

# Video K-Net: A Simple, Strong, and Unified Baseline for Video Segmentation

Xiangtai Li<sup>1,2,4\*</sup> Wenwei Zhang<sup>2\*</sup> Jiangmiao Pang<sup>3,5\*</sup> Kai Chen<sup>4,5</sup>  
 Guangliang Cheng<sup>4</sup> Yunhai Tong<sup>1</sup> Chen Change Loy<sup>2</sup>

<sup>1</sup> Key Laboratory of Machine Perception, MOE, School of Artificial Intelligence, Peking University

<sup>2</sup> S-Lab, Nanyang Technological University <sup>3</sup> CUHK-SenseTime Joint Lab, the Chinese University of Hong Kong

<sup>4</sup> SenseTime Research <sup>5</sup> Shanghai AI Lab

{lxtpku, yhtong}@pku.edu.cn {wenwei001, ccloy}@ntu.edu.sg  
 pangjiangmiao@gmail.com {chenkai, chengguangliang}@sensetime.com

## Abstract

*This paper presents Video K-Net, a simple, strong, and unified framework for fully end-to-end video panoptic segmentation. The method is built upon K-Net, a method that unifies image segmentation via a group of learnable kernels. We observe that these learnable kernels from K-Net, which encode object appearances and contexts, can naturally associate identical instances across video frames. Motivated by this observation, Video K-Net learns to simultaneously segment and track “things” and “stuff” in a video with simple kernel-based appearance modeling and cross-temporal kernel interaction. Despite the simplicity, it achieves state-of-the-art video panoptic segmentation results on Cityscapes-VPS, KITTI-STEP, and VIPSeg without bells and whistles. In particular, on KITTI-STEP, the simple method can boost almost 12% relative improvements over previous methods. On VIPSeg, Video K-Net boosts almost 15% relative improvements and results in 39.8% VPQ. We also validate its generalization on video semantic segmentation, where we boost various baselines by 2% on the VSPW dataset. Moreover, we extend K-Net into clip-level video framework for video instance segmentation, where we obtain 40.5% for ResNet50 backbone and 54.1% mAP for Swin-base on YouTube-2019 validation set. We hope this simple, yet effective method can serve as a new, flexible baseline in unified video segmentation design. Both code and models are released at <https://github.com/lxtGH/Video-K-Net>.*

## 1. Introduction

Video Panoptic Segmentation (VPS) aims at segmenting and tracking every pixel of input video clips [20, 22, 62]. As

\*Equal contribution. Most Work Done when Xiangtai was in SenseTime Research (Beijing).

a fundamental technique to scene understanding, it has received increasing attention in recent years due to its wide applications in many vision systems, including autonomous driving and robot navigation [12, 16]. By definition, VPS is an extension of Panoptic Segmentation (PS) [25] into the video domain with the goal of unifying Video Semantic Segmentation (VSS) [38, 50, 82] and Video Instance Segmentation (VIS) [3, 69] into a single task.

Existing studies for video panoptic segmentation can be mainly divided into *top-down* methods [20, 22] and *bottom-up* approaches [47, 62]. Top-down methods try to solve VPS as a multi-task learning problem via performing semantic segmentation, instance segmentation and multiple object tracking individually. VPSNet [22] is proposed to learn to fuse and warp features. In particular, it applies an optical flow network [14] to align features and attention modules to fuse features [65]. However, these approaches [20, 22] involve many hyper-parameters and ad-hoc designs, leading to a complex system. Bottom-up methods [10, 47, 62] first perform semantic segmentation and then predict near frames’ centers to localize and track each instance where both features are connected by an ASPP module [7] (Fig. 1-(b)). Despite being simpler than top-down approaches, they rely on several post-processing components (i.e., NMS for center grouping, offline mask tracking).

Recently, inspired by the design of object queries in Detection Transformer (DETR) [5], new efforts [11, 56, 74] emerge to explore unified solutions to image panoptic segmentation. In particular, MaskFormer [11] and K-Net [74] unify ‘things’ and ‘stuff’ segmentations by dynamic kernels within the mask classification paradigm. The aforementioned studies in image segmentation motivate us to simplify cumbersome pipelines in video panoptic segmentation. It is noteworthy that several studies [36, 51, 73] in Multi Object Tracking (MOT) also adopt the query design to achieve end-to-end learning. However, most of them introduce *extra tracking queries or boxes* to handle each

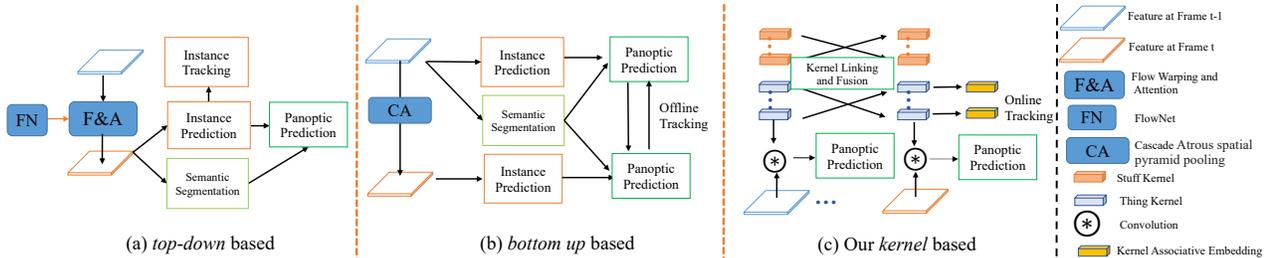


Figure 1. An illustration of previous *top-down* based VPS method (a), *bottom-up* based VPS method (b) and the proposed Video K-Net (c). Unlike previous approaches [22, 47] that perform panoptic segmentation and object tracking with independent modules, our method unifies panoptic segmentation and instance level tracking via kernels in a simpler framework.

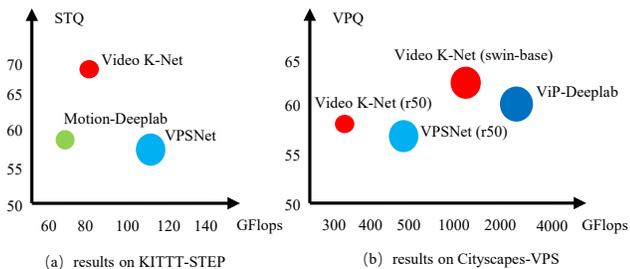


Figure 2. Comparison with previous work. GFLOPS are obtained by using the original image inputs of KITTI-STEP ( $384 \times 1280$ ) and Cityscapes-VPS ( $1024 \times 2048$ ). The number of network parameters are indicated by the radius of circles. Note that our Video K-Net obtains higher accuracy while having lower model complexity on Cityscapes-VPS. Compared to Motion-Deeplab [62], our method achieves a significant gain (10%) with a few extra GFLOPS.

tracked instance.

In this paper, we present Video K-Net, a fully end-to-end framework for video panoptic segmentation. It is observed that the learnable kernels from K-Net [74], which encode object appearances and contexts, are naturally capable of directly associating identical objects across videos. These kernels can be directly used to extract discriminative embeddings for object tracking, which can potentially simplify the MOT framework.

To this end, we design a simple embedding head with a modified contrastive learning loss to obtain better temporal association embeddings. We directly take the dynamic kernels as inputs and output the *kernel association embeddings*. Such design avoids extra tracking query design. The embedding loss forces the kernels to be more distinctive. Then we also link the kernels between adjacent frames for better and more consistent tracking results. Moreover, to jointly learn more consistent temporal mask classification and prediction, we also propose a temporal kernel fusion module to fuse the kernels with corresponding kernel features in a more efficient manner. Fusing at kernel

level avoids computation cost at feature levels (Attention in Fig. 1-(a) and Cascaded ASPP in Fig. 1-(b)). In this way, as shown in Fig. 1-(c), the VPS task can be seen as a kernel linking problem and can significantly reduce the pipeline complexity without extra tracking query and RoI features. The tracking is performed in an online manner. Moreover, it only adds a few extra GFLOPs compared to original K-Net (2.3%GFLOPs). Different from previous works that use RoI heads [22, 41], our kernel based embedding learning avoids box cropping and results in better association performance.

Video K-Net obtains consistent and significant improvements over the strong baseline K-Net with **3.5% STQ** on KITTI-STEP and **4% VPQ** on Cityscapes-VPS. In particular, our method boosts almost **12%** relative improvements over previous bottom up baseline Motion-Deeplab [62]. On VIPSeg dataset [37], our Video K-Net outperforms the co-current baseline CLIP-PanoFCN [37] by **3% VPQ** and boosts almost **16%** relative improvements over CLIP-PanoFCN. Video K-Net achieves new state-of-the-art results on three VPS datasets, including KITTI-STEP [62], Cityscapes-VPS [22], and VIPSeg [37]. Moreover, as shown in Fig. 2, Video K-Net achieves the best trade-off between accuracy and GFLOPs on the two datasets. We further validate the effectiveness of kernel fusion on VSS task with VSPW dataset [38] and we boost previous baselines [8, 76] via considerable margins on two different metrics.

Moreover, we also extend the Video K-Net into clip-level processing via extra temporal kernel fusion module for VIS task [69]. We achieve 40.5% AP using ResNet50 backbone, which lead to **4.0%** improvements over previous VisTR [59] with less GFLOPs. In summary, extensive experiments and analysis on three different types of video segmentation tasks demonstrate that Video K-Net can serve as a new baseline for future research on *unified video segmentation tasks*.

## 2. Related Work

**Panoptic Segmentation.** The goal of this task is to unify the semantic segmentation and instance segmentation into one framework with a single metric named Panoptic Quality (PQ) [25]. Since then, lots of works [9, 10, 24, 27, 30, 44, 57, 66, 67, 70] have been proposed to solve this task with various approaches. However, most methods separate thing and stuff segmentation as individual tasks. Recently, starting from DETR [5], query based approaches [11, 29, 56, 74] unify both things and stuff segmentation as a set prediction problem. In particular, K-Net [74] unifies segmentation tasks in the view of dot product between kernel and feature maps.

**Video Semantic/Instance Segmentation.** Video Semantic Segmentation, a direct extension of semantic segmentation to the video scenario, requires predicting a semantic label to every pixel in each video frame. Several approaches [19, 28, 50, 82] have been proposed in the literature mainly to model the temporal association such as optical flow warping or attention. Recently, the VSPW dataset [38] is proposed to evaluate large scale video semantic segmentation in the wild. Video Instance Segmentation (VIS) [69] extends instance segmentation into video, and it aims to simultaneously classify, segment and track object instances in a given video sequence. Several methods [3, 31, 78] are proposed to link instance-wise feature in video. Recently, several works [21, 59] extend DETR into VIS. Multi-Object Tracking and Segmentation (MOTS) task [55] is proposed to evaluate Multi-Object Tracking along with instance segmentation. However, both VIS and MOTS have limited scale distribution and much fewer objects in the scene.

**Video Panoptic Segmentation.** Video Panoptic Segmentation (VPS) [20, 22] requires generating the instance tracking IDs along with panoptic segmentation results across video clips. Kim *et al.* [22] use Cityscapes video sequences for 6 frames out of each short 30 frame clip and mainly focus on short-term tracks. They proposed Video Panoptic Quality (VPQ) for evaluation. Several works [22, 64] are proposed to solve this task respectively. VIP-Deeplab [47] extends the Panoptic-Deeplab [10] with next frame center map prediction for DVPS task [71]. Then they use such predicted maps for offline tracking. Since Cityscapes VPS only contains short term clips, to allow long term VPS, STEP dataset [62] is proposed. They propose a new metric named Segmentation and Tracking Quality (STQ) that decouples the segmentation and tracking error. They also provide several baselines for reference. Our Video K-Net is verified to work well on both short- and long-term videos.

**Object Tracking.** One of the major tasks in VPS is object tracking. Many studies adopt the tracking-by-detection paradigm [4, 26, 45, 48, 49, 68, 80] and they divide the task into two subtasks where an object detector finds all objects and then a tracking algorithm is employed to associate

Table 1. Toy Experiment results on KITTI-STEP and Cityscapes-VPS set with *STQ* and *VPQ* metrics. Unitrack [60] uses ResNet-50 as the appearance model.

KITTI-STEP		Backbone	STQ	AQ	SQ	-
K-Net	ResNet50		67.5	65.5	68.9	-
K-Net + Unitrack [60]	ResNet50		65.1	64.3	68.9	-
Cityscapes-VPS		Backbone	-	-	-	VPQ
K-Net	ResNet50		-	-	-	54.3
K-Net + Unitrack [60]	ResNet50		-	-	-	53.2

them. There are also several works [2, 41, 43, 61, 75, 79] that detect and track objects at the same time. There are also tracking methods using object queries. For VPS, ViP-DeepLab [47] performs object tracking by clustering all instance pixels. By contrast, our method directly extends and links kernels into kernel association embeddings and avoids such complex post-grouping steps.

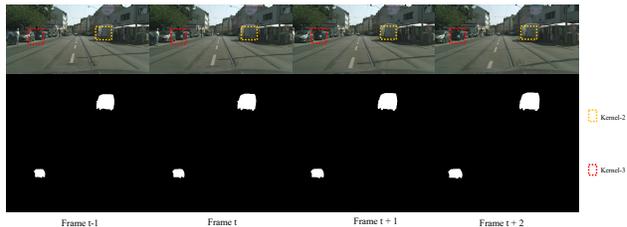


Figure 3. Toy experiment illustration. We use the K-Net directly on Cityscapes video datasets. We find that several instances are originated from **the same kernel** predictions (Red, Yellow boxes, Kernel-2 and Kernel-3). This observation motivates us to use K-Net directly on video. Best view it in color.

## 3. Method

In this section, we will overview the image baseline K-Net. Then we present several toy experiments on using K-Net in video without extra tracking components. Motivated by that, we propose a simple yet effective approach to learn the kernel association embeddings from kernels and introduce two improvements via linking and fusing kernels on K-Net. Finally, we detail the entire framework including both training and inference.

### 3.1. Using K-Net on Video

**Overview of K-Net.** K-Net [74] formulates different image segmentation tasks (semantic, instance, panoptic) into a unified framework via a group of learnable kernels, where each kernel is responsible for generating a mask for either a potential instance or a stuff class. Started with a set of randomly initialized convolutional kernels, the goal of K-Net is to learn kernels in accordance to the segmentation targets (thing or stuff masks). To distinguish and separate

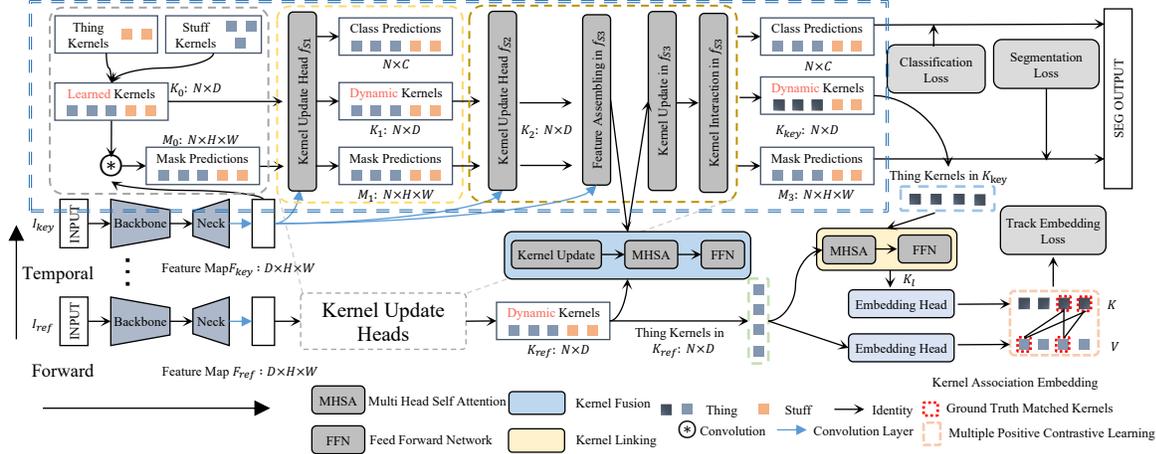


Figure 4. An illustration of our proposed Video K-Net. Our method is based on K-Net [74](in blue dashed box), which is the top-left part of the figure. Video K-Net adds Kernel Fusion at the start phase of the last stage. The Kernel Linking is performed on the output of dynamic kernels. The Embedding Head is appended at the output of kernel linking and takes kernel outputs from both sampled frames.

various instances, a cascaded decoder is proposed to refine the learned kernels iteratively. In particular, the authors propose a kernel update strategy that enables each kernel to be dynamic and conditional on its meaningful area of the input features. As shown in the top-left of Fig. 4, it mainly contains three steps: 1, group feature assembling to obtain corresponding kernel features (in grey color box). 2, adaptive feature update, where the assembled kernel features are used to update the learned kernel (in yellow color box). 3, kernel interaction, where the learned kernels exchange contextual information via self attention layers [53] (in brown color box). Then the final dynamic kernels are used to perform mask classification via fully connected layers. The mask segmentation is done via inner product between the updated kernels and input position-aware features. Then these predictions can be trained in an end-to-end manner with bipartite matching. The above process is repeated several times and the last stage outputs are chosen as the final outputs.

**Toy Experiments on using K-Net on Video.** Before detailing the proposed Video K-Net, we present several toy experiments to show the motivation of our work. As mentioned in the previous section, refining kernels is the key for making K-Net work. We first depict kernel indexes of a trained K-Net on Cityscapes video sequence [12] in Fig. 3. We use K-Net directly into video segmentation task. We argue that the *updated kernels contain discriminative information*, which can be directly used to track the instance in video even *without* adding extra tracking heads. The insights are as follows: 1, Each output instance mask corresponds with one specific kernel. 2, Each kernel absorbs the position-aware features during the adaptive feature update, where the instance aware information has already merged

into each kernel. Thus, the same instance can be decoded from the previous kernel. In Fig. 3, we present one visual examples on Cityscapes-VSP. We train the K-Net directly on both datasets *without* adding tracking components. As shown in that figure, in three frames, several kernels are well linked, such as person and car.

Moreover, we take a further step by directly using the learned kernels as appearance embedding for tracking association. In particular, we use the Quasi-Dense Tracker [41] which is a pure appearance based tracker. We replace the appearance embeddings with our learned kernels from the last stage. As shown in Tab.1, surprisingly, we found the performance of original K-Net is already good enough to perform tracking and even achieves better results than Uni-track [60] that comes with extra appearance network and motion prediction module [63]. These findings motivate us to dive into the design of linking kernels along the temporal dimension.

### 3.2. Extending K-Net into Video

Despite the competitive performance of the original K-Net, it shows several failure cases, for instance, the fast motion object shown in Fig. 3. Thus, we design three improvements on K-Net including the learning of kernel association embeddings via a modified contrastive learning loss, learning to link tracking kernels, and learning to fuse kernels, respectively.

As shown in Fig. 4, given a key image  $I_{key}$  for training, we randomly select a reference image  $I_{ref}$  from its temporal neighborhood. The neighbor distance is constrained by a window size ranging from  $[-2, 2]$  in our experiments by default. Then we use K-Net to obtain the two individual kernels:  $K_{key}$  and  $K_{ref}$ .

**Learning Kernel Association Embeddings.** Our method is motivated by the recent work Quasi-Dense [41] in MOT. Instead of using RoI boxes for feature embedding extraction, we use the kernels directly. We add an extra lightweight embedding head after the original K-Net decoder, to extract embedding features for each kernel. The embedding head is implemented via several fully connected layers. We adopt mask-based assignment for Quasi-Dense learning. A kernel embedding is defined as positive to an object if their corresponding masks have an IoU higher than  $\alpha_1$ , or negative if their corresponding masks have an IoU lower than  $\alpha_2$ . The values of  $\alpha_1$  and  $\alpha_2$  are set as 0.7 and 0.3 in our experiments. The matching of kernels on two sampled frames is positive if the two regions are associated with the same object, and negative otherwise. Instead of following the original design [41] that uses hundreds of proposals during training, we *only optimize kernels that are matched* with ground truth masks. This is because most kernels (unmatched) are not accurate and may cause noise during training, leading to inferior results.

Assume there are  $V$  matched kernels on the key frame as training samples and  $K$  matched kernels on the reference frame as contrastive targets, where both  $V$  and  $K$  are always fewer than kernel number  $N$ . As shown in red boxes of the Fig. 4, these kernels are sampled from  $K_{key}$  and  $K_{ref}$ , respectively.

$$\mathcal{L}_{\text{track}} = - \sum_{\mathbf{k}^+} \log \frac{\exp(\mathbf{v} \cdot \mathbf{k}^+)}{\exp(\mathbf{v} \cdot \mathbf{k}^+) + \sum_{\mathbf{k}^-} \exp(\mathbf{v} \cdot \mathbf{k}^-)}, \quad (1)$$

where  $\mathbf{v}$ ,  $\mathbf{k}^+$ ,  $\mathbf{k}^-$  are kernel embeddings of the training sample, its positive targets, and negative targets in  $K$ . In addition, following previous work [41], we also adopt L2 loss as an auxiliary loss.

$$\mathcal{L}_{\text{aux}} = \left( \frac{\mathbf{v} \cdot \mathbf{k}}{\|\mathbf{v}\| \cdot \|\mathbf{k}\|} - c \right)^2, \quad (2)$$

where  $c$  is 1 if the match of two samples is positive and 0 otherwise. In comparison to the original work, our modification makes the loss calculation process sparse, benefiting the training process and achieving better results. (See the experiment part in Sec. 4.3). In summary, we term this procedure as learning Kernel Association Embedding (KAE).

**Learning to Link Kernels.** Beyond supervision, we take a further step to link kernels  $K_{key}$  and  $K_{ref}$  during training and inference. This forces the kernels to perform interaction along the temporal dimension. We adopt one self-attention layer (*MHSA*: Multi Head Self Attention) with a Feed Forward Network (*FFN*) [53] to learn the correspondence among each query to obtain the updated queries, allowing the full correlation among queries. This process is shown as follows:

$$K_l = FFN(MHSA(K_{key}, K_{ref}, K_{ref}) + K_{key}), \quad (3)$$

where the query, key and value are  $K_{key}$ ,  $K_{ref}$  and  $K_{ref}$ , respectively. In this way, kernels from the reference frame are propagated into key frame via the affinity matrix between kernels. Then the linked kernel  $K_l$  is sent to the tracking head. During the training, this operation is *only optimized with matched thing kernels*. During the inference stage, this operation is performed on *final preserved kernels* in panoptic segmentation map, where other thing kernels are dropped.

**Learning to Fuse Kernels.** The previous linking step may focus on the tracking consistency while ignoring the segmentation consistency. To address this issue, we propose to fuse kernels in between frames on K-Net. In particular, as shown in Fig. 4, we use input kernels of the last stage from two frames. Then we perform kernel updating, where we use kernel features of the current frame to update the previous kernel (Step-2 from original K-Net). Then we adopt one self-attention layer with feed forward layers [53] to fuse the previous kernels into the current frame. The updating step is essential. Directly fusing kernels leads to bad results, since large motion and scale variation often occur in the video scene. More details can be found in Sec. 4.3. In summary, through the above steps, we link kernels of K-Net into the video domain with a little extra computation cost.

### 3.3. Video K-Net Architecture

**Network Architecture.** By default, we adopt the panoptic segmentation setting of K-Net. The kernels are composed of instance kernels and semantic kernels for thing and stuff mask prediction. We adopt semantic FPN [24] for producing high resolution feature map with extra positional encoding as in [5, 53]. The kernel linking operation is appended at the last stage of the K-Net decoder. The kernel fusion is placed at the beginning part. The lightweight tracking head is appended after the decoder. Note that backbone, neck, update head, and embedding head are shared across frames. For VSS task, we directly append the head of Video K-Net at the end of existing semantic segmentation network [8].

**Loss Function and Training.** Compared to the original K-Net, we add the tracking loss at the end of K-Net loss. The total loss function for all kernels is shown as  $L = \lambda_{cls} L_{cls} + \lambda_{ce} L_{ce} + \lambda_{dice} L_{dice} + \lambda_{track} L_{track} + \lambda_{aux} L_{aux}$ , where  $L_{cls}$  is Focal loss [32] for classification, and  $L_{ce}$  and  $L_{dice}$  are Cross Entropy (CE) loss and Dice loss [39, 58] for segmentation, respectively.  $L_{track}$  and  $L_{aux}$  are the tracking loss for instance kernels only. We adopt Hungarian assignment strategy used in [5] for target assignment for end-to-end training. It builds a one-to-one mapping between the predicted instance masks and the ground-truth instances based on the masked-based matching costs. During training, we found that applying the segmentation loss on both key frame and reference frames leads to better results.

**Inference.** We use K-Net to generate panoptic segmenta-

tion results, where we paste thing and stuff masks in a mixed order. We use the kernel embeddings from the learned embedding head as association features. We take these features as the input of Quasi-Dense Tracker [41] where the bidirectional softmax is calculated to associate the instances between two frames in an online manner. Noted that we only track the preserved instance masks from panoptic segmentation maps to save computation cost. We also compare our tracking algorithm with the original Quasi-Dense Tracker [41] using the boxes for appearance modeling.

**Clip-level Training and Inference for VIS.** We also extend our Video K-Net into Clip-level training and inference pipeline for Video Instance Segmentation [69]. In particular, we add three kernel fusion layers on the top of K-Net outputs to jointly fuse both kernels and the corresponding kernel features along the temporal dimension. We use the mean of temporal kernels to represent each object for each clip and assign instance ID to each instance in each frame via kernel index [59].

## 4. Experiment

### 4.1. Experiment Setup

**Dataset.** We carry out experiments on four video-level datasets: KITTI-STEP, Cityscapes-VPS, VSPW, VIPSeg, and YouTube-VIS. KITTI-STEP has 21 and 29 sequences for training and testing, respectively. The training sequences are further split into a training set (12 sequences) and a validation set (9 sequences). Cityscapes-VPS contains 400 training, 50 validation, and 50 test videos. Following previous work [22], all 30 frames are predicted during the inference, and only 6 frames with ground truth are evaluated. Both Cityscapes-VPS and KITTI-STEP have the same class number with Cityscapes dataset [12]. However, the definition of thing and class is different. For VSPW dataset and YouTube-VIS dataset, refer to *the appendix file*.

**Implementation Details.** We implement our models in PyTorch [42] with MMDetection toolbox [6]. We use the distributed training framework with 8 GPUs. Each mini-batch has one image per GPU. Following previous work, we first pretrain the image baseline on Cityscapes dataset [12] where we adopt the same pretraining setting [10, 62] for fair comparison. For STEP pretraining, we change the class distribution of Cityscapes dataset into STEP format to avoid overfitting on STEP. Both ResNet [18] and Swin Transformer [34] are adopted as the backbone networks, and other layers use Xavier initialization [17]. The optimizer is AdamW [35] with weight decay 0.0001. We adopt full image size for random crop in both pretraining and training process. For the large Swin Transformer backbone [34], we also pretrain our model on Mapillary dataset [40]. More details of pretraining and finetuning can be found in the appendix file. Due to diversity of dataset,

Table 2. **Experiment results on KITTI set with both  $STQ$  and  $VPQ$  metric.** OF refers to an optical flow network [52]. The results on validation set are shown in the several top rows, and results on test set are in the bottom rows. P means Panoptic Deeplab [10]. Following [62], we keep two decimal numbers.  $VPQ$  is obtained via average results of window size  $k$  where  $k = 1, 2, 3, 4$  [62]. Top: validation set. Bottom: test set. We find 0.5% noise on this dataset, where we report the average results (three times). Axial means using extra Axial Attention [57].

KITTI-STEP	Backbone	OF	STQ	AQ	SQ	VPQ
P + IoU Assoc.	ResNet50		0.58	0.47	0.71	0.44
P + SORT	ResNet50		0.59	0.50	0.71	0.42
P + Mask Propagation	ResNet50	✓	0.67	0.63	0.71	0.44
Motion-Deeplab [62]	ResNet50		0.58	0.51	0.67	0.40
VPSNet [22]	ResNet50	✓	0.56	0.52	0.61	0.43
TubeFormer-Deeplab [23]	ResNet-50 + axial		0.70	0.64	0.76	0.51
Video K-Net	ResNet50		0.71	0.70	0.71	0.46
Video K-Net	Swin-base		0.73	0.72	0.73	0.53
Video K-Net	Swin-large		0.74	0.73	0.75	0.56
Motion-Deeplab [62]	ResNet50		0.52	0.46	0.60	-
Video K-Net	ResNet50		0.59	0.50	0.62	-
Video K-Net	Swin-base		0.63	0.60	0.65	-

for VSPW, YouTube-VIS and VIP-Seg dataset, refer to *the appendix file*.

**Evaluation Metrics.** For VPS task, as discussed in STEP [62], different metrics result in a significant gap even on the same VPS dataset since different metrics emphasize different properties. We adopt two widely used metrics: Video Panoptic Quality ( $VPQ$ ) and Segmentation and Tracking Quality ( $STQ$ ). The former mainly focuses on mask proposal level as PQ [25] with different window sizes and threshold parameters, while the latter emphasizes pixel level segmentation and tracking without any thresholds.  $STQ$  contains geometric mean of two items: Segmentation Quality ( $SQ$ ) and Association Quality ( $AQ$ ) where  $STQ = (SQ \times AQ)^{\frac{1}{2}}$ . The former evaluates the pixel-level tracking while the latter evaluates the pixel-level segmentation results in a video clip. Due to the interpretability [62] of  $STQ$ , we adopt it for ablation studies. We also report  $VPQ$  for the benchmark comparison on both datasets. For VSS task, we report the Mean Intersection over Union (mIoU) and mean Video Consistency ( $mVC$ ) [38] for reference.

### 4.2. Main Results

**Results on KITTI-STEP.** As shown in Tab. 2, our method achieves 0.71  $STQ$  and 0.46  $VPQ$  using ResNet50 as backbone and achieves the new state-of-the-art results among previous works. After applying a stronger backbone [34], we obtain about 3% gains on  $STQ$  and 7% gains on  $VPQ$ . This suggests strong feature representation leads to better segmentation-level prediction than pixel-level prediction. Compared with VPSNet [22] and Motion-Deeplab [62], our method does not need post-processing steps or extra optical flow to warp features. Compared with Motion-Deeplab [62], our method outperforms it by 13% and 7% on validation set and test set respectively. The results sug-

Table 3. **Results on Cityscapes-VPS validation set.**  $k$  is temporal window size in [22]. All the methods use the single scale inference without other augmentations in the test stage. In each cell, we report  $VPQ$ ,  $VPQ_{thing}$  and  $VPQ_{stuff}$  in order. There is about 0.5% noise on this dataset, where we report the average results (three times).

Method	Backbone	k=0			k=5			k=10			k=15			Average		
VPSNet [22]	ResNet50	65.0	59.0	69.4	57.6	45.1	66.7	54.4	39.2	65.6	52.8	35.8	65.3	57.5	44.8	66.7
SiamTrack [64]	ResNet50	64.6	58.3	69.1	57.6	45.6	66.6	54.2	39.2	65.2	52.7	36.7	64.6	57.3	44.7	55.0
ViP-Deeplab [47]	WideResNet41 [72]	68.2	N/A	N/A	61.3	N/A	N/A	58.2	N/A	N/A	56.2	N/A	N/A	60.9	N/A	N/A
ViP-Deeplab [47]	WideResNet41 [72]+RFP [46] + AutoAug [13]	69.2	N/A	N/A	62.3	N/A	N/A	59.2	N/A	N/A	57.0	N/A	N/A	61.9	N/A	N/A
Video K-Net	ResNet50	65.6	57.4	71.5	57.7	43.4	68.2	54.2	36.5	67.1	52.3	33.1	66.3	57.8	45.0	66.9
Video K-Net	Swin-base [34]	69.2	63.6	73.3	62.0	51.1	70.0	58.4	44.7	68.3	55.8	39.8	67.5	61.2	49.6	69.5
Video K-Net	Swin-base + RFP [46]	70.8	63.2	76.3	63.1	49.3	73.2	59.5	43.4	72.0	56.8	37.0	71.1	62.2	49.8	71.8

Table 4. **Results on VSPW validation set.**  $mVC_c$  means that a clip with  $c$  frames is used. All methods use the same setting for fair comparison.

VPSW	Backbone	mIoU	$mVC_8$	$mVC_{16}$
DeepLabv3+ [8]	ResNet101	35.7	83.5	78.4
PSPNet+ [76]	ResNet101	36.5	84.4	79.8
TCB(PSPNet) [38]	ResNet101	37.5	86.9	82.1
Video K-Net (DeepLabv3+)	ResNet101	37.9	87.0	82.1
Video K-Net (PSPNet)	ResNet101	38.0	87.2	82.3

Table 5. **Results on Video instance segmentation AP (%) on the YouTube-VIS-2019 [69] validation dataset.** \* means using deformable fpn [81]. Axial means using extra Axial Attention [56]. The compared methods are listed by publication date.

Method	backbone	AP	AP <sub>50</sub>	AP <sub>75</sub>	AR <sub>1</sub>	AR <sub>10</sub>
FEELVOS [54]	ResNet50	26.9	42.0	29.7	29.9	33.4
MaskTrack R-CNN [69]	ResNet50	30.3	51.1	32.6	31.0	35.5
MaskProp [3]	ResNet-50	40.0	-	42.9	-	-
MaskProp [3]	ResNet101	42.5	-	45.6	-	-
STEm-Seg [1]	ResNet50	30.6	50.7	33.5	31.6	37.1
STEm-Seg [1]	ResNet101	34.6	55.8	37.9	34.4	41.6
CompFeat [15]	ResNet50	35.3	56.0	38.6	33.1	40.3
VisTR [59]	ResNet50	36.2	59.8	36.9	37.2	42.4
VisTR [59]	ResNet101	40.1	64.0	45.0	38.3	44.9
TubeFormer-DeepLab [23]	ResNet-50 + Axial	38.8	-	-	44.0	51.4
Video K-Net	ResNet50	40.5	63.5	44.5	40.7	49.9
Video K-Net	Swin-base	51.4	77.2	56.1	49.0	58.4
Video K-Net	Swin-base*	54.1	79.0	59.5	49.7	59.9

Table 6. **Results on VIPSeg-VPS [37] validation dataset.** We report VPQ and STQ for reference. Following work [37], we report VPQ score at different window sizes (1,2,4,6).

Method	backbone	$VPQ^1$	$VPQ^2$	$VPQ^4$	$VPQ^6$	VPQ	STQ
ViP-DeepLab [47]	ResNet50	18.4	16.9	14.8	13.7	16.0	22.0
VPSNet [22]	ResNet50	19.9	18.1	15.8	14.5	17.0	20.8
SiamTrack [64]	ResNet50	20.0	18.3	16.0	14.7	17.2	21.1
Clip-PanoFCN [37]	ResNet50	24.3	23.5	22.4	21.6	22.9	31.5
Video K-Net	ResNet50	29.5	26.5	24.5	23.7	26.1	33.1
Video K-Net	Swin-base	43.3	40.5	38.3	37.2	39.8	46.3

gest the suitability of our framework to serve as a new and stronger baseline. We find 0.5% noise on STEP validation set.

**Results on Cityscape-VPS.** Tab. 3 compares our method with previous works on Cityscapes-VPS datasets. For ResNet50, our method outperforms VPSNet [22] by 0.3%. Moreover, we follow the same pretraining procedure of

ViP-Deeplab [47] and achieve 62.2% VPQ. Our method with Swin-Transformer base achieves better results than ViP-Deeplab, which is trained with stronger data augmentation [13]. Note that we find 0.5% noise for this dataset. Moreover, we do not use the kernel fusion on this dataset due to the low frame rate of Cityscapes-VPS.

**Results on VSPW.** We further carry out experiments on VSPW dataset for VSS task to prove the generalization of Video K-Net. All the methods are re-implemented on our codebase and achieve higher results than the original paper [38]. As shown in Tab. 4, Video K-Net boosts Deeplabv3+ [8] and PSPNet [76] by a significant margin on both mIoU (2-3%) and  $mVC$  (3-4%). This proves the generalization ability of both kernel fusion and kernel linking.

**Results on Youtube-VIS-2019.** In Tab. 5, we report the results on Youtube-VIS datasets. Compared with previous work VisTR [59], our method achieves better results (4%) but with less training times (12 epochs vs 300 epochs).

**Results on VIPSeg-VPS.** In Tab. 6, we further report the results on recently challenging VPS dataset VIPSeg [37]. Using ResNet50 backbone, our Video K-Net outperforms Clip-PanoFCN [37] by 3% VPQ and 1.6% STQ with only 12 epoch training (half training time of Clip-PanoFCN). We further use a larger model (Swin-base), which results in new state-of-the-art results and outperforms previous work by 16.6% VPQ and 14.8% STQ.

### 4.3. Ablation Study

**Ablation Study on Each Component.** In Tab. 7a, we first perform ablation studies on the effectiveness of each component. Adding KAE yields 2.8% STQ improvements with more significant gains on AQ (about 2.4%). It shows that directly learning kernel association is a key factor for tracking. Adding Kernel Linking module (KL) further improves AQ by 1.0%. Finally, appending kernel fusing results in 0.7% STQ improvement and 1.4 % improvement on SQ. The AQ decreases a little, mainly because both KL and KF are fighting for tracking quality and segmentation quality, respectively.

**Needs of Extra Appearance Embedding.** In Tab. 7b, we present detailed comparison with recent box based appearance embedding and mask based appearance embedding. In

Table 7. Ablation studies and comparison analysis on KITTI-STEP validation set. All the experiments use ResNet-50 as backbone.

(a) Ablation Study on Each Components.							(b) Needs of Appearance Embeddings			(c) Effect of sampling in association.			
baseline	KAE	KL	KF	STQ	AQ	SQ	Method	AQ	STQ	Method	STQ	AQ	SQ
K-Net				67.5	65.5	68.9	RoI-Align [41]	68.8	69.1	K-Net	67.5	65.5	68.9
	✓			69.3	69.0	69.8	Mask-Emb [64]	67.3	68.1	GT-based (ours)	69.3	69.0	69.8
	✓	✓		70.2	71.2	69.7	Ours	70.8	70.9	sampling in [41]	63.1	62.1	64.3
	✓	✓	✓	70.9	70.8	71.2	Ours + Mask-Emb [64]	70.3	70.8				

(d) Ablation Study on Linking and Fusing Stage.				(e) Ablation Study on Training Settings				(f) Ablation Study on Kernel Fusing			
Stage	STQ	AQ	SQ	Settings	STQ	AQ	SQ	Settings	STQ	AQ	SQ
3	70.9	70.8	71.2	joint training	70.9	70.8	71.2	K-Net	67.5	65.5	68.9
2	68.5	68.2	69.3	only train the key frame	70.1	70.1	69.8	w Update	70.9	70.8	71.2
1	66.9	63.4	67.3					w/o Update	67.1	66.2	68.3

particular, we re-implement their methods on our baseline K-Net. We also use the same embedding loss functions: Equ. 1 and Equ. 2 for fair comparison. From the table, our simple design leads to the best result. From the last row of Tab. 7b, even combing both the appearance embedding and kernels, there is no clear gain. We implement this by adding mask-grouped appearance features from backbone and kernel embeddings together. That indicates kernels have already encoded the discriminative information, which is consistent with our toy experiments results: the pure kernel information is good enough for tracking.

**Effect of Sampling Examples in Association.** Pang *et al.* [41] employ RPN to generate large samples to improve the tracking results. However, we find that directly using this method into our framework leads to inferior results because most kernels are not used to generate masks with Hungarian assignment strategy during the training stage. Adding more samples brings noise to the segmentation prediction shown in Tab. 7c, causing a marked decrease in SQ metric. Thus, we adapt to assign training kernel samples that are matched with ground truth masks. The matched kernels are more robust during training.

**Ablation on Link Stage.** In Tab. 7d, we present ablation on the stage of linking kernels, where we find linking kernels at the last stage achieves the best results. Early fusing of kernels also leads to bad results, since initial kernels are not very accurate. This results in bad misalignment on both kernel and kernel features.

**Effect on Training Settings.** Moreover, in Tab. 7e, we find joint training with both frames  $I_{key}$  and  $I_{ref}$  performs better than only training  $I_{key}$ . This is mainly because both the proposed Kernel Linking and Kernel Fusion need to learn to group similar kernels and separate the different kernels.

**Ablation on Kernel Fusing Design.** In Tab. 7f, we find the importance of Kernel Updating module in Kernel Fusing.

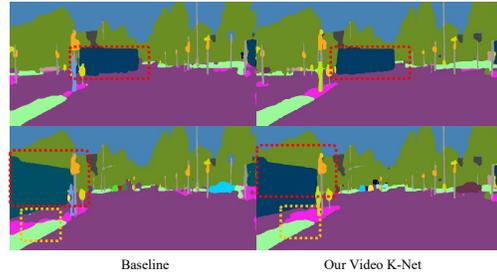


Figure 5. Visualization of improvements over baseline methods. Tracking improvement in red boxes and segmentation improvement in yellow boxes. Best viewed in color and by zooming in.

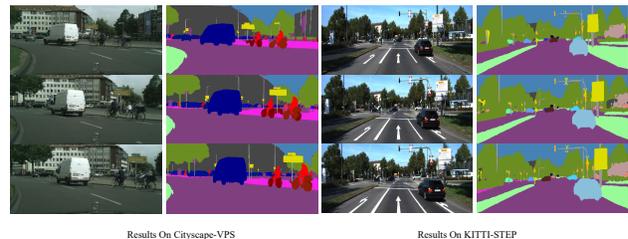


Figure 6. Visual Examples of our Video K-Net. The left is on Cityscapes-VPS validation set, while the right is on KITTI-STEP validation set (video sequences from top to down). The same instances are shown with the same color. Best view it on screen.

Since previous features are not aligned to the current frame, we use current-frame features to re-weight previous kernels to filter out noise and outliers in the previous kernel.

#### 4.4. Visualization and More Analysis

**Visualization over Baseline.** In Fig. 5, we visualize the improvements over the K-Net baseline. The red boxes show the tracking consistency, while the yellow boxes show the

segmentation consistency. Our Video K-Net achieves better results on both tracking and segmentation qualities.

**More Qualitative Results.** In Fig. 6, we present several visual results on Cityscapes-VPS and KITTI-STEP datasets. We find most foreground objects (person and car) can be well tracked and segmented. *More visual examples can be found in the appendix file.*

**Parameter and GFLOPs** Compared with K-Net baseline, our method only adds about **2.3% GFLOPs** and **1.8% Parameters** given  $800 \times 1333$  image inputs.

**Limitation and Future Work.** One limitation of Video K-Net is that we only explore one frame during both training and inference. This setting is mainly for fair comparison with other works [22, 62]. Both exploring the *multiple frames* information via kernels and handling on long video inputs will be our future work.

## 5. Conclusion

We present **Video K-Net** a simple, strong and unified system for fully end-to-end video panoptic segmentation. We propose to learn the kernel association embeddings directly from the encoded ‘thing’ kernels. Then we link and fuse the kernels to jointly improve both tracking and segmentation results. Despite simplicity, our method achieves state-of-the-art results on three VPS datasets including Cityscapes-VPS, KITTI-STEP and recently new proposed VIPSeg datasets. In particular, our method boosts previous methods over 10% improvements on KITTI-STEP and Cityscapes-VPS datasets. In particular, compared to previous works, our method has a much simpler pipeline but achieves better results along with less computation cost. We also show the generalization ability on the VSPW dataset for VSS task and competitive results on YouTube-VIS dataset for VIS task. We hope our proposed Video K-Net would be a simple and strong baseline and benefit unified video segmentation field.

**Acknowledgement.** This study is partly supported under the RIE2020 Industry Alignment Fund Industry Collaboration Projects (IAF-ICP) Funding Initiative, as well as cash and in-kind contribution from the industry partner(s). It is also partially supported by the NTU NAP grant. This research is also supported by the National Key Research and Development Program of China under Grant No. 2020YFB2103402. We also thank for the GPU resource provided by SenseTime Research.

## References

- [1] Ali Athar, Sabarinath Mahadevan, Aljoša Ošep, Laura Leal-Taixé, and Bastian Leibe. Stem-seg: Spatio-temporal embeddings for instance segmentation in videos. In *ECCV*, 2020. 7
- [2] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixé. Tracking without bells and whistles. In *ICCV*, pages 941–951, 2019. 3
- [3] Gedas Bertasius and Lorenzo Torresani. Classifying, segmenting, and tracking object instances in video with mask propagation. In *CVPR*, 2020. 1, 3, 7
- [4] Alex Bewley, Zongyuan Ge, Lionel Ott, Fabio Ramos, and Ben Upcroft. Simple online and realtime tracking. In *ICIP*, pages 3464–3468. IEEE, 2016. 3
- [5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 1, 3, 5, 12
- [6] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, et al. Mmdetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019. 6
- [7] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv:1706.05587*, 2017. 1
- [8] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 2018. 2, 5, 7
- [9] Yifeng Chen, Guangchen Lin, Songyuan Li, Omar Bourahla, Yiming Wu, Fangfang Wang, Junyi Feng, Mingliang Xu, and Xi Li. Banet: Bidirectional aggregation network with occlusion handling for panoptic segmentation. In *CVPR*, 2020. 3
- [10] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and Liang-Chieh Chen. Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. In *CVPR*, 2020. 1, 3, 6, 12
- [11] Bowen Cheng, Alexander G. Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. *arXiv*, 2021. 1, 3
- [12] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016. 1, 4, 6
- [13] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018. 7
- [14] 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 *CVPR*, 2015. 1
- [15] Yang Fu, Linjie Yang, Ding Liu, Thomas S. Huang, and Humphrey Shi. ComFeat: Comprehensive feature aggregation for video instance segmentation. *AAAI*, 2021. 7
- [16] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012. 1

- [17] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010. 6
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6
- [19] Ping Hu, Fabian Caba, Oliver Wang, Zhe Lin, Stan Sclaroff, and Federico Perazzi. Temporally distributed networks for fast video semantic segmentation. *CVPR*, 2020. 3
- [20] Juana Valeria Hurtado, Rohit Mohan, Wolfram Burgard, and Abhinav Valada. Mopt: Multi-object panoptic tracking. *arXiv preprint arXiv:2004.08189*, 2020. 1, 3
- [21] Sukjun Hwang, Miran Heo, Seoung Wug Oh, and Seon Joo Kim. Video instance segmentation using inter-frame communication transformers. *NIPS*, 2021. 3
- [22] Dahun Kim, Sanghyun Woo, Joon-Young Lee, and In So Kweon. Video panoptic segmentation. In *CVPR*, 2020. 1, 2, 3, 6, 7, 9, 12
- [23] Dahun Kim, Jun Xie, Huiyu Wang, Siyuan Qiao, Qihang Yu, Hong-Seok Kim, Hartwig Adam, In So Kweon, and Liang-Chieh Chen. Tubformer-deeplab: Video mask transformer. In *CVPR*, 2022. 6, 7
- [24] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *CVPR*, 2019. 3, 5
- [25] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *CVPR*, 2019. 1, 3, 6
- [26] Laura Leal-Taixé, Cristian Canton-Ferrer, and Konrad Schindler. Learning by tracking: Siamese cnn for robust target association. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 33–40, 2016. 3
- [27] Jie Li, Allan Raventos, Arjun Bhargava, Takaaki Tagawa, and Adrien Gaidon. Learning to fuse things and stuff. *arXiv:1812.01192*, 2018. 3
- [28] Xiangtai Li, Hao He, Yibo Yang, Henghui Ding, Kuiyuan Yang, Guangliang Cheng, Yunhai Tong, and Dacheng Tao. Improving video instance segmentation via temporal pyramid routing. *T-PAMI*, 2022. 3
- [29] Xiangtai Li, Shilin Xu, Yibo Yang, Guangliang Cheng, Yunhai Tong, and Dacheng Tao. Panoptic-partformer: Learning a unified model for panoptic part segmentation. In *ECCV*, 2022. 3
- [30] Yanwei Li, Hengshuang Zhao, Xiaojuan Qi, Liwei Wang, Zeming Li, Jian Sun, and Jiaya Jia. Fully convolutional networks for panoptic segmentation. *CVPR*, 2021. 3
- [31] Huaijia Lin, Ruizheng Wu, Shu Liu, Jiangbo Lu, and Jiaya Jia. Video instance segmentation with a propose-reduce paradigm. *ICCV*, 2021. 3
- [32] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 5
- [33] 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. 12
- [34] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *ICCV*, 2021. 6, 7, 12
- [35] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2017. 6
- [36] Tim Meinhardt, Alexander Kirillov, Laura Leal-Taixe, and Christoph Feichtenhofer. Trackformer: Multi-object tracking with transformers. *arXiv preprint arXiv: 2101.02702*, 2021. 1
- [37] Jiaxu Miao, Xiaohan Wang, Yu Wu, Wei Li, Xu Zhang, Yunchao Wei, and Yi Yang. Large-scale video panoptic segmentation in the wild: A benchmark. In *CVPR*, 2022. 2, 7, 12
- [38] Jiaxu Miao, Yunchao Wei, Yu Wu, Chen Liang, Guangrui Li, and Yi Yang. Vspw: A large-scale dataset for video scene parsing in the wild. In *CVPR*, 2021. 1, 2, 3, 6, 7
- [39] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-Net: Fully convolutional neural networks for volumetric medical image segmentation. In *3DV*, 2016. 5
- [40] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Buló, and Peter Kotschieder. The mapillary vistas dataset for semantic understanding of street scenes. In *ICCV*, 2017. 6, 12
- [41] Jiangmiao Pang, Linlu Qiu, Xia Li, Haofeng Chen, Qi Li, Trevor Darrell, and Fisher Yu. Quasi-dense similarity learning for multiple object tracking. In *CVPR*, June 2021. 2, 3, 4, 5, 6, 8
- [42] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703*, 2019. 6
- [43] Jinlong Peng, Changan Wang, Fangbin Wan, Yang Wu, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, and Yanwei Fu. Chained-tracker: Chaining paired attentive regression results for end-to-end joint multiple-object detection and tracking. In *ECCV*, 2020. 3
- [44] Lorenzo Porzi, Samuel Rota Buló, Aleksander Colovic, and Peter Kotschieder. Seamless scene segmentation. In *CVPR*, 2019. 3
- [45] Lorenzo Porzi, Markus Hofinger, Idoia Ruiz, Joan Serrat, Samuel Rota Buló, and Peter Kotschieder. Learning multi-object tracking and segmentation from automatic annotations. In *CVPR*, pages 6846–6855, 2020. 3
- [46] Siyuan Qiao, Liang-Chieh Chen, and Alan Yuille. Detectors: Detecting objects with recursive feature pyramid and switchable atrous convolution. In *CVPR*, 2021. 7
- [47] Siyuan Qiao, Yukun Zhu, H. Adam, A. Yuille, and Liang-Chieh Chen. Vip-deeplab: Learning visual perception with depth-aware video panoptic segmentation. *CVPR*, 2021. 1, 2, 3, 7, 12
- [48] Samuel Schulter, Paul Vernaza, Wongun Choi, and Manmohan Chandraker. Deep network flow for multi-object tracking. In *CVPR*, pages 6951–6960, 2017. 3
- [49] Sarthak Sharma, Junaid Ahmed Ansari, J Krishna Murthy, and K Madhava Krishna. Beyond pixels: Leveraging ge-

- ometry and shape cues for online multi-object tracking. In *ICRA*, pages 3508–3515, 2018. 3
- [50] Evan Shelhamer, Kate Rakelly, Judy Hoffman, and Trevor Darrell. Clockwork convnets for video semantic segmentation. In *ECCV*, pages 852–868. Springer, 2016. 1, 3
- [51] Peize Sun, Jinkun Cao, Yi Jiang, Rufeng Zhang, Enze Xie, Zehuan Yuan, Changhu Wang, and Ping Luo. Transtrack: Multiple-object tracking with transformer. *arXiv preprint arXiv: 2012.15460*, 2020. 1
- [52] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, pages 402–419. Springer, 2020. 6
- [53] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017. 4, 5
- [54] Paul Voigtlaender, Yuning Chai, Florian Schroff, Hartwig Adam, Bastian Leibe, and Liang-Chieh Chen. Feelvos: Fast end-to-end embedding learning for video object segmentation. In *CVPR*, 2019. 7
- [55] Paul Voigtlaender, Michael Krause, Aljosa Osep, Jonathon Luiten, Berin Balachandar Gnana Sekar, Andreas Geiger, and Bastian Leibe. Mots: Multi-object tracking and segmentation. In *CVPR*, 2019. 3
- [56] Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Max-deeplab: End-to-end panoptic segmentation with mask transformers. *arXiv preprint arXiv:2012.00759*, 2020. 1, 3, 7
- [57] Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Axial-deeplab: Stand-alone axial-attention for panoptic segmentation. In *ECCV*, 2020. 3, 6
- [58] Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. Solo: Segmenting objects by locations. In *ECCV*, 2020. 5
- [59] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *CVPR*, 2021. 2, 3, 6, 7, 12
- [60] Zhongdao Wang, Hengshuang Zhao, Ya-Li Li, Shengjin Wang, Philip HS Torr, and Luca Bertinetto. Do different tracking tasks require different appearance models? *NeurIPS*, 2021. 3, 4
- [61] Zhongdao Wang, Liang Zheng, Yixuan Liu, and Shengjin Wang. Towards real-time multi-object tracking. *arXiv preprint arXiv:1909.12605*, 2019. 3
- [62] M. Weber, J. Xie, M. Collins, Yukun Zhu, P. Voigtlaender, H. Adam, B. Green, A. Geiger, B. Leibe, D. Cremers, Aljosa Osep, L. Leal-Taixé, and Liang-Chieh Chen. Step: Segmenting and tracking every pixel. *NIPS*, 2021. 1, 2, 3, 6, 9, 12
- [63] Nicolai Wojke, Alex Bewley, and Dietrich Paulus. Simple online and realtime tracking with a deep association metric. In *ICIP*, 2017. 4
- [64] Sanghyun Woo, Dahun Kim, Joon-Young Lee, and In So Kweon. Learning to associate every segment for video panoptic segmentation. In *CVPR*, 2021. 3, 7, 8
- [65] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *ECCV*, pages 3–19, 2018. 1
- [66] Yangxin Wu, G. Zhang, Hang Xu, Xiaodan Liang, and Liang Lin. Auto-panoptic: Cooperative multi-component architecture search for panoptic segmentation. *NIPS*, 2020. 3
- [67] Yuwen Xiong, Renjie Liao, Hengshuang Zhao, Rui Hu, Min Bai, Ersin Yumer, and Raquel Urtasun. Upsnet: A unified panoptic segmentation network. In *CVPR*, 2019. 3
- [68] Jiarui Xu, Yue Cao, Zheng Zhang, and Han Hu. Spatial-temporal relation networks for multi-object tracking. In *ICCV*, pages 3988–3998, 2019. 3
- [69] Linjie Yang, Yuchen Fan, and Ning Xu. Video instance segmentation. In *ICCV*, 2019. 1, 2, 3, 6, 7
- [70] Yibo Yang, Hongyang Li, Xia Li, Qijie Zhao, Jianlong Wu, and Zhouchen Lin. Sognet: Scene overlap graph network for panoptic segmentation. In *AAAI*, 2020. 3
- [71] Haobo Yuan, Xiangtai Li, Yibo Yang, Guangliang Cheng, Jing Zhang, Yunhai Tong, Lefei Zhang, and Dacheng Tao. Polyphonicformer: Unified query learning for depth-aware video panoptic segmentation. In *ECCV*, 2022. 3
- [72] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016. 7
- [73] Fangao Zeng, Bin Dong, Tiancai Wang, Cheng Chen, Xiangyu Zhang, and Yichen Wei. Motr: End-to-end multiple-object tracking with transformer. *arXiv preprint arXiv:2105.03247*, 2021. 1
- [74] Wenwei Zhang, Jiangmiao Pang, Kai Chen, and Chen Change Loy. K-net: Towards unified image segmentation. *CoRR*, abs/2106.14855, 2021. 1, 2, 3, 4, 12
- [75] Zheng Zhang, Dazhi Cheng, Xizhou Zhu, Stephen Lin, and Jifeng Dai. Integrated object detection and tracking with tracklet-conditioned detection. *arXiv preprint arXiv:1811.11167*, 2018. 3
- [76] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *CVPR*, 2017. 2, 7
- [77] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ADE20K dataset. *CVPR*, 2017. 12
- [78] Qianyu Zhou, Xiangtai Li, Lu He, Yibo Yang, Guangliang Cheng, Yunhai Tong, Lizhuang Ma, and Dacheng Tao. Transvod: End-to-end video object detection with spatial-temporal transformers, 2022. 3
- [79] Xingyi Zhou, Vladlen Koltun, and Philipp Krähenbühl. Tracking objects as points. *ECCV*, 2020. 3
- [80] Ji Zhu, Hua Yang, Nian Liu, Minyoung Kim, Wenjun Zhang, and Ming-Hsuan Yang. Online multi-object tracking with dual matching attention networks. In *ECCV*, pages 366–382, 2018. 3
- [81] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 7
- [82] Xizhou Zhu, Yuwen Xiong, Jifeng Dai, Lu Yuan, and Yichen Wei. Deep feature flow for video recognition. In *CVPR*, 2017. 1, 3

In this report, we provide the following information in addition to the main paper: more experimental details and more visualization results on Cityscapes-VPS [22] and KITTI-STEP [62]. All the model use AdamW for training.

## 6. More Experimental Details

### Detailed pretraining setting.

For COCO [33] dataset pretraining, all the models are trained following original K-Net settings [74]. We adopt the multi-scale training setting as previous work [5] by resizing the input images such that the shortest side is at least 480 and at most 800 pixels, while the longest size is at most 1333. We also apply random crop augmentations during training, where the train images are cropped with probability 0.5 to a random rectangular patch which is then resized again to 800-1333. All the models are trained for 36 epochs.

For Mapillary [40] dataset pretraining, we mainly follow the Panoptic-Deeplab setting [10]. We adopt the multi-scale training where the the scale ranges from 1.0 to 2.0 of origin image size, then we apply a random crop of  $1024 \times 2048$  patches. The horizontal flip is applied. The pretraining process takes 240 epochs. The Mapillary pretraining is for fair comparison, since the ViP-Deeplab [47] all use the Mapillary pretraining for better results. Note that, we *only* pre-train our largest models with Swin-base [34] as backbone for fair comparison with previous work.

### Training and inference on Cityscapes-VPS.

For Cityscapes-VPS training, we follow previous VP-SNet [22] that we randomly sample one frame from the nearest one frame as the reference frame. We adopt the multi-scale training where the scale ranges from 1.0 to 2.0 of origin images size then we apply a random crop of  $1024 \times 2048$  patches. For the Swin-base model, we apply a random crop of  $800 \times 1600$  patches to save memory and computation cost. The total training epoch is set to 8.

During the inference, the previous frame plays as the reference frame, the kernel information is directly propagated into the next frame in an online manner.

### Training and inference on KITTI-STEP.

For KITTI-STEP training, we follow previous Motion-Deeplab [62] that we randomly sample one frame from near *three* frames as the reference frame. We adopt the multi-scale training where the scale ranges from 1.0 to 2.0 of origin images size, then we apply a random crop of  $384 \times 1248$  patches. The total training epoch is set to 12. The inference procedure is the same as Cityscapes-VPS dataset. Following [62], we also use Cityscapes pretraining before training on STEP, which leads to about 3% STQ gain on the K-Net baseline.

### Training and inference on VSPW.

For VSPW dataset, we adopt the same setting on training K-Net [74] on ADE-datasets [77]. We randomly sample one frame from the nearest three frames as the reference

frame. The inference procedure is the same as Cityscapes-VPS dataset.

### Training and inference on VIPSeg.

For VIPSeg dataset, we adopt the same setting on of KITTI-STEP. In particular, we use COCO-pretrained K-Net following [37]. The entire training time is 12 epochs. We adopt the multi-scale training where the scale ranges from 1.0 to 2.0 of origin images size, then we apply a random crop of  $720 \times 720$  patches. The inference procedure is the same as Cityscapes-VPS dataset.

**Training and inference on YoutubeVIS-2019.** For YoutubeVIS-2019 dataset, we adopt the same training pipeline as VisTR [59] by sampling five frames to train Video K-Net jointly. We adopt COCO-pretrained model for initialization and train the video for only 12 epochs.

## 7. More Visualization Results

**More Visualization on Cityscapes-VPS.** In Fig. 7, we present more visual examples on Cityscapes VPS dataset. Compared with Ground Truth, our method can segment and track well for each pixel for various scale object inputs. We use the Swin-base model for visualization.

**More Visualization on KITTI-STEP.** In Fig. 8, we give more visual results on the KITTI-STEP validation set. On both clips, Video-KNet shows convincing results.

**Failure Cases Analysis.** In Fig. 9, we present three failure cases of our Video K-Net. The first case is the tracking failure case where the car moves fast in remote scene and then there is an ID switch. We believe adding motion cues will improve this case. The last two cases show the segmentation errors. The first is because the color of cars' window is similar to the ground. The second is caused by less training cases: bike on the truck, which makes network hard to predict.

**Border Impact.** Our work pushes the boundary of video segmentation algorithms through simplicity and effectiveness. Since most applications are the video input, this work could also ease and accelerate the model production in real-world applications, such as in autonomous driving. Due to the limited dataset size, we do not evaluate the robustness of the proposed method on corrupted video inputs or extremely long video inputs.

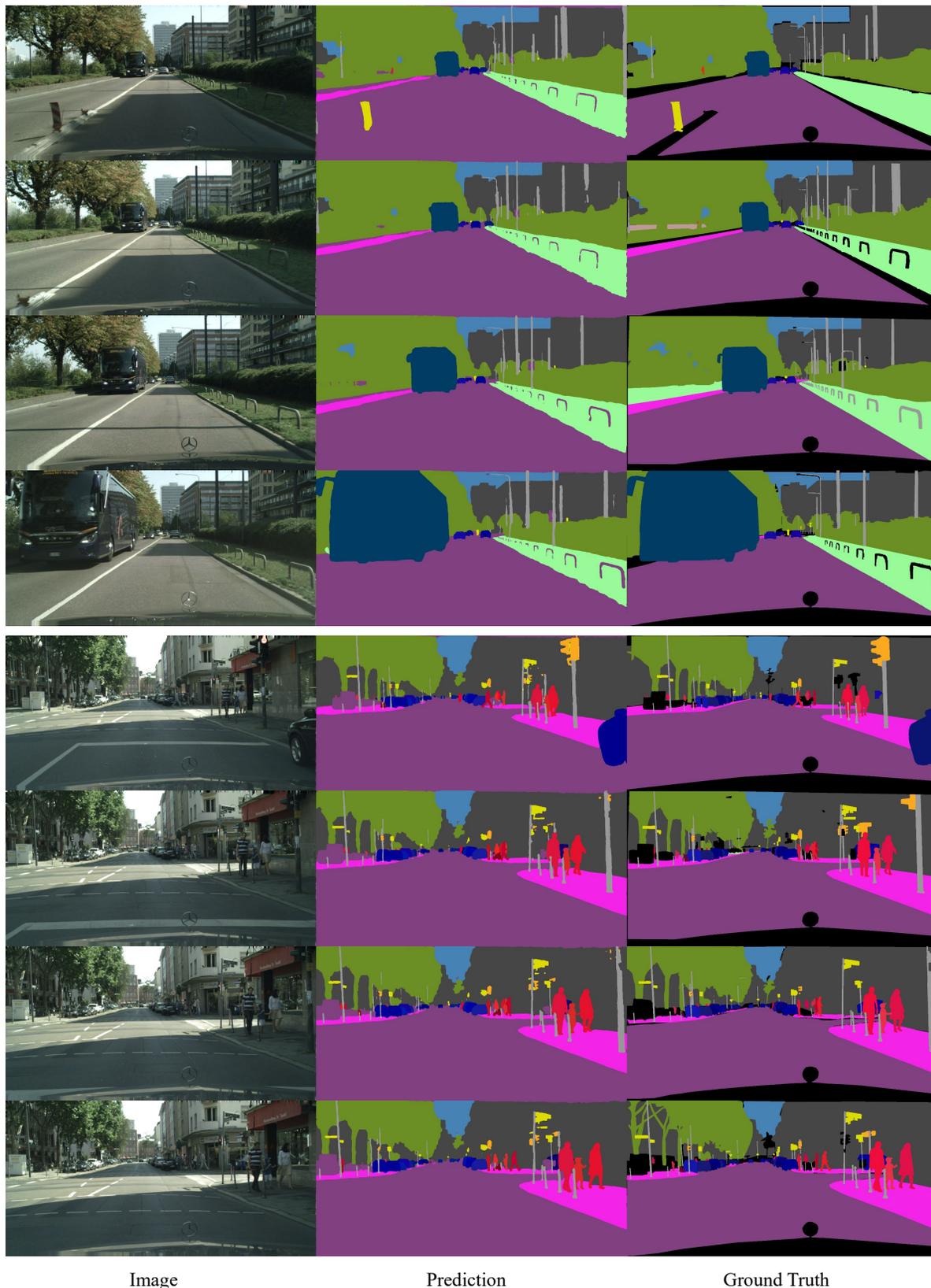


Figure 7. More visual examples of our Video K-Net on CityScapes-VPS dataset. The same instances are shown with the same color. Best view it on screen.

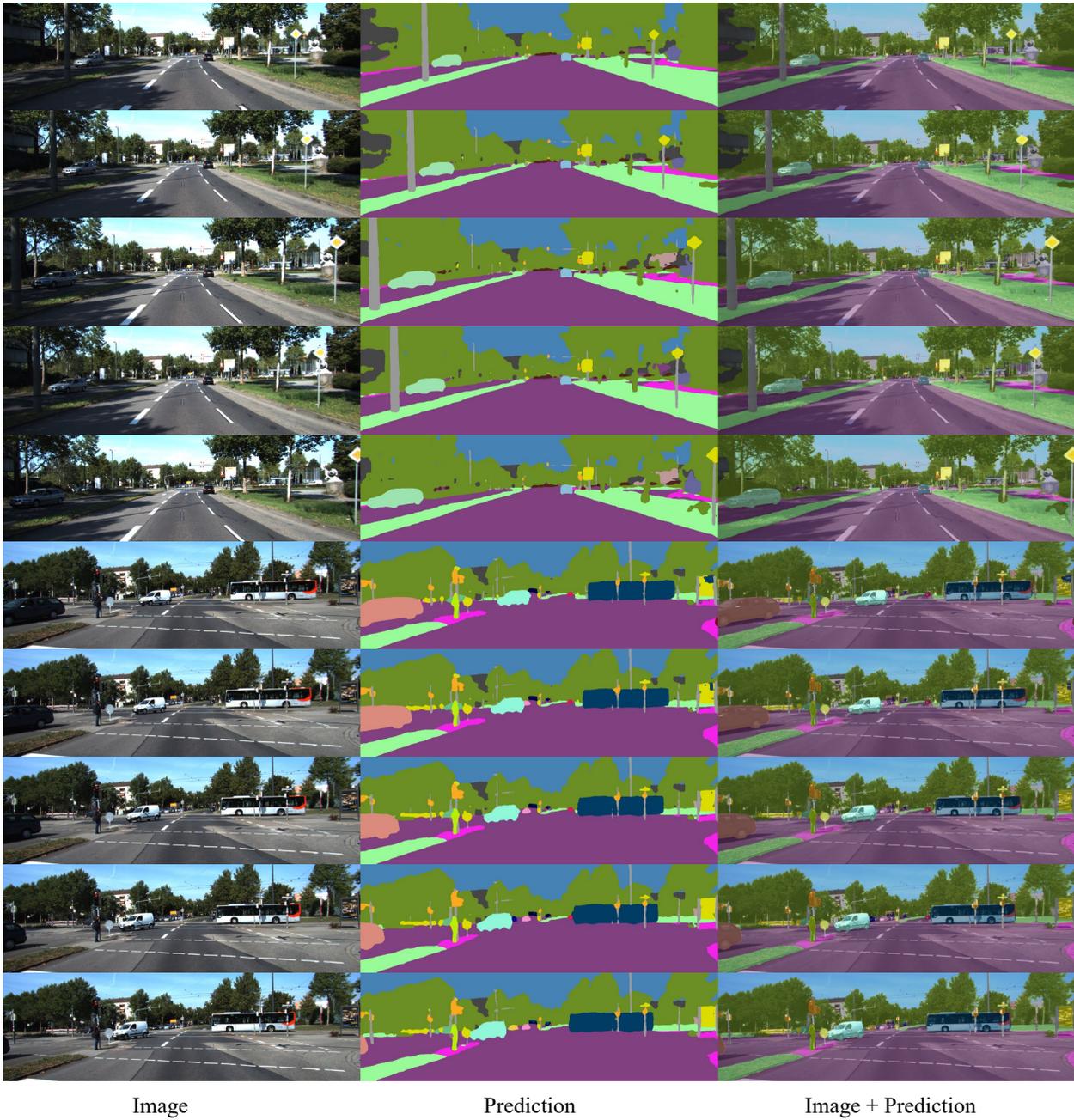


Figure 8. More visual examples of our Video K-Net on STEP dataset. The same instances are shown with the same color. Best view it on screen.

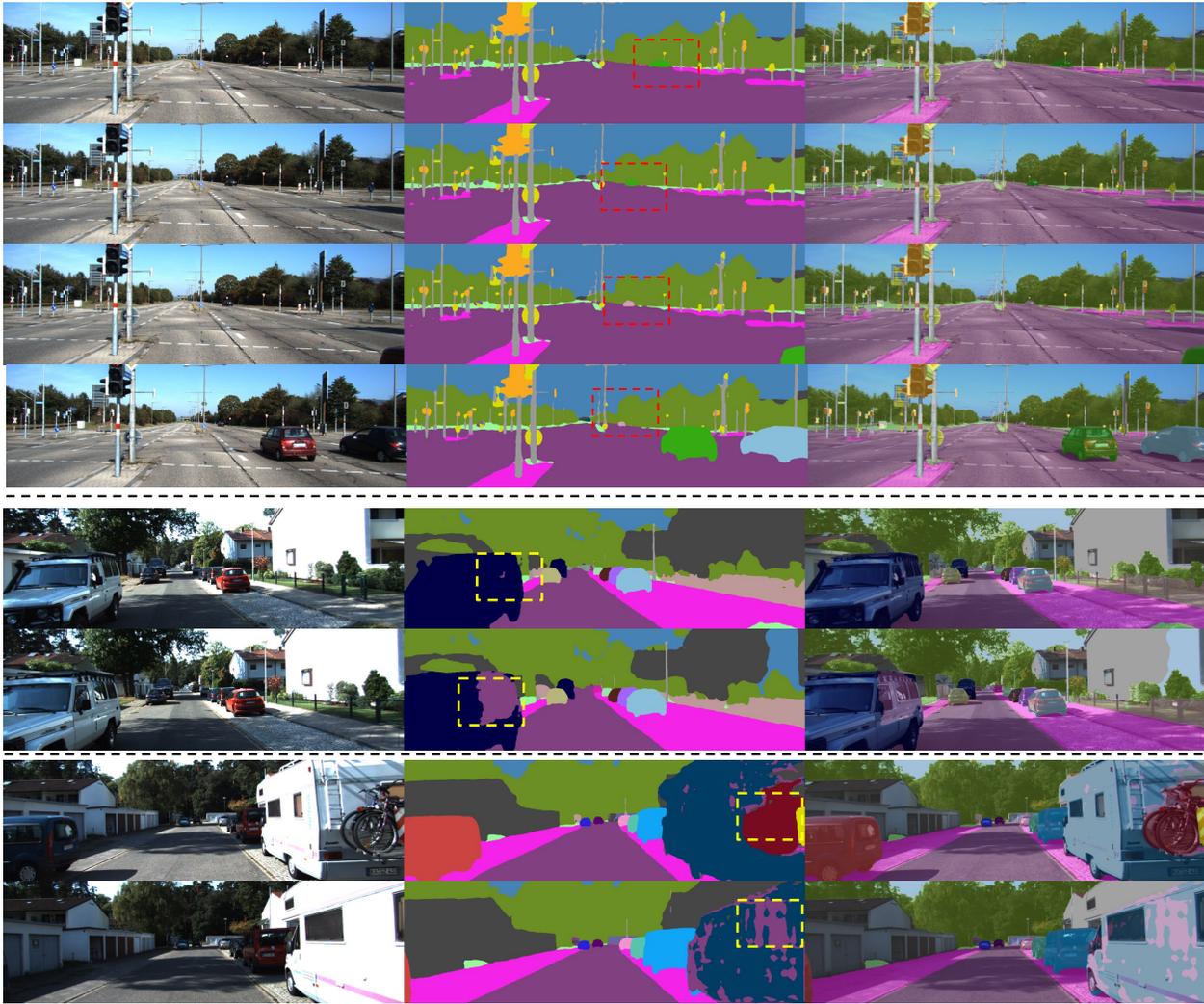


Figure 9. Failure cases of Video K-Net on STEP dataset. The same instances are shown with the same color. The red boxes show the tracking errors, while the yellow boxes show the segmentation errors. Best view it on screen.