

Image Question Answering using Convolutional Neural Network with Dynamic Parameter Prediction

Hyeonwoo Noh

Paul Hongsuck Seo

Bohyung Han

{shgusdngogo, hsseo, bhhhan}@postech.ac.kr

Department of Computer Science and Engineering, POSTECH, Korea

Abstract

We tackle image question answering (ImageQA) problem by learning a convolutional neural network (CNN) with a dynamic parameter layer whose weights are determined adaptively based on questions. For the adaptive parameter prediction, we employ a separate parameter prediction network, which consists of gated recurrent unit (GRU) taking a question as its input and a fully-connected layer generating a set of candidate weights as its output. However, it is challenging to construct a parameter prediction network for a large number of parameters in the fully-connected dynamic parameter layer of the CNN. We reduce the complexity of this problem by incorporating a hashing technique, where the candidate weights given by the parameter prediction network are selected using a predefined hash function to determine individual weights in the dynamic parameter layer. The proposed network—joint network with the CNN for ImageQA and the parameter prediction network—is trained end-to-end through back-propagation, where its weights are initialized using a pre-trained CNN and GRU. The proposed algorithm illustrates the state-of-the-art performance on all available public ImageQA benchmarks.

1. Introduction

One of the ultimate goals in computer vision is holistic scene understanding [30], which requires a system to capture various kinds of information such as objects, actions, events, scene, atmosphere, and their relations in many different levels of semantics. Although significant progress on various recognition tasks [5, 8, 21, 24, 26, 27, 31] has been made in recent years, these works focus only on solving relatively simple recognition problems in controlled settings, where each dataset consists of concepts with similar level of understanding (e.g. object, scene, bird species, face identity, action, texture etc.). There have been less efforts made on solving various recognition problems simultaneously, which is more complex and realistic, even though this is a crucial step toward holistic scene understanding.



Q: What type of animal is this?
Q: Is this animal alone?



Q: Is it snowing?
Q: Is this picture taken during the day?



Q: What kind of oranges are these?
Q: Is the fruit sliced?



Q: What is leaning on the wall?
Q: How many boards are there?

Figure 1. Sample images and questions in VQA dataset [1]. Each question requires different type and/or level of understanding of the corresponding input image to find correct answers.

Image question answering (ImageQA) [1, 17, 23] aims to solve the holistic scene understanding problem by proposing a task unifying various recognition problems. ImageQA is a task automatically answering the questions about an input image as illustrated in Figure 1. The critical challenge of this problem is that different questions require different types and levels of understanding of an image to find correct answers. For example, to answer the question like “how is the weather?” we need to perform classification on multiple choices related to weather, while we should decide between yes and no for the question like “is this picture taken during the day?” For this reason, not only the performance on a single recognition task but also the capability to select a proper task is important to solve ImageQA problem.

ImageQA problem has a short history in computer vision and machine learning community, but there already exist several approaches [10, 16, 17, 18, 23]. Among these methods, simple deep learning based approaches that perform classification on a combination of features extracted from image and question currently demonstrate the state-of-

the-art accuracy on public benchmarks [23, 16]; these approaches extract image features using a convolutional neural network (CNN), and use CNN or bag-of-words to obtain feature descriptors from question. They can be interpreted as a method that the answer is given by the co-occurrence of a particular combination of features extracted from an image and a question.

Contrary to the existing approaches, we define a different recognition task depending on a question. To realize this idea, we propose a deep CNN with a dynamic parameter layer whose weights are determined adaptively based on questions. We claim that a single deep CNN architecture can take care of various tasks by allowing adaptive weight assignment in the dynamic parameter layer. For the adaptive parameter prediction, we employ a parameter prediction network, which consists of gated recurrent units (GRU) taking a question as its input and a fully-connected layer generating a set of candidate weights for the dynamic parameter layer. The entire network including the CNN for ImageQA and the parameter prediction network is trained end-to-end through back-propagation, where its weights are initialized using pre-trained CNN and GRU. Our main contributions in this work are summarized below:

- We successfully adopt a deep CNN with a dynamic parameter layer for ImageQA, which is a fully-connected layer whose parameters are determined dynamically based on a given question.
- To predict a large number of weights in the dynamic parameter layer, we apply hashing trick [3], which reduces the number of parameters significantly with little impact on network capacity.
- We fine-tune GRU pre-trained on a large-scale text corpus [14] to improve generalization performance of our network. Pre-training GRU on a large corpus is natural way to deal with a small number of training data, but no one has attempted it yet to our knowledge.
- This is the first work to report the results on all currently available benchmark datasets such as DAQUAR, COCO-QA and VQA. Our algorithm achieves the state-of-the-art performance on all the three datasets.

The rest of this paper is organized as follows. We first review related work in Section 2. Section 3 and 4 describe the overview of our algorithm and the architecture of our network, respectively. We discuss the detailed procedure to train the proposed network in Section 5. Experimental results are demonstrated in Section 6.

2. Related Work

There are several recent papers to address ImageQA [1, 10, 16, 17, 18, 23]; the most of them are based on deep

learning except [17]. Malinowski and Fritz [17] propose a Bayesian framework, which exploits recent advances in computer vision and natural language processing. Specifically, it employs semantic image segmentation and symbolic question reasoning to solve ImageQA problem. However, this method depends on a pre-defined set of predicates, which makes it difficult to represent complex models required to understand input images.

Deep learning based approaches demonstrate competitive performances in ImageQA [18, 10, 23, 16, 1]. Most approaches based on deep learning commonly use CNNs to extract features from image while they use different strategies to handle question sentences. Some algorithms employ embedding of joint features based on image and question [1, 10, 18]. However, learning a softmax classifier on the simple joint features—concatenation of CNN-based image features and continuous bag-of-words representation of a question—performs better than LSTM-based embedding on COCO-QA [23] dataset. Another line of research is to utilize CNNs for feature extraction from both image and question and combine the two features [16]; this approach demonstrates impressive performance on DAQUAR [17] dataset by allowing to fine-tune the whole parameters.

The prediction of the weight parameters in deep neural networks has been explored in [2] in the context of zero-shot learning. To perform classification of unseen classes, it trains a multi-layer perceptron to predict a binary classifier for class-specific description in text. However, this method is not directly applicable to ImageQA since finding solutions based on the combination of question and answer is a more complex problem than the one discussed in [2], and ImageQA involves a significantly larger set of candidate answers, which requires much more parameters than the binary classification case. Recently, a parameter reduction technique based on a hashing trick is proposed by Chen *et al.* [3] to fit a large neural network in a limited memory budget. However, applying this technique to the dynamic prediction of parameters in deep neural networks is not attempted yet to our knowledge.

3. Algorithm Overview

We briefly describe the motivation and formulation of our approach in this section.

3.1. Motivation

Although ImageQA requires different types and levels of image understanding, existing approaches [1, 10, 18] pose the problem as a flat classification task. However, we believe that it is difficult to solve ImageQA using a single deep neural network with fixed parameters. In many CNN-based recognition problems, it is well-known to fine-tune a few layers for the adaptation to new tasks. In addition, some networks are designed to solve two or more tasks jointly

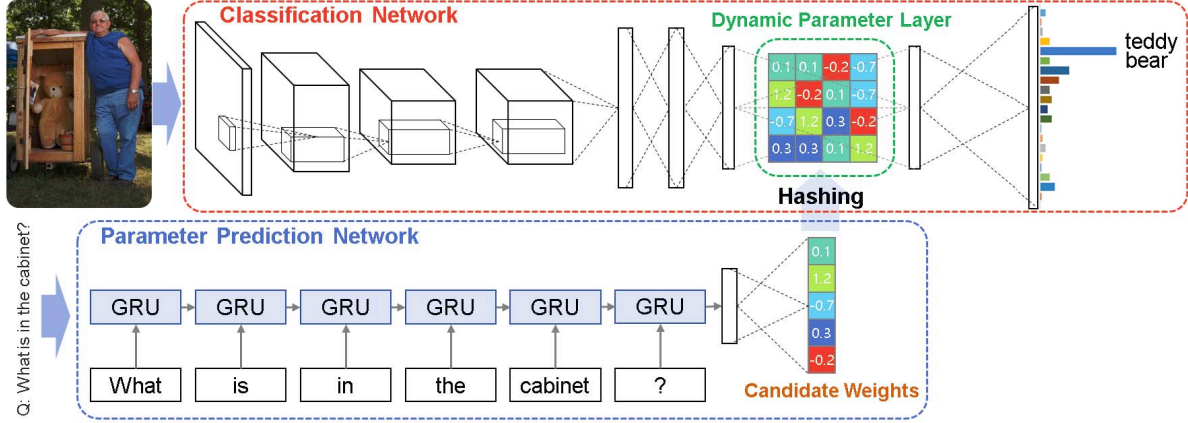


Figure 2. Overall architecture of the proposed Dynamic Parameter Prediction network (DPPnet), which is composed of the classification network and the parameter prediction network. The weights in the dynamic parameter layer are mapped by a hashing trick from the candidate weights obtained from the parameter prediction network.

by constructing multiple branches connected to a common CNN architecture. In this work, we hope to solve the heterogeneous recognition tasks using a single CNN by adapting the weights in the dynamic parameter layer. Since the task is defined by the question in ImageQA, the weights in the layer are determined depending on the question sentence. In addition, a hashing trick is employed to predict a large number of weights in the dynamic parameter layer and avoid parameter explosion.

3.2. Problem Formulation

ImageQA systems predict the best answer \hat{a} given an image I and a question q . Conventional approaches [16, 23] typically construct a joint feature vector based on two inputs I and q and solve a classification problem for ImageQA using the following equation:

$$\hat{a} = \operatorname{argmax}_{a \in \Omega} p(a|I, q; \theta) \quad (1)$$

where Ω is a set of all possible answers and θ is a vector for the parameters in the network. On the contrary, we use the question to predict weights in the classifier and solve the problem. We find the solution by

$$\hat{a} = \operatorname{argmax}_{a \in \Omega} p(a|I; \theta_s, \theta_d(q)) \quad (2)$$

where θ_s and $\theta_d(q)$ denote static and dynamic parameters, respectively. Note that the values of $\theta_d(q)$ are determined by the question q .

4. Network Architecture

Figure 2 illustrates the overall architecture of the proposed algorithm. The network is composed of two sub-networks: classification network and parameter prediction network. The classification network is a CNN. One of the

fully-connected layers in the CNN is the dynamic parameter layer, and the weights in the layer are determined adaptively by the parameter prediction network. The parameter prediction network has GRU cells and a fully-connected layer. It takes a question as its input, and generates a real-valued vector, which corresponds to candidate weights for the dynamic parameter layer in the classification network. Given an image and a question, our algorithm estimates the weights in the dynamic parameter layer through hashing with the candidate weights obtained from the parameter prediction network. Then, it feeds the input image to the classification network to obtain the final answer. More details of the proposed network are discussed in the following subsections.

4.1. Classification Network

The classification network is constructed based on VGG 16-layer net [24], which is pre-trained on ImageNet [6]. We remove the last layer in the network and attach three fully-connected layers. The second last fully-connected layer is the dynamic parameter layer whose weights are determined by the parameter prediction network, and the last fully-connected layer is the classification layer whose output dimensionality is equal to the number of possible answers. The probability for each answer is computed by applying a softmax function to the output vector of the final layer.

We put the dynamic parameter layer in the second last fully-connected layer instead of the classification layer because it involves the smallest number of parameters. As the number of parameters in the classification layer increases in proportion to the number of possible answers, predicting the weights for the classification layer may not be a good option to general ImageQA problems in terms of scalability. Our choice for the dynamic parameter layer can be interpreted as follows. By fixing the classification layer while

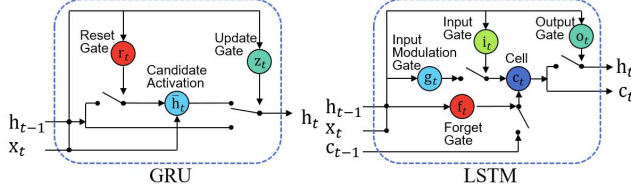


Figure 3. Comparison of GRU and LSTM. Contrary to LSTM that contains memory cell, GRU updates the hidden state directly.

adapting the immediately preceding layer, we obtain the task-independent semantic embedding of all possible answers and use the representation of an input embedded in the answer space to solve an ImageQA problem. Therefore, the relationships of the answers globally learned from all recognition tasks can help solve new ones involving unseen classes, especially in multiple choice questions. For example, when not the exact ground-truth word (*e.g.*, kitten) but similar words (*e.g.*, cat and kitty) are shown at training time, the network can still predict the close answers (*e.g.*, kitten) based on the globally learned answer embedding. Even though we could also exploit the benefit of answer embedding based on the relations among answers to define a loss function, we leave it as our future work.

4.2. Parameter Prediction Network

As mentioned earlier, our classification network has a dynamic parameter layer. That is, for an input vector of the dynamic parameter layer $\mathbf{f}^i = [f_1^i, \dots, f_N^i]^T$, its output vector denoted by $\mathbf{f}^o = [f_1^o, \dots, f_M^o]^T$ is given by

$$\mathbf{f}^o = \mathbf{W}_d(q)\mathbf{f}^i + \mathbf{b} \quad (3)$$

where \mathbf{b} denotes a bias and $\mathbf{W}_d(q) \in \mathbb{R}^{M \times N}$ denotes the matrix constructed dynamically using the parameter prediction network given the input question. In other words, the weight matrix corresponding to the layer is parametrized by a function of the input question q .

The parameter prediction network is composed of GRU cells [4] followed by a fully-connected layer, which produces the candidate weights to be used for the construction of weight matrix in the dynamic parameter layer within the classification network. GRU, which is similar to LSTM, is designed to model dependency in multiple time scales. As illustrated in Figure 3, such dependency is captured by adaptively updating its hidden states with gate units. However, contrary to LSTM, which maintains a separate memory cell explicitly, GRU directly updates its hidden states with a reset gate and an update gate. The detailed procedure of the update is described below.

Let w_1, \dots, w_T be the words in a question q , where T is the number of words in the question. In each time step t , given the embedded vector \mathbf{x}_t for a word w_t , the GRU encoder updates its hidden state at time t , denoted by \mathbf{h}_t ,

using the following equations:

$$\mathbf{r}_t = \sigma(\mathbf{W}_r\mathbf{x}_t + \mathbf{U}_r\mathbf{h}_{t-1}) \quad (4)$$

$$\mathbf{z}_t = \sigma(\mathbf{W}_z\mathbf{x}_t + \mathbf{U}_z\mathbf{h}_{t-1}) \quad (5)$$

$$\bar{\mathbf{h}}_t = \tanh(\mathbf{W}_h\mathbf{x}_t + \mathbf{U}_h(\mathbf{r}_t \odot \mathbf{h}_{t-1})) \quad (6)$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \bar{\mathbf{h}}_t \quad (7)$$

where \mathbf{r}_t and \mathbf{z}_t respectively denote the reset and update gates at time t , and $\bar{\mathbf{h}}_t$ is candidate activation at time t . In addition, \odot indicates element-wise multiplication operator and $\sigma(\cdot)$ is a sigmoid function. Note that the coefficient matrices related to GRU such as \mathbf{W}_r , \mathbf{W}_z , \mathbf{W}_h , \mathbf{U}_r , \mathbf{U}_z , and \mathbf{U}_h are learned by our training algorithm. By applying this encoder to a question sentence through a series of GRU cells, we obtain the final embedding vector $\mathbf{h}_T \in \mathbb{R}^L$ of the question sentence.

Once the question embedding is obtained by GRU, the candidate weight vector, $\mathbf{p} = [p_1, \dots, p_K]^T$, is given by applying a fully-connected layer to the embedded question \mathbf{h}_T as

$$\mathbf{p} = \mathbf{W}_p\mathbf{h}_T \quad (8)$$

where $\mathbf{p} \in \mathbb{R}^K$ is the output of the parameter prediction network, and \mathbf{W}_p is the weight matrix of the fully-connected layer in the parameter prediction network. Note that even though we employ GRU for a parameter prediction network since the pre-trained network for sentence embedding—skip-thought vector model [14]—is based on GRU, any form of neural networks, *e.g.*, fully-connected and convolutional neural network, can be used to construct the parameter prediction network.

4.3. Parameter Hashing

The weights in the dynamic parameter layers are determined based on the learned model in the parameter prediction network given a question. The most straightforward approach to obtain the weights is to generate the whole matrix $\mathbf{W}_d(q)$ using the parameter prediction network. However, the size of the matrix is very large, and the network may be overfitted easily given the limited number of training examples. In addition, since we need quadratically more parameters between GRU and the fully-connected layer in the parameter prediction network to increase the dimensionality of its output, it is not desirable to predict full weight matrix using the network. Therefore, it is preferable to construct $\mathbf{W}_d(q)$ based on a small number of candidate weights using a hashing trick.

We employ the recently proposed random weight sharing technique based on hashing [3] to construct the weights in the dynamic parameter layer. Specifically, a single parameter in the candidate weight vector \mathbf{p} is shared by multiple elements of $\mathbf{W}_d(q)$, which is done by applying a predefined hash function that converts the 2D location in $\mathbf{W}_d(q)$ to the 1D index in \mathbf{p} . By this simple hashing trick, we can reduce

the number of parameters in $\mathbf{W}_d(q)$ while maintaining the accuracy of the network [3].

Let w_{mn}^d be the element at (m, n) in $\mathbf{W}_d(q)$, which corresponds to the weight between m^{th} output and n^{th} input neuron. Denote by $\psi(m, n)$ a hash function mapping a key (m, n) to a natural number in $\{1, \dots, K\}$, where K is the dimensionality of \mathbf{p} . The final hash function is given by

$$w_{mn}^d = p_{\psi(m, n)} \cdot \xi(m, n) \quad (9)$$

where $\xi(m, n) : \mathbb{N} \times \mathbb{N} \rightarrow \{+1, -1\}$ is another hash function independent of $\psi(m, n)$. This function is useful to remove the bias of hashed inner product [3]. In our implementation of the hash function, we adopt an open-source implementation of *xxHash*¹.

We believe that it is reasonable to reduce the number of free parameters based on the hashing technique as there are many redundant parameters in deep neural networks [7] and the network can be parametrized using a smaller set of candidate weights. Instead of training a huge number of parameters without any constraint, it would be advantageous practically to allow multiple elements in the weight matrix to share the same value. It is also demonstrated that the number of free parameter can be reduced substantially with little loss of network performance [3].

5. Training Algorithm

This section discusses the error back-propagation algorithm in the proposed network and introduces the techniques adopted to enhance performance of the network.

5.1. Training by Error Back-Propagation

The proposed network is trained end-to-end to minimize the error between the ground-truths and the estimated answers. The error is back-propagated by chain rule through both the classification network and the parameter prediction network and they are jointly trained by a first-order optimization method.

Let \mathcal{L} denote the loss function. The partial derivatives of \mathcal{L} with respect to the k^{th} element in the input and output of the dynamic parameter layer are given respectively by

$$\delta_k^i \equiv \frac{\partial \mathcal{L}}{\partial f_k^i} \quad \text{and} \quad \delta_k^o \equiv \frac{\partial \mathcal{L}}{\partial f_k^o}. \quad (10)$$

The two derivatives have the following relation:

$$\delta_n^i = \sum_{m=1}^M w_{mn}^d \delta_m^o \quad (11)$$

Likewise, the derivative with respect to the assigned weights in the dynamic parameter layer is given by

$$\frac{\partial \mathcal{L}}{\partial w_{mn}^d} = f_n^i \delta_m^o. \quad (12)$$

¹<https://code.google.com/p/xxhash/>

As a single output value of the parameter prediction network is shared by multiple connections in the dynamic parameter layer, the derivatives with respect to all shared weights need to be accumulated to compute the derivative with respect to an element in the output of the parameter prediction network as follows:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial p_k} &= \sum_{m=1}^M \sum_{n=1}^N \frac{\partial \mathcal{L}}{\partial w_{mn}^d} \frac{\partial w_{mn}^d}{\partial p_k} \\ &= \sum_{m=1}^M \sum_{n=1}^N \frac{\partial \mathcal{L}}{\partial w_{mn}^d} \xi(m, n) \mathbb{I}[\psi(m, n) = k], \end{aligned} \quad (13)$$

where $\mathbb{I}[\cdot]$ denotes the indicator function. The gradients of all the preceding layers in the classification and parameter prediction networks are computed by the standard back-propagation algorithm.

5.2. Using Pre-trained GRU

Although encoders based on recurrent neural networks (RNNs) such as LSTM [11] and GRU [4] demonstrate impressive performance on sentence embedding [19, 25], their benefits in the ImageQA task are marginal in comparison to bag-of-words model [23]. One of the reasons for this fact is the lack of language data in ImageQA dataset. Contrary to the tasks that have large-scale training corpora, even the largest ImageQA dataset contains relatively small amount of language data; for example, [1] contains 750K questions in total. Note that the model in [25] is trained using a corpus with more than 12M sentences.

To deal with the deficiency of linguistic information in ImageQA problem, we transfer the information acquired from a large language corpus by fine-tuning the pre-trained embedding network. We initialize the GRU with the skip-thought vector model trained on a book-collection corpus containing more than 74M sentences [14]. Note that the GRU of the skip-thought vector model is trained in an unsupervised manner by predicting the surrounding sentences from the embedded sentences. As this task requires to understand context, the pre-trained model produces a generic sentence embedding, which is difficult to be trained with a limited number of training examples. By fine-tuning our GRU initialized with a generic sentence embedding model for ImageQA, we obtain the representations for questions that are generalized better.

5.3. Fine-tuning CNN

It is very common to transfer CNNs for new tasks in classification problems, but it is not trivial to fine-tune the CNN in our problem. We observe that the gradients below the dynamic parameter layer in the CNN are noisy since its weights are predicted by the parameter prediction network. Hence, a naïve approach to fine-tune the CNN typically fails

to improve performance, and we employ a slightly different technique for CNN fine-tuning to sidestep the observed problem. We update the parameters of the network using new datasets except the part transferred from VGG 16-layer net at the beginning, and start to update the weights in the subnetwork if the validation accuracy is saturated.

5.4. Training Details

Before training, questions are converted to lower cases and preprocessed by a simple tokenization as in [29]. We also convert answers to lower cases and regard a whole answer in a single or multiple words as a separate class.

The network is trained end-to-end by back-propagation. Adam [13] is used for optimization with initial learning rate 0.01. We clip the gradient to 0.1 to handle the gradient explosion from the recurrent structure of GRU [22]. Training is terminated when there is no progress on validation accuracy for 5 epochs.

Optimizing the dynamic parameter layer is not straightforward since the distribution of the outputs in the dynamic parameter layer is likely to change significantly in each batch. Therefore, we apply batch-normalization [12] to the output activations of the layer to alleviate this problem. In addition, we observe that GRU tends to converge fast and overfit data easily if training continues without any restriction. We stop fine-tuning GRU when the network start to overfit and continue to train the other parts of the network; this strategy improves performance in practice.

6. Experiments

We now describe the details of our implementation and evaluate the proposed method in various aspects.

6.1. Datasets

We evaluate the proposed network on several public ImageQA benchmark datasets such as DAQUAR [17], COCO-QA [23] and VQA [1]. They collected question-answer pairs from existing image datasets and most of the answers are single words or short phrases.

DAQUAR is based on NYUDv2 [20] dataset, and provides two benchmarks. DAQUAR-all consists of 6,795 and 5,673 questions for training and testing respectively, and includes 894 categories in answer. DAQUAR-reduced includes only 37 answer categories for 3,876 training and 297 testing questions. Some questions in this dataset are associated with a set of multiple answers.

The questions in COCO-QA are automatically generated from the image descriptions in MS COCO dataset [15] using the constituency parser with simple question-answer generation rules. The questions in this dataset are typically long and explicitly classified into 4 types depending on the generation rules: object questions, number questions,

color questions and location questions. All answers are with one-words and there are 78,736 questions for training and 38,948 questions for testing.

VQA [1], which is also based on MS COCO dataset [15], contains the largest number of questions: 248,349 for training, 121,512 for validation, and 244,302 for testing, where the testing data is split into test-dev, test-standard, test-challenge and test-reserve. Each question is associated with 10 answers annotated by different people. About 90% of answers have single words and 98% of answers do not exceed three words.

6.2. Evaluation Metrics

DAQUAR and COCO-QA employ both classification accuracy and its relaxed version based on word similarity, WUPS [17]. It uses thresholded Wu-Palmer similarity [28] based on WordNet [9] taxonomy to compute the similarity between words. For predicted answer set \mathcal{A}^i and ground-truth answer set \mathcal{T}^i of the i^{th} example, WUPS is given by

$$\text{WUPS} = \frac{1}{N} \sum_{i=1}^N \min \left\{ \prod_{a \in \mathcal{A}^i} \max_{t \in \mathcal{T}^i} \mu(a, t), \prod_{t \in \mathcal{T}^i} \max_{a \in \mathcal{A}^i} \mu(a, t) \right\}, \quad (14)$$

where $\mu(\cdot, \cdot)$ denotes the thresholded Wu-Palmer similarity between prediction and ground-truth. We use two threshold values (0.9 and 0.0) in our evaluation.

VQA dataset provides open-ended task and multiple-choice task for evaluation. For open-ended task, the answer can be any word or phrase while an answer should be chosen out of 18 candidate answers in the multiple-choice task. In both cases, answers are evaluated by accuracy reflecting human consensus. For predicted answer a_i and target answer set \mathcal{T}^i of the i^{th} example, the accuracy is given by

$$\text{Acc}_{\text{VQA}} = \frac{1}{N} \sum_{i=1}^N \min \left\{ \frac{\sum_{t \in \mathcal{T}^i} \mathbb{I}[a_i = t]}{3}, 1 \right\} \quad (15)$$

where $\mathbb{I}[\cdot]$ denotes an indicator function. In other words, a predicted answer is regarded as a correct one if at least three annotators agree, and the score depends on the number of agreements if the predicted answer is not correct.

6.3. Results

We test three independent datasets, VQA, COCO-QA, and DAQUAR, and first present the results for VQA dataset in Table 1. The proposed Dynamic Parameter Prediction network (DPPnet) outperforms all existing methods non-trivially. We performed controlled experiments to analyze the contribution of individual components in the proposed algorithm—dynamic parameter prediction, use of pre-trained GRU and CNN fine-tuning, and trained 3 additional models, CONCAT, RAND-GRU, and CNN-FIXED.

Table 1. Evaluation results on VQA test-dev in terms of Acc_{VQA}

	Open-Ended				Multiple-Choice			
	All	Y/N	Num	Others	All	Y/N	Num	Others
Question [1]	48.09	75.66	36.70	27.14	53.68	75.71	37.05	38.64
Image [1]	28.13	64.01	00.42	03.77	30.53	69.87	00.45	03.76
Q+I [1]	52.64	75.55	33.67	37.37	58.97	75.59	34.35	50.33
LSTM Q [1]	48.76	78.20	35.68	26.59	54.75	78.22	36.82	38.78
LSTM Q+I [1]	53.74	78.94	35.24	36.42	57.17	78.95	35.80	43.41
CONCAT	54.70	77.09	36.62	39.67	59.92	77.10	37.48	50.31
RAND-GRU	55.46	79.58	36.20	39.23	61.18	79.64	38.07	50.63
CNN-FIXED	56.74	80.48	37.20	40.90	61.95	80.56	38.32	51.40
DPPnet	57.22	80.71	37.24	41.69	62.48	80.79	38.94	52.16

Table 2. Evaluation results on VQA test-standard

	Open-Ended				Multiple-Choice			
	All	Y/N	Num	Others	All	Y/N	Num	Others
Human [1]	83.30	95.77	83.39	72.67	-	-	-	-
LSTM Q+I [1]	54.06	-	-	-	-	-	-	-
DPPnet	57.36	80.28	36.92	42.24	62.69	80.35	38.79	52.79

Table 3. Evaluation results on COCO-QA

	Acc	WUPS 0.9	WUPS 0.0
IMG+BOW [23]	55.92	66.78	88.99
2VIS+BLSTM [23]	55.09	65.34	88.64
Ensemble [23]	57.84	67.90	89.52
ConvQA [16]	54.95	65.36	88.58
DPPnet	61.19	70.84	90.61

CNN-FIXED is useful to see the impact of CNN fine-tuning since it is identical to DPPnet except that the weights in CNN are fixed. RAND-GRU is the model without GRU pre-training, where the weights of GRU and word embedding model are initialized randomly. It does not fine-tune CNN either. CONCAT is the most basic model, which predicts answers using the two fully-connected layers for a combination of CNN and GRU features. Obviously, it does not employ any of new components such as parameter prediction, pre-trained GRU and CNN fine-tuning.

The results of the controlled experiment are also illustrated in Table 1. CONCAT already outperforms LSTM Q+I by integrating GRU instead of LSTM [4] and batch normalization. RAND-GRU achieves better accuracy by employing dynamic parameter prediction additionally. It is interesting that most of the improvement comes from yes/no questions, which may involve various kinds of tasks since it is easy to ask many different aspects in an input image for binary classification. CNN-FIXED improves accuracy further by adding GRU pre-training, and our final model DPPnet achieves the state-of-the-art performance on VQA dataset with large margins as illustrated in Table 1 and 2.

Table 3, 4, and 5 illustrate the results by all algorithms including ours that have reported performance on COCO-QA, DAQUAR-reduced, DAQUAR-all datasets. The proposed algorithm outperforms all existing approaches consistently in all benchmarks. In Table 4 and 5, single answer and multiple answers denote the two subsets of questions divided by the number of ground-truth answers. Also, the numbers

Table 4. Evaluation results on DAQUAR reduced

	Single answer			Multiple answers		
	Acc	0.9	0.0	Acc	0.9	0.0
Multiworld [17]	-	-	-	12.73	18.10	51.47
Askneuron [18]	34.68	40.76	79.54	29.27	36.50	79.47
IMG+BOW [23]	34.17	44.99	81.48	-	-	-
2VIS+BLSTM [23]	35.78	46.83	82.15	-	-	-
Ensemble [23]	36.94	48.15	82.68	-	-	-
ConvQA [16]	39.66	44.86	83.06	38.72	44.19	79.52
DPPnet	44.48	49.56	83.95	44.44	49.06	82.57

Table 5. Evaluation results on DAQUAR all

	Single answer			Multiple answers		
	Acc	0.9	0.0	Acc	0.9	0.0
Human [17]	-	-	-	50.20	50.82	67.27
Multiworld [17]	-	-	-	07.86	11.86	38.79
Askneuron [18]	19.43	25.28	62.00	17.49	23.28	57.76
ConvQA [16]	23.40	29.59	62.95	20.69	25.89	55.48
DPPnet	28.98	34.80	67.81	25.60	31.03	60.77

(0.9 and 0.0) in the second rows are WUPS thresholds.

To understand how the parameter prediction network understand questions, we present several representative questions before and after fine-tuning GRU in a descending order based on their cosine similarities to the query question in Table 6. The retrieved sentences are frequently determined by common subjective or objective words before fine-tuning while they rely more on the tasks to be solved after fine-tuning.

The qualitative results of the proposed algorithm are presented in Figure 4. In general, the proposed network is successful to handle various types of questions that need different levels of semantic understanding. Figure 4(a) shows that the network is able to adapt recognition tasks depending on questions. However, it often fails in the questions asking the number of occurrences since these questions involve the difficult tasks (e.g., object detection) to learn only with image level annotations. On the other hand, the proposed network is effective to find the answers for the same question on different images fairly well as illustrated in Figure 4(b). Refer to our project website² for more detailed results.

7. Conclusion

We proposed a novel architecture for image question answering based on two subnetworks—classification network and parameter prediction network. The classification network has a dynamic parameter layer, which enables the classification network to adaptively determine its weights through the parameter prediction network. While predicting all entries of the weight matrix is infeasible due to its large dimensionality, we relieved this limitation using parameter hashing and weight sharing. The effectiveness of the proposed architecture is supported by experimental results showing the state-of-the-art performances on three different

²<http://cvlab.postech.ac.kr/research/dppnet/>

Table 6. Retrieved sentences before and after fine-tuning GRU

Query question	What body part has most recently contacted the ball?	Is the person feeding the birds?
Before fine-tuning	What shape is the ball? What colors are the ball? What team has the ball? How many times has the girl hit the ball? What number is on the women's Jersey closest to the ball? What is unusual about the ball? What is the speed of the ball?	Is he feeding the birds? Is the reptile fighting the birds? Does the elephant want to play with the birds? What is the fence made of behind the birds? Where are the majority of the birds? What colors are the birds? Is this man feeding the pigeons?
After fine-tuning	What body part is the boy holding the bear by? What body part is on the right side of this picture? What human body part is on the table? What body parts appear to be touching? What partial body parts are in the foreground? What part of the body does the woman on the left have on the ramp? Name a body part that would not be visible if the woman's mouth was closed?	Is he feeding the birds? Is the person feeding the sheep? Is the man feeding the pigeons? Is she feeding the pigeons? Is that the zookeeper feeding the giraffes? Is the reptile fighting the birds? Does the elephant want to play with the birds?



Q: How does the woman feel?
 DPPnet: happy
 Q: What type of hat is she wearing?
 DPPnet: cowboy



Q: Is it raining?
 DPPnet: no
 Q: What is he holding?
 DPPnet: umbrella



Q: What is he doing?
 DPPnet: skateboarding
 Q: Is this person dancing?
 DPPnet: no



Q: How many cranes are in the image?
 DPPnet: 2 (3)
 Q: How many people are on the bench?
 DPPnet: 2 (1)

(a) Result of the proposed algorithm on multiple questions for a single image

Q: What is the boy holding?



DPPnet: surfboard



DPPnet: bat

Q: What animal is shown?



DPPnet: giraffe



DPPnet: elephant

Q: What is this room?



DPPnet: living room

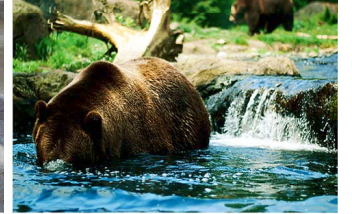


DPPnet: kitchen

Q: What is the animal doing?



DPPnet: resting (relaxing)



DPPnet: swimming (fishing)

(b) Results of the proposed algorithm on a single common question for multiple images

Figure 4. Sample images and questions in VQA dataset [1]. Each question requires a different type and/or level of understanding of the corresponding input image to find correct answer. Answers in blue are correct while answers in red are incorrect. For the incorrect answers, ground-truth answers are provided within the parentheses.

datasets. Note that the proposed method achieved outstanding performance even without more complex recognition processes such as referencing objects. We believe that the proposed algorithm can be extended further by integrating attention model [29] to solve such difficult problems.

Acknowledgements This work was partly supported by IITP grant (B0101-16-0307, Machine Learning Center) and NRF grant (NRF-2011-0031648, Global Frontier R&D Program on Human-Centered Interaction for Coexistence) funded by the Korean government (MSIP).

References

- [1] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh. VQA: visual question answering. In *ICCV*, 2015. 1, 2, 5, 6, 7, 8
- [2] J. Ba, K. Swersky, S. Fidler, and R. Salakhutdinov. Predicting deep zero-shot convolutional neural networks using textual descriptions. In *ICCV*, 2015. 2
- [3] W. Chen, J. T. Wilson, S. Tyree, K. Q. Weinberger, and Y. Chen. Compressing neural networks with the hashing trick. In *ICML*, 2015. 2, 4, 5
- [4] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. In *NIPS Deep Learning Workshop*, 2014. 4, 5, 7
- [5] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild. In *CVPR*, 2014. 1
- [6] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 3
- [7] M. Denil, B. Shakibi, L. Dinh, N. de Freitas, et al. Predicting parameters in deep learning. In *NIPS*, 2013. 5
- [8] 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. In *ICML*, 2014. 1
- [9] C. Fellbaum. Wordnet: An electronic database, 1998. 6
- [10] H. Gao, J. Mao, J. Zhou, Z. Huang, L. Wang, and W. Xu. Are you talking to a machine? dataset and methods for multilingual image question answering. In *NIPS*, 2015. 1, 2
- [11] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. 5
- [12] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 6
- [13] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 6
- [14] R. Kiros, Y. Zhu, R. Salakhutdinov, R. S. Zemel, A. Torralba, R. Urtasun, and S. Fidler. Skip-thought vectors. In *NIPS*, 2015. 2, 4, 5
- [15] 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 *ECCV*, 2014. 6
- [16] L. Ma, Z. Lu, and H. Li. Learning to answer questions from image using convolutional neural network. In *AAAI*, 2016. 1, 2, 3, 7
- [17] M. Malinowski and M. Fritz. A multi-world approach to question answering about real-world scenes based on uncertain input. In *NIPS*, 2014. 1, 2, 6, 7
- [18] M. Malinowski, M. Rohrbach, and M. Fritz. Ask your neurons: A neural-based approach to answering questions about images. In *ICCV*, 2015. 1, 2, 7
- [19] T. Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur. Recurrent neural network based language model. In *INTERSPEECH*, 2010. 5
- [20] P. K. Nathan Silberman, Derek Hoiem and R. Fergus. Indoor segmentation and support inference from rgb-d images. In *ECCV*, 2012. 6
- [21] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *CVPR*, 2014. 1
- [22] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. In *ICML*, 2013. 6
- [23] M. Ren, R. Kiros, and R. S. Zemel. Exploring models and data for image question answering. In *NIPS*, 2015. 1, 2, 3, 5, 6, 7
- [24] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 1, 3
- [25] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *NIPS*, 2014. 5
- [26] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 1
- [27] L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *CVPR*, 2014. 1
- [28] Z. Wu and M. Palmer. Verbs semantics and lexical selection. In *ACL*, 1994. 6
- [29] K. Xu, J. Ba, R. Kiros, A. Courville, R. Salakhutdinov, R. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015. 6, 8
- [30] J. Yao, S. Fidler, and R. Urtasun. Describing the scene as a whole: Joint object detection, scene classification and semantic segmentation. In *CVPR*, 2012. 1
- [31] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning deep features for scene recognition using places database. In *NIPS*, 2014. 1