

Generalized Decoding for Pixel, Image, and Language

Supplementary Material

Xueyan Zou^{*§}, Zi-Yi Dou^{*#}, Jianwei Yang^{*‡♣}, Zhe Gan[†], Linjie Li[†], Chunyuan Li[‡], Xiyang Dai[†], Harkirat Behl[‡]
Jianfeng Wang[†], Lu Yuan[†], Nanyun Peng[#], Lijuan Wang[†], Yong Jae Lee^{¶§}, Jianfeng Gao^{¶‡}

[§] University of Wisconsin-Madison [#] UCLA [‡] Microsoft Research at Redmond [†] Microsoft Cloud & AI

^{*}Equal Technical Contribution [¶]Equal Advisory Contribution [♣]Project Lead

{xueyan,yongjaelee}@cs.wisc.edu {zdou,violetpeng}@cs.ucla.edu {jianwyan,jfgao,zhgan,linjli,chunyl,jianfw,luyuan,lijuanw,hbehl,xidai}@microsoft.com

A. Experiment Settings

A.1. Pretraining

During pretraining, we set a minibatch for segmentation to 32 and image-text pairs to 1024. The image resolution is set to 1024 for segmentation and 224 for image-text data respectively. We follow a similar balanced sampling strategy in [11] to ensure the segmentation data are always observed for a consistent number of epochs, regardless of the total number of image-text pairs. Based on this, we pretrain all models for 50 epochs using AdamW [7] as the optimizer.

Also, in the main paper, all the pre-trained models are trained with 50 epochs of COCO data and roughly 45 epochs of 10 million image-text pairs. And 32 GPUs are used for pre-training with 40-50 training hours for Focal-T model. The AdamW optimizer is used in pretraining with the initial learning rate $1e-4$. A step-wise scheduler is used to decay the learning rate by 0.1 on the fraction [0.88889, 0.96296] of training steps.

A.2. Finetuning

Image-Text Retrieval. For both COCO and Flickr30k image-text retrieval, we finetune the models for 10 epochs using AdamW as the optimizer. We set the image resolution to 384 and the batch size to 2048. The learning rates are $3e-5$ for the X-Decoder part and $3e-6$ for the vision and language backbones.

Image Captioning. Similar to image-text retrieval, we finetune the captioning models for 10 epochs using AdamW as the optimizer. We set the image resolution to 480 and the batch size to 256. The learning rates are $2e-5$ for the X-Decoder part and $2e-6$ for the vision and language backbones. We use beam search during caption generation with the beam size set to 5. We do not use CIDEr optimization for our captioning models.

VQA. For VQA, we add a new classification layer on the top of the model and finetune the models for 10 epochs using AdamW as the optimizer. We set the image resolution

to 640 and the batch size to 256. The learning rates are $1e-4$ for the X-Decoder part, $1e-5$ for the vision and language backbones, and $1e-3$ for the VQA classification layer.

Generic Segmentation. For generic segmentation, we finetune the pretrained checkpoint with 24 epochs with start learning rate $1e-4$. We decay the learning rate by factor 10 at epoch 21 and 23, respectively. The batch size of ADE20k is 64, and 32 for COCO.

Referring Segmentation. For referring segmentation, we also finetune the pretrained checkpoint with 24 epochs. However, as RefCOCO has been used in pretraining, thus the initial learning rate is $1e-5$. It also decays twice at 21 and 23 epochs. We use a batch size of 64 during training. Further, in addition to the normal setting that multiple backbone and language encoder learning rates with 0.1, here we also multiply the transformer encoder learning rate by 0.1.

A.3. Dataset

As mentioned in the main paper, we have exclude the validation sets of Ref-COCOG UMD [13] and COCO Karpathy [12] from COCO2017 [5] in our training setting. This is because the training set of COCO2017 has high overlap with both Ref-COCOG UMD and COCO Karpathy as shown in Fig. 1 below. In addition, we didn't include other validation set of referring segmentation because each split has its own overlap with COCO2017 training set.

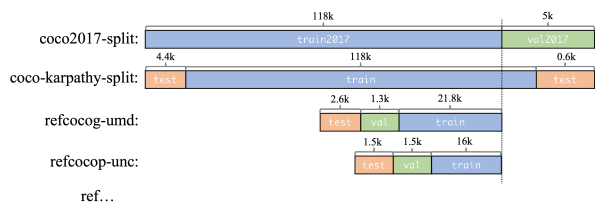


Figure 1. Visualization of dataset overlap. Note that refcocog-umd does also have overlap with refcocog-unc.

B. Open-Vocab Segmentation Benchmark

We propose an open vocabulary segmentation benchmark on 9 datasets with different evaluation metrics. The goal of this benchmark is to provide a comprehensive and standard evaluation protocol for open-vocabulary segmentation on different vocabulary sizes and image domains.

Dataset	Scene	Annotation Format			# Images	# Classes
		Sem	Inst	Pano		
ADE-150	common	✓	✓	✓	2000	150
ADE-847	common	✓	✗	✗	2000	847
Pascal Voc	common	✓	✗	✗	1449	20
Pascal Context-59	common	✓	✗	✗	5105	59
Pascal Context-459	common	✓	✗	✗	5105	459
SUN RGB-D	in-door	✓	✗	✗	5050	37
ScanNet-20	in-door	✓	✗	✓	5436	20
ScanNet-41	in-door	✓	✗	✗	5436	41
Cityscapes	driving	✓	✓	✓	500	19/8/19
BDD	driving	✓	✗	✓	1000	19//40

Table 1. Open-Vocabulary Segmentation Benchmark Statistics.

Table 1 shows the dataset statistics in the benchmark. It supports all generic segmentation tasks including semantic/instance/panoptic segmentation. It covers a variety of scopes ranging from 20 to 847 classes. In addition, the evaluation scene includes common objects, in-door scenes as well as autonomous driving scenarios. To enable a better understanding of the open-vocabulary ability on the training/evaluation datasets. We evaluate the coverage of training datasets captions and evaluation datasets concepts in Fig. 6-12 (we split the caption into single words and phrases to find mappings in categories). The major results of the open-vocabulary segmentation are evaluated in the main paper.

Method	COCO Karpathy				VQAv2 test-dev
	IR@1	TR@1	CIDEr	BLEU	
X-Decoder (T)	49.3	66.7	122.3	37.8	70.6
X-Decoder-VL (T)	44.3 $\downarrow_{5.0}$	60.3 $\downarrow_{6.4}$	113.2 $\downarrow_{9.1}$	34.8 $\downarrow_{3.0}$	69.4 $\downarrow_{1.2}$

Table 2. Compare finetuning result between X-Decoder and X-Decoder-VL which merely uses 4M image-text pairs for pretraining.

C. Extra Ablation Studies

C.1. Complementariness between Vision and VL

In our main paper, we observed that the vision-language pretraining objectives including image-text contrastive learning and image captioning have clear benefits to image segmentation, particularly in the zero-shot setting. Here, we further study the role of segmentation objectives in vision-language understanding. To investigate, we remove the segmentation data (COCO panoptic segmentation and referring segmentation) and only pretrain X-Decoder on the four million image-text pairs, denoted by X-Decoder-VL. Afterwards, we transfer the model to downstream VL tasks. As we can see from Table 2, the performance significantly

drops across all tasks after removing the segmentation data for pretraining. We suspect that segmentation data can help models to learn more fine-grained visual understanding and consequently benefit vision-language tasks. Along with our findings in the main paper, we conclude that *pixel-level segmentation and vision-language learning are complementary to each other for zero-shot and task-specific transfer.*

C.2. Model Architecture Inspection

In Table. 3, we report the results using three different vision backbone architectures, including Swin [6], FocalNet [10] and DaViT [3]. All models in the first block are with tiny size and trained on the combination of image-label and image-text pairs, following the settings in UniCL [11]. In the second block, all the models are initialized with Florence [14] pre-trained DaViT-d5 model. Through the comparisons, we have the following observations: (1) FocalNet and DaViT achieve better performance than Swin across all metrics. Particularly, FocalNet achieves the best performance on generic and referring segmentation, while DaViT is better on the zero-shot vision-language evaluations; (2) After adding the deformable attention, we can see a boost on supervised segmentation but significant (especially large model) degradation on the open-vocabulary segmentation on ADE20K dataset. Based on these experimental results, we make the design choices as mentioned in our main submission: (1) we remove deformable attention in the favor of open-vocabulary segmentation; (2) we use FocalNet as the tiny vision encoder and train it by ourselves using UniCL, while using DaViT [14] as the base and large vision encoder.

C.3. Open-Vocabulary Generic Segmentation Settings Inspection

In Tab. 4, we study the progressive enrichment of data and training settings as well as the pre-trained model usage. X-Decoder-Seg is the baseline of adding a text encoder to Mask2Former [1] with a learnable language encoder. X-Decoder-Seg⁺ takes use of caption nouns for Hungarian matching to enrich the vocabulary size. In addition to the main paper, we add row 3 in Tab. 4 to demonstrate the performance of X-Decoder with only coco image text pairs. Comparing 3rd row and 4th row, we find adding extra image-text pairs for pretraining clearly improve open-vocabulary segmentation performance especially when the vocabulary size is large (e.g. ADE-150, CONTEXT-59/459). The way of pretraining vision backbone also matters. Comparing the last two rows side by side, though the backbone model sizes are similar, using ImageNet-21K for pretraining leads to inferior performance on most of the datasets except for CONTEXT-459 which contains most number of categories. These results demonstrate the benefits of using more image-text pairs for pretraining the vision backbone or our X-Decoder.

Method	Backbone	Deformable Attn.	Generic Segmentation						Referring g-Ref cIoU	Retrieval		Captioning	
			COCO			ADE (open)				COCO-Karpathy		COCO-Karpathy	
			PQ	mAP	mIoU	PQ	mAP	mIoU		IR@1	TR@1	CIDEr	BLEU
X-Decoder (T)	Swin	✗	50.2	38.8	61.9	17.3	9.4	23.7	55.3	28.0	43.7	79.9	24.2
X-Decoder (T)	Swin	✓	52.3	42.7	64.5	17.0	9.3	22.1	59.1	28.1	43.1	87.2	26.9
X-Decoder (T)	Focal	✗	51.4	40.5	62.8	18.8	9.8	25.0	59.8	30.7	48.5	79.9	24.2
X-Decoder (T)	Davit	✗	51.0	39.7	62.4	17.3	9.4	23.6	58.4	31.4	48.8	86.8	26.0
X-Decoder (L)	Davit	✗	56.9	46.7	67.7	21.8	13.1	29.6	64.2	44.7	60.3	111.0	32.6
X-Decoder (L)	Davit	✓	57.4	48.0	69.7	19.1	12.6	26.6	65.1	46.2	61.8	108.2	30.1

Table 3. Model architecture inspection among Swin [6], FocalNet [10] and DaViT [3]. “Deformable Attn.” means multi-scale deformable attention [15] that is used in Mask2Former [1]. All numbers are reported in zero-shot manner without any task-specific finetuning, and the row colored in gray corresponds to the architecture used the main paper.

Model	COCO (p/s)			ITP	ADE-150		VOC	PC-59	PC-459	SUN	SCAN-20	SCAN-41	Cityscapes		BDD			
	m	cls	cap		PQ	mAP	mIoU	mIoU	mIoU	mIoU	mIoU	PQ	mIoU	mIoU	mAP	PQ	mIoU	PQ
X-Decoder-Seg (T)	✓	✓	✗	✗	13.7	6.3	18.0	89.3	59.3	11.5	16.3	6.4	46.6	14.9	30.2	36.9	13.0	
X-Decoder-Seg+(T)	✓	✓	✓	✗	15.0	7.8	21.3	93.1	61.7	10.4	28.7	30.7	30.8	17.1	48.2	16.7	37.1	40.0
X-Decoder (T)	✓	✓	✓	✗	16.6	8.3	22.3	94.4	57.6	11.9	33.1	39.7	26.4	21.9	51.0	15.6	35.5	45.0
X-Decoder (T)	✓	✓	✓	✓	18.8	9.8	25.0	96.2	62.9	12.3	34.5	37.8	30.7	21.7	47.3	16.0	37.2	42.4
X-Decoder (L-IN21K)	✓	✓	✓	✓	19.9	11.7	29.6	95.8	54.2	20.5	42.4	44.9	29.5	27.4	47.2	18.3	33.3	44.9
X-Decoder (L)	✓	✓	✓	✓	21.8	13.1	29.6	97.7	64.0	16.1	43.0	49.5	39.5	29.7	52.0	24.9	38.1	47.2

Table 4. More open-vocabulary segmentation results. We report the results for our X-Decoder pretrained with COCO segmentation and caption annotations only in 3rd row. Additionally, we compare the model initialized with two different pre-trained large vision backbones, FocalNet-Large and DaViT-d5 trained on ImageNet-21K (row 5) and hundreds of millions of image-text pairs (row 6), respectively.

D. Segmentation In the Wild Benchmark

As shown in the main submission, our X-Decoder exhibits a strong generalization ability to segment images in ten settings of seven datasets from different domains, without any dataset-specific finetuning. Inspired by the object detection in the wild setting proposed in GLIP [4], we resort to more domain-specific datasets on the web to further examine the generality of our model. Specifically, we download 55 instance segmentation datasets from Roboflow¹. Afterward, we clean the datasets by excluding those containing visually undetectable categories (e.g. Different species of plant) or categories labeled with other languages. In the end, we compile 25 datasets that are suitable for evaluation into *segmentation in the wild* (*SegInW*) benchmark and report instance segmentation mAP. The dataset meta information is listed in Tab. 5, and exemplar images are shown in Fig. 3.

On the *SegInW* benchmark, we evaluate zero-shot, few-shot, and fine-tuned segmentation for five models (X-Decoder-Seg⁺ as baselines, and X-Decoder with different visual backbone) on three different tuning scales. In Fig. 4, we report the zero-shot instance segmentation performance on 25 datasets separately in a descending order. Accordingly, X-Decoder shows reasonably good generalization ability to a wide range of visual and concept domains. Specifically, it achieves higher mAP on common objects like fruits and animals but lower ones on fine-grained datasets like toolkits and rare concepts like rail and brain tumor. In Fig. 5, we further show the line charts for few-shot learning and fully-finetuning, and observe that:

X-Decoder has privilege on small-scale tuning. As shown in Fig. 5 (a-b), comparing with X-Decoder-Seg⁺ that only extract noun phrase to increase vocabulary size, X-Decoder performs much better with few-shot/finetune setting. Although X-Decoder (B) and X-Decoder-Seg⁺ (B) have similar zero-shot performance, the gap increases with the number of images tuned. However, as the number of parameters tuned increased by a large margin Fig. 5 (c), the performance gap between X-Decoder and X-Decoder-Seg⁺ is shrunk to a small margin.

Zero-Shot gap could be bridged by tuning. X-Decoder (L) and X-Decoder (L-IN21K) are initialized with different pre-trained image backbones. Specifically, X-Decoder (L) is initialized by Florence [14] pre-trained Davit-d5, whereas X-Decoder (L-IN21K) is initialized with FocalNet-L pre-trained on ImageNet-21k [2]. As shown in Fig. 5 (a-c), although the gap between X-Decoder-L and X-Decoder-L-IN21K on the zero-shot setting is relatively large. However, the gap on 5/10/full finetuned settings is much smaller and even cross in some settings.

Tuning class embedding is enough for few-shot settings. As shown in Fig. 5 (e-h), on the smaller scale backbone including (T/B), although tuning the full decoder has a better result, the gap is not obvious on 0-10 shots. And on larger scale models including L/L-IN21K, tuning with class embedding has similar/better results on 0-10 shots.

We show more detailed results in Table 6, Table 7 and Table 8. Similar to main paper, we report the number of parameters tuned in each setting.

¹<https://roboflow.com/>

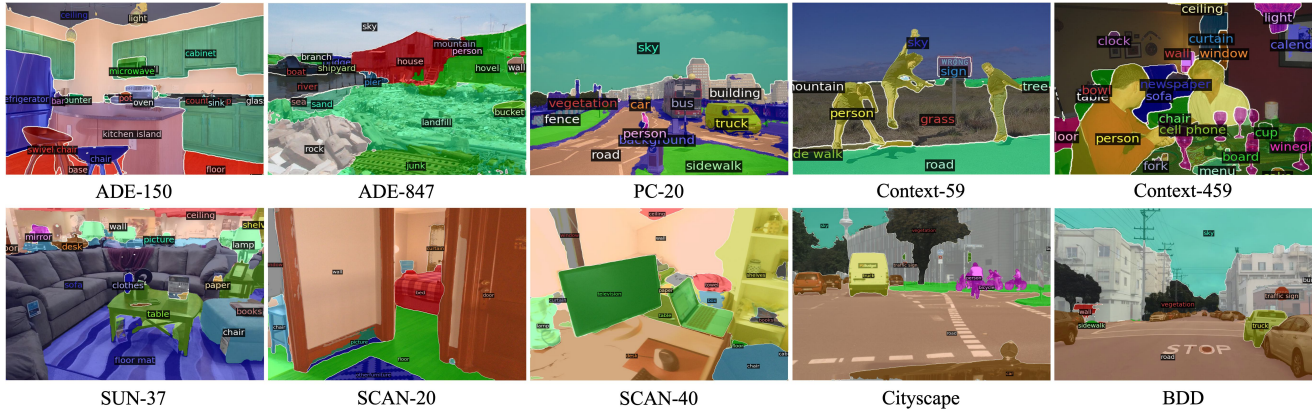


Figure 2. Visualization of zero-shot semantic segmentation on 10 settings of 7 datasets.

Dataset	Categories	# Class	# Images Train Val	URL
Phones	[phone]	1	25 11	https://universe.roboflow.com/workspace-c4esq/phone-xccez/dataset/1
Elephants	[elephant]	1	883 99	https://universe.roboflow.com/ds/4YwrXdlbFy?key=Q5GY9ITu14
Hand-Metal	[hand, metal]	2	504 65	https://universe.roboflow.com/nk950357-gmail-com/lab-k8hyn
Watermelon	[watermelon]	1	65 23	https://universe.roboflow.com/gnous-b5xq6/my-project-38agt/dataset/4
House-Parts	[aluminium door, aluminium window, ...]	22	700 201	https://universe.roboflow.com/testcoco/abc-fqun0/dataset/1
HouseHold-Items	[bottle, mouse, perfume, phone]	4	45 3	https://universe.roboflow.com/maths/household-items-sltdd/dataset/2
Strawberry	[R_strawberry, people]	2	971 87	https://universe.roboflow.com/strawberry-25w7z/strawberry_coco_1/dataset/1
Fruits	[apple, lemon, orange, pear, strawberry]	5	120 9	https://universe.roboflow.com/ds/PEVo9xHLf1?key=PXeJGF0D5q
Butterfly-Squirrel	[butterfly, squirrel]	2	951 237	https://universe.roboflow.com/handwashhygeine/nature-3tkys
Hand	[Hand-Segmentation, hand]	2	210 60	https://universe.roboflow.com/rmutsb-xxgii/hand-segmentation-gqzuh/dataset/1
Garbage	[bin, garbage, pavement, road]	4	325 142	https://universe.roboflow.com/project-blmh9/d2-bj1a0/dataset/1
Chicken	[chicken]	1	19 1	https://universe.roboflow.com/nea-trikic-ljxt3/chickenstf/dataset/1
Rail	[rail]	1	3067 1069	https://universe.roboflow.com/wzk789wzk-gmail-com/rail_dataset/dataset/4
Airplane-Parts	[Airplane, Body, Cockpit, Engine, Wing]	5	39 7	https://universe.roboflow.com/foxehecorp/foxehecorp_airplane_dataset/dataset/4/download
Brain-Tumor	[tumor]	1	236 28	https://universe.roboflow.com/detection-qskiw/segmnetation/dataset/2
Poles	[poles]	1	11 3	https://universe.roboflow.com/ohsee/pole2/dataset/2
Electric-Shaver	[caorau]	1	288 24	https://universe.roboflow.com/fpt-university-1tkhk/caurau
Bottles	[bottle, can, label]	3	357 16	https://universe.roboflow.com/beerup/bottels2/dataset/1
Toolkits	[Allen-key, block, gasket, ...]	8	48 6	https://universe.roboflow.com/mst/mask-2ihnt/dataset/1
Trash	[Aluminium foil, Cigarette, ...]	12	832 92	https://universe.roboflow.com/sara-najafi/trash_segmentation2/dataset/2
Salmon-Fillet	[Salmon, fillet]	1	1991 64	https://universe.roboflow.com/rishik-mishra-rljwe/fl225
Puppies	[puppy]	1	15 3	https://universe.roboflow.com/marcin-bak/puppies-fmoxu/dataset/2
Tablets	[tablets]	1	237 13	https://universe.roboflow.com/detection-qskiw/tablets-instance-segmentation/dataset/1
Cows	[cow]	1	630 60	https://universe.roboflow.com/new-workspace-5abdm/maskrcnn-ofglr/dataset/2
Ginger-Garlic	[garlic, ginger]	2	28 8	https://universe.roboflow.com/george-brown-college-1omrb/ginger-and-garlic-object-segmentation/dataset/1

Table 5. Meta information of *SegInW* benchmark. We list the source links, annotated category names and number of categories for each dataset.

E. Extra Visualization

In this part, we demonstrate the generalization ability to video datasets and flexibility to support task compositions for X-Decoder with more qualitative visualizations.

E.1. Zero-Shot Semantic Segmentation

We visualize zero-shot open vocabulary semantic segmentation on 7 datasets in 10 settings in Fig. 2. The visualization indicates that our model has strong generalization ability on images in different domains as well as categories.

E.2. Zero-Shot Generic Video Segmentation

Open-vocabulary generic segmentation is one of the main advantages of X-Decoder. We also apply generic segmentation in a zero-shot manner to the YoutubeVOS [9] dataset. As shown in Fig. 13, our model can be well generalized to video zero-shot generic segmentation and make predictions that are consistent across frames. As a result, our model can be used in video segmentation directly or a good initialization for further finetuning.

E.3. Zero-Shot Referring Video Segmentation

Besides the generic segmentation on video frames, our X-Decoder can be easily adapted to referring video segmentation as well without any architectural change or finetuning. In Fig. 14, we visualize some examples of referring video segmentation on the YoutubeVOS [9] dataset in a zero-shot manner. We can see that our model can generate rather accurate outputs given various referring phrases. Notably, in addition to the strong segmentation performance for given concepts, the model can also correctly distinguish the spatial locations (e.g., left v.s. right in the first row), and object attributes (e.g., a baby gorilla instead of an adult gorilla in the second row) in these unseen videos.

E.4. Zero-Shot Image Captioning

To test the generalization ability of X-Decoder, we also ask the model generate image captions on the YoutubeVOS [9] dataset, which is in a different domain from the image data. As we can see from the examples in Fig. 15, the model can correctly predict the object, activ-

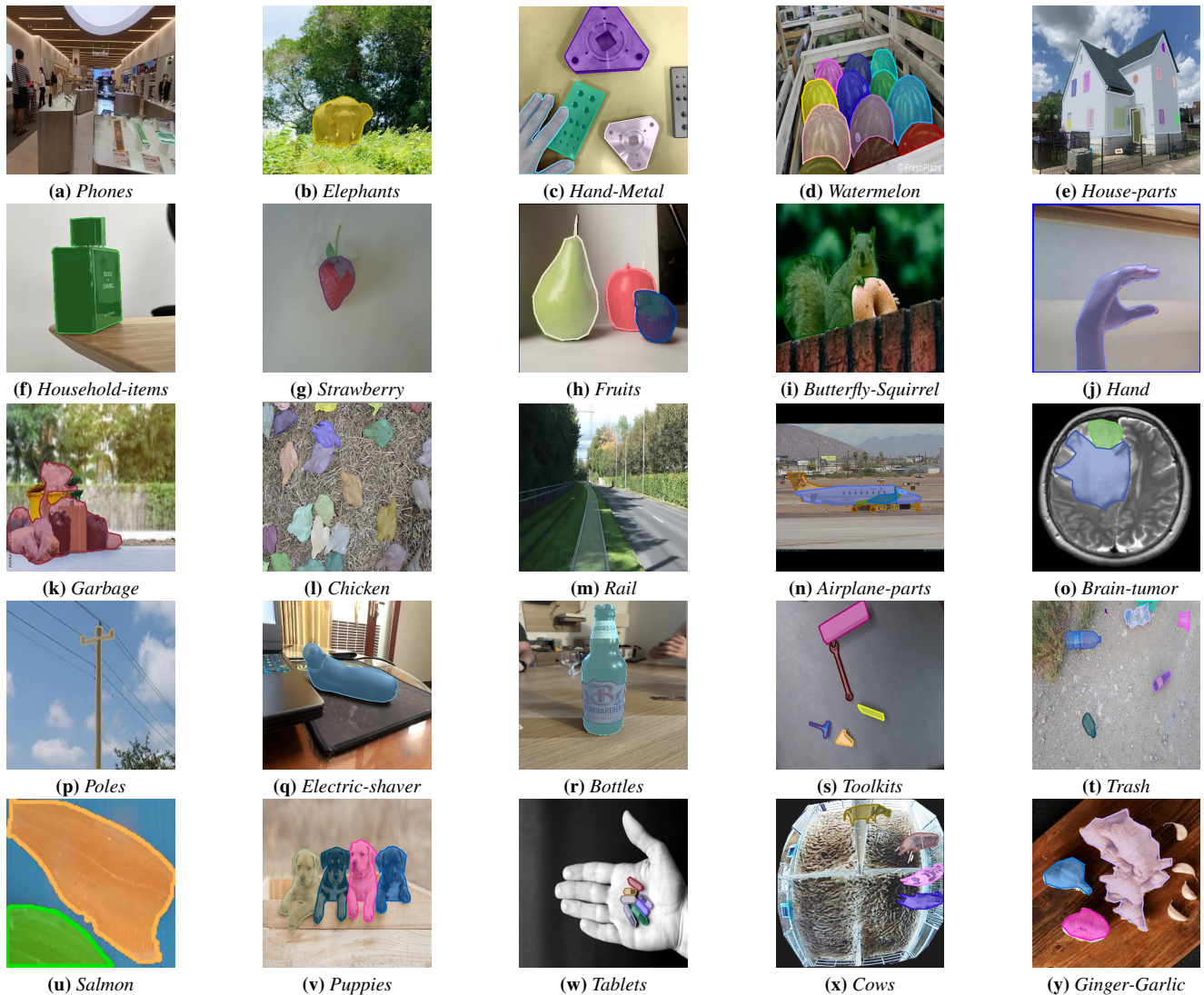


Figure 3. Exemplar images and annotations in *SegInW* benchmark. The benchmark covers a diversity of visual domains and concepts in the daily life.

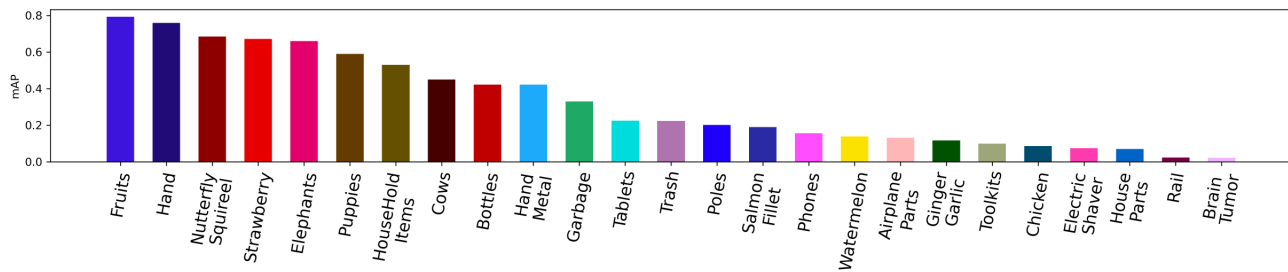


Figure 4. Zero-shot segmentation performance on *SegInW* with X-Decoder-L model. We report the mAP in descending order.

ity, and environment in an image. Interestingly, the captions for the first 6 images sampled from 3 different videos show that our approach can correctly differentiate the movements from similar scenarios (e.g., a man playing vs. a man standing in the first two samples.).

E.5. Zero-Shot Referring Captioning

In compensating for the visualization of the main paper, we add more referring captioning samples in Fig 16. The phrase before “:” is the referring phrase, and the sentence after “:” is the generated caption. The grounding mask of the referring phrase is highlighted in pink. Clearly, our model

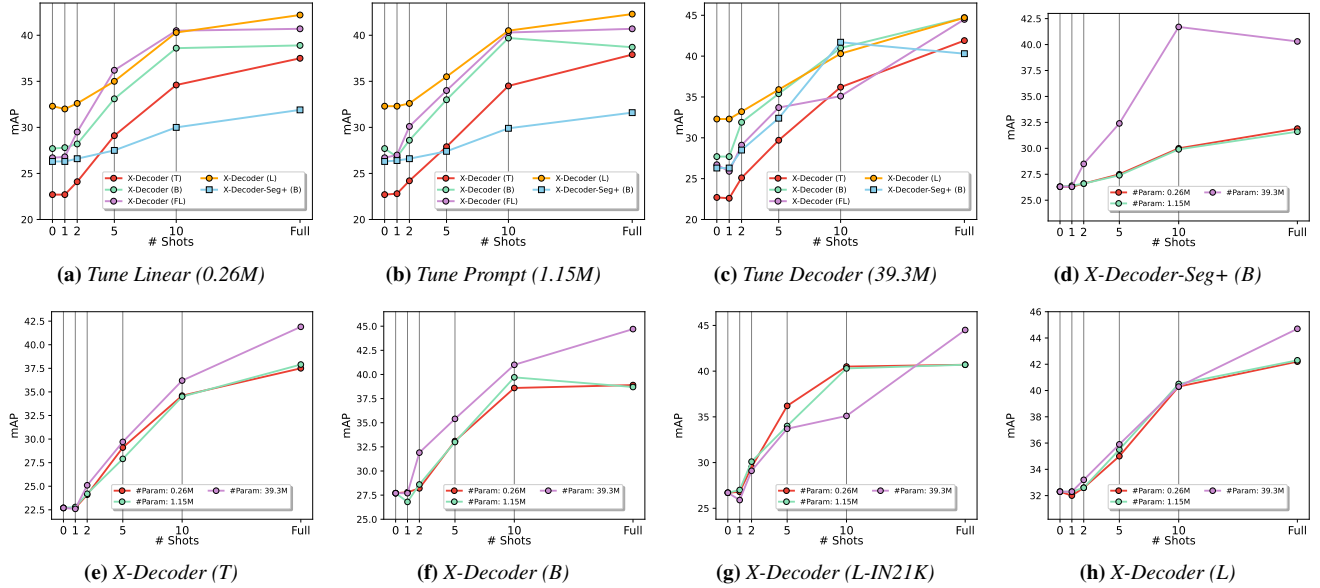


Figure 5. (a-c) Line chart of tuning shot and mAP with different tuning strategies on each backbone architecture. (d-h) Line chart of tuning shot and mAP with different backbone architecture on different number of tuning parameters (strategies).

can simultaneously segments the referred region and generates a region-specific caption. Complementary to regular image captioning systems, such a novel functionality provides a way of interpreting images in a more fine-grained manner. Note that our X-Decoder was never trained to generate such regional captions.

E.6. Zero-Shot Referring Image Editing

Finally, given the high-quality referring segmentation results with X-Decoder, we can effortlessly combine it with off-the-shelf Stable-Diffusion image inpainting model [8] and perform zero-shot referring image editing. As shown in Fig. 17, the model first performs referring segmentation, then the original image and the segmentation mask are fed into the inpainting model to generate the inpainted image. For example, given “change bird to squirrel”, it first extracts the bird segment (blue region) from the input image and then replace the segmented region with a generated squirrel. Likewise in other samples, we can see all the generated images look natural and follow the inpainting instructions very well. These impressive plug-and-play results imply a great potential of combining our X-Decoder and advanced generative AI models for fine-grained precise image editing.

F. Discussions

Future Directions. The extensive quantitative and qualitative results have demonstrated the strong performance and generalization ability of our X-Decoder for a variety of vision and vision-language tasks at different granularities. Upon the current X-Decoder design, we see two directions worth future explorations: (1) *Pretrain the whole model in one stage effectively and efficiently.* Currently, the

model still requires a separate pretraining for the image and text encoders. However, since our model supports large-scale image-text contrastive learning thanks to the decoupled design, we can easily unify the CLIP-style pretraining with the decoder pretraining in an end-to-end manner. (2) *Unify all level of supervisions.* Due to high annotation costs, the pixel-level segmentation annotations by nature are much less than the region-level box and image-level annotations. It is worth building a more unified learning paradigm to jointly learn from pixel-level, region-level and image-level supervision to attain a more powerful unified model.

Social Impact. This work is mainly focused on the design of a generalized decoder for various vision and vision-language tasks. We have used a pretrained image and text encoder and further pretrained the models on a combination of various datasets and tasks. Since the models are trained on large-scale webly-crawled image-text pairs, the negative impact might arise due to the potential offensive or biased content in the data. To mitigate this issue, we need to have a careful sanity check on the training data and model predictions before deploying it in practical scenarios.

Model	Shot #Param	Avg	Airplane-Parts	Bottles	Brain-Tumor	Chicken	Cows	Electric-Shaver	Elephants	Fruits	Garbage	Ginger-Garlic	Hand	Hand-Metal	House-Parts	IHL-Items	Butterfly-Squirrel	Phones	Poles	Puppies	Rail	Salmon-Fillet	Strawberry	Tablets	Toolkits	Trash	Watermelon
X-Decoder (T)	0 0.0M	22.7	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	1 39.3M	22.7	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	3 39.3M	24.1	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	5 39.3M	29.1	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	10 39.3M	34.6	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	3xAll 39.3M	37.5	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder-Seg+ (B)	0 0.0M	26.3	13.2	17.2	0.8	33.0	28.6	4.9	67.9	71.1	28.8	5.2	0.0	0.8	6.8	50.6	53.2	18.8	17.9	68.2	0.7	21.1	86.3	5.8	11.5	12.1	31.7
X-Decoder-Seg+ (B)	1 39.3M	26.3	13.2	17.2	0.8	33.0	28.6	4.9	67.9	71.1	28.8	5.2	0.0	0.8	6.8	50.6	53.2	18.8	17.9	68.2	0.7	21.1	86.3	5.8	11.5	12.1	31.7
X-Decoder-Seg+ (B)	3 39.3M	26.6	13.2	17.2	0.8	33.0	28.6	4.9	67.9	71.1	28.8	5.2	0.0	0.8	6.8	50.6	53.2	18.8	17.9	68.2	0.7	21.1	86.3	5.8	11.5	12.1	31.7
X-Decoder-Seg+ (B)	5 39.3M	27.5	13.2	17.2	0.8	33.0	28.6	4.9	67.9	71.1	28.8	5.2	0.0	0.8	6.8	50.6	53.2	18.8	17.9	68.2	0.7	21.1	86.3	5.8	11.5	12.1	31.7
X-Decoder-Seg+ (B)	10 39.3M	30.0	13.2	17.2	0.8	33.0	28.6	4.9	67.9	71.1	28.8	5.2	0.0	0.8	6.8	50.6	53.2	18.8	17.9	68.2	0.7	21.1	86.3	5.8	11.5	12.1	31.7
X-Decoder-Seg+ (B)	3xAll 39.3M	31.9	13.2	17.2	0.8	33.0	28.6	4.9	67.9	71.1	28.8	5.2	0.0	0.8	6.8	50.6	53.2	18.8	17.9	68.2	0.7	21.1	86.3	5.8	11.5	12.1	31.7
X-Decoder (B)	0 0.0M	27.7	13.0	45.9	0.3	13.6	36.8	4.2	68.0	76.7	30.2	19.4	20.6	18.5	6.7	51.7	53.1	8.9	5.6	55.4	0.8	18.2	81.6	8.0	13.9	27.3	13.0
X-Decoder (B)	1 39.3M	27.8	13.0	45.9	0.3	13.6	36.8	4.2	68.0	76.7	30.2	19.4	20.6	18.5	6.7	51.7	53.1	8.9	5.6	55.4	0.8	18.2	81.6	8.0	13.9	27.3	13.0
X-Decoder (B)	3 39.3M	28.2	13.0	45.9	0.3	13.6	36.8	4.2	68.0	76.7	30.2	19.4	20.6	18.5	6.7	51.7	53.1	8.9	5.6	55.4	0.8	18.2	81.6	8.0	13.9	27.3	13.0
X-Decoder (B)	5 39.3M	31.1	13.0	45.9	0.3	13.6	36.8	4.2	68.0	76.7	30.2	19.4	20.6	18.5	6.7	51.7	53.1	8.9	5.6	55.4	0.8	18.2	81.6	8.0	13.9	27.3	13.0
X-Decoder (B)	10 39.3M	33.9	13.0	45.9	0.3	13.6	36.8	4.2	68.0	76.7	30.2	19.4	20.6	18.5	6.7	51.7	53.1	8.9	5.6	55.4	0.8	18.2	81.6	8.0	13.9	27.3	13.0
X-Decoder (B)	3xAll 39.3M	38.1	13.0	45.9	0.3	13.6	36.8	4.2	68.0	76.7	30.2	19.4	20.6	18.5	6.7	51.7	53.1	8.9	5.6	55.4	0.8	18.2	81.6	8.0	13.9	27.3	13.0

Table 6. SegInW results with tuning on class embedding for different image shots and backbone architectures. (39.3M parameters tuned in the setting.)

Model	Shot #Param	Avg	Airplane-Parts	Bottles	Brain-Tumor	Chicken	Cows	Electric-Shaver	Elephants	Fruits	Garbage	Ginger-Garlic	Hand	Hand-Metal	House-Parts	IHL-Items	Butterfly-Squirrel	Phones	Poles	Puppies	Rail	Salmon-Fillet	Strawberry	Tablets	Toolkits	Trash	Watermelon
X-Decoder (T)	0 0.0M	22.7	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	1 1.15M	22.8	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	3 1.15M	24.1	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	5 1.15M	27.9	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	10 1.15M	32.3	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	3xAll 1.15M	37.3	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder-Seg+ (B)	0 0.0M	26.3	13.2	17.2	0.8	33.0	28.6	4.9	67.9	71.1	28.8	5.2	0.0	0.8	6.8	50.6	53.2	18.8	17.9	68.2	0.7	21.1	86.3	5.8	11.5	12.1	31.7
X-Decoder-Seg+ (B)	1 1.15M	26.4	13.2	17.2	0.8	33.0	28.6	4.9	67.9	71.1	28.8	5.2	0.0	0.8	6.8	50.6	53.2	18.8	17.9	68.2	0.7	21.1	86.3	5.8	11.5	12.1	31.7
X-Decoder-Seg+ (B)	3 1.15M	26.6	13.2	17.2	0.8	33.0	28.6	4.9	67.9	71.1	28.8	5.2	0.0	0.8	6.8	50.6	53.2	18.8	17.9	68.2	0.7	21.1	86.3	5.8	11.5	12.1	31.7
X-Decoder-Seg+ (B)	5 1.15M	27.4	13.2	17.2	0.8	33.0	28.6	4.9	67.9	71.1	28.8	5.2	0.0	0.8	6.8	50.6	53.2	18.8	17.9	68.2	0.7	21.1	86.3	5.8	11.5	12.1	31.7
X-Decoder-Seg+ (B)	10 1.15M	29.9	13.2	17.2	0.8	33.0	28.6	4.9	67.9	71.1	28.8	5.2	0.0	0.8	6.8	50.6	53.2	18.8	17.9	68.2	0.7	21.1	86.3	5.8	11.5	12.1	31.7
X-Decoder-Seg+ (B)	3xAll 1.15M	31.6	13.2	17.2	0.8	33.0	28.6	4.9	67.9	71.1	28.8	5.2	0.0	0.8	6.8	50.6	53.2	18.8	17.9	68.2	0.7	21.1	86.3	5.8	11.5	12.1	31.7
X-Decoder (B)	0 0.0M	27.7	13.0	45.9	0.3	13.6	36.8	4.2	68.0	76.7	30.2	19.4	20.6	18.5	6.7	51.7	53.1	8.9	5.6	55.4	0.8	18.2	81.6	8.0	13.9	27.3	13.0
X-Decoder (B)	1 1.15M	27.8	13.0	45.9	0.3	13.6	36.8	4.2	68.0	76.7	30.2	19.4	20.6	18.5	6.7	51.7	53.1	8.9	5.6	55.4	0.8	18.2	81.6	8.0	13.9	27.3	13.0
X-Decoder (B)	3 1.15M	28.6	13.0	45.9	0.3	13.6	36.8	4.2	68.0	76.7	30.2	19.4	20.6	18.5	6.7	51.7	53.1	8.9	5.6	55.4	0.8	18.2	81.6	8.0	13.9	27.3	13.0
X-Decoder (B)	5 1.15M	32.3	13.0	45.9	0.3	13.6	36.8	4.2	68.0	76.7	30.2	19.4	20.6	18.5	6.7	51.7	53.1	8.9	5.6	55.4	0.8	18.2	81.6	8.0	13.9	27.3	13.0
X-Decoder (B)	10 1.15M	37.3	13.0	45.9	0.3	13.6	36.8	4.2	68.0	76.7	30.2	19.4	20.6	18.5	6.7	51.7	53.1	8.9	5.6	55.4	0.8	18.2	81.6	8.0	13.9	27.3	13.0
X-Decoder (B)	3xAll 1.15M	42.3	13.0	45.9	0.3	13.6	36.8	4.2	68.0	76.7	30.2	19.4	20.6	18.5	6.7	51.7	53.1	8.9	5.6	55.4	0.8	18.2	81.6	8.0	13.9	27.3	13.0

Table 7. SegInW results with tuning on class & mask embeddings and latent queries for different image shots and backbone architectures. (1.15M parameters tuned in the setting.)

Model	Shot #Param	Avg	Airplane-Parts	Bottles	Brain-Tumor	Chicken	Cows	Electric-Shaver	Elephants	Fruits	Garbage	Ginger-Garlic	Hand	Hand-Metal	House-Parts	IHL-Items	Butterfly-Squirrel	Phones	Poles	Puppies	Rail	Salmon-Fillet	Strawberry	Tablets	Toolkits	Trash	Watermelon
X-Decoder (T)	0 0.0M	22.7	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	1 0.26M	22.6	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	3 0.26M	25.1	10.4	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	5 0.26M	29.7	10.5	19.0	1.1	12.0	12.0	1.2	65.6	66.5	28.7	7.9	0.6	22.4	5.5	50.6	62.1	29.9	3.6	48.9	0.7	15.0	41.6	15.2	9.5	19.3	16.2
X-Decoder (T)	10 0.26M	36.2																									

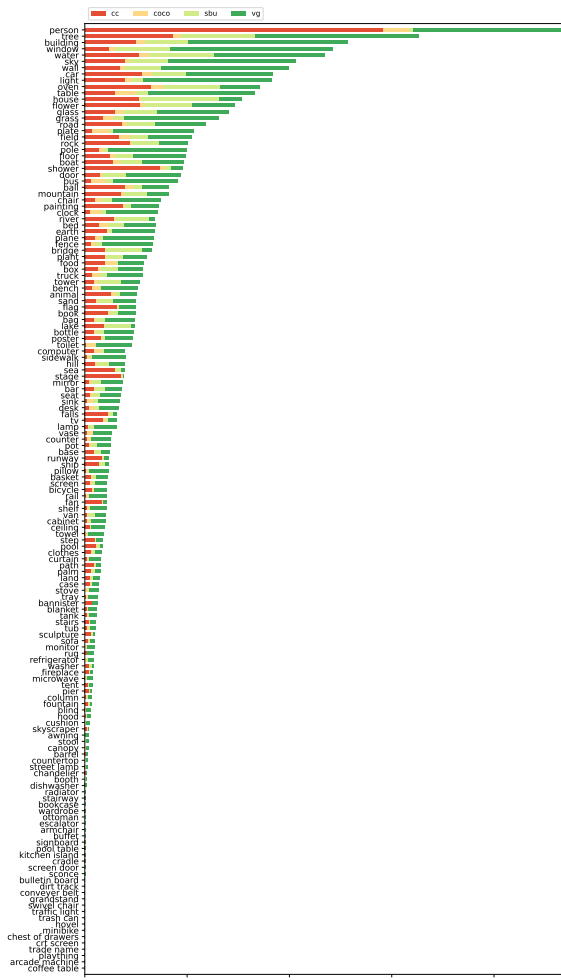


Figure 6. Image captions overlap with ADE20K-150

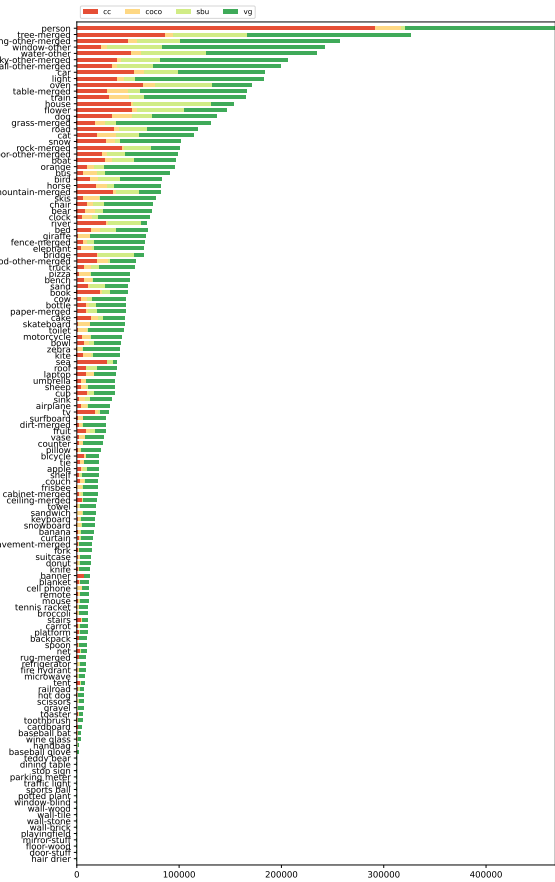


Figure 10. Image captions overlap with COCO

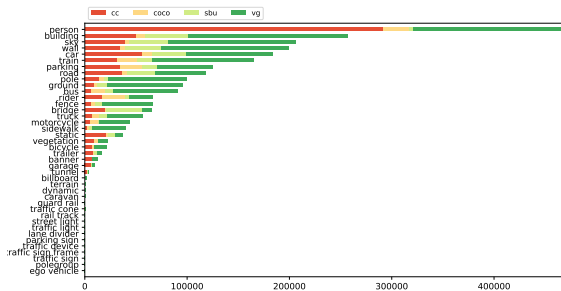


Figure 7. Image captions overlap with BDD-Panoptic

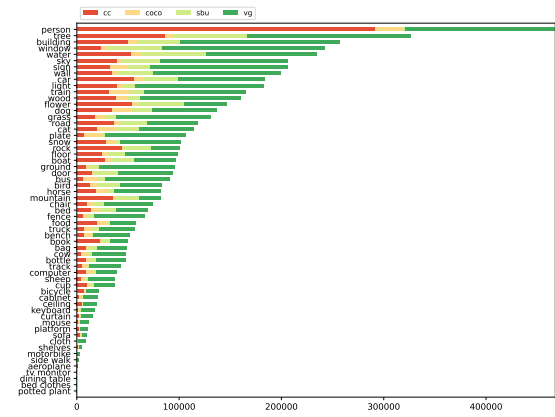


Figure 11. Image captions overlap with Pascal Context-59

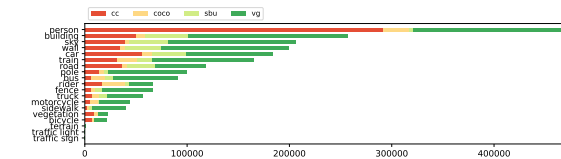


Figure 8. Image captions overlap with BDD-Semantic/Cityscapes

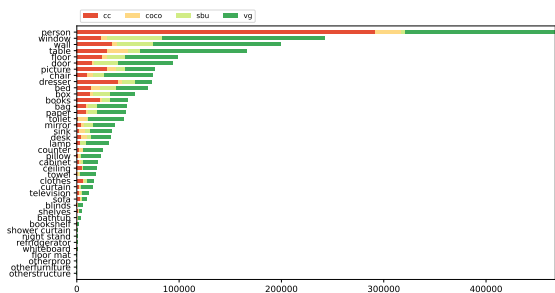


Figure 12. Image captions overlap with ScanNet-40

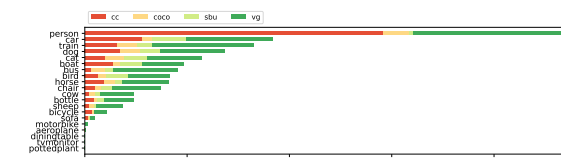


Figure 9. Image captions overlap with Pascal VOC



Figure 13. Zero-Shot Video Generic Segmentation. (Source: YoutubeVOS videos)

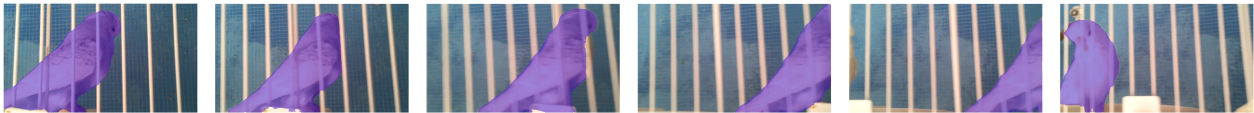
A monkey to the left of another monkey



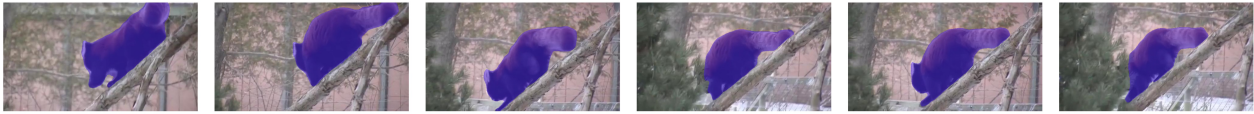
A baby gorilla



A parrot in a cage:



A panda with a bushy tail walking down a branch of a tree



A black bear



A swimming penguin

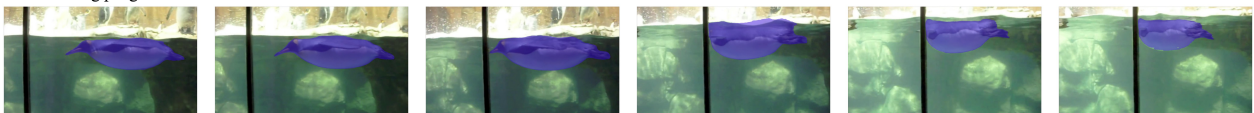


Figure 14. Zero-Shot Referring Video Segmentation. (Source: YoutubeVOS videos)

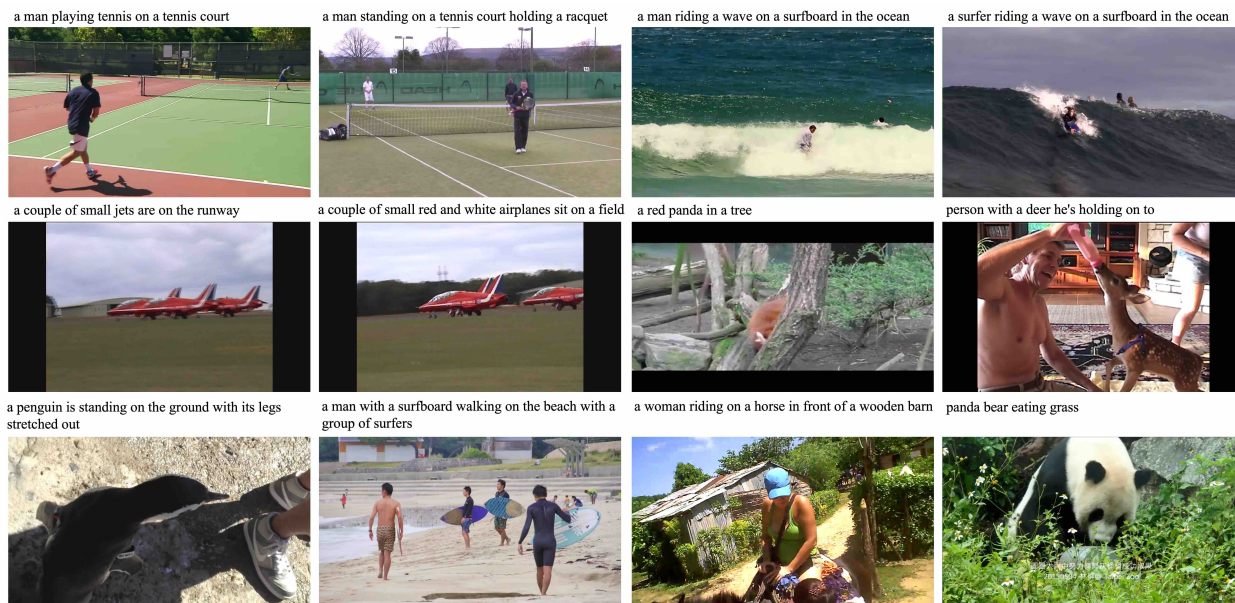


Figure 15. Zero-Shot Image Captioning. (Source: YoutubeVOS videos)

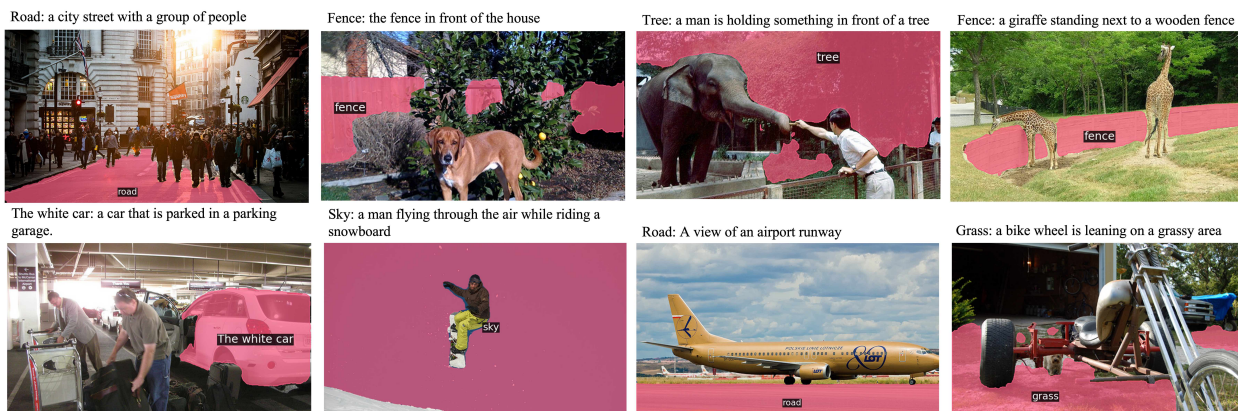


Figure 16. Referring Captioning. (Source: COCO 2017 val images)

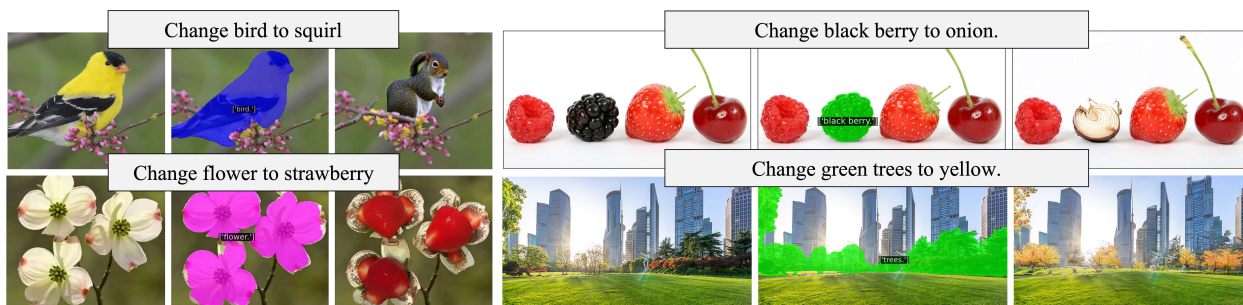


Figure 17. Referring Image Inpainting. (Source: web images)

References

- [1] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1290–1299, 2022.
- [2] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [3] Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jingdong Wang, and Lu Yuan. Davit: Dual attention vision transformers. *arXiv preprint arXiv:2204.03645*, 2022.
- [4] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10965–10975, 2022.
- [5] 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 *ECCV*, 2014.
- [6] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021.
- [7] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [8] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.
- [9] Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas Huang. Youtube-vos: A large-scale video object segmentation benchmark. *arXiv preprint arXiv:1809.03327*, 2018.
- [10] Jianwei Yang, Chunyuan Li, and Jianfeng Gao. Focal modulation networks. *arXiv preprint arXiv:2203.11926*, 2022.
- [11] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Bin Xiao, Ce Liu, Lu Yuan, and Jianfeng Gao. Unified contrastive learning in image-text-label space. In *CVPR*, 2022.
- [12] Xuwang Yin and Vicente Ordonez. Obj2text: Generating visually descriptive language from object layouts. *arXiv preprint arXiv:1707.07102*, 2017.
- [13] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *European Conference on Computer Vision*, pages 69–85. Springer, 2016.
- [14] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, Ce Liu, Mengchen Liu, Zicheng Liu, Yumao Lu, Yu Shi, Lijuan Wang, Jianfeng Wang, Bin Xiao, Zhen Xiao, Jianwei Yang, Michael Zeng, Luwei Zhou, and Pengchuan Zhang. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021.
- [15] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020.