

# **Noise2Void - Learning Denoising from Single Noisy Images**

Alexander Krull, Tim-Oliver Buchholz, Florian Jug  
MPI-CBG/PKS (CSBD), Dresden, Germany

**CVPR 2018**

# **CONTENS**

- **Introduction**
- **Methods**
- **Experiments**

# Introduction

- A noisy image = signal + noise
  - $x = s + n$
- **Traditional CNN:** need pairs of noisy input and clean target images ( $x, s$ )
  - cannot be trained without clean target images
- **N2N:** need pairs of noisy input and noisy target images ( $x, x'$ )
  - $x = s + n, x' = s + n'$
  - the acquisition of such pairs with constant  $s$  is only possible for static scenes
- **N2V:** need single noisy images as input and target ( $x$ )
- **Assumptions:**
  - (i) the signal  $s$  is not pixel-wise independent
  - (ii) the noise  $n$  is conditionally pixel-wise independent given the signal  $s$

# Methods

- **Image Formation:**

- The generation of an image  $\mathbf{x} = \mathbf{s} + \mathbf{n}$  as a draw from the joint distribution

$$p(\mathbf{s}, \mathbf{n}) = p(\mathbf{s})p(\mathbf{n}|\mathbf{s})$$

- The pixels  $s_i$  of the signal are not statistically independent

$$p(\mathbf{s}_i|\mathbf{s}_j) \neq p(\mathbf{s}_i)$$

- Noise are conditionally independent when the signal given

$$p(\mathbf{n}|\mathbf{s}) = \prod_i p(\mathbf{n}_i|\mathbf{s}_i)$$

- We assume the noise to be zero-mean

$$\mathbb{E}[\mathbf{n}_i] = 0, \quad \mathbb{E}[\mathbf{x}_i] = \mathbf{s}_i$$

# Traditional CNN

- We can also see our CNN as a function that takes a patch  $x_{RF}(\mathbf{i})$  as input and outputs a prediction  $\hat{\mathbf{s}}_i$  for the single pixel  $\mathbf{i}$  located at the patch center.

$$f(x_{RF}(\mathbf{i}); \boldsymbol{\theta}) = \hat{\mathbf{s}}_i$$

- Use these pairs to tune the parameters  $\boldsymbol{\theta}$  to minimize pixel-wise loss

$$\arg \min_{\boldsymbol{\theta}} \sum_j \sum_i L \left( f \left( \mathbf{x}_{RF}^j(\mathbf{i}); \boldsymbol{\theta} \right) = \hat{\mathbf{s}}_i^j, \mathbf{s}_i^j \right)$$

- Consider the standard MSE loss

$$L \left( \hat{\mathbf{s}}_i^j, \mathbf{s}_i^j \right) = \left( \hat{\mathbf{s}}_i^j - \mathbf{s}_i^j \right)^2$$

# Noise2Noise

- We start out with noisy image pairs  $(x^j, x'^j)$

$$\mathbf{x}^j = \mathbf{s}^j + \mathbf{n}^j \text{ and } \mathbf{x}'^j = \mathbf{s}^j + \mathbf{n}'^j$$

- As in traditional training, we tune our parameters to minimize a loss

$$\arg \min_{\theta} \sum_j \sum_i L \left( f \left( \mathbf{x}_{\text{RF}}^j(i); \theta \right) = \hat{\mathbf{s}}_i^j, \mathbf{x}_i^j \right)$$

- The training will still converge to the correct solution because the expected value of the noisy input is equal to the clean signal.

$$\mathbb{E} [\mathbf{x}_i] = \mathbf{s}_i$$

# Noise2Void

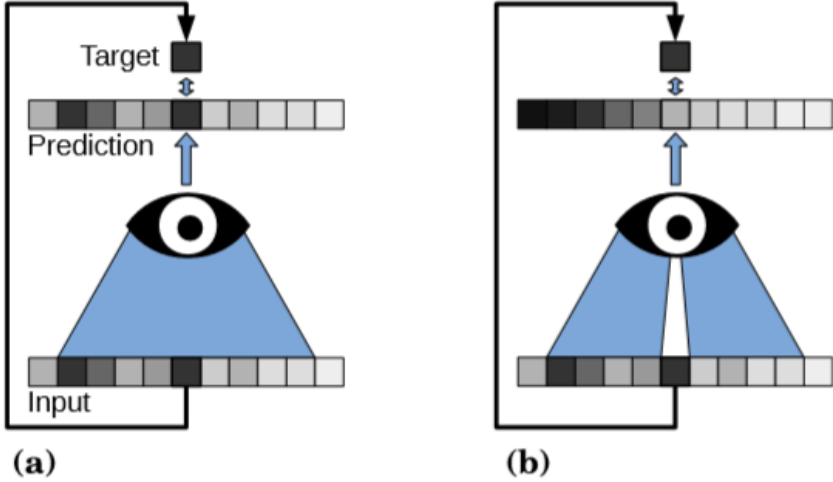


Figure 2: A conventional network versus our proposed blind-spot network. (a) In the conventional network the prediction for an individual pixel depends on a square patch of input pixels, known as a pixel's *receptive field* (pixels under blue cone). If we train such a network using the same noisy image as input and as target, the network will degenerate and simply learn the identity. (b) In a *blind-spot network*, as we propose it, the receptive field of each pixel excludes the pixel itself, preventing it from learning the identity. We show that blind-spot networks can learn to remove pixel wise independent noise when they are trained on the same noisy images as input and target.

- The input and the target are from a single noisy training image  $\mathbf{x}_j$
- The receptive field  $\tilde{\mathbf{x}}_{\text{RF}}^j(i)$  of this network to have a blind-spot in its center.
- We can train it by minimizing the empirical risk

$$\arg \min_{\theta} \sum_j \sum_i L \left( f \left( \tilde{\mathbf{x}}_{\text{RF}}^j(i); \theta \right), \mathbf{x}_i^j \right)$$

- Accuracy will be slightly impaired compared to a normal network

# Noise2Void

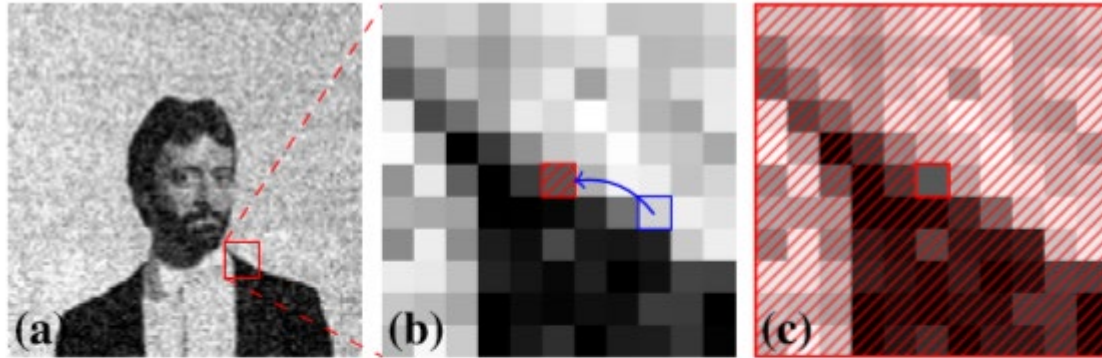


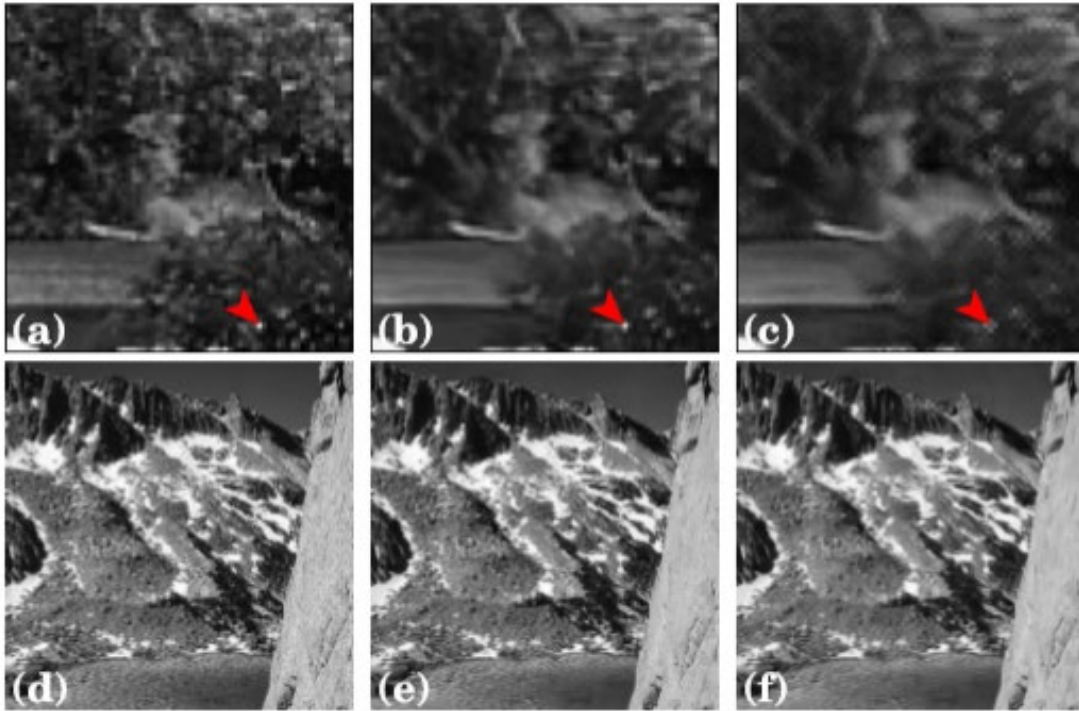
Figure 3: Blind-spot masking scheme used during NOISE2VOID training. (a) A noisy training image. (b) A magnified image patch from (a). During N2V training, a randomly selected pixel is chosen (blue rectangle) and its intensity copied over to create a blind-spot (red and striped square). This modified image is then used as input image during training. (c) The target patch corresponding to (b). We use the original input with unmodified values also as target. The loss is only calculated for the blind-spot pixels we masked in (b).

- If we implement the naïve training scheme, it is not efficient: We have to process an entire patch to calculate the gradients for a single output pixel
- Replace the value in the center of each input patch with a randomly selected value from the surrounding area
- We can now simultaneously calculate the gradients for all of them, while ignoring the rest of the predicted image
- A specialized loss function that is zero for all but the selected pixels

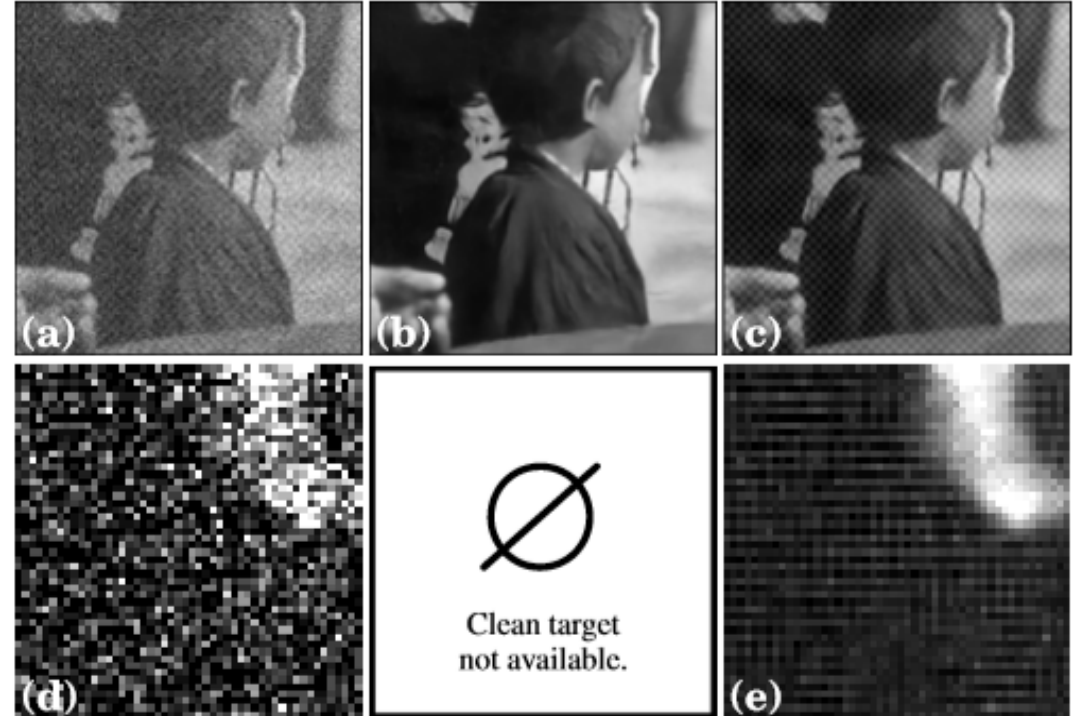
# Experiments - Denoising Results and Costs

	Ground Truth	Input	BM3D	Traditional	NOISE2NOISE	NOISE2VOID
<b>BSD68</b>			 PSNR: 28.59	 PSNR: 29.06	 PSNR: 28.86	 PSNR: 27.71
<b>Simulated Data</b>			 PSNR: 29.96	 PSNR: 32.56	 PSNR: 32.43	 PSNR: 32.28
<b>cryo-TEM</b>	 Does not exist.		 Runtime: ~33.2s	 Clean target not available.	 Runtime: ~1.3s	 Runtime: ~1.3s
<b>CTC-MSC</b>	 Does not exist.		 Runtime: ~4.6s	 Clean target not available.	 Noisy target not available.	 Runtime: ~0.1s
<b>CTC-N2DH</b>	 Does not exist.		 Runtime: ~5.2s	 Clean target not available.	 Noisy target not available.	 Runtime: ~0.1s

# Experiments - Errors and Limitations

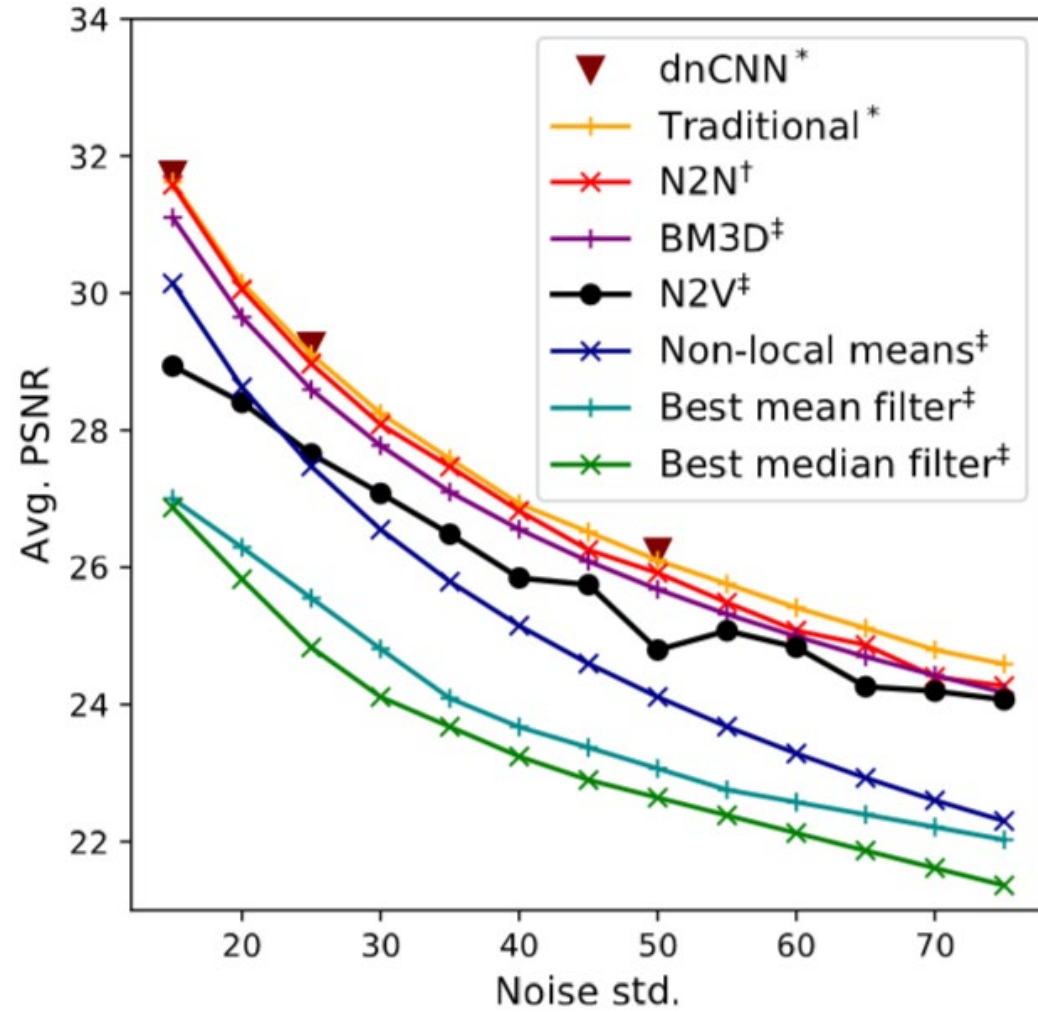


- Images with more high irregularities are difficult to predict
- The more difficult it is to predict a pixel's signal from its surroundings the more errors are expected to appear in N2V predictions



- Structured noise
- The N2V trained CNN removes the unpredictable components of the noise, but reveals the hidden pattern

# Experiments - Performance



Mean filter (5x5)



Median filter (5x5)



N2V

