

2HandedAfforder: Learning Precise Actionable Bimanual Affordances from Human Videos

Marvin Heidinger^{*1}, Snehal Jauhri^{*1}, Vignesh Prasad¹, Georgia Chalvatzaki^{1,2}

* indicates equal contribution

¹Computer Science Department, Technische Universität Darmstadt, Germany

²Hessian.AI, Darmstadt, Germany

{snehal.jauhri, vignesh.prasad, georgia.chalvatzaki}@tu-darmstadt.de

Abstract

When interacting with objects, humans effectively reason about which regions of objects are viable for an intended action, i.e., the affordance regions of the object. They can also account for subtle differences in object regions based on the task to be performed and whether one or two hands need to be used. However, current vision-based affordance prediction methods often reduce the problem to naive object part segmentation. In this work, we propose a framework for extracting affordance data from human activity video datasets. Our extracted 2HANDS dataset contains precise object affordance region segmentations and affordance class-labels as narrations of the activity performed. The data also accounts for bimanual actions, i.e., two hands co-ordinating and interacting with one or more objects. We present a VLM-based affordance prediction model, 2HandedAfforder, trained on the dataset and demonstrate superior performance over baselines in affordance region segmentation for various activities. Finally, we show that our predicted affordance regions are actionable, i.e., can be used by an agent performing a task, through demonstration in robotic manipulation scenarios. Project-website: sites.google.com/view/2handedafforder

1. Introduction

When humans perceive objects, they understand different object regions and can predict which object region *affords* which activities [7], i.e., which object regions can be used for a task. We wish our machines to have this ability, referred to in literature as “affordance grounding”. Affordance grounding has several downstream applications, including building planning agents, VR, and robotics. Affordance grounding is especially important for robotics since robots must reason about various actions that can be performed using different object regions which is a crucial step towards performing useful tasks in everyday, unstructured



Figure 1. **A motivating example:** When labeling affordances for a task ‘Pour into bowl’, typical labeled affordances provided by annotators are not precise and reduce the problem to object part segmentation. Alternatively, our affordance extraction method uses the hand-object interaction sequence to get precise bimanual affordance regions that are not just ‘conceptual’ but also ‘actionable’.

environments. For example, to pour into a bowl, the robot should know that it should hold the bottle in a region close to the center of mass of the bottle (Figure 1), i.e., a region that *affords* pouring. Predicting such affordance regions is challenging since it requires a fine-grained understanding of object regions and their semantic relationship to the task.

Recent advances in large-language and multimodal models have shown impressive visual reasoning capabilities using self-supervised objectives [6, 23, 26]. However, there is still a big gap in their ability to detect accurate object affordance regions in images [18]. Moreover, most existing state-of-the-art affordance detection methods [10, 15, 24, 25, 29] use labeled data [11, 16, 21, 22, 24] that lacks precision and is more akin to object part segmentation rather than *actionable* affordance-region prediction. When humans interact with objects, they are much more *precise* and use specific object regions important in the context of the task. An example is provided in Fig. 1. For the task of pouring into the bowl, part segmentation labels the entire bottom of the

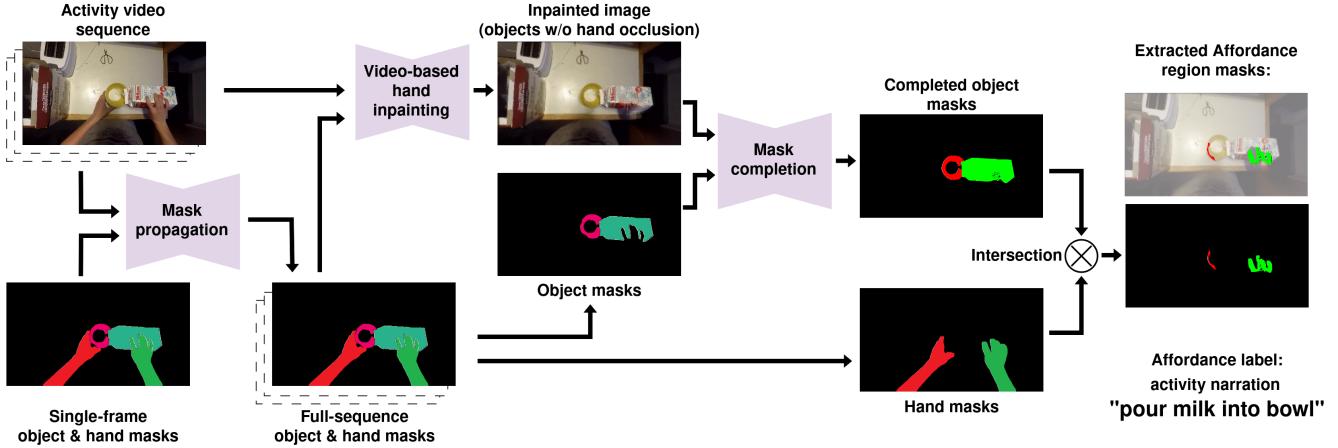


Figure 2. Affordance extraction pipeline. Given a human activity video sequence and a single-frame object and hand masks, we first obtain dense, full-sequence object and hand masks using a video mask-propagation network [3]. We then inpaint out the hands in the RGB images using a video-based hand inpainting model [2]. This gives us an image with the objects reconstructed and un-occluded by the hands. With the inpainted image and the original object masks, we use [27] to “complete” the object masks by again propagating the object masks to the inpainted image. Finally, we can extract the affordance region masks for the given task as the intersection between the completed masks and the hand masks. We also label the affordance class using the narration of the task.

bottle with the affordance ‘pour’. But, to pour correctly, humans leverage the appropriate region of the bottle. Moreover, the affordances are inherently bimanual, i.e., the affordance regions of the bowl and bottle are interconnected.

We argue that affordances should not be labeled but automatically extracted by observing humans performing tasks, e.g. in activity video datasets. We propose a method that uses hand-inpainting and mask completion to extract affordance regions occluded by human hands. This has several advantages. First, by using this procedure, we are able to obtain **bimanual** and **precise** affordances (Figure 1) rather than simply predicting object parts. Second, it makes affordance specification more natural since it is often easier for humans to *show* the object region to interact with, rather than label and segment it correctly in an image. Third, using human activity videos gives us diverse task-specific affordances, with the affordance class label naturally coming from the narration of what task is being done by the human. This makes our affordances **task-oriented** with natural language specification, unlike previous methods focused on predicting task-agnostic interaction hotspots [1, 8].

2. Extraction and Learning of Bimanual Affordances from Human Videos

2.1. Affordance Extraction from Human Videos

We use videos of humans performing tasks to extract precise affordance masks. This involves closely examining the contact regions between the hands and objects. We propose a pipeline to extract affordances that leverages recent advances in hand inpainting [2] and object mask completion [27, 28], providing the first bimanual affordance region

segmentation dataset. Moreover, we use the narration of the task being performed as the affordance text label, obtaining a diverse set of affordance classes for various objects.

We use videos from EPIC-KITCHENS [4], containing ~ 100 hours of egocentric human videos in kitchens. We use VISOR [5] annotations of the dataset, which contain sparse hand-object mask segmentations and binary labels for whether the hand is in contact with the object. To obtain dense hand-object masks for entire video sequences, we use a video-based mask propagation network [3].

With the hand and object masks available over the entire video sequence, we obtain an un-occluded view of the objects by inpainting out the hands. We use a video-based hand inpainting model, VIDM [2], that uses 4 frames from the sequence as input to inpaint the missing regions. This sequence-based inpainting better reconstructs the target objects since the objects may be visible in another frame of the sequence without occlusion. Inpainting provides us with an un-occluded view of the objects. We then precisely segment these un-occluded objects in the inpainted image using mask completion. For this, we use the segmentation masks from the original image and prompt SAM2 [27] to propagate these masks to the new inpainted image. To obtain the final affordance region where the hand and object interact, we can simply compute the intersection of the un-occluded object masks and the hand masks (Fig. 2). For bimanual affordances, we also classify into a bimanual taxonomy [14] of unimanual left, unimanual right, and bimanual actions.

We extract a dataset of 278K images with affordance segmentation masks, narration-based class-labels, and bimanual taxonomy annotations. We call our dataset 2HANDS, i.e., the **2-Handed Affordance + Narration DataSet**.

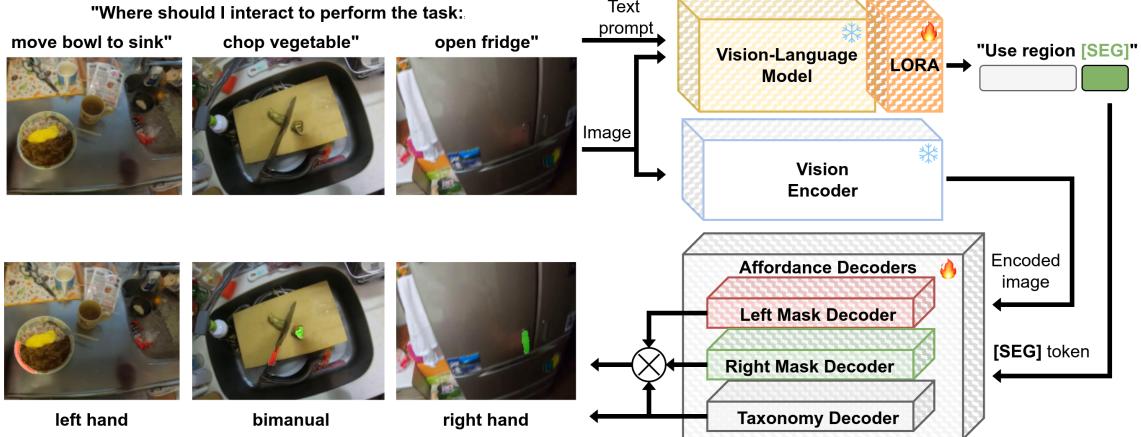


Figure 3. Affordance prediction network. Given an input image and task, we use a question asking where the objects should be interacted for the desired task as a text prompt to a Vision-Language model (VLM). The VLM produces language tokens and a [SEG] token which is passed to the affordance decoders. We also use a SAM [13] vision-backbone to encode the image and pass it to the affordance decoders. The decoders predict the left hand and right hand affordance region masks as well as a taxonomy classification indicating whether the interaction is supposed to be performed with the left hand, right hand, or both hands. The vision encoder is frozen, while the VLM predictions are fine-tuned using LORA [12].

2.2. Task-oriented Bimanual Affordance Prediction

Reasoning segmentation, i.e., text-prompt-based segmentation of full objects, is a difficult task. Segmentation of precise object affordance regions is even more challenging. The complexity is further increased when considering bimanual affordances with multiple objects. To address this challenge, we develop a model for general-purpose bimanual affordance prediction that can process both an input image and any task prompt (e.g., “pour tea from kettle”). We call this model “2HandedAfforder.” We leverage recent developments in reasoning-based segmentation methods [15, 19] and train a VLM-based segmentation model to reason about the required task and predict the relevant affordance region in the input image.

Inspired by reasoning segmentation methods such as by Lai et al. [15], we use a Vision-Language Model (VLM) [20], a LLaVa-13b, to jointly process the input text prompt and image and produce language tokens and a segmentation [SEG] token as output. While VLMs excel at tasks such as visual question answering and image captioning, they are not explicitly optimized for vision tasks like segmentation, where accurately predicting pixel-level information is key. Thus, to have a stronger vision-backbone for our segmentation-related task, we use a modified version of SAM [13]. Given the combined embedding provided by the VLM [SEG] token and SAM image encoder, we use affordance decoders modeled after SAM-style mask decoders to predict the affordances. We use two mask decoders, generating separate affordance masks for the left and right hands, respectively. Furthermore, we add a prediction head to one of the decoders that takes the output token as input and predicts the bimanual taxonomy: ‘uni-

manual left hand’, ‘unimanual right hand’, and ‘bimanual’ using a separate full-connected classifier decoder (Fig. 3).

The VLM is trained to generate a specific output token: a segmentation [SEG] token. Specifically, inspired by LISA [15], we use question-answer templates to encapsulate the narration of the individual tasks in natural language, e.g. ‘**USER:** [IMAGE] Where would you interact with the objects to perform the action {action_narration} in this image? **ANSWER:** Use region: [SEG].’ This [SEG] token encapsulates the general-purpose reasoning information from the VLM for the task which is then used by the affordance decoders. For the left and right hand mask decoders, we initialize the decoders with pre-trained SAM weights and train them to predict segmentation masks using the encoded image and [SEG] token as input. For the taxonomy classifier decoder, as in [24], we pass the left mask decoder output token through an MLP to predict whether the action should be performed with the left hand, right hand, or both hands.

3. Experiments

3.1. ActAffordance Benchmark

To answer the first question of the accuracy of our extracted affordances in the 2HANDS dataset, we evaluate the alignment of our extracted affordance masks with human-annotated affordance regions. As mentioned in Sec. 2.1, when humans label affordances, they often simply label object parts and do not necessarily focus on the precise regions of interaction of the objects [21, 24]. Moreover, the second question regarding the accuracy of 2HandedAfforder is non-trivial. Using only the masks in our 2HANDS dataset as “ground truth” leads to a bias towards our own extracted

Model	ActAffordance Benchmark														
	EPIC-KITCHENS					EGO4D					Combined				
IoU ↑	Precision ↑	HD ↓	Dir. HD ↓	mAP ↑	IoU ↑	Precision ↑	HD ↓	Dir. HD ↓	mAP ↑	IoU ↑	Precision ↑	HD ↓	Dir. HD ↓	mAP ↑	
LISA [15]	0.048	0.056	298	260	0.053	0.038	0.098	336	257	0.084	0.044	0.050	303	255	0.047
LOCATE [17]	0.010	0.014	274	261	0.007	-	-	-	-	-	-	-	-	-	-
AffLLM [25]	0.010	0.010	267	205	0.010	0.015	0.016	229	226	0.014	0.012	0.013	287	225	0.012
2HAffCLIP	0.032	0.077	359	317	0.068	0.023	0.050	306	250	0.047	0.026	0.064	341	292	0.059
2HAff	0.064	0.125	241	185	0.104	0.051	0.137	292	227	0.105	0.058	0.130	262	202	0.104
AffExtract	0.136	0.334	199	169	-	0.253	0.541	163	121	-	0.185	0.420	184	145	-

Table 1. Comparison of our models and baseline methods on the ActAffordance benchmark. Performance is evaluated separately on the EPIC-KITCHENS and EGO4D splits, as well as on the combined benchmark. The reported metrics include IoU (Intersection over Union), Precision, HD (Hausdorff Distance), Dir. HD (Directional Hausdorff Distance), and mAP (Mean Average Precision). For mAP, we average over five different thresholds, and the values for the other metrics correspond to the highest scores obtained across these thresholds. We also run our affordance extraction method, AffExtract, as a measure of data quality and alignment with the benchmark annotations.



Figure 4. Examples of different manipulation tasks executed on a bimanual Tiago++ robot. Red and green masks denote left and right hand affordance mask predictions, respectively. We segment the task-specific object affordance regions, propose grasps for these regions, and use pre-designed motion primitives to execute manipulation tasks. Videos are available at sites.google.com/view/2handedafforder.

affordances. Therefore, we propose a novel benchmark called “ActAffordance” to evaluate both the dataset quality and the predicted affordances. Specifically, we evaluate the alignment of our affordances with the affordances annotated by humans who are shown the *full interaction video sequence*. Annotators predicted ALL possible interaction regions since affordance prediction is inherently multi-modal—for instance, when closing a fridge, a human might choose any point along the door length. The benchmark contains unimanual and bimanual segmentation masks for 400 activities from EPIC-KITCHENS [4] and Ego4D [9], with no overlap with the data used in 2HANDS.

4. Results

Extraction quality & Benchmark performance:

Table 1 shows the quantitative results. The high precision of AffExtract shows a reasonably good alignment with the human-annotated segmentations from the benchmark and meaningful affordance region extraction. The IoU scores are relatively lower, with an average of 0.185, showing the challenge of the task when compared against human-level object understanding. Since ours is the first method to perform bimanual affordance mask detection us-

ing text prompts, we adapt baselines which includes a SOTA text-based reasoning segmentation baseline, LISA [15]. 2HandedAfforder achieves the best results across all metrics. LISA is the next best method since it accurately segments the correct object in the scene, resulting in a natural overlap with the ground truth. This demonstrates the power of reasoning segmentation for the challenging task of prompt-based affordance prediction. Though our models were not trained on any Ego4D data, their performance on Ego4D is still reasonable and often better than the EPIC-KITCHENS split. The IoU scores are low across the board for all methods, indicating further room for improvement on this challenging task.

Real-world Affordance Prediction on a Robot:

We conduct robotic manipulation experiments with various objects using a bimanual Tiago++ robot in a realistic kitchen environment. We deploy our 2HandedAfforder model for affordance region segmentation inference based on task prompts such as ‘pour into cup’. Integrating our affordance prediction into the grasping pipeline leads to greater robot task success. Examples of different manipulation tasks are shown in Figure 4 and in videos at sites.google.com/view/2handedafforder.

References

- [1] Shikhar Bahl, Russell Mendonca, Lili Chen, Unnat Jain, and Deepak Pathak. Affordances from human videos as a versatile representation for robotics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13778–13790, 2023. 2
- [2] Matthew Chang, Aditya Prakash, and Saurabh Gupta. Look ma, no hands! agent-environment factorization of egocentric videos. *Advances in Neural Information Processing Systems*, 36, 2024. 2
- [3] Ho Kei Cheng and Alexander G Schwing. Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model. In *European Conference on Computer Vision*, pages 640–658. Springer, 2022. 2
- [4] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Antonino Furnari, Jian Ma, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100. *International Journal of Computer Vision (IJCV)*, 130:33–55, 2022. 2, 4
- [5] Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma, Amlan Kar, Richard Higgins, Sanja Fidler, David Fouhey, and Dima Damen. Epic-kitchens visor benchmark: Video segmentations and object relations. In *Proceedings of the Neural Information Processing Systems (NeurIPS) Track on Datasets and Benchmarks*, 2022. 2
- [6] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, and Huong Ngo et al. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models, 2024. 1
- [7] James J Gibson. The theory of affordances:(1979). In *The people, place, and space reader*, pages 56–60. Routledge, 2014. 1
- [8] Mohit Goyal, Sahil Modi, Rishabh Goyal, and Saurabh Gupta. Human hands as probes for interactive object understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3293–3303, 2022. 2
- [9] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022. 4
- [10] Andrew Guo, Bowen Wen, Jianhe Yuan, Jonathan Tremblay, Stephen Tyree, Jeffrey Smith, and Stan Birchfield. Handal: A dataset of real-world manipulable object categories with pose annotations, affordances, and reconstructions. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 11428–11435. IEEE, 2023. 1
- [11] Ju He, Shuo Yang, Shaokang Yang, Adam Kortylewski, Xiaoding Yuan, Jie-Neng Chen, Shuai Liu, Cheng Yang, Qihang Yu, and Alan Yuille. Partimagenet: A large, high-quality dataset of parts. In *European Conference on Computer Vision*, pages 128–145. Springer, 2022. 1
- [12] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 3
- [13] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023. 3
- [14] Franziska Krebs and Tamim Asfour. A bimanual manipulation taxonomy. *IEEE Robotics and Automation Letters*, 7(4):11031–11038, 2022. 2
- [15] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589, 2024. 1, 3, 4
- [16] Jaewook Lee, Andrew D. Tjahjadi, Jiho Kim, Junpu Yu, Minji Park, Jiawen Zhang, Yang Li, Sieun Kim, XunMei Liu, Jon E. Froehlich, Yapeng Tian, and Yuhang Zhao. Cookar: Affordance augmentations in wearable ar to support kitchen tool interactions for people with low vision. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*, 2024. 1
- [17] Gen Li, Varun Jampani, Deqing Sun, and Laura Sevilla-Lara. Locate: Localize and transfer object parts for weakly supervised affordance grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10922–10931, 2023. 4
- [18] Gen Li, Deqing Sun, Laura Sevilla-Lara, and Varun Jampani. One-shot open affordance learning with foundation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3086–3096, 2024. 1
- [19] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7061–7070, 2023. 3
- [20] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 3
- [21] Hongchen Luo, Wei Zhai, Jing Zhang, Yang Cao, and Dacheng Tao. Learning affordance grounding from exocentric images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2252–2261, 2022. 1, 3
- [22] Anh Nguyen, Dimitrios Kanoulas, Darwin G Caldwell, and Nikos G Tsagarakis. Object-based affordances detection with convolutional neural networks and dense conditional random fields. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5908–5915. IEEE, 2017. 1
- [23] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, and et al. Gpt-4 technical report, 2024. 1

- [24] Shengyi Qian and David F Fouhey. Understanding 3d object interaction from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 21753–21763, 2023. [1](#) [3](#)
- [25] Shengyi Qian, Weifeng Chen, Min Bai, Xiong Zhou, Zhuowen Tu, and Li Erran Li. Affordancellm: Grounding affordance from vision language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7587–7597, 2024. [1](#) [4](#)
- [26] 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*, pages 8748–8763. PMLR, 2021. [1](#)
- [27] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. [2](#)
- [28] Andranik Sargsyan, Shant Navasardyan, Xingqian Xu, and Humphrey Shi. Mi-gan: A simple baseline for image inpainting on mobile devices. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7335–7345, 2023. [2](#)
- [29] Peize Sun, Shoufa Chen, Chenchen Zhu, Fanyi Xiao, Ping Luo, Saining Xie, and Zhicheng Yan. Going denser with open-vocabulary part segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15453–15465, 2023. [1](#)