# SAM4D: Segment Anything in Camera and LiDAR Streams

Jianyun Xu[1,†]   Song Wang[1,2,†]   Ziqian Ni[1,†]   Chunyong Hu[1]   Sheng Yang[1,✉]   Jianke Zhu[2]   Qiang Li[1]

[1]Unmanned Vehicle Dept., CaiNiao Inc., Alibaba Group   [2]Zhejiang University
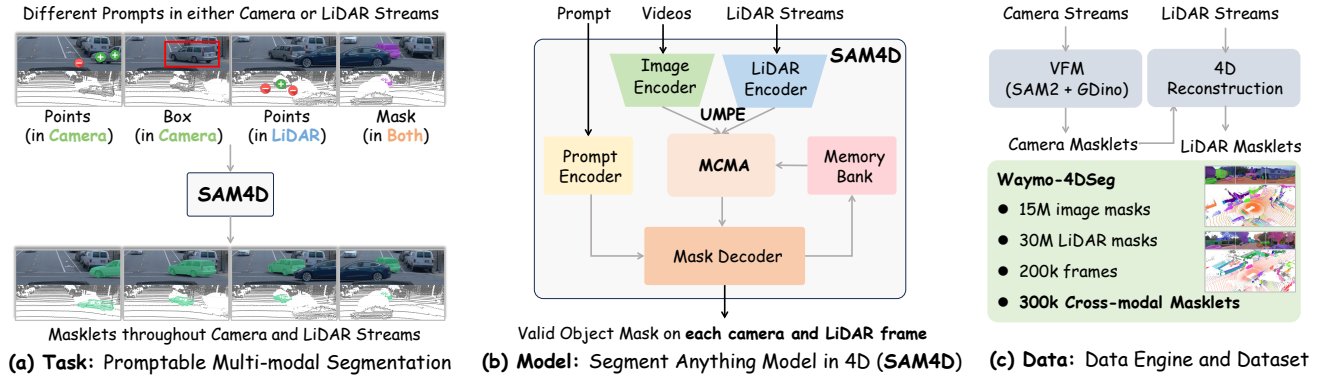
 **Project Page:** SAM4D-Project.github.io

Figure 1. We aim to build a foundation model for 4D segmentation by introducing three interconnected components: (a) a promptable multi-modal segmentation **task**, extending segmentation to both camera and LiDAR streams; (b) a segmentation **model** (SAM4D) that enables 2D-3D joint segmentation with cross-modal prompting and temporal alignment; and (c) an automatic **data** engine for constructing Waymo-4DSeg, a large-scale dataset with over 300k camera-LiDAR associated masklets, providing pseudo labels for SAM4D training.

## Abstract

*We present SAM4D, a multi-modal and temporal foundation model designed for promptable segmentation across camera and LiDAR streams. Unified Multi-modal Positional Encoding (UMPE) is introduced to align camera and LiDAR features in a shared 3D space, enabling seamless cross-modal prompting and interaction. Additionally, we propose Motion-aware Cross-modal Memory Attention (MCMA), which leverages ego-motion compensation to enhance temporal consistency and long-horizon feature retrieval, ensuring robust segmentation across dynamically changing autonomous driving scenes. To avoid annotation bottlenecks, we develop a multi-modal automated data engine that synergizes VFM-driven video masklets, spatiotemporal 4D reconstruction, and cross-modal masklet fusion. This framework generates camera-LiDAR aligned pseudo-labels at a speed orders of magnitude faster than human annotation while preserving VFM-derived semantic fidelity in point cloud representations. We conduct extensive experiments on the constructed Waymo-4DSeg, which demonstrate the powerful cross-modal segmentation ability and*

*great potential in data annotation of proposed SAM4D.*

## 1. Introduction

Segment Anything Model (SAM) [17] has emerged as a foundation model for promptable visual segmentation, demonstrating strong generalization in diverse image domains through user-defined prompts such as points, boxes, and masks. Building on this, SAM2 [37] extends segmentation to videos by incorporating a data engine for large-scale video annotation and a streaming memory mechanism for real-time processing. These advances highlight the potential of promptable segmentation in various downstream tasks [5, 7, 28, 29, 46]. However, existing methods remain limited to the image and video domains without considering other sensor modalities crucial for safety-critical applications such as autonomous driving.

Achieving higher levels of autonomy in driving systems requires robust multi-modal perception [26, 55, 59], where cameras and LiDAR synergistically compensate for each other's limitations, particularly in challenging conditions such as low visibility or poor lighting [11, 52]. Although active depth sensing of LiDAR provides precise geometric

---

†: Equal contribution. ✉: Corresponding author.

1

priors and enables direct temporal feature association, existing segmentation models for LiDAR perception [25, 31, 64] remain largely frame-centric. To the best of our knowledge, no prior work has systematically leveraged cross-modal spatial consistency across synchronized LiDAR scans and camera streams for both 2D and 3D segmentation. These oversights limit the efficacy of joint image and LiDAR annotation, where temporal and cross-modal cues is critical to resolve ambiguities and insufficient observations.

To reduce annotation costs and improve multi-modal segmentation efficiency, we introduce the Promptable Multi-modal Segmentation (PMS) **task**, which enables segmentation across camera and LiDAR sequences based on prompts (*e.g.*, points, boxes, or masks) from both modalities. Furthermore, cross-modal prompting is introduced, allowing a query in one modality (*e.g.*, an image prompt) to guide segmentation in another (*e.g.*, LiDAR).

Based on PMS task, we propose SAM4D, *the first promptable multi-modal segmentation* **model** for camera and LiDAR streams, unifying multi-modal and temporal segmentation within a single framework. Specifically, SAM4D is built upon a multi-modal transformer architecture, integrating Unified Multi-modal Positional Encoding (UMPE) for spatial alignment and Motion-aware Cross-modal Memory Attention (MCMA) for temporal consistency. With UMPE, SAM4D explicitly fuses image and LiDAR features in a shared 3D space, enabling seamless cross-modal prompting and interaction through unified positional encoding. Additionally, MCMA incorporates egomotion compensation, ensuring accurate temporal feature alignment and enhancing long-horizon object tracking in dynamic environments. By integrating multi-modal feature fusion, temporal reasoning, and cross-modal interaction, SAM4D is expected to significantly reduce manual labeling efforts while ensuring robust and temporally consistent segmentation in driving scenarios.

To train SAM4D, we construct Waymo-4DSeg, a large-scale multi-modal segmentation **dataset** based on the Waymo Open Dataset [42], designed to provide high-quality, temporally consistent pseudo-ground truth. Our proposed multi-modal data engine enhances 2D-3D joint annotation by integrating vision foundation model (VFM)-based video masklet generation, 4D reconstruction for LiDAR pseudo-labeling, and cross-modal label fusion. In contrast to previous methods [31, 64] that focus on independent frame annotations, our approach uses a sequence-level propagation strategy, ensuring temporal consistency and cross-modal coherence. This significantly improves annotation efficiency and accuracy, making Waymo-4DSeg a key benchmark for training and evaluating promptable, multi-modal, and temporally aware segmentation models for autonomous driving. Extensive experiments are conducted with Waymo-4DSeg and unseen dataset under different challenging settings, which demonstrate the strong performance and generalizability of SAM4D in promotable multi-modal segmentation.

## 2. Related Work

**Image and LiDAR Segmentation.** Segment Anything (SAM) series [17, 37] introduced a foundation model for image and video segmentation, capable of generating masks based on diverse prompt types while demonstrating strong generalization across various datasets. Subsequent work has focused on improving the granularity of SAM segmentation [16] and efficiency [54, 58, 62], optimizing its performance for finer segmentation tasks. In addition, researchers have explored the applicability of SAM in various 2D downstream tasks [5, 28, 29, 46] and LiDAR point clouds [15, 25, 31, 56]. Liu *et al.* [25] proposed distilling 2D segmentation masks from SAM into LiDAR-based networks to enhance 3D segmentation performance. Other methods [15, 56] project SAM-generated 2D masks into 3D space for further refinement. Approaches such as SAL [31] and PointSAM [64] attempt to build SAM-like promptable segmentation networks directly on point clouds, designing architectures specifically for 3D data. There are also studies [12, 13] that explore the accumulation of sequential LiDAR point clouds for interactive segmentation that leverages temporal consistency by aggregating point cloud information over multiple frames. However, the above methods remain modality-specific, focusing on either image segmentation or LiDAR segmentation in isolation. In contrast, our work is the first to unify image and LiDAR segmentation within a single framework.

**Multi-Modal Perception in Driving.** Recent works have explored multi-modal fusion strategies for Cameras and LiDAR in different spatial representations to improve detection [8, 21, 26, 55, 59], segmentation [4, 18, 30, 41, 50, 51], and occupancy prediction [19, 32, 49, 61] in autonomous systems. Bird's-Eye View (BEV) fusion [22, 26] has gained popularity, where image and LiDAR features are projected into a unified BEV space to facilitate spatial reasoning. Meanwhile, voxel-based fusion [21, 49, 61] operates directly in 3D space, where features of multiple sensors are aggregated into structured voxel grids for fine-grained 3D perception. Sparse representation-based fusion [53, 55, 59] has also gained attention for its efficient feature encoding and wider perceptual coverage, making it a promising direction in multimodal perception research. Despite these advancements, most existing methods output only 3D predictions, lacking exploration of cross-modal interactions and unified 2D-3D segmentation. Our work fills this gap by introducing SAM4D, enabling promptable segmentation across camera and LiDAR streams.

**Lifting Vision Foundation Models for 3D Labeling.** Recent works have leveraged Vision Foundation Models
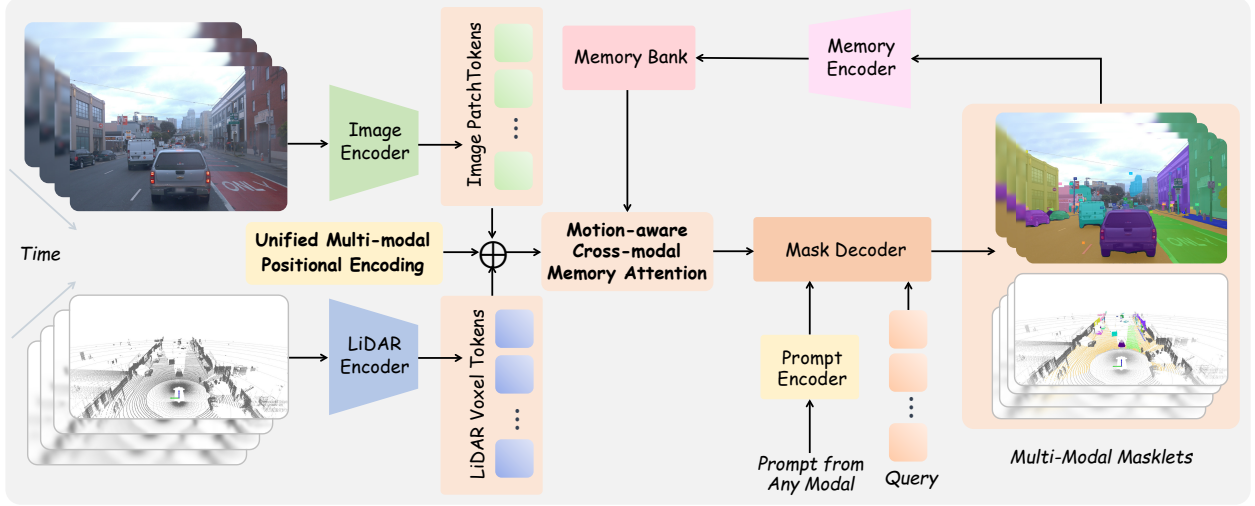
Figure 2. **Overview of the Segment Anything Model in 4D (SAM4D) workflow.** The image and LiDAR encoders generate modality-specific embeddings, which are aligned through the proposed **Unified Multi-modal Positional Encoding**. The **Motion-aware Cross-modal Memory Attention** then processes multi-modal and temporal features, incorporating ego-motion for improved feature interaction. Finally, the updated image and LiDAR features are queried efficiently by mask decoder with diverse input prompts from various modalities.

(VFMs) [17, 36] to enable label-efficient 3D scene understanding by distilling 2D vision priors into 3D representations. Methods like CLIP2Scene [6] and OpenScene [33] transfer CLIP's vision-language embeddings to 3D semantic segmentation, while approaches such as CLIP-FO3D [60] and OVO [43] extend open-vocabulary segmentation to 3D occupancy and point cloud learning. More recent works [1, 48] integrate volume rendering techniques to improve 3D occupancy prediction. Multimodal fusion strategies have also been explored, with VLM2Scene [23] incorporating image, text, and LiDAR representations and VEON [63] enhancing 3D occupancy prediction through vocabulary-driven alignment. Furthermore, methods [27, 57] focus on segmentation of zero-shot point clouds by integrating multimodal visual cues. Although these approaches effectively bridge 2D and 3D feature spaces, they are often limited to individual frames and do not explicitly consider temporal consistency. In contrast, our work introduces a 4D data engine that propagates high-quality labels across entire sequences through temporal reconstruction.

## 3. Promptable Multi-modal Segmentation

The Promptable Multimodal Segmentation (PMS) task is designed to enable *interactive*, *cross-modal*, and *temporal* segmentation across both camera and LiDAR streams. Unlike traditional segmentation tasks that rely on a single modality or frame-by-frame processing, PMS allows prompts in either 2D (images) or 3D (LiDAR point clouds) to guide segmentation across the entire sequence. A prompt can be in the form of positive/negative clicks, boxes, or masks that either define a new object or refine an existing

segmentation. Once a prompt is provided on a specific image frame or LiDAR scan, the model should immediately return valid segmentation masks for both modalities and then propagate the segmentation across the entire sequence, forming masklets that maintain temporal consistency. Additionally, PMS allows users to provide additional prompts at any frame in the sequence, refining segmentation across frames and modalities as needed.

To support PMS, we develop SAM4D, a unified segmentation model capable of processing both videos and LiDAR streams with cross-modal prompting. Additionally, we construct a large-scale dataset based on the Waymo Open Dataset [42] to provide high-quality pseudo-ground truth annotations for PMS. We evaluate SAM4D by simulating interactive multi-modal segmentation scenarios, assessing its ability to segment and track objects across frames and sensor modalities.

## 4. SAM4D Model

### 4.1. Overview

SAM4D extends SAM2 [37] beyond video segmentation to the multimodal domain, addressing the challenges of cross-modal and long-term object segmentation in autonomous driving scenarios. Unified multi-modal positional encoding is proposed to enable the multi-modal feature and prompt interaction. To enhance the ability for long-term object segmentation, we take into account of the ego-motion and design motion-aware cross-modal memory attention. The overview of our proposed SAM4D is illustrated in Fig. 2.

## 4.2. Multi-modal Segmentation Framework

**Multi-modal Feature Embedding.** In the video branch, we follow SAM2 [37] and adopt Hiera [2, 40] with SA-V [37] pre-training to embed each image frame into unconditional patch tokens. In the LiDAR branch, MinkUNet [10], implemented with TorchSparse [44, 45], is utilized to encode sparse point clouds into voxel-level tokens. Throughout the entire interaction process, the image and LiDAR encoder run only once to reduce computational overhead, enabling efficient processing of long-horizon video sequences.

**Motion-aware Cross-modal Memory Attention.** Our memory attention refines feature representations by integrating cross-modal features and previous frame features in memory (see below) to ensure cross-modal and temporal alignment, which is a core component of our method. Unlike SAM2, SAM4D lifts image patches into 3D space via depth estimation, allowing unified positional encoding for image patch tokens and LiDAR voxel tokens (see Sec. 4.3). Furthermore, ego-motion is also embedded in cross-attention with past features and predictions to enable long-term temporal consistency (see Sec. 4.4).

**Prompt Encoder and Mask Decoder.** The prompt encoder supports different input prompts from both the image and LiDAR inputs to define the spatial extent and position of the target. Sparse prompts, such as points and boxes, are represented by positional encoding (see Sec. 4.3) summed with learnable embeddings for each type of prompt, while mask prompts are embedded using convolutions for images and sparse convolutions for LiDAR. The mask decoder processes the prompts from both modalities along with the image and LiDAR features updated by memory attention, simultaneously predicting 2D and 3D segmentation masks.

**Memory Encoder and Memory Bank.** The memory encoder processes both 2D and 3D segmentation masks separately, utilizing convolutions for the image and sparse convolutions for LiDAR to downsample the output. The downsampled masks are then summed element-wise with the initial embeddings from the image and LiDAR encoders, respectively. A lightweight convolutional layer is then applied to fuse the information, generating the final representation in the memory. The memory bank maintains a FIFO queue to store past object features, with up to $N$ unprompted frames being retained. Additionally, a separate FIFO queue stores $M$ prompted frames to preserve keyframes with explicit user input. We store object pointers computed from mask decoder tokens for both the image and LiDAR domains, capturing high-level semantic information of segmented objects and participating in memory attention.

## 4.3. Unified Multi-modal Positional Encoding

To ensure consistent spatial representation for both image and LiDAR modalities, we carefully design a Unified Multi-
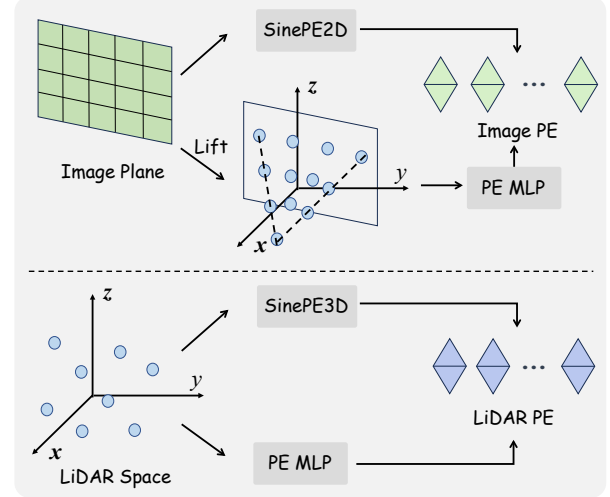


Figure 3. Illustration of the proposed **Unified Multi-modal Positional Encoding**.

modal Positional Encoding (UMPE) scheme. As shown in Fig. 3, this encoding unifies the 2D and 3D features in a shared spatial space, allowing cross-modal interactions while preserving the intrinsic structure of each modality. UMPE consists of two complementary components: (i) a modality-specific positional prior, which encodes features in their native spaces, and (ii) a shared 3D representation, which aligns both modalities in a common spatial domain.

**Positional Encoding for Images.** For a pixel $\mathbf{p} = (u, v)$ in an image feature, we first assign a 2D sinusoidal positional encoding (`SinPE2D`):

$$\mathcal{P}_{\text{img\_sin}} = \texttt{SinPE2D}(u, v), \tag{1}$$

which preserves spatial structure in the image plane. To align image features with LiDAR spatial representations, we estimate a set of depths $D(u, v)$ for each pixel and lift it into 3D space, similar to Lift-Splat-Shoot [34]:

$$\mathbf{x}_{\text{img}} = T_c^l K^{-1} [u * D(u, v), v * D(u, v), D(u, v), 1]^T, \tag{2}$$

where $K \in R^{4 \times 4}$ is the intrinsic matrix of the camera, and $T_c^l \in R^{4 \times 4}$ is the transformation matrix from the camera coordinate to the LiDAR one. This process converts the image into a pseudo-point cloud $\mathbf{x}_{\text{img}}$. We then apply an MLP-based 3D positional encoding:

$$\mathcal{P}_{\text{img\_mlp}} = \texttt{MLP}(\mathbf{x}_{\text{img}}). \tag{3}$$

The final position encoding $\mathcal{P}_{\text{img}}$ is composed of $\mathcal{P}_{\text{img\_sin}}$ and $\mathcal{P}_{\text{img\_mlp}}$, and these two parts ensure that the image features are represented in the same spatial domain as LiDAR while maintaining the original view-based structure.

**Positional Encoding for LiDAR.** For a LiDAR point at $\mathbf{x}_{\text{LiDAR}} = (x, y, z)$, we follow a similar two-stage encoding to obtain $\mathcal{P}_{\text{LiDAR}}$. First, a 3D sinusoidal positional encoding
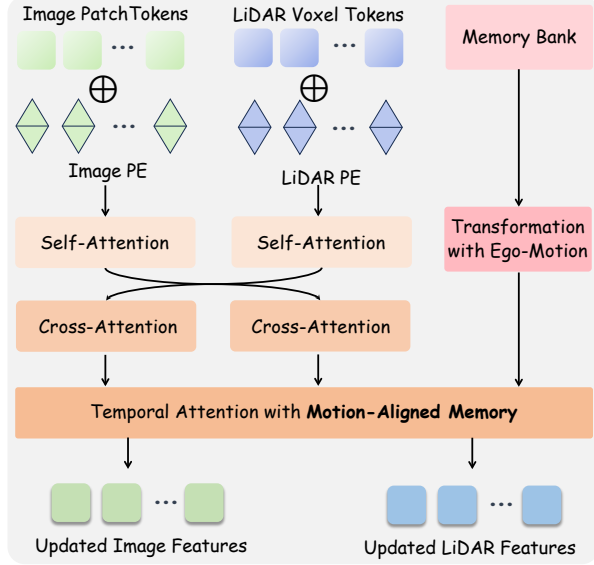
4

Figure 4. Illustration of the proposed **Motion-aware Cross-modal Memory Attention**.

is applied:

$$\mathcal{P}_{\text{LiDAR\_sin}}(x, y, z) = \text{SinPE3D}(x, y, z) \quad (4)$$

which encodes the spatial structure of the point cloud. To ensure consistency with image features lifted into 3D, we utilize the same MLP-based transformation:

$$\mathcal{P}_{\text{LiDAR\_mlp}} = \text{MLP}(\mathbf{x}_{\text{LiDAR}}). \quad (5)$$

For the convenience, we define a symbol $\Phi$ to represent the positional encoding of these two stages for both modalities.

For sparse prompts incluing points or bounding boxes from the image or LiDAR, we apply the same dual-stage positional encoding as used for dense features. The encoded prompts from both modalities are concatenated before being fed into the mask decoder, where missing prompts from one modality are replaced with an empty placeholder. The sparse prompt embeddings are then concatenated to the output tokens, and then applied cross-attention from features updated by the following motion-aware memory attention, enabling the mask decoder to generate both 2D and 3D segmentation masks. By unifying image and LiDAR positional encodings in a shared 3D space, while preserving modality-specific characteristics, UMPE enables further cross-modal feature fusion and interaction in our framework.

## 4.4. Motion-aware Cross-modal Memory Attention

To enhance multi-modal feature representations while ensuring temporal consistency, we introduce Motion-aware Cross-modal Memory Attention (MCMA). This module integrates self-attention, cross-attention across image and LiDAR modalities, and memory-based temporal attention, as

illustrated in Fig. 4. A key distinction from previous approaches [17, 37] is our incorporation of ego-motion compensation, which aligns the features of the past frame to the current coordinate system, allowing for more accurate feature retrieval and reducing errors in dynamically changing autonomous driving scenes.

**Self-Attention for Feature Refinement.** Given the image features $\mathcal{F}_{\text{img}}$ and LiDAR features $\mathcal{F}_{\text{LiDAR}}$ from the respective encoders, along with their positional encodings $\mathcal{P}_{\text{img}}$ and $\mathcal{P}_{\text{LiDAR}}$ obtained from Unified Multi-modal Positional Encoding (UMPE), we first apply self-attention within each modality to refine intra-modal feature representations:

$$\mathcal{F}'_{\text{img}} = \text{SelfAttn}(\mathcal{F}_{\text{img}} + \mathcal{P}_{\text{img}}),$$
$$\mathcal{F}'_{\text{LiDAR}} = \text{SelfAttn}(\mathcal{F}_{\text{LiDAR}} + \mathcal{P}_{\text{LiDAR}}), \quad (6)$$

where $\text{SelfAttn}$ represents the self-attention, allowing each token to attend to others within the same modality.

**Cross-Attention for Multi-modal Fusion.** To facilitate interaction between image and LiDAR features, we perform cross-attention, enabling one modality to incorporate information from the other:

$$\mathcal{F}''_{\text{img}} = \text{CrossAttn}(\mathcal{F}'_{\text{img}}, \mathcal{F}'_{\text{LiDAR}} + \mathcal{P}_{\text{LiDAR}}),$$
$$\mathcal{F}''_{\text{LiDAR}} = \text{CrossAttn}(\mathcal{F}'_{\text{LiDAR}}, \mathcal{F}'_{\text{img}} + \mathcal{P}_{\text{img}}), \quad (7)$$

This step ensures that both modalities share complementary spatial and structural information, enhancing feature expressiveness for segmentation.

**Temporal Attention with Motion-Aligned Memory.** In contrast to SAM2 [37], which only considers short-term object motion, our method explicitly incorporates *ego-motion compensation* to handle large-scale scene changes in autonomous driving scenarios. We maintain a memory bank that stores historical image and LiDAR features $\mathcal{M}_{\text{img}}, \mathcal{M}_{\text{LiDAR}}$ along with their 3D space positions $\mathbf{x}_{\text{img}}$ and $\mathbf{x}_{\text{LiDAR}}$. These features and positions are stored in a FIFO queue, keeping $N$ unprompted frames and $M$ prompted frames for temporal reference.

To correctly align past frame features to the current coordinate frame, we transform stored positions using the ego-motion transformation matrix $T_{t \leftarrow t'}$, which maps historical frame $t'$ to the current frame $t$:

$$\mathcal{M}^{t \leftarrow t'}_{\text{img}} = \mathcal{M}^{t'}_{\text{img}} + \Phi_{\text{img}}(T_{t \leftarrow t'}(\mathbf{x}_{\text{img}})),$$
$$\mathcal{M}^{t \leftarrow t'}_{\text{LiDAR}} = \mathcal{M}^{t'}_{\text{LiDAR}} + \Phi_{\text{LiDAR}}(T_{t \leftarrow t'}(\mathbf{x}_{\text{LiDAR}})), \quad (8)$$

where $T_{t \leftarrow t'} \in SE(3)$ is derived from vehicle odometry, ensuring spatially consistent memory retrieval. $\Phi_{\text{img}}$ and $\Phi_{\text{LiDAR}}$ are the unified multi-modal positional encoding, consisting of a sinusoidal and a MLP embedding.

Similar to SAM2 [37], we also keep the past object pointers $\mathcal{O}^{t'}_{\text{img}}$ and $\mathcal{O}^{t'}_{\text{LiDAR}}$ in memory. Once transformed, previous frame features and object pointers are used in
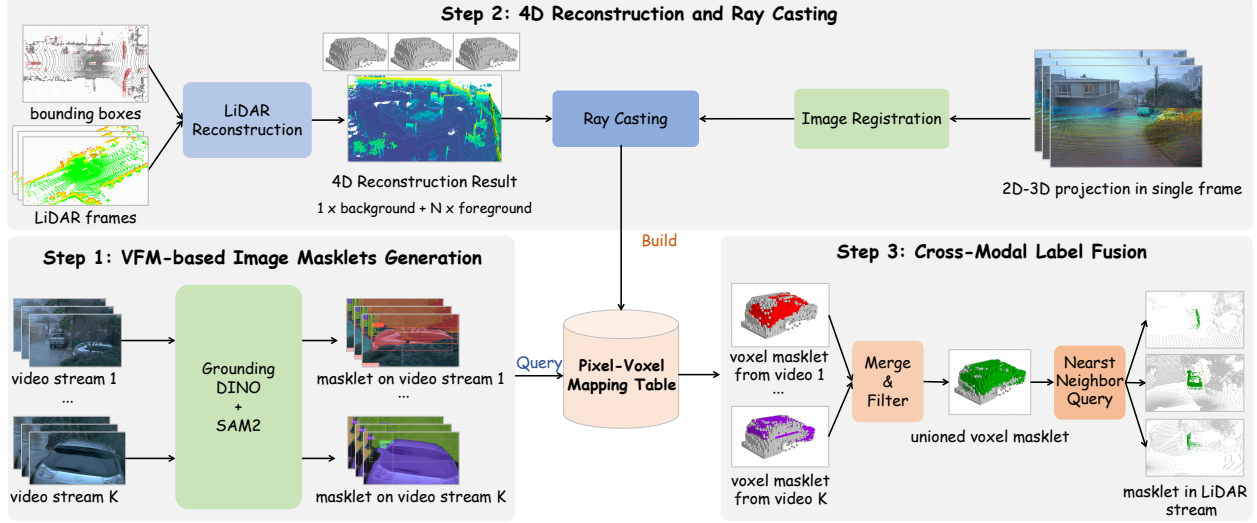
5

Figure 5. **Overview of our data engine**, which is composed of three steps to construct high-quality pseudo labels.

cross-attention to update current features with aligned temporal information:

$$\mathcal{F}_{\text{img}}^{\text{final}} = \text{CrossAttn}(\mathcal{F}_{\text{img}}'', (\mathcal{M}_{\text{img}}^{t \leftarrow t'}, \mathcal{O}_{\text{img}}^{t'})),$$
$$\mathcal{F}_{\text{LiDAR}}^{\text{final}} = \text{CrossAttn}(\mathcal{F}_{\text{LiDAR}}'', (\mathcal{M}_{\text{LiDAR}}^{t \leftarrow t'}, \mathcal{O}_{\text{LiDAR}}^{t'})). \quad (9)$$

By incorporating motion-aware memory alignment, MCMA significantly improves feature consistency across frames, reducing errors in object correspondence caused by large-scale scene changes. This enables SAM4D to perform robust cross-modal and temporal segmentation in dynamic real-world environments.

## 4.5. Training

The SAM4D model is jointly trained on camera and LiDAR sequences with simulated interactive prompting across modalities, following the strategy introduced in the SAM series [17, 37]. Identical loss functions are applied to both image and LiDAR predictions to enforce cross-modal consistency. Additional training details are provided in the supplementary material.

## 5. Data

To the best of our knowledge, there is currently no dataset that simultaneously supports both 2D and 3D segmentation while ensuring instance consistency over time. To quickly establish and expand the training dataset at a low cost, we have carefully designed a multi-modal automatic data engine (Sec. 5.1) to obtain high-quality pseudo-ground-truth data as much as possible. Using this data engine, we construct the Waymo-4DSeg dataset (Sec. 5.2) based on the Waymo Open Dataset [42], providing a large-scale benchmark for multimodal and temporal segmentation.

### 5.1. Data Engine

Our data engine, as shown in Fig. 5, consists of three steps. In Step 1, we leverage vision foundation models (VFM) to generate initial annotations for camera-captured image sequences. Given an image sequence of length $T$, we select keyframes at intervals of $K$ frames. Starting from the first frame, we adopt Grounding-DINO [24, 38], an open-vocabulary detector, and promptable SAM [17, 39] to obtain detections and segmentation masks for common objects in autonomous driving scenes. The generated keyframe masks serve as prompts for SAM2 [37], which propagates the segmentation forward to the next keyframe, producing masklets for the intermediate frames.

In Step 2, we utilize LiDAR frames and pre-annotated 3D bounding boxes of foreground objects to construct a 4D voxel-based reconstruction, serving as an intermediary between image data and LiDAR frames. This 4D reconstruction consists of a single background component and multiple foreground components, each defined in the body coordinate of the object. We also perform exhaustive ray casting from the center of each image toward the voxels to establish a dense pixel-voxel mapping table.

By querying the pixel-voxel mapping table, we can assign the video masklets to the corresponding voxels in Step 3. However, the presence of noise in the labels and masks necessitated the implementation of a filtering step based on a clustering algorithm. We employ the DBSCAN algorithm to cluster voxels according to their BEV positions and selected the main cluster with the highest vote rate while discarding the rest as noise. After filtering, we assessed overlaps between voxel masklets from different videos to merge them into single masklets. Finally, we created a mapping table between points from LiDAR frames and voxels based on their 3D spatial distances, facilitating the transfer of the

6

final voxel masklet to the LiDAR frames. We evaluated the quality of the resulting masklets using cross-modal IoU, which yielded an average score of 0.56.

## 5.2. Constructed Dataset

Our Waymo-4DSeg dataset, derived from the Waymo Open Dataset, follows the original training and validation splits, resulting in 1000 clips (798 for training and 202 for validation), with about 200 frames per clip. On average, we generated 300 masklets for each clip, with each masklet appearing in about 122 frames. This results in an average of 17 masks per image and 170 masks per point cloud. Furthermore, 23.4% of the masklets have been observed in at least two different clips. The semantic categories of our masklets cover nearly all relevant items in autonomous driving scenarios, including dynamic foreground objects (*vehicles*, *pedestrians*), background elements (*buildings*, *trees*) and nearby objects (*curbs*, *lamp posts*, *traffic cones*). The volume of the objects ranges from less than 10 voxels to over 200k voxels (with voxel size of 0.1 meter), occupying an average image area of 1.5k pixels to over 1M pixels. More detailed distribution information and visual results are provided in the supplementary material.

## 6. Experiments

### 6.1. Setup

**Implementation Details.** With the Waymo-4DSeg built, we train our SAM4D model with 6 maximum objects on 16 NVIDIA A100 GPUs for 36 epochs. Unless otherwise mentioned, the experimental results in this section are produced in our default setting, using Hiera-S image encoder [40] with input image in resolution of $768 \times 768$ and Mink-34 LiDAR encoder [10, 44] with input LiDAR points voxelized by size of 0.15m. More details of the implementation can be found in the supplementary material.

**Evaluation Metrics.** In the evaluation, the mean Intersection over Union (mIoU) is adopted to assess segmentation performance for both camera and LiDAR in each single frame. $\mathcal{J}\&\mathcal{F}$ metric [35] in video object segmentation is also reported for image sequences. Additionally, we introduce *Number of Mismatched Predictions* (NMP), which quantify the number of instances where a predicted object fails to match the ground truth with an IoU below 0.01 threshold. This metric captures erroneous associations and misalignments, offering information on the robustness of the model to maintain accurate object correspondence between frames. In the practical evaluation, we randomly selected 48 clips from the validation set.

### 6.2. Main Results

**Promptable Cross-Modal Frame Segmentation.** For objects captured by both the camera and LiDAR, enabling

Table 1. Performance comparison with different prompts on promptable cross-modal frame segmentation.

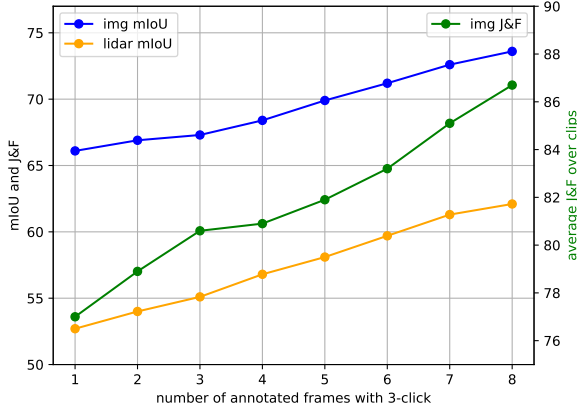| Prompts | Image mIoU (%) ↑ | LiDAR mIoU (%) ↑ |
|---|---|---|
| *Image-Prioritized Prompting* | | |
| 1-click | 68.0 | 42.3 |
| 3-click | 73.6 | **53.1** |
| bounding box | **74.7** | 47.0 |
| *LiDAR-Prioritized Prompting* | | |
| 1-click | 49.6 | 58.8 |
| 3-click | **64.2** | **68.4** |
| bounding box | 46.0 | 63.9 |

Table 2. Performance comparison with different prompts on semi-supervised stream object segmentation.

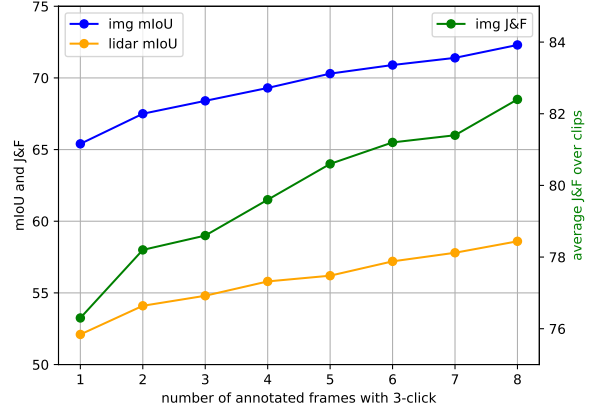| Prompts | Image | | | LiDAR | |
|---|---|---|---|---|---|
| | mIoU (%) ↑ | $\mathcal{J}\&\mathcal{F}$ ↑ (%) | NMP ↓ | mIoU (%) ↑ | NMP ↓ |
| 1-click | 61.4 | 72.2 | 398 | 50.1 | 784 |
| 3-click | 65.6 | 76.3 | 327 | 52.8 | 711 |
| 5-click | 67.1 | 77.7 | 315 | 52.6 | 702 |
| bounding box | 64.5 | 75.4 | 347 | 51.3 | 762 |
| ground-truth mask | **69.8** | **80.1** | **280** | **55.7** | **582** |

segmentation in one modality based on a prompt from the other is crucial for improving efficiency in multimodal annotation. We evaluated this by selecting objects present in both modalities and providing prompts in one single one (image or LiDAR), then measuring the segmentation IoU in both modalities within a single frame. Prompts include single-point (1 click), multi-point (3 clicks), bounding box, and mask inputs. In multiple-click experiments, subsequent clicks are placed on the modality with poorer segmentation after the initial click, simulating human annotation behavior for efficient refinement. As shown in Tab. 1, providing prompts in the image or LiDAR enables the other modality to achieve promising segmentation results, demonstrating the capability of cross-modal prompting in SAM4D.

**Promptable Multi-modal Stream Segmentation.** We further evaluate SAM4D's stream-level promptable segmentation ability, simulating an interactive annotation process. Prompts are given on the first frame where the target appears, with a single-modality prompt if the object is present in one modality, and dual-modality prompts if present in both. Similar to SAM2 [37], experiments are conducted in *offline* and *online* modes. In *offline* mode, segmentation is initialized with 3 point prompts on the first frame, followed by propagation. The frame with the *lowest IoU* is selected for additional prompts, repeating until the prompt limit is reached. In *online* mode, segmentation propagates iteratively, adding prompts to frames where IoU falls below 0.75, until the prompt limit is reached or the sequence ends. As illustrated in Fig. 6, SAM4D achieves stable segmentation performance in both settings, with continuous improvement as additional prompts are introduced.

**Semi-Supervised Stream Object Segmentation.** We extend semi-supervised video object segmentation [9, 35] to

(a) *offline* evaluation (3-click)

(b) *online* evaluation (3-click)

Figure 6. Performance comparison with different promptable frames in interactive offline and online evaluation settings.

Table 3. Performance comparison on nuScenes under semi- supervised stream object segmentation setting.

| nuScenes | Image | | | LiDAR | |
|---|---|---|---|---|---|
| | mIoU (%) ↑ | $\mathcal{J}\&\mathcal{F}$ ↑ (%) | NMP ↓ | mIoU (%) ↑ | NMP ↓ |
| *zero-shot* | 58.4 | 65.8 | 36 | 25.9 | 117 |
| *fine-tuning* | **67.5** | **75.4** | **22** | **44.8** | **70** |

multimodal streams, providing *first frame prompts only* for both image and LiDAR sequences and evaluating segmentation over the full sequence to assess temporal propagation and tracking. The standard mIoU, $\mathcal{J}\&\mathcal{F}$, and NMP are reported for a comprehensive evaluation. As shown in Tab. 2, mask prompts, which encode richer spatial information, achieve the highest segmentation performance across both modalities, outperforming point and box prompts.

**Generalization Experiments on PMS.** We evaluate SAM4D on unseen nuScenes dataset [3, 20] through *zero-shot* transfer and *fine-tuning* which are detailed in the appendix. We adopt the Semi-Supervised Stream Object Segmentation setting, where the mask from the first frame is provided as a prompt to guide the segmentation of subsequent frames. As shown in Tab. 3, SAM4D demonstrates strong *zero-shot* segmentation performance, highlighting its cross-modal generalization to unseen driving scenarios. Further *fine-tuning* on nuScenes enhances segmentation quality, demonstrating the model's ability to adapt and refine predictions in novel environments.

## 6.3. Ablations

We perform ablation studies under the Semi-Supervised Stream Object Segmentation setting to validate the design choices in the SAM4D framework.

**Ablation on Input Modality.** First, we analyze the impact of input modalities by training single-modality variants of SAM4D, where only the image branch (SAM4D-C) or the LiDAR branch (SAM4D-L) is retained, while all other set-

Table 4. Ablation study on the input modality of SAM4D.

| Input | Image | | | LiDAR | |
|---|---|---|---|---|---|
| | mIoU (%) ↑ | $\mathcal{J}\&\mathcal{F}$ ↑ (%) | NMP ↓ | mIoU (%) ↑ | NMP ↓ |
| SAM2+Project | 68.2 | 79.7 | 383 | 32.0 | - |
| SAM4D-C | 68.6 | **80.4** | 301 | - | - |
| SAM4D-L | - | - | - | 47.0 | 799 |
| SAM4D | **69.8** | 80.1 | **280** | **55.7** | **582** |

Table 5. Ablation study on the input resolution of both modalities.

| Resolution | Image | | | LiDAR | |
|---|---|---|---|---|---|
| | mIoU (%) ↑ | $\mathcal{J}\&\mathcal{F}$ ↑ (%) | NMP ↓ | mIoU (%) ↑ | NMP ↓ |
| I-512, V-0.2 | 60.5 | 70.6 | 291 | 48.2 | 985 |
| I-768, V-0.15 | **69.8** | **80.1** | **280** | **55.7** | **582** |

tings are kept the same. Additionally, we introduce a baseline (SAM2+Project), which projects SAM2's video segmentation results onto per-frame point clouds. This method is inherently limited by discrepancies in sensor viewpoint, range, and synchronization between camera and LiDAR. As shown in Tab. 4, the multimodal SAM4D effectively leverages cross-modal interaction and prompting, achieving significantly better segmentation performance compared to its single-modality counterparts.

**Ablation on Input Resolution.** Next, we examine the role of input resolution in the performance of promotable segmentation. Compared to the baseline setting of the image resolution $512 \times 512$ and the resolution of 0.2 m voxels, increasing the resolution to $768 \times 768$ for the images and 0.15 m for the voxels results in a notable performance gain, as presented in Tab. 5. This demonstrates the importance of high-resolution input in dense prediction tasks, where finer spatial details contribute to more accurate segmentation.

**Ablation on Ego-motion in Memory Attention.** Finally, we perform an ablation study on the incorporation of ego-motion in Motion-aware Cross-modal Memory Attention to assess its contribution to temporal feature fusion and object

Table 6. Ablation study on ego-motion in memory attention.

| MCMA | Image | | | LiDAR | |
|------|-------|---|---|-------|---|
| | mIoU (%) ↑ | $\mathcal{J}\&\mathcal{F}$ ↑ (%) | NMP ↓ | mIoU (%) ↑ | NMP ↓ |
| *w/o.* ego-motion | 69.7 | **80.3** | 298 | 52.2 | 746 |
| *w.* ego-motion | **69.8** | 80.1 | **280** | **55.7** | **582** |

tracking. As presented in Tab. 6, ego-motion compensation significantly reduces tracking inconsistencies in stream segmentation, particularly for LiDAR, where NMP decreases from 746 to 592, indicating improved temporal stability. Furthermore, its integration leads to a notable improvement in mIoU, highlighting the importance of ego-motion in enhancing segmentation accuracy over long sequences.

# 7. Conclusion

In this paper, SAM4D is introduced as a multi-modal and temporal model for promptable segmentation across camera and LiDAR streams. Our contributions span **task** (PMS), **model** (SAM4D), and **data** (Waymo-4DSeg). With extensive experiments, SAM4D advances 4D scene understanding, improving segmentation consistency, efficiency, and adaptability in various autonomous driving scenarios. We believe that the insights of the SAM4D model into multi-modal prompting and 4D perception will significantly reduce annotation costs, enabling high-quality scalable 2D-3D joint labeling for large-scale datasets.

## Acknowledgments

## References

[1] Simon Boeder, Fabian Gigengack, and Benjamin Risse. Langocc: Self-supervised open vocabulary occupancy estimation via volume rendering. *arXiv preprint arXiv:2407.17310*, 2024. 3

[2] Daniel Bolya, Chaitanya Ryali, Judy Hoffman, and Christoph Feichtenhofer. Window attention is bugged: how not to interpolate position embeddings. *arXiv preprint arXiv:2311.05613*, 2023. 4

[3] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11621–11631, 2020. 8, 16

[4] Haozhi Cao, Yuecong Xu, Jianfei Yang, Pengyu Yin, Shenghai Yuan, and Lihua Xie. Mopa: Multi-modal prior aided domain adaptation for 3d semantic segmentation. In *IEEE International Conference on Robotics and Automation*, pages 9463–9470, 2024. 2

[5] Keyan Chen, Chenyang Liu, Hao Chen, Haotian Zhang, Wenyuan Li, Zhengxia Zou, and Zhenwei Shi. Rsprompter: Learning to prompt for remote sensing instance segmentation based on visual foundation model. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–17, 2024. 1, 2

[6] Runnan Chen, Youquan Liu, Lingdong Kong, Xinge Zhu, Yuexin Ma, Yikang Li, Yuenan Hou, Yu Qiao, and Wenping Wang. Clip2scene: Towards label-efficient 3d scene understanding by clip. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7020–7030, 2023. 3

[7] Tianrun Chen, Ankang Lu, Lanyun Zhu, Chaotao Ding, Chunan Yu, Deyi Ji, Zejian Li, Lingyun Sun, Papa Mao, and Ying Zang. Sam2-adapter: Evaluating & adapting segment anything 2 in downstream tasks: Camouflage, shadow, medical image segmentation, and more. *arXiv preprint arXiv:2408.04579*, 2024. 1, 13

[8] Xuanyao Chen, Tianyuan Zhang, Yue Wang, Yilun Wang, and Hang Zhao. Futr3d: A unified sensor fusion framework for 3d detection. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 172–181, 2023. 2

[9] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Modular interactive video object segmentation: Interaction-to-mask, propagation and difference-aware fusion. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5559–5568, 2021. 7

[10] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3075–3084, 2019. 4, 7, 12

[11] Yinpeng Dong, Caixin Kang, Jinlai Zhang, Zijian Zhu, Yikai Wang, Xiao Yang, Hang Su, Xingxing Wei, and Jun Zhu. Benchmarking robustness of 3d object detection to common corruptions. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1022–1032, 2023. 1

[12] Ilya Fradlin, Idil Esen Zulfikar, Kadir Yilmaz, Theodora Kontogianni, and Bastian Leibe. Interactive4d: Interactive 4d lidar segmentation. *arXiv preprint arXiv:2410.08206*, 2024. 2

[13] Chenrui Han, Xuan Yu, Yuxuan Xie, Yili Liu, Sitong Mao, Shunbo Zhou, Rong Xiong, and Yue Wang. Scale disparity of instances in interactive point cloud segmentation. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 2660–2667, 2024. 2

[14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 12

[15] Rui Huang, Songyou Peng, Ayca Takmaz, Federico Tombari, Marc Pollefeys, Shiji Song, Gao Huang, and Francis Engelmann. Segment3d: Learning fine-grained class-agnostic 3d segmentation without manual labels. In *European Conference on Computer Vision*, pages 278–295, 2024. 2

[16] Lei Ke, Mingqiao Ye, Martin Danelljan, Yu-Wing Tai, Chi-Keung Tang, Fisher Yu, et al. Segment anything in high quality. *Advances in Neural Information Processing Systems*, 36: 29914–29934, 2023. 2

[17] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 1, 2, 3, 5, 6

[18] Jiale Li, Hang Dai, Hao Han, and Yong Ding. Mseg3d: Multi-modal 3d semantic segmentation for autonomous driving. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21694–21704, 2023. 2

[19] Shiyao Li, Wenming Yang, and Qingmin Liao. Pmafusion: Projection-based multi-modal alignment for 3d semantic occupancy prediction. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3627–3634, 2024. 2

[20] Xiang Li, Junbo Yin, Botian Shi, Yikang Li, Ruigang Yang, and Jianbing Shen. Lwsis: Lidar-guided weakly supervised instance segmentation for autonomous driving. In *AAAI Conference on Artificial Intelligence*, pages 1433–1441, 2023. 8, 16

[21] Yanwei Li, Yilun Chen, Xiaojuan Qi, Zeming Li, Jian Sun, and Jiaya Jia. Unifying voxel-based representation with transformer for 3d object detection. *Advances in Neural Information Processing Systems*, 35:18442–18455, 2022. 2

[22] Tingting Liang, Hongwei Xie, Kaicheng Yu, Zhongyu Xia, Zhiwei Lin, Yongtao Wang, Tao Tang, Bing Wang, and Zhi Tang. Bevfusion: A simple and robust lidar-camera fusion framework. *Advances in Neural Information Processing Systems*, 35:10421–10434, 2022. 2

[23] Guibiao Liao, Jiankun Li, and Xiaoqing Ye. Vlm2scene: Self-supervised image-text-lidar learning with foundation models for autonomous driving scene understanding. In *AAAI Conference on Artificial Intelligence*, pages 3351–3359, 2024. 3

[24] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European Conference on Computer Vision*, pages 38–55, 2024. 6, 13

[25] Youquan Liu, Lingdong Kong, Jun Cen, Runnan Chen, Wenwei Zhang, Liang Pan, Kai Chen, and Ziwei Liu. Segment any point cloud sequences by distilling vision foundation models. *Advances in Neural Information Processing Systems*, 36:37193–37229, 2023. 2

[26] Zhijian Liu, Haotian Tang, Alexander Amini, Xinyu Yang, Huizi Mao, Daniela L Rus, and Song Han. Bevfusion: Multi-task multi-sensor fusion with unified bird's-eye view representation. In *IEEE International Conference on Robotics and Automation*, pages 2774–2781, 2023. 1, 2

[27] Yuhang Lu, Qi Jiang, Runnan Chen, Yuenan Hou, Xinge Zhu, and Yuexin Ma. See more and know more: Zero-shot point cloud segmentation via multi-modal visual data. In *IEEE/CVF International Conference on Computer Vision*, pages 21674–21684, 2023. 3

[28] Jun Ma, Yuting He, Feifei Li, Lin Han, Chenyu You, and Bo Wang. Segment anything in medical images. *Nature Communications*, 15(1):654, 2024. 1, 2

[29] Maciej A Mazurowski, Haoyu Dong, Hanxue Gu, Jichen Yang, Nicholas Konz, and Yixin Zhang. Segment anything model for medical image analysis: an experimental study. *Medical Image Analysis*, 89:102918, 2023. 1, 2

[30] Peizhou Ni, Xu Li, Wang Xu, Dong Kong, Yue Hu, and Kun Wei. Robust 3d semantic segmentation based on multi-phase multi-modal fusion for intelligent vehicles. *IEEE Transactions on Intelligent Vehicles*, 9(1):1602–1614, 2023. 2

[31] Aljoša Ošep, Tim Meinhardt, Francesco Ferroni, Neehar Peri, Deva Ramanan, and Laura Leal-Taixé. Better call sal: Towards learning to segment anything in lidar. In *European Conference on Computer Vision*, pages 71–90, 2024. 2

[32] Jingyi Pan, Zipeng Wang, and Lin Wang. Co-occ: Coupling explicit feature fusion with volume rendering regularization for multi-modal 3d semantic occupancy prediction. *IEEE Robotics and Automation Letters*, 2024. 2

[33] Songyou Peng, Kyle Genova, Chiyu Jiang, Andrea Tagliasacchi, Marc Pollefeys, Thomas Funkhouser, et al. Openscene: 3d scene understanding with open vocabularies. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 815–824, 2023. 3

[34] Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In *European Conference on Computer Vision*, pages 194–210, 2020. 4

[35] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alex Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv preprint arXiv:1704.00675*, 2017. 7

[36] 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, 2021. 3

[37] 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. 1, 2, 3, 4, 5, 6, 7, 12

[38] Tianhe Ren, Qing Jiang, Shilong Liu, Zhaoyang Zeng, Wenlong Liu, Han Gao, Hongjie Huang, Zhengyu Ma, Xiaoke Jiang, Yihao Chen, et al. Grounding dino 1.5: Advance the" edge" of open-set object detection. *arXiv preprint arXiv:2405.10300*, 2024. 6

[39] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, et al. Grounded sam: Assembling open-world models for diverse visual tasks. *arXiv preprint arXiv:2401.14159*, 2024. 6

[40] Chaitanya Ryali, Yuan-Ting Hu, Daniel Bolya, Chen Wei, Haoqi Fan, Po-Yao Huang, Vaibhav Aggarwal, Arkabandhu Chowdhury, Omid Poursaeed, Judy Hoffman, et al. Hiera: A hierarchical vision transformer without the bells-and-whistles. In *International Conference on Machine Learning*, pages 29441–29454, 2023. 4, 7, 12

[41] Inkyu Shin, Yi-Hsuan Tsai, Bingbing Zhuang, Samuel Schulter, Buyu Liu, Sparsh Garg, In So Kweon, and Kuk-Jin Yoon. Mm-tta: multi-modal test-time adaptation for 3d

semantic segmentation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16928–16937, 2022. 2

[42] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2446–2454, 2020. 2, 3, 6, 12

[43] Zhiyu Tan, Zichao Dong, Cheng Zhang, Weikun Zhang, Hang Ji, and Hao Li. Ovo: Open-vocabulary occupancy. *arXiv preprint arXiv:2305.16133*, 2023. 3

[44] Haotian Tang, Zhijian Liu, Xiuyu Li, Yujun Lin, and Song Han. TorchSparse: Efficient Point Cloud Inference Engine. In *Conference on Machine Learning and Systems*, 2022. 4, 7, 12

[45] Haotian Tang, Shang Yang, Zhijian Liu, Ke Hong, Zhongming Yu, Xiuyu Li, Guohao Dai, Yu Wang, and Song Han. TorchSparse++: Efficient Point Cloud Engine. In *Computer Vision and Pattern Recognition Workshops*, 2023. 4, 12

[46] Lv Tang, Haoke Xiao, and Bo Li. Can sam segment anything? when sam meets camouflaged object detection. *arXiv preprint arXiv:2304.04709*, 2023. 1, 2

[47] Ignacio Vizzo, Tiziano Guadagnino, Jens Behley, and Cyrill Stachniss. Vdbfusion: Flexible and efficient tsdf integration of range sensor data. *Sensors*, 22(3):1296, 2022. 13

[48] Antonin Vobecky, Oriane Siméoni, David Hurych, Spyridon Gidaris, Andrei Bursuc, Patrick Pérez, and Josef Sivic. Pop-3d: Open-vocabulary 3d occupancy prediction from images. *Advances in Neural Information Processing Systems*, 36, 2024. 3

[49] Guoqing Wang, Zhongdao Wang, Pin Tang, Jilai Zheng, Xiangxuan Ren, Bailan Feng, and Chao Ma. Occgen: Generative multi-modal 3d occupancy prediction for autonomous driving. In *European Conference on Computer Vision*, pages 95–112, 2024. 2

[50] Song Wang, Jianke Zhu, and Ruixiang Zhang. Metarangeseg: Lidar sequence semantic segmentation using multiple feature aggregation. *IEEE Robotics and Automation Letters*, 7(4):9739–9746, 2022. 2

[51] Song Wang, Wentong Li, Wenyu Liu, Xiaolu Liu, and Jianke Zhu. Lidar2map: In defense of lidar-based semantic map construction using online camera distillation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5186–5195, 2023. 2

[52] Shaoyuan Xie, Lingdong Kong, Wenwei Zhang, Jiawei Ren, Liang Pan, Kai Chen, and Ziwei Liu. Robobev: Towards robust bird's eye view perception under corruptions. *arXiv preprint arXiv:2304.06719*, 2023. 1

[53] Yichen Xie, Chenfeng Xu, Marie-Julie Rakotosaona, Patrick Rim, Federico Tombari, Kurt Keutzer, Masayoshi Tomizuka, and Wei Zhan. Sparsefusion: Fusing multi-modal sparse representations for multi-sensor 3d object detection. In *IEEE/CVF International Conference on Computer Vision*, pages 17591–17602, 2023. 2

[54] Yunyang Xiong, Bala Varadarajan, Lemeng Wu, Xiaoyu Xiang, Fanyi Xiao, Chenchen Zhu, Xiaoliang Dai, Dilin Wang,

Fei Sun, Forrest Iandola, et al. Efficientsam: Leveraged masked image pretraining for efficient segment anything. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16111–16121, 2024. 2

[55] Junjie Yan, Yingfei Liu, Jianjian Sun, Fan Jia, Shuailin Li, Tiancai Wang, and Xiangyu Zhang. Cross modal transformer: Towards fast and robust 3d object detection. In *IEEE/CVF International Conference on Computer Vision*, pages 18268–18278, 2023. 1, 2

[56] Yunhan Yang, Xiaoyang Wu, Tong He, Hengshuang Zhao, and Xihui Liu. Sam3d: Segment anything in 3d scenes. *arXiv preprint arXiv:2306.03908*, 2023. 2

[57] Yihan Zeng, Chenhan Jiang, Jiageng Mao, Jianhua Han, Chaoqiang Ye, Qingqiu Huang, Dit-Yan Yeung, Zhen Yang, Xiaodan Liang, and Hang Xu. Clip2: Contrastive language-image-point pretraining from real-world point cloud data. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15244–15253, 2023. 3

[58] Chaoning Zhang, Dongshen Han, Yu Qiao, Jung Uk Kim, Sung-Ho Bae, Seungkyu Lee, and Choong Seon Hong. Faster segment anything: Towards lightweight sam for mobile applications. *arXiv preprint arXiv:2306.14289*, 2023. 2

[59] Hongcheng Zhang, Liu Liang, Pengxin Zeng, Xiao Song, and Zhe Wang. Sparselif: High-performance sparse lidar-camera fusion for 3d object detection. In *European Conference on Computer Vision*, pages 109–128, 2024. 1, 2

[60] Junbo Zhang, Runpei Dong, and Kaisheng Ma. Clip-fo3d: Learning free open-world 3d scene representations from 2d dense clip. In *IEEE/CVF International Conference on Computer Vision*, pages 2048–2059, 2023. 3

[61] Shuo Zhang, Yupeng Zhai, Jilin Mei, and Yu Hu. Fusionocc: Multi-modal fusion for 3d occupancy prediction. In *ACM International Conference on Multimedia*, pages 787–796, 2024. 2

[62] Xu Zhao, Wenchao Ding, Yongqi An, Yinglong Du, Tao Yu, Min Li, Ming Tang, and Jinqiao Wang. Fast segment anything. *arXiv preprint arXiv:2306.12156*, 2023. 2

[63] Jilai Zheng, Pin Tang, Zhongdao Wang, Guoqing Wang, Xiangxuan Ren, Bailan Feng, and Chao Ma. Veon: Vocabulary-enhanced occupancy prediction. In *European Conference on Computer Vision*, pages 92–108, 2024. 3

[64] Yuchen Zhou, Jiayuan Gu, Tung Yen Chiang, Fanbo Xiang, and Hao Su. Point-SAM: Promptable 3d segmentation model for point clouds. In *International Conference on Learning Representations*, 2025. 2

# SAM4D: Segment Anything in Camera and LiDAR Streams

## Supplementary Material

In this document, we further provide the following materials to support the findings and conclusions drawn in the main body of this paper.

## A. Model and Training Details

### A.1. Model Architecture

**Image Encoder.** The image encoder in SAM4D follows the same architecture as SAM2 [37], but given the smaller scale of our dataset compared to SAM2, we adopt the Hiera-S variant [40] as the default configuration to balance performance and efficiency.

**LiDAR Encoder.** We adopt MinkUNet [10], implemented with TorchSparse [44, 45], as the LiDAR encoder. Inspired by ResNet [14], we define Mink34 and Mink50 as backbone structures, with Mink34 as the default choice. The encoder downsamples the input to stride 32, with feature dimensions $[32, 32, 64, 128, 256]$, and then upsamples to stride 4, extracting voxel features at strides 16, 8, and 4. The stride 16 features are primarily used by the memory module and mask decoder, while the stride 8 and 4 features assist the mask decoder in recovering high-resolution segmentation details, similar to SAM2. To improve the generalization across various datasets, we exclude the original xyz coordinates and do not use intensity and elongation features provided by Waymo Open Dataset [42]. Instead, we assign a binary occupancy value to occupied voxels, ensuring the LiDAR encoder remains dataset-agnostic.

**Mask Decoder.** Sparse prompt tokens from image and LiDAR are concatenated and used as queries for mask prediction, with a shared Transformer module across both modalities. The query token configuration follows SAM2, consisting of mask queries, sparse prompt tokens, IoU tokens, and object pointers stored in memory. This design enables efficient cross-modal segmentation while ensuring consistency between image and LiDAR-based queries.

**Model Parameters.** SAM4D comprises 119.88M parameters, distributed across its core components. The image encoder has 34.32M parameters, while the LiDAR encoder contains 26.94M parameters. The memory module, responsible for temporal feature aggregation and cross-modal attention, is the largest component with 53.96M parameters. The mask decoder, which processes prompt-based queries for segmentation, accounts for 4.66M parameters. This de-

| Config | Value |
|---|---|
| data | Waymo-4DSeg |
| steps | $\sim$44k |
| resolution | Camera: $768 \times 768$, LiDAR: 0.15 m |
| precision | bfloat16 |
| optimizer | AdamW |
| optimizer momentum | $\beta_1 = 0.9, \beta_2 = 0.999$ |
| gradient clipping | type: $\ell_2$, max: 0.1 |
| weight decay | 0.1 |
| learning rate (lr) | OneCycleLR, init: $5e^{-6}$, max: $5e^{-5}$, anneal strategy: cos, pct_start: 0.4 |
| warmup | linear, 7.5k iters |
| layer-wise decay | 0.9 |
| image augmentation | hflip, resize to $768 \times 768$ (square) |
| video augmentation | hflip, affine (deg: 25), colorjitter, grayscale, per-frame colorjitter, mosaic-$2 \times 2$ |
| LiDAR augmentation | rotation-Z (deg: 45), hflip, vflip |
| drop path | 0.2 |
| mask losses (weight) | Focal (20), Dice (1) |
| IoU loss (weight) | $\ell_1$ (1) |
| occlusion loss (weight) | Cross-entropy (1) |
| global attn. blocks | $12 - 16 - 20$ |

Table A1. Hyperparameters for SAM4D full training.

sign balances multimodal fusion, temporal reasoning, and segmentation efficiency.

### A.2. Training Details

Without loss of generality, we use the front-view camera and LiDAR from the Waymo dataset to validate the feasibility of the proposed solution, which can later be extended to multi-camera and LiDAR setups. Since our data is constructed through 4D reconstruction using a 5-camera and LiDAR system, we cannot guarantee that objects appear in both modalities in every frame. Therefore, during the training process, to enable parallel imitation of interaction logic for multiple targets, the following rules are applied when selecting targets for each step in the data pipeline: there is a 0.5 probability that the target exists in both modalities, a 0.25 probability that it appears only in the camera, and a 0.25 probability that it appears only in the LiDAR. During training, the iterative modification logic for mimicking targets is as follows: if the target belongs to both modalities, each prompt randomly selects one modality; if the target belongs to only one modality, the modality in which the target appears is chosen for the prompt.

SAM4D is trained on Waymo-4DSeg for 44k steps with a $768 \times 768$ image resolution and a LiDAR voxel size of 0.15. Specifically, we sample 8-frame sequences and randomly select up to 2 frames to receive prompts. During training, corrective clicks are probabilistically sampled

based on both ground-truth masks and model predictions. The initial prompts are assigned with probabilities of $0.5$ for ground-truth masks, $0.25$ for points, and $0.25$ for bounding boxes, respectively. The loss consists of a combination of focal loss and dice loss for mask prediction and mean absolute error (MAE) loss for IoU prediction. If an object is missing in a given modality, we do not apply supervision to the prediction of that modality. AdamW with OneCycleLR are utilized to optimize the network. Image and video augmentations follow SAM2 (excluding shear), while LiDAR-specific augmentations include Z-axis rotation, hflip, and vflip. Other hyperparameters align with SAM2, which are provided in Tab. A1.

## B. Details on Data Engine and Dataset

### B.1. Data Engine Details

In this section, we provide a more detailed description of the implementation details of the three steps in our data engine to supplement Sec. 5 in our main paper.

**Step 1: Generation of VFM-based image masklets**

In this step, we utilize vision foundation models (VFM) [7, 24] to generate initial annotations including boxes and masks in keyframes firstly. In each new keyframe, we redetect scene objects and match them with the propagated masks from the previous keyframe, merging them as the segmentation result of the current keyframe. Newly detected objects, which were not present in the previous masklets, are first propagated backward to the start of the sequence, and then the merged masks (both new and existing objects) are propagated forward to the next keyframe, continuing this process until the end of the sequence.

This iterative approach produces masklets for the entire image sequence, ensuring consistent object categories and instance IDs across time, laying the foundation for LiDAR ground truth generation in subsequent stages.

**Step 2: 4D Reconstruction and Ray Casting**

Transferring masklets from video to LiDAR frames requires establishing correspondences between pixels and LiDAR points, which is challenging due to the large number of pixels and 3D points in the sequence. To address this challenge, we preemptively perform 4D LiDAR reconstruction using VDBFusion [47], which generates a more efficient representation of spatial occupancy, since the voxel count depends only on scene size and is independent of the number of frames.

The 4D reconstruction comprises multiple foreground components and a single background component. The background consists of static objects and is maintained as a single instance in the world coordinate system. The foreground includes potentially moving entities such as vehicles, cyclists, and pedestrians, each with its own motion trajectory. We leverage the pre-annotated 3D bounding boxes

of these foreground objects to obtain their relative poses in each frame with respect to the world, and subsequently perform individual 4D reconstruction within each object's body coordinate system. This allows the voxels occupied by the foreground objects to remain unchanged even as they move, with only the overall position shifting.

Following this, we generate a dense pixel-voxel mapping table via ray casting. We first compute image poses in the world coordinate system by solving the PnP(Perspective-n-Point) problems, based on the given single-frame point cloud and pixel correspondences. Then, for each image, we construct multiple rays starting from the camera position and ending at the center points of the voxels within the viewing frustum. Each ray falls into an image pixel, and we match the pixel with the voxel that the ray intersects together to build pixel-voxel the mapping table.

**Step 3: Cross-Modal Masklet Fusion**

By querying the pixel-voxel mapping table established in Step 2, we can identify the voxels corresponding to the pixels masked in Step 1, thereby transferring the mask to the voxels. In an ideal scenario, SAM2 ensures consistency of masklets between frames in the video, and we can directly merge the voxel masks from the video stream to obtain an accumulated voxel masklet.

However, we found that both the video masklets and the mapping table are often noisy. A common issue is that SAM2 mis-matches objects, confusing two similar objects appearing in close positions across different frames as the same object. This issue occurs for both background and foreground objects. Additionally, because the pixel-to-voxel correspondence is not always accurate and the edges of 2D masks on images are not perfectly precise, the resulting masks projected onto the voxels are prone to contain noise. Finally, the image segmentation model occasionally misclassifies artifacts such as light spots or other visual anomalies in images as actual objects.

To mitigate these problems, we implemented a clustering approach for noise filtering in the voxel masklets. We employed the DBSCAN algorithm to cluster voxels based on their BEV positions and selected the cluster with the highest average quality as the main cluster, filtering out the rest as noise. Assuming that voxels associated with a single object are adjacent in BEV space, we utilized the DBSCAN algorithm to cluster the voxels based on their BEV positions. We also counted the frequency with which each voxel was mapped to the current object and computed the vote rate as the ratio of this frequency to the total observations in the current image sequence. Ultimately, we selected the cluster with the highest average vote rate as the main cluster, and leave out rest as noise. Fig A2 provides examples of the issues mentioned above and demonstrates the effectiveness of our filter in addressing these problems.

After filtering, we expect significant overlap between

Figure A1. Examples of the Waymo-4Dseg Dataset: In this figure, we visualize the masklets from selected frames of 6 clips. For each clip, we present two frames in two rows, where each row displays the masklets from the side-left, front-left, front, front-right, and side-right images, from left to right, along with the masklets from the LiDAR frames. Through both horizontal and vertical comparisons, the consistency of the masklets across different video streams, modalities, and frames can be observed.
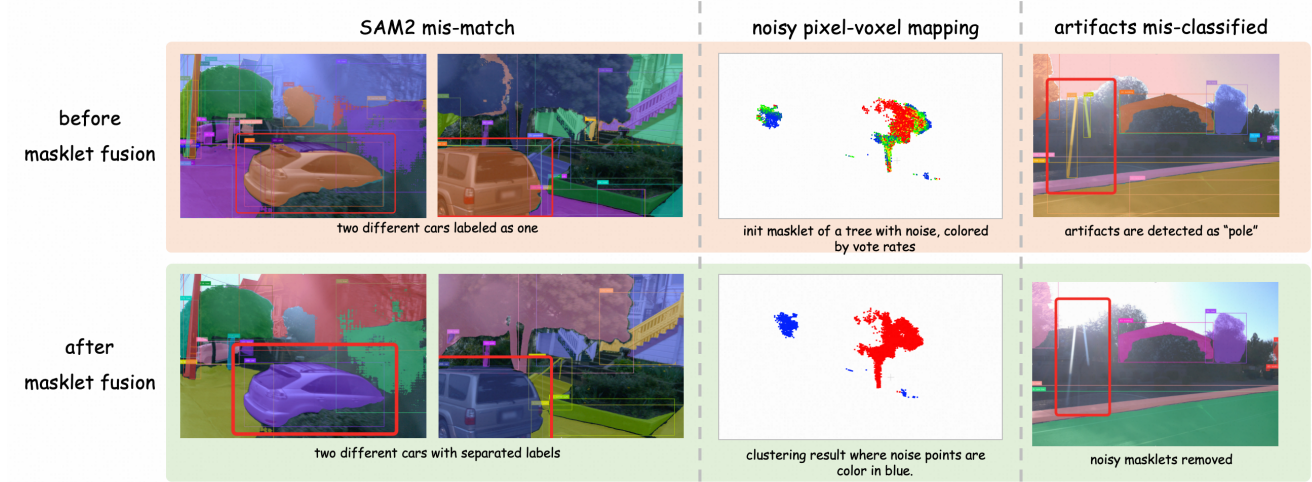
14

**Figure A2.** This figure illustrates three representative problems that may arise before masklet fusion, as well as how our fusion process addresses these issues.

voxel masklets from different videos corresponding to the same object. Overlaps between masklets from two videos were assessed, leading to the merging of those with substantial overlap into a single unified voxel masklet. Finally, we created a mapping table between points from LiDAR frames and voxels based on their 3D spatial distances, facilitating the transfer of the final voxel masklet to the LiDAR frames.

We evaluated the quality of the unified voxel masklets using cross-modal IoU. Assuming that a masklet is visible for image $i$, we calculated the IoU between the voxels mapped by masklet in image $i$ and the visible part of the unified voxel masklet. The average IoU across all images represents one masklet's overall score. The mean score of the masklets in our dataset is 0.56, with a 10th percentile of 0.24. Throughout this process, human annotators play a crucial role in adjusting the parameters based on mask quality and conducting frame-by-frame verification of the final labels in both the image and LiDAR.

### B.2. Dataset Statistical Information

Here, we provide additional information about the dataset and details about our data engine. In Figure A3, we provide detailed statistical information about the masklets in our dataset, including volume distribution, area distribution, the proportion of frames in which cross-modal masklets co-occur, and score distribution. The distributions of volume and area reflect the diversity and richness of the annotated objects in our dataset. Additionally, we calculate the proportion of frames in which cross-modal masklets are present in both modalities, a metric of interest to users. Common scenarios include objects exiting the video frame as the vehicle moves forward while still being detectable by LiDAR, or objects being beyond the scanning range of LiDAR but still visible in the video. Such cases require special han-
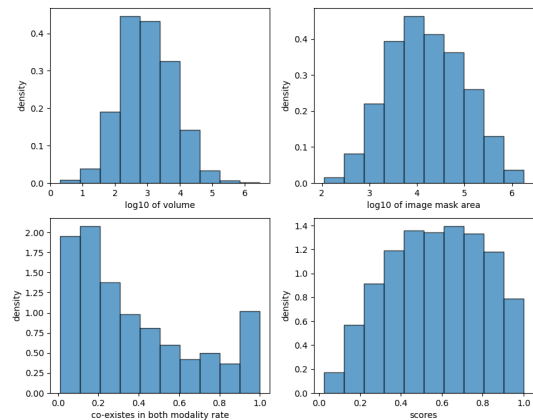


**Figure A3.** Statistical information of masklets in Waymo-4DSeg dataset

dling during model prompting and training, which is why we provide the proportion of frames where cross-modal masklets coexist. We also present the score distribution of the masklets. Although denoising has been applied during processing, a small number of low-quality masklets remain. Fortunately, we can quantify their quality through scores and filter them out during experimentation. In Figure A1, we present the visualized cross-modal masklets from several clips.

## C. Details on the Experimental Settings

### C.1. About the Training and Validation Data.

As mentioned above, when the data engine generates pseudo labels, it creates voxels in 4D for each object, allowing us to obtain the volume of the object in space. Ad-

ditionally, multi-modal consistency checks are performed to calculate the IoU between each frame's image and point cloud. The average IoU over the sequence serves as a reference for the quality of the pseudo-labels, which we refer to as the "score." During SAM4D training, targets with a volume greater than 10 and a score greater than 0.3 are used. For testing, to further ensure the reliability of the ground truth, the volume filtering threshold is increased to 50, and the score threshold is raised to 0.5. Furthermore, there is currently significant ambiguity in the pseudo-labels for ground regions. To ensure better convergence of the LiDAR branch, we temporarily exclude instances near the ground during both training and evaluation. Despite these settings, the number of targets evaluated in each sequence still exceeds 100, resulting in slow evaluation speeds. To accelerate the evaluation, we filter and evaluate only those objects that appear in at least one frame in both the front-view camera and LiDAR.

### C.2. About the Generalization Experiments.

The generalization experiments are conducted on nuScenes dataset [3] with nuInsSeg [20]. The nuInsSeg dataset [20], built on nuScenes [3], provides 2D instance segmentation annotations for foreground objects, with instance IDs corresponding to 3D point cloud segmentation labels.

## D. Limitations and Further Discussions

### D.1. Limitations

While SAM4D effectively integrates multimodal and temporal segmentation, the domain gap across LiDAR sensors remains a challenge, as variations in sensor configurations and point cloud density limit generalization compared to images. Moreover, the spatial representation on point clouds is inherently constrained by single-frame sparsity, occlusions, and blind spots, which may hinder object completeness in certain scenarios. Additionally, while Waymo-4DSeg provides high-quality multimodal labels, the size of the data set can be expanded to cover a broader range of driving conditions, weather variations, and rare long-tail scenarios. Increasing data set diversity would improve the generalizability of the model, particularly in corner cases where data sparsity remains a challenge.

### D.2. Future Work

Currently, SAM4D is trained on pseudo-labels generated by an automated data engine. Although the data labels have undergone multi-modal consistency verification, ambiguities and inaccuracies still persist. Future work will focus on improving SAM4D's data strategy, model adaptability, and scalability. To enhance label quality, we plan to expand dataset scale using our automated data engine and integrate human-annotated subsets for fine-tuning. A confidence-

based filtering mechanism will further refine pseudo-labels iteratively. Additionally, extending SAM4D to incorporate natural language descriptions will enable multimodal segmentation conditioned on text, leveraging LLMs for semantic guidance to reduce reliance on human annotations. Exploring weakly supervised and self-supervised learning will further enhance adaptability while minimizing manual labeling. Beyond data efficiency, improving memory attention and computational efficiency will enable scaling to multi-camera and multi-sensor systems, enhancing 4D spatiotemporal perception in complex environments.

## E. License and Consent with Public Resources

### E.1. Public Datasets

We utilize the Waymo Open Dataset [3] to construct our Waymo-4DSeg dataset. nuScenes [3] and the corresponding nuInstSeg [20] are adopted to further evaluate our model:

- Waymo Open Dataset[1] . . . . . . . . Waymo Dataset License
- nuScenes[2] . . . . . . . . . . . . . . . . . . . . . . . . CC BY-NC-SA 4.0
- nuScenes-devkit[3] . . . . . . . . . . . . . . . . . Apache License 2.0
- nuInsSeg[4] . . . . . . . . . . . . . . . . . . . . . . . . . . . . . MIT License

### E.2. Public Implementation

We leverage publicly available pre-trained models and source codes to investigate the promptable segmentation in the multimodal domain:

- SAM2[5] . . . . . . . . . . . . . . . . . . . . . . . . . Apache License 2.0
- TorchSparse[6] . . . . . . . . . . . . . . . . . . . . . . . . . . . MIT License
- GroundingDINO[7] . . . . . . . . . . . . . . . . Apache License 2.0
- Grounded-SAM-2[8] . . . . . . . . . . . . . . . Apache License 2.0
- VDBFusion[9] . . . . . . . . . . . . . . . . . . . . . . . . . . . . MIT License
- Mask-Propagation[10] . . . . . . . . . . . . . . . . . . . . . MIT License

---

[1]https://waymo.com/open.
[2]https://www.nuscenes.org/nuscenes.
[3]https://github.com/nutonomy/nuscenes-devkit.
[4]https://github.com/Serenos/nuInsSeg.
[5]https://github.com/facebookresearch/sam2.
[6]https://github.com/mit-han-lab/torchsparse.
[7]https://github.com/IDEA-Research/GroundingDINO.
[8]https://github.com/IDEA-Research/Grounded-SAM-2.
[9]https://github.com/PRBonn/vdbfusion.
[10]https://github.com/hkchengrex/Mask-Propagation.