A. Additional Quantitative Results

In this section, we present additional quantitative evaluations for better understanding our experimental outcomes. We first extend our bottom-up grouping experiments nad robustness analysis, followed by additional ablations on aggregation and pretraining strategies.

A.1. Bottom-up Grouping

We run more experiments on bottom-up grouping, following the same evaluation protocol from [10,91] to we compute Jaccard Similarity (JS). JS here measures the average intersection over the union across all segmentation instances regardless of object category. Our results in Tabs. 5 and 6 demonstrate competitive performance by CLIPpy. On two more challenging datasets we note how the model drops in performance, perhaps indicative of more visually cluttered scenes (Tab. 6). We take these results to indicate that CLIPpy perceptually groups semantically related content better than previous work, providing state-of-the-art results in unsupervised segmentation.

A.2. Robustness Analysis

We explore the open-source CLIP model from [68], evaluating its performance on the waterbirds dataset for signs of spurious correlations learned by the model. These results presented in Tab. 7 illustrate how classification performance for the same set of foreground objects drops drastically against different conflicting backgrounds, even in spite of the task involving fine-grained categories of birds.

A.3. Alternate Aggregation Strategies

Having explored two standard visual aggregation techniques as baselines, we ask how aggregation on the textual encoder affects CLIPpy performance. In detail, we first explore maximum pooling for the language modality. Thereafter, we draw further attention to the visual modality, exploring more complex pooling strategies.

Language Modality. In Tab. 8, we explore how pooling on the text embedding affects overall performance. We replace default average pooling with a maximum pooling operation to discover drops in performance across all metrics. This poor performance of maximum pooling based aggregation for language is consistent with prior works [37]. We hypothesize the reasoning as the nature of our task: while we attempt a localization across the visual modality, on the language end we reason with the entire text prompt as a single unit.

Visual Modality. We explored a range of alternate aggregation strategies that performed subpar to spatial max

Method	JS
IIC [42] MDC [9] PiCIE [15] STEGO [36]	6.4 7.1 12.3 21.0
CLIPpy (ours)	22.3

Table 5. **More bottom-up grouping:** CLIPpy achieves competitive Jaccard Similarity (JS) on the Cityscapes Dataset 27 class segmentation setup [36].

	Dataset	ADE20K	COCO
CLIP †	CC-12M	22.9	20.4
CLIPpy	CC-12M	28.9 (+6.0)	26.0 (+5.6)
CLIP †	HOTED 124M	24.2	21.6
CLIPpy	HQITP-134M	29.5 (+5.3)	27.2 (+5.6)

Table 6. More bottom-up grouping: CLIPpy improves Jaccard Similarity datasets.

CLIP [68]	water	land	Δ
waterbird	88.3	70.1	-18.2
landbird	90.5	99.1	-8.6

Table 7. Waterbirds evaluation for OpenAI CLIP model. We demonstrate how even a CLIP model trained on significantly more data (400M image-text pairs) contains stronger sensitivty to spurious correlations than CLIPpy trained with an order of magnitude less data (12M). The first two columns report top-1 classification accuracy (%) and the last column reports difference of diagonal and off-diagonal terms (delta). Higher spurious correlations increase the absolute value of delta.

pooling utilized in CLIPpy. Two noteworthy approaches include Text Similarity Pooling (TSP) and Weighted Maximum Pooling (WMP). In TSP, we measure the similarity of each spatial token (corresponding to different positions) and obtain a normalized distribution using a softmax operation (including a temperature for smoothing). We aggregate the visual modality using a weighted averaging operation where the similarity of each spatial location to the text is the weight. WMP follows the same idea but employs per-channel perlocation embedding values instead of similarity as weights for the aggregation operation. Tab. 10 shows results for TSP, WMP and max pooling. TSP performs poorly across all variations while WMP works better at lower temperatures. Given the common softmax operation, higher temperature values result in smoother weights for both cases, making WMP and TSP more similar to average pooling. In the case of WMP, lower temperatures make the operation similar to simple spatial maximum pooling, which is reflected in the improved results for lower temperature values.

A.4. Pretraining Ablations

Self-supervised pretraining of the vision head of CLIPpy leads to notable performance gains. We also explore how alternate supervised pretraining affects performance. In par-

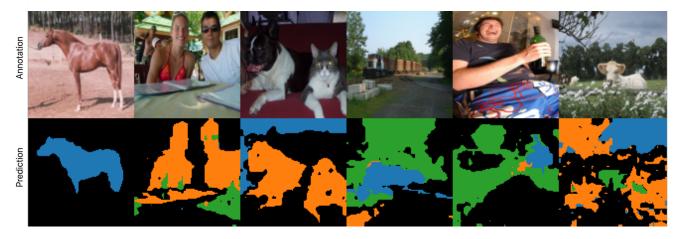


Figure 6. **Qualitative examples of bottom-up unsupervised segmentation with CLIPpy.** We illustrate examples from PASCAL VOC dataset with original image and CLIPpy prediction in the top and bottom rows, respectively. Note that colors correspond to clusters and *not* semantic labels.

Method	IN	VOC	COCO	ADE-20K
Avg	45.3	50.8	23.8	13.1
Max	42.0	42.6	18.9	10.5

Table 8. Alternate Pooling Strategies for Text Modality. Unlike the visual modality, performance drops when replacing the default average pooling (Avg) with maximum pooling (Max).

dataset	image	T5	ImageNet	Pascal VOC	
	init	init?	accuracy	mIoU	Jaccard
	DINO	✓	59.0	50.1	54.6
	IN-1K	/	63.6	21.7	40.2
HQITP	random	/	49.4	37.3	46.3
-134M	DINO		55.2	46.9	53.7
	IN-1K		60.6	20.9	39.1
	random		46.4	37.1	45.8

Table 9. Additional ablation studies with HQITP-134M. Parallel results for Table 4 (center) for ablations on weight initialization. We observe similar trends in results across these experiments.

ticular, we pretrain the image backbone using ImageNet-1K in a fully-supervised setting, and use the penultimate features to initialize CLIPpy. Apart from pretraining, all hyperparameters are unchanged. We present these results in Tab. 9. We observe when training with HQITP-134M that supervised pretraining leads to considerable performance gains for ImageNet-1K top-1 accuracy, but a considerable drop in segmentation performance across all three datasets. However, we note that supervised ImageNet pre-training provides unfair advantage in the case of ImageNet accuracy, in particular given the ability of overfitting those visual concepts. So we focus more on the segmentation results. Interestingly, entirely eliminating visual pre-training, while degrading segmentation performance, outperforms the ImageNet

supervised pre-training initialized model. This reaffirms our hypothesis of better generalization of self-supervised features for segmentation tasks. We also note that results on HQITP-134M indicate trends similar to those with CC-12M.

B. Qualitative examples of localization

In this section, we showcase additional qualitative examples of the success and failures of CLIPpy on bottom-up unsupervised segmentation and top-down semantic segmentation.

Fig. 6 shows examples of bottom-up unsupervised segmentation on PASCAL VOC for CLIPpy. The left three examples highlight the strength of the method for images with less clutter in which there is a single object of interest. The right three examples shows examples highlights the failures in the presence of scene clutter.

Fig. 7 present additional top-down semantic segmentation examples across all three datasets used for evaluation. The examples from PASCAL VOC dataset are the same examples from Fig. 6. For PASCAL VOC, note that the contrast between the top-down and bottom-up segmentations. The top-down segmentation is able to correctly separate the dog and cat classes (column 3) and also improve performance in the more cluttered scenes. On the other hand, in column 4, it missed out on a portion of the train that was segmented properly in the bottom-up setting. Segmentations from CLIPpy may contain discontinuities within a single object region in some cases, especially for the background objects (columns 2-3).

COCO and ADE-20K provide more challenging datasets and highlight several potential failure modes. A notable failure mode for CLIPpy is to reverting to the baseline CLIP behaviour by predicting the salient object class at all locations. This is visible to some extent in column 2 & 3 in the

Method	temp	IN	VOC	COCO	ADE-20K
TSP	0.1	0	2.98	0	0
TSP	1.0	20.15	4.91	1.25	0
TSP	10	24.70	15.83	4.27	2.29
Max		27.05	37.39	17.32	9.69

Method	temp	IN	VOC	COCO	ADE-20K
WMP	0.1	28.36	31.12	13.64	8.12
WMP	1.0	0	0	0	0
WMP	10	0	3.53	0	0
Max		27.05	37.39	17.32	9.69

Table 10. **Alternate Pooling Strategies for Visual Modality.** We report results for TSP and WMP with models trained for lesser steps on CC-12M with no initialization. Max refers to the spatial max operation used in CLIPpy. The temperature of the softmax operation for TSP and WMP is indicated by *temp*.

COCO examples and column 3 in the ADE-20K examples. Additionally, CLIPpy results in false positives for cluttered scenes as visible in some examples of these two datasets.

In summary, CLIPpy is able to localize the salient objects well in less cluttered scenes, even when multiple objects belonging to different classes are present. CLIPpy is also able to coarsely localize some of the salient objects in more cluttered scenes. In terms of limitations, CLIPpy fails to correctly localize some background classes, fails to correctly recognize object boundaries, and may miss smaller objects in cluttered scenes.

C. Details of HQITP-134M dataset

High Quality Image Text Pairs (HQITP-134M) consists of \sim 134 million diverse and high quality images paired with descriptive captions and titles. Images range in spatial resolution from 320 to 2048 pixels on the short side. All images are JPEG format and most are RGB. Each example image is associated with a title, and a list of several captions. A small fraction (\ll 1%) of the examples are missing both captions and title. We favor the associated captions, and find that these tokenize to an average length of 20.1 tokens, although the shortest caption is only one token and the longest is over 1000. This dataset was licensed to our research lab by a third party for commercial use.

To preprocess this dataset for training, we first exclude all pairs for which no valid caption exists. We also perform global exact-byte-match image de-duplication across our full training corpus, meaning that valid examples may be dropped due to appearing in other subsets of our overall training dataset. As we draw each example, we create an image-text pair by sampling from the list of available captions. The text is then tokenized, and the image is resized so that the shortest side is 224 pixels, with a further random crop then applied over the longer dimension to produce a 224 x 224 pixel square image. Lastly, we normalize the image using statistics derived from our full training corpus.

D. Detailed Methodology

In this section, we provide additional details regarding our methodology - in particular the baseline we build off and the modifications we propose.

D.1. Architecture and Training

We provide a quick overview of our architecture in main paper Fig. 2. Consider a batch size N, spatial height H, spatial width W, and depth D. X is a tensor that has a shape of [N, H, W, D] and is the output of an image encoder. Y is a tensor that is shape [N, D] and is the output of a text encoder.

Language Model. We employ a strong language model baseline derived from the transformer architecture [83] and implemented in T5 [69]. T5 models use an encoder-decoder architecture that is trained using a generative span corruption task, and have achieved state-of-the-art on a broad range of NLP tasks including GLUE [85] and Super-Glue [84]. We use the encoder only and discard the decoder part. We employ the T5-base which consists of 12 transformer layers, 12 attention heads, and 768 token channel dimensions.

Image Model. We explore two architectures for image featurization, CNN-based and Vision-Transformers, although we focus the majority of work on the latter. First, we employ the EfficientNet architecture [79] as a high performant CNN architecture, which has been used previously in vision-language models. The specifics of the meta-architecture were derived from considerations based on neural architecture search. Second, we employ the Vision Transformer (ViT) architecture [24]. After some preliminary experiments exploring various ViT architectures across patch sizes and model sizes, we report most of our results based on ViT-B/16 model. This model employs patch sizes of 16×16 pixels. We experiment with this model at two spatial resolutions of 224×224 and 448×448, resulting in 196 and 784 tokens, respectively. We refer the reader to [24, 83] for details. Briefly, ViT is largely inherited from the NLP literature and consists of a hierarchical associative memory. Each layer, termed a transformer, is composed of a Multiheaded Self-Attention (MSA) layer followed by a 2-layer feed-forward multi-layer perceptron (MLP). The primary parameter of ViT is the patch size P specifying the $P \times P$ patch of pixels constituting a token in the architecture.

Contrastive Representation Learning. Let x_i and y_i denote the image and text embeddings (post aggregation) of

the i'th example in the batch. A contrastive loss may be specified as the cross entropy across a batch [43,68]. The cross entropy is calculated between a one-hot encoding specifying the correspondence between the image and text examples, and a softmax-normalized distribution specifying the dot-product similarity between image and text embeddings.

$$L = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp(\boldsymbol{x}_{i}^{\top} \boldsymbol{y}_{i} / \tau)}{\sum_{j=1}^{N} \exp(\boldsymbol{x}_{i}^{\top} \boldsymbol{y}_{j} / \tau)} + \underbrace{-\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp(\boldsymbol{y}_{i}^{\top} \boldsymbol{x}_{i} / \tau)}{\sum_{j=1}^{N} \exp(\boldsymbol{y}_{i}^{\top} \boldsymbol{x}_{j} / \tau)}}_{\text{text-to-image}}$$

The normalization for the image-to-text and text-to-image similarity is computed by summing over the potential matches (indexed by j) to the text and image examples within a batch, respectively. Note that τ is the temperature of the softmax for the normalization.

D.2. Aggregation

We further explore the aggregation techniques explored in this main paper. In the context of contrastive vision language models, **Average pooling** across space has been adopted for CNN-based architectures such as [43]. Borrowed from language modeling [22], the **class token (CLS)**, which is prepended to the image patch tokens [24] to learn an embedding that aggregates information across all patch tokens in order to predict the image label, is als used in ViT based contrastive vision language models [68]. Subsequent work in vision-language models has explored learning pooling strategies [12, 93], heuristically selecting a set of similar neighbors [96] or learning attention-based mechanisms [95].

Our experimentation shows the clear benefits in terms of grouping that maximum pooling provides. We now question why a simple modification in token aggregation results in such behavor. We note how the max pooling operation allows pre-aggregation features (shaped [N, H, W, D]) to determine the spatial location for gradient updates at each step, conditioned on input images. Across different images containing a common object at different spatial locations, the model has to select a conservative and minimal set of spatial locations for gradient updates. At the same time, given the cross-modal contrastive train objective, the aggregated feature of each such image must be aligned towards a common language concept (i.e. related to the common object). We hypothesize that gradient updates at the common object's spatial location is the simplest optimization for the train objective in this case, leading to observed perceptual grouping.

D.3. Pretraining

While we described the pre-training strategies and their effect of emerging grouping behavior in these models, we discuss the suitability of different visual pre-training strategies in terms of inducing such observed grouping. In detail, we compare two types of pre-training: image-level objectives and region-level objectives. The visual encoder rep-

resentation space can be viewed as containing per-image features (post-aggregation) vs per-spatial location features (pre-aggregation). We hypothesize that semantics tied boundaries of this representation space should operate at the latter granularity to induce perceptual grouping. Furthermore, we suggest that initializations facilitating the former will detriment grouping behaviour. In particular, visual pre-training strategies separating image-level representations by semantics (e.g. supervised ImageNet pre-training) will diminish perceptual grouping. Self-supervised pre-training strategies focused on more granular within image representations (e.g. [10]) will tend to enhance perceptual grouping. This hypothesis is empirically validated in ablations (see Table 4).

D.4. Visual Token Sub-Sampling

Motivated by vision transformers' ability to process sequences of length different to train time, we generate higher resolution segmentations during inference by sampling more image patches. In order to increase robustness to such varying resolution, we utilize up to $2\times$ higher resolution images during training but randomly drop 80% of visual tokens to minimize additional compute overhead. While improving inference quality, this also provides training stability possibly due to its regularization effect (see App. E for more details).

D.5. Inference

CLIPpy mainly operates under 3 distinct settings during inference: a) classification, b) bottom-up grouping, c) topdown grouping. On the visual modality, the first utilizes a spatially aggregated single per-image token while the latter two utilize a set of per-region tokens. Classification follows zero-shot analyses from [68] where the model is prompted at inference for a selection of labels (App. I for prompts). Bottom-up grouping follows a form of spectral clustering inspired by [10]. Namely, we compute PCA across the spatial map of image features. Each principal component corresponds to a candidate group and we cluster the proximity of each feature vector to these components. Top-down grouping employs zero shot analysis at each spatial location token given a set of label prompts. This is similar to the ViT-CAM method [30] and generates predictions across space exploiting the transitive property of our aggregation operations.

E. Architecture and Training Details

Experimental Setup. We train our models on two datasets: Conceptual Captions 12M (CC-12M) [11] and High Quality Image Text Pairs (HQITP-134M) consisting of 12 million and 134 million image-text pairs, respectively (App. C for details). For both datasets, text is tokenized, and images resized and center cropped to 224×224 pixels. We report results on EfficientNet-B5 employed by ALIGN [43], and ViT-B/16 employed by CLIP [68] although we focus more on the latter. We train models on 32 GPUs across

4 machines with PyTorch [66]. We evaluate across image classification, localization, and robustness tasks. For image classification, we employ the validation splits of ImageNet [20] and ImageNet-v2 [70], and for robustness we employ the test split of Waterbirds [73]. These datasets contain 1000, 1000, and 2 classes respectively. For segmentation tasks, we employ the validation splits of PASCAL VOC [27], ADE20K [14, 102], COCO [54], COCO (Obj) [54], and Cityscapes [17]. Each of these datasets contain 20, 150, 133, 80, and 27 labels, respectively.

Our implementations of CLIP and ALIGN employ ViT-B/16 [24], and EfficientNet-B5 [79] for the image embedding, respectively, to mirror the primary results presented in each respective vision-language model. For the ViT architecture, we experimented with varying patch sizes P=8 and P=16 in order to leverage open-sourced DINO pretrained weights [10], but report all of our results with P=16.

We train all models on 224×224 images to provide a fair comparison with [68]. Note however that the published version of ALIGN employed a 640×640 resolution. ViT models may operate on images at arbitrary spatial resolution. At inference time we experimented with spatial resolutions of 224×224 and 448×448 , resulting in 196 and 784 tokens, respectively. Results were similar across 224 and 448 resolutions, we report only results at 224 for brevity (except for Fig. 5).

We train models with 32 GPUs across 4 machines with PyTorch [66] using the LAMB optimizer [94] with a cosine decayed learning rate with linear warm-up. We employ an initial learning rate of 1e-3, 2000 warm-up steps, and decay the rate with a single period over the entire training regime (32 epochs for CC12M; 10 epochs for HQITP-134M). We employ a weight decay of 1e-2. All training parameters were determined through moderate hyperparameter tuning.

Patch Sub-Sampling: During training, we also utilize overlapping patch generation with patch sub-sampling as regularization. In particular, the token sequence length is always maintained at the original value of 224 during training through random sub-sampling (patches selected according to uniform distribution). This also allows obtaining a higher resolution feature map from a fixed resolution image during inference.

F. Limitations of CNN-based architectures

The primary focus of our presentation is on the Vision Transformer (ViT) architecture [24]. The reason for this focus is that the Transformer architecture is particularly well suited for multimodal learning tasks because one does not need to craft an architecture for each modality, and tune the training set up for each particular architecture. We also recognize that convolutional neural networks (CNN's) have a long history of providing state-of-the-art CNN results on computer vision related problems. EfficientNet-B5 is a mod-

ern state-of-the-art architecture whose meta-architecture and scaling properties were derived from architecture search considerations [79] (but see [5]), and provides the visual backbone for the ALIGN model [43].

We experimented with this model and find that the image featurization from a CNN-derived backbone achieves favorable results with respect to previously published ViT models on localization problems. Table 3 showcases higher mIoU for semantic segmentation than a model with a ViT backbone (CLIP) on ADE20K, COCO and Pascal VOC when trained on the same dataset. Likewise, previous results published on ALIGN with an EfficientNet-B5 backbone achieved superior results to a ViT model on COCO and ADE20K in terms of semantic segmentation ¹.

Most importantly, in spite of many attempts, we were not able to improve the performance of the EfficientNet-B5 architecture for localization. The best results for a CNNbased architecture we achieved are shown in Table 3, which are notably below previously published results and our best results with a ViT architecture. At best, the addition of maximum pooling at the top layer of a CNN led to marginal gains in terms of mIoU or Jaccard similarity. We suspect that a custom architecture (such as ASPP or FPN) may improve these results further [13,53], but we consider this out of scope as we are attempting to learn a featurization that does not artificially increase the parameterization in order to solve a specialized task of localization. We suspect that this limitation of CNN architectures may reflect the fact that CNN architecture already have learned a representation that is heavily dependent on the spatial geometry derived from a convolutional kernel. Such an inductive bias may be not provide a suitable mechanism for providing global processing in a segmentation task [86].

G. More on Related Work

G.1. Vision-language models for grounding.

Contrastive language image pre-training [68] (CLIP) led to a range of follow up work performing open-vocabulary detection [25,33,44,49,50,98] or segmentation [31,48,100]. While these methods leverage dense human annotations for training, an alternate line of works [19,91–93,103] attempt to learn alignment between regions of images and language with only image level noisy captions for supervision. Their weak supervision allows better scalability (to more data) leading to learning more generic and transferable representations. In fact, multiple such works [23,48,91,92,103] perform zero-shot semantic segmentation. However, un-

 $^{^{1}}$ We note that the previously published ALIGN results were based on a backbone trained with a much larger dataset (1.8B image-text pairs), and it was evaluated at a higher resolution of 640×640 pixels, resulting in a zeroshot image recognition performance of 76.4. In comparison, our baseline and proposed models operate at 224×224 resolution and was trained on a $10\times$ smaller dataset.

	segment label?	segment mask?	ADE20K	COCO	PASCAL VOC
SPNet [90]	✓	✓			18.3
ZS3Net [7]	✓	✓			38.3
LSeg [48]	✓	✓		27.2	47.4
LSeg+ [30]	✓	✓	18.0	55.1 [†]	59.0
ALIGN w/ proposal [30]		✓	12.9	17.9 [†]	22.4
OpenSeg [30]		✓	21.1	36.1 [†]	70.2
OpenSeg + Narr. [30]		✓	24.8	38.1 [†]	72.2

Table 11. Performance of prior zeroshot segmentation models trained on segmentation data. All numbers report the mIoU for semantic segmentation on ADE20K (150 labels), PASCAL-VOC (20 labels) and COCO (50 labels). † indicates models that were trained on image segmentation masks from the COCO dataset.

like [91,92] geared to segment a fixed count of foreground objects, our proposed CLIPpy can better segment arbitrary object counts and background classes. In contrast to [103] using generic image level features, CLIPpy explicitly learns local features during training. Moreover, CLIPpy requires no dense human annotations or task-specific fine-tuning in contrast to [23,48]. We also highlight how [23,91,92] perform grouping independent of language during inference - however CLIPpy can group conditioned on language, allowing to capture variable object boundaries for different language prompts.

G.2. Zero-shot semantic segmentation.

A form of top-down grouping, this relatively new task [3,7,40,41,51,65,75,90,101] attempts to segment unseen classes, usually after a supervised training phase often involving dense annotation based supervision. Following two early representative works [7,90], most later approaches [34,35,45,51,75,81] formulate the task as a pixel-level zero-shot classification problem with a closed set vocabulary. While CLIPpy follows a similar pixel based formulation, in contrast, our method requires no dense human annotations for supervision, no task specific fine-tuning, and is openvocabulary. Recent work [23,49] also explores region-level classification leveraging pre-trained CLIP models [68], but unlike CLIPpy perform grouping independent of language during inference.

G.3. Unsupervised segmentation

. Analogous to bottom-up grouping, these works perform class-agnostic segmentation within the visual modality with no explicit language alignment [10,28,36,42,61]. This topic has a long, rich history in human visual perception [89] and computer vision [58], and has been explored as means of generalizing to new visual domains [59,67]. It is this goal that most closely inspires our work. Early efforts group pixels based on known spatially-local affinities [16,71,76], with subsequent methods leading to region proposal networks for object detection [82] and advances in semantic segmentation [1]. Recent methods employ self-supervision to learn perceptual grouping [15,36] or object-centric groupings [4,26,39,56,88]. Our proposed CLIPpy demonstrates

competitive performance, but additionally aligns groups to the language modality explicitly.

G.4. Learning robust visual representations

. For a long time, ImageNet [20] accuracy was believed to provide a reasonable proxy for quality of learned visual representations [32,47]. However, recent work highlights notable deficiencies in such learned representations [29, 46, 70] including sensitivity to low level textures, failure for domain shifts, and reliance on spurious correlations. These failures inspired a large literature to mitigate learning spurious correlations [2,55,73] by focusing on new optimization techniques. Progress on this issue may address parallel issues in fairness [18]. Resulting methods have largely focused on synthetic data, re-balancing data, and shaping learned embeddings [55, 62]. Nonetheless, theoretical results suggest pessimistic bounds unless additional structure informs the problem (see refs. in [73]). Therein, the structured output predictions of proposed CLIPpy provide another promising solution.

G.5. Prior work on zero shot semantic segmentation

H. Details of robustness analysis

We calculate the zero-shot prediction of the class in a non-standard manner to exploit the spatial reasoning of CLIPpy. We apply the same zero-shot evaluation to the baseline CLIP model. Specifically, we first calculate the embedding for all labels within each category of waterbird, landbird and background. For each of these categories we calculate the average image embeddings across these labels at each spatial location.

To exploit the spatial knowledge of the model, we focus our analysis on all spatial locations which are *not* labeled as background. For all locations which maximally predict a waterbird, we calculate the similarity to the embedding for a waterbird. Likewise, we do the same for all locations maximized by landbird. Finally, our resulting prediction is the class that is closest to its associated embedding.

We find that these results and the corresponding robustness vary substantially due to the selection of prompts for each of the three categories. This matches observations in [68]. In Sec. 3.4 we focus on the results of the model trained on CC-12M using the prompts listed below which contain a minimal amount of prompt engineering. In particular, the prompts for waterbirds and landbirds follow [73]. See also [87].

- background: background
- waterbird: Black footed Albatross, Laysan Albatross, Sooty Albatross, Crested Auklet, Least Auklet, Parakeet Auklet, Rhinoceros Auklet, Brandt Cormorant, Red faced Cormorant, Pelagic Cormorant, Frigatebird, Northern Fulmar, Gadwall, Eared Grebe, Horned Grebe, Pied billed Grebe, Western Grebe, Pigeon Guillemot, California Gull, Glaucous winged Gull, Heermann Gull, Herring Gull, Ivory Gull, Ring billed Gull, Slaty backed Gull, Western Gull, Long tailed Jaeger, Pomarine Jaeger, Red legged Kittiwake, Pacific Loon, Mallard, Hooded Merganser, Red breasted Merganser, Brown Pelican, White Pelican, Horned Puffin, Artic Tern, Black Tern, Caspian Tern, Common Tern, Elegant Tern, Forsters Tern, Least Tern
- landbird: Groove billed Ani, Brewer Blackbird, Red winged Blackbird, Rusty Blackbird, Yellow headed Blackbird, Bobolink, Indigo Bunting, Lazuli Bunting, Painted Bunting, Cardinal, Spotted Catbird, Gray Catbird, Yellow breasted Chat, Eastern Towhee, Chuck will Widow, Bronzed Cowbird, Shiny Cowbird, Brown Creeper, American Crow, Fish Crow, Black billed Cuckoo, Mangrove Cuckoo, Yellow billed Cuckoo, Gray crowned Rosy Finch, Purple Finch, Northern Flicker, Acadian Flycatcher, Great Crested Flycatcher, Least Flycatcher, Olive sided Flycatcher, Scissor tailed Flycatcher, Vermilion Flycatcher, Yellow bellied Flycatcher, American Goldfinch, European Goldfinch, Boat tailed Grackle, Blue Grosbeak, Evening Grosbeak, Pine Grosbeak, Rose breasted Grosbeak, Anna Hummingbird, Ruby throated Hummingbird, Rufous Hummingbird, Green Violetear, Blue Jay, Florida Jay, Green Jay, Dark eyed Junco, Tropical Kingbird, Gray Kingbird, Belted Kingfisher, Green Kingfisher, Pied Kingfisher, Ringed Kingfisher, White breasted Kingfisher, Horned Lark, Western Meadowlark, Mockingbird, Nighthawk, Clark Nutcracker, White breasted Nuthatch, Baltimore Oriole, Hooded Oriole, Orchard Oriole, Scott Oriole, Ovenbird, Western Wood Pewee, Sayornis, American Pipit, Whip poor Will, Common Raven, White necked Raven, American Redstart, Geococcyx, Loggerhead Shrike, Great Grey Shrike, Baird Sparrow, Black throated Sparrow, Brewer Sparrow, Chipping Sparrow, Clay colored Sparrow, House Sparrow, Field Sparrow, Fox Sparrow, Grasshopper Sparrow, Harris Sparrow, Henslow Sparrow, Le Conte Sparrow, Lincoln Sparrow, Nelson Sharp tailed Sparrow, Savannah Sparrow, Seaside Sparrow, Song Sparrow, Tree Sparrow, Vesper Sparrow, White crowned Sparrow, White throated Sparrow, Cape Glossy Starling, Bank Swallow, Barn Swallow, Cliff Swallow, Tree Swallow, Scarlet Tanager, Summer Tanager, Green tailed Towhee, Brown Thrasher, Sage Thrasher, Black capped Vireo, Blue headed Vireo, Philadelphia Vireo, Red eyed Vireo, Warbling Vireo, White eyed Vireo, Yellow throated Vireo, Bay breasted Warbler, Black and white Warbler, Black throated Blue Warbler, Blue winged Warbler, Canada Warbler, Cape May Warbler, Cerulean Warbler, Chestnut sided Warbler, Golden winged Warbler, Hooded Warbler, Kentucky Warbler, Magnolia

Warbler, Mourning Warbler, Myrtle Warbler, Nashville Warbler, Orange crowned Warbler, Palm Warbler, Pine Warbler, Prairie Warbler, Prothonotary Warbler, Swainson Warbler, Tennessee Warbler, Wilson Warbler, Worm eating Warbler, Yellow Warbler, Northern Waterthrush, Louisiana Waterthrush, Bohemian Waxwing, Cedar Waxwing, American Three toed Woodpecker, Pileated Woodpecker, Red bellied Woodpecker, Red cockaded Woodpecker, Red headed Woodpecker, Downy Woodpecker, Bewick Wren, Cactus Wren, Carolina Wren, House Wren, Marsh Wren, Rock Wren, Winter Wren, Common Yellowthroat

I. Prompts for Zero-shot Segmentation

We employed the following prompts for probing our vision-language models for zero-shot semantic segmentation. These prompts were copied from the corresponding label sets of each dataset with some basic considerations, for instance, restoring spaces in compound words. For prompts separated by a comma, the average embedding across all prompts delineated by commas is associated with each label.

Pascal VOC 2012 [27]

- background: background, crops, bush, shrub, tiles, pavement, rug, carpet, box, boxes, speaker, storage, painting, board, panel, poster, clock, cage, drinking glass, park, plaything, toy, fireplace, bag, bag, bed, bench, book, books, building, buildings, cabinet, drawer, ceiling, computer, computer case, cup, cups, door, fence, floor, flower, grass, lawn, turf, ground, soil, dirt, tiles, keyboard, lamp, mountain, hills, mouse, curtain, platform, sign, street, rock, stone, shelf, sidewalk, sky, clouds, snow, track, train track, tree, trees, wall, water, window, wood, woods
- 2. aeroplane: aeroplane, airplane, aeroplanes, airplanes
- 3. bicycle: bicycle, bicycles, bike, bikes
- 4. bird: bird, birds
- 5. boat: boat, boats
- 6. bottle: bottle, bottles, water bottle
- 7. bus: bus, buses
- 8. car: car, cars
- 9. cat: cat, cats, kitties, kitty
- 10. chair: chair, chairs
- 11. cow: cow, cows, calf
- 12. diningtable: diningtable, dining table, diningtables, dining tables, plate, plates
- 13. dog: dog, dogs, puppy, puppies
- 14. horse: horse, horses, foal
- 15. motorbike: motorbike, motorcycle, motorbikes, motorcycles
- person: person, child, girl, boy, woman, man, people, children, girls, boys, women, men, lady, guy, ladies, guys, clothes
- 17. pottedplant: pottedplant, pottedplants, plant pot, plant pots, planter, planters, potted plant
- 18. sheep: sheep
- 19. sofa: sofa, sofas
- 20. train: train, trains, locomotive, locomotives, freight train
- tvmonitor: tvmonitor, monitor, tv, televison, television monitor

COCO 2017 [54]

1. airplane: airplane

- 2. apple: apple
- 3. backpack: backpack
- 4. banana: banana
- 5. baseball bat: baseball bat
- 6. baseball glove: baseball glove
- 7. bear: bear
- 8. bed: bed
- 9. bench: bench
- 10. bicycle: bicycle
- 11. bird: bird
- 12. boat: boat
- 13. book: book
- 14. bottle: bottle
- 15. bowl: bowl
- 16. broccoli: broccoli
- 17. bus: bus
- 18. cake: cake
- 19. car: car
- 20. carrot: carrot
- 21. cat: cat
- 22. cell phone: cell phone
- 23. chair: chair
- 24. clock: clock
- 25. couch: couch
- 26. cow: cow
- 27. cup: cup
- 28. dining table: dining table
- 29. dog: dog
- 30. donut: donut
- 31. elephant: elephant
- 32. fire hydrant: fire hydrant
- 33. fork: fork
- 34. frisbee: frisbee
- 35. giraffe: giraffe
- 36. hair drier: hair drier
- 37. handbag: handbag
- 38. horse: horse
- 39. hot dog: hot dog
- 40. keyboard: keyboard
- 41. kite: kite
- 42. knife: knife
- 43. laptop: laptop
- 44. microwave: microwave
- 45. motorcycle: motorcycle
- 46. mouse: mouse
- 47. orange: orange
- 48. oven: oven
- 49. parking meter: parking meter
- 50. person: person
- 51. pizza: pizza
- 52. potted plant: potted plant
- 53. refrigerator: refrigerator
- 54. remote: remote
- 55. sandwich: sandwich
- 56. scissors: scissors
- 57. sheep: sheep
- 58. sink: sink
- 59. skateboard: skateboard

- 60. skis: skis
- 61. snowboard: snowboard
- 62. spoon: spoon
- 63. sports ball: sports ball
- 64. stop sign: stop sign
- 65. suitcase: suitcase
- 66. surfboard: surfboard
- 67. teddy bear: teddy bear
- 68. tennis racket: tennis racket
- 69. tie: tie
- 70. toaster: toaster
- 71. toilet: toilet
- 72. toothbrush: toothbrush
- 73. traffic light: traffic light
- 74. train: train
- 75. truck: truck
- 76. tv: tv
- 77. umbrella: umbrella
- 78. vase: vase
- 79. wine glass: wine glass
- 80. zebra: zebra

[**ADE-20K** (150 frequent labels) [102]]

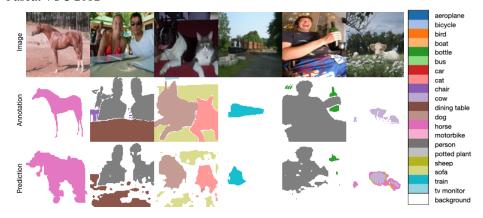
- 1. airplane: airplane, aeroplane, plane
- 2. animal: animal, animate, being, beast, brute, creature, fauna
- 3. apparel: apparel, wearing, apparel, dress, clothes
- 4. arcade: arcade, machine
- 5. armchair: armchair
- 6. ashcan: ashcan, trash, can, garbage, can, wastebin, ash, bin, ash-bin, ashbin, dustbin, trash, barrel, trash, bin
- 7. awning: awning, sunshade, sunblind
- 8. bag: bag
- 9. ball: ball
- bannister: bannister, banister, balustrade, balusters, handrail
- 11. bar: bar
- 12. barrel: barrel, cask
- 13. base: base, pedestal, stand
- 14. basket: basket, handbasket
- 15. bathtub: bathtub, bathing, tub, bath, tub
- 16. bed: bed
- 17. bench: bench
- 18. bicycle: bicycle, bike, wheel, cycle
- 19. blanket: blanket, cover
- 20. blind; blind, screen
- 21. boat: boat
- 22. book: book
- 23. bookcase: bookcase
- 24. booth: booth, cubicle, stall, kiosk
- 25. bottle: bottle
- 26. box: box
- 27. bridge: bridge, span
- 28. buffet: buffet, counter, sideboard
- 29. building: building, edifice
- 30. bulletin: bulletin, board, notice, board
- 31. bus: bus, autobus, coach, charabanc, double-decker, jitney, motorbus, motorcoach, omnibus, passenger, vehicle
- 32. cabinet: cabinet
- 33. canopy: canopy

- 34. car: car, auto, automobile, machine, motorcar
- 35. case: case, display, case, showcase, vitrine
- 36. ceiling: ceiling
- 37. chair: chair
- 38. chandelier: chandelier, pendant, pendent
- 39. chest: chest of drawers, chest, bureau, dresser
- 40. clock: clock
- 41. coffee: coffee, table, cocktail, table
- 42. column: column, pillar
- computer: computer, computing, machine, computing, device, data, processor, electronic, computer, information, processing, system
- 44. conveyer: conveyer, belt, conveyor, belt, conveyer, conveyor, transporter
- 45. counter: counter
- 46. countertop: countertop
- 47. cradle: cradle
- 48. crt: crt, screen
- 49. curtain: curtain, drape, drapery, mantle, pall
- 50. cushion: cushion
- 51. desk: desk
- 52. dirt: dirt, track
- 53. dishwasher: dishwasher, dish, washer, dishwashing, ma-
- 54. door: door, double, door
- 55. earth: earth, ground
- 56. escalator: escalator, moving, staircase, moving, stairway
- 57. fan: fan
- 58. fence: fence, fencing
- 59. field: field
- 60. fireplace: fireplace, hearth, open, fireplace
- 61. flag: flag
- 62. floor: floor, flooring
- 63. flower: flower
- 64. food: food, solid, food
- 65. fountain: fountain
- 66. glass: glass, drinking, glass
- 67. grandstand: grandstand, covered, stand
- 68. grass: grass
- 69. hill: hill
- 70. hood: hood, exhaust, hood
- 71. house: house
- 72. hovel: hovel, hut, hutch, shack, shanty
- 73. kitchen: kitchen, island
- 74. lake: lake
- 75. lamp: lamp
- 76. land: land, ground, soil
- 77. light: light, light, source
- 78. microwave: microwave, microwave, oven
- 79. minibike: minibike, motorbike
- 80. mirror: mirror
- ${\bf 81.\ monitor;\ monitor,\ monitoring,\ device}$
- 82. mountain: mountain, mount
- 83. ottoman; ottoman, pouf, pouffe, puff, hassock
- 84. oven: oven
- 85. painting: painting, picture
- 86. palm: palm, palm, tree
- 87. path: path

- 88. person: person, individual, someone, somebody, mortal, soul
- 89. pier: pier, wharf, wharfage, dock
- 90. pillow: pillow
- 91. plant: plant, flora, plant, life
- 92. plate: plate
- 93. plaything: plaything, toy
- 94. pole: pole
- 95. pool: pool, table, billiard, table, snooker, table
- 96. poster: poster, posting, placard, notice, bill, card
- 97. pot: pot, flowerpot
- 98. radiator: radiator
- 99. railing: railing, rail
- 100. refrigerator: refrigerator, icebox
- 101. river: river
- 102. road: road, route
- 103. rock: rock, stone
- 104. rug: rug, carpet, carpeting
- 105. runway: runway
- 106. sand: sand
- 107. sconce: sconce
- 108. screen: screen, door, screen
- 109. screen: screen, silver, screen, projection, screen
- 110. sculpture: sculpture
- 111. sea: sea
- 112. seat: seat
- 113. shelf: shelf
- 114. ship: ship
- 115. shower: shower
- 116. sidewalk: sidewalk, pavement
- 117. signboard: signboard, sign
- 118. sink: sink
- 119. sky: sky
- 120. skyscraper: skyscraper
- 121. sofa: sofa, couch, lounge
- 122. stage: stage
- 123. stairs: stairs, steps
- 124. stairway: stairway, staircase
- 125. step: step, stair
- 126. stool: stool
- 127. stove: stove, kitchen, stove, range, kitchen, range, cooking, stove
- 128. streetlight: streetlight, street, lamp
- 129. swimming: swimming, pool, swimming, bath, natatorium
- 130. swivel: swivel, chair
- 131. table: table
- 132. tank: tank, storage, tank
- 133. television: television, television, receiver, television, set, tv, tv, set, idiot, box, boob, tube, telly, goggle, box
- 134. tent: tent, collapsible, shelter
- $135. \ \, \texttt{toilet:} \ \, \texttt{toilet,} \ \, \texttt{can,} \ \, \texttt{commode,} \ \, \texttt{crapper,} \ \, \texttt{pot,} \ \, \texttt{potty,} \ \, \texttt{stool,} \ \, \texttt{throne}$
- 136. towel: towel
- 137. tower: tower
- 138. trade: trade, name, brand, name, brand, marque
- 139. traffic: traffic, light, traffic, signal, stoplight
- 140. tray: tray
- 141. tree: tree
- 142. truck: truck, motortruck
- 143. van: van

- 144. vase: vase
- 145. wall: wall
- 146. wardrobe: wardrobe, closet, press
- 147. washer: washer, automatic, washer, washing, machine
- 148. water: water
- 149. waterfall: waterfall, falls
- 150. windowpane: windowpane, window

Pascal VOC 2012



COCO



ADE-20K

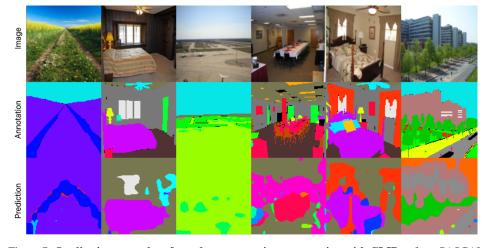


Figure 7. **Qualitative examples of top-down semantic segmentation with CLIPpy** from PASCAL VOC, COCO and ADE-20K. For Pascal VOC, we supply a color legend for the 20 label classes. Note that the Pascal VOC examples correspond to the same examples from Fig. 6.