

## TL;DR

A new framework for generalizing image restoration networks on high-level vision tasks.

```

1 !pip install git+https://github.com/MKFMIKU/d2sm.git
2
3 from d2sm import DeST
4 dest = DeST("54", args.patch_size, args.patch_size, Q=args.dest_q, K=args.dest_k, mean_shift=True)
5
6 dest_loss = dest(self.hr, self.denoised_hr)

```

## Problem Settings

By maximizing the correspondence between each pair of the CNN-restored results and the clear image, a CNN is trained to restore images from the degraded domain into the clear domain.

Recent work called perceptual loss is shown to be able to lead to better visual quality by maximizing the correspondence in the semantic feature space of pre-trained classification network.

GANs, which employ a discriminator to implicitly enforce the distribution of restored images to be consistent with the distribution of clear images in terms of KL- and JS- divergence.

Whether it is possible to combine the pre-trained large-scale networks in an adversarial or statistical manner to bypass their drawbacks and avail their advantages together?

In this work, we exploit semantic features of pretrained classification networks and implicitly matches the probabilistic distribution of clear images at the semantic feature space.

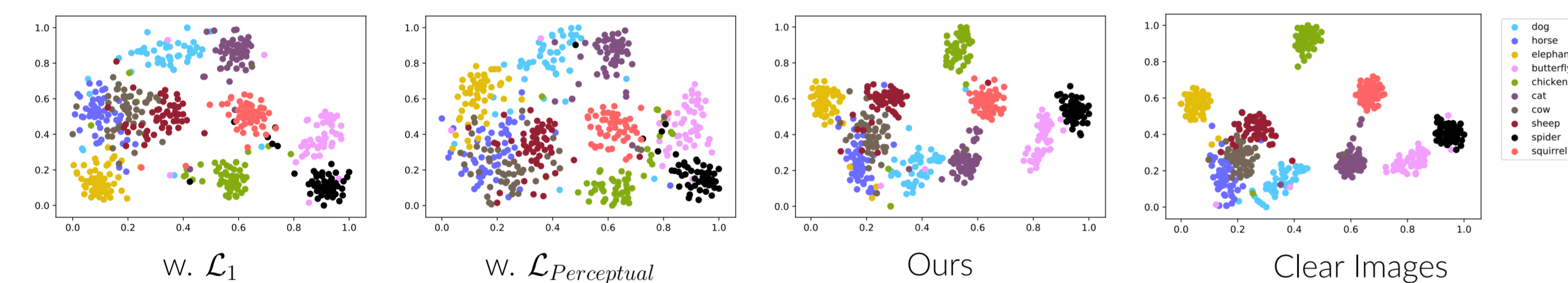


Figure 1. t-SNE of denoised Images in the Semantic Feature Space. Ours preserves most semantics like GT.

## Contributions

- A new image restoration framework that minimizes the distribution divergence instead of the sample-to-sample distance in the semantic feature space.
- A new patch-wise fashion that can decompose complex semantics of images for efficient distribution approximation.
- The method substantially outperforms the original perceptual loss and other SOTA losses, especially in high-level vision tasks that validates D2SM indeed preserve semantics.

## References

- [1] Meng Chang, Qi Li, Huajun Feng, and Zhihai Xu. Spatial-Adaptive Network for Single Image Denoising. In ECCV, 2020.
- [2] Shi Guo, Zifei Yan, Kai Zhang, Wangmeng Zuo, and Lei Zhang. Toward Convolutional Blind Denoising of Real Photographs. In CVPR, 2019.
- [3] Kai Zhang, Wangmeng Zuo, and Lei Zhang. FFDNet: Toward a fast and flexible solution for CNN-based image denoising. TIP, 2018.

## Method

Given  $N$  samples of image pairs that consist of  $T_x = \{x_1, x_2, \dots, x_N\}$  and  $T_y = \{y_1, y_2, \dots, y_N\}$ , we incorporate the mutual information of them in the feature space of  $\Phi(G(\cdot))$  (e.g. pretrained VGG-19 network) into the restoration learning.

By minimizing the divergence of the estimated probability distribution between samples  $T_x$  in  $\Phi(G(\cdot))$  and  $T_y$  in  $\Phi(\cdot)$ , denoted as  $\mathcal{G}'$  and  $\mathcal{G}$ , we force  $G(\cdot)$  to better maintain the geometry of the feature space  $\Phi(\cdot)$  estimated in the clear image domain  $\mathcal{Y}$  as

$$\mathcal{L}(T_x, T_y, G) = \sum_{i=1}^N \sum_{j=1, j \neq i}^N g'_{ji} \log\left(\frac{g'_{ji}}{g_{ji}}\right). \quad (1)$$

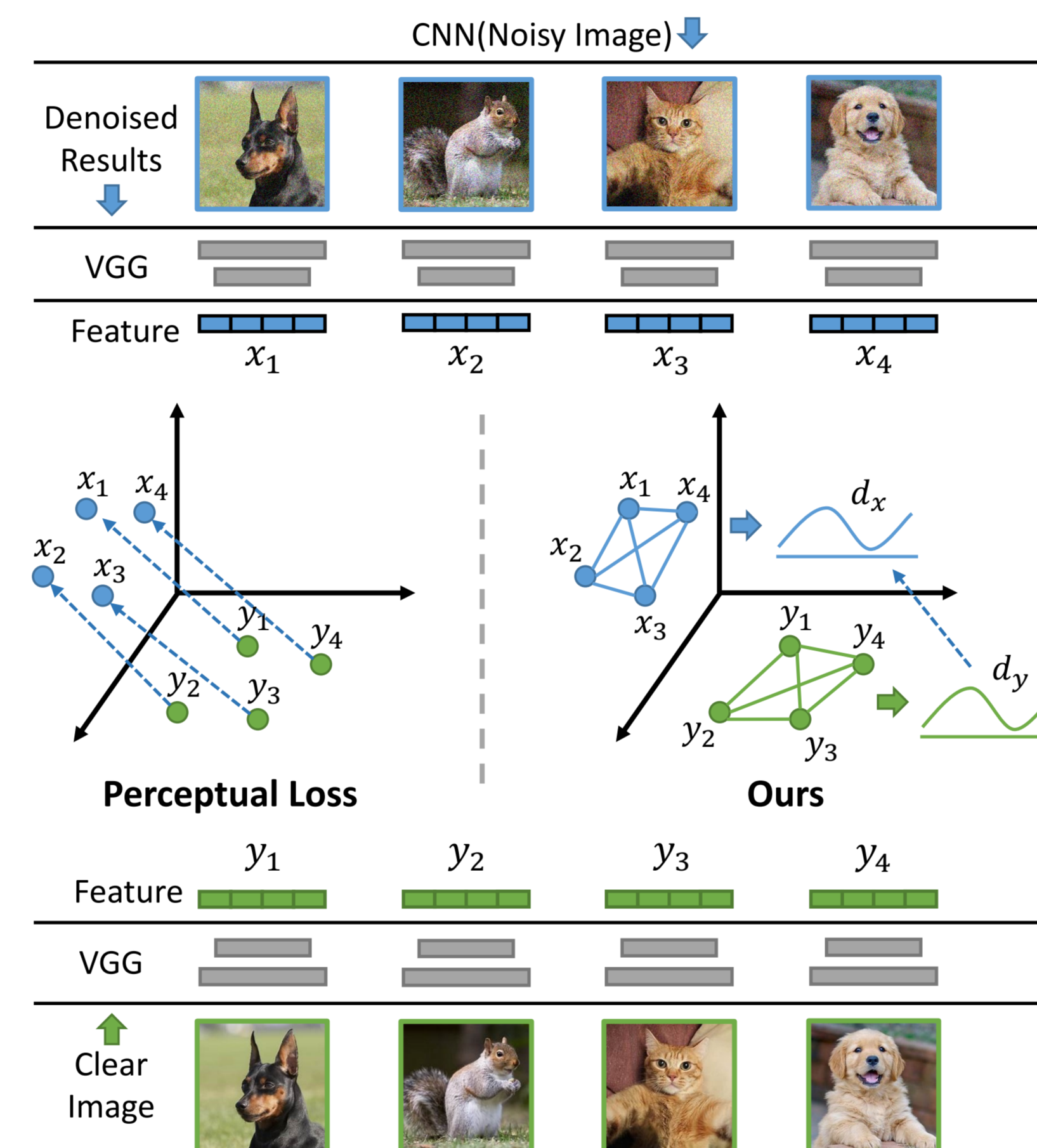


Figure 2. Perceptual loss vs. Ours.

To bypass the memory limitation, we introduce a Memorized Historic Sampling strategy by maintaining two queues of feature samples, i.e.,  $Q^X$  and  $Q^Y$  that can store historic features from previous mini-batches with limited GPU memory cost as

$$g'_{ij} = \frac{K_{\cosine}(Q_i^X, Q_j^X)}{\sum_{k=1, k \neq j}^q K_{\cosine}(Q_k^X, Q_j^X)} \in [0, 1], \quad (2)$$

and

$$Q_{1 \dots N}^X, Q_{N \dots q}^X \leftarrow f_{\{1 \dots N\}}^x, Q_{\{1 \dots q-N\}}^X. \quad (3)$$

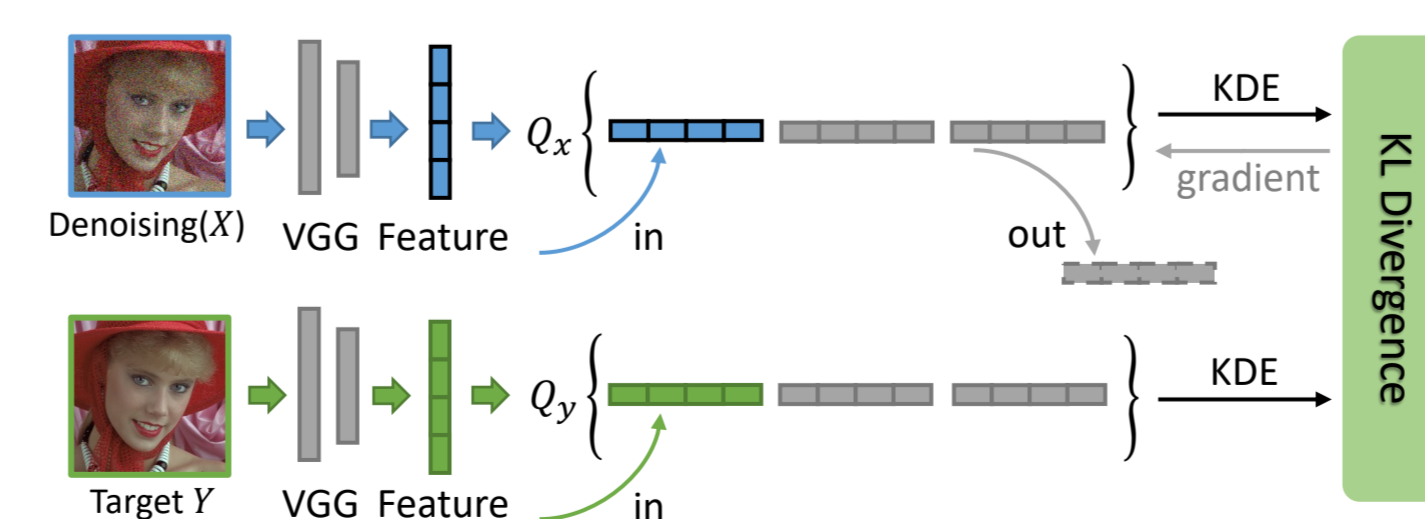


Figure 3. Sampling with Historic Gradients.

## Results

### Cityscape Denoising and Segmentation

We conduct complementary denoising and segmentation experiments on the Cityscapes dataset. We then measure the semantic segmentation accuracy on restored images in the term of Mean Intersection-over-Union (MIoU) in 19 pre-defined semantic classes.

Table 1. Quantitative performance comparison on the cityscape denoising and segmentation.

Method (Backbone)	Objective	Noise-Level $\sigma=25$			Noise-Level $\sigma=35$			Noise-Level $\sigma=50$		
		PSNR $\uparrow$	SSIM $\uparrow$	MIoU (%) $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$	MIoU (%) $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$	MIoU (%) $\uparrow$
FFDNet [3]	$\mathcal{L}_1$	35.033 <sup>(6)</sup>	0.925 <sup>(6)</sup>	0.605 <sup>(8)</sup>	34.074 <sup>(6)</sup>	0.912 <sup>(6)</sup>	0.537 <sup>(8)</sup>	32.845 <sup>(6)</sup>	0.895 <sup>(6)</sup>	0.451 <sup>(7)</sup>
	+ $\mathcal{L}_{SSIM}$	35.567 <sup>(3)</sup>	0.935 <sup>(2)</sup>	0.642 <sup>(2)</sup>	34.469 <sup>(4)</sup>	0.922 <sup>(2)</sup>	0.584 <sup>(2)</sup>	33.180 <sup>(3)</sup>	0.906 <sup>(2)</sup>	0.450 <sup>(8)</sup>
	+ $\mathcal{L}_{Perceptual}$	34.319 <sup>(7)</sup>	0.912 <sup>(7)</sup>	0.629 <sup>(4)</sup>	33.486 <sup>(7)</sup>	0.899 <sup>(7)</sup>	0.582 <sup>(4)</sup>	32.383 <sup>(7)</sup>	0.881 <sup>(7)</sup>	0.509 <sup>(2)</sup>
	+ $\mathcal{L}_{LPIPS}$	35.551 <sup>(4)</sup>	0.929 <sup>(4)</sup>	0.613 <sup>(6)</sup>	34.463 <sup>(5)</sup>	0.916 <sup>(4)</sup>	0.541 <sup>(7)</sup>	33.138 <sup>(5)</sup>	0.899 <sup>(4)</sup>	0.452 <sup>(6)</sup>
	+ $\mathcal{L}_{Contextual}$ + $\mathcal{L}_{CrossEntropy}$	25.115 <sup>(8)</sup>	0.762 <sup>(8)</sup>	0.628 <sup>(8)</sup>	24.938 <sup>(8)</sup>	0.758 <sup>(8)</sup>	0.583 <sup>(3)</sup>	24.775 <sup>(8)</sup>	0.753 <sup>(8)</sup>	0.509 <sup>(2)</sup>
D2SM (Ours)	w/o. Internal	35.913 <sup>(2)</sup>	0.932 <sup>(3)</sup>	0.630 <sup>(3)</sup>	34.800 <sup>(2)</sup>	0.919 <sup>(3)</sup>	0.565 <sup>(5)</sup>	33.477 <sup>(2)</sup>	0.903 <sup>(3)</sup>	0.491 <sup>(4)</sup>
	w/. Internal	<b>35.543<sup>(3)</sup></b>	<b>0.929<sup>(4)</sup></b>	<b>0.612<sup>(7)</sup></b>	<b>34.475<sup>(3)</sup></b>	<b>0.916<sup>(4)</sup></b>	<b>0.546<sup>(6)</sup></b>	<b>33.167<sup>(4)</sup></b>	<b>0.899<sup>(4)</sup></b>	<b>0.463<sup>(5)</sup></b>
CBDNet [2]	-	36.152 <sup>(3)</sup>	0.936 <sup>(2)</sup>	0.655 <sup>(3)</sup>	34.964 <sup>(3)</sup>	0.923 <sup>(3)</sup>	0.599 <sup>(3)</sup>	33.613 <sup>(3)</sup>	0.907 <sup>(3)</sup>	0.539 <sup>(3)</sup>
	w/o. Internal	36.254 <sup>(2)</sup>	0.935 <sup>(3)</sup>	0.679 <sup>(2)</sup>	35.158 <sup>(2)</sup>	0.925 <sup>(2)</sup>	0.631 <sup>(2)</sup>	33.904 <sup>(2)</sup>	0.911 <sup>(2)</sup>	0.550 <sup>(2)</sup>
SADNet [1]	w/. Internal	<b>36.899<sup>(1)</sup></b>	<b>0.941<sup>(1)</sup></b>	<b>0.691<sup>(1)</sup></b>	<b>35.596<sup>(1)</sup></b>	<b>0.929<sup>(1)</sup></b>	<b>0.652<sup>(1)</sup></b>	<b>34.172<sup>(1)</sup></b>	<b>0.914<sup>(1)</sup></b>	<b>0.600<sup>(1)</sup></b>
	-	36.310 <sup>(3)</sup>	0.936 <sup>(3)</sup>	0.674 <sup>(3)</sup>	35.081 <sup>(3)</sup>	0.924 <sup>(2)</sup>	0.637 <sup>(3)</sup>	33.730 <sup>(3)</sup>	0.908 <sup>(3)</sup>	0.581 <sup>(3)</sup>
	w/o. Internal	36.822 <sup>(2)</sup>	0.940 <sup>(2)</sup>	0.691 <sup>(2)</sup>	35.247 <sup>(2)</sup>	0.924 <sup>(2)</sup>	0.655 <sup>(2)</sup>	34.133 <sup>(2)</sup>	0.912 <sup>(2)</sup>	0.600 <sup>(2)</sup>
	w/. Internal	<b>37.130<sup>(1)</sup></b>	<b>0.943<sup>(1)</sup></b>	<b>0.701<sup>(1)</sup></b>	<b>35.839<sup>(1)</sup></b>	<b>0.931<sup>(1)</sup></b>	<b>0.670<sup>(1)</sup></b>	<b>34.440<sup>(1)</sup></b>	<b>0.916<sup>(1)</sup></b>	<b>0.634<sup>(1)</sup></b>

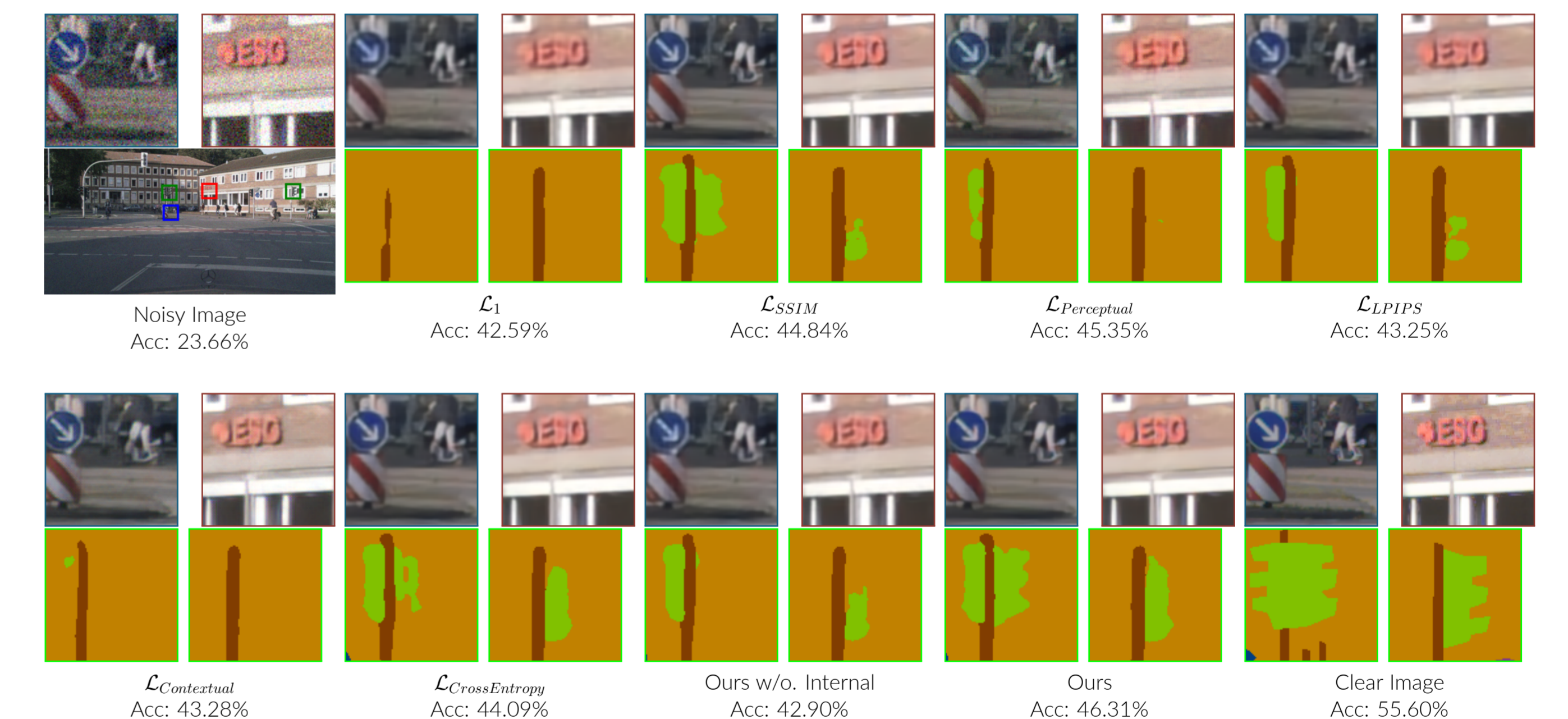


Figure 4. Qualitative comparison on the denoising and segmentation results. Ours preserves most of the semantic details, including the human shape and font edge in the highlighted area. Additionally, in the shown segmentation results, our result is the only one that can be successfully recognized into traffic light.

### Results Analysis

Here we present the quantitative comparison with the three additional divergence estimation or objectives in Cityscapes. The convergence curve visualized below further demonstrates that our proposed method significantly accelerates convergence memorized historic sampling.

Method	$\sigma$	PSNR $\uparrow$	SSIM $\uparrow$	MIoU (%) $\uparrow$
FFDNet	25	35.03	0.925	0.605
+ $\mathcal{L}_{iKLD}$	25	35.97	0.931	0.638
+ $\mathcal{L}_{JSD}$	25	36.31	0.935	0.640
+ $\mathcal{L}_{GAN}$	25	35.55	0.931	0.621
Ours	25	<b>36.45</b>	<b>0.936</b>	<b>0.644</b>

Table 2. Performance comparison with different distribution divergence.

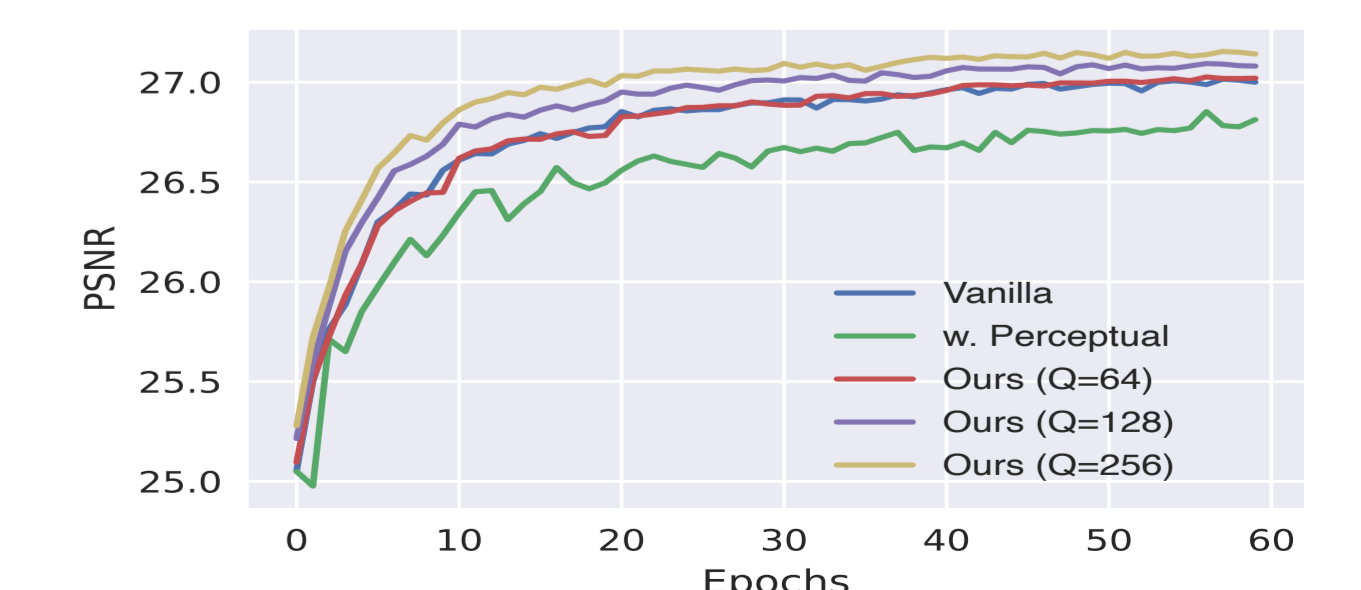


Figure 5. Convergence visualization between different queue size.