

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

Xiangtai Li<sup>1\*</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> School of Artificial Intelligence, Key Laboratory of Machine Perception (MOE), Peking University

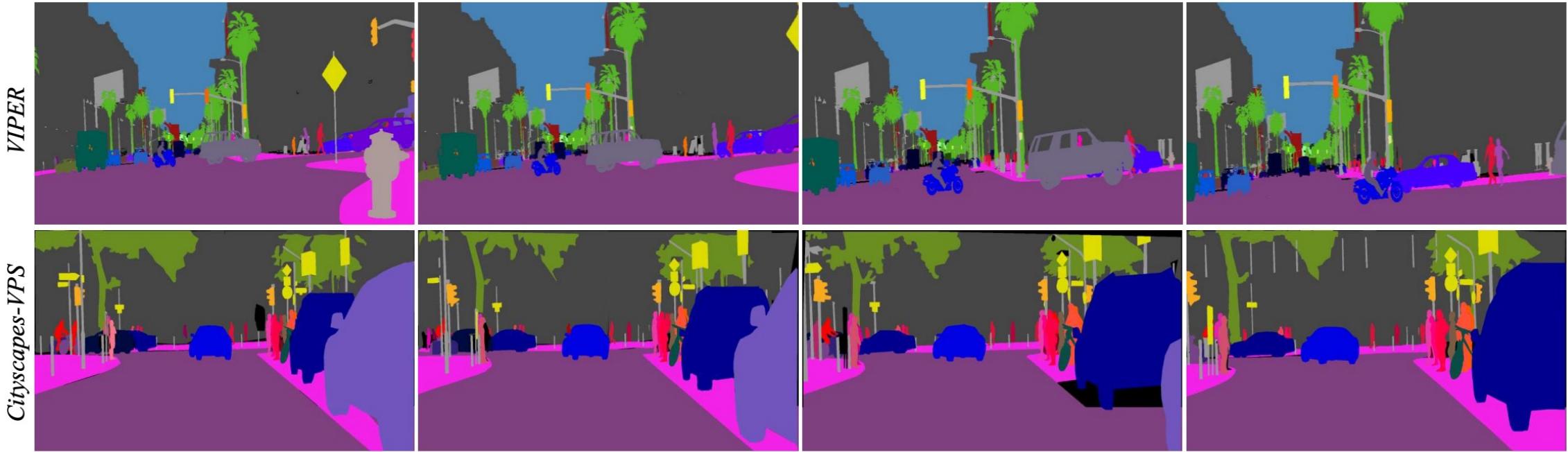
<sup>2</sup> S-Lab, Nanyang Technological University <sup>3</sup> CUHK-SenseTime Joint Lab, the Chinese University of Hong Kong

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

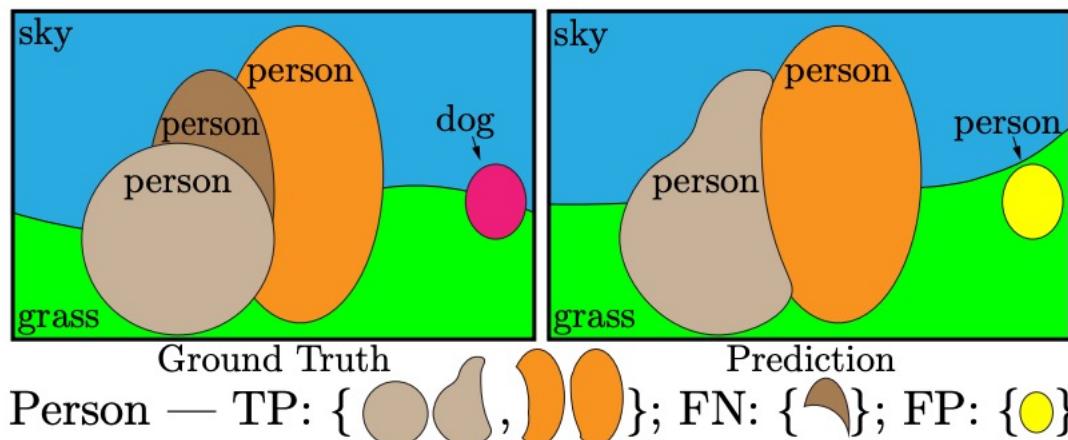
{lxt@pku, yhtong}@pku.edu.cn {wenwei001, ccloy}@ntu.edu.sg

pangjiangmiao@gmail.com {chenkai, chengguangliang}@sensetime.com

# First proposed : CVPR 2020 VPSnet



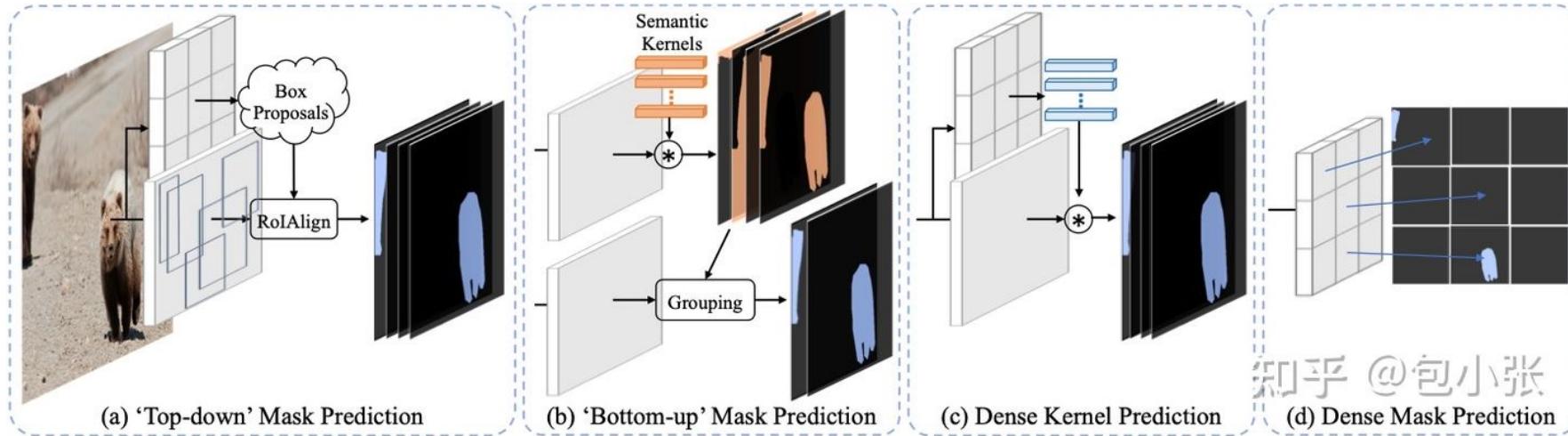
Metric : Video Panoptic Quality (VPQ)、 Segmentation and Tracking Quality (STQ)



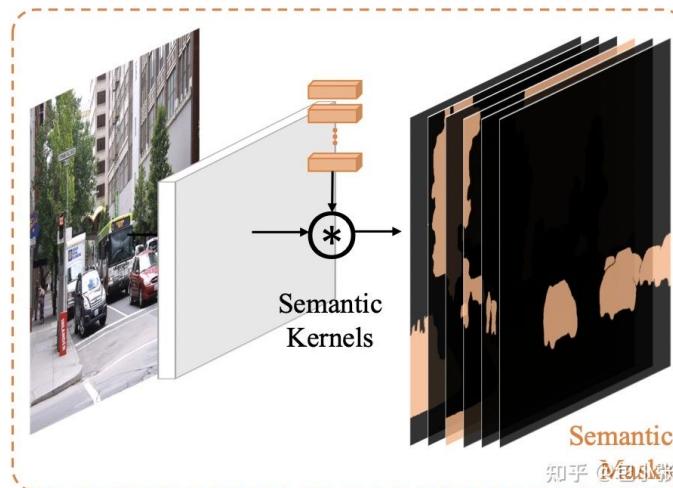
$$PQ = \frac{\sum_{(p,g) \in TP} \text{IoU}(p, g)}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}.$$

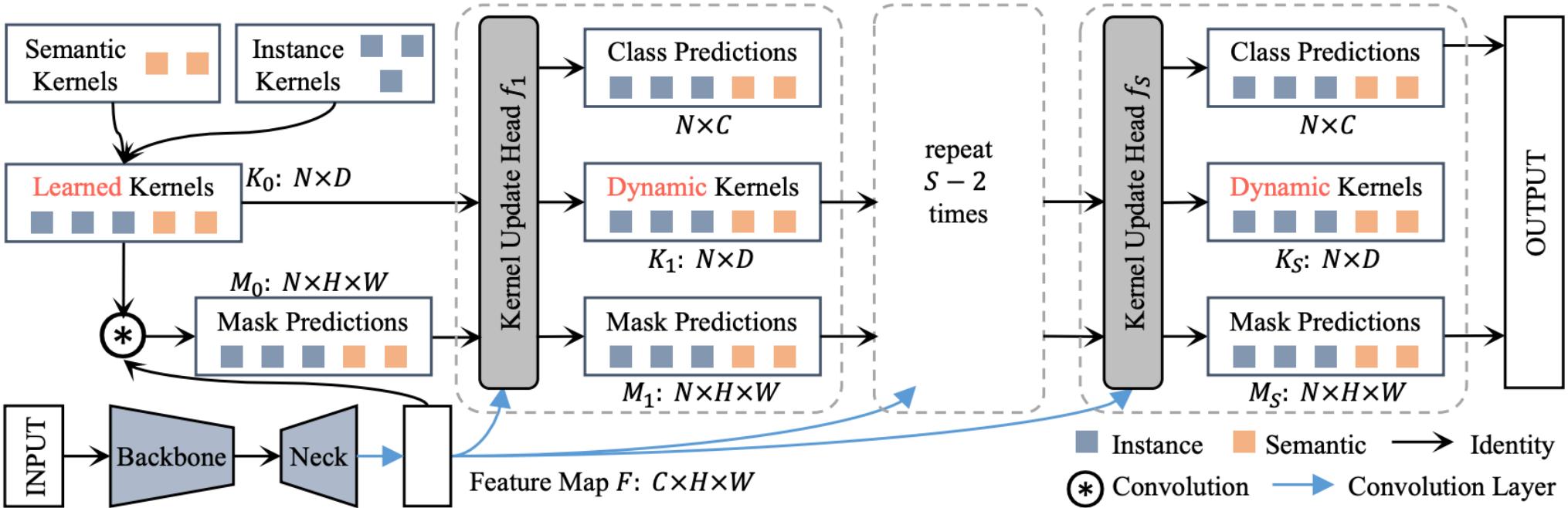
$$VPQ = \frac{1}{K} \sum_k \frac{1}{N_{classes}} \sum_c \frac{\sum_{(p,g) \in TP_c^k} IoU_{2D/3D}(p, g)}{|TP_c^k| + \frac{1}{2}|FP_c^k| + \frac{1}{2}|FN_c^k|}$$

## Different instance segmentation method :



## General Semantic segmentation method :





**Figure 3: K-Net for panoptic segmentation.** A set of learned kernels first performs convolution with the feature map  $F$  to predict masks  $M_0$ . Then the kernel update head takes the mask predictions  $M_0$ , learned kernels  $K_0$ , and feature map  $F$  as input and produce class predictions, group-aware (dynamic) kernels, and mask predictions. The produced mask prediction, dynamic kernels, and feature map  $F$  are sent to the next kernel update head. This process is performed iteratively to progressively refine the kernels and the mask predictions.

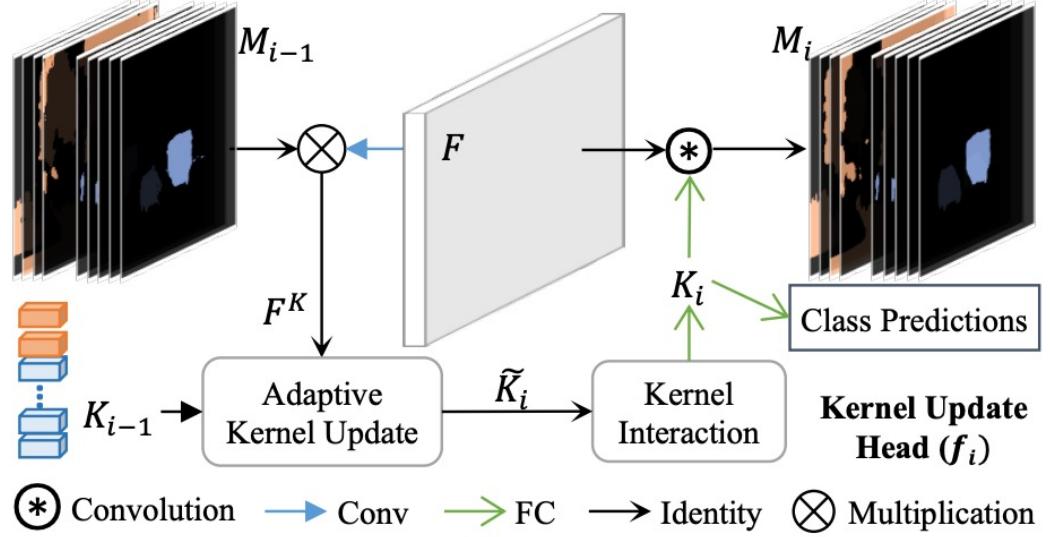


Figure 2: Kernel Update Head.

**Feature assembling:**  $F^K = \sum_u^H \sum_v^W M_{i-1}(u, v) \cdot F(u, v), F^K \in R^{B \times N \times C}$

**Kernel update:**  $F^G = \phi_1(F^K) \otimes \phi_2(K_{i-1}), F^G \in R^{B \times N \times C}$   $\phi_1$  and  $\phi_2$  are linear transformations

$$G^K = \sigma(\psi_1(F^G)), G^F = \sigma(\psi_2(F^G)),$$

$$\tilde{K} = G^F \otimes \psi_3(F^K) + G^K \otimes \psi_4(K_{i-1}),$$

where  $\psi_n, n = 1, \dots, 4$  are different fully connected (FC) layers followed by LayerNorm (LN) and  $\sigma$  is the Sigmoid function.  $\tilde{K}$  is then used in kernel interaction.

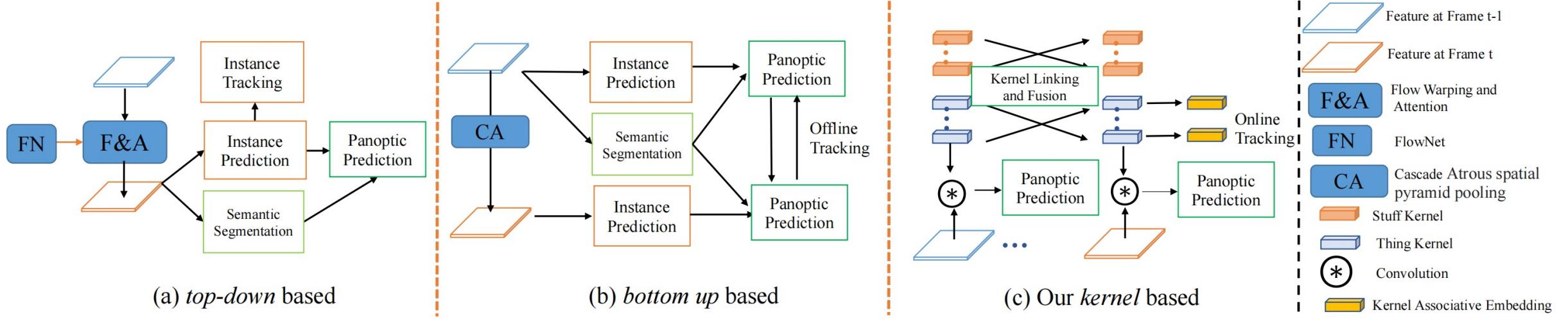


Figure 1. An illustration of previous *top-down* based VPS method (a), *bottom-up* based VPS method (b) and the proposed Video K-Net (c). Unlike previous approaches [22, 43] that perform panoptic segmentation and object tracking with independent modules, our method unifies panoptic segmentation and instance level tracking via kernels in a simpler framework.

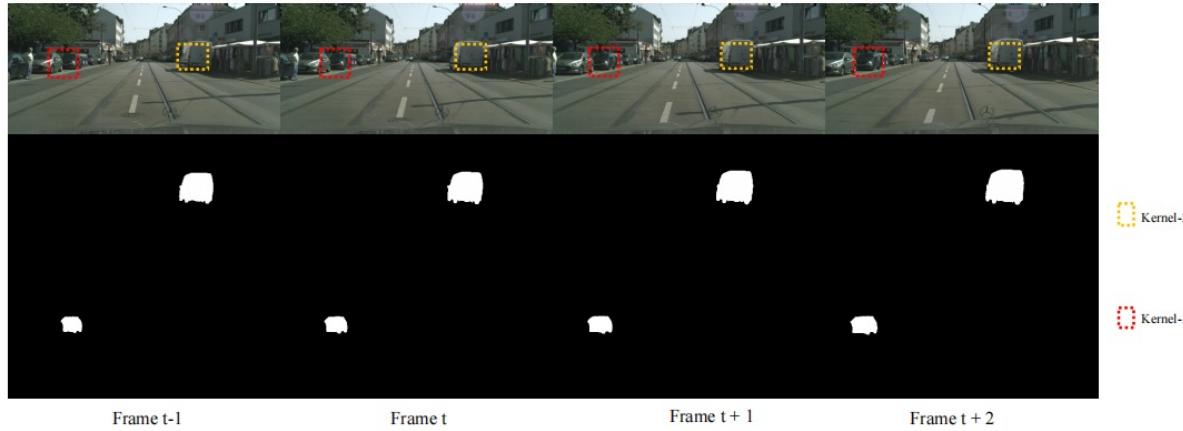


Figure 3. Toy experiment illustration. We use the K-Net directly on Cityscapes video datasets. We find that several instances are originated from **the same kernel** predictions (Red, Yellow boxes, **Kernel-2** and **Kernel-3**). This observation motivates us to use K-Net directly on video. Best view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

Table 1. Toy Experiment results on KITTI-STEP and Cityscape-VPS set with *STQ* and *VPQ* metrics. Unitrack [56] uses ResNet-50 as the appearance model.

<b>KITTI-STEP</b>	Backbone	STQ	AQ	SQ	-
K-Net	ResNet50	67.5	65.5	68.9	-
K-Net + Unitrack [56]	ResNet50	65.1	64.3	68.9	-
<b>Cityscapes-VPS</b>	Backbone	-	-	-	VPQ
K-Net	ResNet50	-	-	-	54.3
K-Net + Unitrack [56]	ResNet50	-	-	-	53.2

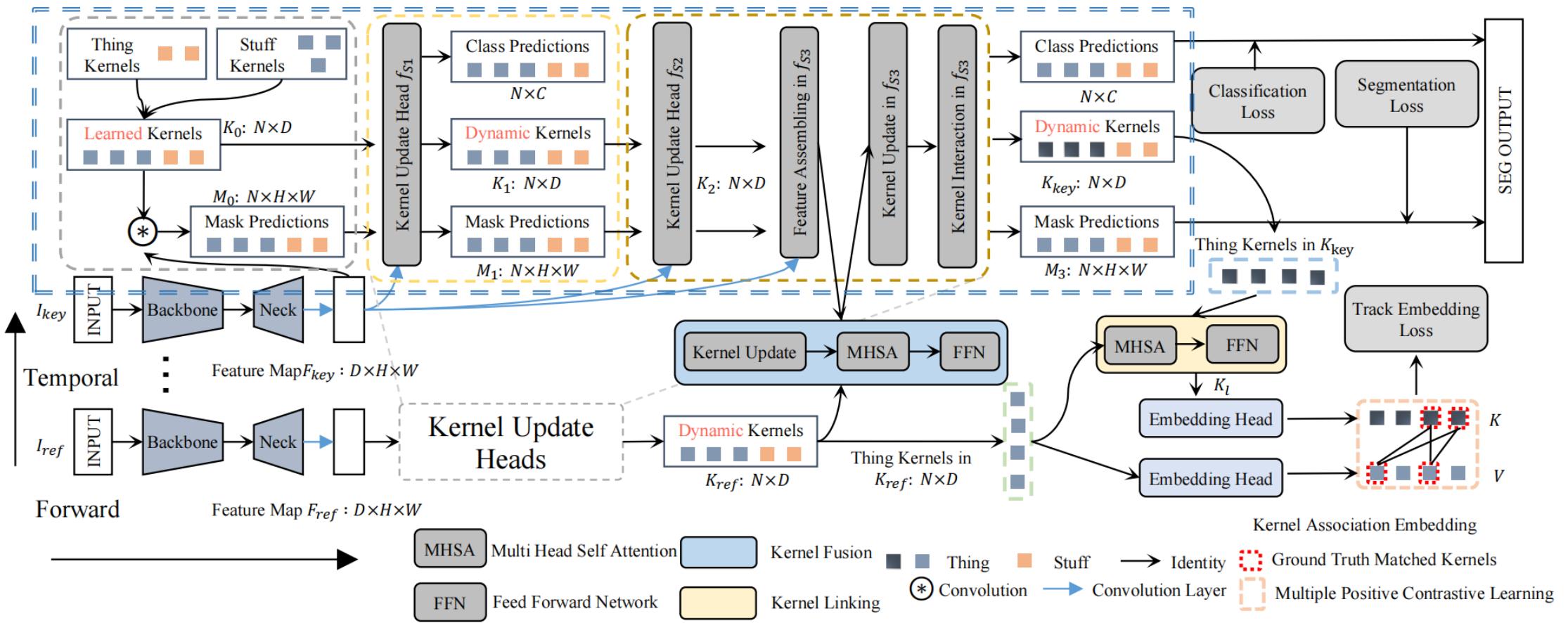
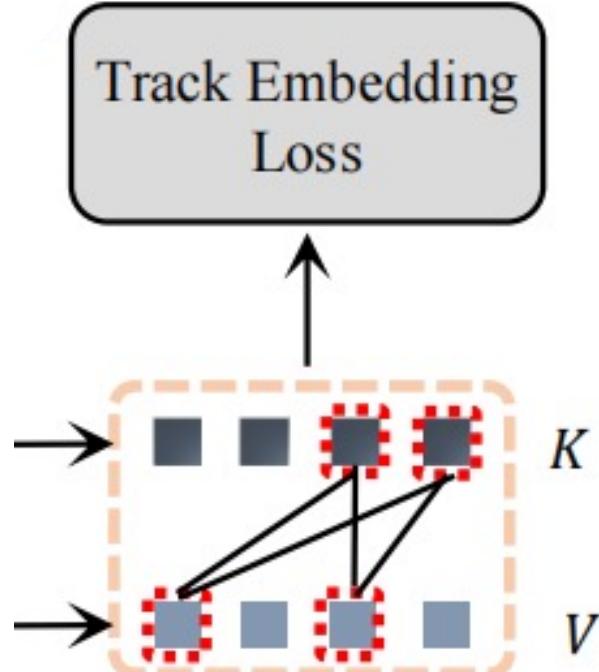


Figure 4. An illustration of our proposed Video K-Net. Our method is based on K-Net [70] (in blue dashed box), which is the top-left part of the figure. Video K-Net adds Kernel Fusion at the start phase of the last stage. The Kernel Linking is performed on the output of dynamic kernels. The Embedding Head is appended at the output of kernel linking and takes kernel outputs from both sampled frames.

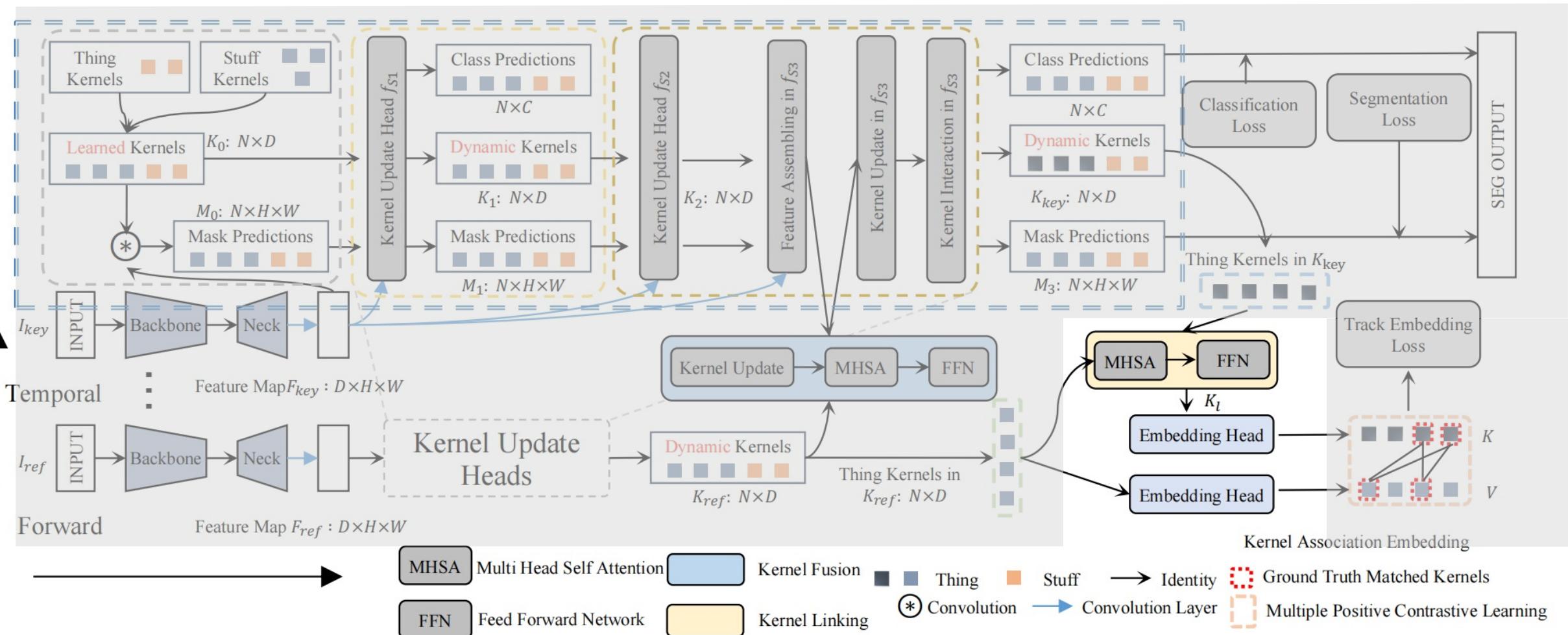
# Learning Kernel Association Embeddings



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

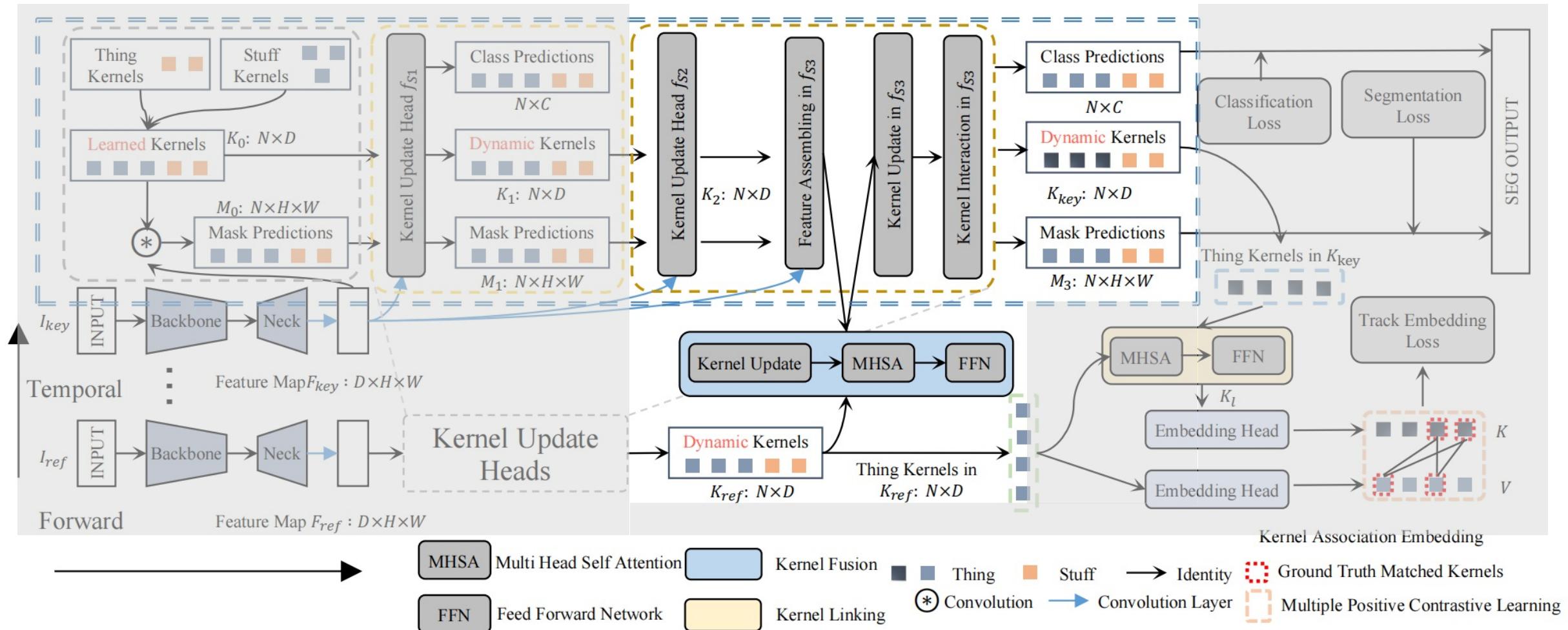
$$\mathcal{L}_{\text{aux}} = \left( \frac{\mathbf{v} \cdot \mathbf{k}}{\|\mathbf{v}\| \cdot \|\mathbf{k}\|} - c \right)^2,$$

# Learning to Link Kernels



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

# Learning to Fuse Kernels



# Inference

Bidirectional softmax:

$$\mathbf{f}(i, j) = \left[ \frac{\exp(\mathbf{n}_i \cdot \mathbf{m}_j)}{\sum_{k=0}^{M-1} \exp(\mathbf{n}_i \cdot \mathbf{m}_k)} + \frac{\exp(\mathbf{n}_i \cdot \mathbf{m}_j)}{\sum_{k=0}^{N-1} \exp(\mathbf{n}_k \cdot \mathbf{m}_j)} \right] / 2$$

# Experiment

KITTI-STEP validation and test set

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 [58]	ResNet50		0.58	0.51	0.67	0.40
VPSNet [22]	ResNet50	✓	0.56	0.52	0.61	0.43
Video K-Net	ResNet50		0.71	0.70	0.71	0.46
Video K-Net	Swin-base		0.74	0.72	0.73	0.53
Motion-Deeplab [58]	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	-

# Experiment

## Cityscapes-VPS validation set

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 [60]	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 [43]	WideResNet41 [68]	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 [43]	WideResNet41 [68]+RFP [42] + 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 [31]	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 [42]	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

## VSPW validation set

VPSW	Backbone	mIoU	$mVC_8$	$mVC_{16}$
DeepLabv3+ [8]	ResNet101	35.7	83.5	78.4
PSPNet+ [72]	ResNet101	36.5	84.4	79.8
TCB(PSPNet) [34]	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

# Experiment

YTVIS-2019 validation set

Method	backbone	AP	AP <sub>50</sub>	AP <sub>75</sub>	AR <sub>1</sub>	AR <sub>10</sub>
FEELVOS [50]	ResNet50	26.9	42.0	29.7	29.9	33.4
MaskTrack R-CNN [65]	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 [55]	ResNet50	34.4	55.7	36.5	33.5	38.9
VisTR [55]	ResNet101	35.3	57.0	36.2	34.3	40.4
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

# Ablation study

(a) Ablation Study on Each Components.

baseline	KAE	KL	KF	STQ	AQ	SQ
K-Net				67.5	65.5	68.9
✓				69.3	69.0	69.8
✓	✓			70.2	71.2	69.7
✓	✓	✓		70.9	70.8	71.2

(b) Needs of Appearance Embeddings

Method	AQ	STQ
RoI-Align [37]	68.8	69.1
Mask-Emb [60]	67.3	68.1
Ours	70.8	70.9
Ours + Mask-Emb [60]	70.3	70.8

(c) Effect of sampling in association.

Method	STQ	AQ	SQ
K-Net	67.5	65.5	68.9
GT-based (ours)	69.3	69.0	69.8
sampling in [37]	63.1	62.1	64.3

(d) Ablation Study on Linking and Fusing Stage.

Stage	STQ	AQ	SQ
3	70.9	70.8	71.2
2	68.5	68.2	69.3
1	66.9	63.4	67.3

(e) Ablation Study on Training Settings

Settings	STQ	AQ	SQ
joint training	70.9	70.8	71.2
only train the key frame	70.1	70.1	69.8

(f) Ablation Study on Kernel Fusing

Settings	STQ	AQ	SQ
K-Net	67.5	65.5	68.9
w Update	70.9	70.8	71.2
w/o Update	67.1	66.2	68.3

