

Learning Multi-Scale Local Conditional Probability Models of Images (top 25%)

Zahra Kadkhodaie¹ Florentin Guth² Stéphane Mallat^{3,4} Eero P. Simoncelli^{4,5}

¹CDS, New York University, zk388@nyu.edu

²DI, ENS, CNRS, PSL University, Paris, France

³Collège de France, Paris, France

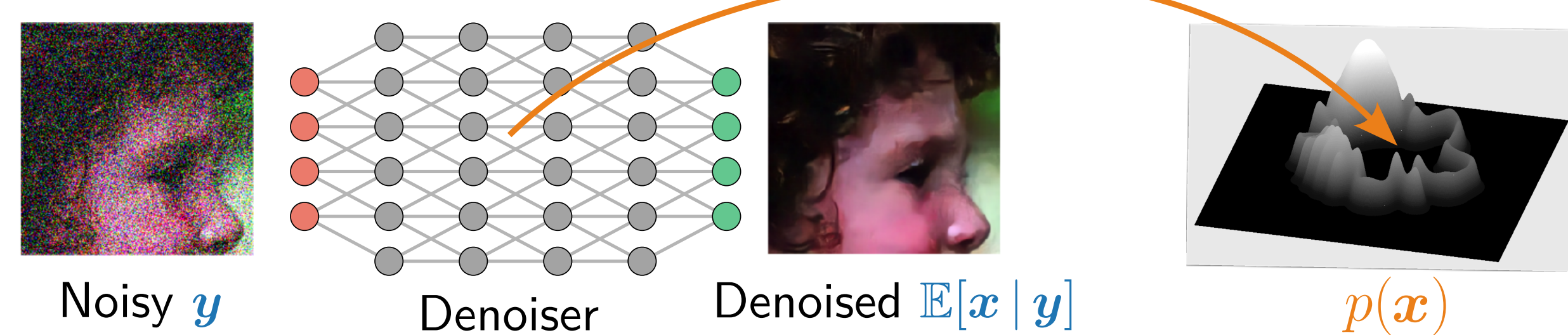
⁴Flatiron Institute, New York, USA

⁵CNS, Courant, and CDS, New York University

Summary

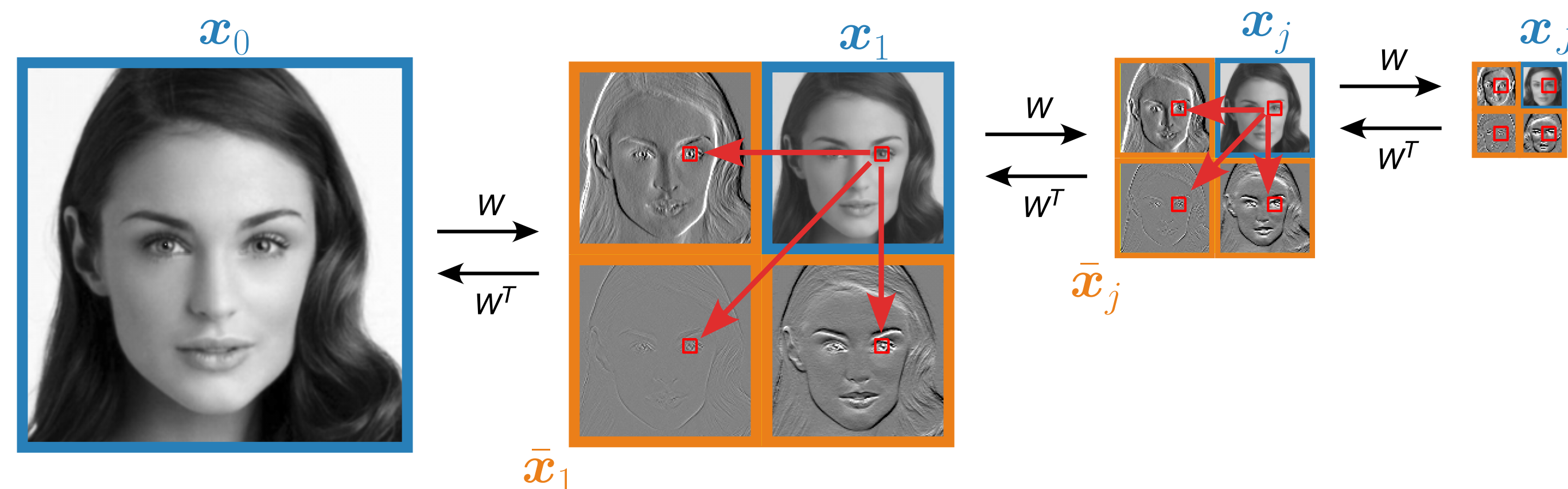
- How do score-based models manage to overcome the **curse of dimensionality**?
- Networks require global receptive fields to capture **global spatial dependencies**
- We show that these **global dependencies become local** after a **multiscale factorization** of the probability distribution
- We obtain high-resolution denoising, super-resolution, and synthesis results with **local conditional networks at each scale**

Denoising and Score-Matching



- A **denoiser** is a tool to learn **priors**: $\mathbb{E}[x|y] = y + \sigma^2 \nabla \log p(y)$ (Miyasawa 1961; Tweedie (via Robbins) 1956; Raphan & Simoncelli 2006; Vincent 2011)
- By using $\nabla \log p(y)$, we can synthesize images by doing **gradient ascent** on the log-probability
- This fails if the denoiser receptive field is smaller than the image size!

Wavelet Conditional Distributions

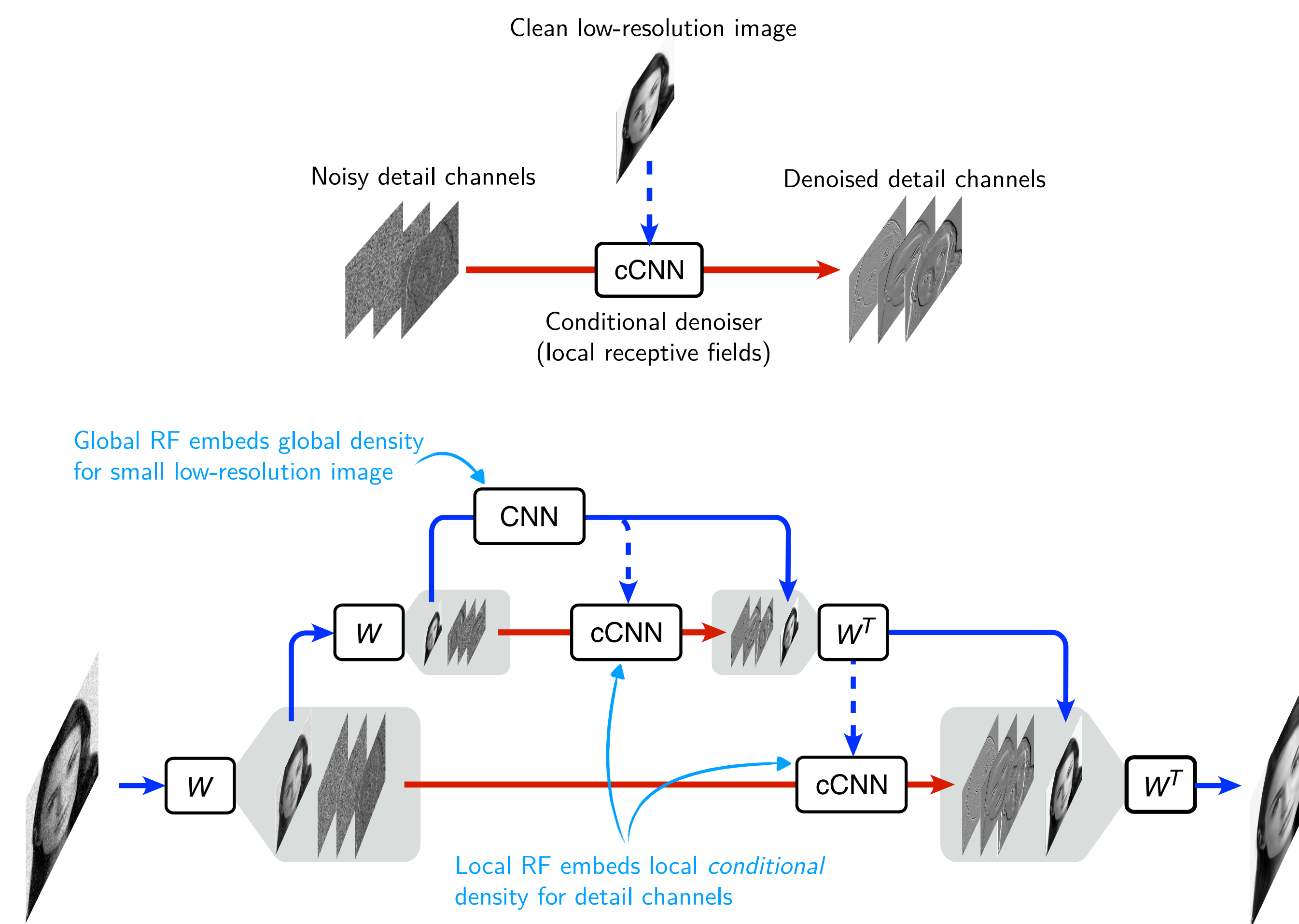


- The probability distribution of the **entire image** x_0 can be factorized as a product of **conditional probabilities at each scale**:

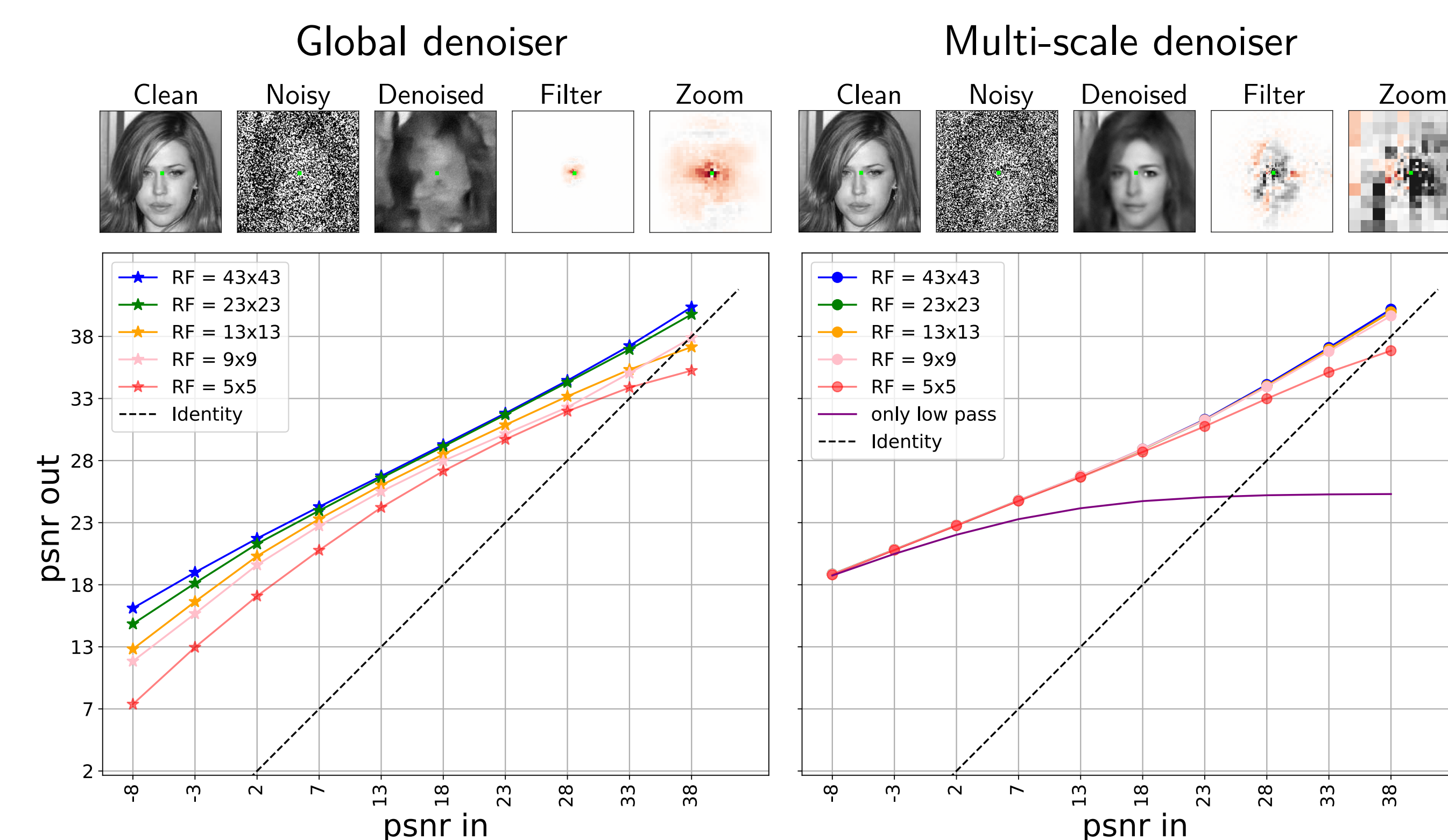
$$p(x_0) = p(x_1, \bar{x}_1) = p(x_1) p(\bar{x}_1|x_1) = p(x_j) \prod_{j=1}^J p(\bar{x}_j|x_j).$$

- This suggests first generating the **lowest-resolution image** x_j and iteratively increasing the resolution by **conditionally generating details** \bar{x}_j
- Theorem:** Restricting the receptive fields of the denoisers is equivalent to enforcing a **Markov property** on \bar{x}_j given x_j

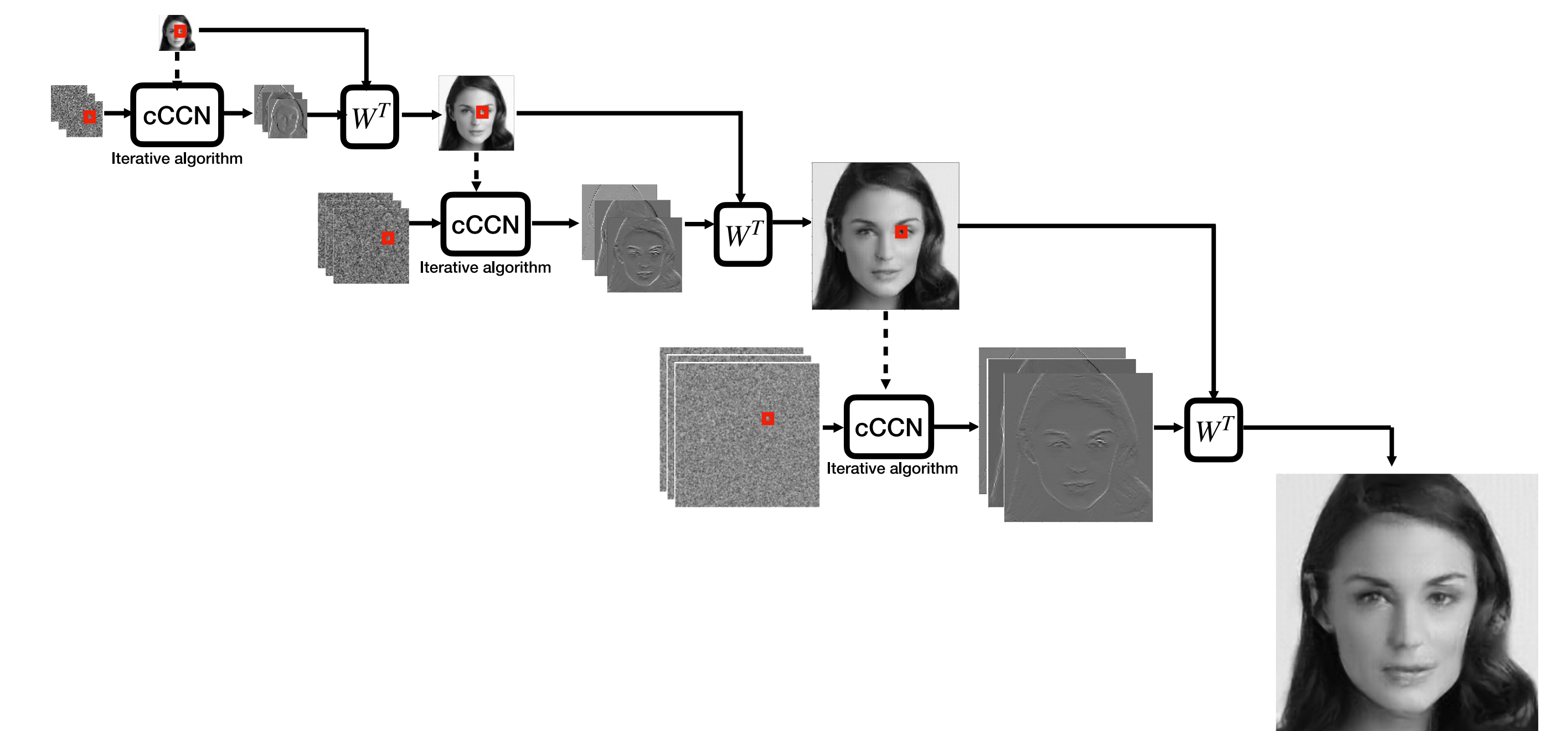
Multi-Scale Local Conditional Denoising



- We evaluate the **denoising performance** of the conditional denoisers
- Their receptive fields can be reduced to 9×9** without harming performance



Super-Resolution and Synthesis



- The conditional denoisers can be used to perform **super-resolution** and **image synthesis**
- The **global structure** is captured with a global prior on the **small low-pass image**

