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

**Anonymous ACL submission**

## Abstract

We propose VALSE (**V**ision **A**nd **L**anguage **S**tructured **E**valuation), a novel benchmark designed for testing general-purpose pretrained vision and language (V&L) models for their *visio-linguistic grounding* capabilities on *specific linguistic phenomena*. VALSE offers a suite of six tests covering various linguistic constructs. Solving these requires models to ground linguistic phenomena in the visual modality, allowing more fine-grained evaluations than hitherto possible. We build VALSE using methods that support the construction of *valid* 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

General-purpose pretrained vision and language (V&L) models have gained notable performance on many V&L tasks (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, V&L research has changed its focus from task-specific architectures to fine-tuning large V&L models.

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*. Specifically, we lack detailed understanding of the extent to which such models are able to ground linguistic phenomena—from morphosyntax to semantics—in the visual modality (Bernardi and Pezzelle, 2021). For example, recent evidence suggests that models are insensitive to linguistic

distinctions of verb-argument structure (Hendricks and Nematzadeh, 2021) and word order (Cirik et al., 2018; Akula et al., 2020).

Our work addresses this gap with VALSE (Vision And Language Structured Evaluation), a benchmark for V&L model evaluation comprising six tasks, or 'pieces', where each piece has the same structure: given a visual input, a model is asked to distinguish real captions from *foils*, where a foil is constructed from a caption by altering a word or phrase that realizes a *specific linguistic phenomenon*, e.g., semantic number of nouns, verb argument structure, or coreference. VALSE uses a resource-lean diagnostic setup that dispenses with 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 gauging the sensitivity of pre-trained V&L models 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* foils that include automatic *and* human validation. We balance *word frequency distributions* between captions and foils, and test against

pretrained models solving the benchmark *unimodally*. We employ *masked language modeling* (MLM) in foil creation and *semantic inference* for validating foils, and finally collect *human annotations* for the entire benchmark.

iv) We establish initial experimental results for pretrained V&L models of diverse architectures on VALSE. These models' overall weak performance indicates that the time is ripe for a novel, reliable foiling dataset targeting 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. 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 pervasive linguistic phenomena that span *multiple tokens*, described in §3.1–§3.6. At the same time, we ensure that our data is robust to perturbations and artefacts by i) controlling for word frequency biases between captions and foils, and ii) testing against *unimodal collapse*, a known issue of V&L models (Goyal et al., 2017; Madhyastha et al., 2018), thereby preventing models from solving the task using a single input modality. 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 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

---

[1]We release our dataset containing all annotators' votes (Prabhakaran et al., 2021) and code upon acceptance.

| pieces | existence | plurality | counting | relations | actions | coreference |
|---|---|---|---|---|---|---|
| **Data collection & metadata** | | | | | | |
| instruments | *existential quantifiers* | *semantic number* | *balanced, adversarial, small numbers* | *prepositions* | *replacement, actant swap* | *standard, clean* |
| #examples† | 505 | 851 | 2,459 | 535 | 1,633 | 812 |
| foil generation method | *nothing ↔ something* | NP replacement (`sg2pl`; `pl2sg`) & quantifier insertion | numeral replacement placement | SpanBERT prediction | action replacement, actant swap | *yes ↔ no* |
| MLM | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ |
| GRUEN | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ |
| NLI | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ |
| src. dataset | Visual7W | MSCOCO | Visual7W | MSCOCO | SWiG | VisDial v1.0 |
| image src. | MSCOCO | MSCOCO | MSCOCO | MSCOCO | SituNet | MSCOCO |
| **Example data** | | | | | | |
| caption (blue) / foil (orange) | *There are no animals / animals shown.* | *A small copper vase with some flowers / exactly one flower in it.* | *There are four / six zebras.* | *A cat plays with a pocket knife on / underneath a table.* | *A man / woman shouts at a woman / man.* | *Buffalos walk along grass. Are they in a zoo? No / Yes.* |
| 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 5 (in the Appendix) we list the number of examples in each individual instrument.

under study. FFs take a *correct caption* as input and change a specific part 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 generated 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.

VALSE data is sourced 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' → '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–

3

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 (pl2sg: 'some flowers' → 'exactly one flower'), or conversely, the caption contains a singular which is foiled by replacing it with a plural (sg2pl). 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 and foil classes are balanced, so that models that exploit class frequency bias are penalized. In **counting adversarial** we ensure that all foils take class $n \in \{0, 1, 2, 3\}$, whereas all correct captions take class $m \in \{m \mid m \geq 4\}$. Biased models are expected to favour more frequent classes. Since small

---

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

numbers are naturally the most frequent, models that resort to such biases should perform poorly on this adversarial test set. **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 with 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 candidates 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 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

4

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 *Reliable* construction of *valid* foils

In VALSE, an instance consisting of an image-caption-foil triple is considered *valid* if: the foil minimally differs from the original caption; the foil does not accurately describe the image; and independent judges agree that the caption, but not the foil, is an accurate description of the image. We consider a *foiling method* to be more *reliable* the more it ensures that a generated foil does not substantially differ from a human caption regarding distributional and plausibility bias, and cannot be easily solved unimodally.

In this section, we discuss automatic and manual means to reliably construct valid foils. In this context, two types of bias are especially worthy of note: distributional bias (§4.1) and plausibility bias (§4.2). In §4.3 we discuss how we apply a natural language inference model to filter examples in our data pipeline, and §4.4 show how we manually validate *all examples* in our benchmark. Random

---

[3]VisDial annotations are not available for the test set.

---

samples from the final version of each instrument are shown in Tab. 6–11.

### 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 procedures yielding foils that are implausible or unnatural, which can 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 and correct captions are similarly plausible, we use language models such as BERT (Devlin et al., 2019) and SpanBERT (Joshi et al., 2020) to suggest 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 originally proposed to automatically score generated sentences based on discourse-level and grammatical properties. We use only the grammaticality component of GRUEN, and retain only foil candidates with a grammaticality score $\geq 0.8$.

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.

5

## 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 valid (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* ($N = 60$, 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]See 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.

[5]By 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.

## 4.4 Manual evaluation of generated foils

As a final step, the data for each instrument was submitted to a manual validation. 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. 2, App. 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 App. E. For the final version of VALSE, we include instances which passed the following validation test: 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%), 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.

We compute inter-annotator agreement for each instrument (Tab. 5). On the valid subset, agreement is low to medium (Krippendorff's $\alpha$: min=0.23, max=0.64, mean=0.42, sd=0.12). We note that there is considerable variation in the number of annotations made by individuals, and $\alpha$ is computed over 5 categories. Hence, this result cannot be straightforwardly interpreted as a ceiling of human performance for VALSE. However, $\alpha$ is higher for pieces on which models also perform better (e.g. existence, Foil-It!; cf. §5).

## 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).

| Metric | Model | Existence quantifiers | Plurality number | Counting balanced | Counting sns.† | Counting adv.† | Sp.rel.‡ relations | Action repl.† | Action actant swap | Coreference standard | Coreference clean | Foil-it! | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 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 | 48.7 | 69.5 | **77.2** | 65.4 | 72.2 | 45.6 | 45.2 | 77.5 | 60.7 |
| | GPT2* | 58.0 | 51.9 | 51.6 | 49.8 | 45.3 | 75.0 | 66.8 | **76.9** | 54.5 | 50.0 | 80.7 | 60.1 |
| | CLIP | 66.9 | 56.2 | 62.1 | 62.5 | 57.5 | 64.3 | **75.6** | 68.6 | 52.1 | 49.7 | **88.8** | 64.0 |
| | LXMERT | 78.6 | 64.4 | 62.2 | 69.2 | 42.6 | 60.2 | 54.8 | 45.8 | 46.8 | 44.2 | 87.1 | 59.6 |
| | ViLBERT | 65.5 | 61.2 | 58.6 | 62.9 | 73.7 | 57.2 | 70.7 | 68.3 | 47.2 | 48.1 | 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 | 39.7 | 45.7 | 48.2 | 48.2 | 50.0 | 39.7 | 49.2 | 44.4 | 49.5 | 47.6 | 48.5 | 46.4 |
| $acc$ | 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 | 2.4 | 50.3 | 50.7 | 50.6 | 51.8 | 49.9 | 52.6 | 50.4 | 50.0 | 50.0 | 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 | 49.3 | 46.5 | 48.3 | 47.8 | 50.0 | 49.3 | 48.8 | 49.7 | 50.0 | 50.0 | 46.6 | 48.8 |
| $\min(p_c, p_f)$ | LXMERT | 41.6 | **42.2** | 50.9 | 50.0 | 37.3 | **28.4** | 35.8 | 36.8 | 18.4 | 17.3 | 69.3 | 38.9 |
| | ViLBERT | 47.9 | 2.1 | 24.4 | 24.7 | 17.5 | 1.5 | 11.9 | 7.1 | 1.3 | 1.9 | 12.9 | 13.9 |
| | 12-in-1 | **85.0** | 33.4 | **64.3** | **61.7** | **59.5** | 13.3 | **47.8** | **37.6** | 15.8 | 13.5 | 48.8 | **43.7** |
| | VisualBERT | 1.3 | 0.3 | 0.0 | 0.0 | 0.0 | 1.3 | 0.0 | 0.0 | 0.0 | 0.0 | 0.2 | 0.3 |
| $AUROC$ ×100 | LXMERT | 60.5 | 57.3 | 53.8 | 57.7 | 50.5 | 51.9 | 52.1 | 47.6 | 49.8 | 49.5 | 76.9 | 55.2 |
| | ViLBERT | 52.5 | 54.1 | 50.8 | 51.6 | 53.5 | 51.2 | 57.2 | **57.8** | 49.9 | 49.9 | 75.2 | 54.9 |
| | 12-in-1 | **96.3** | **67.4** | **72.0** | **77.8** | **75.1** | **55.8** | **61.3** | 55.0 | **59.8** | **59.6** | **81.0** | **69.2** |
| | VisualBERT | 28.9 | 29.0 | 24.5 | 16.5 | 20.9 | 45.2 | 17.7 | 36.3 | 45.3 | 46.3 | 28.5 | 30.8 |

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 naïvely 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).

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 five metrics[6] 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; **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; and **area under the receiver operating characteristic curve** (AUROC), which measures how well models distinguish correct vs. foiled examples across different prediction thresholds. $acc_r$ is more permissive than $acc$ as it accepts model predictions if the score for a foil is lower

---

[6]All metrics are defined in Appendix B.

than the caption's score. Our main metrics are AUROC and $acc_r$. $acc_r$ gives results for a pair ⟨*image, caption*⟩ and ⟨*image, foil*⟩. Both AUROC and $acc_r$ are well suited to evaluate minimally-edited pairs as neither uses a classification threshold. As for $p_c$ and $p_f$, since these are competing metrics where naively increasing one can decrease the other, we report the smaller of the two as an indicator of how informed model predictions are. Since all instruments are implemented as a balanced 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 and are pretrained on a variety of tasks with different training data. We also benchmark two unimodal text-only models, GPT1 (Radford et al., 2018) and GPT2 (Radford et al., 2019). See Appendix D for details on all these models used 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. Also the spatial relations (77.2%), action replacement (66.8%) and actant swap (76.9%) instruments suggest plausibility biases in 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 replacing 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 latter obtains high $p_f$ but very low $p_c$ values—reflected in the $\min(p_c, p_f)$ scores—indicating that VisualBERT learned a heuristic that does not generalise (see Hendricks and Nematzadeh, 2021, for similar observations with other models). We hypothesise that this is due to the way image-sentence alignment is framed in VisualBERT's pretraining: the model expects an image and a correct sentence $c_1$, and predicts whether a second sentence $c_2$ is a match. During pretraining $c_1$ and $c_2$ are likely to differ in many ways, whereas in our setting, they are nearly identical. 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 or counting them in an image; ii) correctly classifying a named spatial relation between objects in an image; iii) distinguishing actions and identifying their participants, even if supported by preference biases; or, iv) tracing multiple references to the same object in an image through the use of pronouns.

**Correct vs. 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 $p_f$ (cf. ViLBERT), or vice-versa (cf. VisualBERT). This suggests that such models are insensitive to VALSE's inputs: models that almost always predict a match will inflate $p_f$ at the expense of $p_c$. $\min(p_c, p_f)$ reveals that VisualBERT and ViLBERT perform poorly and below 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, and actions. These tendencies we see reflected in our main metrics, $acc_r$ and AUROC.

## 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 (as shown by the existence piece), but struggle to ground their interdependence and relationships in visual scenes when forced to respect linguistic indicators. We encourage the community to use VALSE for measuring 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.

---

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

# 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.

Arjun Akula, Spandana Gella, Yaser Al-Onaizan, Song-Chun Zhu, and Siva Reddy. 2020. Words aren't enough, their order matters: On the robustness of grounding visual referring expressions. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6555–6565, Online. Association for Computational Linguistics.

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*.

Volkan Cirik, Louis-Philippe Morency, and Taylor Berg-Kirkpatrick. 2018. Visual referring expression recognition: What do systems actually learn? 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 781–787, New Orleans, Louisiana. Association for Computational Linguistics.

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.

9

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.

Lisa Anne Hendricks and Aida Nematzadeh. 2021. Probing Image-Language Transformers for Verb Understanding. *arXiv*, 2106.09141.

Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension systems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2021–2031, Copenhagen, Denmark. Association for Computational Linguistics.

Robin Jia, Aditi Raghunathan, Kerem Göksel, and Percy Liang. 2019. Certified robustness to adversarial word substitutions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4129–4142, Hong Kong, China. Association for Computational Linguistics.

Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2020. SpanBERT: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics*, 8:64–77.

Kushal Kafle, Robik Shrestha, and Christopher Kanan. 2019. Challenges and prospects in vision and language research. *Frontiers in Artificial Intelligence*, 2:28.

Wonjae Kim, Bokyung Son, and Ildoo Kim. 2021. Vilt: Vision-and-language transformer without convolution or region supervision. *arXiv preprint arXiv:2102.03334*.

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 *International Conference on Learning Representations*.

Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In *Thirteenth International Conference on the Principles of Knowledge Representation and Reasoning*.

Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. 2020a. Unicoder-vl: A universal encoder for vision and language by cross-modal pretraining. In *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*, pages 11336–11344. AAAI Press.

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. *arXiv preprint arXiv:2106.04632*.

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 *Arxiv*.

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 *European Conference on Computer Vision*, pages 121–137. Springer.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer Vision – ECCV 2014*, pages 740–755, Cham. Springer International Publishing.

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 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 13480–13488.

Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems*, pages 13–23.

Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 2020. 12-in-1: Multi-task vision and language representation learning. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Pranava Madhyastha, Josiah Wang, and Lucia Specia. 2019. VIFIDEL: Evaluating the visual fidelity of image descriptions. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6539–6550, Florence, Italy. Association for Computational Linguistics.

Pranava Swaroop Madhyastha, Josiah Wang, and Lucia Specia. 2018. Defoiling foiled image captions. 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 433–438, New

10

Orleans, Louisiana. Association for Computational Linguistics.

Yixin Nie, Haonan Chen, and Mohit Bansal. 2019. Combining fact extraction and verification with neural semantic matching networks. In *Association for the Advancement of Artificial Intelligence (AAAI)*.

Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial NLI: A new benchmark for natural language understanding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4885–4901, Online. Association for Computational Linguistics.

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 *Proceedings of the 1st Workshop on Multimodal Semantic Representations (MMSR)*, pages 32–44, Groningen, Netherlands (Online). Association for Computational Linguistics.

Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.

Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. 2015. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pages 2641–2649.

Vinodkumar Prabhakaran, Aida Mostafazadeh Davani, and Mark Díaz. 2021. On releasing annotator-level labels and information in datasets.

Sarah Pratt, Mark Yatskar, Luca Weihs, Ali Farhadi, and Aniruddha Kembhavi. 2020. Grounded situation recognition. In *Computer Vision - ECCV 2020 - 16th European Conference*, pages 314–332.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*.

Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of NLP models with CheckList. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics.

Daniel Rosenberg, Itai Gat, Amir Feder, and Roi Reichart. 2021. Are VQA systems RAD? Measuring robustness to augmented data with focused interventions. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 61–70, Online. Association for Computational Linguistics.

Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, Melbourne, Australia. Association for Computational Linguistics.

Ravi Shekhar, Sandro Pezzelle, Aurélie Herbelot, Moin Nabi, Enver Sangineto, and Raffaella Bernardi. 2017a. Vision and language integration: Moving beyond objects. In *IWCS 2017 — 12th International Conference on Computational Semantics — Short papers*.

Ravi Shekhar, Sandro Pezzelle, Yauhen Klimovich, Aurélie Herbelot, Moin Nabi, Enver Sangineto, and Raffaella Bernardi. 2017b. FOIL it! find one mismatch between image and language caption. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 255–265, Vancouver, Canada. Association for Computational Linguistics.

Ravi Shekhar, Ece Takmaz, Raquel Fernández, and Raffaella Bernardi. 2019a. Evaluating the representational hub of language and vision models. In *Proceedings of the 13th International Conference on Computational Semantics - Long Papers*, pages 211–222, Gothenburg, Sweden. Association for Computational Linguistics.

Ravi Shekhar, Aashish Venkatesh, Tim Baumgärtner, Elia Bruni, Barbara Plank, Raffaella Bernardi, and Raquel Fernández. 2019b. Beyond task success: A closer look at jointly learning to see, ask, and Guess-What. 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 2578–2587, Minneapolis, Minnesota. Association for Computational Linguistics.

Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2020. Vl-bert: Pretraining of generic visual-linguistic representations.

In *International Conference on Learning Representations*.

Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. 2019. A corpus for reasoning about natural language grounded in photographs. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6418–6428, Florence, Italy. Association for Computational Linguistics.

Hao Tan and Mohit Bansal. 2019. LXMERT: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5100–5111, Hong Kong, China. Association for Computational Linguistics.

Eric Wallace, Matt Gardner, and Sameer Singh. 2020. Interpreting predictions of NLP models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts*, pages 20–23, Online. Association for Computational Linguistics.

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

Junlin Wang, Jens Tuyls, Eric Wallace, and Sameer Singh. 2020b. Gradient-based analysis of NLP models is manipulable. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 247–258, Online. Association for Computational Linguistics.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.

Thomas Wolf, Julien Chaumond, Lysandre Debut, Victor Sanh, Clement Delangue, Anthony Moi, Pierric Cistac, Morgan Funtowicz, Joe Davison, Sam Shleifer, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45.

Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. 2019. Visual Entailment: A Novel Task for Fine-Grained Image Understanding. *arXiv*, 1901.06706.

Mark Yatskar, Luke Zettlemoyer, and Ali Farhadi. 2016. Situation recognition: Visual semantic role labeling for image understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Wanzheng Zhu and Suma Bhat. 2020. GRUEN for evaluating linguistic quality of generated text. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 94–108, Online. Association for Computational Linguistics.

Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. 2016. Visual7W: Grounded Question Answering in Images. In *IEEE Conference on Computer 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' → '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 tokenizer[8] 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 <u>is</u> 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 <u>some flowers</u> in it', arguably still entails the version with the singular '(...) <u>a flower</u>'. 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. '<u>exactly one</u> 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

---

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

[9] NP 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/).

13

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 $\geq 0.8$, 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 4 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.

**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 3 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

---

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

[11]We 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 descibed 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 $\geq 0.8$; 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 3 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 *action 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 metrics[13] 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 4); 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.

We plot caption and foil word frequency distributions for action replacement in Figure 4. We do not plot statistics for the actant swap instrument since by construction it 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

---

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

[13]We 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 valid[14]) 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* →*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 sequences[15] (e.g., 'A woman is driving a mo-

---

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

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

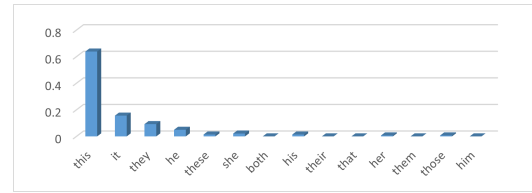Figure 1: Normalized pronoun frequencies in the coreference subset.

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; **pairwise ranking accuracy** $acc_r$, computed using the image-sentence alignment score $\phi$ that the model assigns to correct and foiled image-text pairs; and

---

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

**area under the receiver operating characteristic curve** (AUROC)—a classic metric used in machine learning classification problems—which in our case measures how well models distinguish correct vs. foiled examples across different prediction thresholds. The AUROC has a probabilistic interpretation and can be understood as the probability that a model will assign a higher score to a randomly chosen correct example relative to a randomly chosen foil.

With $acc_r$, 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. In Table 4, we show results for all models investigated according to all above-mentioned metrics.

## 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 leaderboard[17] and it achieves 90% accuracy on MultiNLI devset.

---

[17]github.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 CLIP[18] 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's[19] 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 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

---

[18]github.com/openai/CLIP
[19]github.com/huggingface/transformers
[20]github.com/facebookresearch/vilbert-multi-task

|  | **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.

| Metric | Model | Existence quantifiers | Plurality number | Counting balanced | Counting sns.† | Counting adv.† | Sp.rel.‡ relations | Action repl.† | Action actant swap | Coreference standard | Coreference clean | Foil-it! | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 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 | 48.7 | 69.5 | **77.2** | 65.4 | 72.2 | 45.6 | 45.2 | 77.5 | 60.7 |
|  | GPT2* | 58.0 | 51.9 | 51.6 | 49.8 | 45.3 | 75.0 | 66.8 | **76.9** | 54.5 | 50.0 | 80.7 | 60.1 |
|  | CLIP | 66.9 | 56.2 | 62.1 | 62.5 | 57.5 | 64.3 | **75.6** | 68.6 | 52.1 | 49.7 | **88.8** | 64.0 |
|  | LXMERT | 78.6 | 64.4 | 62.2 | 69.2 | 42.6 | 60.2 | 54.8 | 45.8 | 46.8 | 44.2 | 87.1 | 59.6 |
|  | ViLBERT | 65.5 | 61.2 | 58.6 | 62.9 | 73.7 | 57.2 | 70.7 | 68.3 | 47.2 | 48.1 | 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 | 39.7 | 45.7 | 48.2 | 48.2 | 50.0 | 39.7 | 49.2 | 44.4 | 49.5 | 47.6 | 48.5 | 46.4 |
| $acc$ | 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 | 2.4 | 50.3 | 50.7 | 50.6 | 51.8 | 49.9 | 52.6 | 50.4 | 50.0 | 50.0 | 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 | 49.3 | 46.5 | 48.3 | 47.8 | 50.0 | 49.3 | 48.8 | 49.7 | 50.0 | 50.0 | 46.6 | 48.8 |
| $p_c$ | LXMERT | 41.6 | 68.0 | 50.9 | 50.0 | 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 | 1.3 | 0.3 | 0.0 | 0.0 | 0.0 | 1.3 | 0.0 | 0.0 | 0.0 | 0.0 | 0.2 | 0.3 |
| $p_f$ | LXMERT | 70.1 | 42.2 | 53.0 | 60.8 | 37.3 | 28.4 | 66.4 | 60.2 | 18.4 | 17.3 | 69.3 | 47.6 |
|  | ViLBERT | 47.9 | 2.1 | 24.4 | 24.7 | 17.5 | 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 | 13.3 | 47.8 | 37.6 | 15.8 | 13.5 | 45.9 | 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 |
| $\min(p_c, p_f)$ | LXMERT | 41.6 | **42.2** | 50.9 | 50.0 | 37.3 | **28.4** | 35.8 | 36.8 | 18.4 | 17.3 | 69.3 | 38.9 |
|  | ViLBERT | 47.9 | 2.1 | 24.4 | 24.7 | 17.5 | 1.5 | 11.9 | 7.1 | 1.3 | 1.9 | 12.9 | 13.9 |
|  | 12-in-1 | **85.0** | 33.4 | **64.3** | **61.7** | 59.5 | 13.3 | **47.8** | **37.6** | 15.8 | 13.5 | 48.8 | **43.7** |
|  | VisualBERT | 1.3 | 0.3 | 0.0 | 0.0 | 0.0 | 1.3 | 0.0 | 0.0 | 0.0 | 0.0 | 0.2 | 0.3 |
| $AUROC \times 100$ | LXMERT | 60.5 | 57.3 | 53.8 | 57.7 | 50.5 | 51.9 | 52.1 | 47.6 | 49.8 | 49.5 | 76.9 | 55.2 |
|  | ViLBERT | 52.5 | 54.1 | 50.8 | 51.6 | 53.5 | 51.2 | 57.2 | 57.8 | 49.9 | 49.9 | 75.2 | 54.9 |
|  | 12-in-1 | **96.3** | **67.4** | **72.0** | **77.8** | 75.1 | 55.8 | 61.3 | 55.0 | 59.8 | 59.6 | 81.0 | 69.2 |
|  | VisualBERT | 28.9 | 29.0 | 24.5 | 16.5 | 20.9 | 45.2 | 17.7 | 36.3 | 45.3 | 46.3 | 28.5 | 30.8 |

Table 4: 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 naïvely 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).

ViLBERT[21] and ViLBERT 12-in-1[22].

**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 VisualBERT[23].

**GPT-1 and GPT-2 – Unimodal models** GPT1 (Radford et al., 2018) and GPT2 (Radford et al.,

---

[21] https://dl.fbaipublicfiles.com/vilbert-multi-task/pretrained_model.bin
[22] https://dl.fbaipublicfiles.com/vilbert-multi-task/multi_task_model.bin

[23] github.com/uclanlp/visualbert

| Piece | Instrument | #Inst. | #Valid (%) | #Unan. (%) | #Lex.it. | JS | JS Val. | $\alpha$ | $\alpha$ Valid |
|---|---|---|---|---|---|---|---|---|---|
| **Existence** | *Existential quantifiers* | 534 | 505 (94.6) | 410 (76.8) | 25 | 0.628 | 0.629 | 0.607 | 0.644 |
| **Plurality** | *Semantic Number* | 1000 | 851 (85.1) | 617 (61.7) | 704 | 0.742 | 0.766 | 0.303 | 0.359 |
| **Counting** | *Balanced* | 1000 | 868 (86.8) | 598 (59.8) | 25 | 0.070 | 0.082 | 0.361 | 0.423 |
| | *Small numbers* | 1000 | 900 (90.0) | 637 (63.7) | 4 | 0.059 | 0.071 | 0.417 | 0.473 |
| | *Adversarial* | 756 | 691 (91.4) | 522 (69.0) | 27 | 1.000 | 1.000 | 0.387 | 0.441 |
| **Relations** | *Prepositions* | 614 | 535 (87.1) | 321 (52.3) | 38 | 0.083 | 0.114 | 0.210 | 0.229 |
| **Actions** | *Replacement* | 779 | 648 (83.2) | 428 (54.9) | 262 | 0.437 | 0.471 | 0.229 | 0.318 |
| | *Actant swap* | 1042 | 949 (91.1) | 756 (72.6) | 467 | 0.000 | 0.000 | 0.386 | 0.427 |
| **Coreference** | *standard: VisDial train* | 916 | 708 (77.3) | 499 (54.5) | 2 | 0.053 | 0.084 | 0.291 | 0.360 |
| | *clean: VisDial val* | 141 | 104 (73.8) | 69 (48.9) | 2 | 0.126 | 0.081 | 0.248 | 0.375 |
| **Foil-It!** | *noun replacement* | 1000 | 943 (94.3) | 811 (81.1) | 73 | 0.426 | 0.425 | 0.532 | 0.588 |
| **Overall** | | 8782 | 7702 (87.7) | 5668 (73.6) | | | | | |

Table 5: Manual validation results for each piece in VALSE, as well as for the Foil-it dataset. *#Inst.*: number of instances for linguistic phenomenon. *#Valid (%)*: number (percent) of cases for which at least 2 out of 3 annotators chose the caption; *#Unan. (%)*: number (percent) of cases for which all annotators chose the caption; *#Lex.It.*: number of phrases or lexical items in the vocabulary that differ between foils and captions; *JS*: Jensen-Shannon divergence between foil-caption distributions for all instances in the whole instrument; *JS Val.*: Jensen-Shannon divergence between foil-caption distribution for the valid subset of the instrument, after sub-sampling; $\alpha$: Krippendorff's $\alpha$ coefficient computed over all the instances; $\alpha$ *valid*: Krippendorff's $\alpha$ coefficient computed over the *Valid* instances.

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 2. 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.*

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 5 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
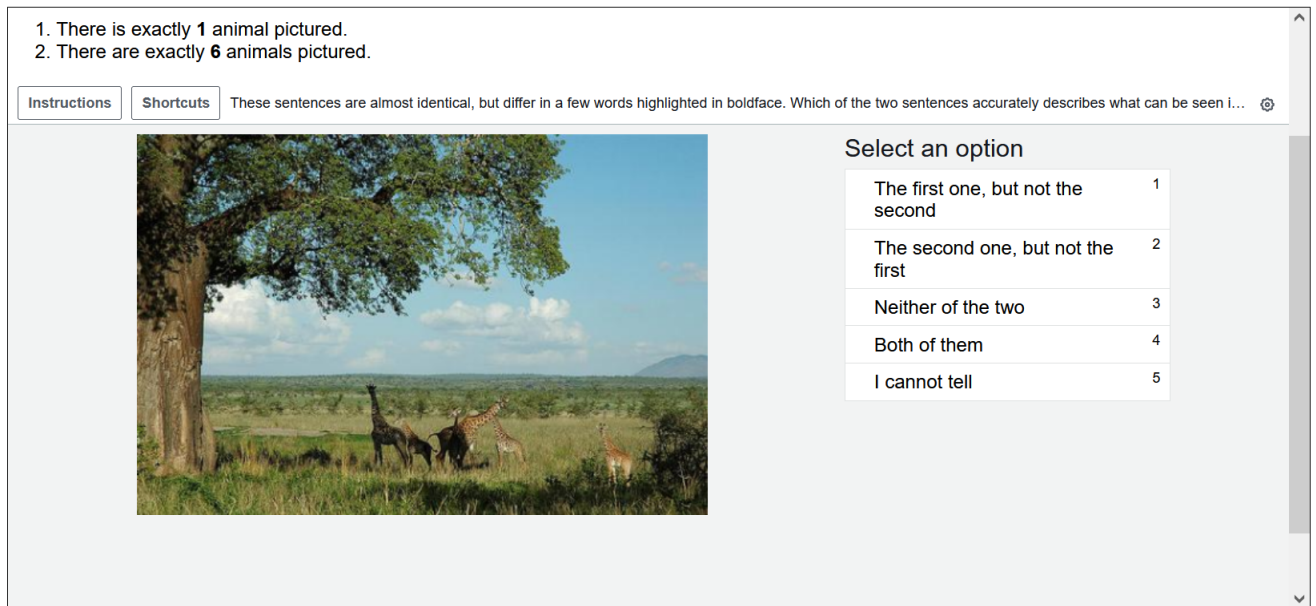
Figure 2: Example of an instance from the validation study. The example is from the Counting piece, *adversarial* instrument (see Section 3.3).

text which accurately describes the image. We also show the number of instances for which annotators unanimously (3/3) chose the caption.

**Annotator agreement**   As shown in Table 5, the proportion of *valid* instances in each instrument was high, ranging from 73.8% to 94.6%, with most instruments having annotators choose the caption well over 80% of the time. The table also shows two inter-annotator agreement statistics, both computed using Krippendorff's $\alpha$: over all the data in a given instrument, and over the valid subset only. On the valid subset, agreement is higher, and ranges from 0.3 to 0.6 (mean = 0.42; sd=0.12). There is a significant positive correlation between the percentage of valid instances per instrument and the $\alpha$ value (Spearman's $\rho = 0.75; p < .05$). The low to medium agreement suggested by the $\alpha$ range is due to two factors: first, the statistic is computed over the entire pool of annotators, of whom there were significant diversions in the amount of annotations they computed (e.g. some workers annotated fewer than 5 HITs); furthermore, the agreement is computed over 5 categories (see above). Given these factors, the inter-annotator agreement results should be treated with caution, and are not straightforwardly interpretable as an index of human performance on VALSE - in particular, the validation task (with 5 categories) was framed differently from the benchmark (which is binary).

**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 5. The foil-caption word frequency distributions can be inspected in Figures 3 and 4. 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 5 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.
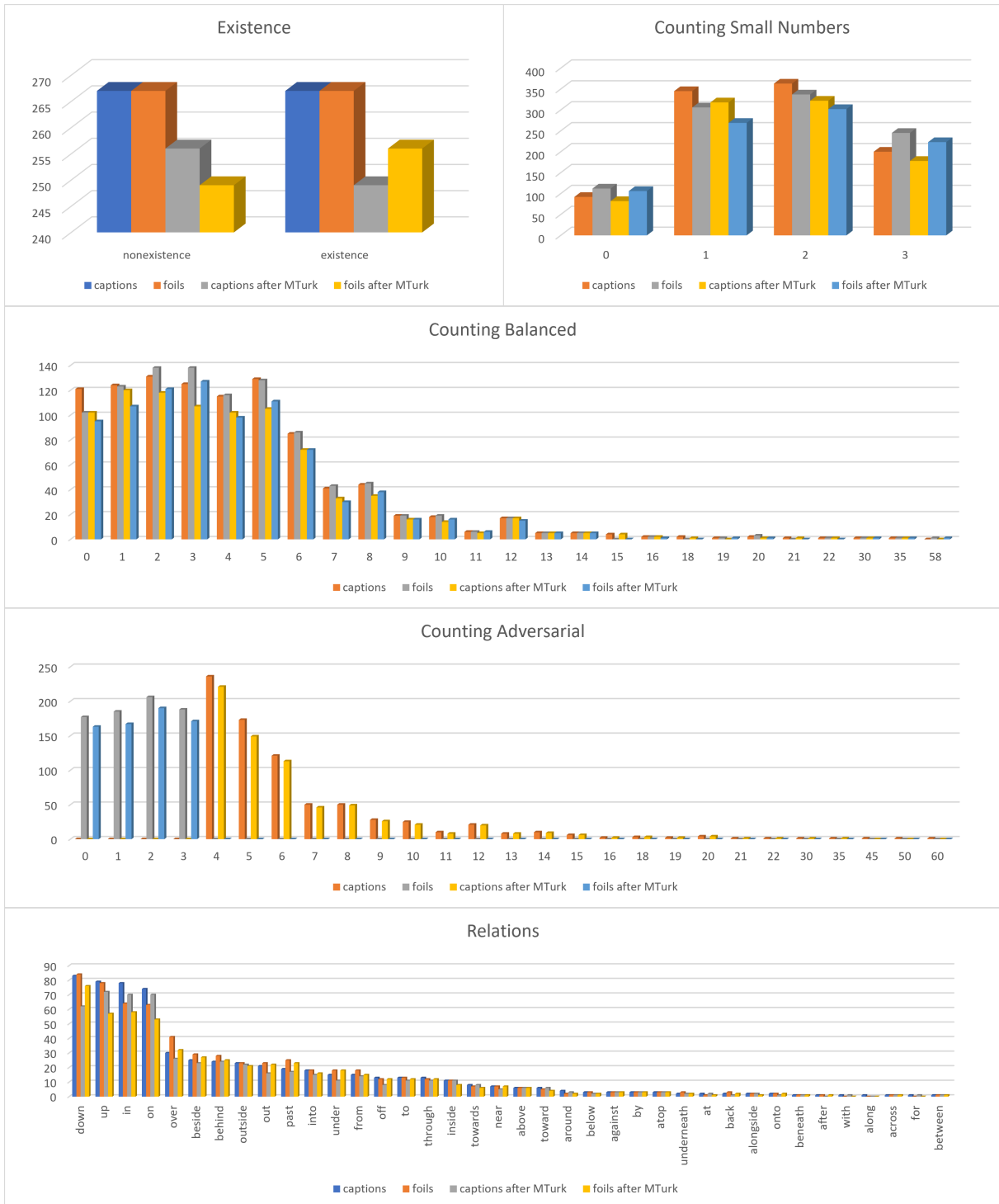
Figure 3: Word frequency distributions for captions and foils before and after the manual validation for existence, counting and relations.
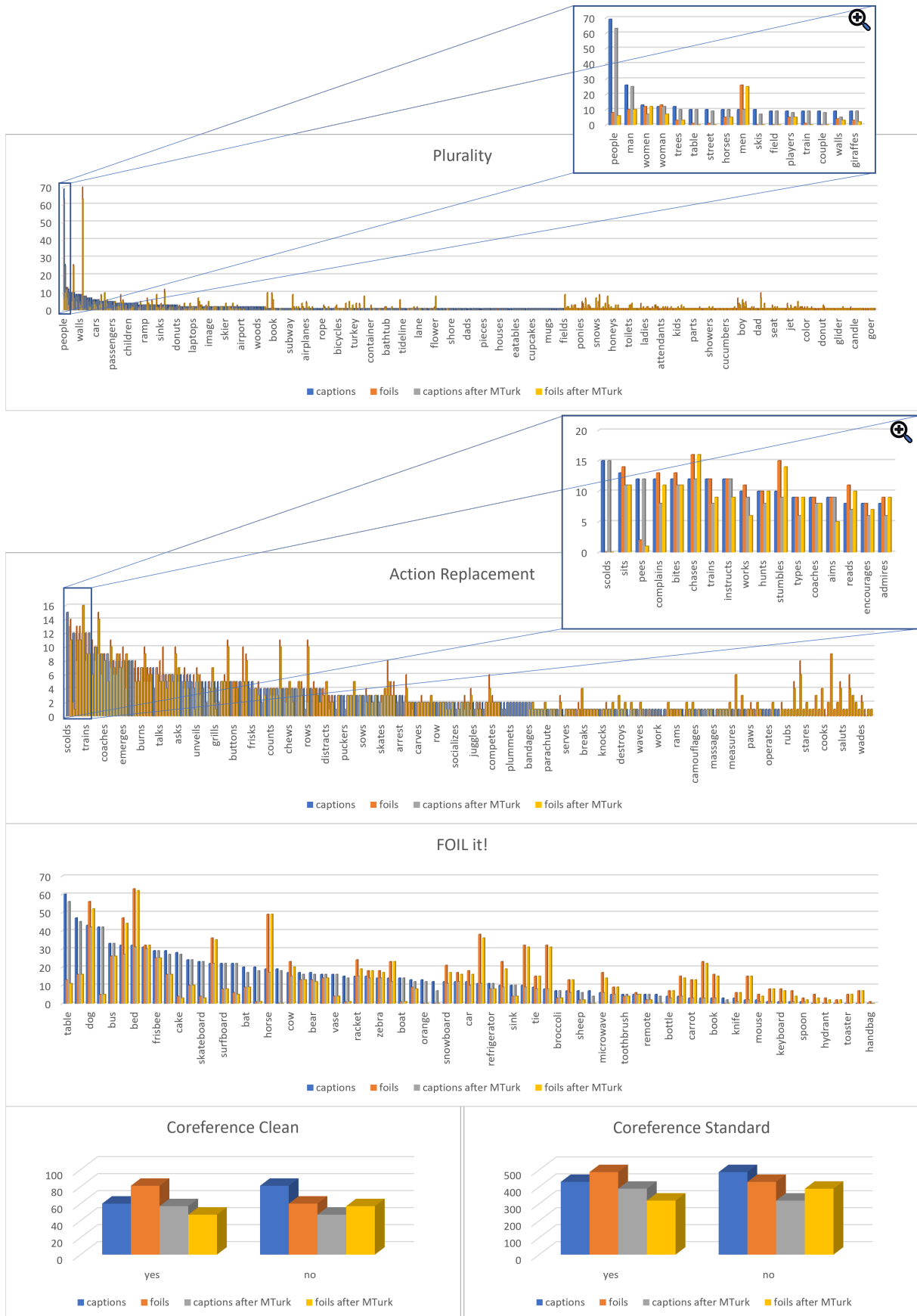
Figure 4: Word frequency distributions for captions and foils before and after the manual validation for plurality, action replacement and FOIL it. The actant swap instrument is not visualised here: By construction, actant swap cannot suffer from distributional bias since caption and foil contain the same words up to a *permutation*.

| piece | image | caption (blue) | foil (orange) |
|-------|-------|----------------|---------------|
| existence |  | There are no people in the picture. | There are people in the picture. |
| |  | There is a truck pictured. | There is no truck pictured. |
| |  | There are no clouds in the sky. | There are clouds in the sky. |
| |  | There are no people riding on elephants. | There are people riding on elephants. |
| |  | There is a kite. | There is no kite. |

Table 6: Randomly selected data examples for existence.

| piece | image | caption (blue) | foil (orange) |
|---|---|---|---|
| plurality |  | Two young men playing frisbee at night on exactly one sports field. | Two young men playing frisbee at night on a number of sports fields. |
| |  | Exactly one row of motorcycles parked together on a grass yard area with a house in the background. | A number of rows of motorcycles parked together on a grass yard area with a house in the background. |
| |  | Two men are looking inside of a single giant barbecue. | Two men are looking inside of a number of giant barbecues. |
| |  | Some children are playing baseball outside in a field. | A single child is playing baseball outside in a field. |
| |  | A number of people riding some motorbikes on the road. | A single person riding some motorbikes on the road. |

Table 7: Randomly selected data examples for plurality.

| piece | image | caption (blue) | foil (orange) |
|---|---|---|---|
| counting |  | There are exactly 8 horses. | There are exactly 5 horses. |
| |  | There is exactly 1 person snowboarding. | There are exactly 4 people snowboarding. |
| |  | There are exactly 6 motorcycles in this photo altogether. | There are exactly 7 motorcycles in this photo altogether. |
| |  | There are exactly 2 banana stalks. | There are exactly 4 banana stalks. |
| |  | There are exactly 12 roman numerals on the clock. | There are exactly 9 roman numerals on the clock. |

Table 8: Randomly selected data examples for counting.

| piece | image | caption (blue) | foil (orange) |
|---|---|---|---|
| relations |  | A baby elephant is walking under a larger elephant. | A baby elephant is walking on a larger elephant. |
| |  | Fruits and vegetables are being sold in a market. | Fruits and vegetables are being sold outside a market. |
| |  | An airplane is letting off white smoke against a blue sky. | An airplane is letting in white smoke against a blue sky. |
| |  | A cow stands on a sidewalk outside a building. | A cow stands on a sidewalk in a building. |
| |  | Three giraffes banding down to drink water with trees in the background. | Three giraffes banding up to drink water with trees in the background. |

Table 9: Randomly selected data examples for relations.

| piece | image | caption (blue) | foil (orange) |
|---|---|---|---|
| actions |  A figure climbs the stairs. | | A figure descends the stairs. |
| |  A woman skips a jump rope. | | A woman releases a jump rope. |
| |  An old man coaches people. | | An old man bothers people. |
| |  The people unveil the prize. | | A prize unveils people. |
| |  A baby drools over clothing. | | A clothing drools over the baby. |

Table 10: Randomly selected data examples for actions.

| piece | image | caption (blue) | foil (orange) |
|---|---|---|---|
| coreference |  | A close up of a hot dog with onions. Is it a big hot dog? Yes. | A close up of a hot dog with onions. Is it a big hot dog? No. |
| |  | A skateboarding man is on a half pipe. Does he wear a helmet? No. | A skateboarding man is on a half pipe. Does he wear a helmet? Yes. |
| |  | 2 women who have painted on mustaches petting a horse. Are they wearing hats? No. | 2 women who have painted on mustaches petting a horse. Are they wearing hats? Yes. |
| |  | Yellow sunflowers are in a blue and white giraffe styled vase. Is it inside? Yes. | Yellow sunflowers are in a blue and white giraffe styled vase. Is it inside? No. |
| |  | An adult giraffe and a child giraffe standing near a fence. Does this look like zoo? Yes. | An adult giraffe and a child giraffe standing near a fence. Does this look like zoo? No. |

Table 11: Randomly selected data examples for coreference.