

Actively Seeking and Learning from Live Data

Damien Teney Anton van den Hengel
Australian Institute for Machine Learning
The University of Adelaide
Adelaide, Australia

{damien.teney, anton.vandenhengel}@adelaide.edu.au

Abstract

One of the key limitations of traditional machine learning methods is their requirement for training data that exemplifies all the information to be learned. This is a particular problem for visual question answering methods, which may be asked questions about virtually anything. The approach we propose is a step toward overcoming this limitation by searching for the information required at test time. The resulting method dynamically utilizes data from an external source, such as a large set of questions/answers or images/captions. Concretely, we learn a set of base weights for a simple VQA model, that are specifically adapted to a given question with the information specifically retrieved for this question. The adaptation process leverages recent advances in gradient-based meta learning and contributions for efficient retrieval and cross-domain adaptation. We surpass the state-of-the-art on the VQA-CP v2 benchmark and demonstrate our approach to be intrinsically more robust to out-of-distribution test data. We demonstrate the use of external non-VQA data using the MS COCO captioning dataset to support the answering process. This approach opens a new avenue for open-domain VQA systems that interface with diverse sources of data.

1. Introduction

One of the ongoing criticisms of modern machine learning methods is that they presume the availability of large volumes of training data [20, 44]. This training data should be representative of the distribution from which the test data will be sampled from, which may be unknowable at training time. These methods usually need constant retraining to accommodate recent data, or to alleviate under-generalizing under a domain shift between the training and test distributions. While there exists a host of approaches to address these limitations (from continuum learning [37, 36] to domain adaptation [9, 30, 42] for example), the information extracted from the training data is typically fixed into the

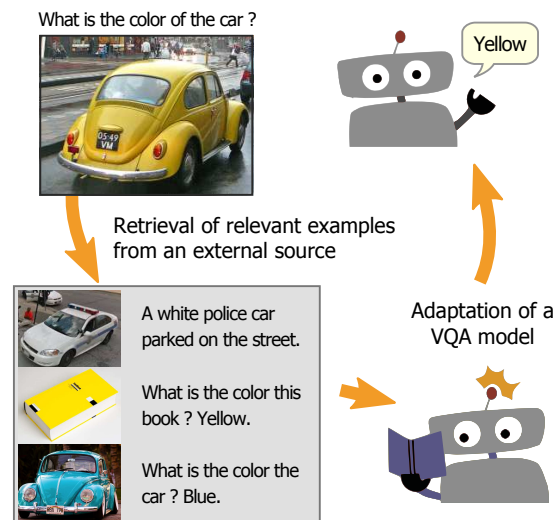


Figure 1. We propose a visual question answering (VQA) system able to retrieve and utilize information from an external source, at test time. The method learns to exploit external information of various forms, and we demonstrate question/answer tuples, but also images and corresponding captions. The method identifies the external information needed to answer a question and adapt its behaviour accordingly. This overcomes limitations of traditional approaches, including overfitting to the training data.

parameters of a model during training, and applied without modification thereafter. The approach we propose here addresses this limitation by exploiting new information as it comes to light, by seeking out relevant data from a large external data source. It actively adapts its behaviour according to the information gained from this data, which represents a fundamental change from pure supervised learning.

This paper demonstrates this novel capability on the task of Visual Question Answering (VQA). The task requires answering a previously unseen question about a previously unseen image. Questions are general and open-ended, and thus require a virtually unlimited array of information and skills to answer. The current approach to VQA is to train

a neural network with end-to-end supervision of questions/answers (QAs). The supervised paradigm has been transformative for most classical tasks of computer vision, but it shows its limits on complex tasks that require more than pixel-processing and pattern recognition alone. VQA models trained in this fashion have revealed to rely mostly on biases and superficial correlations in the training data. For example, questions starting with “How many...” are usually answered with 2 or 3, and those starting with “What sport ...” with the answer tennis, which suffices to obtain high performance on benchmark datasets, where the training and test data are drawn from identical distributions.

The approach proposed in this paper is a step toward *robust* VQA models, *i.e.* capable of reasoning over visual and textual inputs, rather than regurgitating biases learned from a fixed training set. A robust evaluation of these capabilities has recently been made possible. Agrawal *et al.* proposed the VQA-CP (“changing priors”) dataset [1]. In this resampled version of the VQA v2 dataset [15], the training and test sets are drawn from different distributions such that the question type (*i.e.* the first few words such as “What sport ...” or “How many ...”) cannot be relied upon to blindly guess the answer. The performance of existing methods significantly degrades in these conditions.

Our approach borrows ideas from recent research on meta learning [12, 17, 35]. So far, the ubiquitous approach to VQA has attempted to “fit the world” in a neural network, *i.e.* capturing all of the information the method could ever require to answer any question within its weights. In contrast, we train a model to identify and utilize the relevant information from an external source of *support data*. In the simplest instantiation of this principle, the support data is the training set of questions/answers itself, with the major novelty that it does not need to be fixed once the model is trained. The support data can expand at test time and could include data retrieved dynamically from live databases or web searches. The method then adapts itself dynamically using this data. To demonstrate the ability of the model to utilize non-VQA data (*i.e.* other than QA tuples), we use the MS COCO captioning dataset [19, 10] as a source of support data. While VQA data is expensive to acquire, captioned images are omnipresent on the web, and the ability to leverage such data is itself a major contribution.

The evaluation of our approach on VQA-CP demonstrates advantages over classical methods. It generalizes better and obtains state-of-the-art performance on the out-of-distribution test data of VQA-CP. Moreover, the model, once trained on a given distribution of QAs, can successfully adapt to a different distribution of an alternate support set. This is demonstrated with a novel leave-one-out evaluation with VQA-CP. Our experiments clearly demonstrate that the model makes use of the support data at test time, rather than merely capturing biases and priors of a train-

ing set. Consequently, a model trained with our approach could, for example, be reused in another domain-specific application by providing it with a domain-specific support set. This possibility opens the door to systems that reason over vision and language beyond the limited domain covered by any given training set.

The contributions of this paper are summarized as follows.

1. We propose a new approach to VQA in which the model is trained to retrieve and utilize information from an external source, to support its reasoning and answering process. We consider three instantiations of this approach, where the support data is the VQA training set itself (as an evaluation comparable to traditional models), VQA data from a different distribution, and non-VQA image captioning data.
2. We propose an implementation of this approach based on a simple neural network and a gradient-based adaptation, which modifies its weights using selected support data. The method is based on the MAML algorithm [12] with novel contributions for efficient retrieval and cross-domain adaptation.
3. We evaluate the components of our model on the VQA-CP v2 dataset. We demonstrate state-of-the-art performance, benefits in generalization, and the ability to leverage varied sources of support data. The novelty of the approach over existing practices opens the door to multiple opportunities to future research on VQA and vision/language reasoning.

2. Related work

Visual question answering VQA has gathered significant interest in the past few years [5, 39] alongside other tasks combining vision and language such as image captioning [10] or visual dialog [11], for example. The appeal of VQA to the computer vision community is to constitute a practical evaluation of deep visual understanding. Open-domain VQA requires the visual parsing of an image, the comprehension of a text question, and reasoning over multiple pieces of information from these two modalities. See [39] for a survey of modern methods and available datasets.

The ubiquitous approach to VQA is based on supervised learning. It is framed as a classification task over a large set of possible answers, and a machine learning model is optimized over a training set of human-provided questions and answers [5, 15, 18, 48]. Beyond apparent success on VQA benchmarks [14, 33], the approach was revealed to have severe limitations. The models following this approach prove to be overly reliant on superficial statistical regularities in the training sets, and their performance drops dramatically when evaluated on questions drawn from a different distribution [1], or on questions containing words and concepts that appear infrequently in the training data [28, 34]. Popu-

lar benchmarks for VQA [5, 15] have involuntarily encouraged the development of methods that learn and leverage statistical patterns such as biases (*i.e.* the long-tailed distributions of answers) and question-conditioned biases (which make answers easy to guess given a question without the image). These models can essentially bypass the steps of reasoning and image understanding that initially motivated research on VQA.

Robust evaluation of VQA Improved evaluations settings have recently been proposed. In [2, 15, 46], the authors introduced balanced pairs of questions, *i.e.* associating each question with a pair of images that lead to different answers. This procedure, however, had limited benefits. The usual metric of accuracy over individual questions still encouraged to learn and rely on the non-uniform distribution of answers, and the crowd-sourcing procedure used to gather balanced pairs introduced many irrelevant and nonsensical questions to the dataset.

Other recent proposals follow the idea of drawing the training and evaluation questions from different distributions. This discourages overfitting to statistical regularities specific to the training set. In [28, 34], the authors evaluate questions containing words and concepts that appear rarely in the training data. In [1], Agrawal *et al.* propose the VQA-CP dataset (for “changing priors”), in which they enforce different training/test distributions of answers conditioned on the first few words of the question (*e.g.* “What is the color ...” or “How many ...”). Our experiments are conducted on VQA-CP as it represents the most challenging setting currently available.

Robust models for VQA The above robust evaluations have essentially pointed at the inadequacy of current approaches [1, 15, 28, 34, 47]. To address some of these shortcomings, Agrawal *et al.* [1] proposed a modular architecture that prevent it from relying on undesirable biases and priors in the training data. Ramakrishnan *et al.* [27], introduced an information-theoretic regularizer to encourage the model to utilize the image by outperforming a “blind” guesser. In [35], Teney *et al.* proposed a meta learning approach to VQA that improved the recall on rare answers. Their work is the most relevant to this paper, although the methods differ significantly. We use a gradient-based adaptation procedure that update the weights of a whole VQA model, whereas [35] applied existing meta learning algorithms on the final classifier of a simple VQA model. We also formulate the use of support data as a retrieval task, whereas [35] processes the entire support set at every iteration, which is computationally challenging and the evaluation only include small-scale experiments. [35] is also limited to QAs as support data, where our method is much more general.

Meta learning Our central idea is to adapt a VQA model to each given question to incorporate additional informa-

tion from an external source. The adaptation is implemented with the MAML meta learning algorithm [12]. Meta learning or “learning to learn” [21, 31, 37] is a general paradigm to learn to build and/or update machine learning models, *e.g.* to fine-tune the weights of a neural network [7, 6, 32]. Recent works in the area have focused on the adaptation of neural networks for few-shot image recognition [4, 12, 16, 29]. MAML serves to identify a set of weights that can best serve as initial values, before adaptation through one or a few steps of gradient descent. In [13, 43], the authors extended MAML to handle support data from a distinct domain, for robotic imitation learning from demonstration videos. We follow a similar idea to transform the gradients of a loss on captioning data into gradients suitable to update a VQA model. In [17], Huang *et al.* turn the supervised task of language-to-query generation into a meta learning task. They introduce the concept of relevance functions to sample the training set. The approach is similar in spirit to our reformulation of VQA as a meta learning task. However, their aim is to improve accuracy by using specialized adapted models, while our objective is broader, as we also aim to leverage additional (non-VQA) sources of data.

Additional sources of data for VQA The limitations of the mainstream approach to VQA stem from the limited capacity of the training set and of the trained models. Instead of attempting to capture all the training information within the weights of a network, we use an external source of data that is not fixed after training. The capacity and capabilities of the model are thus essentially unbounded. Previous works [40, 38] have interfaced VQA models with knowledge bases, using ad hoc techniques to incorporate external knowledge. In comparison, this paper presents a more general approach, applicable to various types of support data. In [34, 33], the authors used web image search to retrieve visual representations of question and answer words. These representations are however optimized along the other weights of the network and fixed once trained. Recent works on text-based question answering used reinforcement learning to optimize the retrieval of external information [8, 22, 25], which is potentially complementary to our approach.

3. Proposed approach

Our central idea is to learn a VQA model that can subsequently adapt to each particular given question, using additional support data relevant to the question. Intuitively, the adaptation makes the VQA model specialized to the narrow domain of each question. The support data relevant to each question is retrieved dynamically from an external source (Fig. 1), which is assumed to be non-differentiable and/or too large to be processed all at once. Concretely, the support data can be the VQA training set itself (making evaluation

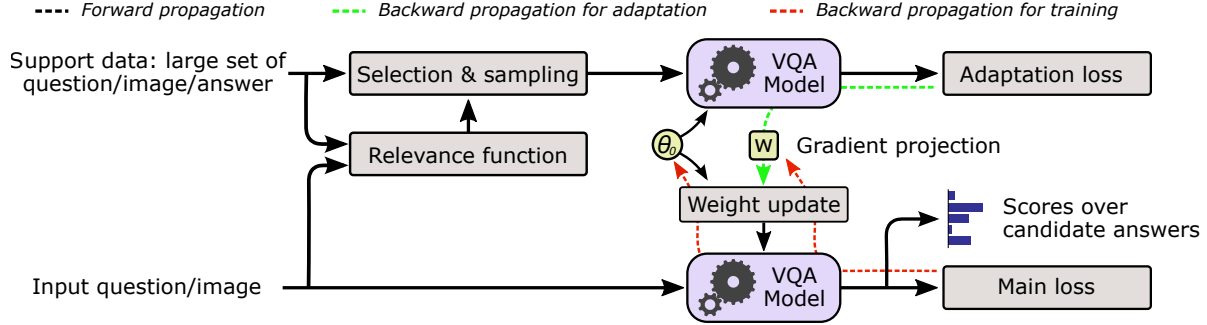


Figure 2. Data flow in the proposed method, using questions/answers as support data (Section 3.2). The input question serves to retrieve pertinent instances from the support data using a relevance function. These instances are passed through the underlying VQA model (Fig. 3) to compute the adaptation loss \mathcal{L}_A , using their ground truth answers. The gradient of the adaptation loss is backpropagated to the weights θ_0 of the VQA model, which are updated, effectively adapting (*i.e.* fine-tuning) the VQA model to the selected support examples. The input question is finally passed through the adapted model to predict scores for the final answer. During training, the gradient of the loss \mathcal{L}_M on the final predictions is backpropagated to optimize the pre-adaptation weights θ_0 and the gradient projection (in yellow).

comparable with traditional methods) but we also demonstrate the use of training QAs from a different distribution (Tables 3–4), and the use of an image captioning dataset (Section 4.1).

3.1. Underlying VQA model

Our approach builds around a standard VQA model that underlies the adaptation procedure. Formally, we denote with $\mathbf{x} = \{\mathbf{q}, \mathbf{v}\}$ the input to the VQA model, made of the question \mathbf{q} (a string of tokens, each corresponding to a word) and of visual features \mathbf{v} pre-extracted from the given image (a feature map produced by a pre-trained convolutional neural network). The VQA model is represented as the function f_θ of parameters θ . It maps \mathbf{x} to a vector of scores with $f_\theta(\mathbf{x}) = \hat{\mathbf{s}}$. The vector $\hat{\mathbf{s}} \in [0, 1]^A$ contains the scores predicted over A candidate answers, typically the few thousands most frequent in the training set. The final answer is the one of largest score, $\arg \max \hat{\mathbf{s}}$. We denote with \mathbf{s} the vector of ground truth scores (which may contain multiple non-zero values when multiple answers are annotated as correct).

The function f is implemented as a neural network and θ denotes the set of all of its weights. Our contributions are not specific to any specific implementation of f . In practice, it corresponds to a classical joint embedding model [33] illustrated in Fig. 3. The network encodes the question as a bag-of-words, taking the average of learned word embeddings. It uses a single-headed, question-guided attention over image locations, a Hadamard product to combine the two modalities, and a non-linear projection followed by a sigmoid to obtain the scores $\hat{\mathbf{s}}$. See Appendix A for details.

3.2. Gradient-based adaptation

The role of the adaptation procedure is to modify the weights of the VQA model to best tailor its capabilities to

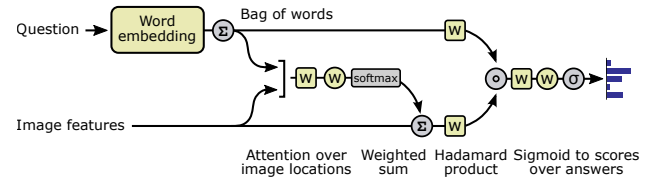


Figure 3. The simple VQA model underlying our method. It implements a classical joint embedding approach [33]. Yellow elements contain learnable weights. Circled and squared ‘w’s represent affine and non-linear projections, respectively. The above network is instantiated twice in the overall diagram of Fig. 2.

a given input question. The motivation for a specialized model is to be potentially be more effective than a general one for a same capacity of the underlying model. Our adaptation procedure is based on MAML [12]. The original MAML algorithm is designed for adaptation using support data of the same form as for task of interest, *i.e.* questions with their ground truth answers. In Section 3.3, we describe an extension to use support data of another task/domain.

The adaptation procedure takes in a set of support elements $\mathcal{S} = \{\mathbf{x}_j\}_j$ and base parameters θ_0 , which it adapts to θ_T over a small number T of updates. The update rule is a gradient descent of step size α :

$$\theta_{i+1} = \theta_i - \alpha \sum_j \nabla_{\theta_i} \mathcal{L}_A(\mathbf{s}_j, f_{\theta_i}(\mathbf{x}_j)) \quad (1)$$

where \mathcal{L}_A is the adaptation loss which evaluates the predictions of the VQA model on the support data. In this case, \mathcal{L}_A is the binary cross-entropy loss typical used to train VQA models [33]. The above adaptation is performed when evaluating a given question at both training and test time. The key to benefit from this approach is to learn base parameters θ_0 that are the most generally and most easily

Algorithm 1: Evaluation of a tr. or test instance.

Input: Test or training instance $\mathbf{x} = (\mathbf{q}, \mathbf{v})$
Support set $\mathcal{S} = \{\mathbf{x}_j\}_j$
with $\mathbf{x}_j = (\mathbf{q}_j, \mathbf{v}_j) \forall j$
Output: Vector of scores \mathbf{s} over candidate answers
// Retrieve support relevant to \mathbf{x} :
 $\mathcal{S}_x \leftarrow \{\mathbf{x}'_j\}_j^K \subset \mathcal{S}$ with max. precomputed $r(\mathbf{x}, \mathbf{x}'_j)$
for $i=0$ to $(T-1)$ **do** *// For each adaptation step*
 $\mathcal{S}'_x \leftarrow K'$ random elements $\in \mathcal{S}_x$
 $\hat{\mathbf{s}}'_j \leftarrow f_\theta(\mathbf{x}'_j) \forall \mathbf{x}'_j \in \mathcal{S}'_x$ *// Forward prop.*
 $\mathbf{d} \leftarrow \sum_j \nabla_{\theta} \mathcal{L}_A(\hat{\mathbf{s}}'_j)$ *// Backprop. adaptation loss*
 $\mathbf{d}' \leftarrow g(\mathbf{d})$ *// Gradient projection*
 $\theta_{i+1} \leftarrow \theta_i - \mathbf{d}'$ *// Update weights of VQA model*
end
 $\hat{\mathbf{s}} \leftarrow f_{\theta_T}(\mathbf{x})$ *// Forward prop. with updated weights*
if training **then**
 $\mathbf{d} \leftarrow \nabla_{\theta_0} \mathcal{L}_M(\hat{\mathbf{s}})$ *// Backprop. main loss*
 $\theta_0 \leftarrow \theta_0 - \alpha \mathbf{d}$ *// Update base weights*
end

adaptable. They are optimized for the following objective:

$$\min_{\theta_0} \sum_{\mathbf{x}_k \sim \mathcal{T}} \mathcal{L}_M(\mathbf{s}_k, f_{\theta_N}(\mathbf{x}_k)) \quad (2)$$

where the elements $(\mathbf{x}_k, \mathbf{s}_k)$ are drawn from a training set \mathcal{T} , and \mathcal{L}_M is the main loss on the VQA model (also called “meta loss” [12]) that corresponds again to a binary cross-entropy. The objective can be optimized with standard backpropagation and stochastic gradient descent [12]. To avoid the expensive differentiation through the T steps of adaptation (Eq. 1), we use a first-order approximation of the gradient as in [23]. The update rule is then

$$\theta_0 \leftarrow \theta_0 - \alpha' \nabla_{\theta_{T-1}} \mathcal{L}_M(\mathbf{s}_k, f_{\theta_T}(\mathbf{x}_k)). \quad (3)$$

where α' is the learning rate. The whole procedure to evaluate any training or test instance is summarized as Algorithm 1. It is worth emphasizing that during training, a support set must be simulated to best mimic the conditions in which the model will be evaluated. If the support set is held constant during training, it would be treated as a static input, and the model is unlikely to generalize to different support data at test time. Therefore, it is crucial to present randomly sampled instances from the support set across the iterations in Algorithm 1.

3.3. Using non-VQA data as support

We now extend the method to use support data other than VQA instances (questions/answers). We apply it to the particular case of images/captions, although the approach is more generally applicable. The challenge is now to produce beneficial updates to the weights θ without access to

a loss on the target VQA model. In practice, the format of captioning data (images with text) facilitates the implementation, as we can use a similar neural network as the VQA model f_θ to process them. We define a model f'_θ similar to f_θ up to the Hadamard product (Fig. 3). The final projection to answers scores is now meaningless for captions.

The adaptation procedure now proceeds as follows. The captions are passed through f' and its output \mathbf{h} (the Hadamard product) is passed to the alternative adaptation loss $\mathcal{L}_{A'} = \|\mathbf{h}\|_2^2$. This squared L2 norm can be interpreted as measuring the compatibility of the embeddings of the caption and of the image. It encourages embedding spaces to align across support images and their captions. Importantly, this loss does not involve ground truth labels or answers, but it allows differentiation with respect to the weights θ ¹. The resulting gradients, however, cannot be assumed to be directly suitable to update the VQA model. We therefore pass them through a learned projection as $g(\nabla_{\theta} \mathcal{L}_{A'})$. This produces gradients that can be plugged into Eq. 1 that now becomes

$$\theta_{i+1} = \theta_i - \alpha \sum_j g(\nabla_{\theta_i} \mathcal{L}_{A'}(g_{\theta_i}(\mathbf{x}_j))) . \quad (4)$$

The projection $g(\cdot)$ is implemented as a non-linear layer that is learned similarly to θ_0 , *i.e.* by backpropagating the gradient of the main loss \mathcal{L}_M as in Eq. 3 (see details in the supplementary material).

3.4. Retrieval of relevant support data

The above descriptions assumed the availability of a set \mathcal{S}_x of support examples relevant to an input question \mathbf{x} . In our experiments, the support data \mathcal{S} is the training split of a large VQA or captioning dataset. The selection of a relevant subset from \mathcal{S} is a crucial step to make the model adaptation both efficient (by processing a much smaller subset $\mathcal{S}_b(x)$) and effective (by focusing the adapted model on a narrow domain around \mathbf{x}). The method described below provides the adaptation algorithm with a subset of the support data of bounded size, and ensures its constant time complexity.

We formalize the retrieval process from \mathcal{S} with a relevance function $r(\mathbf{x}, \mathbf{x}')$. It produces a scalar that reflects the pertinence of a support instance $\mathbf{x}' = (\mathbf{q}', \mathbf{v}') \in \mathcal{S}$ to the input $\mathbf{x} = (\mathbf{q}, \mathbf{v})$. The top- K elements $\{\mathbf{x}'_j\}_j^K \subset \mathcal{S}$ of largest values $r(\mathbf{x}, \mathbf{x}')$ are identified, and then randomly subsampled to the set of K' elements $\mathcal{S}_x = \{\mathbf{x}'_j\}_j^{K'}$.

The relevance function can in principle be learned using the gradient of the main loss $\nabla \mathcal{L}_M$, although we did not explore this option. In our current implementation, we use a static relevance function that allows us to precompute its value between all training elements $\in \mathcal{T}$ and all elements of the simulated support set \mathcal{S}^{tr} . This vastly improves the com-

¹Weights in θ corresponding to the final layers of f_θ and not present in f'_θ receive zero gradients when differentiating through f'_θ .

computational requirements during the training process. Our experiments evaluate conjunctions (products) of the following options:

$$\begin{aligned}
 r_0(\mathbf{x}, \mathbf{x}') &= 1 \quad (\text{Uniform sampling}) \\
 r_1(\mathbf{x}, \mathbf{x}') &= \text{number of common words between } \mathbf{q} \text{ and } \mathbf{q}' \\
 r_2(\mathbf{x}, \mathbf{x}') &= 1 \text{ iff } \mathbf{q}' \text{ contains word matching one of top-5} \\
 &\quad \text{answers from baseline VQA model on } \mathbf{x}. \\
 &= 0 \text{ otherwise} \\
 r_3(\mathbf{x}, \mathbf{x}') &= (\Sigma \mathbf{v} / \|\Sigma \mathbf{v}\|^2) \cdot (\Sigma \mathbf{v}' / \|\Sigma \mathbf{v}'\|^2) \quad (\text{Similarity} \\
 &\quad \text{of globally-pooled, L2-normalized image features}). \tag{5}
 \end{aligned}$$

Note that the retrieval process could alternatively be formulated as a reinforcement learning task. This would allow optimizing the retrieval from “black box” data sources, such as web searches and dynamically-expanding databases [8, 22, 25], which we leave for future work.

4. Experiments

We conducted extensive experiments to evaluate the contribution of the components of our method, and to compare its performance to existing approaches. We use the VQA-CP v2 dataset [1], which is the most challenging benchmark available. Its training and test splits have different distributions of answers conditioned on the first few words of the question, and was built by resampling the VQA v2 dataset [15]. We hold out 8,000 questions from the VQA-CP training data to use as a validation set. All models are trained to convergence (with early stopping) on this validation set. Our underlying VQA model is a reimplementation of [33] (see supplementary material for details). Experiments using captions as support data use the COCO captioning dataset [19]. Since VQA-CP is itself made of images from COCO, we ensure that the captioned images also present in the VQA-CP test set are never used as support (neither during training nor evaluation). Please consult the supplementary material for additional implementation details and results. All results are reported using the standard VQA accuracy metric and broken down into the categories ‘yes/no’, ‘number’, and ‘other’ as in [15].

4.1. Results

Contribution of the proposed components We first evaluate the impact of the proposed components with an ablative study (Table 1). For readability and computational reasons we focus on ‘other’-type questions² with a slightly

²We focus on ‘other’-type questions because random guessing on the yes/no/number questions (or a buggy implementation !) does better than the best model in [1]. We measured that random guessing achieves 72.9% on yes/no questions ([1] gets 65.5%) and random guessing of one/two achieves 34.1% on ‘number’ questions ([1] gets 15.5%). This makes them unreliable for a meaningful analysis.

simplified VQA model. Implementation details are provided in the supplementary material. We examine in Table 1 a series of progressively more elaborate models. Each row corresponds to two different trained models, one trained for QAs as support (evaluated in the first 3 columns), another for captions (evaluated in the last column). All models using adaptation significantly outperform the baseline (first row). Interestingly, the optimal relevance function vary across the models for QAs and captions. The relevance function that includes the image similarity is only moderately useful, while the number of words in common between the question and the support text (QA or caption) proves very effective. Interestingly, in the case of captions, a uniform sampling already gives a clear improvement over the baseline model, but not with QAs, which we explain by the smaller size of the support set of captions.

We report results on both our validation set (of similar distribution as training data) and on the official test set (of different distribution). The overall lower performance on the latter shows the challenge of dealing with out-of-distribution data. The improvement in performance is much clearer on the test set than on the validation set. This demonstrates our contribution to improving generalization – arguably the most challenging aspect of VQA – which is a significant side-effect of our adaptation-based approach.

Using image captions as support data We trained separate models for adaptation to questions/answers and to captions (Table 1 last column). While performance improves over the baseline in both cases, the adaptation using QAs provides a bigger boost, given their direct relevance to the VQA task. The improvement by adaptation to captions demonstrates the ability of the method for picking up relevant information from non-VQA data, which opens a significant avenue for future work. This evaluation currently considers either QAs or captions separately. The combination of the two implies a number of non-trivial design decisions that we will explore in future work.

Amount of retrieved support data In Fig. 4, we examine the performance of the model as a function of the amount of data it is trained with. To make the analysis comparable to the baseline VQA model, the support QAs are the same set of QAs as used for the training (of the baseline and of our model). In the case of captions, we use the same QAs for training, and a similarly subsampled set of captions as support data. We observe that our model is clearly superior to the baseline in all regimes, using both QAs or captions. The gain in performance is maintained even when the model is trained with very little data, in particular when using adaptation with QAs (using as little as 1% of the whole training set).

Unfortunately, the gains in using captions as support data levels off as the amount of support data increases (Fig. 4) and the performance does not surpass that obtained with

	Accuracy on VQA-CP v2 "Other"		
	Val.	Test	
Ours without adaptation	45.46	31.09	
Ours with adaptation and, as support data:	QAs	QAs	Capt.
	Tr.	Tr.	COCO
Uniform sampling $r=r_0$	46.15	31.33	34.00
Relevance function $r=r_1$	44.41	31.79	29.18
Relevance function $r=r_2$	46.49	31.76	33.73
Relevance function $r=r_3$	46.32	31.68	33.51
Relevance function $r=r_2r_3$	46.17	31.09	34.26
Relevance function $r=r_1r_2r_3$	46.79	34.25	33.44

Table 1. Ablative evaluation of the proposed method (see discussion in Section 4.1). Each row corresponds to two different models, trained respectively for QAs (columns 1–3) and for captions (column 4) as support data. Gray cells use additional data during evaluation (QAs from VQA-CP test set in a leave-one-out setting) or during training+evaluation (COCO captions).

QAs. One would rather hope continuing improvement as the model is provided with increasing amounts of support data. We believe that our current results do not prevent this prospect, and that the saturation stems from the particular distribution of captions in COCO. These captions are purely visual and descriptive, and they only cover a limited variety of concepts. In contrast, visual questions often require common sense and knowledge beyond visual descriptions (e.g. Why is the guy wearing such a weird outfit? Is this a healthy breakfast?). Other sources of data, including free-form captions and paired image-text data from the web may be more suitable for this purpose.

Comparison to existing methods Table 2 presents a comparison of our results with existing approaches. We obtain state-of-the-art performance by a large margin over existing models and over our baseline model without adaptation. However, using captions as support data and trained on all question types (*number*, *yes/no*, and *other*), the model performs poorly. We hypothesized that evidence for the *number* and *yes/no* questions was difficult to extract from captions. We therefore trained a model with adaptation using only *other* questions. This model performs significantly better and clearly improves over the baseline. We indeed observed that captions seldom include counts or numbers, which can explain why they do not help on the corresponding questions. In the case of binary questions, it is possible that a different relevance function could address the issue.

Qualitative results Fig. 5 presents results of our best models (using QAs or captions) with visualizations of support data sampled according to the relevance function. We observe that the retrieved support data is both semantically and visually relevant to each question.

Additional experiments and qualitative results are provided in the supplementary material.

	VQA-CP v2 Test			
	Overall	Yes/no	Numbers	Other
SAN [41]	24.96	38.35	11.14	21.74
GVQA [1]	31.30	57.99	13.68	22.14
UpDown [33]	39.06	62.41	15.12	34.47
UpDown + regularizer [27]	42.04	65.49	15.87	36.60
Ours without adaptation	40.71	52.22	11.85	42.88
Ours with adaptation and, as support data:				
QAs (VQA-CP tr.), $r=r_1r_2r_3$	46.00	58.24	29.49	44.33
Captions (COCO), $r=r_1r_3$	39.84	48.78	12.40	42.93
Captions, trained only on 'Other' q.	–	–	–	43.95

Table 2. Comparison with existing methods (accuracy on VQA-CP v2). Our method significantly improves over the comparable baseline (the same VQA model without adaptation) and obtains performance superior to all existing models. Gray cells are not directly comparable to others as they use additional data (as in Table 1).

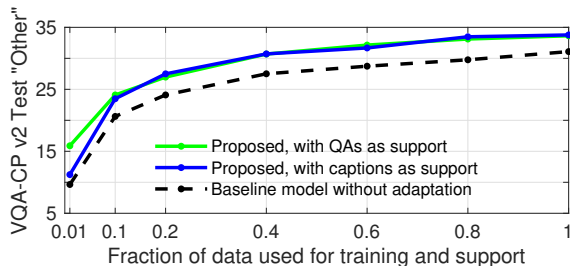


Figure 4. Accuracy as a function of the amount of data used.

5. Conclusions

We presented a new approach to VQA in which the model is trained to interface with an external source of data, and to use it to support its answering process. This is a significant departure from the classical training of a static model on a fixed dataset, which is obviously limited by finite capacity of the model and of the dataset. In contrast, our method retrieves information from the external source specifically for each given question. It then adapts the weights of its underlying VQA model, incorporating information from the external data, and specializing its capabilities to a narrow domain around the input question.

Our experiments demonstrate the benefits of the approach over existing models. It proves intrinsically more robust to out-of-distribution data, and it generalizes to different distributions when provided with novel support data. The model also introduces novel capabilities, in particular for leveraging non-VQA data (image captions) to support the answering process. This presents a number of opportunities to future research, for accessing “black box” data sources, such as web searches and dynamic databases. This opens the door to systems capable of reasoning over vision and language beyond the limited domain covered by any given training set.









Input question	Random selection of retrieved support data			Predicted scores
 Which sport is this ? Correct answer: tennis.	 What sport is taking place ? tennis.	 What sport are they playing ? tennis.	 What sport is this ? tennis.	Without adaptation: <ul style="list-style-type: none"> <input checked="" type="checkbox"/> soccer <input checked="" type="checkbox"/> tennis <input type="checkbox"/> football <input type="checkbox"/> frisbee <input type="checkbox"/> polo After adaptation: <ul style="list-style-type: none"> <input checked="" type="checkbox"/> tennis <input type="checkbox"/> soccer <input type="checkbox"/> frisbee <input type="checkbox"/> polo <input type="checkbox"/> football
 What are two men cutting ? Correct answer: cake.	 What is man cutting pizza with ? knife.	 What object are all four men holding ? knife.	 What are men doing ? cutting cake.	Without adaptation: <ul style="list-style-type: none"> <input type="checkbox"/> knife <input checked="" type="checkbox"/> cake <input type="checkbox"/> candles <input type="checkbox"/> frosting <input type="checkbox"/> cutting cake After adaptation: <ul style="list-style-type: none"> <input checked="" type="checkbox"/> cake <input type="checkbox"/> cutting cake <input type="checkbox"/> yes <input type="checkbox"/> nothing <input type="checkbox"/> knife
 What season is this ? Correct answer: winter.	 What season is it ? fall.	 What season is this ? summer.	 When are these flowers in season ? summer.	Without adaptation: <ul style="list-style-type: none"> <input checked="" type="checkbox"/> winter <input checked="" type="checkbox"/> fall <input type="checkbox"/> spring <input type="checkbox"/> summer <input type="checkbox"/> snow After adaptation: <ul style="list-style-type: none"> <input checked="" type="checkbox"/> winter <input type="checkbox"/> fall <input type="checkbox"/> summer <input type="checkbox"/> spring <input type="checkbox"/> unknown
 Is this breakfast or dinner ? Correct answer: dinner.	 Dinner table with glasses of wine and plates of cheese and crackers.	 Omelet, toast and fruit for breakfast sitting on a table.	 Restaurant table lined for breakfast with breakfast with plates of food.	Without adaptation: <ul style="list-style-type: none"> <input checked="" type="checkbox"/> dinner <input checked="" type="checkbox"/> breakfast <input checked="" type="checkbox"/> dessert <input type="checkbox"/> lunch <input type="checkbox"/> no After adaptation: <ul style="list-style-type: none"> <input checked="" type="checkbox"/> dessert <input type="checkbox"/> cake <input type="checkbox"/> desert <input type="checkbox"/> yes <input type="checkbox"/> lunch
 Food on dinner table in a plate.	 Breakfast plate with egg on toast and greens.	 Table set for breakfast with ham, hashbrowns, croissants and eggs.		

Figure 5. Qualitative results comparing the top-5 answers and their scores predicted by the baseline, and by our model after adaptation. The retrieved support data (random samples are shown) is both visually and semantically relevant to each question.

References

- [1] A. Agrawal, D. Batra, D. Parikh, and A. Kembhavi. Don't just assume; look and answer: Overcoming priors for visual question answering. In *CVPR*, 2018. 2, 3, 6, 7
- [2] A. Agrawal, A. Kembhavi, D. Batra, and D. Parikh. C-vqa: A compositional split of the visual question answering (vqa) v1.0 dataset. *arXiv preprint arXiv:1704.08243*, 2017. 3
- [3] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang. Bottom-up and top-down attention for image captioning and vqa. *arXiv preprint arXiv:1707.07998*, 2017. 11
- [4] M. Andrychowicz, M. Denil, S. Gomez, M. W. Hoffman, D. Pfau, T. Schaul, and N. de Freitas. Learning to learn by gradient descent by gradient descent. In *Advances in Neural Information Processing Systems*, pages 3981–3989, 2016. 3
- [5] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh. VQA: Visual Question Answering. In *Proc. IEEE Int. Conf. Comp. Vis.*, 2015. 2, 3
- [6] S. Bengio, Y. Bengio, J. Cloutier, and J. Gecsei. On the optimization of a synaptic learning rule. In *Preprints Conf. Optimality in Artificial and Biological Neural Networks*, pages 6–8. Univ. of Texas, 1992. 3
- [7] Y. Bengio, S. Bengio, and J. Cloutier. *Learning a synaptic learning rule*. Université de Montréal, Département d'informatique et de recherche opérationnelle, 1990. 3
- [8] C. Buck, J. Bulian, M. Ciaramita, A. Gesmundo, N. Houlsby, W. Gajewski, and W. Wang. Ask the right questions: Active question reformulation with reinforcement learning. *arXiv preprint arXiv:1705.07830*, 2017. 3, 6
- [9] W.-L. Chao, H. Hu, and F. Sha. Cross-dataset adaptation for visual question answering. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 1
- [10] X. Chen, H. Fang, T.-Y. Lin, R. Vedantam, S. Gupta, P. Dollar, and C. L. Zitnick. Microsoft COCO captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. 2
- [11] A. Das, S. Kottur, K. Gupta, A. Singh, D. Yadav, J. M. Moura, D. Parikh, and D. Batra. Visual Dialog. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2
- [12] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. *arXiv preprint arXiv:1703.03400*, 2017. 2, 3, 4, 5, 11
- [13] C. Finn, T. Yu, T. Zhang, P. Abbeel, and S. Levine. One-shot visual imitation learning via meta-learning. In *Conference on Robot Learning (CoRL)*, pages 357–368, 2017. 3
- [14] A. Fukui, D. H. Park, D. Yang, A. Rohrbach, T. Darrell, and M. Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. *arXiv preprint arXiv:1606.01847*, 2016. 2
- [15] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. *arXiv preprint arXiv:1612.00837*, 2016. 2, 3, 6
- [16] S. Hochreiter, A. S. Younger, and P. R. Conwell. Learning to learn using gradient descent. In *International Conference on Artificial Neural Networks*, pages 87–94. Springer, 2001. 3
- [17] P. Huang, C. Wang, R. Singh, W. Yih, and X. He. Natural language to structured query generation via meta-learning. In *HLT-NAACL*, pages 732–738, 2018. 2, 3
- [18] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, M. Bernstein, and L. Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *arXiv preprint arXiv:1602.07332*, 2016. 2
- [19] 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 *Proc. Eur. Conf. Comp. Vis.*, 2014. 2, 6
- [20] G. Marcus. Deep learning: A critical appraisal. *arXiv preprint arXiv:1801.00631*, 2018. 1
- [21] D. K. Naik and R. Mammone. Meta-neural networks that learn by learning. In *Neural Networks, 1992. IJCNN., International Joint Conference on*, volume 1, pages 437–442. IEEE, 1992. 3
- [22] K. Narasimhan, A. Yala, and R. Barzilay. Improving information extraction by acquiring external evidence with reinforcement learning. *arXiv preprint arXiv:1603.07954*, 2016. 3, 6
- [23] A. Nichol, J. Achiam, and J. Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018. 5
- [24] A. Nichol, J. Achiam, and J. Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018. 11
- [25] R. Nogueira and K. Cho. Task-oriented query reformulation with reinforcement learning. *arXiv preprint arXiv:1704.04572*, 2017. 3, 6
- [26] J. Pennington, R. Socher, and C. Manning. Glove: Global Vectors for Word Representation. In *Conference on Empirical Methods in Natural Language Processing*, 2014. 11
- [27] S. Ramakrishnan, A. Agrawal, and S. Lee. Overcoming language priors in visual question answering with adversarial regularization. 2018. 3, 7
- [28] S. K. Ramakrishnan, A. Pal, G. Sharma, and A. Mittal. An empirical evaluation of visual question answering for novel objects. *arXiv preprint arXiv:1704.02516*, 2017. 2, 3
- [29] S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. 2017. 3
- [30] K. Saenko, B. Kulis, M. Fritz, and T. Darrell. Adapting visual category models to new domains. In *European conference on computer vision*, pages 213–226. Springer, 2010. 1
- [31] J. Schmidhuber. *Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook*. PhD thesis, Technische Universität München, 1987. 3
- [32] J. Schmidhuber. Learning to control fast-weight memories: An alternative to dynamic recurrent networks. *Neural Computation*, 4(1):131–139, 1992. 3
- [33] D. Teney, P. Anderson, X. He, and A. van den Hengel. Tips and tricks for visual question answering: Learnings from the 2017 challenge. 2018. 2, 3, 4, 6, 7, 11
- [34] D. Teney and A. van den Hengel. Zero-shot visual question answering. 2016. 2, 3

- [35] D. Teney and A. van den Hengel. Visual question answering as a meta learning task. 2017. [2](#), [3](#)
- [36] S. Thrun. Lifelong learning algorithms. In *Learning to learn*, pages 181–209. Springer, 1998. [1](#)
- [37] S. Thrun and L. Pratt. *Learning to learn*. Springer Science & Business Media, 2012. [1](#), [3](#)
- [38] P. Wang, Q. Wu, C. Shen, A. van den Hengel, and A. Dick. Explicit knowledge-based reasoning for visual question answering. *arXiv preprint arXiv:1511.02570*, 2015. [3](#)
- [39] Q. Wu, D. Teney, P. Wang, C. Shen, A. Dick, and A. van den Hengel. Visual question answering: A survey of methods and datasets. *Computer Vision and Image Understanding*, 2017. [2](#)
- [40] Q. Wu, P. Wang, C. Shen, A. Dick, and A. v. d. Hengel. Ask Me Anything: Free-form Visual Question Answering Based on Knowledge from External Sources. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016. [3](#)
- [41] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola. Stacked Attention Networks for Image Question Answering. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016. [7](#)
- [42] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *Advances in neural information processing systems*, pages 3320–3328, 2014. [1](#)
- [43] T. Yu, C. Finn, A. Xie, S. Dasari, T. Zhang, P. Abbeel, and S. Levine. One-shot imitation from observing humans via domain-adaptive meta-learning. 2018. [3](#)
- [44] A. L. Yuille and C. Liu. Deep nets: What have they ever done for vision? *arXiv preprint arXiv:1805.04025*, 2018. [1](#)
- [45] M. D. Zeiler. ADADELTA: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012. [11](#)
- [46] P. Zhang, Y. Goyal, D. Summers-Stay, D. Batra, and D. Parikh. Yin and yang: Balancing and answering binary visual questions. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016. [3](#)
- [47] B. Zhou, Y. Tian, S. Sukhbaatar, A. Szlam, and R. Fergus. Simple baseline for visual question answering. *arXiv preprint arXiv:1512.02167*, 2015. [3](#)
- [48] Y. Zhu, O. Groth, M. Bernstein, and L. Fei-Fei. Visual7W: Grounded Question Answering in Images. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016. [2](#)