



**HAL**  
open science

## Predicting Deeper into the Future of Semantic Segmentation

Natalia Neverova, Pauline Luc, Camille Couprie, Jakob Verbeek, Yann Lecun

► **To cite this version:**

Natalia Neverova, Pauline Luc, Camille Couprie, Jakob Verbeek, Yann Lecun. Predicting Deeper into the Future of Semantic Segmentation. 2017. hal-01494296v1

**HAL Id: hal-01494296**

**<https://inria.hal.science/hal-01494296v1>**

Preprint submitted on 23 Mar 2017 (v1), last revised 21 Aug 2017 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Predicting Deeper into the Future of Semantic Segmentation

Natalia Neverova\*  
Facebook AI Research  
Paris  
neverova@fb.com

Pauline Luc\*  
Facebook AI Research, Paris  
Université de Grenoble  
paulineluc@fb.com

Camille Couprie  
Facebook AI Research  
Paris  
couprie@fb.com

Jakob Verbeek  
INRIA, Laboratoire Jean Kuntzmann  
Université de Grenoble  
jakob.verbeek@inria.fr

Yann LeCun  
Facebook AI Research, New York  
New York University  
yann@fb.com

## Abstract

*The ability to predict and therefore to anticipate the future is an important attribute of intelligence. It is also of utmost importance in real-time systems, e.g. in robotics or autonomous driving, which depend on visual scene understanding for decision making. While prediction of the raw RGB pixel values in future video frames has been studied in previous work, here we focus on predicting semantic segmentations of future frames. More precisely, given a sequence of semantically segmented video frames, our goal is to predict segmentation maps of not yet observed video frames that lie up to a second or further in the future. We develop an autoregressive convolutional neural network that learns to iteratively generate multiple frames. Our results on the Cityscapes dataset show that directly predicting future segmentations is substantially better than predicting and then segmenting future RGB frames. Our models predict trajectories of cars and pedestrians much more accurately (25%) than baselines that copy the most recent semantic segmentation or warp it using optical flow. Prediction results up to half a second in the future are visually convincing, the mean IoU of predicted segmentations reaching two thirds of the real future segmentations.*

## 1. Introduction

Prediction and anticipation of future events is a key component to intelligent decision-making [27]. Building smarter robotic systems and autonomous vehicles implies making decisions based on the analysis of the current situ-

\*These authors contributed equally

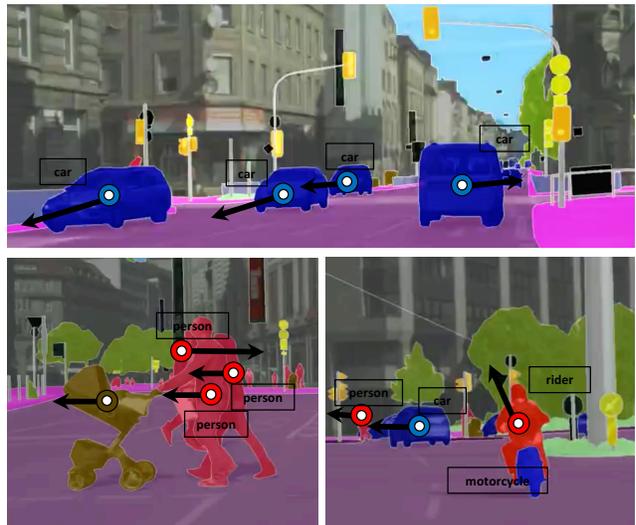


Figure 1: Given the semantic segmentation of past frames obtained by a state of the art model, our goal is to predict segmentations of future frames that have not yet been observed. Our models solve this task by learning the semantic-level scene dynamics.

ation and hypotheses made on what could happen next [7]. While humans can predict vehicle or pedestrian trajectories effortlessly and at the reflex level, it remains an open challenge for current computer vision systems.

The task of predicting future RGB video frames given preceding ones is an interesting one to assess if current vision systems are able to reason about future events, and it has recently received significant attention [25, 26, 22]. Modeling the raw RGB intensities, however, might be a task

that is overly complicated as compared to predicting future high-level scene properties, while the latter is sufficient for many applications. Such high-level future prediction has been studied in various forms, e.g. by explicitly forecasting trajectories of people and other objects in future video frames [1, 17]. In our work we do not explicitly model objects or other scene elements, but instead model the dynamics of semantic segmentation maps of object categories implicitly using convolutional neural networks.

Semantic segmentation is one of the most complete forms of visual scene understanding, where the goal is to label each pixel with the corresponding semantic label (e.g., *tree*, *pedestrian*, *car*, etc.). In our work, we build upon the recent progress in this area [8, 19, 18, 23, 24, 36], and develop models to predict the semantic segmentation of future video frames, given the segmentation of several preceding frames. See Figure 1 for an illustration.

The pixel-level annotations needed for semantic segmentation are expensive to acquire, and this is even worse if we need annotations for each video frame to learn models to predict future semantic segmentations. To alleviate this issue we rely on state-of-the-art semantic image segmentation models to label all frames in videos, and then learn our future segmentation prediction models from these automatically generated annotations.

We systematically study the effect of using RGB frames and/or segmentations as inputs and outputs for our models, and the impact of various loss functions. Our experiments on the Cityscapes dataset [4] suggest that it is advantageous to directly predict future frames at the more abstract semantic-level, rather than to predict at the low-level RGB appearance of future frames [22, 25] and then to apply a semantic segmentation model to these. By moving away from raw RGB predictions and modeling pixel-level object labels instead, we believe that more of the network’s modeling capacity is allocated to learn basic physics and object interaction dynamics. Moreover, artifacts in the predicted RGB frames can lead to poor performance of still image semantic segmentation models.

In this work we make three contributions:

- Our approach is the first to predict future semantic segmentations without requiring extremely costly temporally dense video annotation. Its genericity allows different architectures for still-image segmentation and future segmentation prediction to be swapped in.
- We propose an autoregressive training strategy that outperforms predicting multiple future frames in a single feed-forward batch model.
- Our model convincingly predicts segmentations up to 0.5 seconds into the future, the mean IoU of our predictions reaches two thirds of the one obtained by a state-of-the-art semantic segmentation model [36].

## 2. Related work

Here we discuss the most relevant related work on video forecasting and on disambiguating learning under uncertainty, in particular using adversarial training.

**Video forecasting.** Several authors developed methods to improve the temporal stability of semantic video segmentation. Jin *et al.* [14] train a model to predict the semantic segmentation of the immediate next image from the preceding input frames, and fuse this prediction with the segmentation computed from the next input frame. Nilsson and Sminchisescu [23] use a convolutional RNN model with a spatial transformer component [13] to accumulate the information from past and future frames in order to improve prediction of the current frame segmentation. In a similar spirit, Patraucean *et al.* [24] employ a convolutional RNN to explicitly predict the optical flow, and use these to warp and aggregate per-frame segmentations. In contrast, our work is focused on predicting future segmentations without seeing the corresponding frames. Most importantly, we target a longer time horizon than a single frame.

A second line of related work focuses on generative models for future video frame forecasting. Ranzato *et al.* [25] introduced the first baseline of next video frame prediction. Srivastava *et al.* [26] developed a Long Short Term Memory (LSTM) [12] architecture for the task, and demonstrated a gain in action classification using the learned features. Mathieu *et al.* [22] improved the predictions using a multi-scale convolutional architecture, adversarial training [10], and a gradient difference loss. To reduce the number of parameters to estimate, several authors re-parameterize the problem to predict frame transformations instead of raw pixels [9, 28]. Luo *et al.* [21] employ a convolutional LSTM architecture to predict sequences of up to eight frames of optical flow in RGBd videos. The video pixel network of Kalchbrenner *et al.* [15] also uses LSTMs, and factorizes the temporal and spatial/color dimensions. Rather than predicting pixels or flows, Vondrick *et al.* [29] instead predict features in future frames. They predict the activations in the penultimate layer of the AlexNet in future frames, and use these to predict the presence of human actions in these frames.

**Learning with uncertainty.** Generative Adversarial Networks (GAN) [10] and Variational Auto-Encoders (VAE) [16] are recent unsupervised learning methods that can be used to deal with the inherent uncertainty in future-prediction tasks. GAN training has recently been improved using the Wasserstein GAN approach of Arjovsky and Bottou [2], and the energy based approach of Zhao *et al.* [37]. An interesting approach using GANs for unsupervised image representation learning is proposed in [5], where the generative model is trained along with an inference model that maps images to their latent representations. Vondrick

Model	Input	Output
X2X	$X_{1:t}$	$X_{t+1}$
S2S	$S_{1:t}$	$S_{t+1}$
XS2X	$(X_{1:t}, S_{1:t})$	$X_{t+1}$
XS2S	$(X_{1:t}, S_{1:t})$	$S_{t+1}$
XS2XS	$(X_{1:t}, S_{1:t})$	$(X_{t+1}, S_{t+1})$

Table 1: The 5 models for single-frame prediction.

*et al.* [30] showed that GANs can be applied to video generation. They use a two-stream generative model: one stream generates a static background, while the other generates a dynamic foreground sequence which is pasted on the background. Yang *et al.* [35] use similar ideas to develop an iterative image generation model where objects are sequentially pasted on the image canvas using a recurrent GAN. Xue *et al.* [34] predict future video frames from a single given frame using a VAE approach. Similarly, Walker *et al.* [31] perform forecasting with a VAE, predicting feature point trajectories from still images.

### 3. Predicting future frames and segmentations

This section first presents the different scenarios that we investigated to predict RGB pixel values and/or segmentations of the next frame. Then we describe two extensions of the single frame prediction model to predict further into the future.

#### 3.1. Single frame prediction models

Manual pixel-level supervision is laborious to acquire for still image semantic segmentation, and even more so for all frames in semantic video segmentation. To circumvent the need for datasets with per-frame annotations, we use the state-of-the-art multi-scale Dilated-10 semantic image segmentation network [36] to provide target semantic segmentations for all frames in each video. Then, we use the resulting temporally dense sequences of segmentations to learn our models.

Let us denote with  $X_i$  the  $i$ -th frame of a video sequence, and denote the sequences of frames from  $X_t$  to  $X_T$  as  $X_{t:T}$ . We denote by  $S_i$  the semantic segmentation of frame  $X_i$  given the Dilated-10 network. We represent the segmentations  $S_i$  using log-probabilities rather than the output of the softmax probabilities. This is motivated by recent observations in network distillation that the log-probabilities carry more information when training one network on the outputs of another network [3, 11]. For single-frame future prediction, we consider five different models that differ in whether they take RGB frames and/or segmentations as their input and output. We list these models in Table 1.

**Architectures.** Model X2X is a next-frame prediction

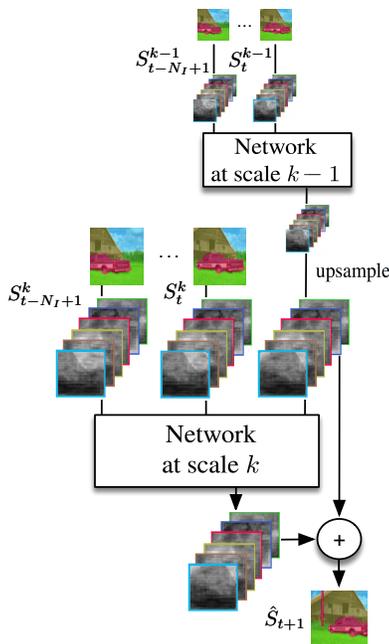


Figure 2: Multi-scale architecture of the S2S model that predicts the semantic segmentation of the next frame given the segmentation maps of the  $N_I$  previous frames.

model, for which we use the multi-scale network of Mathieu *et al.* [22], with two spatial scales. Noting  $C$  the number of output channels, each scale module corresponds to a simple four-layer convolutional network alternating convolutions and ReLU operations, containing feature maps with 128, 256, 128,  $C$  channels each, and convolution kernels of size 3 for the smaller scale, and of 5, 3, 3, 5 for the largest scale. No pooling is employed. The last non-linear function is a hyperbolic tangent, to ensure that the (normalized) RGB values lie in the range  $[-1, 1]$ . The output at a coarser scale is up-sampled, and used in input to the next scale module together with a copy of the input at that scale.

For models that output segmentations, we removed the last hyperbolic tangent non-linearities for the corresponding output channels, since the class log-probabilities are not limited to a fixed range. Apart from this difference, the S2S model, that predicts the future segmentation from past ones, has the same architecture as the X2X model. The multi-scale architecture of the S2S model is illustrated in Figure 2. The other models (XS2X, XS2S, and XS2XS), which take both RGB frames and segmentation maps as input, also use the same internal architecture, and just vary in the number of input and output channels.

**Loss function.** Following [22], for all models the loss function between the model output  $\hat{Y}$  and the target output  $Y$  is the sum of an  $\ell_1$  loss and a gradient difference loss:

$$\mathcal{L}(\hat{Y}, Y) = \mathcal{L}_{\ell_1}(\hat{Y}, Y) + \mathcal{L}_{\text{gdl}}(\hat{Y}, Y). \quad (1)$$

Using  $Y_{ij}$  to denote the pixel elements in  $Y$ , and similarly for  $\hat{Y}$ , the losses are defined as:

$$\mathcal{L}_{\ell_1}(\hat{Y}, Y) = \sum_{i,j} |Y_{ij} - \hat{Y}_{ij}|, \quad (2)$$

$$\begin{aligned} \mathcal{L}_{\text{gdl}}(\hat{Y}, Y) = \sum_{i,j} & \left| |Y_{i,j} - Y_{i-1,j}| - |\hat{Y}_{i,j} - \hat{Y}_{i-1,j}| \right| \\ & + \left| |Y_{i,j-1} - Y_{i,j}| - |\hat{Y}_{i,j-1} - \hat{Y}_{i,j}| \right|, \quad (3) \end{aligned}$$

where  $|\cdot|$  denotes the absolute value function. The  $\ell_1$  loss tries to match all pixel predictions independently to their corresponding target values. The gradient difference loss [22], instead, penalizes errors in the gradients of the target image and the predicted one. This loss is relatively insensitive to low-frequency mismatches between prediction and target (e.g., adding a constant to all pixels does not affect the loss), and is more sensitive to high-frequency mismatches that are perceptually more significant (e.g. blurring the contours of an object). We present a comparison of this loss with a multiclass cross entropy loss in Section 4.

**Adversarial training.** As shown by Mathieu *et al.* [22] in the context of raw images, introducing an adversarial loss term allows the model to disambiguate between modes corresponding to different turns of events, and reduces blur associated with this uncertainty. A recent study by Luc *et al.* [20] has demonstrated the positive influence of the adversarial training for semantic image segmentation, and its effectiveness in detecting higher-order spatial inconsistencies in the produced outputs.

Our formulation of the adversarial loss term is based on the principles of recently introduced Wasserstein GAN [2], with some modifications for the semantic segmentation application. In the case of the S2S model, the discriminator  $\mathcal{D}_\Theta$  is trained to maximize the  $\ell_1$  distance between ground truth sequences  $(S_{1:t}, S_{t+1})$  and sequences  $(S_{1:t}, \hat{S}_{t+1})$  predicted by our model:

$$\max_{\Theta} \left| \sigma(\mathcal{D}_\Theta(S_{1:t}, S_{t+1})) - \sigma(\mathcal{D}_\Theta(S_{1:t}, \hat{S}_{t+1})) \right|, \quad (4)$$

where  $\Theta$  is the set of parameters of the discriminator  $\mathcal{D}$ . The outputs produced by the predictive model are log-probability maps with unbounded values. In the Wasserstein GAN setting, they are encouraged to grow indefinitely. To avoid this and stabilize training, we employ an additional sigmoid non-linearity  $\sigma$  at the output of the discriminator, and set explicit targets for two kinds of outputs: 0 for the sequences produced by the generator and  $\alpha$  for real training sequences  $(S_{1:t}, S_{t+1})$ , where  $\alpha=1^-$  to avoid saturation.

The adversarial regularization term for our predictive model (i.e. the “generator”) then takes the following form:

$$\mathcal{L}_{\text{adv}}(\hat{S}_{t+1}, S_{t+1}) = \lambda |\sigma(\mathcal{D}_\Theta(S_{1:t}, \hat{S}_{t+1})) - \alpha|. \quad (5)$$

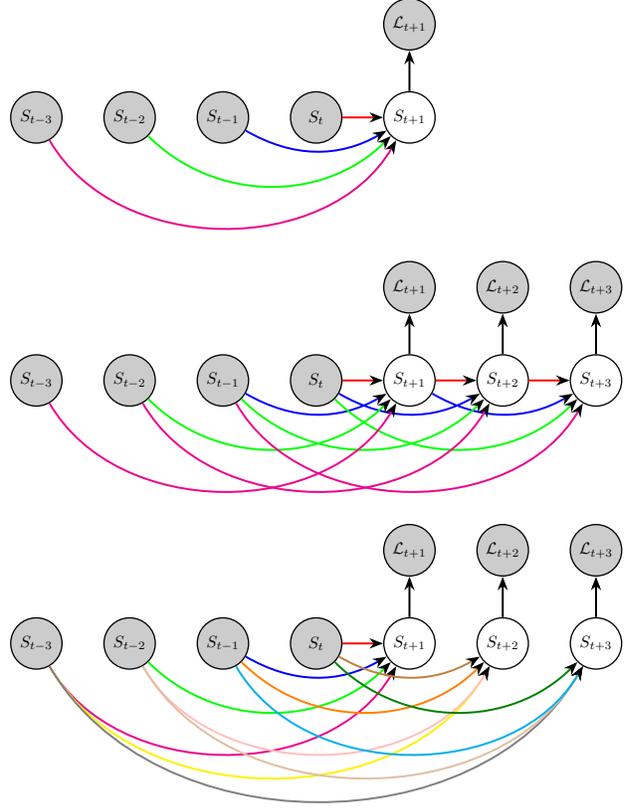


Figure 3: Illustration of the model for single time-step prediction (top), and extensions for multiple time-step prediction: the autoregressive model (middle), and the batch model (bottom). The autoregressive model shares parameters over time; dependency links are colored accordingly.

The structure of the discriminator network is derived from the two-scale architecture described above. The coarse-scale subnetwork has a single convolutional layer  $128 \times 3 \times 3$ , followed by three fully connected layers with 512, 256 and 1 hidden units respectively. The fine-scale subnetwork consists of three convolutional layers ( $128 \times 3 \times 3$ ,  $128 \times 3 \times 3$ ,  $256 \times 3 \times 3$ ) and three fully connected layers with 512, 256, 1 hidden units.

Following [2], we employ clipping of the discriminator weights  $\Theta$  to the range  $[-0.01, 0.01]$  after each gradient update, and set the target coefficient  $\alpha=0.9$  to prevent saturation. In our setting, every iteration of the discriminator training is followed by a single update of the generator parameters. We found that  $\lambda=0.1$  provides the optimal balance between the loss terms.

### 3.2. Predicting deeper into the future

We consider two extensions of the previous models to predict further into the future than a single frame. The first is to expand the output of the network to comprise a batch

of  $m$  frames, i.e. to output  $X_{t+1:t+m}$  and/or  $S_{t+1:t+m}$ . We refer to this as the “batch” approach. The drawback of this approach is that it ignores the recurrence structure of the problem. That is, the fact that  $S_{t+1}$  depends on  $S_{1:t}$  in the same manner as  $S_{t+2}$  depends on  $S_{2:t+1}$ . As a result, the capacity of the model is split to predict the  $m$  output frames, and the number of parameters in the last layer scales linearly with the number of output frames.

In our second approach, we leverage the recurrence property, and iteratively apply a model that predicts a single step into the future, using its prediction for time  $t + 1$  as an input to predict at time  $t + 2$ . This allows us to predict arbitrarily far into the future in an autoregressive manner, without resources scaling with the number of time-steps we want to predict. We refer to this approach as “autoregressive”. See Figure 3 for a schematic illustration of the single time-step model, and the two extensions for longer term prediction.

For the autoregressive mode, we either use the models as trained to predict one time-step ahead as they are, or we fine-tune these models by taking into account the impact the autoregressive approach has on predictions farther away than a single frame. In the latter case, during training we first make a forward pass, predicting one frame ahead at a time, and using the most recent outputs to predict the next time step. We then back-propagate the gradients through time [33], where the gradients w.r.t. the prediction at time  $t$  are based on the loss at time  $t$ , and the impact on the losses at later time steps.

## 4. Experiments

Before presenting our experimental results, we first describe the dataset and evaluation metrics in Section 4.1. We then present results on single-frame prediction, mid-term prediction (0.5 sec.), and long-term prediction (10 sec.).

### 4.1. Dataset and evaluation metrics

The Cityscapes dataset [4] contains 2,975 training, 500 validation and 500 testing video sequences of 1.8 second. Each sequence consists of 30 frames, and a ground-truth semantic segmentation is available for the 20-th frame. We measure the performance of our models on the Cityscape validation set, and refer to the supplementary material for results on the test set.

We assess performance using the standard mean Intersection over Union (IoU) measure, computed w.r.t. the ground truth segmentation of the 20-th frame in each sequence (IoU GT). We also compute the IoU measure w.r.t. the segmentation produced using the Dilated-10 network [36] for the 20-th frame (IoU SEG). The segmentation outputs of the Dilated-10 network are produced at a resolution of  $256 \times 128$  pixels. The IoU SEG metric allows us to validate our models w.r.t. the target segmentations from which they are trained. Finally, we compute the mean IoU

across categories that can move in the scene: *person, rider, car, truck, bus, train, motorcycle, and bicycle* (IoU-MO, for “moving objects”).

To evaluate the quality of the frame RGB predictions, we compute the Peak Signal to Noise Ratio (PSNR) and the Structural Similarity Index Measure (SSIM) measures [32]. The SSIM measures similarity between two images, ranging between -1 for very dissimilar inputs to 1 when the inputs are the same. It is based on comparing local patterns of pixel intensities normalized for luminance and contrast.

Unless specified otherwise, we train our models using a frame rate of 3, and taking 4 frames and/or segmentations as input. That is, the input sequence consists of frames  $\{X_{t-9}, X_{t-6}, X_{t-3}, X_t\}$ , and similarly for segmentations. We performed patch-wise training with  $64 \times 64$  patches for the largest scale resolution, enabling equal class frequency sampling as in [8], using mini-batches of four patches and a learning rate of 0.01.

### 4.2. Short-term prediction

In our first set of experiments we used frames 8, 11, 14, and 17 to predict frame 20, for which we have ground truth segmentation available. In Table 2 we compare our five models. For models that do not directly predict future segmentations, we generate segmentations using the Dilated-10 network based on the predicted RGB frames. We also include two baselines. The first baseline copies the last input frame to the output. For the second baseline we use FlowNet [6] to estimate the optical flow between the last two inputs, and warp the last input using the estimated flow.

From the result we can make several observations. First, in terms of RGB frame prediction (PSNR and SSIM), the performance is comparable for the three models X2X, XS2X, and XS2XS, and substantially better than the two baselines. This shows that our models learn about the scene dynamics in a non-trivial manner, and that adding semantic segmentations either at input and/or output does not have a substantial impact on this ability.

Second, in terms of segmentation prediction (IoU measures), the models that directly predict future segmentations (S2S, XS2S, XS2XS) perform much better than the models that only predict the RGB frames. This suggests that artifacts in the RGB frame predictions propagate into the Dilated-10 network, which then in turn gives degraded segmentations.

Third, the XS2XS model, which predicts both segmentations and RGB frames performs somewhat worse than the models that only predict segmentations (S2S and XS2S), suggesting that some of the modeling capacity is compromised by jointly predicting the RGB frames.

Finally, we find that fine-tuning the S2S model using adversarial training (S2S-adv) gives overall the best results, be it with a small margin over normal training. In Figure 4,

Method	PSNR	SSIM	IoU GT	IoU SEG	IoU-MO GT	IoU-MO SEG
Copy last input	20.6	0.65	49.4	54.6	43.4	48.2
Warp last input	20.9	0.67	50.4	55.5	44.9	49.8
Model X2X	24.0	<b>0.77</b>	23.0	22.3	12.8	11.4
Model S2S	—	—	<b>58.3</b>	64.9	53.8	59.8
Model S2S-adv.	—	—	<b>58.3</b>	<b>65.0</b>	<b>53.9</b>	<b>60.2</b>
Model XS2X	<b>24.2</b>	<b>0.77</b>	22.4	22.5	10.8	10.0
Model XS2S	—	—	58.2	64.6	53.7	59.9
Model XS2XS	24.0	0.76	55.5	61.1	50.7	55.8

Table 2: Single-frame prediction accuracy of our models (taking either RGB frames (X) and/or segmentations (S) as input and output), and baselines which either copy the last input to the output, or warp it to the output using optical flow.

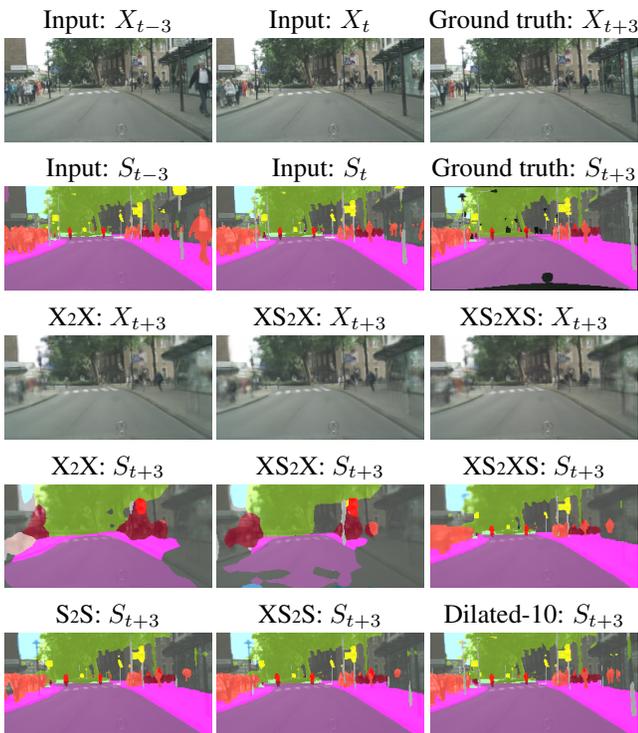


Figure 4: Short-term predictions of RGB frame  $X_{t+3}$  and segmentation  $S_{t+3}$  using our different models, compared to ground truth, and Dilated-10 oracle that has seen  $X_{t+3}$ .

we show qualitative results of the predictions for one of the validation sequences.

Table 3 presents results of an ablation study of the S2S model, assessing the impact of the different loss functions, as well as the impact of using one or two scales. We include the results obtained using the Dilated-10 model as an “oracle”, that predicts the future segmentation based on the future RGB frame, which is not accessible to our other models. This oracle result gives an idea of the maximum performance that could be expected, as it removes the difficulty of

Model	IoU GT	IoU SEG	IoU-MO GT
Dilated-10 oracle	68.8	100	64.7
S2S, 2 scales, $\ell_1$ +gdl	<b>58.3</b>	<b>64.9</b>	<b>53.8</b>
S2S, 1 scale, $\ell_1$ +gdl	57.7	63.9	52.6
S2S, 2 scales, $\ell_1$	57.6	64.0	53.2
S2S, 2 scales, MCE	55.5	60.9	49.7

Table 3: Ablation study with the S2S model, and also comparing to a Dilated-10 oracle that predicts the future segmentation using the future RGB frame as input.

future prediction and reduces to a classic semantic segmentation problem. All variants of the S2S model were trained during about 960,000 iterations, taking about four days of training on a single GPU. The results show that using two scales improves the performance, as does the addition of the gradient difference loss. Training with the  $\ell_1$  and/or gdl loss on the log-probabilities gives better results as compared to training using the Multi-class Cross-Entropy (MCE) loss on the segmentation labels. This is in line with observations made in network distillation [3, 11].

### 4.3. Mid-term prediction

We now address the more challenging task of predicting the mid-term future, *i.e.* the next 0.5 second. In these experiments we take in input frames 2, 5, 8, and 11, and predict outputs for frame 14, 17 and 20. We compare different strategies: batch models, autoregressive models (AR), and models with autoregressive fine-tuning (AR fine-tune). We compare these strategies to our two baselines consisting in copying the last input, and the second one relying on optical flow. For the optical flow baseline, after the first prediction, we also warp the flow field so that the flow is applied to the correct locations at the next time step, and so on. Qualitative prediction results are shown in Figure 5. For models XS2X and XS2S, the auto-regressive training mode is not employed because either the frame or the segmentation in-

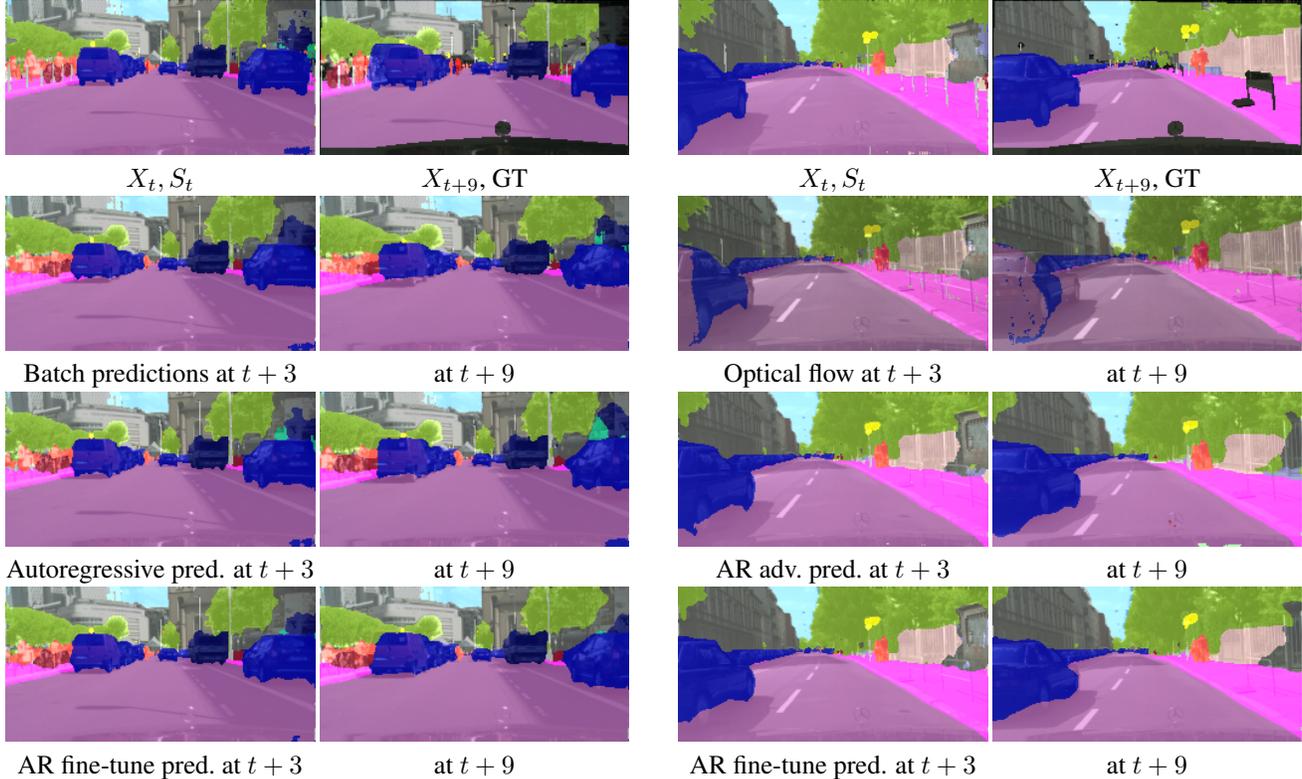


Figure 5: Autoregressive, batch predictions using the S2S model, and optical flow based baseline for two sequences (first sequence left, second sequence right). First row: last input and ground truth. Other rows show predictions overlaid with the true future frames. The full results are provided in the supplementary material.

Model	Frame 14		Frame 20	
	PSNR	SSIM	PSNR	SSIM
Copy last input	20.4	0.64	18.0	0.55
Warp last input	20.7	0.66	18.2	0.57
X2X, AR	<b>23.9</b>	<b>0.76</b>	19.2	0.61
XS2XS, AR	23.8	<b>0.76</b>	19.3	0.61
X2X, batch	23.8	<b>0.76</b>	20.6	<b>0.65</b>
XS2X, batch	<b>23.9</b>	<b>0.76</b>	<b>20.7</b>	<b>0.65</b>
XS2XS, batch	23.8	<b>0.76</b>	<b>20.7</b>	0.64

Table 4: Mid-term RGB frame prediction results.

put are missing for predicting from the second output.

The results for RGB frame prediction in Table 4, show that for frame 14 all models give comparable results, and consistently improve over the copy and warping baseline results. For frame 20, the batch models perform somewhat better than the autoregressive models. When predicting segmentations, *c.f.* Table 5, we find that the autoregressive models are better than the batch models. This is probably due to the fact that the single-step predictions are more accurate for segmentation, which makes them more suitable for autoregressive modeling. For RGB frame prediction, er-

Model	IoU GT	IoU SEG	IoU-MO GT
Copy last input	36.9	39.2	26.8
Warp last input	37.5	39.5	27.9
S2S, AR	45.3	47.2	36.4
S2S-adv, AR	45.1	47.2	37.3
S2S, AR, fine-tune	<b>46.7</b>	<b>49.7</b>	<b>39.3</b>
XS2XS, AR	39.3	40.8	27.4
S2S, batch	42.1	44.2	32.8
XS2S, batch	42.3	44.6	33.1
XS2XS, batch	41.2	43.5	31.4

Table 5: Mid-term segmentation prediction for frame 20 using different models in batch and autoregressive mode.

rors accumulate quickly leading to degraded autoregressive predictions. Among the batch models, using the images as input (XS2S model) helps slightly. Predicting both the images and segmentation (XS2XS model) performs worst, the image prediction task presumably takes up resources otherwise available for the segmentation task.

The S2S is the most effective, as it can be applied in autoregressive mode, and outperforms the XS2XS in this

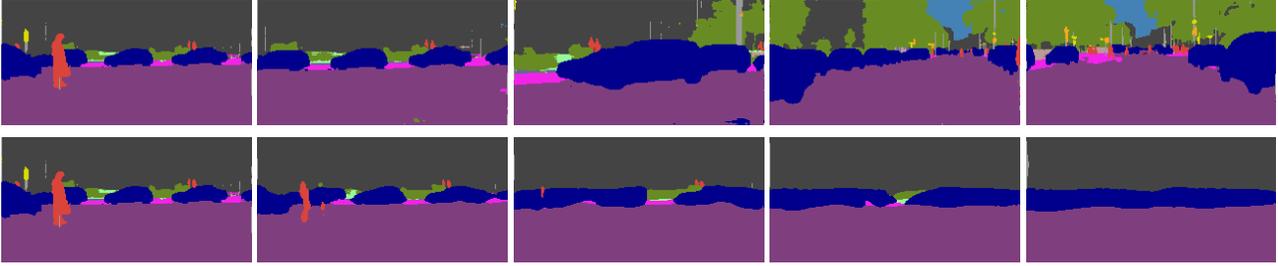


Figure 6: Last input segmentation, and ground truth segmentations at 1, 4, 7, and 10 seconds into the future (top row), and corresponding predictions of the autoregressive S2S model trained with fine-tuning (bottom row).

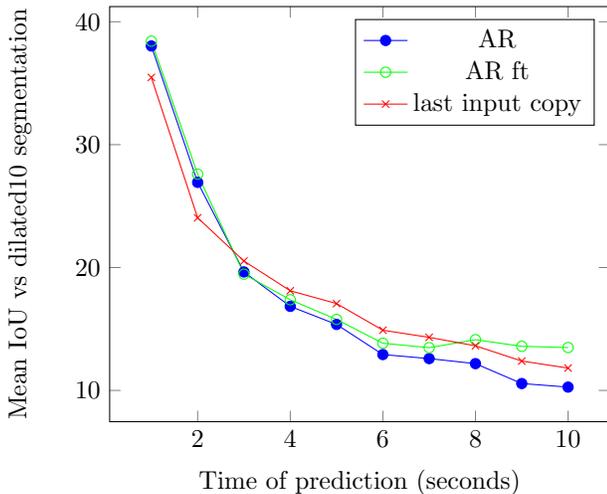


Figure 7: Mean IoU SEG of long-term segmentation prediction for the AR and AR fine-tune S2S models.

setting. In Figure 5 we compare different versions of this model. Visually, the first sequence shows some improvements using the autoregressive fine-tuned model, by more accurately matching contours of the moving cars than the other strategies. The second sequence displays typical failures of the optical flow baseline, as well as some improvements of the adversarial fine-tuning mode on the car contours. More examples are present in the supplementary material, where we can observe that the most difficult cases for our method are dealing with occlusions and situations where the video recording vehicle is turning.

#### 4.4. Long-term prediction

Finally, we consider what happens if we run our best autoregressive model to predict longer sequences of up to 10 seconds into the future. In this experiment we applied our S2S model in autoregressive mode on ten sequences of 238 frames from the Frankfurt long movie of the Cityscapes validation dataset. Given four frame segmentations with a frame rate of 17 images, the model predicts the ten next

ones. Thus, in this setting the images are sampled roughly at 1 Hz. In Figure 7 we report the IoU SEG performance as a function of time. In this extremely challenging setting the predictive performance quickly drops over time. Fine-tuning the model in autoregressive mode improves its performance, but only gives a clear advantage over the input-copy baseline for predictions at one and two seconds ahead. We also applied our model with a frame rate of 3 to predict up to 55 steps ahead, but found this to perform much worse. Figure 6 shows an example of predictions compared to the actual future segmentations. The visualization shows that our models average the different classes into an average future, which is perhaps not entirely surprising. Sampling different possible futures using a GAN or VAE approach could be an interesting way to resolve this issue.

## 5. Conclusion

We introduced a new visual understanding task of predicting future semantic segmentations. We explored five different models for this task relying on RGB and/or segmentations from previous frames. For prediction beyond a single future frame, we considered batch models that predict all future frames at once, and autoregressive models that sequentially predict the future frames. We found that autoregressive training produces the best results for our problem, and that models predicting in the segmentation space work better than those relying on the RGB frames.

While our results are encouraging, there is still room for improvement. Where the Dilated-10 network for semantic image segmentation gives around 69% IoU, this drops to about 58% when predicting 3 frames ahead (0.18 sec.), and to about 47% when for 9 frames (0.54 sec.). Most predicted object trajectories are reasonable, but simply do not always correspond to the actual observed trajectories. To improve results, our relatively basic convolutional network architectures can easily be extended to more sophisticated architectures such as residual nets, LSTMs, and dilated convolutions. Furthermore, GANs or VAE models with stochastic inputs may be useful to address the inherent uncertainty in the prediction of future segmentations.

We plan to open-source our Torch-based implementation upon paper acceptance, and invite the reader to watch videos of the predictions in the supplementary material.

**Acknowledgment.** This work has been partially supported by the grant ANR-16-CE23-0006 “Deep in France” and LabEx PERSYVAL-Lab (ANR-11-LABX-0025-01). We thank Michael Mathieu, Matthijs Douze, Hervé Jegou, Larry Zitnick, Mouhamadou Moustapha Cisse, and Gabriel Synnaeve for their precious comments.

## References

- [1] A. Alahi, V. Ramanathan, and L. Fei-Fei. Socially-aware large-scale crowd forecasting. In *CVPR*, 2014.
- [2] M. Arjovsky and L. Bottou. Towards principled methods for training generative adversarial networks. In *ICLR*, 2017.
- [3] L. Ba and R. Caruana. Do deep nets really need to be deep? In *NIPS*, 2014.
- [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. In *CVPR*, 2016.
- [5] J. Donahue, P. Krahenbuhl, and T. Darrell. Adversarial feature learning. In *ICLR*, 2017.
- [6] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. v.d. Smagt, D. Cremers, and T. Brox. FlowNet: Learning optical flow with convolutional networks. In *ICCV*, 2015.
- [7] A. Dosovitskiy and V. Koltun. Learning to act by predicting the future. In *ICLR*, 2017.
- [8] C. Farabet, C. Couprie, L. Najman, and Y. LeCun. Learning hierarchical features for scene labeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1915–1929, 2013.
- [9] C. Finn, I. Goodfellow, and S. Levine. Unsupervised learning for physical interaction through video prediction. In *NIPS*, 2016.
- [10] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial networks. In *NIPS*, 2014.
- [11] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning Workshop*, 2014.
- [12] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [13] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu. Spatial transformer networks. In *NIPS*, 2015.
- [14] X. Jin, X. Li, H. Xiao, X. Shen, Z. Lin, J. Yang, Y. Chen, J. Dong, L. Liu, Z. Jie, J. Feng, and S. Yan. Video scene parsing with predictive feature learning. *arXiv:1612.00119*, 2016.
- [15] N. Kalchbrenner, A. van den Oord, K. Simonyan, I. Danihelka, O. Vinyals, A. Graves, and K. Kavukcuoglu. Video pixel networks. *arXiv:1610.00527*, 2016.
- [16] D. Kingma and M. Welling. Auto-encoding variational Bayes. In *ICLR*, 2014.
- [17] K. Kitani, B. Ziebart, J. Bagnell, and M. Hebert. Activity forecasting. In *ECCV*, 2012.
- [18] G. Lin, A. Milan, C. Shen, and I. Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *CVPR*, 2017.
- [19] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015.
- [20] P. Luc, C. Couprie, S. Chintala, and J. Verbeek. Semantic segmentation using adversarial networks. In *NIPS Workshop on Adversarial Training*, 2016.
- [21] Z. Luo, B. Peng, D.-A. Huang, A. Alahi, and L. Fei-Fei. Unsupervised learning of long-term motion dynamics for videos. In *CVPR*, 2017.
- [22] M. Mathieu, C. Couprie, and Y. LeCun. Deep multi-scale video prediction beyond mean square error. In *ICLR*, 2016.
- [23] D. Nilsson and C. Sminchisescu. Semantic video segmentation by gated recurrent flow propagation. *arXiv:1612.08871*.
- [24] V. Patraucean, A. Handa, and R. Cipolla. Spatio-temporal video autoencoder with differentiable memory. In *ICLR Workshop*, 2016.
- [25] M. Ranzato, A. Szlam, J. Bruna, M. Mathieu, R. Collobert, and S. Chopra. Video (language) modeling: a baseline for generative models of natural videos. *arXiv:1412.6604*, 2014.
- [26] N. Srivastava, E. Mansimov, and R. Salakhutdinov. Unsupervised learning of video representations using LSTMs. In *ICML*, 2015.
- [27] R. Sutton and A. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- [28] J. Van Amersfoort, A. Kannan, M. Ranzato, A. Szlam, D. Tran, and S. Chintala. Transformation-based models of video sequences. *arXiv:1701.08435*, 2017.
- [29] C. Vondrick, P. Hamed, and A. Torralba. Anticipating the future by watching unlabeled video. In *CVPR*, 2016.
- [30] C. Vondrick, H. Pirsivash, and A. Torralba. Generating videos with scene dynamics. In *NIPS*, 2016.
- [31] J. Walker, C. Doersch, A. Gupta, and M. Hebert. An uncertain future: Forecasting from static images using variational autoencoders. In *ECCV*, 2016.
- [32] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.
- [33] P. Werbos. Generalization of backpropagation with application to a recurrent gas market model. *Neural Networks*, 1(4):339–356, 1988.
- [34] T. Xue, J. Wu, K. Bouman, and W. Freeman. Visual dynamics: Probabilistic future frame synthesis via cross convolutional networks. In *NIPS*, 2016.
- [35] J. Yang, A. Kannan, D. Batra, and D. Parikh. LR-GAN: Layered recursive generative adversarial networks for image generation. In *ICLR*, 2017.
- [36] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. In *ICLR*, 2016.
- [37] J. Zhao, M. Mathieu, and Y. LeCun. Energy-based generative adversarial networks. *ICLR*, 2017.