

Where did I leave my phone ?

Cristian Reyes, Eva Mohedano, Kevin McGuinness and Noel E. O'Connor
Insight Centre for Data Analytics
Dublin, Ireland
cristian.reyes@estudiant.upc.edu
{eva.mohedano, kevin.mcguinness}@insight-centre.org

Xavier Giro-i-Nieto
Universitat Politecnica de Catalunya
Barcelona, Catalonia/Spain
xavier.giro@upc.edu

1. Introduction

The interest of users in having their lives digitally recorded has grown in the last years thanks to the advances on wearable sensors. Wearable cameras are one of the most informative ones, but they generate large amounts of images that require automatic analysis to build useful applications upon them. In this work we explore the potential of these devices to find the last appearance of personal objects among the more than 2,000 images that are generated everyday. This application could help into developing personal assistants capable of helping users when they do not remember where they left their personal objects. We adapt a previous work on instance search [3] to the specific domain of egocentric vision.

2. Methodology

Our goal is to rank the egocentric images captured during a day based on their likelihood to depict the location of a personal object. The whole pipeline is composed of the following stages: ranking by visual similarity, partition between candidate/non-candidate images and temporal-aware reranking within each class.

2.1. Ranking by Visual similarity

Given a certain set of query images Q depicting the object to be found, the algorithm starts by producing a ranking of the images of the day I ordered by their visual similarity score ν . This score is computed according to [3], which uses a bag of visual words model built with local features from a convolutional neural network (CNN).

A feature vector $q = f(Q)$ is generated from the set of images in Q that depict the object to locate. Three different approaches have been explored to define f :

a) No Mask: The q vector is built by averaging the visual words of all the local CNN features from the query images.

b) Mask: The q vector is built by averaging the visual words of the local CNN features that fall inside a query

bounding box that surrounds the object. This allows to consider only the visual words that describe the object.

c) Weighted Mask: The q vector is built by averaging the visual words of the local CNN features of the whole image, but this time weighted depending on their distance to the bounding box. This allows to consider the context in addition to the object.

2.2. Detection of Candidate Moments

As a second step, a thresholding technique is applied to the ranking in order to partition the I set into two subsets named Candidates (C) and Discarded (D) moments.

Two different thresholding techniques were considered in order to create the C and $D = I \setminus C$ sets: TVSS (Threshold on Visual Similarity Scores) and NNDR (Nearest Neighbor Distance Ratio). The TVSS technique builds $C = \{i \in I : \nu_i > \nu_{th}\}$. The NNDR technique is based in the one described by Loewe [2]. Let ν_1 and ν_2 be the two best scores, then it builds $C = \left\{i \in I : \frac{\nu_i}{\nu_1} > p_{th} \frac{\nu_2}{\nu_1}\right\}$.

2.3. Temporal-aware reranking

The temporal-aware reranking step introduces the concept that the lost object is not in the location with the best visual match with the query, but in the last location where it was seen. Image sets R_C and R_D are built by reranking the elements in C and D , respectively, based on their time stamps. The final ranking R is built as the concatenation of $R = [R_C, R_D]$.

We considered two strategies for the temporal reranking: a straightforward sorting from the latest to the earliest timestamp, or a more elaborate one that introduces diversity.

The diversity-aware configuration avoids presenting consecutive images of the same *moment* in the final ranked list. This is especially important in egocentric vision, where sequential images in time often present a high redundancy. Our diversity-based technique is based in the interleaving of samples, which is frequently used in dig-

ital communication. It consists in ordering temporally the images in I but knowing for each image if it belongs to C or D . So we might have something similar to $O = \{i_1^D, \dots, i_{k-1}^D, i_k^C, \dots, i_{l-1}^C, i_l^D, \dots, i_{m-1}^D, i_m^C, \dots, i_{n-1}^C\}$. Then $R_C = \{i_k^C, i_l^C, i_m^C, i_{k+1}^C, i_{l+1}^C, i_{m+1}^C, i_{k+2}^C, \dots\}$ and R_D is built analogously.

3. Experiments

3.1. Dataset annotation

Our work has been developed over the NTCIR Lifelogging Dataset [1] which consists of anonymised images taken every 30 seconds over a period of 30 days. Each day contains around 1,500 images.

This dataset was annotated for this work with five personal objects which could be lost: a phone, headphones, a watch and a laptop. In particular, they were tagged as *relevant* the last appearance of the object within each day.

Queries were defined by considering that the user had a collection of images of the object, not only one. The Q set contained from 3 to 5 images per category. These images showed the objects clearly and were used to build the q vector. This assumption is realistic as the object to be found could be defined from past appearances from the same dataset.

3.2. Training

The proposed system presents some parameters that were learned with the training part of the dataset.

A visual vocabulary for Bag of Words was learned from around 14,000 images of 9 days, generating a total of 25,000 centroids. The thresholds ν_{th} and p_{th} respectively were also learned on the same 9 days used for training. The optimal values found are detailed in Table 1.

	No Mask	Mask	Weighted Mask
ν_{th}	0.04	0.01	0.04
p_{th}	0.17	0.11	0.14

Table 1. Optimal thresholds. In bold those that gave highest mAP

3.3. Test

For evaluating the performance, Mean Average Precision (mAP) was computed for each day, taking into account all the categories. Then these values have been averaged over 15 test days and presented in Table 2.

Applying a thresholding technique has demonstrated to be helpful, as the combination of the object masking and the NNDR thresholding technique has shown the best results.

It must be noticed that mAP is not the best measure in diversity terms, so despite the fact that mAP decreases, the

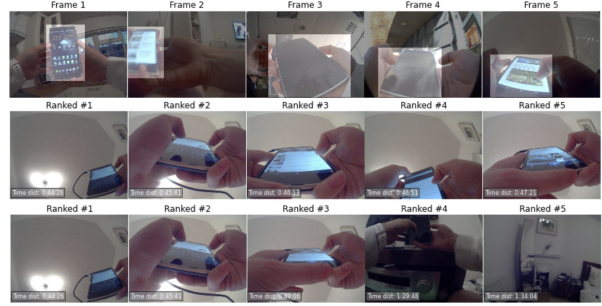


Figure 1. Results obtained for a search in category phone for a certain day. First row are the images that form Q with mask, second row results using NNDR and third results using NNDR + Div.

	No Mask	Mask	Weighted Mask
Temporal Ordering	0.051	0.051	0.051
Visual Similarity	0.102	0.082	0.111
TVSS	0.113	0.111	0.139
NNDR	0.086	0.176	0.093
TVSS + Div	0.096	0.082	0.118
NNDR + Div	0.066	0.166	0.049

Table 2. mAP results obtained when testing over 15 days.

images that form the top of the ranking have shown to be from more diverse scenes as it is shown in Figure 1.

4. Conclusions

This work has presented a good baseline for further research on the problem of finding the last appearance of an object in egocentric images.

Instance search based on bags of convolutional local features has shown promising results on egocentric images. Thresholding and temporal diversity techniques have improved the performance of visual only cues.

We plan to extend the annotations to neighbor images that may also depict relevant information to locate the location where the object was found. This way, not only one image would be considered as relevant, as assumed in the presented experiments.

References

- [1] C. Gurrin, H. Joho, F. Hopfgartner, L. Zhou, and R. Albat. Ntcir lifelog: The first test collection for lifelog research. In *Proceedings of the 39th International ACM SIGIR Conference on Research and Development in Information Retrieval*, Pisa, Italy, July 2016. ACM.
- [2] D. G. Lowe. Distinctive image features from scale-invariant keypoints. page 20, 2004.
- [3] E. Mohamedano, A. Salvador, K. McGuinness, F. Marques, N. E. O’Connor, and X. Giro-i Nieto. Bags of local convolutional features for scalable instance search. In *Proceedings of the ACM International Conference on Multimedia*. ACM, 2016.