

Explaining CLIP’s performance disparities on data from blind/low vision users

Daniela Massiceti[†]Camilla Longden[†]Martin Grayson[†]Agnieszka Słowik[†]Cecily Morrison[†]Samuel Wills[◇][†]Microsoft Research[◇]The World Bank

Abstract

Large multi-modal models (LMMs) hold the potential to usher in a new era of automated visual assistance for people who are blind or low vision (BLV). Yet, these models have not been systematically evaluated on data captured by BLV users. We address this by empirically assessing CLIP, a widely-used LMM likely to underpin many assistive technologies. Testing 25 CLIP variants in a zero-shot classification task, we find that their accuracy is 15 percentage points lower on average for images captured by BLV users than web-crawled images. This disparity stems from CLIP’s sensitivities to 1) image content (e.g. not recognizing disability objects as well as other objects); 2) image quality (e.g. not being robust to lighting variation); and 3) text content (e.g. not recognizing objects described by tactile adjectives as well as visual ones). We delve deeper with a textual analysis of three common pre-training datasets: LAION-400M, LAION-2B and DataComp-1B, showing that disability content is rarely mentioned. We then provide three examples that illustrate how the performance disparities extend to three downstream models underpinned by CLIP: OWL-ViT, CLIPSeg and DALL-E2. We find that few-shot learning with as few as 5 images can mitigate CLIP’s quality-of-service disparities for BLV users in some scenarios, which we discuss alongside a set of other possible mitigations.

1. Introduction

AI-based applications hold the potential to help people who are blind and low vision (BLV) with everyday visual tasks [3, 5]. However, the popularity of video-calling services like Be My Eyes [1] suggest that human assistance is still often required due to the wide set of assistance tasks [44] and varying quality of BLV images [8, 17]. Recent advances in large multi-modal models (LMMs) [19, 49, 52] could potentially address these challenges, enabling a new era of automated visual assistance as highlighted by the early partnership between Open AI and Be My Eyes [2].

Despite the opportunity, little work has evaluated how well LMMs perform on data from BLV users. Performance disparities have been identified for other user groups [6, 36,

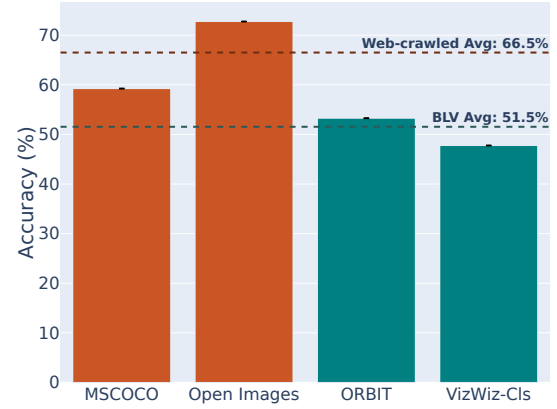


Figure 1. CLIP’s zero-shot object recognition accuracy is 15 percentage points lower in images from BLV users (ORBIT, VizWiz-Classification) versus web-crawled images (MSCOCO, Open Images). Average accuracy (with 95% c.i.) in a standardized zero-shot image classification task is reported over 80-100K images per dataset for 25 CLIP variants.

45, 52, 55, 66] but the evidence for BLV users is either anecdotal [49] or not specific to large multi-modal models [8]. Since BLV users are likely to be one of the biggest beneficiaries of LMMs, often in productivity- and safety-critical situations, it is important to extend studies to this group.

To address this, we systematically evaluate CLIP, a widely used LMM with 8700+ citations and 24M+ downloads¹, on data from BLV users. CLIP’s rich embeddings and strong zero-shot capabilities have led to it underpinning a wide range of downstream tasks including image classification [52], object detection [41, 42], semantic segmentation [37], image captioning [61, 63] and video recognition [35]. It has also been used to create large-scale datasets [23, 34, 57, 58] and evaluation metrics [26, 50]. As CLIP’s pre-trained parameters are often used directly, poor performance can have wide-ranging implications for downstream assistive applications that use them.

We investigate CLIP’s performance on BLV data along three dimensions: image content, image quality, and textual

¹Statistics taken from Google Scholar and OpenAI’s Hugging Face Hub (for CLIP ViT-L/14, ViT-B/32 and ViT-B/16) on 23 October 2023.

content. Visual content considers how well CLIP can recognize BLV-specific objects, such as guide canes. Visual quality assesses robustness to quality variations that characterize BLV images, such as blur and atypical framing [17]. Textual content examines performance on tactile descriptive words used by BLV users in contrast to visual ones, for example “plastic” versus “yellow”. We study each dimension in the context of a zero-shot image classification task, providing a worst-case estimate on how well CLIP will serve downstream assistive applications if used out-of-the box.

Overall, we find that CLIP’s zero-shot classification accuracy is 15 percentage points lower on BLV images compared to web-crawled images across 25 CLIP variants. These variants span architecture size (ViT-B/16 to ViT-g/14), pre-training dataset (WIT [52], LAION [57, 58], DataComp/CommonPool [23]) and pre-training dataset size (80M to 3.8B). On deeper inspection, underperformance stems from CLIP: 1) recognizing disability objects less well than non-disability ones, with 25 percentage points lower accuracy; 2) being sensitive to image quality, particularly occlusion and lighting issues; and 3) recognizing objects described by material less well than color, with discrepancies of 7 percentage points. In all cases, a larger pre-training dataset or architecture does not lead to parity.

To further understand our results, we examine the upstream source and downstream impact of these disparities. First, we conduct a textual analysis of the captions in LAION-400M/2B and DataComp1B and find that disability objects and materials are mentioned $\sim 17\times$ and $\sim 4\times$ less frequently than non-disability objects and colors, respectively. Second, we find performance disparities on BLV data persist in three downstream models that use CLIP: OWL-ViT [41] for object detection, CLIPSeg [37] for semantic segmentation, and DALL-E2 [53] for text-to-image generation. We close by discussing a set of possible mitigations, including few-shot model adaption and application-level solutions, toward making automated visual assistance for BLV users more equitable.

In summary, our work contributes to the literature on how LMMs perform for users in the margins, specifically highlighting how CLIP may underperform for BLV users if integrated into assistive applications. Our contributions are:

- An empirical study of CLIP’s performance on BLV image content, image quality and textual content.
- The first quantification of BLV content representation in LAION-400M, LAION-2B, and DataComp-1B.
- An example-based analysis that illustrates how performance disparities on BLV data persist in three downstream models that use CLIP.

2. Related Works

Large multi-modal models. LMMs now have impressive capabilities in analyzing and synthesizing images [7, 16, 19,

30, 49, 52, 68]. Contrastive models [30, 52, 68], a prominent sub-class, learn joint image and text embeddings by training on massive web-crawled data using a contrastive loss [15, 48]. They are unique in their architecture scale, and in the way they are trained on web-crawled data in an unsupervised manner. Unlike previous models, the rich embeddings they learn are leveraged by a wide range of downstream models – either directly [21, 52], or as part of a larger system [10, 35, 37, 41–43, 46, 53, 61–63, 67].

LMMs and fairness. LMMs are known to have social biases across gender, race, age, and geographic location [6, 36, 45, 66]. CLIP, for example, has been shown to classify people of color as non-human and males as criminal more often than white people and females, respectively [6]. Some works have studied these representational harm for people with disabilities, however only in natural language [28]. Quality-of-service harms arise when an application underperforms or fails for a particular user group [13, 18, 65] – *e.g.* a facial recognition system that does not detect women with darker skin tones [14]. These can be systematically identified and mitigated through disaggregated reporting of a model’s performance [9, 47]. This has not been well studied for people with disabilities generally or BLV people specifically, with the evidence either anecdotal (*e.g.* GPT-4Vision model card [49]) or not specific to LMMs [8].

3. Methodology

Our work investigates CLIP’s robustness to image and text data from BLV users in the context of a zero-shot image classification task. This provides a worst-case estimate of how CLIP will perform out-of-the-box in downstream assistive applications. Here we describe the experimental set-up, CLIP variants, and datasets used in our analyses.

3.1. Episodic zero-shot image classification

An image classifier selects which object $c \in \mathcal{C}$ is present in an image, where \mathcal{C} is the set of possible object classes and $|\mathcal{C}|$ is the task’s “way”. A zero-shot classifier does this without seeing any training images of the classes beforehand.

Our first analysis compares CLIP’s performance on different datasets (rather than the more typical multiple models on a single dataset), requiring our classification task set-up to be standardized across datasets. We take inspiration from the episodic sampling used in meta-learning [22]: for each dataset j annotated with \mathcal{C}_j object classes, we sample T fixed N -way classification tasks, where for each task we randomly sample N classes from \mathcal{C}_j . For each task, we randomly sample M test images per class. The classification accuracy is then computed for all $T \times M \times N$ images and the average (and 95% confidence interval) is reported. We repeat this for each dataset, with T , N and M held constant. We use variations of this to compare CLIP’s performance

between object types (Sec. 4.1) and text prompts (Sec. 4.3) with details provided in each section, respectively.

3.2. Logistic Regression

We also aim to understand which characteristics *within* images and text affect CLIP’s performance. We use logistic regression, a common tool for hypothesis testing, to estimate the marginal effect of each characteristic on the model’s accuracy. This approach avoids the need for careful experimental set-up which controls for all factors except the variable of interest. Logistic regression extends Ordinary Least Squares (OLS) regression to the case when the output variable is binary, as is our case where the model correctly identifies the ground-truth object or not. Formally, we use:

$$p(z_i) = \frac{1}{1 + e^{-z_i}} \quad (1)$$

where $z_i = \alpha_1 + \beta_1 X_i + \alpha_2 D_i + \beta_2 D_i X_i + \epsilon_i$. The output variable is $p(z_i) \in [0, 1]$, the probability that the model correctly identifies the ground-truth object in image i , with 1 for correct, and 0 otherwise. The explanatory variables are X_i , a vector of binary variables that encode whether a particular characteristic is present in image i , and D_i is a binary variable indicating whether the ground-truth object is a disability object (*e.g.* a guide cane). The interaction term $\beta_2 D_i X_i$ measures whether the marginal effect of each characteristic in X_i is compounded or mitigated for disability objects relative to non-disability objects. ϵ_i are residuals which are assumed homoskedastic and uncorrelated.

The coefficients $\alpha_1, \beta_1, \alpha_2, \beta_2$ are estimated through maximum likelihood. In OLS the coefficients directly represent the marginal effect of each X_i variable on the dependent variable. In contrast, here they represent the marginal effect on the log-odds ratio, which is linear in X_i :

$$\ln \left(\frac{p(z_i)}{1 - p(z_i)} \right) = \alpha_1 + \beta_1 X_i + \alpha_2 D_i + \beta_2 D_i X_i + \epsilon_i \quad (2)$$

This makes the coefficients difficult to interpret so we instead report them as $\partial p / \partial x$, the marginal effect of each characteristic $x \in X$ on the model’s probability of being correct, p . We report the average of this marginal effect across all observations in the sample. We interpret each effect through its sign, magnitude, and significance. A negative sign means the model is less likely to be correct when that characteristic is present in an image – on average and holding all other characteristics constant. Its magnitude measures the extent of this impact. Its significance indicates its reliability based on a two-sided t-test that estimates the probability that the marginal effect is different from zero.

3.3. CLIP variants

We study 25 CLIP variants spanning architecture size, pre-training dataset, and pre-training dataset size (see Tab. A.1

for summary). We focus on variants that use a Transformer [64] and Vision Transformer (ViT) [20] as the text and vision encoders respectively as they are most widely used. Specifically, we consider ViT-B/16, ViT-B/32, ViT-L/14, ViT-H/14 and ViT-g/14 vision encoders with associated text encoders. For datasets, we consider OpenAI’s closed-source WIT [52] and open-source LAION (80M/400M/2B) [30, 57, 58], DataComp (S/M/L/XL) [23], and CommonPool (S/M/L/XL) [23] with and without CLIP Score filtering [26]. These span 80M-3.8B image-text pairs.

We use CLIP as a zero-shot classifier by embedding a task’s class labels using its text encoder, and each task image with its vision encoder. An image’s prediction is taken to be the class whose embedding has the highest cosine similarity (after a softmax) with the image’s embedding.

3.4. Datasets

Our analyses are based on two large-scale datasets captured by BLV users: ORBIT [38] and VizWiz-Classification [8]. Both datasets were collected through real-world assistive applications: a personalizable object recognizer app for ORBIT [44]; and a visual question-answering app for VizWiz-Classification [12]. Both are therefore highly representative of typical BLV user data. We contrast these with two common web-crawled datasets – MS-COCO [33] and Open Images [32] – which are typical of the data used to pre-train LMMs, and widely used for benchmarking. We consider only the test and validation sets of these datasets. Below we provide descriptions of the BLV datasets, with the web-crawled datasets described in the appendix.

ORBIT [38] contains 3,822 videos (2.68M frames) of 486 objects collected by 67 BLV users on their mobile phones. For each object, users captured videos which show the object alone, and in a realistic scene alongside other items, which we call the Clean and Clutter datasets, respectively. ORBIT Clean frames are annotated with 6 quality issues (*e.g.* framing, blur) following the categories in [17].

VizWiz-Classification [8] contains 8,900 images from the original VizWiz dataset [25], a dataset of images taken by over 11,000 BLV users via a visual assistance mobile app [12]. All images are annotated with 200 ImageNet object categories and the 6 quality issues of [17] (including an extra “other” quality issue).

4. Experimental Results

Our first finding is that CLIP’s accuracy is 15.0 percentage points lower on BLV datasets (ORBIT and VizWiz-Classification) than web-crawled datasets (MS-COCO and Open Images) (see Fig. 1). We use the standardized zero-shot set-up (see Sec. 3.1) and average the $T \times N \times M$ predictions per dataset from each of the 25 CLIP variants. While the accuracy difference is less for larger CLIP architectures than smaller ones, no model achieves parity (see Fig. B.1).

Table 1. **CLIP underperforms on disability and exclusive disability objects by significant margins compared to non-disability objects.** Zero-shot accuracy is averaged (with 95% c.i.) over 27.5K images of each object type processed by each of the 25 CLIP variants. Experimental details in Sec. 4.1.1.

Object Category	ORBIT Clean	ORBIT Clutter
Excl. disability	36.5% \pm 0.1%	22.6% \pm 0.1%
Disability	41.8% \pm 0.1%	25.8% \pm 0.1%
Non-disability	58.9% \pm 0.1%	50.9% \pm 0.1%

In the best case, the gap is 6.7 percentage points (ViT-g/14, LAION-2B) while in the worst, it is 22.8 percentage points (ViT-B/32, DataComp-M). This preliminary result hints at deeper issues. In the following sections, we aim to identify potential sources of this discrepancy and why it occurs.

4.1. Robustness to image content from BLV users

To understand why accuracy is lower, we first examine BLV image content. The BLV community uses a range of assistive objects, like guide canes and Braille displays [31, 38, 44] (see Fig. 2), which are not included in popular benchmarks [33, 54, 56]. We assess CLIP’s performance on such “disability” objects versus more common objects.

We define disability objects as those that assist BLV people (*e.g.* dog collar); exclusive disability objects as the subset exclusively used by BLV people (*e.g.* guide cane); and non-disability objects as those used by everyone (*e.g.* keys). Three annotators categorized the ORBIT Clean and Clutter datasets² resulting in 55 disability, 42 exclusive disability, and 431 non-disability objects (see App. A.3 for lists).

4.1.1 Disability objects are less well recognized than non-disability objects

We compare zero-shot classification accuracy between disability and non-disability objects using a variant of the episodic set-up described in Sec. 3.1. Specifically, for each disability object we sample two N -way tasks with a “target” object and $N-1$ non-disability “distractor” objects. The first task contains a disability target object and the second task contains a non-disability target. The distractors are randomly sampled from the non-disability objects, each coming from a unique object cluster. We repeat T times for each disability object, sampling a pair of tasks with a different set of distractor objects and non-disability target object. For each task, we randomly sample M frames of the target object, and ask CLIP to classify them from the task’s N possible objects. We report the average accuracy of all frames with a disability and a non-disability object as the target, respectively ($T*55*M$ each). We also report the average accuracy over the subset of frames that are exclusive disability objects. We use $T = 5$, $N = 20$, $M = 100$.

²We do not consider VizWiz-Classification, as none of its 200 ImageNet labels are disability objects.



Figure 2. **Examples from the ORBIT Dataset.** (top) Disability objects: guide canes, liquid level sensor, electronic Braille device. (middle) Quality issues typical in BLV images: underexposure, blur, camera viewpoint, and framing. (bottom) A remote control and a Victor Reader Stream in a clean and clutter frame.

Under this setting, we find that disability and exclusive disability objects have accuracies of 21.1 and 25.3 percentage points less than non-disability objects, respectively, on average across the ORBIT Clean dataset (see Tab. 1). The gap widens by a further 3-4 percentage points when more realistic scenarios are presented from ORBIT Clutter. We find that the worst performing objects include Braille notetakers, talking book devices and liquid level indicators.

We also investigate the role of CLIP’s pre-training dataset size on this finding. We find that accuracy increases with pre-training dataset size generally, but the delta between non-disability and disability objects stays roughly constant (see Fig. B.2). This suggests that web-crawling more data may not be enough to improve performance on potentially long-tailed objects. We see similar trends for increasing architecture sizes (see Fig. B.3).

4.1.2 Disability objects are under-represented in large-scale datasets compared to non-disability objects

To better understand why more pre-training data does not improve performance on disability objects, we analyze the composition of three of CLIP’s large-scale pre-training datasets for the presence of disability content – LAION-400M [57], LAION-2B [58], and DataComp-XL [23] (also called DataComp-1B). These datasets are used for pre-training LMMs more broadly, with DataComp-XL achieving the highest accuracies on ORBIT.

Given the scale of the datasets, we conduct a text-based analysis of their captions as a more computationally tenable approach than analyzing their images. We first extract all noun phrases that contain a physical object³ from the

³A physical object traverses the entity \rightarrow physical-entity \rightarrow object \rightarrow OR(artifact, whole, part, living-thing) hypernym path in WordNet [40].

Table 2. **Disability objects occur 16-17x less frequently in the captions of popular large-scale image-text datasets compared to non-disability objects.** The mentions of 222 disability object synonyms and 312 non-disability synonyms were counted in noun phrases (NPs) extracted from these datasets. Details in Sec. 4.1.2.

	LAION-400M	LAION-2B	DataComp-1B
Captions	401,300,000	2,322,161,808	1,387,173,656
NPs	384,468,921	2,737,763,447	1,342,369,058
Unique NPs	5,984,181	22,657,632	15,071,341
Disability obj. mentions	18,326 (0.0048%)	70,939 (0.0026%)	48,672 (0.0036%)
Non-disability obj. mentions	425,046 (0.1106%)	1,550,043 (0.0566%)	1,126,356 (0.0839%)
Normalized non-dis/dis ratio	16.8	15.6	16.5

captions, referred to as “visual concepts”⁴. We then compute how prevalent ORBIT’s disability and non-disability objects are contained in these visual concepts. We use ORBIT to contextualize our previous results as it is a realistic representation of the types of objects important to BLV users, however, other object lists could be used.

To do this, we first group similar objects from the ORBIT dataset into higher-level clusters (*e.g.* all guide canes). As each cluster could be described in several ways (*e.g.* “symbol canes”, “guide canes”), we assign each two relevant synonyms. This was expanded to 15 synonyms for disability objects based on initial experimentation, resulting in 222 disability object synonyms, and 312 non-disability synonyms overall. We then count how many times each synonym appears within the visual concepts using string matching, allowing partial matches after simple pre-processing (see App. A.5 for details).

We find that disability objects occur 16-17x less frequently than non-disability objects across all three datasets (Tab. 2). We compute this by normalizing the number of mentions by the number of synonyms for disability and non-disability objects, respectively, and taking their ratio. We also see that LAION-2B has 7x the number of noun phrases as LAION-400M, but <4x the unique noun phrases, suggesting that it contains more of the same rather than new visual concepts (see App. A.5 for further statistics).

4.1.3 A few-shot approach can *sometimes* reduce the disability and non-disability accuracy gap

As CLIP is also known to be a good few-shot learner [60], we investigate whether providing several examples of an object can equalize performance between disability and non-disability content. We integrate a ProtoNets approach [59] with the “distractor” set-up described in Sec. 4.1.1, using embeddings directly from CLIP’s vision encoder⁵. Specifically for each disability object, we sample pairs of N -way tasks in the same way, except now we addi-

⁴We release these publicly at [REMOVED FOR REVIEW]

⁵We note that this few-shot set-up does not use CLIP’s text encoder.

Table 3. **A few-shot method using ProtoNets [59] (5-shot) achieves the highest accuracy and lowest accuracy gap between disability and non-disability objects, versus vanilla CLIP (0-shot) and CLIP with LLM-generated object descriptions [39, 51].** Averaged over 25 CLIP variants.

Obj type	ORBIT Clean Acc (%)				ORBIT Clutter Acc (%)			
	0-shot	[39]	[51]	5-shot	0-shot	[39]	[51]	5-shot
Disability	41.8	48.3	50.1	86.2	25.8	32.1	34.2	54.5
Non-disability	58.9	57.0	57.0	88.3	50.9	50.2	49.4	69.1
Accuracy gap	17.1	8.7	6.9	2.1	25.1	18.1	15.2	14.6

tionally sample K training shots of each class which we use to compute the class prototypes. As before, we evaluate the model on M test images for the disability and non-disability target object in each task pair, with the prediction taken to be the closest prototype. We consider $K = [5, 10, 20, 40]$. We compare this to recent methods [39, 51] which improve CLIP’s zero-shot performance by embedding LLM-generated descriptions of objects (rather than just the raw labels). We use GPT-4 as the LLM and the same generation hyperparameters as [39, 51].

We find that augmenting CLIP with LLM-generated object descriptions [39, 51] outperforms vanilla CLIP (0-shot) which just embeds the raw object labels, but not a few-shot approach (5-shot) which embeds a few image examples of each object (see Tab. 3). This holds for both the ORBIT Clean and Clutter datasets. Crucially, the accuracy gap between disability and non-disability objects is lowest with a few-shot approach, though this accuracy gap quickly saturates, with no significant gains coming from more than 5 shots (see Fig. B.4). We also note that while a few-shot approach can reduce the accuracy gap to 2% in the simple images from ORBIT Clean, it is less effective in the more realistic images from ORBIT Clutter, with disability objects performing 14-15% points worse than non-disability objects, even when scaled to 40 shots (see Fig. B.4b).

Furthermore, a few-shot approach is only effective as a mitigation if CLIP is pre-trained on a large enough dataset. We find that for pre-training datasets of less than 100M examples, the accuracy difference is 3-4x larger than that for 100-1000M examples, and 9-10x larger than that for 1B+ examples (see Figs. B.5a and B.5b). These factors are roughly constant across the number of shots. Overall, this speaks to the power of large-scale pre-training, even if a small amount of extra effort is required.

4.2. Robustness to image quality from BLV users

Images captured by BLV users are of more variable quality than those captured by sighted users. These issues include atypical framing, camera blur, camera viewpoint (rotation), occlusion, overexposure, and underexposure [17, 31], which are annotated in the ORBIT Clean and VizWiz Classification datasets. We run the standardized zero-shot set-up (see Sec. 3.1) on these datasets for all CLIP variants. We then use the statistical tools described in Sec. 3.2 to

disentangle the marginal effect of each quality issue on model performance, both in general and for disability objects specifically. For ORBIT, we treat X_i as a binary vector indicating the presence of five quality issues⁶ in image i , D_i as a binary indicating the presence of an exclusive disability object, and $D_i X_i$ as the interactions between them. For VizWiz, we encode seven quality issues (including the “other” category) in X_i , but exclude D_i or $D_i X_i$ as VizWiz labels do not include disability objects.

4.2.1 Blur, viewpoint, occlusion and lighting issues significantly reduce model accuracy.

In Fig. 3, we show that the marginal effects of blur, viewpoint (rotation), occlusion, and lighting issues on model accuracy are negative, large, and statistically significant for most models. All else equal, blur reduces model accuracy by 11 percentage points and 1 percentage point in the ORBIT and VizWiz datasets, respectively, on average. Viewpoint issues by 9 and 8 percentage points on each dataset respectively; occlusion by 9 and 14 percentage points; and lighting issues by 23 and 8 percentage points. We note that these effects are cumulative meaning that the impact on model accuracy is summed if multiple issues occur in the same image. We also note that pre-training on larger datasets, in general, does not guarantee robustness (*e.g.* variants pre-trained on LAION-2B, one of the largest datasets, are negatively affected by viewpoint and occlusion issues by 3-12 and 8-19 percentage points, respectively). We include the raw marginal effects in Tabs. B.3 to B.5.

Framing issues in the ORBIT dataset stand as the exception, with the marginal effect being positive and statistically significant. This can be explained by how the ORBIT videos were collected. To orient the camera, BLV users were instructed to hold it close to the object initially, and then move away. So, the initial frames in the video tend to be at close range – an easier recognition task – but also have framing issues. This is supported by the VizWiz results where framing issues, which occur at further distances from the object, have a negative marginal effect on accuracy.

4.2.2 The impact of quality issues is typically not worse for disability compared to non-disability objects.

Fig. 3 further shows that accuracy is 29 percentage points lower for exclusive disability objects than non-disability objects in the ORBIT Clean dataset, on average across all models, supporting the findings in Sec. 4.1.1. The marginal effect of a quality issue, however, typically affects disability objects no worse than non-disability ones. This can be seen by comparing the net effect of a quality issue on each object type. Let the baseline be the accuracy for non-disability objects. The accuracy for a disability object with no quality

⁶We combine over- and underexposure into a joint “lighting issue” due to low incidence rates of each of these issues.

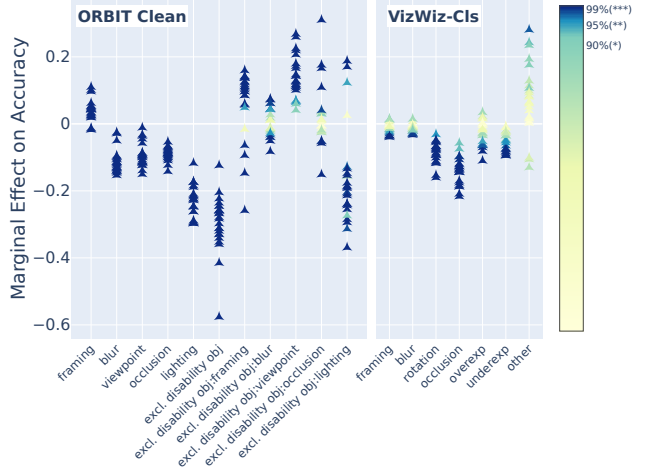


Figure 3. **Blur, viewpoint/rotation, occlusion and lighting issues all have large negative marginal effects on model accuracy, with high statistical significance, but these are not compounded for exclusive disability objects.** Each dot represents a CLIP variant, with its color showing the significance level.

issues will be 29 percentage points lower. Introducing occlusion will reduce the accuracy for non-disability objects by 9 percentage points on average. For disability objects, occlusion will reduce accuracy by this, plus the marginal effect of the interaction term (+2 percentage points), for a net effect of -7. The positive and significant interaction term indicates that having an occlusion issue and being a disability object has an effect that is slightly less than the sum of its parts. The only exception is overexposure issues, which do compound if they co-occur with a disability object.

4.3. Robustness to language used by BLV users

Assistive applications are likely to leverage the multi-modal capabilities of LMMs, so it is important to understand how CLIP performs on the range of language used by BLV people. For example, BLV users commonly use tactile rather than visual words to describe their objects [44]. In this section, we study one instantiation of this – CLIP’s robustness to recognizing objects described by their color, “yellow mug”, versus their material, “plastic cup”.

To do this, three annotators manually labeled the ORBIT validation and test objects (208 objects) with a color and a material⁷. Each adjective was selected from a pre-defined list of 20 colors and 23 materials (see App. A.4). A text prompt was then created for each object using the template “<adjective> <object_name>”, where <adjective> was the object’s color or material, and <object_name> was the noun extracted from the raw object label. We use these templates – referred to as color and material prompts – to examine CLIP’s sensitivity to different object descriptions.

⁷We assigned up to 2 adjectives per object in some cases where objects were multiple colors or materials.

Table 4. **Describing an object by its color (rather than material, or color and material) leads to text embeddings that are most aligned with that object’s image embeddings.** CLIP scores [26] between image and prompt embeddings are averaged (with 95% c.i.) for 100 images per object per prompt type on ORBIT Clean.

Prompt	Obj. name	Material + obj. name	Color + obj. name	Color + material + obj. name
CLIP Score	24.07 \pm 0.02	23.88 \pm 0.02	25.20 \pm 0.02	24.76 \pm 0.02

4.3.1 CLIP classifies objects more accurately when they are described by color rather than material

We compute CLIP scores [26] between an image and four different prompt embeddings, for 100 randomly sampled images of each object in ORBIT Clean. We consider the color and material prompts, a lower bound containing just the object name, and an upper bound adding both color and material adjectives. We expect that the lower bound prompt, which provides the least detail about the object, should align less strongly with the object’s image embedding than the upper bound prompt, which provides the most specific detail. In Tab. 4, however, we see this is not the case. Rather, color prompts have the highest CLIP scores and material prompts the lowest. Interestingly, the upper bound has a lower average CLIP score than the color prompt, suggesting that adding the object’s material is harming alignment.

To quantify the impact of this on accuracy, we run the standard zero-shot set-up (Sec. 3.1), embedding these textual prompts instead of the raw object labels. We see that across all variants, CLIP classifies objects 7.1 percentage points more accurately when they are described by their color rather than their material (see Fig. B.6).

4.3.2 Materials are under-represented in large-scale datasets compared to colors

We further examine this result by measuring how frequently colors versus materials appear in the captions of LAION-400M, LAION-2B and DataComp-1B. We use the extracted noun phrases from Sec. 4.1.2, and count the number of times the 20 material and 23 color annotations are mentioned. In Tab. 5, we see that colors are mentioned $\sim 4x$ more frequently than materials across both datasets, once normalized. This helps to explain some of the results

Table 5. **Materials occur $\sim 4x$ less frequently than colors in the captions of popular large-scale image-text datasets.** The mentions of 20 colors and 23 materials were counted in the noun phrases extracted in Tab. 2.

	LAION-400M	LAION-2B	DataComp-1B
Color mentions	475,060 (0.12%)	1,756,102 (0.06%)	1,165,871 (0.09%)
Material mentions	131,876 (0.03%)	513,014 (0.02%)	354,598 (0.03%)
Norm’d color/ material ratio	4.1	3.9	3.8

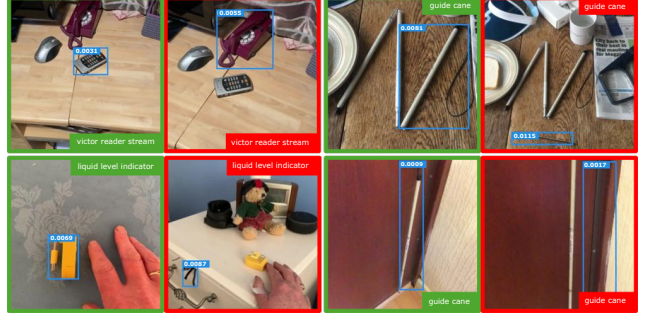


Figure 4. **OWL-ViT [41] detects disability objects less consistently than non-disability objects.** Disability objects are often mistaken for other objects, sometimes with higher confidence.

in Sec. 4.3.1. Taken together, this suggests that models pre-trained on these datasets may perform worse for BLV users who describe their objects by their material, with the potential that this may extend to other tactile-based descriptions.

5. Example-based impact analysis

Sec. 4 broadly shows that CLIP is sensitive to image and textual data provided by BLV users in a zero-shot classification task. We investigate whether these performance disparities persist in three downstream models that use CLIP – OWL-ViT [42], CLIPSeg [37], and DALL-E2 [53]. We run our analysis on 180 BLV images which are systematically selected for 20 objects – the 5 top- and bottom-performing disability and non-disability objects from the ORBIT dataset (see App. C for full protocol). For space reasons, we include CLIPSeg results in App. C.3.

5.1. Object detection with OWL-ViT

Object detection is already widely available in BLV assistive applications [3, 5], and in future, many may rely on models that use CLIP, such as OWL-ViT [41]. OWL-ViT predicts bounding boxes for objects specified in free-form text prompts. It does this by appending a bounding box regression and class-wise layer to CLIP’s (pre-trained) encoders and then fine-tuning on an object detection dataset. We run all 180 images through OWL-ViT (with a ViT-B/32 vision encoder) with the (cleaned) noun phrase extracted from the raw object label as the text prompt. A team of three annotators then manually evaluated the detections. We find:

Disability objects are less consistently detected than non-disability objects. Our results show that 6/10 non-disability objects were correctly detected (taken as the box with the highest confidence) in all 9 frames showing that object, compared to 3/10 disability objects. In many of these failed frames, the model mistook the disability object for another object, often with a higher confidence (Fig. 4). This behavior would have a large negative effect on the user experience of an object detection app.

Table 6. OWL-ViT [41]’s correct bounding box predictions have confidence scores that are $\sim 5\times$ lower for disability than non-disability objects on average. The confidence score of the predicted box per image is averaged (with 95% c.i.) over 90 images for disability and non-disability objects, respectively.

Object	Correct boxes	Incorrect boxes
Dis. objs	0.016 ± 0.008	0.008 ± 0.003
Non-dis. objs	0.084 ± 0.030	0.008 ± 0.003

The model is less confident about disability object detections than non-disability object detections. In Tab. 6, we see that OWL-ViT’s confidence for the correct bounding box is $\sim 5\times$ lower for disability objects compared to non-disability objects. We see that incorrect boxes have similar confidence scores between disability and non-disability objects, which is expected. See examples in Fig. C.1.

5.2. Text-to-image generation with DALL-E2

DALL-E2 [53] also uses CLIP: during training its decoder is conditioned on image embeddings from frozen CLIP. We investigate the downstream impacts of this by examining if DALL-E2 can generate disability content. We create two prompts for each of the 20 objects using the templates: i) “<object_name>” ii) “<object_name> on <surface> next to a <adjacent-object>”. The object name was the object label’s cleaned noun phrase, and the surface/adjacent object was chosen to match a randomly sampled clutter image of that object (see App. C for details). Three annotators then manually evaluated four generations from DALL-E2 per prompt. A generated image was considered correct if it contained the object specified in the prompt. We find:

Generations of disability objects are more likely to be incorrect compared to non-disability objects. DALL-E2 correctly generated the object in the prompt for 18/80 images of disability objects, versus 74/80 images of non-disability objects. For some disability objects, no generations contained a valid representation of the object – including guide canes, electronic Braille devices, and liquid level indicators (see Figs. 5 and C.4). In these cases, the generations either defaulted to a more common object (e.g. a walking stick for “guide cane”) or fabricated an object entirely (e.g. random dot patterns for “Braille sense display”, colorful thermometers for “liquid level sensor”). It also failed to generate specific instances of assistive devices (e.g. “Victor Reader Stream”, a talking book device, resulted in images of books or river streams). In contrast, DALL-E2 generates highly realistic of non-disability objects (see Fig. C.4a).

6. Discussion

Our evaluation of CLIP reveals that it consistently underperforms on BLV data across visual content, visual quality, and textual content, irrespective of architecture size, pre-training dataset, or pre-training dataset size. We discuss mitigation strategies to make LMMs more equitable

for BLV users and marginalized groups more generally.

Our results suggest that the performance disparities come in part from the distribution shift between web-crawled and BLV user data. This highlights the importance of systematic reporting of the contents of large-scale datasets used for pre-training, in the spirit of datasheets for datasets [24]. Our analysis in Secs. 4.1.2 and 4.3.2 provides a starting point, but this should be extended to other datasets and marginalized content. With the data composition known, mitigation strategies can then be developed. For example, assistive device websites and disability dataset platforms like IncluSet [4] could explicitly be crawled.

We also show that a few-shot approach can mitigate performance disparities relating to image content – a more cost-effective alternative than re-training a LMM. The few-shot model adaptation could be done when the LMM is developed, when the application is developed, or by the end-users themselves as part of a teachable paradigm [31, 38]. Each of these options is an open research question with the need to more deeply explore interaction paradigms and light-weight model adaptation techniques [11, 27].

Finally, application-level mitigations should also be considered. For BLV users, auxiliary models could support users to reduce image variance, helping them stabilize the camera or alerting about the lighting conditions, for example. We could also leverage data augmentation techniques that are personalized to individual users or user groups. For BLV users who tend to take blurry images, for example, we could automatically inject blur into the few-shot images so that the model becomes more robust to this quality issue.

The findings in this paper prompt a critical look at the development cycle of current LMMs. Greater transparency and disaggregation in dataset reporting is needed, regardless of the proprietary nature of a dataset. Future work should also explore lightweight model adaption techniques that allow application developers and users to bring equity to their experiences. We must continue to work with marginalized communities – “nothing about us without us” – to equalize the benefit of LMMs and their extraordinary capabilities.

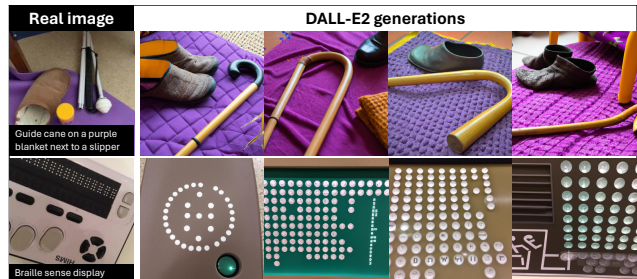


Figure 5. DALL-E2 [53] either defaults to common objects or fabrications when prompted with disability objects like guide canes and electronic Braille devices. Instead, it generates high-quality images of non-disability objects (see Fig. C.4a).

References

- [1] Be My Eyes. <https://www.bemyeyes.com>, . Accessed: 2023-10-11. **1**
- [2] Be My Eyes uses GPT-4 to transform visual accessibility. <https://openai.com/customer-stories/be-my-eyes>, . Accessed: 2023-10-11. **1**
- [3] Google Lookout. <https://play.google.com/store/apps/details?id=com.google.android.apps.accessibility.reveal>. Accessed: 2023-11-06. **1, 7**
- [4] IncluSet. <https://incluset.com/>. Accessed: 2023-11-09. **8**
- [5] Seeing AI. <https://www.microsoft.com/en-us/ai/seeing-ai>. Accessed: 2023-10-11. **1, 7**
- [6] Sandhini Agarwal, Gretchen Krueger, Jack Clark, Alec Radford, Jong Wook Kim, and Miles Brundage. Evaluating clip: towards characterization of broader capabilities and downstream implications. *arXiv preprint arXiv:2108.02818*, 2021. **1, 2**
- [7] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. In *Advances in Neural Information Processing Systems*, 2022. **2**
- [8] Reza Akbarian Bafghi and Danna Gurari. A New Dataset Based on Images Taken by Blind People for Testing the Robustness of Image Classification Models Trained for ImageNet Categories. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. **1, 2, 3**
- [9] Solon Barocas, Anhong Guo, Ece Kamar, Jacquelyn Kroner, Meredith Ringel Morris, Jennifer Wortman Vaughan, W Duncan Wadsworth, and Hanna Wallach. Designing disaggregated evaluations of ai systems: Choices, considerations, and tradeoffs. In *AAAI/ACM Conference on AI, Ethics, and Society*, 2021. **2**
- [10] Manuele Barraco, Marcella Cornia, Silvia Cascianelli, Lorenzo Baraldi, and Rita Cucchiara. The unreasonable effectiveness of CLIP features for image captioning: an experimental analysis. In *IEEE/CVF conference on computer vision and pattern recognition*, 2022. **2**
- [11] Samyadeep Basu, Daniela Massiceti, Shell Xu Hu, and Soheil Feizi. Strong Baselines for Parameter Efficient Few-Shot Fine-tuning. *arXiv preprint arXiv:2304.01917*, 2023. **8**
- [12] Jeffrey P Bigham, Chandrika Jayant, Hanjie Ji, Greg Little, Andrew Miller, Robert C Miller, Robin Miller, Aubrey Tatarowicz, Brandyn White, Samuel White, et al. VizWiz: Nearly real-time answers to visual questions. In *Annual ACM Symposium on User Interface Software and Technology*, 2010. **3**
- [13] Sarah Bird, Richard Edgar, Brandon Horn, Roman Lutz, Vanessa Milan, Mehmoosh Sameki, Hanna Wallach, and Kathleen Walker. Fairlearn: A toolkit for assessing and improving fairness in AI. Technical report, Microsoft, Tech. Rep. MSR-TR-2020-32, 2020. **2**
- [14] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR, 2018. **2**
- [15] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*. PMLR, 2020. **2**
- [16] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. PaLI: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*, 2022. **2**
- [17] Tai-Yin Chiu, Yinan Zhao, and Danna Gurari. Assessing image quality issues for real-world problems. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. **1, 2, 3, 5**
- [18] Sam Corbett-Davies, J Gaebler, Hamed Nilforoshan, Ravi Shroff, and Sharad Goel. The measure and mismeasure of fairness. *Journal Machine Learning Research*, 2023. **2**
- [19] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. InstructBLIP: Towards General-purpose Vision-Language Models with Instruction Tuning. *arXiv preprint arXiv:2305.06500*, 2023. **1, 2**
- [20] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representation*, 2020. **3**
- [21] Han Fang, Pengfei Xiong, Luhui Xu, and Yu Chen. CLIP2Video: Mastering video-text retrieval via image CLIP. *arXiv preprint arXiv:2106.11097*, 2021. **2**
- [22] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, 2017. **2**
- [23] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. DataComp: In search of the next generation of multimodal datasets. *arXiv preprint arXiv:2304.14108*, 2023. **1, 2, 3, 4, 13, 14, 16**
- [24] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. Datasheets for datasets. *Communications of the ACM*, 64(12):86–92, 2021. **8**
- [25] Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. VizWiz Grand Challenge: Answering visual questions from blind people. In *IEEE/CVF conference on computer vision and pattern recognition*, 2018. **3**
- [26] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021. **1, 3, 7, 20**
- [27] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. **8**

- [28] Ben Hutchinson, Vinodkumar Prabhakaran, Emily Denton, Kellie Webster, Yu Zhong, and Stephen Denuyl. Social biases in nlp models as barriers for persons with disabilities. In *Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 2020. 2
- [29] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Open-CLIP, 2021. 12, 13
- [30] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, 2021. 2, 3, 13
- [31] Hernisa Kacorri, Kris M Kitani, Jeffrey P Bigham, and Chieko Asakawa. People with visual impairment training personal object recognizers: Feasibility and challenges. In *CHI Conference on Human Factors in Computing Systems*, 2017. 4, 5, 8
- [32] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International Journal of Computer Vision*, 128(7), 2020. 3, 12
- [33] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*. Springer, 2014. 3, 4, 12
- [34] Zhiqiu Lin, Jia Shi, Deepak Pathak, and Deva Ramanan. The clear benchmark: Continual learning on real-world imagery. In *Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. 1
- [35] Ziyi Lin, Shijie Geng, Renrui Zhang, Peng Gao, Gerard de Melo, Xiaogang Wang, Jifeng Dai, Yu Qiao, and Hongsheng Li. Frozen clip models are efficient video learners. In *European Conference on Computer Vision*, 2022. 1, 2
- [36] Alexandra Sasha Luccioni, Christopher Akiki, Margaret Mitchell, and Yacine Jernite. Stable Bias: Analyzing societal representations in diffusion models. *arXiv preprint arXiv:2303.11408*, 2023. 1, 2
- [37] Timo Lüddecke and Alexander Ecker. Image segmentation using text and image prompts. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 1, 2, 7, 25
- [38] Daniela Massiceti, Luisa Zintgraf, John Bronskill, Lida Theodorou, Matthew Tobias Harris, Edward Cutrell, Cecily Morrison, Katja Hofmann, and Simone Stumpf. Orbit: A real-world few-shot dataset for teachable object recognition. In *IEEE/CVF International Conference on Computer Vision*, 2021. 3, 4, 8
- [39] Sachit Menon and Carl Vondrick. Visual classification via description from large language models. *arXiv preprint arXiv:2210.07183*, 2022. 5
- [40] George A Miller. WordNet: a lexical database for English. *Communications of the ACM*, 38(11):39–41, 1995. 4, 14
- [41] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection. In *European Conference on Computer Vision*, 2022. 1, 2, 7, 8
- [42] Matthias Minderer, Alexey Gritsenko, and Neil Houlsby. Scaling open-vocabulary object detection. *arXiv preprint arXiv:2306.09683*, 2023. 1, 7
- [43] Ron Mokady, Amir Hertz, and Amit H Bermano. CLIP-Cap: CLIP prefix for image captioning. *arXiv preprint arXiv:2111.09734*, 2021. 2
- [44] Cecily Morrison, Rita Marques, Martin Grayson, Daniela Massiceti, Camilla Longden, Linda Yilin Wen, and Edward Cutrell. Understanding Personalized Accessibility through Teachable AI: Designing and Evaluating Find My Things for People who are Blind or Low Vision. In *International ACM SIGACCESS Conference on Computers and Accessibility*, 2023. 1, 3, 4, 6
- [45] Ranjita Naik and Besmira Nushi. Social biases through the text-to-image generation lens. *arXiv preprint arXiv:2304.06034*, 2023. 1, 2
- [46] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 2
- [47] Besmira Nushi, Ece Kamar, and Eric Horvitz. Towards Accountable AI: Hybrid human-machine analyses for characterizing system failure. In *AAAI Conference on Human Computation and Crowdsourcing*, 2018. 2
- [48] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 2
- [49] OpenAI. GPT-4V(ision) System Card, 2023. 1, 2
- [50] Dong Huk Park, Samaneh Azadi, Xihui Liu, Trevor Darrell, and Anna Rohrbach. Benchmark for compositional text-to-image synthesis. In *Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. 1
- [51] Sarah Pratt, Ian Covert, Rosanne Liu, and Ali Farhadi. What does a platypus look like? generating customized prompts for zero-shot image classification. In *International Conference on Computer Vision*, 2023. 5
- [52] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*. PMLR, 2021. 1, 2, 3, 13
- [53] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with CLIP latents. *arXiv preprint arXiv:2204.06125*, 2022. 2, 7, 8
- [54] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. *arXiv preprint arXiv:2104.10972*, 2021. 4
- [55] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image syn-

- thesis with latent diffusion models. In *IEEE/CVF conference on computer vision and pattern recognition*, 2022. [1](#)
- [56] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3): 211–252, 2015. [4](#)
- [57] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. LAION-400M: Open dataset of CLIP-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. [1](#), [2](#), [3](#), [4](#), [13](#), [14](#), [16](#)
- [58] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. LAION-5B: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. [1](#), [2](#), [3](#), [4](#), [13](#), [14](#), [16](#)
- [59] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 2017. [5](#), [14](#), [19](#)
- [60] Haoyu Song, Li Dong, Wei-Nan Zhang, Ting Liu, and Furu Wei. CLIP models are few-shot learners: Empirical studies on VQA and visual entailment. *arXiv preprint arXiv:2203.07190*, 2022. [5](#)
- [61] Yixuan Su, Tian Lan, Yahui Liu, Fangyu Liu, Dani Yogatama, Yan Wang, Lingpeng Kong, and Nigel Collier. Language models can see: Plugging visual controls in text generation. *arXiv preprint arXiv:2205.02655*, 2022. [1](#), [2](#)
- [62] Mingkan Tang, Zhanyu Wang, Zhenhua Liu, Fengyun Rao, Dian Li, and Xiu Li. CLIP4caption: CLIP for video caption. In *ACM International Conference on Multimedia*, 2021.
- [63] Yoad Towel, Yoav Shalev, Idan Schwartz, and Lior Wolf. Zerocap: Zero-shot image-to-text generation for visual-semantic arithmetic. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. [1](#), [2](#)
- [64] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017. [3](#)
- [65] Mingyang Wan, Daochen Zha, Ninghao Liu, and Na Zou. Modeling techniques for machine learning fairness: A survey. *arXiv preprint arXiv:2111.03015*, 2021. [2](#)
- [66] Angelina Wang, Solon Barocas, Kristen Laird, and Hanna Wallach. Measuring representational harms in image captioning. In *ACM Conference on Fairness, Accountability, and Transparency*, 2022. [1](#), [2](#)
- [67] Zhaoqing Wang, Yu Lu, Qiang Li, Xunqiang Tao, Yandong Guo, Mingming Gong, and Tongliang Liu. CRIS: CLIP-driven referring image segmentation. In *IEEE/CVF conference on computer vision and pattern recognition*, 2022. [2](#)
- [68] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021. [2](#)

A. Extended experimental details

A.1. Datasets

Open Images. The Open Images V7 dataset [32] contains 61.4M images with image-level labels spanning 20.6K object classes. The images are web-crawled from Flickr and the classes include items of clothing, food types, animals, vehicles and more. We motivate the choice of this dataset because of its scale and diversity, and because it is widely used for training and benchmarking models within the computer vision community. We only sample images from the validation and test splits that have been verified by humans to contain the labeled object (i.e. all false positives are removed). This is 390,797 validation and 1,319,751 test images, respectively.

MS-COCO. The Microsoft COCO dataset [33] contains 328K images with instance labels spanning 80 object classes. The images are also web-crawled from Flickr and include “common” objects like people, animals, vehicles, furniture, food and more. We motivate this choice of dataset because, like Open Images, it is widely used for training and benchmarking models. We only sample images from the val2017 split (5K images) as the test split does not have ground-truth labels that are publicly available.

A.2. CLIP variants

We include all CLIP variants we study in Tab. A.1, including their pre-training dataset (and size) and model checkpoint. All checkpoints are taken from [open_clip](#) [29].

A.3. Disability and exclusive disability objects

Three annotators manually categorized the 486 ORBIT objects into 55 disability objects, 42 exclusive disability objects (a subset of disability objects) and 431 non-disability objects. Of these, 39, 30 and 310 were unique disability, exclusive disability and non-disability objects, respectively. We include the object lists for each category below:

Unique disability objects [39 objects]: folded cane, solo audiobook player, orbit braille reader and notetaker, victor stream book reader, cane, white cane, digital recorder, magnifier, long cane, pen friend, braille note, dog poo, symbol cane, pocket magnifying glass, glasses, folded long guide cane, insulin pen, dictaphone, white mobility cane, dog lead, retractable dog lead, braille orbit reader, victor reader stream, dogs lead, my hearing aid, water level sensor, braillepen slim braille keyboard, guide dog play cola, black mobility cane, my braille displat, visibility stick, leash, inhaler, liquid level indicator, hearing aid, guide dog harness, orbit reader 20 braille display, folded white cane, my cane

Unique exclusive disability objects [30 objects]: folded cane, solo audiobook player, orbit braille reader and notetaker, victor stream book reader, cane, white cane, digi-

tal recorder, magnifier, long cane, pen friend, braille note, dog poo, symbol cane, pocket magnifying glass, dictaphone, folded long guide cane, white mobility cane, braille orbit reader, victor reader stream, my hearing aid, water level sensor, braillepen slim braille keyboard, black mobility cane, my braille displat, visibility stick, liquid level indicator, hearing aid, orbit reader 20 braille display, folded white cane, my cane

Unique non-disability objects [310 objects]: cushion, tred mill, apple airpods, headphones, ipod stand, wallet for bus pass cards and money, handheld police scanner, shelf unit with things, av tambourine, tea, toothbrush, door, door keys, lotion bottle, pint glass, favourite earrings, proscocco, apple mobile phone, hat, tumble dryer, wall plug, risk watch, green water bottle, apple earpods, hole punch, phone stand, aspirin, tablets, garden shed, desk, knitting basket, dark glasses, headphone case, bin, chap stick, blue headphones, ottawa bus stop, fire stick remote, perfume, hair clip, pink himalayan salt, my purse, yellow marker, ipod in wallet, deodorant, mobile phone, iphone stand, apple phone charger, pencil case, one cup kettle, phone charger, adaptive dryer, skip prep, sunglasses case, eyewear case, apple headphones, front door, cranberry cream tea, backpack, key-chain, 13 measuring cup, microwave, apple wireless keyboard, my tilly hat, dog toy, speaker, water bottle, my airpods, garden table, ruler, journal, stairgate, sleep mask, coffee mug, radar key, lighter, trainer shoe, toaster, vape pen, banana, house keys, winter gloves, cannabis vape battery, my tilly hat upside down, cap, small space screwdriver, dab radio, watering can, wheely bin, litter and dog waste bin, my headphones, my muse s headband, airpods, set of keys, wireless earphones, iphone in case, pink marker, scissors, blue tooth keyboard, remote control, my wraparound sunglasses, finger nail clipper, vagabond ale bottle, face mask, screwdriver, sock, front door to house, my mug, single airpod, back patio gate, earphones, 14 measuring cup, sky q remote, tv unit, lip balm, reptile green marker, coin purse, post box, watch, t-shirts, bus stop sign, buckleys, ladies purse, iphone air pods, recycling bin, black bin, key, black small wallet, table fan, exercise bench, keyboard, hand gel, purse, vase with flowers, white cane, house door, wallet, reading glasses, orange skullcap, baked bean tin, migenta marker, my purple mask, condom box, mediterranean sea salt, my work backpack, personal mug, bottle opener, my slate, clickr, measuring spoon, rice, mug, iphone 6, presentation remote, secateurs, ps4 controller, remote tv, necklace, wardrobe, aspirin vs tylenol, mouse, small screwdriver, socks, eye drops, mustard, hand saw, lipstick, bose wireless headphones, hair brush, hairbrush, pinesol cleaner, memory stick, glasses case, knitting needle, pepper shaker, cup again, bone conducting headset, fridge, usb stick, compact disc, work phone, wine glass, my front door, work bag, headband, airpod pro, walletv, my laptop, money pouch,

Table A.1. All CLIP variants with their pre-training dataset, pre-training dataset size, and checkpoint (taken from [open_clip](#) [29]).

CLIP variant	Pre-training dataset	Dataset size	Checkpoint
ViT-B/16	WIT [52]	400M	openai
ViT-B/16	LAION-80M [30]	80M	Data-80M.Samples-34B.lr-1e-3.bs-88k
ViT-B/16	LAION-400M [57]	400M	laion400m.e32
ViT-B/16	LAION-2B [58]	2B	laion2b.s34b.b88k
ViT-B/16	DataComp-L [23]	140M	datacomp.l.s1b.b8k
ViT-B/16	CommonPool-L [23]	1.28B	commonpool.l.s1b.b8k
ViT-B/16	CommonPool-L (CLIP-Score filt.) [23]	384M	commonpool.l.clip.s1b.b8k
ViT-B/32	WIT [52]	400M	openai
ViT-B/32	LAION-80M [30]	80M	Data-80M.Samples-34B.lr-1e-3.bs-88k
ViT-B/32	LAION-400M [57]	400M	laion400m.e32
ViT-B/32	LAION-2B [58]	2B	laion2b.s34b.b79k
ViT-B/32	DataComp-S [23]	1.4M	datacomp.s.s13m.b4k
ViT-B/32	DataComp-M [23]	14M	datacomp.m.s128m.b4k
ViT-B/32	CommonPool-S [23]	12.8M	commonpool.s.s13m.b4k
ViT-B/32	CommonPool-S (CLIP-Score filt.) [23]	3.8M	commonpool.s.clip.s13m.b4k
ViT-B/32	CommonPool-M [23]	128M	commonpool.m.s128m.b4k
ViT-B/32	CommonPool-M (CLIP-Score filt.) [23]	38M	commonpool.m.clip.s128m.b4k
ViT-L/14	WIT [52]	400M	openai
ViT-L/14	LAION-80M [30]	80M	Data-80M.Samples-34B.lr-1e-3.bs-88k
ViT-L/14	LAION-400M [57]	400M	laion400m.e32
ViT-L/14	LAION-2B [58]	2B	laion2b.s32b.b82k
ViT-L/14	DataComp-XL/1B [23]	1.4B	datacomp.xl.s13b.b90k
ViT-L/14	CommonPool-XL (CLIP-Score filt.) [23]	3.8B	commonpool.xl.clip.s13b.b90k
ViT-H/14	LAION-2B [58]	2B	laion2b.s32b.b79k
ViT-g/14	LAION-2B [58]	2B	laion2b.s34b.b88k

remote, jd whisky bottle, paperclips, pex plumbers pliers, samsung tv remote control, my airpod pros case, portable keyboard, money clip, flat screen television, clear nail varnish, usb c dongle, amazon remote control, digital dab radio, 1 cup, measuring cup, tissue box, baseball cap, earpods, gloves, p939411 white cane, smarttv, skipping rope, back door, i d wallet, bluetooth keyboard, sunglasses, headset, my pill dosette, fridge freezer indicator, usb, apple pencil, black strappy vest, my apple watch, cell phone, apple wath, airpods pro charging case, slippers, dog streetball, corkscrew, airpod case, veg peeler, local post box, brown leather bracelet, pill bottle, my wallet, medication, mayonnaise jar, sofa, bottle, virgin remote control, money, slipper, fish food, styrofoam cup, blue facemask, i phone 11 pro, my keyboard, ipad, nobile phone stand, glasses cleaning wipe, bottle of alcoholic drink, cooker, tv remote, front door keys, tweezers, shed door, kettle, alcohol wipe, make up, battery drill, spanner, apple tv remote, bag, phone case, mini bluetooth keyboard, stylus, shoulder bag, comb, my keys, mirror, my clock, eye glasses, nike trainers, my water bottle, garden wall, sharp knife, my shoes, back pack, grinder, 12 measuring cup, iphone, phone, covid mask, mountain dew can, wheelie bin, car, headphone, keys, large sewing needle, miter saw, apple watch, chicken instant noodles, tv re-

mote control, adaptive tennis ball, embroidery thread cone, washing basket, wrist watch, lime green marker, glass, boot, bed, bose earpods, television remote control, dining table setup, toddler cup, tape measure, adaptive washing machine, pop bottle, electric sanding disc, washing machine, my sennheiser pxc 350-2, ladies silver bracelet

A.4. Colors and materials

Three annotators manually annotated the ORBIT validation and test objects (208 objects) with their color and material. In most cases, each object was labeled with one color and one material, but in some cases up to two labels were selected (*e.g.* a water bottle with a plastic body and metal lid was assigned “plastic metal” as its material). The labels were iterated until all three annotators agreed. All colors and materials were selected from the following lists:

Colors [20 colors]: red, silver, yellow, grey, dark, pink, multicolour, purple, white, beige, burgundy, maroon, blue, green, black, gold, brown, light, transparent, orange

Materials [23 materials]: rubber, crystal, cardboard, denim, material, styrofoam, stone, glass, foam, cloth, leather, ceramic, plastic, wood, paper, embroidered, wooden, suede, canvas, patterned, metal, cotton, lacquered

A.5. Textual analysis of LAION-400M, LAION-2B and DataComp-1B

Our aim is to quantify the prevalence of disability content in large-scale datasets used to pre-train LMMs – specifically, LAION-400M [57], LAION-2B [58] and DataComp-1B (or XL) [23]. To do this, we first extract all visual concepts from the captions of each dataset (see Algorithm 1). We define a visual concept as a noun phrase that contains a physical object (*e.g.* “park bench”). We consider a noun phrase as a phrase that contain a common noun and optional adjectives (*e.g.* “green park bench”). We consider the common noun to be a physical object if it traverses the “entity”, “physical_entity”, and “object” hypernyms and then either the “artifact”, “whole”, “part”, or “living_thing” hypernym in the WordNet tree hierarchy [40] (see Algorithm 2). In Tab. A.2 we report the top ten noun phrases extracted from LAION-400M, LAION-1B and DataComp-1B. We see that many of these are shared across all three datasets, including “image”, “photo”, “man”, and “woman”.

A.5.1 Prevalence of disability vs non-disability objects

In Sec. 4.1.2 of the main paper, we quantify how often disability and non-disability objects occur in the extracted visual concepts. We use the ORBIT objects as a seed set (see lists in App. A.3). Three annotators first grouped the object labels into clusters based on object similarity (*e.g.* all guide canes, all spectacles). Two synonyms were then assigned per cluster to account for different ways objects can be described. Early results showed that the disability clusters’ synonyms occurred extremely rarely in the visual concepts, so we broadened this to 5-16 synonyms per cluster. The 222 and 312 synonyms for the disability and non-disability clusters are provided in Tab. A.3 and Tab. A.4, respectively.

We then count how many times each of these synonyms appears in the extracted visual concepts (see counts in Tabs. A.3 and A.4). We do this using direct string matching allowing for partial matches (*e.g.* “braille note taker” is marked as present in the visual concept “cheap braille note taker”). Before matching, we lower-case and remove punctuation from all synonyms and visual concepts following typical VQA practices (see `processPunctuation` function in the [GT-Vision-Lab/VQA repo](#)). For synonyms that contain multiple words, we also allow for multiple spellings (*e.g.* tread mill and treadmill).

A.5.2 Prevalence of colors vs materials

In Sec. 4.3.2 of the main paper, we quantify how often colors and materials occur in each dataset. We use the color and materials in App. A.4 and directly count their frequency (via partial string matching, as above) in each dataset’s extracted visual concepts (see counts in Tab. A.5).

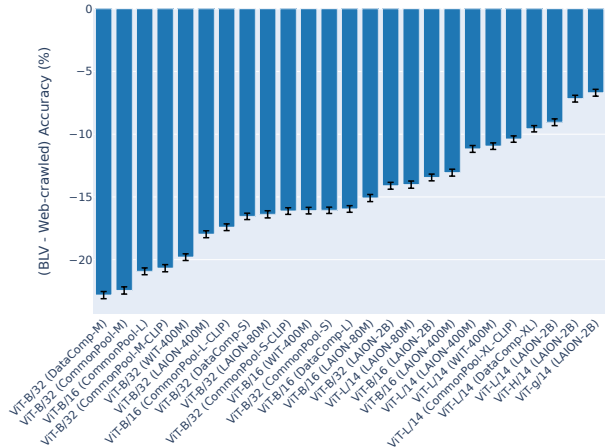


Figure B.1. **No CLIP variant achieves parity between BLV and web-crawled datasets on a standardized zero-shot image classification task (see Sec. 3.1).** Each bar represents the variant’s delta in average accuracy (with 95% c.i.) between all images sampled from BLV versus web-crawled datasets.

B. Extended results

B.1. BLV versus web-crawled data

We extend Fig. 1 in the main paper, with the delta in average accuracy between BLV and web-crawled datasets, reported per CLIP variant in Fig. B.1. We see that no model achieves parity, with smaller architectures (*e.g.* ViT-B/32) pre-trained on smaller datasets (*e.g.* DataComp-M, CommonPool-M) generally having a larger delta than larger architectures (*e.g.* ViT-g/14, ViT-h/14) pre-trained on larger datasets (*e.g.* LAION-2B).

B.2. Robustness to image content from BLV users

B.2.1 Disability objects are less well recognized than non-disability objects

Fig. B.2 shows the difference in average accuracy for disability, exclusive disability and non-disability objects for each CLIP variant, ordered by pre-training dataset size, on the ORBIT Clean (Fig. B.2a) and ORBIT Clutter (Fig. B.2b) datasets. Fig. B.3a and Fig. B.3b show the same, but with the CLIP variants ordered by architecture size. From these figures, we see that the accuracy difference between disability/exclusive disability objects and non-disability objects remains largely constant regardless of test dataset, pre-training dataset size, and architecture size.

B.2.2 A few-shot approach can sometimes reduce the disability and non-disability accuracy gap

In Sec. 4.1.3 of the main paper, we show how a few-shot approach can be effective at reducing the accuracy difference between disability and non-disability objects in some scenarios. We use ProtoNets [59] as the few-shot approach,

Algorithm 1 Pseudocode for `extract_noun_phrases(captions: List[str]) -> List[str]`

```
1: regex_pattern = r""NP: <DT|PRP$\$>?<JJ>*<NN|NNS>"" # a noun phrase (NP) contains a
   singular or plural noun (NN/NNS) which may be prefixed by an article (DT)/possessive
   pronoun (PRP) and/or adjectives (JJ)
2: chunker = nltk.RegexpParser(regex_pattern)
3: noun_phrases = []
4: for caption in captions do
5:     tokens = nltk.word.tokenize(caption.lower()) # tokenize
6:     pos_tags = nltk.pos_tag(tokens) # extract parts of speech
7:     np_tree = chunker(pos_tags) # extract noun phrase tree
8:     noun = extract_noun(np_tree) # extract noun (NN or NNS) from tree
9:     if is_physical_object(noun) then
10:         cleaned_np = clean_np(np_tree) # remove DT/PRP and singularize noun
11:         noun_phrases.extend(cleaned_np)
12:     end if
13: end for
14: return noun_phrases
```

Algorithm 2 Pseudocode for `is_physical_object(word: str) -> bool:`

```
1: synsets = wordnet.synsets(word, "n") # get WordNet noun synsets
2: for synset in synsets do
3:     paths = synset.hypernym_paths() # get hypernym paths
4:     for path in paths do
5:         # path is a list e.g. [Synset("entity.n.01"), ..., Synset("bench.n.01")]
6:         if (path contains "entity.n" AND "physical_entity.n" AND "object.n" \
7:             AND ("artifact.n" OR "whole.n" OR "part.n" OR "living_thing.n")):
8:             return True
9:     end for
10: end for
11: return False
```

which computes an average embedding (or prototype) for each object class by simply averaging the embeddings of K training images for each class. A test image is then classified as the class whose prototype is most similar to the image’s embedding, where similarity is measured by Euclidean distance. We extend Tab. 3 in the main paper with Figs. B.4a and B.4b here. Fig. B.4a shows ProtoNets can reduce the accuracy difference between disability/exclusive disability and non-disability objects on the ORBIT Clean dataset. Fig. B.4b shows ProtoNets’ results on the ORBIT Clutter dataset, however, the few-shot adaptation is less effective, even with 40 shots per object.

In Figs. B.5a and B.5b, we examine how CLIP’s pre-training dataset size influences the few-shot adaptation on ORBIT Clean and Clutter, respectively. We split the CLIP variants into three groups: those pre-trained on 0-100M examples, 100-1000M examples, and 1B+ examples. For each group, we average the delta in accuracy between disability and non-disability objects for all CLIP variants in that group, for each shot setting. For ORBIT Clean, we see that as the pre-training dataset and the number of shots increase,

the delta generally decreases – with 1B+ pre-training examples and a 40-shot setting achieving the lowest delta (-0.08) between disability and non-disability object accuracy. For ORBIT Clutter, however, this trend is less pronounced. Increasing the number of pre-training examples does reduce the delta generally, but the best setting (1B+ pre-training examples, 40 shots) still sees a delta of -10.98 percentage points. Furthermore, for under 100M pre-training examples, the delta remains largely constant (around -18 percentage points) suggesting that a few-shot approach is less effective if the model has not seen enough pre-training data.

B.3. Robustness to image quality from BLV users

We include the raw marginal effects of each quality issue on model accuracy for all CLIP variants in Tabs. B.3 to B.5. These correspond to Fig. 3 in the main paper, with experimental details provided in Sec. 4.2. We note that the same image may be sampled multiple times as a result of the episodic sampling procedure (see Sec. 3.1). No two tasks share the same set of N objects, however, so for a given im-

Table A.2. Top 10 noun phrases extracted from the captions of the LAION-400M [57], LAION-2B [58] and DataComp-1B [23] datasets. See extraction protocol in App. A.5.

LAION-400M		LAION-2B		DataComp-1B	
Noun phrase	Occurrence count	Noun phrase	Occurrence count	Noun phrase	Occurrence count
image	8,930,057	image	69,739,546	image	35,784,938
photo	7,559,650	photo	55,970,047	photo	28,612,087
vector	4,523,146	vector	29,203,691	vector	14,157,166
man	3,458,074	stock	22,209,987	stock	13,062,936
design	3,052,442	man	21,261,979	background	10,829,331
background	2,760,736	background	20,873,562	design	10,473,440
woman	2,513,375	picture	19,322,717	home	8,731,799
home	2,446,557	design	18,335,833	picture	8,523,946
stock	2,236,958	home	17,747,460	man	8,453,533
picture	2,235,123	woman	17,001,793	view	7,595,762

Table A.3. Disability object clusters, their synonyms, and their prevalence in the LAION-400M (L400M), LAION-2B (L2B) and DataComp-1B (DC1B) datasets. Numbers reported are the total number of times each cluster’s synonyms appeared in the dataset’s extracted visual concepts (total visual concepts – LAION-400M: 384,468,921; LAION-2B: 2,737,763,447; and DataComp-1B: 1,342,369,058).

Object cluster	Synonyms	L400M	L2B	DC1B
braille readers	braille note taker, braille reader, braille display, braille notetaker, braille tablet, braille computer, braille keyboard, orbit reader, braillepen slim braille keyboard, braillepen slim keyboard	4	8	6
dictaphones	dictaphone, digital recorder, voice recorder, dictation machine, audio recorder, voice recording device, dictation recorder, audio dictation device, voice transcription device, handheld recorder	40	185	168
digital book readers	digital book player, digital book reader, victor stream, victor reader stream, talking book, humanware reader, solo audiobook player, audiobook player	0	2	1
dog leads	dog lead, dogs lead, dog leash, dogs leash, leash, dog tether, dogs tether	434	1,663	1,402
dog poo	dog poo, dog poop, dog waste, dog scat, dog dung, canine faeces, canine feces, canine faeces	3	11	9
glasses	glass, sight glass, spectacle, eyeglass, reading glass, prescription glass, optical glass, corrective lens, bi focal, eyewear, frame, multi focal, optical, vision aid, spec	17,620	68,169	46,259
guide canes	guide cane, symbol cane, mobility cane, long cane, white cane, blind cane, white mobility cane, vision cane, assistive cane, visibility stick	21	79	50
hearing aids	hearing aid, hearing device, hearing amplifier, assistive listening device, hearing implant, cochlear implant, audio prosthesis, auditory prosthesis	23	122	77
inhalers	inhaler, asthma pump, asthma puffer, aerosol inhaler, inhalant delivery system	81	282	315
insulin pens	insulin pen, insulin injector, insulin delivery pen, insulin auto-injector, insulin syringe pen, insulin dispenser, insulin delivery system, insulin applicator, insulin dosing pen, diabetes pen	2	4	4
liquid level sensors	liquid level sensor, liquid level indicator, liquid level detector, liquid level gauge, water level sensor, water level indicator, water level detector, water level gauge	1	5	8
magnifiers	magnifier, magnifying glass, magnification aid, magnifying lens	89	390	362
audio labelers	penfriend, pen friend, audio labeller, audio labelling device, audio labelling pen, audio labelling tool, voice labeller, voice labelling pen, voice labelling device, voice labelling tool, speech-enabled labeller, speech-enabled labelling device, speech-enabled labelling pen, speech-enabled labelling tool, talking label maker, speech-based label printer	8	19	11

age, the model is always presented a different classification problem. The logistic regression is sensitive to input-output similarities, however, so we filter out all duplicate images to avoid biasing our sample. This resulted in 93,698 images for ORBIT Clean and 6,764 for VizWiz-Classification. We report the prevalence of each quality issue in these datasets in Tabs. B.1a and B.1b.

B.4. Robustness to language used by BLV users

B.4.1 CLIP classifies objects more accurately when they are described by color rather than material

We extend Tab. 4 in the main paper, with Tab. B.2 for the ORBIT Clutter dataset here. We see that the CLIP scores for the lower bound prompt (*i.e.* just the object name) are the highest, followed by the color prompt. Similar to Tab. 4,

Table A.4. Non-disability object clusters, their synonyms, and their prevalence in the LAION-400M (L400M), LAION-2B (L2B) and DataComp-1B (DC1B) datasets. Numbers reported are the total number of times each cluster’s synonyms appeared in the dataset’s extracted visual concepts (total visual concepts – LAION-400M: 384,468,921; LAION-2B: 2,737,763,447; and DataComp-1B: 1,342,369,058).

Object cluster	Synonyms	L400M	L2B	DC1B	Object cluster	Synonyms	L400M	L2B	DC1B
airpods	airpod, ear phone	655	2,077	1,763	make-up	make-up, make up	6,202	23,424	15,361
airpods cases	airpods case, airpods pro case	0	1	2	markers	marker, felt-tip pen	1,254	5,525	3,634
alcohol wipes	alcohol wipe, alcohol pad	1	1	0	measuring spoons	measuring spoon, measuring cup	10	50	55
bags	bag, backpack	15,250	52,351	34,650	medications	medication, pill	689	2,477	2,415
balls	ball, dog toy	6,540	23,289	14,888	mirrors	mirror, looking glass	4,755	18,788	13,327
bananas	banana, fruit	7,631	25,879	21,884	mice	bluetooth mouse, wireless mouse	3	13	16
baskets	basket, crate	3,375	12,380	8,821	mugs	mug, cup	10,360	39,145	30,334
beds	bed, mattress	8,087	35,714	23,575	nail clippers	nail clipper, tweezers	10	48	76
beers	beer, alcohol	458	2,145	1,564	nail polishes	nail polish, nail varnish	916	2,918	1,787
bins	bin, trash can	1,265	4,526	3,370	needles	needle, pin	7,531	27,235	23,560
bottles	bottle, thermos	6,113	21,684	20,062	paper clips	paper clip, paper fastener	106	386	366
bottle openers	bottle opener, cork screw	127	515	683	peelers	peeler, scraper	294	1,085	1,259
bracelets	bracelet, necklace	18,438	63,242	61,734	pencil cases	pencil case, pen case	33	161	152
brushes	brush, comb	3,866	13,624	11,820	phones	phone, iphone	5,217	19,911	11,396
bus stops	bus stop, bus station	36	163	103	phone chargers	phone charger, charging cable	2	13	17
cans	can, tin	2,648	10,906	7,629	phone stands	phone stand, ipad stand	9	32	17
cars	car, vehicle	22,867	84,568	57,283	plugs	plug, socket	3,340	12,519	10,972
CDs	compact disc, cd	9,675	31,511	14,609	police scanners	police scanner, radio scanner	1	2	2
cleaners	cleaner, surface spray	956	3,767	2,850	pops	pop, soda	5,795	20,591	13,506
clocks	clock, timekeeper	2,783	9,560	7,310	post boxes	post box, mail box	608	2,010	1,625
condom boxes	condom box, durex box	0	1	0	purses	purse, wallet	6,221	21,201	15,256
cookers	cooker, air fryer	654	2,612	2,718	radios	radio, receiver	4,638	15,645	12,138
cushions	cushion, pillow	9,081	33,694	24,195	rice	rice, noodle	4,725	16,489	13,833
deodorants	deodorant, perfume	1,564	5,875	5,556	rulers	ruler, tape measure	772	2,963	2,118
dog waste bins	dog waste bin, dog waste container	0	0	0	saucers	mustard, mayonnaise	1,050	3,877	2,941
doors	door, entrance	15,627	54,395	44,418	scissors	secateur, scissor	85	284	268
drills	drill, power tool	1,055	4,115	3,082	screwdrivers	screwdriver, spanner	383	1,466	2,109
electric saws	electric saw, chain saw	368	1,340	1,007	sheds	shed, tool shed	786	3,112	2,270
eye drops	eye drop, eye gel	6	35	18	shoes	shoe, sneaker	11,078	43,179	21,112
face masks	face mask, face covering	43	255	181	skipping ropes	skipping rope, jump rope	7	19	17
fans	fan, air cooler	7,638	24,013	15,585	sleep masks	sleep mask, eye mask	44	158	144
fish foods	fish food, fish flake	4	15	12	socks	sock, sockwear	2,959	8,765	6,443
fridges	fridge, freezer	1,784	5,516	4,926	sofas	sofa, couch	11,588	47,305	31,437
game controllers	wireless controller, game controller	5	29	36	spices	salt, pepper	2,814	9,365	9,068
gates	gate, gateway	3,133	11,797	8,287	styluses	apple pen, stylus	151	688	681
glasses	glass, tumbler	5,022	19,732	13,595	sunglasses	sunglass, shade	6,964	24,927	16,626
glasses cases	glasses case, sunglasses case	2	5	4	tables	desk, table	23,387	104,138	69,152
glasses cleaners	glasses cleaner, lens wipe	0	0	0	tambourines	tambourine, tamborine	89	362	375
gloves	glove, mitten	3,866	13,548	8,228	tea	tea, teabag	7,739	25,126	20,249
grinders	grinder, food processor	421	1,645	1,525	thread cones	thread cone, thread spool	15	49	51
hair clips	hair clip, headband	1,546	5,115	4,381	tissue boxes	tissue box, kleenex	12	29	26
hand sanitizers	hand sanitizer, hand santiser	0	10	5	toothbrushes	toothbrush, dental brush	537	2,329	2,443
hand saws	hand saw, hack saw	46	206	245	tread mills	tread mill, running machine	263	953	627
hats	hat, cap	10,100	34,840	24,185	t-shirts	t-shirt, tee	30,770	100,369	71,408
headphones	headphone, headset	2,241	7,843	6,396	TVs	tv, television	15,350	53,606	35,370
headphone cases	headphone case, headphones case	0	1	1	TV remotes	tv remote, remote control	203	866	658
hole punches	hole punch, paper punch	47	174	142	USB sticks	usb stick, flash drive	221	651	690
iPads	ipad, tablet	6,365	21,157	14,534	vapes	vape, e-cigarette	160	540	469
journals	journal, notebook	6,201	24,084	15,205	vases	vase, jug	4,853	16,954	16,971
kettles	kettle, toaster	1,128	4,879	4,348	walls	wall, fence	14,538	60,387	39,221
keys	key, key chain	8,770	31,370	22,409	wardrobes	wardrobe, cupboard	2,101	8,193	5,926
keyboards	keyboard, keypad	2,491	8,639	9,582	washing machines	washing machine, dryer	767	3,040	2,194
knives	knife, blade	4,012	15,521	17,972	watches	watch, smart watch	9,651	36,043	24,128
laptops	laptop, chromebook	7,522	24,913	17,948	watering cans	watering can, water can	14	110	66
lipsticks	lipstick, lip balm	1,453	5,213	4,635	weight benches	weight bench, gym bench	10	32	33

Table A.5. Colors/materials and their prevalence in the LAION-400M (L400M), LAION-2B (L2B) and DataComp-1B (DC1B) datasets. Numbers reported are the total number of times each color/material appeared in the dataset’s extracted visual concepts (total visual concepts – LAION-400M: 384,468,921; LAION-2B: 2,737,763,447; and DataComp-1B: 1,342,369,058).

		L400M	L2B	DC1B
Colors	beige	2,462	8,857	5,265
	black	87,366	323,959	207,730
	blue	53,947	193,504	131,665
	brown	21,904	84,994	55,553
	burgundy	1,261	4,127	2,446
	dark	10,888	36,735	25,818
	gold	5,882	19,564	12,894
	green	38,876	136,448	92,922
	grey	12,816	47,470	28,847
	light	31,413	120,232	92,203
	maroon	613	2,367	1,495
	multicolour	278	652	404
	orange	10,138	30,733	23,262
	pink	14,925	60,226	35,658
	purple	10,985	37,673	25,270
	red	45,267	165,576	104,289
	silver	8,789	33,614	27,766
	transparent	2,513	10,098	6,807
	white	94,014	366,224	236,456
	yellow	20,723	73,049	49,121
Materials	canvas	897	3,254	1,920
	cardboard	451	1,546	1,154
	ceramic	12,808	48,649	43,339
	cloth	3,498	13,528	8,968
	cotton	12,619	47,188	28,071
	crystal	12,663	46,968	36,329
	denim	5,130	15,265	8,720
	embroidered	2,090	8,026	3,904
	foam	675	2,234	1,894
	glass	3,264	11,420	7,889
	lacquered	351	1,679	1,077
	leather	1,214	4,234	2,505
	material	9,287	44,848	25,692
	metal	3,183	11,780	9,198
	paper	14,615	58,584	40,500
	patterned	1,057	3,841	1,947
	plastic	4,136	14,748	10,997
	rubber	4,620	18,765	12,234
	stone	7,262	28,321	20,570
	styrofoam	83	294	257
	suede	2,787	10,929	5,197
	wood	10,947	43,897	30,791
	wooden	18,239	73,016	51,445

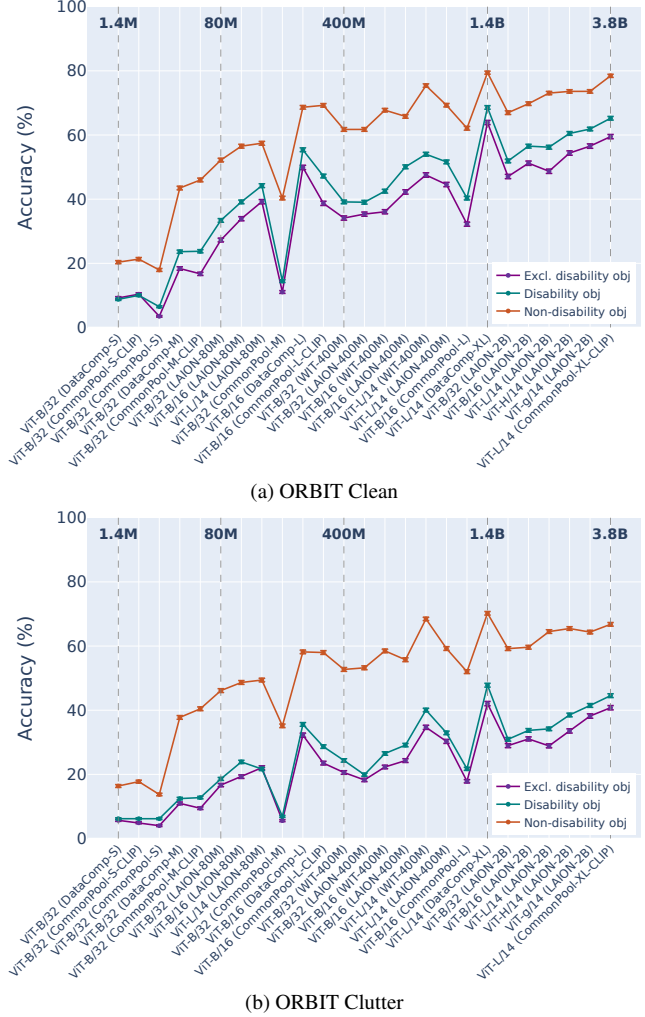


Figure B.2. CLIP’s difference in accuracy between disability and non-disability objects remains largely constant as its pre-training dataset increases. Zero-shot accuracy is averaged (with 95% c.i.) over images from ORBIT Clean/Clutter of each object type. Experimental details in Sec. 4.1.1.

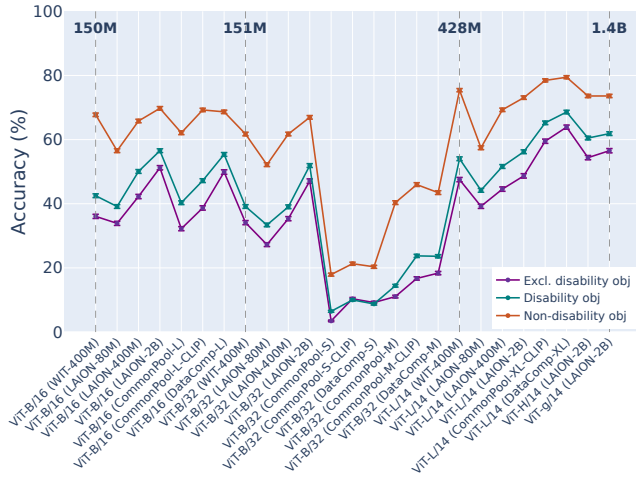
Table B.1. Prevalence of each quality issue in the (a) ORBIT Clean and (b) VizWiz-Classification datasets. Numbers reported as the raw counts of each issue and as a percentage of the total non-disability/disability images (ORBIT Clean) and total images (VizWiz-Classification).

	Total frames	Framing	Blur	Viewpoint	Occlusion	Lighting
Non-disability object	86,185	52,275 (60.7%)	28,235 (32.8%)	14,253 (16.5%)	11,302 (13.1%)	3,788 (4.4%)
Disability object	7,513	3,290 (43.8%)	2,237 (29.8%)	394 (5.2%)	976 (13.0%)	205 (2.7%)

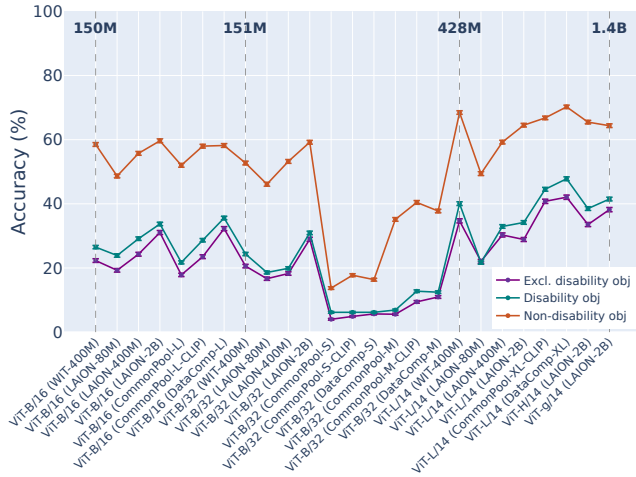
(a) ORBIT Clean

	Total frames	Framing	Blur	Viewpoint	Occlusion	Overexposed	Underexposed	Other
Non-disability object	6,764	3,715 (54.9%)	2,544 (37.6%)	1,118 (16.5%)	142 (2.1%)	327 (4.8%)	288 (4.3%)	16 (0.2%)

(b) VizWiz-Classification

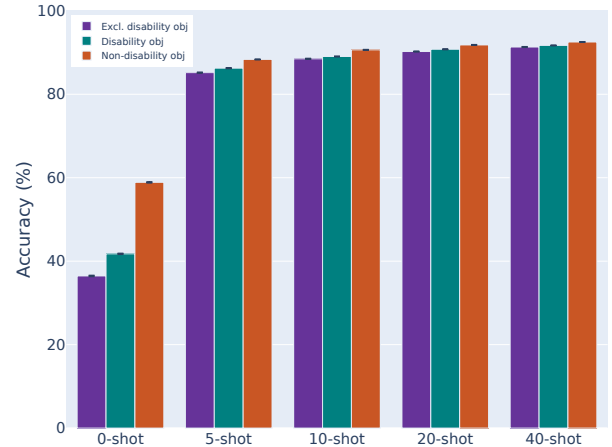


(a) ORBIT Clean

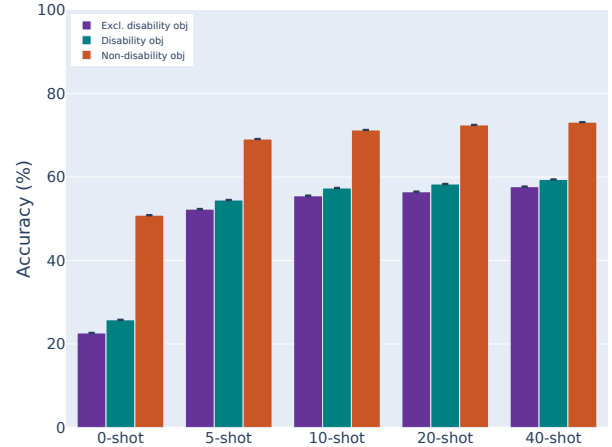


(b) ORBIT Clutter

Figure B.3. CLIP's difference in accuracy between disability and non-disability objects remains largely constant as its architecture size increases. Zero-shot accuracy is averaged (with 95% c.i.) over images from ORBIT Clean/Clutter of each object type. Experimental details in Sec. 4.1.1.



(a) ORBIT Clean



(b) ORBIT Clutter

Figure B.4. A few-shot approach (ProtoNets [59]) can reduce the accuracy gap between disability and non-disability objects, but not for realistic, cluttered images. Bars represent the average accuracy (with 95% c.i.) over all test frames for each shot setting ($K=[5, 10, 20, 40]$). $K=0$ is equivalent to the zero-shot setting described in Sec. 4.1.1. Experimental details in Sec. 4.1.3.

Pre-training dataset size	18+	-14.48	-0.35	-0.36	-0.26	-0.08
	100-1000M	-20.73	-1.45	-1.03	-0.58	-0.49
	0-100M	-15.73	-4.59	-3.45	-2.36	-1.94
		0-shot	5-shot	10-shot	20-shot	40-shot

(a) ORBIT Clean

Pre-training dataset size	18+	-26.16	-12.31	-11.22	-11.39	-10.98
	100-1000M	-28.35	-13.22	-12.19	-12.71	-11.94
	0-100M	-20.29	-18.49	-18.51	-18.55	-18.53
		0-shot	5-shot	10-shot	20-shot	40-shot

(b) ORBIT Clutter

Figure B.5. **The larger the dataset used to pre-train CLIP, the more effective a few-shot approach is at closing the accuracy gap between disability and non-disability objects on ORBIT Clean, but this is less so for ORBIT Clutter especially for pre-training datasets <100M examples.** Each block reports the average delta in accuracy between disability and non-disability objects for the models that fall within that group. Models: 25 CLIP variants.

Table B.2. **Including an object’s material in its prompt leads to text embeddings that are the least aligned with the object’s image embeddings.** CLIP scores [26] between image and prompt embeddings are averaged (with 95% c.i.) for 100 images per object per prompt type on ORBIT Clutter.

Prompt	Obj. name	Material + obj. name	Color + obj. name	Color + material + obj. name
CLIP Score	22.81 \pm 0.02	21.92 \pm 0.02	22.73 \pm 0.02	21.86 \pm 0.02

Table B.3. **Marginal effects of explanatory variables on CLIP’s zero-shot classification accuracy (with ViT-B/16 vision encoders) on the ORBIT Clean and VizWiz-Classification datasets.** Main values are marginal effects, while values in brackets are p-values. */**/* indicates within 90/95/99% confidence interval, respectively. Experimental details in Sec. 4.2. L-80M=LAION-80M, L-400M=LAION-400M, L-2B=LAION-2B, DC-L=DataComp-L, CP-L=CommonPool-L, CP-L-CLIP=CommonPool-L-CLIP.

Dataset	Explanatory variable	ViT-B/16 (WIT)	ViT-B/16 (L-80M)	ViT-B/16 (L-400M)	ViT-B/16 (L-2B)	ViT-B/16 (DC-L)	ViT-B/16 (CP-L)	ViT-B/16 (CP-L-CLIP)
ORBIT Clean	framing	0.044*** (1.0)	0.044*** (1.0)	0.043*** (1.0)	0.042*** (1.0)	0.02*** (1.0)	0.054*** (1.0)	0.033*** (1.0)
	blur	-0.095*** (1.0)	-0.139*** (1.0)	-0.117*** (1.0)	-0.139*** (1.0)	-0.153*** (1.0)	-0.135*** (1.0)	-0.146*** (1.0)
	viewpoint	-0.115*** (1.0)	-0.139*** (1.0)	-0.122*** (1.0)	-0.1*** (1.0)	-0.092*** (1.0)	-0.094*** (1.0)	-0.086*** (1.0)
	occlusion	-0.054*** (1.0)	-0.099*** (1.0)	-0.088*** (1.0)	-0.086*** (1.0)	-0.088*** (1.0)	-0.104*** (1.0)	-0.099*** (1.0)
	lighting	-0.213*** (1.0)	-0.212*** (1.0)	-0.262*** (1.0)	-0.245*** (1.0)	-0.226*** (1.0)	-0.297*** (1.0)	-0.246*** (1.0)
	excl. disability obj	-0.328*** (1.0)	-0.341*** (1.0)	-0.278*** (1.0)	-0.245*** (1.0)	-0.258*** (1.0)	-0.352*** (1.0)	-0.319*** (1.0)
	excl. disability obj:framing	0.143*** (1.0)	0.057*** (0.999)	0.108*** (1.0)	0.097*** (1.0)	0.141*** (1.0)	0.06*** (1.0)	0.097*** (1.0)
	excl. disability obj:blur	0.018 (0.764)	0.074*** (1.0)	-0.015 (0.673)	-0.017 (0.747)	-0.021 (0.843)	0.044*** (0.991)	0.011 (0.551)
	excl. disability obj:viewpoint	0.113*** (1.0)	0.207*** (1.0)	0.173*** (1.0)	0.18*** (1.0)	0.069** (0.969)	0.07* (0.947)	0.065** (0.959)
	excl. disability obj:occlusion	-0.058*** (0.994)	0.007 (0.22)	-0.051** (0.984)	-0.004 (0.169)	-0.006 (0.239)	0.032 (0.837)	0.03 (0.872)
	excl. disability obj:lighting	-0.239*** (1.0)	-0.128** (0.969)	-0.238*** (1.0)	-0.208*** (1.0)	-0.294*** (1.0)	-0.204** (0.986)	-0.214*** (0.999)
VizWiz-Classification	framing	0.001 (0.051)	0.001 (0.093)	-0.009 (0.542)	-0.035*** (0.995)	-0.01 (0.599)	-0.028** (0.979)	-0.018 (0.862)
	blur	-0.028** (0.976)	-0.019 (0.867)	-0.009 (0.52)	0 (0.006)	-0.014 (0.727)	-0.015 (0.761)	-0.02 (0.885)
	rotation	-0.079*** (1.0)	-0.116*** (1.0)	-0.055*** (0.999)	-0.086*** (1.0)	-0.07*** (1.0)	-0.092*** (1.0)	-0.09*** (1.0)
	occlusion	-0.096** (0.978)	-0.144*** (0.999)	-0.138*** (0.999)	-0.187*** (1.0)	-0.142*** (0.999)	-0.209*** (1.0)	-0.186*** (1.0)
	overexposure	-0.034 (0.777)	0.023 (0.592)	-0.033 (0.758)	-0.011 (0.301)	-0.07** (0.987)	-0.064** (0.974)	-0.029 (0.688)
	underexposure	-0.02 (0.498)	-0.088*** (0.996)	-0.084*** (0.995)	-0.065** (0.969)	-0.051* (0.91)	-0.073** (0.984)	-0.074** (0.985)
	other	0.129 (0.668)	0.099 (0.579)	0.009 (0.059)	0.243* (0.911)	0.017 (0.112)	0.115 (0.63)	0.04 (0.251)

Table B.4. **Marginal effects of explanatory variables on CLIP’s zero-shot classification accuracy (with ViT-B/32 vision encoders) on the ORBIT Clean and VizWiz-Classification datasets.** Main values are marginal effects, while values in brackets are p-values. */**/** indicates within 90/95/99% confidence interval, respectively. Experimental details in Sec. 4.2. L-80M=LAION-80M, L-400M=LAION-400M, L-2B=LAION-2B, DC-S=DataComp-S, DC-M=DataComp-M, CP-S=CommonPool-S, CP-S-CLIP=CommonPool-S-CLIP, CP-M=CommonPool-M, CP-M-CLIP=CommonPool-M-CLIP.

Dataset	Explanatory variable	ViT-B/32 (WIT)	ViT-B/32 (L-80M)	ViT-B/32 (L-400M)	ViT-B/32 (L-2B)	ViT-B/32 (DC-S)	ViT-B/32 (DC-M)	ViT-B/32 (CP-S)	ViT-B/32 (CP-S-CLIP)	ViT-B/32 (CP-M)	ViT-B/32 (CP-M-CLIP)
ORBIT Clean	framing	0.062*** (1.0)	0.046*** (1.0)	0.054*** (1.0)	0.051*** (1.0)	0.035*** (1.0)	0.101*** (1.0)	0.046*** (1.0)	0.05*** (1.0)	0.11*** (1.0)	0.097*** (1.0)
	blur	-0.106*** (1.0)	-0.128*** (1.0)	-0.142*** (1.0)	-0.132*** (1.0)	-0.025*** (1.0)	-0.113*** (1.0)	-0.029*** (1.0)	-0.05*** (1.0)	-0.103*** (1.0)	-0.119*** (1.0)
	viewpoint	-0.096*** (1.0)	-0.111*** (1.0)	-0.111*** (1.0)	-0.099*** (1.0)	-0.036*** (1.0)	-0.058*** (1.0)	-0.011*** (0.999)	-0.032*** (1.0)	-0.045*** (1.0)	-0.092*** (1.0)
	occlusion	-0.075*** (1.0)	-0.095*** (1.0)	-0.094*** (1.0)	-0.088*** (1.0)	-0.088*** (1.0)	-0.142*** (1.0)	-0.072*** (1.0)	-0.08*** (1.0)	-0.109*** (1.0)	-0.123*** (1.0)
	lighting	-0.174*** (1.0)	-0.297*** (1.0)	-0.292*** (1.0)	-0.261*** (1.0)	-0.117*** (1.0)	-0.228*** (1.0)	-0.209*** (1.0)	-0.217*** (1.0)	-0.296*** (1.0)	-0.294*** (1.0)
	excl. disability obj	-0.33*** (1.0)	-0.359*** (1.0)	-0.311*** (1.0)	-0.265*** (1.0)	-0.297*** (1.0)	-0.311*** (1.0)	-0.124*** (1.0)	-0.204*** (1.0)	-0.415*** (1.0)	-0.576*** (1.0)
	excl. disability obj:framing	0.137*** (1.0)	0.118*** (1.0)	0.124*** (1.0)	0.113*** (1.0)	-0.016 (0.479)	-0.094*** (1.0)	-0.259*** (1.0)	-0.147*** (1.0)	-0.064*** (0.99)	0.049* (0.949)
	excl. disability obj:blur	0.072*** (1.0)	-0.012 (0.458)	0.011 (0.479)	0.013 (0.593)	0.028 (0.759)	0.045** (0.959)	-0.01 (0.347)	0.043* (0.948)	0.062** (0.987)	0.062** (0.984)
	excl. disability obj:viewpoint	0.128*** (1.0)	0.166*** (1.0)	0.06* (0.911)	0.129*** (1.0)	0.219*** (1.0)	0.143*** (0.999)	0.148** (0.99)	0.225*** (1.0)	0.259*** (1.0)	0.268*** (1.0)
	excl. disability obj:occlusion	-0.151*** (1.0)	0.109*** (1.0)	-0.022 (0.659)	-0.025 (0.779)	0.176*** (1.0)	0.036 (0.758)	-0.013 (0.313)	0.038 (0.791)	0.166*** (1.0)	0.31*** (1.0)
	excl. disability obj:lighting	-0.188*** (1.0)	0.025 (0.312)	-0.255*** (0.998)	-0.285*** (1.0)	0.17*** (1.0)	-0.313** (0.98)	0.188*** (0.998)	0.123* (0.945)	-0.156 (0.768)	-0.273* (0.904)
VizWiz-Classification	framing	0.011 (0.634)	-0.024** (0.963)	-0.002 (0.115)	-0.034*** (0.995)	-0.015* (0.945)	-0.034*** (0.997)	0.014* (0.908)	-0.005 (0.441)	-0.021* (0.936)	-0.003 (0.19)
	blur	-0.009 (0.517)	-0.003 (0.202)	-0.029** (0.98)	0.007 (0.433)	-0.005 (0.455)	-0.026** (0.97)	-0.001 (0.106)	-0.02** (0.97)	-0.033*** (0.995)	-0.03** (0.984)
	rotation	-0.072*** (1.0)	-0.153*** (1.0)	-0.052*** (0.998)	-0.113*** (1.0)	-0.067*** (1.0)	-0.092*** (1.0)	-0.075*** (1.0)	-0.117*** (1.0)	-0.16*** (1.0)	-0.111*** (1.0)
	occlusion	-0.144*** (0.999)	-0.121*** (0.993)	-0.17*** (1.0)	-0.132*** (0.997)	-0.057* (0.907)	-0.175*** (1.0)	-0.098** (0.987)	-0.128*** (0.996)	-0.182*** (1.0)	-0.216*** (1.0)
	overexposure	0.035 (0.775)	-0.063** (0.974)	-0.059** (0.961)	-0.025 (0.618)	-0.004 (0.17)	-0.082*** (0.995)	-0.013 (0.489)	-0.027 (0.778)	-0.11*** (1.0)	-0.063** (0.972)
	underexposure	-0.009 (0.237)	-0.034 (0.754)	-0.095*** (0.998)	-0.046 (0.868)	-0.019 (0.648)	-0.081*** (0.992)	-0.008 (0.307)	-0.049* (0.948)	-0.011 (0.294)	-0.094*** (0.997)
	other	-0.13 (0.705)	0.176 (0.877)	0.06 (0.368)	0.063 (0.388)	0.002 (0.023)	0.084 (0.543)	-0.106 (0.635)	0.107 (0.876)	-0.101 (0.577)	0.016 (0.103)

Table B.5. Marginal effects of explanatory variables on CLIP’s zero-shot classification accuracy (with ViT-L/14, ViT-H/14 and ViT-g/14 vision encoders) on the ORBIT Clean and VizWiz-Classification datasets. Main values are marginal effects, while values in brackets are p-values. */**/* indicates within 90/95/99% confidence interval, respectively. Experimental details in Sec. 4.2. L-80M=LAION-80M, L-400M=LAION-400M, L-2B=LAION-2B, DC-XL=DataComp-XL, CP-XL-CLIP=CommonPool-XL-CLIP.

Dataset	Explanatory variable	ViT-L/14 (WIT)	ViT-L/14 (L-80M)	ViT-L/14 (L-400M)	ViT-L/14 (L-2B)	ViT-L/14 (DC-XL)	ViT-L/14 (CP-XL-CLIP)	ViT-H/14 (L-2B)	ViT-g/14 (L-2B)
ORBIT Clean	framing	-0.013*** (1.0)	0.027*** (1.0)	0.03*** (1.0)	0.019*** (1.0)	-0.016*** (1.0)	-0.018*** (1.0)	0.025*** (1.0)	0.025*** (1.0)
	blur	-0.098*** (1.0)	-0.139*** (1.0)	-0.131*** (1.0)	-0.129*** (1.0)	-0.099*** (1.0)	-0.114*** (1.0)	-0.113*** (1.0)	-0.123*** (1.0)
	viewpoint	-0.117*** (1.0)	-0.15*** (1.0)	-0.114*** (1.0)	-0.124*** (1.0)	-0.087*** (1.0)	-0.09*** (1.0)	-0.115*** (1.0)	-0.106*** (1.0)
	occlusion	-0.065*** (1.0)	-0.091*** (1.0)	-0.085*** (1.0)	-0.078*** (1.0)	-0.068*** (1.0)	-0.068*** (1.0)	-0.078*** (1.0)	-0.082*** (1.0)
	lighting	-0.174*** (1.0)	-0.288*** (1.0)	-0.248*** (1.0)	-0.223*** (1.0)	-0.189*** (1.0)	-0.184*** (1.0)	-0.173*** (1.0)	-0.222*** (1.0)
	excl. disability obj	-0.26*** (1.0)	-0.235*** (1.0)	-0.257*** (1.0)	-0.258*** (1.0)	-0.223*** (1.0)	-0.245*** (1.0)	-0.267*** (1.0)	-0.283*** (1.0)
	excl. disability obj:framing	0.085*** (1.0)	0.138*** (1.0)	0.085*** (1.0)	0.112*** (1.0)	0.16*** (1.0)	0.127*** (1.0)	0.103*** (1.0)	0.138*** (1.0)
	excl. disability obj:blur	-0.012 (0.644)	-0.04** (0.982)	-0.013 (0.635)	-0.032** (0.98)	-0.083*** (1.0)	-0.051*** (1.0)	-0.039*** (0.997)	-0.027** (0.956)
	excl. disability obj:viewpoint	0.101*** (1.0)	0.113*** (0.999)	0.144*** (1.0)	0.108*** (1.0)	0.041 (0.884)	0.119*** (1.0)	0.125*** (1.0)	0.118*** (1.0)
	excl. disability obj:occlusion	-0.052*** (0.998)	-0.023 (0.696)	0.014 (0.56)	-0.005 (0.237)	-0.004 (0.218)	-0.018 (0.765)	0.011 (0.492)	0.041** (0.982)
	excl. disability obj:lighting	-0.192*** (1.0)	-0.369*** (1.0)	-0.185*** (0.999)	-0.147*** (0.998)	-0.133*** (1.0)	-0.153*** (1.0)	-0.176*** (1.0)	-0.243*** (1.0)
VizWiz-Classification	framing	0.007 (0.434)	-0.04*** (0.999)	-0.038*** (0.998)	-0.014 (0.748)	-0.002 (0.138)	0.013 (0.737)	-0.006 (0.402)	-0.008 (0.503)
	blur	-0.005 (0.317)	0.015 (0.773)	-0.014 (0.737)	-0.002 (0.131)	-0.007 (0.404)	-0.009 (0.552)	-0.012 (0.676)	-0.019 (0.87)
	rotation	-0.052*** (0.999)	-0.093*** (1.0)	-0.076*** (1.0)	-0.031* (0.948)	-0.065*** (1.0)	-0.048*** (0.998)	-0.066*** (1.0)	-0.101*** (1.0)
	occlusion	-0.118*** (0.996)	-0.074* (0.915)	-0.143*** (0.999)	-0.126*** (0.998)	-0.134*** (0.999)	-0.126*** (0.999)	-0.136*** (0.999)	-0.106** (0.987)
	overexposure	-0.004 (0.105)	-0.017 (0.449)	-0.052* (0.934)	0.001 (0.015)	0.015 (0.416)	-0.018 (0.499)	0.019 (0.499)	0.004 (0.101)
	underexposure	-0.023 (0.57)	-0.034 (0.743)	-0.092*** (0.998)	-0.076*** (0.99)	-0.023 (0.562)	-0.056* (0.949)	-0.035 (0.761)	-0.028 (0.642)
	other	0.093 (0.525)	0.281** (0.972)	0.08 (0.473)	0.194 (0.83)	0.235 (0.876)	0.051 (0.314)	0.066 (0.395)	0.022 (0.137)

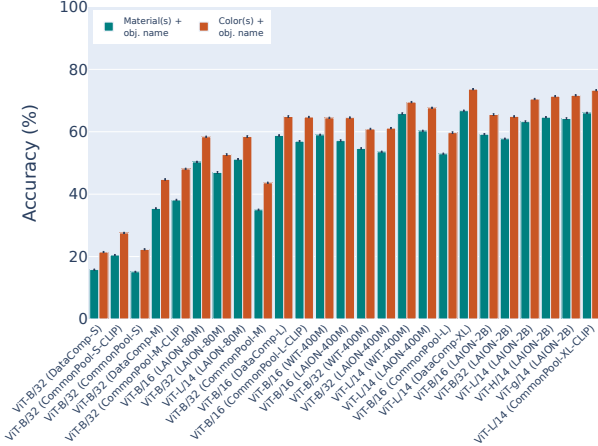


Figure B.6. **All CLIP variants classify objects more accurately when objects are described by their color rather than their material.** Each bar is the average accuracy (with 95% c.i.) over 200K images (100K ORBIT Clean, 100K ORBIT Clutter) for that CLIP variant when given either a material or color prompt. Variants ordered by pre-training dataset size.

we see that both the material prompt and the upper bound prompt which includes the object’s material have the lowest CLIP scores, suggesting that including the object’s material in the prompt harms embedding alignment.

We explore the impact this has on classifier accuracy by combining the textual prompts with the standard zero-shot set-up described in Sec. 3.1. Specifically, rather than embedding the raw ORBIT object labels for each task’s N classes, we instead embed their textual prompts. In the first experiment, we embed all N objects as their color prompts, and in the second as their material prompts. For both experiments, we use $T = 50$, $N = 20$, $M = 100$. In Fig. B.6, we see that across all CLIP variants, objects are classified more accurately when they are described by their color rather than their material – by 7.1 percentage points more, on average. We see that this difference is largely constant regardless of both architecture and pre-training dataset size (see Fig. B.7).

C. Example-based analysis

C.1. Standardized image selection

We run our analysis on 180 images spanning 20 objects which are selected through a standardized process as a way to systematically assess failure cases. Specifically, we select the 5 top- and bottom-performing (disability and non-disability) objects from the ORBIT dataset using the standardized zero-shot classification set-up (see objects in Tab. C.1). We take performance to be the average accuracy per object, computed over all the CLIP variants we considered. For each object, we extract the noun phrase from its raw label and apply simple pre-processing to ensure that it is unambiguous and concise (see cleaned phrases

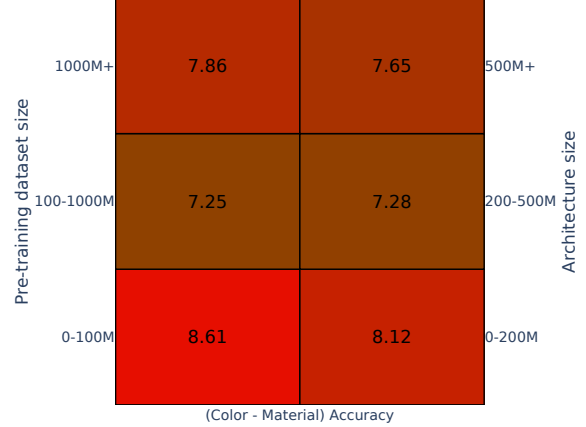


Figure B.7. **Increasing the pre-training dataset and architecture size only marginally reduces the difference in zero-shot accuracy between prompts describing an object by its color versus material.** Numbers reported are the delta in zero-shot accuracy between when a color versus material prompt is used as the text input to CLIP on the ORBIT Clean dataset. Each block averaged the delta for all CLIP variants that fall within that group. Experimental details in Sec. 4.3.1. Models: All 25 CLIP variants Table C.1. **Top- and bottom-performing disability and non-disability objects from the ORBIT dataset sed for the example-based analysis.** The noun phrases are extracted from the raw label and cleaned to ensure they are unambiguous and concise.

		Raw object label	Cleaned noun phrase
Disability objects	Top 5	my braille displat	braille sense display
		dog poo	dog poo
		dog lead	dog lead
		white cane	guide cane
		folded long guide cane	guide cane
	Bottom 5	victor reader stream	victor reader stream
		braille note	braille notetaker
		liquid level indicator	liquid level indicator
		liquid level indicator	liquid level indicator
		dictaphone	dictaphone
Non-disability objects	Top 5	back patio gate	gate
		local post box	post box
		wine glass	wine glass
		tv remote control	remote control
		remote control	remote control
	Bottom 5	digital dab radio	digital radio
		my clock	digital clock
		grinder	tobacco grinder
		dog streetball	ball
		shoulder bag	shoulder bag

in Tab. C.1). These cleaned noun phrases are used as the text prompts for all three downstream models we study. We then sample 9 images for each of the 20 objects – 6 from its clutter videos and 3 from its clean videos. We only sample images where the object is tagged as present. To increase image diversity, we ensure that images are sampled from all videos available for each object, and are sampled

Table C.2. **OWL-ViT’s mean intersection-over-union (IOU) is $\sim 2\times$ lower for disability compared to non-disability objects.** Mean IOU (with 95% c.i.) is computed between the predicted and ground-truth bounding box for each object.

	mean IOU
Disability objects	0.1323 (0.0947)
Non-disability objects	0.2488 (0.1829)

at even intervals. Specifically, for clean videos we sample the 3 frames at 25%, 50% and 75% positions, alternating the video we sample from each time (*e.g.* if an object has 2 clean videos, then we sample 1 frame at 25% of video 1, 1 frame at 50% of video 2, and 1 frame at 75% of video 1). For clutter videos, we sample 6 frames at 25%, 35%, 45%, 55%, 65% and 75%, also alternating the video for each sample. We limit frame sampling to between 25% and 75% of each video as ORBIT data collectors were instructed to start each video with the camera close to the object and then move it further away, so we wanted to exclude frames where the camera might be too close/far from the object.

C.2. Object detection with OWL-ViT

We extend Fig. 4 in the main paper with Fig. C.1 here, where we show one example of OWL-ViT’s bounding box detections for each of the 10 disability and 10 non-disability objects. Specifically, for each object we show the image that had the bounding box with the highest confidence score across all 9 images analyzed for that object. We see that the confidence scores in these images are $\sim 3\times$ lower for disability than non-disability objects, on average. We also see that for 4/10 disability objects, the incorrect object is detected (versus 2/10 non-disability objects).

We also report the mean intersection-over-union (IOU) between OWL-ViT’s predicted and ground-truth bounding box for each object in Tab. C.2. Since ground-truth bounding boxes are only publicly available for the clutter images, we manually annotated the remaining 6 clean images per object. Our results show that the mean IOU is $\sim 2\times$ lower for disability compared to non-disability objects. Taken together, these results suggest that overall, OWL-ViT performs less reliably and confidently for disability content.

C.3. Semantic segmentation with CLIPSeg

Semantic segmentation models are also highly likely to be integrated into assistive applications to help BLV users localize objects. We examine CLIPSeg [37] which trains a decoder on top of CLIP’s frozen vision and text encoders to enable zero-shot image segmentation from text prompts. Unlike OWL-ViT, CLIPSeg does not fine-tune the CLIP encoders, and its pre-trained embeddings are used directly. As before, we run all 180 images through the model with the cleaned noun phrases as text prompts. We find:

Table C.3. **CLIPSeg segments non-disability objects with higher confidence than non-disability objects.** The average confidence (with 95% c.i.) is reported over all pixels with a confidence above 0.1 within the object’s ground-truth bounding box.

	Avg in-box confidence
Disability objects	0.2181 (0.0842)
Non-disability objects	0.4276 (0.1286)

Table C.4. **CLIPSeg incorrectly segments disability objects more often than non-disability objects on the ORBIT Clutter dataset.** Numbers are the average confidence over all pixels *outside* the object’s ground-truth bounding box divided by the average confidence over all pixels in the image (with 95% c.i.), considering only pixels above a 0.1 confidence threshold.

	Avg confusion score
Disability objects	0.2698 (0.1990)
Non-disability objects	0.1292 (0.0749)

Segmentation maps of disability objects are less confident than those of non-disability objects. In Tab. C.3, we compute the average confidence value over all pixels in the segmentation map that fall within the ground-truth bounding box of the target object. To control for the degree of background present across bounding boxes (especially for irregular-shaped objects), we only consider pixels that have a confidence score greater than 0.1. With this, we find that CLIPSeg’s segmentation maps are $\sim 2\times$ more confident for non-disability objects compared to disability objects. In Fig. C.2, we show CLIPSeg’s segmentations for a guide cane versus a TV remote, two objects for which the confidence score difference was most pronounced.

Disability objects are more likely to be segmented as the incorrect object in realistic settings compared to non-disability objects. In Tab. C.4, we compute a confusion score per image: the average confidence score of all the pixels that fall *outside* the object’s ground-truth bounding box, divided by the average confidence score over all pixels in the image. This gives us a measure of how confidently the model is segmenting objects besides the ground-truth object, where a high score indicates the segmentation may be a false positive. Here we also only include confidence scores above a 0.1 threshold. We see that CLIPSeg is $\sim 2\times$ more likely to confuse a disability object with another object compared to a non-disability object in ORBIT Clutter images where multiple objects are present.

We include examples of this in Fig. C.3. Here we see that CLIPSeg fails to segment prominent disability objects (see liquid level indicators, guide canes, and Braille notetakers in Fig. C.3b) but succeeds in segmenting non-disability objects in similarly cluttered scenes (see shoulder bags and wine glasses in Fig. C.3a).

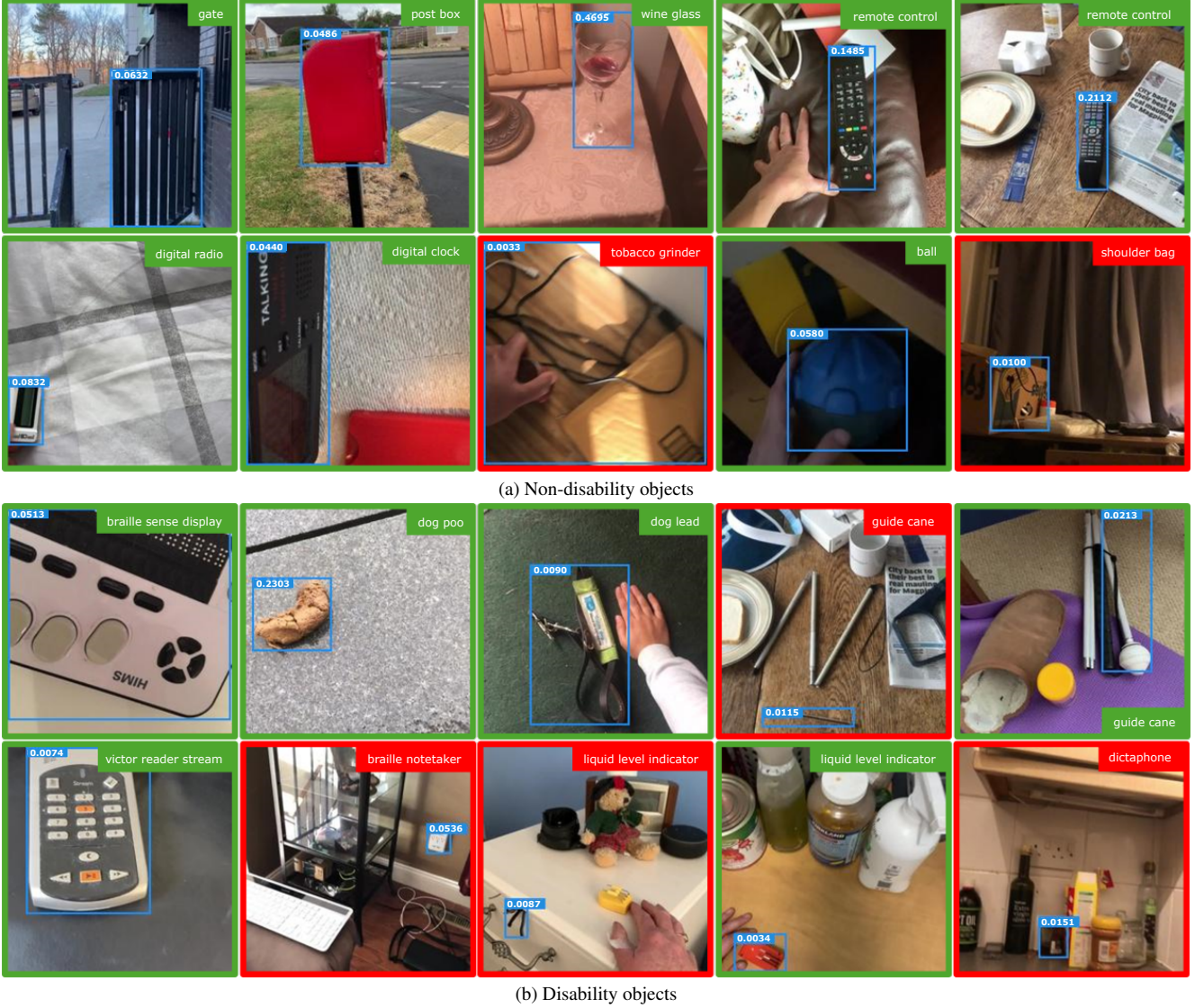


Figure C.1. OWL-ViT detects disability objects less confidently than non-disability objects. For each of the (a) 10 non-disability and (b) 10 disability objects, we show the image with the highest-scoring bounding box out of the 9 images analyzed for that object.

C.4. Text-to-image generations with DALL-E2

Prompt templates. Three annotators manually created two prompts for each of the 20 objects. The first prompt was just the cleaned noun phrase (taken from Tab. C.1). The second prompt combined the cleaned noun phrase with a surface and up to two adjacent objects. The surface and adjacent objects were selected to match an image from the ORBIT Clutter dataset of that object. The image was selected such that it had at least one adjacent object present. The prompt was then created with the template: “<object.name> on <surface> next to <adjacent-object-1> and <adjacent-object-2>” (e.g. “wine glass on a wooden table next to a bottle of wine and a candle”).

We extend Fig. 5 in the main paper with Fig. C.4 here.

We show DALL-E2’s generations for the two prompt types for non-disability (Fig. C.4a) and disability (Fig. C.4b) objects. Overall, we see that DALL-E2 does not generate correct representations for many of the disability objects, either defaulting to a common object or fabricating an object entirely. In contrast, the generations for non-disability objects are highly realistic and mostly correct.

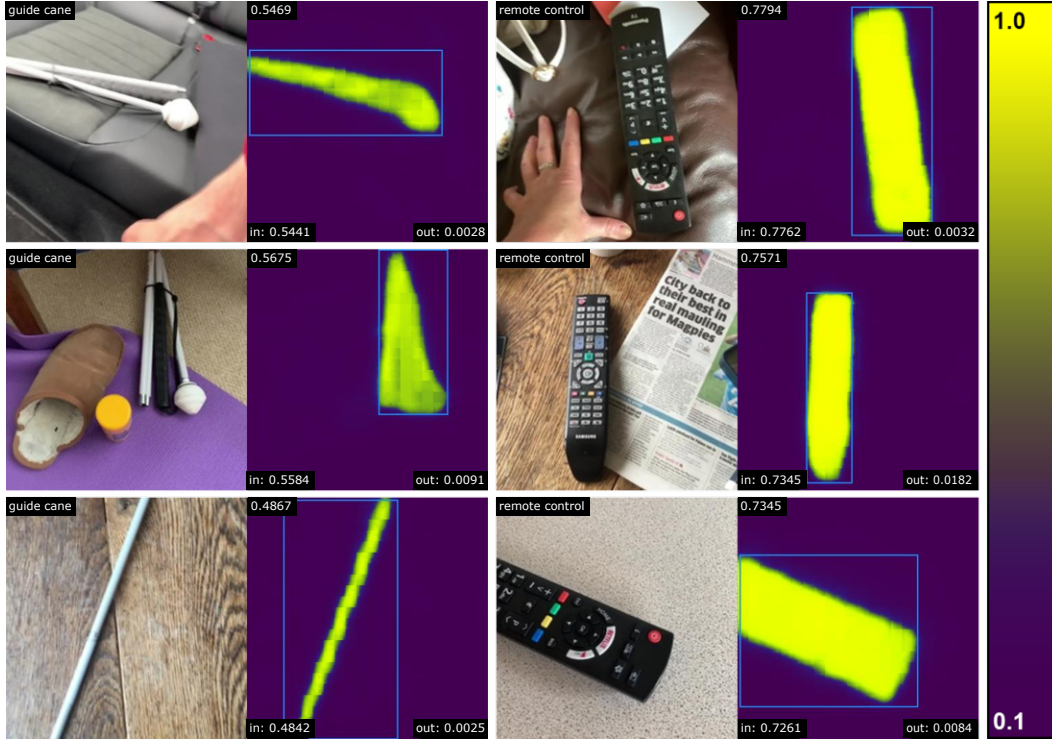


Figure C.2. CLIPSeg segments non-disability objects (right: TV remote) with higher confidence than disability objects (left: guide cane). For each image, we report the average confidence score over all pixels inside and outside the object’s ground-truth bounding box (“in” and “out”, respectively), considering only pixels above a 0.1 confidence threshold. See quantitative results in Tab. C.3.

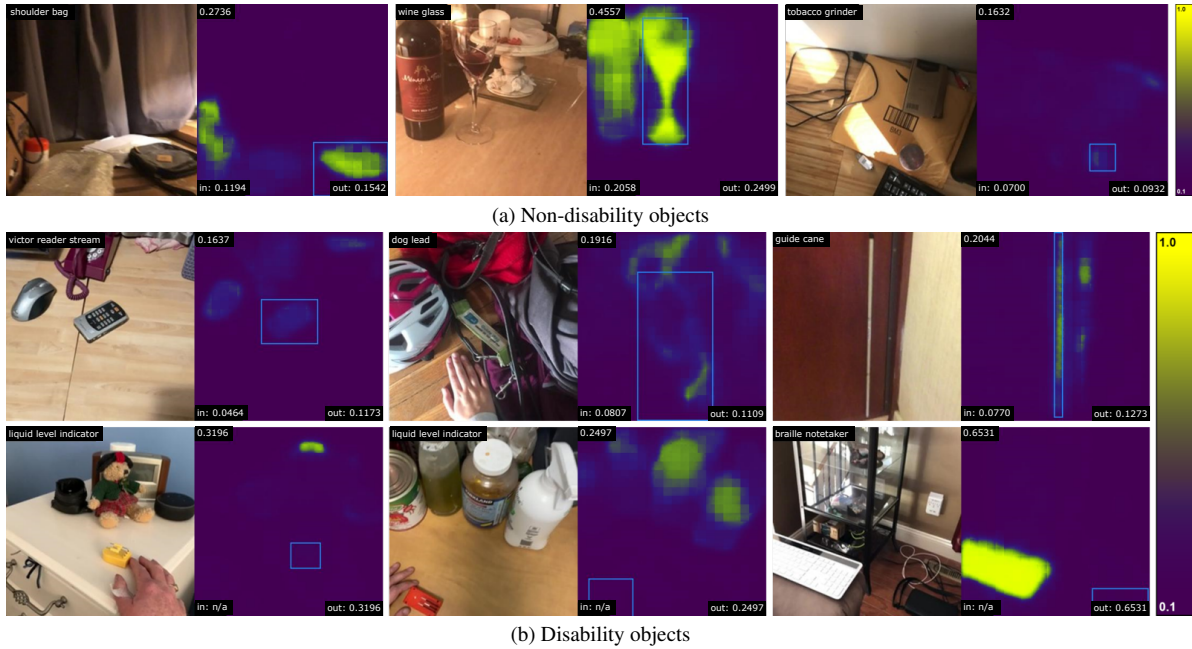


Figure C.3. CLIPSeg is more likely to segment disability objects (bottom) incorrectly in cluttered scenes compared to disability objects (top). For each image, we report the average confidence score over all pixels inside and outside the object’s ground-truth bounding box (“in” and “out”, respectively), considering only pixels above a 0.1 confidence threshold. The correct object is marked by the bounding box. See quantitative results in Tab. C.4.

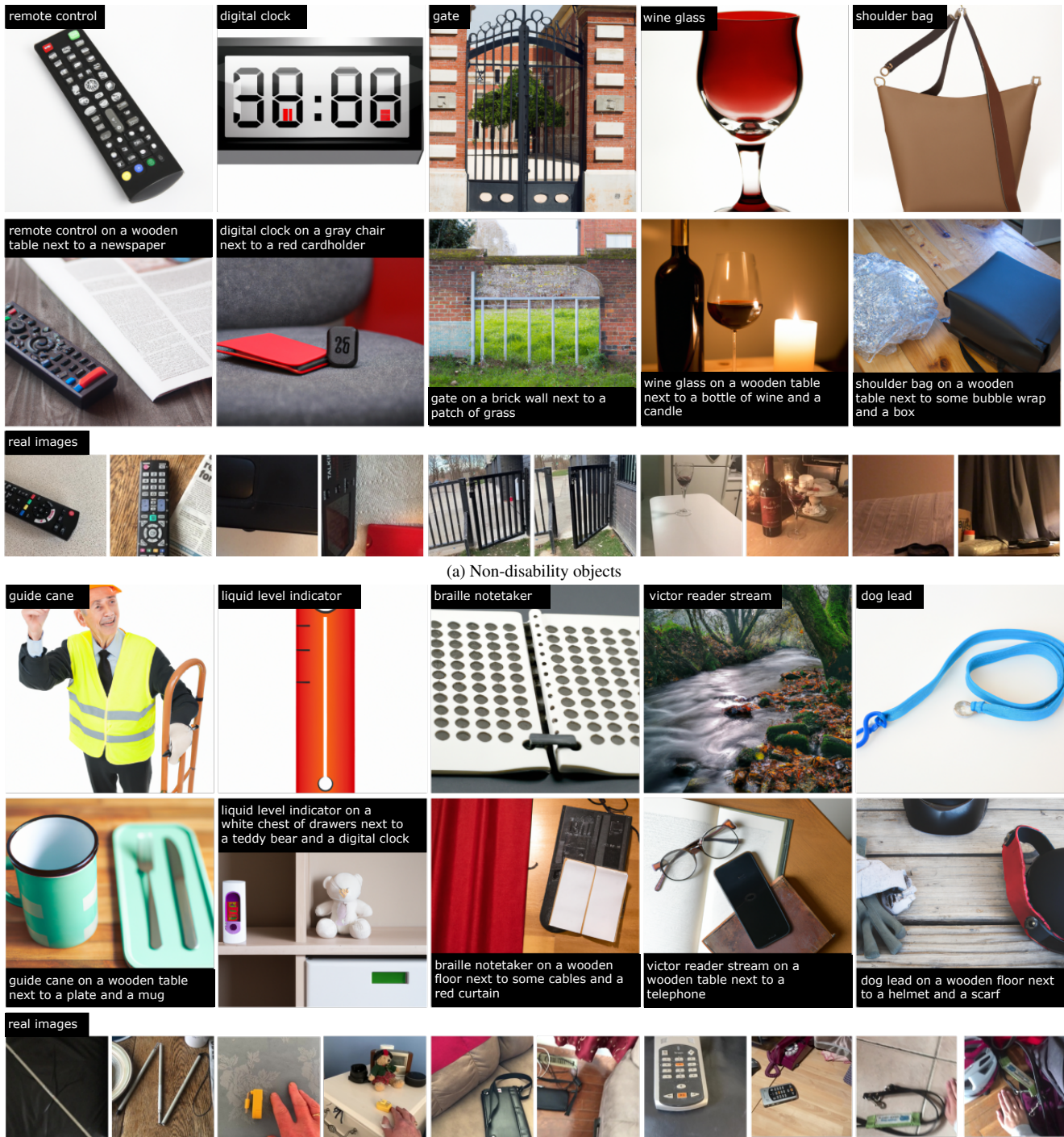


Figure C.4. DALL-E2 generates high-quality images of non-disability objects (a), but defaults to more common objects or fabrications for disability objects (b). For each sub-figure, the top row shows generations for a simple prompt containing just the object name, while the second row shows generations for the richer prompt where a surface and adjacent objects are also specified. The bottom row shows real images of each object.