

# Reassessing Evaluation Practices in Visual Question Answering: A Case Study on Out-of-Distribution Generalization

Aishwarya Agrawal<sup>\*,‡,◇</sup> Ivana Kajić<sup>\*,◇</sup> Emanuele Bugliarello<sup>\*,△</sup>  
Elnaz Davoodi<sup>†,◇</sup> Anita Gergely<sup>†,◇</sup> Phil Blunsom<sup>‡</sup> Aida Nematzadeh<sup>\*,‡,◇</sup>

◇DeepMind △University of Copenhagen ‡University of Oxford  
{aiagrawal,kivana,nematzadeh}@deepmind.com emanuele@di.ku.dk

## Abstract

Vision-and-language (V&L) models pre-trained on large-scale multimodal data have demonstrated strong performance on various tasks such as image captioning and visual question answering (VQA). The quality of such models is commonly assessed by measuring their performance on unseen data that typically comes from the same distribution as the training data. However, when evaluated under out-of-distribution (out-of-dataset) settings for VQA, we observe that these models exhibit poor generalization. We comprehensively evaluate two pretrained V&L models under different settings (i.e. classification and open-ended text generation) by conducting cross-dataset evaluations. We find that these models tend to learn to solve the benchmark, rather than learning the high-level skills required by the VQA task. We also find that in most cases generative models are less susceptible to shifts in data distribution compared to discriminative ones, and that multimodal pretraining is generally helpful for OOD generalization. Finally, we revisit assumptions underlying the use of automatic VQA evaluation metrics, and empirically show that their stringent nature repeatedly penalizes models for correct responses.

## 1 Introduction

Visual Question Answering (VQA) is the task of automatically answering natural language open-ended questions about images. Tackling VQA involves multiple skills, such as language and visual understanding, integrating information between the two (vision and language) modalities, and commonsense and knowledge based reasoning. One of the goals of the VQA research has been fostering the development of systems that are able to answer *any open-ended question about any image*. This

motivation has inspired a fruitful line of research in designing VQA benchmarks (e.g., Malinowski and Fritz, 2014; Antol et al., 2015; Krishna et al., 2017; Goyal et al., 2017; Hudson and Manning, 2019) and models (e.g., Yang et al., 2015; Anderson et al., 2018; Lu et al., 2019; Cho et al., 2021).

In this work, we investigate if recent pretrained VQA models can indeed answer any open-ended question about images or if they are mostly suitable for answering questions from the VQA benchmarks they are optimized for. In other words, *are models learning to solve the task or learning to solve the datasets?* We believe the former is more aligned with the goal of building real-world VQA systems.

To measure whether models learn to solve the task of VQA, we believe we need to examine their *out-of-distribution (OOD)* generalization capabilities: how they perform on examples drawn from a distribution other than that of the training set. In this work, we extensively evaluate OOD generalization of current pretrained V&L models by conducting cross-dataset evaluations (without any adaptation to the test domain).

Through our extensive experiments, we provide in-depth discussion on the following questions:

- *How well do recent models generalize under OOD settings?* We observe a notable drop in performance from IID to OOD settings across models and benchmarks, demonstrating that models mostly learn to solve specific benchmarks as opposed to learning general skills for answering questions about images. This result is not simply due to a mismatch between the set of answers between the training and test VQA datasets, nor due to poor representation of test answers in VQA training data.
- *Is multimodal pretraining beneficial for OOD generalization?* We find that while image-text pretraining is helpful in most OOD settings, it is not always more useful than in IID ones. Moreover, it is least useful for OOD evalua-

\*denotes equal first author contribution. † denotes equal contribution. ‡ denotes equal senior contribution. Detailed contributions follow at the end of the paper.

tion on the VIZWIZ benchmark, highlighting the challenges of a real-world benchmark.

- *Is generative modeling more robust to distribution shifts?* In most cases, we observe that generative models—which are not bound to predictions over a fixed set of answers curated from the training data—are more robust to OOD evaluation than discriminative (*i.e.*, classification-based) ones. Moreover, we quantify what the limitations of discriminative models are for real-world VQA applications (*e.g.*, answering questions of visually-impaired users), where the answers a deployed model needs to produce cannot be predetermined.
- *Are current automatic VQA metrics too stringent for OOD evaluation?* We examine if the performance of our pretrained models is negatively impacted by the current standard VQA accuracy metrics, which match predicted answer strings to a limited number of ground-truth answers. Human evaluation reveals the stringent nature of such accuracy metrics, which is especially pronounced in the OOD settings. Nevertheless, while the IID-to-OOO performance gap is reduced after human evaluation, models still exhibit poor generalization to OOD VQA benchmarks.

We believe our OOD evaluations and supporting analyses expose the shortcomings of current models, and recommend future work to adopt these evaluation practices to provide real-world, robust assessment of VQA systems.

## 2 Related Work

**Beyond IID evaluation in VQA.** Previous work has evaluated VQA models beyond the IID setting for robustness to *specific and controlled* aspects – novel compositions of seen concepts (Agrawal et al., 2017; Johnson et al., 2017; Hudson and Manning, 2019), change in prior distributions of answers per question type (Agrawal et al., 2018; Gokhale et al., 2020; Niu et al., 2021), adversarial examples provided by humans (Sheng et al., 2021; Li et al., 2021b), consistency, negation, and simple perturbation in questions (Jimenez et al., 2022), counter-examples (Dancette et al., 2021), and controlled shifts in language and vision modalities (Akula et al., 2021). Our focus, however, is to evaluate for *overall* robustness to OOD data without controlling for specific aspects, by testing our

models on different OOD benchmarks. We believe our experimental setting more closely emulates the expected experience of deployed VQA systems. Moreover, when the exact nature of distribution shift between train and test splits is known (such as in (Agrawal et al., 2018)), approaches developed to tackle such shifts tend to rely on the explicit knowledge of construction of such OOD splits resulting in inflated sense of progress (Teney et al., 2020).

Similar to us, Zhang et al. (2021); Hudson and Manning (2019) also present some experimental results on VQA OOD evaluation, however they do it in limited manner (*e.g.*, do not consider all pairs of datasets, do not evaluate the effect of multimodal pretraining, etc.). To our best knowledge, ours is the first work to extensively quantifying the extent of IID to OOD performance drops in current VQA models and study the effect of several factors: answer overlap, multimodal pretraining, generative vs. discriminative modeling, and stringent evaluation metric.

**Domain adaptation in VQA.** Some studies (Jabri et al., 2016; Chao et al., 2018; Zhang et al., 2021) have explored domain adaptation of VQA models from one VQA benchmark to another. Our focus, instead, is on evaluating *zero-shot* cross-benchmark generalization *without* any adaptation. This allows us to assess the robustness of current models towards unforeseen distribution shifts. Our work is similar to that of Torralba and Efros (2011) and Hendrycks et al. (2020), who study OOD generalization in vision and text.

**Zero-shot VQA with pretrained models.** In an emerging line of research (Tsimpoukelli et al., 2021; Alayrac et al., 2022; Song et al., 2022; Piergiovanni et al., 2022), large-scale pretrained unimodal models (Brown et al., 2020; Radford et al., 2021) are repurposed to tackle VQA in zero-shot or few-shot fashion. While such zero-shot VQA evaluations are a better test of generalization than IID evaluations, our focus, differently, is on investigating whether models can generalize to unseen datasets *upon being taught the task by showing examples from one dataset*. Moreover, this line of work does not focus on a thorough analysis of models in OOD settings (which is hard to define for these models due to the massive amount of data they are pretrained on).

### 3 Experimental Setup

In this section, we present our framework to examine OOD generalization in VQA. We examine two pretrained Transformers across five benchmarks.

#### 3.1 Models

We evaluate the performance of two representative, widely-used pretrained models that have achieved strong performance in various V&L tasks in the last few years: ViLBERT (Lu et al., 2019) and ALBEF (Li et al., 2021a). We evaluate these models in a broad range of settings (generative/discriminative, w/wo pretraining, and multiple benchmarks), resulting in 128 experiments. We chose these models as they include components shown to be important in the literature: cross-attention (ViLBERT and ALBEF), and contrastive learning (ALBEF). We note that our goal is to study trends that hold across different models, and we leave for future work controlled comparisons across architectures.

**ViLBERT** is one of the first, yet strong models in the recent pretrain–fine-tune paradigm in V&L. Its inputs are a sequence of sub-word tokens (Wu et al., 2016), and a set of regions of interest given by a Faster R-CNN (Ren et al., 2015; Anderson et al., 2018). The authors fine-tune it on VQAV2 by learning a classifier over the most frequent answers. We first re-implement this model successfully, and then extend it to a generative setting by pretraining and fine-tuning a Transformer decoder (more details in App. A). We denote the discriminative/generative version as ViLBERT<sub>DISC</sub>/ViLBERT<sub>GEN</sub>. Unless otherwise specified, results for ViLBERT<sub>DISC</sub> are from our code base for direct comparison with ViLBERT<sub>GEN</sub>.

**ALBEF** is a state-of-the-art V&L encoder whose visual inputs are image patches encoded by a vision Transformer (Dosovitskiy et al., 2021; Touvron et al., 2021) that is jointly trained with the rest of the model. Li et al. (2021a) fine-tune ALBEF on VQAV2 by adding a 6-layer Transformer decoder to generate answers (ALBEF<sub>GEN</sub>). We use the official implementation,<sup>1</sup> and furthermore train a discriminative variant (ALBEF<sub>DISC</sub>) by learning a multi-answer classifier, as in ViLBERT<sub>DISC</sub>.

In our analysis, we investigate the role of multi-modal pretraining. ViLBERT was pretrained on 3M image–text pairs from Conceptual Captions (CC; Sharma et al. 2018). Li et al. (2021a) re-

<sup>1</sup><https://github.com/salesforce/ALBEF>.

Dataset	# Train (imgs / qns)	# Val (imgs / qns)	# Classes	Coverage [%]
VQAV2	82,783 / 443,757	40,504 / 214,354	3,129	98.07 / 98.07
GQA	72,140 / 943,000	10,234 / 132,062	1,533	99.78 / 99.79
VG	59,635 / 868,259	39,645 / 577,063	3,449	76.55 / 76.55
VIZWIZ	20,523 / 20,523	4,319 / 4,319	3,112	96.76 / 97.01

Table 1: Datasets statistics. #classes is the number of classes we use for the discriminative models; coverage is the percentage of questions that can be answered with our selected classes in train/validation splits.

leased two checkpoints for ALBEF: one pretrained on 4M images from CC, MS-COCO (Lin et al., 2014), SBU (Ordonez et al., 2011) and Visual Genome (Krishna et al., 2017); the other one is further pretrained on Conceptual 12M (Changpinyo et al., 2021) for a total of 14M images.<sup>2</sup>

#### 3.2 Datasets and Evaluation Metrics

**Datasets.** We ground our analysis on five diverse VQA datasets: VQAV2 (Goyal et al., 2017), GQA (Hudson and Manning, 2019), VISUAL GENOME (VG; Krishna et al. 2017), VIZWIZ (Gurari et al., 2018) and VQA-CP (Agrawal et al., 2018). VQAV2 is the most commonly used VQA dataset to date. VQA-CP re-splits it such that, for every question type, train and test sets have different prior distributions of answers. VG includes questions centered around either the full image or a specific region. GQA is a large-scale dataset that focuses on compositionality of template-generated questions. Finally, VIZWIZ is the only real-world VQA dataset, collected from visually-impaired people. VG and GQA have one answer per question, while the other datasets include 10 answers per question. See Tab. 1 and App. A for more details.

There are several differences among these datasets. Both VQAV2 and GQA mostly have one-word answers (89% and 81%, respectively) whilst there are fewer in VG (57%) and VIZWIZ (67%). The type of questions also varies: VG does not contain binary ‘yes/no’ questions, but rather spans 6 WH-questions. By design, GQA questions require more compositional skills but do not test for counting; while VIZWIZ questions are more conversational as they were collected through a speech interface and has a significant proportion of OCR questions (21%). Moreover, a significant number of VIZWIZ questions (28%) are *unanswerable* due to the challenges faced by the visually-impaired users in taking pictures, resulting in poor focus,

<sup>2</sup>We also conducted experiments with ViLBERT pretrained on same datasets as the 4M ALBEF checkpoint. We found no significant difference compared to the results presented throughout this paper.

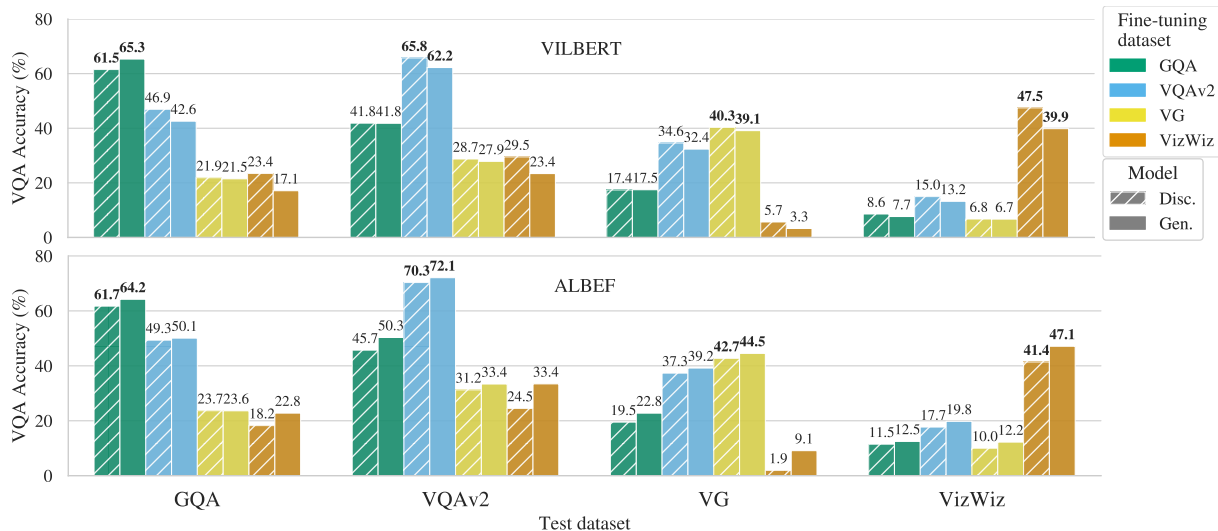


Figure 1: IID (highlighted in bold) vs. OOD performance. Top: ViLBERT pretrained on CC. Bottom: ALBEF pretrained on CC, VG, SBU, MS-COCO and C12M datasets. All models are initialized with BERT weights.

poor lighting or entirely missing the entity of interest. As such, the distribution of images in VIZWIZ is different from other datasets.

**Evaluation metrics.** These VQA benchmarks compute model accuracy between its prediction and the ground-truth answer(s) by string matching (after simple pre-processing). VQAV2 and VIZWIZ, each with 10 answers per question, account for diversity in ground-truth answers by scoring a given model answer as  $\min\{1.0, 0.3 \times \text{count}\}$ , where count is the number of annotators that used that answer. For GQA and VG, both with one answer per question, we use top-1 accuracy.<sup>3</sup>

### 3.3 Training Details

Following common practice, for discriminative models, we select the top- $k$  most frequent answers as the set of answer classes to perform classification over. Here  $k$  is a dataset-dependent variable, chosen to cover most of the questions (see Tab. 1). All models are trained on the respective training sets and evaluated on the validation sets. For VG, we randomly split the data into training and validation (60%/40%) with no image is in both splits.

## 4 Out-of-Distribution Generalization

We examine to what extent our models learn to solve a specific VQA benchmark by latching on dataset-specific correlations, as opposed to learning

<sup>3</sup>We note that GQA and VG propose top-5 accuracy. We, instead, opt for top-1 accuracy to keep a consistent setup with VQAV2 and VIZWIZ. And we believe top-5 accuracy is impractical for many applications, such as answering questions for visually-impaired users.

more general skills required in VQA. We fine-tune a pretrained model on the train split of one benchmark (e.g., GQA) and evaluate it on the validation split of a different one (e.g., VG). Overall, we evaluate models by fine-tuning them on each benchmark and testing them against all benchmarks. If pretrained models are indeed learning the VQA skill, we expect to see a small drop in performance between the IID and OOD settings.

The results are presented in Fig. 1, with different evaluation benchmarks grouped on the  $x$ -axis. First, across all models and for each benchmark, we see a notable drop in the VQA accuracy from the IID to the OOD setting. While such a drop might be anticipated, we found the extent of the drop surprising given the impressive performance of current pretrained VL models. For all models shown, the largest drops are observed when evaluating models on the VIZWIZ benchmark. Moreover, even the smallest performance drop, which happens when fine-tuning models on VQAV2 and evaluating them on VG, remains relatively large (i.e., 5.3 points for ALBEF<sub>GEN</sub>). These results show that *pretrained models are largely learning the fine-tuning benchmark without learning to solve the VQA task.*

Second, we observe that fine-tuning on VQAV2 results in the lowest drop in IID to OOD performance across all conditions—the VQAV2 bar (shown in blue in Fig. 1) is the closest to the IID one for GQA, VG, and VIZWIZ. We conclude that fine-tuning on VQAV2 yields a model that best generalizes to OOD settings in our benchmarks. This result is not simply due to the size of the fine-tuning benchmark as VG is larger than VQAV2.

Similarly, all the models achieve highest OOD performance on VQAV2. We conjecture that VQAV2 is the most diverse benchmark of our selection.

#### 4.1 Evaluating on Shared Answer Sets

Discriminative models treat VQA as a multi-answer classification task over the set of top- $k$  most frequent answers in the fine-tuning data. This limits their performance: if a certain answer is not frequent in the fine-tuning data, a discriminative model will perform poorly for such an answer during test time. While this limitation also affects IID evaluation, we expect it to have a stronger effect in OOD generalization (due to potentially different answer distributions between the fine-tuning and test sets). We next examine to what extent this limitation affects OOD performance by controlling for the mismatch in answer sets between the fine-tuning and test sets. We do so by considering only the test questions whose answers are included in the top- $k$  answers of a given fine-tuning dataset (for more details, see App. B).

Fig. 2 shows the improvement in the VQA accuracy over the IID and OOD evaluation accuracy (in Fig. 1) when controlling for the shared answer set. For IID evaluation, only one intersection of answer sets is reported, corresponding to the smallest gap between IID and OOD evaluation, with remaining numbers reported in Tab. 10 (App. B). Thus, the difference between the height of the IID bar (#) and the OOD bar (\*) with respect to which answer intersection between IID and OOD is computed, represents the best case scenario for OOD generalization, *i.e.*, the least drop from IID to OOD.

We observe a similar pattern across the models: in most cases, using a shared answer set improves the performance. *Overall, we still observe a notable gap between the OOD and IID settings for the best case OOD generalization scenario, showing that a shared answer set does not circumvent the difficulty of OOD generalization for these models.* A few cases where IID evaluations with a shared answer set hurt performance are discussed in App. B. When evaluating on the shared answer set, we further examine if the drop in accuracy from IID to OOD is due to the low frequency of the test answer classes in the OOD fine-tuning set. The details of the correlation computation and the results are explained in App. B and Tab. 9, respectively. This result indicates that frequency of the answer class is a contributing factor to the weak OOD generalization, but we also explore other causes in Sec. 7.

	VQAV2	GQA	VG	VIZWIZ
VQAV2	92.9	96.7	65.1	43.6
GQA	73.5	99.9	44.8	36.6
VG	52.7	62.4	74.2	32.3
VIZWIZ	79.4	82.5	40.9	86.2

Table 2: Maximum achievable accuracy for all test answers based on the top- $k$  answers present in the respective fine-tuning sets. Rows correspond to fine-tuning datasets, columns correspond to the test benchmarks

#### 4.2 The Case for the Generative Evaluation

A discriminative model cannot correctly answer questions for which the answers lie outside the predefined top- $k$  classes; therefore, by treating VQA as a classification task, we can define the upper-bound performance of discriminative models on VQA by computing the accuracy given all answers in the test set being answered correctly. The upper-bound VQA accuracy is shown in Tab. 2; we observe a large drop from IID to OOD evaluations for most conditions. VIZWIZ has the lowest achievable accuracies in OOD evaluation.

However, our ALBEF<sub>DISC</sub> and ViLBERT<sub>DISC</sub> models still perform notably worse than maximum achievable accuracy in all settings (smallest gap of 21.5% across all conditions, see Fig. 7 in App. B); *as a result, the poor OOD performance in the discriminative setting is not simply due to the low maximum achievable accuracy.* We conclude that the common practice of modeling VQA as a classification task severely limits the generalization capability of models to new datasets. On the other hand, generative models do not suffer from a fixed class set. They can generate a larger set of answers—all words for which the tokens occur in the pretraining data, including those that are out-of-vocabulary for the given VQA fine-tune datasets. We argue that generative modeling is a more promising solution for real-world application of VQA; similarly, recent work has identified text generation as a way to unify various V&L tasks (*e.g.*, Cho et al., 2021; Wang et al., 2022; Alayrac et al., 2022).

We next ask *whether our ViLBERT<sub>GEN</sub> and ALBEF<sub>GEN</sub> models are more successful in OOD generalization compared to their discriminative counterparts.* For each model (*i.e.*, generative/discriminative ALBEF/ViLBERT), we first calculate the gap between the IID setting and each OOD setting (*i.e.*,  $\Delta$  OOD), resulting in three values per benchmark. For instance, for the VQAV2 benchmark,  $\Delta$  OOD numbers are calculated between the model fine-tuned on VQAV2 and those

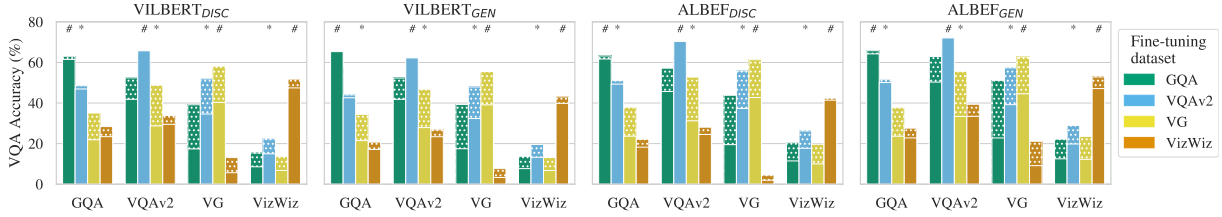


Figure 2: IID (#) vs OOD performance when controlling for the shared answer set. Solid bars are as in Fig. 1; stacked dotted bars are improvements when evaluating on questions with shared answer sets between IID and OOD settings. For IID, the shared answer set is computed with respect to a dataset denoted with \*.

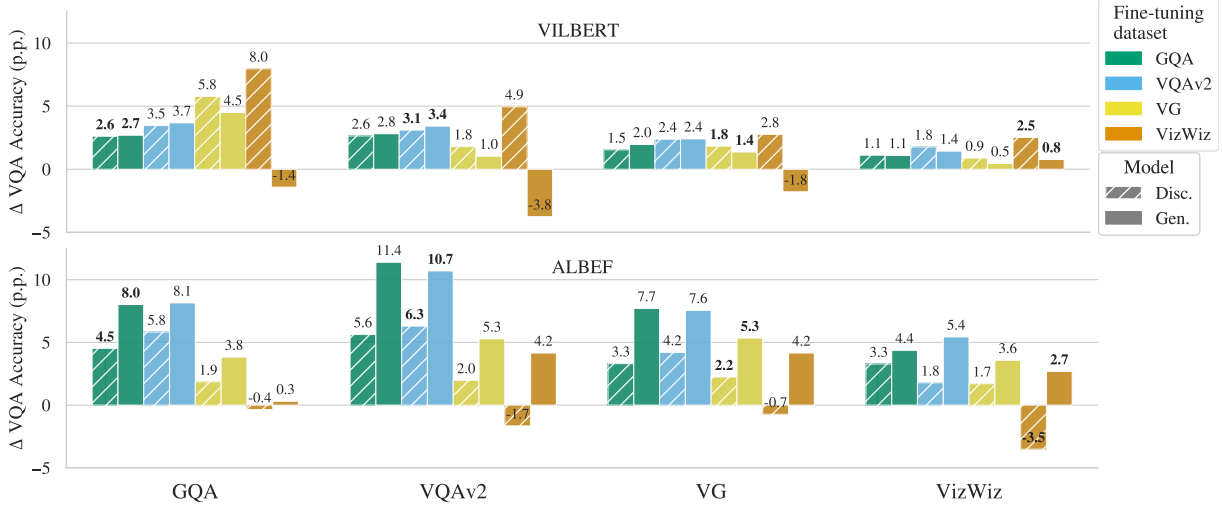


Figure 3: Percentage point difference in VQA accuracy between models with and without multimodal pretraining, for OOD and IID (highlighted in bold) evaluations. All models are initialized with BERT weights.

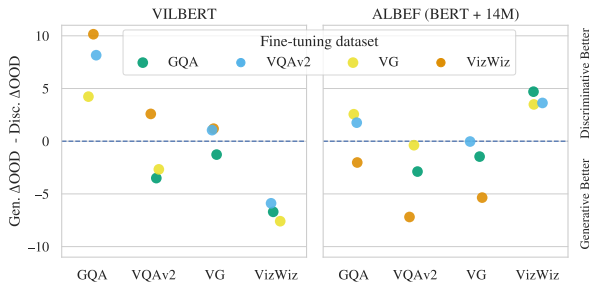


Figure 4: Difference in  $\Delta$  OOD values between discriminative and generative models.

fine-tuned on VG, GQA, and VIZWIZ. Note that the higher the  $\Delta$  OOD value, the poorer a model is in OOD generalization. We then compute the difference between the  $\Delta$  OOD values of the generative and discriminative models. Fig. 4 visualizes this result; the benchmarks are shown on the  $x$ -axis and each circle represents the difference in  $\Delta$  OOD values between the generative and discriminative model for a given fine-tuning dataset. If a generative model is more robust to OOD evaluation, we expect to see smaller  $\Delta$  OOD value for that model compared to its discriminative counterpart: when the circles are below the  $x$ -axis (depicting negative values), the generative model is more robust than the discriminative one. We observe ALBEF<sub>GEN</sub>

models often outperform their discriminative counterparts with respect to OOD generalization.

## 5 The Effect of Multimodal Pretraining

Previous work has shown that pretraining on multimodal (*i.e.*, image–text) data improves IID performance (*e.g.*, Lu et al., 2019; Li et al., 2021a); here, we ask if multimodal pretraining can help in OOD settings as well. We repeat the experiments in Sec. 4 without pretraining our models on multimodal data; instead we train the models on the train split of one benchmark and test it on the validation split of another. Fig. 3 shows the difference between the VQA accuracy of models with and without multimodal pretraining: each bar shows the gap between a bar in Fig. 1 and the equivalent experiment without multimodal pretraining.

We observe that multimodal pretraining is helpful in almost all conditions, since the majority of values displayed in Fig. 3 are positive. Pretraining is improving OOD performance likely because it can reduce the gap between the train and OOD test data by potentially exposing the model to a more diverse set of data points during pretraining. In our

Model	MM PT	VQAv2	VQA-CP	Drop
CF-VQA	–	53.6	63.5	9.9
VILBERT <sub>DISC</sub>	no	66.7	42.5	24.2
VILBERT <sub>DISC</sub>	yes	67.0	42.9	24.1
ALBEF <sub>DISC</sub>	no	64.0	40.1	23.9
ALBEF <sub>DISC</sub>	yes (4M)	70.0	44.4	25.6
ALBEF <sub>DISC</sub>	yes (14M)	70.3	45.2	25.1
ALBEF <sub>GEN</sub>	no	61.4	36.6	24.8
ALBEF <sub>GEN</sub>	yes (4M)	71.0	49.2	21.8
ALBEF <sub>GEN</sub>	yes (14M)	72.1	49.6	22.5

Table 3: Performance of models on VQAv2 (IID) and VQA-CP (OOD). The last column shows drop in performance from VQAv2 to VQA-CP. MM PT: Multimodal Pretraining.

experiments, the maximum gain from multimodal pretraining is indeed observed in OOD settings for both VILBERT (fine-tune on VIZWIZ; test on GQA) and ALBEF (fine-tune on GQA; test on VQAv2); however, *multimodal pretraining is not always more useful in OOD settings compared to IID ones*. For example, when evaluating VILBERT on VQAv2, pretraining helps the IID setting more than some of the OOD ones. Lastly, multimodal pretraining is detrimental for some cases where models are fine-tuned on VIZWIZ.

We observe that multimodal pretraining is more effective for the generative ALBEF compared to the discriminative ALBEF (cf. the shaded and solid bar with the same color in Fig. 3 bottom). For the VILBERT model, we generally do not observe such a pattern—discriminative and generative models mostly show comparable improvements due to multimodal pretraining. We observe only small improvements when increasing the size of the multimodal pretraining dataset for the ALBEF model (see Fig. 8 in App. B for more details).

## 6 Evaluation on VQA-CP

In this section, we evaluate the models<sup>4</sup> on the VQA under Changing Priors dataset (VQA-CP; Agrawal et al. 2018). This dataset is designed such that, for every question type, train and test splits have different prior distributions of answers. Thus, models that overfit to answer priors in training data and lack sufficient visual grounding show poor generalization on the VQA-CP test set. For comparison, we also evaluate the performance of Counterfactual VQA (CF-VQA; Niu et al. 2021), a state-of-art method on VQA-CP, which does not use either the Transformer architecture nor multi-

<sup>4</sup>ALBEF and VILBERT<sub>DISC</sub> (using the official codebase).

modal pretraining. However, it explicitly tackles the language (*i.e.*, question and answer) biases in VQA-CP.

Tab. 3 shows that for all the Transformer-based models, there is a large drop in the performance (at least 22%) from VQAv2 to VQA-CP. Thus, in spite of advances in the Transformer architecture and pretraining on diverse datasets, models are still overfitting to answer priors in the training data and lack sufficient visual grounding (Agrawal et al., 2018). However, the drop is much less for CF-VQA (10%), suggesting that *incorporating inductive biases specific to the generalization problem (i.e., modeling language bias) helps more than the Transformer architecture or scaling up the amount of pretraining data*. We also observe that the drop from VQAv2 to VQA-CP is often lower for the generative ALBEF than the discriminative ALBEF (except for ALBEF without any multimodal pretraining). Thus, *generative models are more robust than the discriminative ones*, especially when they are pretrained (similarly to the observations made in Sec. 4.2). As for the effect of pretraining, for generative ALBEF, pretraining helps reduce the drop from VQAv2 to VQA-CP. However, for discriminative models, pretraining does not seem to help generalization (in fact, it worsen ALBEF).

## 7 Qualitative Analysis

To dig deeper into the potential causes of the poor OOD generalization of our pretrained models, we perform a qualitative study. To this end, we randomly sample and manually examine failure cases in top-30 answer classes with the highest performance drop when moving from IID to OOD evaluation. We only focus on answer classes that are present in both the train and test splits, ensuring that performance drop is not due to the absence of answer classes in the training set. We report the top-5 classes that contribute the most to the drop in performance for each OOD setting in Tab. 11 (App. C). We notice that the following answer classes appear frequently across different OOD settings: yes/no answers, directions (left/right), colors, and numbers. In the following, we discuss a few major potential causes for the poor OOD generalization, and mention VILBERT<sub>DISC</sub> responses as examples in the discussion, although similar observations hold for other models.

**Poor reasoning skills.** Models evaluated on GQA, but fine-tuned on another dataset, show the

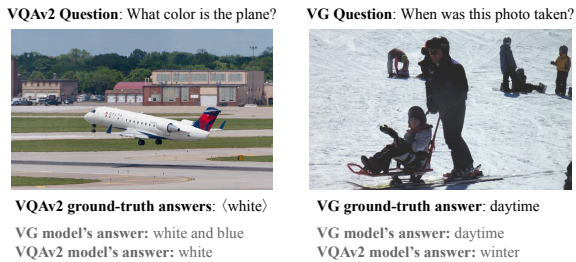


Figure 5: Examples where models’ prediction are correct but not accounted for in the ground-truth set. ⟨ ⟩ denotes a list of unique (out of 10) ground-truth answers. VG (VQAV2) model refers to a ViLBERT<sub>DISC</sub> fine-tuned on VG (VQAV2).

highest performance drop on classes such as “yes”, “no”, “right”, “left”, “top”, and “bottom”. For instance, ViLBERT<sub>DISC</sub> fine-tuned on VQAV2, and evaluated on GQA underperforms ViLBERT<sub>DISC</sub> that has been both fine-tuned and evaluated on GQA by 24% for the answer class “no.” Upon qualitative examination, we find that for many of such failure cases, the GQA questions are more compositional and hence require more complex reasoning (*e.g.*, “Are there both bison and zebras in the image?”, “Is the cheese to the right or to the left of the empty plate?”) than the questions for the same answer classes in other datasets (*e.g.*, from VQAV2 train set: “Is the TV turned on?”, “Which hand is the man holding up?”). This study re-affirms previous findings that VQA models lack sufficient logical, spatial, and compositional reasoning skills (Johnson et al., 2017; Hudson and Manning, 2019) but for the more recent, pretrained Transformer models.

**Overfitting to the answer priors.** Previous studies have shown that VQA models tend to be biased towards the prior distribution of answers in the training set (per question type) (Agrawal et al., 2018). We find that this limitation exists in the more recent pretrained models as well, and it is especially hurtful in the OOD settings because the priors need not be the same across train and test sets, unlike in the IID settings. For instance, ViLBERT<sub>DISC</sub> fine-tuned on VQAV2 predicts “2” for a lot of questions with target answer “1” in the VG test set. Similarly, sometimes ViLBERT<sub>DISC</sub> fine-tuned on VG incorrectly predicts “helmet” for VQAV2 test questions such as “What is the skateboarder wearing to protect his head?”, “What protective gear is he wearing?” when the skateboarder is not wearing anything. This indicates that the model is relying on answer priors rather than visual grounding. Our experimental results on VQA-CP

(Sec. 6) directly quantify the extent of such limitations in current models.

**Overfitting to the question format.** We observe instances of models failing to correctly answer questions when the format of the questions changes between the fine-tuning and test sets. For instance, questions about “chair” in the VQAV2 fine-tuning set are mostly of the form “What is ... sitting on?” whereas in the GQA test set, they are mostly of the form “What kind of furniture is ...?”. Thus, the “chair” class accuracy of ViLBERT<sub>DISC</sub> fine-tuned on VQAV2 drops from 48% when tested on VQAV2 to 38% on the GQA test set. Similarly, ViLBERT<sub>DISC</sub> fine-tuned on GQA fails terribly for “dog” and “cat” classes on the VG test set (accuracy drops of 47% and 43% respectively between GQA–GQA (fine-tuned on GQA, tested on GQA) and GQA–VG). GQA questions are mostly of the form “What animal ...?” or “What kind of animal ...?” whereas VG questions often do not mention the word “animal” and are of the form “Who is ...?” or “What is ...?” (*e.g.*, “Who is holding the Frisbee?”, “What is on the leash?”). To the best of our knowledge, no previous work has reported such behavior of VQA models (*i.e.*, they tend to overfit to the question format).

Finally, we observe cases where correct model responses are evaluated as incorrect by the VQA evaluation metric, as such responses differ from the ground-truth answers. In the next section, we provide examples of such cases and examine the impact of **stringent evaluation metric** on poor OOD generalization by engaging human raters to evaluate responses.

## 8 Human Evaluation

In our qualitative study, we observed that the stringent nature of the standard VQA evaluation metrics (*i.e.*, performing string matching of model responses with a small set of ground-truth answers) repeatedly penalizes models for correct responses because those responses do not exist in the set of ground-truth answers (Fig. 5). For example, the evaluation metric fails to take into account differences (between model response and ground-truth) due to specificity of the answers (*e.g.*, “on table” vs. “table”, “pizza slices” vs. “pizza”), synonyms, and different interpretations of the question (*e.g.*, Fig. 5 right).

In this section, we aim to quantify how robust standard VQA metrics are by performing human



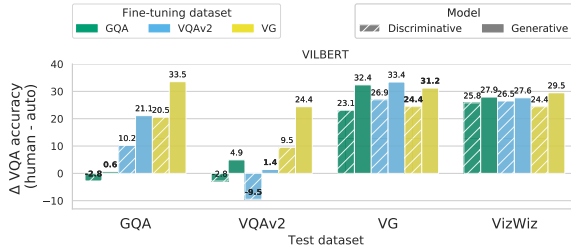


Figure 6: Difference in human and automatic accuracy of  $VILBERT_{DISC}$  (shaded bars) and  $VILBERT_{GEN}$  (plain bars) for GQA, VQAV2, VG and VIZWIZ. Accuracies in bold denote the IID settings.

evaluation of our models for both IID and OOD settings. The details of the setup and in-depth results are provided in App. D. Below we present our main findings.

**Human evaluation yields notably higher accuracies than the automatic evaluation.** This is shown in Fig. 6, where the increase can be up to 33.5% when moving from automatic to human evaluation.<sup>5</sup> This implies the current automatic metrics miss out on a lot of correct responses due to their stringent nature. Interestingly, this increase in model accuracy from automatic to human evaluation is higher for  $VILBERT_{GEN}$  than  $VILBERT_{DISC}$  for all the benchmarks. This is expected because the generative model is more likely to produce longer, more varied answers, which might not be awarded using automatic metric but are still correct responses. Moreover, human evaluation helps OOD settings more than the IID settings for most of the benchmarks (*e.g.*, GQA, VQAV2). This is also expected, because in the OOD settings, a model might not learn the format of the test answer (“on table” vs. “table”, “clear vs. sunny”) from the train set (unlike in the IID settings) and hence it is more likely to be penalized by the automatic accuracy metric. Thus, we conclude that the currently used accuracy metrics for VQA are not robust, especially for generative models and OOD evaluation settings. Hence, to more accurately evaluate the goodness of our models, *we need to develop better evaluation metrics for VQA.*

**Even after human evaluation, models still exhibit poor OOD generalization.** Although human evaluation improves the models’ accuracies and more so for the OOD than the IID settings, we observe that the models’ performance in OOD settings is still worse compared to that of IID set-

<sup>5</sup>In some cases, human evaluation yields lower accuracy than the automatic evaluation. We discuss this in App. D.

tings, albeit with reduced margin (see App. D for quantitative results). We also note that while  $VILBERT_{DISC}$  usually outperformed  $VILBERT_{GEN}$  with the automatic evaluation,  $VILBERT_{GEN}$  outperforms  $VILBERT_{DISC}$  for all the test sets under human evaluation. This reinforces the observations in Sec. 4.2 regarding stronger OOD generalization of generative models over discriminative ones.

## 9 Conclusion

In this study, we show that, despite their impressive performance when evaluated on test data drawn from the same distribution as the training data, recent V&L models perform poorly in out-of-distribution (OOD) settings. We conclude that these models learn to solve specific benchmarks as opposed to the skill of visual question answering (VQA). Interestingly, in most cases, we observe that the generative models are more robust to OOD generalization compared to the discriminative ones. Moreover, pretraining the models on large image-text data often helps in OOD generalization. Our results also highlight the importance of human evaluation for a more accurate assessment of model performance: we find that the current VQA automatic metrics miss out on a notable number of correct model responses. Human evaluation is especially important as the community is shifting towards generative VQA models which, unlike discriminative ones, can produce answers that go beyond those seen in a training/fine-tuning dataset. Finally, to make progress towards more capable models, we need more rigorous evaluation protocols that shed light on models’ strengths and short-comings. We believe testing models in OOD settings is a step towards this direction as it helps evaluate models for general skills required to solve the task as opposed to benchmark-specific correlations.

## Limitations

We list some limitations of our work which could benefit from future investigations.

First, when exploring potential factors for poor OOD generalization, our quantitative analysis focused only on differences in answer distributions between fine-tuning and test datasets. However, future work should investigate differences in question distribution, image distribution and combinations of these three variables.

Second, it would be interesting to conduct further investigation to understand why multimodal

pretraining does not help in certain cases. A correlation analysis between improvement in accuracy (due to multimodal pretraining), and between pretraining and fine-tuning/test data could be useful.

Third, the models investigated in our study (ViLBERT and ALBEF) are pretrained on a relatively small number (millions) of data points compared to language-only pretrained transformers, such as BERT, trained on billions of tokens. Such large-scale pretraining has been shown to improve OOD robustness for language-only models (Hendrycks et al., 2020). Hence, we leave for future work to investigate multimodal models trained on billions of image–text pairs (for instance, LAION-5B; Schuhmann et al. 2021).

Lastly, in this study, we only focus on standard VQA evaluation metrics for each benchmark. However, it would be interesting to also evaluate the robustness of metrics such as WUPS (Malinowski and Fritz, 2014) that compute answer similarities based on the distance between them in the WordNet (Miller, 1995) tree and are expected to be more robust than the standard metrics.

## Ethics Statement

Below we present some considerations related to the ethical and broader impact of our work.

First, all datasets used in our study are from published work and are publicly available, including the VIZWIZ data (Gurari et al., 2018) which has been curated from visually impaired users and released publicly after proper filtering to preserve the privacy of the users.

Second, for human evaluation of our models, we collected human data via the Amazon Mechanical Turk platform. We detail the data collection process and measures taken to control the quality of collected data in App. D. As for the ethical considerations related to collecting data from human subjects, our data collection campaign was approved by an ethics review board in our institution. Human subjects were paid at the rate of 0.15 USD per HIT (Human Intelligence Task) resulting in an hourly payment well above minimum wage.

Third, by testing models on a data distribution different from the training one, the OOD evaluation setting studied in our work has the following broader impacts: it highlights (1) the challenges of generalizing to real-world VQA datasets such as VIZWIZ, and (2) the kind of biases learned (and also potentially amplified) by the models.

Lastly, we discuss both potentially beneficial and harmful applications of the task of Visual Question Answering studied in our work. VQA has many potential applications beneficial for society:

- Aiding visually impaired users in understanding their surroundings (Human: What is on the shelf above the microwave? AI: Canned containers.)
- Teaching children through interactive demos (Kid: What animal is that? AI: That is Dall Sheep. You can find those in Alaska.)
- Aiding analysts in processing large quantities of visual surveillance data (Analyst: What kind of car did the man in red shirt leave in? AI: Blue Toyota Prius.)
- Interacting with in-home physical robots (Human: Is my laptop in my bedroom upstairs? AI: Yes. Human: Is the charger plugged in?)
- Making visual social media content more accessible (AI: Your friend Bob just uploaded a picture from his Hawaii trip. Human: Great, is he at the beach? AI: No, on a mountain.)

But like most other technology, VQA could also be used for potentially harmful applications such as:

- Invasion of individual’s privacy by using VQA to query streams of video data being recorded by CCTV cameras at public places.
- Visually impaired users often need assistance with parsing data containing personal information (Ahmed et al., 2015), such as credit cards, personal mails, etc. Such VQA systems could be configured to leak/retain personally identifiable information.

## Contributions

*Aishwarya Agrawal* initiated and designed the project, ran experiments and analyses on the official codebase of ViLBERT, contributed significantly to paper writing, and provided project support and advice. *Ivana Kajić* was responsible for the project’s technical infrastructure for the re-implementation of ViLBERT, ran experiments and analyses, and contributed significantly to paper writing. *Emanuele Bugliarello* co-led the data preparation, ran experiments and analyses on ALBEF, and contributed significantly to paper writing.

*Elnaz Davoodi* co-lead the data preparation, and helped setting up and running re-implementation experiments. *Anita Gergely* led the human evaluation experiments, and contributed to paper writing. *Phil Blunsom* provided project advice. *Aida Nematzadeh* provided significant project support and advice, helped running experiments, and contributed significantly to paper writing.

## Acknowledgements

We are grateful to Antoine Miech, Lisa Anne Hendricks and Chris Dyer for their constructive feedback. ■ During this project, *Emanuele Bugliarello* was supported by the funding from the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 801199.

## References

- Abien Fred Agarap. 2018. [Deep learning using rectified linear units \(relu\)](#). *arXiv preprint arXiv:1803.08375*.
- Aishwarya Agrawal, Dhruv Batra, Devi Parikh, and Aniruddha Kembhavi. 2018. [Don’t just assume; look and answer: Overcoming priors for visual question answering](#). In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Aishwarya Agrawal, Aniruddha Kembhavi, Dhruv Batra, and Devi Parikh. 2017. [C-VQA: A compositional split of the visual question answering \(vqa\) v1.0 dataset](#). *ArXiv*, abs/1704.08243.
- Tousif Ahmed, Roberto Hoyle, Kay Connelly, David Crandall, and Apu Kapadia. 2015. [Privacy concerns and behaviors of people with visual impairments](#). In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, pages 3523–3532.
- Arjun Akula, Soravit Changpinyo, Boqing Gong, Piyush Sharma, Song-Chun Zhu, and Radu Soricut. 2021. [CrossVQA: Scalably generating benchmarks for systematically testing VQA generalization](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2148–2166, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Saahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. 2022. [Flamingo: a visual language model for few-shot learning](#). In *Advances in Neural Information Processing Systems*.
- 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. 2015. [VQA: Visual question answering](#). In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Emanuele Bugliarello, Ryan Cotterell, Naoaki Okazaki, and Desmond Elliott. 2021. [Multimodal pretraining unmasked: A meta-analysis and a unified framework of vision-and-language BERTs](#). *Transactions of the Association for Computational Linguistics*, 9:978–994.
- Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. 2021. [Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3558–3568.
- Wei-Lun Chao, Hexiang Hu, and Fei Sha. 2018. [Cross-dataset adaptation for visual question answering](#). In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. [UNITER: Universal image-text representation learning](#). In *European Conference on Computer Vision (ECCV)*.
- Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. 2021. [Unifying vision-and-language tasks via text generation](#). In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 1931–1942. PMLR.

- Corentin Dancette, Rémi Cadène, Damien Teney, and Matthieu Cord. 2021. [Beyond question-based biases: Assessing multimodal shortcut learning in visual question answering](#). In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1574–1583.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. [Imagenet: A large-scale hierarchical image database](#). In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. [An image is worth 16x16 words: Transformers for image recognition at scale](#). In *International Conference on Learning Representations*.
- Tejas Gokhale, Pratyay Banerjee, Chitta Baral, and Yezhou Yang. 2020. [MUTANT: A training paradigm for out-of-distribution generalization in visual question answering](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 878–892, Online. Association for Computational Linguistics.
- 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Danna Gurari, Qing Li, Abigale J. Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P. Bigham. 2018. [VizWiz grand challenge: Answering visual questions from blind people](#). In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzić, Rishabh Krishnan, and Dawn Song. 2020. [Pretrained transformers improve out-of-distribution robustness](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2744–2751, Online. Association for Computational Linguistics.
- Drew A. Hudson and Christopher D. Manning. 2019. [GQA: A new dataset for real-world visual reasoning and compositional question answering](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Allan Jabri, Armand Joulin, and Laurens van der Maaten. 2016. [Revisiting visual question answering baselines](#). In *Computer Vision – ECCV 2016*, pages 727–739, Cham. Springer International Publishing.
- Carlos E. Jimenez, Olga Russakovsky, and Karthik Narasimhan. 2022. [CARETS: A consistency and robustness evaluative test suite for VQA](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6392–6405, Dublin, Ireland. Association for Computational Linguistics.
- Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross Girshick. 2017. [CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning](#). In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. [Visual Genome: Connecting language and vision using crowdsourced dense image annotations](#). *International journal of computer vision*, 123(1):32–73.
- Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. 2021a. [Align before fuse: Vision and language representation learning with momentum distillation](#). In *Advances in Neural Information Processing Systems*, volume 34, pages 9694–9705. Curran Associates, Inc.
- Linjie Li, Jie Lei, Zhe Gan, and Jingjing Liu. 2021b. [Adversarial VQA: A new benchmark for evaluating the robustness of vqa models](#). In *International Conference on Computer Vision (ICCV)*.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. [Microsoft COCO: Common objects in context](#). In *Computer Vision – ECCV 2014*, pages 740–755, Cham. Springer International Publishing.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *International Conference on Learning Representations*.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. [ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks](#). In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Mateusz Malinowski and Mario Fritz. 2014. [A multi-world approach to question answering about real-world scenes based on uncertain input](#). In *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc.

- George A. Miller. 1995. [Wordnet: A lexical database for english](#). *Commun. ACM*, 38(11):39–41.
- Yulei Niu, Kaihua Tang, Hanwang Zhang, Zhiwu Lu, Xian-Sheng Hua, and Ji-Rong Wen. 2021. [Counterfactual VQA: A cause-effect look at language bias](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12700–12710.
- Vicente Ordonez, Girish Kulkarni, and Tamara Berg. 2011. [Im2text: Describing images using 1 million captioned photographs](#). In *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc.
- AJ Piergiovanni, Wei Li, Weicheng Kuo, Mohammad Saffar, Fred Bertsch, and Anelia Angelova. 2022. [Answer-Me: Multi-task open-vocabulary visual question answering](#). *arXiv preprint arXiv:2205.00949*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. [Learning transferable visual models from natural language supervision](#). In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR.
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. [Faster R-CNN: Towards real-time object detection with region proposal networks](#). In *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.
- Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. [LAION-400M: Open dataset of clip-filtered 400 million image-text pairs](#). *arXiv preprint arXiv:2111.02114*.
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. [Conceptual Captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, Melbourne, Australia. Association for Computational Linguistics.
- Sasha Sheng, Amanpreet Singh, Vedanuj Goswami, Jose Magana, Tristan Thrush, Wojciech Galuba, Devi Parikh, and Douwe Kiela. 2021. [Human-adversarial visual question answering](#). In *Advances in Neural Information Processing Systems*, volume 34, pages 20346–20359. Curran Associates, Inc.
- Haoyu Song, Li Dong, Weinan Zhang, Ting Liu, and Furu Wei. 2022. [CLIP models are few-shot learners: Empirical studies on VQA and visual entailment](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6088–6100, Dublin, Ireland. Association for Computational Linguistics.
- Damien Teney, Ehsan Abbasnejad, Kushal Kafle, Robik Shrestha, Christopher Kanan, and Anton van den Hengel. 2020. [On the value of out-of-distribution testing: An example of goodhart's law](#). 33:407–417.
- Antonio Torralba and Alexei A. Efros. 2011. [Unbiased look at dataset bias](#). In *CVPR 2011*, pages 1521–1528.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. 2021. [Training data-efficient image transformers & distillation through attention](#). In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 10347–10357. PMLR.
- Maria Tsimpoukelli, Jacob L. Menick, Serkan Cabi, S. M. Ali Eslami, Oriol Vinyals, and Felix Hill. 2021. [Multimodal few-shot learning with frozen language models](#). In *Advances in Neural Information Processing Systems*, volume 34, pages 200–212. Curran Associates, Inc.
- Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. 2022. [SimVLM: Simple visual language model pretraining with weak supervision](#). In *International Conference on Learning Representations*.
- 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*.
- Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. 2019. [Visual entailment: A novel task for fine-grained image understanding](#). *arXiv preprint arXiv:1901.06706*.
- Saining Xie, Ross Girshick, Piotr Dollar, Zhuowen Tu, and Kaiming He. 2017. [Aggregated residual transformations for deep neural networks](#). In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alexander J. Smola. 2015. [Stacked attention networks for image question answering](#). *CoRR*, abs/1511.02274.
- Mingda Zhang, Tristan D. Maidment, Ahmad Diab, Adriana Kovashka, and Rebecca Hwa. 2021. [Domain-robust VQA with diverse datasets and methods but no target labels](#). *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7042–7052.

## A Experimental Setup Details

In this section, we report additional details regarding our experimental setup.

### A.1 Models

We evaluate the performance of two strong models, ViLBERT (Lu et al., 2019) and ALBEF (Li et al., 2021a). These models belong to the family of pretrained Transformer that has recently achieved state-of-the-art performance on several V&L tasks, and are specifically instances of dual-stream architectures (Bugliarello et al., 2021). In this paradigm, models are first pretrained on a large collection of image–caption pairs, and then fine-tuned to solve specific downstream tasks. ViLBERT is pretrained using three objectives, masked language modeling (MLM; Devlin et al. 2019), masked region modeling and image–text matching (ITM; Chen et al. 2020). ALBEF is pretrained using MLM, ITM and an image–text contrastive loss (Li et al., 2021a). We refer the reader to Sec. 3 for an overall description of these models. Tab. 4 lists pretraining and architecture details for both models. All the models were fine-tuned using the AdamW (Loshchilov and Hutter, 2019) optimizer, with model-specific hyperparameters in Tab. 6.

**ViLBERT.** In this model, the textual inputs are first processed through 6 Transformer layers, before being combined with visual inputs through inter- and intra-modal attention layers. We re-implement this architecture, and confirm comparable performance by reproducing its results with the ones obtained through the official codebase<sup>6</sup> (see Tab. 5 for IID performance of both implementations). A key difference between the two implementations is in the image features: while the official model uses 10–100 regions of interest (RoIs) from a ResNeXt-152 (Xie et al., 2017), our re-implementation relies on 100 RoIs extracted by Faster R-CNN (Ren et al., 2015) trained on VG.

We then extend our codebase to implement a generative version of ViLBERT by replacing the discriminative decoder with an autoregressive decoding head.<sup>7</sup> The decoder is trained with teacher-forcing, and in datasets with several responses, the

<sup>6</sup><https://github.com/facebookresearch/vilbert-multi-task/>.

<sup>7</sup>Our implementation of the generative decoder follows that of TransformerDecoder available at <https://github.com/pytorch/pytorch/blob/master/torch/nn/modules/transformer.py>.

most frequent response is used as the ground truth response. Pretraining on CC, when used, is done by generating text. We also examine the effect of pretraining on both the encoder and decoder, and find the learning to be more stable when using only pretrained encoder, although further hyperparameter exploration could mitigate this difference. We study the effects of multimodal pretraining on the Conceptual Captions dataset (Sharma et al., 2018) with 3M images.

**ALBEF.** Like ViLBERT, ALBEF is a dual-stream encoder but with two main differences: first, the visual inputs are image patches that are processed through a vision Transformer (Dosovitskiy et al., 2021); and second, the cross-modal interactions happen through standard Transformer cross-attention at each layer (whereas ViLBERT uses co-attention layers specifically designed for intra- and inter-modal interactions). In addition, the model is trained with pseudo-targets that are generated from a moving-average version of its weights. We run experiments using the official codebase.<sup>8</sup> The visual backbone is a DeiT-B/16 (Touvron et al., 2021) pretrained on ImageNet-1k (Deng et al., 2009) at resolution  $224 \times 224$ , and further trained during the multimodal pretraining phase. For the downstream VQA benchmarks, we follow the authors and resize input images to  $384 \times 384$  and apply random augmentation during fine-tuning. Li et al. (2021a) formulated the VQA task as generative by adding a 6-layer Transformer decoder initialized from the pretrained encoder. We follow this approach and also evaluate a discriminative version by learning a two-layer MLP with ReLU (Agarap, 2018) non-linearity in between, following the authors’ setup for the Visual Entailment benchmark (Xie et al., 2019). We found the hyperparameters proposed by Bugliarello et al. (2021) to work better. During inference, we evaluate ALBEF in two ways: first, following the authors, we rank the in-domain candidate answers based on their likelihood; second, we let the model generate any possible answer in an open-ended fashion through greedy decoding. We found these two approaches to minimally affect final performance (see Tab. 7). Unless otherwise specified, we report results given by generation as it reflects open-ended question answering.

<sup>8</sup><https://github.com/salesforce/ALBEF>.

Model	# Params	Pretrain data	# Images	# Captions
VILBERT	240M	CC	3.3M	3.3M
ALBEF (4M)	450M	+COCO+SBU+VG	4M	5M
ALBEF (14M)	450M	+C12M	14M	15M

Table 4: Pretrained models details. ALBEF size includes both the main model and its moving average. Pretraining data: CC (Sharma et al., 2018), COCO (Lin et al., 2014), SBU (Ordonez et al., 2011), VG (Krishna et al., 2017), C12M (Changpinyo et al., 2021).

Model	VQAv2	GQA	VG	VIZWIZ
Official codebase	67.04	66.78	40.69	44.46
Re-implementation	65.75	61.51	40.31	47.46

Table 5: Comparison between the official and our codebases for VILBERT<sub>DISC</sub> in the IID setting.

## A.2 Datasets

Tab. 1 lists statistics for each dataset in our study. VQAv2 (Goyal et al., 2017) is the most commonly used VQA dataset to date, it consists of 265K images and 1.1M question-image pairs, each with 10 ground-truth answers. VQA-CP (Agrawal et al., 2018) re-splits the VQAv2 dataset such that, for every question type, train and test sets have different prior distributions of answers. VG (Krishna et al., 2017) includes 108K images and 1.7M questions, each paired with a single answer, centered around either the full image or a specific region within it. GQA (Hudson and Manning, 2019) is another large-scale effort (22M questions, each with one answer) that focuses on compositionality of template-generated questions for real-world images (from VG). Following prior work, we use the GQA *balanced* subset (1.5M questions). Finally, VIZWIZ (Gurari et al., 2018) is the only real-world VQA dataset as it was collected from visually impaired people. It consists of 31K image-question pairs, each paired with 10 answers.

## A.3 Training Details

Following common practice, for discriminative models, we select the top- $k$  most frequent answers from the fine-tuning dataset, as the set of answer classes to perform classification over. Here  $k$  is a dataset-dependent variable. For VQAv2 and GQA, we use the same answer sets as VILBERT (3,129 and 1,533, respectively). For VIZWIZ, we select the answers that appear at least 8 times in training and validation sets, for a total of 3,112 answers that cover 97% of the data. For VG, we select the answers that appear at least 29 times in the dataset, for a total of 3,449 answers that cover 76.5% of the data. Importantly, combined with the VQA accu-

Benchmark	Qn len	# Classes	BS	LR	# Epochs
VQAv2	16	3,129	256	4e-5	20
GQA	26	1,533	256	4e-5	20
VG	16	3,449	256	4e-5	20
VIZWIZ	40	3,112	256	4e-5	20

(a) Parameters used for VILBERT models. The internal codebase uses LAMB optimizer with the initial LR of 1e-3, with the best checkpoint selected on eval dataset.

Benchmark	Qn len	# Classes	BS	LR	# Epochs
VQAv2	16	3,129	256	1e-4	20
GQA	26	1,533	256	1e-4	20
VG	16	3,449	256	1e-4	20
VIZWIZ	40	3,129	256	1e-4	20

(b) Discriminative ALBEF.

Benchmark	Qn len	Answer len	BS	LR	# Epochs
VQAv2	16	6	256	2e-5	20
GQA	26	5	256	2e-5	20
VG	16	8	256	2e-5	20
VIZWIZ	40	11	256	2e-5	20

(c) Generative ALBEF.

Table 6: Hyperparameters used in our experiments. Question and answer lengths are in tokens, BS is the batch size, LR is the learning rate.

racy metric defined above, this results in an upper-bound to the accuracy that discriminative models can achieve in each dataset (see Tab. 2).

All models are trained on the respective training sets and evaluated on the validation sets, which lets us conduct in-depth analyses that would otherwise be impossible to carry out on private test sets. As there is no official split of VG, we randomly sample the data into training (60%) and validation (40%) such that no image appears in both splits.

## B Additional Results

**Evaluation with Shared Answer Sets** While different answer sets are an apparent issue for discriminative models, they also impact the performance of generative models, as the number of data points for each answer class seen by the generative model during fine-tuning varies: data-points in top- $k$  answer set are more frequent than others (by definition of top- $k$ ). In other words, even though a tokenizer used to produce an answer could generate it, it is unlikely (or less likely) to do so if it has not seen (or seen rarely) that combination of tokens during fine-tuning. Thus, even for generative models, we consider performance on top- $k$  most frequent classes for each benchmark.

Thus, we report the accuracy on the subset of

	VQAV2	GQA	VG	VIZWIZ
VQAV2	72.37	50.56	38.94	19.81
GQA	50.32	64.26	22.80	12.51
VG	33.40	24.99	43.35	12.60
VIZWIZ	34.17	22.73	8.89	48.44

(a) Ranking-based evaluation of ALBEF<sub>GEN</sub>.

	VQAV2	GQA	VG	VIZWIZ
VQAV2	72.09	50.10	39.20	19.81
GQA	50.33	64.24	22.79	12.50
VG	33.39	23.64	44.55	12.25
VIZWIZ	34.44	22.80	9.13	47.14

(b) Generation-based evaluation of ALBEF<sub>GEN</sub>.

Table 7: Performance of ALBEF<sub>GEN</sub> (14M) when tested via ranking (top) and generation (bottom). Rows correspond to the fine-tuning datasets, columns correspond to the test benchmarks. The model performs similarly in both setups. We found similar results for ALBEF<sub>GEN</sub> (4M).

test questions whose answers are shared between *both* the IID and the OOD models. For instance, when comparing the performance of the VQAV2 and VG fine-tuned models on the VQAV2 test set, we compute the average accuracy on those VQAV2 questions whose ground truth answers are present in the top- $k$  answers from VQAV2 as well as the top- $k$  answers from VG: we extract the common answer labels (between VQAV2 and VG top- $k$  answers) and compute performance on test questions belonging to these shared answer labels only.

For IID evaluations, there are several possible ways to define shared answer sets based on OOD vocabs. While a subset is shown in Fig. 2, Tab. 10 lists the VQA accuracy of each model in the IID settings when evaluated on the questions in the test sets whose answers are shared between the top- $k$  answers in both the IID and the OOD settings (see Sec. 4.1 for more details).

In some IID cases, restricting the answer set to common answers hurts the performance (indicated as a lack of dotted bar in Fig. 2). Interestingly, this pattern is observed across all models for some IID evaluations where the shared answer set is computed with respect to the VG benchmark only: VQAV2 for ViLBERT<sub>DISC</sub> (-7.65 pp drop), VQAV2 (-8.65 pp drop) and GQA (-8.00 pp drop) for ViLBERT<sub>GEN</sub>, VQAV2 (-6.86 pp drop) for ALBEF<sub>DISC</sub>, and VQAV2 (-6.89 pp drop) for ALBEF<sub>GEN</sub>. This seems to indicate that the GQA and VQAV2 questions corresponding to shared answer set with VG are more difficult than the average difficulty of these test sets.

	VQAV2	GQA	VG	VIZWIZ
VQAV2	–	0.41	0.51	0.11 <sup>^</sup>
GQA	0.25	–	0.44	0.14 <sup>^</sup>
VG	0.28	0.38	–	0.03 <sup>^</sup>
VIZWIZ	0.46	0.54	0.48	–

(a) Discriminative ALBEF.

	VQAV2	GQA	VG	VIZWIZ
VQAV2	–	0.45	0.48	0.26
GQA	0.26	–	0.45	0.22
VG	0.27	0.35	–	0.09 <sup>^</sup>
VIZWIZ	0.48	0.56	0.52	–

(b) Generative ViLBERT.

	VQAV2	GQA	VG	VIZWIZ
VQAV2	–	0.49	0.50	0.27
GQA	0.30	–	0.45	0.34
VG	0.25	0.38	–	0.16
VIZWIZ	0.50	0.57	0.51	–

(c) Discriminative ViLBERT.

Table 8: Spearman’s rank correlation between drops in test accuracy (from IID to OOD) and the differences in proportion of answer classes between IID and OOD fine-tune sets. Unless otherwise specified with a <sup>^</sup> character,  $\rho$  values are significant with  $p < .05$ . Rows correspond to the fine-tuning datasets, columns correspond to the test benchmarks.

	VQAV2	GQA	VG	VIZWIZ
VQAV2	–	0.43	0.51	0.25
GQA	0.27	–	0.43	0.19
VG	0.26	0.36	–	0.13
VIZWIZ	0.47	0.55	0.48	–

Table 9: Spearman’s rank correlation between drops in test accuracy (from IID to OOD) and the differences in proportion of answer classes between IID and OOD fine-tuning sets for ALBEF<sub>GEN</sub>.  $p < .05$  for all  $\rho$ . Rows correspond to the fine-tuning datasets, columns correspond to the test benchmarks.

**Answer Frequency Correlation** In order to examine the relationship between accuracy drop for less frequent classes, we first compute per answer-class accuracy (average accuracy of all test questions belonging to the same answer class) for answers in shared answer set. We then sort the shared answer classes based on their weighted drop in per-class accuracy from IID to OOD (IID accuracy - OOD accuracy), *i.e.* absolute drop in per-class accuracy weighted by number of data points belonging to that class in the test set. We then compute the Spearman’s rank correlation of these weighted drop in per-class accuracies with difference in percentage frequencies of the answer classes between IID



and OOD fine-tuning sets (percentage frequency of an answer class in IID minus its percentage frequency in OOD).

Tab. 8 list Spearman’s rank correlations of IID-to-OOD drops in test accuracy vs. proportion of answer classes in respective (IID and OOD) fine-tuning sets for ALBEF<sub>DISC</sub>, ViLBERT<sub>GEN</sub> and ViLBERT<sub>DISC</sub> (see Sec. 4.1 for more details). As a simple baseline test, we also compute correlations and p-values for a permuted dataset to confirm their lack of significance, or correlation values close to zero.

**Maximum Achievable Scores** Table 2 lists maximum achievable accuracies, and Figure 7 shows the difference between those scores and bar values shown in Figure 1. In our analyses, we also noted that differences in answer pre-processing strategies can result in slightly different numbers than those reported in Tab. 2. However, those differences did not change the conclusion of our findings.

**Effect of pretraining data size on ALBEF** For the ALBEF model, while we often observe improvements by increasing the size of the multi-modal pretraining dataset (4M vs. 14M), the improvements are small. When pretraining on the smaller dataset (4M, see Fig. 8), we observe a median improvement (over no pretraining) of 1.9% for the discriminative and 4.9% for the generative ALBEF, while the median additional improvements due to larger pretraining dataset (14M) are 0.1% and 0.6% respectively (refer to Fig. 3). Surprisingly, there are also dataset pairs for which larger pretraining has a negative effect when compared to the performance with a smaller pretraining set (e.g., ALBEF model fine-tuned on VIZWIZ and tested on VQAV2).

## C Potential Causes of Poor OOD Generalization: A Qualitative Study

In section 4, we observe that our pretrained models exhibit poor OOD generalization for the task of VQA. We also noted that this poor generalization is not entirely explained by the absence or poor representation of test answer classes in the training data. Here, we perform a qualitative study to dig deeper into the potential causes of the poor OOD generalization. We manually examine 20 randomly-sampled qualitative examples of failure cases on top-30 answer classes contributing the most to the drop in performance from IID to OOD.

We only focus on answer classes that are shared between the train and test splits to make sure the performance drop is not due to the absence of answer classes in the training dataset. We report the top-5 classes that contribute the most to the drop in performance for each OOD setting in Tab. 11. Below, we describe four major potential causes<sup>9</sup> for the poor OOD generalization that we can infer from our qualitative study on ViLBERT<sub>DISC</sub><sup>10</sup> and ALBEF<sub>GEN</sub>. The specific examples reported below are for ViLBERT<sub>DISC</sub>.

**Poor reasoning skills.** In Tab. 11, we can see that a model fine-tuned on VQAV2, VG, or VIZWIZ and evaluated on GQA shows the highest performance drop on classes such as “yes”, “no”, “right”, “left”, “top”, and “bottom”. For instance, VQAV2–GQA (fine-tuned on VQAV2, evaluated on GQA) model underperforms GQA–GQA model by 24% for “no.” Upon qualitative examination, we find that for many of such failure cases, the GQA questions are more compositional and hence require more complex reasoning (e.g., “Are there both bison and zebras in the image?”, “Is the cheese to the right or to the left of the empty plate?”) than the questions for the same answer classes in other datasets (e.g., from VQAV2 train set: “Is the TV turned on?”, “Which hand is the man holding up?”). This study re-affirms previous findings (Johnson et al., 2017; Hudson and Manning, 2019) – VQA models lack sufficient logical, spatial, and compositional reasoning skills – for the more recent, pretrained Transformer models.

**Overfitting to the answer priors.** Previous studies have shown that VQA models tend to be biased towards the prior distribution of answers in the training set (per question type) (Agrawal et al., 2018). We find that this limitation exists in the more recent pretrained models as well, and it is especially hurtful in the OOD settings because the priors need not be the same across train and test sets, unlike in the IID settings. For instance, ViLBERT<sub>DISC</sub> fine-tuned on VQAV2 predicts “2” for a lot questions with target answer “1” in the VG test set. Similarly, sometimes ViLBERT<sub>DISC</sub> fine-tuned on VG incorrectly predicts “helmet” for

<sup>9</sup>For poor OOD generalization on the VIZWIZ benchmark, one of the reasons could be difference in image distributions between VIZWIZ (that contains many blurry pictures, or pictures with poor lighting conditions) and other three datasets (that contain clear pictures).

<sup>10</sup>We use the model trained with the official codebase.

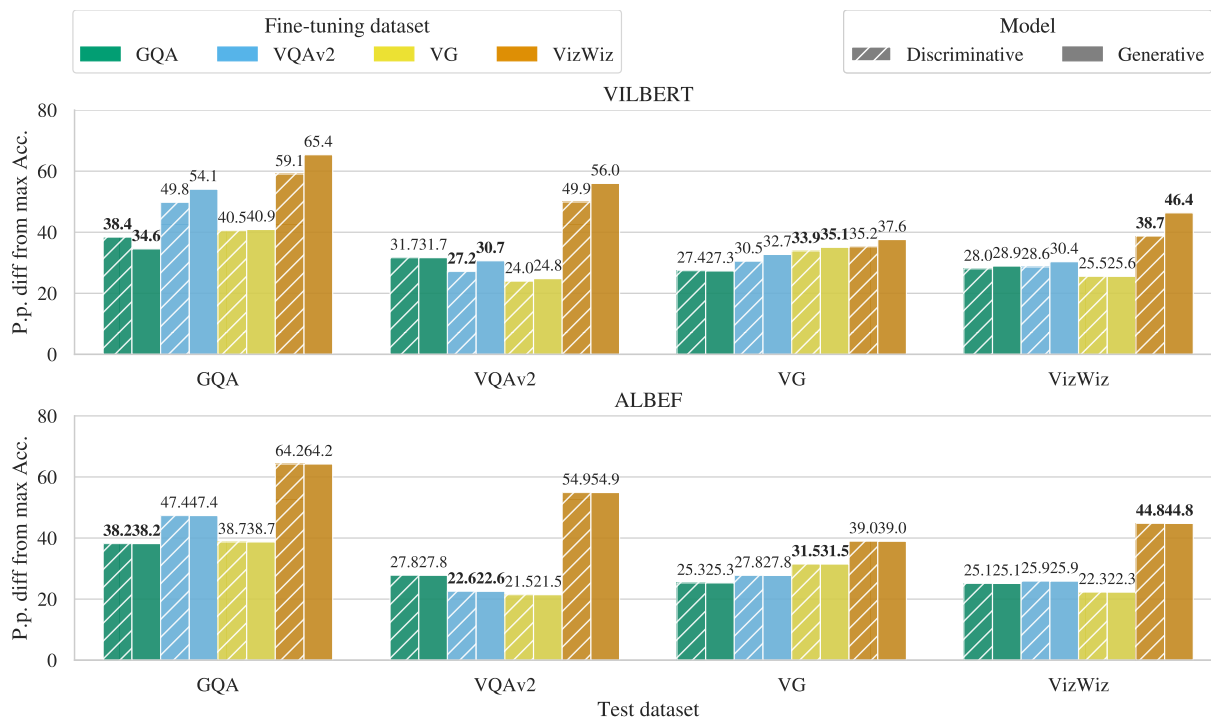


Figure 7: Percentage point difference between maximum achievable accuracies in Table 2 and accuracies in Figure 1. Results for VILBERT pretrained on the same data as ALBEF 4M are also shown.

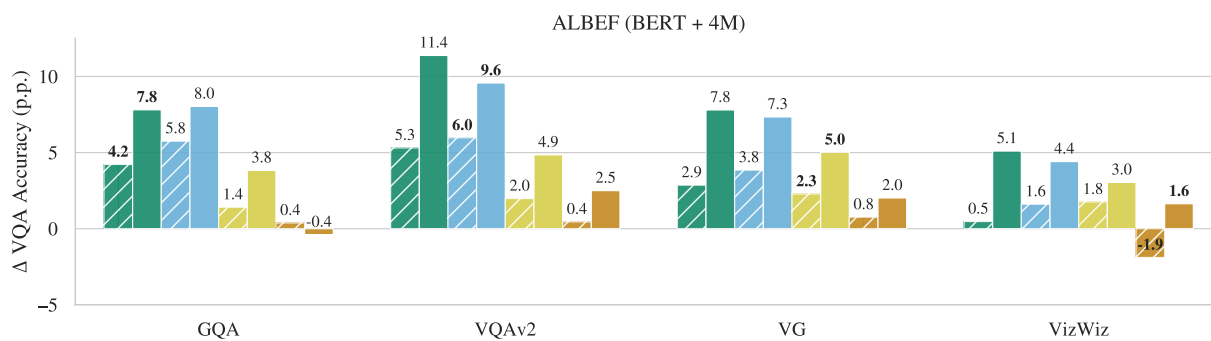


Figure 8: Difference in VQA accuracy (p.p.) for ALBEF that has and has not been pretrained on the 4M dataset.

VQAv2 test questions such as “What is the skateboarder wearing to protect his head?”, “What protective gear is he wearing?” when the skateboarder is not wearing anything. This indicates that the model is relying on answer priors rather than visual grounding. Our experimental results on VQA-CP (Sec. 6) directly quantify the extent of such limitations in current models.

**Overfitting to the question format.** For each answer class, there is usually a limited variation in the format of questions in the fine-tuning set. For some of the answer classes showing poor OOD generalization, we found that certain question formats are quite dominant in the fine-tuning set, and that these dominant formats are different between the OOD fine-tuning and test sets. Thus, we conjecture

that models are likely overfitting to such dominant formats in fine-tuning data and hence fail to generalize at test time when the format changes. For instance, questions about “chair” in the VQAv2 fine-tuning set are mostly of the form “What is ... sitting on?” whereas in the GQA test set, they are mostly of the form “What kind of furniture is ...?”. Thus, the “chair” class accuracy of VILBERT<sub>DISC</sub> fine-tuned on VQAv2 drops from 48% when tested on VQAv2 to 38% on the GQA test set. As another example, VILBERT<sub>DISC</sub> trained on GQA fails terribly for “dog” and “cat” classes on VG test set (accuracy drops of 47% and 43% respectively, where drop is between GQA–GQA and GQA–VG). GQA questions are mostly of the form “What animal ...?” or “What kind of animal

Model	Test	Fine-tune	Answer Set	VQA Acc. (IID)	VQA Acc. (OOD)	Difference
ViLBERT <sub>DISC</sub>	GQA	GQA	GQA $\cap$ VQAv2	63.05	48.53	14.52
ViLBERT <sub>DISC</sub>	GQA	GQA	GQA $\cap$ VG	52.32	35.08	17.24
ViLBERT <sub>DISC</sub>	GQA	GQA	GQA $\cap$ VizWiz	64.91	28.39	36.52
ViLBERT <sub>DISC</sub>	VG	VG	VG $\cap$ VQAv2	58.18	51.97	6.21
ViLBERT <sub>DISC</sub>	VG	VG	VG $\cap$ GQA	59.57	39.15	20.42
ViLBERT <sub>DISC</sub>	VG	VG	VG $\cap$ VizWiz	60.65	13.23	47.42
ViLBERT <sub>DISC</sub>	VizWiz	VizWiz	VizWiz $\cap$ VQAv2	51.78	22.37	29.41
ViLBERT <sub>DISC</sub>	VizWiz	VizWiz	VizWiz $\cap$ GQA	51.05	15.40	35.65
ViLBERT <sub>DISC</sub>	VizWiz	VizWiz	VizWiz $\cap$ VG	50.07	13.59	36.48
ViLBERT <sub>DISC</sub>	VQAv2	VQAv2	VQAv2 $\cap$ GQA	71.74	52.42	19.32
ViLBERT <sub>DISC</sub>	VQAv2	VQAv2	VQAv2 $\cap$ VG	58.10	48.84	9.26
ViLBERT <sub>DISC</sub>	VQAv2	VQAv2	VQAv2 $\cap$ VizWiz	68.20	33.87	34.33
ViLBERT <sub>GEN</sub>	GQA	GQA	GQA $\cap$ VQAv2	67.01	44.03	22.98
ViLBERT <sub>GEN</sub>	GQA	GQA	GQA $\cap$ VG	57.35	34.40	22.95
ViLBERT <sub>GEN</sub>	GQA	GQA	GQA $\cap$ VizWiz	69.63	20.73	48.90
ViLBERT <sub>GEN</sub>	VG	VG	VG $\cap$ VQAv2	55.52	47.93	7.59
ViLBERT <sub>GEN</sub>	VG	VG	VG $\cap$ GQA	57.65	39.26	18.39
ViLBERT <sub>GEN</sub>	VG	VG	VG $\cap$ VizWiz	58.84	7.67	51.17
ViLBERT <sub>GEN</sub>	VizWiz	VizWiz	VizWiz $\cap$ VQAv2	43.06	19.58	23.48
ViLBERT <sub>GEN</sub>	VizWiz	VizWiz	VizWiz $\cap$ GQA	42.57	13.71	28.86
ViLBERT <sub>GEN</sub>	VizWiz	VizWiz	VizWiz $\cap$ VG	41.11	13.20	27.91
ViLBERT <sub>GEN</sub>	VQAv2	VQAv2	VQAv2 $\cap$ GQA	68.01	52.35	15.66
ViLBERT <sub>GEN</sub>	VQAv2	VQAv2	VQAv2 $\cap$ VG	53.59	46.78	6.81
ViLBERT <sub>GEN</sub>	VQAv2	VQAv2	VQAv2 $\cap$ VizWiz	64.69	26.82	37.87
ALBEF <sub>DISC</sub>	GQA	GQA	GQA $\cap$ VQAv2	63.24	51.06	12.18
ALBEF <sub>DISC</sub>	GQA	GQA	GQA $\cap$ VG	53.09	37.97	15.12
ALBEF <sub>DISC</sub>	GQA	GQA	GQA $\cap$ VizWiz	64.85	22.14	42.71
ALBEF <sub>DISC</sub>	VG	VG	VG $\cap$ VQAv2	61.38	55.83	5.55
ALBEF <sub>DISC</sub>	VG	VG	VG $\cap$ GQA	63.76	43.82	19.94
ALBEF <sub>DISC</sub>	VG	VG	VG $\cap$ VizWiz	64.12	4.52	59.60
ALBEF <sub>DISC</sub>	VizWiz	VizWiz	VizWiz $\cap$ VQAv2	42.42	26.23	16.19
ALBEF <sub>DISC</sub>	VizWiz	VizWiz	VizWiz $\cap$ GQA	40.49	20.44	20.05
ALBEF <sub>DISC</sub>	VizWiz	VizWiz	VizWiz $\cap$ VG	38.15	19.68	18.47
ALBEF <sub>DISC</sub>	VQAv2	VQAv2	VQAv2 $\cap$ GQA	76.64	57.18	19.46
ALBEF <sub>DISC</sub>	VQAv2	VQAv2	VQAv2 $\cap$ VG	63.47	52.78	10.69
ALBEF <sub>DISC</sub>	VQAv2	VQAv2	VQAv2 $\cap$ VizWiz	72.84	28.08	44.76
ALBEF <sub>GEN</sub>	GQA	GQA	GQA $\cap$ VQAv2	65.81	51.72	14.09
ALBEF <sub>GEN</sub>	GQA	GQA	GQA $\cap$ VG	56.08	37.71	18.37
ALBEF <sub>GEN</sub>	GQA	GQA	GQA $\cap$ VizWiz	67.05	27.61	39.44
ALBEF <sub>GEN</sub>	VG	VG	VG $\cap$ VQAv2	62.71	57.33	5.38
ALBEF <sub>GEN</sub>	VG	VG	VG $\cap$ GQA	65.48	51.17	14.31
ALBEF <sub>GEN</sub>	VG	VG	VG $\cap$ VizWiz	66.20	21.13	45.07
ALBEF <sub>GEN</sub>	VizWiz	VizWiz	VizWiz $\cap$ VQAv2	52.85	28.96	23.89
ALBEF <sub>GEN</sub>	VizWiz	VizWiz	VizWiz $\cap$ GQA	52.58	22.21	30.37
ALBEF <sub>GEN</sub>	VizWiz	VizWiz	VizWiz $\cap$ VG	51.94	23.56	28.38
ALBEF <sub>GEN</sub>	VQAv2	VQAv2	VQAv2 $\cap$ GQA	78.03	62.93	15.10
ALBEF <sub>GEN</sub>	VQAv2	VQAv2	VQAv2 $\cap$ VG	65.20	55.64	9.56
ALBEF <sub>GEN</sub>	VQAv2	VQAv2	VQAv2 $\cap$ VizWiz	74.38	39.37	35.01

Table 10: VQA accuracy of each model in the IID settings (see column VQA Acc. (IID)) when evaluated on the questions in the test sets whose answers are shared between the top- $k$  answers in both the IID and the OOD settings. Please refer to Sec. 4.1 for more details. Answer Set: OOD benchmarks with respect to which IID shared answer set accuracy is computed. VQA Acc. (OOD): OOD accuracy on questions corresponding to the shared answer set, *i.e.* when fine-tuned on the OOD dataset mentioned in Answer Set column and tested on the benchmark mentioned in the Test column. Difference: VQA Acc. (IID) - VQA Acc. (OOD). Gray bands highlight the OOD benchmarks with respect to which IID shared answer set accuracy is computed in Fig. 2.

...?” whereas VG questions often do not mention the word “animal” and are of the form “Who is ... ?” or “What is ... ?” (*e.g.*, “Who is holding the Frisbee?”, “What is on the leash?”). Similarly, for the answer class “pizza”, ViLBERT<sub>DISC</sub> fine-tuned

on VG has mostly seen questions of the format “What food is this?”, “What is the man eating?”, “What is on the plate?”, “What’s in the box?” in VG fine-tuning set. However, when evaluated on the VQAv2 test set, the model fails to respond cor-

rectly for questions about “pizza” such as, “What snack is this?” (model response: “pineapple”), “What recipe this will become?” (model response: “cheese”), “What’s in the bowl” (model response: “tomato sauce”). For the last example, perhaps the model is not expecting pizza to be in a bowl.

Related to above, we observed that sometimes `VILBERTDISC` fails to produce the correct answer type for a given question. For instance, `VILBERTDISC` fine-tuned on VG responds with “woman” to the question “Is the person who is cutting these carrots right handed or left handed?”. So it appears as if the VG model does not understand the question structure in this example, *i.e.* the response is expected to be either “right” or “left”. Similarly, for the question “Are there more blue or black shirts?”, VG model responds with “rolled up”. Similarly, it answers “1 on right” to the question “What type of apple is shown?”, instead of describing some attribute of apple such as “green”.

**Stringent evaluation metric.** We notice that sometimes the models’ responses are correct but they are evaluated as incorrect because those responses do not exist in the ground-truth answers. For instance, VQAV2–VG model gets penalized for answering “table” instead of “on table”<sup>11</sup> (Q: “Where is . . . ?”) or “sunny” instead of “clear” (Q: “How is the weather?”). More examples in Fig. 5. This effect is expected to be more pronounced for the OOD evaluation than IID, because in IID a model can learn the format of the test answer (“on table” vs. “table”, “clear vs. sunny”) from the train set, whereas in OOD the format in the train set can be different from the test set. Also, such stringent evaluation (*i.e.*, performing string matching with a small set of ground-truth answers) is expected to hurt generative models more than discriminative ones because they show more variations in the form of the answers as they are not limited by a fixed answer vocabulary (*e.g.*, “pizza slices” instead of “pizza” (Q: “What are these?”), “pizzeria” instead of “pizza” (Q: “What kind of restaurant is this?”)). We observed that, VG model (model fine-tuned on VG) evaluated on GQA answers questions about “man” with “snowboarder”, “man on left” (*i.e.* more descriptive referring expressions) than just saying “man” but it does not get any credit

<sup>11</sup>Note that before computing the accuracy, both the predicted and the ground truth answers are pre-processed for answer normalization but such pre-processing is very basic. More details of the pre-processing can be found at <https://visualqa.org/evaluation.html>

because GQA ground truth is “man”. To quantify the extent of this issue and measure its effect on discriminative vs. generative models, IID vs. OOD settings, we perform human evaluation of machine generated answers and provide additional insights in Sec. 8.

**Poor performance of GQA model on color questions (both IID and OOD):** `VILBERTDISC` fine-tuned on GQA does not seem to be transferring well to color questions in the VQAV2 and VG test sets (and even in IID GQA test set). In Tab. 11, we can see that the top-5 answer classes with highest drop in IID-to-OOO performance for GQA model have quite a few colors. For instance, for the answer class “red” in the VG test set, GQA model fails to correctly answer simple questions (given the kind of questions GQA model is fine-tuned on) such as “What is the primary color of the sign on the right?”, “What is the main color of the strawberry?”, “What color is the pull luggage of the woman?”, “What color are the pepperonis?”. It is not clear why GQA model does not perform well on color questions.

## D Human Evaluation

**Method.** We used Amazon Mechanical Turk to collect human judgment about model responses on a random subset of 10K questions for each of the test sets—VQAV2, GQA and VG. Since the size of VIZWIZ test set is less than 10K, we collected human judgment on all the VIZWIZ test questions. However, we dropped the questions that were tagged as “unanswerable” or “unsuitable” (more details are provided below under “Filtering VizWiz data”). The total number of VIZWIZ test questions for which we collected human judgment is 1440 (per model). We performed human evaluation of the responses from the following models – `VILBERTDISC`<sup>12</sup> and `VILBERTGEN` trained on the VQAV2, GQA, VG datasets. We did not collect human judgements for models *fine-tuned* on VIZWIZ, because a significant proportion of the responses from these models tend to be “unanswerable” or “unsuitable” (35% on VQAV2, 39% on GQA, 65% on VG, and 64% on VIZWIZ). Collecting human feedback about such responses would

<sup>12</sup>For `VILBERTDISC`, we had initially collected human judgements for the version trained using the official codebase, and we did not collect annotations again for our reimplementation due to time constraints. Given our results above, we do not expect significant differences between the two versions.

Train data	Test data	Model	Answer classes	
VQAv2	GQA	Discriminative ViLBERT	no, yes, left, right, top	
		Discriminative ViLBERT (in-house)	no, yes, right, bottom, color	
		Generative ViLBERT (in-house)	no, yes, right, left, bottom	
		Discriminative ALBEF	no, left, yes, bottom, chair	
	VG	Generative ALBEF	left, no, yes, bottom, top	
		Discriminative ViLBERT	1, no 1, daytime, on table, in sky	
		Discriminative ViLBERT (in-house)	daytime, 1, white, 2, black	
		Generative ViLBERT (in-house)	daytime, white, 2, black, 1	
	VizWiz	Discriminative ALBEF	1, daytime, in sky, on table, white	
		Generative ALBEF	1, daytime, black, in sky, clear	
		Discriminative ViLBERT	no, blue, yes, white, black	
		Discriminative ViLBERT (in-house)	yes, black, water bottle, corn, soup	
GQA	VQAv2	Generative ViLBERT (in-house)	pink, brown, corn, wine, keys	
		Discriminative ALBEF	keyboard, no, soup, cake, samsung	
		Generative ALBEF	soup, lotion, black, brown, corn	
		Discriminative ViLBERT	yes, no, white, red, black	
	VG	Discriminative ViLBERT (in-house)	no, yes, white, red, tennis	
		Generative ViLBERT (in-house)	no, yes, white, red, tennis	
		Discriminative ALBEF	no, yes, white, red, right	
		Generative ALBEF	no, yes, right, red, black and white	
	VizWiz	Discriminative ViLBERT	white, trees, green, black, black and white	
		Discriminative ViLBERT (in-house)	white, black, trees, green, blue	
		Generative ViLBERT (in-house)	white, trees, black, green, brown	
		Discriminative ALBEF	white, trees, black and white, grass, green	
VG	VQAv2	Generative ALBEF	trees, green, black and white, black, grass	
		Discriminative ViLBERT	no, blue, yes, white, laptop	
		Discriminative ViLBERT (in-house)	blue, white, black, dog, laptop	
		Generative ViLBERT (in-house)	white, blue, laptop, black, dog	
	GQA	Discriminative ALBEF	no, keyboard, soup, red, cake	
		Generative ALBEF	no, dog, keyboard, laptop, blue	
		Discriminative ViLBERT	0, white, nothing, gray, red	
		Discriminative ViLBERT (in-house)	0, 3, left, nothing, brown	
	VizWiz	Generative ViLBERT (in-house)	0, 1, gray, left, wii	
		Discriminative ALBEF	0, nothing, left, brown, 2	
		Generative ALBEF	0, 3, nothing, right, gray	
		Discriminative ViLBERT	right, left, bottom, top, gray	
VizWiz	GQA	Discriminative ViLBERT (in-house)	bottom, left, top, color, large	
		Generative ViLBERT (in-house)	left, bottom, color, top, gray	
		Discriminative ALBEF	left, bottom, top, black, chair	
		Generative ALBEF	left, bottom, color, top, gray	
	VG	Discriminative ViLBERT	blue, black, grey, red, soup	
		Discriminative ViLBERT (in-house)	grey, black, blue, white, computer screen	
		Generative ViLBERT (in-house)	grey, blue, black, pink, computer screen	
		Discriminative ALBEF	grey, soup, remote, cake, samsung	
	VQAv2	GQA	Generative ALBEF	grey, blue, soup, wine, pink
			Discriminative ViLBERT	no, yes, 1, 2, white
			Discriminative ViLBERT (in-house)	no, 1, 2, 0, white
			Generative ViLBERT (in-house)	no, yes, 1, 2, white
VG		Discriminative ALBEF	no, 1, 2, yes, blue	
		Generative ALBEF	yes, no, 1, 2, 0	
		Discriminative ViLBERT	no, right, left, man, bottom	
		Discriminative ViLBERT (in-house)	no, right, bottom, man, top	
VizWiz		GQA	Generative ViLBERT (in-house)	no, yes, left, right, man
			Discriminative ALBEF	no, left, bottom, top, man
			Generative ALBEF	yes, left, no, bottom, top
			Discriminative ViLBERT	1, white, green, 2, black
VG	VG	Discriminative ViLBERT (in-house)	1, 2, white, green, man	
		Generative ViLBERT (in-house)	1, 2, white, green, black	
		Discriminative ALBEF	1, green, white, 1, 2, blue	
		Generative ALBEF	1, 2, black, man, green	

Table 11: Top-5 answer classes with highest performance drop from IID to OOD (for the same test set) for all OOD configurations. The answer classes are sorted by drop in weighted (wtd) accuracy, i.e. drop in absolute (abs) accuracy weighted by the # test questions for that answer class.

not provide useful insights, because all questions in VQAv2, GQA and VG should be answerable, therefore all cases of “unanswerable” should be incorrect. Such responses are just a side effect of a model’s priors caused by all the unanswerable training points in the VIZWIZ fine-tuning set.

For each response, we asked 5 raters to evaluate the question, image, and a given model response, and indicate through a binary choice whether they considered the model response a correct answer to the question or not. To control the quality of the data, we filtered out low quality data using different heuristics such as distribution of yes/no answers for each worker, their mean submission times, average agreement with their fellow workers, or average alignment with the automatic accuracy.<sup>13</sup> In each of these cases, we looked at random samples from the outliers to qualitatively confirm our hypothesis. More details about the human evaluation interface are presented in the next paragraph.

To compute human accuracy of a model response (for a given question and image), we considered a response correct if at least 4 raters voted it is correct, and incorrect otherwise. We decided so in order to decrease noise introduced by cases where there was low agreement between raters.

**Data collection interface.** Fig. 10 shows a sample of the interface the MTurk raters used to submit their responses. The workers were shown some examples, but in order not to bias them, we did not give them detailed guidance as to what should be considered correct for not - rather we asked them to rely on common sense, and consider an answer correct if it seems both factually accurate and natural in the context. See Fig. 11 for details.

**Filtering VizWiz data.** For human evaluation we filtered out all image-question pairs for which the ground truth answer indicates it is *unanswerable*. That is, we have not collected human feedback for questions for which the ground truth answer appears in the following list:

- "unanswerable", "unsuitable"
- "insufficient image"
- "unknown", "unsure", "not clear"
- "blurry", "too blurry"

<sup>13</sup>How frequently a worker’s response (yes/no) aligns with the automatic accuracy computed (100.0/0.0) More specifically, we equate the worker’s *yes* response with 100.0 and *no* with 0.0 and look at the average difference between the worker’s response and the automatic accuracy

- "i don’t know", "don’t know", "i don no", "no idea"
- "unusable image", "unsuitable image", "unstable image", "insufficient image quality", "unreadable"
- "i can’t tell", "can’t tell", "can’t see"
- "0" <sup>14</sup>

In particular, this left us with 1440 questions for the VIZWIZ dataset.

**Results.** We report the human accuracies for VILBERT<sub>DISC</sub> and VILBERT<sub>GEN</sub> in Fig. 9 (bottom). We also report the accuracies obtained using automatic metrics (please see Sec. 3.2 for description of automatic metrics for each dataset) computed over the same random subset of test questions as that used for human evaluation in Fig. 9 (top). Please refer to Sec. 8 for discussion of results.

**Qualitative examples of questions being incorrectly penalized by automatic evaluation** Tab. 12 shows some examples for responses which were awarded 0.0 accuracy using automatic metrics but were marked as correct by all 5 raters during human evaluation.

**Discussion on VQA data quality** For the collected human judgement data, we find that for a significant number of questions (32%) there was low agreement between the 5 raters, i.e. either 3/5 answered *correct* while the remaining 2/5 answered *incorrect*, or vice-versa. Note that this is after we already filtered out low quality judgements. We have to recognize that, despite our best efforts to control data quality using our heuristics, there might still be low quality data in our dataset. Low quality annotations can be misleading and might distort the results of our analysis. Yet, we believe that we have collected a large enough sample to dampen the effect of these on the reported results. Upon examining some examples from such low agreement questions, we find that many such cases highlight the quality of the VQA data. For instance, questions not being sufficiently objective but up for interpretation, questions phrased poorly that make it difficult to understand what the question is asking about, *etc.* We discuss these further below.

- **Low agreement due to ambiguity.** One reason why human raters might give different feedback stems from ambiguity and subjectiv-

<sup>14</sup>Qualitative examples have suggested that often this was used to indicate *unanswerable*.

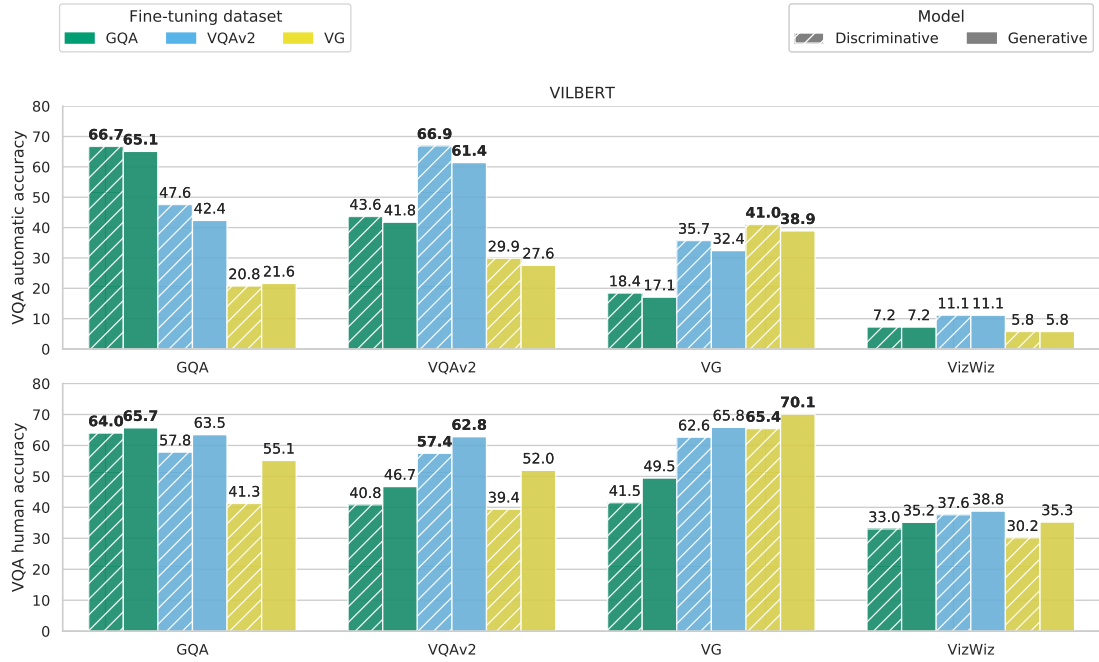


Figure 9: Human evaluation of  $\text{ViLBERT}_{\text{DISC}}$  (shaded bars) and  $\text{ViLBERT}_{\text{GEN}}$  (plain bars) and comparison with automatic evaluation on a random subset (10K) of test questions for each test dataset – GQA, VQAV2, VG and VIZWIZ. Accuracies in bold denote the IID settings. Top: accuracies obtained using automatic evaluation. Bottom: accuracies obtained using human evaluation.

ity around the question and the contents of the image. In these cases it is up to the raters subjective opinion to judge whether the answer is acceptable or not. Find some qualitative examples in Tab. 13.

- **Low agreement due to poor quality question.** Some questions in the original dataset are of rather poor quality which makes it near impossible for the rater to provide a valuable response. Find some qualitative examples of such questions in Tab. 14.

Also surprisingly, from Fig. 6, we see that for models fine-tuned on VQAV2 or GQA and tested on VQAV2, and models fine-tuned on GQA and tested on GQA, human evaluation yields lower accuracy than automatic evaluation! This is not as expected. Upon examining some examples for responses with 100.0 automatic accuracy but marked as *incorrect* by at least 4 human raters, we again find some noise in the ground-truth answers. Tab. 15 shows some examples. Below we report the number of questions where at least 4 human raters voted incorrect even though the automatic metric indicated  $\geq 90.0$  accuracy. Generative case: {GQA  $\rightarrow$  GQA (fine-tuned on GQA, tested on GQA): 86, GQA  $\rightarrow$  VQAV2: 49, VQAV2  $\rightarrow$  VQAV2: 48}, Discriminative case: {GQA  $\rightarrow$  GQA: 128, GQA  $\rightarrow$  VQAV2: 52, VQAV2  $\rightarrow$  VQAV2: 76}.



dataset: VG  
img\_id: 2413078  
q\_id: 151766  
Q: What are they wearing on their heads?  
A: helmet  
GT: helmets  
accuracy: 0.0  
votes: 5 yes, 0 no



dataset: VQAV2  
img\_id: 546983  
q\_id: 546983002  
Q: What is flying in the sky?  
A: kite  
GT: kites  
accuracy: 0.0  
votes: 5 yes, 0 no



dataset: GQA  
img\_id: 2413903  
q\_id: 5199731  
Q: Which kind of device is on the table?  
A: laptop  
GT: cell phone  
accuracy: 0.0  
votes: 5 yes, 0 no

Table 12: A few examples of questions to which the model gave a response that was objectively correct, yet the automatic evaluation metric has marked these data points as 0% accurate. (*votes* here refers to how many raters selected *yes* (i.e. correct) or *no* (i.e. incorrect) when asked about this data point, while GT stands for *ground truth*.)



dataset: VG  
img\_id: 2358330  
q\_id: 700783  
Q: Where is he riding?  
A: park  
GT: in street  
accuracy: 0.0  
votes: 3 yes, 2 no



dataset: VQAV2  
img\_id: 254750  
q\_id: 254750003  
Q: Where is the toilet paper?  
A: bathroom  
GT: on sink  
accuracy: 0.0  
votes: 3 yes, 2 no



dataset: GQA  
img\_id: 2338989  
q\_id: 17319928  
Q: What is on the green sign?  
A: word  
GT: flag  
accuracy: 0.0  
votes: 3 yes, 2 no

Table 13: Low agreement due to ambiguity. In many cases, whether an answer is correct could be up to interpretation. (*votes* here refers to how many raters selected *yes* (i.e. correct) or *no* (i.e. incorrect) when asked about this data point, while GT stands for *ground truth*)





dataset: VG  
img\_id: 2396675  
q\_id: 1453804  
Q: What is the kitchen dresser?  
A: cabinet  
GT: brown  
accuracy: 0.0  
votes: 2 yes, 3 no



dataset: VQAv2  
img\_id: 503518  
q\_id: 503518006  
Q: What is happening?  
A: phone  
GT: watching videos, showing phone  
accuracy: 0.0  
votes: 2 yes, 3 no



dataset: GQA  
img\_id: 2346071  
q\_id: 5863992  
Q: What kind of furniture is playing a game?  
A: table  
GT: couch  
accuracy: 0.0  
votes: 2 yes, 3 no

Table 14: Low agreement due poor question quality. Some questions have poor phrasing that make it difficult to understand what exactly is being asked. In these cases even the humans are not sure what to answer. (*votes* here refers to how many raters selected *yes* (i.e. correct) or *no* (i.e. incorrect) when asked about this data point, while GT stands for *ground truth*)



dataset: VQAv2  
img\_id: 367228  
q\_id: 367228001  
Q: Is the kite flying high enough?  
A: yes  
GT: [no, yes, yes, no, no, no, yes, no, no, yes]  
accuracy: 100.0  
votes: 1 yes, 4 no




dataset: VQAv2  
img\_id: 197745  
q\_id: 197745007  
Q: How many spots are on this animal?  
A: 100  
GT: [70, 100, 100, numerous, 200, 100, 100, 100, 20, lots]  
accuracy: 100.0  
votes: 1 yes, 4 no



dataset: VQAv2  
img\_id: 264737  
q\_id: 264737002  
Q: How many animals are in the picture?  
A: 6  
GT: [7, 6, 6, 9, 6, 6, 6, 7, 7, 6]  
accuracy: 100.0  
votes: 1 yes, 4 no

Table 15: Examples of the few cases where humans considered the response incorrect despite 100.0 automatic accuracy. (*votes* here refers to how many raters selected *yes* (i.e. correct) or *no* (i.e. incorrect) when asked about this data point, while GT stands for *ground truth*)



**Question: What are the people standing on?**

**Response: sand**

Is the displayed **response** a correct answer to the above **question** about the image?

Correct

Incorrect

Figure 10: Sample of the MTurk interface the raters used to annotate data.

### Instructions

You will see a question and a corresponding, short response displayed next to the image. The question and response both concern the contents of the image. Based on the image, is the provided response a correct answer to the question?

**Your task:** Please indicate whether the response is correct by selecting **Correct** or **Incorrect**.

A response should be considered **correct** if:

1.  It is relevant to the question and answers the question in a direct, grammatically correct way, AND
2.  It accurately reflects the contents of the displayed image.

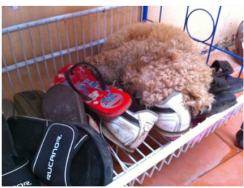
A response should be considered **incorrect** if:

1.  It is not a coherent, natural response to the question, OR
2.  It does not accurately reflect the contents of the displayed image.

**Note:** We expect the responses to be short and concise, not full sentences. Thus a response should **not** be considered incorrect solely on the basis that it is brief.

Please see the examples below to understand the task better:

- **Example 1:**  This response is correct.
 




**Question: What colour is the flip-flop?**

**Response: red**

Is the displayed **response** a correct answer to the above **question** about the image?

Correct

Incorrect
- **Example 2:**  This response is correct.
 




**Question: Is this a cat?**

**Response: no**

Is the displayed **response** a correct answer to the above **question** about the image?

Correct

Incorrect
- **Example 3:**  This response is incorrect. Although the response addresses the question, it does not accurately reflect the image's contents. A correct response would be 'baseball'.
 



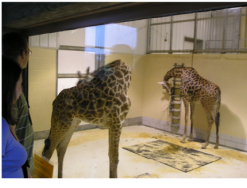
**Question: What are they playing?**

**Response: basketball**

Is the displayed **response** a correct answer to the above **question** about the image?

Correct

Incorrect
- **Example 4:**  This response is incorrect. Although it is true that the scene is at the zoo, the response does not directly address the question. A correct response would be 'no'.
 



**Question: Is this at a museum?**

**Response: at zoo**

Is the displayed **response** a correct answer to the above **question** about the image?

Correct

Incorrect

Figure 11: Instructions given to MTurk raters.