

Food Recognition: A New Dataset, Experiments, and Results

Gianluigi Ciocca, Paolo Napoletano, and Raimondo Schettini

Abstract—We propose a new dataset for the evaluation of food recognition algorithms that can be used in dietary monitoring applications. Each image depicts a real canteen tray with dishes and foods arranged in different ways. Each tray contains multiple instances of food classes. The dataset contains 1027 canteen trays for a total of 3616 food instances belonging to 73 food classes. The food on the tray images has been manually segmented using carefully drawn polygonal boundaries. We have benchmarked the dataset by designing an automatic tray analysis pipeline that takes a tray image as input, finds the regions of interest, and predicts for each region the corresponding food class. We have experimented with three different classification strategies using also several visual descriptors. We achieve about 79% of food and tray recognition accuracy using convolutional-neural-networks-based features. The dataset, as well as the benchmark framework, are available to the research community.

Index Terms—Algorithm benchmarking, convolutional neural networks (CNN), food dataset, food recognition.

I. INTRODUCTION

HEALTH care on food and good practices in dietary behavior are drawing people's attention recently. Nowadays, technology can support the users to keep tracks of their food consumption, and to increase the awareness in their daily diet by monitoring their food habits. In the recent years many research works have demonstrated that machine learning and computer vision techniques can help to build systems to automatically recognize diverse foods and to estimate the food quantity [1]–[5]. To be useful for dietary monitoring, food recognition systems should also be able to operate in “wild” environments such as restaurants, canteens, and such. Obviously, a fair benchmarking of these systems, requires the availability of suitable datasets that actually pose the challenges of the food recognition task in unconstrained environments.

A. Food Recognition Systems

Research works in the literature have often focused on different aspects of the food recognition problem. Many works

address the challenges in the recognition of food by developing recognition strategies that differ in terms of features and classification methodologies. With respect to the features, the work of He *et al.* [6] describes the food image by combining both global and local features, while the work of Farinella *et al.* [7] uses a vocabulary built on textons. SIFT and local binary patterns (LBP) are used in [8], while in [9], the context of where the pictures are taken is also exploited along with the visual features. With respect to the classification strategies, the most widely used are k -nearest neighbor (k -NN) classifiers [6], [10], and support vector machines (SVMs) [7], [8]. An evaluation of different classification methodologies is reported in [5] where SVM, artificial neural networks, and random forest classification methods are analyzed. Recently, convolutional neural network (CNN) are used in the context of food recognition [11]–[13].

Other works in the literature focus on the design of a complete system for diet monitoring in real scenario. Often these systems exploit mobile application for food recognition, assessment, and logging. Examples of such systems are FoodLog [14], DietCam [15], Menu-Match [16], FoodCam [17], and those described in [10], [18], and [19].

Food quantity estimation is very important in the context of a dietary monitoring applications since on it depends the assessment of the food intakes. Works that tackle this problem are, for example, [20]–[27]. All these works require a reference information to be able to estimate the quantity of food on the plate. This information may come from markers or tokens for camera calibration, the size of reference objects, e.g., thumb or eating tools, or from the specific location where the food is consumed, e.g., canteen. Other works, instead of estimating the amount of food from 2-D images, use 3-D techniques coupled with template matching or shape reconstruction algorithms [20], [28], [29].

Very few works specifically consider the problem of left-over estimation. Often the problem is theoretically treated as a special case of the problem of food recognition and quantity estimation [18], [23]. Only one work to date explicitly tackles the problem with assessment experiments on a dedicated dataset [10].

B. Food Datasets

Regardless of the objective, a dataset of food images is required to evaluate the performance of the different feature extraction and classification algorithms proposed. To this end, the above research works either use existing datasets or introduce new ones.

Manuscript received April 13, 2016; revised August 22, 2016 and November 16, 2016; accepted November 30, 2016. Date of publication December 7, 2016; date of current version May 3, 2017.

The authors are with the Department of Informatics, System and Communication, University of Milano-Bicocca, Milano 20126, Italy (e-mail: ciocca@disco.unimib.it; paolo.napoletano@disco.unimib.it; schettini@disco.unimib.it).

Digital Object Identifier 10.1109/JBHI.2016.2636441

TABLE I
LIST OF FOOD DATASETS USED IN THE LITERATURE

Name	Year	#Images	#Classes	Type	Acquisition	Task	Annotation	Availability	Reference
Food50	2009	5000	50	Single	Wild	Food Recognition	Label	Proprietary	[30]
PFID	2009	1098 ^a	61 ^a	Single	Wild and Lab	Food Recognition	Label	Public	[63]
TADA	2009	50/256	–	Single and Multi	Lab	Food Recognition	–	Proprietary	[22]
Food85 ^b	2010	8500	85	Single	Wild	Food Recognition	Label	Proprietary	[31]
Chen	2012	5000	50	Single	Wild	Food Recognition	Label	Public	[34]
UEC FOOD-100	2012	9060	100	Single and Multi	Wild	Food Recognition	BBox	Public	[32], [64]
Food-101	2014	101 000	101	Single	Wild	Food Recognition	Label	Public	[35]
UEC FOOD-256 ^c	2014	31 397	256	Single and Multi	Wild	Food Recognition	BBox	Public	[33], [65]
UNICT-FD889	2014	3583	889	Single	Wild	Near Duplicate Food Retrieval	Label	Public	[66]
Diabetes	2014	4868	11	Single	Wild	Food Recognition	Label	Public	[5]
UNIMIB2015	2015	1000 × 2	15	Multi	Wild/Canteen	Food Recognition and Leftover Estimation	Poly	Public ^d	[10]
UNIMIB2016	2016 ^d	1027	73	Multi	Wild/Canteen	Food Recognition	Poly	Public ^d	–

^aNumbers refer to the baseline dataset.

^bIncludes Food50.

^cIncludes UECFOOD-100.

^d<http://www.ivl.disco.unimib.it/activities/food-recognition/>

One of the first food dataset was introduced in [30]. It contains 50 food categories (mostly Japanese food) and the images, gathered from the Web, depict a close-up of a single food. Using MKL-based feature fusion, they obtained a recognition accuracy of 61.34%. This dataset was enlarged to 85 food categories in a subsequent work [31]. Using a similar approach to the previous work, the authors achieved a classification accuracy of 62.85%. These two datasets are proprietary. Other proprietary datasets are the ones introduced in [4] and [22]. These datasets have been acquired in a lab settings and use markers to help the recognition phase. Differently from the previous datasets, the TADA dataset [22] contains images of real foods (256 images) as well as food replica (50 images). Also, the images can have multiple food depicted. This makes the dataset more challenging since it requires the segmentation of each food in the image. Another dataset that contains images with multiple foods is the UECFOOD-100 dataset [32]. It is public and contains more than 9000 images of 100 food categories. For the recognition, SVM classifiers with color histogram and SURF features are used, achieving a classification rate of 81.55% for the top five category candidates when the ground-truth (GT) bounding boxes are given. The dataset was extended to 256 food categories in [33] and the classification rate in this case was 74.4% for the top five categories. Chen *et al.* [34] published a dataset of 5000 images of 50 foods. Using multilabel SVMs trained on SIFT, LBP, color, and Gabor features, they achieved a food recognition overall accuracy of 68.3%.

Currently, the largest dataset available is Food-101 [35]. It contains 101 000 images divided into 101 food categories. Random forest are used to mine discriminant superpixel-grouped parts in the food images. These parts are then classified with SVM achieving an average accuracy of 50.76% on the 101 classes. If the Food-101 is the largest dataset available, the UNICT889 [7] is the dataset with most food categories. It contains 889 classes on a total of 3583 images. Given these numbers, each class contains few instances of a given food. However, the goal of the authors is the near duplicate food retrieval, and not food recognition. Different features are tested and the best

results for near duplicate retrieval was achieved by color Bag-of-Textons with a mean average precision of 67.5%.

Anthimopoulos *et al.* [5] uses a dataset of 4868 food images organized into 11 classes to evaluate a food recognition system based on the bag-of-features model. The system is designed to help diabetic patients in controlling their carbohydrates daily consumption. Different visual features and classification strategies are tested and the best combination shows a classification accuracy of slightly less than 78% using a 10 000 words dictionary.

In [10], we presented a dataset used for testing a system that recognizes foods and estimates food leftovers. The dataset contains 2000 images of 15 classes of foods placed on trays. The images were acquired in a real canteen location, and are paired with the corresponding leftover images acquired after the meals. The images are associated to a given canteen customer by using a QR code automatically generated by the dietary monitoring system on the customer's mobile. In [10], the first dataset is explicitly designed for both food recognition and leftover estimation.

Table I summarizes the characteristics of the different datasets of food images available in the literature. For each dataset, we report its size and the number of food categories it contains. The datasets have been categorized according to the type of images considered (i.e., images containing a single food or a set of foods), the acquisition procedure (e.g., in-the-wild for unconstrained acquisitions, or in-the lab for constrained acquisitions), the task for which it is used or created, the annotation type (label only, bounding boxes, or polygonal areas), and the availability (i.e., either public, or proprietary). Fig. 1 shows some examples of the images contained in each dataset.

As it can be seen from Table I and Fig. 1, most of the existing datasets depict single instance foods with only three dataset having multiple instance foods in the images. Not all the environments (and cultures) are characterized by a single food plate. For example, Asian food usually are placed in different small plates and are usually brought on the table at the same time (UECFOOD-100 is an example). Moreover, in all the canteen

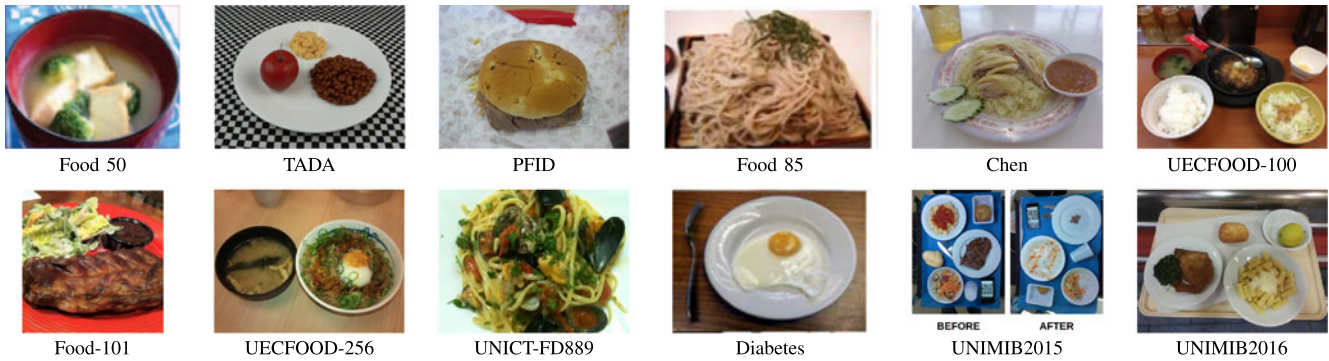


Fig. 1. Food dataset examples.

environments different plates, for the first course, main course, side dishes, and desserts, are placed on the same tray. In these cases, it is more convenient to take a single picture of the whole meal than separate pictures for each food. To date, only the UNIMIB2015 dataset is specifically designed for the canteen environments.

Canteens and cafeterias are very important in everyday life because they are often the preferred (or only) choice for workers, employees, or students. Cafeterias and canteens are receiving more and more attention with respect to customers' health and wellness. The problem of healthy and balanced meal in schools is seriously tackled by the different health agencies with the aim at reducing obesity and unbalanced nutrition. For example, the Department of Health of the Italian Government promoted an extensive campaign for food and nutrition education.¹ The Department of Health of the Australian Government, compiled a very detailed report with guidelines for healthy foods in school canteens.² Similar actions can be found across many other countries (e.g., U.K.,³ USA,⁴ etc...).

Also corporations are addressing the dietary wellness of their employees. For example, Google re-engineered its cafeterias to drive people toward healthier food choices by changing food disposition and using color coding to highlight food calories.⁵ Other corporate dining services are following a similar approach to provide healthier food and to educate their employees to a correct diet.⁶

For these reasons, we believe that datasets of food images acquired in canteen or cafeteria environments are very important for the problem of food recognition and dietary monitoring, and large and representative datasets are thus required.

In this paper, we introduce a new food dataset named UNIMIB2016. This dataset is similar to our previous dataset UNIMIB2015. Both contains images taken in a canteen environment where different foods are placed on a tray to be taken to the dining table. Differently from the UNIMIB2015 dataset,

here we have much more classes and the dishes are more difficult to locate due to the similar color of plates, tray and placemat. In fact, in UNIMIB2015, the placemat being dark blue is clearly distinguishable from the other items. In UNIMIB2016, the placemat is white as the plates. This could make it more difficult in the location and segmentation of the plates. Moreover, the higher number of food classes with respect to UNIMIB2015 makes this dataset more representative of the typical foods found in canteens. As it can be seen from Fig. 3, many food classes have a very similar appearance. For example, we have four different "Pasta al sugo," each with different ingredients added (i.e., fish, vegetables, or meat). Finally, on the tray there can be other "noisy" objects that must be ignored during the recognition. For example, we may find cell phones, wallets, id cards, and other personal items. For these reasons we need to design of a very accurate recognition algorithm.

These differences make this dataset more challenging than the previous one for the task of food recognition. Finally, as in the UNIMIB2015 dataset, here we have conducted a careful annotation of the food regions using polygonal shapes. This will allow the design of food quantity estimation algorithms using a very precise GT. However, the UNIMIB2015 dataset is the only dataset available that contains images and annotations of canteen trays taken before and after the meal (see Fig. 1), and therefore, can be used for leftover estimation. Also the two dataset are both publicly available for research purposes.

II. UNIMIB2016 FOOD DATASET

The dataset has been collected in a real canteen environment. The particularities of this setting are that each image depicts different foods on a tray, and some foods (e.g., fruit, bread, and dessert) are placed on the placemats rather than on plates. Sides are often served in the same plate as the main dish making it difficult to separate the two. Moreover, the acquisition of the images has been performed in a semicontrolled settings so the images present visual distortions as well as illumination changes due to shadows. These characteristics make this dataset challenging requiring both the segmentation of the trays for food localization, and a robust way to deal with multiple foods.

Fig. 2 shows the location where the images have been acquired. It is a canteen situated within the University of Milano-Bicocca Campus that serves students and faculty members.

¹http://www.salute.gov.it/imgs/c_17_pubblicazioni_1248_allegato.pdf

²<https://education.nt.gov.au/policies/canteen-nutrition-and-healthy-eating>

³<http://www.schoolfoodplan.com/actions/school-food-standards/>

⁴<http://www.fns.usda.gov/school-meals/child-nutrition-programs>

⁵<http://www.fastcodesign.com/1669355/6-ways-google-hacks-its-cafeterias-so-googlers-eat-healthier>

⁶<http://www.timesfreepress.com/news/business/aroundregion/story/2015/jan/20/todays-company-cafeterias-offer-healthier-brighter-fare/283592/>

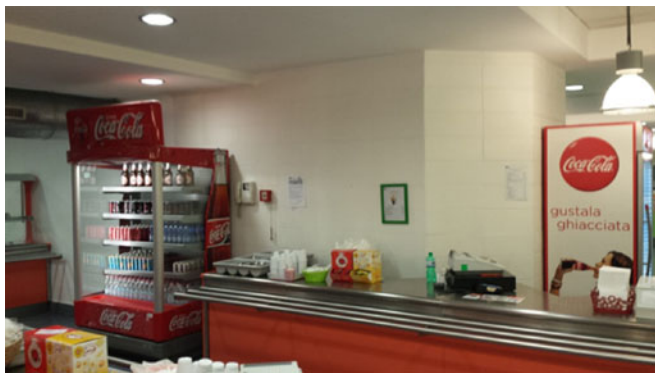


Fig. 2. Canteen situated within the University of Milano-Bicocca campus where we have acquired the images in the UNIMIB2016 dataset.

Images have been acquired using a hand-held Samsung Galaxy S3 (GT-i9300) smart phone. The acquisition station is located at the end of the tray line after the cashier. Customers place the tray on the acquisition station and the images are taken by an operator. Unfortunately, due to privacy issues and the intense affluence of customers, we have been unable to take pictures of the trays after the meal. This has prevented us to include leftover information in this dataset as we have done in the UNIMIB2015 dataset.

We have collected a total of 1442 images that went through a quality check phase where we removed excessively blurred images and duplicated photos. After this phase, we obtained a final dataset of 1027 tray images, 73 food categories, and a total of 3616 food instances. Fig. 3 shows a sample of each food category of the UNIMIB2016 dataset, while Fig. 4 shows some examples of the acquired images.

To create the GT, we have annotated the dataset using an improved version of our image annotation tool [36]–[38]. The modifications include the support of touchscreens, the drawing of freehand shapes, and the automatic approximation these shapes to polygon using the Ramer–Douglas–Peucker algorithm [39], [40]. These modification allowed us a significant speed up in the annotation process with respect to the standard point and click mouse. Fig. 4 shows some examples of annotations superimposed to the acquired images. Using our tool, to each image we have associated an annotation file containing the list of food identities, and the segmentation region of each food in terms of points of the polygon surrounding it.

Most of the existing food databases are characterized by images that contain a single food (often in a close-up setting), and in most of the cases, the food annotations are provided in terms of bounding boxes around the food (see Table I). The UNIMIB2016 dataset is characterized by images that contain multiple foods with accurate segmentation (see Fig. 4). The annotations will allow researchers to work on methods for food segmentation, as well as food quantity estimation.

III. TRAY ANALYSIS

In Fig. 5, we show the schema of our tray analysis method. The segmentator module takes the tray image as input. The output of this module is a list of regions of interest. For benchmarking,

we also consider the regions of interest obtained from the GT annotations. The regions of interest are then processed by the food class predictor. The output of the predictor is a list of recognized foods. Given a region of interest, we investigate the use of three different approaches for predicting the food class. The first approach is a global one that extracts the visual features from the whole region of interest. The second approach is a local one that extracts the visual features from local patches of the region of interest. The third approach combines the posterior probabilities computed by the global and local classifiers with the sum and product operators [41]. Given a region of interest r_i , the probability that a region is of class m is calculated in two ways:

- 1) sum rule: $P(m|r_i) = P_G(m|r_i) + P_L(m|r_i)$;
- 2) product rule: $P(m|r_i) = P_G(m|r_i) \cdot P_L(m|r_i)$

where $P_G(m|r_i)$ and $P_L(m|r_i)$ are the probability that a region of interest r_i is of class m with respect to the global and local approach, respectively. The sum rule is expected to produce reliable results when the two approaches catch information that is highly correlated, while the product rule is expected to be effective when the two approaches catch independent information.

A. Tray Segmentation

Fig. 6 shows the segmentation pipeline of the segmentator module used to detect the regions of interest. It is composed of four main steps. First, in order to speed up the computation without losing relevant information, the input RGB image is resized to a height of 320 pixels. The resized image undergoes two separate processing pipeline: a saturation-based one, and a color texture one. In the first one, the image is first gamma corrected and, then, the RGB values are converted to HSV to extract the saturation channel (see Step 1a of Fig. 6). These values are automatically thresholded and morphological operations are applied to clean up the obtained binary image (see Step 1b of Fig. 6). We have noticed that the saturation channel contain good cues for the localization of food regions since they have saturation values higher than the plate regions, the tray, and the cutlery. Of course, other regions may have saturation values comparable to those of the food and, thus, we have introduced a second processing based on the segmentation algorithm JSEG [42] that works on both color and texture features (see Step 2a of Fig. 6). We use the standard implementation of the authors with the default parameters (i.e., automatic segmentation) and we found that it works well in most cases. The segmentation is able to detect the regions having similar visual characteristics.

The segmented image is then processed in order to remove nonrelevant regions (see Step 2b of Fig. 6). For instance, the regions that touch the border of the image do not belong to the food regions, and thus, can be eliminated. Also, regions larger or smaller than predefined thresholds can be discarded as well (e.g., the placemat, the tray, highlights). The final segmented image contains with high probability the food regions and few nonrelevant ones. In order to retain only the food regions, the outputs of the saturation-based processing and the output of the color and texture processing are combined together (see Step 3 of Fig. 6). The combination performs a cross analysis be-



Fig. 3. Segmented images of the 73 food categories in the proposed UNIMIB2016 dataset. On the right, the Italian names of the classes. Note that in some cases foods slightly differ in the ingredients, and thus, are named as “FoodName 1,” “FoodName 2,” etc.



Fig. 4. Examples of acquired trays. The black polygon around the food represents the manual annotations.

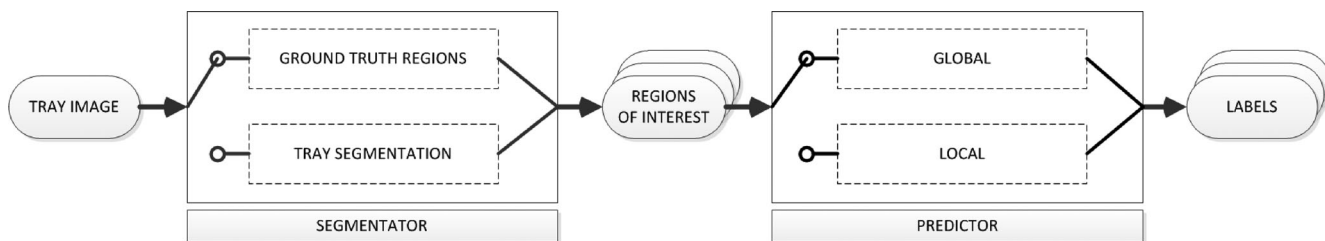


Fig. 5. Tray analysis pipeline.

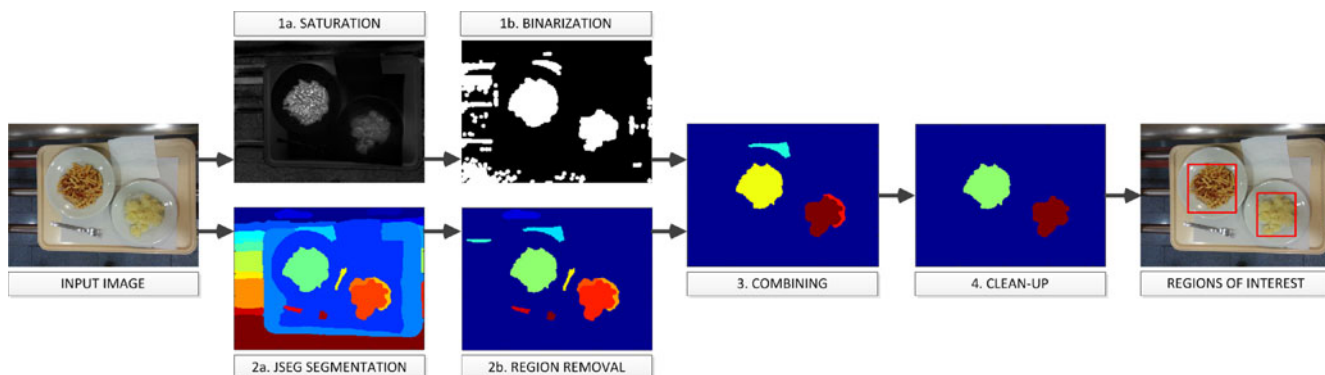


Fig. 6. Processing pipeline for the food segmentation.

TABLE II
REGION-BASED AND BOUNDARY-BASED SEGMENTATION
PERFORMANCE RESULTS

	Region-based			Boundary-based		
	Covering	PRI	VI	Recall	Precision	F-Measure
JSEG	0.385	0.389	3.106	0.870	0.198	0.323
Proposed	0.916	0.931	0.429	0.714	0.734	0.724

tween the two outputs with the aim to retain only the segmented regions that have a large percentage of saturated pixels. With this analysis, we are able to remove most of the regions of the cutlery and the spurious ones while retaining the food regions. To further ensure that only few, relevant regions are retained for the classification phase, geometric constraints are used to clean-up the output of the combining step (see step 4 of Fig. 6). The bounding boxes of all the regions of interest are passed to the prediction phase.

In order to assess the proposed segmentation pipeline, we applied the evaluation benchmarks suggested in [43]. Specifically, we computed the following region-based measures: covering of GT (Covering), the probabilistic rand index (PRI), and the variation of information (VI). Moreover, following the same work, we also computed the boundary-based precision, recall, and F1 measures. We compare the final results obtained by the proposed segmentation pipeline against the segmentation initially obtained by the JSEG algorithm. Results are reported in Table II. As it can be seen, the proposed strategy obtains the best segmentation results by all the measures considered. The region-based measures shows the highest improvements: 0.916 against 0.385, and 0.931 against 0.389 for Covering and PRI,

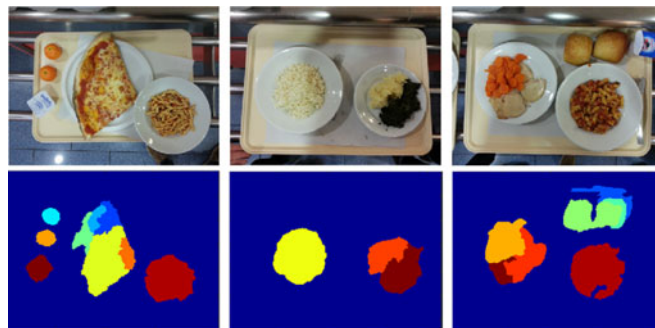


Fig. 7. Examples of segmentation results of some UNIMIB2016 images.

respectively, while the obtained VI is 0.429 against the initial 3.106 (in this case the lower the better). With respect to the boundary-based measures, we see that the initial segmentations have a high recall but with a very low precision, while the proposed one has a more balanced precision-recall values. On the overall, the proposed segmentation pipeline outperforms the JSEG one with an F-Measure of 0.724 against 0.323. The results shows that the proposed segmentation strategy is able to effectively locate the food regions.

Fig. 7 shows some results of our segmentation pipeline. As it can be seen, we are able to separate different food on the same plate. We still have some spurious regions that we hope to classify as nonfood regions in the next phase. Moreover, the JSEG algorithm often oversegments foods that shows heterogeneous regions such as the pizza slice or very textured foods such as the salads and vegetables. Each one of these regions will be independently classified. Before being passed to the classification phase, the coordinates of the bounding boxes of the

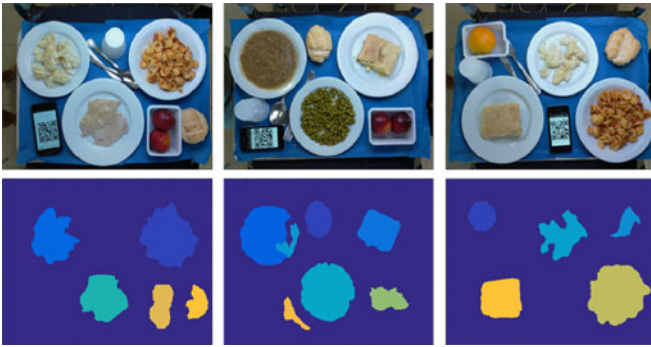


Fig. 8. Examples of segmentation results of some UNIMIB2015 images.

food regions are transformed back to match the image's original size.

To further prove the effectiveness of the proposed segmentation strategy, we have experimented it on the UNIMIB2015 database achieving a boundary-based F-Measure of 0.782 and a region-based covering of 0.945. Some UNIMIB2015 segmentation examples are shown in Fig. 8.

B. Classification of the Regions of Interest

Fig. 9 shows the processing pipeline for the food classification used in the predictor module. As we discussed earlier, here we compare three different classification strategies: a global strategy (top path in Fig. 9), a local one (bottom path), and a combination of them. The classification module works as follows. Depending on the classification strategy, from each region of interest one subimage (global strategy) or several, nonoverlapping, and image patches (local strategy) are extracted. These images are then fed to a feature extractor where several visual descriptors are computed. Specifically, we have evaluated the following visual descriptors: color histogram (HIST) [44], Gabor features (Gabor) [45], Opponent Gabor features (OG) [46], Local Color Contrast (LCC) [47], [48], chromaticity moments (CM) [46], complex wavelet features (CWT) [46], [49], color and edge directivity descriptor (CEDD) [50], nonuniform invariant LBP on the RGB channels [51], CNN [52], [53], and bag of convolutional filter responses (BoCFR) [54]–[56].

The visual descriptors are independently evaluated by pre-trained classifiers for predicting the corresponding food label. We experimented the use of two classifiers as predictor: the k-NN and SVM. The training of the classifiers is done by considering a suitable split of the UNIMIB2016 that will be described in Section IV. In the case of the local classification strategy, for each region of interest, we have several food labels, one for each image patch. Thus, a postprocessing phase to merge all these labels into a final classification decision is necessary. The local strategy is similar to the one presented in our previous work [10], and it should be useful when the food region contains part of different foods as often happens in the case of the side dishes.

IV. EXPERIMENTAL SETUP

For comparison, we evaluate the different visual features and classification strategies. In order to evaluate how much the

segmentation process influences the classification process, we also present experiments considering the ideal food segmentation provided by the GT.

A. Visual Descriptors

In this paper, we compare several visual descriptors. All feature vectors are L^2 normalized⁷:

- 1) 768-dimensional RGB [44];
- 2) 32-dimensional *Gabor* features composed of mean and standard deviation of six orientations extracted at four frequencies for each color channel [46];
- 3) 264-dimensional *OG* feature vector extracted as Gabor features from several inter/intrachannel combinations: monochrome features extracted from each channel separately and opponent features extracted from couples of colors at different frequencies [46];
- 4) 256-dimensional *LCC* feature vector. The LCC vector is obtained by comparing the color vectors at a given location with the average color in a surrounding neighborhood in terms of angular difference [47], [57];
- 5) 10-dimensional feature vector composed of normalized *CM* as defined in [46];
- 6) 8-dimensional *dual tree complex wavelet transform* (CWT) features obtained considering four scales, mean and standard deviation, and three color channels [46], [49];
- 7) 144-dimensional *CEDD* features. This descriptor uses a fuzzy version of the five digital filters proposed by the MPEG-7 edge histogram descriptor, forming six texture areas. CEDD uses two fuzzy systems that map the colors of the image in a 24-color custom palette [46], [49];
- 8) 18-dimensional *LBP* feature vector for each channel. We consider LBP applied to color images represented in RGB [58], [59]. We select the LBP with a circular neighborhood of radius 2 and 16 elements, and 18 uniform and rotation invariant patterns;
- 9) 4096-dimensional *CNNs* features (CNN4096). The CNN-based features are obtained as the intermediate representations of trained deep CNNs. The networks are used to generate a visual descriptor by removing the final softmax nonlinearity and the last fully connected layer. The network used in this paper is the BVLC AlexNet trained on ILSVRC 2012 [60];
- 10) 128-dimensional *CNNs* features (CNN128). Features are extracted in the same way as in the case of CNN4096. Here the network is the Vgg M [53] that is similar to the one presented in [61] with a reduced number of filters in the convolutional layer four. The last fully connected layer is 128-dimensional. Also this network is trained on ILSVRC 2012;
- 11) 1024-dimensional *BoCFR*: we consider the BoCFR of the first convolutional layer of the BVLC AlexNet trained for ILSVRC 2012. We built a codebook of 1024 visual words by exploiting images from external sources.

⁷The feature vector is divided by its L^2 -norm.

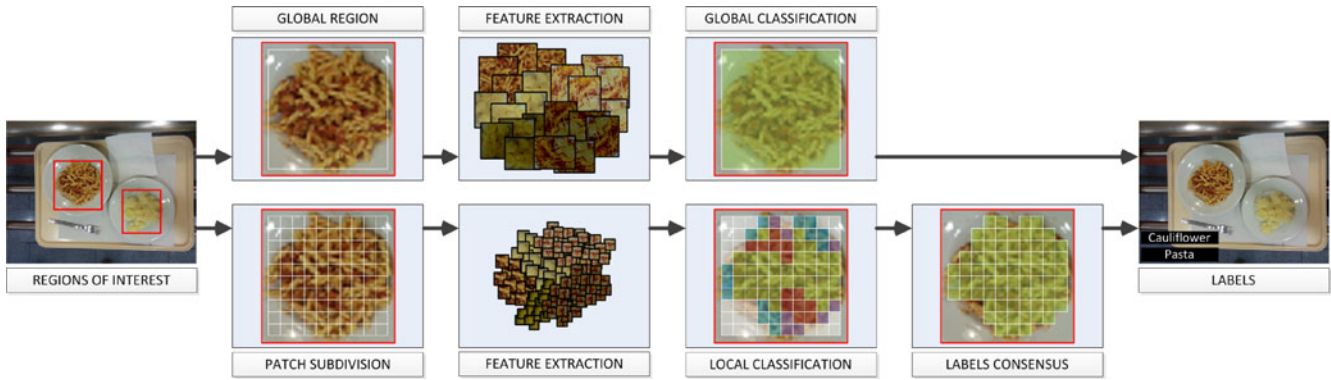


Fig. 9. Processing pipeline for the food classification.

B. Training Process

Since we collected our dataset in a real canteen scenario, and with different daily menus, the number of occurrences of each food is highly variable. This number ranges from a maximum of 479 instances for the “Pane” class down to one for some other classes (e.g., “Strudel” and “Rucola”). We have removed from the dataset the images with foods having fewer than four instances. The final dataset used in the experiments thus contains 1010 tray images and 65 foods. We split the 1010 tray images into a training set and a test set such that the sets contain about 70% and 30% of each food instances, respectively. This resulted in a training set of 650 tray images, and a test set of 360 images.

For training the global and local food classifiers, we extracted the visual descriptors from the regions of interest provided by the GT segmentation of the training trays.

Regarding the k -NN classifier, we have evaluated different values of k ranging from 1 to 11 and we have selected the value that gave the best results across visual descriptors and classification strategies, that is $k = 3$. For what concerns the SVM classifier, we have adopted the radial basis function kernel with width and regularization parameters found after a cross-validation procedure.

During the prediction process, in the case of the local classification approach, the region of interest is subdivided in patches of size 140×140 . The resulting patches may contain both food and nofood classes. This is quite clear looking at the bottom part of the Fig. 9. For both the global and local classifiers, during the training process, we added the class nofood to the classifier by choosing randomly samples from the portion of the tray images that do not overlap with the regions of interest. Once the prediction of each patch is obtained, the class with the maximum number of patches predicted is assigned to the region of interest.

C. Evaluation Measures

To cope with the class imbalance problem of the test set, we jointly used two assessment metrics for food recognition: the *standard accuracy* (SA) and the *macro average accuracy* (MAA) [62]. Denoting NP_c the number of positives, i.e., the number of times the class c occurs in the dataset; TP_c the number of *true positives* for class c , i.e., the number of times that the

system recognizes the dish c ; and C the number of classes, for each class, the metrics can be defined as follows:

$$SA = \frac{\sum_{c=1}^C TP_c}{\sum_{c=1}^C NP_c}; \quad MAA = \frac{1}{C} \sum_{c=1}^C A_c = \frac{1}{C} \sum_{c=1}^C \frac{TP_c}{NP_c}.$$

The metric for the evaluation of the error in the tray analysis is the *Tray accuracy*. This is defined as the percentage of trays correctly analyzed. A tray is correctly analyzed when all the foods contained are correctly recognized.

V. RESULTS

Results are presented in Table III. It’s quite clear that the CNN-based visual descriptors achieve better results than others in all the classification strategy. In particular, the CNN4096 features coupled with the combination of posterior probability strategy obtains the best performance. It is quite interesting to note that, apart some exceptions, the combination strategy, with both k -NN and SVM classifiers, reduces the performance with respect to the use of global and patch-based approaches. It happens in all cases when global and patch-based approaches are coupled with visual descriptors that are not good performing. It also interesting to note that the patch-based approach outperforms the global approach only when it is coupled with the SVM classifier. This is due to the fact that the radial basis function used in the SVM classifier is more suitable than the linear k -NN to separate the food classes in the feature space when the number of samples increases. Moreover, the patch-based strategy greatly outperforms the global one when coupled with traditional visual descriptors (no CNN-based). This suggest that the lower discriminant power of these features, compared to the CNN-based ones, is somewhat compensated by the larger amount of information obtained by aggregating the classification results from the local patches. For example, among the non-CNN visual descriptors, the HIST RGB combined with the local classification approach achieves a performance that is very close to some of the CNN-based descriptors. This is due to the fact that the local approach in some way takes into account the spatial variability of the food. In fact, the local approach, when applied to the UNIMIB2015 dataset, has demonstrated to be very useful for food quantity estimation [10]: the number of

TABLE III
FOOD RECOGNITION RESULTS USING THE PROPOSED TRAY ANALYSIS PIPELINE AND k -NN OR SVM CLASSIFIER

Classifier	Segment.	Approach	Measure	LBP	CEDD	Hist	Gabor	OG	LCC	CM	CWT	CNN128	CNN4096	BoCFR		
k-NN	Proposed	G	Food SA	0.343	0.423	0.555	0.397	0.463	0.320	0.439	0.276	0.656	0.728	0.689		
			Food MAA	0.139	0.184	0.356	0.168	0.253	0.127	0.210	0.079	0.467	0.585	0.490		
			Tray Accuracy	0.353	0.383	0.561	0.367	0.446	0.306	0.409	0.231	0.676	0.732	0.689		
		P	Food SA	0.488	0.594	0.689	0.597	0.667	0.492	0.608	0.624	0.679	0.697	0.697	0.689	
			Food MAA	0.202	0.315	0.474	0.318	0.443	0.201	0.326	0.387	0.453	0.473	0.490	0.490	
			Tray Accuracy	0.438	0.560	0.685	0.563	0.673	0.433	0.573	0.621	0.674	0.692	0.694	0.689	
		$G \oplus P$	Food SA	0.490	0.608	0.673	0.612	0.684	0.489	0.593	0.591	0.742	0.742	0.764	0.729	
			Food MAA	0.193	0.298	0.470	0.329	0.453	0.160	0.299	0.334	0.509	0.509	0.561	0.539	
			Tray Accuracy	0.399	0.515	0.636	0.540	0.655	0.367	0.509	0.536	0.715	0.715	0.738	0.711	
	$G \otimes P$	Food SA	0.436	0.477	0.637	0.461	0.515	0.350	0.511	0.313	0.714	0.714	0.763	0.716		
		Food MAA	0.198	0.235	0.428	0.238	0.331	0.150	0.285	0.137	0.504	0.504	0.601	0.554		
		Tray Accuracy	0.360	0.402	0.592	0.398	0.497	0.301	0.454	0.274	0.696	0.696	0.747	0.709		
	GT	G	Food SA	0.394	0.446	0.628	0.427	0.536	0.358	0.518	0.289	0.748	0.748	0.820	0.761	
			Food MAA	0.171	0.219	0.380	0.192	0.299	0.151	0.255	0.085	0.555	0.555	0.652	0.559	
			Tray Accuracy	0.434	0.492	0.662	0.470	0.570	0.408	0.534	0.313	0.783	0.783	0.842	0.782	
		P	Food SA	0.543	0.656	0.719	0.682	0.719	0.557	0.651	0.723	0.745	0.745	0.774	0.734	
			Food MAA	0.221	0.312	0.505	0.367	0.458	0.201	0.346	0.420	0.464	0.464	0.500	0.510	0.510
			Tray Accuracy	0.501	0.625	0.720	0.648	0.721	0.499	0.632	0.681	0.738	0.738	0.762	0.743	
		$G \oplus P$	Food SA	0.504	0.629	0.732	0.641	0.752	0.518	0.631	0.623	0.814	0.814	0.855	0.811	
			Food MAA	0.210	0.313	0.493	0.377	0.492	0.176	0.332	0.360	0.586	0.586	0.631	0.577	
			Tray Accuracy	0.431	0.529	0.686	0.580	0.701	0.391	0.557	0.565	0.787	0.787	0.826	0.777	
	$G \otimes P$	Food SA	0.437	0.536	0.705	0.518	0.619	0.412	0.586	0.330	0.805	0.805	0.858	0.791		
		Food MAA	0.222	0.273	0.457	0.291	0.380	0.183	0.327	0.143	0.611	0.611	0.685	0.614		
		Tray Accuracy	0.389	0.461	0.650	0.475	0.580	0.368	0.552	0.295	0.791	0.791	0.840	0.785		
SVM	Proposed	G	Food Accuracy	0.398	0.465	0.610	0.396	0.434	0.320	0.432	0.297	0.694	0.694	0.715	0.666	
			Food MAA	0.185	0.215	0.346	0.203	0.234	0.098	0.211	0.093	0.479	0.479	0.546	0.449	
			Tray Accuracy	0.440	0.440	0.575	0.394	0.403	0.313	0.408	0.255	0.703	0.703	0.738	0.669	
		P	Food SA	0.607	0.645	0.721	0.627	0.732	0.515	0.606	0.650	0.742	0.742	0.783	0.708	
			Food MAA	0.332	0.356	0.483	0.377	0.498	0.168	0.330	0.428	0.496	0.496	0.560	0.479	
			Tray Accuracy	0.585	0.605	0.705	0.630	0.729	0.421	0.570	0.655	0.720	0.720	0.767	0.708	
		$G \oplus P$	Food SA	0.640	0.628	0.703	0.670	0.713	0.382	0.612	0.646	0.777	0.777	0.798	0.702	
			Food MAA	0.387	0.399	0.452	0.446	0.518	0.100	0.261	0.453	0.616	0.616	0.632	0.465	
			Tray Accuracy	0.596	0.610	0.690	0.638	0.712	0.304	0.469	0.640	0.768	0.768	0.789	0.702	
	$G \otimes P$	Food SA	0.489	0.555	0.612	0.529	0.640	0.414	0.498	0.504	0.746	0.746	0.789	0.698		
		Food MAA	0.281	0.354	0.367	0.322	0.443	0.114	0.277	0.228	0.626	0.626	0.636	0.442		
		Tray Accuracy	0.465	0.513	0.580	0.499	0.630	0.322	0.461	0.441	0.756	0.756	0.777	0.689		
	GT	G	Food SA	0.480	0.520	0.643	0.456	0.533	0.425	0.495	0.326	0.774	0.774	0.825	0.756	
			Food MAA	0.231	0.249	0.375	0.234	0.298	0.134	0.274	0.106	0.552	0.552	0.644	0.489	
			Tray Accuracy	0.525	0.562	0.667	0.502	0.560	0.416	0.538	0.347	0.798	0.798	0.842	0.753	
		P	Food SA	0.646	0.718	0.759	0.711	0.795	0.609	0.650	0.718	0.816	0.816	0.857	0.763	
			Food MAA	0.346	0.405	0.518	0.388	0.538	0.180	0.360	0.449	0.541	0.541	0.575	0.505	
			Tray Accuracy	0.659	0.694	0.762	0.694	0.788	0.489	0.646	0.726	0.804	0.804	0.838	0.763	
		$G \oplus P$	Food SA	0.672	0.700	0.721	0.698	0.769	0.385	0.581	0.702	0.872	0.872	0.891	0.734	
			Food MAA	0.419	0.444	0.505	0.470	0.545	0.092	0.263	0.454	0.677	0.677	0.684	0.508	
			Tray Accuracy	0.641	0.658	0.723	0.665	0.745	0.279	0.459	0.670	0.845	0.845	0.871	0.702	
	$G \otimes P$	Food SA	0.565	0.619	0.642	0.576	0.711	0.418	0.528	0.551	0.816	0.816	0.858	0.722		
		Food MAA	0.322	0.370	0.434	0.359	0.471	0.125	0.310	0.248	0.670	0.670	0.687	0.557		
		Tray Accuracy	0.530	0.567	0.634	0.546	0.669	0.324	0.498	0.464	0.814	0.814	0.843	0.691		

Proposed: our automatic segmentation. GT: ground-truth segmentation. G: global approach. P: local, patch-based approach. $G \oplus P$, combination exploiting the sum of posteriors. $G \otimes P$, combination exploiting the product of posteriors. For each row, the best result is reported in bold.

patches labeled as food X suggests the quantity of the food X itself.

Overall, the SVM classifier performs slightly better than k-NN with a tray accuracy of 78.9% obtained using the sum of posteriors combination strategy. The Table III contains also the results achieved using the ideal GT as a perfect segmentation algorithm. The differences between the results obtained using the proposed segmentation pipeline and GT, allow us to evaluate

the influences of the automatic segmentation on the classification performance of the entire pipeline. It is quite clear that when the ideal segmentation is used we achieve a gain of about 10% with a maximum of 86% accuracy for the food recognition.

The best food recognition accuracy obtained by using our patch-based approach on UNIMIB2015 and measured with the SA and MAA, is 99.05% and 99.03%, respectively. Using the same patch-based approach, the best food recognition SA and

MAA on UNIMIB2016 are 78.3% and 56%, respectively. However, it should be noted that the UNIMIB2015 dataset consists of 15 food classes while the UNIMIB2016 dataset consists of 65 food classes.

VI. CONCLUSION

In recent years, it has been demonstrated that visual recognition and machine learning methods can be used to develop systems that keep tracks of human food consumption. The actual usefulness of these system heavily depends on the capability of recognizing foods in unconstrained environments. In this paper, we proposed a new dataset for the evaluation of food recognition algorithms. The images have been acquired in a real canteen and depict a real canteen tray with foods arranged in different ways. Each tray contains multiple instances of food classes. We collected a set of 1027 canteen trays for a total of 3616 food instances belonging to 73 food classes. The tray images have been manually segmented using carefully drawn polygonal boundaries. We designed a suitable automatic tray analysis pipeline that takes a tray image as input, finds the regions of interest, and predicts the corresponding food class for each region. We evaluated three different classification strategies using several visual descriptors. The best performance has been obtained by using CNNs-based features. The dataset, as well as the benchmark framework, are made available to the research community. Thanks to the way it has been annotated, this database along with the UNIMIB2015 can be used for food segmentation, recognition, and quantity estimation.

ACKNOWLEDGMENT

The authors would like to thank A. Albé for his invaluable help in the collection of the images during the acquisition campaign.

REFERENCES

- [1] W. Wu and J. Yang, "Fast food recognition from videos of eating for calorie estimation," in *Proc. IEEE Int. Conf. Multimedia Expo.*, 2009, pp. 1210–1213.
- [2] N. Yao *et al.*, "A video processing approach to the study of obesity," in *Proc. IEEE Int. Conf. Multimedia Expo.*, 2007, pp. 1727–1730.
- [3] S. Yang, M. Chen, D. Pomerleau, and R. Sukthankar, "Food recognition using statistics of pairwise local features," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2010, pp. 2249–2256.
- [4] M. Bosch, F. Zhu, N. Khanna, C. Boushey, and E. Delp, "Combining global and local features for food identification in dietary assessment," in *Proc. 18th IEEE Int. Conf. Image Process.*, 2011, pp. 1789–1792.
- [5] M. M. Anthimopoulos, L. Gianola, L. Scarnato, P. Diem, and S. G. Mougiakakou, "A food recognition system for diabetic patients based on an optimized bag-of-features model," *IEEE J. Biomed. Health Informat.*, vol. 18, no. 4, pp. 1261–1271, Jul. 2014.
- [6] Y. He, C. Xu, N. Khanna, C. Boushey, and E. Delp, "Analysis of food images: Features and classification," in *Proc. IEEE Int. Conf. Image Process.*, 2014, pp. 2744–2748.
- [7] G. Farinella, M. Moltisanti, and S. Battiato, "Classifying food images represented as bag of textons," in *Proc. IEEE Int. Conf. Image Process.*, 2014, pp. 5212–5216.
- [8] D. T. Nguyen, Z. Zong, P. O. Ogunbona, Y. Probst, and W. Li, "Food image classification using local appearance and global structural information," *Neurocomputing*, vol. 140, pp. 242–251, 2014.
- [9] V. Bettadapura, E. Thomaz, A. Parnami, G. Abowd, and I. Essa, "Leveraging context to support automated food recognition in restaurants," in *Proc. IEEE Winter Conf. Appl. Comput. Vis.*, 2015, pp. 580–587.
- [10] G. Ciocca, P. Napoletano, and R. Schettini, "Food recognition and leftover estimation for daily diet monitoring," in *Proc. New Trends Image Anal. Process. Workshops*, 2015, vol. 9281, pp. 334–341.
- [11] H. Kagaya, K. Aizawa, and M. Ogawa, "Food detection and recognition using convolutional neural network," in *Proc. ACM Int. Conf. Multimedia*, 2014, pp. 1085–1088.
- [12] Y. Kawano and K. Yanai, "Food image recognition with deep convolutional features," in *Proc. ACM Int. Joint Conf. Pervasive Ubiquitous Comput.*, 2014, pp. 589–593.
- [13] W. Zhang, D. Zhao, W. Gong, Z. Li, Q. Lu, and S. Yang, "Food image recognition with convolutional neural networks," in *Proc. IEEE 12th Int. Conf. Ubiquitous Intell. Comput. IEEE 12th Int. Conf. Auton. Trusted Comput. IEEE 15th Int. Conf. Scalable Comput. Commun. Assoc. Workshops*, 2015, pp. 690–693.
- [14] K. Kitamura, T. Yamasaki, and K. Aizawa, "Foodlog: Capture, analysis and retrieval of personal food images via web," in *Proc. ACM Multimedia Workshop Multimedia Cooking Eating Act.*, 2009, pp. 23–30.
- [15] F. Kong and J. Tan, "Dietcam: Automatic dietary assessment with mobile camera phones," *Pervasive Mobile Comput.*, vol. 8, no. 1, pp. 147–163, 2012.
- [16] O. Bejbom, N. Joshi, D. Morris, S. Saponas, and S. Khullar, "Menu-match: Restaurant-specific food logging from images," in *Proc. IEEE Winter Conf. Appl. Comput. Vis.*, 2015, pp. 844–851.
- [17] Y. Kawano and K. Yanai, "Foodcam: A real-time food recognition system on a smartphone," *Multimedia Tools Appl.*, vol. 74, pp. 5263–5287, 2015.
- [18] F. Zhu *et al.*, "The use of mobile devices in aiding dietary assessment and evaluation," *IEEE J. Sel. Topics Signal Process.*, vol. 4, no. 4, pp. 756–766, Aug. 2010.
- [19] Z. Ahmad, N. Khanna, D. A. Kerr, C. J. Boushey, and E. J. Delp, "A mobile phone user interface for image-based dietary assessment," *Proc. IS&T/SPIE*, vol. 9030, pp. 903007–903007, 2014.
- [20] M. Puri, Z. Zhu, Q. Yu, A. Divakaran, and H. Sawhney, "Recognition and volume estimation of food intake using a mobile device," in *Proc. Workshop Appl. Comput. Vis.*, 2009, pp. 1–8.
- [21] M. Sun *et al.*, "Determination of food portion size by image processing," in *Proc. 30th Annu. Int. Conf. Eng. Med. Biol. Soc.*, 2008, pp. 871–874.
- [22] A. Mariappan *et al.*, "Personal dietary assessment using mobile devices," in *IS&T/SPIE Electron. Imag.*, International Society for Optics and Photonics, vol. 7246, 2009, pp. 72460Z–72460Z.
- [23] P. Pouladzadeh, S. Shirmohammadi, and R. Al-Maghrabi, "Measuring calorie and nutrition from food image," *IEEE Trans. Instrum. Meas.*, vol. 63, no. 8, pp. 1947–1956, Aug. 2014.
- [24] P. Pouladzadeh, G. Villalobos, R. Almaghrabi, and S. Shirmohammadi, "A novel SVM based food recognition method for calorie measurement applications," in *Proc. IEEE Int. Conf. Multimedia Expo Workshops*, 2012, pp. 495–498.
- [25] G. Villalobos, R. Almaghrabi, P. Pouladzadeh, and S. Shirmohammadi, "An image processing approach for calorie intake measurement," in *Proc. IEEE Int. Symp. Med. Meas. Appl.*, 2012, pp. 1–5.
- [26] E. A. Akpro Hippocrate, H. Suwa, Y. Arakawa, and K. Yasumoto, "Food weight estimation using smartphone and cutlery," in *Proc. 1st Workshop IoT-Enabled Healthcare Wellness Technol. Syst.*, 2016, pp. 9–14.
- [27] P. Pouladzadeh, P. Kuhad, S. V. B. Peddi, A. Yassine, and S. Shirmohammadi, "Food calorie measurement using deep learning neural network," in *Proc. IEEE Int. Instrum. Meas. Technol. Conf.*, 2016, pp. 1–6.
- [28] J. Chae *et al.*, "Volume estimation using food specific shape templates in mobile image-based dietary assessment," in *IS&T/SPIE Electron. Imag.*, International Society for Optics and Photonics, vol. 7873, pp. 78730K–78730K, 2011.
- [29] Y. He, C. Xu, N. Khanna, C. Boushey, and E. Delp, "Food image analysis: Segmentation, identification and weight estimation," in *Proc. IEEE Int. Conf. Multimedia Expo*, 2013, pp. 1–6.
- [30] T. Joutou and K. Yanai, "A food image recognition system with multiple kernel learning," in *Proc. 16th IEEE Int. Conf. Image Process.*, 2009, pp. 285–288.
- [31] H. Hoashi, T. Joutou, and K. Yanai, "Image recognition of 85 food categories by feature fusion," in *Proc. IEEE Int. Symp. Multimedia*, 2010, pp. 296–301.
- [32] Y. Kawano and K. Yanai, "Real-time mobile food recognition system," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog. Workshops*, 2013, pp. 1–7.
- [33] Y. Kawano and K. Yanai, "Automatic expansion of a food image dataset leveraging existing categories with domain adaptation," in *Proc. ECCV Workshop Transferring Adapting Source Knowl. Comput. Vis.*, 2014, pp. 3–17.

- [34] M.-Y. Chen *et al.*, "Automatic chinese food identification and quantity estimation," in *Proc. SIGGRAPH Asia Tech. Briefs*, 2012, Art. no. 29.
- [35] L. Bossard, M. Guillaumin, and L. Van Gool, "Food-101—mining discriminative components with random forests," in *Proc. Comput. Vis.*, 2014, pp. 446–461.
- [36] G. Ciocca, P. Napoletano, and R. Schettini, "IAT—image annotation tool: Manual," 2015, to be published, arXiv:1502.05212.
- [37] S. Bianco, G. Ciocca, P. Napoletano, and R. Schettini, "An interactive tool for manual, semi-automatic and automatic video annotation," *Comput. Vis. Image Underst.*, vol. 131, pp. 88–99, 2015.
- [38] S. Bianco *et al.*, "Cooking action recognition with IVAT: An interactive video annotation tool," in *Proc. Int. Conf. Image Anal. Process.*, 2013, pp. 631–641.
- [39] U. Ramer, "An iterative procedure for the polygonal approximation of plane curves," *Comput. Graphic Image Process.*, vol. 1, no. 3, pp. 244–256, 1972.
- [40] D. H. Douglas and T. K. Peucker, "Algorithms for the reduction of the number of points required to represent a digitized line or its caricature," *Cartographica: Int. J. Geographic Inf. Geovisualization*, vol. 10, no. 2, pp. 112–122, 1973.
- [41] J. Kittler, M. Hatef, R. P. Duin, and J. Matas, "On combining classifiers," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 20, no. 3, pp. 226–239, Mar. 1998.
- [42] Y. Deng and B. S. Manjunath, "Unsupervised segmentation of color-texture regions in images and video," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 23, no. 8, pp. 800–810, Aug. 2001.
- [43] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, "Contour detection and hierarchical image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 5, pp. 898–916, May 2011.
- [44] C. L. Novak *et al.*, "Anatomy of a color histogram," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recog.*, 1992, pp. 599–605.
- [45] B. S. Manjunath and W.-Y. Ma, "Texture features for browsing and retrieval of image data," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 18, no. 8, pp. 837–842, Aug. 1996.
- [46] F. Bianconi, R. Harvey, P. Southam, and A. Fernández, "Theoretical and experimental comparison of different approaches for color texture classification," *J. Electron. Imaging*, vol. 20, no. 4, 2011, Art. no. 043006.
- [47] C. Cusano, P. Napoletano, and R. Schettini, "Combining local binary patterns and local color contrast for texture classification under varying illumination," *J. Opt. Soc. Amer. A*, vol. 31, no. 7, pp. 1453–1461, 2014.
- [48] C. Cusano, P. Napoletano, and R. Schettini, "Intensity and color descriptors for texture classification," in *Proc. IS&T/SPIE*, vol. 8661, 2013.
- [49] M. Barilla and M. Spann, "Colour-based texture image classification using the complex wavelet transform," in *Proc. 5th Int. Conf. Electr. Eng. Comput. Sci. Autom. Control*, Nov. 2008, pp. 358–363.
- [50] S. A. Chatzichristofis and Y. S. Boutalis, "Cedd: Color and edge directivity descriptor: A compact descriptor for image indexing and retrieval," in *Proc. Comput. Vis. Syst.*, 2008, pp. 312–322.
- [51] T. Ojala, M. Pietikäinen, and T. Mäenpää, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 7, pp. 971–987, Jul. 2002.
- [52] A. S. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "Cnn features off-the-shelf: An astounding baseline for recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, 2014, pp. 512–519.
- [53] A. Vedaldi and K. Lenc, "Matconvnet: Convolutional neural networks for MATLAB," in *Proc. 23rd ACM int. conf. Multimedia*, ACM, Oct. 2015, pp. 689–692.
- [54] J. Ng, F. Yang, and L. Davis, "Exploiting local features from deep networks for image retrieval," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, 2015, pp. 53–61.
- [55] J. Sivic and A. Zisserman, "Video google: A text retrieval approach to object matching in videos," in *Proc. 9th IEEE Int. Conf. Comput. Vis.*, 2003, pp. 1470–1477.
- [56] Y. Yang and S. Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," in *Proc. Int. Conf. Adv. Geographic Inf. Syst.*, 2010, pp. 270–279.
- [57] C. Cusano, P. Napoletano, and R. Schettini, "Evaluating color texture descriptors under large variations of controlled lighting conditions," *J. Opt. Soc. Amer. A*, vol. 33, no. 1, pp. 17–30, 2016.
- [58] T. Mäenpää and M. Pietikäinen, "Classification with color and texture: Jointly or separately?" *Pattern Recognit.*, vol. 37, no. 8, pp. 1629–1640, 2004.
- [59] G. Ciocca, C. Cusano, and R. Schettini, "Image orientation detection using lbp-based features and logistic regression," *Multimedia Tools Appl.*, vol. 74, no. 9, pp. 3013–3034, 2015.
- [60] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2012, pp. 1097–1105.
- [61] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 818–833.
- [62] H. He and Y. Ma, *Imbalanced Learning: Foundations, Algorithms, and Applications*. New York, NY, USA: Wiley, 2013.
- [63] M. Chen, K. Dhingra, W. Wu, L. Yang, R. Sukthankar, and J. Yang, "Pfid: Pittsburgh fast-food image dataset," in *Proc. 16th IEEE Int. Conf. Image Process.*, 2009, pp. 289–292.
- [64] Y. Matsuda, H. Hoashi, and K. Yanai, "Recognition of multiple-food images by detecting candidate regions," in *Proc. IEEE Int. Conf. Multimedia Expo*, 2012, pp. 25–30.
- [65] Y. Kawano and K. Yanai, "Foodcam-256: A large-scale real-time mobile food recognitionsystem employing high-dimensional features and compression of classifier weights," in *Proc. ACM Int. Conf. Multimedia*, 2014, pp. 761–762.
- [66] G. M. Farinella, D. Allegra, and F. Stanco, "A benchmark dataset to study the representation of food images," in *Proc. ECCV Eur. Conf. Comput. Vis. Workshop Assist. Comput. Vis. Robot.*, 2014, pp. 584–599.

Gianluigi Ciocca received the M.Sc. degree in computer science from the University of Milano, Milano, Italy, in 1998. Since 2003, he has been with the Department of Informatics, Systems and Communication (DISCo), University of Milano-Bicocca, where he received the Ph.D. degree in computer science in 2006.

He is currently an Associate Professor of computer science at DISCo. His current research interests include image and video analysis, pattern recognition, classification.

He is a member of NeuroMi: Milan Center for Neuroscience and a Fellow at the Institute of Multimedia Information Technologies of the Italian National Research Council, where his research focused on the development of systems for the management of image and video databases, and the development of new methodologies and algorithms for automatic indexing.

Paolo Napoletano received the Master's degree in telecommunications engineering from the University of Naples Federico II, Naples, Italy, with a thesis focused on transmission of electromagnetic fields, in 2003. In 2007, he received the Ph.D. degree in information engineering from the University of Salerno, Fisciano, Italy, with a thesis focused on computational vision and pattern recognition.

He is currently a Postdoc Researcher in the Department of Informatics, Systems, and Communication, University of Milano-Bicocca, Milano, Italy. His current research interests include image and video analysis and information retrieval.

Raimondo Schettini obtained the Laurea degree in physics from the University of Milano, Milano, Italy, in 1986.

He is currently a Full Professor at the University of Milano-Bicocca, Milano, Italy. He is the Vice-Director of the Department of Informatics, Systems, and Communication, and the Head of Imaging and Vision Laboratory (www.ivl.disco.unimib.it). Since 1987, he has been associated with the Italian National Research Council, where he has led the Color Imaging Lab from 1990 to 2002. He has been a Team Leader in several research projects and published more than 300 refereed papers and several patents about color reproduction, image processing, analysis, and classification.

Mr. Schettini is a Fellow of the International Association of Pattern Recognition for his contributions to pattern recognition research and color image analysis.