

# Visual Commonsense in Pretrained Unimodal and Multimodal Models

Chenyu Zhang   Benjamin Van Durme   Zhuowan Li\*   Elias Stengel-Eskin\*

Johns Hopkins University

{czhan105, vandurme, zli110, elias}@jhu.edu

## Abstract

Our commonsense knowledge about objects includes their typical visual attributes; we know that bananas are typically yellow or green, and not purple. Text and image corpora, being subject to reporting bias, represent this world-knowledge to varying degrees of faithfulness. In this paper, we investigate to what degree unimodal (language-only) and multimodal (image and language) models capture a broad range of visually salient attributes. To that end, we create the Visual Commonsense Tests (ViComTe) dataset covering 5 property types (color, shape, material, size, and visual co-occurrence) for over 5000 subjects. We validate this dataset by showing that our grounded color data correlates much better than ungrounded text-only data with crowdsourced color judgments provided by Paik et al. (2021). We then use our dataset to evaluate pretrained unimodal models and multimodal models. Our results indicate that multimodal models better reconstruct attribute distributions, but are still subject to reporting bias. Moreover, increasing model size does not enhance performance, suggesting that the key to visual commonsense lies in the data.<sup>1</sup>

## 1 Introduction

The observation that human language understanding happens in a rich multimodal environment has led to an increased focus on visual grounding in natural language processing (NLP) (Baltrusaitis et al., 2019; Bisk et al., 2020), driving comparisons between traditional unimodal text-only models and multimodal models which take both text and image inputs. In this work, we explore to what extent unimodal and multimodal models are able to capture commonsense visual concepts across five types of relations: color, shape, material, size, and visual co-occurrence (cf. Fig. 1). We further explore how this

ability is influenced by reporting bias (Gordon and Van Durme, 2013), the tendency of large corpora to over- or under-report events. We define visual commonsense as knowledge about generic visual concepts, e.g. “knobs are usually round”, and we measure this knowledge via frequency distributions over potential properties (e.g. round, square, etc). A visually-informed language model should be able to capture such properties. Our color, shape, material, and co-occurrence data are mined from Visual Genome (Krishna et al., 2016), and our size data are created from object lists. They contain a large number of examples of per-object attribute distributions and “object-attribute” pairs.

Paik et al. (2021) evaluate language models’ color perception using a human-annotated color dataset (CoDa), finding that reporting bias negatively influences model performance and that multimodal training can mitigate those effects. In this work, we confirm those findings while extending the evaluation to a broader range of visually salient properties, resulting in a more comprehensive metric for visual commonsense. In order to elicit visual commonsense from language models, we utilize soft prompt tuning (Qin and Eisner, 2021), which trains optimal templates by gradient descent for each model and relation type that we explore. We also utilize knowledge distillation to enhance a text-only model’s visual commonsense ability, where the vision-language model serves as the teacher.

The major contributions of this work are: (1) we design a comprehensive analytic dataset, ViComTe, for probing English visual commonsense, that is applicable to any language model; (2) we use ViComTe to study models’ ability to capture empirical distributions of visually salient properties. We examine unimodal language models, multimodal vision-language (VL) models, and a knowledge-distilled version of a VL model; and (3) we analyze the effects of reporting bias on the visually-grounded vs. ungrounded datasets and models.

\*Joint Advising

<sup>1</sup>The dataset and code is available at <https://github.com/ChenyuHeidiZhang/VL-commonsense>.

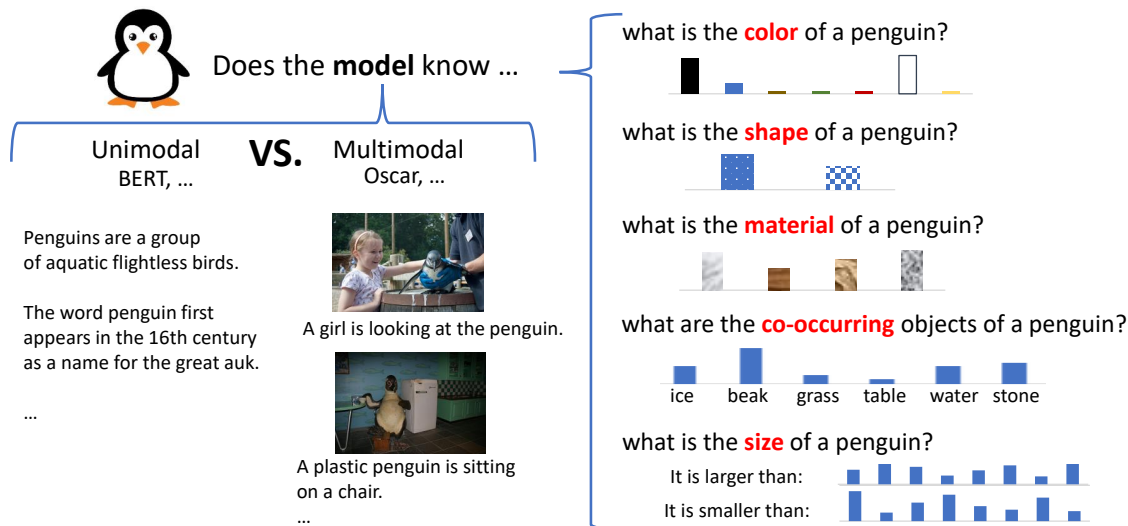


Figure 1: We compare unimodal and multimodal models’ abilities to capture visual commonsense knowledge. The commonsense knowledge is evaluated on five relation types: color, shape, material, size, and visual co-occurrence. We compare the model outputs with the gold distribution from ViComTe, which is mined from Visual Genome.

## 2 Related Work

### 2.1 Vision-Language Modeling

Recent advances in vision-language (VL) modeling have led to increased success on benchmark tasks. Most VL models learn joint image and text representations from cross-modal training of transformers with self-attention, including LXMERT (Tan and Bansal, 2019), ViLBERT (Lu et al., 2019), VisualBERT (Li et al., 2019), UNITER (Chen et al., 2020), etc. Oscar (Li et al., 2020) additionally uses object tags in images as anchor points to facilitate the learning of image-text alignments and VinVL (Zhang et al., 2021) presents an improved object detection model. CLIP (Radford et al., 2021) learns by predicting caption-image alignment from a large internet corpus of (image, text) pairs.

While our work uses textual prompt tuning techniques, there have also been work on visual prompt engineering to enhance the performance of pre-trained vision-language models. Zhou et al. (2021) model context in prompts as continuous representations and learn to optimize that context. Yao et al. (2021) develop a cross-modal prompt tuning framework that reformulates visual grounding as a fill-in-the-blank problem for both image and text.

### 2.2 Visual Commonsense

In one of the early attempts at learning visual commonsense, Vedantam et al. (2015) measure the plausibility of a commonsense assertion in the form of (obj1, relation, obj2) based on its similarity

to known plausible assertions, using both visual scenes and accompanying text. Zellers et al. (2021) learn physical commonsense via interaction, and use this knowledge to ground language. Frank et al. (2021) probe whether VL models have learned to construct cross-modal representations from both modalities via cross-modal input ablation.

Note that our definition of visual commonsense differs from that of Zellers et al. (2019), where the model is required to perform commonsense reasoning based on an image. Our definition of visual commonsense is more similar to the idea of stereotypic tacit assumptions (Prince, 1978) – the propositional beliefs that humans hold about generic concepts, such as “dogs have to be walked”. Weir et al. (2020) probe neural language models for such human tacit assumptions and demonstrate the models’ success. We extend this intuition to visual concepts and explore how visual information may help language models to capture such assumptions.

There has also been earlier work on the McRae feature norms (McRae et al., 2005), in which human annotators wrote down attributes that describe the meaning of words. For instance, “car” can be labeled as “has four wheels” and “apple” can be labeled as “is green”. Silberer et al. (2013) expand the McRae dataset into a set of images and their visual attributes and construct visually grounded distributional models that can represent image features with visual attributes.

Zhu et al. (2020) examine the “language prior” problem in Visual Question Answering models,

where models tend to answer based on word frequencies in the data, ignoring the image contents. In this work, we explore to what extent such a language prior is recruited absent a visual input.

### 2.3 Reporting Bias

Pretrained language models such as BERT (Devlin et al., 2019) are trained on billions of tokens of text, capturing statistical regularities present in the training corpora. However, their textual training data can suffer from reporting bias, where the frequency distribution of specific events and properties in text may not reflect the real-world distribution of such properties (Gordon and Van Durme, 2013). For example, while grass is typically green, this may be under-reported in web corpora (as it is assumed to be true), and while motorcycle crashes may be more common in the real world, plane crashes are mentioned far more in news text (Gordon and Van Durme, 2013). Misra et al. (2016) highlight the reporting bias in “human-centric” image annotations and find that the noise in annotations exhibits a structure that can be modeled.

## 3 Dataset: ViComTe

### 3.1 Dataset Mining

For each relation color, shape, material, size, and object co-occurrence, our data take the form of (subject, object) tuples extracted from object distributions per subject. The goal is to predict the object and its distribution from the subject and relation. Table 1 summarizes the number of classes and subject-object pairs for each relation.<sup>2</sup>

**Color, Shape, Material** For color, shape, and material, the subject is a noun and the object is the color, shape, or material property of the noun, mined from attributes of Visual Genome (VG) (Krishtna et al., 2016).<sup>3</sup> We manually create a list of single-word attributes for each relation, and only VG subjects that are matched with a specific attribute for more than a threshold number of times are recorded, in order to avoid noise in the dataset. The thresholds for color, material, and shape are 5, 2, and 1, respectively, chosen based on the availability of attributes of each relation in VG. VG attributes are filtered with the following steps: (1) attribute “Y colored / made / shaped” is treated as “Y”; (2) select only the last word for compound

attributes (e.g. treat “forest green” as “green”); (3) similar attributes are merged into a main attribute class (e.g. “maroon” and “crimson” become “red”).

The above procedure produces a distribution over the set of attributes for each subject noun. From that distribution, a (subject, object) data instance is generated for each subject where the object is the attribute that associates with it the most. See the first three rows of Table 1 for examples.

**Size** Size is separated into `size_smaller` and `size_larger`, where the subject is a noun and the object is another noun that is smaller or larger, respectively, than the subject. To form the size dataset, we obtain a set of concrete nouns that appears in VG, which we manually classify into 5 size categories (`tiny`, `small`, `medium`, `large`, and `huge`). Typical objects in each category includes *pill*, *book*, *table*, *lion*, *mountain*, respectively. We randomly pick two nouns from different categories to form a (subject, object) pair.

**Visual Co-occurrence** The visual co-occurrence dataset is generated in a similar way to the color, shape, and material datasets. Co-occurrence distribution is extracted from Visual Genome where two objects that occur in the same scene graph together for more than 8 times are recorded, and a (subject, object) instance is generated for each subject, where the object is the noun that co-occurs with the subject the most.

### 3.2 Data Grouping

Following Paik et al. (2021), we split the color, shape, and material datasets each into three groups: SINGLE, MULTI, and ANY. The SINGLE group is for subjects whose most common attribute covers more than 80% of the probability, e.g., the color of *snow* is almost always white. The MULTI group is defined as subjects not in the SINGLE group where more than 90% of the probability falls in the top 4 attribute classes, e.g., the color of a penguin in Fig. 1. The rest of the subjects are in the ANY group. Lower model performance for the SINGLE group would indicate the influence of reporting bias. For example, if the model is unable to correctly capture the distribution of the color of *snow*, it is likely because the color of snow has low probability of being reported in the training corpus, as people know it is white by default.

<sup>2</sup>See Appendix A.1 for more information on the object classes.

<sup>3</sup>Licensed under CC-BY 4.0.

Relation	# Classes	# (subj, obj) Pairs	Ex Template	Ex (subj, obj) Pair
color	12	2877	[subj] <i>can be of color</i> [obj]	( <i>sky, blue</i> )
shape	12	706	[subj] <i>has shape</i> [obj] .	( <i>egg, oval</i> )
material	18	1423	[subj] <i>is made of</i> [obj] .	( <i>sofa, cloth</i> )
size (smaller)	107	2000	[subj] <i>is smaller than</i> [obj] .	( <i>book, elephant</i> )
size (larger)	107	2000	[subj] <i>is larger than</i> [obj] .	( <i>face, spoon</i> )
co-occurrence	5939	2108	[subj] <i>co-occurs with</i> [obj] .	( <i>fence, horse</i> )

Table 1: Summary of the ViComTe dataset and the manual templates, including the number of classes, (subject, object) pairs, and an example pair for each relation.

Source	Group	Spearman $\rho$	# Subjs	Avg # Occ	Top5 # Occ	Btm5 # Occ	Acc@1
VG	All	64.3 $\pm$ 23.9	355	1252.6	64.6	308.6	
	SINGLE	62.2 $\pm$ 24.0	131	494.9	64.6	1181.6	80.2
	MULTI	69.3 $\pm$ 20.7	136	1156.1	2062.2	347.0	
	ANY	58.4 $\pm$ 27.1	88	2529.6	8452.4	1213.4	
Wikipedia	All	33.4 $\pm$ 30.6	302	543.6	1758.0	49.8	
	SINGLE	29.6 $\pm$ 29.9	110	352.2	345.8	35.0	35.5
	MULTI	33.9 $\pm$ 30.9	119	500.8	1242.0	27.6	
	ANY	38.2 $\pm$ 30.4	73	902.0	3000.2	161.2	

Table 2: Evaluation of ViComTe (mined from VG) and Wikipedia-mined color datasets by comparing with the human-annotated dataset CoDa. Reported are the average Spearman correlation ( $\times 100$ ), number of common subjects, average number of occurrences of the common subjects, average number of occurrences of subjects with top- and bottom-5 Spearman correlations, and the percentage of top-1 attributes being matched for the single group. ViComTe has higher correlations with human annotations.

### 3.3 Templates

In order to elicit model response and extract target objects and distributions from text, we manually design a set of templates for each relation. There are 7 templates for color, shape, and material each, 8 for size, and 4 for visual co-occurrence. See Table 1 for example templates.

### 3.4 Wikipedia Data

In order to compare text-based and visually-grounded data, we mine the color, shape, and material datasets from Wikipedia data, which is typically used in model pretraining. To mine these text-based datasets, we combine the sets of subjects in VG, take the manual list of attributes as objects again, and extract (subject, object) pairs if the pair matches any of the pre-defined templates. In Section 3.5 we will show the advantages of the VG-mined dataset over this text-based dataset.

### 3.5 Dataset Evaluation

To ensure the validity of ViComTe, we compare our color dataset with the human-annotated CoDa dataset (Paik et al., 2021), which we assume is close to real-world color distributions and has minimal reporting bias. We see a reasonably strong correlation with CoDa, indicating that the ViComTe dataset is a good and cost-effective approximation to human annotations.

**Metrics** We report the Spearman’s rank-order correlation between the two distributions in comparison, averaged across all subjects. The Spearman correlation is used instead of the Pearson correlation since for our purpose the rank of the object distributions is more important than the exact values, which may change due to data variability. The top-1 accuracy (Acc@1) is the percentage of the objects with the highest probability in the source distributions matching those in the target distributions. These two metrics are also used in later sections when evaluating model distributions.

**Analysis** Table 2 shows the detailed results of the evaluation of the ViComTe and Wikipedia color datasets by comparing with the human-annotated dataset, CoDa. We can see that ViComTe has much higher Spearman correlation with CoDa, as well as substantially higher top-1 accuracy for the SINGLE group. The correlation is expected to be low for the ANY group, because objects in the ANY group can have many possible colors.

Reporting bias is present in both datasets, as the average number of occurrences of SINGLE group subjects are much fewer than that of the MULTI and ANY group subjects. Counter-intuitively, for ViComTe, the highly-correlated SINGLE group subjects have fewer average occurrences than the ones with low correlations. This is contrary to our expectation that more frequent objects would better

reflect the human-perceived distribution and can be explained by SINGLE subjects being easier to represent even without a large amount of data.

One example where the Wikipedia distribution diverges from the CoDa distribution is “penguin”, whose most likely color in CoDa is black, followed by white and gray; however, its top color in Wikipedia is blue, because “blue penguin” is a specific species with an entry in Wikipedia, even if it is not as common as black and white penguins. One example where the VG distributions diverge from CoDa is “mouse”, because in VG, most occurrences of “mouse” are computer mice, which are most commonly black, whereas when asked about the word “mouse”, human annotators typically think about the animal, so that the most likely colors in CoDa are white and gray.<sup>4</sup>

### 3.6 Dataset splits

Each of the color, shape, material, size, and co-occurrence datasets is split into 80% training data and 20% test data. All evaluation metrics are reported on the test set. The training set is used for the logistic regression and the soft prompt tuning algorithm (Section 4.2).

## 4 Probing Visual Commonsense

### 4.1 Models

We examine 7 pretrained transformer-based models and 2 variations of them, trained on a variety of data. BERT (Devlin et al., 2019), ALBERT (Lan et al., 2020), and RoBERTa (Liu et al., 2019) are trained on text only using a masked language modeling objective (MLM). Oscar (Li et al., 2020) is a vision-language model based on the BERT architecture, trained with an combined MLM and contrastive loss on text-image pairs. VisualBERT (Li et al., 2019) is another vision-language model based on BERT that learns joint representation of images and text. Tan and Bansal (2020) introduce the “vokenization” method, which aligns language tokens to their related images, mitigating the shortcomings of models trained on visually-grounded datasets in text-only tasks. Since our task is purely text-based, we also experiment with a pretrained vokenization model (BERT + VLM on Wiki). Finally, we use representations from CLIP (ViT-B/32) (Radford et al., 2021), which is trained with a contrastive image-caption matching loss.

<sup>4</sup>Additional examples are provided in Appendix A.3.

**Distilled Oscar** As our experiments involve exclusively textual inputs, we develop a knowledge-distilled version of Oscar (“Distilled”) which corrects for the lack of image input in our task. Knowledge distillation (Hinton et al., 2015; Sanh et al., 2019) is the process of transferring knowledge from one model to another, where the student model is trained to produce the output of the teacher model. Here, we use Oscar as the teacher and BERT as the student. The training data is part of the Oscar pretraining corpus: COCO (Lin et al., 2014), Flickr30k (Young et al., 2014), and GQA (Hudson and Manning, 2019), and the Distilled Oscar model has access to the text data only. We use the Kullback-Leibler loss to measure the divergence between the output logits of BERT and Oscar, and optimize the pretrained BERT on that loss to match the outputs of Oscar. Configurable parameters are set the same as for Oscar pretraining.

**CaptionBERT** Since VL models are trained largely on caption data, it could be that the differences between a text-only model and a VL model come not from a difference in modalities – text vs. images and text – but from a difference in domain – webtext vs. image captions. In order to disentangle the effects of the domain difference from those of visual inputs, we train a BERT model from scratch (“CaptionBERT”) on Oscar’s caption-based text data (the same data as for the Distilled model). If CaptionBERT, which does not have exposure to visual inputs, performs better than BERT and similarly to VL models (which are trained with visual inputs), it would suggest that the training domain matters more than the modality. If, on the other hand, CaptionBERT performs worse than VL models, it would highlight the importance of modality.

### 4.2 Evaluation Methods

We compare the visual commonsense abilities of pretrained unimodal and multimodal models. Given a list of prompts and a subject word, each model outputs the distribution of the target word. Following Paik et al. (2021), we apply zero-shot probes to models that are trained on a language modeling objective, and conduct representation probes for those that are not. We report the prediction accuracy and the Spearman correlation of the output distribution with the true distribution.

We use models trained with an MLM objective (BERT, Distilled, etc) directly for zero-shot predic-

tion of masked tokens.<sup>5</sup> For Oscar we add a word-prediction head on top of it. The results across templates are aggregated in two modes. In the “best template” mode, for each example, the highest Spearman correlation among all templates is reported, and the top-1 result is regarded as correct if the true target object is the same as the top-1 result of any of the templates. In the “average template” mode, the output distribution is the mean of the distributions across all templates.

Since CLIP is not trained on a token-prediction objective, we implement logistic regression on top of the frozen encoder output, to predict the target attribute or object. The input is each of the templates with the subject [X] filled with an input in the dataset. Like Paik et al. (2021), to give the model ample chance of success, we take the template that results in the best test accuracy score, report that accuracy and the Spearman correlation associated with that template. For the classification head, we use the Scikit-Learn implementation of Logistic Regression (random\_state=0, C=0.316, max\_iter=2000) (Pedregosa et al., 2011).

**Soft prompt tuning** In order to overcome the limitation of self-designed prompts, we incorporate prompt tuning technique that learns soft prompts by gradient descent, from Qin and Eisner (2021).<sup>6</sup> The algorithm minimizes the log loss:

$$\sum_{(x,y) \in E_r} -\log \sum_{\mathbf{t} \in T_r} p(y|\mathbf{t}, x)$$

for a set of example pairs  $E_r$  and template set  $T_r$ .

### 4.3 Size Evaluation

The size dataset differs from the other datasets in that we use relative sizes (X is larger/smaller than Y), as absolute size information is hard to obtain. Thus, we use two evaluation strategies for size.

**Rank partition** First, as in the previous prediction task, given a template such as “[X] is larger than [Y]” and an object [X], we ask the model to predict the distribution of [Y], taking only the distribution  $D$  of nouns in the size dataset. For the current object [X], we take the nouns in size categories that are smaller than the category of [X] ( $N_{sm}$ ), and those that are in larger categories ( $N_{lg}$ ).

Let the length of  $N_{sm}$  be  $m$  and the length of  $N_{lg}$  be  $n$ . Then for the “larger” templates, we compute the average percentage of overlap between the top  $n$  objects in  $D$  and  $N_{lg}$  and that between the bottom  $m$  objects in  $D$  and  $N_{sm}$ . For the “smaller” templates, the “top” and “bottom” are reversed.

**Adjective projection** The second approach follows that of van Paridon et al. (2021), which projects the word to be evaluated onto an adjective scale. In this case, we compute the word embeddings of the adjectives “small” and “large” and the nouns from models, so the scale is  $\vec{\text{large}} - \vec{\text{small}}$  and the projection is calculated by cosine similarity. For instance, for the example noun “bear”, the projection score is given by:

$$\text{cos\_sim}(\vec{\text{large}} - \vec{\text{small}}, \vec{\text{bear}})$$

With good word embeddings, larger nouns are expected to have higher projection scores. The validity of the adjective scales from word representations is shown by Kim and de Marneffe (2013).

### 4.4 Measuring Model Reporting Bias

We measure the reporting bias of our models by comparing model performance on datasets with different levels of reporting bias and on the SINGLE, MULTI, ANY groups of the ViComTe dataset.

We assume that CoDa contains no reporting bias, in which case we can interpret Table 2 as showing that ViComTe contains a relatively small amount of it, and Wikipedia contains a relatively large amount. Thus, a larger correlation of model outputs with ViComTe and a smaller one with Wikipedia would indicate less model reporting bias.

Also, since the SINGLE group subjects are those whose attribute distribution concentrates on a single attribute, these subject-attribute pairs are less likely to be reported in text corpora or even image annotations. Therefore, lower model correlation on the SINGLE group than the MULTI and the ANY groups would be a sign of model reporting bias.

## 5 Results

The experimental results show that multimodal models outperform text-only models, suggesting their advantage in capturing visual commonsense. However, all models are subject to the influence of reporting bias, as they correlate better with the distributions from Wikipedia than those from CoDa

<sup>5</sup>For the target words that contain more than one subword tokens, we use the first token as the target.

<sup>6</sup><https://github.com/hiaoxui/soft-prompts>

Tune	Model	Color		Shape		Material		Cooccur
		Spearman $\rho$	Acc@1	Spearman $\rho$	Acc@1	Spearman $\rho$	Acc@1	Spearman $\rho$
No	BERT <sub>b</sub>	26.1 ± 31.0*	11.7	38.7 ± 15.1	6.7	33.7 ± 19.6	30.0	4.7 ± 3.5
	Oscar <sub>b</sub>	26.4 ± 30.7*	24.0	45.9 ± 14.1	<b>53.0</b>	38.6 ± 17.5	<b>43.3</b>	9.8 ± 6.9
	Distilled	34.8 ± 27.3	27.5	<b>46.2 ± 14.2</b>	37.3	36.1 ± 20.2	37.7	<b>10.1 ± 7.5</b>
	BERT <sub>l</sub>	<b>37.6 ± 30.3</b>	<b>30.3</b>	42.7 ± 17.1	28.4	36.6 ± 19.1	35.7	5.2 ± 3.8
	Oscar <sub>l</sub>	31.8 ± 28.3	17.1	40.0 ± 16.9	38.1	<b>39.2 ± 17.1</b>	40.5	9.7 ± 6.7
Yes	BERT <sub>b</sub>	48.0 ± 22.9	47.4	49.2 ± 12.7*	76.1	41.2 ± 15.3	45.2	11.3 ± 7.9
	Oscar <sub>b</sub>	<b>58.1 ± 21.1</b>	<b>67.9</b>	50.4 ± 11.5*	81.3	45.3 ± 14.3	<b>66.2</b>	12.7 ± 9.3
	Distilled	57.1 ± 21.9	64.6	<b>50.5 ± 12.3</b>	<b>82.8</b>	<b>45.4 ± 14.8</b>	<b>66.2</b>	<b>13.0 ± 10.1</b>
	BERT <sub>l</sub>	37.6 ± 30.3	30.3	49.2 ± 12.6	78.4	43.7 ± 15.1	53.3	11.4 ± 8.0
	Oscar <sub>l</sub>	57.6 ± 21.6	65.3	50.1 ± 12.2	81.3	45.2 ± 15.2	65.8	12.8 ± 9.6

Table 3: Spearman correlation and top-1 accuracy (both  $\times 100$ ) of zero shot probing, before and after soft prompt tuning (“N” and “Y” for the “Tune” column). This is the “average template” case where the output distribution is the mean of distributions across all templates. The Spearman correlation reported is the mean across all subjects  $\pm$  standard deviation, comparing the output distribution and the Visual Genome distribution. The subscripts *b* and *l* indicate the size of the model, and Distilled is the BERT model after distilling from Oscar. Asterisk indicates where there is *no* significant difference between BERT<sub>b</sub> and Oscar<sub>b</sub> (t-test p-value > 0.05).

and ViComTe. Prompt tuning and knowledge distillation substantially enhance model performance, while increasing model size does not.

### 5.1 Results with MLM Objective

**Color, Shape, Material** The resulting model performance for the “average template” mode is shown in Table 3. Prompt tuning is done in this mode only. Note that because the top-1 accuracy is taken among all possible classes of each relation, it should be interpreted together with the number of classes (Table 1).

We can see from Table 3 that Oscar does better than BERT in almost all cases. Significant difference between Oscar (base) and BERT (base) is seen in most cases. Also, after soft prompt tuning, both the Spearman correlation and the accuracy substantially improved. Although there is considerable variation of the Spearman correlations, we find consistent improvement per example with both prompt tuning and multimodal pretraining (Appendix A.2).

Table 3 also shows that knowledge distillation helps improve the performance of BERT in all cases, and the distilled model can sometimes even outperform the teacher model, Oscar. Moreover, the large version of each model does not always outperform its base counterpart, suggesting that increasing the size of the model does not enhance the model’s ability to understand visual commonsense. Instead, training with visually grounded data does.

Fig. 2 illustrates the Spearman correlations of different models with the color distributions from CoDa, ViComTe and Wikipedia, under the “best template” mode.<sup>7</sup> All models correlate moderately

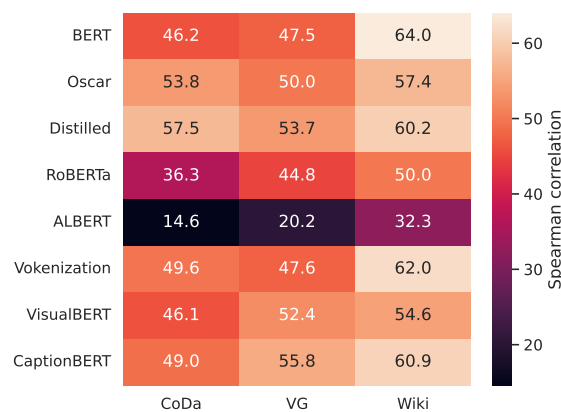


Figure 2: Spearman correlations ( $\times 100$ ) for color, under the “best template” case, for base models on CoDa, VG, and Wikipedia. While all models correlate the best with Wikipedia, BERT is the most biased.

with all three datasets, with the highest correlations to Wikipedia, indicating text-based reporting bias in all model types. BERT has the largest correlation gap between Wikipedia and CoDa, whereas the visually-grounded models have smaller gaps, indicating less reporting bias in VL models.

**Visual Co-occurrence** Table 3 also contains the results on visual co-occurrence before and after prompt tuning. Only the Spearman correlations are reported, because the top-1 accuracy is meaningless due to the large number of possible co-occurring objects with any noun.

Before prompt tuning, BERT has small Spearman correlations, suggesting that it may contain little knowledge about the visual co-occurrence relationship. Oscar demonstrates more such knowledge under the zero-shot setting. After prompt tuning, all model performances improve.

<sup>7</sup>Appendix A.2 contains further details.

Model	Color		Shape		Material		Co-occur
	Spearman $\rho$	Acc@1	Spearman $\rho$	Acc@1	Spearman $\rho$	Acc@1	Spearman $\rho$
BERT <sub>b</sub>	48.0 $\pm$ 21.6	51.4	53.2 $\pm$ 13.4	78.4	41.3 $\pm$ 15.6	51.1	30.2 $\pm$ 11.7
Oscar <sub>b</sub>	<b>52.5 <math>\pm</math> 20.8</b>	63.1	54.4 $\pm$ 14.8	<b>80.6</b>	<b>43.2 <math>\pm</math> 14.4</b>	<b>63.0</b>	31.2 $\pm$ 12.1
CLIP	51.9 $\pm$ 20.8	<b>63.8</b>	<b>54.5 <math>\pm</math> 13.9</b>	79.9	42.9 $\pm$ 15.0	<b>63.0</b>	<b>31.3 <math>\pm</math> 11.6</b>

Table 4: Spearman correlation and top-1 accuracy (both  $\times 100$ ) with a logistic regression head on model encoder outputs. Oscar and CLIP have comparable performance, both slightly better than BERT.

Group	Model	Color		Shape		Material	
		Spearman $\rho$	Acc@1	Spearman $\rho$	Acc@1	Spearman $\rho$	Acc@1
SINGLE	BERT <sub>b</sub>	36.8 $\pm$ 19.0	54.8	48.3 $\pm$ 12.3	83.0	35.9 $\pm$ 14.3	51.6
	Oscar <sub>b</sub>	39.9 $\pm$ 15.3	60.3	<b>49.3 <math>\pm</math> 11.6</b>	87.0	<b>38.5 <math>\pm</math> 12.8</b>	<b>65.1</b>
	CLIP	<b>41.0 <math>\pm</math> 15.2</b>	<b>66.3</b>	49.2 $\pm$ 14.5	<b>90.0</b>	38.1 $\pm$ 12.8	64.1
MULTI	BERT <sub>b</sub>	49.7 $\pm$ 21.2	42.3	<b>65.9 <math>\pm</math> 16.9</b>	59.5	53.8 $\pm$ 16.2	51.3
	Oscar <sub>b</sub>	<b>51.2 <math>\pm</math> 19.9</b>	50.6	65.2 $\pm$ 17.4	64.9	<b>56.2 <math>\pm</math> 13.0</b>	53.9
	CLIP	50.5 $\pm$ 21.1	<b>55.4</b>	64.6 $\pm$ 18.9	<b>67.6</b>	56.2 $\pm$ 14.3	<b>59.2</b>
ANY	BERT <sub>b</sub>	56.5 $\pm$ 19.5	46.1	100.0 $\pm$ 0	–	58.7 $\pm$ 15.2	<b>35.7</b>
	Oscar <sub>b</sub>	<b>62.5 <math>\pm</math> 18.9</b>	<b>58.4</b>	100.0 $\pm$ 0	–	60.4 $\pm$ 17.1	<b>35.7</b>
	CLIP	60.3 $\pm$ 18.2	55.8	100.0 $\pm$ 0	–	<b>63.5 <math>\pm</math> 20.5</b>	21.4

Table 5: Per-group Spearman correlation and top-1 accuracy (both  $\times 100$ ) with a logistic regression head on model encoder outputs. Note that the ANY group for shape only has one example, so the accuracy is less meaningful and is omitted. All models have higher correlations in the MULTI and ANY groups than the SINGLE group, which is a sign of reporting bias.

## 5.2 Results with Classification Head

Table 4 shows the results of BERT, CLIP, and Oscar when topped with a classification head. We observe that Oscar and CLIP achieve similar performance and both outperform BERT. Note that, while Visual Genome is part of Oscar’s pretraining corpus and one might suspect that that gives it an advantage, CLIP is trained on a large corpus from web search that is unrelated to Visual Genome. Therefore, we can conclude that multimodal models pretrained on both images and text outperform text-only models.

Table 5 breaks down the results in Table 4 into three subject groups. Oscar and CLIP outperform BERT in almost all cases. The top-1 accuracy is higher for the SINGLE group than for the MULTI and ANY groups, perhaps because the SINGLE group subjects have only one most likely target attribute, which may be easier to predict. Note that the Spearman correlations for all three models become higher from group SINGLE to MULTI to ANY. Paik et al. (2021) argue that higher correlation for the ANY and MULTI groups is a sign of model reporting bias, as objects in those two groups are more often reported. Thus, the results here indicate that reporting bias is still present in multimodal models.

## 5.3 Results: Size Relation

Table 6 shows results of the rank partition method (Section 4.3), before and after prompt tuning. Sur-

Tune	Model	Larger	Smaller
N	BERT <sub>b</sub>	80.0	67.1
	Oscar <sub>b</sub>	79.5	67.7
	Distilled	<b>84.6</b>	60.7
	BERT <sub>t</sub>	80.9	66.1
	Oscar <sub>t</sub>	79.4	<b>70.7</b>
Y	BERT <sub>b</sub>	69.9	55.7
	Oscar <sub>b</sub>	70.6	57.3
	Distilled	70.6	57.3
	BERT <sub>t</sub>	70.0	55.7
	Oscar <sub>t</sub>	70.6	57.3

Table 6: Percent correct for size relation, for “larger” and “smaller” templates, before and after soft prompt tuning. Interestingly, tuning does not help with size.

prisingly, prompt tuning does not help in this case. Moreover, the performance for the “larger” templates is higher than that of the “smaller” templates, suggesting that the models contain inherent preference towards the “larger” templates.

Fig. 3 shows the results of the adjective projection method.<sup>8</sup> For BERT and Oscar, we use the average embedding of the subword tokens of the nouns projected onto that of the adjectives “large” and “small”. For CLIP, we take the textual encoder outputs as the embeddings, resulting in a different score range from that of BERT and Oscar. The results show the following trend: larger objects are projected onto the “large” end of the spectrum, although the trend is sometimes broken towards the “huge” end. This may be due to the “huge” group

<sup>8</sup>Appendix A.2 contains per-object plot for BERT vs Oscar.



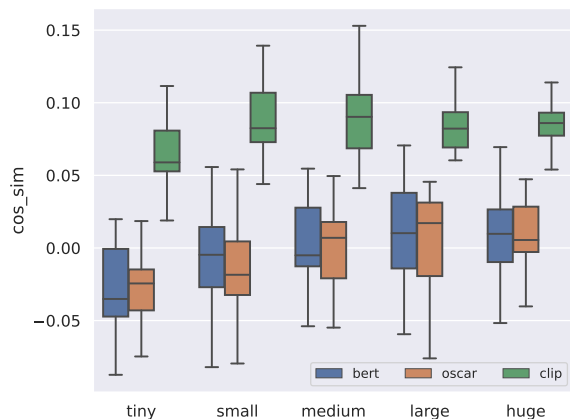


Figure 3: The size projection scores, where the x-axis indicates the object groups. Outliers are omitted. All three models perform reasonably well, as larger objects have higher cosine similarities in general.

including nouns such as “pool” and “house” which can be modified by a relative size indicator “small”.

#### 5.4 Analysis and Limitations

In Table 3, the accuracy of BERT for shape is particularly low (only 6.7%), despite that shape has only 12 classes. We hypothesize that this is due to reporting bias on shape in the text corpora that BERT is trained on. This hypothesis is supported by mining sentences from Wikipedia that contain (noun, attribute) pairs, where we see that the relation shape has fewer number of occurrences than material and color (Appendix A.3).

We also investigate whether the advantage of the visually-grounded models over pure-language models comes from the domain difference between web corpora and image captions, or the presence of actual visual input. Although its teacher is trained with visual inputs, the Distilled model is trained only on captions data and its performance matches that of Oscar, so we hypothesize that grounded training data enhance models’ ability to capture visual commonsense. The CaptionBERT results support the hypothesis in favor of domain difference, since it performs better than BERT in both CoDa and VG (Fig. 2). Nevertheless, the visual inputs also have an effect, as Oscar has a higher correlation than CaptionBERT on CoDa. Thus, it seems that both domain and modality affect the ultimate model performance.

Finally, although multimodal models show improvement on the task, sometimes the improvement is not significant and the resulting correlations are still weak. Further work is needed to enhance the visual commonsense abilities of the models and

mitigate reporting bias, and our datasets can serve as an evaluation method.

## 6 Conclusion

In this paper, we probe knowledge about visually salient properties from pretrained neural networks. We automatically extract dataset of five visual relations: color, shape, material, size, and co-occurrence, and show that our ViComTe dataset has a much higher correlation with human perception data for color than data mined from Wikipedia. We then apply several probing techniques and discover that visually-supervised models perform better than pure language models, which indicates that they can better capture such visual properties. Distilling the knowledge from a visually-supervised model into a pure language model results in comparable performance with the teacher model.

We also observe less reporting bias in both visually-grounded text (VG-mined datasets) than Wikipedia text and visually-grounded models (Oscar, DistilledOscar, VisualBERT, and CLIP) than pure language models. However, visually-grounded models are still subject to the influence of reporting bias, as seen in the per-group analysis, where both types of models perform better for the MULTI group than the SINGLE group.

## Acknowledgments

We would like to thank the reviewers for their comments and suggestions. Chenyu Zhang is supported by the Pistritto Research Fellowship. Elias Stengel-Eskin is supported by an NSF Graduate Research Fellowship. Zhuowan Li is supported by NSF 1763705.

## References

- Tadas Baltrusaitis, Chaitanya Ahuja, and Louis-Philippe Morency. 2019. [Multimodal machine learning: A survey and taxonomy](#). *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(2):423–443.
- Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, Nicolas Pinto, and Joseph Turian. 2020. [Experience grounds language](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8718–8735, Online. Association for Computational Linguistics.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. UNITER: Universal image-text representation learning. In *European Conference on Computer Vision (ECCV)*.
- 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.
- Stella Frank, Emanuele Bugliarello, and Desmond Elliott. 2021. Vision-and-language or vision-for-language? on cross-modal influence in multimodal transformers. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Jonathan Gordon and Benjamin Van Durme. 2013. Reporting bias and knowledge acquisition. In *Proceedings of the 2013 Workshop on Automated Knowledge Base Construction*.
- Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the knowledge in a neural network. *ArXiv*, abs/1503.02531.
- Drew A. Hudson and Christopher D. Manning. 2019. GQA: A new dataset for compositional question answering over real-world images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (ICCV)*, page 6700–6709.
- Joo-Kyung Kim and Marie-Catherine de Marneffe. 2013. [Deriving adjectival scales from continuous space word representations](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1625–1630, Seattle, Washington, USA. Association for Computational Linguistics.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. 2016. [Visual Genome: Connecting language and vision using crowdsourced dense image annotations](#).
- 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 (ICLR)*.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. In *Arxiv*.
- Xiujun Li, Xi Yin, Chunyuan Li, Xiaowei Hu, Pengchuan Zhang, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *European Conference on Computer Vision (ECCV)*.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft COCO: Common objects in context. In *European Conference on Computer Vision (ECCV)*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#).
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. [ViLBERT: Pretraining task-agnostic visual-linguistic representations for vision-and-language tasks](#). In *Annual Conference on Neural Information Processing Systems (NeurIPS)*.
- Ken McRae, George S. Cree, Mark S. Seidenberg, and Chris McNorgan. 2005. Semantic feature production norms for a large set of living and nonliving things. *Behavior Research Methods*, 37:547–559.
- Ishan Misra, C. Lawrence Zitnick, Margaret Mitchell, and Ross Girshick. 2016. [Seeing through the human reporting bias: Visual classifiers from noisy human-centric labels](#). In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2930–2939.
- Cory Paik, Stéphane Aroca-Ouellette, Alessandro Roncone, and Katharina Kann. 2021. [The world of an octopus: How reporting bias influences a language model’s perception of color](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 823–835, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jeroen van Paridon, Qiawen Liu, and Gary Lupyan. 2021. [How do blind people know that blue is cold? distributional semantics encode color-adjective associations](#). In *Proceedings of the Annual Meeting of the Cognitive Science Society*.

- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Ellen F. Prince. 1978. On the function of existential presupposition in discourse. In *Chicago Linguistic Society (Vol. 14, pp. 362–376)*.
- Guanghui Qin and Jason Eisner. 2021. Learning how to ask: Querying lms with mixtures of soft prompts. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, page 5203–5212, Online. Association for Computational Linguistics.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of bert: smaller, faster, cheaper and lighter. *ArXiv*, abs/1910.01108.
- Carina Silberer, Vittorio Ferrari, and Mirella Lapata. 2013. Models of semantic representation with visual attributes. In *ACL*.
- 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 (EMNLP)*.
- Hao Tan and Mohit Bansal. 2020. **Vokenization: Improving language understanding with contextualized, visual-grounded supervision**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2066–2080, Online. Association for Computational Linguistics.
- Ramakrishna Vedantam, Xiao Lin, Tanmay Batra, C. Lawrence Zitnick, and Devi Parikh. 2015. Learning common sense through visual abstraction. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 2542–2550.
- Nathaniel Weir, Adam Poliak, and Benjamin Van Durme. 2020. On the existence of tacit assumptions in contextualized language models. *Proceedings of the Annual Meeting of the Cognitive Science Society*.
- Yuan Yao, Ao Zhang, Zhengyan Zhang, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. 2021. CPT: colorful prompt tuning for pre-trained vision-language models. *ArXiv*, abs/2109.11797.
- Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. **From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions**. In *Transactions of the Association for Computational Linguistics*, volume 2, pages 67–78.
- Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. From recognition to cognition: Visual commonsense reasoning. In *2019 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Rowan Zellers, Ari Holtzman, Matthew Peters, Roozbeh Mottaghi, Aniruddha Kembhavi, Ali Farhadi, and Yejin Choi. 2021. **PIGLeT: Language grounding through neuro-symbolic interaction in a 3D world**. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2040–2050.
- Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. 2021. **VinVL: Making visual representations matter in vision-language models**. In *2021 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2021. Learning to prompt for vision-language models. *arXiv preprint arXiv:2109.01134*.
- Xi Zhu, Zhendong Mao, Chunxiao Liu, Peng Zhang, Bin Wang, and Yongdong Zhang. 2020. **Overcoming language priors with self-supervised learning for visual question answering**. In *International Joint Conference on Artificial Intelligence (IJCAI)*, page 1083–1089.

## A Appendix

### A.1 List of Objects

Table 7 shows the list of all possible attributes for relations color, shape, and material. Table 8 shows the list of objects in the five categories of relation size. Visual co-occurrence has a large number of objects that are not listed here for space reasons.

Relation	Classes
Color	<b>black, blue</b> (aqua, azure, cyan, indigo, navy), <b>brown</b> (khaki, tan), <b>gray</b> (grey), <b>green</b> (turquoise), <b>orange</b> (amber), <b>pink</b> (magenta), <b>purple</b> (lavender, violet), <b>red</b> (burgundy, crimson, maroon, scarlet), <b>silver, white</b> (beige), <b>yellow</b> (blond, gold, golden)
Shape	<b>cross, heart, octagon, oval, polygon</b> (heptagon, hexagon, pentagon), <b>rectangle, rhombus</b> (diamond), <b>round</b> (circle), <b>semicircle, square, star, triangle</b>
Material	<b>bronze</b> (copper), <b>ceramic, cloth, concrete, cotton, denim, glass, gold, iron, jade, leather, metal, paper, plastic, rubber, stone</b> (cobblestone, slate), <b>tin</b> (pewter), <b>wood</b> (wooden)

Table 7: List of all objects for relation color, shape, and material. Inside the parentheses are the attributes that are grouped into the object class.

Size	Objects
Tiny	ant, leaf, earring, candle, lip, ear, eye, nose, pebble, shrimp, pendant, spoon, dirt, pill, bee
Small	bird, tomato, pizza, purse, bowl, cup, mug, tape, plate, potato, bottle, faucet, pot, knob, dish, book, laptop, menu, flower, pillow, clock, teapot, lobster, duck, balloon, helmet, hand, face, lemon, microphone, foot, towel, shoe
Medium	human, door, dog, cat, window, lamp, chair, tire, tv, table, desk, sink, guitar, bicycle, umbrella, printer, scooter, pumpkin, monitor, bag, coat, vase, deer, horse, kite
Large	elephant, car, tree, suv, pillar, stairway, bed, minivan, fireplace, bus, boat, cheetah, wall, balcony, bear, lion
Huge	building, airplane, plane, clocktower, tower, earth, pool, mountain, sky, road, house, hotel, tank, town, city, dinosaur, whale, school

Table 8: List of objects in five size categories.

### A.2 Additional Probing

**Best template mode** Table 9 contains zero-shot results under the “best template” mode, for BERT (base), Oscar (base), BERT distilled from Oscar, RoBERTa (base), ALBERT (base), Vokenization, and VisualBERT (base). These results demonstrate

similar trends as the ones in the “average template” mode.

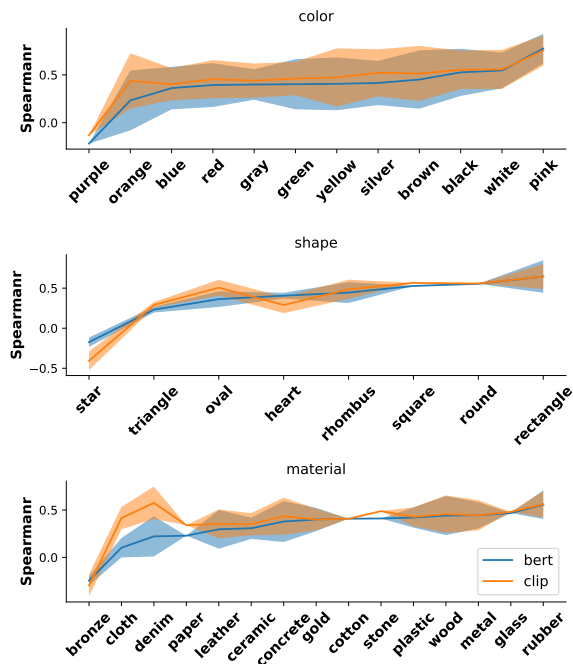


Figure 4: Spearman correlation per object class for BERT and CLIP with the logistic regression head, for color, shape, and material. The error margins are the standard deviations.

**Per-object analysis** Fig. 4 illustrates the fine-grained Spearman correlation  $\pm$  standard deviation per object group for BERT and CLIP.

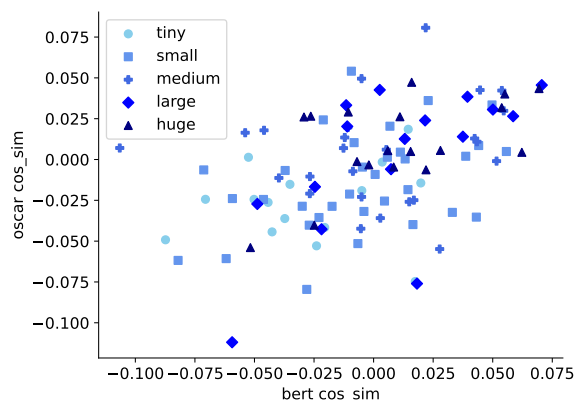


Figure 5: The size projection scores from BERT and Oscar, where each point is one object. Cosine similarities are correlated between Oscar and BERT.

**Size per-object** Fig. 5 shows how the per-object projection scores on the size spectrum from BERT and Oscar are correlated.

Model	Color		Shape		Material		Cooccur
	Spearman $\rho$	Acc@1	Spearman $\rho$	Acc@1	Spearman $\rho$	Acc@1	Spearman $\rho$
BERT <sub>b</sub>	47.5 ± 21.6	41.8	48.2 ± 12.0	64.3	41.9 ± 15.4	55.3	6.1 ± 4.0
Oscar <sub>b</sub>	50.0 ± 19.8	59.8	<b>52.7 ± 10.0</b>	89.3	<b>46.5 ± 13.7</b>	<b>74.6</b>	10.1 ± 7.2
Distilled	53.7 ± 21.3	57.7	51.4 ± 11.1	74.3	46.0 ± 13.6	<b>74.6</b>	10.4 ± 7.8
RoBERTa <sub>b</sub>	44.8 ± 19.8	41.6	45.4 ± 12.4	69.3	33.0 ± 15.5	39.1	1.1 ± 1.4
ALBERT <sub>b</sub>	20.2 ± 24.8	13.4	29.8 ± 15.7	13.6	25.0 ± 17.9	27.8	6.6 ± 5.1
Vokenization	47.6 ± 20.9	51.6	49.8 ± 13.1	72.9	39.4 ± 16.0	52.5	6.0 ± 3.7
VisualBERT	52.4 ± 19.8	65.3	48.7 ± 12.9	66.4	43.4 ± 15.5	59.5	<b>10.7 ± 8.1</b>
CaptionBERT	<b>55.8 ± 20.6</b>	<b>70.0</b>	51.3 ± 11.8	<b>91.4</b>	42.6 ± 15.4	54.6	10.2 ± 7.5

Table 9: Spearman correlation and top-1 accuracy (both  $\times 100$ ) of zero shot probing. This is the “best template” case discussed in Section 4.1.

**Per-Subject Comparison** Fig. 6 and Fig. 7 show how the Spearman correlations of 10 individual subjects improve after soft prompt tuning and after multimodal pretraining. Consistent improvement can be seen in color, material, and cooccurrence. Although we report average Spearman correlations in Table 3 and there are large standard deviations, here we show that when improvement is observed collectively, it is also consistent across subjects. With shape, the improvement is less obvious (45.9 to 50.4 for prompt tuning and 49.2 to 50.4 for multimodal pretraining).

### A.3 Error Analysis

**Data** The three subjects with the highest and lowest Spearman correlation are shown in Fig. 8 and Fig. 9.

**Wikipedia** Table 10 shows the number of (noun, attribute) pairs of the three relation types in Wikipedia. Shape has fewer occurrences than material and color.

	Color	Shape	Material
Total	331480	195921	307879
Avg 12	27623.3	16326.8	24634.7

Table 10: First row is the total number of occurrences of (noun, attribute) pairs for relations shape, material, and color in Wikipedia. Second row is the average number of occurrences across the top 12 attributes for each relation. Shape has the fewest number of occurrences.

**Model** Table 11 shows the errors made by BERT and Oscar in the “average template” mode before prompt tuning. Overall, subjects with low correlation are those that are less often reported in Visual Genome as well as in textual data.

### A.4 Resources

**BERT, RoBERTa, ALBERT** We use the Huggingface implementations of BERT, RoBERTa, and

ALBERT.

**Oscar** See the GitHub repository for the code and pretrained Oscar: <https://github.com/microsoft/Oscar>.

**CLIP** We use the CLIP model released by OpenAI: <https://github.com/openai/CLIP>.

**Vokenization** See the GitHub repository for the pretrained model: <https://github.com/airsplay/vokenization>.

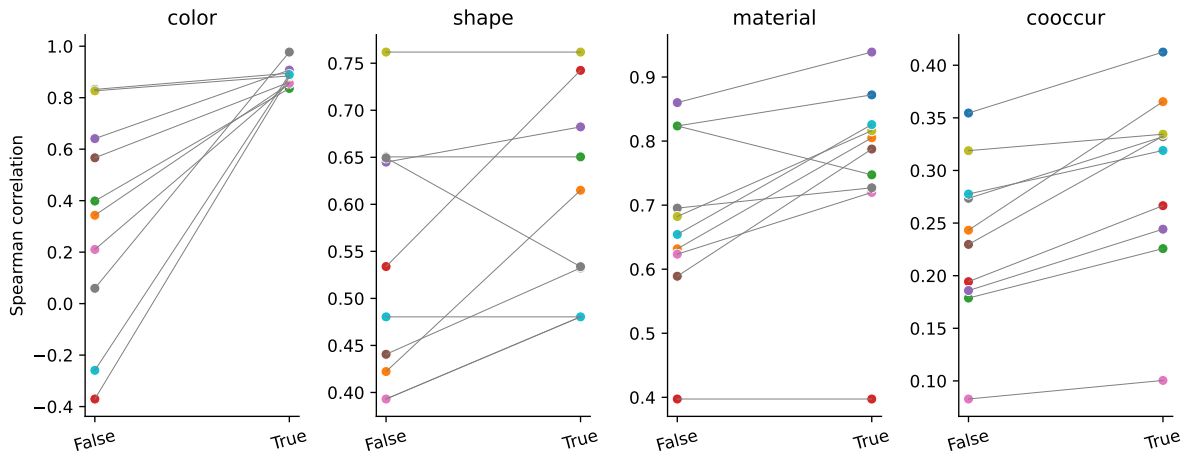


Figure 6: Spearman correlation of 10 subjects for each relation type before and after soft prompt tuning, with Oscar (base). Almost all individual subject has increased correlation after prompt tuning, except in relation shape.

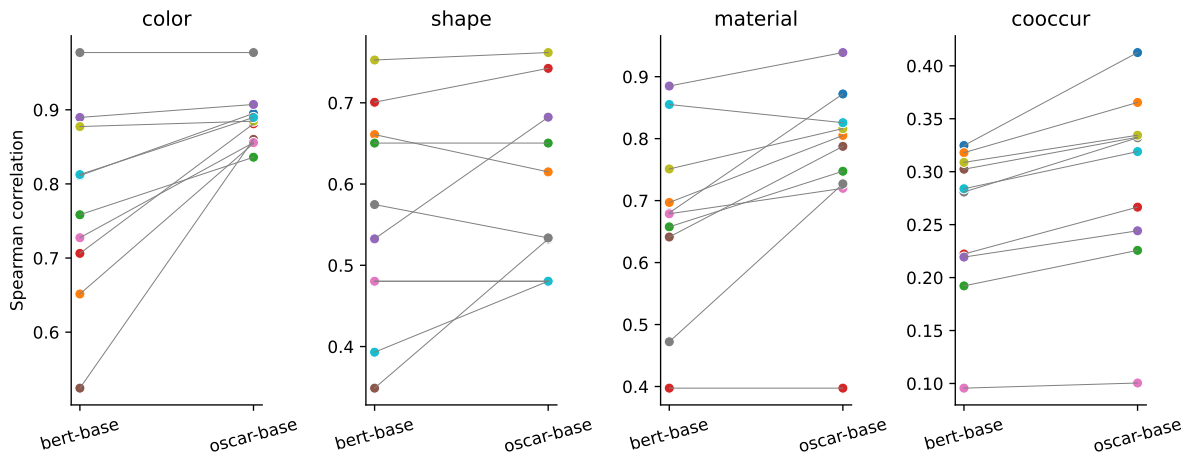


Figure 7: Spearman correlation of 10 subjects for each relation type with BERT (base) and Oscar (base), after soft prompt tuning. Almost all individual subject has higher correlation with Oscar than with BERT, except in relation shape.

Relation	High Corr Subjs		Low Corr Subjs	
	BERT <sub>b</sub>	Oscar <sub>b</sub>	BERT <sub>b</sub>	Oscar <sub>b</sub>
Color	lace, jacket, design	balloon, jacket, apple	flush, water faucet, muffler	hinge, leg, slack
Shape	mirror, vase, container	chair, pizza, vase	connector, log, knot	banana, toast, phone
Material	wall, tray, board	fence, wall, shelf	sheep, fabric, patch	elephant, rug, patch

Table 11: Three subjects each with high and low correlations for relations color, shape, and material.

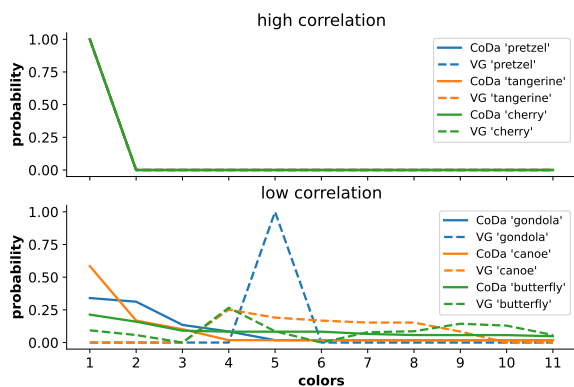


Figure 8: VG vs. CoDa distribution of 3 subjects with the lowest and highest correlation, ordered by probability of colors in CoDa.

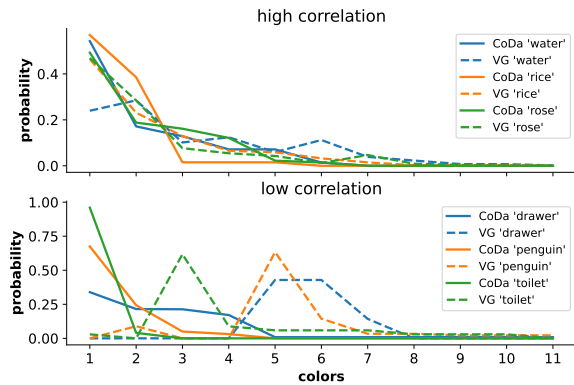


Figure 9: Wikipedia vs. CoDa distribution of 3 subjects with the lowest and highest correlation, ordered by probability of colors in CoDa.