inspired slam system,” *IEEE Transactions on Robotics*, vol. 24, no. 5, pp. 1038–1053, 2008.
- [6] P. Gronát, G. Obozinski, J. Sivic, and T. Pajdla, “Learning and calibrating per-location classifiers for visual place recognition,” in *2013 IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 907–914.
- [7] J. Knopp, J. Sivic, and T. Pajdla, “Avoiding confusing features in place recognition,” in *Computer Vision – ECCV 2010*, K. Daniilidis, P. Maragos, and N. Paragios, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 748–761.
- [8] F. Radenović, G. Toliás, and O. Chum, “Fine-tuning cnn image retrieval with no human annotation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 7, pp. 1655–1668, 2019.
- [9] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems*, F. Pereira, C. Burges, L. Bottou, and K. Weinberger, Eds., vol. 25. Curran Associates, Inc., 2012. [Online]. Available: <https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>
- [10] S. Liu and W. Deng, “Very deep convolutional neural network based image classification using small training sample size,” in *2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR)*, 2015, pp. 730–734.
- [11] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.

-
- [12] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, “Netvlad: Cnn architecture for weakly supervised place recognition,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 5297–5307.
- [13] H. J. Kim, E. Dunn, and J.-M. Frahm, “Learned contextual feature reweighting for image geo-localization,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 3251–3260.
- [14] Y. Ge, H. Wang, F. Zhu, R. Zhao, and H. Li, “Self-supervising fine-grained region similarities for large-scale image localization,” in *Computer Vision – ECCV 2020*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds. Cham: Springer International Publishing, 2020, pp. 369–386.
- [15] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu, “Cosface: Large margin cosine loss for deep face recognition,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5265–5274.
- [16] J. Johnson, M. Douze, and H. Jégou, “Billion-scale similarity search with GPUs,” *IEEE Transactions on Big Data*, vol. 7, no. 3, pp. 535–547, 2019.
- [17] A. Babenko and V. Lempitsky, “The inverted multi-index,” in *2012 IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 3069–3076.
- [18] H. Jégou, R. Tavenard, M. Douze, and L. Amsaleg, “Searching in one billion vectors: re-rank with source coding,” in *ICASSP 2011 - International Conference on Acoustics, Speech and Signal Processing*. Prague, Czech Republic: IEEE, May 2011, pp. 861–864. [Online]. Available: <https://hal.inria.fr/inria-00566883>
- [19] T. Ge, K. He, Q. Ke, and J. Sun, “Optimized product quantization,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 4, pp. 744–755, 2014.
- [20] P. Schönemann, “A generalized solution of the orthogonal procrustes problem,” *Psychometrika*, vol. 31, no. 1, pp. 1–10, 1966. [Online]. Available: <https://EconPapers.repec.org/RePEc:spr:psycho:v:31:y:1966:i:1:p:1-10>
- [21] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>