

# XVO: Generalized Visual Odometry via Cross-Modal Self-Training

Lei Lai\* Zhongkai Shangguan\* Jimuyang Zhang Eshed Ohn-Bar  
Boston University

{leilai, sgzk, zhangjim, eohnbar}@bu.edu

## Abstract

We propose XVO, a semi-supervised learning method for training generalized monocular Visual Odometry (VO) models with robust off-the-self operation across diverse datasets and settings. In contrast to standard monocular VO approaches which often study a known calibration within a single dataset, XVO efficiently learns to recover relative pose with real-world scale from visual scene semantics, i.e., without relying on any known camera parameters. We optimize the motion estimation model via self-training from large amounts of unconstrained and heterogeneous dash camera videos available on YouTube. Our key contribution is twofold. First, we empirically demonstrate the benefits of semi-supervised training for learning a general-purpose direct VO regression network. Second, we demonstrate multi-modal supervision, including segmentation, flow, depth, and audio auxiliary prediction tasks, to facilitate generalized representations for the VO task. Specifically, we find audio prediction task to significantly enhance the semi-supervised learning process while alleviating noisy pseudo-labels, particularly in highly dynamic and out-of-domain video data. Our proposed teacher network achieves state-of-the-art performance on the commonly used KITTI benchmark despite no multi-frame optimization or knowledge of camera parameters. Combined with the proposed semi-supervised step, XVO demonstrates off-the-shelf knowledge transfer across diverse conditions on KITTI, nuScenes, and Argoverse without fine-tuning.

## 1. Introduction

Monocular Visual Odometry (VO) methods for recovering ego-motion from a sequence of images have mostly been studied within a *restricted scope*, where a single dataset, such as KITTI [30], may be used for both training and evaluation under a fixed pre-calibrated camera [39, 48, 57, 82, 116, 119, 123, 131, 133, 135]. However, very few

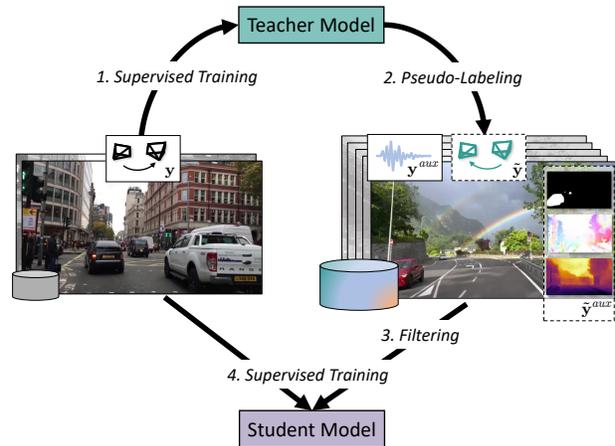


Figure 1: **Learning General-Purpose Monocular Visual Odometry (VO) Models from Multi-Modal and Pseudo-Labeled Videos.** Our proposed XVO framework first trains an ego-motion prediction *teacher model* over a small initial dataset, e.g., nuScenes [7]. We then expand the original dataset through pseudo-labeling of in-the-wild videos. Motivated by how humans learn general representations through observation of large amounts of multi-modal data, we employ multiple auxiliary prediction tasks, including segmentation, flow, depth, and audio, as part of the semi-supervised training process. Finally, we leverage *uncertainty-based filtering* of potentially noisy pseudo-labels and train a robust student model.

studies have analyzed the task of *generalized VO*, i.e., relative pose estimation with real-world scale across differing scenes and capture setups.

For instance, consider an autonomous robot or vehicle deployed at a large scale. The robot is highly likely to encounter environments for which no ground truth ego-motion data was previously collected. In such novel settings, current VO methods will quickly exhibit poor ego-motion estimation and drift [25, 26, 31, 49, 50, 68, 103, 135]. Moreover, our robot may be required to adjust its camera setup over its lifetime (e.g., to a new camera) or leverage data from a fleet of robots with varying or perhaps unknown camera configurations. Yet, existing VO methods generally

\* Equally contributed.

assume carefully calibrated camera parameters during training [25, 26, 68, 82, 103, 116, 135]. Specifically, to simplify the ill-posed monocular pose recovery task, researchers often resort to relying on knowledge of the camera intrinsic parameters to incorporate various geometric or photometric consistency-based mechanisms [39, 48, 57, 96, 119, 123, 131, 133]. In this work, we do not make such an assumption as we are concerned training VO models that can learn from and operate under diverse unconstrained videos in the wild. Specifically, we pursue an orthogonal direction to prior work based on semi-supervised learning and explore more scalable and camera-agnostic deployment settings.

Our key hypothesis is that neural network models can learn to circumvent issues related to pose and scale ambiguity in generalized VO settings through observation of ample amounts of diverse scene and motion video data. This approach is motivated by humans’ ability to flexibly estimate motion in arbitrary conditions through a general understanding of salient scene properties (e.g., object sizes) [77]. This general understanding is developed over large amounts of perceptual data, often multi-modal in nature [75, 88, 88]. For instance, cross-modal information processing between audio and video has been shown to play a role in spatial reasoning and proprioception [60, 61, 71, 79, 106]. Indeed, collected online videos often have audio, which can be used as a further source of cross-modal supervision. As further discussed in Sec. 3, we find ambient audio to correlate at times with scenarios where monocular VO tends to fail, such as determining ego-speed when the vehicle is stopped at a dense intersection or as context for the current traffic scenario when estimating translation (e.g., highway driving). Extracting information related to flow, segmentation, or depth can also further guide learning generalized representations. To fully explore the utility of self-training VO models, we analyze a unified multi-modal framework and its impact on guiding semi-supervised VO learning from large amounts of unconstrained sources.

As far as we are aware, we are the first to study the feasibility of self-training for direct, calibration-free, ego-motion pose regression with an absolute real-world scale. Specifically, we find that incorporating additional modalities via simple multi-task learning can significantly enhance model robustness and generalization. When paired with an uncertainty-based filtering module, we achieve state-of-the-art VO performance with a single broadly usable model which we validate for the autonomous driving use-case. Moreover, our training and inference is highly efficient as the auxiliary learning formulation does not alter the two-frame input, i.e., in contrast to methods relying on extracting rich intermediate representations [4, 103, 119, 131]. We demonstrate state-of-the-art results on KITTI using the proposed two-frame VO model structure without requiring elaborate long-term memory, computationally expensive it-

erative refinement steps, or knowledge of camera parameters. Our code is available at <https://genxvo.github.io/>.

## 2. Related Work

**Monocular Visual Odometry:** Despite recent advances, both geometrical and learning-based VO approaches are still mostly evaluated over limited datasets under similar training and testing conditions [9, 27, 42, 67, 76, 97, 99, 101]. For instance, training and evaluation are both conducted on KITTI [30], which contains a fixed camera setup with limited diversity and density of scenarios [103]. More recently, learning-based approaches leveraging unsupervised learning for VO [48, 50, 80, 119, 122, 133], have shown state-of-the-art performance on KITTI. Notably, UnDeepVO [48] utilizes stereo imagery for training to recover real-world scale without the need for labeled datasets. GeoNet [119] combines depth, optical flow, and camera pose to holistically learn a VO prediction model. TartanVO [103] conditions the VO model on the intrinsic parameters in order to achieve robust generalized performance. However, the aforementioned methods all require known camera intrinsics [48, 119] in inference resulting in a restricted use-case and cannot leverage data with unknown camera parameters. In light of these challenges, our work develops mechanisms to enable VO models to learn from and operate over diverse datasets without known calibration. Specifically, our method leverages semi-supervised and multi-modal learning techniques to learn robust generalized representations for estimating motion and real-world scale. Therefore, our approach is orthogonal to most related methods which emphasize self- and un-supervised learning of models based on warping and consistency tasks which rely on precise camera calibration [39, 48, 119].

**Semi-Supervision for Computer Vision Tasks:** Semi-supervised learning approaches have been extensively studied within the computer vision and machine learning communities [18, 47, 58, 89, 94, 118]. However, prior works have not yet focused on the monocular VO task, instead emphasizing object detection [8], semantic segmentation [105], 3D reconstruction [109], action recognition [108] or low-level computer vision tasks [6, 10, 17, 35, 40, 62, 64, 74, 91, 114, 117]. Consequently, fundamental research questions related to the impacts of improving VO model generalization remain unanswered, e.g., whether semi-supervised learning can be used to enhance reasoning over real-world scale [46, 81, 104, 127] and long-tail scenarios [126] or how uncertainty mechanisms can contribute to more robust training from heterogeneous video data. Specifically, we explore the role of multi-task and multi-modal learning in order to improve semi-supervised VO model training.

**Auxiliary Learning:** Our proposed method primarily aims

to enhance the performance of a VO model through auxiliary tasks within a semi-supervised learning framework. Auxiliary learning [52, 54, 129] aims to use auxiliary tasks to enhance the performance of the primary tasks. It has been effectively employed in diverse domains, including computer vision [15, 43, 63, 111], natural language processing [16, 19, 45], and robotics [38, 70, 72, 90, 100, 124, 125]. Xu et al. [111] applied auxiliary image classification and saliency detection to improve the performance of the semantic segmentation. Song et al. [90] leverages an auxiliary task of velocity estimation to enhance the ability to avoid obstacles of an indoor mobile robot. While several related studies employ auxiliary supervision derived from the ground-truth depth and optical flow [95, 96], in this work our goal is to explore the use of such supervision from potentially noisy pseudo-labels, e.g., as regularization for learning robust internal representations for VO.

**Cross-Modal Learning:** Cross-modal learning is inspired by how biological systems learn by incorporating complementary information from multiple modalities, such as vision, sound, and touch. Prior research in computational cross-modal learning has focused on learning a shared representation space where samples from distinct modalities i.e., image, audio, text, can be aligned [3, 36, 132]. Moreover, the addition multiple tasks and modalities have been shown to benefit generalization for various machine perception and learning tasks [2, 36, 84, 120, 121, 121]. For instance, audio generation [21, 93, 134], image captioning [51, 110], speech recognition [1, 86], navigation [13], and multimedia retrieval [11, 29, 37] have all shown improved performance due to cross-modal training. However, such studies tend to focus on simplified domains, e.g., restricted acoustic or haptic environments, whereas we analyze dense and dynamic scenes in the wild.

### 3. Method

Our proposed framework comprises three main steps: (1) uncertainty-aware training of an initial (i.e., teacher) VO model (Sec. 3.2); (2) pseudo-labeling with the removal of low-confidence and potentially noisy samples (Sec. 3.3); (3) self-training with pseudo-labeled and auxiliary prediction tasks of a robust VO student model (Sec. 3.4).

#### 3.1. Problem Setting

**Direct Pose Regression:** Our goal is to learn a general function for mapping two observed image frames  $\mathbf{x}_i = \{\mathbf{I}_{i-1}, \mathbf{I}_i\}$ , with  $\mathbf{I} \in \mathcal{R}^{W \times H \times 3}$ , to a relative camera pose with real-world scale  $\mathbf{y}_i = [\mathbf{R}_i | \mathbf{t}_i] \in \text{SE}(3)$  with rotation  $\mathbf{R}_i \in \text{SO}(3)$  and translation  $\mathbf{t}_i \in \mathbb{R}^3$ . Given a dataset comprising annotated labels of pose ground-truth,  $\mathcal{D}_L = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ , learning-based approaches for VO often optimize for a regression loss [5, 80, 96, 101, 119]. In

practice, the direct pose regression task often exhibits drift due to issues with absolute scale ambiguity and compound errors, particularly in cases of dense and dynamic scenes. For instance, small errors in rotation estimation can result in large errors over multiple time steps which impact the evaluation. While we formulate a two-frame regression task, prior methods have relied on longer-term memory in order to improve model robustness [39, 50, 101, 113], however, this comes at a computational and memory cost. Moreover, most monocular methods only produce up-to-scale predictions [49, 103, 119], as will be further discussed in Sec 4. Instead, we rely on a semi-supervised training process to mitigate issues in absolute scale recovery while enabling a simple two-frame model to achieve state-of-the-art results.

**Self-Training with Auxiliary Tasks:** In addition to a labeled odometry dataset  $\mathcal{D}_L$ , our framework assumes access to a large dataset that is not annotated with respect to the ego-motion task but potentially other complementary tasks that are auxiliary to the main VO task, i.e.,  $\mathcal{D}_U = \{(\mathbf{x}_i, \mathbf{y}_i^{aux})\}_{i=1}^M$ . Moreover, we assume access to a set of models for generating a *pseudo-labeled* dataset [8, 46, 81, 114, 127], i.e.,  $\mathcal{D}_{PL} = \{(\mathbf{x}_i, \tilde{\mathbf{y}}_i, \mathbf{y}_i^{aux}, \tilde{\mathbf{y}}_i^{aux})\}_{i=1}^M$  which can be joined with the original dataset  $\mathcal{D} = \mathcal{D}_L \cup \mathcal{D}_{PL}$  for supervised training (Sec. 3.4). We note that this is a practical assumption as there are abundant computer vision models for obtaining various pseudo-labels. As will be discussed in Sec. 3.3, these pseudo-labels may be filtered by removing high-uncertainty samples. Overall, the cross-modal self-training objective can be defined as

$$\mathcal{L}_{xvo} = \mathcal{L}_{vo}(\mathbf{y}) + \lambda_u \mathcal{L}_{unc}(\mathbf{y}) + \mathcal{L}_{aux}(\mathbf{y}^{aux}, \tilde{\mathbf{y}}^{aux}) \quad (1)$$

where  $\mathcal{L}_{vo}$  is a main VO task loss,  $\mathcal{L}_{unc}$  is an uncertainty estimation loss,  $\mathcal{L}_{aux}$  is defined over the auxiliary prediction tasks, and  $\lambda_u$  is a scalar hyper-parameter. We demonstrate our semi-supervised formulation to benefit various known issues with VO, e.g., improving scale recovery. Moreover, our formulation is kept efficient during inference as it does not alter the two-frame input  $\mathbf{x}$ , i.e., in contrast to methods relying on extracting intermediate representations as input, such as flow [103] or depth [131]. Next, we define our network structure and training.

#### 3.2. Ego-Motion Network Training

Our approach first trains a direct ego-motion teacher model (shown as the main encoder and middle branch in Fig. 1) over the labeled dataset  $\mathcal{D}_L$ . To enable learning from an unconstrained video, we do not incorporate any dependency on intrinsic parameters, i.e., either as input [25, 26, 103] or for computing a supervisory loss [48, 82, 96, 116, 119]. We find our network design to provide an efficient but surprisingly strong baseline, matching state-of-the-art on the KITTI benchmark despite no elaborate multi-frame optimization.

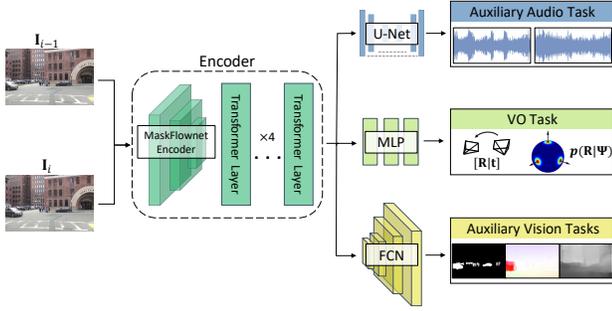


Figure 2: **Network Architecture.** Our initial teacher model (used for pseudo-labeling and filtering) encodes two concatenated image frames and predicts relative camera pose and its uncertainty. The complete cross-modal architecture leverages a similar architecture but with added auxiliary prediction branches with complementary tasks that can further guide self-training, e.g., prediction branches for audio reconstruction, dynamic object segmentation, optical flow, and depth.

**Encoder:** We employ a high-capacity feature extractor for effectively leveraging the rich multi-task supervision in later stages (Sec. 3.4). The feature extractor is a MaskFlowNet encoder [130], which was found to outperform the commonly used PWC-Net [92, 103], followed by four transformer self-attention layers [14, 23]. The patch size is  $12 \times 16$ , with each layer comprising four heads and 256 hidden parameters. The encoder structure for the initial teacher model and cross-modal student is kept the same.

**VO Decoder:** The VO decoder branch consists of three Fully Connected (FC) layers that regress relative pose  $\mathbf{y} = [\mathbf{R}|\mathbf{t}]$  and an uncertainty estimate for the prediction. The VO task optimizes a Mean Squared Error (MSE) loss over predicted translation  $\hat{\mathbf{t}} \in \mathcal{R}^3$  and Euler angle rotations  $\hat{\boldsymbol{\theta}} \in \mathcal{R}^3$ ,

$$\mathcal{L}_{vo} = \|\mathbf{t} - \hat{\mathbf{t}}\|_2^2 + \lambda_{\theta} \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|_2^2 \quad (2)$$

**Uncertainty Estimation:** To account for the difficulty in the absolute scale pose regression task, we propose to also model prediction uncertainty. We adopt a matrix Fisher distribution [65], which provides a framework for modeling rotation distribution on  $\text{SO}(3)$ . The probability density function of the matrix Fisher distribution is as follows:

$$p(\mathbf{R}|\Psi) = \frac{1}{c(\Psi)} \exp(\text{tr}(\Psi^T \mathbf{R})) \quad (3)$$

where  $\Psi \in \mathbb{R}^{3 \times 3}$  are the distribution parameters,  $\mathbf{R} \in \text{SO}(3)$  is the pose rotation matrix, and  $c(\Psi)$  is a normalization constant [59]. Given the estimated parameters  $\hat{\Psi}$

we use the negative log likelihood of  $\mathbf{R}$  in the predicted distribution as a loss, i.e.,

$$\mathcal{L}_{unc} = -\log(p(\mathbf{R}|\hat{\Psi})) \quad (4)$$

As a proxy for prediction (i.e., pseudo-label) quality, we find it is sufficient to model uncertainty in rotation prediction, however more elaborate estimation methods can also be used [41, 81, 115]. The confidence predictions will be used to remove potentially noisy pseudo-labels prior to the self-training process, as discussed next.

### 3.3. Pseudo-Label Selection

The VO model from Sec. 3.2 can be used to obtain pseudo-labels over an unlabeled (i.e., with respect to the main VO task) data  $\mathcal{D}_U$ . However, incorrect predictions can introduce noise and heavily degrade model training [81, 89]. Hence, it is crucial to remove low-confidence samples prior to the cross-modal self-training.

In our regression problem, we measure the confidence of a pseudo-label based on the entropy of the predicted matrix Fisher distribution (i.e., a lower entropy represents increased confidence),

$$H(p(\mathbf{R}|\hat{\Psi})) < \tau_u \quad (5)$$

where we set a fixed threshold  $\tau_u$  to ensure the network prediction is sufficiently certain to be selected. To generate pseudo-labels, the VO model is tested on out-of-domain data with highly diverse and dynamic scenes. Based on our analysis in Sec. 4, we find the uncertainty-aware selection mechanism to be crucial for robust self-training irrespective of the auxiliary training tasks.

### 3.4. Self-Training with Auxiliary Tasks

To learn effective representations for generalized VO at scale, we propose to incorporate supervision from auxiliary but potentially complementary prediction tasks in addition to the generated VO pseudo-labels on  $\mathcal{D}_U$ . The introduced auxiliary tasks regularize the self-training process, particularly in cases where VO pseudo-labels may be inaccurate but other modalities may contain relevant information for reducing ambiguity. Our approach is motivated by the success of multi-task frameworks for computer vision tasks [2, 32, 43, 120, 121, 128]. However, we emphasize that related studies often leverage high-quality annotated labels and not noisy pseudo-labels based on model predictions. We sought to incorporate useful auxiliary tasks, as unrelated or noisy supervision can impede the learning process and result in a detrimental effect on the main task model. We set the auxiliary labeled task as an audio prediction task  $\mathbf{y}_i^{aux} := \mathbf{A}_i \in \mathcal{R}^{2 \times L}$ , and the auxiliary pseudo-labeled tasks  $\tilde{\mathbf{y}}_i^{aux} := [\tilde{\mathbf{S}}_i, \tilde{\mathbf{D}}_i, \tilde{\mathbf{F}}_i] \in \mathcal{R}^{W \times H \times C}$  as segmentation, depth, and flow prediction, respectively. Subsequently, we leverage multi-task learning (as shown in Fig. 3)



Figure 3: **Illustration of the Importance of Audio.** The frame is consistent with the red arrow marked on the waveform. Left: audio amplitude decreases and maintains a low level when the vehicle is going to wait for traffic lights. Right: audio experiences many ups and downs representing acceleration and brake in a narrow urban area.

and minimize a loss composed of four terms,

$$\mathcal{L}_{aux} = \lambda_a \mathcal{L}_{audio} + \lambda_s \mathcal{L}_{seg} + \lambda_f \mathcal{L}_{flow} + \lambda_d \mathcal{L}_{depth} \quad (6)$$

over the entire dataset  $\tilde{\mathcal{D}}$ . We note that we drop the explicit label source to avoid clutter. Next, we define each term and corresponding decoder. We empirically observe the additional tasks to improve generalization in evaluation, both within and across VO datasets.

**Audio Decoder:** We utilize audio labels, generally available for online videos, as an auxiliary prediction task. We note that prior work often studies such cross-modal reasoning for basic navigation scenarios [12, 13, 28] and not for in-the-wild videos where dense dynamic objects may generate significant ambient noise. In our settings, an audio signal can provide complementary information to visual information regarding the overall traffic scenario as well as ego-speed. This insight will be affirmed by our findings in Sec. 4, where the audio task is shown to provide synergistic supervision, both for the main VO task and when combined with other auxiliary tasks. For instance, Fig. 3 depicts how an idling ego-vehicle may generate lower audio levels, which, in conjunction with the visual scene features, can help disambiguate ego-motion from surrounding motion when stopped at intersections. As drift due to surrounding motion is a common failure mode for VO models, we further incorporate a segmentation task for dynamic objects below.

Our audio decoder is based on a 1D U-Net architecture [85], consisting of a residual 1D convolutional block [34] and an attention block [98], and reconstructs the dual-channel raw audio of the two input frames using the encoder features. We employ a two-term MSE and spectral loss,

$$\mathcal{L}_{audio} = \|\mathbf{A}_i - \hat{\mathbf{A}}_i\|_2^2 + \|\text{FT}(\mathbf{A}_i) - \text{FT}(\hat{\mathbf{A}}_i)\|_2^2 \quad (7)$$

where FT represents the short-time Fourier transform [20].

**Segmenting Dynamic Objects:** The relative motion caused by dynamic objects can often lead to inaccurate

pose predictions, e.g., when stopped at a traffic light with oncoming traffic. To facilitate disambiguating potentially dynamic objects from the static background, we incorporate a segmentation prediction for pedestrian and vehicle classes [4]. As this task involves extensive manual annotation, we intentionally do not assume it is provided as part of the originally labeled dataset  $\mathcal{D}_L$  and instead leverage an off-the-shelf model based on Mask R-CNN [33]. The model is pre-trained on the COCO dataset [53]. We use the detector to construct a pseudo-label semantic segmentation  $\tilde{\mathbf{S}} \in \mathcal{R}^{W \times H \times 2}$  of foreground and background in the two input frames. We leverage an FCN [55] decoder, consisting of 11 transposed convolutional layers followed by a convolutional layer and a final sigmoid activation function, and minimize a Dice loss,

$$\mathcal{L}_{seg} = 1 - 2 \frac{\sum_{j,k,c} \tilde{\mathbf{S}}_i \circ \hat{\mathbf{S}}_i}{\sum_{j,k,c} \tilde{\mathbf{S}}_i^2 + \sum_{j,k,c} \hat{\mathbf{S}}_i^2} \quad (8)$$

where  $\hat{\mathbf{S}}_i \in \mathbb{R}^{H \times W \times 2}$  is the decoder predicted segmentation, and  $j \in [1, H], k \in [1, W], c \in [1, 2]$ . As dynamic objects often cause ego-motion estimation drift, the prediction task can regularize self-training by providing a useful invariant prior (i.e., across datasets and settings) of background and foreground knowledge. Moreover, the segmentation task complements the audio task in many scenarios as dynamic objects may also generate ambient audio.

**Depth and Flow Tasks:** We explore two additional auxiliary tasks based on depth and optical flow estimation, as they potentially offer valuable information about the structure of the surroundings and the camera motion and are frequently employed in VO tasks [56, 66, 73, 103]. We utilize an MSE as the loss function for both depth and flow tasks,

$$\begin{aligned} \mathcal{L}_{flow} &= \|\tilde{\mathbf{F}}_i - \hat{\mathbf{F}}_i\|_2^2 \\ \mathcal{L}_{depth} &= \|\tilde{\mathbf{D}}_i - \hat{\mathbf{D}}_i\|_2^2 \end{aligned} \quad (9)$$

To simplify the model, we maintain the identical decoder structure used as in the dynamic object segmentation task (see Fig. 2), with the exception of eliminating the final Sigmoid layer.

### 3.5. Implementation Details

Our models are trained using three NVIDIA RTX 3090 GPUs using a batch size of six. The learning rate is set to 0.001 and with decay 0.99. Given the main VO objective, we set  $\lambda_\theta = 1$  and  $\lambda_u = 0.1$ . Remaining auxiliary loss hyper-parameters, i.e.,  $\lambda_a, \lambda_s, \lambda_f, \lambda_d$ , are set to 0.01. For our semi-supervised training, we obtain a diverse set of 59,000 unlabeled samples across different geographical locations, times of day, and environmental conditions. We split the nuScenes benchmark [7] into training, validation, and evaluation sets, to train an initial teacher model for 15

epochs. The student model is trained for 15 epochs on a mix of labeled nuScenes and pseudo-labeled YouTube data. We note that we do not employ careful ratio optimization [8] when mixing the datasets without and instead solely rely on the uncertainty-based selection mechanism. We leverage data augmentation strategies, including random cropping and resizing, for improving generalization and simulating varying camera intrinsics [103]. During inference, the model runs at 77 FPS on a single NVIDIA GTX 3090 GPU. Additional details regarding the training and experiments can be found in the supplementary.

## 4. Experiments

In this section, we comprehensively analyze our XVO framework. As our goal is to build generalized VO systems, we emphasize generalization ability across different datasets with various camera setups, specifically in the context of varying autonomous driving settings.

**Datasets:** To understand the role of cross-modal self-training on model generalization, we evaluate our proposed XVO method using three commonly employed datasets, KITTI [30], nuScenes [7] and Argoverse 2 [107]. Out of the three, KITTI is the most popular VO benchmark, consisting of 11 sequences 00-10 with ground truth. As KITTI is an older benchmark (2012), its camera intrinsics vary significantly from the other two benchmarks. nuScenes consists of about 15 hours of driving data (totaling 197,000 images) from four regions in Boston and Singapore: Boston-Seaport, Singapore-OneNorth, Singapore-Queenstown, and Singapore-HollandVillage. In contrast to KITTI which was captured in sunny driving with mostly static objects, nuScenes incorporates complex real-world driving maneuvers in dense streets and various conditions, e.g., nighttime, difficult illumination conditions with low visibility, as well as artifacts on the camera lens, such as rain droplets or dirt. Finally, Argoverse 2 is a large dataset with 1,000 driving sequences across six US cities. We leverage a test dataset that includes 150 sequences and 48,022 images.

**Procedure and Baselines:** We generally train within one region on nuScenes (HollandVillage) and evaluate the remaining regions and datasets. This is in contrast to prior evaluation procedures where models can learn to memorize the scale and camera setup without generalization through training and testing on the same camera setup and similar environments. We also directly compare with prior state-of-the-art using the standard KITTI protocol [48, 119]. As our approach does not leverage known intrinsics, we separate approaches that do assume such knowledge in their pipeline to ensure meaningful analysis. We further indicate whether methods predict pose with absolute scale, as some methods output up-to-scale estimates and use the ground-truth scale to align and evaluate their model, e.g., TartanVO [103].

Nonetheless, TartanVO is one of the few approaches that have shown generalization across datasets without the need to fine-tune or perform online adaptation strategies and is therefore our main baseline.

**Metrics:** We follow standard evaluation metrics of average translation  $t_{rel}$  (in %) and rotation  $r_{rel}$  (in degrees per 100 meters) errors, computed over all possible subsequences within a test sequence of lengths (100, ..., 800) meters [30, 103]. We refer the reader to the KITTI leaderboard for more details regarding the metric. However, we observe prior measures to only provide a proxy evaluation of real-world scale predictions as the errors could potentially vary along the trajectory independently of trajectory-level measures (our supplementary contains additional details). To explore the benefits of semi-supervised training on real-world scale estimation, we sought to directly quantify scale recovery within consecutive frames in a single metric. We therefore also report the *average scale error (se)* over predicted and ground-truth translation,

$$se = 1 - \min \left( \frac{\|\hat{\mathbf{t}}\|_2}{\max(\|\mathbf{t}\|_2, \epsilon)}, \frac{\|\mathbf{t}\|_2}{\max(\|\hat{\mathbf{t}}\|_2, \epsilon)} \right) \quad (10)$$

where  $\epsilon$  prevents dividing by zero.

### 4.1. Results

We examine the role of the main components in the proposed framework below. Complete ablation, e.g., across modality combinations and training settings, can be found in the supplementary.

**Teacher Model Performance:** Table 1 compares our proposed encoder architecture for the main VO task with prior methods. When trained in a supervised learning manner on KITTI, our teacher model achieves the lowest translation error of 3.4% even without access to camera intrinsics or multi-step optimization. This suggests that basic modifications to underlying network structure, e.g., through an improved encoder and attention-based mechanism, can result in significant gains for the monocular VO task. Given the effective network structure, we now turn to analyzing the benefits of the proposed semi-supervised framework.

**Semi-Supervised VO Training:** Table 1 also analyzes the generalization performance of the proposed semi-supervised learning framework on KITTI. Specifically, we show our initial teacher model that is trained on the nuScenes (HollandVillage) dataset to not generalize well to the KITTI testing set (25.27% translation and 8.17° rotation error) due to domain shift and differing camera settings. However, after semi-supervised training, the errors for the student model are reduced by 40% and 50% in translation and rotation errors, respectively. The best self-trained model with auxiliary tasks (complete ablation can be found

Table 1: **Analysis on the KITTI Benchmark.** We abbreviate ‘intrinsic-free’ as I (i.e., a method which does not assume the intrinsics) and ‘real-world scale’ as S (i.e., a method is able to recover real-world scale). To ensure meaningful comparison, we categorize models based on supervision type. Firstly, we present unsupervised learning methods, followed by supervised learning methods, then generalized VO methods, and finally our XVO ablation. In the case of TartanVO, we analyze robustness to noise applied to the intrinsics. We train two teacher models: one based on KITTI (as shown in supervised learning approaches) and the other on nuScenes (as displayed at the end of the Table with ablations).

Method	I	S	Seq 03		Seq 04		Seq 05		Seq 06		Seq 07		Seq 10		Avg	
			$t_{rel}$	$r_{rel}$	$t_{rel}$	$r_{rel}$	$t_{rel}$	$r_{rel}$	$t_{rel}$	$r_{rel}$	$t_{rel}$	$r_{rel}$	$t_{rel}$	$r_{rel}$	$t_{rel}$	$r_{rel}$
<i>Unsupervised Methods:</i>																
SfMLearner [133]	✗	✗	10.78	3.92	4.49	5.24	18.67	4.10	25.88	4.80	21.33	6.65	14.33	3.30	15.91	4.67
GeoNet [119]	✗	✗	19.21	9.78	9.09	7.55	20.12	7.67	9.28	4.34	8.27	5.93	20.73	9.10	14.45	7.40
Zhan et al. [122]	✗	✓	15.76	10.62	3.14	2.02	4.94	2.34	5.80	2.06	6.49	3.56	12.82	3.40	8.16	4.00
UnDeepVO [48]	✗	✓	5.00	6.17	4.49	2.13	3.40	1.50	6.20	1.98	3.15	2.48	10.63	4.65	5.48	3.15
<i>Supervised Methods:</i>																
DeepVO [101]	✓	✓	8.49	6.89	7.19	4.97	2.62	3.61	5.42	5.82	3.91	4.60	8.11	8.83	5.96	5.79
ESP-VO [102]	✓	✓	6.72	6.46	6.33	6.08	3.35	4.93	7.24	7.29	3.52	5.02	9.77	10.2	6.16	6.66
GFS-VO [112]	✓	✓	5.44	3.32	2.91	1.30	3.27	1.62	8.50	2.74	3.37	2.25	6.32	2.33	4.97	2.26
Xue et al. [113]	✓	✓	<b>3.32</b>	2.10	2.96	1.76	2.59	1.25	4.93	1.90	<b>3.07</b>	<b>1.76</b>	3.94	1.72	3.47	<b>1.75</b>
<b>Our Teacher (KITTI)</b>	✓	✓	3.46	<b>2.00</b>	<b>1.67</b>	<b>0.70</b>	<b>2.12</b>	<b>0.92</b>	<b>3.92</b>	<b>1.46</b>	5.93	3.96	<b>3.31</b>	<b>1.52</b>	<b>3.40</b>	1.76
<i>Baseline Generalized VO Methods:</i>																
TartanVO [103]	✗	✗	4.20	2.80	6.19	4.35	5.84	3.24	4.21	2.51	7.11	4.96	8.00	3.21	5.93	3.51
TartanVO (10% Noise)	✗	✗	9.33	3.12	10.88	4.71	11.77	5.39	11.88	4.52	14.70	10.74	11.76	3.61	11.72	5.35
TartanVO (20% Noise)	✗	✗	17.79	4.42	21.58	5.04	20.12	8.54	18.80	6.26	21.34	16.27	17.45	5.03	19.51	7.59
TartanVO (30% Noise)	✗	✗	25.89	7.06	34.91	4.54	22.48	10.17	19.32	5.23	19.40	13.33	25.06	8.43	24.51	8.13
<i>Proposed Generalized VO Methods:</i>																
Our Teacher (nuScenes)	✓	✓	26.78	4.92	26.02	2.42	23.65	8.85	23.97	6.47	30.66	20.32	20.57	6.01	25.27	8.17
Student w/o Filter	✓	✓	26.98	9.68	22.56	2.15	14.77	5.83	<b>11.38</b>	<b>1.62</b>	16.45	9.35	20.23	8.99	18.73	6.27
Student	✓	✓	20.30	3.97	16.33	1.57	11.12	4.19	15.60	5.69	7.77	3.48	19.91	5.59	15.17	4.08
<b>XVO</b>	✓	✓	<b>14.53</b>	<b>3.93</b>	<b>16.29</b>	<b>0.96</b>	<b>8.31</b>	<b>2.76</b>	15.31	5.49	<b>5.86</b>	<b>3.00</b>	<b>12.17</b>	<b>3.45</b>	<b>12.08</b>	<b>3.27</b>

Table 2: **Average Quantitative Results across Datasets.** We test on KITTI (sequences 00-10), Argoverse 2, and the unseen regions in nuScenes. All results are the average over all scenes. We present translation error, rotation error and scale error. Approaches such as TartanVO do not estimate real-world scale but may be aligned with ground truth (GT) scale in evaluation. A, S, F, D are the abbreviation of Audio, Seg, Flow, Depth.

Method	KITTI 00-10			Argoverse 2			nuScenes		
	$t_{err}$	$r_{err}$	$se$	$t_{err}$	$r_{err}$	$se$	$t_{err}$	$r_{err}$	$se$
<i>Baseline Generalized VO Methods:</i>									
TartanVO w/ GT Align	6.37	3.32	/	8.55	5.77	/	9.61	6.83	/
TartanVO w/o GT Align	21.67	3.33	0.29	41.11	5.77	0.40	28.23	6.83	0.29
<i>Proposed Generalized VO Methods:</i>									
Teacher (nuScenes)	26.16	6.84	0.25	10.89	3.40	0.16	15.93	6.73	0.20
Student w/o Filter	20.64	5.68	0.21	10.80	7.33	0.14	9.32	4.60	0.14
Student	17.04	4.02	0.16	9.16	3.40	0.14	10.54	3.94	0.13
Student+Seg	16.31	3.77	0.16	9.17	<b>3.18</b>	0.13	11.35	4.05	0.14
Student+Flow	15.60	3.19	0.19	9.04	4.45	0.13	<b>9.13</b>	4.06	<b>0.13</b>
Student+Depth	17.49	3.89	0.20	9.25	4.11	0.13	11.86	6.46	0.15
Student+Audio	<b>14.37</b>	<b>3.06</b>	<b>0.16</b>	<b>8.00</b>	<b>3.08</b>	<b>0.12</b>	<b>9.26</b>	<b>3.20</b>	<b>0.12</b>
Student+Audio+Seg	<b>14.20</b>	<b>3.02</b>	<b>0.16</b>	8.67	3.63	0.13	11.29	3.70	0.14
Student+S+F+D	18.23	3.88	0.21	8.79	4.89	0.13	8.93	<b>3.44</b>	0.13
Student+A+S+F+D	16.74	4.40	0.18	<b>7.89</b>	3.54	<b>0.12</b>	9.98	4.36	0.15

in the supplementary) results in further student performance gains, e.g., a further reduction in translation error by 20%.

We also compare with the most related TartanVO [103] baseline which utilizes the ground-truth for scale alignment and has access to camera intrinsics. However, even with the ground-truth alignment, TartanVO exhibits quick degradation with minimal noise in the intrinsics (enabling a more fair comparison as our method is not provided these as input). Moreover, we explore the generalization of our training framework by evaluating on various datasets in Table 2. We emphasize that none of the trained models have access to samples from Argoverse 2 or KITTI dataset during training. By predicting real-world scale, our student model with all auxiliary tasks outperforms the baseline TartanVO in all three datasets, e.g., by 80% in translation and 70% scale error on Argoverse 2, without any ground-truth alignment. This indicates the proposed method to improve reasoning over scale and scene semantics across arbitrary conditions.

**Impact of Uncertainty-Aware Sample Selection:** When inspecting the various pseudo-labels, we observed many cases of drift and incorrect predictions due to the harsh generalization settings. Hence, the uncertainty-aware pseudo-label selection mechanism plays a crucial role in the semi-supervised learning process. As shown in Table 1 and Table 2, discarding pseudo-labels with low confidence consistently improves performance, both with and without multi-modal supervision. We notice how a student model without the uncertainty-aware sample removal (i.e., ‘Student w/o

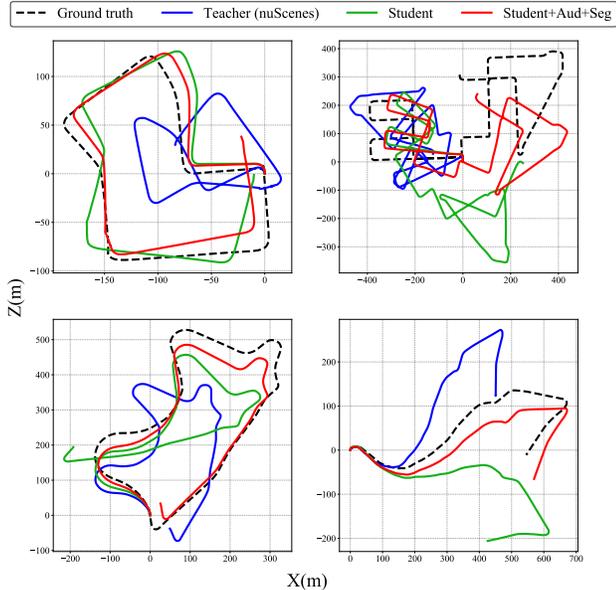


Figure 4: **Qualitative Analysis on KITTI.** We find that incorporating audio and segmentation tasks as part of the semi-supervised learning process significantly improves ego-pose estimation on KITTI.

Filter’) provides only mild improvements compared to the teacher. Once noisy samples are filtered out of the dataset, the performance on KITTI and nuScenes improve significantly, e.g., from 26.16 to 17.04 and 15.93 to 10.54 translation error respectively.

**Ablation on Auxiliary Tasks:** We sought to understand the role of the various explored auxiliary tasks, i.e., audio, segmentation, depth, and flow. We first analyze the impact of adding an audio reconstruction task for the VO problem. Although extracted audio includes some ambient noise, we can see that XVO consistently benefits from the proposed audio supervision across the evaluation datasets. This can be explained by the consistent quality of the ground-truth audio labels, i.e., when compared to the noise in pseudo-labels generated by the auxiliary prediction models on our unconstrained videos. In general, we find that audio, segmentation, and flow tasks result in better performance when compared to the depth prediction task. While prior research often leverages monocular depth prediction for improving VO on KITTI, this is a significantly challenging task in more general settings which results in noisier pseudo-labels. We also investigate various combinations of auxiliary branches and find the combination of segmentation and audio branch performs better than a single auxiliary task on KITTI. While this is encouraging, KITTI contains simpler scenarios with relatively few dynamic traffic participants. In such simpler settings, our segmentation branch can be used to obtain reliable pseudo-labels and learn effi-

cient generalized features. However, this finding does not extend to nuScenes and Argoverse 2 which frequently contain dense and dynamic scenes. We also find that simply adding prediction tasks does not provide further gains due to the pseudo-label noise and a more brittle and difficult optimization process. Complete ablations on auxiliary tasks can be seen in our supplementary.

## 4.2. Qualitative Results

Fig. 4 depicts the prediction of driving trajectories on KITTI sequences 7, 8, 9, and 10. The trajectory predicted using the teacher model that is trained on nuScenes is not able to recover scale accurately. Due to the semi-supervised training process, the student model is shown to have better scale recovery and generalization despite the lack of calibration knowledge. Nonetheless, the student model fails to estimate accurate rotation in more challenging scenes on KITTI, e.g., top right and bottom left scenarios. Finally, the cross-modal trained model is shown to robustly estimate translation, rotation, and scale, even in the most complicated route in Fig. 4-top right. Additional qualitative examples are provided in the supplementary.

## 5. Conclusion

In this paper, we present XVO, a novel method for generalized visual odometry estimation via cross-modal self-training. Our efficient network structure achieves state-of-the-art results on KITTI, despite having no knowledge of camera parameters or multi-frame optimization as in prior methods. Moreover, our framework leverages a mixed-label semi-supervised setting over a large dataset of internet videos to further enhance generalization performance. Specifically, we show that additional auxiliary segmentation and audio reconstruction tasks can significantly impact cross-dataset generalization. Our trained VO models can be used across platforms and settings without fine-tuning, i.e., due to general reasoning over semantic visual characteristics of scenes. Moreover, our training settings of improving the performance of a model that is initially trained on a small and restricted dataset are broadly applicable to various robotics use-cases. We hope our work can inspire future researchers to explore scalable VO models that can benefit a broad range of applications. Given the limited utility of combining multiple auxiliary tasks in our settings, a future direction would be to study better methods for learning from noisy and diverse auxiliary pseudo-labels. Moreover, while we achieved state-of-the-art results with a two-frame approach, multi-frame optimization could provide further benefits by alleviating drift.

**Acknowledgments:** We thank the Red Hat Collaboratory and National Science Foundation (IIS-2152077) for supporting this research.

## References

- [1] Triantafyllos Afouras, Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman. Deep audio-visual speech recognition. In *PAMI*, 2018. 3
- [2] Roman Bachmann, David Mizrahi, Andrei Atanov, and Amir Zamir. Multimaec: Multi-modal multi-task masked autoencoders. In *ECCV*, 2022. 3, 4
- [3] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. In *PAMI*, 2018. 3
- [4] Aseem Behl, Kashyap Chitta, Aditya Prakash, Eshed Ohn-Bar, and Andreas Geiger. Label-efficient visual abstractions for autonomous driving. In *IROS*, 2020. 2, 5
- [5] Aseem Behl, Despoina Paschalidou, Simon Donné, and Andreas Geiger. Pointflownet: Learning representations for rigid motion estimation from point clouds. In *CVPR*, 2019. 3
- [6] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. In *NeurIPS*, 2019. 2
- [7] 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 *CVPR*, 2020. 1, 5, 6
- [8] Benjamin Caine, Rebecca Roelofs, Vijay Vasudevan, Jiquan Ngiam, Yuning Chai, Zhifeng Chen, and Jonathon Shlens. Pseudo-labeling for scalable 3d object detection. *arXiv preprint arXiv:2103.02093*, 2021. 2, 3, 6
- [9] Carlos Campos, Richard Elvira, Juan J Gómez Rodríguez, José MM Montiel, and Juan D Tardós. Orb-slam3: An accurate open-source library for visual, visual-inertial, and multimap slam. *T-RO*, 37(6):1874–1890, 2021. 2
- [10] Yair Carmon, Aditi Raghunathan, Ludwig Schmidt, John C Duchi, and Percy S Liang. Unlabeled data improves adversarial robustness. In *NeurIPS*, 2019. 2
- [11] Angel Chang, Will Monroe, Manolis Savva, Christopher Potts, and Christopher D Manning. Text to 3d scene generation with rich lexical grounding. *arXiv preprint arXiv:1505.06289*, 2015. 3
- [12] Changan Chen, Ziad Al-Halah, and Kristen Grauman. Semantic audio-visual navigation. In *CVPR*, 2021. 5
- [13] Changan Chen, Unnat Jain, Carl Schissler, Sebastia Vicenc Amengual Gari, Ziad Al-Halah, Vamsi Krishna Ithapu, Philip Robinson, and Kristen Grauman. Soundspaces: Audio-visual navigation in 3d environments. In *ECCV*, 2020. 3, 5
- [14] Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L Yuille, and Yuyin Zhou. Transunet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306*, 2021. 4, 13
- [15] Amirhossein Dadashzadeh, Alan Whone, and Majid Mirmehdi. Auxiliary learning for self-supervised video representation via similarity-based knowledge distillation. In *CVPR*, 2022. 3
- [16] Keqi Deng, Songjun Cao, Yike Zhang, Long Ma, Gaofeng Cheng, Ji Xu, and Pengyuan Zhang. Improving ctc-based speech recognition via knowledge transferring from pre-trained language models. In *ICASSP*, 2022. 3
- [17] Zhun Deng, Linjun Zhang, Amirata Ghorbani, and James Zou. Improving adversarial robustness via unlabeled out-of-domain data. In *AISTAT*, 2021. 2
- [18] Christophe Denis and Mohamed Hebiri. Consistency of plug-in confidence sets for classification in semi-supervised learning. *J. Nonparametric Stat.*, 32(1):42–72, 2020. 2
- [19] Lucio M Dery, Paul Michel, Mikhail Khodak, Graham Neubig, and Ameet Talwalkar. Aang: Automating auxiliary learning. *arXiv preprint arXiv:2205.14082*, 2022. 3
- [20] Prafulla Dhariwal, Heewoo Jun, Christine Payne, Jong Wook Kim, Alec Radford, and Ilya Sutskever. Jukebox: A generative model for music. *arXiv preprint arXiv:2005.00341*, 2020. 5
- [21] Shangzhe Di, Zeren Jiang, Si Liu, Zhaokai Wang, Leyan Zhu, Zexin He, Hongming Liu, and Shuicheng Yan. Video background music generation with controllable music transformer. In *ACMMM*, 2021. 3
- [22] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *arXiv preprint arXiv:2010.11929*, 2020. 13
- [23] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 4
- [24] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *ICCV*, 2015. 13
- [25] Jakob Engel, Vladlen Koltun, and Daniel Cremers. Direct sparse odometry. In *PAMI*, 2017. 1, 2, 3
- [26] Jakob Engel, Thomas Schöps, and Daniel Cremers. Lsd-slam: Large-scale direct monocular slam. In *ECCV*, 2014. 1, 2, 3
- [27] Mahdi Abolfazli Esfahani, Han Wang, Keyu Wu, and Shenghai Yuan. Aboldeepio: A novel deep inertial odometry network for autonomous vehicles. *IEEE Trans. Intell. Transp. Syst.*, 21(5):1941–1950, 2019. 2
- [28] Chuang Gan, Yiwei Zhang, Jiajun Wu, Boqing Gong, and Joshua B Tenenbaum. Look, listen, and act: Towards audio-visual embodied navigation. In *ICRA*, 2020. 5
- [29] Ralph Gasser, Luca Rossetto, and Heiko Schuldt. Multi-modal multimedia retrieval with vitrivr. In *ICMR*, 2019. 3
- [30] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012. 1, 2, 6, 14
- [31] Ariel Gordon, Hanhan Li, Rico Jonschkowski, and Anelia Angelova. Depth from videos in the wild: Unsupervised monocular depth learning from unknown cameras. In *CVPR*, 2019. 1
- [32] Daniel Gordon, Abhishek Kadian, Devi Parikh, Judy Hoffman, and Dhruv Batra. Splitnet: Sim2sim and task2task

- transfer for embodied visual navigation. In *ICCV*, 2019. 4
- [33] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 5
- [34] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5
- [35] Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. Using self-supervised learning can improve model robustness and uncertainty. In *NeurIPS*, 2019. 2
- [36] Zanming Huang, Zhongkai Shanguan, Jimuyang Zhang, Gilad Bar, Matthew Boyd, and Eshed Ohn-Bar. ASSIS-TER: Assistive navigation via conditional instruction generation. In *ECCV*, 2022. 3
- [37] Bogdan Ionescu, Henning Müller, Renaud Péteri, Johannes Rückert, Asma Ben Abacha, Alba G Seco de Herrera, Christoph M Friedrich, Louise Bloch, Raphael Brüngel, Ahmad Idrissi-Yaghir, et al. Overview of the imageclef 2022: Multimedia retrieval in medical, social media and nature applications. In *CLEF*, 2022. 3
- [38] Mobarakol Islam, Daniel Anojan Atputharuban, Ravikiran Ramesh, and Hongliang Ren. Real-time instrument segmentation in robotic surgery using auxiliary supervised deep adversarial learning. *RA-L*, 4(2):2188–2195, 2019. 3
- [39] Ganesh Iyer, J Krishna Murthy, Gunshi Gupta, Madhava Krishna, and Liam Paull. Geometric consistency for self-supervised end-to-end visual odometry. In *CVPRW*, 2018. 1, 2, 3
- [40] Jisoo Jeong, Seungeui Lee, Jeesoo Kim, and Nojun Kwak. Consistency-based semi-supervised learning for object detection. In *NeurIPS*, 2019. 2
- [41] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *CVPR*, 2018. 4
- [42] Alex Kendall, Matthew Grimes, and Roberto Cipolla. Posenet: A convolutional network for real-time 6-dof camera relocalization. In *ICCV*, 2015. 2
- [43] Iasonas Kokkinos. Ubertnet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory. In *CVPR*, 2017. 3, 4
- [44] Rainer Kümmerle, Bastian Steder, Christian Dornhege, Michael Ruhnke, Giorgio Grisetti, Cyrill Stachniss, and Alexander Kleiner. On measuring the accuracy of slam algorithms. *Autonomous Robots*, 27, 2009. 14
- [45] Po-Nien Kung, Sheng-Siang Yin, Yi-Cheng Chen, Tse-Hsuan Yang, and Yun-Nung Chen. Efficient multi-task auxiliary learning: selecting auxiliary data by feature similarity. In *EMNLP*, 2021. 3
- [46] Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *ICMLW*, 2013. 2, 3
- [47] Qimai Li, Xiao-Ming Wu, Han Liu, Xiaotong Zhang, and Zhichao Guan. Label efficient semi-supervised learning via graph filtering. In *CVPR*, 2019. 2
- [48] Ruihao Li, Sen Wang, Zhiqiang Long, and Dongbing Gu. Undeepvo: Monocular visual odometry through unsupervised deep learning. In *ICRA*, 2018. 1, 2, 3, 6, 7
- [49] Shunkai Li, Xin Wu, Yingdian Cao, and Hongbin Zha. Generalizing to the open world: Deep visual odometry with online adaptation. In *CVPR*, 2021. 1, 3
- [50] Shunkai Li, Fei Xue, Xin Wang, Zike Yan, and Hongbin Zha. Sequential adversarial learning for self-supervised deep visual odometry. In *ICCV*, 2019. 1, 2, 3
- [51] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *ECCV*, 2020. 3
- [52] Lukas Liebel and Marco Körner. Auxiliary tasks in multi-task learning. *arXiv preprint arXiv:1805.06334*, 2018. 3
- [53] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 5
- [54] Shikun Liu, Andrew Davison, and Edward Johns. Self-supervised generalisation with meta auxiliary learning. In *NeurIPS*, 2019. 3
- [55] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 5, 13
- [56] Shing Yan Loo, Ali Jahani Amiri, Syamsiah Mashohor, Sai Hong Tang, and Hong Zhang. Cnn-svo: Improving the mapping in semi-direct visual odometry using single-image depth prediction. In *ICRA*, 2019. 5
- [57] Reza Mahjourian, Martin Wicke, and Anelia Angelova. Un-supervised learning of depth and ego-motion from monocular video using 3d geometric constraints. In *CVPR*, 2018. 1, 2
- [58] Gideon S Mann and Andrew McCallum. Simple, robust, scalable semi-supervised learning via expectation regularization. In *ICML*, 2007. 2
- [59] Kanti V Mardia, Peter E Jupp, and KV Mardia. *Directional statistics*, volume 2. Wiley Online Library, 2000. 4
- [60] Viorica Marian, Sayuri Hayakawa, and Scott R Schroeder. Cross-modal interaction between auditory and visual input impacts memory retrieval. *Front. Neurosci.*, 15:661477, 2021. 2
- [61] Sarah Maslen. “hearing” ahead of the sound: How musicians listen via proprioception and seen gestures in performance. *Senses Soc.*, 2022. 2
- [62] Alexander Mey and Marco Loog. Improved generalization in semi-supervised learning: A survey of theoretical results. In *PAMI*, 2022. 2
- [63] Ishan Misra, Abhinav Shrivastava, Abhinav Gupta, and Martial Hebert. Cross-stitch networks for multi-task learning. In *CVPR*, 2016. 3
- [64] Takeru Miyato, Andrew M Dai, and Ian Goodfellow. Adversarial training methods for semi-supervised text classification. *arXiv preprint arXiv:1605.07725*, 2016. 2
- [65] David Mohlin, Josephine Sullivan, and Gérald Bianchi. Probabilistic orientation estimation with matrix fisher distributions. In *NeurIPS*, 2020. 4
- [66] Peter Muller and Andreas Savakis. Flowdometry: An optical flow and deep learning based approach to visual odometry. In *WACV*, 2017. 5
- [67] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *T-RO*, 31(5):1147–1163, 2015. 2
- [68] Raul Mur-Artal and Juan D Tardós. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cam-

- eras. *T-RO*, 33(5):1255–1262, 2017. 1, 2
- [69] Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *ECCV*, 2012. 13
- [70] Ishan Nigam. Learning from auxiliary supervision. Master’s thesis, Carnegie Mellon University, 2018. 3
- [71] Casey O’Callaghan. Seeing what you hear: Cross-modal illusions and perception. *Philos.*, 18:316–338, 2008. 2
- [72] Eshed Ohn-Bar, Aditya Prakash, Aseem Behl, Kashyap Chitta, and Andreas Geiger. Learning situational driving. In *CVPR*, 2020. 3
- [73] Tejas Pandey, Dexmont Pena, Jonathan Byrne, and David Moloney. Leveraging deep learning for visual odometry using optical flow. *J. Sens.*, 21(4):1313, 2021. 5
- [74] Mohammad Peikari, Sherine Salama, Sharon Nofech-Mozes, and Anne L Martel. A cluster-then-label semi-supervised learning approach for pathology image classification. *Sci. Rep.*, 8(1):1–13, 2018. 2
- [75] Jean Piaget, Margaret Cook, et al. *The origins of intelligence in children*, volume 8. IUP Press NY, 1952. 2
- [76] Jin-Chun Piao and Shin-Dug Kim. Real-time visual-inertial slam based on adaptive keyframe selection for mobile ar applications. *TMM*, 21(11):2827–2836, 2019. 2
- [77] Sabrina Pitzalis, Stefano Sdoia, Alessandro Bultrini, Giorgia Committeri, Francesco Di Russo, Patrizia Fattori, Claudio Galletti, and Gaspare Galati. Selectivity to translational egomotion in human brain motion areas. *PLoS one*, 8(4), 2013. 2
- [78] David Prokhorov, Dmitry Zhukov, Olga Barinova, Konushin Anton, and Anna Vorontsova. Measuring robustness of visual slam. In *ICMVA*, 2019. 14
- [79] Monique Radeau. Auditory-visual spatial interaction and modularity. *CPC*, 13(1):3–51, 1994. 2
- [80] Anurag Ranjan, Varun Jampani, Lukas Balles, Kihwan Kim, Deqing Sun, Jonas Wulff, and Michael J Black. Competitive collaboration: Joint unsupervised learning of depth, camera motion, optical flow and motion segmentation. In *CVPR*, 2019. 2, 3
- [81] Mamshad Nayeem Rizve, Kevin Duarte, Yogesh S Rawat, and Mubarak Shah. In defense of pseudo-labeling: An uncertainty-aware pseudo-label selection framework for semi-supervised learning. In *ICLR*, 2021. 2, 3, 4
- [82] Chris Rockwell, Justin Johnson, and David F Fouhey. The 8-point algorithm as an inductive bias for relative pose prediction by vits. *arXiv preprint arXiv:2208.08988*, 2022. 1, 2, 3
- [83] Sebastian Ruder. An overview of gradient descent optimization algorithms. In *arXiv preprint arXiv:1609.04747*, 2016. 13
- [84] Alexander Sax, Bradley Emi, Amir R Zamir, Leonidas Guibas, Silvio Savarese, and Jitendra Malik. Mid-level visual representations improve generalization and sample efficiency for learning active tasks. In *CoRL*, 2019. 3
- [85] Flavio Schneider, Zhijing Jin, and Bernhard Schölkopf. Mousai: Text-to-music generation with long-context latent diffusion. *arXiv preprint arXiv:2301.11757*, 2023. 5
- [86] Dmitriy Serdyuk, Otavio Braga, and Olivier Siohan. Audio-visual speech recognition is worth  $32 \times 32 \times 8$  voxels. In *ASRU*, 2021. 3
- [87] Zhongkai Shanguan, Lei Lin, Wencheng Wu, and Beilei Xu. Neural process for black-box model optimization under bayesian framework. *arXiv preprint arXiv:2104.02487*, 2021. 13
- [88] Linda Smith and Michael Gasser. The development of embodied cognition: Six lessons from babies. *Artif. Life*, 11(1-2):13–29, 2005. 2
- [89] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. In *NeurIPS*, 2020. 2, 4
- [90] Hailuo Song, Ao Li, Tong Wang, and Minghui Wang. Multimodal deep reinforcement learning with auxiliary task for obstacle avoidance of indoor mobile robot. *J. Sens.*, 21(4):1363, 2021. 3
- [91] Nasim Souly, Concetto Spampinato, and Mubarak Shah. Semi supervised semantic segmentation using generative adversarial network. In *ICCV*, 2017. 2
- [92] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *CVPR*, 2018. 4
- [93] Xiaodong Tan, Mathis Antony, and H Kong. Automated music generation for visual art through emotion. In *ICCC*, 2020. 3
- [94] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *NeurIPS*, 2017. 2
- [95] Zachary Teed and Jia Deng. Droid-SLAM: Deep visual slam for monocular, stereo, and rgb-d cameras. *NeurIPS*, 2021. 3
- [96] Zachary Teed, Lahav Lipson, and Jia Deng. Deep patch visual odometry. *arXiv preprint arXiv:2208.04726*, 2022. 2, 3
- [97] Alexander Vakhtov, Victor Lempitsky, and Yinqiang Zheng. Stereo relative pose from line and point feature triplets. In *ECCV*, 2018. 2
- [98] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 5
- [99] Ke Wang, Siyuan Zhang, Junlan Chen, Fan Ren, and Lei Xiao. A feature-supervised generative adversarial network for environmental monitoring during hazy days. *Sci. Total Environ.*, 748:141445, 2020. 2
- [100] Lirui Wang, Yu Xiang, Wei Yang, Arsalan Mousavian, and Dieter Fox. Goal-auxiliary actor-critic for 6d robotic grasping with point clouds. In *CoRL*, 2022. 3
- [101] Sen Wang, Ronald Clark, Hongkai Wen, and Niki Trigoni. Deepvo: Towards end-to-end visual odometry with deep recurrent convolutional neural networks. In *ICRA*, 2017. 2, 3, 7
- [102] Sen Wang, Ronald Clark, Hongkai Wen, and Niki Trigoni. End-to-end, sequence-to-sequence probabilistic visual odometry through deep neural networks. *IJRR*, 37(4-5):513–542, 2018. 7
- [103] Wenshan Wang, Yaoyu Hu, and Sebastian Scherer. Tartanvo: A generalizable learning-based vo. In *CoRL*, 2021. 1, 2, 3, 4, 5, 6, 7, 15, 16
- [104] Yuxi Wang, Junran Peng, and ZhaoXiang Zhang.

- Uncertainty-aware pseudo label refinery for domain adaptive semantic segmentation. In *ICCV*, 2021. 2
- [105] Yuchao Wang, Haochen Wang, Yujun Shen, Jingjing Fei, Wei Li, Guoqiang Jin, Liwei Wu, Rui Zhao, and Xinyi Le. Semi-supervised semantic segmentation using unreliable pseudo-labels. In *CVPR*, 2022. 2
- [106] Shams Watkins, Ladan Shams, Sachiyo Tanaka, J-D Haynes, and Geraint Rees. Sound alters activity in human v1 in association with illusory visual perception. *Neuroimage*, 31(3):1247–1256, 2006. 2
- [107] Benjamin Wilson, William Qi, Tanmay Agarwal, John Lambert, Jagjeet Singh, Siddhesh Khandelwal, Bowen Pan, Ratnesh Kumar, Andrew Hartnett, Jhony Kaesemodel Pontes, et al. Argoverse 2: Next generation datasets for self-driving perception and forecasting. In *NeurIPS*, 2021. 6
- [108] Zhen Xing, Qi Dai, Han Hu, Jingjing Chen, Zuxuan Wu, and Yu-Gang Jiang. Svformer: Semi-supervised video transformer for action recognition. In *CVPR*, 2023. 2
- [109] Zhen Xing, Hengduo Li, Zuxuan Wu, and Yu-Gang Jiang. Semi-supervised single-view 3d reconstruction via prototype shape priors. In *ECCV*, 2022. 2
- [110] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015. 3
- [111] Lian Xu, Wanli Ouyang, Mohammed Bennamoun, Farid Boussaid, Ferdous Sohel, and Dan Xu. Leveraging auxiliary tasks with affinity learning for weakly supervised semantic segmentation. In *ICCV*, 2021. 3
- [112] Fei Xue, Qiuyuan Wang, Xin Wang, Wei Dong, Junqiu Wang, and Hongbin Zha. Guided feature selection for deep visual odometry. In *ACCV*, 2019. 7
- [113] Fei Xue, Xin Wang, Shunkai Li, Qiuyuan Wang, Junqiu Wang, and Hongbin Zha. Beyond tracking: Selecting memory and refining poses for deep visual odometry. In *CVPR*, pages 8575–8583, 2019. 3, 7
- [114] Pengxiang Yan, Guanbin Li, Yuan Xie, Zhen Li, Chuan Wang, Tianshui Chen, and Liang Lin. Semi-supervised video salient object detection using pseudo-labels. In *ICCV*, 2019. 2, 3
- [115] Jihan Yang, Shaoshuai Shi, Zhe Wang, Hongsheng Li, and Xiaojuan Qi. ST3D: Self-training for unsupervised domain adaptation on 3d object detection. In *CVPR*, 2021. 4
- [116] Nan Yang, Lukas von Stumberg, Rui Wang, and Daniel Cremers. D3vo: Deep depth, deep pose and deep uncertainty for monocular visual odometry. In *CVPR*, 2020. 1, 2, 3
- [117] Xiangli Yang, Zixing Song, Irwin King, and Zenglin Xu. A survey on deep semi-supervised learning. *arXiv preprint arXiv:2103.00550*, 2021. 2
- [118] Yingda Yin, Yingcheng Cai, He Wang, and Baoquan Chen. Fishermatch: Semi-supervised rotation regression via entropy-based filtering. In *CVPR*, 2022. 2
- [119] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *CVPR*, 2018. 1, 2, 3, 6, 7
- [120] Amir R Zamir, Alexander Sax, Nikhil Cheerla, Rohan Suri, Zhangjie Cao, Jitendra Malik, and Leonidas J Guibas. Robust learning through cross-task consistency. In *CVPR*, 2020. 3, 4, 15
- [121] Amir R Zamir, Alexander Sax, William Shen, Leonidas J Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *CVPR*, 2018. 3, 4, 15
- [122] Huangying Zhan, Ravi Garg, Chamara Saroj Weerasekera, Kejie Li, Harsh Agarwal, and Ian Reid. Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction. In *CVPR*, 2018. 2, 7
- [123] Huangying Zhan, Chamara Saroj Weerasekera, Jia-Wang Bian, and Ian Reid. Visual odometry revisited: What should be learnt? In *ICRA*, 2020. 1, 2
- [124] Jimuyang Zhang, Zanming Huang, and Eshed Ohn-Bar. Coaching a teachable student. In *CVPR*, 2023. 3
- [125] Jimuyang Zhang and Eshed Ohn-Bar. Learning by watching. In *CVPR*, 2021. 3
- [126] Jimuyang Zhang, Minglan Zheng, Matthew Boyd, and Eshed Ohn-Bar. X-World: Accessibility, vision, and autonomy meet. In *ICCV*, 2021. 2
- [127] Jimuyang Zhang, Ruizhao Zhu, and Eshed Ohn-Bar. Selfd: self-learning large-scale driving policies from the web. In *CVPR*, 2022. 2, 3
- [128] Yu Zhang and Qiang Yang. A survey on multi-task learning. *TKDE*, 34(12):5586–5609, 2021. 4
- [129] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. Facial landmark detection by deep multi-task learning. In *ECCV*, 2014. 3
- [130] Shengyu Zhao, Yilun Sheng, Yue Dong, Eric I Chang, Yan Xu, et al. Maskflownet: Asymmetric feature matching with learnable occlusion mask. In *CVPR*, 2020. 4, 13
- [131] Wang Zhao, Shaohui Liu, Yezhi Shu, and Yong-Jin Liu. Towards better generalization: Joint depth-pose learning without posenet. In *CVPR*, 2020. 1, 2, 3
- [132] Liangli Zhen, Peng Hu, Xu Wang, and Dezhong Peng. Deep supervised cross-modal retrieval. In *CVPR*, 2019. 3
- [133] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, 2017. 1, 2, 7
- [134] Yipin Zhou, Zhaowen Wang, Chen Fang, Trung Bui, and Tamara L Berg. Visual to sound: Generating natural sound for videos in the wild. In *CVPR*, 2018. 3
- [135] Yuliang Zou, Zelun Luo, and Jia-Bin Huang. Df-net: Unsupervised joint learning of depth and flow using cross-task consistency. In *ECCV*, 2018. 1, 2

# Supplementary Material for XVO: Generalized Visual Odometry via Cross-Modal Self-Training

## Abstract

*In this supplementary document, we first discuss additional implementation details (Sec. 1), including our network architecture, training settings, and employed evaluation metrics. Next, we provide ablative analysis across different baselines, datasets, and auxiliary tasks (Sec. 2). Finally, we include additional qualitative examples (Sec. 3, also shown in our supplementary video).*

## 1. Implementation Details

### 1.1. Network Architecture

As mentioned in the main paper, we learn a unified feature extractor based on a MaskFlowNet encoder [130] that is followed by four self-attention layers [14, 22]. The encoder structure for the initial teacher model and the subsequent cross-modal student model is kept the same, and provides state-of-the-art results on KITTI. Task decoders are then added over the shared visual encoder and supervised by auxiliary audio, visual, and motion prediction tasks in order to encourage the model to extract information that is relevant across a wide range of conditions, including potentially unseen environments.

**Task Decoders:** Due to their similar output representations, we incorporate segmentation, depth, flow auxiliary tasks using a similar decoder structure (also outlined in the main paper). The goal of these three auxiliary tasks is to encourage the visual encoder to identify geometric and motion cues that can also be relevant for the primary VO task. The segmentation decoder employs an FCN [55] decoder, consisting of 11 transposed convolutional layers followed by a convolutional layer and a final sigmoid activation function. The output size is  $289 \times 296$ . The optical flow decoder shares the same structure but without a sigmoid and pretrained on the Flying Chairs [24] dataset. The depth decoder was pretrained using NYU Depth V2 [69].

### 1.2. Data

**Augmentation Strategies:** To generalize model performance, we use several data augmentation techniques including random cropping and resizing which simulates varying camera intrinsics.

**YouTube Videos:** YouTube has ample navigation videos, often with corresponding audio. We searched for high-resolution driving dash camera videos and downloaded diverse data with different times of day (both daylight and nighttime), weathers (rain, snow, sun), environments (urban, rural, suburban), and location (Boston, Washington, London, Singapore, Paris, Switzerland, Milan, Ireland, Tokyo). The original video data is about 100 minutes at 30FPS, which we then subsample by three to 10FPS to obtain the final dataset.

### 1.3. Training Details

All models are trained for 15 epochs using Stochastic Gradient Descent with a batch size of six and input image size of  $640 \times 384$  [83, 87]. Training takes approximately two days. We perform extensive ablation on the role of the various auxiliary tasks for model training in Sec. 2. We empirically set the entropy threshold in Eqn. 5 of the main paper as  $\tau_u = -5.668$ . As we do not find it beneficial to iterate over semi-supervised training (i.e., with the new student as the teacher), we only perform one iteration of teacher-student training.

### 1.4. Evaluation Metrics

We sought to accurately measure the ability of our model to predict general real-world pose, including absolute scale, throughout a trajectory. In this section, we define our metrics, including a proposed scale error metric which can more accurately account for scale errors in prediction along a trajectory. While our best performing model is shown to outperform other models in terms of all three defined metrics, we demonstrate the semi-supervised learning process to consistently improve scale estimation as well (discussed further in Sec. 2).

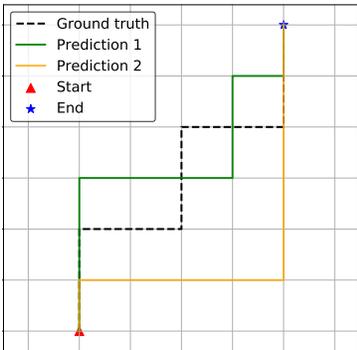


Figure 5: **Low Translation and Rotation Error Despite Inaccurate Predictions.** Given our emphasis on estimating real-world pose, including absolute scale, across diverse platforms, we highlight issues with the definition of standard translation and rotation errors. As standard VO metrics are computed over trajectory end-points, the trajectory itself can vary arbitrarily between the end-points while maintaining low error. In the example shown, all predicted trajectories will have zero translation and rotation error. However, using an average scale error will demonstrate a difference. The scale error for Prediction 1 and Prediction 2 is 0.33 and 0.30 respectively. Therefore, introducing the scale error provides a more holistic evaluation of model performance while directly measuring scale.

**Translation and Rotation Errors at Trajectory End-Points:** Standard visual odometry metrics average relative pose errors at fixed distances along the trajectory [30, 44, 78]. We follow common metrics to compute translation and rotation errors which are computed over all possible trajectory subsequences of length  $(100, \dots, 800)$  meters. Given a trajectory of length  $l$ , standard error metrics measure relative pose between an estimated pose and the ground-truth pose of the last frame with respect to the first frame,  $[\mathbf{R}'|\mathbf{t}']$ . Translation and rotation errors are computed as averaged translational (%) drift and averaged rotational drift ( $^\circ/100m$ ):

$$t_{rel} = \|\mathbf{t}'\|_2 * \frac{100}{l} \quad (11)$$

$$r_{rel} = \arccos(\max(\min(d, 1), -1)) * \frac{180}{\pi} * \frac{100}{l} \quad (12)$$

where  $d = \frac{\text{tr}(\mathbf{R}') - 1}{2}$ .

**Issues with Standard Metrics:** Relative translation and rotation errors can measure an error drift over the trajectory subsequence. However, these measures can neglect instantaneous errors along the trajectory subsequence. Here, multiple differing trajectories can all have similar rotation and translation errors, given similar translation and rotation at their end-points. Critically, the *translation error can be low even if the estimated trajectory itself is far from the ground-truth trajectory* prior to the end-point used for evaluation. To emphasize, as long as the estimated translation and rotation at the end of the trajectory is same with ground-truth, translation error and rotation error can be low despite potentially high two-frame inaccuracies.

We present a toy example in Fig. 5 and a trajectory of six consecutive time steps with ground-truth relative rotations of  $[0^\circ, 90^\circ, 0^\circ]$ ,  $[0^\circ, -90^\circ, 0^\circ]$ ,  $[0^\circ, 90^\circ, 0^\circ]$ ,  $[0^\circ, -90^\circ, 0^\circ]$ ,  $[0^\circ, 0^\circ, 0^\circ]$ , and ground-truth relative translations is kept fixed at each time step  $[0, 0, 20]$ . In Prediction 1, assume our predicted relative rotations are similar to the ground-truth and our predicted relative translation are  $[0, 0, 30], [0, 0, 30], [0, 0, 20], [0, 0, 10], [0, 0, 10]$ , respectively. In Prediction 2, assume our predicted relative rotations are  $[0^\circ, 90^\circ, 0^\circ]$ ,  $[0^\circ, -90^\circ, 0^\circ]$ ,  $[0^\circ, 0^\circ, 0^\circ]$ ,  $[0^\circ, 0^\circ, 0^\circ]$ ,  $[0^\circ, 0^\circ, 0^\circ]$  and our predicted relative translations are  $[0, 0, 10], [0, 0, 40], [0, 0, 10], [0, 0, 20], [0, 0, 20]$ , respectively. Once we recover the trajectories based on the ground-truth relative pose and predicted relative pose, all trajectories share the same endpoint with same direction. Here, while the model fails at two-frame prediction, the translation error and rotation error are both zero. Although this example may seem contrived, we sought a more accurate evaluation of our model which directly estimates real-world scale using two-frame input, as discussed next.

**Scale Error:** Towards a more accurate evaluation of real-world scale along the trajectory, we also compute and average the two-frame scale error (defined in the Eqn. 10 in the main paper). We therefore also report the *average scale error (se)*, which

is an error computed over predicted and ground-truth translation,

$$se = 1 - \min \left( \frac{\|\hat{\mathbf{t}}\|_2}{\max(\|\mathbf{t}\|_2, \epsilon)}, \frac{\|\mathbf{t}\|_2}{\max(\|\hat{\mathbf{t}}\|_2, \epsilon)} \right). \quad (13)$$

Here, the scale error is always computed over two sequential frames, as opposed to over end-points of a multi-step trajectory, and therefore provides a more holistic evaluation of model performance while directly measuring scale. The scale error of Prediction 1 and Prediction 2 in Fig. 5 are then 0.33 and 0.30, respectively.

## 2. Ablation Studies

This section provides comprehensive analysis with various modalities and training settings using our cross-dataset evaluation. In Sec. 2.1 we discuss comparison with our baseline, TartanVO [103]. Next, we explore the impact of various auxiliary prediction tasks in Sec. 2.2 as well as the uncertainty-based filtering process for the pseudo-labels in Sec. 2.3 and the impact of additional semi-supervised learning iterations (Sec. 2.4).

### 2.1. TartanVO Baseline

Our main comparative baseline is TartanVO [103] as it is one of the few VO methods that can perform cross-dataset generalization without fine-tuning. Moreover, TartanVO does not leverage long-term memory to correct for drift (i.e., the input is two-frame, as with our model). However, TartanVO *leverages the intrinsic parameters* to adapt across camera settings as input to the network, while our method does not. Moreover, in evaluation, the approach aligns the predicted pose to the ground-truth, i.e., by only estimating pose up to a scale factor which is determined using the ground-truth in evaluation. In contrast, our study emphasizes accurate real-world scale prediction. Since up-to-scale prediction can limit the applicability of the method in practice, we analyze performance with and without the ground-truth scale alignment step in Table 3. As shown in the table, the removal of the ground-truth alignment step allows us to better evaluate and compare with our models. On average across the dataset, our top performing models (e.g., Student+Audio) outperform TartanVO both with and without ground-truth alignment by (e.g., 10.54 vs. 14.84 and 3.11 vs. 5.31 for the translation and rotation errors, respectively). We note that we cannot compute a scale error for methods that estimate pose up to a scale.

### 2.2. Modality Ablation

Table 3 shows the result of training a teacher model on nuScenes and evaluating on the three benchmarks, including unseen portions of nuScenes. We then train various student models with different auxiliary tasks. For legibility, we also color the top three performing methods within each dataset and metric, as well as summarized performance across all settings (however, while KITTI only has 11 sequences, nuScenes has 765 testing sequences). We also note that our results are with *pseudo labels of depth and motion cues and not ground-truth* which is often employed by related studies in cross-task [120, 121]. In our semi-supervised settings, the pseudo-labels can still contain noise despite our filtering mechanism.

**Results on KITTI:** On KITTI, we find adding depth and audio-based models (both audio, audio+segmentation, and audio+segmentation+flow) to perform best. This finding is consistent both with the overall averages and the per-sequence results for KITTI (shown in Table 4). As depth provides a general cue regarding scene geometry, we find it to benefit pose estimation performance on KITTI (translation error with the student model drops from 17.04 to 13.09), but to a lesser extent on the other datasets where the audio-based approaches perform best. This can be explained by inspecting the KITTI scenes, which contain simpler layouts with few dynamic objects compared to Argoverse and nuScenes (where flow and segmentation-based prediction is also shown to be helpful). In such scenarios, depth cues can be more complex in contrast with the others, e.g., audio or motion cues.

**Results on All Benchmarks:** Although we do find fluctuations in the benefit of auxiliary tasks within each dataset, the audio-based model is shown to perform consistently best across all datasets (10.54 translation, 3.11 rotation, and 0.13 scale errors). While other combinations of auxiliary tasks also perform well, we do not find them to benefit generalized representation learning as much as the audio task itself. Overall, the results affirm our hypothesis that the audio task can act as a useful regularizer for the VO task. While the audio can have significant ambient noise, the pseudo-labels with the depth, flow, and segmentation tasks can also contain noise. Moreover, many approaches for such tasks are often developed and tuned for KITTI, potentially incorporating various implicit biases towards KITTI-specific settings. Nonetheless, the addition of such tasks can benefit the audio-only model at times, e.g., on KITTI with the audio+segmentation model (14.20 vs. 14.37

Table 3: **Model Ablation.** We train models on nuScenes and test our approach using various model settings and datasets (KITTI sequences 00-10, Argoverse 2, and unseen areas in nuScenes). We also show overall average results. TartanVO [103], our main baseline, leverages the intrinsics and does not predict absolute scale directly, but instead evaluates with scale alignment using the ground-truth scale at each time step. Results are shown for translation, rotation, and scale errors. ‘w/o Filter’ refers to removal of the proposed uncertainty-aware pseudo-label sample selection mechanism. Lowest three errors are highlighted within each column (when two numbers are identical to two decimal places, we compare their original precision).

Method	KITTI 00-10			Argoverse 2			nuScenes			Average		
	$t_{err}$	$r_{err}$	$se$	$t_{err}$	$r_{err}$	$se$	$t_{err}$	$r_{err}$	$se$	$t_{err}$	$r_{err}$	$se$
TartanVO w/ GT Alignment [103]	6.37	3.32	/	8.55	5.77	/	9.61	6.83	/	8.17	5.31	/
TartanVO w/o GT Alignment	21.67	3.33	0.29	41.11	5.77	0.40	28.23	6.83	0.29	30.34	5.31	0.33
Teacher (nuScenes)	26.16	6.84	0.25	10.89	3.40	0.16	15.93	6.73	0.20	17.66	5.66	0.20
Student w/o Filter	20.64	5.68	0.21	10.80	7.33	0.14	9.32	4.60	0.14	13.59	5.87	0.16
Student	17.04	4.02	0.16	9.16	3.40	0.14	10.54	3.94	0.13	12.24	3.79	0.14
Student+Seg w/o Filter	16.65	3.42	0.20	9.83	4.63	0.13	8.33	3.75	0.12	11.60	3.93	0.15
Student+Seg	16.31	3.77	0.16	9.17	3.18	0.13	11.35	4.05	0.14	12.28	3.67	0.14
Student+Flow w/o Filter	20.47	5.65	0.22	10.74	5.95	0.13	10.02	4.87	0.13	13.74	5.49	0.16
Student+Flow	15.60	3.19	0.19	9.04	4.45	0.13	9.13	4.06	0.13	11.26	3.90	0.15
Student+Depth w/o Filter	19.15	4.71	0.18	8.64	3.70	0.13	11.48	3.80	0.14	13.09	4.07	0.15
Student+Depth	13.09	3.03	0.16	9.91	3.08	0.14	10.86	5.55	0.14	11.29	3.89	0.15
Student+Audio w/o Filter	16.45	3.88	0.19	9.83	5.38	0.14	9.02	4.03	0.14	11.77	4.43	0.16
Student+Audio	14.37	3.06	0.16	8.00	3.08	0.12	9.26	3.20	0.12	10.54	3.11	0.13
Student+Audio+Seg	14.20	3.02	0.16	8.67	3.63	0.13	11.29	3.70	0.14	11.39	3.45	0.14
Student+Audio+Depth	15.45	3.86	0.18	8.43	4.51	0.13	10.44	4.94	0.14	11.44	4.44	0.15
Student+Seg+Depth	15.07	3.98	0.15	10.03	3.10	0.15	9.58	3.40	0.13	11.56	3.49	0.14
Student+Audio+Flow	15.84	3.09	0.19	7.76	3.59	0.12	8.44	3.45	0.12	10.68	3.38	0.14
Student+Flow+Depth	17.08	3.07	0.20	7.84	3.85	0.12	10.01	3.33	0.14	11.64	3.42	0.15
Student+Audio+Seg+Flow	16.95	2.98	0.19	8.32	3.59	0.12	8.20	3.35	0.12	11.16	3.31	0.26
Student+Audio+Flow+Depth	18.19	3.70	0.21	7.95	4.43	0.13	9.80	3.71	0.14	11.88	3.95	0.16
Student+Audio+Seg+Depth	17.51	3.49	0.20	9.33	4.26	0.13	9.07	3.61	0.13	11.97	3.79	0.15
Student+Seg+Flow+Depth	18.23	3.88	0.21	8.79	4.89	0.13	8.93	3.44	0.13	11.98	4.07	0.16
Student+A+S+F+D	16.74	4.40	0.18	7.89	3.54	0.12	9.98	4.36	0.15	11.53	4.10	0.15

Table 4: **Results for KITTI Sequences.** Model results for each of the KITTI sequences (used to compute averages in Table 3). We train various models using nuScenes, and evaluate without fine-tuning on KITTI. ‘w/o Filter’ refers to removal of the proposed uncertainty-aware pseudo-label sample selection mechanism. Lowest three errors are highlighted within each column.

Method	Seq 00			Seq 01			Seq 02			Seq 03			Seq 04			Seq 05			Seq 06			Seq 07			Seq 08			Seq 09			Seq 10		
	$t_{err}$	$r_{err}$	$se$																														
Our Teacher (nuScenes)	24.62	7.51	0.24	40.02	2.44	0.41	25.19	5.14	0.25	26.78	4.92	0.25	26.02	2.42	0.26	23.65	8.85	0.22	23.97	6.47	0.24	30.66	20.32	0.26	33.03	6.69	0.20	23.23	4.44	0.22	20.57	6.01	0.18
Student w/o Filter	18.16	6.45	0.19	41.24	1.84	0.41	19.05	4.48	0.22	26.98	9.68	0.23	22.56	2.15	0.22	14.77	5.83	0.17	11.38	1.62	0.18	16.45	9.35	0.20	19.36	7.95	0.18	16.82	4.14	0.20	20.23	8.99	0.15
Student	14.71	5.42	0.12	39.79	2.56	0.38	12.54	2.82	0.14	20.30	3.97	0.18	16.33	1.57	0.16	11.12	4.19	0.11	15.60	5.69	0.11	7.77	3.48	0.18	13.15	4.41	0.13	16.17	4.50	0.12	19.91	5.59	0.14
Student+Seg w/o Filter	12.41	3.32	0.17	41.20	3.28	0.39	15.70	2.56	0.19	22.16	3.09	0.22	22.39	3.23	0.22	10.92	3.51	0.15	10.41	1.86	0.15	12.88	6.57	0.19	10.53	3.14	0.15	14.14	2.88	0.16	10.36	4.21	0.15
Student+Seg	13.88	5.09	0.12	40.73	1.94	0.40	13.22	3.11	0.14	17.18	5.14	0.15	18.40	1.74	0.18	11.71	4.27	0.12	15.26	5.03	0.12	5.97	3.52	0.17	11.79	3.88	0.12	15.85	3.77	0.14	15.40	4.02	0.15
Student+Flow w/o Filter	17.92	6.19	0.18	42.25	1.45	0.43	20.56	5.30	0.22	21.55	4.34	0.22	24.03	4.52	0.23	17.00	6.82	0.19	13.26	3.19	0.18	18.68	10.79	0.20	17.07	7.33	0.17	21.11	6.03	0.21	11.73	6.19	0.17
Student+Flow	9.52	2.68	0.15	38.94	2.40	0.40	13.02	2.82	0.16	23.31	3.09	0.23	19.40	1.17	0.19	8.41	1.94	0.16	19.60	6.61	0.13	9.24	4.68	0.21	8.85	2.81	0.15	12.63	2.87	0.16	9.72	4.00	0.15
Student+Depth w/o Filter	20.11	7.56	0.14	40.46	2.72	0.40	15.38	3.32	0.17	18.92	4.83	0.16	21.33	2.91	0.21	12.19	5.05	0.13	22.01	7.49	0.14	13.77	6.38	0.18	13.82	4.36	0.14	15.84	2.97	0.17	16.86	4.26	0.16
Student+Depth	7.57	2.13	0.13	38.66	2.96	0.37	11.38	2.75	0.13	16.12	3.29	0.17	11.71	0.91	0.12	9.96	3.00	0.12	13.58	5.51	0.11	8.05	4.48	0.18	8.58	2.19	0.12	10.58	2.84	0.13	7.84	3.27	0.15
Student+Audio w/o Filter	13.01	4.00	0.17	39.97	3.58	0.38	15.71	3.26	0.18	21.48	4.80	0.21	21.15	2.05	0.21	10.84	3.83	0.14	10.09	1.50	0.16	12.69	6.76	0.17	12.37	4.39	0.15	13.76	3.16	0.16	9.89	5.37	0.14
Student+Audio	9.09	2.01	0.12	37.46	3.32	0.35	13.87	2.46	0.16	16.06	3.97	0.15	19.27	1.78	0.19	9.77	2.98	0.13	10.91	3.16	0.12	11.26	6.61	0.17	7.91	1.84	0.12	13.42	2.77	0.14	9.01	2.83	0.14
Student+Audio+Seg	10.49	3.71	0.12	42.46	2.73	0.40	11.61	2.38	0.14	14.53	3.93	0.14	16.29	0.96	0.17	8.31	2.76	0.14	15.31	5.49	0.11	5.86	3.00	0.17	7.05	2.03	0.12	11.64	2.78	0.13	12.17	3.45	0.14
Student+Audio+Depth	12.99	4.71	0.16	37.79	1.63	0.39	15.47	3.94	0.16	17.06	3.37	0.18	13.42	2.63	0.14	9.53	2.77	0.14	20.02	7.20	0.13	8.95	3.70	0.19	10.18	3.47	0.15	15.12	4.57	0.16	9.40	4.42	0.14
Student+Seg+Depth	14.73	5.22	0.12	39.28	2.39	0.37	13.72	3.29	0.11	17.98	5.28	0.13	11.38	1.91	0.11	14.84	5.76	0.11	8.47	2.05	0.10	10.68	7.20	0.19	12.25	3.96	0.11	10.71	3.12	0.10	11.78	3.58	0.16
Student+Audio+Flow	11.55	2.23	0.18	36.98	2.89	0.37	15.22	3.09	0.18	20.41	3.35	0.20	19.83	2.33	0.20	10.85	2.63	0.16	12.24	2.96	0.16	11.43	4.56	0.22	8.83	2.16	0.15	15.12	3.74	0.15	11.78	4.08	0.15
Student+Flow+Depth	11.90	2.15	0.18	41.76	2.44	0.41	15.51	2.56	0.18	22.31	3.81	0.22	20.00	1.95	0.20	12.58	3.16	0.16	12.15	2.85	0.16	14.35	6.93	0.20	11.21	2.32	0.15	15.95	2.56	0.18	10.18	3.06	0.15
Student+Audio+Seg+Flow	11.51	3.19	0.15	41.72	3.31	0.42	15.72	2.70	0.19	22.06	3.58	0.22	20.36	1.76	0.20	10.28	2.78	0.13	13.98	3.51	0.16	8.84	3.44	0.18	10.47	2.33	0.14	16.64	2.87	0.18	14.85	3.33	0.15
Student+Audio+Flow+Depth	13.02	2.69	0.20	42.33	2.15	0.42	15.99	2.82	0.19	22.36	4.39	0.23	21.46	2.27	0.21	12.64	3.24	0.18	12.83	2.80	0.16	15.04	8.07	0.21	12.14	3.19	0.17	18.42	3.81	0.20	13.89	5.26	0.15
Student+Seg+Depth	13.27	2.78	0.18	42.21	2.22	0.41	15.79	2.81	0.19	20.24	3.93	0.19	24.15	3.42	0.24	12.75	3.61	0.16	12.86	2.78	0.15	14.22	7.30	0.22	10.68	2.30	0.15	16.67	3.61	0.19	9.78	3.61	0.16
Student+Seg+Flow+Depth	13.12	2.58	0.19	41.47	1.87	0.41	17.48	3.06	0.21	20.40	3.83	0.20	21.89	3.73	0.21	13.61	3.87	0.18	13.43	2.54	0.17	15.98	8.15	0.21	12.49	3.37	0.17	18.39	3.90	0.21	12.28	5.76	0.17
Student+A+S+F+D	14.98	4.59	0.16	37.72	2.28	0.37	16.23	4.57	0.14	16.76	5.43	0.14	12.10	3.17	0.13	16.70	5.28	0.15	10.80	3.59	0.17	15.69	8.95	0.24	13.98	4.07	0.14	12.52	3.93	0.13	16.66	4.55	0.17

translation error and 3.02 vs. 3.06 rotation error). As combining multiple modalities can potentially interfere with the main VO task as well as introduce noise, effective combinations with multiple auxiliary tasks require further study in the future.

**Scale Error Discussion:** Based on our analysis in Table 3, we do show the scale error to differ from the other metrics. For instance, the segmentation+depth model has the lowest scale error on KITTI but not translation nor rotation errors. Similarly, the audio+segmentation+flow model achieves low scale error on Argoverse. Moreover, while the translation error between the depth and audio-based models on KITTI varies, their scale error is similar (audio slightly lower). As the audio-based model outperforms the depth in terms of overall performance, we can see how the *scale error provides a better predictor* in this case for overall performance. The scale error is shown to consistently select the audio-based student model on the other datasets, while translational error exhibits more variability. For instance, on Argoverse, the translation errors of audio+flow,

flow+depth, and all tasks are lower than the audio model, yet the scale error selects the audio model.

### 2.3. Uncertainty-Aware Pseudo-Label Filter

Table 3 also analyzes the role of the uncertainty-based removal of pseudo-labeled samples prior to training the student model. We find our sample removal mechanism to generally benefit the student models regardless of the task. The impact is more pronounced for the student, flow, and depth models. The segmentation model shows a slight increase in average overall translation error (from 11.60 to 12.28), but reduction in rotation and scale errors. Moreover, the scale error for the audio-based model significantly drops after filtering out potentially noisy pseudo-labels, from 0.16 to 0.13. These findings affirm the benefits of our proposed approach as well as evaluation with the scale error.

### 2.4. Additional Pseudo-Labeling Iterations

We investigated the benefit of additional iterations of pseudo-labeling, where the trained student model is used to re-label the YouTube videos with potentially higher quality pseudo-labels. However, we did not find additional pseudo-labeling iterations and re-training to benefit model performance. For instance, in the case of the strongest audio-based model, the scale error increases slightly (from 0.13 and 0.14) following another iteration. The translation and rotation errors also slightly increase by 3% and 2%, respectively, and plateau with subsequent iterations.

## 3. Additional Qualitative Results

To further motivate our use of the audio modality, Fig. 6 depicts additional examples from our dataset with overlaid audio. In addition to ego-motion speed (e.g., stopped or driving), the audio cues also contain general context regarding the traffic scenario (e.g., highway, urban). We also visualize additional success and failure examples for Argoverse (Fig. 7), nuScenes (Fig. 8), and KITTI (Fig. 9).

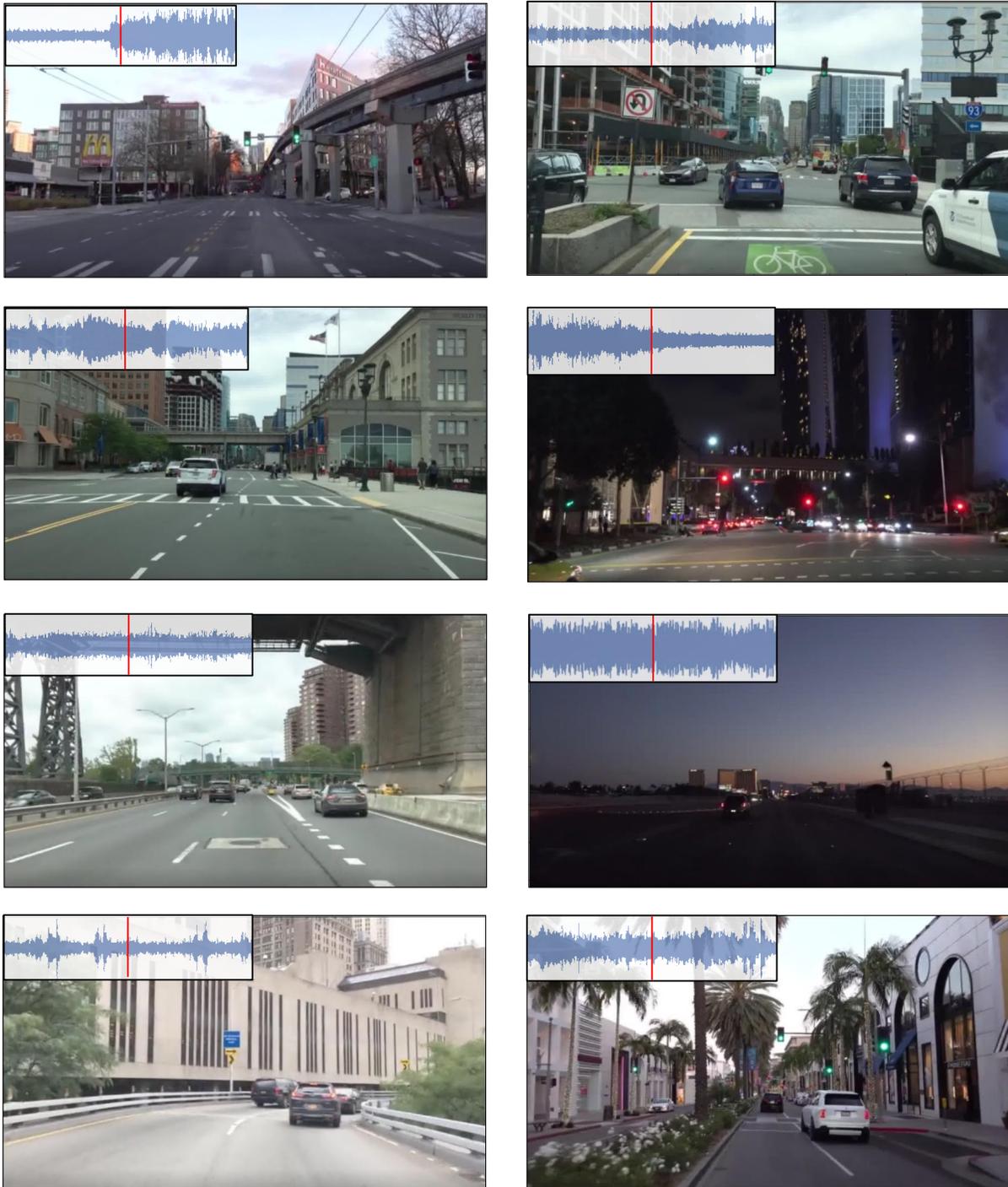


Figure 6: **Role of Audio.** First row: The audio amplitude is shown to increase when the vehicle speeds up from a stopped state; Second row: The audio amplitude decreases when the vehicle begins to slow down; Third row: High and stable audio amplitude on a freeway; Row 4: Various urban events, e.g., stop-and-go traffic, presence of surrounding vehicles, or slowing down to turn.

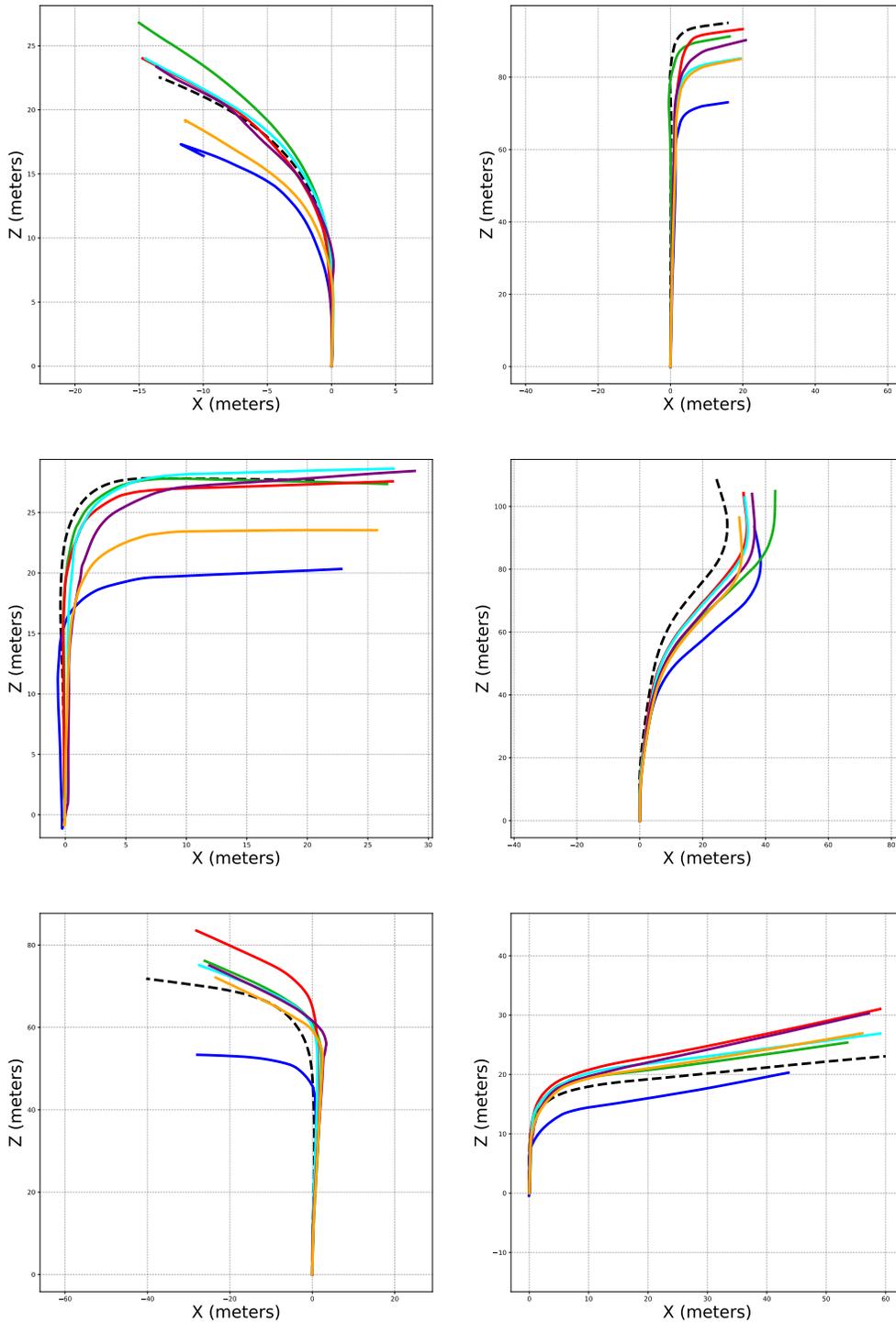


Figure 7: **Recovered Trajectories on Argoverse.** First two rows depict **success cases** where the proposed approach is shown to improve the predicted trajectory. Third row (left) depicts a **failure case** due to incorrect estimation of rotation during a turn.

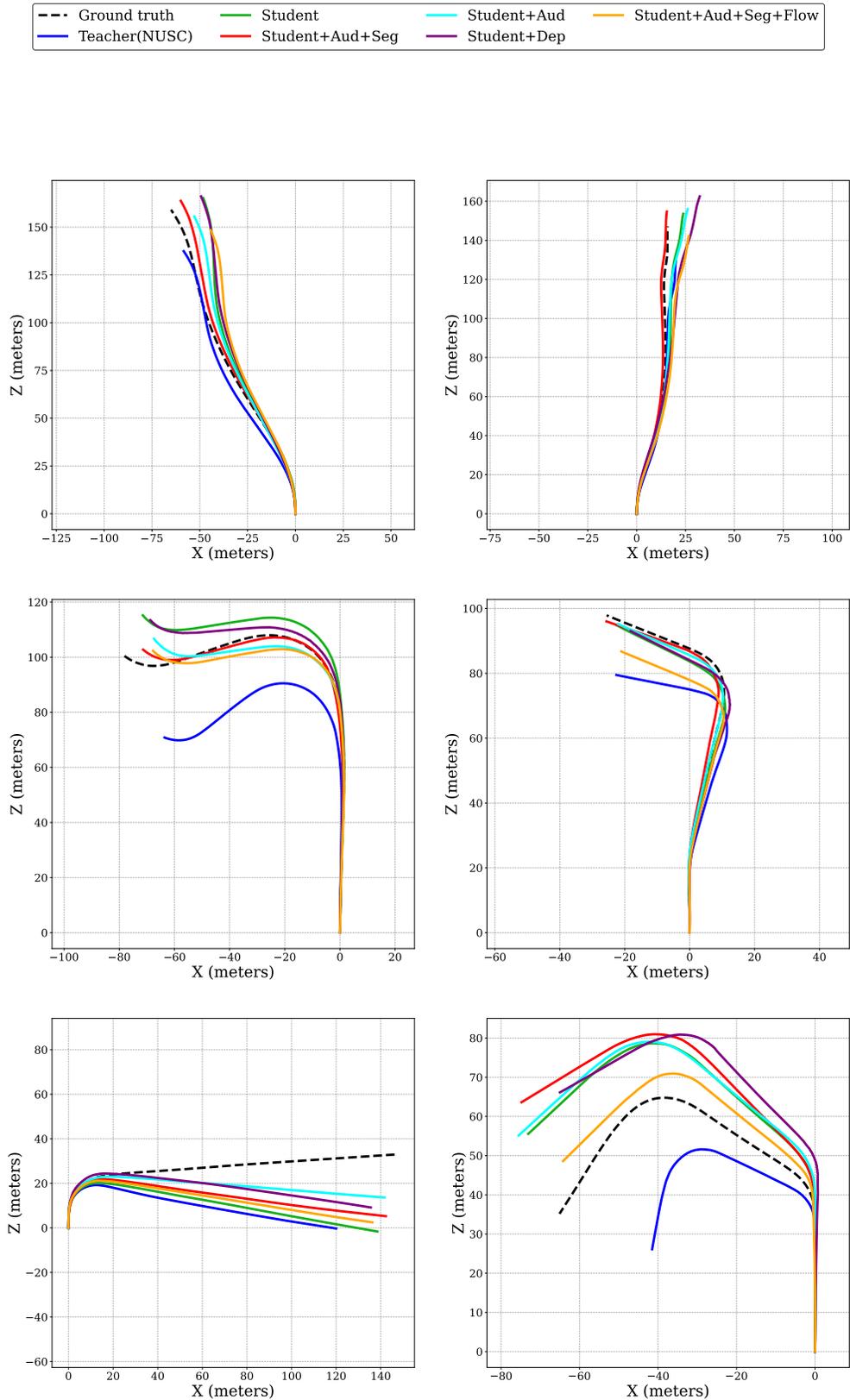


Figure 8: **Recovered Trajectories on nuScenes.** First two rows depict **success cases** where the proposed approach is shown to improve the predicted trajectory. The third row (left) depicts a **failure case**, due to a challenging turn in a busy intersection where most approaches under-shoot estimated pose.

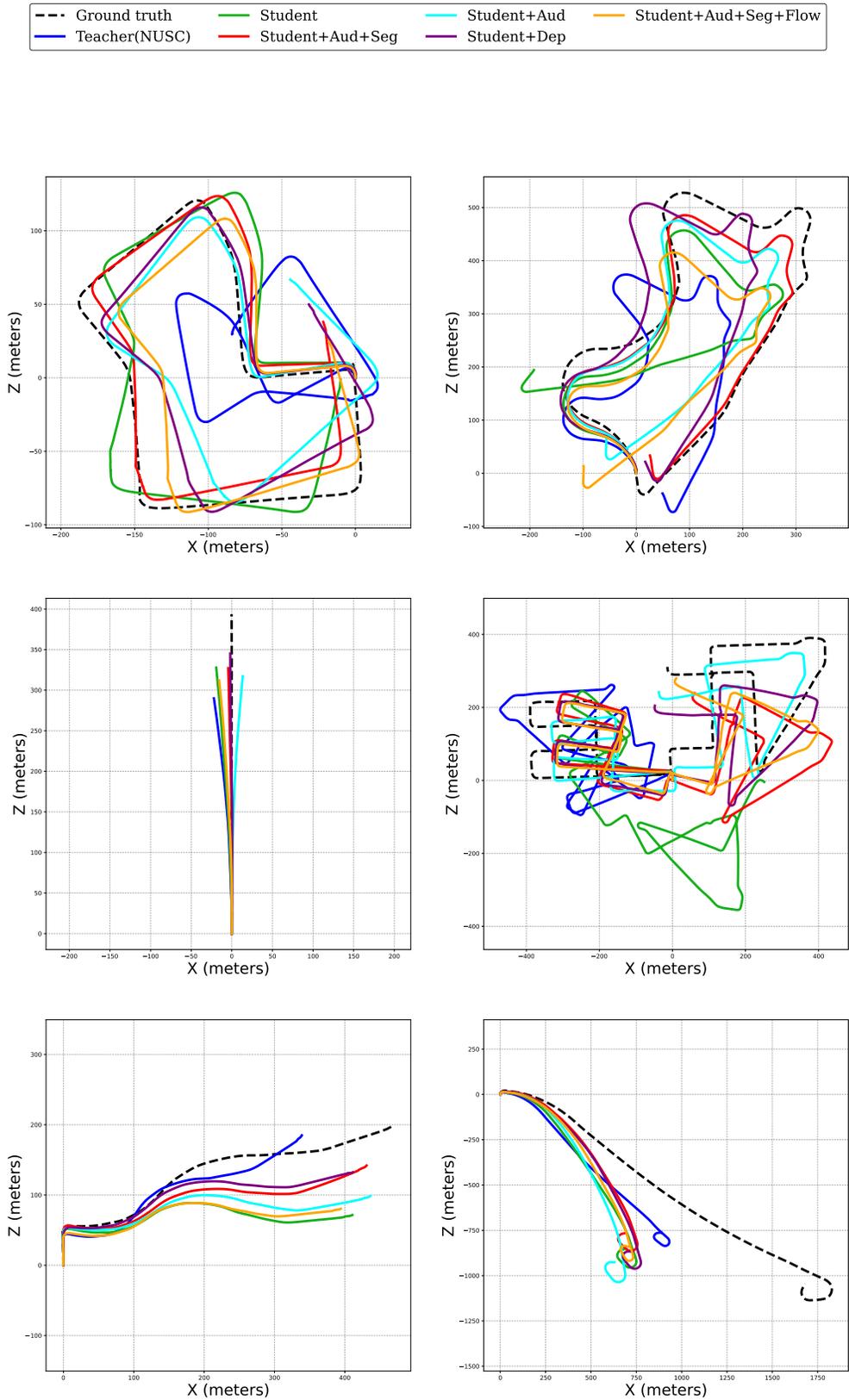


Figure 9: **Recovered Trajectories on KITTI**. First two rows depict **success cases** where the proposed approach is shown to improve the predicted trajectory. The third row depicts two **failure cases** where most of the approaches fail, either due to a lengthy forward motion on a stretch minimal clear visual landmarks or along a curved route.