

DualNet: Domain-Invariant Network for Visual Question Answering

Kuniaki Saito, Andrew Shin, Yoshitaka Ushiku, Tatsuya Harada

The University of Tokyo
7 Chome-3-1 Hongo, Bunkyo
Tokyo 113-8654, Japan

Abstract

Visual question answering (VQA) task not only bridges the gap between images and language, but also requires that specific contents within the image are understood as indicated by linguistic context of the question, in order to generate the accurate answers. Thus, it is critical to build an efficient embedding of images and texts. We implement DualNet, which fully takes advantage of discriminative power of both image and textual features by separately performing two operations. Building an ensemble of DualNet further boosts the performance. Contrary to common belief, our method proved effective in both real images and abstract scenes, in spite of significantly different properties of respective domain. Our method was able to outperform previous state-of-the-art methods in real images category even without explicitly employing attention mechanism, and also outperformed our own state-of-the-art method in abstract scenes category, which recently won the first place in VQA Challenge 2016.

Introduction

Recent rise of deep learning methods including convolutional neural networks (CNN) and recurrent neural networks (RNN) has escalated a large number of artificial intelligence tasks to an unprecedented stage, where the performance frequently rivals that of humans. Tasks such as object classification, scene classification, and object detection demonstrated the ability to correctly recognize and locate the images both holistically and regionally, whereas tasks such as caption generation or object retrieval demonstrated that deep learning methods can successfully bridge the gap between images and language. Visual question answering (VQA) task further promotes the boundary of deep learning applicability and complicates the problem by necessitating multiple prerequisites, potentially encompassing all of the above-mentioned capabilities; as it needs to understand the question, locate or classify the objects/scenes mentioned in the question, and generate appropriate answers.

In this paper, we introduce DualNet, which attempts to fully exploit the discriminative information provided by the images and textual features, by separately performing addition and multiplication of input features to form a common embedding space. As we shall see in Experiment Section, it shows clear advantage over performing only one operation, and outperforms many recent state-of-the-art methods, without using any attention mechanism. Furthermore, it turns

out that building an ensemble of DualNets with varying dimensions leads to even more superior performances, despite feeding identical set of input features to all DualNet units.

Another advantage of our DualNet is that it is applicable to both real images and abstract scenes categories. So far, it has widely been considered that successful methods for real images cannot be directly ported to abstract scenes domain, as they have fundamentally different characteristics. In fact, applying the basic setting of fc7 features for images and long short-term memory (LSTM) with one hidden layer for questions, which results in 58.16 for real images, yields only about 55 in abstract scenes domain. Indeed, most of the previous papers on VQA have tackled only one domain, presumably due to such reason. Our DualNet, however, results in superior performances in both domains, demonstrating that it is applicable to a wider domain, provided that features are plausible.

It is also noteworthy that we do not employ any attention mechanism, which has become one of the most common approaches in VQA. While useful, building attention mechanism necessitates a separate stage of training to map language to specific regions of image, and complicates the procedure. Instead, our method demonstrates that basic set of features can provide rich amount of information without building attention mechanism, provided a network is designed in a way that fully exploits the features' discriminative capacity.

This paper is hereafter structured as follows; In Related Works, we review the recent innovations and trends in VQA, and briefly discuss how our method diverts from them. In Method, we describe both the motivation behind and the implementation details of our DualNet architecture. In Experiment, we apply our proposed model to actual VQA dataset and discuss the results with examples and comparisons to other methods. Finally, we conclude the paper and discuss future work in Conclusion.

Related Work

Visual question answering (VQA) task itself has only recently been introduced with the advent of dataset provided by (Antol et al. 2015), consisting of 0.25M images, 0.76M questions, and 10M answers. They also report baseline results from methods with multi-layer perceptron and LSTM (Hochreiter and Schmidhuber 1997).

VQA: Real Images Real images category is currently by far the more popular and competitive task in VQA. (Malinowski, Rohrbach, and Fritz 2016) introduced Ask Your Neurons. Unlike the baseline provided by (Antol et al. 2015), in which image features and question features are embedded to common space at the last stage prior to classification, they built a system where image features are shared at each LSTM unit for processing question features. They also performed comparison of different operations for fusing input features, and concluded that summation performs better than multiplication. In our work, however, both summation and multiplication are performed, which demonstrates significant improvements.

Many recent papers reporting competitive results have relied heavily on various types of attention mechanism. (Yang et al. 2016) introduced stacked attention networks (SANs), which relies on semantic representation of each question to search for relevant regions in the image. More specifically, they built multiple-layer attention mechanism, which locates the relevant region multiple times so that more accurate region of interest can be retrieved.

In a similar manner, (Shih, Singh, and Hoiem 2015) attempts to locate relevant regions in the image. They map the textual queries to features from different regions by embedding them to a common space and comparing their relevance via inner product.

(Xiong, Merity, and Socher 2016) proposed a number of improvements to dynamic memory network (DMN). Their proposed DMN+ model introduced a novel input module based on a two-level encoder with sentence reader and input fusion layer, and implemented memory based on gated recurrent units (GRU).

(Ilievski, Yan, and Feng 2016) proposed focused dynamic attention (FDA) model, which exploits an object detector to determine regions of interest. LSTM is used to embed the region features and global features into common space.

(Xu and Saenko 2015) proposed spatial memory network in which neuron activations of different spatial regions are stored in memory, and regions with high relevance are chosen depending on the question. The latter step was made possible by their novel spatial attention architecture designed to align words with patches.

Unlike most of the works mentioned above, our work does not employ any attention mechanism, yet demonstrates superior performance by fully exploiting features provided to the network.

VQA: Abstract Scenes Relatively few results have been reported on abstract scene categories compared to real images.

(Zhang et al. 2016) converted the questions to a tuple containing essential clues to the visual concept of the images. Each tuple (P, R, S) consists of a primary object (P), secondary object (S), and their relation (R). Mutual information was employed to determine which object corresponds to primary object and secondary object. They also augmented the dataset using crowd-sourcing in order to balance the biases in the dataset. Their visual features included histogram-like vectors for primary and secondary objects, as well as abso-

lute and relative locations of the objects modeled by GMMs. We show that this model’s performance is enhanced by addition of deep features, both holistically and regionally, and applying our DualNet further improves the performance.

Method

In this section, we describe the details of our proposed network architecture “DualNet”, which we demonstrate to work well both on real images and abstract image. Furthermore, we demonstrate that it performs well on various combinations of image features when combined with sentence features from encoders such as LSTM.

Motivation In the VQA task, it is necessary to determine how to combine visual features with sentence features because a network cannot answer correctly unless they have enough knowledge about what the questions are asking and which features are necessary to answer them correctly. Figure 1 shows an example of fusing features which employs element-wise multiplication (Antol et al. 2015). There are other options to fuse the features such as element-wise summation. Some of the previous works have examined and compared the behaviors of network depending on the fusing mechanisms (Malinowski, Rohrbach, and Fritz 2016). According to them, the performance of network varies depending on the way the image features and sentence features are fused. This indicates that, even with non-linearity of network, the information can vary according to the fusing methods. Most architectures only used one method to fuse the features; for example, summation or multiplication only.

However, the features combined with different fusing methods should contain different information. For example, some information should be preserved (or lost) only by summation, whereas some are preserved only by multiplication. For this reason, we propose to integrate two kinds of operations, namely element-wise summation and element-wise multiplication. Moreover, we propose to use different kinds of image features. The motivation is to fully take advantage of different information present in different kinds of features. For example, holistic features used in abstract scenes (Zhang et al. 2016) display completely different characteristics from CNN features. Likewise, for CNN features, different network structures also result in different characteristics of the extracted features. Thus, our DualNet benefits further by exploiting a combination of features from different networks and different methods.

Implementation We now go through the theoretical background of our network. We skip notations for bias parameters in the following equations for clarity.

$$Q = LSTM(x_1, x_2, x_3, \dots, x_t) \quad (1)$$

First, we input one-hot vector of words sequentially and obtain question vector from the last hidden layer of LSTM.

$$I_{M_1}' = \tanh(W_{M_1} I_1) \quad (2)$$

$$I_{M_2}' = \tanh(W_{M_2} I_2) \quad (3)$$

$$Q_M' = \tanh(W_{M_q} Q) \quad (4)$$

$$F_M = I_{M_1}' \circ I_{M_2}' \circ Q_M' \quad (5)$$

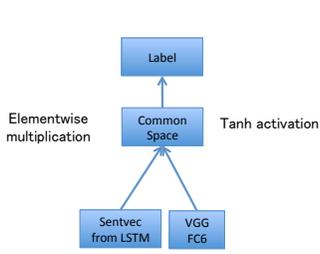


Figure 1: Basic network architecture for VQA

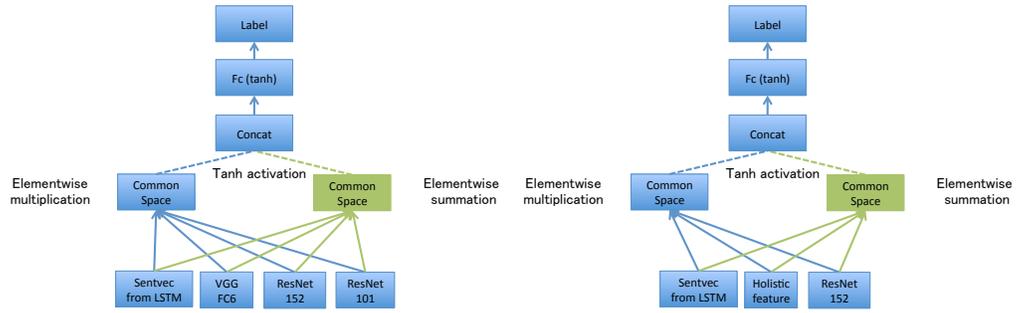


Figure 2: DualNet for real images

Figure 3: DualNet for abstract scenes

Eq. (2) to (5) correspond to the fusing of image features and text features by multiplication. \circ refers to the element-wise multiplication.

$$I_{S_1}' = \tanh(W_{S_1}I_1) \quad (6)$$

$$I_{S_2}' = \tanh(W_{S_2}I_2) \quad (7)$$

$$Q_S' = \tanh(W_{S_q}Q) \quad (8)$$

$$F_S = I_{S_1}' + I_{S_2}' + Q_S' \quad (9)$$

Eq. (6) to (9) correspond to the fusing of the features by summation. Our proposed network does not share the weight between multiplication and summation because we expect each operation to extract different kinds of information.

$$F = \text{Concat}(F_M, F_S) \quad (10)$$

$$\text{Output} = W_{f_2} \tanh(W_{f_1}F) \quad (11)$$

We concatenate the features from element-wise multiplication and element-wise summation. In this example, we have shown the case where we use two kinds of image features. We can change the number of image features depending on the needs, but the overall workflow will remain the same regardless of the number of features.

The proposed model architecture for real image is described in Fig.2. It uses L2-normalized features from the first fully-connected layer (fc6) of VGG-19 (Simonyan and Zisserman 2014) extracted using Caffe (Jia et al. 2014) trained on ImageNet (Deng et al. 2009), and the uppermost fully-connected layers from ResNet-152 and ResNet-101 (He et al. 2015) for image features, in order to construct DualNet. The proposed model architecture for abstract image is described in Fig.3. It uses L2-normalized holistic feature, and fully-connected layer of ResNet-152 for image features.

Experiment

In this section, we describe and discuss the results from the experiments using our DualNet architecture on VQA dataset.

Real Images

The dataset consists of 82,783 training images, 40,504 validation images and 81,434 test images. 3 questions are attached to each image. We can evaluate the model by using a

subset of test split called test-dev split on the VQA evaluation server. In this experiment, we used both train and validation splits for training, and tested on both test-dev and test splits. Since the number of test submissions for the complete test split is limited, the evaluation on the complete test split was restricted to selected key methods.

LSTM in our model consists of 2 layers with 512 hidden units. We used 2,000 most frequent answers as labels, and relied on rms-prop to optimize our model. The batch size was 300 and learning rate was set to 0.0004. We optimized the hyper-parameters based on the evaluations on test-dev split.

The model performance slightly changed according to the dimension of common space. We show the result of 1024 dimension as single DualNet's result. For the ensemble of DualNets, we set the common space dimensions differently for each unit. We changed common space dimension from 500 to 3000 for each DualNet unit in our ensemble, which consists of 19 DualNet units. We tuned the weight for each unit in the ensemble based on their result on test-dev split.

Abstract Scenes

Abstract scenes category contains 20,000 images for training, 10,000 images for validation, and 20,000 images for test images, where each image is accompanied by 3 questions. Unlike real images, there is no test-dev split.

So far, it has widely been believed that abstract scenes possess fundamentally different properties from those of real images, and thus successful methods for real images cannot be directly ported to abstract scenes, necessitating a significantly different approach.

Baseline 1 As our first baseline, we follow Zhang et al. (Zhang et al. 2016) whose method was described in Related Works.

Baseline 2 We now describe the second baseline also implemented by ourselves, which recently won the first place in VQA Challenge 2016, and is currently the state-of-the-art method in the abstract scenes category.

On top of the features described in (Zhang et al. 2016), we added features from the uppermost fully-connected layer from ResNet with 152 layers, and fc7 layer of VGG with 19 layers for holistic features. We alternated between two different setups for regional features as following:

Table 1: Performances of each method on test-dev split of real images category

	Open-Ended				Multiple-Choice			
	All	Y/N	Num	Others	All	Y/N	Num	Others
DPPnet (Noh, Seo, and Han 2015)	57.22	80.7	37.2	41.7	62.50	80.8	38.9	52.2
deeper LSTM Q+norm (Lu et al. 2015)	57.75	80.5	36.8	43.1	62.70	80.5	38.2	53.0
SAN (Yang et al. 2016)	58.70	79.3	36.6	46.1	-	-	-	-
FDA (Ilievski, Yan, and Feng 2016)	59.24	81.1	36.2	45.8	64.01	81.5	39.0	54.7
DMN+(Xiong, Merity, and Socher 2016)	60.30	80.5	36.8	48.3	-	-	-	-
Sum only	56.81	78.4	35.2	43.3	-	-	-	-
Mul only	59.15	80.6	37.0	45.8	-	-	-	-
DualNet	60.47	81.0	37.1	48.2	65.80	80.8	39.8	58.9
DualNet (ensembled)	61.47	82.0	37.9	49.2	66.66	82.1	39.8	59.5

Table 2: Performances of each method on test-std split of real images category

	Open-Ended				Multiple-Choice			
	All	Y/N	Num	Others	All	Y/N	Num	Others
DPPnet (Noh, Seo, and Han 2015)	57.36	80.28	36.92	42.24	62.69	80.35	38.79	52.79
D-NMN (Andreas et al. 2016)	58.0	-	-	-	-	-	-	-
deeper LSTM Q+norm (Lu et al. 2015)	58.16	80.56	36.53	43.73	63.09	80.59	37.70	53.64
AYN (Malinowski, Rohrbach, and Fritz 2016)	58.43	78.24	36.27	46.32	-	-	-	-
SAN (Yang et al. 2016)	58.90	-	-	-	-	-	-	-
FDA (Ilievski, Yan, and Feng 2016)	59.54	81.34	35.67	46.10	64.18	81.25	38.3	55.20
DMN+ (Xiong, Merity, and Socher 2016)	60.36	80.43	36.82	48.33	-	-	-	-
DualNet (ensembled)	61.72	81.92	37.84	49.66	66.72	81.95	39.72	59.55

1) Avg. Softmax of Top Regions: we first extract 10 regions from each image using Deep Proposal (Ghodrati et al. 2015), which proposes regions based on objectness measure and applies non-maximum suppression to filter out overlaps. We then extract softmax probabilities for each region, which correspond to 201 classes used in ILSVRC object detection task. We used Fast-RCNN (Girshick 2015) and VGG-16 trained for the task. Finally, we average the softmax probabilities of all 10 regions to obtain one 201-dimensional vector.

2) VLAD Coding of CNN with Coordinates: The general procedure is similar to (Shin et al. 2016) except we do not employ spatial pyramid. We run selective search for each image, which returns approximately 1,000 region proposals for each image. Using Fast-RCNN, we extract fc7 features from all regions. Dimensionality of fc7 features is reduced to 256 using PCA. We then concatenate 8-dimensional coordinate vector (x_min, y_min, x_max, y_max, x_center, y_center, width, height) as in (Hu et al. 2016) so that each region is 264-dimensional. Finally, we apply VLAD coding to all regions of an image with one cluster to obtain the final one 264-dimensional vector for each image.

It turns out that 1) performs better on yes/no and number questions, while 2) performs better on others category. We thus alternated between the two methods depending on the type of question, which was predicted by key phrase extraction; e.g., ‘how many’ indicating number category, etc. We had batch size of 400, and 500 possible answers, and set number of word embeddings for questions as 1,000. LSTM with one hidden layer of 256 hidden units was employed.

Training was performed for 100 epochs.

Results & Analysis

Table 1 shows the results of each method on test-dev split, and Table 2 reports results on test-std split. As shown in the tables, we outperformed the previous state-of-the-art methods on real image category published prior to the 2016 VQA Challenge. We can clearly see the effectiveness of our network structure through comparison to the summation-only network and multiplication-only network. The multiplication network obtained 59.15 and summation network obtained only 56.81. The performance of summation network is much poorer than multiplication network. However, when combining two paths, we were able to improve the performance significantly. This indicates that the two paths extract different kinds of information from more than three kinds of features, reminiscent of the way ‘‘And’’ and ‘‘Or’’ gates behave in electronic circuits.

Comparing with the methods such as DMN (Xiong, Merity, and Socher 2016), SAN (Yang et al. 2016) and FDA (Ilievski, Yan, and Feng 2016), which used the attention mechanism, our model still achieves better performance. This indicates that we can construct an efficient model without explicitly including spatial information from local features. It also suggests that the image features from VGG and ResNet must contain spatial information to a useful extent, since our model demonstrates high performance on the questions that require the model to have knowledge about particular image regions in order to answer correctly. Figure 4 shows examples of questions and

Table 3: Performances of each method on test data of abstract scenes category

	Open-Ended				Multiple-Choice			
	All	Y/N	Num	Others	All	Y/N	Num	Others
Baseline1 (Zhang et al. 2016)	65.02	77.5	52.5	56.4	69.21	77.5	52.9	66.7
Baseline 2	67.39	79.6	57.1	58.2	71.18	79.6	56.2	67.9
MRN (Kim et al. 2016)	62.56	79.1	51.6	48.9	67.99	79.1	52.6	62.0
DualNet	68.87	80.0	57.9	61.1	73.29	80.0	58.5	71.8
DualNet (ensembled)	69.73	80.7	58.8	62.1	74.02	80.8	59.2	72.4



(a) Q: What fruit is yellow and brown?
A: banana



(b) Q: Is this a laptop? A: yes



(c) Q: How many screens are there? A: 2

Figure 4: Examples of question and generated answers in real images



(a) Q: What is the boy playing with?
A: teddy bear



(b) Q: Are there any animals swimming in the pond? A: No



(c) Q: How many trees? A: 1

Figure 5: Examples of question and generated answers in abstract scenes

generated answers in real images along with the images.

As for abstract image, our DualNet method significantly outperformed the result of Baseline2 which won the first place VQA Challenge 2016. The idea of combining two paths proves to be effective when using different kinds of image features for abstract images as well. Although many works have been published for VQA, few works have tackled the task with abstract image dataset. We suspect that this is because the abstract scenes dataset is too small to construct a large network architecture, which frequently includes attention mechanism. Due to the small number of training samples, training complex network can decrease the performance. On the other hand, in our model, the architecture is so simple that our model is not influenced by the limitation of samples. Figure 5 shows examples of questions and generated answers in abstract scenes along with the images.

Conclusion

We implemented DualNet to efficiently and fully account for discriminative information in images and textual features by performing separate operations for input features and build-

ing ensemble with varying dimensions. Experiment results demonstrate that DualNet outperforms many previous state-of-the-art results and that it is applicable to both real images and abstract scenes despite their fundamentally different characteristics. In particular, we were able to outperform our own previous state-of-the-art results on abstract scenes category, which recently won the first place at VQA Challenge 2016. Since our method was able to perform well even without attention mechanism, it will be an interesting future work to examine the combination of DualNet and attention mechanism.

Acknowledgments

This work was funded by ImpACT Program of Council for Science, Technology and Innovation (Cabinet Office, Government of Japan).

References

[Andreas et al. 2016] Andreas, J.; Rohrbach, M.; Darrell, T.; and Klein, D. 2016. Learning to compose neural networks for question answering. *arXiv preprint arXiv:1601.01705*.

- [Antol et al. 2015] Antol, S.; Agrawal, A.; Lu, J.; Mitchell, M.; Batra, D.; Zitnick, C. L.; and Parikh, D. 2015. Vqa: Visual question answering. *The IEEE International Conference on Computer Vision (ICCV)*.
- [Deng et al. 2009] Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*.
- [Ghodrati et al. 2015] Ghodrati, A.; Diba, A.; Pedersoli, M.; Tuytelaars, T.; and Gool, L. V. 2015. Deepproposal: Hunting objects by cascading deep convolutional layers.
- [Girshick 2015] Girshick, R. 2015. Fast r-cnn.
- [He et al. 2015] He, K.; Zhang, X.; Ren, S.; and Sun, J. 2015. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*.
- [Hochreiter and Schmidhuber 1997] Hochreiter, S., and Schmidhuber, J. 1997. Long short-term memory. *Neural Comput.* 9(8):1735–1780.
- [Hu et al. 2016] Hu, R.; Xu, H.; Rohrbach, M.; Feng, J.; Saenko, K.; and Darrell, T. 2016. Natural language object retrieval.
- [Ilievski, Yan, and Feng 2016] Ilievski, I.; Yan, S.; and Feng, J. 2016. A focused dynamic attention model for visual question answering. *arXiv preprint arXiv:1604.01485*.
- [Jia et al. 2014] Jia, Y.; Shelhamer, E.; Donahue, J.; Karayev, S.; Long, J.; Girshick, R.; Guadarrama, S.; and Darrell, T. 2014. Caffe: Convolutional architecture for fast feature embedding. In *ACMMM*.
- [Kim et al. 2016] Kim, J.-H.; Lee, S.-W.; Kwak, D.-H.; Heo, M.-O.; Kim, J.; Ha, J.-W.; and Zhang, B.-T. 2016. Multimodal residual learning for visual qa. *arXiv preprint arXiv:1606.01455*.
- [Lu et al. 2015] Lu, J.; Lin, X.; Batra, D.; and Parikh, D. 2015. Deeper lstm and normalized cnn visual question answering model. https://github.com/VT-vision-lab/VQA_LSTM_CNN.
- [Malinowski, Rohrbach, and Fritz 2016] Malinowski, M.; Rohrbach, M.; and Fritz, M. 2016. Ask your neurons: A deep learning approach to visual question answering. *arXiv preprint arXiv:1605.02697*.
- [Noh, Seo, and Han 2015] Noh, H.; Seo, P. H.; and Han, B. 2015. Image question answering using convolutional neural network with dynamic parameter prediction. *arXiv preprint arXiv:1511.05756*.
- [Shih, Singh, and Hoiem 2015] Shih, K.; Singh, S.; and Hoiem, D. 2015. Where to look: Focus regions for visual question answering. In *arXiv:1511.07394*.
- [Shin et al. 2016] Shin, A.; Yamaguchi, M.; Ohnishi, K.; and Harada, T. 2016. Dense image representation with spatial pyramid vlad coding of cnn for locally robust captioning.
- [Simonyan and Zisserman 2014] Simonyan, K., and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *CoRR* abs/1409.1556.
- [Xiong, Merity, and Socher 2016] Xiong, C.; Merity, S.; and Socher, R. 2016. Dynamic memory networks for visual and textual question answering. *arXiv preprint arXiv:1603.01417*.
- [Xu and Saenko 2015] Xu, H., and Saenko, K. 2015. Ask, attend, and answer: Exploring question-guided spatial attention for visual question answering. In *arXiv:1511.05234*.
- [Yang et al. 2016] Yang, Z.; He, X.; Gao, J.; Deng, L.; and Smola, A. 2016. Stacked attention networks for image question answering.
- [Zhang et al. 2016] Zhang, P.; Goyal, Y.; Summers-Stay, D.; Batra, D.; and Parikh, D. 2016. Yin and yang: Balancing and answering binary visual questions.