

# High-Quality Self-Supervised Deep Image Denoising

Samuli Laine

NVIDIA

Tero Karras

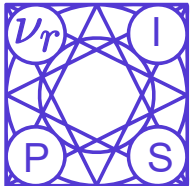
NVIDIA

Jaakko Lehtinen

NVIDIA, Aalto University

Timo Aila

NVIDIA



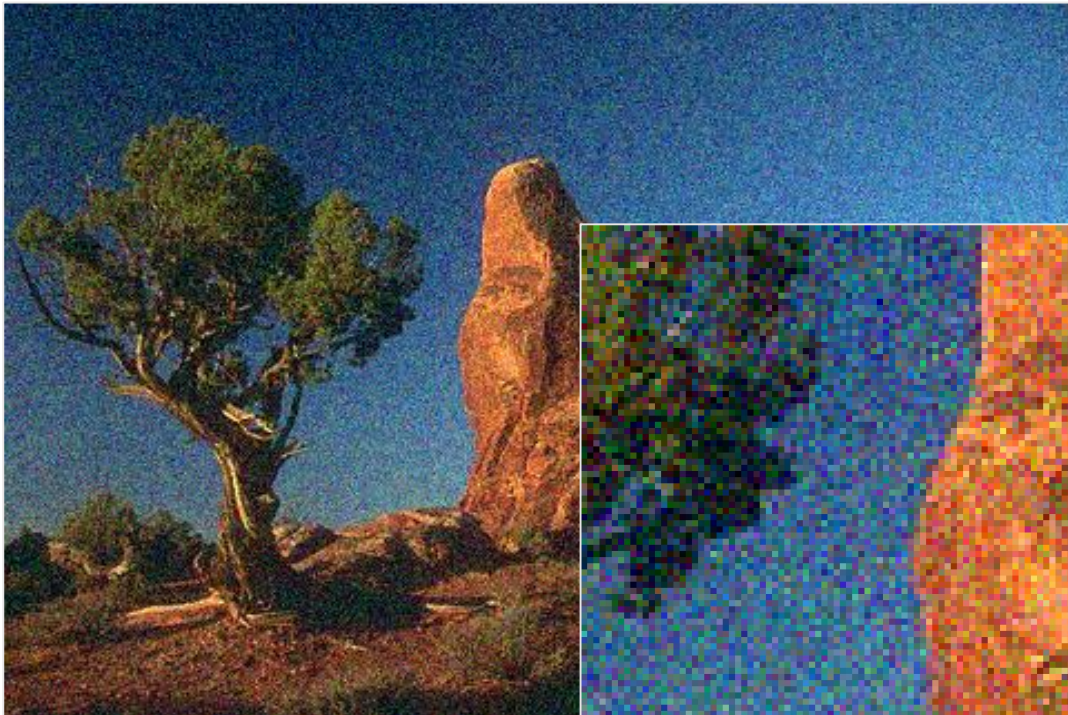
NeurIPS | 2019



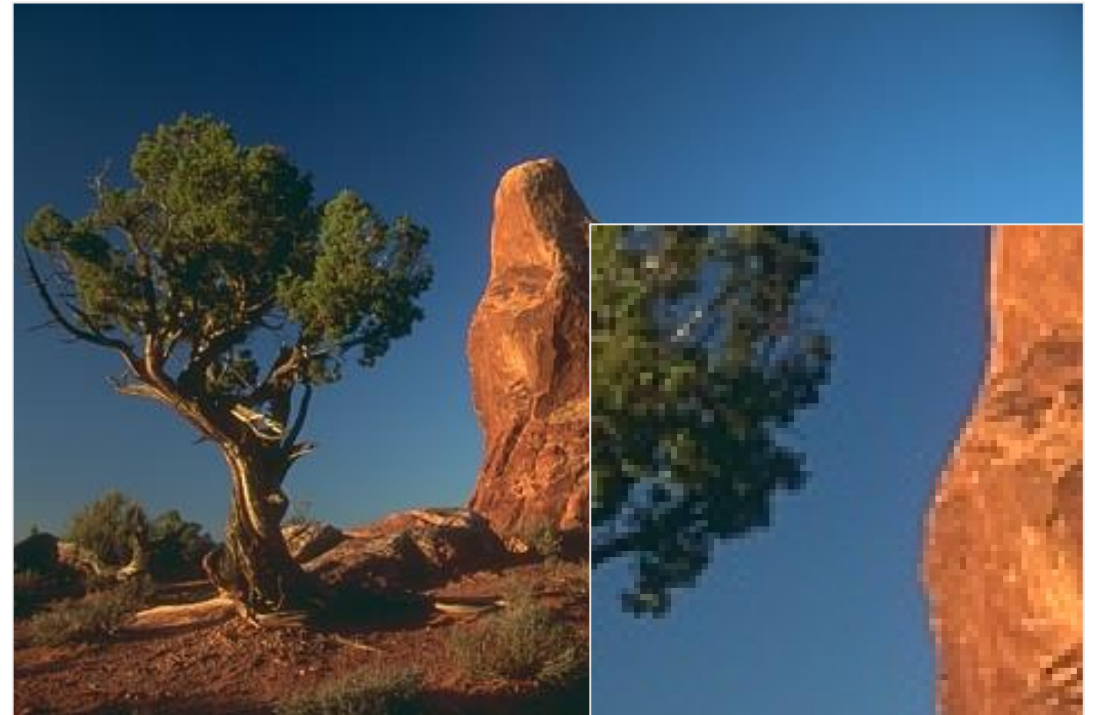
# Goals and contributions

- Train a deep denoiser network
  - Why is deep learning needed? So that denoiser adapts to underlying data!
- Remove the need for separate training data with self-supervision
- **Contribution 1:** Bayesian approach for high-quality denoising results
  - Target is to match a denoiser trained with clean reference data
- **Contribution 2:** Improve training performance with an efficient blind-spot network architecture

# Background: Traditional training



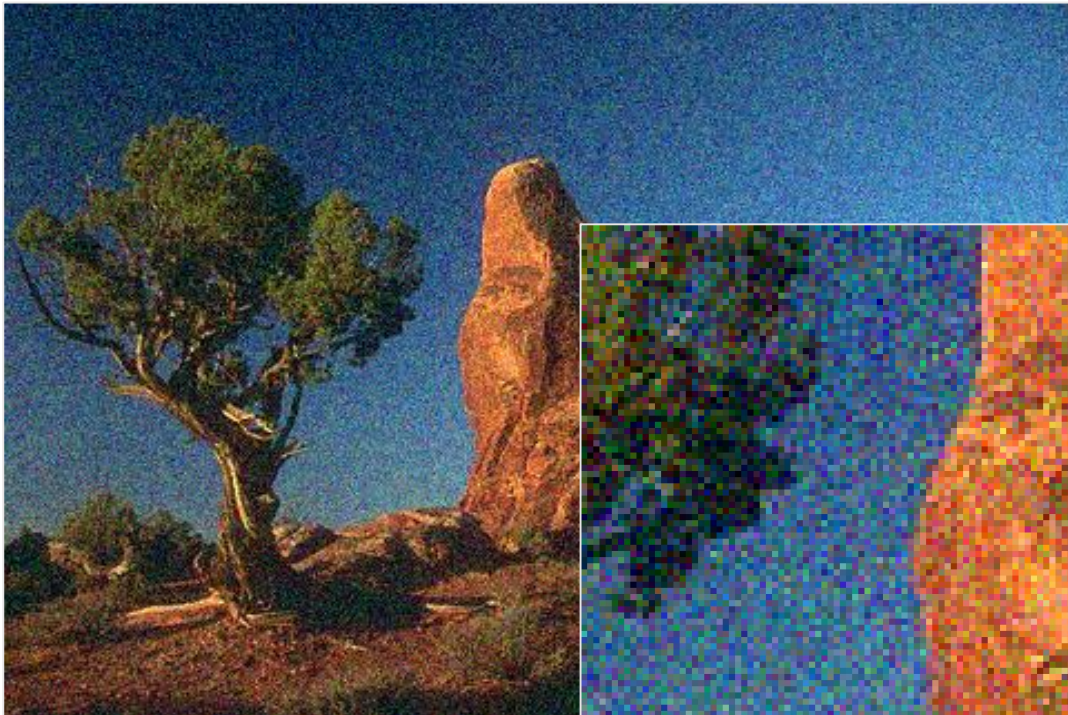
Input: **Noisy image**



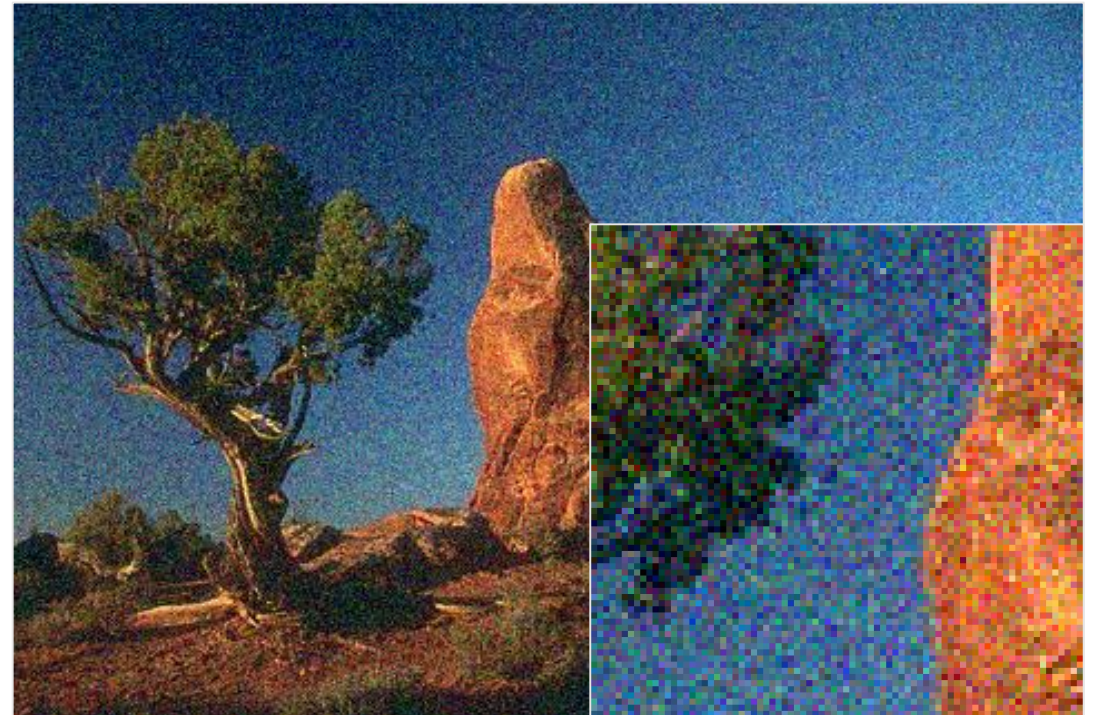
Target: **Clean image**

# Background: Noise2Noise training

[Lehtinen et al., 2018]



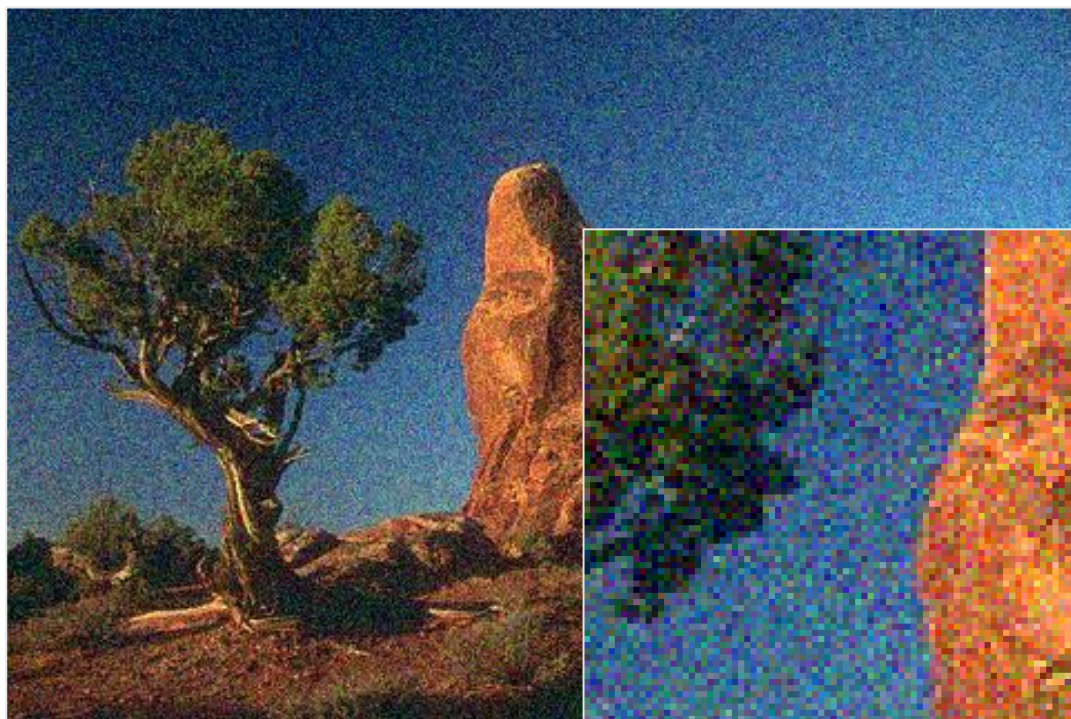
Input: **Noisy image**



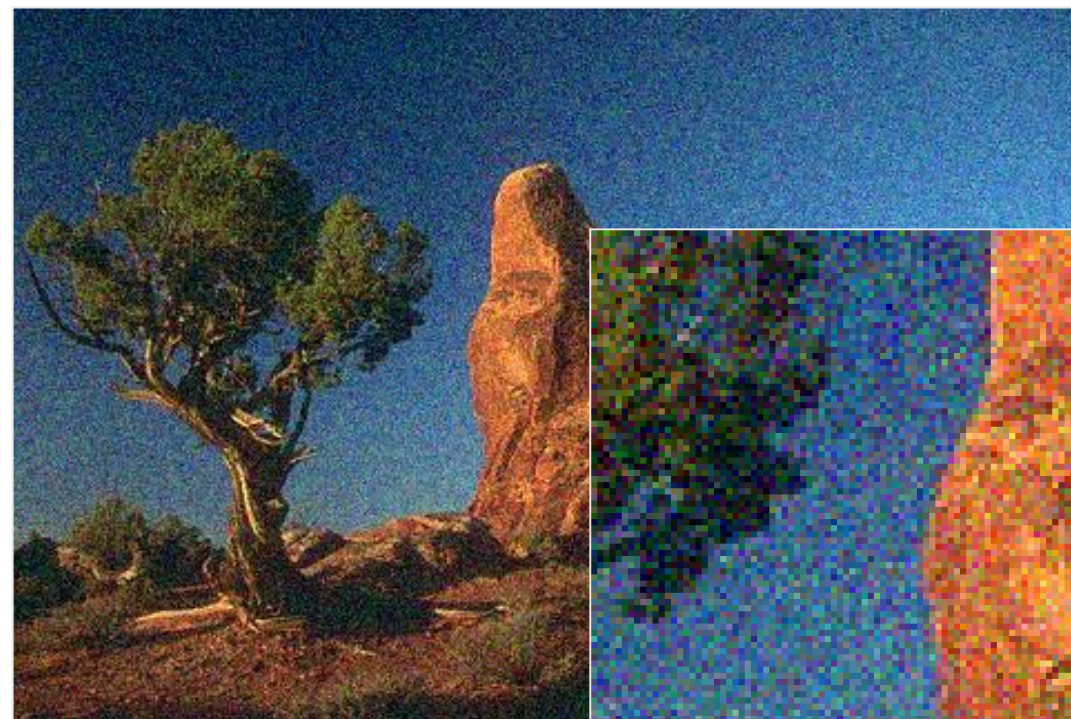
Target: **A different noisy image**

# Background: Noise2Void training

[Krull et al., 2018]



Input: **Noisy image**



Target: **The same noisy image**

# Concepts and notation

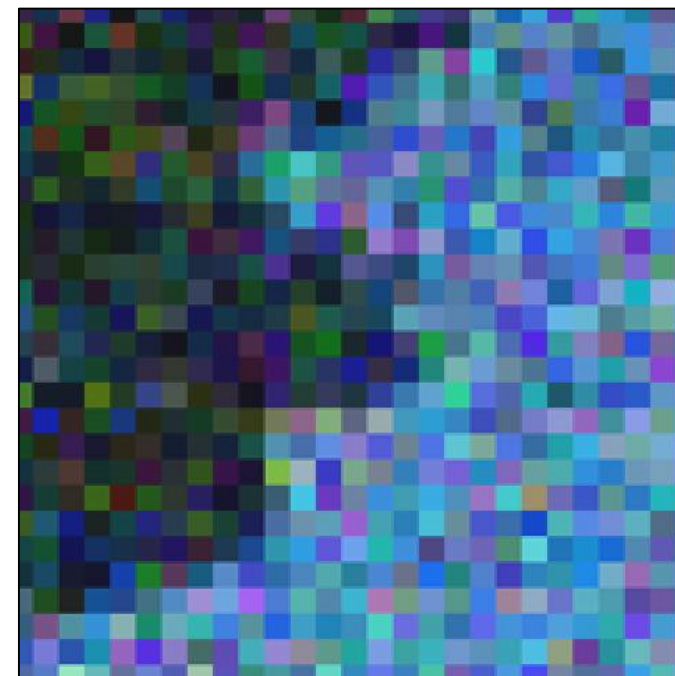
$x$  clean pixel value

$y$  noisy pixel value

$\Omega_y$  noisy *context*, i.e.,  
noisy image except  
pixel  $y$



Clean image



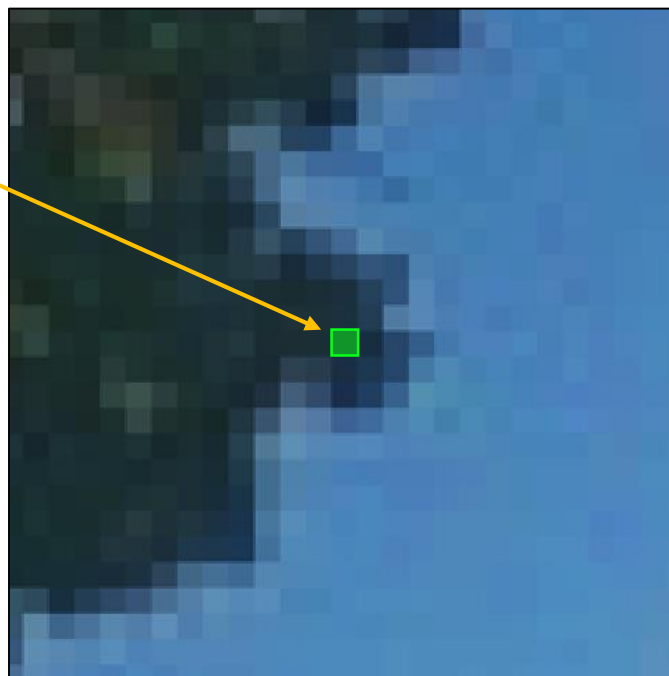
Noisy image

# Concepts and notation

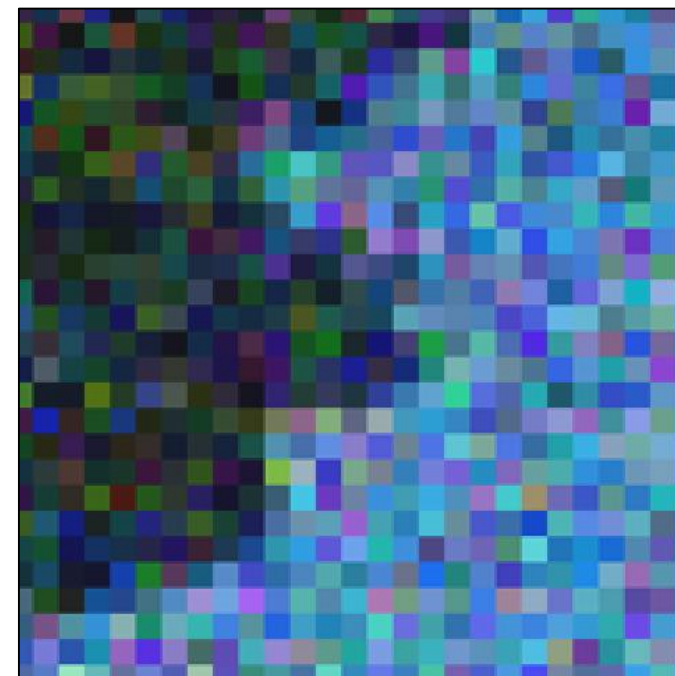
$x$  clean pixel value

$y$  noisy pixel value

$\Omega_y$  noisy *context*, i.e.,  
noisy image except  
pixel  $y$



Clean image



Noisy image

# Concepts and notation

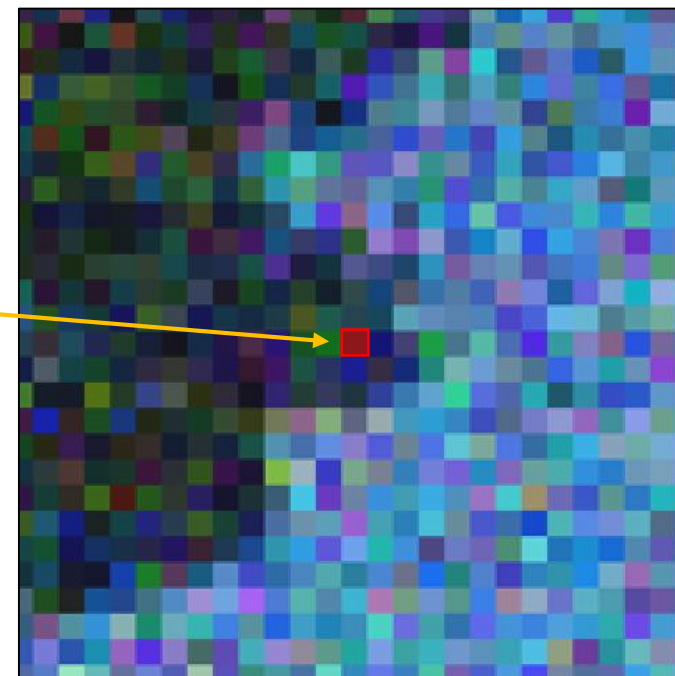
$x$  clean pixel value

$y$  noisy pixel value

$\Omega_y$  noisy *context*, i.e.,  
noisy image except  
pixel  $y$



Clean image



Noisy image



# Concepts and notation

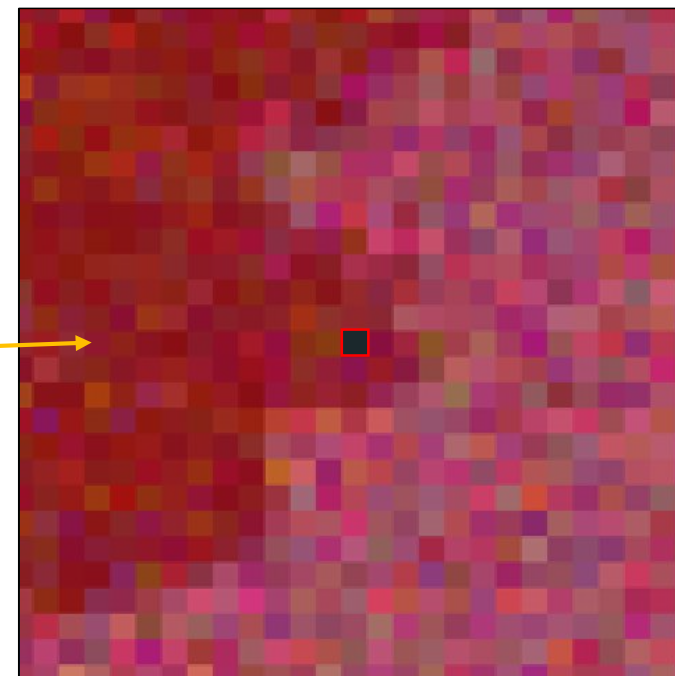
$x$  clean pixel value

$y$  noisy pixel value

$\Omega_y$  **noisy context**, i.e.,  
noisy image except  
pixel  $y$



Clean image



Noisy image



# What is supervised training?

- Lump noisy pixel  $\mathbf{y}$  and context  $\mathbf{\Omega}_y$  together
- Learn to infer clean pixel  $\mathbf{x}$  as  $\mathbb{E}_{\mathbf{x}}[p(\mathbf{x}|\mathbf{y}, \mathbf{\Omega}_y)]$
- I.e., train  $f_\theta: \mathbf{y}, \mathbf{\Omega}_y \rightarrow \mathbf{x}$  by optimizing  $\operatorname{argmin}_\theta \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{\Omega}_y} [L(f_\theta(\mathbf{y}, \mathbf{\Omega}_y), \mathbf{x})]$ 
  - Simplifying assumptions made here: L2 loss, zero-mean noise

# What is Noise2Void training?

- Only use context  $\Omega_y$  for inference [Krull et al., 2018]
  - Thus, approximate clean pixel  $x$  as  $\mathbb{E}_x[p(x|\Omega_y)]$
- Can replace  $x$  with  $y$  if noise is zero-mean [Lehtinen et al., 2018]
  - Optimize  $\operatorname{argmin}_{\theta} \mathbb{E}_y[L(f_{\theta}(\Omega_y), y)]$  – no clean  $x$  is needed
- This is equivalent if corruption is independent between pixels!
  - See [Batson and Royer, 2019] for further analysis

# Limitations of Noise2Void

- Ignoring  $\mathbf{y}$  when denoising clearly leaves useful information unused
- While we can regress  $f_{\theta}: \Omega_{\mathbf{y}} \rightarrow \mathbf{y}$ , we cannot regress  $f_{\theta}: \mathbf{y}, \Omega_{\mathbf{y}} \rightarrow \mathbf{y}$ 
  - Trivial solution is to pass pixel value through as-is  $\rightarrow$  no denoising
  - Hence, at training time we cannot use  $\mathbf{y}$  as an input
- Our solution is to bring in  $\mathbf{y}$  via Bayesian inference at test time
  - Concurrent work by [Krull et al., 2019]

# A more complete view

- Assume a known noise model  $p(\mathbf{y}|\mathbf{x})$  that is independent of  $\Omega_y$
- Observed noisy data (training data) now relates to clean data as

$$\underbrace{p(\mathbf{y}|\Omega_y)}_{\text{Training data}} = \int \underbrace{p(\mathbf{y}|\mathbf{x})}_{\text{Noise model}} \underbrace{p(\mathbf{x}|\Omega_y)}_{\text{Unobserved}} d\mathbf{x}$$

- This lets us learn to predict a parametric model for  $p(\mathbf{x}|\Omega_y)$  that we represent as a multivariate Gaussian  $\mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$  over color components

# Test-time inference

- The (unnormalized) posterior probability of  $\mathbf{x}$ , given observations of  $\mathbf{y}$  and  $\Omega_{\mathbf{y}}$ , is given by Bayes' rule as

$$\underbrace{p(\mathbf{x}|\mathbf{y}, \Omega_{\mathbf{y}})}_{\text{Posterior}} \propto \underbrace{p(\mathbf{y}|\mathbf{x})}_{\text{Noise model}} \underbrace{p(\mathbf{x}|\Omega_{\mathbf{y}})}_{\text{Prior}}$$

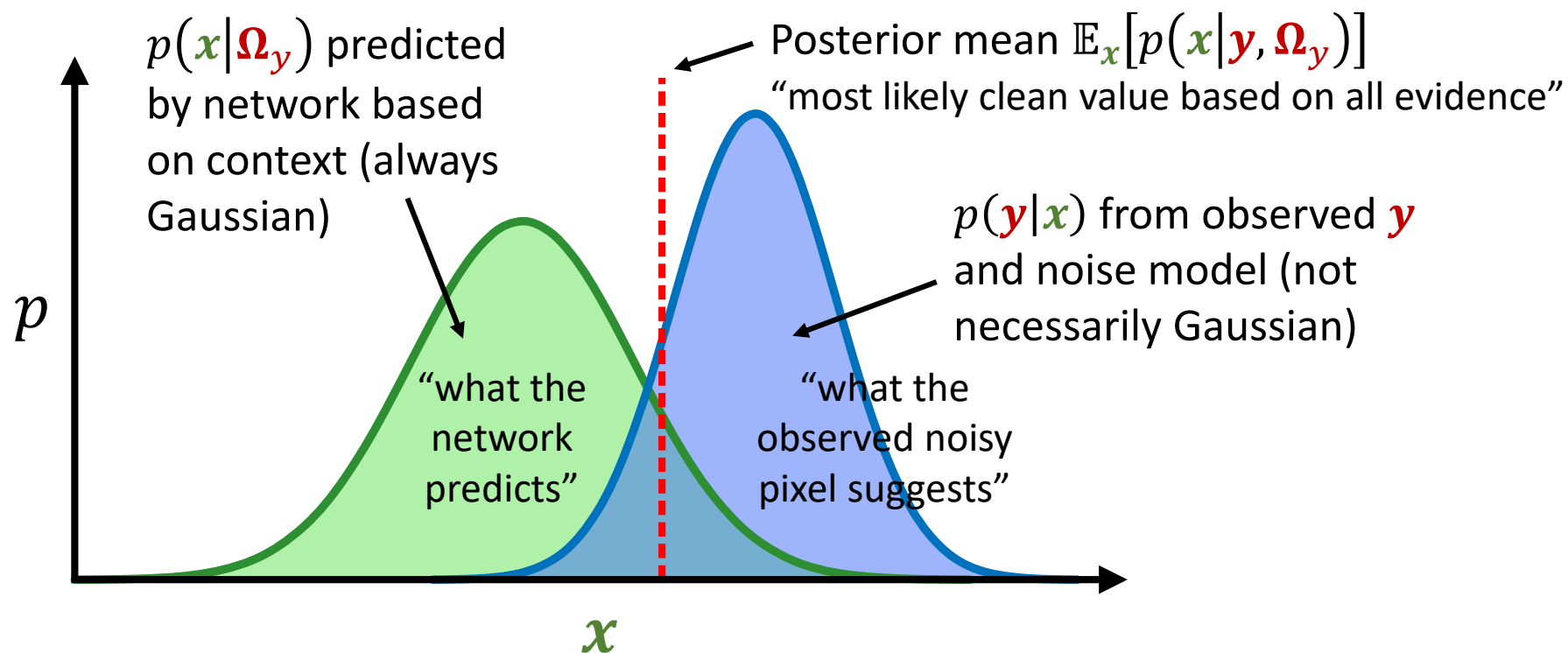
← Predicted by the network

← Known

- We can make our best guess of  $\mathbf{x}$  based on the posterior distribution
- Concretely, we output the posterior mean  $\mathbb{E}_{\mathbf{x}}[p(\mathbf{x}|\mathbf{y}, \Omega_{\mathbf{y}})]$  because it minimizes MSE and therefore maximizes PSNR

# Test-time inference – a sketch

- Simplified view of a 1D (monochromatic) case



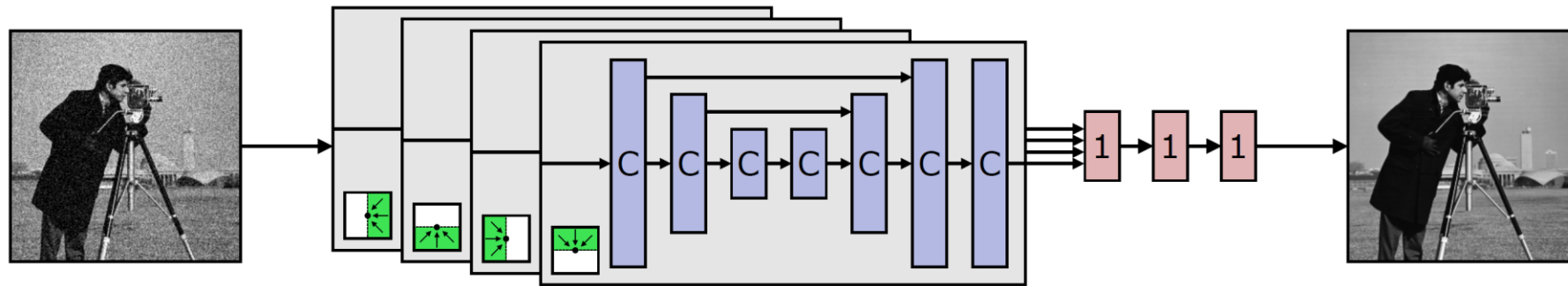
# Summary of our approach

- In training phase, train neural network  $f_\theta$  to map context  $\Omega_y$  to mean  $\mu_x$  and variance  $\Sigma_x$  to approximate prior  $p(\mathbf{x}|\Omega_y)$ 
  - Known noise model maps  $\mathcal{N}(\mu_x, \Sigma_x) \rightarrow \mathcal{N}(\mu_y, \Sigma_y)$  so training can be done using standard Gaussian process regression (see e.g., [Nix and Weigend, 1994])
- At test time, evaluate  $f_\theta(\Omega_y)$  and compute posterior mean  $\mathbb{E}_x[p(\mathbf{x}|\mathbf{y}, \Omega_y)]$  by closed-form integration



# Implementing blind-spot network efficiently

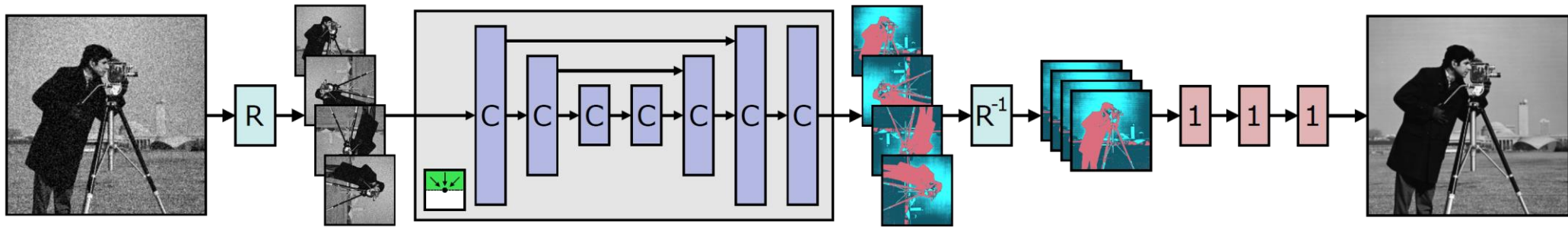
- Our solution: Combine information from four branches, each having its receptive field restricted to one direction only



- Restricting the receptive field to one half-space is easier than removing just one pixel

# Optimizing a bit

- Roll the four branches into one, rotate image data instead



- Implicitly shares weights between branches
- Implementation details in the paper

# Unknown noise parameters

- What if the noise model has an unknown parameter? What if the parameter varies for every image?
  - E.g., standard deviation  $\sigma$  in Gaussian noise  $\mathcal{N}(\mathbf{0}, \sigma^2 I)$
- We show that these can be estimated from the data as well, so that each image in training and test data can have a different, unknown amount of noise
  - Requires regularization in certain cases to break ambiguity (is the image actually noisy vs. is the clean signal hard to predict) – see paper for details

# Results: Gaussian noise ( $\sigma = 25$ )

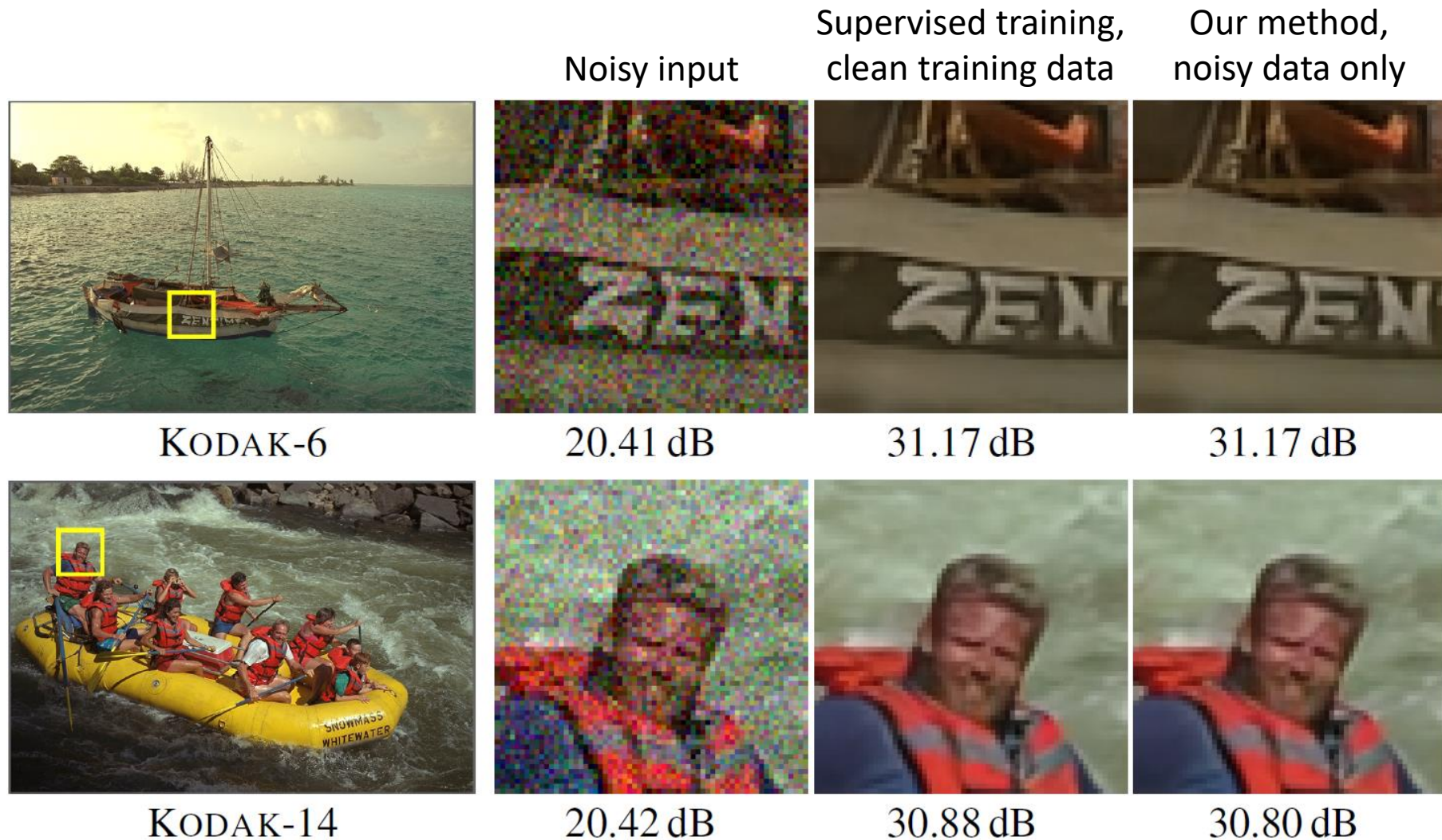


Table 1: Image quality results for Gaussian noise. Values of  $\sigma$  are shown in 8-bit units.

Noise type	Method	$\sigma$ known?	KODAK	BSD300	SET14	Average
Gaussian $\sigma = 25$	Baseline, N2C	no	32.46	31.08	31.26	31.60
	Baseline, N2N	no	32.45	31.07	31.23	31.58
	Our	yes	32.45	31.03	31.25	31.57
	Our	no	32.44	31.02	31.22	31.56
	Our ablated, diag. $\Sigma$	yes	31.60	29.91	30.58	30.70
	Our ablated, diag. $\Sigma$	no	31.55	29.87	30.53	30.65
	Our ablated, $\mu$ only	no	30.64	28.65	29.57	29.62
	CBM3D	yes	31.82	30.40	30.68	30.96
	CBM3D	no	31.81	30.40	30.66	30.96
Gaussian $\sigma \in [5, 50]$	Baseline, N2C	no	32.57	31.29	31.27	31.71
	Baseline, N2N	no	32.57	31.29	31.26	31.70
	Our	yes	32.47	31.19	31.21	31.62
	Our	no	32.46	31.18	31.13	31.59
	Our ablated, diag. $\Sigma$	yes	31.59	30.06	30.54	30.73
	Our ablated, diag. $\Sigma$	no	31.58	30.05	30.45	30.69
	Our ablated, $\mu$ only	no	30.54	28.56	29.41	29.50
	CBM3D	yes	31.99	30.67	30.78	31.15
	CBM3D	no	31.99	30.67	30.72	31.13

Supervised training with clean targets

Our result when  $\sigma$  is known vs. estimated from data

Noise2Void: Ignore  $y$  and predict based on context only

Our results are within 0.04 dB from supervised training

Table 1: Image quality results for Gaussian noise. Values of  $\sigma$  are shown in 8-bit units.

Noise type	Method	$\sigma$ known?	KODAK	BSD300	SET14	Average
Gaussian $\sigma = 25$	Baseline, N2C	no	32.46	31.08	31.26	31.60
	Baseline, N2N	no	32.45	31.07	31.23	31.58
	Our	yes	32.45	31.03	31.25	31.57
	Our	no	32.44	31.02	31.22	31.56
	Our ablated, diag. $\Sigma$	yes	31.60	29.91	30.58	30.70
	Our ablated, diag. $\Sigma$	no	31.55	29.87	30.53	30.65
	Our ablated, $\mu$ only	no	30.64	28.65	29.57	29.62
	CBM3D	yes	31.82	30.40	30.68	30.96
	CBM3D	no	31.81	30.40	30.66	30.96
Gaussian $\sigma \in [5, 50]$	Baseline, N2C	no	32.57	31.29	31.27	31.71
	Baseline, N2N	no	32.57	31.29	31.26	31.70
	Our	yes	32.47	31.19	31.21	31.62
	Our	no	32.46	31.18	31.13	31.59
	Our ablated, diag. $\Sigma$	yes	31.59	30.06	30.54	30.73
	Our ablated, diag. $\Sigma$	no	31.58	30.05	30.45	30.69
	Our ablated, $\mu$ only	no	30.54	28.56	29.41	29.50
	CBM3D	yes	31.99	30.67	30.78	31.15
	CBM3D	no	31.99	30.67	30.72	31.13

Supervised training with clean targets

Our result when  $\sigma$  is known vs. estimated from data

Noise2Void: Ignore  $y$  and predict based on context only

Our results are within 0.04 dB from supervised training

Table 1: Image quality results for Gaussian noise. Values of  $\sigma$  are shown in 8-bit units.

Noise type	Method	$\sigma$ known?	KODAK	BSD300	SET14	Average
Gaussian $\sigma = 25$	Baseline, N2C	no	32.46	31.08	31.26	31.60
	Baseline, N2N	no	32.45	31.07	31.23	31.58
	Our	yes	32.45	31.03	31.25	31.57
	Our	no	32.44	31.02	31.22	31.56
	Our ablated, diag. $\Sigma$	yes	31.60	29.91	30.58	30.70
	Our ablated, diag. $\Sigma$	no	31.55	29.87	30.53	30.65
	Our ablated, $\mu$ only	no	30.64	28.65	29.57	29.62
	CBM3D	yes	31.82	30.40	30.68	30.96
	CBM3D	no	31.81	30.40	30.66	30.96
Gaussian $\sigma \in [5, 50]$	Baseline, N2C	no	32.57	31.29	31.27	31.71
	Baseline, N2N	no	32.57	31.29	31.26	31.70
	Our	yes	32.47	31.19	31.21	31.62
	Our	no	32.46	31.18	31.13	31.59
	Our ablated, diag. $\Sigma$	yes	31.59	30.06	30.54	30.73
	Our ablated, diag. $\Sigma$	no	31.58	30.05	30.45	30.69
	Our ablated, $\mu$ only	no	30.54	28.56	29.41	29.50
	CBM3D	yes	31.99	30.67	30.78	31.15
	CBM3D	no	31.99	30.67	30.72	31.13

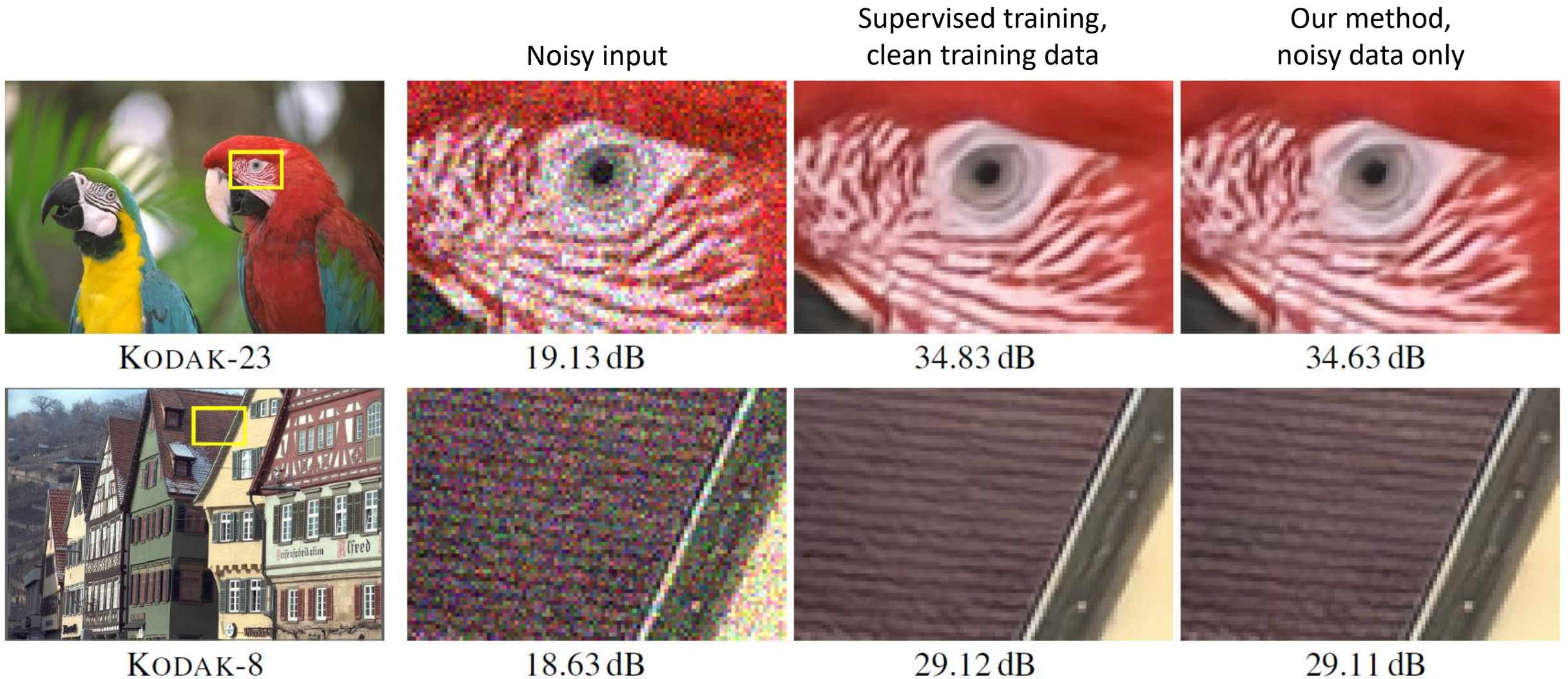
Supervised training with clean targets

Our result when  $\sigma$  is known vs. estimated from data

Noise2Void: Ignore  $y$  and predict based on context only

Close to baseline with variable noise ( $\sigma \in [5, 50]$ ) as well

# Results: Poisson noise ( $\lambda = 30$ )





# Results: Impulse noise ( $\alpha = 0.5$ )



KODAK-20



Noisy input

9.30 dB



Supervised training,  
clean training data

34.90 dB



Our method,  
noisy data only

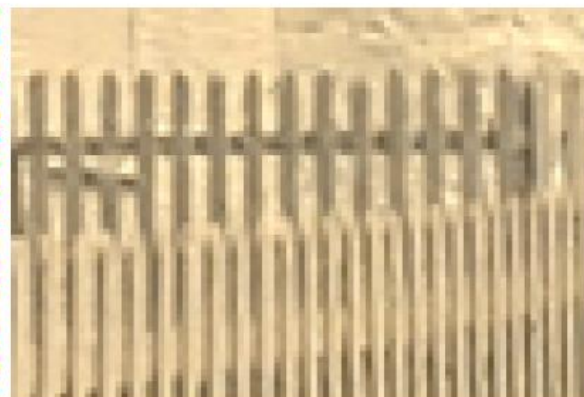
34.55 dB



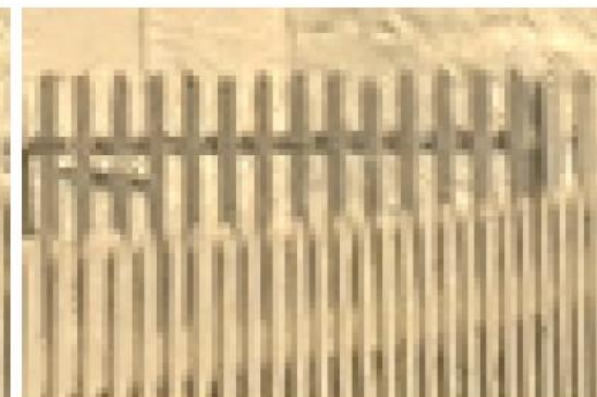
KODAK-19



12.09 dB



33.62 dB

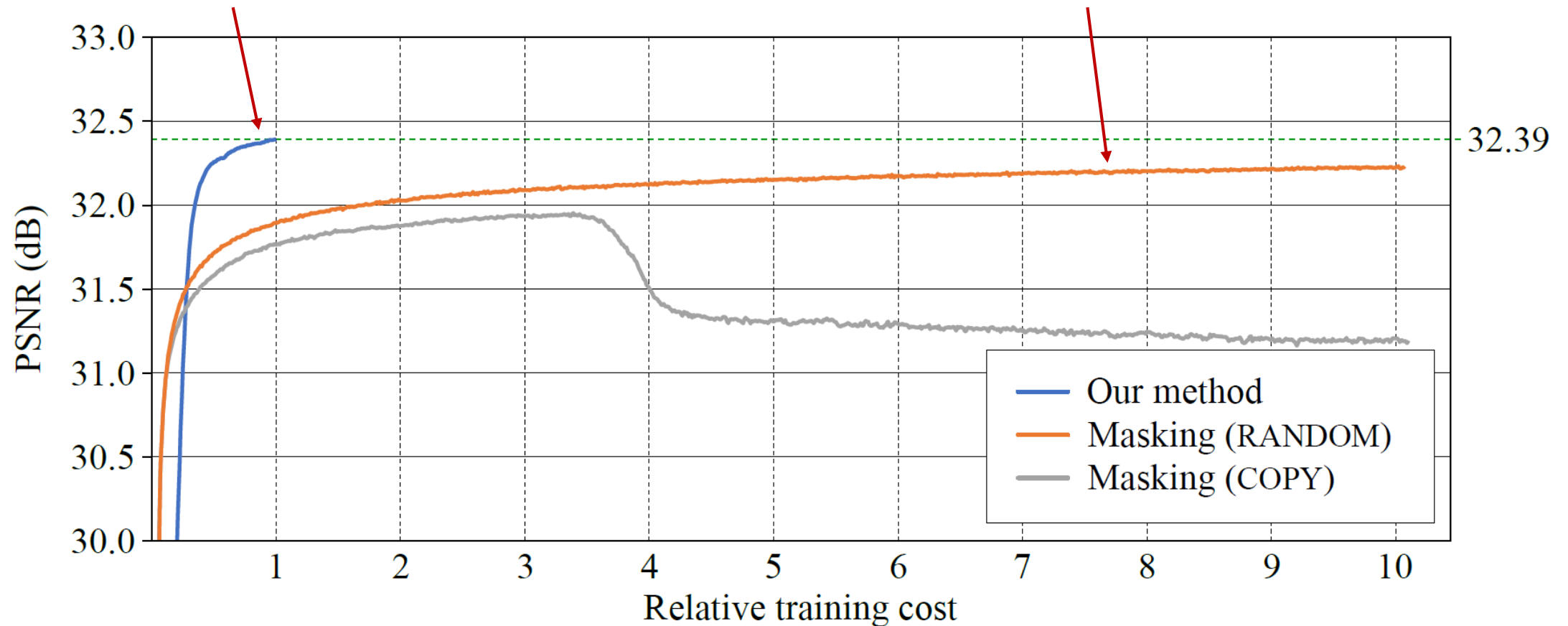


33.35 dB

# Evaluation of network architecture

Our blind-spot network architecture converges quickly

Standard network architecture with masking-based training [Krull et al., 2018]



# Conclusions

- Training high-quality denoisers is possible **with noisy data only**, when we have **just one** noisy realization of each training image
  - Can train a denoiser from a corpus of noisy data – no separate training set is required
- Result quality is comparable to traditionally trained networks
- Future work: Extend to more general corruptions?
  - Can we relax the assumption that noise is independent between pixels?

# Thank you

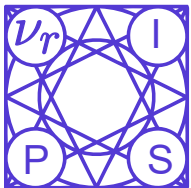
Paper: <https://arxiv.org/abs/1901.10277>

Code: <https://github.com/NVlabs/selfsupervised-denoising>

Feel free to contact with any questions: [slaine@nvidia.com](mailto:slaine@nvidia.com)

## References

- J. Batson and L. Royer. Noise2Self: Blind denoising by self-supervision. In *Proc. International Conference on Machine Learning (ICML)*, pages 524–533, 2019.
- A. Krull, T.-O. Buchholz, and F. Jug. Noise2Void – Learning denoising from single noisy images. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2129–2137, 2019.
- A. Krull, T. Vicar, and F. Jug. Probabilistic Noise2Void: Unsupervised content-aware denoising. *CoRR*, abs/1906.00651, 2019.
- J. Lehtinen, J. Munkberg, J. Hasselgren, S. Laine, T. Karras, M. Aittala, and T. Aila. Noise2Noise: Learning image restoration without clean data. In *Proc. International Conference on Machine Learning (ICML)*, 2018.
- D. A. Nix and A. S. Weigend. Estimating the mean and variance of the target probability distribution. *Proc. IEEE International Conference on Neural Networks (ICNN)*, pages 55–60, 1994.



**NeurIPS | 2019**

