# Measuring and Predicting Tag Importance for Image Retrieval

Shangwen Li, Sanjay Purushotham, Chen Chen, Yuzhuo Ren, *Student Member, IEEE,*
and  C.-C. Jay Kuo, *Fellow, IEEE*

**Abstract**—Textual data such as tags, sentence descriptions are combined with visual cues to reduce the semantic gap for image retrieval applications in today's Multimodal Image Retrieval (MIR) systems. However, all tags are treated as equally important in these systems, which may result in misalignment between visual and textual modalities during MIR training. This will further lead to degenerated retrieval performance at query time. To address this issue, we investigate the problem of tag importance prediction, where the goal is to automatically predict the tag importance and use it in image retrieval. To achieve this, we first propose a method to measure the relative importance of object and scene tags from image sentence descriptions. Using this as the ground truth, we present a tag importance prediction model to jointly exploit visual, semantic and context cues. The Structural Support Vector Machine (SSVM) formulation is adopted to ensure efficient training of the prediction model. Then, the Canonical Correlation Analysis (CCA) is employed to learn the relation between the image visual feature and tag importance to obtain robust retrieval performance. Experimental results on three real-world datasets show a significant performance improvement of the proposed MIR with Tag Importance Prediction (MIR/TIP) system over other MIR systems.

**Index Terms**—Multimodal image retrieval (MIR), image retrieval, semantic gap, tag importance, importance measure, importance prediction, cross-domain learning.

✦

## 1 INTRODUCTION

IMAGE retrieval [1] is a long-standing problem in the computer vision and information retrieval fields. Current image search engines in service rely heavily on the text data. It attempts to match user-input keywords (tags) to accompanying texts of an image. It is thus called Tag Based Image Retrieval (TBIR) [2], [3]. TBIR is effective in achieving semantic similarity [4] due to the rich semantic information possessed by the text data. However, with the exponential growth of web images, one cannot assume all images on the Internet have the associated textual data. Clearly, TBIR is not able to retrieve untagged images even if they are semantically relevant to the input text query.

On the other hand, Content Based Image Retrieval (CBIR) [4], [5], [6], [7] takes an image as the query and searches for relevant images based on the visual similarities between the query image and the database images. Despite a tremendous amount of effort in the last two decades, the performance of CBIR system is bounded by the semantic gap between low-level visual features and high-level semantic concepts.

To bridge the semantic gap, one idea is to leverage well annotated Internet images. Due to the rich information over the Internet, numerous images are well annotated with text information such as tags, labels, sentences or even paragraphs. Their visual content provides low-level visual features while their associated textual information provides meaningful semantic information. The complementary nature of texts and images provides more complete

descriptions of underlying content. It is intuitive to combine both image and text modalities to boost the image retrieval performance, leading to Multimodal Image Retrieval (MIR).
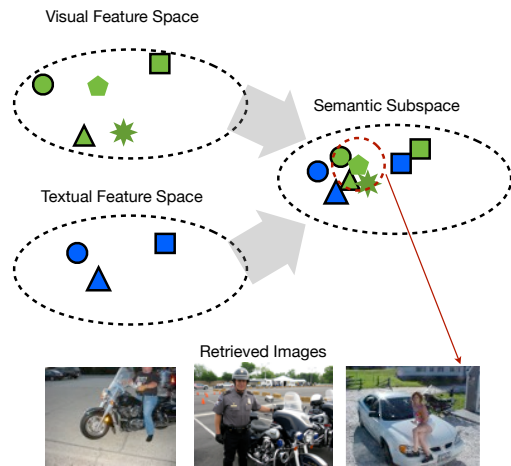
When a large number of database images with their textual information are available, one could use MIR to find a common subspace for the visual and textual features of these tagged images using machine learning algorithms such as the Canonical Correlation Analysis (CCA) [8], [9], [10], [11], [12], [13] or more recently the Deep Learning [14], [15], [16], [17], [18] techniques. Once the common subspace is found, the visual similarity between the query and tagged database images in the same subspace can indicate the semantic similarity, thus reducing the semantic gap to some extent. Moreover, untagged database images can be retrieved if their projected visual features in the subspace are sufficiently close to that of the query image. An example of MIR applied to image-to-image retrieval is shown in Fig. 1. Another advantage of the MIR framework is that it can accommodate different retrieval tasks at the same time. Since the subspace between the visual and the textual domains has been built, cross-modality information search, such as text-to-image and image-to-text search [19], [20], [21], can be achieved by leveraging the subspace [8], [9], [10].

The performance of an MIR system is highly dependent on the quality of tags. Unfortunately, as pointed out in [22], [23], tags provided by people over the Internet are often noisy, and they might not have strong relevance to image content. Even if a tag is relevant, the content it represents might not be perceived as important by humans. Take Fig. 2 as an example. It has two images with the same tags: "car", "motorbike", and "person". While the tag "motorbike" is perceived as having higher importance in the left image,

- *S. Li, S. Purushotham, C. Chen, Y. Ren, and C. Kuo are with the Ming Hsieh Department of Electrical Engineering, University of Southern California, Los Angeles, CA 90089. E-mail: {shangwel, yuzhuore}@usc.edu, {sanjayp2005, ohyline}@gmail.com, cckuo@sipi.usc.edu*

(a) A toy MIR database with tagged and untagged images.



(b) Image-to-image search.

Fig. 1. Image-to-image search under the MIR framework. (a) A toy database consisting of both tagged and untagged images. (b) Illustration of the image-to-image search. Tagged database images are used to learn the common semantic subspace between the visual and the textual domains in the training stage. Visual features of both the query image and untagged database images will be projected into the common semantic subspace to calculate the visual similarity during the query time.



Fig. 2. Two images with the same object tags ("car", "motorbike", and "person") but substantially different visual content.

it is not as important as "car" and "person" in the right one. Thus, if the left one is the query image, the right one will not be a good retrieved result. Therefore, capturing and incorporating human perceived tag importance can significantly improve the performance of automatic MIR systems, which will be clearly demonstrated in this work.

When people describe an image, they tend to focus on important content in the image. Thus, sentence descriptions

serve as a natural indicator of tag importance. Berg *et al.* [24] modeled tag importance as a binary value, i.e. a tag is important if and only if it appears in a sentence. However, when multiple tags are present in the sentences, they might not be equally important. Thus, the binary-valued tag importance cannot capture the relative tag importance, which degrades the retrieval performance of MIR systems. To address this deficiency, we study the relative tag importance prediction problem and incorporate the predicted importance in an MIR system, leading to the MIR with Tag Importance Prediction (**MIR/TIP**) system. To the best of our knowledge, this is the first work that embeds predicted tag importance into the semantic subspace learning in an MIR system.

It is worthwhile to mention that tags can be rather versatile in the wild. On one hand, it may refer to object ("person") or scene ("beach"), which is closely related to the image visual content. On the other hand, it may also refer to more abstractive content such as attribute("joy") or action ("surfing"). Last but not least, it may refer to photography techniques or camera parameters. The different properties of these content make it hard to study all types of tags at the same time. Thus, in this paper, we focus on studying the importance of object and scene tags. Specifically, we assume all tags are relevant to image visual content, and do not consider irrelevant (or noisy) tags since tag cleaning is a research topic on its own. Based on this assumption, we treat image object/scene category labels as the tags in our experiments.

To build an MIR/TIP system, we need to address the following questions:

1) How to define tag importance?
2) How to predict the defined tag importance?
3) How to embed tag importance in MIR?

For the first question, we propose a method to measure the object and scene tag importance from human provided sentence descriptions based on natural language processing (NLP) tools. A subjective test is conducted to validate its benefits in image retrieval. For the second question, we present a novel prediction model that integrates visual, semantic, and context cues. While the first two cues were explored before in [24], [25], context cue has not been considered for tag importance prediction. The context cue contributes to tag importance prediction significantly, and results in improved image retrieval performance as demonstrated in our experiments. To train the prediction model, we use the Structural Support Vector Machine (SSVM) [41] formulation. For the third question, Canonical Correlation Analysis (CCA) [11] is adopted to incorporate the predicted tag importance in the proposed MIR/TIP system.

The rest of this paper is organized as follows. Related previous work and an overview of the proposed MIR/TIP system are presented in Section 2 and Section 3, respectively. A technique to measure tag importance based on human sentence descriptions is discussed in Section 4. Tag importance prediction is studied in Section 5, where the measured tag importance is used as the ground truth. The MIR/TIP system is described in Section 6. Experimental results are shown in Section 7. Finally, concluding remarks are given and future research directions are pointed out in Section 8.

## 2 RELATED WORK

In this section, we review recent work on tag importance prediction and MIR systems that are closely related to our work.

**Importance Prediction.** Importance in images is a concept that has recently gained attention in visual research community. Elazary and Itti [26] used the naming order of objects as the interestingness indicator and saliency to predict their locations in an image. A formal study of object importance was conducted by Spain and Perona [25], who developed a forgetful urn model to measure object importance from ordered tag lists, and then, used visual cues to predict the object importance value. Berg *et al.* [24] used human sentence descriptions to measure the importance of objects, scenes and attributes in images, and proposed various visual and semantic cues for importance prediction. To better understand user defined content importance, Yun *et al.* [27] studied the relationship between the human gaze and descriptions. Parikh *et al.* [28], [29] proposed a ranking method for the same attribute across different images to capture its relative importance in multiple images. Instead of predicting tag importance for images directly, some works focus on image tags reranking to achieve better retrieval performance. For example, Liu *et al.* [23] developed a random walk-based approach to rerank tags according to their relevance to image content. Similarly, Tang *et al.* [30] proposed a two-stage graph-based relevance propagation approach. Zhuang and Hoi [31] proposed a two-view tag weighting approach to exploit correlation between tags and visual features. Lan and Mori [32] proposed a Max-Margin Riffled Independence Model to rerank object and attribute tags in an image in order of decreasing relevance or importance. More recently, to address a very large tag space, Feng *et al.* [33] cast tag ranking as a matrix recovery problem. However, all the previous approaches used ranked tag list, human annotators, or binary labels for defining tag importance and failed to capture the continuous-valued relative tag importance which is addressed in our paper.

**Multimodal Image Retrieval (MIR).** The current state-of-the-art MIR systems aim at finding a shared latent subspace between image visual and textual features so that the information in different domains can be represented in a unified subspace. Several learning methods have been developed for this purpose, including the canonical correlation analysis (CCA) [34] and its extension known as the kernel CCA (KCCA) [35]. The main idea of CCA is to find a common subspace for visual and textual features so that their projections into this lower dimensional representation are maximally correlated. Hardoon *et al.* [11] adopted KCCA to retrieve images based on their content using the text query. Rasiwasia *et al.* [12] replaced the textual modality with an abstract semantic feature space in KCCA training. More recently, Gong *et al.* [10] proposed a three-view CCA that jointly learns the subspace of visual, tag and semantic features. Hwang and Grauman [8], [9] adopted human provided ranked tag lists as object tag importance and used them in KCCA learning, which is most relevant to our work. Deep learning based multimodal methods adopt the same idea of learning the shared representation of multi-modal data. However, they are based on recently developed deep learning techniques such as stacked autoencoder [14], [16], [17] and deep convolutional neural network [17], [18]. The quality of the learned semantic subspace with shared representations of visual and textual modalities highly depends on the quality of tags. A noisy or unimportant tag may lead to misalignment between the visual and the textual domains, resulting in degenerated retrieval performance. Thus, there is a need to systematically build an MIR system by incorporating the learned tag importance model, which is one of the focuses of this paper.

## 3 SYSTEM OVERVIEW

A high-level description of the proposed MIR/TIP system is given in Fig. 3. It consists of the following three stages (or modules):

1) Tag importance measurement;
2) Tag importance prediction; and
3) Multimodal retrieval assisted by predicted tag importance.

It is assumed in our experiments that there are three types of images on the web:

A images with human provided sentences and tags;
B images with human provided tags only; and
C images without any textual information.

In the tag importance measurement stage, images in Type A are used to obtain the measured tag importance, which will serve as the ground truth tag importance. In the tag importance prediction stage, images in Type A are first used as the training images to create a tag importance prediction model. Then, tag importance for images in Type B will be predicted based on the learned model. Finally, images in Types A and B will be used as training images to learn the CCA semantic subspace in the multimodal retrieval stage. Images in Type C will serve as test images to validate the performance of our MIR/TIP system. Table 1 summarizes the textual and visual features used in different stages of MIR/TIP. Details of each module will be described in the following sections.

## 4 MEASURING TAG IMPORTANCE

Measuring human-perceived importance of a tag is a critical yet challenging task in MIR. Researchers attempted to measure the importance of tags associated with images from two human provided sources: 1) ranked tag lists [8], [9], [25], and 2) sentence descriptions [24]. The major drawback of ranked tag lists is their unavailability. Tags are rarely ranked according to their importance, but rather listed randomly. Obtaining multiple ranked tag lists from human (using the Amazon Turk) is labor intensive and, thus, not a feasible solution. In contrast, human sentence descriptions are easier to obtain due to the rich textual information on the Internet. For this reason, we adopt sentence descriptions as the source to measure tag importance.

Clearly, the binary-valued tag importance as proposed in [24] cannot capture relative tag importance in an image. For example, both "person" and "motorbike" in Fig. 4 are important since they appear in multiple sentences. However, as compared with "person" that appears in all five

TABLE 1
Visual and textual features used for tag importance measurement, prediction, and MIR.

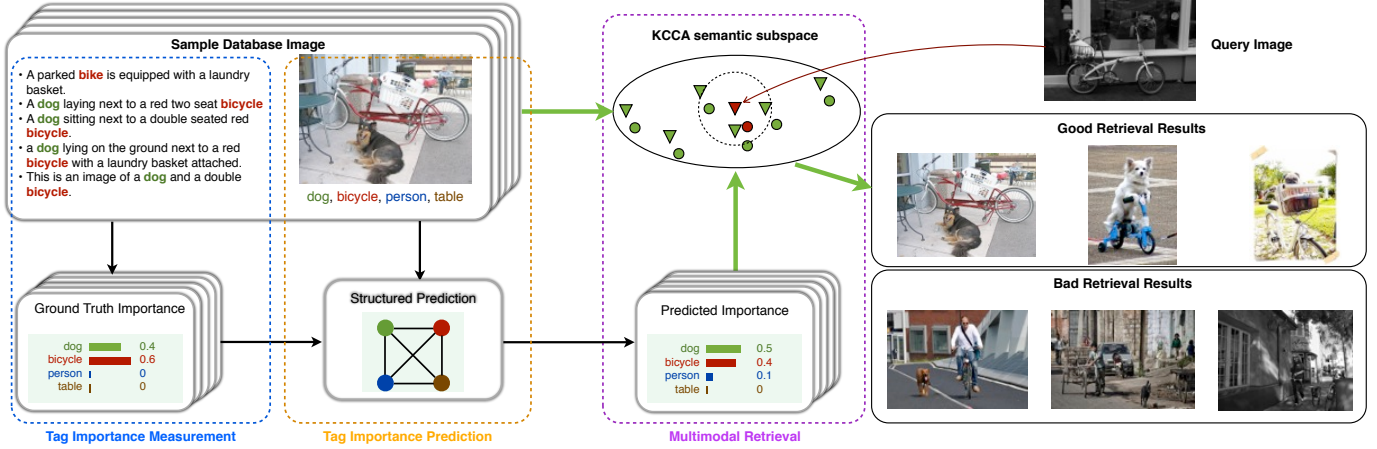| | | Textual | Visual |
|---|---|---|---|
| Tag Importance Measurement (Sec. 4) | | Tags, Sentences | NA |
| Tag Importance Prediction (Sec. 5) | Training Data Image Type A | Tags, Measured Importance | R-CNN Detected Object Bounding Box, Places VGG16 FC7, Saliency |
| | Testing Data Image Type B | Tags | R-CNN Detected Object Bounding Box, Places VGG16 FC7, Saliency |
| MIR System (Sec. 6) | Training Data Image Type A & B | Tags, (A) Measured Importance, (B) Predicted Importance | ImageNet VGG16 FC7, Places VGG16 FC7 |
| | Testing Data Image Type C | NA | ImageNet VGG16 FC7, Places VGG16 FC7 |



Fig. 3. An overview of the proposed MIR/TIP system. Given a query image with important "dog" and "bicycle", the MIR/TIP system will rank the good retrieval examples with important "dog" and "bicycle" ahead of bad retrieval ones with less important "dog" and "bicycle".



Fig. 4. An example of object importance measurement using sentence descriptions, where object tags "person" and "motorbike" appear in sentences in synonyms such as "man" and "scooter", respectively.



Fig. 5. An example for comparison between probability importance and discounted probability importance, where the "bicycle" in both images are equally important with probability importance but not with the discounted probability importance (left "bicycle": 0.6; right "bicycle": 1.0).

sentences, "motorbike" only appears twice. This shows that humans perceive "person" as more important than "motorbike" in Fig. 4. Thus, tag importance should be quantified in a finer scale rather than a binary value.

The desired tag importance should serve the following two purposes.

1) **Within-image comparison.** Tag importance should teach the retrieval system to ignore unimportant content within an image.
2) **Cross-image comparison.** Given two images with the same tag, tag importance should identify the image in which the tag has a more important role.

In the following paragraphs, we will first introduce the idea of measuring object tag importance and, then, extend it to account for scene tag importance.

**Object Tag Importance.** To achieve within-image comparison, one heuristic way is to define the importance of an

object tag in an image as the probability of it being mentioned in a sentence. This is called **probability importance**. To give an example, for the left image in Fig. 5 (i.e. the sample database image in Fig. 3), the importance of "dog" and "bicycle" are 0.8 and 1 since they appear in four and five sentences, respectively. While this notion can handle within-image comparison, it fails to model cross-image comparison. For instance, as compared with the right image in Fig. 5, where the "bicycle" is the only tag appearing in all five sentences, it is clear that the "bicycle" in the left image is less important. However, its probability importance has the same value, 1, in both images.

To better handle cross-image comparison, we propose a measure called **discounted probability importance**. It is based on the observation that different people describe an image in different levels of detail. Obviously, tags mentioned by detail-oriented people should be discounted accordingly. Mathematically, discounted probability impor-

tance of object tag $t$ in the $n$th image $\mathbf{I}_n(t)$ is defined as

$$\mathbf{I}_n(t) = \sum_{k=1}^{K_n} \frac{\mathbb{I}\left\{t \in \mathcal{T}_n^{(k)}\right\}}{\left|\mathcal{T}_n^{(k)}\right| K_n}, \qquad (1)$$

where $K_n$ is the total number of sentences for the $n$th image, $\mathbb{I}$ is the indicator function, and $\mathcal{T}_n^{(k)}$ is the set of all object tags in the $k$th sentence of the $n$th image. An example of measured tag importance using Eq. (1) is shown in Fig. 4. Also, for the bicycles in Fig. 5, the measured tag importance using discounted probability are 0.6 (left) and 1 (right).

To identify the appearance of an object tag in a particular sentence, we need to map the object tag to the word in that sentence. For example, for the first sentence in Fig. 4, we need to know the word "man" corresponds to tag "person" and "scooter" corresponds to tag "motorbike". We use the WordNet based Semantic distance [36] to measure the similarity between concepts. The effectiveness of this synonymous term matching method has been demonstrated in [24].

**Scene Tag Importance.** To jointly measure scene and object tag importance, we need to consider two specific properties of the scene tag in sentence descriptions of an image.

1) An image usually has fewer scene tags than object tags.
2) The grammatical role of the scene tag in a sentence is a strong indicator of its importance.

The first property results in an imbalance between the scene and object tags. Even if both scene and object tags appear in one sentence, scene importance can be discounted to a lower value if there are many object tags in the sentence. To understand the second property, we consider the following two sentences:

- Sentence 1: A sandy beach covered in white surfboards near the ocean.
- Sentence 2: Surfboards sit on the sand of a beach.

Clearly, the scene tag "beach" is the major constituent of the first sentence because it appears as the main subject of the whole sentence. On the other hand, it becomes the minor constituent in the second sentence because it appears in the modifier phrase of subject "surfboard".

To infer the grammatical role of a scene tag in a sentence, we first leverage the Stanford Lexicalized Probabilistic Context Free Grammar Parser [37], [38] to obtain the sentence parse tree [39]. Two examples are shown in Fig. 6. Then, we identify whether the scene tag appears in a prepositional phrase by checking whether there is a "PP" node in the path from the root to the scene tag leaf.

**Joint Object/Scene Tag Importance.** Based on the above analysis, we propose an algorithm to jointly measure the importance of object and scene tags in one sentence. It is summarized in Algorithm I. The final scene and object tag importance values are obtained by averaging the importance vector $\mathbf{I}$ over all sentences associated with the same image.

Parameters $\alpha$ and $\beta$ in Algorithm I are the scene weights when the scene tag appears as the modifier and the subject of a sentence, respectively. They are used to account for the
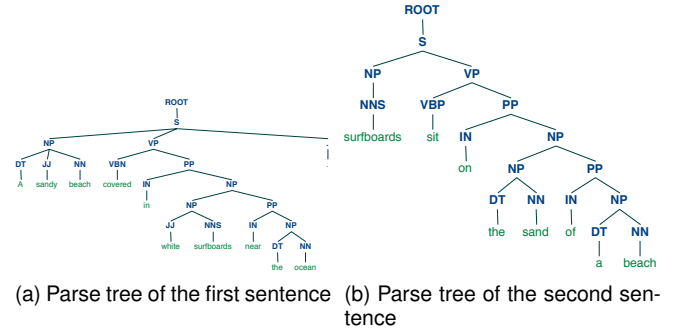


(a) Parse tree of the first sentence  (b) Parse tree of the second sentence

Fig. 6. The parse trees of two sentences. The acronyms in the trees are: S (Sentence), NP (Noun Phrase), VP (Verb Phrase), PP (Preposition Phrase), DT (Determiner), JJ (Adjective), NN (Singular Noun), NNS (Plural Noun), VBN (Verb, past participle), VBP (Verb, non-3rd person singular present), IN (Preposition or subordinating conjunction).

different grammatical roles of the scene tag in a sentence. We set $\alpha = 1$ and $\beta = 2$ in our experiments. The object and scene tag importance values are computed in lines 7 and 8 of Algorithm I, respectively. The formula in line 7 is essentially the discounted probability defined in Eq. (1). The scene tag is adjusted to account for imbalance of scene and object tag numbers.

**Algorithm I: Measuring object/scene tag importance in a sentence description.**

1: Input: The set of all object tags $\mathcal{T}_o$, the set of all object tags mentioned in sentence $\mathcal{T}_o^s$, current sentence $s$, scene tag $t_s$.
2: **Set** scene factor $c_s$ to default value 0;
3: **Set** sentence tree $\mathbf{T}$ to parseSentence($s$);
4: **Set** path $\mathbf{p}$ to findPath($\mathbf{T}, \mathbf{T}.root, t_s$).
5: **If** "PP" in $\mathbf{p}$ and "PP".$left = prep$: **Set** $c_s = \alpha$ ;
6: **Else if** $\mathbf{p} \neq NULL$: **Set** $c_s = \beta$;
7: **For** $t \in \mathcal{T}_o$: $\mathbf{I}(t) = \frac{\mathbb{I}\{t \in \mathcal{T}_o^s\}}{|\mathcal{T}_o^s|(1+c_s)}$;
8: $\mathbf{I}(t_s) = \frac{c_s}{1+c_s}$.

In the following, the measured tag importance will serve as the ground truth for our experiments in tag importance prediction in Section 7.

## 5 PREDICTING TAG IMPORTANCE

The problem of tag importance prediction is studied in this section. First, we discuss three feature types used for prediction. They are semantic, visual and context cues. Then, we describe a prediction model, in which interdependency between tag importance is characterized by the Markov Random Field (MRF) [40]. The model parameters are learned using the Structural Support Vector Machine (SSVM) [41].

### 5.1 Three Feature Types

**Semantic Features.** Some object categories are more attractive to humans than others [24]. For example, given tags "cat", "chair" and "potted plant" for the image in Fig. 7a, people tend to describe the "cat" more often than the "chair" and the "potted plant". The same observation applies to scene importance. For example, in Fig. 7d, people often mention (as observed in the dataset) the whole image as

(a) Object semantic  (b) Object visual

cat: 0.8; chair: 0.2; potted plant: 0     bus: 0.2; car: 0; person: 0.8     bus: 0.9; car: 0; person: 0.1

(c) Object context

dog: 0.6; sofa: 0.4     sofa: 0.8; tv monitor: 0.2

(d) Scene semantic  (e) Scene visual

bathroom: 0.9;     beach: 0.8;     beach: 0.2
sink:0; toilet:0;     person: 0; dog: 0; bench: 0     cow: 0.6; boat: 0.3

(f) Object scene context

living room: 0.3; clock: 0; tv: 0; chair: 0;     living room: 0.8; clock: 0; chair: 0;     living room: 0.9; laptop:0; vase: 0; chair: 0;
dining table: 0; couch: 0.1; person 0.5;     dining table: 0; couch: 0.1; cat: 0.2;     dining table: 0; couch: 0.1; cell phone: 0
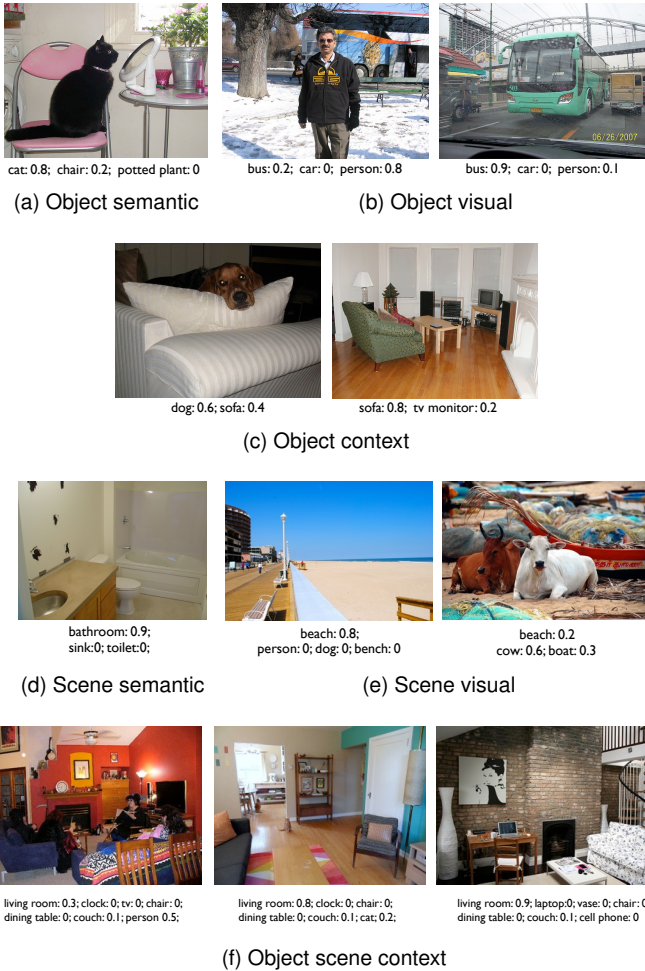
Fig. 7. Examples of various cues for predicting object and scene tag importance. The texts below images give the ground truth tag importance.

a "bathroom" image rather than describing objects "toilet" and "sink" in the image. This is because objects "toilet" and "sink" are viewed as necessary components of a bathroom. When all component objects are combined together to form a scene, people tend to mention the scene as a whole rather than describe each individual object.

The semantic cue can be modeled as a categorical feature. It is a $|\mathcal{C}|$-dimensional binary vector with value one in the $i$th index indicating the $i$th object/scene category, where $|\mathcal{C}|$ denotes the number of different object/scene categories.

**Visual Features.** Human do not consider an object/scene important just because of its category. For the case of object tags, we show a case in Fig. 7b, where both images have tags "bus" and "person". However, their importance differs because of their visual properties. To capture visual cues, we first apply Faster R-CNN [42] (with RPN as object proposal network and Fast R-CNN with VGG16 as detector network) to extract object tags' corresponded bounding boxes, and then calculate the following properties using the detected bounding boxes: 1) the area and log(area) as the size features; 2) the max-, min- and mean-distances to the image center, the vertical mid-line, the horizontal mid-line, and the third box [25] as location features; and 3) relative saliency. For the last item, we use the spectral residual approach

in [43] to generate the saliency map. Even though false detection will affect the tag importance prediction, it can be reduced significantly by simply removing the proposal whose object category is not in the tag list of current image. The performance loss of tag importance prediction caused by object detection error will be studied in Section 7. By concatenating all above features, we obtain a 15-D visual feature vector. An object tag may correspond to multiple object instances in an image. For this case, we add the size and saliency features of all object instances to obtain the corresponding tag size and saliency features, but take the minimum value among all related object instances to yield the tag location feature.

For the scene tag importance visual feature, global patterns such as "openness" and color property of an image are useful for predicting scene tag importance. This is shown in Fig. 7e, where two images have important "beach" and unimportant "beach", respectively. We use the FC7 layer (Fully Connected layer 7) features extracted using VGG16 trained on the Places dataset [44] to model the scene property.

**Object Context Features.** The object context features are used to characterize how the importance of an object tag is affected by the importance of other object tags. When two object tags coexist in an image, their importance is often interdependent. Consider the sample database image in Fig. 3. If the "dog" did not appear, the "bicycle" would be of great importance due to its large size and centered location, and the discounted probability importance of "bicycle" would be 1 based on Eq. (1). To model interdependency, we should consider not only relative visual properties between two object tags but also their semantic categories. Fig. 7c shows two object context examples. For the left image, although the "sofa" has a larger size and a better location, people tend to describe the "dog" more often since the "dog" gets more attention. On the other hand, for the right image, people tend to have no semantical preference between the "sofa" and the "TV monitor". However, due to the larger size of the "sofa", it is perceived as more important by humans.

To extract an object context feature, we conduct two tasks: 1) analyze the relative visual properties within an object tag pair, and 2) identify the tag pair type (i.e. semantic categories of two object tags that form a pair). To model the difference of visual properties for tag pair $(t_i, t_j)$, we use $s_i - s_j$ and $d_i - d_j$ as the relative size and location, respectively, where $s$ and $d$ denote the bounding box area and the corresponding mean distance to the image center as described in the visual feature section. The final object context feature $\mathbf{g}_{ij}^o$ for tag pair $(t_i, t_j)$ is defined as

$$\mathbf{g}_{ij}^o = \left[ (s_i - s_j) \cdot \mathbf{p}_{ij}^\mathsf{T} \quad (d_i - d_j) \cdot \mathbf{p}_{ij}^\mathsf{T} \right]^\mathsf{T}, \qquad (2)$$

where $\mathbf{p}_{ij}$ is the tag pair type vector for tag pair $(t_i, t_j)$.

**Object Scene Context Features.** It is intuitive that scene tag importance and object tag importance are also interdependent. Consider three images in Fig. 7f with the same scene type - "living room". The right image is a classical "living room" scene and all objects are components of the scene structure. For the middle image, object "cat" only takes out a bit of importance of the "living room" due to its semantic interestingness but relative small size. At last,

object tag: car; dog; potted plant
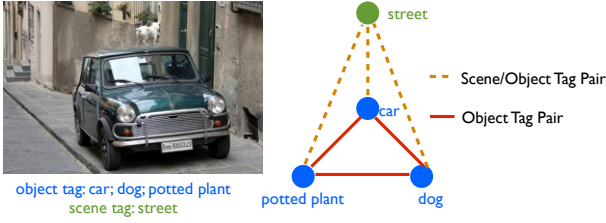scene tag: street

Fig. 8. A sample image and its corresponding joint MRF model.

in the left most image, people become the dominant objects due to their prominent size. Thus, the importance of "living room" has been suppressed by the tag "person". Clearly, by removing the "person" tag in the left image and the "cat" tag in the middle image, the "living room" tags in the three images would be equally important.

To model the context cues between scene and object tag importance, we use a similar approach described in object context features. That is, we define a tag pair type vector $\mathbf{p}_{is}$ to indicate the object and scene tags that form edge $(t_i, t_s)$. It is a categorical vector that models the semantic of a tag pair. Moreover, as discussed above, the size of an object may affect the interaction between the scene and object. Thus, the context feature is defined as

$$\mathbf{g}_{is}^s = s_i \mathbf{p}_{is}, \tag{3}$$

where $s_i$ is the total size of the $i$th object tag's bounding boxes.

## 5.2 Tag Importance Prediction Model

The interdependence of tag importance in an image defines a structured prediction problem [45]. It can be mathematically formulated using the MRF model [40]. Formally, each image is represented by an MRF model denoted by $(V, E)$, where $V = V_o \bigcup v_s$, and where $V_o$ is the set of object tags of the current image and $v_s$ is the current scene tag (assuming scene tag is presented in the image). Consider edge set $E = E_o \bigcup E_{os}$, where $E_o = \{(v_i, v_j) : v_i, v_j \in V_o\}$ is the edge set for object-object tag pairs and $E_{os} = \{(v_i, v_s) : v_i \in V_o\}$ is the edge set for object-scene tag pairs. An exemplary MRF model built for a sample image with scene tag "street" and object tags "car", "dog", "potted plant" is shown in Fig. 8. Specifically, each tag located in a vertex has its own visual and semantic cues to predict importance while each edge enforces the output to be compatible with the relative importance between tags.

Under a log-linear MRF model, the energy function to be minimized can be expressed as

$$E(\mathbf{X}_o, \mathbf{x}_s, \mathbf{G}_o, \mathbf{G}_s, \mathbf{y}; \mathbf{w}) =$$
$$\sum_{i \in V_o} \underbrace{\mathbf{w}_{V_o}^\mathsf{T} \boldsymbol{\varphi}_V(\mathbf{x}_i, y_i)}_{\text{object tag visual \& semantic}} + \sum_{(i,j) \in E_o} \underbrace{\mathbf{w}_{E_o}^\mathsf{T} \boldsymbol{\varphi}_E(\mathbf{g}_{ij}^o, y_i, y_j)}_{\text{object tag pair context}}$$
$$+ \underbrace{\mathbf{w}_{V_s}^\mathsf{T} \boldsymbol{\varphi}_V(\mathbf{x}_s, y_s)}_{\text{scene tag visual \& semantic}} + \sum_{(i,s) \in E_{os}} \underbrace{\mathbf{w}_{E_{os}}^\mathsf{T} \boldsymbol{\varphi}_E(\mathbf{g}_{is}^s, y_i, y_s)}_{\text{object scene tag pair context}}, \tag{4}$$

where $\mathbf{y} = \{y_i\}$ is the predicted tag importance output vector, $\mathbf{X}_o = \{\mathbf{x}_i\}$ and $\mathbf{x}_s$ are the concatenation of object and scene tag visual and semantic feature vectors as described

in Section 5.1, respectively. $\mathbf{G}_o = \left\{\mathbf{g}_{ij}^o\right\}$ is the object context feature vector calculated using Eq. (2), $\mathbf{G}_s = \{\mathbf{g}_{is}^s\}$ is the object scene context feature vector calculated using Eq. (3).

The weight vector $\mathbf{w} = \left[\mathbf{w}_{V_o}^\mathsf{T}, \mathbf{w}_{E_o}^\mathsf{T}, \mathbf{w}_{v_s}^\mathsf{T}, \mathbf{w}_{E_{os}}^\mathsf{T}\right]^\mathsf{T}$ in Eq. (4) will be learned from training data. $\boldsymbol{\varphi}_V$ and $\boldsymbol{\varphi}_E$ are joint kernel maps [46] defined as $\boldsymbol{\varphi}_V(\mathbf{x}_i, y_i) \equiv \mathbf{x}_i \otimes \boldsymbol{\delta}(y_i), \boldsymbol{\varphi}_E(\mathbf{g}_{ij}, y_i, y_j) \equiv \mathbf{g}_{ij} \otimes \boldsymbol{\delta}(y_i - y_j)$ where $\otimes$ is the Kronecker product,

$$\boldsymbol{\delta}(y_i) \equiv [\mathbb{I}\{y_i = 0\}, \mathbb{I}\{y_i = 0.1\}, \cdots, \mathbb{I}\{y_i = 1\}],$$
$$\boldsymbol{\delta}(y_i - y_j) \equiv [\mathbb{I}\{y_i - y_j = -1\}, \cdots, \mathbb{I}\{y_i - y_j = 1\}],$$

where $\mathbb{I}$ is the indicator function. It is worthwhile to mention that the ground truth tag importance usually takes certain discrete values in experiments, and it does not affect the retrieval performance if it is rounded to the nearest tenth. Thus, the ground truth tag importance is quantized into 11 discrete levels, starting from 0 to 1 with intervals of 0.1. This leads to an 11-D $\boldsymbol{\delta}(y_i)$ vector and a 21-D $\boldsymbol{\delta}(y_i - y_j)$ vector. There are two reasons to define $\boldsymbol{\delta}(y_i - y_j)$ such a form. First, relative importance can be quantified. Second, dimensionality of the vector $\mathbf{w}$ can be greatly reduced.

The above model can be simplified to yield the binary-valued tag importance as done in [24]. This is achieved by treating $y_i$ as a binary class label and redefining $\boldsymbol{\delta}(y_i)$ and $\boldsymbol{\delta}(y_i - y_j)$. The performance of this simplified model will be reported in the first half of Section 7.4.

**Learning.** To learn model parameter $\mathbf{w}$, one straightforward way is to apply the probabilistic parameter learning approach [45], i.e., treating the energy function in Eq. (4) as the negative of the log likelihood of data and applying gradient-based optimization. However, this learning approach ignores the ordinal nature of the output importance label. For example, if a tag has ground truth importance value 1, the predicted importance 0 will be penalized the same as the predicted importance 0.9. This clearly deviates from intuition. On the other hand, the Loss Minimizing Parameter Learning approach such as the Structural Support Vector Machine (SSVM) [41] allows a customizable loss function for different prediction tasks. It can be exploited by taking the ordinal nature of output importance label into account. As a result, we use the SSVM to learn weight vector $\mathbf{w}$ and adopt one slack variable with the margin rescaling formulation in [47]. The optimization problem becomes

$$\min_{\mathbf{w}, \xi} \frac{1}{2}\|\mathbf{w}\|^2 + C\xi$$
$$\text{s.t.} \forall (\bar{\mathbf{y}}_1, \cdots, \bar{\mathbf{y}}_N) \in \mathcal{Y}^N \tag{5}$$
$$\frac{1}{N}\sum_{n=1}^{N} \mathbf{w}^\mathsf{T} \boldsymbol{\delta}_n(\bar{\mathbf{y}}_n) \geq \frac{1}{N}\sum_{n=1}^{N} \Delta(\hat{\mathbf{y}}_n, \bar{\mathbf{y}}_n) - \xi$$

where $\hat{\mathbf{y}}_n$ is the ground truth tag importance vector for the $n$th training image, $\mathcal{Y}^N$ is the set of all possible output for the training dataset, $C$ is the regularization parameter and $\xi$ is the slack variable, $\boldsymbol{\delta}_n(\bar{\mathbf{y}}_n)$ is the $\boldsymbol{\delta}(\bar{\mathbf{y}})$ for $n$th image where

$$\boldsymbol{\delta}(\bar{\mathbf{y}}) = \Psi(\mathbf{X}_o, \mathbf{x}_s, \mathbf{G}_o, \mathbf{G}_s, \hat{\mathbf{y}}) - \Psi(\mathbf{X}_o, \mathbf{x}_s, \mathbf{G}_o, \mathbf{G}_s, \bar{\mathbf{y}}),$$

and

$$\Psi\left(\mathbf{X}_o, \mathbf{x}_s, \mathbf{G}_o, \mathbf{G}_s, \mathbf{y}\right) = -\begin{bmatrix} \sum_{i \in V_o} \boldsymbol{\varphi}_V(\mathbf{x}_i, y_i) \\ \sum_{(i,j) \in E_o} \boldsymbol{\varphi}_E(\mathbf{g}_{ij}^o, y_i, y_j) \\ \boldsymbol{\varphi}_V(\mathbf{x}_s, y_s) \\ \sum_{(i,s) \in E_{os}} \boldsymbol{\varphi}_E(\mathbf{g}_{is}^s, y_i, y_s) \end{bmatrix}.$$

Then, we get

$$\mathbf{w}^\mathsf{T}\Psi\left(\mathbf{X}_o, \mathbf{x}_s, \mathbf{G}_o, \mathbf{G}_s, \mathbf{y}\right) = -E(\mathbf{X}_o, \mathbf{x}_s, \mathbf{G}_o, \mathbf{G}_s, \mathbf{y}; \mathbf{w}).$$

In the weight learning process, we define the following loss function:

$$\Delta(\hat{\mathbf{y}}, \bar{\mathbf{y}}) \equiv \frac{1}{|V|}\sum_{i \in V} |\hat{y}_i - \bar{y}_i|, \tag{6}$$

which is the Mean Absolute Difference (MAD) between the ground truth and the predicted tag importance values of one image. Finally, we applied the standard cutting plane algorithm [47] to optimize and obtain the final weight vector $\mathbf{w}$.

**Inference.** After learning the weight vector $\mathbf{w}$, we can determine the vector $\mathbf{y}$ that minimizes Eq. (4). Moreover, finding the maximum violated constraint in the cutting plane training also needs inference. Despite the fully connected graph structure in the MRF model, the number of tags in an image is usually limited. As a result, even if we try all possible outputs, the computational complexity is still acceptable. Our experimental results show that inference takes approximately only 0.2s per image in C on a 2.4GHz CPU 4GB RAM PC. For this reason, we adopt the exact inference approach in this work, and we will investigate fast inference techniques in future work.

## 6 MULTIMODAL IMAGE RETRIEVAL

In this section, we discuss our MIR/TIP system by employing CCA/KCCA. First, we will review the CCA and KCCA. Then we will describe the visual and textual features used in our MIR/TIP experiments. Note that other learning methods introduced in Section 2 can also be used in place of CCA/KCCA.

### 6.1 CCA and KCCA

In multimodal image retrieval, an image is associated with both visual feature vector $\mathbf{f}_v$ and textual feature vector $\mathbf{f}_t$ (e.g. the tag vector). Given these feature pairs $(\mathbf{f}_v^{(i)}, \mathbf{f}_t^{(i)})$ for $N$ images, two design matrices $\mathbf{F}_v \in \mathbb{R}^{N \times D_v}$ and $\mathbf{F}_t \in \mathbb{R}^{N \times D_t}$ can be generated, where the $i$th row in $\mathbf{F}_v$ and $\mathbf{F}_t$ correspond to $\mathbf{f}_v^{(i)}$ and $\mathbf{f}_t^{(i)}$ respectively. CCA aims at finding a pair of matrices $\mathbf{P}_v \in \mathbb{R}^{D_v \times c}$ and $\mathbf{P}_t \in \mathbb{R}^{D_t \times c}$ that project visual and textual features into a common $c$ dimensional subspace with maximal normalized correlation:

$$\max_{\mathbf{P}_v, \mathbf{P}_t} \text{trace}\left(\mathbf{P}_v^\mathsf{T}\mathbf{F}_v^\mathsf{T}\mathbf{F}_t\mathbf{P}_t\right)$$
$$\text{s.t.} \quad \mathbf{P}_v^\mathsf{T}\mathbf{F}_v^\mathsf{T}\mathbf{F}_v\mathbf{P}_v = \mathbf{I}, \ \mathbf{P}_t^\mathsf{T}\mathbf{F}_t^\mathsf{T}\mathbf{F}_t\mathbf{P}_t = \mathbf{I}. \tag{7}$$

The above optimization problem can be reduced to a generalized eigenvalue problem [11], and the eigenvectors corresponding to the largest $c$ eigenvalues are stacked horizontally to form $\mathbf{P}_v$ and $\mathbf{P}_v$.

To measure the similarity of projected features in subspace to achieve cross-modality retrieval, we adopted the Normalized CCA metric proposed in [10], [48]. After solving the CCA problem in Eq. (7), the similarity between visual features $\mathbf{F}_v$ and textual features $\mathbf{F}_t$ will be computed as:

$$\frac{\left(\mathbf{F}_v\mathbf{P}_v\text{diag}\left(\lambda_1^t, \cdots, \lambda_c^t\right)\right)\left(\mathbf{F}_t\mathbf{P}_t\text{diag}\left(\lambda_1^t, \cdots, \lambda_c^t\right)\right)^\mathsf{T}}{\|\mathbf{F}_v\mathbf{P}_v\text{diag}\left(\lambda_1^t, \cdots, \lambda_c^t\right)\|_2 \|\mathbf{F}_t\mathbf{P}_t\text{diag}\left(\lambda_1^t, \cdots, \lambda_c^t\right)\|_2}, \tag{8}$$

where $\lambda_1, \cdots, \lambda_c$ correspond to the top $c$ eigenvalues, and $t$ is the power of the eigenvalues (we set $t$=4 as in [10], [48]).

To model nonlinear dependency between visual and textual feature vectors, a pair of nonlinear transforms, $\Phi_v$ and $\Phi_t$, are used to map visual and textual features into high dimensional spaces, respectively. With kernel functions

$$\mathbf{K}_m(\mathbf{f}_m^{(i)}, \mathbf{f}_m^{(j)}) = \Phi_m(\mathbf{f}_m^{(i)})^\mathsf{T}\Phi_m(\mathbf{f}_m^{(j)}) \quad m = v, t$$

$\Phi_v$ and $\Phi_t$ are only computed implicitly. This kernel trick lead to KCCA, which attempts to find the maximally correlated subspace with the two transformed spaces [11]. However, since the time and space complexity of KCCA is $O(N^2)$, it is not practically applicable to large scale image retrieval. We thus adopt the CCA for the retrieval experiments.

We also tried out the scalable KCCA proposed in [10] by constructing the approximate kernel mapping but find almost no performance improvement in our experiment setting.

### 6.2 Retrieval Features

Features used in CCA for MIR/TIP and experimental settings are given below.

**Visual Features.** To capture both object and scene properties, we use the VGG16 trained on the ImageNet [49] and the Places [44] to extract visual features. The output of the FC7 (Fully Connected layer 7) for both networks are concatenated together to form a 8192-D visual feature vector.

**Textual Features.** We consider 5 types of textual features, i.e. 1) the tag vector, 2) the predicted binary-valued tag importance vector, 3) the true binary-valued tag importance vector, 4) the predicted continuous-valued tag importance vector, and 5) the true continuous-valued tag importance vector, each of which is used in a retrieval experimental setting as discussed in Section 7.

## 7 EXPERIMENTAL RESULTS

In this section, we first discuss the datasets and our experimental settings. We proceed to present our subjective test results to justify that the ground truth tag importance based on descriptive sentences is consistent with human perception. Then, we compare the performance of different tag importance prediction models. Finally, experiments on three retrieval tasks are conducted to demonstrate the superior performance of the proposed MIR/TIP system.

### 7.1 Datasets

To test the proposed system, we need datasets that have many annotated data available, including sentence descriptions, object tags, object bounding boxes and scene tags. Table 2 lists the profiles of several image datasets in the public domain. Among them, the UIUC, COCO, and VisualGenome datasets appear to meet our need the most since

they have descriptive sentences for each image. However, the Visual Genome dataset aims to be an "open vocabulary" dataset with 80,138 object categories but on average only 49 instances for each object category. While the large tag vocabulary size makes it challenging to learn CCA subspace, limited number of instances per category makes Faster R-CNN training difficult. We thus adopt the UIUC and COCO datasets to conduct our experiments, and leave the Visual Genome dataset to our future work. To test our proposed full system with scene tag, we enrich COCO with scene tags and call it **COCO Scene dataset** (will be discussed later in this section). In the following, we consider 2 types of datasets: 1) the datasets with only object tags (e.g. COCO, UIUC) and 2) the dataset with both object and scene tags (e.g. COCO Scene).

TABLE 2
Comparison of major image datasets.

| Dataset | Sentences | Object Tag | Bounding Box | Scene Tag |
|---|---|---|---|---|
| UIUC [50] | Yes | Yes | Yes | No |
| COCO [51] | Yes | Yes | Yes | No |
| SUN2012 [52] | No | Yes | Yes | Yes |
| LabelMe [53] | No | Yes | Yes | No |
| ImageNet [54] | No | Yes | Yes | No |
| Visual Genome [55] | Yes | Yes | Yes | No |

For the datasets with only object tags, we adopt UIUC and full COCO datasets for small and large scale experiments respectively. Moreover, the UIUC [50] dataset was used for image importance prediction in [24] and can serve as a convenient benchmarking dataset. Each image in these two datasets is annotated by 5 sentence descriptions and each object instance in an image is labeled with a bounding box. The UIUC dataset consists of 1000 images with objects from 20 different categories. The COCO dataset has in total 123,287 images with objects from 80 categories. We use object categories as object tags. Thus, UIUC has approximately 1.8 tags per image while COCO has on average 3.4 tags per image. We mainly use UIUC dataset for importance prediction experiment since it is not a challenging dataset for retrieval.

For dataset with both object and scene tags, we generated an experimental dataset based on images in the COCO dataset. Specifically, we first identified 30 common scene types in the COCO dataset. Then, 50 human workers were invited to manually classify 60,000 images randomly drawn from the COCO dataset into one of 31 groups, which include 30 scene types mentioned above and an extra group indicated as "Not sure/None of above". This is necessary as the COCO dataset contains a large amount of object centric images, whose scene types are hard to identify even for human. This resulted in a dataset consisting of 25,124 images with a tag vocabulary of 110. It has on average 4.3 tags per image. We will refer to this dataset as the COCO Scene dataset for the rest of this section. (For the statistics of COCO Scene dataset, please refer to supplementary material.)

## 7.2 Retrieval Experiment Settings

As introduced in Section 6, different tag features correspond to different MIR experiment settings because the semantic subspace is determined by applying CCA to visual and textual features. We thus compared the following 5 MIR settings:

- **Traditional MIR**: Textual features are the binary-valued tag vectors. This is the benchmark method used in [8], [9], [10].
- **MIR/PBTI**: Textual features are **P**redicted **B**inary-valued **T**ag **I**mportance vectors. This corresponds to the predicted importance proposed in [24].
- **MIR/PCTI**: Textual features are **P**redicted **C**ontinuous-valued **T**ag **I**mportance vectors. This is *our proposed system*.
- **MIR/TBTI**: Textual features are **T**rue **B**inary-valued **T**ag **I**mportance vectors. This serves as the upper bound for the binary-valued tag importance proposed in [24].
- **MIR/TCTI** : Textual features are **T**rue **C**ontinuous-valued **T**ag **I**mportance vectors. This gives the best retrieval performance, which serves as the performance bound.

Among the above five systems, the last two are not achievable since they assume the tag importance prediction to be error free.

Moreover, we evaluate our system in terms of 3 retrieval tasks:

- **I2I** (Image to Image retrieval): Given a query image, the MIR systems will project the visual features into the CCA subspace and rank the database images according to Eq. (8). We also test a baseline retrieval system (**Visual Only**) that ranks the database images using visual features' Euclidean distance.
- **T2I** (Tag to Image retrieval): Given a tag list, the MIR systems will project the tag feature into the CCA subspace and rank the database images according to Eq. (8). Note our system can support weighted tag list as query as in [10], in which the weights represent the importance of tags.
- **I2T** (Image annotation): Given a query image, the MIR systems will find 50 nearest neighbors in the CCA subspace and use their textual features to generate an average textual feature vectors, based on which the tags in tag vocabulary will be ranked. We also test a baseline tagging system using deep features to find nearest neighbors and their corresponding tag vectors to rank tags.

For all retrieval tasks, we adopt the Normalized Discounted Cumulative Gain (NDCG) as the performance metric since it is a standard and commonly used metric [8], [9], [32]. Moreover, it helps quantify how an MIR system performs. The NDCG value for the top k results is defined as $\text{NDCG@k} = \frac{1}{Z} \sum_{i=1}^{k} \frac{2^{r_i}-1}{\log_2(i+1)}$, where $r_i$ is a relevance index (or function) between the query and the $i$th ranked image, and $Z$ is a query-specific normalization term that ensures the optimal ranking with the NDCG score of 1. The relevance index measures the similarity between retrieved results and the query in terms of ground truth continuous-valued tag importance, i.e. whether an MIR system can preserve important content of the query in retrieval results or not. For I2T retrieval task, the relevance of a tag to the

query image is set as its ground truth continuous-valued tag importance. The choice of the relevance index for the other two tasks will be discussed in detail in the next section.

## 7.3 Subjective Test Performance of Measured Tag Importance

Since ground truth continuous-valued tag importance is used to measure the degree of object/scene importance in an image as perceived by a human, it is desired to design a subjective test to evaluate its usefulness. Here, we would like to evaluate it by checking how much it will help boost the retrieval performance. Specifically, we compare the performance of two different relevance functions and see whether the defined ground truth continuous-valued tag importance correlates human experience better. The two relevance functions are given below.

1) The relevance function with measured ground truth tag importance:

$$r_g(p, q) = \frac{\langle \mathbf{I}_p, \mathbf{I}_q \rangle}{\|\mathbf{I}_p\| \|\mathbf{I}_q\|}, \qquad (9)$$

where $\mathbf{I}_k$ denotes the ground truth continuous-valued importance vector for image $k$ ($k = p$ or $q$).

2) The relevance function with binary-valued importance [24] (whether appeared in sentences or not):

$$r_b(p, q) = \frac{\langle \mathbf{t}_p, \mathbf{t}_q \rangle}{\|\mathbf{t}_p\| \|\mathbf{t}_q\|}, \qquad (10)$$

where $\mathbf{t}_k$ denotes the binary-valued tag importance vector for image $k$ ($k = p$ or $q$).

In the experiment, we randomly selected 500 image queries from the COCO Scene dataset and obtained the top two retrieved results with the max relevance scores using two relevance functions mentioned above. In the subjective test, we presented the two retrieved results of the same query in pairs to the subject, and asked him/her to choose the better one among the two. We invited five subjects (one female and four males with their ages between 25 and 30) to take the test. There were 1500 pairwise comparisons in total. We randomized the order of two relevance functions in the GUI to minimize the bias. Moreover, each subject viewed each query at most once. We made the following observation from the experiment. As compared with the results using the relevance function with binary-valued importance $r_b$, the results using our relevance function $r_g$ were favored in 1176 times (out of 1500 or 78.4%). This indicates that the relevance function with ground truth importance $r_g$ does help improve the retrieval performance and it also demonstrates the validity of the proposed methodology in extracting sentenced-based ground truth tag importance.

## 7.4 Performance of Tag Importance Prediction

To evaluate tag importance prediction performance, we first compare the performance of the state-of-the-art method with that of the proposed tag importance prediction method on UIUC dataset. Then, under continuous-valued tag importance setting, we study the effect of different feature types on the proposed tag importance prediction model,

along with the loss introduced by binary-valued importance for all datasets introduced in Section 7.1.

For the purpose of performance benchmarking, we simplified our structured model as discussed in Sec. 5.2 to achieve binary-valued tag importance prediction and compared with [24]. Here we use accuracy as the evaluation metric. Same as in [24], accuracy is defined as the percentage of correctly classified object instances against the total number of object instances in test images. For fair comparison, we also ran 10 simulations of 4-fold cross validation, and compared the mean and standard deviation of estimated accuracy. The performance comparison results are shown in Table 3. The baseline method simply predicts "yes" (or important) for every object instance while the next column refers to the best result obtained in the work of Berg *et al.* [24]. Our simplified structured model can further improve the prediction accuracy of [24] by 5.6%.

TABLE 3
Performance comparison of tag importance prediction.

| Methods | Baseline | Berg *et al.* [24] | Proposed |
|---|---|---|---|
| Accuracy Mean | 69.7% | 82.0% | 87.6% |
| Accuracy STD | 1.3 | 0.9 | 0.7 |

Next, we evaluate the continuous-valued tag importance prediction performance of 7 different models. Here we use prediction error as the evaluation metric, which is defined as the average MAD in Eq. (6) across all test images. These 7 models are:

1) Equal importance of all tags (called the Baseline);
2) Visual features only (denoted by Visual);
3) Visual and semantic features (denoted by Visual+Semantic);
4) Visual, semantic and context features (denoted by our Model);
5) Visual, semantic and context features with ground truth bounding boxes (denoted by our Model/True bbox);
6) Equal importance of tags that are mentioned in any sentence. This corresponds to the true binary-valued tag importance computed as in [24] (denoted by binary true);
7) Equal importance of tags that are predicted as "important" using the model proposed in [24] (denoted by binary predicted).

For the 2nd and 3rd models, we adopt the ridge regression models trained by visual features and visual plus semantic features, respectively. The 4th model is the proposed structured model as described in Section 5.2. These 3 models help us to understand the impact of different feature types as described in Sec. 5.1. The 5th model differs from the 4th model in that it uses ground truth bounding boxes to compute the visual features. It enables us to identify how object detection error will affect the tag importance prediction. For the 6th and 7th models, they are used to quantify how the binary-valued tag importance (i.e. treating important tags as equally important) results in tag importance prediction error, which serves as an indicator of performance loss in retrieval. Specifically, the 6th and 7th models will generate

the true and predicted important tags using the method proposed in [24], respectively. Then, the important tags within the same image will be treated as equally important and assigned the same continuous-valued importance. Note that the 6th model is not achievable but only serves as the best case for the 7th model.

We used 5-fold cross validation to evaluate these prediction models. The prediction errors for UIUC and COCO datasets are shown in Figs. 9(a) and (b), respectively. For COCO Scene dataset, we show the prediction error for all, object, and scene tags in Figs. 9(c).

These figures show that our proposed structured prediction model (4th) can achieve approximately **40%** performance gain with respect to the baseline (1st). For COCO and COCO Scene dataset, we observe performance gain of all 3 feature types, among which visual, semantic, and context features result in approximately 28%, 11%, and 11% prediction error reduction respectively. By comparing the 4th to the 5th model, we find the object detection error only results in approximately 3% performance loss. Moreover, it is noted that even true binary-valued tag importance lead to non negligible prediction error by ignoring relative importance between tags, and this error will propagate to predicted binary-valued tag importance model, resulting in 48% and 45% performance loss over our proposed model on COCO and COCO Scene dataset, respectively. Lastly, for COCO Scene dataset, it is observed that scene tag importance is more difficult to predict as compared to object tag importance. Thus, the overall error (average over both object and scene tag importance) is higher than the average error of object tag importance but lower than that of scene tag importance.

The tag importance prediction performance on UIUC dataset differs from that of COCO and COCO Scene in two parts: 1) the visual features results in 1% performance loss compared with baseline; 2) the binary-valued tag importance has less performance loss compared to COCO and COCO Scene. The above phenomena were caused by the bias of the UIUC dataset, in which 454 out of 1000 images have only one tag, and 415 of them have tag with importance value 1. This bias makes modeling the relative importance between tags within the same image insignificant. Thus, the baseline and binary-valued tag importance based models can achieve reasonable performance on UIUC dataset but not on COCO and COCO Scene datasets.

### 7.5 Performance of Multimodal Image Retrieval

In this subsection, we show the retrieval experimental results on COCO and COCO Scene datasets using the settings given in Section 7.2. For both datasets, we randomly sampled 10% of images as queries and the other as database images. Among the database images, 50% were used as MIR training images. For I2I and I2T experiments, we directly used the image as the query. For T2I experiment, weighted tag list based on ground truth tag importance vector was used as query.

**I2I Results.** The NDCG curves of COCO and COCO Scene datasets are shown in Fig. 10(a), (b), respectively. We have the following observations from the plots. First, for all datasets, we find significant improvement of all MIR



(a) UIUC dataset.     (b) COCO dataset.
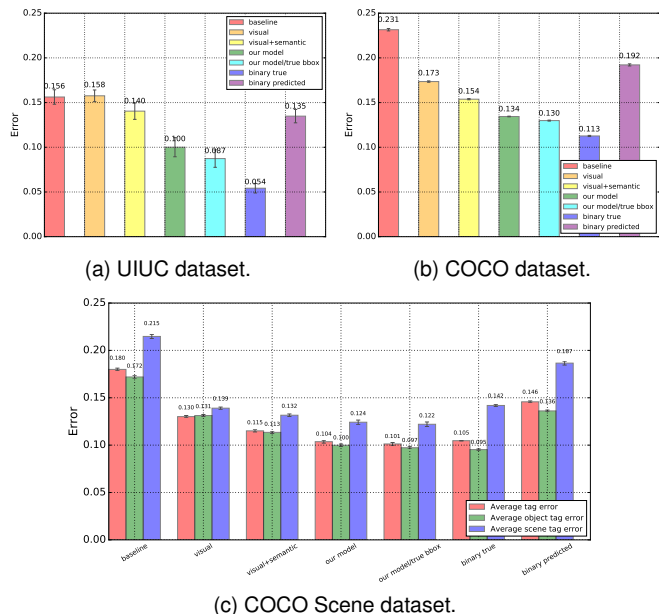
(c) COCO Scene dataset.

Fig. 9. Comparison of continuous-valued tag importance prediction errors of seven models: (a) the UIUC dataset, (b) the COCO dataset, and (c) the COCO Scene dataset.

systems over visual baseline. It seems deep features can give reasonable performance on retrieving the most similar image, but its efficiency lags behind MIR systems as $K$ becomes larger. Second, our proposed MIR/PCTI system exhibits considerable improvements over other practical MIR systems, including Traditional MIR and MIR/PBTI. Specifically, for $K = 50$ (the typical number of retrieved images user is willing to browse for one query), the MIR/PCTI can achieve approximately 14% and 10% gain over visual baseline, 4% and 2% over Traditional MIR, and 2% and 2% over MIR/PBTI on COCO and COCO Scene datasets, respectively. Moreover, our proposed MIR/PCTI system can even match the upper bound of MIR system using binary-valued importance, and it only has 1% and 2% performance gap with its upper bound. Finally, by associating Fig. 10 to Fig. 9, we can identify that the tag importance prediction performance roughly correlates to retrieval performance. Thus, better tag importance prediction leads to better I2I retrieval performance. Some qualitative I2I retrieval results are shown in Fig. 11. Generally speaking, our proposed system can capture the overall semantic of queries more accurately, such as "person playing wii in the living room" for the 1st query and "person playing frisbee in yard area" for the 2nd query, while the remaining 3 systems fail to preserve some important objects such as "remote" or "frisbee".

**T2I Results.** We show the NDCG curves of T2I results on COCO and COCO Scene datasets in Fig. 12 (a) and (b), respectively. Our proposed MIR/PCTI model shows consistent superior performance over Traditional MIR and MIR/PBTI on both datasets. Particularly, for $K$=50, the MIR/PCTI outperforms the Traditional MIR by 4% and 11%, and the MIR/PBTI by 2% and 5% on COCO and COCO Scene datasets, respectively. Moreover, the proposed MIR/PCTI system only has 1% and 2% performance gap with its upper bound on the two datasets. Fig. 13 shows

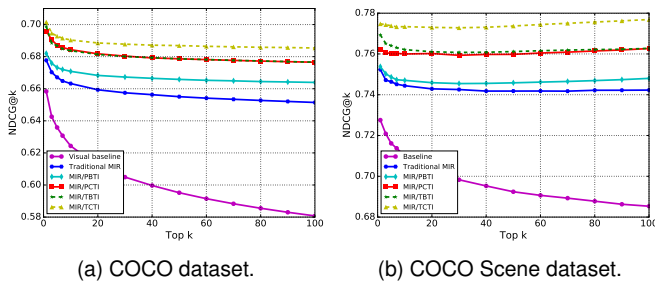(a) COCO dataset.  (b) COCO Scene dataset.

Fig. 10. The NDCG curves for the image-to-image retrieval on the (a) COCO and (b) COCO Scene datasets. The dashed lines are upper bounds for importance based MIR systems.
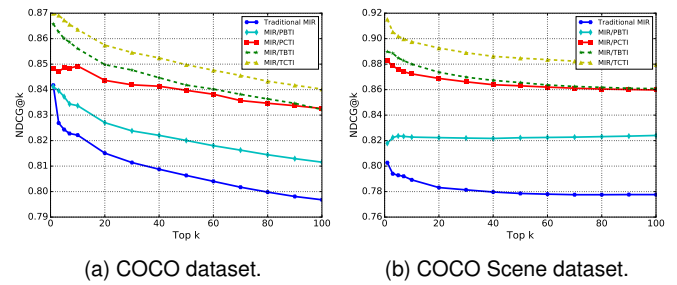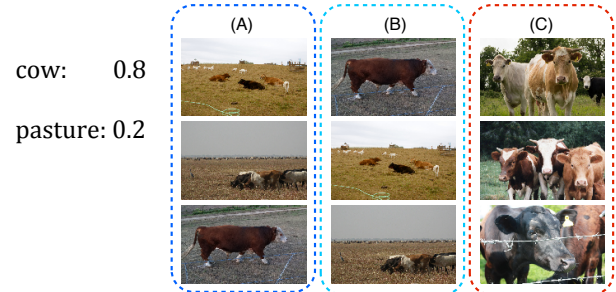


(a) COCO dataset.  (b) COCO Scene dataset.

Fig. 12. The NDCG curves for the tag-to-image retrieval on the (a) COCO and (b) COCO Scene datasets. The dashed lines are upper bounds for importance based MIR systems.
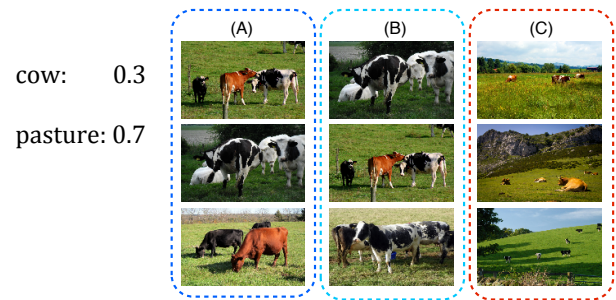


Fig. 11. Top three I2I retrieved results for two exemplary queries, where the four columns show four retrieval systems: (A) Visual Baseline, (B) Traditional MIR, (C) MIR/PBTI, and (D) MIR/PCTI.



(a) A input tag query seeking object centric image.



(b) A input tag query seeking scene centric image.

Fig. 13. Tag-to-Image retrieval results for two exemplary query with different focus, where the three columns correspond to the top three ranked tags of three MIR systems: (A)Traditional MIR, (B) MIR/PBTI, and (C) MIR/PCTI.
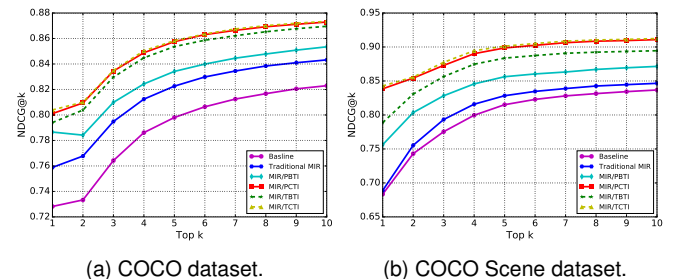


(a) COCO dataset.  (b) COCO Scene dataset.

Fig. 14. The NDCG curves for auto ranked tag list generation on the (a) COCO and (b) COCO Scene datasets. The dashed lines are upper bounds for importance based MIR systems.

the two qualitative results of T2I retrieval, where the two input queries consist of the same tag pair but have different focus. It is observed that our proposed system can correctly retrieve scene/object centric images as indicated by the importance value, while the other two systems can not.

**I2T Results.** The results of tagging on COCO and COCO Scene datasets are shown in Fig. 14 (a) and (b), respectively. Again, we observe consistent improvements of MIR/PCTI over Traditional MIR and MIR/PBTI. Specifically, for $K = 3$ (the typical number of tags each image have in these two datasets), the MIR/PCTI can achieve approximately 8% and 13% gain over baseline, 5% and 10% over Traditional MIR, and 3% and 5% over MIR/PBTI on COCO and COCO Scene datasets, respectively. More surprisingly, its performance can outperform the upper bound of MIR/PBTI and match that of MIR/TCTI. This suggests that our proposed system can not only generate tags but also rank them according to their importance. Sample qualitative tagging results are shown in Fig. 15, in which we can see that our proposed

model will rank more important tags ("cow" and "cat") ahead of unimportant/wrong ones ("person" and "bathroom").
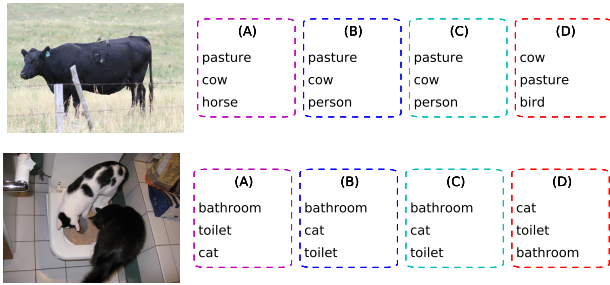
Fig. 15. Tagging results for two exemplary images, where the four columns correspond to the top three ranked tags of four MIR systems: (A) Baseline, (B) Traditional MIR, (C) MIR/PBTI, and (D) MIR/PCTI.
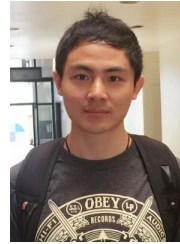
# 8 CONCLUSION AND FUTURE WORK

A multimodal image retrieval scheme based on tag importance prediction (MIR/TIP) was proposed in this work. Both object and scene tag importance were measured from human sentence descriptions and used as the ground truth based on a Natural Language Processing methodology. Three types of features (namely, semantic, visual and context) were identified and a structured model was proposed for object and scene tag importance prediction. Model parameters were trained using the Structural Support Vector Machine formulation. It was shown by experimental results that the proposed sentence-based tag importance measure and the proposed tag importance prediction can significantly boost the performance of various retrieval tasks.

To make the system more practical to real world applications, it is worthwhile to extend the importance idea to other types of tags. Specifically, it would be interesting to extend the idea to scene graph based Visual Genome dataset [55] with more densely annotated object bounding boxes, attributes and relationship information. Furthermore, it is promising to improve retrieval performance by training Convolutional Neural Network and Canonical Correlation Analysis jointly using end-to-end learning. Finally, how to jointly optimize tag importance prediction and retrieval system remains an open question and will be explored in our future work.

## REFERENCES

[1] Y. Rui, T. S. Huang, and S.-F. Chang, "Image retrieval: Current techniques, promising directions, and open issues," *Journal of visual communication and image representation*, vol. 10, no. 1, pp. 39–62, 1999.

[2] L. Chen, D. Xu, I. W. Tsang, and J. Luo, "Tag-based web photo retrieval improved by batch mode re-tagging," in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. IEEE, 2010, pp. 3440–3446.

[3] Y. Liu, D. Xu, I. W.-H. Tsang, and J. Luo, "Textual query of personal photos facilitated by large-scale web data," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 33, no. 5, pp. 1022–1036, 2011.

[4] R. Datta, D. Joshi, J. Li, and J. Z. Wang, "Image retrieval: Ideas, influences, and trends of the new age," *ACM Computing Surveys (CSUR)*, vol. 40, no. 2, p. 5, 2008.

[5] A. W. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain, "Content-based image retrieval at the end of the early years," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 22, no. 12, pp. 1349–1380, 2000.

[6] R. Datta, J. Li, and J. Z. Wang, "Content-based image retrieval: approaches and trends of the new age," in *Proceedings of the 7th ACM SIGMM international workshop on Multimedia information retrieval*. ACM, 2005, pp. 253–262.

[7] M. S. Lew, N. Sebe, C. Djeraba, and R. Jain, "Content-based multimedia information retrieval: State of the art and challenges," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 2, no. 1, pp. 1–19, 2006.

[8] S. J. Hwang and K. Grauman, "Accounting for the relative importance of objects in image retrieval." in *BMVC*, 2010, pp. 1–12.

[9] S. J. Hwang and K. Grauman, "Learning the relative importance of objects from tagged images for retrieval and cross-modal search," *International journal of computer vision*, vol. 100, no. 2, pp. 134–153, 2012.

[10] Y. Gong, Q. Ke, M. Isard, and S. Lazebnik, "A multi-view embedding space for modeling internet images, tags, and their semantics," *International journal of computer vision*, vol. 106, no. 2, pp. 210–233, 2014.

[11] D. Hardoon, S. Szedmak, and J. Shawe-Taylor, "Canonical correlation analysis: An overview with application to learning methods," *Neural computation*, vol. 16, no. 12, pp. 2639–2664, 2004.

[12] N. Rasiwasia, J. Costa Pereira, E. Coviello, G. Doyle, G. R. Lanckriet, R. Levy, and N. Vasconcelos, "A new approach to cross-modal multimedia retrieval," in *Proceedings of the international conference on Multimedia*. ACM, 2010, pp. 251–260.

[13] J. Costa Pereira, E. Coviello, G. Doyle, N. Rasiwasia, G. R. Lanckriet, R. Levy, and N. Vasconcelos, "On the role of correlation and abstraction in cross-modal multimedia retrieval," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 36, no. 3, pp. 521–535, 2014.

[14] H. Zhang, Y. Yang, H. Luan, S. Yang, and T.-S. Chua, "Start from scratch: Towards automatically identifying, modeling, and naming visual attributes," in *Proceedings of the ACM International Conference on Multimedia*. ACM, 2014, pp. 187–196.

[15] G. Andrew, R. Arora, J. Bilmes, and K. Livescu, "Deep canonical correlation analysis," in *Proceedings of the 30th International Conference on Machine Learning*, 2013, pp. 1247–1255.

[16] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, "Multimodal deep learning," in *Proceedings of the 28th international conference on machine learning (ICML-11)*, 2011, pp. 689–696.

[17] W. Wang, X. Yang, B. C. Ooi, D. Zhang, and Y. Zhuang, "Effective deep learning-based multi-modal retrieval," *The VLDB Journal*, pp. 1–23, 2015.

[18] J. Johnson, L. Ballan, and F.-F. Li, "Love thy neighbors: Image annotation by exploiting image metadata," *arXiv preprint arXiv:1508.07647*, 2015.

[19] K. Barnard, P. Duygulu, D. Forsyth, N. De Freitas, D. M. Blei, and M. I. Jordan, "Matching words and pictures," *The Journal of Machine Learning Research*, vol. 3, pp. 1107–1135, 2003.

[20] Y. Gong, Y. Jia, T. Leung, A. Toshev, and S. Ioffe, "Deep convolutional ranking for multilabel image annotation," *arXiv preprint arXiv:1312.4894*, 2013.

[21] B. Zhou, V. Jagadeesh, and R. Piramuthu, "Conceptlearner: Discovering visual concepts from weakly labeled image collections," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1492–1500.

[22] L. Wu, R. Jin, and A. K. Jain, "Tag completion for image retrieval," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 35, no. 3, pp. 716–727, 2013.

[23] D. Liu, X.-S. Hua, L. Yang, M. Wang, and H.-J. Zhang, "Tag ranking," in *Proceedings of the 18th international conference on World wide web*. ACM, 2009, pp. 351–360.

[24] A. C. Berg, T. L. Berg, H. Daume, J. Dodge, A. Goyal, X. Han, A. Mensch, M. Mitchell, A. Sood, K. Stratos *et al.*, "Understanding and predicting importance in images," in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2012, pp. 3562–3569.

[25] M. Spain and P. Perona, "Measuring and predicting object importance," *International Journal of Computer Vision*, vol. 91, no. 1, pp. 59–76, 2011.

[26] L. Elazary and L. Itti, "Interesting objects are visually salient," *Journal of vision*, vol. 8, no. 3, p. 3, 2008.

[27] K. Yun, Y. Peng, D. Samaras, G. J. Zelinsky, and T. Berg, "Studying relationships between human gaze, description, and computer vision," in *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*. IEEE, 2013, pp. 739–746.

[28] D. Parikh and K. Grauman, "Relative attributes," in *Computer Vision (ICCV), 2011 IEEE International Conference on*. IEEE, 2011, pp. 503–510.

[29] N. Turakhia and D. Parikh, "Attribute dominance: What pops out?" in *Computer Vision (ICCV), 2013 IEEE International Conference on*. IEEE, 2013, pp. 1225–1232.

[30] M. Li, J. Tang, H. Li, and C. Zhao, "Tag ranking by propagating relevance over tag and image graphs," in *Proceedings of the 4th International Conference on Internet Multimedia Computing and Service*. ACM, 2012, pp. 153–156.

[31] J. Zhuang and S. C. Hoi, "A two-view learning approach for image tag ranking," in *Proceedings of the fourth ACM international conference on Web search and data mining*. ACM, 2011, pp. 625–634.

[32] T. Lan and G. Mori, "A max-margin riffled independence model for image tag ranking," in *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*. IEEE, 2013, pp. 3103–3110.

[33] S. Feng, Z. Feng, and R. Jin, "Learning to rank image tags with limited training examples," *Image Processing, IEEE Transactions on*, vol. 24, no. 4, pp. 1223–1234, 2015.

[34] H. Hotelling, "Relations between two sets of variates," *Biometrika*, pp. 321–377, 1936.

[35] P. L. Lai and C. Fyfe, "Kernel and nonlinear canonical correlation analysis," *International Journal of Neural Systems*, vol. 10, no. 05, pp. 365–377, 2000.

[36] Z. Wu and M. Palmer, "Verbs semantics and lexical selection," in *Proceedings of the 32nd annual meeting on Association for Computational Linguistics*. Association for Computational Linguistics, 1994, pp. 133–138.

[37] D. K. C. D. Manning, "Natural language parsing," in *Advances in Neural Information Processing Systems 15: Proceedings of the 2002 Conference*, vol. 15. MIT Press, 2003, p. 3.

[38] C. D. Manning, M. Surdeanu, J. Bauer, J. Finkel, S. J. Bethard, and D. McClosky, "The Stanford CoreNLP natural language processing toolkit," in *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, 2014, pp. 55–60. [Online]. Available: http://www.aclweb.org/anthology/P/P14/P14-5010

[39] S. Bird, E. Klein, and E. Loper, *Natural Language Processing with Python*. O'Reilly Media, 2009.

[40] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 32, no. 9, pp. 1627–1645, 2010.

[41] I. Tsochantaridis, T. Hofmann, T. Joachims, and Y. Altun, "Support vector machine learning for interdependent and structured output spaces," in *Proceedings of the twenty-first international conference on Machine learning*. ACM, 2004, p. 104.

[42] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems*, 2015, pp. 91–99.

[43] X. Hou and L. Zhang, "Saliency detection: A spectral residual approach," in *Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on*. IEEE, 2007, pp. 1–8.

[44] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva, "Learning deep features for scene recognition using places database," in *Advances in neural information processing systems*, 2014, pp. 487–495.

[45] S. Nowozin and C. H. Lampert, "Structured learning and prediction in computer vision," *Foundations and Trends® in Computer Graphics and Vision*, vol. 6, no. 3–4, pp. 185–365, 2011.

[46] G. H. Bakir, T. Hofmann, B. Schölkopf, A. J. Smola, B. Taskar, and S. V. N. Vishwanathan, *Predicting Structured Data (Neural Information Processing)*. The MIT Press, 2007.

[47] T. Joachims, T. Finley, and C.-N. J. Yu, "Cutting-plane training of structural svms," *Machine Learning*, vol. 77, no. 1, pp. 27–59, 2009.

[48] Y. Gong, L. Wang, M. Hodosh, J. Hockenmaier, and S. Lazebnik, "Improving image-sentence embeddings using large weakly annotated photo collections," in *Computer Vision–ECCV 2014*. Springer, 2014, pp. 529–545.

[49] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[50] C. Rashtchian, P. Young, M. Hodosh, and J. Hockenmaier, "Collecting image annotations using amazon's mechanical turk," in *Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon's Mechanical Turk*. Association for Computational Linguistics, 2010, pp. 139–147.

[51] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer Vision–ECCV 2014*. Springer, 2014, pp. 740–755.

[52] J. Xiao, J. Hays, K. Ehinger, A. Oliva, A. Torralba *et al.*, "Sun database: Large-scale scene recognition from abbey to zoo," in *Computer vision and pattern recognition (CVPR), 2010 IEEE conference on*. IEEE, 2010, pp. 3485–3492.

[53] B. C. Russell, A. Torralba, K. P. Murphy, and W. T. Freeman, "Labelme: a database and web-based tool for image annotation," *International journal of computer vision*, vol. 77, no. 1-3, pp. 157–173, 2008.

[54] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. IEEE, 2009, pp. 248–255.

[55] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, M. Bernstein, and L. Fei-Fei, "Visual genome: Connecting language and vision using crowdsourced dense image annotations," 2016. [Online]. Available: http://arxiv.org/abs/1602.07332

**Shangwen Li** received the B.S. and M.S. degrees, both in electrical engineering, from the Zhejiang University, Hangzhou, China in 2008 and 2011 respectively. Since August 2012, he has been been with Media Communications Lab at University of Southern California (USC), Los Angeles. His research interests include image retrieval, scene understanding and object detection using machine learning techniques.

**Sanjay Purushotham** received his PhD in Electrical Engineering at the University of Southern California (USC) in 2015, where he was a research assistant in the Media Communications Lab (MCL) advised by Prof. C.-C. Jay Kuo. In Sept. 2015, he joined the Department of Computer Science at USC as a Post-doc scholar. His research interests are in machine learning and computer vision.

**Chen Chen** received his EE bachelor degree from the Beijing University of Posts and Telecoms (BUPT) in 2010. He is now working on computer vision researches as a PhD candidate in the Media and Creative Lab of the Viterbi Engineer School in USC. His research focuses on scene image segmentation, classification and retrieval using advanced machine learning technologies.

**Yuzhuo Ren** received her B.S. degree from Hebei University of Technology, China, in 2011 and the M.S. degree from University of Southern California, USA, in 2013, both in electrical engineering. She is currently a Ph.D student in USC Media Communications Lab supervised by Prof. C.-C. Jay Kuo. Her research interest is the field of scene understanding, image segmentation, 3D layout estimation using computer vision and machine learning techniques.

**C.-C. Jay Kuo** received the B.S. degree from National Taiwan University, Taipei, Taiwan, in 1980, and the M.S. and Ph.D. degrees from the Massachusetts Institute of Technology, Cambridge, MA, USA, in 1985 and 1987, respectively, all in electrical engineering. He is currently the Director of the Multimedia Communications Laboratory and Deans Professor of Electrical Engineering at the University of Southern California, Los Angeles, CA, USA. His research interests include digital image/video analysis and modeling, multimedia data compression, communication and networking, computer vision and machine learning. He has co-authored about 230 journal papers, 870 conference papers, and 13 books. He is a fellow of the Institute of Electrical and Electronics Engineers (IEEE), the American Association for the Advancement of Science (AAAS) and the International Society for Optical Engineers (SPIE).