

# Video Scene Parsing with Predictive Feature Learning

Xiaojie Jin<sup>1</sup> Xin Li<sup>2</sup> Huaxin Xiao<sup>2</sup> Xiaohui Shen<sup>3</sup> Zhe Lin<sup>3</sup> Jimei Yang<sup>3</sup>  
Yunpeng Chen<sup>2</sup> Jian Dong<sup>4</sup> Luoqi Liu<sup>4</sup> Zequn Jie<sup>2</sup> Jiashi Feng<sup>2</sup> Shuicheng Yan<sup>4,2</sup>

<sup>1</sup>NUS Graduate School for Integrative Science and Engineering, NUS

<sup>2</sup>Department of ECE, NUS

<sup>3</sup>Adobe Research

<sup>4</sup>360 AI Institute

## Abstract

*In this work, we address the challenging video scene parsing problem by developing effective representation learning methods given limited parsing annotations. In particular, we contribute two novel methods that constitute a unified parsing framework. (1) **Predictive feature learning** from nearly unlimited unlabeled video data. Different from existing methods learning features from single frame parsing, we learn spatiotemporal discriminative features by enforcing a parsing network to predict future frames and their parsing maps (if available) given only historical frames. In this way, the network can effectively learn to capture video dynamics and temporal context, which are critical clues for video scene parsing, without requiring extra manual annotations. (2) **Prediction steering parsing** architecture that effectively adapts the learned spatiotemporal features to scene parsing tasks and provides strong guidance for any off-the-shelf parsing model to achieve better video scene parsing performance. Extensive experiments over two challenging datasets, Cityscapes and Camvid, have demonstrated the effectiveness of our methods by showing significant improvement over well-established baselines.*

## 1. Introduction

Video scene parsing (VSP) aims to predict per-pixel semantic labels for every frame in scene videos recorded in unconstrained environments. It has drawn increasing attention as it benefits many important applications like drones navigation, autonomous driving and virtual reality.

In recent years, remarkable success has been made by deep convolutional neural network (CNN) models in image parsing tasks [3, 5, 21, 28, 29, 34, 43, 44]. Some of those CNN models are thus proposed to be used for parsing scene videos frame by frame. However, as illustrated in Figure 1, naively applying those methods suffers from noisy and inconsistent labeling results across frames, since the important temporal context cues are ignored. For example, in the second row of Figure 1, the top-left region of

class *building* in the frame  $T+4$  is incorrectly classified as *car*, which is temporally inconsistent with the parsing result of its preceding frames. Besides, for current data-hungry CNN models, finely annotated video data are rather limited as collecting pixel-level annotation for long videos is very labor-intensive. Even in the very recent scene parsing dataset Cityscapes [4], there are only 2,975 finely annotated training samples vs. overall 180K video frames. Deep CNN models are prone to over-fitting the small training data and thus generalize badly in real applications.

To tackle these two problems, we propose a novel Parsing with prEdictive feAtuRe Learning (**PEARL**) approach which is both annotation-efficient and effective for VSP. By enforcing them to predict future frames based on historical ones, our approach guides CNNs to learn powerful spatiotemporal features that implicitly capture video dynamics as well as high-level context like structures and motions of objects. Attractively, such a learning process is nearly annotation-free as it can be performed using any unlabeled videos. After this, our approach further adaptively integrates the obtained temporal-aware CNN to steer any image scene parsing models to learn more spatial-temporally discriminative frame representations and thus enhance video scene parsing performance substantially.

Concretely, there are two novel components in our proposed approach: *predictive feature learning* and *prediction steering parsing*. As shown in Figure 1, given frames  $T$  to  $T+3$ , predictive feature learning aims to learn discriminative spatiotemporal features by enforcing a CNN model to predict the future frame  $T+4$  as well as the parsing map of  $T+4$  if available. Such predictive learning enables the CNN model to learn features capturing the cross-frame object structures, motions and other temporal cues, and provide better video parsing results, as demonstrated in the third row of Figure 1. To further adapt the obtained CNN along with its learned features to the parsing task, our approach introduces a prediction steering parsing architecture. Within this architecture, the temporal-aware CNN (trained by frame prediction) is utilized to guide an image-parsing CNN model to parse the current frame by providing tem-

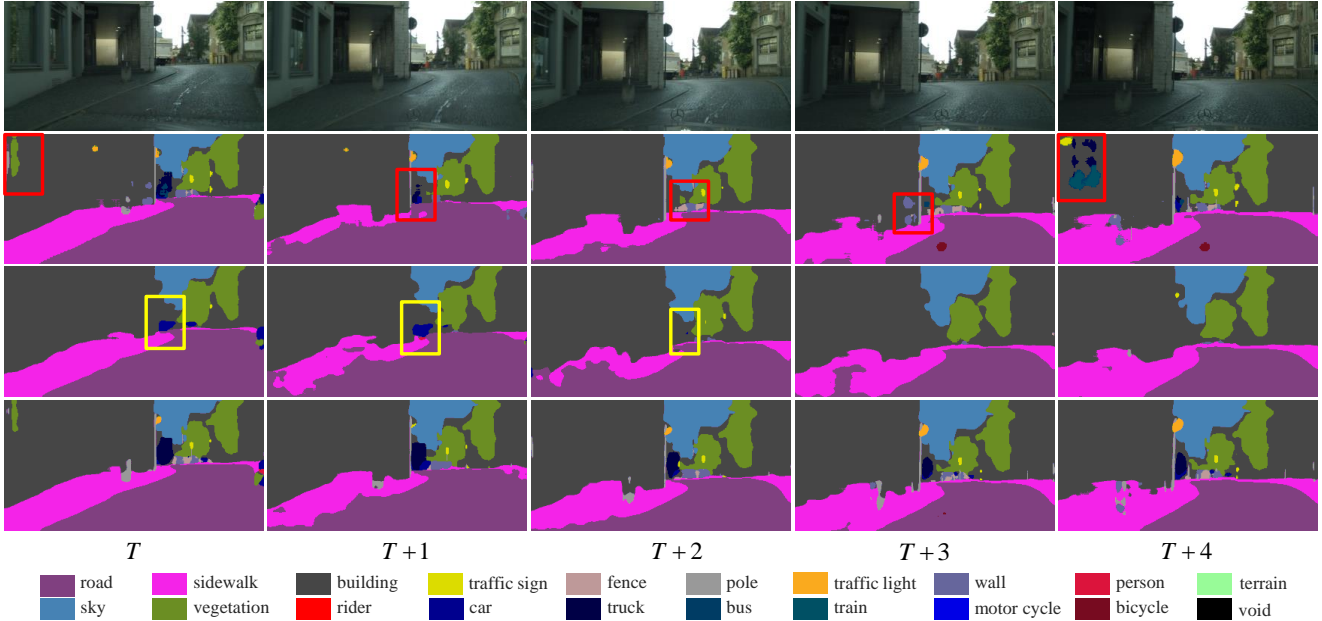


Figure 1: Illustration on core ideas of predictive feature learning for video scene parsing. **Top**: a five-frame sequence from the Cityscapes video dataset. **Second row**: frame parsing produced by the conventional VGG16 model [3]. Since it is unable to model temporal context, severe noise and inconsistency across frames can be observed within the red boxes. **Third row**: results from predictive feature learning and parsing. The regions showing inconsistency in the second row are classified consistently across frames by predictive feature learning. Besides, the motion trajectories of cars are correctly captured (see yellow boxes). **Bottom**: parsing maps produced by PEARL with better accuracy and temporal consistency. PEARL combines the advantages of traditional image parsing model (the second row) and predictive parsing model (the third row). Best viewed in color and zoomed pdf.

poral cues implicitly. The two parsing networks are jointly trained end-to-end and provide parsing results with strong cross-frame consistency and richer local details (as shown in the bottom row of Figure 1).

We conduct extensive experiments over two challenging datasets and compare our approach with strong baselines, *i.e.*, state-of-the-art VGG16 [33] and Res101 [10] based parsing models. Our approach achieves the best results on both datasets. In the comparative study, we demonstrate its superiority to other methods that model temporal context, *e.g.*, using optical flow [27].

To summarize, we make the following contributions to video scene parsing:

- A novel predictive feature learning method is proposed to learn the spatiotemporal features and high-level context from a large amount of unlabeled video data.
- An effective prediction steering parsing architecture is presented which utilizes the temporal consistent features to produce temporally smooth and structure preserving parsing maps.
- Our approach achieves state-of-the-art performance on

two challenging datasets, *i.e.*, Cityscapes and Camvid.

## 2. Related Work

Recent image scene parsing progress is mostly stimulated by various new CNN architectures, including the fully convolutional architecture (FCN) with multi-scale or larger receptive fields [5, 21, 34] and the combination of CNN with graphical models [3, 28, 43, 44, 29]. There are also some recurrent neural networks based models [12, 17, 26, 32, 39]. However, without incorporating the temporal information when directly applying them to every frame of videos, the parsing results commonly lack cross-frame consistency and the quality is not good.

To utilize temporal consistency across frames, the motion and structure features in 3D data are employed by [6, 35, 42]. In addition, [9, 14, 16, 23] use CRF to model spatiotemporal context. However, those methods suffer from high computation cost as they need to perform expensive inference of CRF. Some other methods employ optical flow to capture the temporal consistency explored in [11, 30]. Different from above works that heavily depend on labeled data for supervised learning, our proposed approach takes ad-

vantage of both the labeled and unlabeled video sequences to learn richer temporal context information.

Generative adversarial networks were firstly introduced in [8] to generate natural images from random noises, and have been widely used in many fields including image synthesis [8], frame prediction [22, 24] and semantic inpainting [25]. Our approach also uses adversarial loss to learn more robust spatiotemporal features in frame predictions. Our approach is more related to [22, 24] by using adversarial training for frame prediction. However, different from [22, 24], PEARL tackles the VSP problem by utilizing spatiotemporal features learned in frame prediction.

### 3. Predictive Feature Learning for VSP

#### 3.1. Motivation and Notations

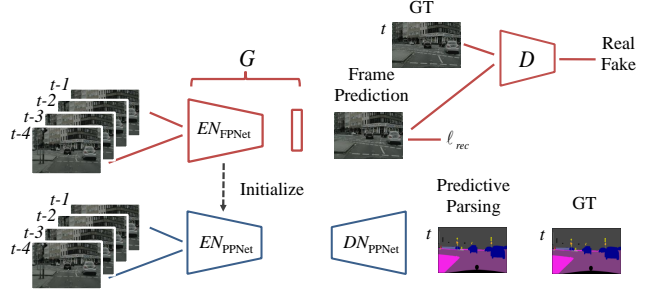
The proposed approach is motivated by two challenging problems of video scene parsing: first, how to leverage the temporal context information to enforce cross-frame smoothness and produce structure preserving parsing results; second, how to build effective parsing models even in presence of insufficient training data.

Our approach solves these two problems through a novel predictive feature learning strategy. We consider the partially-labeled video collection used for predictive feature learning, denoted as  $\{\mathcal{X}, \mathcal{Y}\}$ , where  $\mathcal{X} = \{X_1, \dots, X_N\}$  denotes the raw video frames and  $\mathcal{Y} = \{Y_{r_1}, \dots, Y_{r_M}\}$  denotes the provided dense annotations for a subset of  $\mathcal{X}$ . Here  $M \ll N$  as collecting large-scale annotation is not easy.  $Y_{r_j}(p, q) \in \{1, \dots, C\}$  denotes the ground truth category at location  $(p, q)$  in which  $C$  is the number of semantic categories. Correspondingly, let  $\hat{\mathcal{X}} = \{\hat{X}_i, i = 1, \dots, N\}$  and  $\hat{\mathcal{Y}} = \{\hat{Y}_i, i = 1, \dots, N\}$  denote predicted frames and predicted parsing maps, respectively. We use  $P_i^s = \{X_{i-k}\}_{k=1}^s$  to denote the  $s$  preceding frames ahead of  $X_i$ . For the first several frames in a video, we define their preceding set as  $X_{i-k} = X_1$  if  $i \leq k$ .

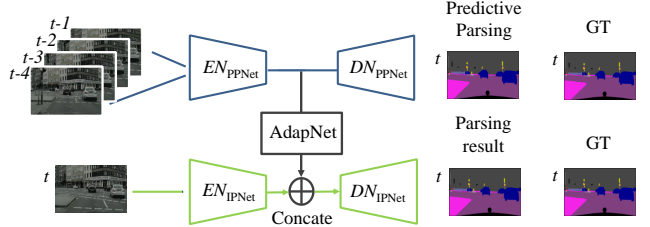
Video scene parsing can be formulated as seeking a parsing function  $\mathcal{F}$  that maps any frame sequence to the parsing map of the most recent frame:

$$\hat{Y}_i = \mathcal{F}(X_i, P_i^s). \quad (1)$$

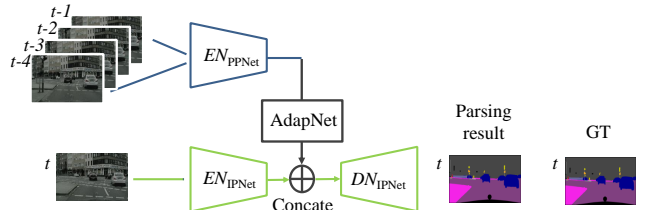
The above definition reveals the difference between static scene image parsing and video scene parsing — the video scene parsing model  $\mathcal{F}$  has access to historical/temporal information for parsing the current target. We also want to highlight an important difference between our problem setting and some existing works [9, 23]: instead of using the whole video (including both past and future frames *w.r.t.*  $X_i$ ) to parse the frame  $X_i$ , we aim to perform parsing based on causal inference (or online inference) where only past frames are observable. This setting aligns better with



(a) The framework of predictive feature learning in PEARL



(b) The architecture of prediction steering parsing network in PEARL



(c) A variant of PEARL

Figure 2: (a) The framework of predictive feature learning. The terms *EN* and *DN* represent encoder/decoder network respectively. First, FPNet (highlighted in red) learns to *predict frame t* given only frames  $t-4$  to  $t-1$  via its generator  $G$  and discriminator  $D$ . Second, PPNet (highlighted in blue) performs *predictive parsing* on frame  $t$  without seeing it, based on its  $EN_{PPNet}$  which is initialized by  $EN_{FPNet}$  and connected to  $DN_{PPNet}$ . (b) The architecture of prediction steering parsing network (PSPNet). Given a single input frame  $t$ , the image parsing network (IPNet) (highlighted in green) parses it by integrating the learned features from  $EN_{PPNet}$  plus a shallow AdapNet. PPNet and IPNet are jointly trained. (c) An important variant of PEARL to verify effectiveness of features learned by predictive feature learning. Only the  $EN_{FPNet}$  or  $EN_{PPNet}$  is concatenated with  $EN_{IPNet}$  through AdapNet. The weights of  $EN_{FPNet}/EN_{PPNet}$  are fixed during training. Best viewed in color.

real-time applications like autonomous driving where future frames cannot be seen in advance.

#### 3.2. Predictive Feature Learning

Predictive feature learning aims to learn spatiotemporal features capturing high-level context like object mo-

tions and structures from two consecutive predictive learning tasks, *i.e.*, the frame prediction and the predictive parsing. Figure 2 gives an overview. In the first task, we train an FPNet for future frame prediction given several past frames by a new generative adversarial learning framework developed upon GANs [8, 24]. Utilizing a large amount of unlabeled video sequence data, the FPNet is capable of learning rich spatiotemporal features to model the variability of content and dynamics of videos. Then we further augment FPNet through Predictive Parsing, *i.e.* predicting parsing results of the future frame given previous frames. This adapts FPNet to another model (called PPNet) suitable for parsing.

**Frame Prediction** The architecture of FPNet is illustrated in Figure 2a. It consists of two components, *i.e.*, the generator (denoted as  $G$ ) which generates future frame  $\hat{X}_i = G(P_i^s)$  based on its preceding frames  $P_i^s$ , and the discriminator (denoted as  $D$ ) which plays against  $G$  by trying to identify predicted frame  $\hat{X}_i$  and the real one  $X_i$ . There is an  $Encoder_{FPNet}$  ( $EN_{FPNet}$  in Figure 2a) which maps the input video sequence to spatiotemporal features and a followed output layer to produce the RGB values of the predicted frame using the learned features. Note that  $Encoder_{FPNet}$  can choose any deep networks with various architectures, *e.g.* VGG16 and Res101. We adapt them to be compatible with video inputs by using group convolution [13] for the first convolutional layer, where the group number is equal to the number of input past frames.

FPNet alternatively trains  $D$  and  $G$  for predicting frames with progressively improved quality. Denote learnable parameters of  $D$  and  $G$  as  $W_D$  and  $W_G$  respectively. The objective for training  $D$  is to minimize the following binary cross-entropy loss where  $G$  is fixed:

$$\min_{W_D} \ell_D \triangleq -\log(1 - D(G(P_i^s))) - \log D(X_i). \quad (2)$$

Minimizing the above loss gives  $D$  a stronger discriminative ability to distinguish the predicted frames  $G(P_i^s)$  from real ones  $X_i$ , enforcing  $G$  to predict future frames with higher quality. Towards this target,  $G$  learns to predict future frames more like real ones through

$$\min_{W_G} \ell_G = \ell_{rec} + \lambda_{adv} \ell_{adv}, \quad (3)$$

where

$$\ell_{rec} = \|X_j - \hat{X}_j\|_2, \text{ and } \ell_{adv} = -\log D(G(\hat{P}_i^s)).$$

Minimizing the combination of reconstruction loss and adversarial loss supervises  $G$  to predict the frame that looks both similar to its corresponding real frame and sufficiently authentic to fool the strong competitor  $D$ . Our proposed frame prediction model is substantially different from vanilla GAN and more tailored for VSP problems. The

key difference lies in the generator  $G$  that takes past frame sequence  $P_i^s$  as input to predict the future frame, instead of crafting new samples completely from random noise as vanilla GANs. Therefore, the future frames coming from such “temporally conditioned” FPNet should present certain temporal consistency with past frames. On the other hand, FPNet can learn representations containing implicit temporal cues desired for solving VSP problems.

As illustrated by Figure 3, FPNet produces real-looking frame predictions by learning both the content and dynamics in videos. By comparing with the ground truth frame, the prediction frame resembles both the structures of objects/stuff like building/vegetation and the motion trajectories of cars, demonstrating that FPNet learns robust and generalized spatiotemporal features from video data.

In our experiments, we use a GoogLeNet [36] as  $D$  and we try both Res101 and VGGNet for  $G$ . More details are given in Section 4.1.

**Predictive Parsing** The features learned by FPNet so far are trained for video frame generation. To adapt these spatiotemporal features to VSP problems, FPNet performs the second predictive learning task, *i.e.*, predicting the parsing map of one frame  $X_i$  given only its preceding frames  $P_i^s$  (without seeing the frame to parse). This predictive parsing task is very challenging as we do not have any information of the current video frame.

Also, directly training FPNet for this predictive parsing task from scratch will not succeed. There are not enough data with annotations for training a good model free of overfitting. Thus, training FPNet for frame prediction at first gives a good starting model for accomplishing the second task. In this perspective, frame prediction training is important for both spatiotemporal feature learning and feature adaption.

Details on how FPNet performs the predictive parsing task are given in Figure 2a. For predicting parsing maps, we modify the architecture of FPNet as follows. We remove the  $D$  component from FPNet as well as the output layer in  $G$ . Then we add a deconvNet (*i.e.*  $DN_{PPNet}$  in the figure) on top of modified  $G$ , which produces the parsing map sharing the same size with frames. We call this new FPNet as PPNet, short for Predictive Parsing Net.

More details about the structure of PPNet are given in Section 4.1. Based on the notations given in Section 3.1, the objective function for training PPNet is defined as

$$\ell_{PPNet} = - \sum_{(p,q) \in X_{r_j}} \omega_{Y_{r_j}(p,q)} h_{Y_{r_j}(p,q)}(W_{PPNet}, P_{r_j}^s), \quad (4)$$

where  $h_{Y_{r_j}(p,q)}$  denotes the per-pixel logarithmic probability predicted by PPNet for the category label  $Y_{r_j}(p,q)$ . We introduce the weight vector  $\omega$  to balance scene classes with



Figure 3: Example frame predictions of FPNet on Cityscape val set. **Top:** ground truth video sequence. **Bottom:** frame prediction of FPNet. FPNet produces visually similar frames with the ground truth, demonstrating that FPNet learns robust spatiotemporal features to model the structures of objects (*building*) and stuff (*vegetation*), and motion information of moving objects (*cars*), both of which are critical for VSP problems.

different frequency in the training set. We will further discuss its role in the experiments.

Examples of predicted parsing maps from PPNet are shown in Figure 1 (the third row). Compared with parsing results from a traditional CNN parsing model on the single frame, the parsing maps of PPNet present two distinct properties. First, the parsing maps are temporally smooth which are reflected in the temporally consistent parsing results of regions like building where CNN models produce noisy and inconsistent parsing results. This demonstrates the PPNet indeed learns the temporally continuous dynamics from the video data. Secondly, PPNet tends to miss objects of small sizes, *e.g.*, transport signs and poles. One reason is the inevitable blurry prediction [24] since the high frequency spectrum is prone to being smoothed. This problem can be mitigated by parsing at multiple scales [2] and we will investigate in the future.

In contrast, the conventional image parsing network relies on locally discriminative features such that it can capture small objects. However, due to lacking temporal information, its produced parsing maps are noisy and lack temporal consistency with past frames. Above observations motivate us to combine the strengths of PPNet and the CNN-based image parsing model to improve the overall VSP performance. Therefore, we develop the following prediction steering parsing architecture.

### 3.3. Prediction Steering Parsing

To take advantage of the temporally consistent spatiotemporal features learned by PPNet, we propose a novel architecture to integrate PPNet and the traditional image parsing network (short for IPNet) into a unified framework, called PSPNet, short for Prediction Steering Parsing Net.

As illustrated in Figure 2b, PSPNet has two interconnected branches: one is the PPNet for predictive feature learning and the other is IPNet for frame-by-frame parsing. Similar to FPNet, the IPNet can also be chosen freely from any existed image parsing networks, *e.g.* FCN [21] and DeepLab [2]. As a high-level description, IPNet consists of two components, a feature encoder  $Encoder_{IPNet}$  ( $EN_{IPNet}$ ) which transforms the input frame to dense pixel features

and a deconvNet ( $DN_{IPNet}$ ) that produces per-pixel parsing map. Through an AdapNet (a shallow CNN), PPNet communicates its features to IPNet and steers the overall parsing process. In this way, the integrated features within IPNet gain two complementary properties, *i.e.*, descriptiveness for the temporal context and discriminability for different pixels within a single frame. Therefore the overall PSPNet model is capable of producing more accurate video scene parsing results than both PPNet and IPNet. Formally, the objective function of training PSPNet end-to-end is defined as

$$L_{PSPNet} = - \sum_{(p,q) \in X_{r_j}} \omega_{Y_{r_j}(p,q)} \left( h_{Y_{r_j}(p,q)}(W_{PPNet}, P_{r_j}^s) + \lambda_{IP} f_{Y_{r_j}(p,q)}(W_{IPNet}, X_{r_j}) \right), \quad (5)$$

where  $f_{Y_{r_j}(p,q)}(W_{IPNet}, X_{r_j})$  denotes the per-pixel logarithmic probability produced by the  $DN_{IPNet}$  and  $\lambda_{IP}$  balances the effect of PPNet and IPNet. We start training PSPNet by taking the trained PPNet in predictive parsing as initialization. We find this benefits the convergence of PSPNet training. In Section 4.2.1, we give more empirical studies.

Now we proceed to explain the role of AdapNet. There are two disadvantages by naively combining the intermediate features of PPNet and IPNet. Firstly, since the output features from those two feature encoders generally have different distributions, naively concatenating features harms the final performance as the “large” features dominate the “smaller” ones. Although during training, the followed weights in deconvNet may adjust accordingly, it requires careful tuning of parameters thus is subject to trial and error to find the best settings. Similar observations have been made in previous literature [19]. However, different from [19] which uses a normalization layer to tackle the scale problem, we use a more powerful AdapNet to transform the features of PPNet to be with proper norm and scale. Secondly, the intermediate features of PPNet and IPNet are with different semantic meanings, which means they reside in different manifold spaces. Therefore naively combining them increases the difficulty of learning the transformation from feature space to parsing map in the

followed  $DN_{IPNet}$ . By adding the AdapNet to convert the feature space in advance, it eases the training of  $DN_{IPNet}$ . Detailed explorations of the architecture of AdapNet follows in Section 4.2.1.

### 3.4. Discussion

**Unsupervised Feature Learning** Currently we train FPNet in a pseudo semi-supervised way, *i.e.* initialize  $G$  and  $D$  with ImageNet pre-trained models for faster training. Without using the pre-trained models, our approach becomes an unsupervised feature learning one. We also investigate the fully unsupervised learning strategy of FPNet in the experiments. The resulting FPNet is denoted as  $FPNet_{VGG11}^*$ . As shown in Table 1,  $FPNet_{VGG11}^*$  performs similarly well as  $FPNet_{VGG16}$  using the pre-trained VGG16 model. In the future, we will perform unsupervised learning on FPNet using deeper architectures and further improve its ability.

**Proactive Parsing** Our predictive learning approach recalls another challenging but attractive task, *i.e.*, to predict the future parsing maps for a few seconds without seeing them, only based on past frames. Achieving this allows autonomous vehicles or other parsing-demanded devices to receive parsing information ahead and get increased buffer time for decision making. Our approach indeed has the potential to accomplish this proactive parsing task. As one can observe from Figure 1, the predicted parsing maps capture the temporal information across frames, such as motions of cars, and the predicted parsing map reflects such dynamics and shows roughly correct prediction. In the future, we will investigate how to enhance the performance of our approach on predictive parsing to get higher-quality and longer-term future results.

## 4. Experiments

### 4.1. Settings and Implementation Details

**Datasets** Since PEARL tackles the scene parsing problem with temporal context, we choose Cityscapes [4] and Camvid [1] for evaluation. Both datasets provide annotated frames as well as adjacent frames, suitable for testing the temporal modeling ability. Cityscapes is a large-scale dataset containing finely pixel-wise annotations on 2,975/500/1,525 train/val/test frames with 19 semantic classes and another 20,000 coarsely annotated frames. Each finely annotated frame is sampled from the 20th frame of a 30-frame video clip in the dataset, giving in total 180K frames. Since there are no video data provided for the coarsely annotated frames, we only use finely annotated ones for training PEARL. Every frame in Cityscapes has a resolution of  $1024 \times 2048$  pixels.

The Camvid dataset contains 701 color images with annotations on 11 semantic classes. These images are ex-

tracted from driving videos captured at daytime and dusk. Each video contains 5,000 frames on average, with a resolution of  $720 \times 960$  pixels, giving in total 40K frames.

**Baselines** We conduct experiments to compare PEARL with two baselines which use different deep network architectures.

- **VGG16-baseline** Our VGG16-baseline is built upon DeepLab [3]. We make the following modifications. We add three deconvolutional layers after  $f_{c7}$  to learn better transformations to label maps. The architectures of three added deconvolutional layers in VGG16-baseline are  $O(256)-K(4)-S(2)-P(1)$ ,  $O(128)-K(4)-S(2)-P(1)$  and  $O(64)-K(4)-S(2)-P(1)$  respectively. Here  $O(n)$  denotes  $n$  output feature maps,  $K(n)$  denotes the kernel size of  $n \times n$ ,  $S(n)$  denotes a stride of length  $n$  and  $P(n)$  denotes the padding size of  $n$ . The layers before  $f_{c7}$  (included) constitute the encoder network ( $EN$  in Figure 2) and the other layers form the decoder network ( $DN$  in Figure 2).
- **Res101-baseline** Our Res101-baseline is modified from [10] by adapting it to a fully convolutional network following [21]. Specifically, we replace the average pooling layer and the 1000-way classification layer with a fully convolutional layer to produce dense label maps. Also, we modify  $conv5\_1$ ,  $conv5\_2$  and  $conv5\_3$  to be dilated convolutional layers by setting their dilation size to be 2 to enlarge the receptive fields. As a result, the output feature maps of  $conv5\_3$  have a stride of 16. Following [21], we utilize high-frequency features learned in bottom layers by adding skip connections from  $conv1$ ,  $pool1$ ,  $conv3\_3$  to corresponding up-sampling layers to produce label maps with the same size as input frames. The layers from  $conv1$  to  $conv5\_3$  belong to  $EN$  while the other layers belong to  $DN$ . Following [40], we also use hard training sample mining to reduce over-fitting.

Note that IPNets in PEARL share the same network architectures as baseline models. The encoder networks in FPNet/PPNet and the decoder network in PPNet share the same network architectures of the encoder network and the decoder network in baseline models, respectively.

**Evaluation Metrics** Following previous practice, we use the mean IoU (mIoU) for Cityscapes, and per-pixel accuracy (PA) and average per-class accuracy (CA) for Camvid. In particular, mIoU is defined as the pixel intersection-over-union (IOU) averaged across all categories; PA is defined as the percentage of all correctly classified pixels; and CA is the average of all category-wise accuracies.

**Implementation Details** Throughout the experiments, we set the number of preceding frames of each frame as

4, *i.e.*,  $s = 4$  in  $P_i^s$  (ref. Section 3.1 in the main text). When training FPNet, we randomly select frame sequences from those with a length of 4 and also enough movement (the  $\ell_2$  distance between the raw frames is larger than a threshold 230). In this way, we finally obtain 35K/8.8K sequences from Cityscapes and Camvid respectively. The input frames for training FPNet are all normalized such that values of their pixels lie between -1 and 1. For training PPNet and PSPNet, we only perform mean value subtraction on the frames. For training PPNet, we select the 4 frames before the annotated images to form the training sequences, where the frames are required to have sufficient motion, consistent with FPNet.

We perform random cropping and horizontal flipping to augment the training frames for FPNet, PPNet and PSPNet. In addition, for training FPNet, the temporal order of a sequence (including the to-predict frame and 4 preceding frames) is randomly reversed to model various dynamics in videos. The hyperparameters in PEARL are fixed as  $\lambda_{adv} = 0.2$  in Eq. (3) and  $\lambda_{IP} = 0.3$  in Eq. (5) and the probability threshold of hard training sample mining in Res101-baseline as 0.9. The values are chosen through cross-validation.

Since the class distribution is extremely unbalanced, we increase the weight of rare classes during training, similar to [5, 12, 32]. In particular, we adopt the re-weighting strategy in [32]. The weight for the class  $y_{i,j}$  is set as  $\omega_{y_{i,j}} = 2^{\lceil \log_{10}(\eta/f_{y_{i,j}}) \rceil}$  where  $f_{y_{i,j}}$  is the frequency of class  $y_{i,j}$  and  $\eta$  is a dataset-dependent scalar, which is defined according to the 85%/15% frequent/rare classes rule.

All of our experiments are carried out on NVIDIA Titan X and Tesla M40 GPUs using Caffe library.

**Computational Efficiency** Since no post-processing is required in PEARL, the running time of PEARL<sub>Res101</sub> for obtaining the parsing map of a frame with resolution  $1,024 \times 2,048$  is only 0.8 seconds on a modern GPU, among the fastest methods in existing works.

## 4.2. Results

Examples of parsing maps produced by PEARL are illustrated in Figure 5, where VGG16 is used in both the baseline model and PEARL. As can be seen from Figure 5, compared with the baseline model, PEARL produces smoother parsing maps, *e.g.* for the class *vegetation*, and stronger discriminability for small objects, *e.g.* for the classes *pole*, *pedestrian*, *traffic sign* as highlighted in red boxes. Such improvements are attributed to the capability of PEARL to learn the temporally consistent features and discriminative features for local pixel variance simultaneously.

Table 1: Comparative study of effects of FPNet and PPNet on final performance of PEARL over Cityscapes val set. The front models of FPNet and PPNet are VGG16 for fair comparison.  $feat_{Net}$  means only the output features of “Net” are combined with IPNet. FPNet\* means FPNet is trained from random initialization.

Methods	mIoU
VGG16-baseline	63.4
$feat_{FPNet_{VGG16}} + IPNet_{VGG16}$	68.6
$feat_{FPNet^*_{VGG16}} + IPNet_{VGG16}$	68.4
$feat_{PPNet_{VGG16}} + IPNet_{VGG16}$	69.1
PEARL <sub>VGG16</sub>	<b>69.8</b>

### 4.2.1 Cityscapes

**Ablation Analysis** We investigate contribution of each component of our approach.

(1) *Predictive Feature Learning.* To investigate the contribution of the two predictive feature learning networks, FPNet and PPNet, we conduct three experiments. The comparison results are listed in Table 1, where the front-CNN of FPNet and PPNet both use the VGG16 architecture.

First, we would like to verify the effectiveness of the features learned by FPNet for VSP. We concatenate the output features (denoted as  $feat_{FPNet}$ ) of  $Encoder_{FPNet}$  with the output features of  $Encoder_{IPNet}$  in IPNet, as shown in Figure 2c, and fix the  $Encoder_{FPNet}$  during training. In this way, the FPNet only extracts spatiotemporal features for the IPNet. As can be seen from Table 1, by combining  $feat_{FPNet}$ , the mIoU increases from 63.4 (of the VGG16-baseline) to 68.6, demonstrating FPNet indeed learns useful spatiotemporal features through frame prediction.

Similarly, we replace  $Encoder_{FPNet}$  in the above experiment with  $Encoder_{PPNet}$  to investigate the influence of PPNet on final performance. In the experiment, the per-pixel cross-entropy loss layer of PPNet is removed and the weights of  $Encoder_{PPNet}$  are fixed. As illustrated in Table 1, combining IPNet with  $feat_{PPNet}$  further increases the mIoU by 0.5 compared to combining  $feat_{FPNet}$ , demonstrating features learned by PPNet from predictive parsing are useful for VSP.

Finally, we look into the effectiveness of joint training of PPNet and IPNet. It is observed from Table 1 that the resulting model, *i.e.* PEARL<sub>VGG16</sub> achieves the best performance, benefiting from the joint end-to-end training strategy.

(2) *Comparison with Optical Flow Methods.* To verify the superiority of PEARL on learning the temporal information specific for VSP, we compare PEARL with other temporal context modeling methods. First, we naively pass each frame in  $P_i^s$  and  $X_i$  through baseline models (both VGG16-baseline and Res101-baseline) and merge

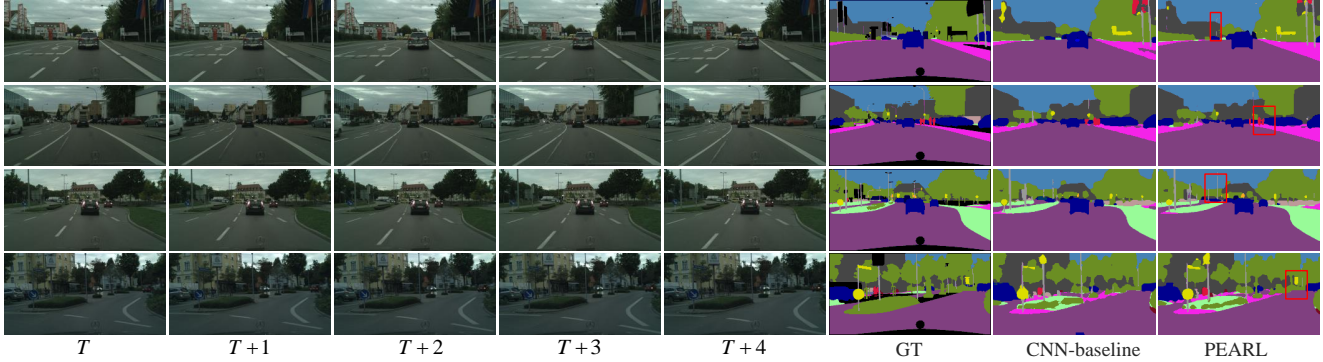


Figure 5: Examples of parsing results of PEARL on Cityscape val set. The first five images in each row represents a different video sequence, which are followed by ground truth annotations, the parsing map of the baseline model and the parsing map of our proposed PEARL, all for frame  $T+4$ . It is observable that PEARL produces more smooth label maps and shows stronger discriminability for small objects (highlighted in red boxes) compared to the baseline model. Best viewed in color and zoomed in pdf.

Table 3: Performance comparison of PEARL with state-of-the-arts on Cityscapes *test* set. Note for fast inference, single scale testing is used in PEARL without any post-processing like CRF.

Methods	road	sidewalk	building	wall	fence	pole	traffic light	traffic sign	vegetation	terrain	sky	person	rider	car	truck	bus	train	motorcycle	bicycle	mIoU
FCN_8s [21]	97.4	78.4	89.2	34.9	44.2	47.4	60.1	65.0	91.4	69.3	93.9	77.1	51.4	92.6	35.3	48.6	46.5	51.6	66.8	65.3
DPN [20]	97.5	78.5	89.5	40.4	45.9	51.1	56.8	65.3	91.5	69.4	94.5	77.5	54.2	92.5	44.5	53.4	49.9	52.1	64.8	66.8
Dilation10 [41]	97.6	79.2	89.9	37.3	47.6	53.2	58.6	65.2	91.8	69.4	93.7	78.9	55.0	93.3	45.5	53.4	47.7	52.2	66.0	67.1
DeepLab [2]	97.9	81.3	90.3	48.8	47.4	49.6	57.9	67.3	91.9	69.4	94.2	79.8	59.8	93.7	56.5	67.5	57.5	57.7	68.8	70.4
Adelaide [18]	98.0	82.6	90.6	44.0	50.7	51.1	65.0	71.7	92.0	<b>72.0</b>	94.1	81.5	61.1	94.3	61.1	65.1	53.8	61.6	70.6	71.6
LRR-4X [7]	97.9	81.5	91.4	<b>50.5</b>	52.7	59.4	<b>66.8</b>	<b>72.7</b>	92.5	70.1	95.0	81.3	60.1	94.3	51.2	67.7	54.6	55.6	69.6	71.8
PEARL <sup>1</sup> (ours)	<b>98.3</b>	<b>83.9</b>	<b>91.6</b>	47.6	<b>53.4</b>	<b>59.5</b>	<b>66.8</b>	72.5	<b>92.7</b>	70.9	<b>95.2</b>	<b>82.4</b>	<b>63.5</b>	<b>94.7</b>	<b>57.4</b>	<b>68.8</b>	<b>62.2</b>	<b>62.6</b>	<b>71.5</b>	<b>73.4</b>

Table 2: Comparative study of PEARL with optical flow based method on two deep networks: VGG16 and Res101 to verify the superiority of PEARL on modeling temporal information.  $feat_{OF}$  means the optical flow maps calculated by epic flow [27].

Methods	mIoU
VGG16-baseline	63.4
$feat_{OF}$ + IPNet <sub>VGG16</sub>	64.5
PEARL <sub>VGG16</sub>	<b>69.8</b>
Res101-baseline	72.5
$feat_{OF}$ + IPNet <sub>Res101</sub>	72.7
PEARL <sub>Res101</sub>	<b>74.9</b>

their probability maps to obtain the final label map of  $X_i$ . It is verified by experiments that such methods achieve worse performance than baseline models due to their weakness of utilizing temporal information and the noisy probability map produced for each frame. Since optical flow

Table 4: Comparison with state-of-the-arts on Cityscapes val set. Single scale testing is used in PEARL w/o post-processing as CRF.

Methods	mIoU
VGG16-baseline (ours)	63.4
FCN [21]	61.7
Pixel-level Encoding [38]	64.3
DPN [20]	66.8
Dilation10 [41]	67.1
DeepLab-VGG16 [2]	62.9
Deep Structure [18]	68.6
Clockwork FCN [31]	64.4
PEARL <sub>VGG16</sub> (ours)	<b>69.8</b>
Res101-baseline (ours)	72.5
DeepLab-Res101 [2]	71.4
PEARL <sub>Res101</sub> (ours)	<b>74.9</b>

is naturally capable of modeling the temporal information in videos, we use it as a strong baseline to compete with PEARL. We employ the epic flow [27] for computing all

<sup>1</sup><https://www.cityscapes-dataset.com/method-details/?submissionID=328>



optical flows. Then we warp the parsing map of the frame  $X_{i-1}$  and merge it with that of the frame  $X_i$ , according to the optical flow calculated between these two frames. In this way, the temporal context is modeled explicitly via optical flow. This method performs better than the last method but its performance is still inferior to baseline models. It is because the CNN models produce the parsing map without knowing the temporal information during training.

Then we conduct the third experiment by concatenating the optical flows calculated from  $X_{i-1}$  to  $X_i$  with the frame  $X_i$ , which forms 5-channel raw data (RGB plus X/Y channels of optical flow). Based on optical flow augmented data, we re-train baseline models. During training, each kernel in the first convolutional layer of baseline models is randomly initialized for the weights corresponding to the X/Y channels of optical flow. This method is referred to as “feat<sub>OF</sub> + IPNet”. The comparative results of “feat<sub>OF</sub> + IPNet” and PEARL using VGG16 and Res101 are displayed in Table 2. From the results, one can observe “feat<sub>OF</sub> + IPNet” achieves higher performance than baselines models as it uses temporal context during training. Notably, PEARL significantly beats “feat<sub>OF</sub> + IPNet” on both network architectures, proving its superior ability to model temporal information for VSP problems.

(3) *Ablation Study of AdapNet* As introduced in Section 3.3, AdapNet improves the performance of PEARL by learning the latent transformations from  $Encoder_{PPNet}$  to  $Encoder_{IPNet}$ . In our experiments, the AdapNet contains one convolutional layer followed by ReLU. The kernel size of the convolutional layer is 1 and the number of kernels is equal to that of output channels of  $Encoder_{PPNet}$ . Compared to the PEARL w/o AdapNet, adding AdapNet brings 1.1/0.3 mIoU improvements for PEARL<sub>VGG16</sub> and PEARL<sub>Res101</sub>, respectively. We also conduct experiments by increasing convolutional layers of AdapNet, but only observe marginal improvements. Since deeper AdapNet brings more computation cost, we use AdapNet with one convolutional layer.

**Comparison with State-of-the-arts** The comparison of PEARL with other state-of-the-arts on Cityscapes val set is listed in Table 4, from which one can observe PEARL achieves the best performance among all compared methods on both network architectures. Note loss re-weighting is not used on this dataset.

Specifically, PEARL<sub>VGG16</sub> and PEARL<sub>Res101</sub> significantly improve the corresponding baseline models by 6.4/2.4 mIoU, respectively. Notably, compared with [31] which proposed a temporal skip network based on VGG16 for video scene parsing, PEARL<sub>VGG16</sub> beats it by 5.4 in terms of mIoU. We also note that different from other methods which extensively modify VGG16 networks to enhance its discriminative power for image parsing, e.g. [2, 18], our PEARL<sub>VGG16</sub> is built on vanilla VGG16 architecture. Thus

it is reasonable to expect further improvement on the VSP performance by using more powerful base network architectures. Furthermore, we submit PEARL<sub>Res101</sub> to the online evaluation server of Cityscapes to compare with other state-of-the-arts on Cityscapes test set. As shown in Table 3, our method achieves the best performance among all top methods which have been published till the time of submission. Note in inference, PEARL only uses single-scale testing without CRF post-processing for the sake of fast inference.

#### 4.2.2 Camvid

We further investigate the effectiveness of PEARL on Camvid. Its result and the best results ever reported on this dataset are listed in Table 5. Following [12, 32], loss re-weighting is used on this dataset. One can observe that PEARL performs much better than all competing methods — significantly improving the PA/CA of the baseline model (Res101-baseline) by 1.6%/2.5% respectively, once again demonstrating its strong capability of improving VSP performance. Notably, compared to the optical flow based methods [16] and [20] which utilize CRF to model temporal information, PEARL shows great advantages in performance, verifying its superiority in modeling the spatiotemporal context for VSP problems.

Table 5: Comparison with the state-of-the-art methods on CamVid.

Methods	PA(%)	CA(%)
Res101-baseline (ours)	92.6	80.0
Ladicky <i>et al.</i> [15]	83.8	62.5
SuperParsing [37]	83.9	62.5
DAG-RNN [32]	91.6	78.1
MPF-RNN [12]	92.8	82.3
Liu <i>et al.</i> [20]	82.5	62.5
RTDF [16]	89.9	80.5
PEARL (ours)	<b>94.2</b>	<b>82.5</b>

## 5. Conclusion

We proposed a predictive feature learning method for effective video scene parsing. It contains two novel components: predictive feature learning and prediction steering parsing. The first component learns spatiotemporal features by predicting future frames and their parsing maps without requiring extra annotations. The prediction steering parsing architecture then guides the single frame parsing network to produce temporally smooth and structure preserving results by using the predictive feature learning outputs. Extensive experiments on Cityscapes and Camvid fully demonstrated the effectiveness of our approach.

## References

- [1] G. J. Brostow, J. Shotton, J. Fauqueur, and R. Cipolla. Segmentation and recognition using structure from motion point clouds. In *ECCV*. 2008. 6
- [2] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *CoRR*, abs/1606.00915, 2016. 5, 8, 9
- [3] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*, 2015. 1, 2, 6
- [4] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. *arXiv preprint arXiv:1604.01685*, 2016. 1, 6
- [5] C. Farabet, C. Couprie, L. Najman, and Y. LeCun. Learning hierarchical features for scene labeling. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(8):1915–1929, 2013. 1, 2, 7
- [6] G. Floros and B. Leibe. Joint 2d-3d temporally consistent semantic segmentation of street scenes. In *CVPR*, pages 2823–2830. IEEE, 2012. 2
- [7] G. Ghiasi and C. C. Fowlkes. Laplacian reconstruction and refinement for semantic segmentation. *CoRR*, abs/1605.02264, 2016. 8
- [8] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014. 3, 4
- [9] B. L. . X. H. . S. Gould. Multi-class semantic video segmentation with exemplar-based object reasoning. In *WACV*, 2016. 2, 3
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015. 2, 6
- [11] J. Hur and S. Roth. Joint optical flow and temporally consistent semantic segmentation. In *ECCV*, pages 163–177. Springer, 2016. 2
- [12] X. Jin, Y. Chen, J. Feng, Z. Jie, and S. Yan. Multi-path feedback recurrent neural network for scene parsing. *CoRR*, abs/1608.07706, 2016. 2, 7, 9
- [13] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, pages 1097–1105, 2012. 4
- [14] A. Kundu, V. Vineet, and V. Koltun. Feature space optimization for semantic video segmentation. In *CVPR*, 2016. 2
- [15] L. Ladickỳ, P. Sturges, K. Alahari, C. Russell, and P. H. Torr. What, where and how many? combining object detectors and crfs. In *ECCV*. 2010. 9
- [16] P. Lei and S. Todorovic. Recurrent temporal deep field for semantic video labeling. In *ECCV*, pages 302–317. Springer, 2016. 2, 9
- [17] M. Liang, X. Hu, and B. Zhang. Convolutional neural networks with intra-layer recurrent connections for scene labeling. In *NIPS*, 2015. 2
- [18] G. Lin, C. Shen, A. v. d. Hengel, and I. Reid. Exploring context with deep structured models for semantic segmentation. *arXiv preprint arXiv:1603.03183*, 2016. 8, 9
- [19] W. Liu, A. Rabinovich, and A. C. Berg. Parsenet: Looking wider to see better. *arXiv preprint arXiv:1506.04579*, 2015. 5
- [20] Z. Liu, X. Li, P. Luo, C.-C. Loy, and X. Tang. Semantic image segmentation via deep parsing network. In *ECCV*, pages 1377–1385, 2015. 8, 9
- [21] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 1, 2, 5, 6, 8
- [22] W. Lotter, G. Kreiman, and D. Cox. Unsupervised learning of visual structure using predictive generative networks. *arXiv preprint arXiv:1511.06380*, 2015. 3
- [23] B. Mahasseni, S. Todorovic, and A. Fern. Approximate policy iteration for budgeted semantic video segmentation. *CoRR*, abs/1607.07770, 2016. 2, 3
- [24] M. Mathieu, C. Couprie, and Y. LeCun. Deep multi-scale video prediction beyond mean square error. *arXiv preprint arXiv:1511.05440*, 2015. 3, 4, 5
- [25] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature learning by inpainting. *arXiv preprint arXiv:1604.07379*, 2016. 3
- [26] P. H. Pinheiro and R. Collobert. Recurrent convolutional neural networks for scene parsing. *arXiv preprint arXiv:1306.2795*, 2013. 2
- [27] J. Revaud, P. Weinzaepfel, Z. Harchaoui, and C. Schmid. Epicflow: Edge-preserving interpolation of correspondences for optical flow. In *CVPR*, pages 1164–1172, 2015. 2, 8, 9
- [28] A. Roy and S. Todorovic. Scene labeling using beam search under mutex constraints. In *CVPR*, 2014. 1, 2
- [29] A. G. Schwing and R. Urtasun. Fully connected deep structured networks. *arXiv preprint arXiv:1503.02351*, 2015. 1, 2
- [30] L. Sevilla-Lara, D. Sun, V. Jampani, and M. J. Black. Optical flow with semantic segmentation and localized layers. *arXiv preprint arXiv:1603.03911*, 2016. 2

- [31] E. Shelhamer, K. Rakelly, J. Hoffman, and T. Darrell. Clockwork convnets for video semantic segmentation. *arXiv preprint arXiv:1608.03609*, 2016. 8, 9
- [32] B. Shuai, Z. Zuo, G. Wang, and B. Wang. Dag-recurrent neural networks for scene labeling. *arXiv preprint arXiv:1509.00552*, 2015. 2, 7, 9
- [33] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2
- [34] R. Socher, C. C. Lin, C. Manning, and A. Y. Ng. Parsing natural scenes and natural language with recursive neural networks. In *ICML*, 2011. 1, 2
- [35] P. Sturgess, K. Alahari, L. Ladicky, and P. H. Torr. Combining appearance and structure from motion features for road scene understanding. In *BMVC*, 2009. 2
- [36] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. *arXiv preprint arXiv:1409.4842*, 2014. 4
- [37] J. Tighe and S. Lazebnik. Superparsing: scalable non-parametric image parsing with superpixels. In *ECCV*. 2010. 9
- [38] J. Uhrig, M. Cordts, U. Franke, and T. Brox. Pixel-level encoding and depth layering for instance-level semantic labeling. *arXiv preprint arXiv:1604.05096*, 2016. 8
- [39] F. Visin, K. Kastner, A. C. Courville, Y. Bengio, M. Matteucci, and K. Cho. Reseg: A recurrent neural network for object segmentation. *CoRR*, abs/1511.07053, 2015. 2
- [40] Z. Wu, C. Shen, and A. van den Hengel. High-performance semantic segmentation using very deep fully convolutional networks. *CoRR*, abs/1604.04339, 2016. 6
- [41] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015. 8
- [42] C. Zhang, L. Wang, and R. Yang. Semantic segmentation of urban scenes using dense depth maps. In *ECCV*, pages 708–721. Springer, 2010. 2
- [43] Y. Zhang and T. Chen. Efficient inference for fully-connected crfs with stationarity. In *CVPR*, 2012. 1, 2
- [44] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. Torr. Conditional random fields as recurrent neural networks. In *ICCV*, 2015. 1, 2

# Appendices

## A. Qualitative Evaluation of PEARL

### A.1. More Results of Frame Prediction from FPNet in PEARL

Please refer to Figure 6;

### A.2. More Results of Video Scene Parsing

Please refer to Figure 9.

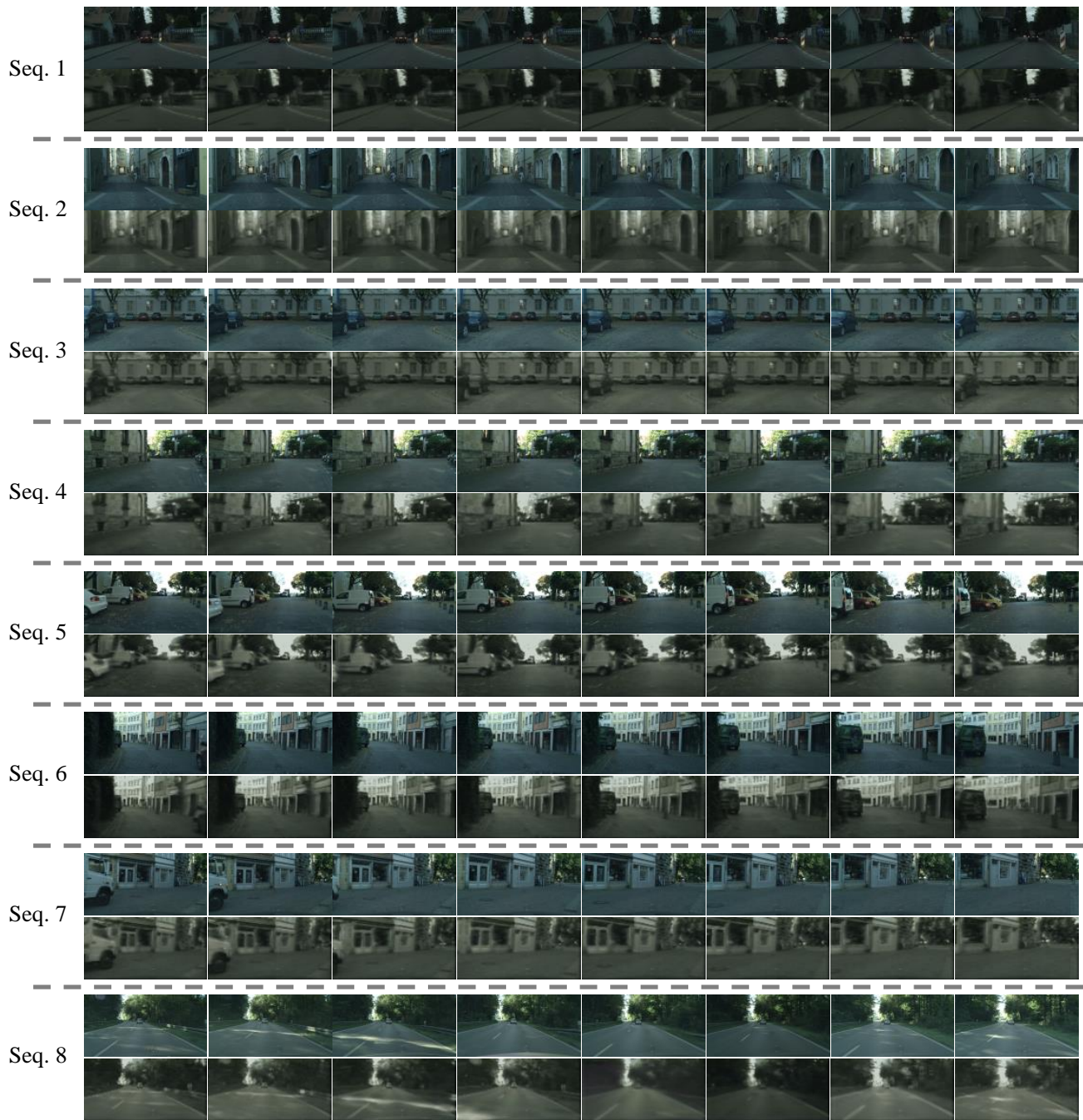
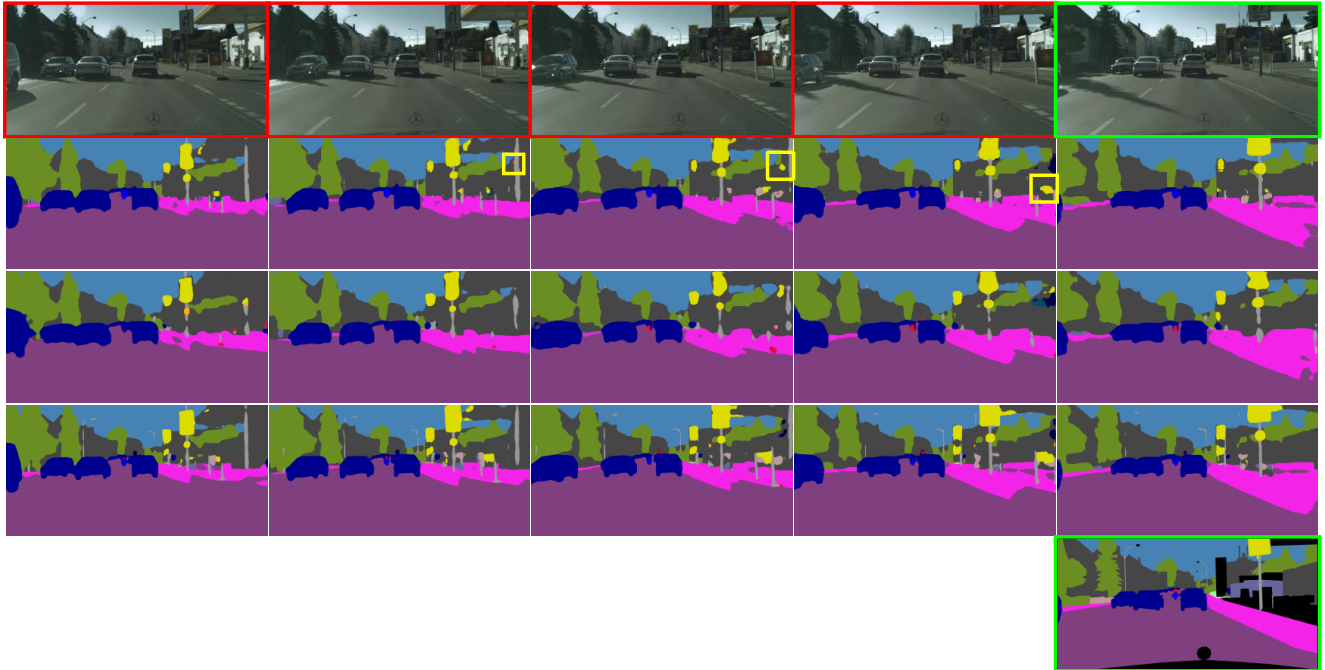
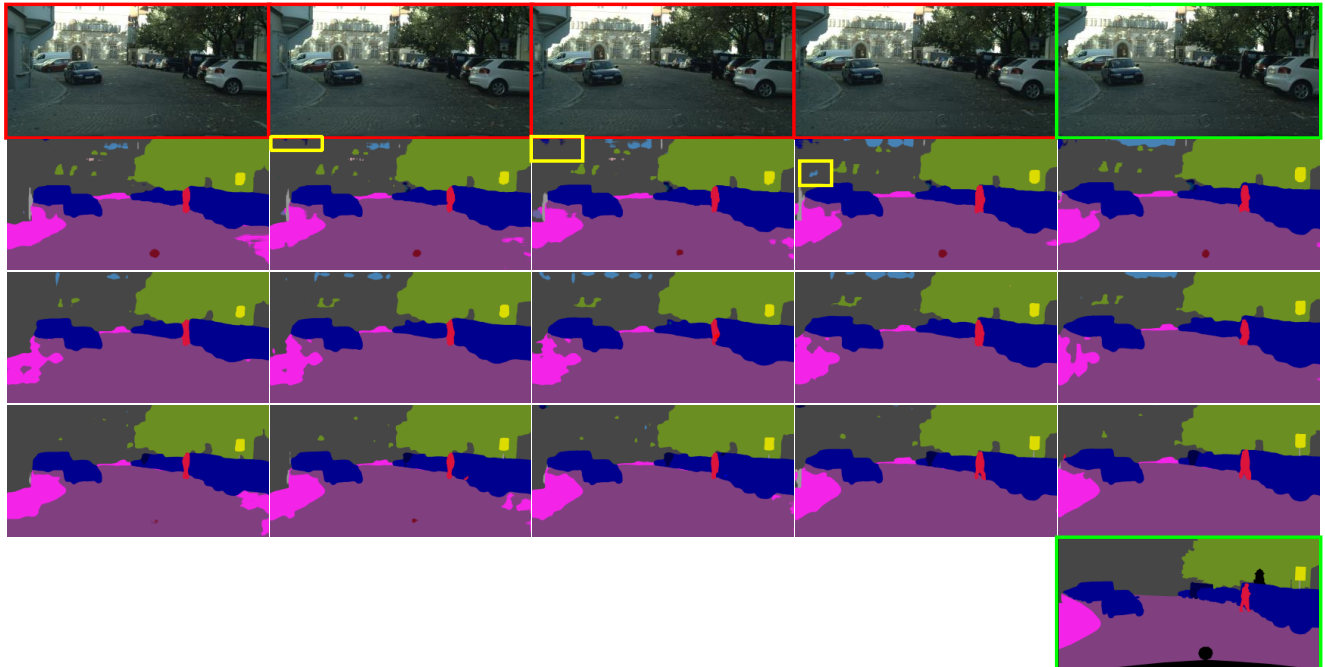


Figure 6: Eight sequences of videos in Cityscapes val set with frame prediction results. For each sequence, the upper row contains eight ground truth frames and the bottom row contains frame predictions produced by FPNet in PEARL. It is observed that FPNet is able to model the structures of objects and stuff as well as the motion information of moving objects in videos. Best viewed in color and zoomed pdf.

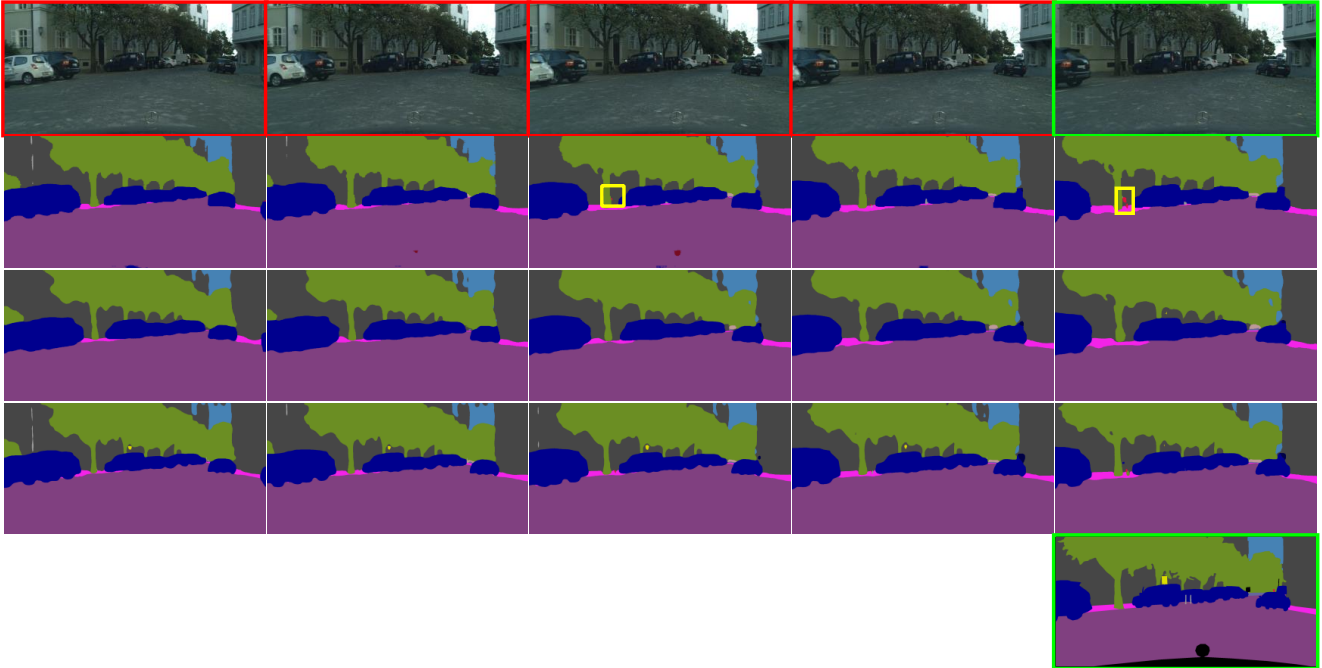
Seq. 1



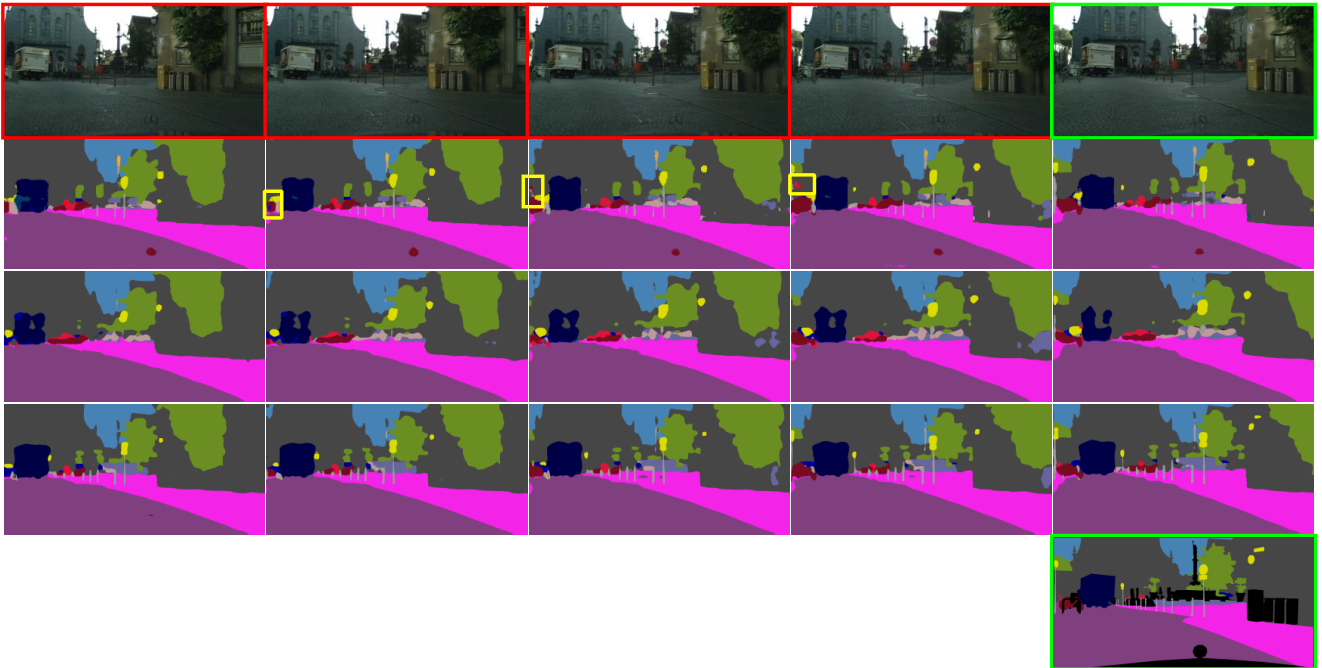
Seq. 2



Seq. 3



Seq. 4



Seq. 5

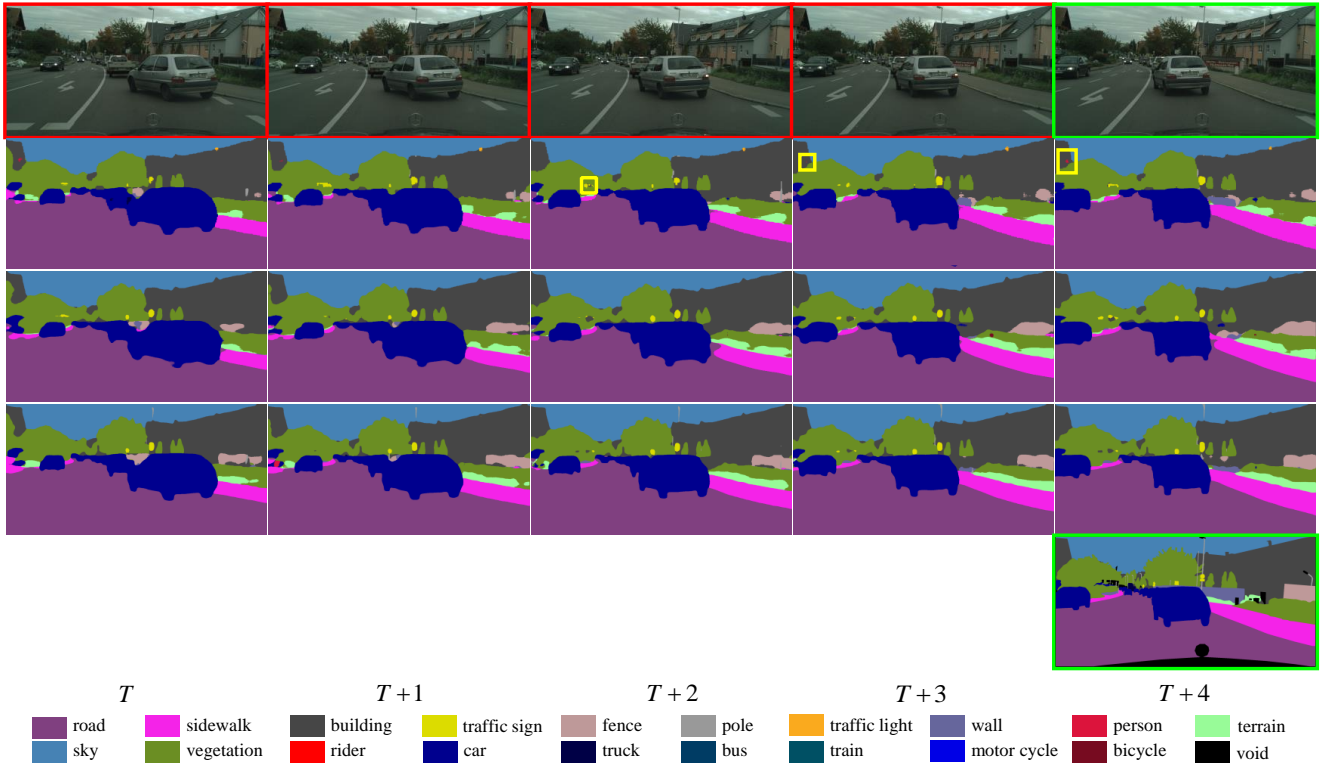


Figure 9: Examples of parsing results of PEARL on Cityscape val set. All parsing maps have the same resolution as input frames. **Top row**: a five-frame sequence (the four preceding frames are highlighted by red and the target frame to parse has a green boundary). **Second row**: frame parsing maps produced by the VGG16-baseline model. Since it cannot model temporal context, the baseline model produces parsing results with undesired inconsistency across frames as in yellow boxes. **Third row**: predictive parsing maps output by PEARL. The inconsistent parsing regions in the second row are classified consistently across frames. **Fourth row**: parsing maps produced by PEARL with better accuracy and temporal consistency due to the combination of advantages of traditional image parsing model (the second row) and predictive parsing model (the third row). **Bottom row**: the ground truth label map (with green boundary) for the frame  $T+4$ . Best viewed in color and zoomed pdf.