

Semantic Image Search for Robotic Applications

Tomas Kulvicius^a, Irene Markelic^a, Minija Tamosiunaite^a and Florentin Wörgötter^a

^a*Georg-August-Universität Göttingen
Bernstein Center for Computational Neuroscience
Department for Computational Neuroscience
III Physikalisches Institut - Biophysik
Friedrich-Hund Platz 1, DE-37077 Göttingen, Germany
E-mail: {tomas,irene,minija,worgott}@physik3.gwdg.de*

Abstract. Generalization in robotics is one of the most important problems. New generalization approaches use internet databases in order to solve new tasks. Modern search engines can return a large amount of information according to a query within milliseconds. However, not all of the returned information is task relevant, partly due to the problem of polysemes. Here we specifically address the problem of object generalization by using image search. We suggest a bi-modal solution, combining visual and textual information, based on the observation that humans use additional linguistic cues to demarcate intended word meaning. We evaluate the quality of our approach by comparing it to human labelled data and find that, on average, our approach leads to improved results in comparison to Google searches, and that it can treat the problem of polysemes.

Keywords. Internet-based Knowledge, Polysemy, Semantic Search, Image Database Cleaning

1. Introduction

Humans can generalize to new tasks very quickly whereas for robots this is still not an easy task which makes it one of the most important and relevant problems in robotics. One of the most common approaches in generalization is learning from previous experiences (Ude et al., 2010; Nemec et al., 2011; Kober et al., 2012; Kronander et al., 2011). Some new approaches use internet databases in order to generalize to new situations (Tenorth et al., 2011; Beetz et al., 2011; Tamosiunaite et al., 2011). In particular, here we are interested in generalization in object domain by using image search. Although modern search engines like Google or Yahoo do an amazing job in returning a large number of images according to a query within milliseconds, not all of the returned images are task/context-relevant. A reason for spurious results is that most image searches rely on text-based queries, which is justified, since visual and textual information are dual to some degree. An *image* of a cup can be interpreted as the visual representation of the concept cup, whereas the *word* cup can be seen as a linguistic handle to the concept cup as represented in the human mind (Grush, 2004). Therefore, existing tools for text-based information retrieval applied to image search can lead to relatively good results (Brin and Page,

1998). Problems arise mainly due to ambiguities: 1) The same linguistic handle can map to several, different concepts, e.g., homonyms and polysemes. An example is the word “jaguar” which can refer to a car or an animal. Without any further information, e.g., contextual information, it is not possible to infer which domain is actually referred to. 2) Text-based image search relies on the assumption that textual information that is somehow related to an image, e.g., text placed close-by an image on a web page refers to the image content (Brin and Page, 1998). This assumption is reasonable, however not always correct, e.g., not every web-page creator names images according to their content.

A lot of effort has been spent on trying to resolve the problem of obtaining unclear image search results, often with the goal of object detection or image categorization, by making additional use of image content in form of visual cues, e.g., features like local image patches, edges, texture, color, deformable shapes, etc. (Fergus et al., 2003, 2004, 2005; Guillaumin et al., 2010; Berg and Forsyth, 2006; Schroff et al., 2007; Khan et al., 2011; jia Li et al., 2007; Wang et al., 2009; Jing and Baluja, 2008; A.D. Holub, 2008). All these approaches use textual information, too. Either implicitly by using the results of text-based image search engines e.g., Fergus et al. (2005, 2004), or

constructing their own image search (Schroff et al., 2007; Berg and Forsyth, 2006; A.D. Holub, 2008), or explicitly, by making use of image tags and labels as found in photo-sharing websites like Flickr (Guillaumin et al., 2010; Wang et al., 2009; Berg and Forsyth, 2006). An interesting work is Wang et al. (2009), because it is somewhat inverse to the standard procedure: Instead of using images with similar text labels to obtain image features for classification, they reverse the problem and use similar images to obtain textual features.

To our knowledge all of the aforementioned approaches achieve an improved precision of the result set, however, none can automatically cope with the problem of polysemes. For example in Fergus et al. (2004) a re-ranking of images obtained from Google searches was proposed, based on the observation that images related to the search are visually similar while unrelated images differed. This “visual consistency”, what we will here call inter-image similarity, was measured using a probabilistic, generative image model, and the EM-algorithm was used for estimating the model parameters from image features. Naturally, due to the underlying assumption, this will not work well for homonyms, since for these many images that are actually closely related to the search can have a very different appearances. A similar problem was faced in Fergus et al. (2005), where an extended version of pLSA (probabilistic Latent Semantic Analysis) was used to learn a clustering of images obtained from a Google search. A solution suggested in Berg and Forsyth (2006) copes with the polysemes problem but requires human supervision for this stage. Google text search is used to collect webpages for 10 animals. Then LDA (Latent Dirichlet Allocation) is applied to text from these pages to discover a set of latent topics. Images extracted from the webpages are then assigned to the identified topics, according to their nearby word likelihood. The problem of polysemes is tackled by a human user who manually selects or rejects these image sets.

Here, we present our approach which we call SIMSEA (Semantic Image SEArch) which also aims at increasing the precision of Internet image search results. Its most prominent advantage is that it can cope near-to automatically with polysemes. This is achieved by exploiting the fact that also humans need to resolve ambiguities in every-day speech, e.g., we may say “the bank - that you can *sit* on” to distinguish it from the bank that deals with money. Thus, we give additional cues to demarcate our intended meaning of a word. Here, we combine this linguistic refinement with the image-level in the following way: We conduct several different image searches, where we pair the basic search term with an additional linguistic cue. E.g., if interested in the category “cup”, (the basic search term), we search for “coffee cup”, “tea cup”, etc. The expectation is that images that are retrieved by more

than one of these subsearches are more likely to be of interest, than those that are retrieved only once. Note that for simplicity, in this paper we defined additional cues manually. In general, automated extraction of object descriptors (cues) can be done using methods of natural language processing (Cimiano, 2006; Olivie et al., 2011; McAuley et al., 2012), however, this is out of the scope of the current paper.

To compute the similarity between images from different subsearches, we use a “Bag-of-Words” (or “Bag-of-Features”) representation as often used in image classification. More precisely, we compute a codebook based on PHOW (Bosch et al., 2007a,b) features. However, also other, or additional, features are possible.

In addition to this procedure for achieving cleaner search results, we propose a ranking of the retrieved images, which is simply based on the idea that an image is the more relevant the more subsearches it, or a very similar match, is contained in.

We evaluate the quality of the obtained image set, and our proposed ranking, by comparing it to human labelled data.

The paper is structured as follows: We give a detailed description of our procedure in section 2.1, followed by the explanation of how we evaluated our method and the presentation of the achieved results in section 3. Finally we discuss and conclude our work in section 4.

2. Methods

2.1. SIMSEA Overview

The approach is summarized in Fig. 1 and an overview on its stages, which are enumerated in the figure, are described below, followed by a more detailed description of each stage in the paragraphs 2.2 and 2.3.

The goal is to find “clean” results for image searches with respect to given task/context. For that we conduct several image searches to which we refer as *subsearches* (4), see Fig. 1. A subsearch is conducted using the *basic search term* (1) with an additional *linguistic cue* (2+3). E.g., if interested in the category “cup”, we search for “coffee cup”, “tea cup”, etc. (using Google). The set of images retrieved by a subsearch is consequently referred to as *subsearch results* (5). The expectation is that images that are retrieved by more than one subsearch are more likely to be task/context-relevant than those that do not, they form the final *result set* (6+7). We do not consider only images that have exact copies in other subsearch result sets, but instead relax this demand and also consider images as relevant if merely a similar image is returned by another subsearch.

Finally, we suggest to rank the retrieved result set (8). The ranking is supposed to indicate how relevant a given image is, e.g., a glass-image with a high

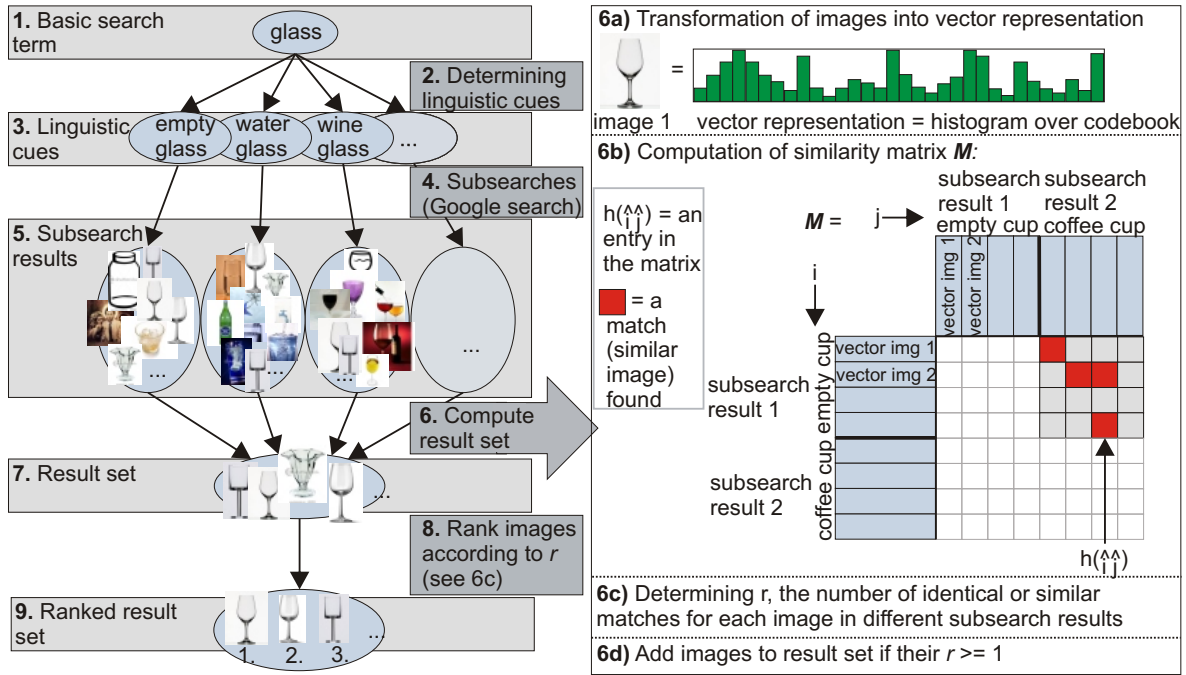


Fig. 1. The different stages of SIMSEA sketched for the category “glass”. Note that M , for clarity, is only shown for the first two subsearch results.

ranking factor should be considered to be very likely a true representative of the category glass, whereas an image with a low ranking factor can be considered to be very likely not a good representative of its class. As stated, we assume that images that have similar counterparts in other subsearch results are more likely to be relevant. This measurement can be used as a simple relevance ranking of the resulting images: The more often an image (or a similar counterpart) occurred in other subsearches the higher its relevance.

2.2. Linguistic Cues and Subsearches

We investigate four different categories (basic search terms) taken from a kitchen scenario: “cup”, “glass”, “milk” and “apple”. Glass and apple are homonyms (vision, drinking, and material; or brand and fruit). Milk is another special case, because as a liquid it usually comes in a more or less characteristic container. For each of the four categories we conduct a varying number of subsearches in which we combine the basic search term with an additional linguistic cue. For the category milk we conduct six subsearches, namely: “cold milk”, “fresh milk”, “healthy milk”, “tasty milk”, and “hot milk”. We also conduct a query with the basic search term “milk” without any additional cue. The linguistic cues we use for apples are: “delicious”, “green”, “red”, “ripe”, “sour”, “sweet”, and “unripe”. In addition we search for the basic terms “apple” and “apples”. For glass we use: “empty”, “full”, “juice”, “milk”, “water”, “wine”, and the word “glass” by itself. For cup: “coffee”, “full”, “tea”, and simply “cup”. The strategy for selecting the cues was

to select those that restrict the domain to the desired kitchen domain.

2.3. Computing the Result Set

As explained, the expectation is that images that have similar counterparts (or matches) in the other subsearch results are likely to meet the user expectations. To be able to measure inter-image similarity we use a “Bag-of-Words” approach. In such an approach each image is represented by a histogram over a fixed number of so-called “visual words” which are also often referred to as “codebook”. These visual words are usually created based on local image features. In our case we use PHOW features (which are explained below) but other features can be used, too.

First, the codebook needs to be generated. For that we take a small, randomly chosen subset of images, we use 40, from each category. We compute PHOW features for all these 160 (40×4 categories) images which we then quantize into k vectors - the visual words - using k -means clustering. We set k to 200.

After having created the codebook, we can represent each image by a vector, a histogram over the codebook, which is computed as follows (step 6a in Fig. 1): We determine the PHOW features of each image i , map these to the k visual words that make up the codebook, and compute the histogram which counts which visual word occurred how many times in the given image. This histogram, called “vector image i ” in Fig. 1, is then used to represent image i .

To compute the similarity matrix M (step 6b), we consider all images i, j , which are contained in differ-

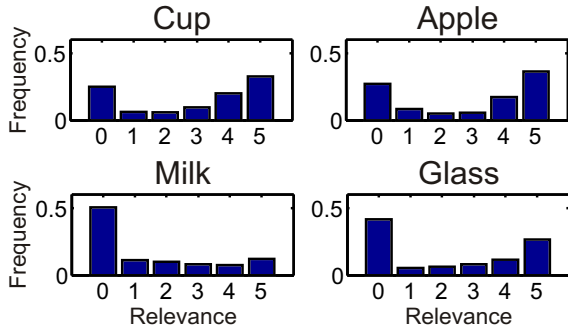


Fig. 2. Histogram of image category membership assigned by the five human subjects from which we derive the image relevance.

ent subsearch result sets. For each such image pair we compute the Hellinger distance $h(\hat{i}\hat{j})$, described in detail below. Here, \hat{i} and \hat{j} refer to the vector representation of the two images. This leads to an upper diagonal matrix M , as indicated in Fig. 1.

In step 6c we then determine for each image i the number of similar matches or counterparts for that image which we refer to as the ranking factor r_i . We consider two images to be matches if their Hellinger distance is above a certain threshold. Determining r_i is simply counting the number of matches for that image. For example, the first row of matrix M in Fig. 1 shows all matches for the first image in the subsearch for “empty cup” in orange. Since there is only one match, r for this image is one. For the second image of the same subsearch there are two matches, thus r for this image is two, etc.

Finally, we select those images to be part of the final result set whose ranking factor r is larger than one (step 6d).

The ranking, step 8, is trivial, it merely consists of ordering the images from the computed result set according to their ranking factor. The ranking is supposed to indicate how relevant a given image is, e.g., a glass-image with a high ranking factor should be considered to be very likely a true representative of the category glass, whereas an image with a low ranking factor can be considered to be very likely not a good representative of its class.

Pyramid Histogram of Visual Words (Bosch et al., 2007a,b) are state-of-the-art image descriptors based on a variant of dense SIFT (Lowe, 2004). A grid with a self-defined spacing (here we use 5 pixels) is laid over an image and at each grid point four SIFT descriptors, varying in radii to allow for scale variations, are computed. This can be done on various levels, hence “Pyramid”, but here we suffice with the first level, thus, to be precise we are actually using HOW descriptors. We use the VLFeat library (Vedaldi and Fulkerson, 2010) to compute the PHOW descriptors and the subsequent vector representation of the images.

To compute the similarity between image pairs we use the Hellinger distance¹. The Hellinger distance between two distributions P and Q is denoted $H(P, Q)$ and satisfies $0 \leq H(P, Q) \leq 1$ (where 1 denotes large distances and 0 no distances, i.e., identical images). It is defined as follows.

$$H(P, Q) = \sqrt{1 - BC(P, Q)}, \quad (1)$$

where BC denotes the Bhattacharyya coefficient which, in the discrete case, is defined as:

$$BC(P, Q) = \sum_{x \in X} \sqrt{P(x)Q(x)}. \quad (2)$$

Here X denotes the common domain over which the two distributions are defined. We define two images to be similar if their Hellinger distance is above a fixed threshold (we use 0.15, experimentally chosen). For n subsearches, where s_i denotes the i ’th subsearch, we compare each image from each subsearch to all images from all other subsearches except to its own. Thus, if the total number of images is $N = |s_1| + |s_2| + \dots + |s_n|$ (where the vertical bars denote the number of elements in the set), we have $C = \binom{N}{2} - \sum_{i=1}^n \binom{|s_i|}{2}$, where C is the total number of comparisons that need to be computed. This is depicted in the visualization of M in Fig. 1. Note, we do not compare images from the same subsearch to each other. This is because we are not interested in intra-subsearch similarity due to the following reason: We may receive many images of the same topic during one search but which are unrelated to what we are interested in. If we counted the intra-subsearch similarity these images would be evaluated as highly relevant to our search interest which they are not.

3. Results

Since the goal is to find a subset of images which meets the semantic expectation of the user, we need some “ground truth”, i.e., a set of true samples, to evaluate our algorithm. For this issue we let several human subjects classify the same data that was input to the algorithm according to the given categories. This way we can gather various subjective human opinions and determine those images that get assigned the same labels by all subjects and also those where opinions differed. In the following we describe the ground truth retrieval procedure.

3.1. Ground-Truth Retrieval

We asked five human subjects to aid in retrieving the ground-truth data to which we compare our algorithm. Each human was instructed to decide for each image from the subsearches for milk (hot, tasty, cold, ...) if it belonged, in his or her opinion, to the category milk.

¹We also used the χ^2 -distance, which gave very similar results. The Hellinger distance has the advantage to be bounded.

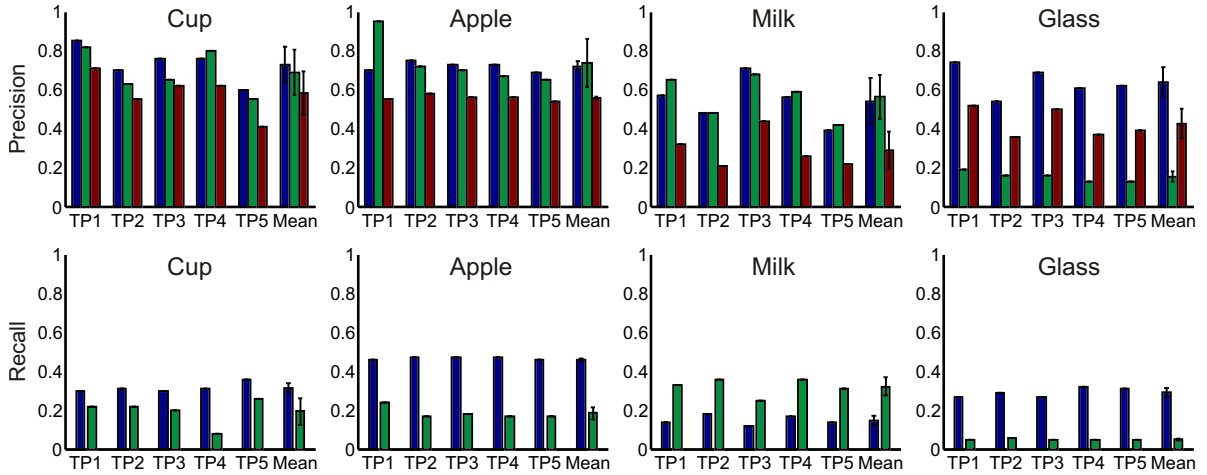


Fig. 3. Precision and recall of SIMSEA (blue), a standard Google search (Google, in green) and the cumulative data from all subsearches for a given category (SumGoogle, in red) with respect to the data obtained from each test person (TP1-5) for the categories. The vertical errorbar for the mean indicates the variance.

The same was done for the three other categories, cup, glass and apples. To make this evaluation as fair as possible, all humans were given precisely the same information by means of an instruction. Basically, the subjects were told that there are four categories and that they are from a kitchen scenario, thus, glass was supposed to be for drinking, and not for aiding vision, etc.

As explained, we suggest a ranking procedure which is supposed to reflect the relevance of a given image. So far we loosely defined relevance as “goodness of an image as class representative”. To evaluate our ranking result, we again require some “ground truth” data that we can compare it to, which we obtain as follows: In Fig. 2 we show a histogram indicating for each category how many of the test persons considered a given image as being member of a category. Since there were five test persons each image can be selected as category member between zero and five times. We assume that images which were considered by none of the test persons as category member should be assigned the lowest relevance, and vice versa, images considered by all test persons should be assigned the highest relevance. Thus, for each image we can compute a measure based on the five human subjects decisions to which we refer to as *relevance* whereas the equivalent measure of SIMSEA is called *ranking*. We will use these measurements in Section 3.2.

3.2. Evaluation

We assess the quality of the algorithm by computing precision and recall on its output, see Eq. 3, with respect to the ground truth data from each human subject.

$$\begin{aligned} \text{precision} &:= (A \cap B) / |A| \\ \text{recall} &:= (A \cap B) / |B|, \end{aligned} \quad (3)$$

where A is the set of retrieved samples and B is the set of true samples, i.e., in our case A is the set of samples retrieved by SIMSEA and B is the set of samples belonging to a given category selected by each human subject. Since there were five human subjects, there are five true sample sets, with respect to which we compute precision and recall. The results are given in Fig. 3. We compare these results to (i) standard Google searches and also to (ii) the union of all subsearches of a given category. For (i), we conduct standard Google searches with the basic search terms for each category, e.g., for the category milk, the set A is the set of images returned by a Google search using the search term “milk”, and again compute precision and recall the way described before. In Fig. 3 we refer to this evaluation as “Google”. For (ii) we set A to the union of the images from all subsearches of a given category, i.e., if we conducted n subsearches for the category milk, we have s_1, s_2, \dots, s_n subsearch result sets and we set $A = s_1 \cup s_2 \dots \cup s_n$. In Fig. 3 we refer to this evaluation as “SumGoogle”. Note, that for SumGoogle the recall is always one. This is because the ground truth set from all human subjects is a subset of the union of subsearches for a category, in other words $B \subset A$.

To be useful, precision and recall of SIMSEA should be higher than those of the standard Google search and SumGoogle. In other words, most human subjects should find that the output of SIMSEA gives more relevant results than the Google standard search and SumGoogle (precision), and also that SIMSEA returns more of the overall available relevant samples (recall). It can be seen from Fig. 3 that except for the category “milk” SIMSEA indeed outperforms the standard Google search and SumGoogle.

For the category “cup”, test person 4 (TP4) agrees more with the results of the Google search, but all



Fig. 4. The first 19 glasses from a Google search for the word “glass” (upper panel) and the best ten glasses according to the ranking obtain by using SIMSEA algorithm (bottom panel).

other human subjects consider more images retrieved by the automatic routine to be important. It can also be seen that the values for precision and recall differ between the subjects which shows, what we had already expected, that assigning images to a certain category also depends on subjective opinions. For the category “apple”, TP1 shows a very clear preference for the Google search results. Due to TP1 also the precision is higher for the Google search than the automatic routine. However, TP1’s opinion is not in accordance with the that of the other subjects, which all have a precision value around 0.7 and therefore we consider this to be an outlier. Without TP1’s influence SIMSEA outperforms the Google search for “apple”, too.

For the category “milk” we can observe a different case, most human subjects are more in accordance with the results of the Google standard search. A possible reason for that can be found in Fig. 2. We see that for all categories there are clear peaks for images that all human subjects consider as category member and for those that all human subjects consider to not be category members. Except for the category “milk”. Here, the relation between images with full voting or relevance (all five human subjects) and ambiguous decisions, i.e., where for example only two out of three subjects considered an image as relevant, is higher than compared to the other categories. In other words, there are many images in the milk category for which even humans find a clear decision difficult. This might be due to the fact that we have already stated that milk as a liquid is depicted to be contained in more or less characteristic containers. We can assume that for this reason SIMSEA is not performing well for this category either.

To select “clean” images from a Google search we

use a ranking r as described above, i.e., frequency with which they occur in the different subsearches. To visualize the effect of the ranking we show the best ten glasses in Fig. 4, bottom panel. In the upper panel we show the first 19 images returned by a Google search for the word “glass”. We can see that Google search results include images of glasses from domains others than the desired kitchen domain (in this case $\approx 42\%$). SIMSEA in contrast was successful in eliminating those.

4. Discussion

We proposed a method based on the combination of linguistic cues with the image domain that is useful for retrieving cleaner results in image searches, in particular it is able to tackle the problem of polysemes. This is a novel approach and we have given the proof of principle by showing that it indeed leads to cleaner search results.

In addition we suggested a ranking, based on the occurrence frequency of images between different subsearches. We could show that this roughly reflects a human based relevance measure.

Although we have introduced the notion of linguistic cues, we have not tackled the issue where these cues might come from, or how they should best be chosen. Automated extraction of object descriptors (cues) can be done using methods of natural language processing (Cimiano, 2006; Olivie et al., 2011; McAuley et al., 2012). However, this is an issue falling in the domain of linguistics and is not the core of this paper.

It is obvious that our method can only be as good as the subsearch results which depend on the “right” linguistic cues. If unrelated images occur in many of

the subsearches, these images will erroneously be part of the result set.

Similar to the effectiveness of human linguistic refinement to distinguish intended meaning from other, our method has its strength when dealing with polysemes or homonyms. For example, the result for the category “glass” is very good, where it had to distinguish drinking glasses from the material glass and the vision aid. In contrast the method did not perform well for the liquid “milk”.

Another critical issue is the similarity function between images. Here we used PHOW features and the Hellinger distance, which works satisfactory, but also sometimes lead to artefacts, i.e., images that do not appear similar to humans can sometimes turn out to be very similar when using PHOW features. Here, different features and metrics may lead to an improvement of the method. Another option can be to follow the idea of (A.D. Holub, 2008) and to learn an appropriate metric by solving a constrained optimization problem.

In summary, we believe that this a novel and promising idea for data “cleaning” which can be used to automatically form training data sets using Internet search which later can be used for object classification/recognition and generalization. In future work we are going to include more classes and make such image search completely automatic by augmenting it with an automated extraction of object descriptors from language.

5. Acknowledgements

The research leading to these results has received funding from the European Community’s Seventh Framework Programme FP7/2007-2013 (Programme and Theme: ICT-2011.2.1, Cognitive Systems and Robotics) under grant agreement no. 600578, ACAT.

6. References

P. Perona A.D. Holub, P. Moreels. Unsupervised clustering for google searches of celebrity images. *8th IEEE Int. Conf. Automatic Face and Gesture Recognition*, 2008.

Michael Beetz, Ulrich Klank, Ingo Kresse, Alexis Maldonado, Lorenz Mösenlechner, Dejan Pangercic, Thomas Rühr, and Moritz Tenorth. Robotic Roommates Making Pancakes. In *11th IEEE-RAS Int. Conf. on Humanoid Robots*, pages 529–536, Bled, Slovenia, October, 26–28 2011.

Tamara L. Berg and David A. Forsyth. Animals on the web. In *CVPR*, pages 1463–1470, 2006.

A. Bosch, A. Zisserman, and X. Munoz. Representing shape with a spatial pyramid kernel. In *ACM Int. Conf. Image and Video Retrieval*, 2007a.

Anna Bosch, Andrew Zisserman, and Xavier Muoz. Image classification using random forests and ferns. In

ICCV, pages 1–8, 2007b.

Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual web search engine. *Comput. Netw. ISDN Syst.*, 30:107–117, April 1998.

P. Cimiano. *Ontology Learning and Population from Text: Algorithms, Evaluation and Applications*. Springer Verlag, 2006.

R. Fergus, P. Perona, and A. Zisserman. Object class recognition by unsupervised scale-invariant learning. In *CVPR*, pages 264–271, 2003.

R. Fergus, P. Perona, and A. Zisserman. A visual category filter for google images. In *8th Europ. Conf. Computer Vision*, pages 242–256, May 2004.

R. Fergus, L. Fei-Fei, P. Perona, and A. Zisserman. Learning object categories from google’s image search. In *10th IEEE Int. Conf. Computer Vision*, volume 2, pages 1816–1823, oct. 2005.

Rick Grush. The emulation theory of representation: Motor control, imagery, and perception. *Behavioral and Brain Sciences*, 27:377442, 2004.

Matthieu Guillaumin, Jakob Verbeek, and Cordelia Schmid. Multimodal semi-supervised learning for image classification. In *CVPR*, 2010.

Li jia Li, Gang Wang, and Li Fei-fei. Optimol: automatic online picture collection via incremental model learning. In *CVPR*, 2007.

Yushi Jing and S. Baluja. Visualrank: Applying pagerank to large-scale image search. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 30(11):1877–1890, Nov. 2008.

Inayatullah Khan, Peter M. Roth, and Horst Bischof. Learning object detectors from weakly-labeled internet images. In *35th OAGM/AAPR Workshop*, 2011.

Jens Kober, Andreas Wilhelm, Erhan Oztop, and Jan Peters. Reinforcement learning to adjust parametrized motor primitives to new situations. *Auton. Robots*, 33(4):361–379, 2012.

K. Kronander, M.S.M. Khansari-Zadeh, and A. Billard. Learning to control planar hitting motions in a minigolf-like task. In *2011 IEEE/RSJ Int. Conf. Intelligent Robots and Systems*, pages 710–717, sept. 2011.

David G. Lowe. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vision*, 60:91–110, Nov. 2004.

J. J. McAuley, J. Leskovec, and D. Jurafsky. Learning attitudes and attributes from multi-aspect reviews. In *International Conference on Data Mining*, 2012.

B. Nemeč, R. Vuga, and A. Ude. Exploiting previous experience to constrain robot sensorimotor learning. In *11th IEEE-RAS Int. Conf. Humanoid Robots*, pages 727–732, oct. 2011.

J. Olivie, C. Christianson, and J. McCarry. *Handbook of natural Language Processing and Machine Translation*. Springer, 2011.

F. Schroff, A. Criminisi, and A. Zisserman. Harvesting image databases from the web. In *11th IEEE Int. Conf. on Computer Vision*, pages 1–8, Oct. 2007.

M. Tamosiunaite, I. Markelic, T. Kulvicius, and F. Worgotter. Generalizing objects by analyzing language. In *11th IEEE-RAS Int. Conf. Humanoid Robots*, pages 557–563, oct. 2011.

Moritz Tenorth, Ulrich Klank, Dejan Pangercic, and Michael

- Beetz. Web-enabled Robots – Robots that Use the Web as an Information Resource. *Rob. & Automat. Magazine*, 18(2):58–68, 2011.
- A. Ude, A. Gams, T. Asfour, and J. Morimoto. Task-specific generalization of discrete and periodic dynamic movement primitives. *IEEE Trans. Rob.*, 26(5):800–815, oct. 2010.
- A. Vedaldi and B. Fulkerson. Vlfeat – an open and portable library of computer vision algorithms. In *18th annual ACM Int. Conf. Multimedia*, 2010.
- Gang Wang, D. Hoiem, and D. Forsyth. Building text features for object image classification. In *IEEE Int. Conf. Computer Vision and Pattern Recognition*, pages 1367–1374, Jun. 2009.