# Learning What Makes a Difference from Counterfactual Examples and Gradient Supervision

Damien Teney, Ehsan Abbasnedjad, Anton van den Hengel

Australian Institute for Machine Learning
University of Adelaide
North Terrace, SA 5005 Adelaide, Australia

**Abstract.** One of the primary challenges limiting the applicability of deep learning is its susceptibility to learning spurious correlations rather than the underlying mechanisms of the task of interest. The resulting failure to generalise cannot be addressed by simply using more data from the same distribution. We propose an auxiliary training objective that improves the generalization capabilities of neural networks by leveraging an overlooked supervisory signal found in existing datasets. We use pairs of minimally-different examples with different labels, a.k.a counterfactual or contrasting examples, which provide a signal indicative of the underlying causal structure of the task. We show that such pairs can be identified in a number of existing datasets in computer vision (visual question answering, multi-label image classification) and natural language processing (sentiment analysis, natural language inference). The new training objective orients the gradient of a model's decision function with pairs of counterfactual examples. Models trained with this technique demonstrate improved performance on out-of-distribution test sets.

## 1 Introduction

Most of today's machine learning methods rely on the assumption that the training and testing data are drawn from a same distribution [69]. One implication is that models are susceptible to poor real-world performance when the test data differs from what is observed during training. This limited capability to generalise partly arises because supervised training essentially amounts to identifying correlations between given examples and their labels. However, correlations can be spurious, in the sense that they may reflect dataset-specific biases or sampling artifacts, rather than intrinsic properties of the task of interest [41,67]. When spurious correlations do not hold in the test data, the model's predictive performance suffers and its output becomes unreliable and unpredictable. For example, an image recognition system may rely on common co-occurrences of objects, such as people together with a dining table, rather visual evidence for each recognized object. This system could then hallucinate people when a table is observed (Fig. 4).
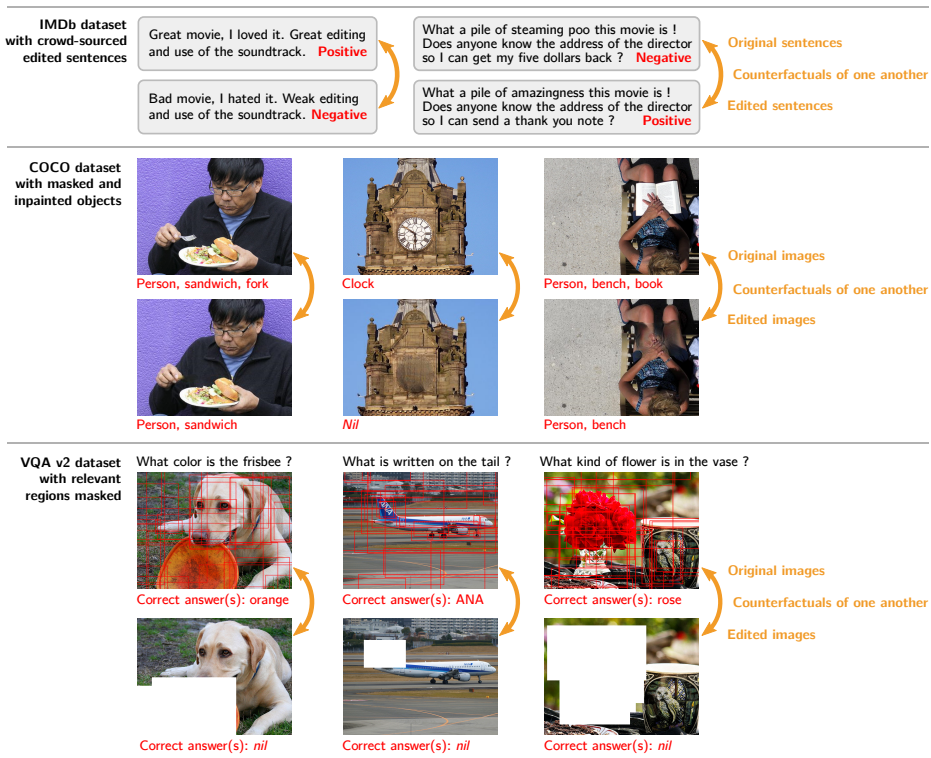
**Fig. 1.** The traditional supervised learning process treats each training example individually. We use counterfactual *relations* between examples as an additional supervisory signal. In some datasets, the relations are provided explicitly, as in sentiment analysis (top) with sentences edited by annotators to flip their label between positive and negative. In other datasets (COCO, middle, and VQA v2, bottom), we show that counterfactual examples can be generated from existing annotations, by masking relevant regions. Our method then leverages the *relations* between the original and edited examples, which proves superior to simple data augmentation.

A model capable of generalization and extrapolation beyond its training distribution should ideally capture the causal mechanisms at play behind the data. Acquiring additional training examples from the same distribution cannot help in this process [49]. Rather, we need either to inject strong prior assumptions in the model, such as inductive biases encoded in the architecture of a neural network, or a different type of training information. Ad hoc methods such as data augmentation and domain randomization fall in the former category, and they only defer the limits of the system by hand-designed rules.

In this paper, we show that many existing datasets contain an overlooked signal that is informative about their causal data-generating process. This information is present in the form of groupings of training examples, and it is often discarded by the shuffling of points occurring during stochastic training. We show that this information can be used to learn a model that is more faith-

ful to the causal model behind the data. This training signal is fundamentally different and complementary to the labels of individual points. We use pairs of minimally-dissimilar, differently-labeled training examples, which we interpret as counterfactuals of one another. In some datasets, such pairs are provided explicitly [7,12,47,48]. In others, they can be identified from existing annotations [18,25,36,61,62,70].

The intuition for our approach is that relations between pairs of counterfactual examples indicate what changes in the input space map to changes in the space of labels. In a classification setting, this serves to constrain the geometry of a model's decision boundary between classes. Loosely speaking, we complement the traditional "curve fitting" to individual training points of standard supervised learning, with "aligning the curve" with pairs of counterfactual training points.

We describe a novel training objective (gradient supervision) and its implementation on various architectures of neural networks. The vector difference in input space between pairs of counterfactual examples serves to supervise the orientation of the gradient of the network. We demonstrate the benefits of the method on four tasks in computer vision and natural language processing (NLP) that are notoriously prone to poor generalization due to dataset biases: visual question answering (VQA), multi-label image classification, sentiment analysis, and natural language inference. We use annotations from existing datasets that are usually disregarded, and we demonstrate significant improvements in generalization to out-of-distribution test sets for all tasks.

In summary, the contributions of this paper are as follows.

1. We propose to use relations between training examples as additional information in the supervised training of neural networks (Section 3.1). We show that they provide a fundamentally different and complementary training signal to the fitting of individual examples, and explain how they improve generalization (Section 3.3).
2. We describe a novel training objective (gradient supervision) to use this information and its implementation on multiple architectures of neural networks (Section 4).
3. We demonstrate that the required annotations are present in a number of existing datasets in computer vision and NLP, although they are usually discarded. We show that our technique brings improvements in out-of-distribution generalization on VQA, multi-label image classification, sentiment analysis, and natural language inference.

## 2   Related work

This work proposes a new training objective that improves the generalization capabilities of models trained with supervision. This touches a number of core concepts in machine learning.

The predictive performance of machine learning models rests on the fundamental assumption of statistical similarity of the distributions of training and

test data. There is a growing interest for evaluating and addressing the limits of this assumption. Evaluation on **out-of-distribution data** is increasingly common in computer vision [2,8,27] and NLP [33,78]. These evaluations have shown that some of the best models can be right for the wrong reasons [2,21,25,67,68]. This happens when they rely on dataset-specific biases and artifacts rather than intrinsic properties of the task of interest. When these biases do not hold in the test data, the predictive performance of the models can drop dramatically [2,8].

When poor generalization is viewed as a deficiency of the training data, it is often referred to as **dataset biases**. They correspond to correlations between inputs and labels in a dataset that can be exploited by a model to exhibit strong performance on a test set containing these same biases, without actually solving the task of interest. Several popular datasets used in vision-and-language [24] and NLP [78] have been shown to exhibit strong biases, leading to an inflated sense of progress on these tasks.

Recent works have discussed generalization from a **causal perspective** [6,29,50,54]. This sheds light on the possible avenues for addressing the issue. In order to generalize perfectly, a model should ideally capture the real-world causal mechanisms at play behind the data. The limits of identifiability of causal models from observational data have been well studied [49]. In particular, additional data from a single biased training distribution can not solve the problem. The alternative options are to use strong assumptions (*e.g.* inductive biases, engineered architectures, hand-designed data augmentations), or additional data, collected in controlled conditions and/or of a different type than labeled examples. This work uses the latter option, using pairings of training examples that represent counterfactuals of one another. Recent works that follow this line include the principle of invariant risk minimization (IRM [6]). IRM uses multiple training environments, *i.e.* non-IID training distributions, to discover generalizable invariances in the data. Teney *et al.* [17] showed that existing datasets could be automatically partitioned to create these environments, and demonstrated improvements in generalization for the task of visual question answering (VQA).

Generalization is also related to the wide area of **domain adaptation** [22]. Our objective in this paper is not to adapt to a particular new domain, but rather to learn a model that generalizes more broadly by using annotations indicative of the causal mechanisms of the task of interest. In domain adaptation, the idea of finding a data representation that is invariant across domains is limiting, because the true causal factors that our model should rely on may differ in their distribution across training domains. We refer the reader to [6] for a formal discussion of these issues.

The growing popularity of high-level tasks in **vision-and-language** [4,5,19] has brought the issue of dataset biases to the forefront. In VQA, language biases cause models to be overly reliant on the presence of particular words in a question. Improving the data collection process can help [24,80] but it only addresses precisely identified biases and confounders. Controlled evaluations for VQA now include out-of-distribution test sets [2,65]. Several models and training methods [11,15,16,26,27,40,52] have been proposed with significant improvements.

They all use strong prior knowledge about the task and/or additional annotations (question types) to improve generalization. Some methods also supervise the model's attention [37,51,57] with ground truth human attention maps [18]. All of these methods are specific to VQA or to captioning [30,37] whereas we describe a much more general approach.

Evaluating generalization overlaps with the growing interest in **adversarial examples** for evaluation [9,13,31,33,42,46,71]. The term has been used to refer both to examples purposefully generated to fool existing models [43,44], but also to hard natural examples that current models struggle with [9,13,71,79]. Our method is most related to the use of these examples for **adversarial training**. Existing methods focus mostly on the generation of these examples then mix them with the original data in a form of data augmentation [34,58,75]. We argue that this shuffling of examples destroys valuable information. In many datasets, we demonstrate that relations between training points contain valuable information. The above methods also aim at improving robustness to targeted adversarial attacks, which often use inputs outside the manifold of natural data. Most of them rely on prior knowledge and unsupervised regularizers [35,32,74] whereas we seek to exploit additional supervision to improve generalization on natural data.

## 3   Proposed approach

We start with an intuitive motivation for our approach, then describe its technical realization. In Section 3.3, we analyze more formally how it can improve generalization. In Section 4, we demonstrate its application to a range of tasks.

### 3.1   Motivation

Training a machine learning model amounts to fitting a function $f(\cdot)$ to a set of labeled points. We consider a binary classification task, in which the model is a neural network $f$ of parameters $\boldsymbol{\theta}$ such that $f_{\boldsymbol{\theta}} : \mathbb{R}^d \rightarrow \{0, 1\}$, and a set of training points[1] $\mathcal{T} = \{(\boldsymbol{x}_i, y_i)\}_i$, with $\boldsymbol{x}_i \in \mathbb{R}^d$ and $y_i \in [0, 1]$. By training the model, we typically optimize $\boldsymbol{\theta}$ such that the output of the network $\tilde{y}_i = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$ minimizes a loss $\mathcal{L}_{\mathrm{Main}}(\tilde{y}_i, y_i)$ on the training points. However, this does not specify the behaviour of $f$ between these points, and the decision boundary could take an arbitrary shape (Fig. 3). The typical practice is to restrain the space of functions $\mathcal{F} \supset f$ (e.g. a particular architecture of neural networks) and of parameters $\Theta \supset \boldsymbol{\theta}$ (e.g. with regularizers [32,35,74]). The capability of the model to interpolate and extrapolate beyond the training points depends on these choices.

Our motivating intuition is that many datasets contain information that is indicative of the shape of an ideal $f$ (in the sense of being faithful to the data-generating process, see Section 3.3) between training points. In particular, we are

---

[1] By *input space*, we refer to a space of feature representations of the input, *i.e.* vector representations ($\boldsymbol{x}$) obtained with a pretrained CNN or text encoder.
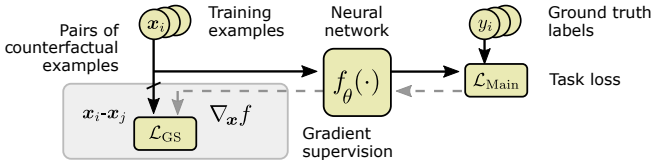
**Fig. 2.** The proposed gradient supervision (GS) is an auxiliary loss on the gradient of a neural network with respect to its inputs, which is simply computed by backpropagation (dashed lines). Supervision for this gradient is generated from pairs of training examples identified as counterfactuals of one another. The loss $\mathcal{L}_{GS}$ is a cosine distance that encourages the gradient of the network to align with the vector between pairs of counterfactual examples.

interested in pairs of training examples that are **counterfactuals** of one another. Given a labeled example $(\boldsymbol{x}_1, y_1)$, we define its counterfactuals as examples such as $(\boldsymbol{x}_2, y_2)$ that represents an alternative premise $\boldsymbol{x}_2$ ("counter to the facts") that lead to different outcome $y_2$. These points represent "minimal changes" ($\|\boldsymbol{x}_1 - \boldsymbol{x}_2\| \ll$, in a semantic sense) such that their label $y_1 \neq y_2$. All possible counterfactuals to a given example $\boldsymbol{x}_1$ constitute a distribution. We assume the availability of samples from it, forming pairs such as $\{(\boldsymbol{x}_1, y_1), (\boldsymbol{x}_2, y_2)\}$. The counterfactual relation is undirected.

**Obtaining pairs of counterfactual examples.**   Some existing datasets explicitly contain pairs of counterfactual examples [7,12,34,47,48]. For example, [34] contains sentences (movie reviews) with positive and negative labels. Annotators were instructed to edit a set of sentences to flip the label, thus creating counterfactual pairs (see examples in Fig. 1). Existing works simply use these as additional training point. Our contribution is to use the *relation* between these pairs, which is usually discarded. In other datasets, counterfactual examples can be created by masking parts of the input, thus creative negative examples. In Section 4, we apply this approach to the COCO and VQA v2 datasets.

## 3.2   Gradient supervision

To exploit relations between counterfactual examples, we introduce an auxiliary loss that supervises the gradient of the network $f_{\boldsymbol{\theta}}$. We denote the gradient of the network with respect to its input at a point $\boldsymbol{x}_i$ with $\boldsymbol{g}_i = \nabla_{\boldsymbol{x}} f(\boldsymbol{x}_i)$. Our new gradient supervision (GS) loss encourages $\boldsymbol{g}_i$ to align with a "ground truth" gradient vector $\hat{\boldsymbol{g}}_i$:

$$\mathcal{L}_{GS}(\boldsymbol{g}_i, \hat{\boldsymbol{g}}_i) \;=\; 1 \;-\; (\boldsymbol{g}_i . \hat{\boldsymbol{g}}_i) \,/\, (\|\boldsymbol{g}_i\| \, \|\hat{\boldsymbol{g}}_i\|) \;. \tag{1}$$

This definition is a cosine distance between $\boldsymbol{g}_i$ and $\hat{\boldsymbol{g}}_i$. Assuming $\{(\boldsymbol{x}_i, y_i), (\boldsymbol{x}_j, y_j)\}$ is a pair of counterfactual examples, a "ground truth" gradient at $\boldsymbol{x}_i$ is obtained as $\hat{\boldsymbol{g}}_i = \boldsymbol{x}_j - \boldsymbol{x}_i$. This represents the translation in the input space that should change the network output from $y_i$ to $y_j$. Minimizing Eq. 1 encourages the network's gradient to align with this vector at the training points. Assuming $f$ is
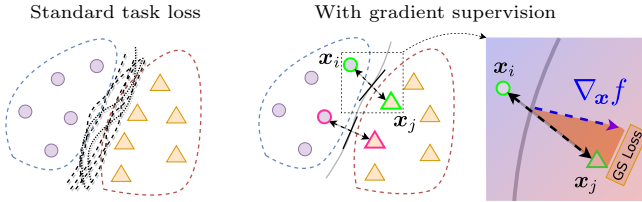
**Fig. 3.** The proposed gradient supervision constrains the geometry of a model's decision boundary between classes. (Left) We show possible decision boundaries consistent with a conventional supervised loss for two classes (circles and triangles representing training points of each). (Right) The gradient supervision uses pairs of counterfactual examples $(\boldsymbol{x}_i, \boldsymbol{x}_j)$ to constrain classifier $f$ such that its local gradient aligns with the vector between these points. When the gradient supervision counterfactuals are included with the GS loss, the boundary is clearer between two classes. On the right, we show the GS loss for a pair of counterfactuals.

continuously differentiable, it also constrains the shape of $f$ *between* training points. This makes $f$ more faithful to the generating process behind the training data (see Section 3.3). Also note that the GS loss uses a local linearization of the network. Although deep networks are highly non-linear globally, first-order approximations have found multiple uses, for example in providing explanations [53,56] and generating adversarial examples [23]. In our application, this approximation is reasonable since pairs of counterfactual examples lie close to one another and to the classification boundary, by definition.

The network is optimized for a combination of the main and GS losses, $\mathcal{L} = \mathcal{L}_{\mathrm{Main}} + \lambda \mathcal{L}_{\mathrm{GS}}$, where $\lambda$ is a scalar hyperparameter. The optimization of the GS loss requires backpropagating second-order derivatives through the network. The computational cost over standard supervised training is of two extra backpropagations through the whole model for each mini-batch.

**Multiclass output.** In cases where the network output $\boldsymbol{y}$ is a vector, a ground truth gradient is only available for classes for which we have positive examples. Denoting such a class $gt$, we apply the GS loss only on the gradient of this class, using $\boldsymbol{g}_i = \nabla_{\boldsymbol{x}} f_i(\boldsymbol{x}_i)$. If a softmax is used, the output for one class depends on that of the others, so the derivative of the network is taken on its logits to make it dependent on one class only.

### 3.3   How gradient supervision improves generalization

By training a machine learning model $f_\theta$, we seek to approximate an ideal $\mathcal{F}$ that represents the real-world process attributing the correct label $y = \mathcal{F}(\boldsymbol{x})$ to any possible input $\boldsymbol{x}$. Let us considering the Taylor expansion of $f$ at a training point $\boldsymbol{x}_j$:

$$f(\boldsymbol{x}_j) \;=\; f(\boldsymbol{x}_i) \;+\; f'(\boldsymbol{x}_i)\,(\boldsymbol{x}_i - \boldsymbol{x}_j) \;+\; \underbrace{\frac{1}{2}\, f''(\boldsymbol{x}_i)\,(\boldsymbol{x}_i - \boldsymbol{x}_j)^2 \;+\; \ldots}_{\approx 0} \quad (2)$$

Our definition of a pair of counterfactual examples $(\boldsymbol{x}_i, \boldsymbol{x}_j)$ (Section 3.1) implies that $(\boldsymbol{x}_i\text{-}\boldsymbol{x}_j)^m$ approaches 0 $(m > 1)$. For such a pair of nearby points, the terms beyond the first order virtually vanish. It follows that the distance between $f(\boldsymbol{x}_j)$ and $f(\boldsymbol{x}_i)$ is maximized when the dot product $\nabla_{\boldsymbol{x}} f(\boldsymbol{x}_i) . (\boldsymbol{x}_i\text{-}\boldsymbol{x}_j)$ is maximum. This is precisely the desired behavior of $f$ in the vicinity of $\boldsymbol{x}_i$ and $\boldsymbol{x}_j$, since their ground truth labels $y_i$ and $y_j$ are different by our definition of counterfactuals. This leads to the definition of the GS loss in Eq. 1. Geometrically, it encourages the gradient of $f$ to align with the vector pointing from a point to its counterfactual, as illustrated in Fig. 3.

The conventional empirical risk minimization with non-convex functions leads to large numbers of local minimas. They correspond to multiple plausible decision boundaries with varied capability for generalization. Our approach essentially modifies the optimization landscape for the parameters $\theta$ of $f$ such that the minimizer found after training is more likely to reflect the ideal $\mathcal{F}$.

## 4   Applications

The proposed method is applicable to datasets with counterfactual examples in the training data. They are sometimes provided explicitly [7,12,34,47,48]. Most interestingly, we show that they can be also be generated from existing annotations [18,25,36,61,62,70].

We selected four classical tasks in vision and language that are notoriously subject to poor generalization due to dataset biases. Our experiments aim (1) to measure the impact of gradient supervision on performance for well-known tasks, and (2) to demonstrate that the necessary annotations are available in a variety of existing datasets. We therefore prioritized the breadth of experiments and the use of simple models (details in supp. mat.) rather than chasing the state of the art on any particular task. The method should readily apply to more complex models for any of these tasks.

### 4.1   Visual question answering

The task of visual question answering (VQA) involves an image and a related question, to which the model must determine the correct answer among a set of approximately 2,000 candidate answers. Models trained on existing datasets (*e.g.* VQA v2 [24]) are notoriously poor at generalization because of dataset biases. These models rely on spurious correlations between the correct answer and certain words in the question. We use the training/test splits of VQA-CP [2] that were manually organized such that the correlation between the questions' prefixes (first few words) and answers differ at training/test time. Most methods evaluated on VQA-CP **use the explicit knowledge of this fact** [2,11,15,26,52,66,17,76] or even of the ground truth set of prefixes, which defeats the purpose of evaluating generalization. As discussed in the introduction, strong background assumptions are one of the two options to improve generalization beyond a set of labels. Our method, however, follows the other

| Test data → | Val. | | | | Test | | | | Test "focused" | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | All | YesNo | Nb | Other | All | YesNo | Nb | Other | All | YesNo | Nb | Other |
| SAN [76] | – | – | – | – | 25.0 | 38.4 | 11.1 | 21.7 | – | – | – | – |
| GVQA [2] | – | – | – | – | 31.3 | 58.0 | 13.7 | 22.1 | – | – | – | – |
| UpDown [64] | – | – | – | – | 39.1 | 62.4 | 15.1 | 34.5 | | | | |
| Ramakrishnan et al., 2018 [52] | – | – | – | – | 42.0 | **65.5** | 15.9 | 36.6 | – | – | – | – |
| Grand and Belinkov, 2019 [26] | – | – | – | – | 42.3 | 59.7 | 14.8 | 40.8 | – | – | – | – |
| RUBi [11] | – | – | – | – | 47.1 | 68.7 | 20.3 | 43.2 | – | – | – | – |
| Teney et al., 2019 [66] | – | – | – | – | 46.0 | 58.2 | **29.5** | 44.3 | – | – | – | – |
| Unshuffling [17] | – | – | – | – | 42.39 | 47.72 | 14.43 | 47.24 | – | – | – | – |
| Strong baseline [17] + CF data | 63.3 | 79.4 | 45.5 | 53.7 | 46.0 | 61.3 | 15.6 | **46.0** | 44.2 | 57.3 | 9.2 | **42.2** |
| **+ CF data + GS** | 62.4 | 77.8 | 43.8 | 53.6 | **46.8** | 64.5 | 15.3 | 45.9 | **46.2** | 63.5 | 10.5 | 41.4 |
| Weak baseline (BUTD [64]), trained on 'Other' only | – | – | – | 54.7 | – | – | – | 43.3 | – | – | – | 40.6 |
| + CF data | – | – | – | 55.9 | – | – | – | 45.0 | – | – | – | 40.6 |
| **+ CF data + GS** | – | – | – | 56.1 | – | – | – | 44.7 | – | – | – | 38.3 |

**Table 1.** Application to VQA-CP v2. Existing methods all rely on built-in knowledge of the construction procedure of the dataset, defeating some of the claimed improvements in robustness. Using counterfactual data with the proposed gradient supervision (GS) improve performance on most question types on the out-of-distribution test sets (see text for discussion).

option of using a different type of data, and **does not rest on the knowledge of the construction of VQA-CP**.

**Generating counterfactual examples.** We build counterfactual examples for VQA-CP using annotations of human attention from [18]. Given a question/image/answer triple $(q, I, a)$, we build its counterfactual counterpart $(q, I', a')$ by editing the image and answer. The image $I$ is a set of features pre-extracted with a bottom-up attention model [3] (typically a matrix of dimensions $N \times 2048$). We build $I'$ ($N' \times 2048$, $N' \leq N$) by masking the features whose bounding boxes overlap with the human attention map past a certain threshold (details in supp. mat.). The vector $a$ is a binary vector of correct answers over all candidates. We simply set all entries in $a'$ to zero.

**Experimental setting.** For training, we use the training split of VQA-CP, minus 8,000 questions held out as an "in-domain" validation set (as in [17]). We generate counterfactual versions of the training examples that have a human attention map (approx. 7% of them). For evaluation, we use (1) our "in-domain" validation set (held out from the training set), (2) the official VQA-CP test set (which has a different correlation between prefixes and answers), and (3) a new *focused* test set.

The *focused* test set contains the questions from VQA-CP test from which we only keep image features of regions looked at by humans to answer the questions. We essentially perform the opposite of the building of counterfactual examples, and mask regions where the human attention is **below** a low threshold. Answering questions from the *focused* test set should intuitively be easier, since the background and distracting image regions have been removed. However, a model that relies on context (question or irrelevant image regions) rather than strictly on the relevant visual evidence will do poorly on the focused test set. This serves

to measure robustness beyond the question biases that VQA-CP was specifically designed for.

**Results**   We present results of our method applied on top of two existing models. The first (*weak baseline*) is the popular BUTD model [3,64]. The second (*strong baseline*) is the "unshuffling" method of [17], which was specifically tuned to address the language biases evaluated with VQA-CP. We compare the baseline model with the same model trained with the additional counterfactual data, and then with the additional GS loss. The performance improves on most question types with each of these additions. The "focused" provides an out-of-distribution evaluation complementary to the VQA-CP test set (which only accounts for language biases). It shows the improvements expected from our method to a larger extent that the VQA-CP test set. This suggests that evaluating generalization in VQA is still not completely addressed with the current benchmarks. Importantly, the improvements over both the weak and strong baselines indicate that **the proposed method is not redundant with existing methods that specifically address the language biases measured by VQA-CP**, like the strong baseline. Additional details are provided in the supplementary material.

## 4.2   Multi-label image classification

We apply our method to the COCO dataset [36]. Its images feature objects from 80 classes. They appear in common situations such that the patterns of co-occurrence are highly predictable: a bicycle often appears together with a person, and a traffic light often appears with cars, for example. These images serve as the basis of a number of benchmarks for image detection [36], captioning [14], visual question answering [5], etc. They all inherit the biases inherent to the COCO images [2,30,72] which is an increasing cause of concern. A method to improve generalization in this context has a wide potential impact.

**Experimental setting.**   We consider a simple multi-label classification task that captures the core issue of dataset biases that affect higher-level tasks (captioning for example [30]). Each image is associated with a binary vector of size 80 that represents the presence of at least one object of the corresponding class in the image. The task is to predict this binary vector. Performance is measured with the mean average precision (mAP) over all classes. The model is a feed-forward neural network that performs an 80-class binary classification with sigmoid outputs, over pre-extracted ResNet-based visual features. We pre-extract these features with the bottom-up attention model of Anderson *et al.* [3]. They are spatially pooled into a 2048-dimensional vector. The model is trained with a standard binary cross-entropy loss (details in the supplementary material).

**Generating counterfactual examples.**   Counterfactual examples can be generated using existing annotation in COCO. Agarwal *et al.* [1] used the inpainter GAN [59] to edit images by masking selected objects. This only requires the original labels and bounding boxes. The edited images represent a "minimal change" that makes the corresponding label negative, which agrees with our definition of

| | COCO Multi-label classification | | |
|---|---|---|---|
| | Original | Edited images | Hard edited |
| Test data → | images | images | images |
| Random predictions (chance) | 5.1 | 3.9 | 7.8 |
| Baseline w/o edited tr. examples | 71.8 | 58.1 | 54.8 |
| Baseline w/ edited tr. examples | 72.1 | 64.0 | 56.0 |
| + **GS, counterfactual relations** | **72.9** | **65.2** | **57.7** |
| + GS, random relations | 71.8 | 63.9 | 56.1 |

**Table 2.** Application to multi-label classification on COCO. We use counterfactual examples generated by masking objects with the inpainter GAN [1,59]. Our method allows to train a model that is less reliant less on common object co-occurrences of the training set. The most striking improvements are measurable with images that feature sets of objects that appear rarely ("Edited") or never ("Hard edited") during training.
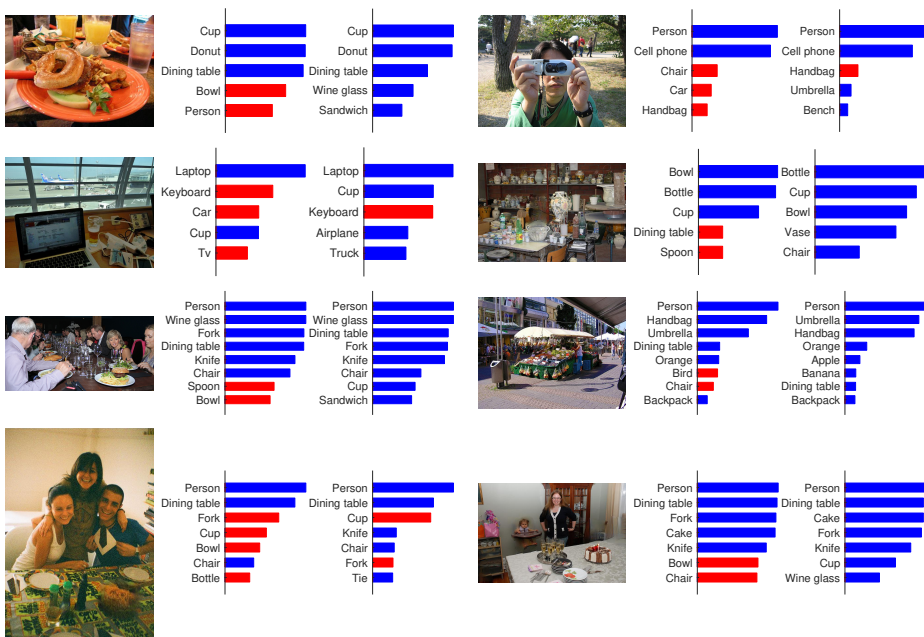


**Fig. 4.** Qualitative examples of multi-label classification on COCO. We show the input image and the scores of the top predicted labels by the baseline and by our method (blue: correct, red: incorrect). The baseline can erroneously predict common co-occurring objects, such as a *person* with food items (top left) even though there is no visual evidence for the former. Our method is better at predicting unusual combinations, such as as a *donut* with a *wineglass* (first row, left) or a *laptop* with an *airplane* (second row, left).

counterfactuals. The vector of ground truth labels for edited images are edited accordingly. For training, we use all images produced by [1] from the COCO *train2014* split (original and edited versions). For evaluation, we use their images from the *val2014* split (original and edited version, evaluated separately).

We also create an additional evaluation split named "Hard edited images". It contains a subset of edited images with patterns of classes that **never** appear in the training set.

**Results.**   We first compare the baseline model trained with the original images only, and then with the original and edited images (Table 2). The performance improves (71.8→72.1%), which is particularly clear when evaluated on edited images (58.1→64.0%). This is because the patterns of co-occurrence in the training data cannot blindly relied on with the edited images. The images in this set depict situations that are unusual in the training set, such as a surfboard without a person on top, or a man on a tennis court who is not holding a racquet. A model that relies on common co-occurrences in the training set rather than strictly on visual evidence can do well on the original images, but not on edited ones. An improvement from additional data is not surprising. It is still worth emphasizing that the edited images were generated "for free" using existing annotations in COCO.

Training the model with the proposed gradient supervision (GS) further improves the precision (72.1→72.9%). This is again more significant on the edited images (64.0→65.2%). The improvement is highest on the set of "hard edited images" (56.0→57.7%). As an ablation, we train the GS model with random pairwise relations instead of relations between counterfactual pairs. The performance is clearly worse, showing that the value of GS is in leveraging an additional training signal, rather than setting arbitrary constraints on the gradient like existing unsupervised regularizers [35,32,74]. In Fig. 4, we provide qualitative examples from the evaluation sets where the predictions of our model improve over the baseline.

### 4.3   NLP Tasks: sentiment analysis and natural language inference

The task of **sentiment analysis** is to assign a positive or negative label to a text snippet, such as a movie or restaurant review. For training, we use the extension of the IMDb dataset [39] of movie reviews by Kaushik *et al.* [34]. They collected counterfactual examples by instructing crowdworkers to edit sentences from the original dataset to flip their label. They showed that a standard model trained on the original data performs poorly when evaluated on edited data, indicating that it relies heavily on dataset biases (*e.g.* the movie genre being predictive of the label). They then used edited data during training (simply mixing it with the original data) and showed much better performance in all evaluation settings, even when controlling for the amount of additional training examples. Our contribution is to use GS to leverage the relations between the pairs of original/edited examples.

The task of **natural language inference (NLI)** is to classify a pair of sentences, named the premise and the hypothesis, into {*entailment, contradiction, neutral*} according to their logical relationship. We use the extension of the SNLI dataset [10] by Kaushik *et al.* [34]. They instructed crowdworkers to edit original examples to change their labels. Each original example is supplemented

| Test data → | IMDb with counterfactuals | | | Zero-shot transfer | | |
|---|---|---|---|---|---|---|
| | Val. | Test original | Test edited | Amazon | Twitter | Yelp |
| Random predictions (chance) | 51.4 | 47.7 | 49.2 | 47.3 | 53.3 | 45.4 |
| Baseline w/o edited tr. data | 71.2 | 82.6 | 55.3 | 78.6 | 61.0 | 82.8 |
| Baseline w/ edited tr. data | 85.7 | 82.0 | 88.7 | 80.8 | 63.1 | 87.4 |
| + **GS**, counterfactual rel. | **89.8** | **83.8** | **91.2** | **81.6** | **65.4** | **88.8** |
| + GS, random relations | 50.8 | 49.2 | 52.0 | 47.4 | 61.2 | 57.4 |

| Test data → | SNLI with counterfactuals | | | Zero-shot transfer |
|---|---|---|---|---|
| | Val. | Test original | Test edited | MultiNLI dev. |
| Random predictions (chance) | 30.8 | 34.6 | 32.9 | 31.9 |
| Baseline w/o edited tr. data | 61.8 | 42.0 | 59.0 | 46.0 |
| Baseline w/ edited tr. data | 61.3 | 39.1 | 57.8 | 42.4 |
| + **GS**, counterfactual relations | **64.8** | **44.4** | **61.2** | **46.8** |
| + GS, random relations | 58.5 | 40.4 | 58.6 | 45.7 |

**Table 3.** Application to sentiment analysis (top) and natural language inference (bottom), trained resp. on the subsets of the IMDb and SNLI datasets augmented with "edited" counterfactual examples [34]. Our technique brings clear improvements over mere data augmentation baseline (accuracy in %), in particular when evaluated on the edited test data (on which biases from the original training data cannot be relied on) and on test data from other datasets (no fine-tuning is used).

with versions produced by editing the premise or the hypothesis, to either of the other two classes. The original and edited data together are therefore four times as large as the original data alone.

**Results on sentiment analysis.**   We first compare a model trained with the original data, and with the original and edited data as simple augmentation (Table 3). The improvement is significant when tested on edited data (55.3→88.7%). We then train the model with our GS loss. The added improvement is visible on both the original data (82.0→83.8%) and on the edited data (88.7→91.2%). The evaluation on edited examples is the more challenging setting, because spurious correlations from the original training data cannot be relied on. The ablation that uses GS with random relations completely fails, confirming the value of the supervision with relations between pairs of related examples.

We additionally evaluate the model on out-of-sample data with three additional test sets: Amazon Reviews [45], Semeval 2017 (Twitter data) [55], and Yelp reviews [77]. The model trained on IMDb is applied **without any fine-tuning** to these, which constitutes a significant challenge in terms of generalization. We observe a clear gain over the data augmentation baseline on all three.

**Results on NLI.**   We perform the same set of experiments on NLI. The fairest point of comparison is again the model trained with the original and edited data. Using the GS loss on top of it brings again a clear improvement (Table 3), both when evaluated on standard test data and on edited examples. As an additional measure of generalization, we also evaluate the same models on the dev. set of MultiNLI [73] without any fine-tuning. There is a significant domain shift
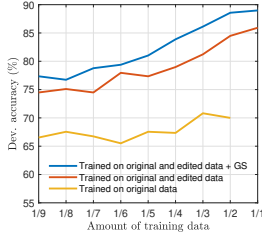
**Fig. 5.** Training on different amounts of data for sentiment analysis. The proposed gradient supervision (GS) brings a clear improvement over a model trained with edited examples for simple data augmentation, and even more so over a model trained with the same number of original training examples.

between the datasets. Using the edited examples for data augmentation actually hurts the performance here, most likely because they constitute very "unnatural" sentences, such that easy-to-pick-up language cues cannot be relied on. Using GS (with always uses the edited data as augmentations as well) brings back the performance higher, and above the baseline trained only on the original data.

**Limitations.**   Our NLP experiments were conducted with simple models and relatively little data. The current state of the art in sentiment analysis and NLI is achieved by transformer-based models [20] trained on vastly more data. Kaushik *et al.* [34] showed that counterfactual examples are much more valuable than the same amount of standard data, including for fine-tuning a BERT model for NLI. The application of our technique to the extremely-large data regime, including with large-scale language models, is an exciting direction for future work.

## 5    Conclusions

We proposed a new training objective that improves the generalization capabilities of neural networks by supervising their gradient, and using an unused training signal found in many datasets. While most machine learning models rely on identifying correlations between inputs and output, we showed that relations between counterfactual examples provide a fundamentally different, complementary type of information. We showed theoretically and empirically that our technique can shape the decision boundary of the model to be more faithful to the causal mechanisms that generated the data. Practically speaking, the model is then more likely to be "right for the right reasons". We showed that this effect brings significant improvements on a number of tasks when evaluated with out-of-distribution test data. We demonstrated that the required annotations can be extracted from existing datasets for a number of tasks.

There is a number of additional tasks and datasets on which our method can readily apply [7,12,47,48,61,62,70]. Scaling up the technique to state-of-the-art models in vision and NLP [20,28,38,60,63] is another exciting direction for future work.

# Supplementary material

## A      Application to VQA

**Data.**   We generate the counterfactual examples by masking image features on-the-fly, during training, according the the human attention maps of [18]. We use image features from [3], which correspond to bounding boxes in the image. We mask the features whose boxes overlap with a fraction of the the human attention map above a fixed threshold. We use the precomputed overlap score from [57], which is a scalar in $[0, 1]$, and set the threshold at 0.2 (setting it at 0 would mask the occasional boxes that encompass nearly the whole image, which is not desirable). This value was set manually by verifying for the intended effect on a few training examples (that is, masking most of the relevant visual evidence). See Fig. 6 for examples of original questions and their counterfactual versions.

**Experimental setting.**   Our experiments use a validation set (8,000 questions chosen at random) held out from the original VQA-CP training set. Note that most existing methods evaluated in VQA-CP use the extremely unsanitary practice of using the VQA-CP test split for model selection. This is extremely concerning since the whole purpose of VQA-CP is to evaluate generalization to an out-of-distribution test set. The variance in evaluating the 'number' and 'yes/no' questions is moreover extremely high, because the number of reasonable answers on each of these types is very limited. For example, a model that answers *yes* or *no* at random, or produces constantly either answer, can fare extremely well (upwards of 62% accuracy) on these questions. This can very well result from a buggy implementation or a "lucky" random seed, identified by model selection on the test set (!). This is the reason why we include an evaluation on the 'other' type of questions in isolation. All of these issues have been pointed out by a few authors [26,16,66].

Our *focused* test set is a subset of the official VQA-CP test set. It is created in a similar manner as the counterfactual examples. We mask features that overlap with human attention maps *below* (instead of above) a threshold of 0.8. This value was set manually by verifying for the intended effect on a few examples (masking the background but not the regions necessary to answer the question). The *focused* test set is much smaller than the official test set since it only comprises questions for which a human attention map is available.

**Models.**   Our baseline model follows the general description of Teney *et al.* [64]. We use the features of size $36 \times 2048$ provided by Anderson *et al.* [3]. Our 'strong baseline' uses the additional procedure described in [17] on top of this baseline, using the code provided by the authors.

**Existing methods.**   The method presented in [57] could have constituted an ideal point of comparison with ours, as it was evaluated on VQA-CP and used human attention maps. However, after extensive discussions with the authors,

What is floating in the sky ?
GT Answer(s): kites, kite, sail

Where is the woman sitting ?
GT Answer(s): stairs, steps

What team is the batter on ?
GT Answer(s): white, yankees, mets, giants

Where is the baby looking ?
GT Answer(s): laptop, screen, monitor

What is the sex of rider ?
GT Answer(s): female, male

What kind of boat is on the water ?
GT Answer(s): canoe, paddle

What sport is the person participating in ?
GT Answer(s): surfing

What is this person standing on ?
GT Answer(s): skateboard

What is the person in photo holding ?
GT Answer(s): surfboard

**Fig. 6.** Application to VQA. Examples of original examples (with their ground truth answer) and their counterfactual version. Red boxes indicate regions that were candidates for masking when generating the counterfactual versions.

we still have not been able to replicate any of the performance claimed in the paper. We found a number of errors in the paper, as well as inconsistencies in the reported results, and an extreme sensitivity to a single hyperparameter (their reported results were obtained with a single run on a single random seed). We chose not to mention this work in our main paper until these issues have been resolved.

**Why not use the same technique for the VQA and COCO experiments ? Inpainting in pixel space *vs* masking image features.**  The two approaches are applicable in both cases. The only reason was to showcase the use of multiple techniques to generate counterfactual examples. The human attention map are specific to VQA and not applicable to the COCO experiments.

# B    Application to image classification with COCO

**Data.**    We use the edited images released by [1] together with the corresponding original images from COCO. The edited images were created with the inpainter GAN [59] to mask ground truth bounding boxes of specific objects. The images come from the COCO splits *train2014* and *val2014*. We keep this separation for our experiments as follows. Images from *train2014* (323,116 counting original and edited ones) are used for training, except a random subset (1,000 images) that we hold out for validation (model selection, early stopping). Images from *val2014* (3,361 original and 3,361 edited) are used exclusively for testing.

We identified a subset (named *Hard edited*) of the edited images from *val2014* whose ground truth vector (which indicated the classes appearing in the image) is never seen during training (614 images).

The set of edited images provided by [1] is a non-standard subset of COCO, so no directly-comparable results have been published for the multi-label classification task that we consider.

**Model.**    We pre-extract image features from all images with the ResNet-based, bottom-up attention model [3]. These features are averaged across spatial locations, giving a single vector of dimensions 2048 to represent each image. Our model is a 3-layer ReLU MLP of size 64, followed by a linear/sigmoid output layer of size 80 (corresponding to the 80 COCO classes). This baseline model was first tuned for best performance on the validation set (tuning the number of a layers and their size, the batch size, and learning rate), before adding the proposed GS loss. The model is optimized with AdaDelta, mini-batches of size 512, and a binary cross-entropy loss.

Performance is measured with a standard mean average precision (mAP) (as defined in the Pascal VOC challenge) over all 80 classes.

The Fig. 4 in the paper shows the input image with the scores of the top-$k$ predicted labels by the baseline and by our method. The $k$ corresponds to the number of ground truth labels of each image.

Masked object: **car**
(left and right, behind the truck)

Masked object: **person**

Masked object: **skateboard**

Masked object: **surfboard**

Masked object: **boat**

Masked object: **tie**
(on both persons in the foreground)

Masked object: **bicycle**
(against the railing on the right)

Masked object: **person**

Masked object: **horse**

Masked object: **tie**

**Fig. 7.** Application to multi-label image classification with COCO. Examples of original and edited images.

**Random baseline.**  In our ablations, this model is identical to the standard baseline, but it is trained with a randomly shuffled training set. We shuffle the inputs $\{x_i\}_i$ and the ground truth labels $\{y_i\}_i$ of all training examples. The model is thus not getting any relevant training signal from any example. It can only leverage static dataset biases (*i.e.* a class imbalance).

## C    Application to NLP tasks

**Sentiment analysis data.**  We use the subset of the IMDb dataset [39] for which Kaushik *et al.* [34] obtained counterfactual examples. We use their 'paired' version of the data, which only contains original examples that do have an edited version. For **training**, we use the 'train' split of original and edited data (3414 examples). For **validation** (model selection, early stopping), we use the 'dev' set of paired examples. For **testing**, we use the 'test' split, reporting accuracy over the original and edited examples separately. For testing on other datasets, we use a random subset (2000 examples) of the test sets of Amazon Reviews [45], Semeval 2017 (Twitter data) [55], and Yelp reviews [77] similarly to [34].

**Sentiment analysis model.**  We first optimized a simple baseline model on the validation set (tuning the number of a layers, embedding sizes, batch size, and learning rate). We then added the proposed gradient supervision, tuned its hyperparameters on the validation set (regularizer weight) then reported the performance on the test sets at the epoch of best performance on the validation set. The sentences are tokenized and trimmed to a maximum of 32 tokens. The model encodes a sentence as a bag of words, using word embeddings of size 50, averaged to the exact length of each sentence (*i.e.* not including the padding of the shorter sentences). The vocabulary is limited to the 20,000 most frequent words in the dataset. The averaged vector is passed to a simple linear classifier with a sigmoid output. All weights, including word embeddings, are initialized from random values, and optimized with AdaDelta, in mini-batches of size 32, with a binary cross-entropy loss. The best weight for the GS regularizer was found to be $\lambda=20$. To reduce the noise in the evaluation due to the small size of the training set, we use an ensemble of 6 identical models trained in parallel. The reported results uses the output of the ensemble, that is the average of the logits of the 6 models.

**NLI data.**  The experiments on NLI follow a similar procedure to those on sentiment analysis. We use the subset of the SNLI dataset [10] for which Kaushik *et al.* [34] collected counterfactual examples. We use their biggest version of the data, named 'all combined', that contains counterfactual examples with edited premises and edited hypotheses. For testing, we evaluate accuracy separately on original and edited examples (edited premises and edited hypotheses combined). For testing transfer, we use the 'dev' set of MultiNLI [73]. Whereas the SNLI dataset contains sentence pairs derived from image captions, MultiNLI is more diverse. It contains sentences from transcribed speech, popular fiction, and government reports. Compared to SNLI, it contains more linguistic diversity and complexity.

| Test data → | Yelp |
|---|---|
| Random predictions (chance) | 45.4 |
| Baseline w/o edited tr. data | 82.8 |
| Baseline w/ edited tr. data | 87.4 |
| + **GS**, counterfactual rel. | **88.8** |
| + GS, random relations | 57.4 |

**Table 4.** Application to sentiment analysis. Results on the Yelp dataset. This column was missing in Table 3 in the paper (a code-generating error replicated the values from the *Amazon* column into the *Yelp* column).

**NLI model.**    The premise and hypothesis sentences are tokenized and trimmed to a maximum of 32 tokens. They are encoded separately as bags of words, using frozen Glove embeddings (dimension 300), then a learned linear/ReLU projection to dimension 50, and an average to the length of each sentence (without using the padding). They are then passed through a batch normalization layer, then concatenated, giving a vector of size 100. The vector is passed through 3 linear/ReLU layers, then a final linear/sigmoid output layer. The model is trained with AdaDelta, with mini-batches of size 512, and a binary cross-entropy loss. The best weight for the GS regularizer was found to be $\lambda=0.01$. Similarly to our experiments on sentiment analysis, we evaluate an ensemble of 6 copies of the model described above.

# References

1. Agarwal, V., Shetty, R., Fritz, M.: Towards causal vqa: Revealing and reducing spurious correlations by invariant and covariant semantic editing. arXiv preprint arXiv:1912.07538 (2019)
2. Agrawal, A., Batra, D., Parikh, D., Kembhavi, A.: Don't just assume; look and answer: Overcoming priors for visual question answering. In: Proc. IEEE Conf. Comp. Vis. Patt. Recogn. pp. 4971–4980 (2018)
3. Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., Zhang, L.: Bottom-up and top-down attention for image captioning and vqa. CVPR (2018)
4. Anderson, P., Wu, Q., Teney, D., Bruce, J., Johnson, M., Sünderhauf, N., Reid, I., Gould, S., van den Hengel, A.: Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3674–3683 (2018)
5. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh, D.: VQA: Visual Question Answering. In: Proc. IEEE Int. Conf. Comp. Vis. (2015)
6. Arjovsky, M., Bottou, L., Gulrajani, I., Lopez-Paz, D.: Invariant risk minimization. arXiv preprint arXiv:1907.02893 (2019)
7. Baradel, F., Neverova, N., Mille, J., Mori, G., Wolf, C.: Cophy: Counterfactual learning of physical dynamics. arXiv preprint arXiv:1909.12000 (2019)
8. Barbu, A., Mayo, D., Alverio, J., Luo, W., Wang, C., Gutfreund, D., Tenenbaum, J., Katz, B.: Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. In: Proc. Advances in Neural Inf. Process. Syst. pp. 9448–9458 (2019)

9. Bartolo, M., Roberts, A., Welbl, J., Riedel, S., Stenetorp, P.: Beat the ai: Investigating adversarial human annotations for reading comprehension. arXiv preprint arXiv:2002.00293 (2020)

10. Bowman, S.R., Angeli, G., Potts, C., Manning, C.D.: A large annotated corpus for learning natural language inference. In: Proc. Conf. Empirical Methods in Natural Language Processing (2015)

11. Cadene, R., Dancette, C., Ben-younes, H., Cord, M., Parikh, D.: Rubi: Reducing unimodal biases in visual question answering. arXiv preprint arXiv:1906.10169 (2019)

12. Camburu, O.M., Rocktäschel, T., Lukasiewicz, T., Blunsom, P.: e-snli: Natural language inference with natural language explanations. In: Proc. Advances in Neural Inf. Process. Syst. pp. 9539–9549 (2018)

13. Chen, M., DArcy, M., Liu, A., Fernandez, J., Downey, D.: Codah: An adversarially-authored question answering dataset for common sense. In: Proceedings of the 3rd Workshop on Evaluating Vector Space Representations for NLP. pp. 63–69 (2019)

14. Chen, X., Fang, H., Lin, T.Y., Vedantam, R., Gupta, S., Dollar, P., Zitnick, C.L.: Microsoft COCO captions: Data collection and evaluation server. arXiv preprint arXiv:1504.00325 (2015)

15. Clark, C., Yatskar, M., Zettlemoyer, L.: Don't take the easy way out: Ensemble based methods for avoiding known dataset biases. arXiv preprint arXiv:1909.03683 (2019)

16. Damien Teney, Ehsan Abbasnejad, A.v.d.H.: On incorporating semantic prior knowledge in deep learning through embedding-space constraints. arXiv preprint arXiv:1909.13471 (2019)

17. Damien Teney, Ehsan Abbasnejad, A.v.d.H.: Unshuffling data for improved generalization. arXiv preprint arXiv:2002.11894 (2020)

18. Das, A., Agrawal, H., Zitnick, C.L., Parikh, D., Batra, D.: Human attention in visual question answering: Do humans and deep networks look at the same regions? In: Proc. Conf. Empirical Methods in Natural Language Processing (2016)

19. Das, A., Kottur, S., Gupta, K., Singh, A., Yadav, D., Moura, J.M., Parikh, D., Batra, D.: Visual Dialog. In: CVPR (2017)

20. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)

21. Feng, S., Wallace, E., Boyd-Graber, J.: Misleading failures of partial-input baselines. arXiv preprint arXiv:1905.05778 (2019)

22. Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., Lempitsky, V.: Domain-adversarial training of neural networks. The Journal of Machine Learning Research (2016)

23. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572 (2014)

24. Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., Parikh, D.: Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. arXiv preprint arXiv:1612.00837 (2016)

25. Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., Parikh, D.: 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. pp. 6904–6913 (2017)

26. Grand, G., Belinkov, Y.: Adversarial regularization for visual question answering: Strengths, shortcomings, and side effects. arXiv preprint arXiv:1906.08430 (2019)

27. Guo, Y., Cheng, Z., Nie, L., Liu, Y., Wang, Y., Kankanhalli, M.: Quantifying and alleviating the language prior problem in visual question answering. arXiv preprint arXiv:1905.04877 (2019)

28. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proc. IEEE Conf. Comp. Vis. Patt. Recogn. (2016)

29. Heinze-Deml, C., Meinshausen, N.: Conditional variance penalties and domain shift robustness. arXiv preprint arXiv:1710.11469 (2017)

30. Hendricks, L.A., Burns, K., Saenko, K., Darrell, T., Rohrbach, A.: Women also snowboard: Overcoming bias in captioning models. In: European Conference on Computer Vision. pp. 793–811. Springer (2018)

31. Iyyer, M., Wieting, J., Gimpel, K., Zettlemoyer, L.: Adversarial example generation with syntactically controlled paraphrase networks. arXiv preprint arXiv:1804.06059 (2018)

32. Jakubovitz, D., Giryes, R.: Improving dnn robustness to adversarial attacks using jacobian regularization. In: Proc. Eur. Conf. Comp. Vis. pp. 514–529 (2018)

33. Jia, R., Liang, P.: Adversarial examples for evaluating reading comprehension systems. arXiv preprint arXiv:1707.07328 (2017)

34. Kaushik, D., Hovy, E., Lipton, Z.C.: Learning the difference that makes a difference with counterfactually-augmented data. arXiv preprint arXiv:1909.12434 (2019)

35. Li, Y., Cohn, T., Baldwin, T.: Learning robust representations of text. In: Proc. Conf. Empirical Methods in Natural Language Processing. pp. 1979–1985. Association for Computational Linguistics (2016)

36. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: Common objects in context. In: Proc. Eur. Conf. Comp. Vis. (2014)

37. Liu, C., Mao, J., Sha, F., Yuille, A.: Attention correctness in neural image captioning. In: Proc. Conf. AAAI (2017)

38. Liu, Y., Wang, Y., Wang, S., Liang, T., Zhao, Q., Tang, Z., Ling, H.: Cbnet: A novel composite backbone network architecture for object detection. arXiv preprint arXiv:1909.03625 (2019)

39. Maas, A.L., Daly, R.E., Pham, P.T., Huang, D., Ng, A.Y., Potts, C.: Learning word vectors for sentiment analysis. In: Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies-volume 1. pp. 142–150. Association for Computational Linguistics (2011)

40. Mahabadi, R.K., Henderson, J.: Simple but effective techniques to reduce biases. arXiv preprint arXiv:1909.06321 (2019)

41. Mitchell, T.M.: The need for biases in learning generalizations. Department of Computer Science, Laboratory for Computer Science Research  (1980)

42. Miyato, T., Dai, A.M., Goodfellow, I.: Adversarial training methods for semi-supervised text classification. arXiv preprint arXiv:1605.07725 (2016)

43. Moosavi-Dezfooli, S.M., Fawzi, A., Fawzi, O., Frossard, P.: Universal adversarial perturbations. In: Proc. IEEE Conf. Comp. Vis. Patt. Recogn. pp. 1765–1773 (2017)

44. Nguyen, A., Yosinski, J., Clune, J.: Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In: Proc. IEEE Conf. Comp. Vis. Patt. Recogn. pp. 427–436 (2015)

45. Ni, J., Li, J., McAuley, J.: Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In: Proc. Conf. Empirical Methods in Natural Language Processing. pp. 188–197. Association for Computational Linguistics

46. Nie, Y., Williams, A., Dinan, E., Bansal, M., Weston, J., Kiela, D.: Adversarial nli: A new benchmark for natural language understanding. arXiv preprint arXiv:1910.14599 (2019)
47. Park, D.H., Darrell, T., Rohrbach, A.: Robust change captioning. In: Proc. IEEE Conf. Comp. Vis. Patt. Recogn. pp. 4624–4633 (2019)
48. Park, D.H., Darrell, T., Rohrbach, A.: Viewpoint invariant change captioning. arXiv preprint arXiv:1901.02527 (2019)
49. Pearl, J.: Causality: models, reasoning and inference, vol. 29. Springer (2000)
50. Peters, J., Bühlmann, P., Meinshausen, N.: Causal inference by using invariant prediction: identification and confidence intervals. Journal of the Royal Statistical Society: Series B (Statistical Methodology) (2016)
51. Qiao, T., Dong, J., Xu, D.: Exploring human-like attention supervision in visual question answering. In: Proc. Conf. AAAI (2018)
52. Ramakrishnan, S., Agrawal, A., Lee, S.: Overcoming language priors in visual question answering with adversarial regularization. In: Proc. Advances in Neural Inf. Process. Syst. pp. 1541–1551 (2018)
53. Ribeiro, M.T., Singh, S., Guestrin, C.: "why should i trust you ?" explaining the predictions of any classifier. In: Proc. ACM SIGKDD Int. Conf. Knowledge discovery & data mining (2016)
54. Rojas-Carulla, M., Schölkopf, B., Turner, R., Peters, J.: Invariant models for causal transfer learning. The Journal of Machine Learning Research (2018)
55. Rosenthal, S., Farra, N., Nakov, P.: SemEval-2017 task 4: Sentiment analysis in twitter. In: Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017). pp. 502–518. Association for Computational Linguistics (2017)
56. Selvaraju, R.R., Das, A., Vedantam, R., Cogswell, M., Parikh, D., Batra, D.: Gradcam: Why did you say that? arXiv preprint arXiv:1611.07450 (2016)
57. Selvaraju, R.R., Lee, S., Shen, Y., Jin, H., Ghosh, S., Heck, L., Batra, D., Parikh, D.: Taking a hint: Leveraging explanations to make vision and language models more grounded. In: Proc. IEEE Int. Conf. Comp. Vis. (2019)
58. Shafahi, A., Najibi, M., Ghiasi, M.A., Xu, Z., Dickerson, J., Studer, C., Davis, L.S., Taylor, G., Goldstein, T.: Adversarial training for free! In: Proc. Advances in Neural Inf. Process. Syst. pp. 3353–3364 (2019)
59. Shetty, R.R., Fritz, M., Schiele, B.: Adversarial scene editing: Automatic object removal from weak supervision. In: Proc. Advances in Neural Inf. Process. Syst. pp. 7706–7716 (2018)
60. Su, W., Zhu, X., Cao, Y., Li, B., Lu, L., Wei, F., Dai, J.: Vl-bert: Pre-training of generic visual-linguistic representations. arXiv preprint arXiv:1908.08530 (2019)
61. Suhr, A., Lewis, M., Yeh, J., Artzi, Y.: A corpus of natural language for visual reasoning. In: Proc. Conf. Association for Computational Linguistics. vol. 2, pp. 217–223 (2017)
62. Suhr, A., Zhou, S., Zhang, A., Zhang, I., Bai, H., Artzi, Y.: A corpus for reasoning about natural language grounded in photographs. arXiv preprint arXiv:1811.00491 (2018)
63. Tan, H., Bansal, M.: Lxmert: Learning cross-modality encoder representations from transformers. arXiv preprint arXiv:1908.07490 (2019)
64. Teney, D., Anderson, P., He, X., van den Hengel, A.: Tips and tricks for visual question answering: Learnings from the 2017 challenge. CVPR (2018)
65. Teney, D., van den Hengel, A.: Zero-shot visual question answering. arXiv preprint arXiv:1611.05546 (2016)

66. Teney, D., van den Hengel, A.: Actively seeking and learning from live data. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
67. Torralba, A., Efros, A.A., et al.: Unbiased look at dataset bias. In: CVPR. vol. 1, p. 7 (2011)
68. Vapnik, V., Izmailov, R.: Rethinking statistical learning theory: learning using statistical invariants. Machine Learning (2019)
69. Vapnik, V.N.: An overview of statistical learning theory. IEEE Transactions on Neural Networks **10**(5), 988–999 (1999)
70. Vo, N., Jiang, L., Sun, C., Murphy, K., Li, L.J., Fei-Fei, L., Hays, J.: Composing text and image for image retrieval-an empirical odyssey. In: Proc. IEEE Conf. Comp. Vis. Patt. Recogn. pp. 6439–6448 (2019)
71. Wallace, E., Boyd-Graber, J.: Trick me if you can: Adversarial writing of trivia challenge questions. In: ACL Student Research Workshop (2018)
72. Wang, T., Zhao, J., Yatskar, M., Chang, K.W., Ordonez, V.: Balanced datasets are not enough: Estimating and mitigating gender bias in deep image representations. In: Proc. IEEE Int. Conf. Comp. Vis. pp. 5310–5319 (2019)
73. Williams, A., Nangia, N., Bowman, S.R.: A broad-coverage challenge corpus for sentence understanding through inference. arXiv preprint arXiv:1704.05426 (2017)
74. Woods, W., Chen, J., Teuscher, C.: Adversarial explanations for understanding image classification decisions and improved neural network robustness. Nature Machine Intelligence **1**(11), 508–516 (2019)
75. Xie, C., Tan, M., Gong, B., Wang, J., Yuille, A., Le, Q.V.: Adversarial examples improve image recognition. arXiv preprint arXiv:1911.09665 (2019)
76. Yang, Z., He, X., Gao, J., Deng, L., Smola, A.: Stacked Attention Networks for Image Question Answering. In: Proc. IEEE Conf. Comp. Vis. Patt. Recogn. (2016)
77. Yelp: Yelp dataset challenge. `http://www.yelp.com/dataset_challenge`
78. Zellers, R., Bisk, Y., Schwartz, R., Choi, Y.: Swag: A large-scale adversarial dataset for grounded commonsense inference. arXiv preprint arXiv:1808.05326 (2018)
79. Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., Choi, Y.: Hellaswag: Can a machine really finish your sentence? arXiv preprint arXiv:1905.07830 (2019)
80. Zhang, P., Goyal, Y., Summers-Stay, D., Batra, D., Parikh, D.: Yin and yang: Balancing and answering binary visual questions. In: Proc. IEEE Conf. Comp. Vis. Patt. Recogn. (2016)