

A VISUAL ANALYSIS ON RECOGNIZABILITY AND DISCRIMINABILITY OF ONOMATOPOEIA WORDS WITH DCNN FEATURES

Wataru Shimoda Keiji Yanai

Department of Informatics, The University of Electro-Communications, Tokyo, Japan
{shimoda-k,yanai}@mm.inf.uec.ac.jp

ABSTRACT

In this paper, we examine the relation between onomatopoeia and images using a large number of Web images. The objective of this paper is to examine if the images corresponding to Japanese onomatopoeia words which express the feeling of visual appearance can be recognized by the state-of-the-art visual recognition methods. In our work, first, we collect the images corresponding to onomatopoeia words using an Web image search engine, and then we filter out noise images to obtain clean dataset with automatic image re-ranking method. Next, we analyze the recognizability of various kinds of onomatopoeia images using improved Fisher vector (IFV) and deep convolutional neural network (DCNN) features. In addition, we collect images corresponding to the pairs of nouns and onomatopoeia words, and we examine if the images associated with the same nouns and the different onomatopoeia words are visually discriminable or not. By the experiments, it has been shown that the DCNN features extracted from the layer 7 of Overfeat’s network pre-trained with the ILSVRC 2013 data have prominent ability to represent onomatopoeia images, and most of the onomatopoeia words have visual characteristics which can be recognized.

Index Terms— onomatopoeia, Web images, DCNN features

1. INTRODUCTION

In general, an “onomatopoeia” is a word that phonetically imitates, resembles or suggests the source of the sound that it describes such as “tic tac” and “quack”. In English language, an onomatopoeia is commonly used only for expressing sounds in everyday life. However, onomatopoeia words in Japanese language are commonly used in the boarder purpose such as expressing feeling of visual appearance or touch of objects or materials. Figure 1 shows a “fuwa-fuwa” object, which means being very softy like very soft cotton. In Japanese language, there are so many onomatopoeia words like “fuwa-fuwa” expressing some kinds of feeling of appearance or touch.

The relation between images and onomatopoeia has not been never explored in the context of multimedia research,



Fig. 1. An example photo of “fuwa-fuwa” object.

although many works related to words and images have been done so far. Then, in this paper, we try to analyze the relation between images and onomatopoeia by using a large number of tagged images on the Web. Especially, we examine if onomatopoeia images can be recognized by the state-of-the-art visual recognition method. As a case study on onomatopoeia images, we focus on onomatopoeia in Japanese language, because Japanese language has much more onomatopoeia words which are used in the more broader context compared to other languages such as English.

In this paper, we collect images corresponding to Japanese onomatopoeia words representing feeling of appearance or touch of objects from the Web, and then analyze the relation between onomatopoeia words and images corresponding to them in terms of recognizability using two kinds of state-of-the-art image representations, Improved Fisher Vector [1] and Deep Convolutional Neural Network Features (DCNN features) [2].

In the experiments on collecting images associated with onomatopoeia words, we found that some onomatopoeia words are strongly related to specific kinds of objects. Then, in addition, we collect images corresponding to pairs of nouns and onomatopoeia words, and we examine if the images associated with the same nouns and the different onomatopoeia words are discriminable or not.

2. RELATED WORKS

In this section, we mention some works on material recognition as related works on onomatopoeia.

Since Japanese onomatopoeia represents feeling of appearance, recognition of onomatopoeia image is more related

to material recognition than generic object recognition. As works on material recognition, the work on Flickr Material Database (FMD) [3] is the most representative. They constructed FMD which consists of ten kinds of material photos, “Fabric”, “Foliage”, “Glass”, “Leather”, “Metal”, “Paper”, “Plastic”, “Stone”, “Water” and “Wood”. Each of these material classes has unique visual characteristics which enables people to estimate which material class a given material photo belongs to. However, it was unexplored what kinds of visual features are effective for it. The situation was different from object recognition where local features and bag-of-features representation were proved to be effective. Liu et al. [3] proposed a method to classify material photos based on topic modeling with various kinds of image features. They achieved 44.6% classification accuracy. Cimpoi et al. [4] proposed to represent material images with state-of-the-art image representations, Improved Fisher Vector [1] and Deep Convolutional Neural Network Features (DCNN features) extracted by DeCAF [5], and achieved 67.1% for 10 class material photo classification of FMD. They also created the larger-scale textured photo database, Describable Textures Dataset (DTD), which consists of 47 classes, texture attributes. Inspired by their work, we also use IFV and DCNN features in this paper.

Both FMD [3] and DTD [4] are constructed by gathering images from the Web and selecting good images by hand. Since DTD is relatively a large-scale dataset, they used crowd-sourcing service, Amazon Mechanical Turk (AMT), to select good images out of the images gathered from the Web. Nowadays, AMT is commonly used to image filtering. However, it costs more than a little expense. In this work, we adopt fully automatic image gathering method to built an onomatopoeia image dataset based on the method on automatic Web image gathering and re-ranking with pseudo-positive training samples [6, 7]. An automatic method is helpful to prevent human’s prejudice from getting into the process of image selection.

3. METHODS

In this paper, first we construct an onomatopoeia image database automatically, and next analyze the relation between onomatopoeia words and the corresponding images in terms of visual recognizability of onomatopoeia words. In addition, we carry out the same image gathering process and analysis for the pairs of nouns and onomatopoeia words.

3.1. Gathering onomatopoeia images

To gather onomatopoeia image, we use Bing Image Search API by providing Japanese onomatopoeia words as query words. Most of the upper-ranked images in the search results can be regarded as the images which correspond to the given onomatopoeia word. However, some images irrelevant

to the given word are expected to be included even in the upper-ranked results. Therefore, we re-rank the results obtained from Bing Image Search API so that only relevant images are ranked in the upper rank. To re-rank images, we use the similar approach as [7, 6] where no human supervision is needed. We regard the upper-ranked images in the search result as pseudo-positive training samples and random images as negative samples, and train SVM with them. Then, we apply the trained SVM to the images in the original search results, and sort images in the descending order of the SVM output values to obtain re-ranked results. In our work, we repeat this re-ranking process twice. The detail of the procedure of image collection is as follows:

- (1) Prepare Japanese onomatopoeia words.
- (2) Gather 1000 images corresponding to each onomatopoeia word using Bing Image Search API.
- (3) Extract an image feature vector from each of the gathered images using Improved Fisher Vector [1] and Deep Convolutional Neural Network Features (DCNN features) [2].
- (4) Regard the top-10 images in the search result as pseudo-positive samples and random images as negative samples, and train a linear SVM with them.
- (5) Apply the trained SVM to the images in the original search results, and sort images in the descending order of the SVM output values.
- (6) Carry out the second re-ranking step. Train a linear SVM with the top-20 images in the re-ranked results as pseudo-positive samples, apply it, and sort images in the descending order of the SVM output values again.
- (7) Finally regard the top-50 images as the images corresponding to the given onomatopoeia word.

3.2. Evaluation of recognizability of onomatopoeia words

After gathering onomatopoeia images, we evaluate to what extent the images corresponding to an onomatopoeia word can be recognized by state-of-the-art object recognition methods.

We mix onomatopoeia images and random noise images and discriminate onomatopoeia images from noise images, and examine if we can separate onomatopoeia images from noise images for these mixed images by visual recognition methods regarding each of the onomatopoeia image sets.

To evaluate it fairly, we adopt 5-fold cross validation. We prepare 50 onomatopoeia images selected in the previous step and 5000 random images. In each fold, we select 40 onomatopoeia images as positive samples and 4000 random images as negative samples, and train a linear SVM. Then, we apply the trained SVM into the mixed image set containing

1010 images and rank all the images in the descending order of the SVM output values, and evaluate the result with average precision. We repeat this for five times changing the training samples. In our work, we regard that the obtained mean average precision over the five fold means the recognizability of the corresponding onomatopoeia word.

The average precision is calculated in the following equation:

$$AP = \frac{1}{m} \sum_{k=1}^m Precision_{true}(k)$$

, where m is the number of positive sample (50), and $Precision_{true}(k)$ means the precision value within the k -th positive samples.

3.3. Analysis on pairs of nouns and onomatopoeia words

In the experiments on gathering images associated with onomatopoeia words, we found that some onomatopoeia words are strongly related to specific kinds of objects. For example, most of the image associated with “fuwa-fuwa”, “gotsu-gotsu” and “toro-toro” were cloud images, mountain images and food images, respectively. Then, in addition, we collect images corresponding to pairs of nouns and onomatopoeia words, and we examine if the images associated with the same nouns and the different onomatopoeia words are discriminable or not.

Firstly, we prepare some noun words (20 in the experiments), and we collect snippet texts including the given nouns via Bing Text Search API. Then, we extract 120 adjective words including onomatopoeia words from the collected snippet texts for each given noun.

Secondly, we list up the compound words consisting of the given nouns and one of 120 adjective words, and collect 100 images for each compound words with Bing Image Search API.

Thirdly, we evaluate recognizability of each compound words, and for each noun we select around 10 compound words which have relatively higher recognizability.

Finally, we carry out multi-class classification over the compound words including the same noun and evaluate classification rate by 5-fold cross validation. We regard the result of this multi-class classification as “visual discriminability” of onomatopoeia words within the specific kinds of objects.

3.4. Image Features

As image representation in both the image collection step and the evaluation step, we use two kinds of state-of-the-art features, Improved Fisher Vector [1] and Deep Convolutional Neural Network Features (DCNN features) [2].



Fig. 2. The structure of the Deep Convolutional Neural Network (DCNN) for the ILSVRC 2013 dataset in Over feat[2]. We extracted feature vectors from the Layer-5, the Layer-6 and the Layer-7.

3.4.1. Improved Fisher Vector(IFV)

To encode an image to IFV, we follow the method proposed by Perronnin et al. [1]. First, we extract SIFT local features randomly from a given image, and apply PCA to reduce their dimension from 128 to 64. Next, we code them into a Fisher Vector with the GMM consisting of 64 Gaussian, and obtain IFV after L2-normalizing the Fisher Vector. Since the dimension of the IFV is $2DK$ where D is the number of dimension of the local features and K is the number of elements of MGM, totally it is $2 \times 64 \times 64 = 8192$.

3.5. Deep Convolutional Neural Network (DCNN)

Recently, it has been proved that Deep Convolutional Neural Network (DCNN) is very effective for large-scale object recognition. However, it needs a lot of training images. In fact, one of the reasons why DCNN won the Image Net Large-Scale Visual Recognition Challenge (ILSVRC) 2013 is that the ILSVRC dataset contains one thousand training images per category [8]. This situation does not fit common visual recognition tasks. Then, to make the best use of DCNN for common image recognition tasks, Donahue et al. [5] proposed the pre-trained DCNN with the ILSVRC 1000-class dataset was used as a feature extractor.

Following Donahue et al. [5], we extract the network signals from the middle layers (layer 5, 6 and 7) in the pre-trained DCNN as a DCNN feature vector. We use the pre-trained deep convolutional neural network in Overfeat [2] as shown in Figure 2. This is slight modification of the network structure proposed by Krizhevsky et al. [8] at the LISVRC 2012 competition. In the experiments, we extract raw signals from layer-5, layer-6 or layer-7, where the dimension of the signals are 36864, 3072 and 4096, respectively, and L2-normalize them to use them as DCNN feature vectors.

3.6. Support Vector Machine(SVM)

For classification in both the image collection step and the evaluation step on recognizability of onomatopoeia images, we use a linear SVM which is commonly used as a classifier for IFV and DCNN, since they are relatively higher dimensional. In the experiments, we used LIBLINEAR [9] as an implementation of SVM.

Table 1. Twenty kinds of Japanese onomatopoeia words used in the experiments.

onomatopoeia	meaning	onomatopoeia	meaning
pika-pika	shining gold	mofu-mofu	softly
basha-basha	splashing water	mock-mock	mountainous smoke or clouds
fuwa-fuwa	softly; spongy	kara-kara	hanging many metals
nyoki-nyoki	shooting up one after another	bou-bou	overgrown
kira-kira	shining stars	fuwa-fuwa	well-roasted
gune-gune	winding	shiwa-shiwa	wrinkled; crumpled
toge-toge	thorny; prickly	zara-zara	sandy; gritty
butsu-butsu	a rash	kari-kari	crispy; crunch
puru-puru	fresh and juicy	guru-guru	whirling
gotsu-gotsu	rugged; angular; hard; stiff	giza-giza	notched; corrugated



Fig. 3. Examples of onomatopoeia images gathered from the Web.

4. EXPERIMENTS

In the experiments, we collected images related to twenty onomatopoeia words and examined their recognizability with Fisher Vector and DCNN features. The twenty Japanese onomatopoeia words we used in the experiments and their meanings are shown in Table 1, the visual recognizability of which we will examine in the experiments. In addition, we examined visual discriminability of onomatopoeia words with images associated with noun/onomatopoeia(adjective) pairs.

4.1. Data Collection

We gathered 1000 images for each of the twenty Japanese onomatopoeia words using Bing Image Search API, and repeated re-ranking twice using four kinds of image features. Finally we obtained an onomatopoeia image dataset containing twenty onomatopoeia categories where each category has fifty images without any human supervision. Figure 3 shows some images corresponding to ten onomatopoeia words.

We evaluated the precision of the onomatopoeia datasets constructed with four different kinds of image representations by subjective evaluation. Figure 4 shows the precision value of the selected fifty images on each of the twenty given onomatopoeia words in case of using IFV, DCNN Layer-5,

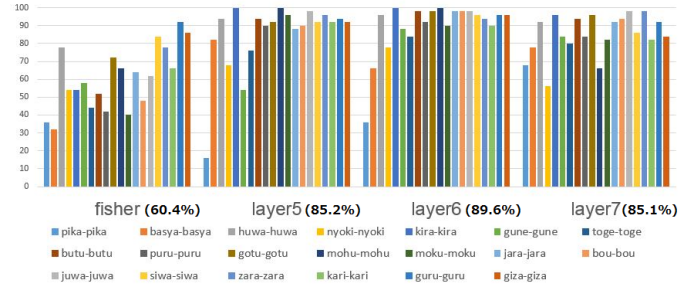


Fig. 4. Precision of the collected images corresponding to the 20 given onomatopoeia words.

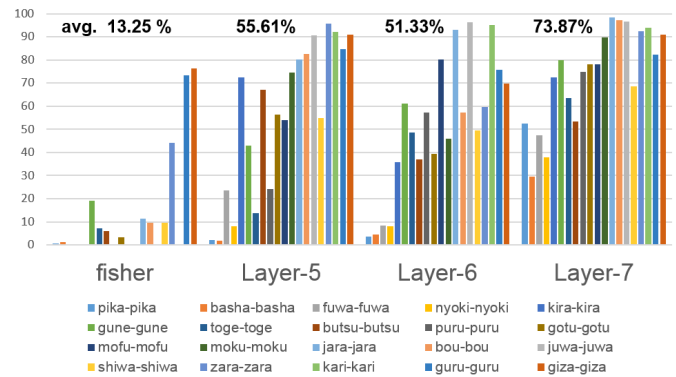


Fig. 5. Evaluation results of the recognizability of each onomatopoeia word.

DCNN Layer-6 and DCNN Layer-7 as a feature, respectively. The mean values of the precision over twenty words for four kinds of features were 60.4%, 85.2%, 89.6% and 85.1%, respectively. As a result, DCNN features outperformed IFV clearly, and DCNN Layer-6 achieved the best performance.

4.2. Evaluation of recognizability

Figure 5 shows the results on recognizability of each of the twenty onomatopoeia words represented by the average precision of the results of separation of 50 onomatopoeia images from 5000 noise images with five-fold cross-validation in case of using IFV, DCNN Layer-5, DCNN Layer-6 and DCNN Layer-7 as a feature, respectively. As a result, DCNN Layer-7 achieved the best result on the average over twenty onomatopoeia words. We guess the reason why DCNN Layer-7 achieved the best is that the activation signals extracted from the layer-7 reflects the semantics of the images more, which helps represent visual characteristics of onomatopoeia.

Figure 6 and 7 shows the top-20 “gotsu-gotsu” (which means being stiff or hard) images in the descending order of the SVM output values. In case of IFV, separation was failed, because the top-20 images contains many images irrelevant



Fig. 6. The top-20 images of “gotsu-gotsu” classified with IFV features.



Fig. 7. The top-20 images of “gotsu-gotsu” classified with DCNN Layer-7 features.

to “gotsu-gotsu”. On the other hand, the results by DCNN Layer-7 does not contain prominent irrelevant images. This result shows that DCNN features has high ability to express visual onomatopoeia elements in images.

4.3. Evaluation of discriminability of onomatopoeia words within the same target objects

In addition to evaluation of recognizability of single onomatopoeia words, we carried out evaluation of discriminability of onomatopoeia words regarding the same target objects. We made multi-class classification in 5-fold cross validation using DCNN Layer-7 features within the same nouns. As the nouns of target objects, we used four nouns: “dog”, “shoes”, “cake” and “flower”. Note that although we collected images associated with pairs of nouns and adjectives including onomatopoeia words for twenty kinds of noun words, the nouns more than six compound words of which are evaluated as having “high recognizability” was only four kinds.

The gathered images and evaluation results of recognizability and discriminability of nouns/adjective(onomatopoeia) compound words on “dog”, “shoes”, “cake” and “flower” are shown in Figure 8, 9, 10 and 11. The adjective words which are paired with the given nouns includes some non-onomatopoeia words. This is because the number of the detected onomatopoeia words for each nouns was not enough for multi-class classification experiments, and we detected some adjectives which have high recognizability as additional “onomatopoeia” words for multi-class classification. The several categories in the upper rows in the each table corresponds to nouns/onomatopoeia pairs, while the rest categories in the lower rows corresponds to nouns/adjective pairs.

As shown in the table, “recognizability” of most of the compound words are high, since we selected the words having high recognizability. Regarding the confusion matrix of the multi-class classification, the values in the lower right corner represents the classification rate of 8-class “dog”, 6-class “shoes”, 7-class “cake” and 7-class “flower”, 52.5%, 85.7%, 72.3% and 84.6%, respectively. All the classification results were much more than the random rate. This results shows

the selected onomatopoeia words have visual characteristics which can be discriminated from other onomatopoeia words and other non-onomatopoeia adjectives even within the same object category. The proposed method will help mine the onomatopoeia words having discriminative visual property corresponding to the specific object categories.

5. CONCLUSIONS

In this paper, we examined if the images corresponding to Japanese onomatopoeia words which express the feeling of visual appearance or touch of objects can be recognized by the state-of-the-art visual recognition methods. In our work, first, we collect the images corresponding to onomatopoeia words using an Web image search engine, and then we filter out noise images to obtain clean dataset with automatic image re-ranking method. Next, we analyze recognizability of various kinds of onomatopoeia images by using improved Fisher vector (IFV) and deep convolutional neural network (DCNN) features. In addition, we collect images corresponding to pairs of nouns and onomatopoeia words, and we examine if the images associated with the same nouns and the different onomatopoeia words are discriminable or not.

By the experiments, it has been shown that the DCNN features extracted from the layer 7 of Overfeat’s network pre-trained with the ILSVRC 2013 data have prominent ability to represent onomatopoeia images, and most of the onomatopoeia words have visual characteristics which can be recognized.

6. REFERENCES

- [1] F. Perronnin, J. Sanchez, and T. Mensink, “Improving the fisher kernel for large-scale image classification,” in *Proc. of European Conference on Computer Vision*, 2010.
- [2] P. Sermanet, D. Eigen, X. Zhang, M.I Mathieu, R. Fergus, and Y. LeCun, “Overfeat: Integrated recognition, localization and detection using convolutional networks,” in *Proc. of International Conference on Learning Representations*, 2014.
- [3] C. Liu, L. Sharan, E. Adelson, and R. Rosenholtz, “Exploring features in a bayesian framework for material recognition,” in *Proc. of IEEE Computer Vision and Pattern Recognition*, 2010.
- [4] M. Cimpoi, S. Maji, I Kokkinos, S. Mohamed, and A. Vedaldi, “Describing textures in the wild,” in *Proc. of IEEE Computer Vision and Pattern Recognition*, 2014.
- [5] J Donahue, Y Jia, O Vinyals, J Hoffman, N Zhang, E Tzeng, and T Darrell, “DeCAF: A deep convolutional activation feature for generic visual recognition,” *arXiv preprint arXiv:1310.1531*, 2013.
- [6] K. Yanai and K. Barnard, “Probabilistic Web image gathering,” in *Proc. of ACM SIGMM International Workshop on Multimedia Information Retrieval*, pp. 57–64, 2005.
- [7] F. Schroff, A. Criminisi, and A. Zisserman, “Harvesting image databases from the web,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 4, pp. 754–766, 2011.
- [8] A. Krizhevsky, I. Sutskever, and G E. Hinton, “Imagenet classification with deep convolutional neural networks..” in *Advances in Neural Information Processing Systems*, 2012.
- [9] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, “LIBLINEAR: A library for large linear classification,” *The Journal of Machine Learning Research*, vol. 9, pp. 1871–1874, 2008.



	recognizability		confusion matrix									
	25	50	25	4	4	2	3	6	2	1	56.0	
"fuwa-fuwa" dog	33.8	43.0	28	4	4	2	3	6	2	1	56.0	
cute dog	75.9	86.5	1	28	2	13	2	3	1	0	56.0	
small dog	51.8	53.2	6	3	19	5	5	4	7	1	38.0	
pretty dog	88.3	88.6	0	13	2	25	0	9	0	1	50.0	
scary dog	20.1	43.5	2	8	6	2	20	1	7	4	40.0	
white dog	61.5	70.9	5	5	1	3	0	29	1	6	58.0	
long dog	27.0	39.0	5	2	3	1	2	4	27	6	54.0	
strong dog	44.4	61.3	0	1	0	2	0	11	2	34	68.0	
			59.6	43.8	51.4	47.2	62.5	43.3	57.4	64.2	52.5	

Fig. 8. Evaluation results of pairs of "dog" and onomatopoeia (adjectives).



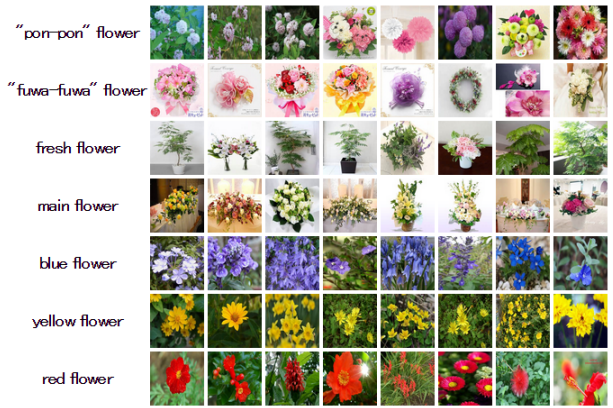
	recognizability		confusion matrix									
	25	50	25	3	4	1	4	3	70.0			
"guru-guru" shoes	80.8	89.3	35	3	4	1	4	3	70.0			
"pika-pika" shoes	100.0	98.6	0	49	0	1	0	0	98.0			
casual shoes	100.0	100.0	0	0	48	1	1	0	96.0			
clean shoes	96.0	99.6	0	1	13	34	1	1	68.0			
white shoes	87.8	93.7	1	0	2	3	44	0	88.0			
red shoes	11.8	53.9	1	2	0	0	0	47	94.0			
			94.6	89.1	71.6	85.0	88.0	92.2	85.7			

Fig. 9. Evaluation results of pairs of "shoes" and onomatopoeia.



	recognizability		confusion matrix									
	25	50	25	3	10	2	3	1	0 <th>62.0</th>	62.0		
"goro-goro" cake	87.5	96.6	31	3	10	2	3	1	0	62.0		
"pasa-pasa" cake	84.4	92.0	3	26	3	7	5	6	0	52.0		
"saku-saku" cake	86.0	96.9	8	4	24	8	2	3	1	48.0		
"fuwa-fuwa" cake	95.0	99.2	1	2	2	42	3	0	0	84.0		
smooth cake	72.5	96.3	3	2	3	3	37	2	0	74.0		
deep cake	90.9	96.8	1	0	2	1	0	46	0	92.0		
light cake	55.6	89.9	0	0	0	0	3	0	47	94.0		
			66.0	70.3	54.5	66.7	69.8	79.3	97.9	72.3		

Fig. 10. Evaluation results of pairs of "cake" and onomatopoeia.



	recognizability		confusion matrix									
	25	50	25	3	2	2	4	0	4 <th>70.0</th>	70.0		
"pon-pon" flower	49.0	94.9	35	3	2	2	4	0	4	70.0		
"fuwa-fuwa" flower	63.0	80.9	7	36	2	3	2	0	0	72.0		
fresh flower	74.0	91.7	2	1	38	7	1	0	1	76.0		
main flower	83.2	97.3	0	2	3	44	1	0	0	88.0		
blue flower	48.1	93.9	1	0	0	0	49	0	0	98.0		
yellow flower	100.0	100.0	0	0	0	0	0	50	0	100.0		
red flower	84.1	93.7	3	0	0	0	3	0	44	88.0		
			72.9	85.7	84.4	78.6	81.7	100.0	89.8	84.6		

Fig. 11. Evaluation results of pairs of "flower" and onomatopoeia.