

# TOWARDS EMBEDDED WASTE SORTING

## *Using Constellations of Visual Words*

Toon Goedemé

*De Nayer Technical University, Embedded System Design (EmSD)*

*Jan De Nayerlaan 5, 2860 Sint-Katelijne-Waver, Belgium*

*Katholieke Universiteit Leuven, VISICS, ESAT/PSI, Kasteelpark Arenberg 10, 3001 Heverlee, Belgium*

**Keywords:** Waste sorting, local image features, SURF, SIFT, visual words.

**Abstract:** In this paper, we present a method for fast and robust object recognition, especially developed for implementation on an embedded platform. As an example, the method is applied to the automatic sorting of consumer waste. Out of a stream of different thrown-away food packages, specific items — in this case beverage cartons — can be visually recognised and sorted out. To facilitate and optimise the implementation of this algorithm on an embedded platform containing parallel hardware, we developed a voting scheme for constellations of visual words, i.e. clustered local features (SURF in this case). On top of easy implementation and robust and fast performance, even with large databases, an extra advantage is that this method can handle multiple identical visual features in one model.

## 1 INTRODUCTION

We do not live in a world with unlimited resources, therefore the principle of the *TetraPak* company is '*a package should save more than it costs*'. One key issue in their recycling process is sorting the beverage carton fraction out of the consumer waste stream. Although sometimes beverage cartons are separately collected, at most places a mixed 'recyclable' fraction is separately collected, which has to be sorted out afterwards. Sorting out some subfractions is easy, e.g. by using magnets for ferrometals. Some other subfractions are less easily automated and have to be sorted manually. This is the case with beverage cartons also. In waste processing plants, people have to pick out the beverage cartons from a stinking never-ending stream of waste on conveyor belts ...

Although techniques such as the measurement of UV light reflection can help the automated sorting process, we present in this work a reliable visual method. The system's input consists of images from a camera which is placed above the conveyor belt. These images are rapidly matched with a database of beverage carton photos. In real-time, a large fraction of all beverage cartons can be identified and picked out. Missing items in the database can be quickly added, on the basis of a photograph of the beverage carton.

The remainder of this text is organised as follows.

Section 2 gives an overview of relevant related work. In section 3, our algorithm is described. Some *real-waste* experiments are presented in section 4. The paper ends with a conclusion in section 5.

## 2 RELATED WORK

Since long, general object recognition is one of the core research subjects in computer vision. Numerous techniques are proposed, traditionally mainly based on the template matching technique (Rosenfeld and Kak, 1976). A few years ago, a major revolution in the field was the appearance of the idea of local image features (Tuytelaars et al., 1999; Lowe, 1999). Indeed, looking at local parts instead of the entire pattern to be recognised has the inherent advantage of robustness to partial occlusions. In both template and query image, local regions are extracted around interest points, each described by a descriptor vector for comparison. The development of robust local feature descriptors, like e.g. Mindru's generalised colour moment based ones (Mindru et al., 1999), added robustness to illumination and changes in viewpoint.

Many researchers proposed algorithms for local region matching. The differences between approaches lie in the way in which interest points, local image regions, and descriptor vectors are extracted.

An early example is the work of Schmid and Mohr (Schmid et al., 1997), where geometric invariance was still under image rotations only. Scaling was handled by using circular regions of several sizes. Lowe et al. (Lowe, 1999) extended these ideas to real scale-invariance. More general affine invariance has been achieved in the work of Baumberg (Baumberg, 2000), that uses an iterative scheme and the combination of multiple scales, and in the more direct, constructive methods of Tuytelaars & Van Gool (Tuytelaars et al., 1999; Tuytelaars and Gool, 2000), Matas et al. (Matas et al., 2002), and Mikolajczyk & Schmid (Mikolajczyk and Schmid, 2002). Although these methods are capable to find very qualitative correspondences, most of them are too slow for use in a real-time application as the one we envision here. Moreover, none of these methods are especially suited for the implementation on an embedded computing system, where both memory and computing power must be as low as possible to ensure reliable operation at the lowest cost possible.

The classic recognition scheme with local features, presented in (Lowe, 1999; Tuytelaars and Gool, 2000), and used in many applications such as in our previous work on robot navigation (Goedemé et al., 2005; Goedemé et al., 2006), is based on finding one-on-one matches. Between the query image and a model image of the object to be recognised, bijective matches are found. For each local feature of the one image, the most similar feature in the other is selected.

This scheme contains a fundamental drawback, namely its disability to detect matches when multiple identical features are present in an image. In that case, no guarantee can be given that the most similar feature is the correct correspondence. Such pattern repetitions are quite common in the real world, though, especially in man-made environments. To reduce the number of incorrect matches due to this phenomenon, in classic matching techniques a criterium is used such as comparing the distance to the most and the second most similar feature (Lowe, 1999). Of course, this practice throws away a lot of good matches in the presence of pattern repetitions.

In this paper, we present a possible solution to this problem by making use of the *visual word* concept. Visual words are introduced (Sivic and Zisserman, 2003; Li and Perona, 2005; Zhang and Schmid, 2005) in the context of object classification. Local features are grouped into a large number of clusters with those with similar descriptors assigned into the same cluster. By treating each cluster as a *visual word* that represents the specific local pattern shared by the keypoints in that cluster, we have a visual word vocabu-

lary describing all kinds of such local image patterns. With its local features mapped into visual words, an image can be represented as a *bag of visual words*, as a vector containing the (weighted) count of each visual word in that image, which is used as feature vector in the classification task.

In contrast to the in categorisation often used bag-of-words concept, in this paper we present the *constellation-of-words* model. The main difference is that not only the presence of a number of visual words is tested, but also their relative positions.

### 3 ALGORITHM

Figure 1 gives an overview of the algorithm. It consists of two phases, namely the model construction phase (upper row) and the matching phase (bottom row).

First, in a model photograph (*a*), local features are extracted (*b*). Then, a vocabulary of visual words is formed by clustering these features based on their descriptor. The corresponding visual words on the image (*c*) are used to form the model description. The relative location of the image centre (the *anchor*) is stored for each visual word instance (*d*).

The bottom row depicts the matching procedure. In a query image, local features are extracted (*e*). Matching with the vocabulary yields a set of visual words (*f*). For each visual word in the model description, a vote is cast at the relative location of the anchor location (*g*). The location of the object can be found based on these votes as local maxima in a voting Hough space (*h*). Each of the following subsections describes one step of this algorithm in detail.

**Local Feature Extraction.** We chose to use SURF as local feature detector, instead of the often used SIFT detector. SURF (Bay et al., 2006; Fasel and Gool, 2007) is developed to be substantially faster, but at least as performant as SIFT.

**Interest Point Detector.** In contrast to SIFT (Lowe, 1999), which approximates Laplacian of Gaussian (LoG) with Difference of Gaussians (DoG), SURF approximates second order Gaussian derivatives with box filters, see figure 2. Image convolutions with these box filters can be computed rapidly by using integral images as defined in (Viola and Jones, 2001). Interest points are localised in scale and image space by applying a non-maximum suppression in a  $3 \times 3$  neighbourhood. Finally, the found maxima of the determinant of the approximated Hessian matrix are interpolated in scale and image space.

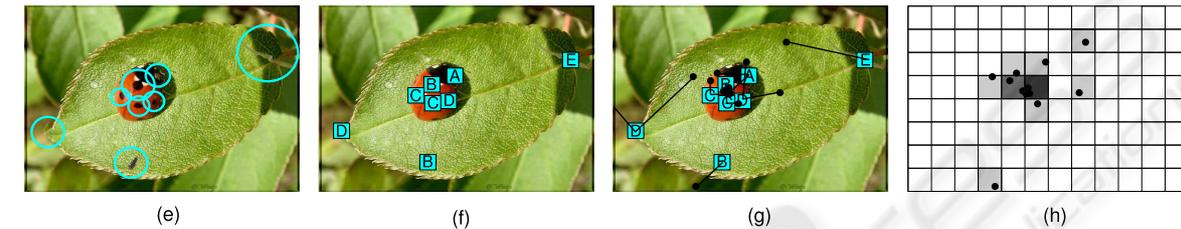
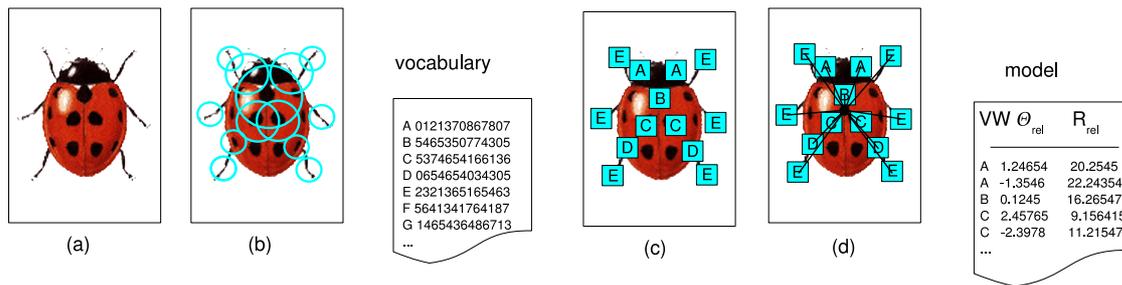


Figure 1: Overview of the algorithm. Top row (model building): (a) model photo, (b) extracted local features, (c) features expressed as visual words from the vocabulary, (d) model description with relative anchor positions for each visual word. Bottom row (matching): (e) query image with extracted features, (f) visual words from the vocabulary, (g) anchor position voting based on relative anchor position, (h) Hough voting space.



Figure 2: Left: two filters based on Gaussian derivatives. Right: their approximation using box filters.

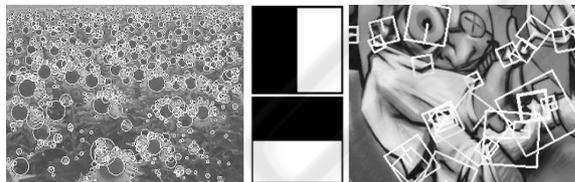


Figure 3: Middle: Haar wavelets. Left and right: examples of extracted SURF features.

**Descriptor.** In a first step, SURF constructs a circular region around the detected interest points in order to assign a unique orientation to the former and thus gain invariance to image rotations. The orientation is computed using Haar wavelet responses in both x and y direction as shown in the middle of figure 3. The Haar wavelets can be easily computed via integral images, similar to the Gaussian second order approximated box filters. Once the Haar wavelet responses are computed, they are weighted with a Gaussian centred at the interest points. In a next step the dominant orientation is estimated by summing the horizontal and vertical wavelet responses within a rotating wedge, covering an angle of  $\frac{\pi}{3}$  in the wavelet re-

sponse space. The resulting maximum is then chosen to describe the orientation of the interest point descriptor. In a second step, the SURF descriptors are constructed by extracting square regions around the interest points. These are oriented in the directions assigned in the previous step. Some example windows are shown on the right hand side of figure 3. The windows are split up in  $4 \times 4$  sub-regions in order to retain some spatial information. In each sub-region, Haar wavelets are extracted at regularly spaced sample points. In order to increase robustness to geometric deformations and localisation errors, the responses of the Haar wavelets are weighted with a Gaussian, centred at the interest point. Finally, the wavelet responses in horizontal  $d_x$  and vertical directions  $d_y$  are summed up over each sub-region. Furthermore, the absolute values  $|d_x|$  and  $|d_y|$  are summed in order to obtain information about the polarity of the image intensity changes. The resulting descriptor vector for all  $4 \times 4$  sub-regions is of length 64. See figure 4 for an illustration of the SURF descriptor for three different image intensity patterns. More details about SURF can be found in (Bay et al., 2006) and (Fasel and Gool, 2007).

**Visual Words.** As explained before, the next step is forming a vocabulary of visual words. This is accomplished by clustering a big set of extracted SURF features. It is important to build this vocabulary using a large number of features, in order to be representa-

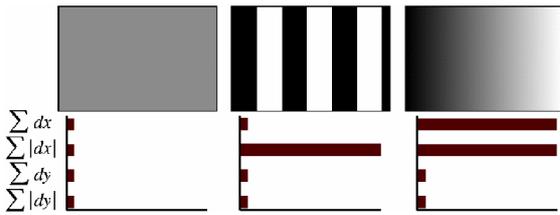


Figure 4: Illustrating the SURF descriptor.

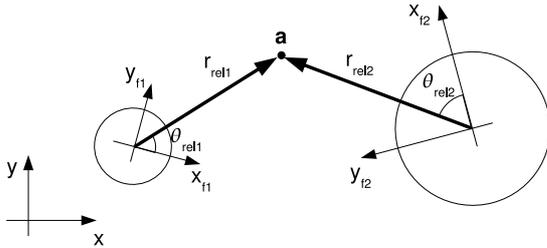


Figure 5: The position of the anchor point is stored in the model as polar coordinates relative to the visual word scale and orientation.

tive for all images to be processed.

The clustering itself is easily carried out with the *k*-means algorithm. Distances between features are computed as the Euclidean distance between the corresponding SURF descriptors. Keep in mind that this model-building phase can be processed off-line, the real-time behaviour is only needed in the matching step.

In the fictive ladybug example of figure 1, each visual word is symbolically presented as a letter. It can be seen that the vocabulary exists of a file linking each visual word symbol with a mean descriptor vector of the corresponding cluster.

### 3.1 Model Construction

All features found on a model image are matched with the visual word vocabulary, as shown in fig. 1 (c). In addition to the popular bag-of-words models, which consist of a set of visual words, we add the relative constellation of all visual words to the model description.

Each line in the model description file consists of the symbolic name of a visual word, and the relative coordinates  $(r_{rel}, \theta_{rel})$  to the anchor point of the model item. As anchor point, we chose for instance the centre of the model picture. These coordinates are expressed as polar coordinates, relative to the individual axis frame of the visual word. Indeed, each visual word in the model photograph has a scale and an orientation because it is extracted as a SURF feature. Figure 5 illustrates this. The resulting model is a very compact description of the appearance of the

model photo. Many of these models, based on the same visual word vocabulary, can be saved in a compact database. In our beverage carton sorting application, we build a database of all different carton prints to be recognised.

### 3.2 Matching

Once a database of objects to be recognised is built, these objects can be detected in a query image. In our application, a camera overviews a section of the conveyor belt. The object detection algorithm here described gives cues where beverage cartons are located. With this information, a mechanical device can sort out the beverage cartons.

This part of the algorithm is time-critical. We are spending lots of efforts in speeding up the matching procedure, in order to be able to implement it on an embedded system.

The first operation carried out on incoming images is extracting SURF features, exactly as described in section 3. After local feature extraction, matching is performed with the visual words in the vocabulary. We used Mount's ANN (Approximate Nearest Neighbour) (Arya et al., 1998) algorithm for this, which is very performant. As seen in fig. 1 (f), some of the visual words of the object are recognised, amidst other visual words.

**Anchor Location Voting.** Because each SURF feature has a certain scale and rotation, we can reconstruct the anchor pixel location by using the feature-relative polar coordinates of the object anchor. For each instance in the object model description, this yields a vote for a certain anchor location. In figure 1 (g), this is depicted by the black lines ending with a black dot at the computed anchor location.

Ideally, all these locations would coincide at the correct object centre. Unfortunately, this is not the case due to mismatches and noise. Moreover, if there are two identical visual words in the model description of an object (as is the case in the ladybug example for words *A*, *C* and *D*), each detected visual word of that kind in the query image will cast to different anchor location votes, of which only one can be correct.

**Object Detection.** For all different models in the database, anchor location votes can be quickly computed. Next task is to decide where a certain object is detected. Because a certain object can be present more than once in the query image, it is clear that a simple average of the anchor position votes is not a sufficient technique, even if robust estimators like

RANSAC are used to eliminate outliers. Therefore, we construct a *Hough space*, a matrix which is initiated at zero and incremented at each anchor location vote, fig. 1 (*h*). The local maxima of the resulting Hough matrix are computed and interpreted as detected object positions.

## 4 EXPERIMENTS

For preliminary experiments, we implemented this algorithm using Octave and an executable of the SURF extractor. Figure 6 shows some typical results of different phases of the algorithm. The test images were made by pouring out a 'recyclable fraction' garbage bag and taking  $640 \times 480$  photographs of it from about 1 meter distance.

In fig. 6, first two model photographs are shown, for two types of beverage cartons. Each of such images, having a resolution of about  $100 \times 150$  pixels, yielded a thorough description of the carton print in a model description containing on the average 65 features, what boils down to a model file size of only 3.5 KB.

In the middle of the top row, the anchor position voting output is shown for the milk carton detection step. From matched visual words, black lines are drawn towards the anchor position. It is clearly visible that many lines point at the centres of both milk cartons. In the Hough voting space, next to it, this leads to two black spots at the positions of the milk cartons. The bottom row shows comparable experimental results for other query and model images.

The cartons were detected by finding local maxima in the Hough space. We performed experiments on 25 query images, containing in total 189 milk cartons. We were able to detect 84% of the trained types. Detection failures were mostly due to a large occlusion of one carton by another object.

## 5 CONCLUSIONS AND FUTURE WORK

In this paper, we presented an algorithm for object detection based on the concept of visual word constellation voting. The preliminary experiments proved the performance of this approach. The method has the advantages that it is computing-power and memory efficient and that it can handle pattern repetitions in the models.

We applied this method on the vision-based sorting process of consumer waste, by detecting the

beverage cartons based on a database of previously trained beverage carton prints.

As told before, our aim in this work is an embedded implementation of this algorithm. The Octave implementation presented here is only a first step towards that. But we believe the proposed approach has a lot of advantages. The SURF extraction phase can mostly be migrated to a parallel hardware implementation on FPGA. Visual word matching is sped up using the ANN-libraries, making use of Kd-trees. Of course a large part of the memory is used by the (mostly sparse) hough space. A better description of the voting space will lead to a great memory improvement of the algorithm.

## REFERENCES

- Arya, S., Mount, D., Netanyahu, N., Silverman, R., and Wu, A. (1998). An optimal algorithm for approximate nearest neighbor searching. In *J. of the ACM*, vol. 45, pp. 891-923.
- Baumberg, A. (2000). Reliable feature matching across widely separated views. In *Computer Vision and Pattern Recognition, Hilton Head, South Carolina*, pp. 774-781.
- Bay, H., Tuytelaars, T., and Gool, L. V. (2006). Speeded up robust features. In *ECCV*.
- Fasel, B. and Gool, L. V. (2007). Interactive museum guide: Accurate retrieval of object descriptions. In *Adaptive Multimedia Retrieval: User, Context, and Feedback, Lecture Notes in Computer Science, Springer, volume 4398*.
- Goedemé, T., Nuttin, M., Tuytelaars, T., and Gool, L. V. (2006). Omnidirectional vision based topological navigation. In *International Journal of Computer Vision and International Journal of Robotics Research, Special Issue: Joint Issue of IJCV and IJRR on Vision and Robotics*.
- Goedemé, T., Tuytelaars, T., Vanacker, G., Nuttin, M., and Gool, L. V. (2005). Feature based omnidirectional sparse visual path following. In *IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2005, pp. 1003-1008, Edmonton*.
- Li, F.-F. and Perona, P. (2005). A bayesian hierarchical model for learning natural scene categories. In *Proc. of the 2005 IEEE Computer Society Conf. on Computer Vision and Pattern Recognition*, pages 524531.
- Lowe, D. (1999). Object recognition from local scale-invariant features. In *International Conference on Computer Vision*.
- Matas, J., Chum, O., Urban, M., and Pajdla, T. (2002). Robust wide baseline stereo from maximally stable extremal regions. In *British Machine Vision Conference, Cardiff, Wales, pp. 384-396*.
- Mikolajczyk, K. and Schmid, C. (2002). An affine invariant interest point detector. In *ECCV, vol. 1, 128-142*.



Figure 6: Some experimental results. Top row: model photos of milk and juice cartons, query image with matching visual words (white) and relative anchor locations (black) for the milk carton, hough space. Bottom row: Two query images with detected milk cartons, one with detected juice carton.

- Mindru, F., Moons, T., , and Gool, L. V. (1999). Recognizing color patterns irrespective of viewpoint and illumination. In *Computer Vision and Pattern Recognition, vol. 1*, pp. 368-373.
- Rosenfeld, A. and Kak, A. (1976). Digital picture processing. In *Computer Science and Applied Mathematics*, Academic Press, New York.
- Schmid, C., Mohr, R., and Bauckhage, C. (1997). Local grey-value invariants for image retrieval. In *International Journal on Pattern Analysis and Machine Intelligence, Vol. 19, no. 5*, pp. 872-877.
- Sivic, J. and Zisserman, A. (2003). Video google: A text retrieval approach to object matching in videos. In *Proc. of 9th IEEE Intl Conf. on Computer Vision, Vol. 2*.
- Tuytelaars, T. and Gool, L. V. (2000). Wide baseline stereo based on local, affinely invariant regions. In *British Machine Vision Conference, Bristol, UK*, pp. 412-422.
- Tuytelaars, T., Gool, L. V., D'haene, L., , and Koch, R. (1999). Matching of affinely invariant regions for visual servoing. In *Intl. Conf. on Robotics and Automation*, pp. 1601-1606.
- Viola, P. and Jones, M. (2001). Rapid object detection using a boosted cascade of simple features. In *Computer Vision and Pattern Recognition*.
- Zhang, M., Marszalek, S. L. and Schmid, C. (2005). Local features and kernels for classification of texture and object categories: An in-depth study. In *Technical report, INRIA*.