

EGO-CH: Dataset and Fundamental Tasks for Visitors Behavioral Understanding using Egocentric Vision

Francesco Ragusa^{a,b}, Antonino Furnari^a, Sebastiano Battiato^a, Giovanni Signorello^c, Giovanni Maria Farinella^{a,c,*}

^aDMI-IPLab, University of Catania

^bXGD - XENIA s.r.l., Acicastello, Catania, Italy

^cCUTGANA, University of Catania

Abstract

Equipping visitors of a cultural site with a wearable device allows to easily collect information about their preferences which can be exploited to improve the fruition of cultural goods with augmented reality. Moreover, egocentric video can be processed using computer vision and machine learning to enable an automated analysis of visitors' behavior. The inferred information can be used both online to assist the visitor and offline to support the manager of the site. Despite the positive impact such technologies can have in cultural heritage, the topic is currently understudied due to the limited number of public datasets suitable to study the considered problems. To address this issue, in this paper we propose EGOcentric-Cultural Heritage (EGO-CH), the first dataset of egocentric videos for visitors' behavior understanding in cultural sites. The dataset has been collected in two cultural sites and includes more than 27 hours of video acquired by 70 subjects, with labels for 26 environments and over 200 different Points of Interest. A large subset of the dataset, consisting of 60 videos, is associated with surveys filled out by real visitors. To encourage research on the topic, we propose 4 challenging tasks (room-based localization, point of interest/object recognition, object retrieval and survey prediction) useful to understand visitors' behavior and report baseline results on the dataset.

Keywords: Egocentric Vision, First Person Vision, Localization, Object Detection, Object Retrieval

1. INTRODUCTION

Cultural sites receive many visitors every day. For a cultural site manager, it is hence paramount to 1) provide services able to assist the visitors, and 2) analyze their behavior to measure the performance of the site and understand what can be improved. For example using indicators [1] such as: a) Attraction index: to measure how much a point of interest attracts the visitors, b) Retention index: to measure the average time spent observing information element (e.g., a caption, a video a panel, etc.), c) Sweep Rate Index (SRI): it is used to calculate if visitors move slowly or quickly through the exhibition, d) Diligent Visitor Index (DVI): the percentage of visitors who

stopped in front of more than half of the points of interest. Classic approaches addressed the former task through the delivery of printed material (e.g., maps of the museum), the use of audio-guides and the installation of informative panels. Similarly, the analysis of visitors' behavior has generally been performed through the administration of questionnaires. It should be noted that such approaches often require manual intervention and are limited especially when the number of visitors is large. Recent works [2, 3, 4] have highlighted that the use of wearable devices such as smart glasses can provide a convenient platform to tackle the considered tasks in an automated fashion. Using such technology, it is possible to provide to the user services such as automated localization (e.g., to help visitors navigating the site) and recognition of currently observed

*Corresponding author

Email address: gfarinella@dmf.unict.it (Giovanni Maria Farinella)

Ponts Of Interest (POIs)¹ to provide more information on relevant objects and suggest what to see next. Conveniently, localization and POI recognition can be used by the manager of the cultural site to obtain information about the visitors and understand their behavior by inferring where they have been, how much time they have spent in a specific environment and what POIs have been liked most.

Despite the aforementioned technologies can have a significant impact on cultural heritage, they are currently under-explored due to the lack of public benchmark datasets. To address this issue, in this paper we propose EGOcentric-Cultural Heritage (EGO-CH), the first large dataset of egocentric videos for visitors behavioral understanding in cultural sites. The dataset has been collected in two cultural sites located in Sicily, Italy: Galleria Regionale di Palazzo Bellomo² and Monastero dei Benedettini³. The overall dataset contains more than 27 hours of video, including 26 environments, over 200 Points of Interest and 70 visits. We release EGO-CH with a set of annotations useful to tackle fundamental tasks related to visitors behavior understanding in cultural sites, and specifically, temporal labels specifying the location of the visitor as well as the currently observed POI, bounding box annotations around POIs, surveys filled out by visitors at the end of each tour in the cultural site. Figure 1 reports some sample frames from the proposed dataset. The dataset can be publicly accessed upon request to the authors from our webpage <http://iplab.dmi.unict.it/EGO-CH/>.

We propose 4 fundamental tasks for visitors behavioral understanding using egocentric vision: 1) *room-based localization*, consisting in recognizing the environment in which the visitor is located in each frame of the video, 2) *Point of Interest recognition*, which consists in correctly detecting and localizing all objects in the image frames, 3) *object retrieval*, which consists in matching an observed object from the egocentric point of view to a reference image contained in the museum catalogue

¹In this work, we refer to the definition of Point Of Interest (POI) given in [5], as an element which can attract the attention of visitors. Most POIs are objects such as paintings and statues, but architectural elements such as pavements can qualify as POIs, despite not being objects. Therefore, in this paper the notations “Point Of Interest” and “object” are not used interchangeably.

²<http://www.regione.sicilia.it/beniculturali/palazzobellomo/>.

³<http://www.monasterodeibenedettini.it/>

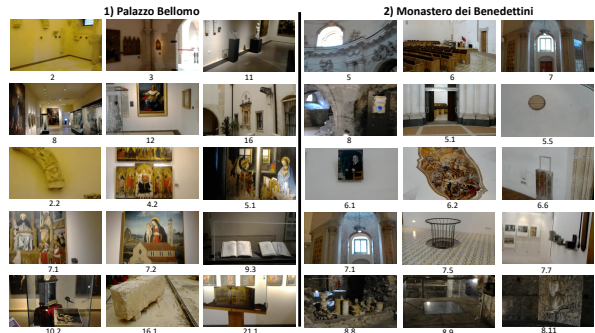


Figure 1: Sample frames from the two cultural sites belonging to EGO-CH: 1) Palazzo Bellomo, 2) Monastero dei Benedettini. The first two rows show frames extracted from the training videos and related to the environments, whereas the remaining rows show frames of the training videos related to POIs. See Section 5.2 for more details.

of all artworks, 4) *survey prediction*, which consists in generating the survey associated to a visit from video. We also provide baseline results for each task on the proposed dataset. The experimental results suggest that the proposed dataset is a challenging benchmark for visitors behavioral understanding using egocentric vision.

In sum, the contributions of this work are: 1) we present EGO-CH, a new challenging dataset of egocentric videos acquired in two cultural sites, 2) the dataset has been labeled to tackle 4 main tasks useful to understand visitors behavior, 3) we report baseline results for each task.

2. RELATED WORK

Visitors Behavioural Understanding and Site Manager Assistance in Cultural Sites

Several works investigated the use of wearable systems to augment the fruition in cultural sites [2]. Razavian et al. [6] proposed a method to estimate the attention of the visitors of an exhibition, whereas in [7] a CNN to perform localization and object recognition is introduced in order to develop a context aware audio guide. Raptis et al. [8] studied the design of mobile applications in museum environments and highlighted that context influences interaction. In [3, 4], the problem of localizing the visitors of a museum from egocentric videos is considered. The inferred localization can be used to provide behavioral information to the manager of the site. Past works investigated specific applications, generally relying on data collected on purpose and not publicly released. In this work, we

aim at standardizing the fundamental problems of visitors behavioral understanding in cultural sites by proposing a public dataset and a series of tasks.

Datasets on Cultural Heritage Few image-based datasets focusing on cultural heritage have been proposed in past works. Koniusz et al. [9] proposed the OpenMIC dataset containing photos captured in ten different exhibition spaces of several museums and explored the problem of artwork identification. DelChiaro et al. [10] proposed NoisyArt, a dataset composed of artwork images collected from Google Images and Flickr correlated by metadata gathered from DBpedia. In contrast with the aforementioned works, we propose the first dataset composed of egocentric videos, and release it publicly. The dataset can be used to address different tasks related to visitors behavioral understanding in cultural sites. A significative part of the proposed dataset has been collected by real visitors (i.e., 60 visits) and hence it is a realistic set of data for benchmarking.

Localization Ahmetovic et al. [11] presented NavCog, a system to navigate with a smartphone in complex indoor and outdoor environments exploiting Bluetooth Low Energy beacons. Kendall et al. [12] proposed to infer the 6 Degrees of Freedom pose of a camera from egocentric images using a CNN. In [3], it has been considered the problem of localizing a visitor in a cultural site from egocentric images to provide behavioral information to the site manager. In this work, we consider the work presented in [3] as a baseline for the localization task.

Point Of Interest/Object Recognition Seidenari et. al [7] and Taverriti et al. [13] proposed to perform object classification and artwork recognition to assist tourists with additional information about the observed objects. In general, object detectors (e.g., YOLOv3 [14]) have been used to detect artworks in cultural sites. However, it should be noted that, as pointed out in [5], depending on the cultural site, not all Points Of Interest are objects. For instance, a point of interest can be an architectural element such as a pavement, or even a corridor. In this case, it should be considered that object detectors can be limited. In this work, we consider the YOLOv3 object detector [14] as baseline for Point Of Interest/Object recognition.

Object Retrieval Many previous works investigated approaches to image retrieval. Rubhasy et al. [15] used an ontology-based approach to retrieval in multimedia cultural heritage collections.

Environment	#video	#frame
1 Sala1	1	3721
2 Sala2	4	7968
3 Sala3	4	8285
4 Sala4	5	11497
5 Sala5	5	11461
6 Sala6	1	2630
7 Sala7	4	10613
8 Sala8	2	6910
9 Sala9	4	10505
10 Sala10	3	5830
11 Sala11	3	7343
12 Sala12	1	2463
13 Sala13	1	3040
14 Cortile degli Stemmi	2	5853
15 Sala delle Carrozze	1	3259
16 Cortile Parisio	4	10374
17 Biglietteria	2	4099
18 Portico	2	5701
19 Scala Catalana	2	6399
20 Loggetta	2	4169
21 Box Sala8	2	4661
22 Area Sosta	2	3505
Total	57	140286

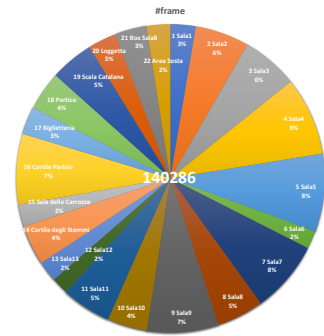


Figure 2: Number of training videos collected in each environment and corresponding number of frames for the cultural site “Palazzo Bellomo” (left), along with a pie chart representation of the same data (right).

The goal is to enable the integration of different types of cultural heritage media and to retrieve relevant heritage media given a query. Kwan et al. [16] proposed matrix of visual perspectives to address Content-based Image Retrieval (CBIR) of cultural heritage symbols, whereas Iakovidis et al. [17] perform pattern-based Content-based Image Retrieval. The work of [18] focused on discarding image outliers using Content-based Image Retrieval. Despite the availability of advanced approaches, for generality and ease of comparison, in this paper we consider simple baselines based on image representation and nearest neighbor search to address the object retrieval task.

3. THE EGO-CH DATASET

3.1. Data Collection

The dataset has been acquired using a head-mounted Microsoft HoloLens device in two cultural sites located in Sicily, Italy: 1) Palazzo Bellomo (Table 1), located in Siracusa⁴, and 2) Monastero dei Benedettini (Table 2), located in Catania⁵.

Palazzo Bellomo This cultural site is composed of 22 environments and contains 191 Points of Interest (e.g., statues, paintings, etc.).⁶ Training videos have been collected by operators instructed to walk around in order to capture images of each environment from different points of view. To simplify labeling, each training video contains

⁴<http://www.regione.sicilia.it/beniculturali/palazzobellomo/>
⁵<http://monasterodeibenedettini.it/>

⁶See the supplementary material for the list of environments and POIs.

Table 1: Details regarding the cultural site "Palazzo Bellomo".

Subset	Resolution	FPS	AVG Time (min)	# POIs	#environments	bbox annotations	temporal segments
Training	1280x720	29.97	1.4	191	22	56686	57
Test	1280x720	29.97	31.27	191	22	13402	340

Table 2: Details regarding the cultural site "Monastero dei Benedettini".

Subset	Resolution	FPS	AVG Time (min)	# POIs	#environments	bbox annotations	temporal segments
Training	1216x684	24.00	2.2	35	4	33366	48
Validation	1216x684	24.00	3.5	35	4	2235	20
Test	1408x792	30.03	21	35	4	71310	455

only frames from a given environment. At least one training video has been collected per environment. In the case of outdoor environments (e.g., courtyards), we collected multiple videos to include different lighting conditions. We have collected a total of 57 training video in this cultural site. Figure 1(left) shows some frames acquired in the considered cultural site, whereas Figure 2 reports the number/percentage of frames acquired in each environment. Ten test videos have been collected separately asking 10 volunteers to visit the cultural site. One of the 10 videos (i.e., "Test 3") was selected randomly and used as validation set, whereas the remaining 9 videos are used for evaluation purposes. No specific instructions on where to go, what to look at and how much time to spend in a specific environment/POI has been provided to the visitors. Most of the subjects had limited confidence with the cultural site. This provided a natural means to collect realistic data of visitors exploring the environments and observing Points of Interest. All the videos have a resolution of 1280×720 pixels and a frame-rate of 29.97 fps. The average duration of test videos is 31.27 min, with the longest one being 50.23 min. See the supplementary material for more details about training/test videos. We also include 191 reference images related to the considered POIs to be used for one-shot image retrieval. The images are akin to the images generally included in museum catalogs.⁷

Monastero dei Benedettini This dataset is composed of 4 environments and contains 35 Points Of Interest.⁸ Differently from "Palazzo Bellomo", the POIs belonging to this cultural site include both objects such as paintings and statues as well as architectural elements, such as pavements, which



Figure 3: Some example bounding box annotations from the cultural site "Monastero dei Benedettini".

cannot be easily recognized using object detection techniques as noted in [5]. See Figure 1(right) for some qualitative examples of the considered objects. Training videos have been collected with the same acquisition modality considered for the "Palazzo Bellomo" cultural site. Figure 4 reports the number/percentage of frames acquired in each environment. Training and validation videos have a resolution of 1216×684 pixels and a frame-rate of 24 fps. Five validation videos have been collected by asking volunteers to visit the cultural site following the same protocol used for "Palazzo Bellomo". Additionally, we collected 60 test videos by asking real visitors inexperienced with both the research project and its goals and the HoloLens device to freely visit the cultural site. No specific instructions have been given to the visitors, who were free to explore the 4 environments and the 35 POIs. This allowed us to obtain realistic data of how a visitor would move in a cultural site. Test videos have been collected over a period of three months. Moreover, at the end of the visit, we administered the visitor a survey, the content of which is described in Section 5.2.2. The 60 test videos have a resolution of 1408×792 pixels and a frame-rate of 30.03 fps. The average video length is 21 min, with the maximum

⁷Examples reference images for both cultural sites are included in the supplementary material.

⁸See the supplementary material for the list of environments and POIs.

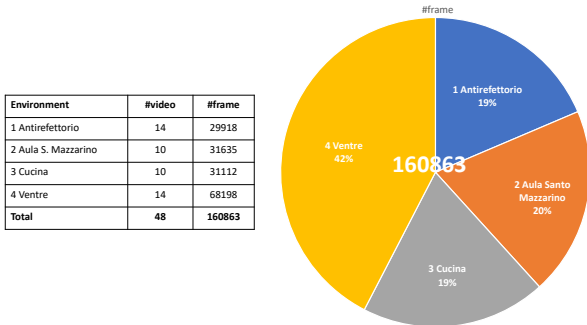


Figure 4: Number of training videos collected in each environment and corresponding number of frames for the cultural site “Monastero dei Benedettini” (left), along with a pie chart representation of the same data (right).

length being 42 *min.* See the supplementary material for more details about training/validation/test videos. Similarly to “Palazzo Bellomo”, we include 35 reference images related to the considered POIs for one-shot image retrieval⁷. Please note that this set of data is adapted from and extends significantly the dataset proposed in [3], introducing 60 new labelled videos collected by real visitors. Specifically, the overall dataset presented in this work contains +1600 minutes of video, data from +70 more subjects, +91369 bounding box annotations and an additional cultural site “Palazzo Bellomo” comprising 22 environments and 191 points of interest.

3.2. Annotations

Temporal Labels All test and validation videos have been temporally labeled to indicate in every frame the environment in which the visitor is located and the observed point of interest, if any. If the visitor is not located in one of the considered environment (e.g., a stair), the frame is marked as “negative”⁹. It is worth noting that there are no negative frames in “Palazzo Bellomo” since all environments are part of the museum, whereas negative frames are contained in “Monastero dei Benedettini”. This is due to the different nature of the two sites: “Palazzo Bellomo” is a museum, consisting in a limited set of rooms, whereas “Monastero dei Benedettini” is a much more complex environment including many corridors and stairs which have not been labeled as locations of interest for visitors. Similarly, we mark as “negative” all frames in which

⁹ Examples of “negative” frames are reported in the supplementary material.

the visitor is not observing any of the considered POIs. Each location is identified by a number that denotes a specific environment (1 – 22 for “Palazzo Bellomo” and 1 – 4 for “Monastero dei Benedettini”). Each point of interest is denoted by a code in the form X.Y (e.g., 3.5) where “X” denotes the environment in which the point of interest is located and “Y” identifies the point of interest. See Figure 1 for some examples.

Bounding Box Annotations A subset of frames from the dataset (sampled at at 1 fps) has been labeled with bounding boxes indicating the presence and locations of all POIs. Specifically, each POI has been labeled with a tuple $(class, x, y, w, h)$ indicating the class of the POI and its bounding box information. It is worth mentioning that, as noted in [5], a POI can be an object (e.g., a painting or a statue) or a different element (e.g., a pavement or a specific location), which cannot be strictly defined as an object. Indeed, the kind of POIs contained in a cultural site depends on the nature of the site itself. In EGO-CH, “Palazzo Bellomo” contains only objects as POIs, whereas “Monastero dei Benedettini” contains both objects and other elements. Nevertheless, all elements are labeled with class type and bounding box annotations. Figure 3 shows examples of labeled frames from the 60 visits of “Monastero dei Benedettini”.

Surveys The 60 test videos collected in the “Monastero dei Benedettini” are associated with surveys which have been administered to the visitors at the end of the visits. Specifically, the visitors are asked to rate a subset of 33 out of the 35 Points Of Interest (a picture of each point is shown) or specify if any of them had not been seen it during the visit. The rating is expressed as a number ranging from -7 (not liked) to $+7$ (liked).

The EGO-CH dataset is publicly available at our website: <http://iplab.dmi.unict.it/EGO-CH/>. The reader is referred to the supplementary material for more details about the dataset and the experiments. The dataset can be used only for research purposes and is available upon the acceptance of an agreement.

4. PROPOSED TASKS AND BASELINES

In this Section, we propose four tasks which can be addressed using the proposed dataset. The tasks are related to problems investigated in previous works on cultural heritage [4, 3, 7, 13]. We believe

that solving these tasks can bring useful information about the behavior of the visitors of a cultural site.

4.1. Room-based Localization

Task: The task consists in determining the room in which the visitor of a cultural site is located from egocentric images collected using a wearable device. Localization information can be used both to provide a “where am I” service to the visitor and to collect behavioral information useful for the site manager to understand what paths do visitors prefer and where they spend more time in the cultural site.

Baseline: As a baseline for this task, we consider the approach proposed in [3, 19]. This approach is selected as a baseline due to the limited work on room-based localization in the cultural heritage domain [3] and due to the state-of-the-art performance of the approach shown in [19]. Given a set of locations, the considered approach allows to segment a given video into video shots related to the specified locations. If a given shot is not related to any of the locations, the algorithm automatically labels it as a “negative segment” through a “negative rejection” stage. The method is composed by three steps, as illustrated in Figure 5. For each cultural site, we trained a VGG-19 CNN to discriminate between locations (“Discrimination” stage). The “Negative Rejection” step has been considered only for the data of “Monastero dei Benedettini”, since “Palazzo Bellomo” does not contain negative locations. The “Sequential Modeling” stage allows to obtain a temporal segmentation of the input video where each segment is associated to one of the considered environments. This algorithm is chosen as it achieves state-of-the-art performance in the task of location-based egocentric video segmentation [19, 3]. Two hyper-parameters are involved in the algorithm: K , related to the “negative rejection” stage and ϵ , which regulates the amount of temporal smoothing applied to the predictions. The reader is referred to [3] for more details.

Implementation Details and Evaluation

Measures: We evaluated our method following [3] using FF_1 score and ASF_1 score. Specifically, the FF_1 score is the F_1 score applied to individual frames and, as such, it does not evaluate the ability of the methods to produce a temporally coherent segmentation. ASF_1 is the F_1 score applied to temporal segments rather than frames and measures the ability to detect video segments coherent

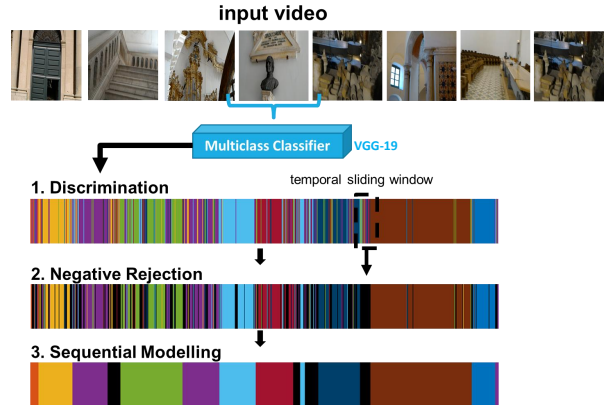


Figure 5: The method used to perform room-based localization. The method is composed by three steps: 1) Discrimination, 2) Negative Rejection, 3) Sequential Modeling. See [19, 3] for more details.

with the ground truth. Both scores are normalized between 0 and 1. The hyper-parameters of the algorithm K and ϵ are tuned on the validation sets of the proposed dataset. Specifically, $\epsilon = 10^{-273}$ is found by optimizing the validation ASF_1 score with a grid search in the range $[10^{-1} : 10^{-299}]$ on “Palazzo Bellomo” (see [3] for details). Since no negative locations are contained in “Palazzo Bellomo”, the “negative rejection” stage is not performed and hence the parameter K is not optimized. Similarly, we find $\epsilon = 10^{-89}$ and $K = 100$ on “Monastero dei Benedettini”.¹⁰

Results: Table 3 reports the results obtained by the baseline in the two cultural sites¹¹. On “Palazzo Bellomo”, the baseline achieves good FF_1 scores for most rooms, obtaining an average value of 0.81. Much lower results are observed when the ASF_1 score is considered. In this case, an average value of 0.59 is reached. Lower results equal to 0.68 and 0.40 are obtained in the “Monastero dei Benedettini”. This is partly due to the presence of negatives, which are not included in “Palazzo Bellomo” and to the more challenging nature of the test set of “Monastero dei Benedettini”, which contains 60 videos collected by real visitors within 3 months with different lighting condition and blur as shown in Figure 6. The overall results highlight that addressing the considered task on the proposed

¹⁰The supplementary material reports more implementation details.

¹¹Extended tables, qualitative results and confusion matrix are included in the supplementary material.

Table 3: Room-based localization results. For each cultural site, the last row reports the Average (AVG) of the FF_1 and ASF_1 scores.

1) Palazzo Bellomo		
Room	FF_1 score	ASF_1 score
Sala1	0.71	0.48
Sala2	0.92	0.79
Sala3	0.84	0.50
Sala4	0.92	0.59
Sala5	0.94	0.64
Sala6	0.77	0.52
Sala7	0.94	0.61
Sala8	0.89	0.64
Sala9	0.91	0.47
Sala10	0.84	0.69
Sala11	0.84	0.58
Sala12	0.80	0.66
Sala13	0.80	0.66
Cortile degli Stemmi	0.85	0.64
Sala Carrozze	0.91	0.67
Cortile Parisio	0.75	0.50
Biglietteria	0.65	0.44
Portico	0.69	0.51
Scala Catalana	0.76	0.63
Loggetta	0.71	0.51
Box Sala8	0.94	0.79
Area Sosta	0.43	0.47
AVG	0.81	0.59

2) Monastero dei Benedettini		
Class	FF_1 score	ASF_1 score
Antirefettorio	0.75	0.54
Aula S. Mazzarino	0.33	0.12
Cucina	0.79	0.34
Ventre	0.97	0.60
Negative	0.54	0.33
AVG	0.68	0.40

dataset is challenging. In particular, issues such as varying lighting conditions and the presence of negatives need to be addressed in task-specific investigations.

4.2. Point of Interest/Object Recognition

Task: This task consists in recognizing the points of interest which the user is looking at. This can be useful to understand the visitor’s behavior and answer questions as “What are the most viewed points of interest?” and “How long have they been observed?”. Moreover, a system able to recognize points of interest could suggest the visitor what to see next, as well as provide information with Augmented Reality. The dataset could be used to perform standard object detection task.

Baseline: Due to its real-time performance and to its popularity in the cultural heritage domain [5, 7, 13], we consider a YOLOv3 [14] object detector as a baseline for the task. The detector has been trained on the training sets of “Palazzo Bellomo” and “Monastero dei Benedettini”.



Figure 6: Some sample frames from different visits acquired within 3 months. Each row represents similar positions in the same environment with different lighting conditions.

Implementation Details and Evaluation

Metrics: We trained YOLOv3 using the standard anchors provided by the authors for the COCO dataset. We use mean Average Precision (mAP) with threshold on IoU equal to 0.5 for the evaluations. In order to use YOLOv3 to detect artworks, a detection threshold is specified to discard detections with low confidence scores. For each cultural site, we tuned this threshold on the validation sets by choosing the value which maximizes mAP in the range [5⁻⁴; 1⁻³; 5⁻³; 1⁻²; 3⁻²; 5⁻²; 0.10; 0.15; 0.2; 0.25; 0.3; 0.35; 0.40]. To train the detector on “Palazzo Bellomo”, we set the initial learning rate to 0.001 and the detection threshold to 0.01. On “Monastero dei Benedettini”, we set the initial learning rate to 0.01 and the detection threshold to 0.001.

Results: Table 4 reports the results obtained in the two cultural sites. The results obtained on “Palazzo Bellomo” are much lower than the ones obtained on “Monastero dei Benedettini” mainly because of the larger set of POIs contained in the former site (191) versus the lower number of POIs contained in the latter (35). In both cases, the results are in general very low, which highlights the challenging nature of the proposed dataset and tasks. Among the challenges of the dataset, as previously discussed, it should be considered that some of the points of interest represent architectural elements such as corridors or pavements, which might

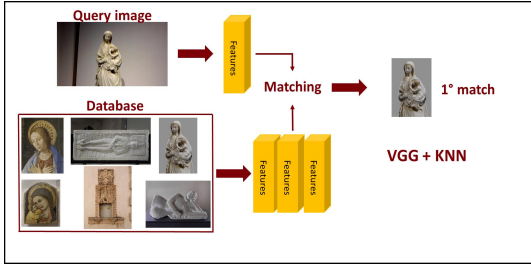


Figure 7: Diagram of the baseline for the object retrieval task.

Table 4: Object detection results. The reported mean Average Precision (mAP) is averaged over all test videos. Per-class Average Precision (AP) values are reported in the supplementary material.

Cultural Site	mAP
1) Palazzo Bellomo	10.59%
2) Monastero dei Benedettini	15.45%

be challenging to detect with a simple object detector, as pointed out in [5]. Moreover, differently from other object detection tasks, POIs here need to be recognized at the instance level. For instance, the dataset contains multiple paintings which should be recognized as separate objects. We leave the investigation of more specific approaches to future investigations.

4.3. Object Retrieval

Task: Given a query image containing an object, the task consists in retrieving an image of the same object from a database. This task can be useful to perform automatic recognition of artworks when detection can be bypassed, i.e., when the user places the artwork in the center of the field of view using a wearable or mobile device. Moreover, the task is particularly of interest especially considering that artwork detection is a hard task, as highlighted in the previous section. We obtain a set of query images by extracting image patches from the bounding boxes annotated in the test set and consider two variants of the task. This accounts to 23727 image patches for “Palazzo Bellomo” and 44978 image patches for “Monastero dei Benedettini”.¹² We consider two variants of this task. In the first variant, object retrieval is framed as a one-shot retrieval problem. In this case, the database contains

¹²The supplementary material reports examples of extracted image patches.

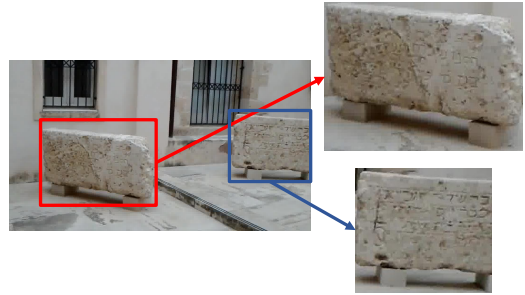


Figure 8: Example of patches extracted using bounding boxes annotations.

only the reference images associated to each POI, whereas the whole set of image patches is used as the test set, i.e., only a single labeled sample is assumed to be available for each object. In the second variant, we split the set of image patches into a training set (70% - used as DB) and a test set (30%). It should be noted that the first variant of the task is much more challenging both due to the presence of few labeled samples and to the domain shift which affects the two sets of images: reference images for the POIs and image patches cropped from egocentric images. Figure 8 shows an example of image patches cropped from the egocentric images using bounding box annotations.

Baseline: Given the lack of investigation of approaches for retrieval in the scenario of First-Person vision in the cultural heritage domain, we consider a simple image-retrieval pipeline for both variants of the task. The pipeline uses VGG19 CNN pre-trained on ImageNet to represent image patches, while matching is performed and matched using a K-NN. A scheme of the considered baseline is shown in Figure 7.

Implementation Details and Evaluation Measures: We have extracted all features from the FC7 layer of the VGG19 network pre-trained on ImageNet. When the second variant of the task is considered, we perform K-NN using $K = [1; 3; 5]$. We evaluated the performance of our baseline using standard metrics for image-retrieval: precision, recall and F_1 score.

Results: Table 5 shows the results of the baseline on the image retrieval variants. In both cultural sites, one-shot retrieval does not achieve good results. This is probably due to the fact that one-shot retrieval relies on a limited number of training samples, which are drawn from a different distribution as compared to test samples. This suggests

Table 5: Object retrieval results for the two variant of the task.

Points of Interest Retrieval				
1) Palazzo Bellomo				
Variant	K	Precision	Recall	F ₁ score
1 - One Shot	1	0.004	0.007	0.001
	1	0.69	0.66	0.67
	3	0.69	0.62	0.62
2 - Many Shots	5	0.69	0.62	0.62
	7	0.68	0.62	0.62
	9	0.67	0.61	0.62
	11	0.67	0.61	0.61
2) Monastero dei Benedettini				
Variant	K	Precision	Recall	F ₁ score
1 - One shot	1	0.29	0.07	0.08
	1	0.87	0.87	0.87
	3	0.88	0.87	0.87
2 - Many Shots	5	0.88	0.88	0.88
	7	0.88	0.87	0.87
	9	0.87	0.87	0.87
	11	0.87	0.86	0.86

that dedicated methodologies should be considered to tackle one-shot retrieval and the domain shift problem. Better results are obtained on both sites in the second variant of the task, when the effect of one-shot retrieval and domain shift is reduced. Best results are obtained in “Palazzo Bellomo” for $K = 1$ (F_1 score of 0.67) and in “Monastero dei Benedettini” for $K = 5$ (F_1 score of 0.88).

4.4. Survey Prediction

Task: Each test video of the “Monastero dei Benedettini” is associated to a survey collected from visitors at the end of the visit. We define this task as predicting the content of a survey from the analysis of the related egocentric video. We deem this to be possible as the egocentric video contains information on what the visitor has seen during the visit. In particular, the task consists in predicting for each POI 1) if the POI has been remembered by the visitor and 2) how the POI would be rated by the visitor in a $[-7, 7]$ scale. This task investigates automatic algorithms for automatically “filling in” surveys from videos.

Baseline: Since the proposed task is novel and very challenging, as a proof of concept, we propose a baseline which takes as input the temporal annotations indicating the objects observed by the visitors in the 60 visits. To obtain fixed-length descriptors for each video, we accumulate the number of frames in which a given POI has been observed in a Bag Of Word representation. In such representation, each component of the fixed-length vector indicates the total time in which a

specific point of interest has been observed by the visitor. The vector is hence sum-normalized to reduce the influence of videos with different lengths. The whole training set is normalized with z-scoring and classification is performed using K-NN. We consider two baselines. The first one simply performs a binary classification to predict whether a POI has been remembered by the visitor or not. The second one predicts both if the POI has been seen and what score has been assigned to it. This is tackled as a 15-class classification problem, where class -8 indicates that the POI has not been remembered, whereas the other 14 classes represent the scores from -7 to 7 assigned by the visitors to POIs. We would like to note that we treat the problem as a classification task, as the scores assigned by the visitors are discrete integer numbers. Also, the dataset contains a limited set of data-points, which would prevent the algorithm from generalizing beyond the discrete set of labels available at training time.

Implementation Details and Evaluation Measures: We perform our experiments using a leave-one-out strategy. We tested different values for k ranging from 1 to 9 and chose $K = 9$ which resulted to be optimal in our experiments. We evaluate results with weighted precision, recall and F_1 score.

Results: Table 6 reports the results obtained in the case of binary classification (remembered vs not remembered)¹³. The number of instances belonging to each class is reported in the last column. The results suggest that this task is very challenging. Indeed, even if a POI appears in some frames, this does not imply that the visitor remembers it. Table 7 shows that the multi-class task¹³ is even more challenging, with classes containing fewer examples (e.g., $-6, -5, -4, -3$) hard to recognize. As a final remark, it is worth noting that the results suggest that the task can be addressed to some degree. We expect that more complex approaches leveraging the analysis of the semantics of the input videos and the estimation of the attention of the visitor can achieve much better performance.

The code of our baselines is public available. See our web-page for the details: <https://iplab.dmi.unict.it/EGO-CH/#code>.

¹³ See the supplementary material for the extended tables.

Table 6: Survey prediction results - binary classification task.

Class	Precision	Recall	F ₁ score	support
Not Remembered	0,43	0,2	0,27	561
Remembered	0,74	0,89	0,81	1419
AVG	0,65	0,7	0,66	1980

Table 7: Survey prediction results - multi-class classification. “Weighted AVG” reports the average scores weighted by the number of samples in each class.

Class	Precision	Recall	F ₁ score	Support
Not Remem.	0,32	0,63	0,43	561
-7	0,52	0,24	0,33	49
-6	0	0	0	8
-5	0	0	0	8
-4	0	0	0	5
-3	0	0	0	5
-2	0,09	0,08	0,08	13
-1	0	0	0	10
0	0,18	0,15	0,17	104
1	0	0	0	36
2	0,02	0,02	0,02	65
3	0,12	0,02	0,04	91
4	0,1	0,04	0,06	181
5	0,13	0,07	0,09	213
6	0,14	0,09	0,11	248
7	0,33	0,29	0,31	383
weighted AVG	0,23	0,27	0,23	1980

5. CONCLUSION

We presented EGO-CH, a dataset for visitors behavioral understanding using egocentric vision. The dataset includes more than 27 hours of video, 70 visits acquired by real visitors, 26 environments and over 200 different points of interest related to two different cultural sites. We publicly release the dataset along with temporal labels for locations and observed points of interest, bounding box annotations for objects, and surveys associated to 60 visits. Baseline results on the challenging tasks of Room-based Localization, Point of Interest/Object Recognition, Object Retrieval and Survey Prediction show the potential of the dataset for visitors behavioral understanding. We believe that EGO-CH can be a valuable benchmark to tackle the proposed tasks, as well as others not investigated in this paper. Future works can address the evaluation considering more advanced baselines and investigate specialized approach to the four proposed tasks.

Acknowledgment

This research is part of the project VALUE - Visual Analysis for Localization and Understanding of Environments (N. 08CT6209090207, CUP G69J18001060007) supported by PO FESR 2014/2020 - Azione 1.1.5. - “Sostegno

all’avanzamento tecnologico delle imprese attraverso il finanziamento di linee pilota e azioni di validazione precoce dei prodotti e di dimostrazioni su larga scala” del PO FESR Sicilia 2014/2020, and Piano della Ricerca 2016-2018 linea di Intervento 2 of DMI, University of Catania. The authors would like to thank Regione Siciliana Assessorato dei Beni Culturali dell’Identit Siciliana - Dipartimento dei Beni Culturali e dell’Identit Siciliana and Polo regionale di Siracusa per i siti culturali - Galleria Regionale di Palazzo Bellomo.

References

- [1] A. Bollo, L. Pozzolo, Analysis of visitor behaviour inside the museum: An empirical study, in: Arts & Cultural Management, 2005.
- [2] R. Cucchiara, A. Del Bimbo, Visions for augmented cultural heritage experience, IEEE MultiMedia 21 (1) (2014) 74–82.
- [3] F. Ragusa, A. Furnari, S. Battiato, G. Signorello, G. M. Farinella, Egocentric visitors localization in cultural sites, J. Comput. Cult. Herit. 12 (2) (2019) 11:1–11:19.
- [4] F. Ragusa, L. Guarnera, A. Furnari, S. Battiato, G. Signorello, G. M. Farinella, Localization of visitors for cultural sites management, in: International Joint Conference on e-Business and Telecommunications - Volume 2: ICETE, 2018, pp. 407–413.
- [5] F. Ragusa, A. Furnari, S. Battiato, G. Signorello, G. M. Farinella, Egocentric point of interest recognition in cultural sites, in: VISAPP, 2019.
- [6] A. S. Razavian, O. Aghazadeh, J. Sullivan, S. Carlsson, Estimating attention in exhibitions using wearable cameras, ICPR (2014) 2691–2696.

- [7] L. Seidenari, C. Baecchi, T. Uricchio, A. Ferracani, M. Bertini, A. D. Bimbo, Deep artwork detection and retrieval for automatic context-aware audio guides, TOMM 13 (3s) (2017) 35.
- [8] D. Raptis, N. K. Tselios, N. M. Avouris, Context-based design of mobile applications for museums: a survey of existing practices, in: Mobile HCI, 2005.
- [9] P. Koniusz, Y. Tas, H. Zhang, M. Harandi, F. Porikli, R. Zhang, Museum exhibit identification challenge for domain adaptation and beyond (2018). [arXiv:1802.01093](https://arxiv.org/abs/1802.01093).
- [10] R. D. Chiaro, A. Bagdanov, A. D. Bimbo, {NoisyArt}: A Dataset for Webly-supervised Artwork Recognition, in: International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications, 2019.
- [11] D. Ahmetovic, C. Gleason, K. M. Kitani, H. Takagi, C. Asakawa, Navcog: Turn-by-turn smartphone navigation assistant for people with visual impairments or blindness, in: Web for All Conference, W4A '16, 2016, pp. 9:1–9:2.
- [12] A. Kendall, M. Grimes, R. Cipolla, Posenet: A convolutional network for real-time 6-dof camera relocation, in: ICCV, 2015, pp. 2938–2946.
- [13] G. Taverriti, S. Lombini, L. Seidenari, M. Bertini, A. Del Bimbo, Real-time wearable computer vision system for improved museum experience, in: ACM Multimedia, 2016, pp. 703–704.
- [14] J. Redmon, A. Farhadi, Yolov3: An incremental improvement, CoRR abs/1804.02767 (2018). [arXiv:1804.02767](https://arxiv.org/abs/1804.02767).
- [15] A. Rubhasy, A. A. G. Y. Paramartha, I. Budi, Z. A. Hasibuan, Management and retrieval of cultural heritage multimedia collection using ontology, International Conference on Information Technology, Computer, and Electrical Engineering (2014) 255–259.
- [16] P. Kwan, K. Kameyama, J. Gao, K. Toraichi, Content-based image retrieval of cultural heritage symbols by interaction of visual perspectives., IJPRAI 25 (2011) 643–673.
- [17] D. K. Iakovidis, E. E. Kotsifakos, N. Pelekis, H. Karanikas, I. Kopanakis, T. Mavroudakos, Y. Theodoridis, Pattern-based retrieval of cultural heritage images, 2007.
- [18] K. Makantasis, A. Doulamis, N. Doulamis, M. Ioannides, In the wild image retrieval and clustering for 3d cultural heritage landmarks reconstruction, Multimedia Tools and Applications 75 (7) (2016) 3593–3629.
- [19] A. Furnari, S. Battiato, G. M. Farinella, Personal-location-based temporal segmentation of egocentric videos for lifelogging applications, Journal of Visual Communication and Image Representation 52 (2018) 1–12.
- [20] Galleria regionale di palazzo bellomo (2007).
URL <http://www.regione.sicilia.it/beniculturali/palazzobellomo/>
- [21] Monastero dei benedettini.
URL <http://www.monasterodeibenedettini.it/>
- [22] G. Huang, Z. Liu, L. van der Maaten, K. Q. Weinberger, Densely connected convolutional networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017.

SUPPLEMENTARY MATERIAL

This supplementary material complements the submitted paper. It reports additional details on the EGO-CH dataset and experiments.

5.1. INTRODUCTION

This document is intended for the convenience of the reader and reports additional information about the proposed dataset and the performed experiments. This supplementary material is related to the following submission:

- F. Ragusa, A. Furnari, S. Battiato, G. Signorello, G. M. Farinella, “EGO-CH: Dataset and Challenges for Visitors Behavioral Understanding from Egocentric Vision”, Pattern Recognition Letters, DOI: 10.1016/j.patrec.2019.12.016

The reader is referred to the manuscript and to our web page <http://iplab.dmi.unict.it/EGO-CH/> for further information.

5.2. THE EGO-CH DATASET

5.2.1. Data Collection

The dataset has been acquired using a head-mounted Microsoft HoloLens device in two cultural sites located in Sicily, Italy: 1) “Palazzo Bellomo”, located in Siracusa [20], and 2) “Monastero dei Benedettini”, located in Catania [21].

1) *Palazzo Bellomo*. This cultural site is composed by 22 environments and contains 191 Points of Interest (e.g., statues, paintings, etc.). Figure 9 and Figure 10 report some frames related to the different environments and some points of interest. Table 14 details the list of the acquired training video. Some of the videos are related to the 22 rooms of the cultural site, whereas other are related to specific points of interest. For each video we report its total duration, the amount of required storage, the number of frames, as well as the percentage of frames with respect to the whole training set.

Table 15 reports the list of the 10 test videos acquired by volunteers visiting the cultural site. For each video, we report its total duration, the amount of required storage, the number of frames, the number of environments encountered in the video, as well as the sequence of environments, as visited by the subject acquiring the video. All the training/test videos have a resolution of 1280×720 pixels

Table 8: FF_1 Results for Room- Based Localization on Palazzo Bellomo.

FF_1 score										
Class	Test1	Test2	Test4	Test5	Test6	Test7	Test8	Test9	Test10	AVG
1_Sala1	0,00	0,96	0,96	0,95	0,99	0,94	0,97	0,97	0,91	0,85
2_Sala2	0,82	0,91	0,96	0,95	0,99	0,99	0,99	0,97	0,95	0,95
3_Sala3	0,65	0,84	0,84	0,83	0,94	0,88	0,61	0,85	0,83	0,81
4_Sala4	0,84	0,96	0,98	0,93	/	0,96	0,78	0,97	0,93	0,92
5_Sala5	0,98	0,88	0,96	0,91	0,65	0,91	0,97	0,97	0,99	0,91
6_Sala6	0,84	0,80	0,43	0,00	0,86	0,00	0,18	0,83	0,73	0,52
7_Sala7	0,67	0,93	0,22	0,00	0,38	0,88	0,88	0,40	0,96	0,59
8_Sala8	0,64	0,75	0,60	0,52	0,56	0,56	0,64	0,60	0,77	0,63
9_Sala9	0,90	0,89	0,73	0,88	0,41	0,92	0,98	0,99	0,88	0,84
10_Sala10	0,96	0,98	0,92	0,72	0,00	0,98	0,85	0,97	0,78	0,80
11_Sala11	0,93	0,96	0,96	0,97	0,00	0,97	0,97	0,98	0,95	0,85
12_Sala12	0,82	0,88	0,87	0,90	0,00	0,87	0,85	0,92	0,87	0,78
13_Sala13	0,96	0,95	0,95	0,76	0,00	0,93	0,95	0,94	0,96	0,82
14_CortiledegliStemmi	0,78	0,68	0,74	0,88	0,17	0,90	0,79	0,92	0,00	0,65
15_SalaCarrozze	0,84	0,89	0,95	0,93	0,91	0,84	0,97	0,95	/	0,91
16_CortileParisio	0,33	0,51	0,52	0,45	0,35	0,60	0,81	0,91	0,80	0,59
17_Biglietteria	0,25	0,56	0,00	0,26	0,00	0,00	0,00	0,36	0,00	0,16
18_Portico	0,68	0,67	0,78	0,54	0,69	0,69	0,62	0,78	0,62	0,67
19_ScalaCatalana	0,00	0,00	0,54	0,55	0,78	0,67	0,58	0,85	0,49	0,50
20_Loggetta	0,00	0,50	0,80	0,44	0,72	0,53	0,55	0,83	0,68	0,56
21_BoxSala8	0,96	0,97	0,97	/	0,00	0,97	0,99	0,81	0,94	0,83
22_AreaSosta	0,74	0,85	0,45	0,83	0,00	0,56	0,72	0,85	0,00	0,56
mFF_1	0,66	0,79	0,73	0,68	0,45	0,75	0,76	0,85	0,72	0,71

Table 9: ASF_1 Results for Room- Based Localization on Palazzo Bellomo.

ASF_1 score										
Class	Test1	Test2	Test4	Test5	Test6	Test7	Test8	Test9	Test10	AVG
1_Sala1	0,00	0,93	0,93	0,91	0,97	0,89	0,94	0,94	0,84	0,82
2_Sala2	0,64	0,83	0,91	0,90	0,97	0,97	0,97	0,94	0,90	0,89
3_Sala3	0,36	0,53	0,51	0,53	0,61	0,63	0,29	0,59	0,72	0,53
4_Sala4	0,39	0,92	0,96	0,86	/	0,92	0,31	0,94	0,87	0,77
5_Sala5	0,94	0,64	0,66	0,62	0,34	0,63	0,54	0,95	0,97	0,70
6_Sala6	0,46	0,43	0,18	0,00	0,42	0,00	0,09	0,59	0,51	0,30
7_Sala7	0,50	0,86	0,13	0,00	0,10	0,78	0,28	0,25	0,92	0,42
8_Sala8	0,36	0,44	0,42	0,35	0,35	0,41	0,49	0,53	0,57	0,44
9_Sala9	0,81	0,80	0,51	0,78	0,22	0,84	0,95	0,97	0,78	0,74
10_Sala10	0,92	0,96	0,85	0,19	0,00	0,95	0,27	0,94	0,64	0,64
11_Sala11	0,53	0,92	0,91	0,93	0,00	0,93	0,65	0,95	0,91	0,75
12_Sala12	0,38	0,79	0,76	0,81	0,00	0,77	0,30	0,85	0,76	0,60
13_Sala13	0,93	0,90	0,90	0,62	0,00	0,87	0,91	0,89	0,92	0,77
14_CortiledegliStemmi	0,57	0,53	0,54	0,77	0,08	0,81	0,49	0,67	0,00	0,50
15_SalaCarrozze	0,72	0,80	0,91	0,87	0,83	0,72	0,94	0,90	/	0,84
16_CortileParisio	0,30	0,63	0,43	0,65	0,38	0,47	0,59	0,83	0,66	0,55
17_Biglietteria	0,26	0,43	0,00	0,11	0,00	0,00	0,00	0,39	0,00	0,13
18_Portico	0,51	0,47	0,60	0,45	0,48	0,52	0,42	0,64	0,47	0,51
19_ScalaCatalana	0,00	0,00	0,43	0,51	0,61	0,55	0,42	0,75	0,39	0,41
20_Loggetta	0,00	0,31	0,54	0,28	0,46	0,38	0,48	0,72	0,44	0,40
21_BoxSala8	0,92	0,95	0,95	/	0,00	0,93	0,98	0,67	0,88	0,79
22_AreaSosta	0,59	0,73	0,29	0,71	0,00	0,75	0,57	0,74	0,00	0,49
$mASF_1$	0,50	0,67	0,61	0,56	0,32	0,67	0,54	0,76	0,63	0,58

Table 10: FF_1 Results for Room- Based Localization on Monastero dei Benedettini.

ID_visit	1	2	3	4	Neg.	AVG	ID_visit	1	2	3	4	Neg.	AVG
4805	0,79	0,82	0,74	0,95	0,45	0,75	2043	0,52	0,44	0,62	0,91	0,33	0,564
1804	0,32	0,36	0,79	0,98	0,55	0,6	3996	0,17	0,68	0,57	0,95	0,72	0,618
4377	0,46	0,53	0,8	0,96	0,32	0,614	3455	0,78	0,65	0,8	0,99	0,55	0,754
1669	0,75	0,72	0,9	0,99	0,45	0,762	4785	0,02	0	0,64	0,95	0,33	0,388
1791	0,49	0,5	0,84	0,95	0,43	0,642	2047	0,95	0,84	0,84	0,98	0,38	0,798
3948	0	/	0,59	0,98	0,49	0,515	1912	0,66	0,51	0,69	0,94	0,46	0,652
3152	0,39	0,72	0,87	0,97	0,76	0,742	3232	0,66	0,53	0,85	0,99	0,49	0,704
4361	0,55	0,78	0,82	0,97	0,28	0,68	4442	0,8	0,77	0,82	0,98	0,23	0,72
3976	0,87	0,43	0,81	0,93	0,66	0,74	3646	0,63	0,29	0,66	0,97	0,4	0,59
3527	0,86	0,79	0,85	0,99	0,56	0,81	4833	0,67	0,82	0,67	0,88	0,24	0,656
4105	0,62	0	0,77	0,97	0,09	0,49	3478	0,71	0,75	0,84	0,98	0,37	0,73
1399	0,46	0,23	0,79	0,98	0,43	0,578	4396	0,74	0,69	0,89	0,96	0,26	0,708
3836	0,51	0	0,72	0,99	0,55	0,554	2894	0,83	0,8	0,66	0,92	0,48	0,738
4006	0,57	0,78	0,78	0,99	0,37	0,698	4414	0,75	0,7	0,91	0,95	0,22	0,706
4415	0,86	0,61	0,82	0,97	0,35	0,722	4639	0,57	0,52	0,85	0,99	0,39	0,664
3008	0,69	0	0,7	0,99	0,51	0,578	1004	0,08	0,72	0,49	0,97	0,26	0,504
4660	0,62	0,81	0,83	0,98	0,57	0,762	1917	0,41	0,71	0,36	0,87	0,29	0,528
2826	0,31	0,52	0,76	0,97	0,61	0,634	1153	0,62	0,7	0,71	0,91	0,32	0,652
1099	0,62	0,42	0,85	0,98	0,49	0,672	2244	0,94	0,56	0,74	0,99	0,43	0,732
4391	0,74	0,72	0,72	0,98	0,33	0,698	2614	0,88	0,56	0,83	0,99	0,4	0,732
3929	0,26	0	0,8	0,99	0,43	0,496	1624	0,33	0,8	0,62	0,99	0,33	0,614
3362	0,34	0,68	0,68	0,95	0,21	0,572	3441	0,61	0,25	0,84	0,99	0,41	0,62
1379	0,41	0	0,81	0,96	0,45	0,526	4793	0,52	/	0,68	0,99	0,33	0,63
2600	0,26	0	0,63	0,96	0,3	0,43	4083	0,93	/	0,73	0,99	0,71	0,84
1430	0,94	0	0,69	0,87	0,58	0,616	4906	0,46	0,36	0,59	0,94	0,31	0,532
2956	0,32	0,45	0,33	0,91	0,65	0,532	1160	0,88	0,84	0,72	0,98	0,46	0,776
4742	0,13	0,58	0,59	0,84	0,5	0,528	3416	0,56	/	0,5	0,82	0,2	0,52
3651	0,77	0,41	0,88	0,98	0,41	0,69	1051	0,78	0,76	0,64	0,95	0,57	0,74
1064	0,77	0,23	0,83	0,99	0,22	0,608	2580	0,71	0,45	0,63	0,96	0,44	0,638
3818	0,68	0,64	0,73	0,99	0,47	0,702	1109	0,91	0,28	0,85	0,97	0,36	0,674
							mFF1	0,59	0,51	0,73	0,96	0,42	0,64

Table 11: ASF_1 Results for Room- Based Localization on Monastero dei Benedettini.

ID_visit	1	2	3	4	Neg.	AVG	ID_visit	1	2	3	4	Neg.	AVG
4805	0,55	0,62	0,4	0,12	0,36	0,41	2043	0,4	0,27	0,28	0,13	0,27	0,27
1804	0,32	0,1	0,39	0,05	0,21	0,214	3996	0,23	0,3	0,22	0,05	0,38	0,236
4377	0,31	0,12	0,26	0,08	0,27	0,208	3455	0,63	0,38	0,65	0,99	0,48	0,626
1669	0,68	0,34	0,67	0,9	0,33	0,584	4785	0,06	0	0,25	0,45	0,07	0,166
1791	0,41	0,41	0,25	0,06	0,22	0,27	2047	0,9	0,72	0,72	0,37	0,34	0,61
3948	0	/	0,33	0,26	0,49	0,27	1912	0,42	0,25	0,2	0,08	0,3	0,25
3152	0,39	0,43	0,64	0,44	0,58	0,496	3232	0,52	0,1	0,34	0,13	0,35	0,288
4361	0,26	0,35	0,32	0,08	0,2	0,242	4442	0,61	0,22	0,33	0,06	0,12	0,268
3976	0,68	0,26	0,35	0,03	0,5	0,364	3646	0,62	0,3	0,26	0,19	0,22	0,318
3527	0,65	0,37	0,38	0,28	0,44	0,424	4833	0,54	0,55	0,26	0,1	0,21	0,332
4105	0,49	0	0,49	0,16	0,06	0,24	3478	0,57	0,4	0,47	0,14	0,24	0,364
1399	0,55	0,1	0,33	0,23	0,29	0,3	4396	0,41	0,1	0,2	0,07	0,12	0,18
3836	0,47	0	0,35	0,98	0,34	0,428	2894	0,55	0,67	0,36	0,13	0,31	0,404
4006	0,47	0,49	0,49	0,55	0,26	0,452	4414	0,59	0,26	0,83	0,11	0,08	0,374
4415	0,74	0,07	0,35	0,14	0,23	0,306	4639	0,5	0,35	0,67	0,39	0,27	0,436
3008	0,57	0	0,5	0,99	0,41	0,494	1004	0,22	0,47	0,17	0,11	0,17	0,228
4660	0,45	0,52	0,17	0,05	0,27	0,292	1917	0,39	0,19	0,03	0,03	0,21	0,17
2826	0,31	0,29	0,33	0,06	0,25	0,248	1153	0,47	0,33	0,21	0,14	0,16	0,262
1099	0,52	0,24	0,55	0,96	0,17	0,488	2244	0,77	0,35	0,49	0,35	0,32	0,456
4391	0,54	0,13	0,14	0,2	0,24	0,25	2614	0,43	0,06	0,39	0,99	0,26	0,426
3929	0,15	0	0,31	0,22	0,39	0,214	1624	0,3	0,33	0,24	0,29	0,29	0,29
3362	0,2	0,35	0,38	0,05	0,22	0,24	3441	0,49	0,11	0,44	0,88	0,42	0,468
1379	0,22	0	0,43	0,05	0,35	0,21	4793	0,35	/	0,36	0,28	0,37	0,34
2600	0,32	0	0,39	0,31	0,26	0,256	4083	0,88	/	0,41	0,25	0,48	0,505
1430	0,43	0	0,1	0,02	0,21	0,152	4906	0,36	0,1	0,31	0,06	0,21	0,208
2956	0,41	0,2	0,16	0,29	0,31	0,274	1160	0,79	0,5	0,39	0,14	0,28	0,42
4742	0,15	0,18	0,25	0,11	0,17	0,172	3416	0,32	/	0,25	0,42	0,34	0,3325
3651	0,32	0,25	0,74	0,48	0,28	0,414	1051	0,5	0,31	0,21	0,05	0,4	0,294
1064	0,63	0,09	0,52	0,97	0,19	0,48	2580	0,36	0,29	0,33	0,05	0,28	0,262
3818	0,43	0,49	0,2	0,06	0,29	0,294	1109	0,82	0,17	0,46	0,11	0,28	0,368
							mASF1	0,46	0,26	0,37	0,28	0,28	0,40

Table 12: Point of Interest Retrieval related to Palazzo Bellomo cultural site.

1) Palazzo Bellomo				
Variant	K	Precision	Recall	F_1 score
1 - One Shot	1	0.02	0.01	0.00
	3	0.62	0.59	0.6
	5	0.62	0.56	0.56
2 - Many Shots	3	0.62	0.56	0.56
	5	0.62	0.56	0.56
	7	0,61	0,56	0,56
	9	0,61	0,55	0,56
11	0,61	0,55	0,55	

Table 13: Point of Interest Retrieval related to Monastero dei Benedettini cultural site.

2) Monastero dei Benedettini				
Variant	K	Precision	Recall	F_1 score
1 - One shot	1	0.38	0.07	0.09
	3	0,83	0,83	0,83
	5	0,84	0,83	0,83
2 - Many Shots	3	0,84	0,83	0,83
	5	0,84	0,84	0,83
	7	0,84	0,83	0,83
	9	0,84	0,83	0,83
11	0,83	0,83	0,82	



Figure 9: Sample frames for each of the 22 considered environments of “Palazzo Bellomo”.

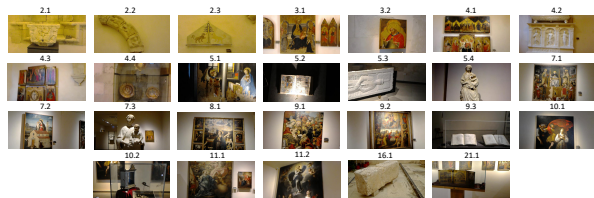


Figure 10: Sample frames of points of interest of “Palazzo Bellomo”.



Figure 11: Sample reference images related to the cultural site “Palazzo Bellomo”.

and a frame-rate of 29.97 fps. We also include 191 reference images related to the considered POIs to be used for one-shot image retrieval. The images are akin to pictures generally included in museums catalog. Figure 11 shows some examples of such reference images.

2) *Monastero dei Benedettini*. This dataset is composed by 4 environments and contains 35 Points of Interest. Figure 12 and Figure 13 report some frames related to the 4 different environments and some of the points of interest. Table 16 reports details on the acquired training videos, highlighting the total duration of the videos, the required storage, the number of frames and the percentage of frames with respect to the whole training set. Training and validation videos have a resolution of 1216×684 pixels and a frame-rate of 24 fps. Five validation videos have been collected by asking volunteers to visit the cultural site with rules similarly to the one used for “Palazzo Bellomo”. Table 17 shows the number of frames belonging to each video (left) and the number of frames belonging for each class (right).

Additionally, we collected 60 test videos by asking real visitors to freely visit the cultural site. No specific instructions have been given to the visitors, who were free to explore the 4 environments and the 35 POIs. The 60 test videos have a resolution of 1408×792 pixels and a frame-rate of 30.03 fps. The average number of frames for each video is 39296. We also include 35 reference images related to the considered POIs to be used for one-shot image retrieval. Figure 14 shows some example of reference images.

Table 14: List of training videos of “Palazzo Bellomo”.

Name	Time (s)	Storage (MB)	#frame	%frame
1.0.Sala1	124	156.229	3721	3,13%
2.0.Sala2	117	148.480	3525	2,96%
3.0.Sala3	100	125.924	3000	2,52%
3.0.Sala3_S	73	92.589	2200	1,85%
4.0.Sala4	97	122.941	2914	2,45%
5.0.Sala5	99	126.213	2992	2,51%
6.0.Sala6	87	110.451	2630	2,21%
7.0.Sala7	113	143.257	3402	2,86%
8.0.Sala8	147	186.470	4427	3,72%
9.0.Sala9	143	180.971	4298	3,61%
10.0.Sala10	71	90.697	2154	1,81%
11.0.Sala11	104	131.983	3145	2,64%
12.0.Sala12	82	103.785	2463	2,07%
13.0.Sala13	101	128.013	3040	2,55%
14.0.CortiledegliStemmi	104	131.962	3131	2,63%
14.0.CortiledegliStemmi_S	90	113.822	2722	2,29%
15.0.SalaCarrozze	108	136.968	3259	3,12%
16.0.CortileParisio	124	156.605	3718	3,12%
16.0.Cortile.Parisio_S	68	86.646	2062	1,73%
17.0.Biglietteria	83	104.532	2489	2,09%
17.0.Biglietteria_S	53	68.071	1610	1,35%
18.0.Portico	126	159.800	3791	3,18%
18.0.Portico_S	63	80.044	1910	1,60%
19.0.ScalaCatalana	97	123.010	2918	2,45%
19.0.ScalaCatalana_S	116	110.063	3481	2,92%
20.0.Loggetta	80	101.584	2425	2,04%
20.0.Loggetta_S	58	73.722	1744	1,46%
21.0.BoxSala8	85	107.540	2562	2,15%
22.0.AreaSosta	64	81.934	1945	1,63%
22.0.Area.Sosta_S	52	65.340	1560	1,31%
2.1.Sala2_Acquasantiera	54	68.445	1623	2,94%
2.2.Sala2_FrammentiArchitett.	46	58.265	1393	2,53%
2.3.Sala2_LastraconLeoni	47	60.199	1427	2,59%
3.1.Sala3_MadonnainTrono	65	83.198	1972	3,58%
3.2.Sala3_FrammentoS.Leonardo	37	47.061	1113	2,02%
4.1.Sala4_MadonnainTrono	75	94.767	2252	4,08%
4.2.Sala4_MonumentoE.d'Aragona	86	108.296	2580	4,68%
4.3.Sala4_TrasfigurazioneCristo	76	96.106	2277	4,13%
4.4.Sala4_Piatti	49	62.281	1474	2,67%
5.1.Sala5_Annunciazione	76	96.952	2295	4,16%
5.2.Sala5_LibroD'OreMiniato	46	59.011	1406	2,55%
5.3.Sala5_LastraG.Cabastida	100	127.023	3017	5,47%
5.4.Sala5_MadonnadelCardillo	61	77.568	1829	3,32%
7.1.Sala7_DisputaS.Tommaso	74	94.188	2234	4,05%
7.2.Sala7_TraslazioneSantaCasa	76	96.045	2281	4,14%
7.3.Sala7_MadonnacolBambino	90	113.202	2696	4,89%
8.1.Sala8_ImmacolataConcezione	82	104.570	2483	4,50%
9.1.Sala9_AdorazionedeiMagi	60	76.171	1803	3,27%
9.2.Sala9_S.ElenaCostantinoeMadonna	76	96.227	2283	4,14%
9.3.Sala9_TaccuinidiDisegni	70	89.647	2121	3,85%
10.1.Sala10_MartirioS.Lucia	58	74.248	1759	3,19%
10.2.Sala10_VoltodiCristo	64	80.896	1917	3,48%
11.1.Sala11_MiracolodiS.Orsola	66	84.297	2002	3,63%
11.2.Sala11_Immacolata	73	92.424	2196	3,98%
16.1.CortileParisio_LapidiEbraiche	85	108.098	2563	4,65%
16.1.CortileParisio_LapidiEbraiche_S	67	85.173	2031	3,68%
21.1.BoxSala8_StoriedellaGenesi	70	88.300	2099	3,81%
AVG	81.72	103.02	2462.53	2.07%

Table 15: List of test videos of “Palazzo Bellomo”.

Name	Time (s)	MB	#Frame	%Frame	#Environments	Environments - Temporal sequence
Test1	1906	2.400.360	57123	11,13%	22	16->17->18->1->18->3->2->3->18->4->18->14->15->14->19->20->6->5->6->7->21->8->9->22->10->11->12->13->5->6->20->19->18->17
Test2	1413	1.435.096	42348	8,25%	22	16->17->18->1->18->3->2->3->18->4->18->14->15->14->19->20->6->5->6->7->8->21->8->9->22->10->11->12->13->5->6->20->19->18->17
Test3	1830	2.304.410	54845	10,69%	22	16->17->18->1->18->3->2->3->18->4->18->14->15->14->19->20->6->5->13->12->11->10->22->9->8->21->8->7->6->20->19->18->17
Test4	1542	1.942.200	46214	5,49%	22	16->17->18->1->18->3->2->3->18->4->18->14->15->14->19->20->6->5->6->7->8->21->8->9->22->10->11->12->13->5->6->20->19->18->17
Test5	1034	1.302.612	30989	9,00%	22	16->17->18->1->18->3->2->3->18->4->18->14->15->14->19->20->6->5->6->7->8->9->22->10->11->12->13->5->6->20->19->18->17
Test6	1949	2.273.926	58411	11,38%	22	16->17->8->1->18->3->2->3->18->14->15->14->19->20->6->5->13->5->6->7->8->21->8->9->22->10->11->12->11->10->22->9->8->7->6->20->19->18->17
Test7	1332	1.677.047	39920	7,78%	22	16->17->14->15->14->18->1->18->3->2->3->18->4->18->19->20->6->5->6->7->8->21->8->9->22->10->11->12->13->5->6->20->19->18->17->16
Test8	3023	3.806.383	90599	17,65%	22	16->17->14->5->14->18->4->18->3->2->3->18->1->18->19->20->6->5->13->12->11->10->22->9->8->21->8->7->6->20->19->14
Test9	2236	2.815.878	67013	13,05%	22	16->17->18->1->18->3->2->3->18->4->18->14->15->14->19->20->6->5->13->12->11->10->22->9->8->21->8->7->6->20->19->14->17
Test10	858	1.080.389	25714	5,01%	22	16->17->14->19->20->6->7->8->21->8->9->22->10->11->12->13->5->6->20->19->18->4->18->1->18->3->2->3->18



Figure 12: Sample frames from the 4 considered environments of “Monastero dei Benedettini”.

5.2.2. Annotations

Temporal Labels. All test and validation videos have been temporally labeled to indicate in every frame the environment in which the visitor is located and the currently observed point of interest, if any. If the visitor is not located in any of the considered environments or they are not observing any of the considered POIs we mark as that frame as “negative” (Figure 15).

Bounding Box Annotations. A subset of frames from the dataset (sampled at at 1 fps) has been labeled with bounding boxes indicating the presence and locations of all POIs. Figure 13 shows some example of labeled frames from the training set of the “Monastero dei Benedettini”.

The EGO-CH dataset is publicly available at our

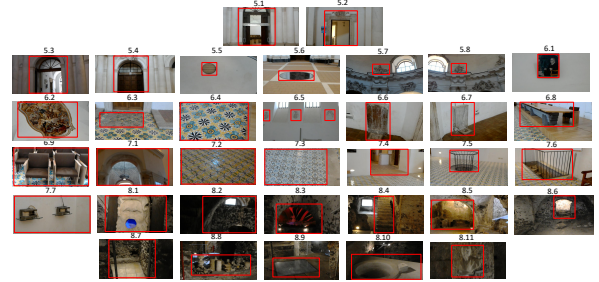


Figure 13: Sample frames from the 35 considered POIs of “Monastero dei Benedettini”, with the related bounding box annotations.

website: <http://iplab.dmi.unict.it/EGO-CH/>. The dataset can be used only for research purposes and is available upon request to the authors.

5.3. EXPERIMENTAL DETAILS

In this section, we present additional details on baseline experiments related to the 4 tasks proposed in the paper.

5.3.1. Room-based Localization

1) Palazzo Bellomo. We split the Training Set into two subsets to train and validate the VGG-19 for the “Discrimination” stage (no “negative” frames are used for training). Table 18 reports the number of frames belonging to the two subsets for each of the 22 considered environments.

We report the results obtained by the baseline on the 9 test videos (“Test3” has been used for validation) in Table 19 considering the *FF1* score metric,

Table 16: List of training videos of “Monastero dei Benedettini”

Name	Time (s)	Storage (MB)	#frame	%frame
5.0_Antirefettorio	247	244.500	5933	1,21%
5.0_Antirefettorio1	263	260.395	6315	1,29%
6.0_SantoMazzarino	241	239.007	5800	1,18%
7.0_Cucina	239	237.268	5753	1,17%
8.0_Ventre	679	832.844	20385	4,15%
5.1_Antirefettorio_PortaA.S.Mazz.Ap.	67	67.001	1630	0,33%
5.1_Antirefettorio_PortaA.S.Mazz.Ch.	74	73.589	1785	0,36%
5.2_Antirefettorio_PortaMuseoFab.Ap.	50	49.840	1211	0,25%
5.2_Antirefettorio_PortaMuseoFab.Ch.	62	61.846	1503	0,31%
5.3_Antirefettorio_PortaAntiref.	51	58.557	1537	0,31%
5.4_Antirefettorio_PortaRef.Piccolo	54	53.767	1306	0,27%
5.5_Antirefettorio_Cupola	57	56.089	1377	0,28%
5.6_Antirefettorio_AperturaPavimento	55	54.586	1322	0,27%
5.7_Antirefettorio_S.Agata	48	47.820	1165	0,24%
5.8_Antirefettorio_S.Scolastica	58	57.919	1407	0,29%
5.9_Antirefettorio_ArcoconFirma	62	76.700	1864	0,38%
5.10_Antirefettorio_BustoVaccarini	65	64.298	1563	0,32%
6.1_SantoMazzarino_QuadroS.Mazz.	71	70.401	1716	0,35%
6.2_SantoMazzarino_Affresco	213	211.424	5124	1,04%
6.3_SantoMazzarino_PavimentoOr.	99	98.691	2397	0,49%
6.4_SantoMazzarino_PavimentoRes.	69	69.148	1675	0,34%
6.5_SantoMazzarino_BassorilieviManc.	117	115.348	2823	0,57%
6.6_SantoMazzarino_LavamaniSx	151	149.882	3637	0,74%
6.7_SantoMazzarino_LavamaniDx	93	92.928	2256	0,46%
6.8_SantoMazzarino_TavoloRelatori	150	148.661	3603	0,73%
6.9_SantoMazzarino_Poltrone	108	107.374	2604	0,53%
7.1_Cucina_Edicola	369	437.219	11086	2,26%
7.2_Cucina_PavimentoA	52	52.163	1268	0,26%
7.3_Cucina_PavimentoB	52	52.244	1266	0,26%
7.4_Cucina_PassavivandePavim.Orig.	81	80.733	1961	0,40%
7.5_Cucina_AperturaPav.1	57	57.170	1385	0,28%
7.5_Cucina_AperturaPav.2	53	52.875	1280	0,26%
7.5_Cucina_AperturaPav.3	62	62.320	1509	0,31%
7.6_Cucina_Scala	77	76.587	1856	0,38%
7.7_Cucina_SalaMetereologica	156	154.394	3748	0,76%
8.1_Ventre_Doccione	103	102.683	2492	0,51%
8.2_Ventre_VanoRacc.Cenere	126	124.837	3026	0,62%
8.3_Ventre_SalaRossa	300	296.792	7202	1,47%
8.4_Ventre_ScalaCucina	214	212.372	5152	1,05%
8.5_Ventre_CucinaProv.	148	146.379	3553	0,72%
8.6_Ventre_Ghiacciaia	69	68.562	1668	0,34%
8.6_Ventre_Ghiacciaia1	266	263.817	6398	1,30%
8.7_Ventre_Latrina	102	100.542	2468	0,50%
8.8_Ventre_OssaScarti	154	152.861	3713	0,76%
8.8_Ventre_OssaScarti1	36	36.345	886	0,18%
8.9_Ventre_Pozzo	300	297.167	7209	1,47%
8.10_Ventre_Cisterna	57	56.572	1384	0,28%
8.11_Ventre_BustoP.Tacchini	110	109.520	2662	0,54%
AVG	133.06	137.38	3351.31	1.04%

Table 17: List of validation videos of “Monastero dei Benedettini”.

Name	#frame	Class	#frame
Test1	4141	1 Antirefettorio	88613
Test3	18678	2 Aula S. Mazzarino	8395
Test4	13731	3 Cucina	9712
Test5	15958	4 Ventre	20513
Test7	1124	Negatives	6399
Total	53632	Total	53632



Figure 14: Sample references images related to the “Monastero dei Benedettini”.

and in in Table 20 considering the ASF_1 score. As example, Figure 16 illustrate qualitatively the segmentation result of the baseline on “Test7”. Figure 17 reports the confusion matrix of the baseline on the test set.

2) *Monastero dei Benedettini*. Similarly to “Palazzo Bellomo”, we split the Training Set into two subsets to train and validate the VGG-19 (no “negative” frames are used for training on this cultural site). Table 23 reports the number of frames belonging to the two subsets for each of the 4 considered environments. We report the results obtained by the baseline method over the 60 test videos in Table 24 and Table 25 considering the FF_1 score metric, and in Table 26 and Table 27 considering the ASF_1 score.

DenseNet Backbone. We performed experiments using another backbone in the same pipeline to solve the first task. We used DenseNet[22], a densely convolutional connect network which connects each layer to every other layer in a feed-



Figure 15: Sample frames from “Monastero dei Benedettini” marked as “negative locations”.

forward fashion. Table 21 and Table 22 report the results obtained with DenseNet at the end of the Sequential Modeling step. We have evaluated the model using F_1 score and Asf_1 score. As shown the tables we did not obtain an improvement respect the results obtained with the backbone VGG. In particular for the “Palazzo Bellomo” we obtained a FF_1 score of 0.71 and a ASF_1 score of 0.58 which are lower respect the scores obtained with VGG ($FF_1 = 0.82$, $ASF_1 = 0.59$). From Table 24 to Table 27 we report the results obtained in the “Monastero dei Benedettini” using the backbone DenseNet. As shown, neither considering this cultural site we obtained an improvement of FF_1 and ASF_1 scores.

5.3.2. Points of Interest Recognition

1) *Palazzo Bellomo*. We used the subset of the Training set considered in the first task annotated with bounding box for a total of 56686 frames. We split this subset in training/validation sets to train and validate the object detector. In particular, we used 41111 frames as training set and 15575 as validation set. To find the optimal detection threshold, we used “Test3” as validation video. We obtaining the value of 0.05 maximizing mAP over the valida-

Table 18: Number of frames belonging to the two subsets (Training/Validation) to train the CNN for “Palazzo Bellomo”.

	Training	Validation
1 Sala1	2605	1116
2 Sala2	5578	2390
3 Sala3	5800	2486
4 Sala4	8048	3449
5 Sala5	8023	3438
6 Sala6	1841	789
7 Sala7	7429	3184
8 Sala8	4837	2073
9 Sala9	7354	3152
10 Sala10	4081	1749
11 Sala11	5140	2203
12 Sala12	1724	739
13 Sala13	2128	912
14 Cortile degli Stemmai	4097	1756
15 Sala delle Carrozze	2281	978
16 Cortile Parisio	7262	3112
17 Biglietteria	2869	1230
18 Portico	3991	1710
19 Scala Catalana	4479	1920
20 Loggetta	2918	1251
21 Box Sala8	3263	1398
22 Area Sosta	2454	1052
Total	98202	42084

tions et. In Table 32 we report the results on the 9 test videos (excluding “Test3”) of the object detector when using the mean Average Precision (mAP) as evaluation metric. The table also reports the number of frames annotated with bounding box for each test video. Per-class AP values are reported in Table 33.

2) *Monastero dei Benedettini*. A subset of frames from the dataset (sampled at at 1 fps) has been labeled with bounding boxes. The annotated frames with bounding box of the Training set used in the first task are 33366. We split this set into training/validation sets to train and validate the object detector. In particular, we used 23363 frames as training set and 10003 as validation set. We used the 10 validation videos (2235 frames) to find the optimal threshold of the object detector. Table 34 reports the AP values obtained for each of the 35 considered points of interest belonging to the validation set, using the optimal threshold found through validation (0.001). We annotated the 60 real visits with bounding boxes for a total of 71310

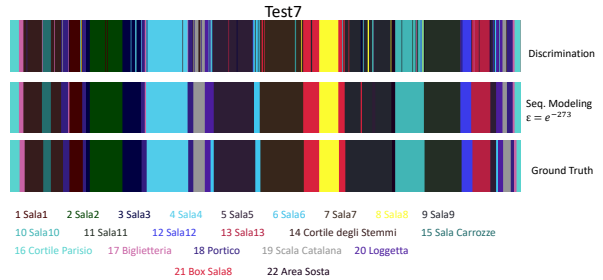


Figure 16: Color-coded segmentations for the test video “Test7” of “Palazzo Bellomo”.

images. We tested the object detector on these frames. The results are shown in Table 35. Also, for each test video, we report the number of frames annotated with bounding boxes. Per-class AP values are reported in Table 36.

5.3.3. Object Retrieval

To address this task, we extracted image patches from the bounding box annotations of the dataset.

1) *Palazzo Bellomo*. In Table 37, we report the number of image patches which have been extracted for each test video. For one-shot learning, we have used reference images for training and all the image patches for testing, whereas to perform many-shot learning, we used the patches belonging to test videos 1 – 7 for training (15185 patches) and the others to test (8542 patches).

2) *Monastero dei Benedettini*. In Table 38, we report the number of image patches extracted from the 60 test videos. One-shot learning has been performed using reference images from training and all extracted patches for testing. For many-shot learning, we used 30497 image patches belonging to the visits with IDs from 100 to 147 for training, and 14551 patches belonging to the visits with ID from 148 to 166 for testing.

Object Retrieval with DenseNet. We tried to extract features using a different backbone respect to the proposed baseline based on VGG-19. We extracted the features from the FC7 layer of DenseNet[22]. In this way, we obtained for each image a fixed-size vector of 1024 values. We have followed the same pipeline used with VGG to perform object retrieval. Table 39 shows the results

Table 19: Detailed results of the 9 test videos of “Palazzo Bellomo” using the FF_1 score. The “/” sign indicates that no samples from that class was present in the test video.

Class	FF_1 score									
	Test1	Test2	Test4	Test5	Test6	Test7	Test8	Test9	Test10	AVG
1_Sala1	0,16	0,00	0,81	0,96	0,86	0,96	0,90	0,85	0,92	0,71
2_Sala2	0,78	0,67	0,96	0,96	0,99	0,99	0,97	0,97	0,96	0,92
3_Sala3	0,68	0,75	0,97	0,87	0,72	0,96	0,83	0,89	0,91	0,84
4_Sala4	0,89	0,96	0,93	0,91	/	0,98	0,86	0,91	0,95	0,92
5_Sala5	0,90	0,94	0,98	0,95	0,89	0,95	0,95	0,95	0,97	0,94
6_Sala6	0,84	0,86	0,80	0,59	0,76	0,57	0,93	0,93	0,68	0,77
7_Sala7	0,99	0,95	0,99	0,88	0,84	0,99	0,92	0,93	0,97	0,94
8_Sala8	0,85	0,95	0,96	0,91	0,67	0,90	0,84	0,95	0,93	0,89
9_Sala9	0,86	0,97	0,95	0,90	0,76	0,93	0,94	0,94	0,90	0,91
10_Sala10	0,87	0,96	0,96	0,97	0,00	0,98	0,95	0,97	0,90	0,84
11_Sala11	0,86	0,96	0,97	0,97	0,00	0,97	0,94	0,96	0,96	0,84
12_Sala12	0,82	0,91	0,86	0,94	0,00	0,88	0,96	0,94	0,88	0,80
13_Sala13	0,74	0,94	0,91	0,85	0,00	0,95	0,97	0,92	0,94	0,80
14_CortiledegliStemmi	0,92	0,80	0,93	0,89	0,73	0,92	0,97	0,94	0,57	0,85
15_SalaCarrozze	0,91	0,89	0,96	0,93	0,90	0,90	0,85	0,96		0,91
16_CortileParisio	0,74	0,50	0,92	0,71	0,48	0,91	0,99	0,74	0,72	0,75
17_Biglietteria	0,64	0,79	0,81	0,49	0,74	0,55	0,66	0,61	0,53	0,65
18_Portico	0,71	0,42	0,77	0,70	0,70	0,73	0,71	0,75	0,72	0,69
19_ScalaCatalana	0,70	0,78	0,80	0,77	0,39	0,86	0,83	0,95	0,76	0,76
20_Loggetta	0,62	0,39	0,75	0,58	0,67	0,77	0,84	0,94	0,81	0,71
21_BoxSala8	0,97	0,97	0,98	/	0,79	0,97	0,99	0,94	0,94	0,94
22_AreaSosta	0,24	0,81	0,56	0,46	0,00	0,87	0,56	0,78	0,77	0,56
mFF1	0,76	0,78	0,89	0,82	0,57	0,89	0,88	0,90	0,84	0,81

Table 20: Detailed results of the 9 test videos of “Palazzo Bellomo” using the ASF_1 score. The “/” sign indicates that no samples from that class was present in the test video.

Class	ASF_1									
	Test1	Test2	Test4	Test5	Test6	Test7	Test8	Test9	Test10	AVG
1_Sala1	0,07	0,00	0,25	0,92	0,15	0,92	0,62	0,55	0,84	0,48
2_Sala2	0,32	0,46	0,92	0,92	0,97	0,98	0,65	0,94	0,92	0,79
3_Sala3	0,18	0,34	0,66	0,63	0,19	0,88	0,16	0,66	0,83	0,50
4_Sala4	0,58	0,92	0,41	0,45	/	0,95	0,21	0,28	0,89	0,59
5_Sala5	0,57	0,56	0,84	0,77	0,27	0,72	0,49	0,64	0,94	0,64
6_Sala6	0,47	0,61	0,44	0,35	0,32	0,40	0,87	0,80	0,44	0,52
7_Sala7	0,97	0,91	0,97	0,28	0,07	0,97	0,21	0,16	0,94	0,61
8_Sala8	0,40	0,86	0,84	0,59	0,13	0,68	0,70	0,74	0,85	0,64
9_Sala9	0,19	0,93	0,90	0,45	0,32	0,18	0,64	0,14	0,49	0,47
10_Sala10	0,28	0,92	0,92	0,93	0,00	0,96	0,42	0,94	0,81	0,69
11_Sala11	0,21	0,92	0,94	0,94	0,00	0,94	0,13	0,21	0,92	0,58
12_Sala12	0,37	0,82	0,76	0,89	0,00	0,79	0,91	0,63	0,78	0,66
13_Sala13	0,32	0,88	0,83	0,74	0,00	0,90	0,94	0,40	0,89	0,66
14_CortiledegliStemmi	0,78	0,61	0,83	0,79	0,16	0,84	0,72	0,69	0,39	0,64
15_SalaCarrozze	0,49	0,80	0,92	0,87	0,37	0,82	0,17	0,91	/	0,67
16_CortileParisio	0,31	0,22	0,65	0,41	0,16	0,80	0,98	0,47	0,54	0,50
17_Biglietteria	0,41	0,67	0,69	0,32	0,45	0,29	0,38	0,37	0,36	0,44
18_Portico	0,54	0,28	0,60	0,43	0,51	0,58	0,56	0,49	0,58	0,51
19_ScalaCatalana	0,61	0,48	0,66	0,64	0,41	0,72	0,65	0,91	0,62	0,63
20_Loggetta	0,43	0,31	0,37	0,46	0,23	0,63	0,70	0,87	0,64	0,51
21_BoxSala8	0,65	0,94	0,96	/	0,34	0,94	0,98	0,64	0,88	0,79
22_AreaSosta	0,24	0,69	0,42	0,47	0,00	0,76	0,41	0,64	0,62	0,47
mASF1	0,43	0,64	0,72	0,63	0,24	0,76	0,57	0,59	0,72	0,59

Table 21: Detailed results of the 9 test videos of “Palazzo Bellomo” using the FF_1 score. The backbone used is DenseNet. The “/” sign indicates that no samples from that class was present in the test video.

FF_1 score										
Class	Test1	Test2	Test4	Test5	Test6	Test7	Test8	Test9	Test10	AVG
1_Sala1	0,00	0,96	0,96	0,95	0,99	0,94	0,97	0,97	0,91	0,85
2_Sala2	0,82	0,91	0,96	0,95	0,99	0,99	0,99	0,97	0,95	0,95
3_Sala3	0,65	0,84	0,84	0,83	0,94	0,88	0,61	0,85	0,83	0,81
4_Sala4	0,84	0,96	0,98	0,93	/	0,96	0,78	0,97	0,93	0,92
5_Sala5	0,98	0,88	0,96	0,91	0,65	0,91	0,97	0,97	0,99	0,91
6_Sala6	0,84	0,80	0,43	0,00	0,86	0,00	0,18	0,83	0,73	0,52
7_Sala7	0,67	0,93	0,22	0,00	0,38	0,88	0,88	0,40	0,96	0,59
8_Sala8	0,64	0,75	0,60	0,52	0,56	0,56	0,64	0,60	0,77	0,63
9_Sala9	0,90	0,89	0,73	0,88	0,41	0,92	0,98	0,99	0,88	0,84
10_Sala10	0,96	0,98	0,92	0,72	0,00	0,98	0,85	0,97	0,78	0,80
11_Sala11	0,93	0,96	0,96	0,97	0,00	0,97	0,97	0,98	0,95	0,85
12_Sala12	0,82	0,88	0,87	0,90	0,00	0,87	0,85	0,92	0,87	0,78
13_Sala13	0,96	0,95	0,95	0,76	0,00	0,93	0,95	0,94	0,96	0,82
14_CortiledegliStemmi	0,78	0,68	0,74	0,88	0,17	0,90	0,79	0,92	0,00	0,65
15_SalaCarrozze	0,84	0,89	0,95	0,93	0,91	0,84	0,97	0,95	/	0,91
16_CortileParisio	0,33	0,51	0,52	0,45	0,35	0,60	0,81	0,91	0,80	0,59
17_Biglietteria	0,25	0,56	0,00	0,26	0,00	0,00	0,00	0,36	0,00	0,16
18_Portico	0,68	0,67	0,78	0,54	0,69	0,69	0,62	0,78	0,62	0,67
19_ScalaCatalana	0,00	0,00	0,54	0,55	0,78	0,67	0,58	0,85	0,49	0,50
20_Loggetta	0,00	0,50	0,80	0,44	0,72	0,53	0,55	0,83	0,68	0,56
21_BoxSala8	0,96	0,97	0,97	/	0,00	0,97	0,99	0,81	0,94	0,83
22_AreaSosta	0,74	0,85	0,45	0,83	0,00	0,56	0,72	0,85	0,00	0,56
mFF1	0,66	0,79	0,73	0,68	0,45	0,75	0,76	0,85	0,72	0,71

Table 22: Detailed results of the 9 test videos of “Palazzo Bellomo” using the ASF_1 score. The backbone used is DenseNet. The “/” sign indicates that no samples from that class was present in the test video.

ASF_1 score										
Class	Test1	Test2	Test4	Test5	Test6	Test7	Test8	Test9	Test10	AVG
1_Sala1	0,00	0,93	0,93	0,91	0,97	0,89	0,94	0,94	0,84	0,82
2_Sala2	0,64	0,83	0,91	0,90	0,97	0,97	0,97	0,94	0,90	0,89
3_Sala3	0,36	0,53	0,51	0,53	0,61	0,63	0,29	0,59	0,72	0,53
4_Sala4	0,39	0,92	0,96	0,86	/	0,92	0,31	0,94	0,87	0,77
5_Sala5	0,94	0,64	0,66	0,62	0,34	0,63	0,54	0,95	0,97	0,70
6_Sala6	0,46	0,43	0,18	0,00	0,42	0,00	0,09	0,59	0,51	0,30
7_Sala7	0,50	0,86	0,13	0,00	0,10	0,78	0,28	0,25	0,92	0,42
8_Sala8	0,36	0,44	0,42	0,35	0,35	0,41	0,49	0,53	0,57	0,44
9_Sala9	0,81	0,80	0,51	0,78	0,22	0,84	0,95	0,97	0,78	0,74
10_Sala10	0,92	0,96	0,85	0,19	0,00	0,95	0,27	0,94	0,64	0,64
11_Sala11	0,53	0,92	0,91	0,93	0,00	0,93	0,65	0,95	0,91	0,75
12_Sala12	0,38	0,79	0,76	0,81	0,00	0,77	0,30	0,85	0,76	0,60
13_Sala13	0,93	0,90	0,90	0,62	0,00	0,87	0,91	0,89	0,92	0,77
14_CortiledegliStemmi	0,57	0,53	0,54	0,77	0,08	0,81	0,49	0,67	0,00	0,50
15_SalaCarrozze	0,72	0,80	0,91	0,87	0,83	0,72	0,94	0,90	/	0,84
16_CortileParisio	0,30	0,63	0,43	0,65	0,38	0,47	0,59	0,83	0,66	0,55
17_Biglietteria	0,26	0,43	0,00	0,11	0,00	0,00	0,00	0,39	0,00	0,13
18_Portico	0,51	0,47	0,60	0,45	0,48	0,52	0,42	0,64	0,47	0,51
19_ScalaCatalana	0,00	0,00	0,43	0,51	0,61	0,55	0,42	0,75	0,39	0,41
20_Loggetta	0,00	0,31	0,54	0,28	0,46	0,38	0,48	0,72	0,44	0,40
21_BoxSala8	0,92	0,95	0,95	/	0,00	0,93	0,98	0,67	0,88	0,79
22_AreaSosta	0,59	0,73	0,29	0,71	0,00	0,75	0,57	0,74	0,00	0,49
mFF1	0,50	0,67	0,61	0,56	0,32	0,67	0,54	0,76	0,63	0,58

Table 24: Detailed results on the 60 test videos of “Monastero dei Benedettini”, considering the evaluation measure FF_1 score. The “/” sign indicates that no samples from that class were present in the test video. The four classes are: 1) Antirefettorio, 2) Aula S. Mazzarino, 3) Cucina, 3) Ventre, whereas Neg. represents the negatives.

ID_Visit	1	2	3	4	Neg.	AVG
4805	0,91	0,33	0,75	0,99	0,43	0,682
1804	0,69	0,43	0,8	0,99	0,43	0,668
4377	0,77	0,47	0,84	0,99	0,5	0,714
1669	0,8	0,51	0,92	0,99	0,67	0,778
1791	0,59	0,22	0,78	0,98	0,35	0,584
3948	0,8	/	0,66	0,99	0,61	0,765
3152	0,71	0,25	0,86	0,98	0,75	0,71
4361	0,89	0,32	0,81	0,099	0,62	0,5478
3976	0,97	0,62	0,92	0,98	0,68	0,834
3527	0,85	0	0,82	0,99	0,63	0,658
4105	0,81	0	0,77	0,99	0,66	0,646
1399	0,65	0	0,76	0,99	0,62	0,604
3836	0,65	0	0,76	0,99	0,62	0,604
4006	0,81	0,82	0,92	0,99	0,75	0,858
4415	0,87	0,49	0,77	0,98	0,73	0,768
3008	0,82	0,2	0,63	0,99	0,23	0,574
4660	0,82	0,2	0,63	0,99	0,23	0,574
2826	0,79	0,57	0,81	1	0,41	0,716
1099	0,77	0,27	0,64	0,98	0,5	0,632
4391	0,8	0,03	0,79	0,98	0,55	0,63
3929	0,94	0	0,84	0,99	0,67	0,688
3362	0,46	0,25	0,76	0,99	0,46	0,584
1379	0,84	0	0,75	0,98	0,33	0,58
2600	0,74	0	0,78	0,99	0,59	0,62
1430	0,57	0,22	0,72	0,96	0,47	0,588
2956	0,53	0,33	0,73	0,96	0,8	0,67
4742	0,14	0,43	0,92	0,99	0,67	0,63
3651	0,94	0,66	0,82	0,99	0,49	0,78
1064	0,77	0,12	0,88	0,99	0,31	0,614
3818	0,93	0,14	0,71	0,99	0,51	0,656

Table 25: Continued from Table 24

ID_Visit	1	2	3	4	Neg.	AVG
2043	0,72	0,47	0,79	0,98	0,66	0,724
3996	0,47	0,4	0,72	0,98	0,76	0,666
3455	0,85	0,33	0,88	0,99	0,61	0,732
4785	0,03	0	0,73	0,96	0,65	0,474
2047	0,95	0,55	0,86	1	0,6	0,792
1912	0,79	0,44	0,67	0,96	0,62	0,696
3232	0,89	0,33	0,82	0,99	0,73	0,752
4442	0,9	0,37	0,82	0,99	0,54	0,724
3646	0,67	0,09	0,8	1	0,34	0,58
4833	0,78	0,59	0,86	0,95	0,66	0,768
3478	0,82	0,51	0,94	1	0,67	0,788
4396	0,81	0,53	0,89	0,99	0,38	0,72
2894	0,67	0	0,87	0,97	0,55	0,612
4414	0,88	0,51	0,9	0,98	0,58	0,77
4639	0,73	0,21	0,83	0,99	0,61	0,674
1004	0,15	0,44	0,6	0,98	0,48	0,53
1917	0,44	0,48	0,51	1	0,42	0,57
1153	0,78	0,54	0,81	0,98	0,57	0,736
2244	0,86	0,3	0,77	0,99	0,32	0,648
2614	0,97	0,59	0,84	0,99	0,39	0,756
1624	0,91	0,71	0,69	0,99	0,48	0,756
3441	0,82	0,45	0,79	0,99	0,46	0,702
4793	0,82	/	0,74	0,99	0,42	0,7425
4083	0,84	/	0,72	0,99	0,73	0,82
4906	0,77	0,2	0,74	0,99	0,42	0,624
1160	0,84	0,66	0,74	1	0,42	0,732
3416	0,77	/	0,82	0,99	0,68	0,815
1051	0,79	0,43	0,78	0,98	0,4	0,676
2580	0,73	0,18	0,89	0,99	0,48	0,654
1109	0,81	0,28	0,89	0,99	0,43	0,68
AVG	0,75	0,33	0,79	0,99	0,54	0,68

Table 26: Detailed results on the 60 test videos of “Monastero dei Benedettini”, considering the evaluation measure ASF_1 score. The “/” sign indicates that no samples from that class were present in the test video. The four classes are: 1) Antirefettorio, 2) Aula S. Mazzarino, 3) Cucina, 3) Ventre, whereas Neg. represents the negatives.

ID_Visit	1	2	3	4	Neg.	AVG
4805	0,71	0,06	0,2	0,31	0,22	0,3
1804	0,55	0,06	0,26	0,15	0,14	0,232
4377	0,55	0,05	0,18	0,98	0,25	0,402
1669	0,4	0,22	0,27	0,93	0,49	0,462
1791	0,46	0,14	0,21	0,25	0,19	0,25
3948	0,67	/	0,26	0,98	0,51	0,605
3152	0,54	0,07	0,66	0,88	0,47	0,524
4361	0,47	0,12	0,17	0,27	0,29	0,264
3976	0,94	0,24	0,45	0,22	0,56	0,482
3527	0,72	0	0,44	0,46	0,55	0,434
4105	0,54	0	0,49	0,39	0,51	0,386
1399	0,15	0	0,44	0,98	0,26	0,366
3836	0,15	0	0,44	0,98	0,26	0,366
4006	0,56	0,69	0,53	0,98	0,5	0,652
4415	0,48	0,15	0,24	0,49	0,48	0,368
3008	0,68	0,1	0,26	0,98	0,1	0,424
4660	0,68	0,09	0,26	0,98	0,09	0,42
2826	0,71	0,14	0,24	0,99	0,41	0,498
1099	0,54	0,15	0,16	0,49	0,31	0,33
4391	0,67	0,05	0,32	0,96	0,53	0,506
3929	0,91	0	0,43	0,98	0,46	0,556
3362	0,36	0,07	0,45	0,57	0,36	0,362
1379	0,49	0	0,21	0,18	0,23	0,222
2600	0,64	0	0,41	0,98	0,38	0,482
1430	0,17	0,09	0,03	0,08	0,1	0,094
2956	0,33	0,06	0,46	0,37	0,48	0,34
4742	0,1	0,09	0,39	0,62	0,17	0,274
3651	0,63	0,19	0,54	0,65	0,28	0,458
1064	0,63	0,05	0,3	0,65	0,13	0,352
3818	0,78	0,11	0,17	0,25	0,35	0,332

Table 27: Continued from Table 26

ID_Visit	1	2	3	4	Neg.	AVG
2043	0,4	0,22	0,2	0,27	0,25	0,268
3996	0,45	0,12	0,39	0,4	0,49	0,37
3455	0,69	0,16	0,39	0,96	0,5	0,54
4785	0,05	0,11	0,53	0,13	0,33	0,205
2047	0,9	0,38	0,22	0,99	0,36	0,57
1912	0,54	0,1	0,16	0,31	0,39	0,3
3232	0,73	0,04	0,49	0,98	0,53	0,554
4442	0,59	0,285	0,22	0,27	0,25	0,323
3646	0,4	0,13	0,2	0,66	0,34	0,346
4833	0,54	0,15	0,53	0,19	0,36	0,354
3478	0,71	0,29	0,54	0,99	0,46	0,598
4396	0,35	0,02	0,05	0,18	0,08	0,136
2894	0,43	0	0,46	0,35	0,18	0,284
4414	0,75	0,11	0,31	0,24	0,3	0,342
4639	0,61	0,12	0,44	0,39	0,43	0,398
1004	0,22	0,11	0,27	0,21	0,38	0,238
1917	0,41	0,11	0,15	0,31	0,26	0,248
1153	0,54	0,11	0,33	0,59	0,26	0,366
2244	0,58	0,16	0,45	0,97	0,26	0,484
2614	0,65	0,05	0,39	0,98	0,3	0,474
1624	0,71	0,29	0,35	0,39	0,23	0,394
3441	0,59	0,13	0,26	0,88	0,25	0,422
4793	0,68	/	0,52	0,98	0,38	0,64
4083	0,65	/	0,31	0,35	0,59	0,475
4906	0,57	0,06	0,33	0,59	0,33	0,376
1160	0,7	0,15	0,34	0,99	0,29	0,494
3416	0,45	/	0,24	0,97	0,19	0,4625
1051	0,55	0,28	0,33	0,4	0,27	0,366
2580	0,37	0,1	0,68	0,21	0,37	0,346
1109	0,66	0,16	0,63	0,98	0,36	0,558
AVG	0,54	0,12	0,34	0,60	0,33	0,40

Table 28: Detailed results on the 60 test videos of “Monastero dei Benedettini”, considering the evaluation measure FF_1 score. We used the backbone DenseNet. The “/” sign indicates that no samples from that class were present in the test video. The four classes are: 1) Antirefettorio, 2) Aula S. Mazzarino, 3) Cucina, 3) Ventre, whereas Neg. represents the negatives.

ID_Visit	1	2	3	4	Neg.	AVG
4805	0,79	0,82	0,74	0,95	0,45	0,75
1804	0,32	0,36	0,79	0,98	0,55	0,6
4377	0,46	0,53	0,8	0,96	0,32	0,614
1669	0,75	0,72	0,9	0,99	0,45	0,762
1791	0,49	0,5	0,84	0,95	0,43	0,642
3948	0	/	0,59	0,98	0,49	0,515
3152	0,39	0,72	0,87	0,97	0,76	0,742
4361	0,55	0,78	0,82	0,97	0,28	0,68
3976	0,87	0,43	0,81	0,93	0,66	0,74
3527	0,86	0,79	0,85	0,99	0,56	0,81
4105	0,62	0	0,77	0,97	0,09	0,49
1399	0,46	0,23	0,79	0,98	0,43	0,578
3836	0,51	0	0,72	0,99	0,55	0,554
4006	0,57	0,78	0,78	0,99	0,37	0,698
4415	0,86	0,61	0,82	0,97	0,35	0,722
3008	0,69	0	0,7	0,99	0,51	0,578
4660	0,62	0,81	0,83	0,98	0,57	0,762
2826	0,31	0,52	0,76	0,97	0,61	0,634
1099	0,62	0,42	0,85	0,98	0,49	0,672
4391	0,74	0,72	0,72	0,98	0,33	0,698
3929	0,26	0	0,8	0,99	0,43	0,496
3362	0,34	0,68	0,68	0,95	0,21	0,572
1379	0,41	0	0,81	0,96	0,45	0,526
2600	0,26	0	0,63	0,96	0,3	0,43
1430	0,94	0	0,69	0,87	0,58	0,616
2956	0,32	0,45	0,33	0,91	0,65	0,532
4742	0,13	0,58	0,59	0,84	0,5	0,528
3651	0,77	0,41	0,88	0,98	0,41	0,69
1064	0,77	0,23	0,83	0,99	0,22	0,608
3818	0,68	0,64	0,73	0,99	0,47	0,702

Table 29: Continued from Table 28

ID_Visit	1	2	3	4	Neg.	AVG
2043	0,52	0,44	0,62	0,91	0,33	0,564
3996	0,17	0,68	0,57	0,95	0,72	0,618
3455	0,78	0,65	0,8	0,99	0,55	0,754
4785	0,02	0	0,64	0,95	0,33	0,388
2047	0,95	0,84	0,84	0,98	0,38	0,798
1912	0,66	0,51	0,69	0,94	0,46	0,652
3232	0,66	0,53	0,85	0,99	0,49	0,704
4442	0,8	0,77	0,82	0,98	0,23	0,72
3646	0,63	0,29	0,66	0,97	0,4	0,59
4833	0,67	0,82	0,67	0,88	0,24	0,656
3478	0,71	0,75	0,84	0,98	0,37	0,73
4396	0,74	0,69	0,89	0,96	0,26	0,708
2894	0,83	0,8	0,66	0,92	0,48	0,738
4414	0,75	0,7	0,91	0,95	0,22	0,706
4639	0,57	0,52	0,85	0,99	0,39	0,664
1004	0,08	0,72	0,49	0,97	0,26	0,504
1917	0,41	0,71	0,36	0,87	0,29	0,528
1153	0,62	0,7	0,71	0,91	0,32	0,652
2244	0,94	0,56	0,74	0,99	0,43	0,732
2614	0,88	0,56	0,83	0,99	0,4	0,732
1624	0,33	0,8	0,62	0,99	0,33	0,614
3441	0,61	0,25	0,84	0,99	0,41	0,62
4793	0,52	/	0,68	0,99	0,33	0,63
4083	0,93	/	0,73	0,99	0,71	0,84
4906	0,46	0,36	0,59	0,94	0,31	0,532
1160	0,88	0,84	0,72	0,98	0,46	0,776
3416	0,56	/	0,5	0,82	0,2	0,52
1051	0,78	0,76	0,64	0,95	0,57	0,74
2580	0,71	0,45	0,63	0,96	0,44	0,638
1109	0,91	0,28	0,85	0,97	0,36	0,674
mFF1	0,59	0,51	0,73	0,96	0,42	0,64

Table 30: Detailed results on the 60 test videos of “Monastero dei Benedettini”, considering the evaluation measure ASF_1 score. We used the backbone DenseNet. The “/” sign indicates that no samples from that class were present in the test video. The four classes are: 1) Antirefettorio, 2) Aula S. Mazzarino, 3) Cucina, 3) Ventre, whereas Neg. represents the negatives.

ID_visit	1	2	3	4	Neg.	AVG
4805	0,55	0,62	0,4	0,12	0,36	0,41
1804	0,32	0,1	0,39	0,05	0,21	0,214
4377	0,31	0,12	0,26	0,08	0,27	0,208
1669	0,68	0,34	0,67	0,9	0,33	0,584
1791	0,41	0,41	0,25	0,06	0,22	0,27
3948	0	/	0,33	0,26	0,49	0,27
3152	0,39	0,43	0,64	0,44	0,58	0,496
4361	0,26	0,35	0,32	0,08	0,2	0,242
3976	0,68	0,26	0,35	0,03	0,5	0,364
3527	0,65	0,37	0,38	0,28	0,44	0,424
4105	0,49	0	0,49	0,16	0,06	0,24
1399	0,55	0,1	0,33	0,23	0,29	0,3
3836	0,47	0	0,35	0,98	0,34	0,428
4006	0,47	0,49	0,49	0,55	0,26	0,452
4415	0,74	0,07	0,35	0,14	0,23	0,306
3008	0,57	0	0,5	0,99	0,41	0,494
4660	0,45	0,52	0,17	0,05	0,27	0,292
2826	0,31	0,29	0,33	0,06	0,25	0,248
1099	0,52	0,24	0,55	0,96	0,17	0,488
4391	0,54	0,13	0,14	0,2	0,24	0,25
3929	0,15	0	0,31	0,22	0,39	0,214
3362	0,2	0,35	0,38	0,05	0,22	0,24
1379	0,22	0	0,43	0,05	0,35	0,21
2600	0,32	0	0,39	0,31	0,26	0,256
1430	0,43	0	0,1	0,02	0,21	0,152
2956	0,41	0,2	0,16	0,29	0,31	0,274
4742	0,15	0,18	0,25	0,11	0,17	0,172
3651	0,32	0,25	0,74	0,48	0,28	0,414
1064	0,63	0,09	0,52	0,97	0,19	0,48
3818	0,43	0,49	0,2	0,06	0,29	0,294

Table 31: Continued from Table 30

ID_Visit	1	2	3	4	Neg.	AVG
2043	0,4	0,27	0,28	0,13	0,27	0,27
3996	0,23	0,3	0,22	0,05	0,38	0,236
3455	0,63	0,38	0,65	0,99	0,48	0,626
4785	0,06	0	0,25	0,45	0,07	0,166
2047	0,9	0,72	0,72	0,37	0,34	0,61
1912	0,42	0,25	0,2	0,08	0,3	0,25
3232	0,52	0,1	0,34	0,13	0,35	0,288
4442	0,61	0,22	0,33	0,06	0,12	0,268
3646	0,62	0,3	0,26	0,19	0,22	0,318
4833	0,54	0,55	0,26	0,1	0,21	0,332
3478	0,57	0,4	0,47	0,14	0,24	0,364
4396	0,41	0,1	0,2	0,07	0,12	0,18
2894	0,55	0,67	0,36	0,13	0,31	0,404
4414	0,59	0,26	0,83	0,11	0,08	0,374
4639	0,5	0,35	0,67	0,39	0,27	0,436
1004	0,22	0,47	0,17	0,11	0,17	0,228
1917	0,39	0,19	0,03	0,03	0,21	0,17
1153	0,47	0,33	0,21	0,14	0,16	0,262
2244	0,77	0,35	0,49	0,35	0,32	0,456
2614	0,43	0,06	0,39	0,99	0,26	0,426
1624	0,3	0,33	0,24	0,29	0,29	0,29
3441	0,49	0,11	0,44	0,88	0,42	0,468
4793	0,35	/	0,36	0,28	0,37	0,34
4083	0,88	/	0,41	0,25	0,48	0,505
4906	0,36	0,1	0,31	0,06	0,21	0,208
1160	0,79	0,5	0,39	0,14	0,28	0,42
3416	0,32	/	0,25	0,42	0,34	0,3325
1051	0,5	0,31	0,21	0,05	0,4	0,294
2580	0,36	0,29	0,33	0,05	0,28	0,262
1109	0,82	0,17	0,46	0,11	0,28	0,368
mASF1	0,46	0,26	0,37	0,28	0,28	0,40

Table 32: Detailed results obtained using the YoloV3 object detector YoloV3 on the 9 test videos of “Palazzo Bellomo”. The last column reports the number of frames belonging to each test video. The last row indicates the average of mAP score obtained for each test video.

	mAP	#images
Test1	12,72	1644
Test2	13,61	1238
Test4	12,31	1398
Test5	8,65	848
Test6	8,9	1453
Test7	10,29	1200
Test8	10,98	2826
Test9	9,97	2004
Test10	7,85	791
AVG	10,59	13402

Table 33: Per-class AP values obtained on the 9 test videos. The “/” sign indicates that no samples from that class were present in the test videos.

Class	AP	Class	AP	Class	AP	Class	AP
0	10,53	50	6,67	100	5,34	150	0,07
1	41,45	51	0,00	101	9,45	151	2,18
2	49,60	52	0,60	102	5,44	152	17,34
3	41,79	53	0,00	103	10,69	153	15,24
4	13,27	54	0,00	104	10,77	154	35,10
5	66,73	55	1,85	105	2,71	155	0,00
6	66,97	56	0,00	106	2,67	156	0,71
7	72,62	57	0,00	107	0,00	157	3,90
8	52,41	58	0,00	108	3,80	158	1,92
9	68,21	59	0,00	109	6,11	159	0,00
10	2,69	60	1,85	110	16,50	160	0,63
11	14,79	61	0,13	111	0,00	161	18,92
12	2,19	62	6,05	112	12,32	162	11,44
13	44,65	63	0,00	113	0,00	163	18,31
14	35,27	64	2,34	114	0,00	164	12,00
15	16,58	65	2,22	115	21,75	165	26,78
16	61,05	66	0,74	116	9,50	166	11,97
17	28,68	67	18,36	117	4,98	167	11,18
18	46,37	68	9,19	118	1,97	168	1,04
19	9,68	69	5,70	119	1,64	169	23,41
20	51,04	70	1,14	120	33,36	170	11,95
21	11,11	71	2,64	121	0,39	171	0,52
22	45,00	72	9,19	122	9,74	172	2,82
23	48,80	73	11,26	123	3,62	173	5,12
24	10,40	74	0,11	124	17,70	174	37,13
25	0,00	75	34,85	125	0,00	175	30,37
26	17,47	76	6,56	126	0,96	176	18,87
27	0,00	77	0,75	127	1,05	177	/
28	14,01	78	10,02	128	4,66	178	0,00
29	0,00	79	5,16	129	12,58	179	0,00
30	3,11	80	16,57	130	15,71	180	3,43
31	16,71	81	17,89	131	6,43	181	0,00
32	0,61	82	/	132	2,74	182	0,43
33	2,76	83	25,18	133	0,00	183	0,10
34	0,99	84	0,86	134	0,00	184	0,00
35	0,00	85	1,00	135	0,00	185	0,04
36	0,56	86	0,00	136	5,02	186	0,00
37	4,60	87	15,17	137	4,49	187	0,00
38	16,01	88	12,64	138	0,30	188	3,86
39	3,38	89	8,99	139	0,00	189	0,00
40	14,20	90	25,49	140	0,83	190	10,42
41	0,00	91	0,29	141	11,66	mAP	10.66
42	0,00	92	0,00	142	0,00		
43	25,33	93	8,84	143	1,25		
44	0,11	94	11,98	144	6,56		
45	0,00	95	4,79	145	25,27		
46	0,00	96	0,00	146	2,01		
47	0,00	97	0,00	147	31,66		
48	20,07	98	0,00	148	2,15		
49	12,90	99	0,00	149	0,23		

Table 34: Per-class AP values obtained on the validation set using the optimal threshold of 0.001.

Class	AP
5.1 PortaAulaS.Mazzarino	36,05
5.2 PortaIngressoMuseoFabbrica	37,99
5.3 PortaAntirefettorio	20,44
5.4 PortaIngressoRef.Piccolo	26,07
5.5 Cupola	73,98
5.6 AperturaPavimento	80,52
5.7 S.Agata	74,89
5.8 S.Scolastica	66,84
6.1 QuadroSantoMazzarino	76,89
6.2 Affresco	60,39
6.3 PavimentoOriginale	37,41
6.4 PavimentoRestaurato	13,11
6.5 BassorilieviMancanti	25,7
6.6 LavamaniSx	41,55
6.7 LavamaniDx	25,8
6.8 TavoloRelatori	14,59
6.9 Poltrone	23,48
7.1 Edicola	42,57
7.2 PavimentoA	6,03
7.3 PavimentoB	0,93
7.4 PassavivandePavimentoOriginale	44,44
7.5 AperturaPavimento	33,12
7.6 Scala	46,49
7.7 SalaMetereologica	30,83
8.1 Doccione	46,65
8.2 VanoRaccoltaCenere	75,6
8.3 SalaRossa	18,52
8.4 ScalaCucina	42,48
8.5 CucinaProv. v.	48,26
8.6 Ghiacciaia	47,7
8.7 Latrina	72,86
8.8 OssaeScarti	50,36
8.9 Pozzo	51,74
8.10 Cisterna	17,03
8.11 BustoPietroTacchini	41,15
Negatives	5,8
mAP	40.51

Table 35: Detailed results of the YoloV3 object detector on the the 60 real visits. The second column reports the number of frames annotated with bounding box contained in each visit.

ID_Visit	#images	mAP	ID_Visit	#images	mAP
156	1288	19,89	117	1349	19,89
154	2443	18,73	115	1511	18,73
153	1620	20,34	135	1396	20,34
155	786	18,58	137	2484	18,58
110	1671	18,36	136	566	18,36
109	679	19,96	132	1233	19,96
108	1065	13,84	134	1177	13,84
107	1728	15,24	130	1401	15,24
158	1660	13,37	105	1434	13,37
157	874	12,40	124	936	12,40
160	660	19,20	123	1571	19,20
159	654	16,08	103	2411	16,08
129	751	14,88	104	1492	14,88
125	544	19,51	122	1794	19,51
126	892	21,75	120	939	21,75
163	683	12,75	140	1050	12,75
165	1689	17,91	139	1454	17,91
161	1563	22,52	138	1736	22,52
162	979	20,02	145	1343	20,02
166	597	10,80	146	1370	10,80
164	1197	13,98	114	868	13,98
142	1161	17,27	112	840	17,27
144	868	11,61	111	726	11,61
143	824	9,64	113	851	9,64
101	1894	13,52	149	1612	13,52
102	824	10,50	148	847	10,50
100	1343	19,07	147	450	19,07
119	564	17,43	152	1519	17,43
118	740	16,74	150	1437	16,74
116	1618	16,55	151	942	16,55
			Tot./AVG	71310	

Table 36: Per-class AP values obtained on the 60 real visits.

Class	AP
5.1 PortaAulaS.Mazzarino	37,86
5.2 PortaIngressoMuseoFabbrica	27,18
5.3 PortaAntirefettorio	2,23
5.4 PortaIngressoRef.Piccolo	15,44
5.5 Cupola	65,80
5.6 AperturaPavimento	0,89
5.7 S.Agata	41,82
5.8 S.Scolastica	31,64
6.1 QuadroSantoMazzarino	13,22
6.2 Affresco	53,69
6.3 PavimentoOriginale	6,95
6.4 PavimentoRestaurato	4,35
6.5 BassorilieviMancanti	16,57
6.6 LavamaniSx	0,83
6.7 LavamaniDx	0,58
6.8 TavoloRelatori	4,99
6.9 Poltrone	9,39
7.1 Edicola	34,01
7.2 PavimentoA	1,59
7.3 PavimentoB	3,26
7.4 PassavivandePavimentoOriginale	9,79
7.5 AperturaPavimento	22,77
7.6 Scala	20,44
7.7 SalaMeteoreologica	11,71
8.1 Doccione	13,78
8.2 VanoRaccoltaCenere	16,47
8.3 SalaRossa	17,35
8.4 ScalaCucina	13,36
8.5 CucinaProv. .	16,33
8.6 Ghiacciaia	3,98
8.7 Latrina	22,29
8.8 OssaeScarti	29,39
8.9 Pozzo	13,43
8.10 Cisterna	5,45
8.11 BustoPietroTacchini	23,13
Negatives	16,81
mAP	17,47

Table 37: Number of patches extracted from each of the 10 test videos of “Palazzo Bellomo”.

Test video	#images
Test1	2568
Test2	2048
Test3	2672
Test4	2224
Test5	1439
Test6	2086
Test7	2148
Test8	4108
Test9	2914
Test10	1520
Total	23727

Table 38: Number of image patches extracted from the 60 test videos of “Monastero dei Benedettini”.

ID_Visit	#images	ID_Visit	#images
100	770	135	765
101	696	136	414
102	613	137	354
103	1733	138	824
104	768	139	707
105	929	140	494
107	1011	142	770
108	659	143	536
109	234	144	598
110	918	145	897
111	365	146	1307
112	727	147	173
113	288	148	692
114	561	149	954
115	623	150	609
116	968	151	652
117	810	152	699
118	907	153	1244
119	545	154	1691
120	669	155	666
121	774	156	709
122	1156	157	515
123	957	158	846
124	652	159	544
125	491	160	327
126	702	161	985
129	587	162	902
130	820	164	693
132	771	165	1133
134	884	166	690
Total		Total	44978

Table 39: Results using Densenet.

Points of Interest Retrieval				
1) Palazzo Bellomo				
Variant	K	Precision	Recall	F1 score
1 - One Shot	1	0,02	0,01	0,00
2 - Many Shots	1	0,62	0,59	0,6
	3	0,62	0,56	0,56
	5	0,62	0,56	0,56
	7	0,61	0,56	0,56
	9	0,61	0,55	0,56
11	0,61	0,55	0,55	
2) Monastero dei Benedettini				
Variant	K	Precision	Recall	F1 score
1 - One shot	1	0,38	0,07	0,09
2 - Many Shots	1	0,83	0,83	0,83
	3	0,84	0,83	0,83
	5	0,84	0,84	0,83
	7	0,84	0,83	0,83
	9	0,84	0,83	0,83
11	0,83	0,83	0,82	

Table 40: Results of the binary classifier obtained using a KNN with different values of K .

K	Precision	Recall	F1 score	Support
1	0,62	0,61	0,62	1980
3	0,62	0,65	0,63	
5	0,63	0,67	0,64	
7	0,64	0,69	0,65	
9	0,65	0,7	0,66	

Table 41: Results of the multi-class classifier obtained using a KNN with different values of K .

K	Precision	Recall	F1 score	Support
1	0,2	0,2	0,2	1980
3	0,2	0,23	0,19	
5	0,2	0,24	0,21	
7	0,22	0,24	0,22	
9	0,23	0,27	0,23	