

Deep Image Prior

Supplementary Material

In this supplemental material we provide more details about the setup in each experiment. We also extend the figures from the main paper with more qualitative comparisons. The code for all experiments is published at <https://github.com/DmitryUlyanov/deep-image-prior>. Interactive comparisons are available at the project page https://dmitryulyanov.github.io/deep_image_prior.

1. Architectures

In general we found that deep architectures *without* too many short paths from input to output define very good deep priors. Having too many short paths prevent learning regular image patterns, resulting in memorization at the pixel level and deficient image priors.

While other options are possible, we mainly experimented with fully-convolutional architectures, where the input $z \in \mathbb{R}^{C' \times W \times H}$ has the same spatial resolution as the the output of the network $f_\theta(z) \in \mathbb{R}^{3 \times W \times H}$.

We use encoder-decoder (“hourglass”) architecture (possibly with skip-connections) for f_θ in all our experiments except noted otherwise (fig. 1), varying small number of hyperparameters. Although the best results can be achieved by carefully tuning an architecture for a particular task (and potentially for a particular image), we found that wide range of hyperparameters and architectures give acceptable results.

We use LeakyReLU [4] as a non-linearity. As a down-sampling technique we simply use strides implemented within convolution modules. We also tried average/max pooling and downsampling with Lanczos kernel, but did not find a consistent difference between any of them. As an up-sampling operation we choose between bilinear upsampling and nearest neighbour upsampling. An alternative upsampling method could be to use transposed convolutions, but the results we obtained using them were worse. We use reflection padding instead of zero padding in convolution layers everywhere except for the feature inversion experiment.

We considered two ways to create the input z : 1. *random*, where the z is filled with uniform noise between zero and 0.1, 2. *meshgrid*, where we initialize $z \in \mathbb{R}^{2 \times W \times H}$ using `np.meshgrid` (see fig. 2). Such initialization serves as an additional smoothness prior to the one imposed by the

structure of f_θ itself. We found such input to be beneficial for large-hole inpainting, but not for other tasks.

During fitting of the networks we sometimes use a noise-based regularization. I.e. at each iteration we perturb the input z with an additive normal noise. While we have found such regularization to impede optimization process, we also observed that the network was able to optimize its objective to zero no matter the variance of the additive noise (i.e. the network was always able to adapt to any reasonable variance for sufficiently large number of optimization steps).

We found the optimization process tends to destabilize as the loss goes down and approaches a certain value. Destabilization is observed as significant loss increase and the blur in generated image $f_\theta(z)$. From such destabilization point the loss goes down again till destabilized one more time. To remedy this issue we simply track the optimization loss and return to parameters from the previous iteration if the loss difference between two consecutive iterations is higher than a certain threshold.

Finally, we use *ADAM* optimizer [5] in all our experiments. All experiments are implemented in PyTorch.

Below, we provide the remaining details of the network architectures. We use the notation introduced in fig. 1.

After that, a large number of additional experimental results are provided. *Electronic zoom-in is recommended for almost all figures.*

Super-resolution (fig. 1 and 5 of main paper) (default architecture).

$$z \in \mathbb{R}^{32 \times W \times H} \sim U(0, \frac{1}{10})$$

$$n_u = n_d = [128, 128, 128, 128, 128]$$

$$k_u = k_d = [3, 3, 3, 3, 3]$$

$$n_s = [4, 4, 4, 4, 4]$$

$$k_s = [1, 1, 1, 1, 1]$$

$$\sigma_p = \frac{1}{30}$$

$$\text{num_iter} = 2000$$

$$\text{LR} = 0.01$$

$$\text{upsampling} = \text{bilinear}$$

The decimation operator d is composed of low pass filtering operation using Lanczos2 kernel (see [9]) and resampling, all implemented as a single (fixed) convolutional layer.

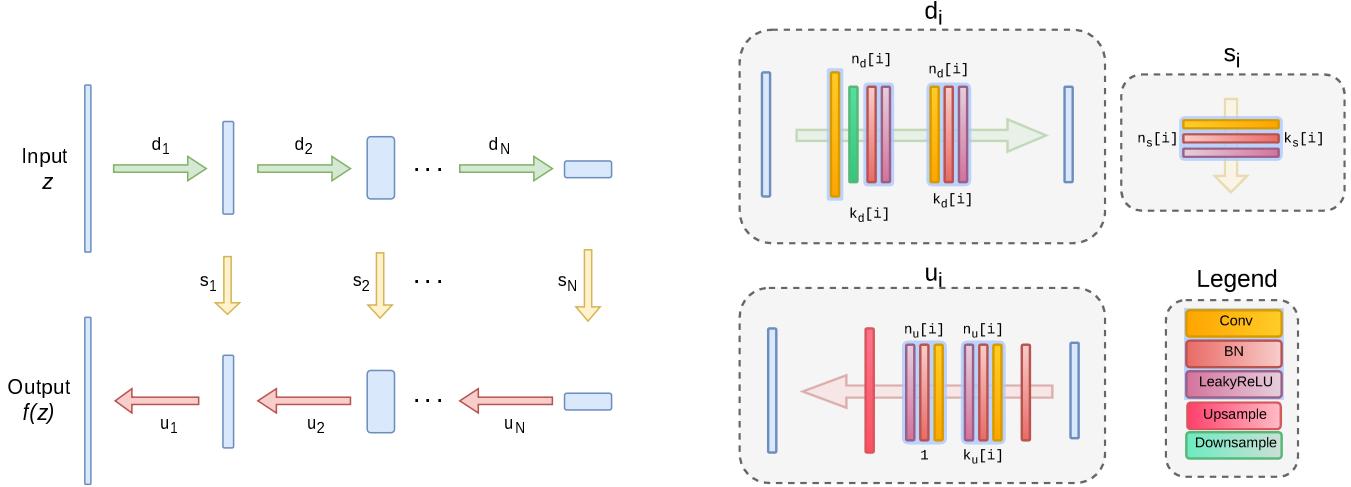


Figure 1: **The architecture used in the experiments.** We use “hourglass” (also known as “decoder-encoder”) architecture. We sometimes add skip connections (yellow arrows). $n_u[i]$, $n_d[i]$, $n_s[i]$ correspond to the number of filters at depth i for the upsampling, downsampling and skip-connections respectively. The values $k_u[i]$, $k_d[i]$, $k_s[i]$ correspond to the respective kernel sizes.

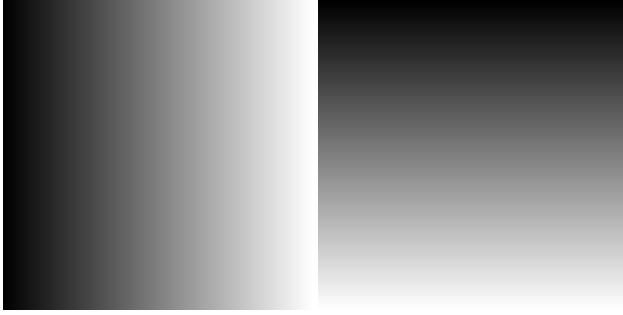


Figure 2: **“Meshgrid” input z** used in some inpainting experiments. These are two channels of the input tensor; in BCWH layout: $z[0, 0, :, :]$, $z[0, 1, :, :]$. The intensity encodes the value: from zero (black) to one (white). Such type of input can be regarded as a part of the prior which enforces smoothness.

For 8x super-resolution (9) we have changed the standard deviation of the input noise to $\sigma_p = \frac{1}{20}$ and the number of iterations to 4000.

Text inpainting (fig. 7 top row of main paper). We used the same hyperparameters as for super-resolution but optimized the objective for 6000 iterations.

Large hole inpainting (fig. 6 of main paper). We used the same hyperparameters as for super-resolution, but used meshgrid as an input, removed skip connections and optimized for 5000 iterations.

Large hole inpainting (fig. 8 of main paper). We used the following hyperparameters:

```

 $z \in \mathbb{R}^{32 \times W \times H} \sim U(0, \frac{1}{10})$ 
 $n_u = n_d = [16, 32, 64, 128, 128, 128]$ 
 $k_d = [3, 3, 3, 3, 3, 3]$ 
 $k_u = [5, 5, 5, 5, 5, 5]$ 
 $n_s = [0, 0, 0, 0, 0, 0]$ 
 $k_s = [\text{NA}, \text{NA}, \text{NA}, \text{NA}, \text{NA}, \text{NA}]$ 
 $\sigma_p = 0$ 
num_iter = 5000
LR = 0.1
upsampling = nearest

```

In figures 8(c) and 8(d) we simply sliced off last layers to get smaller depth.

Denoising. Hyperparameters were set to be the same as in the case of super-resolution with only difference in iteration number, which was set to 1800. We used the following implementations of referenced denoising methods: [6] for CBM3D and [1] for NLM. We used exponential sliding window with weight $\gamma = 0.99$. **JPEG artifacts removal (fig. 3).** Although we could use the same setup as in other denoising experiments, the hyperparameters we used to generate the image in figure 3 were the following:

```

 $z \in \mathbb{R}^{3 \times W \times H} \sim U(0, \frac{1}{10})$ 
 $n_u = n_d = [8, 16, 32, 64, 128]$ 
 $k_u = k_d = [3, 3, 3, 3, 3]$ 
 $n_s = [0, 0, 0, 4, 4]$ 
 $k_s = [\text{NA}, \text{NA}, \text{NA}, 1, 1]$ 
 $\sigma_p = \frac{1}{30}$ 
num_iter = 2400
LR = 0.01
upsampling = bilinear

```

Image reconstruction (fig. 7 bottom). We used the same setup as in the case of superresolution and denoising, but set num_iter = 11000, LR = 0.001. **Alexnet inversion.**

```

 $z \in \mathbb{R}^{32 \times W \times H} \sim U(0, \frac{1}{10})$ 
 $n_u = n_d = [16, 32, 64, 128, 128, 128]$ 
 $k_u = k_d = [7, 7, 5, 5, 3, 3]$ 
 $n_s = [4, 4, 4, 4, 4]$ 
 $k_s = [1, 1, 1, 1, 1]$ 
num_iter = 3100
LR = 0.001
upsampling = nearest

```

We used num_iter = 10000 for the VGG inversion experiment (16)

References

- [1] A. Buades. NLM demo. [http://demo.ipol.im/
demo/bcm_non_local_means_denoising/](http://demo.ipol.im/demo/bcm_non_local_means_denoising/). 2
- [2] A. Dosovitskiy and T. Brox. Inverting convolutional networks with convolutional networks. In *Proc. CVPR*, 2016. 14
- [3] D. Glasner, S. Bagon, and M. Irani. Super-resolution from a single image. In *ICCV*, pages 349–356. IEEE Computer Society, 2009. 6, 7, 8
- [4] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015. 1
- [5] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014. 1
- [6] M. Lebrun. BM3D code. [https://github.com/
gfaccioli/bm3d](https://github.com/gfaccioli/bm3d). 2
- [7] A. Mahendran and A. Vedaldi. Understanding deep image representations by inverting them. In *Proc. CVPR*, 2015. 14
- [8] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. 14
- [9] K. Turkowski. Filters for common resampling-tasks. *Graphics gems*, pages 147–165, 1990. 1

	Baboon	Barbara	Bridge	Coastguard	Comic	Face	Flowers	Foreman	Lenna	Man	Monarch	Pepper	Ppt3	Zebra
No prior	22.24	24.89	23.94	24.62	21.06	29.99	23.75	29.01	28.23	24.84	25.76	28.74	20.26	21.69
Bicubic	22.44	25.15	24.47	25.53	21.59	31.34	25.33	29.45	29.84	25.7	27.45	30.63	21.78	24.01
TV prior	22.34	24.78	24.46	25.78	21.95	31.34	25.91	30.63	29.76	25.94	28.46	31.32	22.75	24.52
Glasner et al.	22.44	25.38	24.73	25.38	21.98	31.09	25.54	30.4	30.48	26.33	28.22	32.02	22.16	24.34
Ours	22.29	25.53	24.38	25.81	22.18	31.02	26.14	31.66	30.83	26.09	29.98	32.08	24.38	25.71
SRResNet-MSE	23.0	26.08	25.52	26.31	23.44	32.71	28.13	33.8	32.42	27.43	32.85	34.28	26.56	26.95
LapSRN	22.83	25.69	25.36	26.21	22.9	32.62	27.54	33.59	31.98	27.27	31.62	33.88	25.36	26.98

Table 1: Detailed super-resolution PSNR comparison on Set14 dataset with factor $4x$. Note that our method does not aim to maximize PSNR.

	Baboon	Barbara	Bridge	Coastguard	Comic	Face	Flowers	Foreman	Lenna	Man	Monarch	Pepper	Ppt3	Zebra
No prior	21.09	23.04	21.78	23.63	18.65	27.84	21.05	25.62	25.42	22.54	22.91	25.34	18.15	18.85
Bicubic	21.28	23.44	22.24	23.65	19.25	28.79	22.06	25.37	26.27	23.06	23.18	26.55	18.62	19.59
TV prior	21.3	23.72	22.3	23.82	19.5	28.84	22.5	26.07	26.74	23.53	23.71	27.56	19.34	19.89
SelfExSR	21.37	23.9	22.28	24.17	19.79	29.48	22.93	27.01	27.72	23.83	24.02	28.63	20.09	20.25
Ours	21.38	23.94	22.2	24.21	19.86	29.52	22.86	27.87	27.93	23.57	24.86	29.18	20.12	20.62
LapSRN	21.51	24.21	22.77	24.1	20.06	29.85	23.31	28.13	28.22	24.2	24.97	29.22	20.13	20.28

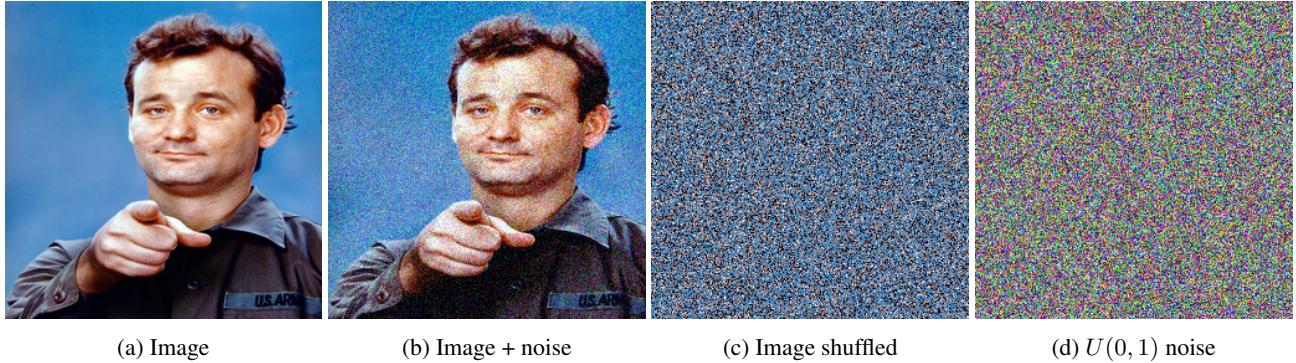
Table 2: Detailed super-resolution PSNR comparison on Set14 dataset with factor $8x$.

	Baby	Bird	Butterfly	Head	Woman
No prior	30.16	27.67	19.82	29.98	25.18
Bicubic	31.78	30.2	22.13	31.34	26.75
TV prior	31.21	30.43	24.38	31.34	26.93
Glasner et al.	32.24	31.1	22.36	31.69	26.85
Ours	31.49	31.8	26.23	31.04	28.93
LapSRN	33.55	33.76	27.28	32.62	30.72
SRResNet-MSE	33.66	35.1	28.41	32.73	30.6

Table 3: Detailed super-resolution PSNR comparison on Set5 dataset with factor $4x$.

	Baby	Bird	Butterfly	Head	Woman
No prior	26.28	24.03	17.64	27.94	21.37
Bicubic	27.28	25.28	17.74	28.82	22.74
TV prior	27.93	25.82	18.4	28.87	23.36
SelfExSR	28.45	26.48	18.8	29.36	24.05
Ours	28.28	27.09	20.02	29.55	24.5
LapSRN	28.88	27.1	19.97	29.76	24.79

Table 4: Detailed super-resolution PSNR comparison on Set5 dataset with factor $8x$.



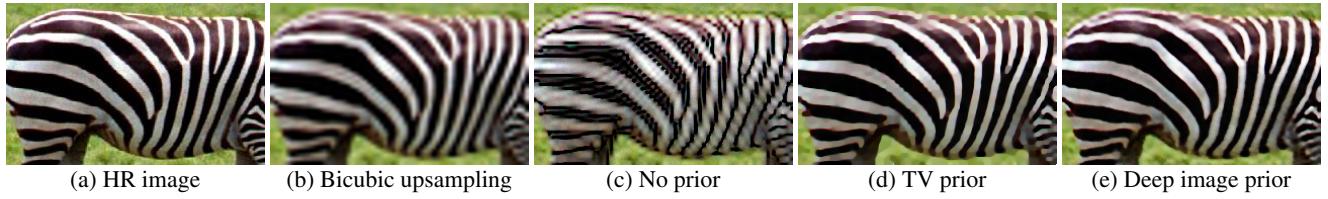
(a) Image

(b) Image + noise

(c) Image shuffled

(d) $U(0, 1)$ noise

Figure 3: Images used in fig. 2 of the main paper to demonstrate the noise impedance effect.



(a) HR image

(b) Bicubic upsampling

(c) No prior

(d) TV prior

(e) Deep image prior

Figure 4: **Prior effect in super-resolution.** (c):Direct optimization of data term $E(x; x_0)$ with respect to the pixels leads to ringing artifacts. TV prior removes ringing artifacts but introduces cartoon effect. Deep prior leads to the result that is both clean and sharp.Figure 5: **Inpainting diversity.** Left: original image (black pixels indicate holes). The remaining four images show results obtained using deep prior corresponding to different input vectors z .

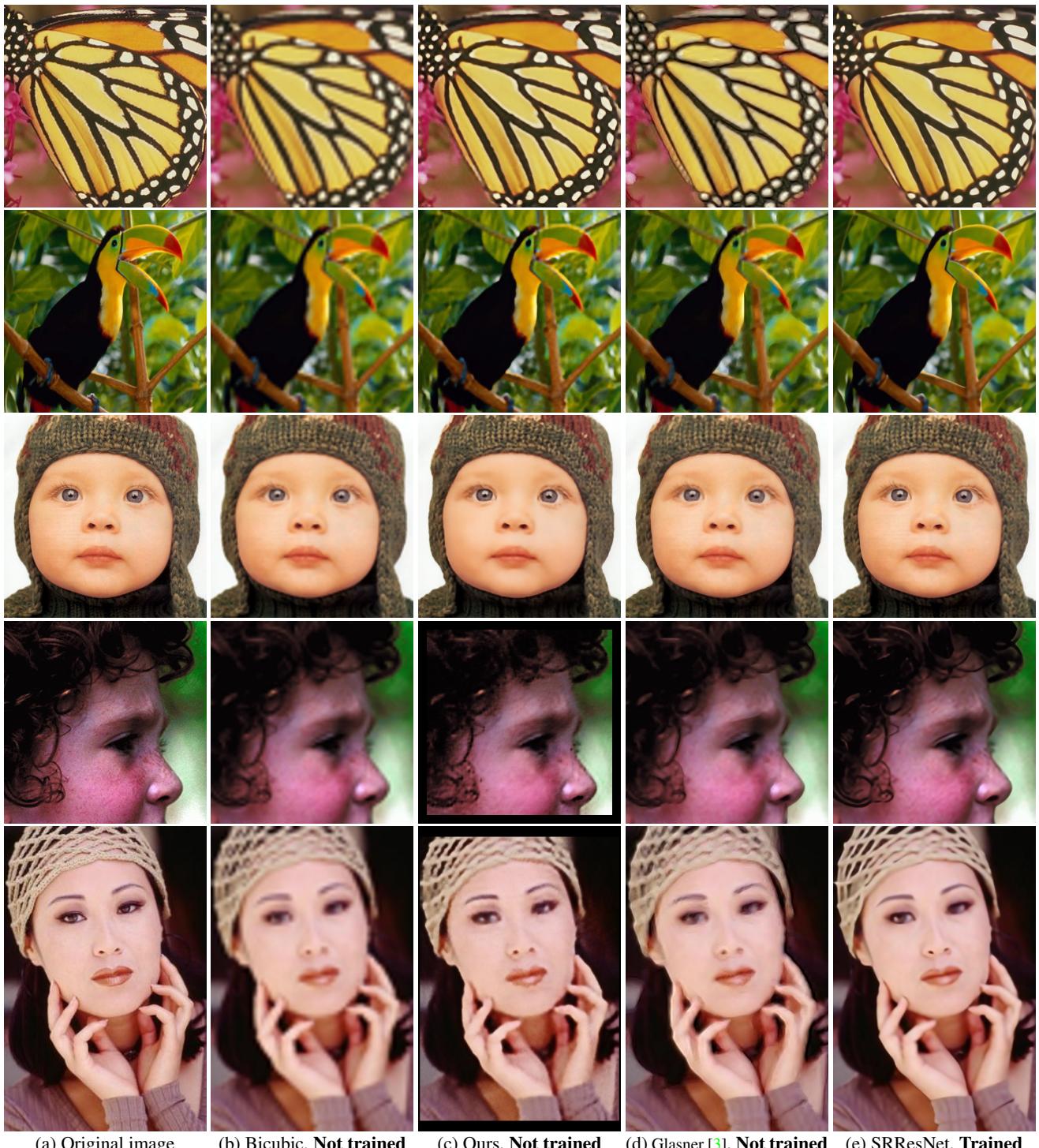


Figure 6: **4x image super-resolution.** Comparison on the Set5 dataset. In our experiments we center cropped the images to make spatial dimensions divisible by 32 (resulting in the black frames in the plots that can be removed with more sophisticated padding schemes).

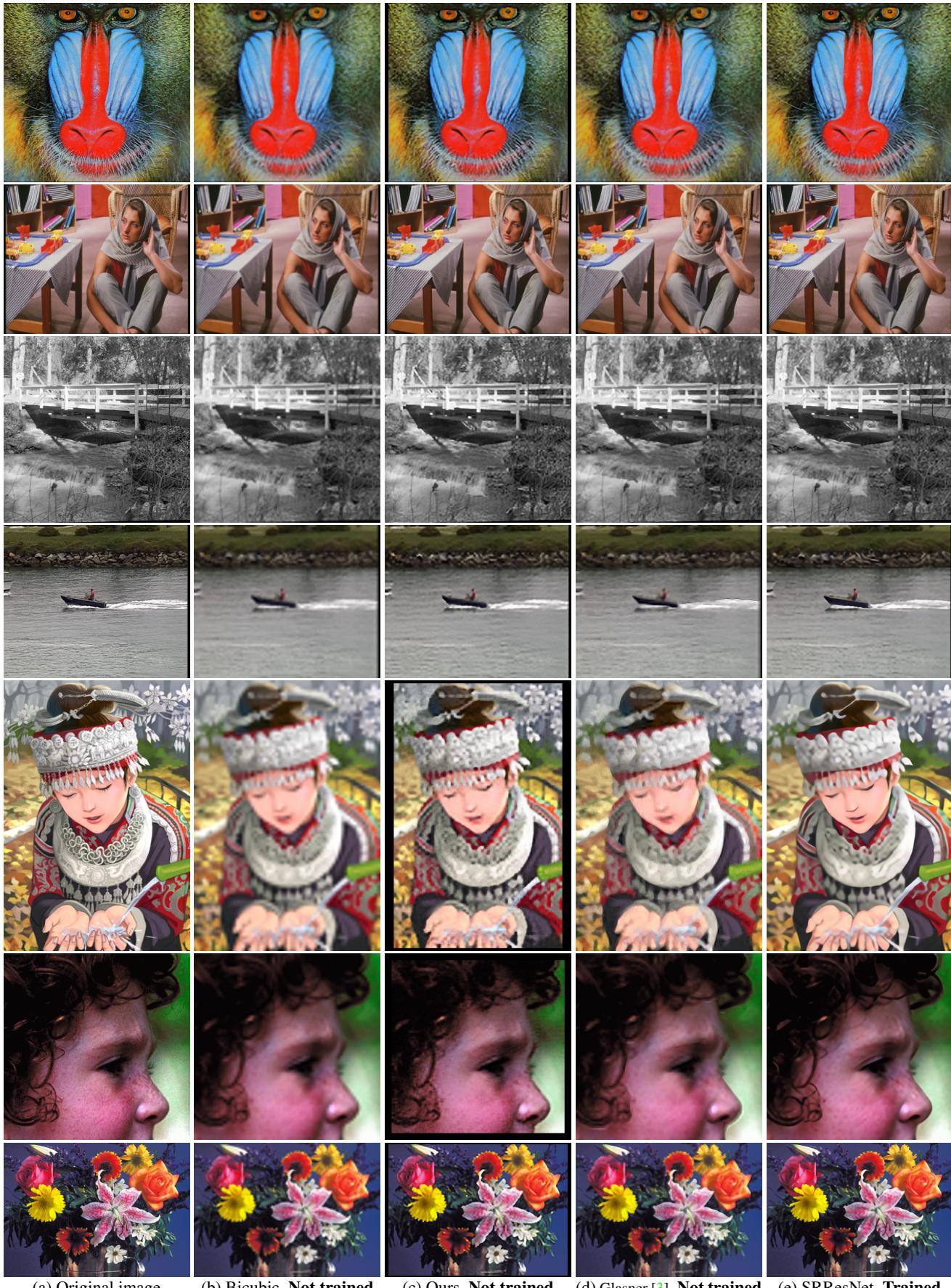


Figure 7: **4x image super-resolution.** Comparison on the Set14 dataset. Part 1.

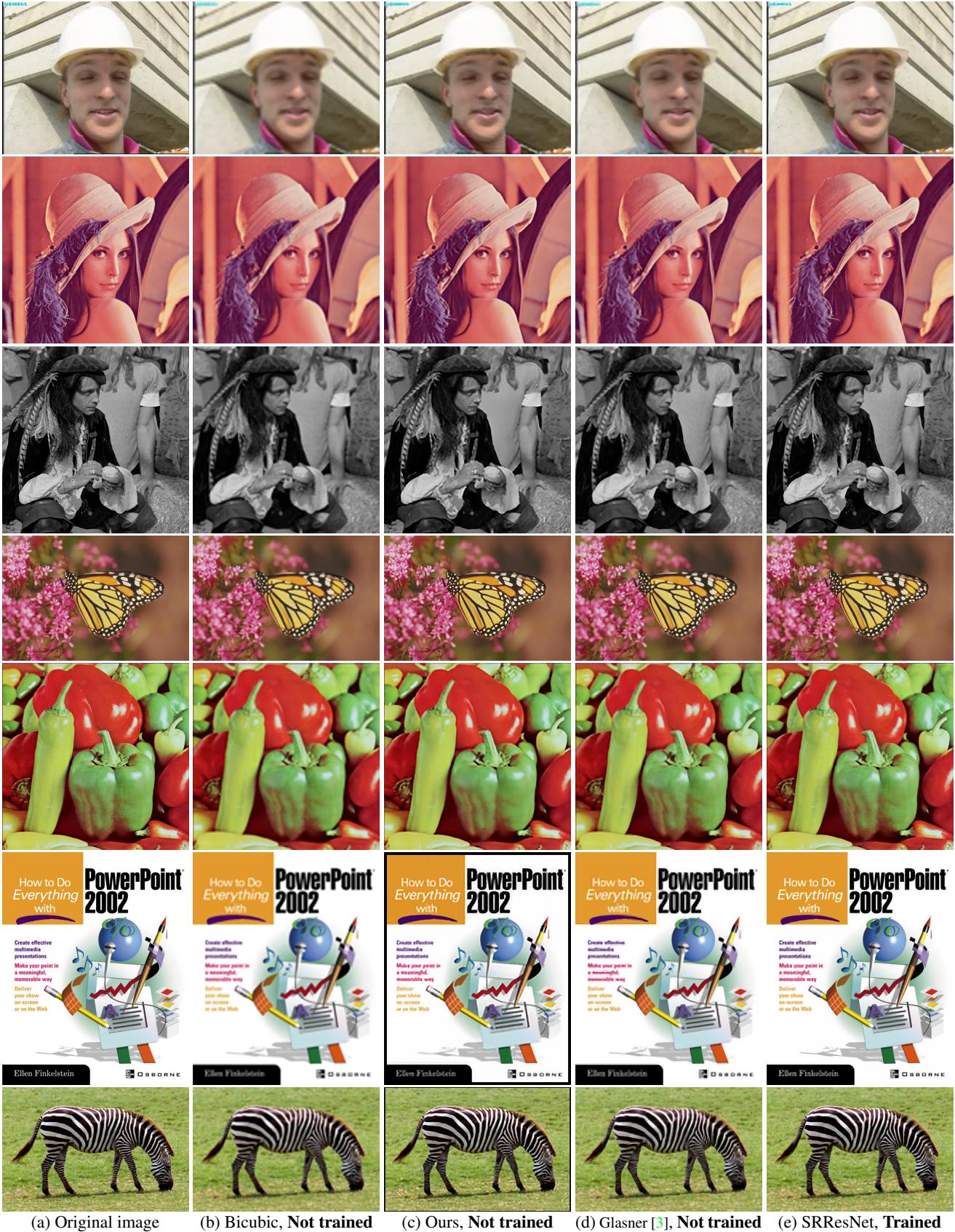


Figure 8: **4x image super-resolution.** Comparison on the Set14 dataset. Part 2.

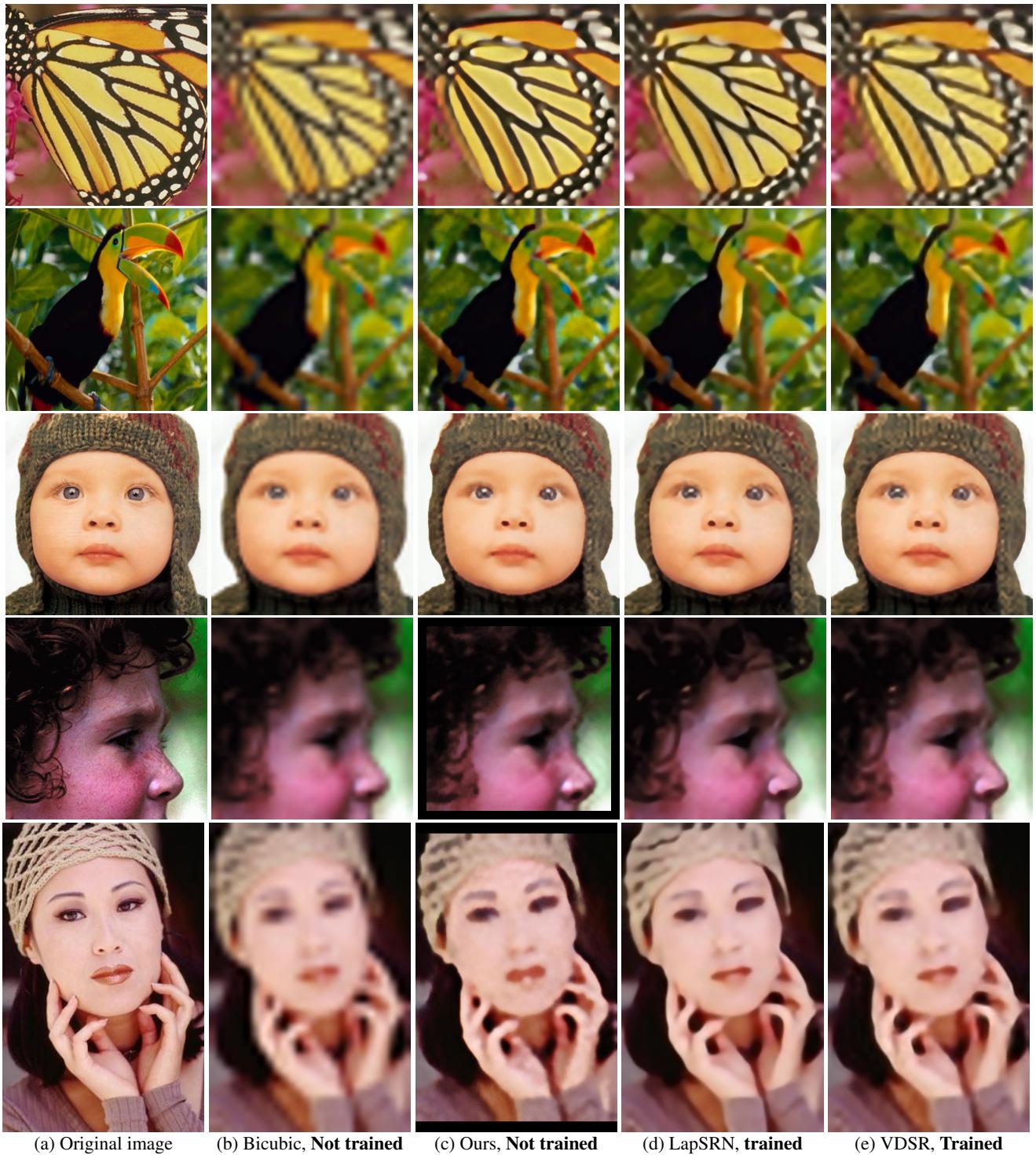


Figure 9: **8x image super-resolution.** Comparison on the Set5 dataset.

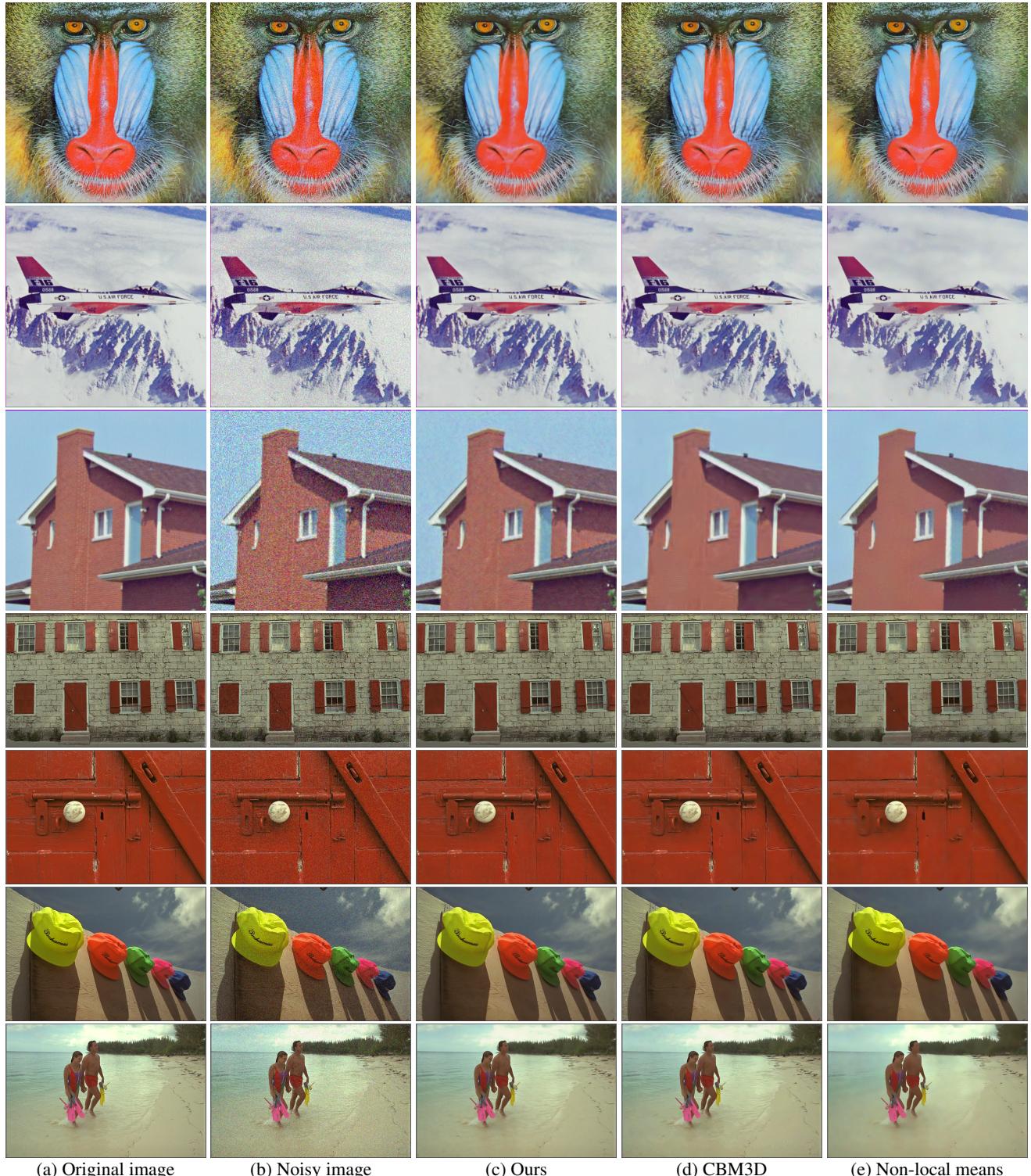


Figure 10: **Denoising**, $\sigma = 25$. Part 1.



Figure 11: **Denoising, $\sigma = 25$.** Part 2.

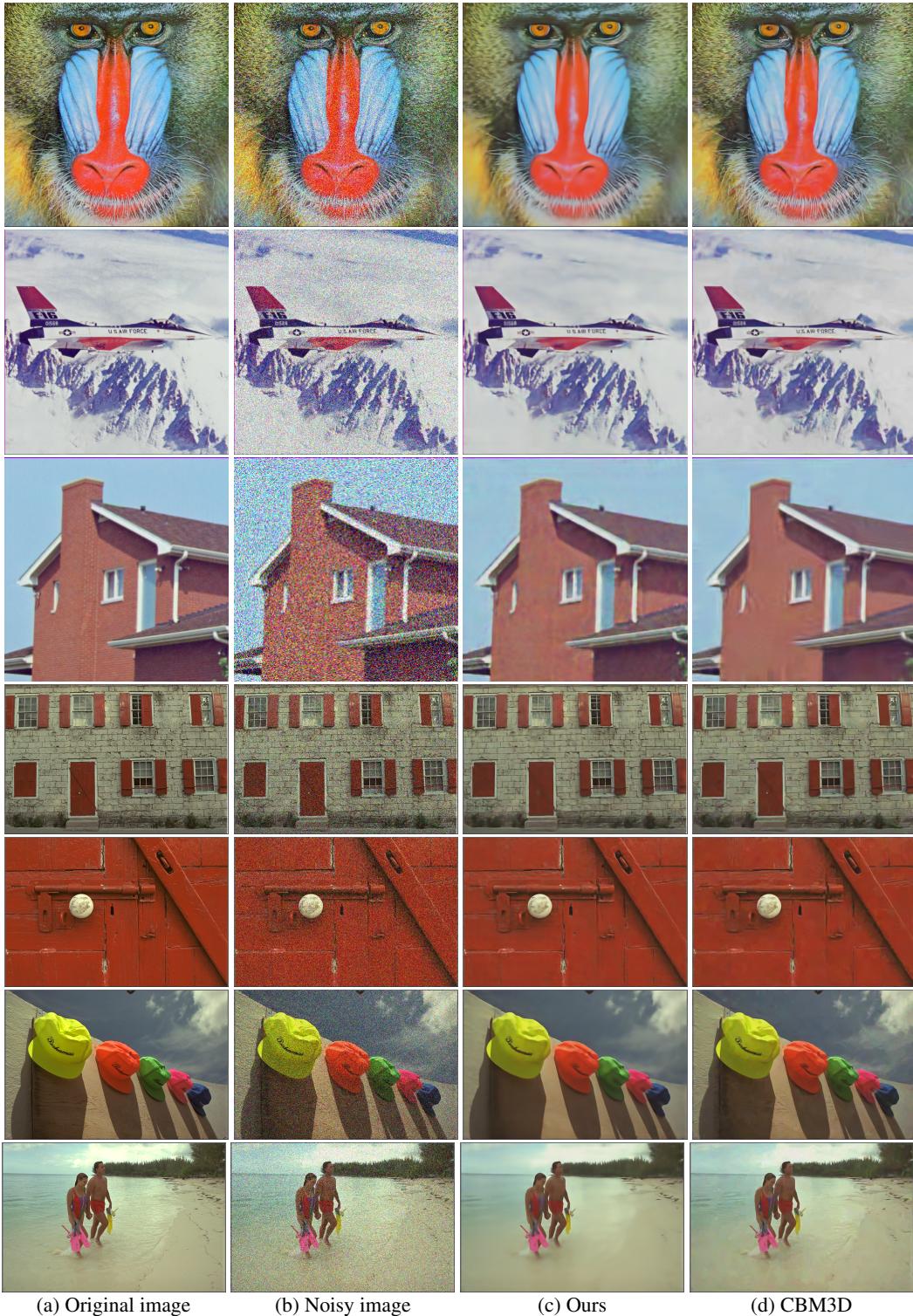


Figure 12: Denoising results on the standard benchmark set, $\sigma = 50$. Part 1.

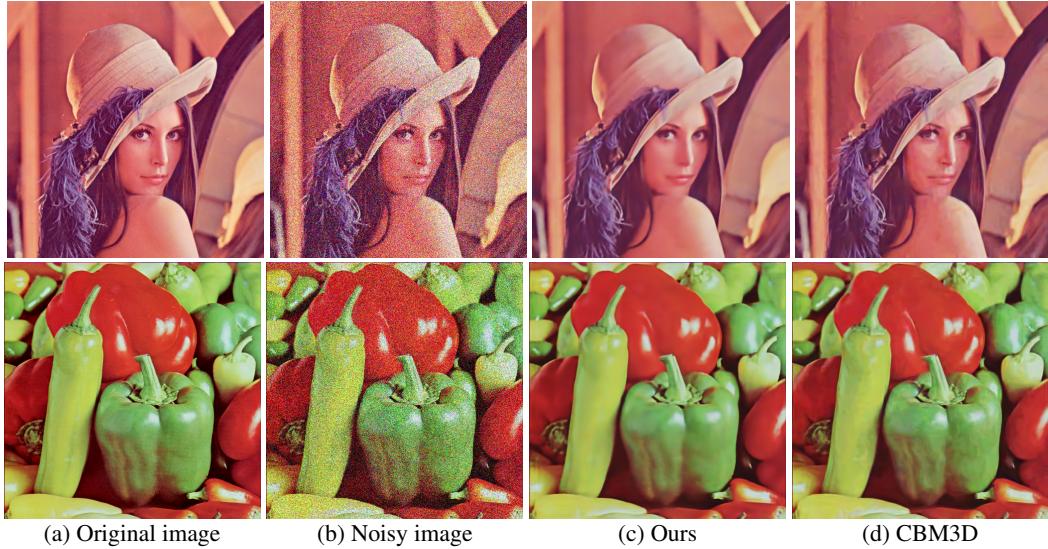
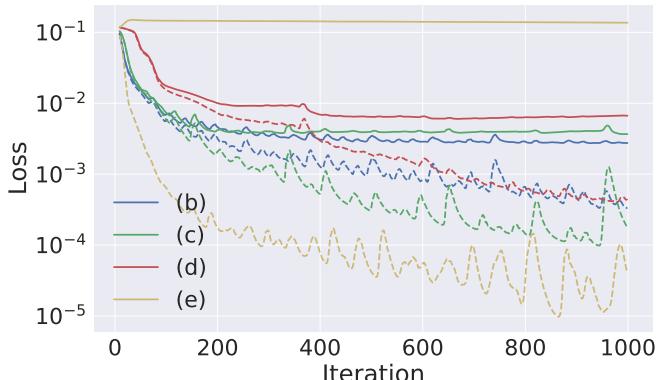
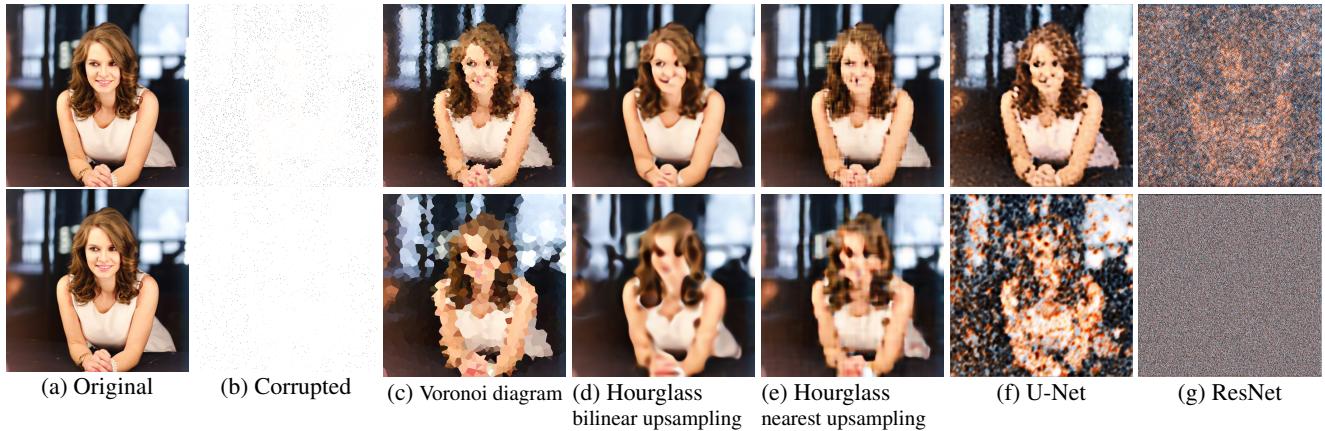
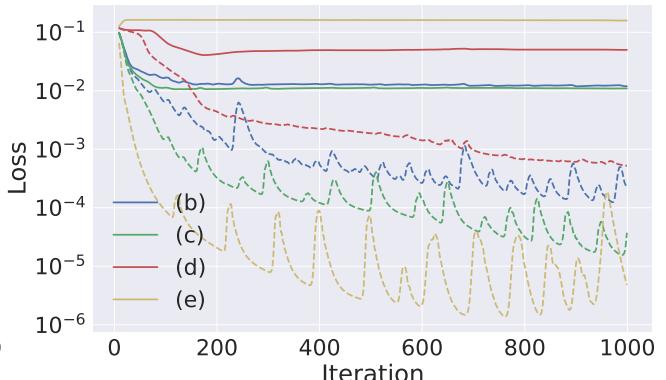


Figure 13: Denoising results on the standard benchmark set, $\sigma = 50$. Part 2.



(h) Loss/MSE curves for the top row (2% pixels known)



(i) Loss/MSE curves for the bottom row (0.05% pixels known)

Figure 14: Reconstruction using small number of pixels. We reconstruct ‘Kate’ image from 2% and 0.05% of pixels in the first and the second rows respectively. Different architectures lead to different inpainting results, yet each converges to a local minima of the data term (dashed line). Solid line shows MSE error of the generated image compared to the original image.

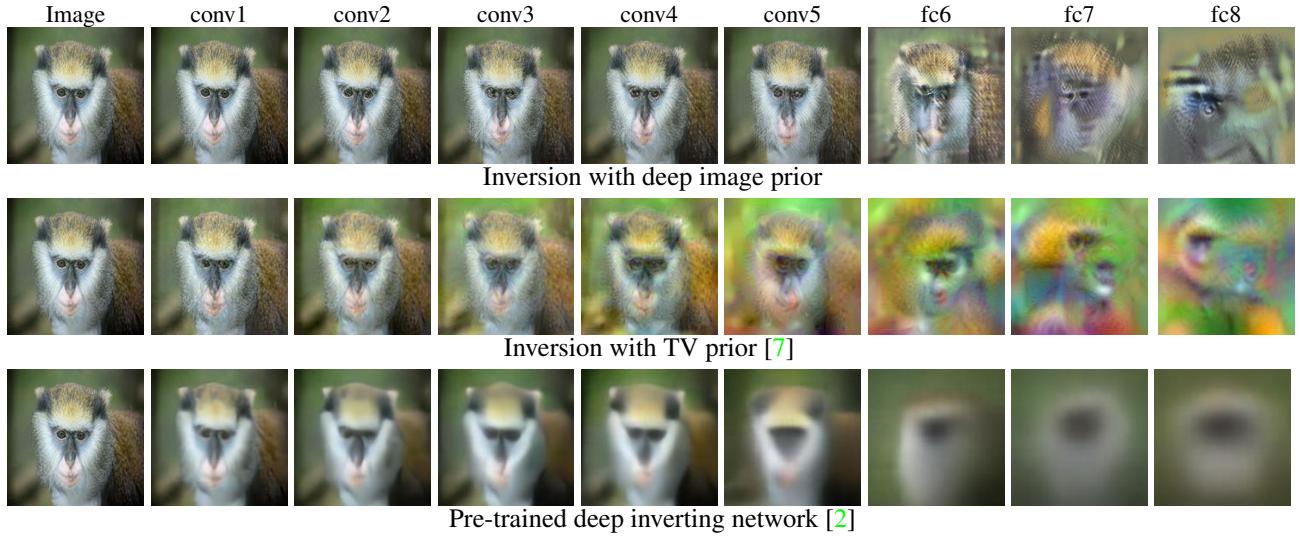


Figure 15: **Inversion of AlexNet activations at different layers with different priors** (see main text for the discussion of the problem).

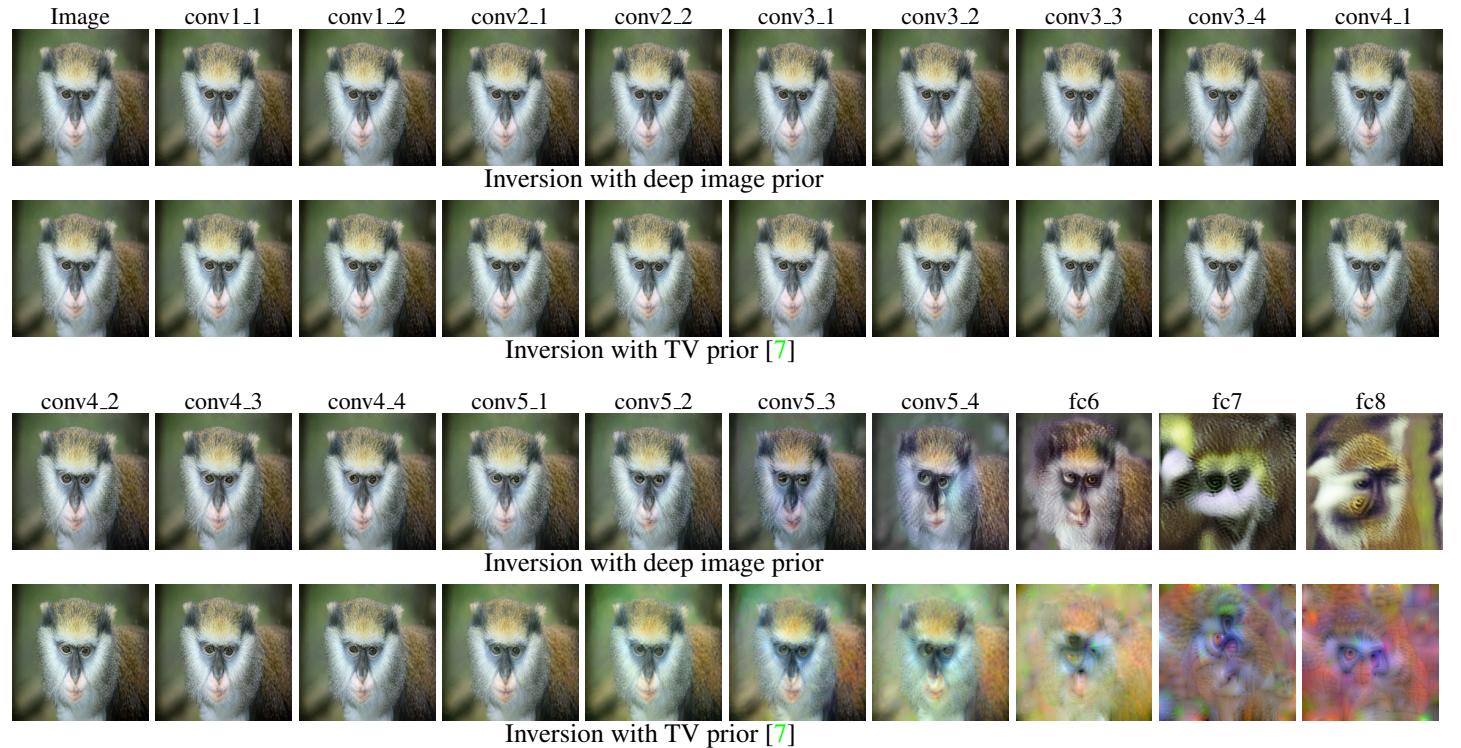


Figure 16: Inversion of VGG-19 [8] network activations at different layers with different priors (see main text for the discussion of the problem).