

# Training for X-Ray Vision: Amodal Segmentation, Amodal Content Completion, and View-Invariant Object Representation from Multi-Camera Video

Moore, Alexander \*  
moore278@llnl.gov

Saini, Amar \*  
saini5@llnl.gov

Cancilla, Kylie  
cancilla5@llnl.gov

Poland, Doug  
poland1@llnl.gov

Carrano, Carmen  
carrano2@llnl.gov

July 2, 2025

## Abstract

Amodal segmentation and amodal content completion require using object priors to estimate occluded masks and features of objects in complex scenes. Recent amodal segmentation work has introduced using temporal features to enrich object representations for improved amodal segmentation and enforce the concept of object permanence and temporal consistency in amodal video segmentation models which modal video object segmentation lacks. Until now, no data has provided an additional dimension for object context: the possibility of multiple cameras sharing a view of a scene. We introduce *MOVi-MC-AC: Multiple Object Video with Multi-Cameras and Amodal Content*, the largest amodal segmentation and first amodal content dataset to date. Cluttered scenes of generic household objects are simulated in multi-camera video. MOVi-MC-AC contributes to the growing literature of object detection, tracking, and segmentation by including two new contributions to the deep learning for computer vision world. Multiple Camera (MC) settings where objects can be identified and tracked between various unique camera perspectives are rare in both synthetic and real-world video. We introduce a new complexity to synthetic video by providing consistent object ids for detections and segmentations between both frames and multiple cameras each with unique features and motion patterns on a single scene. Amodal Content (AC) is a reconstructive task in which models predict the appearance of target objects through occlusions. In the amodal segmentation literature, some datasets have been released with amodal detection, tracking, and segmentation labels. However, to date no dataset has provided ground-truth amodal content labels. While other methods rely on slow cut-and-paste schemes to generate amodal content pseudo-labels, they do not account for natural occlusions present in the modal masks. MOVi-MC-AC provides labels for 5.8 million object instances, setting a new maximum in the amodal dataset literature, along with being the first to provide ground-truth amodal content. The full dataset is available at <https://huggingface.co/datasets/Amar-S/MOVi-MC-AC>

## 1 Introduction

<sup>1</sup>The ability to conceive of whole objects from glimpses at parts of objects is called gestalt psychology [1]. The shape and size of object bounding boxes and masks in video may rapidly change as objects undergo changes in position or occlusion through time. Tracking [2, 3], video object segmentation [4, 5, 6], object retrieval and re-identification [7, 8], and video inpainting [9, 10] could benefit from consistent object representations which maintain a cohesive object view invariant of occlusion, representation, or perspective change [11]. Amodal segmentation and content completion are vital in real-world applications of machine learning requiring consistent object understanding and object permanence through complex video such as robots and autonomous driving [12, 13, 14]. Monocular image amodal segmentation models [15, 16, 17, 18, 19] rely on object priors to estimate occluded object size and shape through obscurations. Recent monocular

\*Equal Contribution.

<sup>0</sup>This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344. Release number: LLNL-JRNL-2007808

	MOVi-MC-AC (Ours)	MOVi-Amodal (Amazon)	SAIL-VOS 3D	SAIL-VOS	COCOA	COCOA-cls	D2S	DYCE
<b>Statistics</b>								
Image or Video	Video	Video	Video	Video	Image	Image	Image	Image
Synthetic or Real	Synthetic	Synthetic	Synthetic	Synthetic	Real	Real	Real	Synthetic
Number of Video Scenes	2041	838	203	201	-	-	-	-
Number of Scene Images	293,904	20,112	237,611	111,654	5,073	3,499	5,600	5,500
Number of Classes	1,033	930	178	162	-	80	60	79
Number of Instances	5,899,104	295,176	3,460,213	1,896,296	46,314	10,562	28,720	85,975
Number of Occluded Instances	4,089,229	247,565	-	1,653,980	28,106	5,175	16,337	70,766
Average Occlusion Rate	45.2%	52.0%	-	56.3%	18.8%	10.7%	15.0%	27.7%
<b>Provided Modalities</b>								
Scene-Level RGB Frames	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Modal Object Masks	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Model Object RGB Content	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Scene-Level (Modal) Depth Masks	Yes	Yes	Yes	Yes	No	No	No	No
Amodal Object Masks	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Amodal Object RGB Content	Yes	No	No	No	No	No	No	No
Amodal Object Depth Masks	Yes	No	No	No	No	No	No	No
Multiple Cameras (multi-view)	Yes	No	No	No	No	No	No	No
Scene-object descriptors (instance re-id)	Yes	Yes	No	No	No	No	No	No

Table 1: A comparison of contemporary datasets for amodal video object segmentation, tracking, and amodal content completion. While MOVi-MC-AC does not contain articulating objects, it is by far the largest dataset which provides complete modal masks, amodal masks, and amodal content for all objects in every scene as well as six distinct camera views with unique camera extrinsics and motion patterns with a united object identification label enabling new directions of research in computer vision.

video amodal segmentation models [20, 21] use context from temporally distant video features to estimate amodal segmentation masks across time by exploiting object priors of video diffusion models trained on synthetic data. So far, no existing research has investigated using multi-view images and video to generate consistent object representations for the purpose of amodal segmentation. We further develop this research area to introduce multi-view video amodal content completion, a new task in which object visuals are estimated through occlusion using both temporal context as well as multi-view information. We release the first dataset to contain ground-truth amodal segmentation masks for all objects in the scene as well as ground-truth amodal content (or the visible "x-ray" view) of all objects in every scene.

## 1.1 Contributions

We make the following contributions to the amodal segmentation and amodal content completion challenge in deep learning for computer vision:

1. We release MOVi Multi-Camera Amodal Content, the first dataset to include complete ground-truth annotations for amodal content, amodal masks, amodal detections, masklets and tracks of obscured objects with over 5 million instances.
2. We introduce the task of multicamera video amodal content completion, including new metrics adapted from the image reconstruction literature to measure the ability of content completion models to correctly predict the shape and appearance of occluded object regions.

## 2 Related Work

MOVi-MC-AC is built to fulfill the needs of many tasks in computer vision and object perception for robotics, as well as introduce new tasks to push the field into view-invariant, object-centric representations of objects

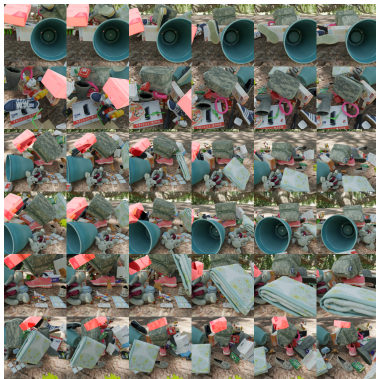


(a) Camera 1.

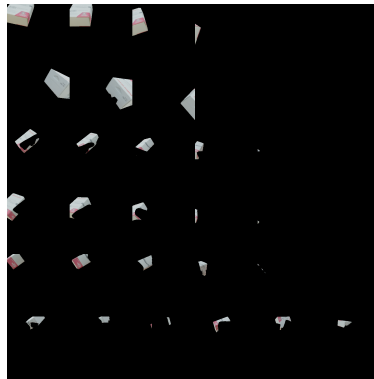


(b) Camera 6.

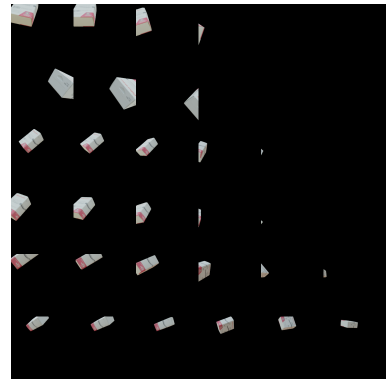
Figure 1: Amodal content completion from multiple cameras must leverage temporal information from one camera view as well as multiple camera perspectives to most accurately predict the visual features of highly-occluded objects from many perspectives simultaneously. MOVi-MC-AC is the first dataset to include ground truth amodal content of occluded objects in video as well as the option to utilize information from multiple cameras on the same scene. Amodal video segmentation and amodal content completion models share information between temporal and camera features to estimate the unoccluded view of objects. Rows: (1) RGB, (2) Modal/Amodal Masks, (3) Amodal Content (4) Overlay. Each column is a new frame in the video.



(a) Up to six cameras observe the same scene with different perspectives and motion characteristics.



(b) Object visible features are partially observed by each camera from one perspective.



(c) Amodal object content is the ground-truth unobscured view of the object generated.

Figure 2: Amodal content completion from multiple cameras must leverage temporal information from one camera view as well as multiple camera perspectives to most accurately predict the visual features of highly-occluded objects from many perspectives simultaneously. MOVi-MC-AC is the first dataset to include ground truth amodal content of occluded objects in video as well as the option to utilize information from multiple cameras on the same scene. Amodal video segmentation and amodal content completion models share information between temporal and camera features to estimate the unoccluded view of objects.

in cluttered video scenes. We introduce a brief survey of relevant tasks in segmentation and amodal content completion.

**Image amodal segmentation.** Amodal image segmentation predicts object masks invariant of occlusion [22, 19, 16, 17, 23, 18, 13, 19, 24, 25, 26]. Image amodal segmentation models require strong object priors to complete object shapes in the absence of temporal context cues. In the presence of temporal context such as real-world video, image amodal segmentation models underperform video amodal segmenters [20].

**Amodal content completion.** Amodal content completion is relevant to neuroscience in researching how the human mind completes incomplete patterns [27]. Amodal completion uses object and shape priors in trained diffusion models to estimate occluded content of various kinds [28, 9, 10]. Scene reconstruction [25] can be performed by decomposing scenes into a collection of de-occluded amodal objects. Extending the object shape priors of image amodal segmentation models, amodal content completion of image and videos requires object appearance priors or temporal context to accurately estimate visual characteristics of occluded objects.

**Amodal video segmentation.** Contemporary state-of-the-art video object segmentation (VOS) [4, 5, 6] segments objects from frame-to-frame through affinity or writing to a memory buffer of references but are prone to identity switching and track loss at low frame rates due to occlusions causing rapid changes in the object mask which modifies the memory encoding [29]. Amodal video segmentation benefits spatiotemporal stability of object tracking by maintaining consistent object mask representations through video regardless of occlusion [11, 30, 31].

**Amodal video content completion.** To the author’s knowledge, Diffusion-VAS [20] and TACO [21] are the only video amodal content completion models currently published. Diffusion-VAS [20] uses a three-stage process to estimate depth from monocular video before an amodal segmentation model estimates the amodal mask of the target object using the modal mask and depth before a final pass uses the amodal mask estimate and RGB video to estimate the object amodal content, all powered by a video diffusion model. TACO [21] trains a latent video diffusion model to denoise representations of video amodal content in a VAE latent.

**Multi-camera deep learning.** Multiple-camera paradigms are common in autonomous driving [32], robotics [33], and person re-identification [34, 35, 36, 37] settings. Until now, no dataset has provided full segmentation masks on a multi-camera video corpus. This emphasis on consistent object representations through multiple cameras on the same cluttered scene of generic objects may benefit downstream tasks needing persistence of object representations. Multiple cameras introduces new challenges in object detection, tracking, re-identification as multiple camera perspectives can learn shared representations of objects with view invariance [38].

### 3 Data

MOVi Multi-Camera Amodal Content (MOVi-MC-AC) is a collection of 2041 scenes split into a 1651-scene training set and a 390-scene testing set. A scene is a collection of videos containing 6 cameras with unique motion characteristics sampled from static, linear motion, or a moving arc with a camera tracking the middle of the scene [39]. Videos are 2-second simulations in which 24 frames are collected. Each scene contains 2 to 40 objects. Of these, 1 to 20 objects are static on the floor, while 1 to 20 objects are dynamically thrown through the air, sometimes causing extreme obscurations of camera views. The training and testing set contain disjoint sets of objects, enabling re-identification of generic objects to be testable with unseen object classes. Considering all scenes, cameras, and objects, there are approximately 20 million files in this dataset.

**Annotations.** We provide two kinds of annotations: scene-level and object-level. Scene-level annotations for each camera include the RGB video data collected by the camera, the instance segmentation mask for all objects seen by the camera, and the ground-truth scene depth image. Object-level annotations include the unoccluded amodal RGB content used as the target in amodal content completion, the unoccluded amodal segmentation mask used as the target in amodal image segmentation [1], amodal video segmentation [21], amodal instance segmentation [17], amodal panoptic segmentation [16], and the unoccluded depth of the target object which is currently not used as a feature or target in any amodal task. We also provide scene-object descriptors, which associates objects to all relevant scenes that they exist in. Visibility and obscuration rates are also provided within these descriptors.

MOVi-MC-AC is unique among other segmentation datasets where full annotations are provided for **each**

object in each scene. This means MOVi-MC-AC is appropriate for instance and semantic segmentation as well as training model and amodal detectors with information from multiple cameras, whereas VOS and amodal segmentation datasets generally are not (see Section 5) [6, 13].

Enabled tasks include image segmentation, video object segmentation (VOS), object detection and classification, video object tracking, object re-identification across views or between cameras, and object re-identification across scenes. Along with these common tasks in computer vision research, we also enable new **amodal** tasks by providing ground-truth labels which no dataset has released until now, including object-based retrieval using amodal content, amodal object detection, amodal 3D detection, amodal video object tracking, multicamera amodal segmentation, and amodal content completion.

Our provided scene-object descriptors enables development of object re-id & retrieval pipelines. Finally, from the MOVi dataset engine, we also have access to object names and meta-class categories with descriptions, further supporting research in grounded/referring tracking, and video object segmentation as in language-based detection [12].

Table 1 compares contemporary datasets relevant to amodal segmentation. We note that no dataset until now supports amodal content prediction natively. Contemporary models use synthetic dataset generation through layering occluded object masks over the target object until a desired occlusion level is achieved [20], **MOVi-MC-AC is the first dataset to provide ground-truth amodal content.**

## 4 Metrics

To accompany our new proposed task of multiple-camera video object amodal content prediction, we introduce the following metrics derived from contemporary amodal video object segmentation and computer vision image reconstruction literature.

### 4.1 Amodal Segmentation Metrics

Amodal segmentation accuracy can be quantified by segmentation metrics between the predicted and ground truth amodal mask of an occluded object.

Following common practice in amodal segmentation, we propose using the mIoU and mIoU<sub>occ</sub> as evaluation metrics for amodal mask predictions [15, 40, 41]. Given videos including the modal and amodal masks for the target object, where the ground-truth modal mask is  $M_i$ , and the predicted and ground-truth amodal masks are  $\hat{A}_i$  and  $A_i$ , the mIoU metric is given as:

$$\text{IoU} = \frac{\hat{A}_i \cap A_i}{\hat{A}_i \cup A_i}$$

and mIoU<sub>occ</sub> as:

$$\text{mIoU}_{\text{occ}} = \frac{(\hat{A}_i - M_i) \cap (A_i - M_i)}{(\hat{A}_i - M_i) \cup (A_i - M_i)}$$

following Diffusion-VAS [20].

### 4.2 Image Reconstruction Metrics

MOVi-MC-AC is the first dataset to provide ground-truth amodal content labels for direct training on occluded content reconstruction. In order to measure content completion accuracy of models, we will adapt the following existing metrics from image reconstruction:

**PSNR metric.** The Peak Signal-to-Noise Ratio is a log-distance between the mean squared error of a target image  $x_1$  and its predicted reconstruction  $x_2$  for images scaled to a range such that the maximum is  $v$ :

$$\text{PSNR}(x_1, x_2) = 10 \log_{10} \frac{v^2}{\text{MSE}(x_1, x_2)} \quad (1)$$

**LPIPS metric.** The Learned Perceptual Image Patch Similarity [42] uses VGG or AlexNet as embeddings to extract features at multiple scales then measures the distance between the target  $x_1$  and reconstruction  $x_2$  for activations at layer  $l$ :

$$LPIPS(x_1, x_2) = \sum_l |f_l(x_1) - f_l(x_2)|^2 \quad (2)$$

**SSIM metric.** Structural Similarity Index Measure [43] compares the mean luminance, contrast, and structure between two images, where  $x$  and  $y$  are the two image patches being compared (e.g., a reference image and a distorted image). SSIM uses  $x$  and  $y$  are the two image patches being compared (e.g., a reference image and a distorted image),  $\mu_x$  and  $\mu_y$  are the mean luminance of  $x$  and  $y$ , respectively,  $\sigma_x^2$  and  $\sigma_y^2$  are the variances of  $x$  and  $y$ ,  $\sigma_{xy}$  is the covariance between  $x$  and  $y$ , and  $C_1$  and  $C_2$  are small constants to stabilize the division when the denominator is close to zero:

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)} \quad (3)$$

### 4.3 Occluded Segmentation and Reconstruction Metrics

For each image reconstruction metric, we apply the function on the ground-truth amodal mask and estimated amodal mask, or the ground truth amodal content and the estimated amodal content. In addition to the entire amodal mask and content, amodal segmentation literature also introduces  $mIOU_{occ}$ , the segmentation performance taken only in the occluded regions to prevent influence of the easy-to-predict visible mask. The suffix “ $_{occ}$ ” applied to each of our amodal segmentation and amodal content completion metrics means to first subset the amodal prediction to only the occluded region using the modal mask, before applying the metric on the occluded regions of the ground truth and prediction.

Given amodal mask  $A$  and modal mask  $V$ , the occluded mask  $O$  is given by  $A \setminus V$ . Using  $O$  we can generate subimages  $I'_1$  and  $I'_2$  by applying the ground-truth occluded mask over the ground-truth amodal object content:

$$\begin{aligned} I'_1 &= I_1 \cdot O \\ I'_2 &= I_2 \cdot O \end{aligned} \quad (4)$$

Using these two masked images, we can compute ground-truth amodal content quantitative metrics on occluded regions:

$$Metric_{occ} = Metric(I'_1, I'_2) \quad (5)$$

## 5 Future Work

MOVi-MC-AC enables a wide range of new tasks in computer vision for the detection, tracking, and segmentation of multiple objects across camera views in cluttered scenes. We propose the following tasks as open challenges enabled by MOVi-MC-AC, which could not be directly trained for until now:

1. **Multi-camera object detection and tracking.** Multiple Object Tracking (MOT) algorithms could be adapted to detect an object in multiple camera views simultaneously with the benefit of greater spatial and viewpoint contexts and assign it a consistent unique object id. Then persistently track the object through the multiple views through video.
2. **Multi-scene object retrieval.** Given an attention prompt (such as a bounding box or segmentation mask) in one video or one multi-camera scene, retrieve the object as a detection in a new camera view of the same scene to unite the camera perspectives. One could further detect the object in a new scene with new cameras and object clutter context. This equates to a view-invariant object representation learnable through multiple cameras to retrieve the object from a gallery of new view points in cluttered scenes.

## References

- [1] Ege Ozguroglu, Ruoshi Liu, Dídac Surís, Dian Chen, Achal Dave, Pavel Tokmakov, and Carl Vondrick. pix2gestalt: Amodal segmentation by synthesizing wholes, 2024.
- [2] Nir Aharon, Roy Orfaig, and Ben-Zion Bobrovsky. Bot-sort: Robust associations multi-pedestrian tracking, 2022.
- [3] Yifu Zhang, Peize Sun, Yi Jiang, Dongdong Yu, Fucheng Weng, Zehuan Yuan, Ping Luo, Wenyu Liu, and Xinggang Wang. Bytetrack: Multi-object tracking by associating every detection box, 2022.
- [4] Ho Kei Cheng and Alexander G. Schwing. Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model, 2022.
- [5] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Joon-Young Lee, and Alexander Schwing. Putting the object back into video object segmentation, 2024.
- [6] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos, 2024.
- [7] Mang Ye, Shuoyi Chen, Chenyue Li, Wei-Shi Zheng, David Crandall, and Bo Du. Transformer for object re-identification: A survey, 2024.
- [8] Dongchen Han, Baodi Liu, Shuai Shao, Weifeng Liu, and Yicong Zhou. Feature aggregation and connectivity for object re-identification. *Pattern Recognition*, 157:110869, 2025.
- [9] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models, 2022.
- [10] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis, 2023.
- [11] Douglas Poland and Amar Saini. Seeing objects in a cluttered world: Computational objectness from motion in video, 2024.
- [12] Chilam Cheang, Haitao Lin, Yanwei Fu, and Xiangyang Xue. Learning 6-dof object poses to grasp category-level objects by language instructions, 2022.
- [13] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3354–3361, 2012.
- [14] Xuelin Qian, Li Wang, Yi Zhu, Li Zhang, Yanwei Fu, and Xiangyang Xue. Impdet: Exploring implicit fields for 3d object detection, 2022.
- [15] Ke Fan, Jingshi Lei, Xuelin Qian, Miaopeng Yu, Tianjun Xiao, Tong He, Zheng Zhang, and Yanwei Fu. Rethinking amodal video segmentation from learning supervised signals with object-centric representation, 2023.
- [16] Rohit Mohan and Abhinav Valada. Amodal panoptic segmentation, 2022.
- [17] Minh Tran, Khoa Vo, Kashu Yamazaki, Arthur Fernandes, Michael Kidd, and Ngan Le. Aisformer: Amodal instance segmentation with transformer, 2024.
- [18] Yuting Xiao, Yanyu Xu, Ziming Zhong, Weixin Luo, Jiawei Li, and Shenghua Gao. Amodal segmentation based on visible region segmentation and shape prior, 2020.
- [19] Lei Ke, Yu-Wing Tai, and Chi-Keung Tang. Deep occlusion-aware instance segmentation with overlapping bilayers, 2021.

- [20] Kaihua Chen, Deva Ramanan, and Tarasha Khurana. Using diffusion priors for video amodal segmentation, 2024.
- [21] Ruijie Lu, Yixin Chen, Yu Liu, Jiaxiang Tang, Junfeng Ni, Diwen Wan, Gang Zeng, and Siyuan Huang. Taco: Taming diffusion for in-the-wild video amodal completion. *arXiv preprint arXiv:2503.12049*, 2025.
- [22] Patrick Follmann, Rebecca König, Philipp Härtinger, and Michael Klostermann. Learning to see the invisible: End-to-end trainable amodal instance segmentation, 2018.
- [23] Minh Tran, Khoa Vo, Tri Nguyen, and Ngan Le. Amodal instance segmentation with diffusion shape prior estimation, 2024.
- [24] Ke Li and Jitendra Malik. Amodal instance segmentation, 2016.
- [25] Xiaohang Zhan, Xingang Pan, Bo Dai, Ziwei Liu, Dahua Lin, and Chen Change Loy. Self-supervised scene de-occlusion, 2020.
- [26] Lu Qi, Li Jiang, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Amodal instance segmentation with kins dataset. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3009–3018, 2019.
- [27] Jordy Thielen, Sander E. Bosch, Tessa M. van Leeuwen, Marcel A. J. van Gerven, and Rob van Lier. Neuroimaging findings on amodal completion: A review. *i-Perception*, 10(2):2041669519840047, 2019. PMID: 31007887.
- [28] Katherine Xu, Lingzhi Zhang, and Jianbo Shi. Amodal completion via progressive mixed context diffusion, 2023.
- [29] Clayton Bromley, Alexander Moore, Amar Saini, Douglas Poland, and Carmen Carrano. Addressing issues with working memory in video object segmentation, 2024.
- [30] N Dinesh Reddy, Robert Tamburo, and Srinivasa G. Narasimhan. Walt: Watch and learn 2d amodal representation from time-lapse imagery. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9346–9356, 2022.
- [31] Jasmin Breitenstein, Franz Jünger, Andreas Bär, and Tim Fingscheidt. Foundation models for amodal video instance segmentation in automated driving, 2024.
- [32] Ahmed Rida Sekkat, Yohan Dupuis, Varun Ravi Kumar, Hazem Rashed, Senthil Yogamani, Pascal Vasseur, and Paul Honeine. Synwoodscape: Synthetic surround-view fisheye camera dataset for autonomous driving. *IEEE Robotics and Automation Letters*, 7(3):8502–8509, July 2022.
- [33] Aron Schmied, Tobias Fischer, Martin Danelljan, Marc Pollefeys, and Fisher Yu. R3d3: Dense 3d reconstruction of dynamic scenes from multiple cameras, 2023.
- [34] Tatjana Chavdarova, Pierre Baqué, Stéphane Bouquet, Andrii Maksai, Cijo Jose, Louis Lettry, Pascal Fua, Luc Van Gool, and François Fleuret. The wildtrack multi-camera person dataset, 2017.
- [35] Xiaotian Han, Quanzeng You, Chunyu Wang, Zhizheng Zhang, Peng Chu, Houdong Hu, Jiang Wang, and Zicheng Liu. Mmptrack: Large-scale densely annotated multi-camera multiple people tracking benchmark, 2021.
- [36] Temitope Ibrahim Amosa, Patrick Sebastian, Lila Iznita Izhar, Oladimeji Ibrahim, Lukman Shehu Ayinla, Abdulrahman Abdullah Bahashwan, Abubakar Bala, and Yau Alhaji Samaila. Multi-camera multi-object tracking: A review of current trends and future advances. *Neurocomputing*, 552:126558, 2023.
- [37] James M. Ferryman and Ali Shahrokni. Pets2009: Dataset and challenge. *2009 Twelfth IEEE International Workshop on Performance Evaluation of Tracking and Surveillance*, pages 1–6, 2009.



- [38] Yang You, Yixin Li, Congyue Deng, Yue Wang, and Leonidas Guibas. Multiview equivariance improves 3d correspondence understanding with minimal feature finetuning, 2025.
- [39] Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J Fleet, Dan Gnanapragasam, Florian Golemo, Charles Herrmann, Thomas Kipf, Abhijit Kundu, Dmitry Lagun, Issam Laradji, Hsueh-Ti (Derek) Liu, Henning Meyer, Yishu Miao, Derek Nowrouzezahrai, Cengiz Oztireli, Etienne Pot, Noha Radwan, Daniel Rebain, Sara Sabour, Mehdi S. M. Sajjadi, Matan Sela, Vincent Sitzmann, Austin Stone, Deqing Sun, Suhani Vora, Ziyu Wang, Tianhao Wu, Kwang Moo Yi, Fangcheng Zhong, and Andrea Tagliasacchi. Kubric: a scalable dataset generator. 2022.
- [40] Jianxiong Gao, Xuelin Qian, Yikai Wang, Tianjun Xiao, Tong He, Zheng Zhang, and Yanwei Fu. Coarse-to-fine amodal segmentation with shape prior, 2023.
- [41] Jian Yao, Yuxin Hong, Chiyu Wang, Tianjun Xiao, Tong He, Francesco Locatello, David Wipf, Yanwei Fu, and Zheng Zhang. Self-supervised amodal video object segmentation, 2022.
- [42] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric, 2018.
- [43] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image Quality Assessment: From Error Visibility to Structural Similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, April 2004.