

CAPWAP: Captioning with a Purpose

Adam Fisch^{1*} Kenton Lee² Ming-Wei Chang² Jonathan H. Clark² Regina Barzilay¹

¹Massachusetts Institute of Technology, ²Google Research
{fisch, regina}@csail.mit.edu,
{kentonl, mingweichang, jhclark}@google.com

Abstract

The traditional image captioning task uses generic reference captions to provide textual information about images. Different user populations, however, will care about different visual aspects of images. In this paper, we propose a new task, **Captioning with A Purpose** (CAPWAP). Our goal is to develop systems that can be *tailored* to be useful for the information needs of an intended population, rather than merely provide generic information about an image. In this task, we use question-answer (QA) pairs—a natural expression of information need—from users, instead of reference captions, for both training and post-inference evaluation. We show that it is possible to use reinforcement learning to directly optimize for the intended information need, by rewarding outputs that allow a question answering model to provide correct answers to sampled user questions. We convert several visual question answering datasets into CAPWAP datasets, and demonstrate that under a variety of scenarios our purposeful captioning system learns to anticipate and fulfill specific information needs better than its generic counterparts, as measured by QA performance on user questions from unseen images, when using the caption alone as context.

1 Introduction

The image captioning task typically selects for captions having high similarity with generic human references. While this task definition has driven much of the research in the field, the end-purpose of these captions is not always clearly articulated. We argue that (1) generic annotations may not be representative of users’ information needs, (2) user questions are a more natural way of articulating information needs, and (3) optimizing captions to provide correct answers to those questions allows

*Work primarily completed while interning at Google.



Task	Caption	Information Need
Captioning	There is a green bus.	(Unspecified)
Visual QA	(Unspecified)	Where’s it headed?
CAPWAP	At least three people are boarding the #14 bus to Bembridge.	Which bus is this? Where’s it headed? How many people are boarding?

Figure 1: The informational purpose of generic captioning is not clearly defined, and VQA provides only reactionary information. The objective of the CAPWAP task is ultimately to provide more informative captions that specifically *anticipate* and *satisfy* users’ potential needs. In CAPWAP, we use QA as an *implicit* signal for information need: e.g., in the image above, a good caption that has been generated in advance should be able to be used to answer, *Where is this bus headed?*

training to focus on information need. For example, in the VizWiz mobile application (Bigham et al., 2010), visually impaired users upload images from their everyday lives, along with questions about them that need to be answered. These questions serve as a powerful signal for aspects of the image that they find important.

Consider the image in Figure 1, where an annotator might provide a generic caption such as *There is a green bus*. This may be used to answer: *What color is the bus?* However, it would provide no utility to a user asking: *Where is this bus headed?* In fact, examples from VizWiz demonstrate a clear disconnect between the type of information provided by







Task	Training Data	Prediction Function	Evaluation
Captioning	{  , ref.caption}	{  } → pred.caption	SIMILARITY(pred.caption, ref.caption)
Visual QA	{  , ref.question, ref.answer}	{  , ref.question} → pred.answer	ACCURACY(pred.answer, ref.answer)
CAPWAP	{  , ref.question, ref.answer}	{  } → pred.caption	ACCURACY(QA(ref.question, pred.caption), ref.answer)

Table 1: The CAPWAP task combines elements of generic image captioning with visual question answering. Training consists of images paired with visual questions and answers. A CAPWAP model should directly *anticipate* user information needs by outputting *captions* that can be used to answer *future* questions drawn from a distribution similar to the training data. Accordingly, the “QA” function represents the inference of an answer to a question using the generated caption as context. We approximate this with a strong automatic question answering model.

today’s systems (e.g., arbitrary descriptions of entities and actions) versus what visually-impaired users need to know (e.g., fine-grained details to help make decisions).

Here, we propose an alternative framing for captioning: Captioning with A Purpose (CAPWAP). We do not assume the existence of a *universal* caption distribution. A good caption is highly subjective; different users will care about different aspects of a given image. Instead, we assume a distribution of visual question-answer pairs that are representative of population’s information needs. Here we aim to map images to text that can serve as context to answer likely questions under this distribution. At test time, the goal is to anticipate similar user questions for a new image, and implicitly answer them *before* they even need to be asked.

We use image-question-answer triplets as supervision, and require the model to generate from the latent space of captions that provide contextual support for the answer (Table 1). Within our task definition, any sampled caption that can be used to answer these questions is considered useful. Under this formulation, very different captions may be scored identically if they deliver the same content—regardless of word choice. Note that this is different from either standard visual question answering (VQA) or query-focused summarization: the target questions are not available prior to generation; at test time, they are used *only* for evaluation.

Existing approaches cannot be readily applied in this setting, as there are no gold reference captions for training—and off-the-shelf captioning systems transfer quite poorly (§6). To address the new learning challenge that arises in CAPWAP, we propose a novel model-in-the-loop reinforcement learning (RL) approach that acts as a strong baseline for

this task. Our approach assumes a fixed question answering (QA) system that predicts an answer to a question using some input context. The captioning model receives a reward if it generates text which the QA system can use to predict the correct answer. Applying RL, however, is nontrivial. A naïve exploration of the caption generation space can lead to sparse rewards—resulting in long training times and disappointing quality. We show that our approach can be significantly improved by using a novel, synthetic pre-training routine to push the initial policy towards areas of high-reward.

We repurpose four VQA datasets for CAPWAP: VQA (Goyal et al., 2017), GQA (Hudson and Manning, 2019), Visual7W (Zhu et al., 2016), and VizWiz (Gurari et al., 2018). These datasets range in style from synthetic QA pairs (GQA) to natural information-seeking questions asked by visually-impaired users (VizWiz). We find that our method produces significantly more informative captions with respect to the given questions (up to $3.8\times$ exact match), compared to models trained on generic captions from COCO (Lin et al., 2014).

Our key contributions are as follows:

1. We define a new task (CAPWAP) that generates image captions for the *purpose* of fulfilling specific information needs expressed by different target user populations.
2. We demonstrate that our information-need-driven model can generate much higher quality captions on this task than those of state-of-the-art traditional *generic* captioning systems.
3. We propose a novel synthetic pre-training routine that greatly improves the performance of reinforcement learning under this new paradigm.

2 Related Work

Since the early days of the field, human-written references have been used for the supervised training and evaluation of text generation systems, including image captioning, summarization, and other related applications (Edmundson, 1969; Lin and Hovy, 2003; Ordonez et al., 2011; Vinyals et al., 2015). Recently, researchers have begun to consider a multitude of different objectives for reference comparison (Böhm et al., 2019; Gao et al., 2019), or even parametric regressions trained on human judgements (Louis and Nenkova, 2013; Peyrard and Gurevych, 2018). Though diverse in approach, each ultimately relies on designing a robust general-purpose metric. In practice, engineering such a metric is challenging—if at all possible (Spärck Jones, 1994, 1999). Here we take a more empirical approach by relying on the information need expressed by users’ questions.

Many studies have observed that reference-trained captioning models suffer from systematic usability issues—including being rigid, neglecting relevant image aspects, and regurgitating frequent phrases (Wang et al., 2017; Dai et al., 2017). As a result, much effort has been focused on developing secondary, corrective objectives—for instance, “discriminability” losses encouraging captions to be unique (Dai and Lin, 2017; Liu et al., 2018; Luo et al., 2018). While these measures provide some fixes, they do not necessarily reflect user information needs—a central concept in CAPWAP.

The idea of using QA for assessing information quality has been proposed in recent work for text summarization (Arumae and Liu, 2019; Eyal et al., 2019; Scialom et al., 2019). The primary distinctions with our work are both the domain (images) and how questions are obtained—both of which impact the task objective and learning procedure. In this prior work, questions are generated programmatically (e.g., following Hermann et al., 2015). Such “questions” may not necessarily reflect real user preferences. Our work focuses on QA not as just another method to improve standard reference-based metrics, but as a key, flexible way of formulating user information need—and as such we focus on challenging, *real* QA datasets. Furthermore, we train on this signal, rather than rely on it solely for evaluation (Wang et al., 2020).

Efforts to leverage VQA resources to drive image

captioning, and vice-versa, via variations of transfer learning, have also received extensive interest in recent years (Li et al., 2018; Wu et al., 2019; Yang and Xu, 2019). As opposed to optimizing metrics for specific VQA or supervised captioning benchmarks, the primary focus in CAPWAP is on modeling the target user population in order to anticipate the correct information-need.

In a similar vein, VQA and textual QA resources have also been leveraged for active learning (Shen et al., 2019; Li et al., 2017), where the model learns to query its environment for information *it* is uncertain about to help improve its performance on the given task. The key distinction with our work is the *directionality* of the questions. In CAPWAP, the model uses questions posed by the users to infer *their* latent information need—which is a distinctly different, and quite challenging, setting.

3 Problem Formulation

We begin by formulating the CAPWAP task. In our setting, questions and answers are the only source of direct supervision assumed during training. At test time, the model is not given questions in advance, but rather must *anticipate* the information need of the user, and generate captions that answer the forthcoming questions in expectation.

Task Setting: Given an image \mathbf{x} the model must output a caption \mathbf{y} , such that \mathbf{y} entails the answer \mathbf{a} for a question-answer pair (\mathbf{q}, \mathbf{a}) sampled from some underlying distribution \mathcal{D} . Examples from \mathcal{D} are given during training, but are not known in advance by the generation model at test time.

Information Need: We assume that the QA data from \mathcal{D} is derived by the following process:

1. an image \mathbf{x} is drawn from distribution $p(\mathbf{x})$;
2. a question-answer pair (\mathbf{q}, \mathbf{a}) targeting an informative detail of \mathbf{x} perceived as important to a user in \mathcal{D} is drawn from distribution $p(\mathbf{q}, \mathbf{a}|\mathbf{x})$.

The operating assumption is that the marginal distribution over (\mathbf{q}, \mathbf{a}) pairs represents the visual interests of the *typical* user. In other words, answers to common questions represent the type of information that is often considered important. This is comparable to *content selection* (Peyrard, 2019).

Question Anticipation: We do not assume the existence of a “gold” caption. Rather, the caption \mathbf{y} is assumed to be a latent variable, and $G_\theta(\mathbf{y}|\mathbf{x})$ is a

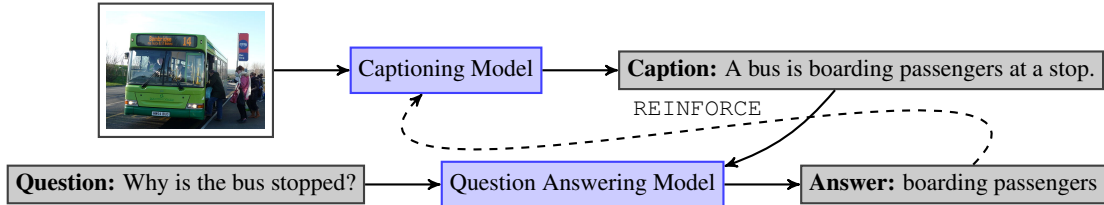


Figure 2: Overview of our proposed approach to the CAPWAP task. The captioning model $G_\theta(\mathbf{y}|\mathbf{x})$ is learned using supervision from question-answer-image triples. Generated text that can be used to answer the question correctly, according to an extractive question answering model, is rewarded in our model-in-the-loop reinforcement learning framework. The questions, answers, and the question answering system are discarded after training.

stochastic generator that we must learn. A sample $\mathbf{y} \sim G_\theta(\mathbf{y}|\mathbf{x})$ should provide *contextual support* for a new, randomly sampled question-answer pair. We estimate this using the accuracy of a pre-trained QA model $\mathcal{M}(\mathbf{q}, \mathbf{y})$, when using \mathbf{y} as context for \mathbf{q} . CAPWAP requires maximizing the expectation:

$$\operatorname{argmax}_\theta \mathbb{E}_{G_\theta(\mathbf{y}|\mathbf{x})} [\mathbb{E}_{p(\mathbf{q}, \mathbf{a}|\mathbf{x})} [\mathcal{R}(\mathbf{y}, \mathbf{q}, \mathbf{a})]] \quad (1)$$

where θ parameterizes $G_\theta(\mathbf{y}|\mathbf{x})$, and we choose our reward to be $\mathcal{R}(\mathbf{y}, \mathbf{q}, \mathbf{a})$, any appropriate accuracy metric for comparing the output of $\mathcal{M}(\mathbf{q}, \mathbf{y})$ with \mathbf{a} (expressed as ACCURACY in Table 1).

CAPWAP vs. Other Tasks: Table 1 compares our setting to those of both standard (generic) captioning and visual question answering. Both standard captioning and CAPWAP models output a single caption per image, but CAPWAP does not compare to references. Both VQA and CAPWAP models are trained and evaluated with QA data, but CAPWAP does not provide the question prior to generation. VQA models output *single answers*, whereas CAPWAP models output anticipatory *contexts*.

4 An Approach to CAPWAP

Given that we only have access to question-answer pairs during training, but not during inference, how can we learn a model for this task? Eq. 1 naturally lends itself to a reinforcement learning (RL) framework where the model receives a reward \mathbf{r} (e.g., $\mathbf{r} = \mathcal{R}(\mathbf{y}, \mathbf{q}, \mathbf{a})$) for each generated caption \mathbf{y} and training QA pair (\mathbf{q}, \mathbf{a}) . $G_\theta(\mathbf{y}|\mathbf{x})$ can be cast as a policy, and updated with policy gradients.

Optimizing such a policy, however, poses a technical challenge because the model is only rewarded for correct (or partially correct) answers, which is initially a rare event. Transferring $G_\theta(\mathbf{y}|\mathbf{x})$ from *generic* captioning data can be a useful starting point. Our method then follows this recipe:

1. Initialize $G_\theta(\mathbf{y}|\mathbf{x})$ using fully-supervised off-the-shelf captioning data, $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \sim \mathcal{D}_{\text{generic}}$;
2. Fine-tune $G_\theta(\mathbf{y}|\mathbf{x})$ using policy gradient on targeted visual QA data, $(\mathbf{x}, \mathbf{q}, \mathbf{a}) \sim \mathcal{D}_{\text{target}}$.

In Sections 4.1 and 4.2 we detail our model for $G_\theta(\mathbf{y}|\mathbf{x})$, and the above training procedure.

Note that $\mathcal{D}_{\text{generic}}$ is assumed to be *out-of-domain* for our intended captioning purpose, $\mathcal{D}_{\text{target}}$. Since we are interested in diverse user-generated questions and information needs, the generic captioning data can often diverge dramatically from our end goal. To improve transfer, in Section 4.3 we further develop a novel mechanism for automatically generating in-domain synthetic data that can be used as pre-training for guiding $G_\theta(\mathbf{y}|\mathbf{x})$ towards balanced areas of high reward in $\mathcal{D}_{\text{target}}$.

4.1 Model Architecture

We briefly describe our base captioning model, which consists of a Faster R-CNN and Transformer-based encoder-decoder, following the sequence-to-sequence framework common in state-of-the-art image captioning systems (Anderson et al., 2018; Vinyals et al., 2015; Zhou et al., 2019). See Appendix A for full technical details. Given an image \mathbf{x} , we first represent it as a sequence of detected object bounding box embeddings, computed from a pre-trained Faster R-CNN model (Anderson et al., 2018). We then generate caption word-pieces $\mathbf{y} = (y_1, \dots, y_n)$ using a Transformer-based architecture (Vaswani et al., 2017).

4.2 Policy Training

We describe our RL framework for training our captioning model using QA data. See Appendix B for additional technical details, including hyperparameter settings and optimization choices.

Initialization: We initialize $G_\theta(\mathbf{y}|\mathbf{x})$ using maximum likelihood estimation (MLE) on a corpus of out-of-domain generic captions $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$, as common practice (Ranzato et al., 2016). This warm-starts our policy with an initial set of grounded image concepts, albeit not necessarily the ones we ultimately care about. Given the generic reference $\tilde{\mathbf{y}} = (\tilde{y}_1, \dots, \tilde{y}_n)$, we minimize the cross-entropy:

$$\mathcal{L}_{XE}(\theta) = - \sum_{i=1}^n \log G_\theta(\tilde{y}_i | \tilde{\mathbf{x}}, \tilde{y}_{j < i}) \quad (2)$$

QA Model: We implement the QA model \mathcal{M} using a BERT_{LARGE} extractive model fine-tuned on SQuAD 2.0 (Rajpurkar et al., 2018)—which contains unanswerable questions. As an extractive model, \mathcal{M} predicts a span $\mathbf{y}_{i\dots j}$. Important for our use-case, \mathcal{M} is both able to be accurate when predicting the answer \mathbf{a} when \mathbf{a} is present in \mathbf{y} , and also able to abstain from answering when \mathbf{a} is *not* logically entailed (i.e., predict “no answer”).

QA Reward: We take $\mathcal{R}(\mathbf{y}, \mathbf{q}, \mathbf{a})$ from Eq. 1 as the F1 score of the predicted answer with the gold answer. We control for reward noise with a confidence threshold for predicting “no answer.”

Policy Gradient: We use REINFORCE with a baseline (Williams, 1992) to compute the policy gradient $\nabla_\theta \mathcal{L}_{QA}(\theta)$ of the QA reward:

$$-\mathbb{E}_{G_\theta(\mathbf{y}|\mathbf{x})} [(\mathcal{R}(\mathbf{y}, \mathbf{q}, \mathbf{a}) - b) \nabla_\theta \log G_\theta(\mathbf{y}|\mathbf{x})] \quad (3)$$

We take b as $\mathcal{R}(\hat{\mathbf{y}}, \mathbf{q}, \mathbf{a})$, where $\hat{\mathbf{y}}$ is the argmax (test-time prediction) of G_θ , following the self-critical method of Rennie et al. (2017).

4.3 Synthetic Policy Pre-Training

In the beginning of training, the generated captions typically do not correctly answer many questions, leading to almost no reward signal. More formally, the reward is sparse if the policy $G_\theta(\mathbf{y}|\mathbf{x})$ is not well-initialized. As a result, REINFORCE becomes extremely sample-inefficient. When the target distribution is strikingly divergent from the one present in the generic captioning data—a key setting in this work—supervised pre-training on the out-of-domain data does not yield a usable initialization. As a substitute, we derive a method for generating a synthetic dataset of captions $\mathcal{D}_{\text{synthetic}}$ with high-reward as a form of guided policy search (Levine and Koltun, 2013). The full method then consists of three stages that train on the three datasets: $\mathcal{D}_{\text{generic}} \rightarrow \mathcal{D}_{\text{synthetic}} \rightarrow \mathcal{D}_{\text{target}}$.

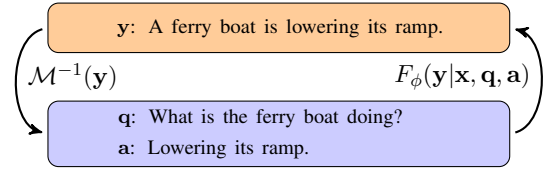


Figure 3: A demonstration of reverse engineering the connections between question generation (\mathcal{M}^{-1}) and context generation (F_ϕ). \mathbf{x} is the image (not shown). See Algorithm 1 in Appendix C for full details.

For the extractive QA model to possibly yield a positive reward, the answer must be a span of the caption. When the question and answer are known in advance, it is typically fairly simple to reverse engineer a candidate caption that meets this constraint (e.g., by inverting *wh*-movement). Figure 3 demonstrates this concept. If we have an auxiliary model $F_\phi(\mathbf{y}|\mathbf{x}, \mathbf{q}, \mathbf{a})$ that can automate this reverse engineering step, we can synthetically generate captions to use for pre-training, as in Eq. 2.¹

QA Conditional Model: Motivated by this, we learn $F_\phi(\mathbf{y}|\mathbf{x}, \mathbf{q}, \mathbf{a})$ by explicitly conditioning on QA pairs when generating a caption that supports the answer span by design. Concretely, we include the word-pieces of the question $\mathbf{q} = (q_1, \dots, q_l)$ and answer $\mathbf{a} = (a_1, \dots, a_m)$ as inputs when decoding \mathbf{y} , while \mathbf{y} satisfies $\mathcal{M}(\mathbf{y}, \mathbf{q}) = \mathbf{a}$.

How do we train $F_\phi(\mathbf{y}|\mathbf{x}, \mathbf{q}, \mathbf{a})$ effectively without access to any paired data $(\mathbf{x}, \mathbf{q}, \mathbf{a}, \mathbf{y})$? We create automatic $(\tilde{\mathbf{x}}, \hat{\mathbf{q}}, \hat{\mathbf{a}}, \tilde{\mathbf{y}})$ examples from the out-of-domain generic captioning data used in Section 4.2 by using the (text-based) question generation model \mathcal{M}^{-1} of Alberti et al. (2019). At a high level, given a generic caption $\tilde{\mathbf{y}} = (\tilde{y}_1, \dots, \tilde{y}_n)$, (1) picks an answer span $\hat{\mathbf{a}} \subseteq \tilde{\mathbf{y}}$, (2) generates a question $\hat{\mathbf{q}}$ following some inferred distribution $p(\mathbf{q}|\hat{\mathbf{a}}, \tilde{\mathbf{y}})$, and (3) confirms that the sample obeys “round-trip filtering”, i.e., that the original QA model answers the synthetic example correctly ($\mathcal{M}(\hat{\mathbf{q}}, \tilde{\mathbf{y}}) = \hat{\mathbf{a}}$). We then train the model for F_ϕ using conditional cross-entropy:

$$\mathcal{L}_{CXE}(\phi) = - \sum_{i=1}^n \log F_\phi(\tilde{y}_i | \tilde{\mathbf{x}}, \hat{\mathbf{q}}, \hat{\mathbf{a}}, \tilde{y}_{j < i}) \quad (4)$$

Synthetic Data Generation: After training, we transfer $F_\phi(\mathbf{y}|\mathbf{x}, \mathbf{q}, \mathbf{a})$ with fixed weights to generate reverse engineered captions $\hat{\mathbf{y}}$ using the

¹Methods for constrained decoding (Anderson et al., 2016; Hokamp and Liu, 2017, *inter alia*) that enforce $\mathbf{a} \subseteq \mathbf{y}$ are related, yet complementary, and can be incorporated into any F_ϕ . It is more important to ensure that not only is the answer contained in the caption, but also that it is logically supported.

true $(\mathbf{x}, \mathbf{q}, \mathbf{a}, \text{null})$ examples from our target QA datasets. For each example, we decode the top- k captions using beam search, and keep those with $R(\hat{\mathbf{y}}, \mathbf{q}, \mathbf{a}) \geq c$, where c is a threshold (e.g., $c = 1.0$ for an exact match F1 score). These examples paired with the high-scoring captions are used to create the synthetic captions dataset $\mathcal{D}_{\text{synthetic}}$. We then use $\mathcal{D}_{\text{synthetic}}$ as further weak supervision for initializing $G_{\theta}(y|\mathbf{x})$, again following Eq. 2. See Appendix C for additional technical details.

5 Experimental Setup

Evaluation: Our primary evaluation assumes a dataset of questions and answers about images. Conceptually, if the correct answers are supported by the generated caption in expectation, then we consider it to be sufficiently informative.²

Automatic Evaluation: Our automatic proxy of informativeness utilizes the state-of-the-art extractive question answering model (\mathcal{M}) described in Section 4.2 that is trained on SQuAD 2.0.³ \mathcal{M} is applied to the given QA pair, with the generated caption as the “context.” We report EM, measuring exact match with the gold answer, and F1, measuring word overlap. If there are multiple answers, then we take the maximum score over all.

Human Evaluation: For human evaluation we ask raters to judge whether a caption is less, equally, or more informative than another caption with respect to the question-answer pair. We also gather human ratings for two properties that are desirable regardless of the target audience: (1) fluency (whether the caption is grammatical and coherent) and (2) fidelity (whether the caption makes any false assertions regarding what is in the image).

5.1 Datasets

We evaluate our method on four converted visual question answering datasets. We filter questions that are unanswerable, or have ‘yes/no’ or non-alphabetic answers.⁴ Appendix D gives additional size, splitting, and pre-processing details.

COCO (Lin et al., 2014): For all experiments,

²Note that traditional captioning metrics such as ROUGE, BLEU, and CIDEr rely on gold references, which are not available in our new setting (in fact, by our definition, there is no one “gold” caption). Thus, we cannot include them.

³For evaluation we turn off the “no answer” option.

⁴Numerical answers that are written out (e.g., *two* vs. 2) are not disqualified. This requirement simplifies evaluation.

we use COCO as the source of out-of-domain generic captions for pre-training. COCO contains images covering 80 object categories and various scenes gathered from Flickr, paired with five human-written reference captions.

CapVQA (Goyal et al., 2017): VQA v2.0 originally contains questions written by crowd-workers where the prompt was to write queries that are easy for humans to answer, but challenging for a hypothetical robot that mainly knows only about objects. VQA is the only dataset we consider that fully covers the same images as COCO.

CapGQA (Hudson and Manning, 2019): GQA contains challenging compositional questions derived from scene graphs of everyday images using various human-specified grammars.

CapVisual7W (Zhu et al., 2016): Visual7W contains questions written by crowd-workers about objects that, in general, require richer and longer answers than those in VQA. We use only the “telling” split of the dataset (i.e., the questions that require open-ended natural language answers).

CapVizWiz (Gurari et al., 2018): VizWiz consists of natural visual questions asked by visually-impaired users of a mobile application who were seeking answers to their daily visual needs. Each question is answered by a remote assistant.

5.2 Generic Captioning Models

In addition to our baseline captioning model trained only to maximize the likelihood of COCO references (MLE in the tables), we compare to two state-of-the-art *generic* image captioning methods (also trained on COCO data). Huang et al. (2019) directly optimizes the CIDEr metric with policy gradients, while Luo et al. (2018) optimizes both CIDEr and a “discrimination” loss intended to encourage models to describe each image’s uniquely identifying aspects. These models are included in order to highlight the differences in applicability between off-the-shelf models trained for generic image captioning versus those for CAPWAP.

6 Results

In the following, we address several key research questions relating to our approach to CAPWAP, and the broader assumptions, strengths, and limitations of using QA to drive the process.

Model	CapVQA		CapGQA		CapVisual7W		CapVizWiz	
	EM	F1	EM	F1	EM	F1	EM	F1
Human reference (from COCO)	16.5	25.7	-	-	-	-	-	-
Luo et al. (2018)	12.0	20.1	9.6	13.9	6.5	11.8	4.7	11.8
Huang et al. (2019)	16.0	25.0	9.9	14.9	6.9	14.0	6.0	13.4
Our generic baseline: MLE	16.8	25.2	8.0	11.1	6.9	13.2	4.9	12.6
Our CAPWAP model: RL	23.1	32.3	15.7	19.3	10.5	18.4	22.5	28.5
Our CAPWAP model: RL + SYN	24.2	33.2	16.6	19.8	9.2	15.4	19.5	27.8

Table 2: *Does the proposed approach better fulfill information needs?* We show question answering test performance when applying an extractive question answering model on predicted captions (see Table 1). Existing captioning models trained on generic references (rows 2-4)—or even the generic references themselves (row 1)—do not capture the information requested by different QA datasets. Applying our RL method for tailoring towards CAPWAP (row 5) leads to more informative captions with respect to those questions (and by extension, for the assumed end-users). Adding synthetic pre-training data (+ SYN) improves results on several datasets (row 6).

A = Purposeful (RL + SYN) vs. B = Generic		
Dataset	Informativeness	
	A > B	B > A
CapVQA	27%	20%
CapGQA	31%	20%
CapVisual7W	38%	22%
CapVizWiz	37%	20%

Table 3: *Do raters think that the proposed approach provides more informative captions?* Human evaluation of the informativeness of captions with respect to our QA datasets agrees with our automatic evaluation—finding our model to have better information coverage than MLE, our baseline without QA rewards.

Evaluation of Generic Captions: We begin by empirically verifying our introductory claim that training on generic reference captions can poorly reflect the varying, user-specific information need. Table 2 presents the results of the baseline generic captioning systems when evaluated in terms of how well the predicted captions support QA over different distributions. Though they are strong methods as measured on the COCO benchmark, unsurprisingly, they still fail to capture all the information necessary to answer diverse visual questions. Performance on CapVizWiz is exceptionally poor; the visually-impaired users ask for information strikingly different than what is represented in COCO. The causes of this poor performance go beyond simple limitations in the current state-of-the-art models; the target references themselves are insufficient. For example, on CapVQA, where the images overlap with COCO and thus human captions are directly available, the average performance of these “gold” references is only slightly better—supporting our conjecture that good captions for one purpose are not necessarily good for

Train	Evaluation			
	CapVQA	CapGQA	CapV7W	CapVizWiz
CapVQA	33.8	14.3	15.6	16.4
CapGQA	24.8	20.2	12.6	13.5
CapV7W	26.0	11.4	15.2	14.0
CapVizWiz	23.9	12.0	12.2	31.3
Generic	25.5	11.7	12.8	14.1

Table 4: *Does the proposed approach tailor to specific information need?* We show transfer performance (dev F1) of the RL + SYN policies learned on different QA datasets. The in-domain F1 peaks indicate that the model is producing distribution-specific captions.

another, even on the same images.

Adaptation to Information Need: We next test the effectiveness of our proposed approach at tailoring captions to meet the specific information need stipulated by our datasets. Our results in Table 2 demonstrate significant improvements by our QA-driven models (RL and RL + SYN) across all four datasets—achieving an average gain of 8.0 absolute F1. Notably, we improve by 7.5 EM over the average human caption on CapVQA, and by 16.5 EM over the best generic model on CapVizWiz. Table 4 further illustrates that the adaptation process is indeed tailored to the respective QA datasets. The improvements on our automatic QA-based metrics (using the proxy model \mathcal{M}) also translate to human judgements. Table 3 presents the results of our human A/B test of our proposed model vs. the MLE baseline. Relative to MLE, we find that our method is significantly more informative with respect to unseen QA pairs across all datasets. As expected, the largest improvements are on the datasets whose questions deviate significantly from the generic COCO content (e.g., CapVizWiz).





Dataset	Input Image	Generic Caption Output	Purposeful Caption (RL + SYN) Output	Unseen Question	Unseen Answer
CapVQA		a man playing tennis	a man in white shirt and white shorts playing a game of tennis on a grass court	what color is he wearing?	white shirt and shorts
CapGQA		a man riding on the back of a horse	a man riding a horse next to a woman on the right side	what color is the horse the man is to the right of?	brown
CapVisual7W		a couple of people that are standing in the snow	two people posing for a picture taken at a ski slope	why are the kids wearing coats ?	it is cold
CapVizWiz		a plate of food that is on a table	this is a picture of a sweet corn frozen dinner	what kind of tv dinner is this?	lean cuisine

Figure 4: Example outputs⁵ comparing our model trained for CAPWAP (RL + SYN) with the baseline model trained on generic COCO references (MLE). These examples are representative of the way in which the various datasets ask about image content. Tendencies include colors for CapVQA, spatial relations for CapGQA, higher-level concepts in CapVisual7W, and OCR in CapVizWiz. Note that our approach to CAPWAP does not (and likely cannot ever perfectly) anticipate *all* unseen questions—but is distributionally closer in terms of content selection.

A = RL + SYN vs. B = RL				
Dataset	Fluency		Fidelity	
	A > B	B > A	A > B	B > A
CapVQA	57%	26%	55%	29%
CapGQA	84%	6%	76%	11%
CapVisual7W	69%	19%	66%	22%
CapVizWiz	38%	37%	24%	47%

Table 5: Does synthetic data improve secondary measures of caption quality? Raters find that this strategy dramatically improves fluency and fidelity (§5) when compared to a model with only on-policy sampling.

Importance of Synthetic Pre-training: A deficiency of QA-based rewards is that they neither explicitly enforce text fluency, nor penalize the system when content is produced that is either *not* relevant or *not* true. On the other hand, when reference captions are available, it is easy to learn a fluent language model. Table 5 shows that incorporating synthetic, guiding “silver” samples from our auxiliary QA conditional model $F_\phi(\mathbf{y}|\mathbf{x}, \mathbf{q}, \mathbf{a})$ (§4.3) to bridge the gap between $\mathcal{D}_{\text{generic}}$ and each considered $\mathcal{D}_{\text{target}}$ dramatically reduces the fluency and fidelity issues that arise from training solely with QA rewards. Ultimately, however, Table 6 shows that our model still suffers on these secondary metrics as compared to the reference-trained MLE baseline. This is a challenge that is shared with nearly all other comparable RL-based methods for text generation (e.g., Guo et al., 2018; Paulus et al., 2018, *et cetera*). Incorporating complementary fluency re-

⁵Examples are chosen to highlight model differences.

A = Purposeful (RL + SYN) vs. B = Generic		
Dataset	Fluency	
	A ≥ B	Fidelity A ≥ B
CapVQA	37%	33%
CapGQA	30%	57%
CapVisual7W	37%	53%
CapVizWiz	34%	72%

Table 6: Do more informative captions come at a cost? Raters find that our tailored “purposeful” approach is less fluent more than half the time (left). Still, this system has greater or equal fidelity to the image content on most datasets (right). At a high level, while the system does give more *relevant* information, it may do so in a less fluent way, a direction we leave for future work.

wards (e.g., via pre-trained language model perplexity) is a valuable direction for future work.

Qualitative Discussion: The qualitative effect of our method is quite intuitive (see Figure 4 as well as Table 7). For example, many questions in CapGQA ask about spatial relations, which is reflected in the generated captions. On the other hand, CapVizWiz users often ask about detailed information about meals, and the adapted model attempts to provide a more useful description beyond a “plate of food.” Of the above datasets, note that only CapVizWiz consists of questions asked by genuinely interested users. Interestingly, this property unearths yet another challenge: CapVizWiz questions can be long-form and quite different from SQuAD, and reverse engineering them (using F_ϕ) for pre-training is noisier (as evidenced by the performance of +SYN

Model	Dataset											
	CapVQA			CapGQA			CapVisual7W			CapVizWiz		
	y	adj%	V ²	y	adj%	V ²	y	adj%	V ²	y	adj%	V ²
Generic: MLE	10.8	6.9	3.6	11.1	7.7	2.5	10.7	6.5	3.1	10.7	9.1	1.3
CAPWAP: RL	16.2	15.8	4.8	15.6	17.5	2.0	20.1	15.5	3.9	5.0	12.2	1.4
CAPWAP: RL + SYN	13.9	11.1	4.8	13.8	8.7	2.0	13.6	7.6	4.4	9.6	7.0	2.6

Table 7: *How do the generated captions differ qualitatively?* We present a number of automatic qualitative measures of caption content calculated over the dev sets: average caption length ($|y|$), adjective production rate (adj%), and the total vocabulary size of the unique unigrams and bigrams emitted ($|V^2|$). Captions are measured in tokens (PTB-style), adjectives are identified using NLTK (Loper and Bird, 2002), and vocabulary size is measured in thousands. Both CAPWAP methods tend to produce longer captions, presumably with more descriptions (higher number of adjectives). Notably, RL + SYN manages to maintain more “natural” adjective production rates and richer language usage (in terms of bigram usage) than RL only, supporting the human quality ratings in Table 5.

Model	Reward	Dataset											
		CapVQA			CapGQA			CapVisual7W			CapVizWiz		
		IN	EM	F1	IN	EM	F1	IN	EM	F1	IN	EM	F1
RL + SYN	answer supported	42.5	24.7	33.8	25.1	16.9	20.2	16.7	9.1	15.2	34.3	22.0	31.3
w/o “no answer”	answer most likely	40.8	23.3	32.3	29.6	19.1	23.8	16.3	8.9	15.1	33.8	20.8	30.9
w/o QA model	answer present	46.4	22.4	31.5	40.2	13.9	21.6	17.4	8.9	15.2	39.0	13.5	26.6
Generic MLE	None	32.9	17.1	25.5	17.0	8.4	11.7	14.2	6.7	12.8	18.5	5.9	14.1

Table 8: *What is the impact of using the QA model to provide rewards?* We present an ablation study across our different datasets when using the QA model with the “no answer” option or not, as well as a simple indicator reward, $\mathbf{1}\{a \subseteq y\}$, that simply measures if the answer string is present at all (without running the expensive QA model). Our results show that while the indicator reward increases the indicator metric (IN) the most, these are likely mostly spurious or disfluent generations. Using the QA model improves the F1 and EM scores across all datasets—and in all cases except one improves further when confidence is used.

in Table 5). While the artificial settings of the other datasets are not ideal, their diversity serves to demonstrate the flexibility of our approach.

Ablation Studies: Tables 7 and 8 show the effects of different design choices in our RL and RL + SYN models. A significant challenge for CAPWAP systems, as previously discussed and illustrated in Table 6, is learning information need while maintaining fluency. Table 7 shows how synthetic pre-training regularizes the model to stay closer to human-level production patterns. Similarly, Table 8 shows how using the QA model to provide rewards (as opposed to a simple keyword search) helps the model avoid spurious rewards.

Future Work: The CAPWAP paradigm introduces new challenges for learning effective systems, some of which our approach solves, and others which it still leaves open (e.g., maintaining fluency and fidelity). While some may be addressed by large-scale multi-modal models (Li et al., 2019; Tan and Bansal, 2019), it is still unclear whether they would fully cover the diversity of information

that real users are interested in (e.g., OCR).

7 Conclusion

We defined and studied the CAPWAP task, where question-answer pairs provided by users are used as a source of supervision for learning their visual information needs. Our results indicate that measuring caption content by its ability to logically support the answers to typical QA pairs from a target audience is (1) not only feasible, but also (2) a good proxy for uncovering information need. We hope this work will motivate the image captioning field to learn to anticipate and provide for the information needs of specific user communities.

Acknowledgements

We thank the MIT NLP group, the Google AI Language team, and the anonymous EMNLP reviewers for their valuable feedback. AF is supported in part by an NSF Graduate Research Fellowship.

References

- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. TensorFlow: Large-scale machine learning on heterogeneous systems. Software available from tensorflow.org.
- Chris Alberti, Daniel Andor, Emily Pitler, Jacob Devlin, and Michael Collins. 2019. Synthetic QA corpora generation with roundtrip consistency. In *ACL*.
- Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. 2016. Guided open vocabulary image captioning with constrained beam search. In *EMNLP*.
- Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-up and top-down attention for image captioning and visual question answering. In *CVPR*.
- Kristjan Arumae and Fei Liu. 2019. Guiding extractive summarization with question-answering rewards. In *NAACL*.
- Jeffrey P. Bigham, Chandrika Jayant, Hanjie Ji, Greg Little, Andrew Miller, Robert C. Miller, Robin Miller, Aubrey Tatarowicz, Brandyn White, Samuel White, and Tom Yeh. 2010. VizWiz: Nearly real-time answers to visual questions. In *23rd Annual ACM Symposium on User Interface Software and Technology*.
- Florian Böhm, Yang Gao, Christian M. Meyer, Ori Shapira, Ido Dagan, and Iryna Gurevych. 2019. Better rewards yield better summaries: Learning to summarise without references. In *EMNLP-IJCNLP*.
- Bo Dai, Sanja Fidler, Raquel Urtasun, and Dahua Lin. 2017. Towards diverse and natural image descriptions via a conditional gan. In *ICCV*.
- Bo Dai and Dahua Lin. 2017. Contrastive learning for image captioning. In *NeurIPS*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*.
- H. P. Edmundson. 1969. New methods in automatic extracting. *J. ACM*, 16(2):264–285.
- Matan Eyal, Tal Baumel, and Michael Elhadad. 2019. Question answering as an automatic evaluation metric for news article summarization. In *NAACL*.
- Yang Gao, Christian M. Meyer, Mohsen Mesgar, and Iryna Gurevych. 2019. Reward learning for efficient reinforcement learning in extractive document summarization. In *IJCAI*.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *CVPR*.
- Tszhang Guo, Shiyu Chang, Mo Yu, and Kun Bai. 2018. Improving reinforcement learning based image captioning with natural language prior. In *EMNLP*.
- Danna Gurari, Qing Li, Abigale Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey Bigham. 2018. Vizwiz grand challenge: Answering visual questions from blind people. In *CVPR*.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *NeurIPS*.
- Chris Hokamp and Qun Liu. 2017. Lexically constrained decoding for sequence generation using grid beam search. In *ACL*.
- Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. 2019. Attention on attention for image captioning. In *ICCV*.
- Drew A. Hudson and Christopher D. Manning. 2019. GQA: A new dataset for real-world visual reasoning and compositional question answering. *CVPR*.
- Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2019. Spanbert: Improving pre-training by representing and predicting spans. *arXiv preprint arXiv:1907.10529*.
- Andrej Karpathy and Fei-Fei Li. 2015. Deep visual-semantic alignments for generating image descriptions. In *CVPR*.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *ICLR*.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yanis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. 2016. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *arXiv preprint arXiv:1602.07332*.
- Sergey Levine and Vladlen Koltun. 2013. Guided policy search. In *ICML*.
- Jiwei Li, Alexander H. Miller, S. Chopra, Marc’Aurelio Ranzato, and Jason Weston. 2017. Learning through dialogue interactions by asking questions. In *ICLR*.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*.
- Qing Li, Jianlong Fu, Dongfei Yu, Tao Mei, and Jiebo Luo. 2018. Tell-and-answer: Towards explainable vi-

- sual question answering using attributes and captions. In *EMNLP*.
- Chin-Yew Lin and Eduard Hovy. 2003. Automatic evaluation of summaries using n-gram co-occurrence statistics. In *NAACL*.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft COCO: Common objects in context. In *ECCV*.
- Xihui Liu, Hongsheng Li, Jing Shao, Dapeng Chen, and Xiaogang Wang. 2018. Show, tell and discriminate: Image captioning by self-retrieval with partially labeled data. In *ECCV*.
- Edward Loper and Steven Bird. 2002. Nltk: The natural language toolkit. In *In Proceedings of the ACL Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational Linguistics. Philadelphia: Association for Computational Linguistics*.
- Annie Louis and Ani Nenkova. 2013. Automatically assessing machine summary content without a gold standard. *Computational Linguistics*, 39(2).
- Ruotian Luo, Brian Price, Scott Cohen, and Gregory Shakhnarovich. 2018. Discriminability objective for training descriptive captions. In *CVPR*.
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Ranking sentences for extractive summarization with reinforcement learning. In *NAACL*.
- Vicente Ordonez, Girish Kulkarni, and Tamara L. Berg. 2011. Im2text: Describing images using 1 million captioned photographs. In *NeurIPS*.
- Romain Paulus, Caiming Xiong, and Richard Socher. 2018. A deep reinforced model for abstractive summarization. In *ICLR*.
- Maxime Peyrard. 2019. A simple theoretical model of importance for summarization. In *ACL*.
- Maxime Peyrard and Iryna Gurevych. 2018. Objective function learning to match human judgements for optimization-based summarization. In *NAACL*.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for SQuAD. In *ACL*.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *EMNLP*.
- Marc'Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2016. Sequence level training with recurrent neural networks. In *ICLR*.
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NeurIPS*.
- Steven J. Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2017. Self-critical sequence training for image captioning. In *CVPR*.
- Thomas Scialom, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano. 2019. Answers unite! unsupervised metrics for reinforced summarization models. In *EMNLP-IJCNLP*.
- Tingke Shen, Amlan Kar, and Sanja Fidler. 2019. Learning to caption images through a lifetime by asking questions. In *ICCV*, pages 10392–10401.
- Karen Spärck Jones. 1994. Towards better nlp system evaluation. In *Human Language Technologies Workshop*.
- Karen Spärck Jones. 1999. Automatic summarizing: Factors and directions. In I. Mani and M. T. Maybury, editors, *Advances in Automatic Text Summarization*. MIT Press, Cambridge, MA.
- Hao Tan and Mohit Bansal. 2019. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NeurIPS*.
- Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. CIDEr: Consensus-based image description evaluation. In *CVPR*.
- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *CVPR*.
- Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. Asking and answering questions to evaluate the factual consistency of summaries. In *ACL*.
- Liwei Wang, Alexander Schwing, and Svetlana Lazebnik. 2017. Diverse and accurate image description using a variational auto-encoder with an additive gaussian encoding space. In *NeurIPS*.
- Ronald J. Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3-4).
- Jialin Wu, Zeyuan Hu, and Raymond Mooney. 2019. Generating question relevant captions to aid visual question answering. In *ACL*.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Xiaoshan Yang and Changsheng Xu. 2019. Image captioning by asking questions. *ACM Trans. Multimedia Comput. Commun. Appl.*, 15(2s).
- Luwei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason J. Corso, and Jianfeng Gao. 2019. Unified

vision-language pre-training for image captioning and vqa. *arXiv preprint arXiv:1909.11059*.

Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. 2016. Visual7W: Grounded question answering in images. In *CVPR*.

A Model Architecture Details

The captioning architecture we use is a standard Transformer sequence-to-sequence model. As the model is not the main focus, we did not do any extensive hyper-parameter tuning or ablations beyond ensuring that we had a reasonable baseline model on COCO (114 CIDEr on COCO captions).

Image Encoder: For each image \mathbf{x} , we take represent the region embeddings $\mathbf{o}_i \in \mathbb{R}^{2048}$ of the bounding boxes for the k most confident object detections. We use the pre-trained Faster R-CNN (Ren et al., 2015) model of Anderson et al. (2018).⁶ We then map each region embedding to $\tilde{\mathbf{o}}_i \in \mathbb{R}^d$ using a single dense layer with a ReLU. Inspired by the positional token embeddings in the BERT model (Devlin et al., 2019), we then augment $\tilde{\mathbf{o}}_i$ with learned position (the rasterized coordinate of the bounding box center), segment (a constant “image” component identifier), and confidence (the detection rank of the object) embeddings to obtain the full object representation:

$$\hat{\mathbf{o}}_i = \tilde{\mathbf{o}}_i + \mathbf{p}_i + \mathbf{s}_i + \mathbf{c}_i.$$

Text Decoder: We decode the caption autoregressively—starting with the [CLS] token and terminating on [SEP]. At each time-step t we concatenate the image embeddings with special delimiters and the word-pieces decoded thus far, to obtain a joint context:

$$\mathbf{h} = \{ [\text{IMG}], \hat{\mathbf{o}}_1, \dots, \hat{\mathbf{o}}_k, [\text{CLS}], \mathbf{w}_1, \dots, \mathbf{w}_t \},$$

where $\mathbf{w}_i \in \mathbb{R}^d$ is the word piece embedding (the sum of token, position, and segment embeddings). We then encode \mathbf{h} using multi-layer Transformer, and compute the probability of generating w_{t+1} using a softmax over the 30,522 word-piece vocabulary (the BERT vocabulary). For efficiency, we encode whole sequences at a time with a left-to-right attention mask: image regions may attend to all other image regions, and tokens may attend to all previous tokens and image regions.

⁶<https://github.com/peteanderson80/bottom-up-attention>

Hyperparameters: In our experiments we use the top 64 object regions and a 6-layer Transformer with 512 hidden input units, 8 attention heads, and 2,048 hidden units in the intermediate feed-forward layer. During inference we do beam search with a beam size of 3 and a length penalty α of 0.6 (Wu et al., 2016). We implemented our model in Tensorflow (Abadi et al., 2015).

B Model Training Details

QA Model Threshold: During inference, the QA model $\mathcal{M}(\mathbf{q}, \mathbf{y})$ computes the probability of the “no answer” option $p_{\mathcal{M}}(\text{NONE}|\mathbf{q}, \mathbf{y})$ and the probability of the most likely answer span $p_{\mathcal{M}}(\mathbf{y}_{i\dots j}|\mathbf{q}, \mathbf{y})$. We adjust how precise this reward is by treating the log odds ratio c of the “no answer” vs. span options as a hyper-parameter when choosing the prediction $\hat{\mathbf{a}}$:

$$\log \left(\frac{p_{\mathcal{M}}(\mathbf{y}_{i\dots j}|\mathbf{q}, \mathbf{y})}{p_{\mathcal{M}}(\text{NONE}|\mathbf{q}, \mathbf{y})} \right) \begin{cases} > c, & \hat{\mathbf{a}} = \mathbf{y}_{i\dots j} \\ \leq c, & \hat{\mathbf{a}} = \text{NONE} \end{cases}$$

Depending on the value of c , we may only answer if we are confident the answer is supported—not just the most probable (e.g., based on answer-type)—to avoid potentially spurious rewards obtained by guessing or elimination. Table 8 shows an ablation over some choices of c .

Answer Normalization: For both training and evaluation, we normalize the gold and predicted answers by removing articles and punctuation when comparing them (see Rajpurkar et al., 2016).

Policy Gradient: We approximate the policy gradient (Eq. 3) using a single Monte-Carlo sample $\mathbf{y} = (y_1, \dots, y_n)$ from $G_{\theta}(\mathbf{y}|\mathbf{x})$. We accelerate training by restricting samples to be from a set of high-probability candidates with non-zero reward (cf. Anderson et al., 2018; Narayan et al., 2018, *inter alia*). We decode using beam search and sample from the top- k beams ($k = 16$).

Training: For MLE pre-training on $\mathcal{D}_{\text{generic}}$ (Eq. 2), all examples are shuffled and divided into mini-batches of 256 examples each. For RL adaptation to $\mathcal{D}_{\text{target}}$ (Eq. 3), we use a mini-batch size of 128. To help regularize the fluency of the model, during RL training we continue to multi-task on the supervised generic captions, as in MIXER (Ranzato et al., 2016). For both settings, we train for a maximum of 120K steps and choose the best

Algorithm 1 Synthetic data generation procedure for policy pre-training.

Definitions: $\mathcal{D}_{\text{generic}}$ is assumed out-of-domain generic captioning data with input images $\tilde{\mathbf{x}}$ and supervised reference captions $\tilde{\mathbf{y}}$. $\mathcal{D}_{\text{target}}$ is the in-domain target QA data with input images \mathbf{x} , questions \mathbf{q} , and answers \mathbf{a} . \mathcal{M} is the automatic QA model used in this paper for evaluation (Table 1). $\mathcal{M}^{-1}(\mathbf{y})$ is a pre-trained QA generation model, that takes in some context \mathbf{y} and outputs a predicted QA pair $(\hat{\mathbf{q}}, \hat{\mathbf{a}})$.

```
1: function TRAIN( $\mathcal{D}_{\text{generic}}, \mathcal{M}^{-1}, T$ )
2:    $\phi \leftarrow$  random
3:   for  $i = 1$  to  $T$  do
4:      $\tilde{\mathbf{x}}, \tilde{\mathbf{y}} \sim \mathcal{D}_{\text{generic}}$ 
5:      $\hat{\mathbf{q}}, \hat{\mathbf{a}} \leftarrow \mathcal{M}^{-1}(\mathbf{y})$ 
6:      $\mathcal{L} \leftarrow -\log F_{\phi}(\mathbf{y}|\mathbf{x}, \mathbf{q}, \mathbf{a})$ 
7:      $\phi \leftarrow$  MINIMIZE( $\mathcal{L}, \phi$ )
8:   return  $F_{\phi}$ 
9:
10: function GENERATE( $\mathcal{D}_{\text{target}}, \mathcal{M}, F_{\phi}$ )
11:    $\mathcal{D}_{\text{synthetic}} \leftarrow []$ 
12:   for  $(\mathbf{x}, \mathbf{q}, \mathbf{a}) \in \mathcal{D}_{\text{target}}$  do
13:      $\tilde{\mathbf{y}} \leftarrow \operatorname{argmax} F_{\phi}(\mathbf{y}|\mathbf{x}, \mathbf{q}, \mathbf{a})$ 
14:     if  $\mathcal{M}(\mathbf{q}, \tilde{\mathbf{y}}) = \mathbf{a}$  then
15:       APPEND( $\mathcal{D}_{\text{synthetic}}, (\mathbf{x}, \tilde{\mathbf{y}})$ )
16:   return  $\mathcal{D}_{\text{synthetic}}$ 
```

model based on the dev set performance (using COCO CIDEr (Vedantam et al., 2015) for MLE pre-training and QA F1 for RL). For optimization, we use Adam (Kingma and Ba, 2015) with a linear warm-up and decay schedule. Training was performed on a 4×4 TPU, and took about 1-2 hours per experiment.

C Synthetic Pre-training Details

QA Conditional Model: We use the same basic architecture for $F_{\phi}(\mathbf{y}|\mathbf{x}, \mathbf{q}, \mathbf{a})$ as for the main captioning model $G_{\theta}(\mathbf{y}|\mathbf{x})$, and only introduce two new “question” segment and “answer” segment embeddings that we add to differentiate the conditional text from the generated text in the Transformer. The full input then becomes:

$$\mathbf{h} = \{ [\text{IMG}], \hat{\mathbf{o}}_1, \dots, \hat{\mathbf{o}}_k, [\text{Q}], \mathbf{q}_1, \dots, \mathbf{q}_m, [\text{A}], \mathbf{a}_1, \dots, \mathbf{a}_n, [\text{CLS}], \mathbf{w}_1, \dots, \mathbf{w}_t \},$$

where the segment delimiter, \mathbf{q}_i , and \mathbf{a}_j vectors are defined the same way as—and shared with—the caption’s input word-piece embeddings \mathbf{w}_i (see §A). We decode auto-regressively as before.

Training and Generating: We train $F_{\phi}(\mathbf{y}|\mathbf{x}, \mathbf{q}, \mathbf{a})$ using both the corpus of generic captions used

for MLE pre-training in Section 4.2 (i.e., COCO) and additional Wikipedia text. We create automatic $(\mathbf{x}, \mathbf{q}, \mathbf{a}, \mathbf{y})$ and $(\text{null}, \mathbf{q}, \mathbf{a}, \mathbf{y})$ examples for COCO and Wikipedia sentences, respectively. To offset biases present in the question generation model (which is out-of-domain for caption-styled text as it is trained on SQuAD), we add $(\mathbf{x}, \text{null}, \mathbf{a}, \mathbf{y})$ examples from the generic captions by selecting random spans of \mathbf{y} , using the sampler of Joshi et al. (2019) (random spans with Poisson distributed lengths). Algorithm 1 illustrates the full synthetic data generation process.

D Converted Dataset Details

Splits: For the COCO and VQA datasets we use the ‘Karpathy’ splits from Karpathy and Li (2015). For GQA, we use the ‘balanced’ splits, but limit to $\sim 5\text{K}$ images each for the new test and dev sets from the original GQA dev set. Both the Visual7W and GQA datasets have images from Visual Genome (Krishna et al., 2016), and thus some (partially) overlap with the ‘Karpathy’ COCO images. Since we use COCO for pre-training ($\mathcal{D}_{\text{generic}}$), we avoid data leakage by mapping Visual Genome IDs to COCO IDs and either filter questions about images that are in the COCO train or dev sets, or re-

	Dataset	Train		Development		Test	
		Images	QA	Images	QA	Images	QA
I	COCO	113,287	–	5,000	–	5,000	–
II	CapVQA	104,311	297,484	4,617	13,081	4,615	12,847
	CapGQA	69,450	611,102	5,000	43,015	4,739	41,398
	CapVisual7W	20,268	93,878	3,448	16,314	4,892	22,769
	CapVizWiz	10,027	10,027	960	960	1,905	1,905

Table D.1: Statistics of the datasets used in this paper. Type I: generic/no QA. Type II: target/QA.

assign the data to match COCO. Finally, on VizWiz, we combine all of the original data and randomly re-partition it into test and dev sets of $\sim 1\text{K}$ and $\sim 2\text{K}$ images each, keeping the rest for training.