

VALSE: A Task-Independent Benchmark for Vision and Language Models Centered on Linguistic Phenomena

Anonymous ACL submission

Abstract

We propose VALSE (Vision And Language Structured Evaluation), a novel benchmark designed for testing general-purpose pretrained vision and language (V&L) models for specific *visio-linguistic grounding* capabilities. Currently, V&L models are evaluated on *tasks* such as visual question answering or visual reasoning, which do not address their fine-grained linguistic capabilities. VALSE addresses this gap by offering a suite of six tests targeting *specific linguistic phenomena*. Solving these tests requires models to ground these phenomena in the visual modality, allowing more fine-grained evaluations than hitherto possible. We build VALSE using methods that support the construction of *reliable* foils, and report results from evaluating five widely-used V&L models. Our experiments suggest that current models have considerable difficulty addressing most phenomena. Hence, we expect VALSE to serve as an important benchmark to measure future progress of pretrained V&L models from a *linguistic perspective*, complementing the canonical task-centred V&L evaluations.

1 Introduction

Recently, general-purpose pretrained vision and language (V&L) models have gained notable performance on all V&L tasks they are finetuned on, e.g. visual question answering (VQA), visual commonsense reasoning, phrase grounding or image retrieval (Lu et al., 2019; Tan and Bansal, 2019; Li et al., 2019; Chen et al., 2020; Li et al., 2020a; Su et al., 2020). As a result, the focus of V&L research has broadened beyond neural architectures designed for specific tasks, to large V&L models that are fine-tuned on several V&L tasks.

Current benchmarks give a good perspective on model performance on a wide range of V&L tasks (Cao et al., 2020; Lourie et al., 2021; Li et al., 2021), but the field is only starting to assess *why* models perform so well and whether models learn

specific capabilities that span multiple V&L tasks. In particular, we currently lack understanding of the extent to which such models are able to ground specific linguistic phenomena—at the level of morphosyntax and semantics—in the visual modality (Bernardi and Pezzelle, 2021).

In this paper, we address this gap with VALSE (Vision And Language Structured Evaluation): a benchmark for V&L model evaluation made up of six different tasks, or ‘pieces’. Each piece has the same structure: Given a visual input, a V&L model is required to distinguish real captions from *foils*, where a foil is constructed from a caption by altering a word or phrase corresponding to a *specific linguistic phenomenon*, for example semantic number in noun phrases; verb argument structure; discourse-level coreference, etc. VALSE uses a resource-lean diagnostic setup that does not require large-scale annotation (e.g., of bounding boxes), and builds on existing high-quality image captioning and VQA data. VALSE is designed to leverage the existing prediction heads in pretrained (or finetuned) V&L models; for that reason, our benchmark does not include any re-training and can be interpreted as a *zero-shot* evaluation. We build *test* data for each piece so as to safeguard against the possibility of models exploiting artefacts or statistical biases in the data, a well-known issue with highly parameterised neural models pretrained on large amounts of data (Goyal et al., 2017; Madhyastha et al., 2018; Kafle et al., 2019). With this in view, we propose novel methods to guard against the emergence of *artefacts* during foiling.

Our main contributions are:

- i) We introduce VALSE, a novel benchmark aimed at testing the multimodal capacities of pre-trained V&L models by gauging their sensitivity to *foiled* instances.
- ii) We cover a wide spectrum of basic linguistic phenomena affecting the linguistic *and* visual

- modalities: existence, plurality, counting, spatial relations, actions, and entity coreference.
- iii) We investigate novel strategies to build *valid* and *reliable* foils that include automatic and human validation. We balance the *word frequency distributions* between caption and foil data, and test against the capabilities of pre-trained models to solve the benchmark *unimodally*. We make use of *masked language modeling* (MLM) predictions in foil creation and *semantic inference* predictions for validating foils, and finally collect *human annotations* for the entire benchmark.
 - iv) We establish initial experimental results using a variety of publicly available pretrained V&L models with diverse architectures. The overall weak performance of V&L models on VALSE indicates that time is ripe for a more detailed and reliable foiling dataset targeted at the visual grounding capabilities of V&L models through the lens of linguistic constructs.¹

2 Background and Related work

Pretrained V&L models learn to combine vision and language through self-supervised multitask learning. Tasks include *multimodal masked modeling*—where words in the text and object labels or regions in the image are masked out, then predicted—and *image-sentence alignment*, whereby a model learns to predict whether an image and a text correspond to each other. Major architectures are single- and dual-stream multimodal transformers: *single-stream models* concatenate word and image features, and encode the resulting sequence with a single transformer stack; *dual-stream models* use distinct transformer stacks to handle visual and textual inputs, and additional layers (e.g. co-attention) to fuse these into multimodal features.

Benchmarking V&L models V&L models (Li et al., 2019; Lu et al., 2019; Tan and Bansal, 2019; Lu et al., 2020; Li et al., 2020b; Kim et al., 2021) are commonly evaluated on V&L *tasks* such as VQA (Goyal et al., 2017), visual reasoning (Suhr et al., 2019), or image retrieval (Lin et al., 2014; Plummer et al., 2015).

Given how well transformer-based models perform across unimodal and multimodal tasks, research efforts have recently started to address what makes them so effective, and to what extent they

learn generalisable representations. Techniques to address these questions in unimodal and multimodal V&L contexts include: adversarial examples (Jia and Liang, 2017; Jia et al., 2019); investigation of the impact of bias, be it linguistic (Gururangan et al., 2018), visual semantic (Agarwal et al., 2020), or socio-economic (Garg et al., 2019); and the use of linguistically-informed counterfactual and minimally-edited examples (Levesque et al., 2012; Gardner et al., 2020). A trend within the latter research line that is specific to V&L models is *vision-and-language foiling* (Shekhar et al., 2017b; Gokhale et al., 2020; Bitton et al., 2021; Parcalabescu et al., 2021; Rosenberg et al., 2021), where the idea is to create counterfactual (i.e., *foiled*) and/or minimally edited examples by performing data augmentation on captions (Shekhar et al., 2017b,a) or images (Rosenberg et al., 2021).

Since most V&L models are pretrained on some version of the image-text alignment task, it is possible to test their ability to distinguish correct from foiled captions (in relation to an image) in a zero-shot setting. The construction of foils can serve many investigation purposes. With VALSE, we target the linguistic grounding capabilities of V&L models, focusing on complex phenomena that encompass multiple tokens (i.e., coreference chains, verb-argument structure, or full noun phrases with diverse reference properties such as plurality, existence or counting). At the same time, we ensure that our data is robust to known perturbations and artifacts by i) controlling for word frequency biases between captions and foils, and ii) testing against *unimodal collapse*, thereby preventing models to solve the task by concentrating on a single input modality. This is especially important as it has been shown that V&L models are prone to such problems (Goyal et al., 2017; Madhyastha et al., 2018). The issue of neural models exploiting data artefacts is well-known (Gururangan et al., 2018; Jia et al., 2019; Wang et al., 2020b; He et al., 2021) and methods have been proposed to uncover such effects, including gradient-based, adversarial perturbations or input reduction techniques (cf. Wallace et al., 2020). Yet, these methods are still not fully understood (He et al., 2021) and can be unreliable (Wang et al., 2020b).

Our work is related to Gardner et al. (2020), who construct *task-specific contrast sets* for NLU. However, our focus is on modelling *linguistic phenomena* instead of tasks, and we construct carefully

¹We release our dataset and code upon acceptance.







	pieces	existence	plurality	counting	relations	actions	coreference
	instruments	<i>existential quantifiers</i>	<i>semantic number</i>	<i>balanced, adversarial, small numbers</i>	<i>prepositions</i>	<i>replacement, actant swap</i>	<i>standard, clean</i>
Data collection & metadata	#examples [†]	505	851	2, 459	535	1, 633	812
	foil generation method	<i>nothing ↔ something</i>	NP replacement (sg2pl; pl2sg) & quantifier insertion	numeral placement	re- SpanBERT prediction	action replacement, actant swap	<i>yes ↔ no</i>
	MLM	✗	✗	✗	✓	✓	✗
	GRUEN	✗	✓	✗	✓	✗	✗
	NLI	✗	✓	✗	✓	✗	✗
	src. dataset	Visual7W	MSCOCO	Visual7W	MSCOCO	SWiG	VisDial v1.0
Example data	image src.	MSCOCO	MSCOCO	MSCOCO	MSCOCO	SituNet	MSCOCO
	caption (blue) / foil (orange)	<i>There are no animals / animals shown.</i>	<i>A small copper vase with some flowers / exactly one flower in it.</i>	<i>There are four / six zebras.</i>	<i>A cat plays with a pocket knife on / underneath a table.</i>	<i>A man / woman shouts at a woman / man.</i>	<i>Buffalos walk along grass. Are they in a zoo? No / Yes.</i>
	image						

Table 1: Overview of pieces and instruments in VALSE, with number of examples per piece; the foil generation method used; whether masked language modelling (MLM), GRUEN, and NLI filtering are used; dataset and image sources; and image-caption-foil examples. [†]The number of examples is the sum of the examples available for each instrument in the piece. In Table 4 (in the Appendix) we list the number of examples in each individual instrument.

curated, balanced, single foils from valid instances that we select from multiple multimodal datasets.

3 Constructing the VALSE benchmark

We resort to a musical analogy to describe VALSE: Vision And Language Structured Evaluation is composed of 6 *pieces*, each corresponding to a specific linguistic phenomenon (see Table 1 for an overview). Each piece consists of one or more *instruments* designed to evaluate a model’s ability to ground that specific linguistic phenomenon.

All instruments are built by applying *foiling functions* (FFs) specific to the linguistic phenomenon under study. FFs take a *correct caption* as input and change a specific part of it to produce a *foiled caption* (or *foil*). We design FFs such that the sentences they produce fail to describe the image, while still being grammatical and otherwise valid sentences.

Of course, a *foiled caption* may be less likely than the original caption from which it was produced, and such unwarranted biases can be easily picked up by overparameterised V&L models. Moreover, an automatic FF may fail to produce a foil that contradicts the image, for example by altering the original caption to yield a near-synonymous one, or one that is entailed by the original caption. For phenomena that make it difficult to control these crucial properties of foils, we apply additional filters: i) some FFs make use of strong LMs to propose changes to captions, so that the gener-

ated foils are still high-probability sentences; ii) we use state-of-the-art natural language inference (NLI) methods to detect cases where there is an *entailment* between caption and foil, and filter out such foils from the dataset (see §4 for discussion). As a final measure, we employ human annotators to validate all generated testing data in VALSE.

We build VALSE by sourcing data from existing V&L datasets. Below, we describe each piece and its instruments, and the corresponding task setup in VALSE. For each instrument, we follow the same procedure: i) we identify captions that contain instances of the targeted linguistic phenomenon; ii) we apply a FF that automatically replaces the expression with a variant that contradicts the original expression’s visual content, thereby constructing one or more foils from each target instance in the original caption; as discussed in §4; we then iii) subject the obtained foils to various filters, with the aim of distilling a subset of *valid* and *reliable* foils that cannot be easily tricked by a new generation of highly parameterised pretrained V&L models.

3.1 Existence

The **existence** piece has a single instrument and targets instances with **existential quantifiers**. Models need to differentiate between examples i) where *there is no entity* of a certain type or ii) where *one or more of these entities* are visible in an image.

We use the Visual7W visual question answering

dataset (Zhu et al., 2016) and source its ‘how many’ examples, building a pool of those whose answers are numerals (0, 1, 2, etc.). We use templates to transform question and answer fields into a declarative statement that correctly describes what can be seen in the image, e.g. ‘Q: How many animals are shown? A: 0’ \rightarrow ‘There are 0 animals shown’. We then transform these statements into an existential statement. In the example above, we replace the numeral by the word ‘no’ to create a correct caption (‘There are no animals shown’) and remove the numeral altogether to create a foil (‘There are animals shown’). The existence piece has 505 image–caption–foil tuples after manual validation out of 534 candidates (cf. §4), and captions/foils are balanced: 50% of the (correct) captions originally have answer 0, and the remaining have answer 1 or greater. Full details are provided in A.1.

3.2 Plurality

The **plurality** piece has a single instrument, concerned with **semantic number**. It is intended to test whether a model is able to distinguish between noun phrases denoting a single entity in an image (‘exactly one flower’), versus multiple entities (‘some flowers’). The dataset consists of 851 instances from 1000 generated candidates (cf. §4), evenly divided between cases where the caption contains a plural NP, foiled by replacing it with a singular (p12sg: ‘some flowers’ \rightarrow ‘exactly one flower’), or conversely, the caption contains a singular which is foiled by replacing it with a plural (sg2p1). Foil candidates were generated from the COCO 2017 validation set (Chen et al., 2015). Full details are provided in A.2.

3.3 Counting

The **counting** piece has three instruments: **balanced**, **adversarial** and **small numbers**. All instances are *statements about the number of entities visible in an image*. The model needs to differentiate between examples where *the specific number of entities in the associated image* is correct or incorrect, given the statement. Similarly to the existence piece, we use the Visual7W VQA dataset (Zhu et al., 2016) and source its ‘how many’ examples whose answers are numerals (0, 1, 2, etc.). We use templates to transform question and answer fields into a declarative statement describing the image and create foils by replacing the numeral in the correct statement by another numeral.

All three instruments are designed to show

whether models learn strategies that generalize beyond the training distribution, and to what extent a model exploits class frequency bias.² In **counting balanced** we cap the number of examples to a maximum per class and make sure correct/foil classes are balanced, so that models that exploit class frequency bias are penalized. In **counting adversarial** we make sure that all foils take class $n \in \{0, 1, 2, 3\}$, whereas all correct captions take class $n \in \{n \mid n \geq 4\}$. Biased models are expected to favour more frequent classes. Since small numbers are naturally the most frequent, models that resort to such biases should perform poorly on this adversarially built test. **Counting small numbers** is a sanity check where all correct captions and foils have class $n \in \{0, 1, 2, 3\}$, and caption/foil classes are balanced. Since models likely have been exposed to many examples in this class set and all such classes are high-frequency, with this instrument we disentangle model performance from class exposure. Counting balanced, adversarial, and small numbers have 868 (1000), 691 (756), and 900 (1000) instances after (before) manual validation, respectively (cf. §4). For details, see A.3.

3.4 Spatial relations

The **relations** piece has a single instrument and focuses on the ability of models to distinguish between different spatial relations. Foils differ from the original caption only by the replacement of a spatial preposition. As for plurals, the data was sourced from the COCO 2017 validation split. To create foils, we first identified all preposition sequences in captions (e.g., ‘in’, ‘out of’). Foils were created by masking the prepositions and using SpanBERT (Joshi et al., 2020) to generate replacements of between 1–3 words in length. We keep SpanBERT candidates which differ from the original preposition sequence, but exist in the dataset. There are 535 instances after manual validation out of 614 proposed instances (cf. §4), and we ensure that prepositions are similarly distributed among captions and foils. Full details are provided in A.4.

3.5 Actions

The **actions** piece has two instruments: i) **action replacement** and ii) **actant swap**. They test a V&L model’s capability to i) identify whether an *action* mentioned in the text matches the action

²We take the original answer in Visual7W as the example class: e.g., in ‘There are 0 animals shown’, the class is 0.

seen in the image (e.g., ‘a man [shouts](#) / [smiles](#) at a woman’), and ii) correctly identify the *participants* of an action and the *roles* they play (e.g., is it the man who is shouting or is it the woman, given the picture in Table 1?).

The SWiG dataset (Pratt et al., 2020) contains 504 action verbs, and we generate captions and foils from SWiG annotations of semantic roles and their fillers. For the action replacement piece, we exchange action verbs with other verbs from SWiG that fit the context as suggested by BERT. For the actant swap, we swap role fillers in the role annotations, hence generating action descriptions with inverted roles. Action replacement and actant swap have 648 (779) and 949 (1042) instances after (before) manual validation, respectively (cf. §4). See A.5 for full details.

3.6 Coreference

The **coreference** piece aims to uncover whether V&L models are able to perform pronominal coreference resolution. It encompasses cases where i) the pronoun has a noun (phrase) antecedent and pronoun and (noun) phrase are both grounded in the visual modality (‘A woman is driving a motorcycle. Is she wearing a helmet?’), and cases where ii) the pronoun refers to a region in the image or even to the entire image (‘Is this outside?’).

We create foils based on VisDial v1.0 (Das et al., 2017) with images from MSCOCO (Lin et al., 2014). VisDial captions and dialogues are Q&A sequences. We select image descriptions of the form [Caption. Question? Yes/No.] where the question contains at least one pronoun. When foiling, we exchange the answer from *yes* to *no* and vice-versa (see Table 1). We ensure a 50-50% balance between *yes* / *no* answers.

The coreference piece consists of two instruments: **coreference standard** originating from the VisDial train set and a small **coreference clean** set from the validation set, containing 708 (916) and 104 (141) examples after (before) manual validation, respectively (cf. §4).³ See A.6 for full details.

4 Constructing *valid* and *reliable* foils

In the context of VALSE, instances consisting of image-caption-foil triples are *valid* if: foils minimally differ from the original caption; foils do not accurately describe the image; and independent judges agree that the captions, but not the foils, are

accurate descriptions of the image. As for reliability, a foiling method is more *reliable* the more it ensures that generated foils do not substantially differ from human captions regarding distributional and plausibility bias, and cannot be easily solved unimodally.

In this section, we discuss automatic and manual means to ascertain validity and reliability. Two types of bias are especially worthy of note when constructing a foiling benchmark: distributional bias (see §4.1) and plausibility bias (see §4.2). In §4.3 we discuss how we apply a natural language inference model to filter examples in our data pipeline, and in §4.4 we discuss how we manually validate *all the examples* used in our benchmark.

4.1 Mitigating distributional bias

A first form of bias is related to distributional imbalance between captions and foils (e.g., certain words or phrases having a high probability only in foils). Previous foiling datasets exhibit such imbalance, enabling models to solve the task disregarding the image (Madhyastha et al., 2019). To mitigate this problem, for each phenomenon and throughout our data creation process, we ensure that the token *frequency distributions* in correct and foiled captions are approximately the same (cf. App. A and E).

4.2 Countering plausibility bias

A second form of bias may arise from automatic foil construction procedures yielding foils that are implausible or unnatural, and thereby can act as signals that facilitate their detection. Often, VALSE pieces can be safely foiled by simple rules (e.g., switching from existence to non-existence, or from singular to plural or vice versa). However, with *spatial relations* and *actions*, a foil could be deemed unlikely given only the textual modality and independently of the image, e.g., ‘a man stands [under](#) / [on](#) a chair’. Such **plausibility biases** may be detected by large language models that incorporate commonsense knowledge (Petroni et al., 2019; Wang et al., 2020a), and we expect future V&L models to exhibit similar capabilities.

To ensure that foiled captions are deemed as plausible as correct captions by LMs, we use language models such as BERT (Devlin et al., 2019) and SpanBERT (Joshi et al., 2020) to suggest plausible replacements in our foiling functions. Additionally, in the case of spatial relations and plurals, we also apply a grammaticality filter using GRUEN (Zhu and Bhat, 2020). GRUEN was orig-

³VisDial annotations are not available for the test set.

inally proposed as a method to assign automatically generated sentences a composite score which reflects discourse-level and grammatical properties. We use only the grammaticality component of GRUEN, and retain only foil candidates with a grammaticality score ≥ 0.8 .

Furthermore, we evaluate unimodal, language-only models on VALSE to verify whether our benchmark could be solved by a multimodal model with strong linguistic capacities in **unimodal collapse**, whereby a model silently relies on a single modality within which biases are easier to exploit (Goyal et al., 2017; Shekhar et al., 2019a). By evaluating VALSE with unimodal models, we establish a baseline that V&L models should exceed if we are to expect true multimodal integration.

4.3 Filtering foils with NL Inference

When constructing foils, we need to ensure that they *fail* to describe the image. To test this automatically, we apply natural language inference (NLI) with the following rationale: We consider an image and its caption as a premise and its entailed hypothesis, respectively (a similar rationale is applied in the visual entailment task; Xie et al., 2019). In addition, we consider the *caption as premise* and the *foil as its hypothesis*. If a NLI model predicts the foil to be entailed (E) by the caption, it cannot be a good foil since by transitivity it will give a truthful description of the image. By contrast, if the foil is predicted to contradict (C) or to be neutral (N) with respect to the caption, we take this as an indicator of a good (C) or a plausible (N) foil.⁴

We use the NLI model ALBERT (Lan et al., 2020) finetuned on the task (see Appendix C for details). Filtering with NLI was initially applied to *relations*, *plurals* and *actions*, on the grounds that foils in these pieces may induce substantive changes to lexical content.⁵ Following automatic labelling of caption-foil pairs, we manually validated a sample labelled as E, C or N. For *relations* ($N = 30$), labels were found to be near 100% accurate with only 2 (0.06%) errors overall. For *plurals*

⁴See the following examples from action replacement:
P: A mother scolds her son.

H1: A mother encourages her son. (C; good foil);
H2: A mother camps with her son. (N; needs image control);
H3: A mother talks to her son. (E; not a suitable foil)

If the NLI prediction is N, we still need to check the image, since the description might happen to fit the image content.

⁵By contrast, existence and counting foils involve a more straightforward swap (e.g., between numerical quantities); similarly, coreference foils simply involve the replacement of a positive with a negative answer.

($N = 60$, 50% *sg2pl* and 50% *pl2sg*), the error rate was also low, with 0 errors for C, 33% errors for E and 11% errors for N. Here, a number of entailment errors were due to odd formulations arising from the automatic foiling process, whereas no such oddities were observed for C. We therefore include only foils labelled C in the final relations and plurals pieces. For *actions*, the model labelled contradictions very accurately (0% error) but was erroneous up to 97.1% for E, meaning that a large number of valid foils would be spuriously excluded. To avoid reducing the dataset too much, we did not use NLI filtering for actions, but relied on human annotation as a final validity check.

4.4 Manual evaluation of generated foils

Each instance in VALSE comprises an image, its caption and a foiled caption (cf. Table 1). As shown above, we take various automatic measures to ensure the quality of foils in each piece. As a final step, the entire data for each instrument was submitted to a manual validation, which took the following form: for each instance, annotators were shown the image, the caption and the foil. Caption and foil were numbered and displayed above each other to make differences more apparent, with differing elements highlighted in boldface (Fig. 4, E). Annotators were not informed which text was the caption and which was the foil, and captions appeared first (numbered 1) 50% of the time. The task was to determine which of the two texts accurately described what could be seen in the image. In each case, annotators had a forced choice between five options: a) the first, but not the second; b) the second, but not the first; c) both of them; d) neither of the two; and e) I cannot tell.

Each item was annotated by three individuals. The validation was conducted on Amazon Mechanical Turk with a fixed set of annotators who had qualified for the task. For details see Appendix E. We consider an instance to have passed the validation test if at least two out of three annotators identified the caption, but not the foil, as the text which accurately describes the image. Across all instruments, 87.7% of the instances satisfied this criterion (min 77.3%; max 94.6%; full details in Appendix E), with 73.6% of instances overall having a unanimous (3/3) decision that the caption, but not the foil, was an accurate description. We consider these figures high, suggesting that the automatic construction and filtering procedures yield

foils which are likely to be valid, in the sense discussed in §4 above.

5 Benchmarking with VALSE

We propose VALSE as a task-independent, *zero-shot* benchmark to assess the extent to which models learn to ground specific linguistic phenomena as a consequence of their pretraining (or fine-tuning). VALSE is built in the spirit of approaches such as Checklist (Ribeiro et al., 2020), including pairs consisting of captions and minimally edited foils.

The only requirement to evaluate a model on our benchmark is: *i*) to have a binary classification head to predict whether an image-sentence pair is foiled, or *ii*) to predict an image-sentence matching score between the image and the caption vs. the foil, returning the pair with the highest score. Systems reporting results on VALSE are expected to report any data used in model training prior to testing on VALSE, for comparability.

5.1 Benchmark Metrics

We employ four metrics⁶ for evaluation: overall **accuracy** (acc) on all classes (foil and correct); **precision** (p_c) measuring how well models identify the *correct* examples; **foil precision** (p_f) measuring how well *foiled* cases are identified; and **pairwise ranking accuracy** (acc_r), which measures whether the image-sentence alignment score is greater for a correct image-text pair than for its foiled pair. acc_r is more permissive than acc because it accepts model predictions if the score for a foil is lower than the caption’s score. Our main metric is acc_r , which gives results for a pair $\langle image, caption \rangle$ and $\langle image, foil \rangle$ and is better suited to evaluate minimally-edited pairs as it does not need a classification threshold. Since all instruments are implemented as a binary classification, the random baseline is always 50%.

5.2 V&L models

We benchmark five V&L models on VALSE: CLIP (Radford et al., 2021), LXMERT (Tan and Bansal, 2019), ViLBERT (Lu et al., 2019), ViLBERT 12-in-1 (Lu et al., 2020), and VisualBERT (Li et al., 2019). These models have different architectures, are pretrained on a variety of tasks and using different training data. We also benchmark two unimodal text-only models, GPT1 (Radford et al., 2018) and GPT2 (Radford et al., 2019), discussed below. See

Appendix D for details about all unimodal and V&L models we use in our evaluation.

Unimodal models GPT1 and GPT2 are autoregressive language models pretrained on English text. We test whether VALSE is solvable by these unimodal models by computing the perplexity of the correct and foiled caption and *predicting the entry with the lowest perplexity*. If the perplexity is higher for the foil, we take this as an indication that the foiled caption may suffer from **plausibility bias** or other linguistic biases (cf. §4.2).

5.3 Experiments and Results

We test V&L and unimodal models on VALSE in a zero-shot setting, and also evaluate on a number of correct captions and foils from the *FOIL it!* dataset (Shekhar et al., 2017b) (cf. App. A.7 for details). All results are listed in Table 2.

Unimodal results For most instruments, unimodal results are close to random and hence do not signal strong linguistic or plausibility biases. One exception is the original *FOIL it!* dataset, in line with Madhyastha et al. (2019)’s findings. Spatial relations (77.2%), action replacement (66.8%) and actant swap (76.9%) instruments also suggest plausibility biases in the foils. Such biases are hard to avoid in automatic foil generation for actions due to the verb arguments’ selectional restrictions, which are easily violated when flipping role fillers, or exchanging the verb. Similar considerations hold for relations: though SpanBERT proposals are intended to aid selection of likely replacements for prepositions, plausibility issues arise with relatively rare argument-preposition combinations.

While these might be the first instruments in VALSE to be solved in the future, current V&L models struggle to detect even blatant mismatches of actant swap, e.g., ‘A ball throws a tennis player.’ For VALSE, the unimodal scores will serve as a baseline for the pairwise accuracy of V&L models.

Multimodal results The best zero-shot results are achieved by ViLBERT 12-in-1 with the highest scores across the board, followed by ViLBERT, LXMERT, CLIP,⁷ and finally VisualBERT. The high p_f values for the latter indicate that the model is accurate at predicting foils, but far less so at predicting correct captions. We hypothesise that this is due to the way image-sentence alignment is framed

⁶All metrics are defined in Appendix B.

⁷CLIP works in a contrastive fashion, therefore we report only acc_r (cf. Appendix D for details).

Metric	Model	Existence quantifiers	Plurality number	Counting			Sp.rel. [‡] relations	Action repl. [†]	Action actant swap	Coreference		Foil-it!	Avg.
				balanced	sns. [†]	adv. [†]				standard	clean		
	Random	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0
acc_r	GPT1*	61.8	53.1	51.2	<u>48.7</u>	69.5	77.2	65.4	72.2	<u>45.6</u>	<u>45.2</u>	77.5	60.7
	GPT2*	58.0	51.9	51.6	<u>49.8</u>	<u>45.3</u>	75.0	66.8	76.9	54.5	<u>50.0</u>	80.7	60.1
	CLIP	66.9	56.2	62.1	62.5	57.5	64.3	75.6	68.6	52.1	<u>49.7</u>	88.8	64.0
	LXMERT	78.6	64.4	62.2	69.2	<u>42.6</u>	60.2	54.8	<u>45.8</u>	<u>46.8</u>	<u>44.2</u>	87.1	59.6
	ViLBERT	65.5	61.2	58.6	62.9	73.7	57.2	70.7	68.3	<u>47.2</u>	<u>48.1</u>	86.9	63.7
	12-in-1	95.6	72.4	76.7	80.2	77.3	67.7	65.9	58.9	75.7	69.2	86.9	75.1
	VisualBERT	<u>39.7</u>	<u>45.7</u>	<u>48.2</u>	<u>48.2</u>	<u>50.0</u>	<u>39.7</u>	<u>49.2</u>	<u>44.4</u>	<u>49.5</u>	<u>47.6</u>	<u>48.5</u>	<u>46.4</u>
acc	LXMERT	55.8	55.1	52.0	55.4	49.9	50.8	51.1	48.5	49.8	49.0	70.8	53.5
	ViLBERT	<u>2.4</u>	50.3	50.7	50.6	51.8	<u>49.9</u>	52.6	50.4	<u>50.0</u>	<u>50.0</u>	55.9	51.3
	12-in-1	89.0	62.0	64.9	69.2	66.7	53.4	57.3	52.2	54.4	54.3	71.5	63.2
	VisualBERT	<u>49.3</u>	<u>46.5</u>	<u>48.3</u>	<u>47.8</u>	<u>50.0</u>	<u>49.3</u>	<u>48.8</u>	<u>49.7</u>	<u>50.0</u>	<u>50.0</u>	<u>46.6</u>	<u>48.8</u>
p_c	LXMERT	41.6	68.0	50.9	<u>50.0</u>	61.5	73.1	35.8	36.8	81.2	80.8	72.3	59.3
	ViLBERT	56.8	98.5	77.0	76.6	86.1	98.3	93.2	93.7	98.7	98.1	98.8	88.7
	12-in-1	85.0	90.7	64.3	76.7	59.5	93.5	66.7	66.8	92.9	95.2	94.3	80.5
	VisualBERT	<u>1.3</u>	<u>0.3</u>	<u>0.0</u>	<u>0.0</u>	<u>0.0</u>	<u>1.3</u>	<u>0.0</u>	<u>0.0</u>	<u>0.0</u>	<u>0.0</u>	<u>0.2</u>	<u>0.3</u>
p_f	LXMERT	70.1	<u>42.2</u>	53.0	60.8	<u>37.3</u>	28.4	66.4	60.2	<u>18.4</u>	<u>17.3</u>	69.3	<u>47.6</u>
	ViLBERT	47.9	2.1	24.4	24.7	<u>17.5</u>	1.5	11.9	7.1	1.3	1.9	12.9	13.9
	12-in-1	93.1	33.4	65.6	61.7	74.0	<u>13.3</u>	<u>47.8</u>	<u>37.6</u>	15.8	<u>13.5</u>	<u>48.8</u>	45.9
	VisualBERT	97.3	92.8	96.7	95.7	100.0	97.3	97.6	99.4	100.0	100.0	93.0	97.3

Table 2: Performance of unimodal and multimodal models on the VALSE benchmark according to different metrics. We bold-face the best overall result per metric, and underscore all results below (or at) the random baseline. acc_r is a pairwise ranking accuracy where a prediction is considered correct if $p(caption, img) > p(foil, img)$. Precision p_c and foil precision p_f are *competing* metrics where naively increasing one can decrease the other: therefore looking at the smaller number among the two gives a good intuition of how informed is a model prediction. [†]sns. Counting small numbers. **adv.** Counting adversarial. **repl.** Action replacement. [‡] **Sp.rel.** Spatial relations. *Unimodal text-only models that do not use images as input. CLIP is only tested in pairwise ranking mode (fn. 6).

in VisualBERT’s pretraining: the model expects an image and a (correct) sentence c_1 , and predicts whether a second sentence c_2 is correct or a foil. During pretraining c_1 and c_2 are likely to differ in many ways, whereas in our setting, they are nearly identical, modulo a word or phrase replaced by the foiling procedure. This may bias the model against predicting foils, which would raise the value p_f .

Instruments centered on individual objects like existence and the *FOIL it!* dataset are almost solved by ViLBERT 12-in-1, highlighting that models are capable of identifying named objects and their presence in images. However, none of the remaining pieces can be reliably solved in our adversarial foiling settings: i) distinguishing references to single vs. multiple objects of a given type or counting them in an image; ii) correctly classifying a named spatial relation between objects in an image; iv) distinguishing actions and reliably identifying their participants, even if supported by preference biases; or, v) tracing references to the same object in an image through the use of pronouns.

Correct and foil precision p_c and p_f show that V&L models struggle to solve the phenomena in VALSE. When a model achieves high precision on correct captions p_c this is often at the expense of very low precision on foiled captions p_f (e.g., ViLBERT), or vice-versa (e.g., VisualBERT). This

suggests that such models are insensitive to the inputs in VALSE: models that almost always predict a match will inflate p_f at the expense of p_c . Considering $\min(p_c, p_f)$ reveals that VisualBERT and ViLBERT perform poorly and below the random baseline, and LXMERT close to or below it. ViLBERT 12-in-1 performs strongly on existence, well on counting, but struggles on plurality, spatial relations, coreference, actions. These tendencies we see reflected in our main pairwise metric, acc_r .

6 Conclusions and Future Work

We present the VALSE benchmark to help the community improve V&L models by hard-testing their visual grounding capabilities through the lens of linguistic constructs. Our experiments show that V&L models identify named objects and their presence in images well, but struggle to ground objects, their interdependence and relationships in visual scenes when forced to respect refined linguistic indicators. We encourage the community to use VALSE for measuring the progress towards V&L models capable of true language grounding.

VALSE is designed as a living benchmark. As future work we plan to extend it to further linguistic phenomena, and to source data from diverse V&L datasets to cover more linguistic variability and image distributions.

References

- Vedika Agarwal, Rakshith Shetty, and Mario Fritz. 2020. Towards causal vqa: Revealing and reducing spurious correlations by invariant and covariant semantic editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9690–9698.
- Raffaella Bernardi and Sandro Pezzelle. 2021. [Linguistic issues behind visual question answering](#). *Language and Linguistics Compass*, 15(6):1–25.
- Yonatan Bitton, Gabriel Stanovsky, Roy Schwartz, and Michael Elhadad. 2021. [Automatic generation of contrast sets from scene graphs: Probing the compositional consistency of GQA](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 94–105, Online. Association for Computational Linguistics.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Jize Cao, Zhe Gan, Yu Cheng, Licheng Yu, Yen-Chun Chen, and Jingjing Liu. 2020. Behind the scene: Revealing the secrets of pre-trained vision-and-language models. *arXiv preprint arXiv:2005.07310*.
- Xinlei Chen, Hao Fang, Tsung-yi Lin, Ramakrishna Vedantam, C Lawrence Zitnick, Saurabh Gupta, and Piotr Doll. 2015. [Microsoft COCO Captions: Data Collection and Evaluation Server](#). *arXiv*, 1504.00325:1–7.
- 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 *ECCV*.
- Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José M.F. Moura, Devi Parikh, and Dhruv Batra. 2017. Visual Dialog. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- 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.
- Christiane Fellbaum. 1998. *WordNet: An Electronic Lexical Database*. Bradford Books.
- Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hannaneh Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfaty, Eric Wallace, Ally Zhang, and Ben Zhou. 2020. [Evaluating models’ local decision boundaries via contrast sets](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1307–1323, Online. Association for Computational Linguistics.
- Sahaj Garg, Vincent Perot, Nicole Limtiaco, Ankur Taly, Ed H. Chi, and Alex Beutel. 2019. [Counterfactual fairness in text classification through robustness](#). In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society, AIES ’19*, page 219–226, New York, NY, USA. Association for Computing Machinery.
- Albert Gatt and Ehud Reiter. 2009. [SimpleNLG: A realisation engine for practical applications](#). In *Proceedings of the 12th European Workshop on Natural Language Generation (ENLG 2009)*, pages 90–93, Athens, Greece. Association for Computational Linguistics.
- Gabriel Goh, Nick Cammarata †, Chelsea Voss †, Shan Carter, Michael Petrov, Ludwig Schubert, Alec Radford, and Chris Olah. 2021. [Multimodal neurons in artificial neural networks](#). *Distill*. <https://distill.pub/2021/multimodal-neurons>.
- 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*, pages 6904–6913.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. 2018. [Annotation artifacts in natural language inference data](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112, New Orleans, Louisiana. Association for Computational Linguistics.
- Feijuan He, Yaxian Wang, Xianglin Miao, and Xia Sun. 2021. [Interpretable visual reasoning: A survey](#). *Image and Vision Computing*, 112:104194.
- Robin Jia and Percy Liang. 2017. [Adversarial examples for evaluating reading comprehension systems](#).

787	In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 2021–2031, Copenhagen, Denmark. Association for Computational Linguistics.	
788		
789		
790		
791	Robin Jia, Aditi Raghunathan, Kerem Göksel, and Percy Liang. 2019. Certified robustness to adversarial word substitutions . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 4129–4142, Hong Kong, China. Association for Computational Linguistics.	
792		
793		
794		
795		
796		
797		
798		
799	Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2020. SpanBERT: Improving pre-training by representing and predicting spans . <i>Transactions of the Association for Computational Linguistics</i> , 8:64–77.	
800		
801		
802		
803		
804	Kushal Kafle, Robik Shrestha, and Christopher Kanan. 2019. Challenges and prospects in vision and language research . <i>Frontiers in Artificial Intelligence</i> , 2:28.	
805		
806		
807		
808	Wonjae Kim, Bokyung Son, and Ildoo Kim. 2021. Vilt: Vision-and-language transformer without convolution or region supervision. <i>arXiv preprint arXiv:2102.03334</i> .	
809		
810		
811		
812	Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. Albert: A lite bert for self-supervised learning of language representations . In <i>International Conference on Learning Representations</i> .	
813		
814		
815		
816		
817	Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In <i>Thirteenth International Conference on the Principles of Knowledge Representation and Reasoning</i> .	
818		
819		
820		
821	Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. 2020a. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training . In <i>The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020</i> , pages 11336–11344. AAAI Press.	
822		
823		
824		
825		
826		
827		
828		
829		
830		
831	Linjie Li, Jie Lei, Zhe Gan, Licheng Yu, Yen-Chun Chen, Rohit Pillai, Yu Cheng, Luowei Zhou, Xin Eric Wang, William Yang Wang, et al. 2021. Value: A multi-task benchmark for video-and-language understanding evaluation . <i>arXiv preprint arXiv:2106.04632</i> .	
832		
833		
834		
835		
836		
837	Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. In <i>Arxiv</i> .	
838		
839		
840		
	Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. 2020b. Oscar: Object-semantics aligned pre-training for vision-language tasks . In <i>European Conference on Computer Vision</i> , pages 121–137. Springer.	841
		842
		843
		844
		845
		846
	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 <i>Computer Vision – ECCV 2014</i> , pages 740–755, Cham. Springer International Publishing.	847
		848
		849
		850
		851
		852
	Nicholas Lourie, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Unicorn on rainbow: A universal commonsense reasoning model on a new multitask benchmark. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pages 13480–13488.	853
		854
		855
		856
		857
		858
	Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks . In <i>Advances in Neural Information Processing Systems</i> , pages 13–23.	859
		860
		861
		862
		863
	Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 2020. 12-in-1: Multi-task vision and language representation learning. In <i>The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> .	864
		865
		866
		867
		868
	Pranava Madhyastha, Josiah Wang, and Lucia Specia. 2019. VIFIDEL: Evaluating the visual fidelity of image descriptions . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 6539–6550, Florence, Italy. Association for Computational Linguistics.	869
		870
		871
		872
		873
		874
	Pranava Swaroop Madhyastha, Josiah Wang, and Lucia Specia. 2018. Defoiling foiled image captions . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 433–438, New Orleans, Louisiana. Association for Computational Linguistics.	875
		876
		877
		878
		879
		880
		881
		882
	Yixin Nie, Haonan Chen, and Mohit Bansal. 2019. Combining fact extraction and verification with neural semantic matching networks . In <i>Association for the Advancement of Artificial Intelligence (AAAI)</i> .	883
		884
		885
		886
	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial NLI: A new benchmark for natural language understanding . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4885–4901, Online. Association for Computational Linguistics.	887
		888
		889
		890
		891
		892
		893
	Letitia Parcalabescu, Albert Gatt, Anette Frank, and Iacer Calixto. 2021. Seeing past words: Testing the cross-modal capabilities of pretrained v&l models on counting tasks . In <i>Proceedings of the 1st</i>	894
		895
		896
		897

898	Workshop on Multimodal Semantic Representations	954
899	(MMSR), pages 32–44, Groningen, Netherlands (On-	955
900	line). Association for Computational Linguistics.	956
901	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel,	957
902	Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and	
903	Alexander Miller. 2019. Language models as knowl-	Ravi Shekhar, Sandro Pezzelle, Aurélie Herbelot, Moin
904	edge bases? In <i>Proceedings of the 2019 Confer-</i>	Nabi, Enver Sangineto, and Raffaella Bernardi.
905	<i>ence on Empirical Methods in Natural Language</i>	2017a. Vision and language integration: Moving be-
906	<i>Processing and the 9th International Joint Confer-</i>	yond objects . In <i>IWCS 2017 — 12th International</i>
907	<i>ence on Natural Language Processing (EMNLP-</i>	<i>Conference on Computational Semantics — Short</i>
908	<i>IJCNLP)</i> , pages 2463–2473, Hong Kong, China. As-	<i>papers</i> .
909	sociation for Computational Linguistics.	963
910	Bryan A Plummer, Liwei Wang, Chris M Cervantes,	964
911	Juan C Caicedo, Julia Hockenmaier, and Svetlana	Ravi Shekhar, Sandro Pezzelle, Yauhen Klimovich, Au-
912	Lazebnik. 2015. Flickr30k entities: Collecting	rémie Herbelot, Moin Nabi, Enver Sangineto, and
913	region-to-phrase correspondences for richer image-	Raffaella Bernardi. 2017b. FOIL it! find one mis-
914	to-sentence models. In <i>Proceedings of the IEEE</i>	match between image and language caption . In <i>Pro-</i>
915	<i>international conference on computer vision</i> , pages	<i>ceedings of the 55th Annual Meeting of the Associa-</i>
916	2641–2649.	<i>tion for Computational Linguistics (Volume 1: Long</i>
917	Sarah Pratt, Mark Yatskar, Luca Weihs, Ali Farhadi,	<i>Papers)</i> , pages 255–265, Vancouver, Canada. Asso-
918	and Aniruddha Kembhavi. 2020. Grounded situa-	ciation for Computational Linguistics.
919	tion recognition . In <i>Computer Vision - ECCV 2020 -</i>	971
920	<i>16th European Conference</i> , pages 314–332.	972
921	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	Ravi Shekhar, Ece Takmaz, Raquel Fernández, and
922	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish	Raffaella Bernardi. 2019a. Evaluating the represen-
923	Sastry, Amanda Askell, Pamela Mishkin, Jack Clark,	tational hub of language and vision models . In <i>Pro-</i>
924	et al. 2021. Learning transferable visual models	<i>ceedings of the 13th International Conference on</i>
925	from natural language supervision. <i>arXiv preprint</i>	<i>Computational Semantics - Long Papers</i> , pages 211–
926	<i>arXiv:2103.00020</i> .	222, Gothenburg, Sweden. Association for Compu-
927	Alec Radford, Karthik Narasimhan, Tim Salimans, and	tational Linguistics.
928	Ilya Sutskever. 2018. Improving language under-	973
929	standing by generative pre-training.	974
930	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	975
931	Dario Amodei, and Ilya Sutskever. 2019. Language	976
932	models are unsupervised multitask learners. <i>OpenAI</i>	977
933	<i>blog</i> , 1(8):9.	978
934	Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin,	979
935	and Sameer Singh. 2020. Beyond accuracy: Be-	Ravi Shekhar, Aashish Venkatesh, Tim Baumgärtner,
936	havioral testing of NLP models with CheckList . In	Elia Bruni, Barbara Plank, Raffaella Bernardi, and
937	<i>Proceedings of the 58th Annual Meeting of the Asso-</i>	Raquel Fernández. 2019b. Beyond task success: A
938	<i>ciation for Computational Linguistics</i> , pages 4902–	closer look at jointly learning to see, ask, and Guess-
939	4912, Online. Association for Computational Lin-	What . In <i>Proceedings of the 2019 Conference of</i>
940	guistics.	<i>the North American Chapter of the Association for</i>
941	Daniel Rosenberg, Itai Gat, Amir Feder, and Roi Re-	<i>Computational Linguistics: Human Language Tech-</i>
942	ichart. 2021. Are VQA systems RAD? Measuring	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages
943	robustness to augmented data with focused interven-	2578–2587, Minneapolis, Minnesota. Association
944	tions . In <i>Proceedings of the 59th Annual Meeting of</i>	for Computational Linguistics.
945	<i>the Association for Computational Linguistics and</i>	988
946	<i>the 11th International Joint Conference on Natu-</i>	Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu,
947	<i>ral Language Processing (Volume 2: Short Papers)</i> ,	Furu Wei, and Jifeng Dai. 2020. VI-bert: Pre-
948	pages 61–70, Online. Association for Computational	training of generic visual-linguistic representations .
949	Linguistics.	In <i>International Conference on Learning Represen-</i>
950	Piyush Sharma, Nan Ding, Sebastian Goodman, and	<i>tations</i> .
951	Radu Soricut. 2018. Conceptual captions: A	993
952	cleaned, hypernymed, image alt-text dataset for au-	994
953	tomatic image captioning . In <i>Proceedings of the</i>	Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang,
		Huajun Bai, and Yoav Artzi. 2019. A corpus for
		reasoning about natural language grounded in pho-
		tographs . In <i>Proceedings of the 57th Annual Meet-</i>
		<i>ing of the Association for Computational Linguistics</i> ,
		pages 6418–6428, Florence, Italy. Association for
		Computational Linguistics.
		1000
		Hao Tan and Mohit Bansal. 2019. LXMERT: Learning
		cross-modality encoder representations from trans-
		formers . In <i>Proceedings of the 2019 Conference on</i>
		<i>Empirical Methods in Natural Language Processing</i>
		<i>and the 9th International Joint Conference on Natu-</i>
		<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages
		5100–5111, Hong Kong, China. Association for
		Computational Linguistics.
		1008
		Eric Wallace, Matt Gardner, and Sameer Singh. 2020.
		Interpreting predictions of NLP models . In <i>Proceed-</i>
		<i>ings of the 2020 Conference on Empirical Methods</i>
		1009
		1010
		1011

1012 *in Natural Language Processing: Tutorial Abstracts*,
1013 pages 20–23, Online. Association for Computational
1014 Linguistics.

1015 Chenguang Wang, Xiao Liu, and Dawn Song. 2020a.
1016 Language models are open knowledge graphs.
1017 *arXiv preprint arXiv:2010.11967*.

1018 Junlin Wang, Jens Tuyls, Eric Wallace, and Sameer
1019 Singh. 2020b. [Gradient-based analysis of NLP mod-](#)
1020 [els is manipulable](#). In *Findings of the Association*
1021 *for Computational Linguistics: EMNLP 2020*,
1022 pages 247–258, Online. Association for Computa-
1023 tional Linguistics.

1024 Adina Williams, Nikita Nangia, and Samuel Bowman.
1025 2018. [A broad-coverage challenge corpus for sen-](#)
1026 [tence understanding through inference](#). In *Proceed-*
1027 *ings of the 2018 Conference of the North American*
1028 *Chapter of the Association for Computational Lin-*
1029 *guistics: Human Language Technologies, Volume*
1030 *1 (Long Papers)*, pages 1112–1122, New Orleans,
1031 Louisiana. Association for Computational Linguis-
1032 tics.

1033 Thomas Wolf, Julien Chaumond, Lysandre Debut, Vic-
1034 tor Sanh, Clement Delangue, Anthony Moi, Pier-
1035 ric Cistac, Morgan Funtowicz, Joe Davison, Sam
1036 Shleifer, et al. 2020. Transformers: State-of-the-
1037 art natural language processing. In *Proceedings of*
1038 *the 2020 Conference on Empirical Methods in Nat-*
1039 *ural Language Processing: System Demonstrations*,
1040 pages 38–45.

1041 Ning Xie, Farley Lai, Derek Doran, and Asim Kadav.
1042 2019. [Visual Entailment: A Novel Task for Fine-](#)
1043 [Grained Image Understanding](#). *arXiv*, 1901.06706.

1044 Mark Yatskar, Luke Zettlemoyer, and Ali Farhadi.
1045 2016. Situation recognition: Visual semantic role
1046 labeling for image understanding. In *Proceedings of*
1047 *the IEEE Conference on Computer Vision and Pat-*
1048 *tern Recognition (CVPR)*.

1049 Wanzheng Zhu and Suma Bhat. 2020. [GRUEN for](#)
1050 [evaluating linguistic quality of generated text](#). In
1051 *Findings of the Association for Computational Lin-*
1052 *guistics: EMNLP 2020*, pages 94–108, Online. As-
1053 sociation for Computational Linguistics.

1054 Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-
1055 Fei. 2016. Visual7W: Grounded Question Answer-
1056 ing in Images. In *IEEE Conference on Computer*
1057 *Vision and Pattern Recognition*.

A Benchmark creation

A.1 Existence

The **existence** piece has a single instrument and targets instances with **existential quantifiers**. Models need to differentiate between examples i) where *there is no entity* of a certain type or ii) where *there is one or more of these entities* visible in an image.

Data sources We use the Visual7W visual question answering dataset (Zhu et al., 2016) to source examples, starting with the ‘how many’ questions in Visual7W and building a pool of those whose answers are numerals (e.g., 0, 1, 2, etc.). We use the templates from Parcalabescu et al. (2021) to transform question and answer fields into a declarative statement that correctly describes what can be seen in the image, e.g., ‘Q: How many animals are shown? A: 0’ \rightarrow ‘There are 0 animals shown’.

Foiling method Let us use $x =$ ‘There are N animals shown’ as a running example for a correct caption, where N is a number. If $N > 0$, we simply remove N from the sentence, effectively creating the statement $\exists x$ or ‘There are animals shown’. If $N = 0$, we replace N by ‘no’, creating the statement $\neg \exists x$ or ‘There are no animals shown’. If necessary, we fix singular-plural agreement. To create data with balanced correct and foil classes, we select 50% of our examples from those where the correct answer is originally 0, and the remaining 50% from those where the correct answer is any other number (e.g., 1, 2, etc.). To create foils, we then simply convert the statement from $\exists x$ to $\neg \exists x$, and vice-versa.

A.2 Plurality

The **plurality** piece has a single instrument, concerned with **semantic number**, that is, the distinction between single entities in an image (‘exactly one flower’) and multiple instances of the same type (‘some flowers’). In this piece, foil candidates are created either by converting a singular NP and its coreferents to a plural, or vice versa.

Data sources The data was sourced from the validation split of the COCO 2017 dataset (Chen et al., 2015). Captions are only foiled if their length after tokenization with the pretrained BERT tokenizer⁸ is of 80 tokens or less. This is done to minimise the risk that captions and foils need to be truncated

to accommodate the input specifications of current pretrained V&L models.

Foiling method Foiling is done in two directions: singular-to-plural (*sg2pl*) or plural-to-singular (*pl2sg*). Given a caption, NP chunking is applied to identify all non-pronominal NPs. In the *sg2pl* case, a foiled version of a caption containing a singular NP is created by pluralising the head noun. We automatically identify anaphoric expressions coreferring to the singular NP within the caption and pluralise them in the same way. For NPs which are subjects of copular VPs or VPs with an auxiliary requiring subject-verb number agreement (e.g. ‘ N *is* V ’), we also pluralise the verb. Note that this procedure creates a potential foil for every singular NP in the caption; thus, more than one foil candidate can be created for each instance in the source dataset.⁹ In the *pl2sg* case, the same procedure is carried out, but turning a plural NP, as well as its coreferents, into a singular. We generate all foil candidates using the Checklist framework (Ribeiro et al., 2020), within which we implement our procedures for data perturbation.

An important consideration, especially in the *pl2sg* case, is that singularising an NP in a foil can still be truth-preserving. Specifically, a caption with a plural NP, such as ‘A small copper vase with some flowers in it’, arguably still entails the version with the singular ‘(...) a flower’. As a result, the singular version may still correctly be judged to match the image. One way around this problem is to insert a quantifier in the singular NP which makes it explicit that exactly one instance and no more is intended (e.g. ‘exactly one flower’). This may however result in a biased dataset, with such singular quantifiers acting as signals for singular foils and enabling models to solve the task with no grounding in the visual information. We avoid this by adopting a uniform strategy for both *sg2pl* and *pl2sg*. We determine two singular quantifiers (‘exactly one N ’ and ‘a single N ’) and two plural quantifiers (‘some N ’, ‘a number of N ’). When a foil candidate is generated, we alter the *original* NP by inserting one of the two quantifiers matching its semantic number, and generate a foil with one

⁸We use the `bert-large-cased` pretrained tokenizer distributed as part of the `transformers` python library.

⁹NP chunking is performed using the Spacy v.3 pipeline for English using the `en_core_web_md` pretrained models. Coreference chains are detected using the pretrained English model for Coreferee (github.com/msg-systems/coreferee). Pluralisation of head nouns is carried out using the `inflect` engine (github.com/jaraco/inflect/).

of the two quantifiers for the other number. In the foregoing example, we end up with ‘A small copper vase with some flowers / exactly one flower in it.’

After generating all candidate foils, in both directions, we use the GRUEN pretrained model (Zhu and Bhat, 2020) to score the foils for grammaticality. We only keep foils with a score ≥ 0.8 , and run each foil-caption pair through the NLI model described in Section 4.3, keeping only pairs whose predicted label is *contradiction*, for an initial candidate set of 1000 cases (500 *sg2pl* and 500 *pl2sg*), of which 851 (85.1%) are considered valid following manual validation (see §4.4. Figure 3 shows the distribution of nouns in captions and foils, before and after the validation. Note that the validation process does not result in significant change to the distributions.

A.3 Counting

The **counting** piece comes in three instruments: **balanced**, **adversarial** and **small numbers**. All three instruments include instances with *statements about the number of entities visible in an image*. The model needs to differentiate between examples where *the specific number of entities in the associated image* is correct or incorrect, given the statement.

All three instruments are designed to show whether models learn strategies that generalize beyond the training distribution, and to what extent a model exploits class frequency bias.¹⁰ In **counting balanced** we cap the number of examples to a maximum per class and make sure correct/foil classes are balanced, so that models that exploit class frequency bias are penalized. In **counting adversarial** we make sure that all foils take class $n \in \{0, 1, 2, 3\}$, whereas all correct captions take class $n \in \{n \mid n \geq 4\}$. Biased models are expected to favour more frequent classes and these correspond to smaller numbers, therefore models that resort to such biases should perform poorly on this adversarially built test. Instrument **counting small numbers** is a sanity check where all correct captions and foils have class $n \in \{0, 1, 2, 3\}$, and caption/foil classes are balanced. Models likely have been exposed to many examples in this class set, so with this instrument we assess model performance certain it does not suffer from (class) exposure bias.

¹⁰We take the original answer in Visual7W as the example class. E.g., in *There are four zebras*, the class is 4.

Data sources We use the Visual7W visual question answering dataset (Zhu et al., 2016) and source its ‘how many’ examples, building a pool of those whose answers are numerals (e.g., 0, 1, 2, etc.). We use the templates from Parcalabescu et al. (2021) to transform question and answer fields into a declarative statement that correctly describes what can be seen in the image.

Foiling method We create foils by directly replacing the numeral in the correct caption by another numeral. When creating foils we make sure that the class distribution for correct and foiled captions are approximately the same, i.e., there are a similar number of correct and foiled examples in each class in each instrument. The only exception is the counting adversarial instrument, where the classes used in correct and foiled captions are disjoint, i.e., $n \in \{0, 1, 2, 3\}$ and $n \in \{n \mid n \geq 4\}$, respectively. See Figure 2 for a visualisation of these distributions.

A.4 Spatial relations

The **relations** piece has one instrument and focuses on the ability of models to distinguish between different spatial relations, as expressed by prepositions. Foils therefore consist of captions identical to the original except for the replacement of a spatial preposition.

Data sources Data was sourced from the COCO 2017 validation split (Chen et al., 2015). To generate foil candidates, we first extracted from the original COCO captions all the sequences consisting of one or more consecutive prepositions (e.g., ‘on’ or ‘out of’). Foils are generated by detecting these preposition spans, and replacing them with another preposition span attested in the list.

Foiling method To generate foils, we mask the preposition span in an original caption, and use SpanBERT (Joshi et al., 2020), a pretraining method based on BERT (Devlin et al., 2019).¹¹ The advantage of SpanBERT over BERT is that in a masked language modelling context, with masks spanning more than a single word, SpanBERT predicts sequences and takes into account their joint probability, whereas BERT trained with standard Masked Language Modelling can only predict single tokens independently. With SpanBERT, we

¹¹We use SpanBERT with the pretrained bert-large-cased model distributed as part of the transformers Python library.

generate replacements of between 1 and 3 tokens in length, in each case retaining only the best prediction out of the top k which matches one of the preposition sequences in the pre-extracted list.

After all candidates are generated, we apply GRUEN (Zhu and Bhat, 2020) to score the foils for grammaticality, and further apply the NLI model described in Section 4.3 to label the entailment relationship between caption and foil pairs. From the resulting data, we sample as follows: i) we keep only caption-foil pairs labelled as *contradiction*, where the GRUEN grammaticality score is ≥ 0.8 ; ii) for every caption-foil pair sampled where p is replaced with q , we search for another caption-foil pair where q is replaced with p , if present. This strategy yields a roughly balanced dataset, where no single preposition or preposition sequence is over-represented in captions or foils.

These processes result in an initial set of 614 cases, of which 535 (87.1%) are selected following manual validation described in §4.4.

Figure 2 shows proportions in captions and foils of the prepositions. E.g.: ‘A cat plays with a pocket knife [on](#) / [underneath](#) a table.’

As with plurals, we implement procedures for foil candidate generation by extending the *perturb* functionality in Checklist (Ribeiro et al., 2020).

A.5 Actions

The **action** piece consists of two instruments: i) **action replacement** and ii) **actant swap**. They are testing a V&L model’s capability of i) identifying whether an *action* mentioned in the textual modality matches the action seen in the image or not (e.g. ‘a man [shouts](#) / [smiles](#) at a woman’) and ii) correctly identifying the *participants* of an action and the *roles* they are playing in it (e.g., given the picture in Table 1: is it the man or the woman who shouts?).

Data source For creating interesting foils with *diverse* actions, we focus on the SWiG dataset (Pratt et al., 2020) that comprises 504 action verbs annotated with semantic roles and their fillers, which are grounded in images of the *imSitu* dataset (Yatskar et al., 2016). We generate English captions for the images using SimpleNLG (Gatt and Reiter, 2009)¹². For generation we use the specified *ac-*

tion verb, the realized FrameNet semantic roles and their annotated filler categories (see Table 1 for *shout*: AGENT: man, ADDRESSEE: woman), and generate short captions, with realization of two roles in active form. We apply various filters to ensure high quality of the generated captions using diverse metrics¹³ and manual checks through AMT crowdsourcing.

Foiling method When creating the **action replacement** instrument, we need to make sure that the action replacement suits the context. We propose action replacements with BERT (Devlin et al., 2019) that need to satisfy three conditions: 1) the proposed action verbs originate from the SWiG dataset – otherwise new verbs are introduced on the foil side only, which may induce biases; 2) the frequency distribution of action verbs on the caption and on the foil side is approximately the same (cf. Figure 3); 3) we constrain the replacement verbs to be either antonyms of the original verb or at least not synonyms, hyponyms or hypernyms to the original, according to WordNet (Fellbaum, 1998) in order to avoid situations where replacements are almost synonymous to the original action. The **actant swap** instrument is based on the original image annotations, but swaps the two role fillers (e.g., ‘A woman shouts at the man.’ for the image in Table 1). To avoid agreement mistakes, we *generate* these foils using the inverted role fillers as input.

The frequency distributions of words in which captions and foils differ, are plotted in Figure 3 for action replacement. The actant swap instrument is not visualised: By construction, actant swap cannot suffer from distributional bias since caption and foil contain the same words up to a *permutation*.

A.6 Coreference

The **coreference** piece consists of two pieces: **coreference standard** and **coreference clean**. It aims to uncover whether V&L models are able to perform pronoun coreference resolution. The coreference phenomenon encompasses both cases where i) the pronoun refers to a noun (phrase) and both the pronoun and the (noun) phrase are grounded

¹²SimpleNLG is a surface realization engine that – given some content and crucial syntactic specifications – performs surface generation including morphological adjustments.

¹³We use the GRUEN metric (Zhu and Bhat, 2020) that scores grammaticality, naturalness and coherence of generations and compute perplexity with GPT-2 to rank alternative outputs. We determined appropriate thresholds based on manual judgements of acceptability and chose the highest-ranked candidates. The final data quality is controlled by crowdsourced annotation with AMT.

in the visual modality (e.g. ‘A woman is driving a motorcycle. Is she wearing a helmet?’), and cases where ii) the pronoun refers directly to a region in the image or even to the whole image (e.g. ‘A man is sitting on a bench. Is this outside?’).

Data source We source the data from VisDial v1.0 (Das et al., 2017), which contains images from MSCOCO (Lin et al., 2014), their captions and dialogues about the images in form of Q&A sequences. To ensure that the coreference phenomenon is present in the [Caption. Question? Yes/No.] formulations, we check whether pronouns are present in the *question*. The list of pronouns and their frequencies in our train-val-test splits are represented in Figure 1.

The **coreference standard** instrument contains 916 data samples (708 are valid¹⁴) from the VisDial’s training set. The data of **coreference clean** instrument consisting of 141 samples (104 are valid), originates from VisDial’s validation set. With models that have been trained on VisDial, we would be in the situation where models are tested on their training data. Therefore we also have the *coreference clean instrument* based on the validation set of VisDial to test models safely. Unfortunately, we cannot use VisDial’s test set because the required question-answers annotations necessary for foiling are withheld.

Foiling method When foiling, we take the image description of the form [Caption. Question? Yes/No.] and exchange the answer: *yes* \rightarrow *no* and vice-versa (see example in Table 1). This way, we keep the full textual description including pronoun and noun (phrase) intact, hence ensuring that the coreference phenomenon is present and valid in the foil too, and rely on the model to interpret affirmation and negation correctly. Note that we rely on the capability of models to correctly interpret negation also in the existence piece (cf. §3.1).

Arguably, coreference is the most difficult phenomenon to foil in VALSE. Especially in cases where pronouns refer to a noun (phrase) (e.g., ‘A woman is driving a motorcycle. Is she wearing a helmet? Yes.’), exchanging the pronoun with another pronoun would generate incoherent and unlikely sequences¹⁵ (e.g., ‘A woman is driving a mo-

¹⁴The majority of manual annotators validated that the caption describes the image but the foil does not.

¹⁵Even more, the possibilities of exchanging pronouns with pronouns in grammatical ways are very limited: *she* – *he* but not *she* – *they* / *her* / *their*.

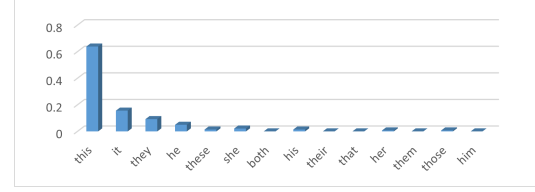


Figure 1: Normalized pronoun frequencies in the coreference subset.

torcycle. Is he wearing a helmet?’), and exchanging it with a noun phrase would furthermore break the pronoun coreference phenomenon because there would be no pronoun anymore (e.g., ‘A woman is driving a motorcycle. Is the man wearing a helmet?’). Therefore when foiling the coreference piece, we aim to keep the original description intact for ensuring the preservation of the coreference phenomenon. Hence we rely on the answers containing *yes* or *no*¹⁶ and exchange affirmative to negative answers and vice-versa.

A.7 FOIL it! data

We include an additional piece in VALSE consisting of 1000 randomly sampled entries from the *FOIL it!* dataset (Shekhar et al., 2017b). Each entry in *FOIL it!* consists of an MSCOCO (Lin et al., 2014) image and a foiled caption where a noun phrase depicting an object visible in the image was replaced by a semantically related noun phrase. Since examples in the *FOIL it!* dataset are linked to MSCOCO, we use these links to retrieve one correct caption from the five captions available for the image, and create an image–caption–foil triple. From the original 1000 entries, 943 have been validated by our manual annotation procedure (in Appendix E). Please refer to Shekhar et al. (2017b) for more details.

B Evaluation metrics

We evaluate pretrained V&L models on VALSE using **accuracy** (*acc*), the overall accuracy on all classes; **precision** or *positive predictive value* (*p_c*), which measures the proportion of correctly identified *correct captions*; and **foil precision** or *negative predictive value* (*p_f*), which measures the proportion of correctly identified *foiled examples*.

The **pairwise ranking accuracy** *acc_r* is computed using the image-sentence alignment score ϕ that the model assigns to correct and foiled image-

¹⁶If the answer is longer than just *yes/no* (e.g., ‘Yes, she is’) we shorten it to *yes/no*.

	CLIP (Radford et al., 2021)	LXMERT (Tan and Bansal, 2019)	ViLBERT (Lu et al., 2019)	ViLBERT 12-in-1 (Lu et al., 2020)	VisualBERT (Li et al., 2019)
model type	separate image and text encoders	dual stream	dual stream	dual stream	single stream
pretraining data	400M image-text pairs	MSCOCO	Conceptual Captions	Conceptual Captions	MSCOCO
pretraining tasks	ISA	ISA, MLM, MOP, VQA	ISA, MLM, MOP	ISA, MLM, MOP	ISA, MLM, MOP
finetuning	–	VQA	–	12 V&L tasks	–

Table 3: V&L models evaluated with VALSE in our experiments. **ISA**: image-sentence alignment; **MLM**: masked language modelling; **MOP**: masked object prediction; **VQA**: visual question answering.

text pairs. A prediction is considered successful, if given an image (i) paired with a correct (c) versus a foil (f) text, the score of the positive/correct pair is greater than that of the foiled pair.

$$acc_r = \frac{\sum_{(i,c) \in C} \sum_{f \in F} s(i, c, f)}{|C| + |F|},$$

$$s(i, c, f) = \begin{cases} 1, & \text{if } \phi(i, f) \leq \phi(i, c), \\ 0, & \text{otherwise,} \end{cases}$$

where C is the set of correct image-caption pairs (i, c), and F is the set of foils for the pair (i, c).

The **pairwise accuracy** acc_r is important for two reasons: First, it enables V&L models to be evaluated on VALSE without a binary classification head for classifying image-sentence pairs as correct or foiled. For example, CLIP (Radford et al., 2021) is a model that computes a score given an image-sentence pair. This score can be used to compare the scores of a correct image-sentence pair and the corresponding foiled pair. By contrast, a model like LXMERT (Tan and Bansal, 2019) has a binary image-sentence classification head and can predict a correct pair independently of the foiled pair (and vice-versa). Second, acc_r enables the evaluation of unimodal models on VALSE, as motivated in §4.2.

C Filtering methods

NLI filtering For NLI filtering we make use of the *HuggingFace* (Wolf et al., 2020) implementation of ALBERT (xxlarge-v2) that was already finetuned on the concatenation of SNLI (Bowman et al., 2015), MultiNLI (Williams et al., 2018), FEVER-NLI (Nie et al., 2019) and ANLI datasets (Nie et al., 2020). The model is the best performing on the ANLI benchmark leaderboard¹⁷ and it achieves 90% accuracy on MultiNLI devset.

¹⁷github.com/facebookresearch/anli

D Vision & Language and Unimodal Models

In Table 3 we summarise the five V&L models used in our experiments, their architecture, pretraining tasks and data, and finetuning tasks (if any).

CLIP CLIP (Radford et al., 2021) is composed of two transformer-based text and an image encoders. These are jointly trained on 400M image-text pairs through contrastive learning for predicting high scores for paired image-text examples and low scores when image-text samples are not paired in the dataset. CLIP has shown zero-shot capabilities in e.g. object classification, OCR, activity recognition (Radford et al., 2021). Goh et al. (2021) have shown the existence of multimodal neurons in CLIP, responding to the same topic regardless of whether it is represented in an image, drawing or handwritten text. We use CLIP’s image-text alignment scores for benchmarking on VALSE: Given an image, we compare whether CLIP¹⁸ predicts higher image-text similarity for the correct or for the foiled caption.

LXMERT LXMERT (Tan and Bansal, 2019) is a dual-stream transformer model combining V&L through cross-modal layers. It is pretrained on MSCOCO (Lin et al., 2014) and on multiple VQA datasets for (i) multimodal masked word and object prediction, (ii) image-sentence alignment, i.e., determining whether a text corresponds to an image or not, and (iii) question-answering. For benchmarking on VALSE, we use LXMERT’s¹⁹ image-sentence alignment head.

ViLBERT and ViLBERT 12-in-1 ViLBERT (Lu et al., 2019) is a BERT-based transformer architecture that combines V&L on two separate streams

¹⁸github.com/openai/CLIP

¹⁹github.com/huggingface/transformers

Piece	Instrument	#Instances	#Valid (%)	#Unanimous (%)	#Lexical Items	JS-div	JS-div valid
Existence	<i>Existential quantifiers</i>	534	505 (94.6)	410 (76.8)	25	0.628	0.629
Plurality	<i>Semantic Number</i>	1000	851 (85.1)	617 (61.7)	704	0.742	0.766
Counting	<i>Balanced</i>	1000	868 (86.8)	598 (59.8)	25	0.070	0.082
	<i>Small numbers</i>	1000	900 (90.0)	637 (63.7)	4	0.059	0.071
	<i>Adversarial</i>	756	691 (91.4)	522 (69.0)	27	1.000	1.000
Relations	<i>Prepositions</i>	614	535 (87.1)	321 (52.3)	38	0.083	0.114
Actions	<i>Replacement</i>	779	648 (83.2)	428 (54.9)	262	0.437	0.471
	<i>Actant swap</i>	1042	949 (91.1)	756 (72.6)	467	0.000	0.000
Coreference	<i>standard: VisDial train</i>	916	708 (77.3)	499 (54.5)	2	0.053	0.084
	<i>clean: VisDial val</i>	141	104 (73.8)	69 (48.9)	2	0.126	0.081
Foil-It!	<i>noun replacement</i>	1000	943 (94.3)	811 (81.1)	73	0.426	0.425
Overall		8782	7702 (87.7)	5668 (73.6)			

Table 4: Manual validation results for each piece in VALSE, as well as for the Foil-it dataset. *Valid*: number (percent) of cases for which at least 2 out of 3 annotators chose the caption; *Unanimous*: number (percent) of cases for which all annotators chose the caption; *Lexical Items*: number of phrases or lexical items in the vocabulary that differs between foils and captions; *JS-div*: Jensen-Shannon divergence between foil-caption distributions for the whole instrument; *JS-div valid*: Jensen-Shannon divergence between foil-caption distribution for the valid subset of the instrument, after sub-sampling.

by co-attention layers. It is pretrained on Google Conceptual Captions (Sharma et al., 2018) on (i) multimodal masked word and object prediction; and (ii) image-sentence alignment. ViLBERT 12-in-1 (Lu et al., 2020) further finetuned a ViLBERT model checkpoint on 12 different tasks including VQA, image retrieval, phrase grounding and others.²⁰ We use the image-sentence alignment head of the publicly available model checkpoints for ViLBERT²¹ and ViLBERT 12-in-1²².

VisualBERT VisualBERT (Li et al., 2019) is also a BERT-based transformer. Its single-stream architecture encodes image regions and linguistic features via a transformer stack, using self-attention to discover the alignments between the two modalities. VisualBERT is pretrained on MSCOCO captions (Chen et al., 2015) on two tasks: (i) masked language modelling, and (ii) sentence-image prediction. The latter is framed as an extension of the next sentence prediction task used with BERT. Inputs consist of an image and a caption, with a second caption which has a 50% probability of being random. The goal is to determine if the second caption is also aligned to the image. In our experiments, we use the publicly

available implementation of VisualBERT²³.

GPT-1 and GPT-2 – Unimodal models GPT1 (Radford et al., 2018) and GPT2 (Radford et al., 2019) are transformer-based autoregressive language models pretrained on English data through self-supervision. We test whether our benchmark is solvable by these unimodal models by computing the perplexity of the correct sentence and compare it to the perplexity of the foiled sentence. In case the computed perplexity is higher for the foil than for the correct sentence, we assume that the correctly detected foiled caption may possibly suffer from a **plausibility bias** (as described in section 4.2) or from other biases (e.g. a model’s preference towards affirmative or negative sentences).

E Mechanical Turk Annotation and Evaluation

Setup The validation study was conducted on all the data for each instrument in VALSE, as well as for the FOIL it! data (Shekhar et al., 2019b). Each instance consisted of an image, a caption and a foiled version of the caption, as shown in Figure 4. Annotators received the following general instructions:

You will see a series of images, each accompanied by two short texts. Your task is to judge which of the two texts accurately describes what can be seen in the image.

²⁰github.com/facebookresearch/vilbert-multi-task

²¹https://dl.fbaipublicfiles.com/vilbert-multi-task/pretrained_model.bin

²²https://dl.fbaipublicfiles.com/vilbert-multi-task/multi_task_model.bin

²³github.com/uclanlp/visualbert

Each instance was accompanied by the caption and the foil, with the ordering balanced so that the caption appeared first 50% of the time. In each instance, the caption and foil were placed above each other, with the differing parts highlighted in bold. Annotators were asked to determine *which of the two sentences accurately describes what can be seen in the image?* In each case, they had to choose between five options: (a) the first, but not the second; (b) the second, but not the first; (c) both of them; (d) neither of the two; and (e) I cannot tell. We collected three annotations for each instance, from three independent workers.

Annotator selection We recruited annotators who had an approval rating of 90% or higher on Amazon Mechanical Turk. We ran an initial, pre-selection study with 10 batches of 100 instances each, in order to identify annotators who understood the instructions and performed the task adequately. The pre-selection batches were first manually annotated by the authors, and we identified ‘good’ annotators based on the criterion that they preferred the caption to the foil at least 70% of the time. Based on this, we selected a total of 63 annotators. Annotators were paid \$0.05 per item (i.e. per HIT on Mechanical Turk).

Results Table 4 shows, for each instrument, the number of instances in total, as well as the proportion of instances which we consider *valid*, that is, those for which at least two out of three annotators chose the caption, *but not the foil*, as the text which accurately describes the image. We also show the number of instances for which annotators unanimously (3/3) chose the caption.

Bias check While measures were taken to control for distributional bias between captions and foils in the different pieces of VALSE (cf. §4.1), it is possible that sub-sampling after manual validation could reintroduce such biases. To check that this is not the case, we compare the *word frequency distributions between captions and foils* in the original pieces, and the word frequency distribution of the manually validated set. We report the Jensen-Shannon divergence and the number of words that differ between caption and foil in Table 4. The foil-caption word frequency distributions can be inspected in Figures 2 and 3. The Jensen-Shannon (JS) divergence is defined as:

$$JS(f \parallel c) = \sqrt{\frac{KL(f \parallel m) + KL(c \parallel m)}{2}}$$

where f is the normalized word frequency for foils, c the normalized word frequency for captions, m is the point-wise mean of f and c , and KL is the Kullback-Leibler divergence.

As Table 4 shows, the JS-divergence between caption and foil distributions remains the same, or changes only marginally (compare columns *JS-div* and *Js-div valid*, where *#Lexical Items* indicates the number of lexical/phrasal categories in the relevant distributions). This indicates that no significant bias was introduced as a result of subsampling after manual validation.

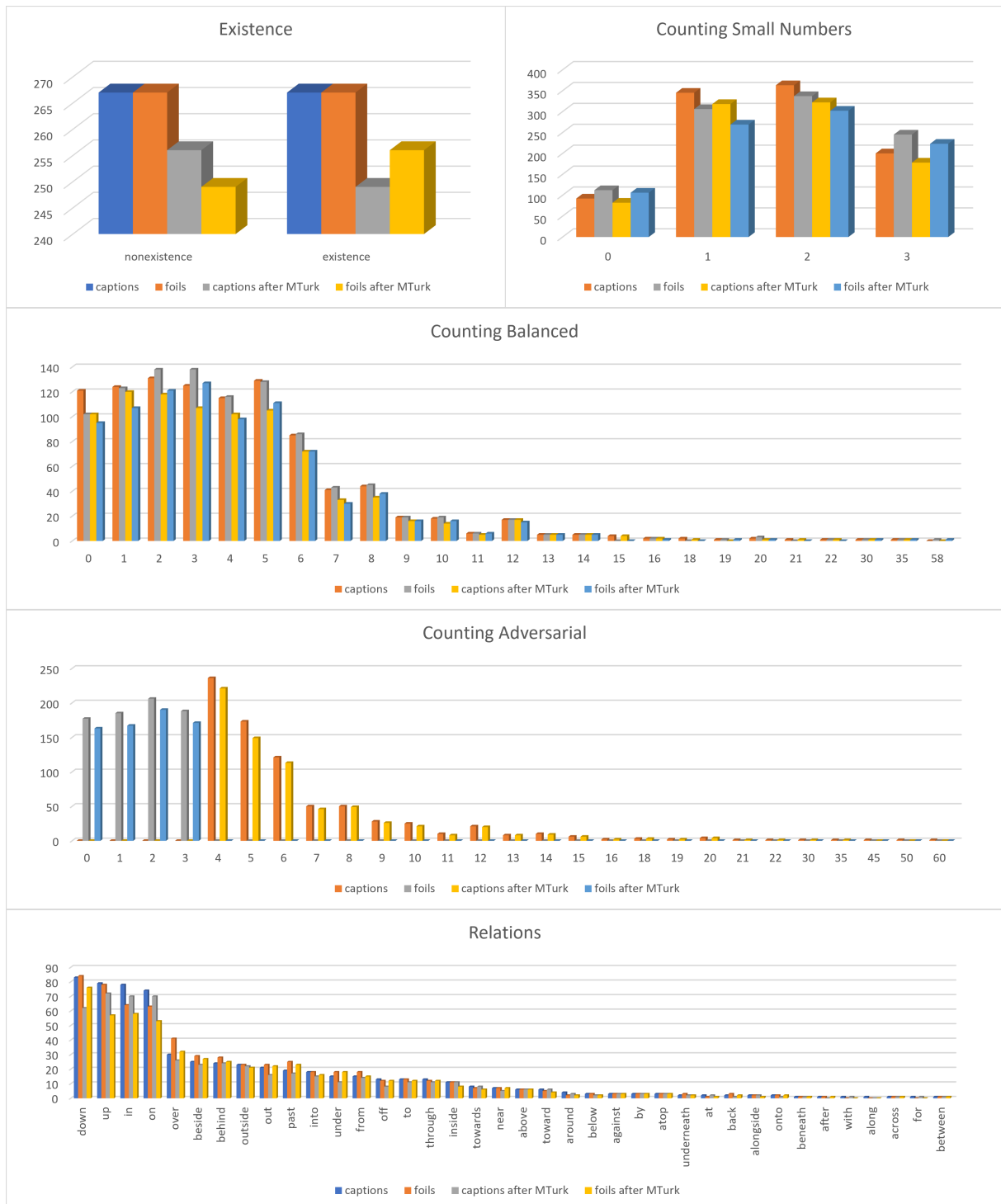


Figure 2: Word frequency distributions for captions and foils before and after the manual validation for existence, counting and relations.

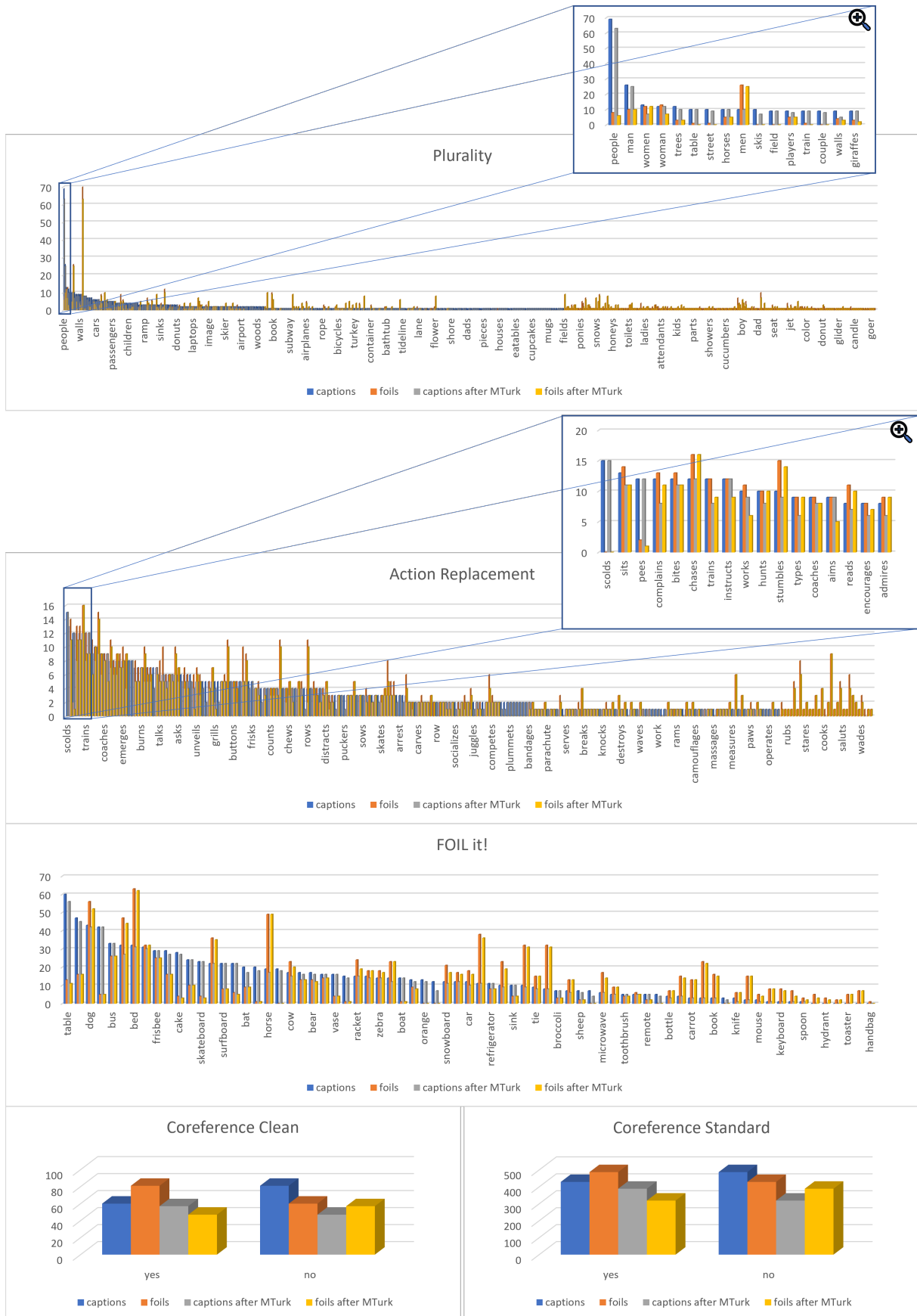


Figure 3: Word frequency distributions for captions and foils before and after the manual validation for plurality, action replacement and FOIL it. The actant swap instrument is not visualized here: By construction, actant swap cannot suffer from distributional bias since caption and foil contain the same words up to a *permutation*.


1. There is exactly **1** animal pictured.

2. There are exactly **6** animals pictured.

Instructions

Shortcuts

These sentences are almost identical, but differ in a few words highlighted in boldface. Which of the two sentences accurately describes what can be seen i...



Select an option

The first one, but not the second	1
The second one, but not the first	2
Neither of the two	3
Both of them	4
I cannot tell	5

Figure 4: Example of an instance from the validation study. The example is from the Counting piece, *adversarial* instrument (see Section 3.3).