

# Image Restoration

Giacomo Boracchi

CVPR USI, May 22 2020

[giacomo.boracchi@polimi.it](mailto:giacomo.boracchi@polimi.it)

<https://boracchi.faculty.polimi.it/>

# Inverse Problems

# Direct Problem





# The Inverse Problem



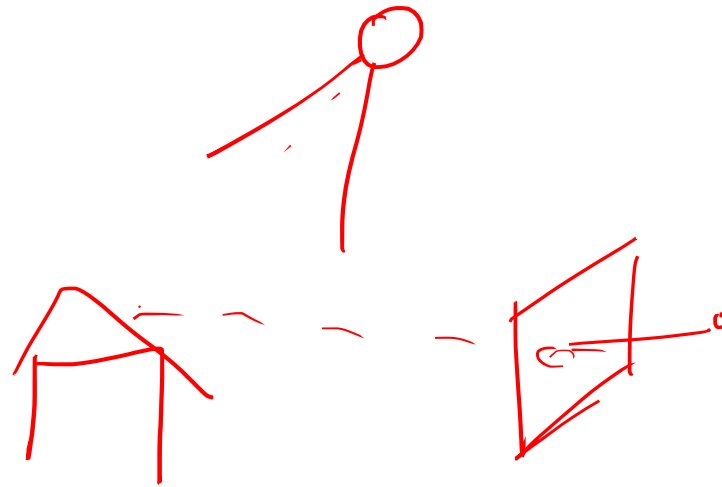


# Direct and Inverse problems

Often, direct problems can be addressed through simulation with respect to some known physical (natural) law.

An example related to CV:

- Image rendering



# Direct and Inverse problems

Inverse problems goes in the opposite direction and are much more difficult to solve.

They are often ill posed as they admit infinite many solutions

An example:

- 3D reconstruction from a single image



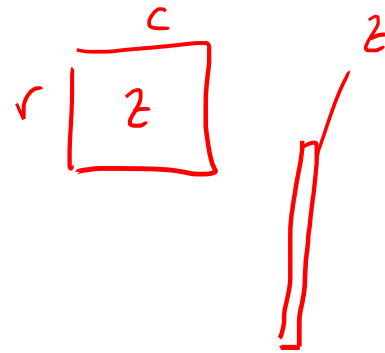


# Inverse problems in Imaging

One typically observe

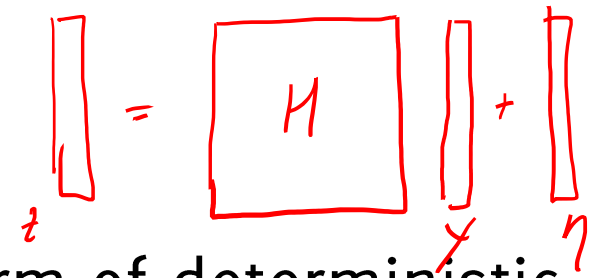
$$\underline{z} = H\underline{y} + \underline{\eta}$$

*Handwritten annotations:*  $\in \mathbb{R}^d$  above  $\underline{z}$ ,  $\in \mathbb{R}^d$  above  $\underline{y}$ ,  $\in \mathbb{R}^d$  above  $\underline{\eta}$ , and  $\mathbb{R}^d$  to the right of  $H$ .



Where  $y, z, \eta \in \mathbb{R}^d$  are images arranged as vectors that corresponds to

- $y$  : the unknown input image, which is to be recovered
- $z$  : the observation
- $\eta$  : the noise corrupting image acquisition



And  $H \in \mathbb{R}^{d \times d}$  is a linear operator describing some form of deterministic distortion during the image acquisition process

# Blur

Blurry images can be modeled

$$z = Hy + \eta$$

$\mathbb{R}^d$   $\mathbb{R}^{r \times c}$

Where  $Hy$  in its simplest form is spatially invariant  $Hy = (y \circledast h)$

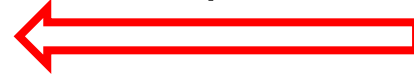
$y$



Forward problem



Inverse problem



*Deblurring*

$z$





# Blur

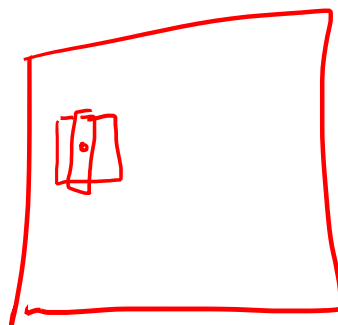
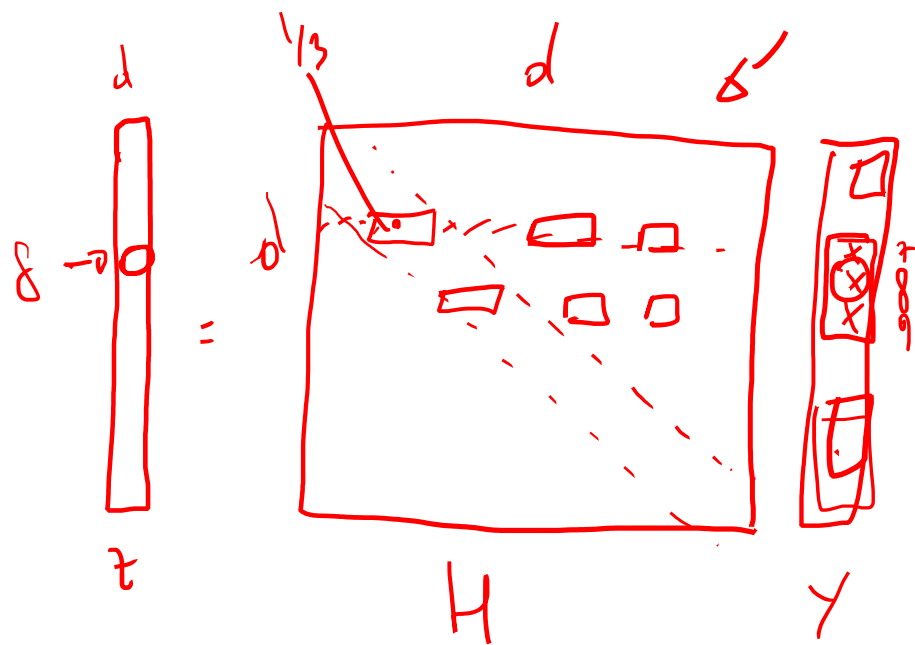
Blurry images can be modeled

$$z = Hy \quad y = H^{-1}z$$

$$z = Hy + \eta$$

Where  $Hy$  in its simplest form is spatially invariant  $Hy = (y \circledast h)$

What  $H$  corresponds to a spatially invariant blur operator?



$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

$$z = y * h$$

$$\underline{\underline{F(z)}} = \underline{\underline{F(y)}} \cdot \underline{\underline{F(h)}}$$

$$F(z) = \frac{F(y) + \eta}{F(h) \sigma}$$

# The Deblurring Problem

The convolution theorem states that

$$\text{Fourier}((f * g)) = FG$$

- Convolution can be computed by performing
  - Fourier transform,
  - Element-wise product
  - Inverse Fourier transform
- (and if FFT is possible, this is less expensive than computing convolutions in space domain)
- Blur can be easily inverted in Fourier domain

Note that the convolution theorem holds when the signal is periodic, and thus the circular convolution has to be computed



# The Deblurring Problem

Given the image  $y$  and the kernel  $h$  a blurred observation is given by

$$z = (y \circledast h)$$

then, when the filter  $h$  is exactly known, the convolutional blur can be inverted in Fourier Domain

$$Z = YH$$

thus,

$$\hat{y} = \text{Fourier}^{-1} \left( \frac{Z}{H} \right)$$

in order to define this ratio in frequencies where  $H = 0$

$$\hat{y} = \text{Fourier}^{-1} \left( \frac{Z\bar{H}}{\|H\|^2 + \epsilon} \right)$$

being  $\epsilon > 0$  a regularization parameter (difficult to choose in practice).

# The Deblurring Problem

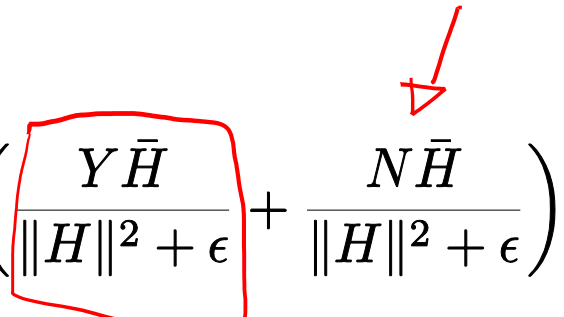
but what's up when noise appears?

$$z = (y * h) + \eta$$

Then, when computing the Fourier transform of the observation we have

$$Z = \underline{YH} + \underline{N}$$

Thus

$$\hat{y} = \text{Fourier}^{-1} \left( \frac{Y\bar{H}}{\|H\|^2 + \epsilon} + \frac{N\bar{H}}{\|H\|^2 + \epsilon} \right)$$


.. thus even unperceptible amount of noise, may become problematic in the second term of the sum.



Not all blurs are spatially invariant



# Inpainting

Blurry images can be modeled

$$z = Hy + \eta$$

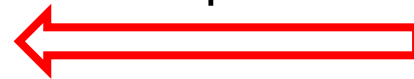
$H$



Forward problem



Inverse problem



# Inpainting

Blurry images can be modeled

$$z = Hy + \eta$$

What  $H$  corresponds to an inpainting operator?



# Denoising

The simplest of Inverse Problems

# Denoising: The Issue

A Detail in Camera Raw Image z



# Denoising: The Issue

Denoised  $\hat{y}$



# Denoising: The Issue

A Detail in Camera Raw Image z



# Denoising: The Issue

Denoised  $\hat{y}$





# Image Formation Model

$$z = \mathcal{H}x + \eta$$

*inverse problems*

Observation model is

$$z(x) = y(x) + \eta(x),$$

$$x \in \mathcal{X}$$

*$\mathcal{H}$  - colatht,*



# Image Formation Model

Observation model is

$$z(x) = y(x) + \underline{\eta(x)}, \quad x \in \mathcal{X}$$

Where

- $x$  denotes the pixel coordinates in the domain  $\mathcal{X} \subset \mathbb{Z}^2$
- $y$  is the original (noise-free and unknown) image
- $z$  is the noisy observation
- $\eta$  is the noise realization

For the sake of simplicity we assume AWG:  $\eta \sim N(0, \sigma^2)$  and  $\eta(x)$  independent realizations.

The noise standard deviation  $\sigma$  is also assumed as known.

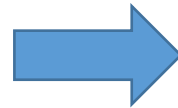
# Goal of Image Denoising

The goal of **image denoising** is to compute  $\hat{y}$  *realistic* estimate of the original image  $y$ , given the noisy observation  $z$

Denoising is an **ill posed problem** and requires some form of **regularization** to promote outputs that are close to natural images



$z$



$\hat{y} \approx y$

$$z = y + \eta$$

# Image Denoising

Denoising is a fundamental step in image processing pipelines

- Improves the quality of digital images to the standard we are used to
- Eases the following algorithms in imaging pipelines from those solving low-level (e.g., edge detection), till high-level (recognition) problems
- It is also a tool to quantitatively assess the performance of a descriptive model for images.

$$z = \underset{\uparrow}{\bar{y}} + \underset{\uparrow}{\bar{\eta}}$$

# Denoising as a Regression Problem

# Image Denoising

Estimating  $y(x)$  from  $z(x)$  can be statistically treated as regression of  $z$  given  $x$

$$\hat{y}(x) = E[z | x]$$

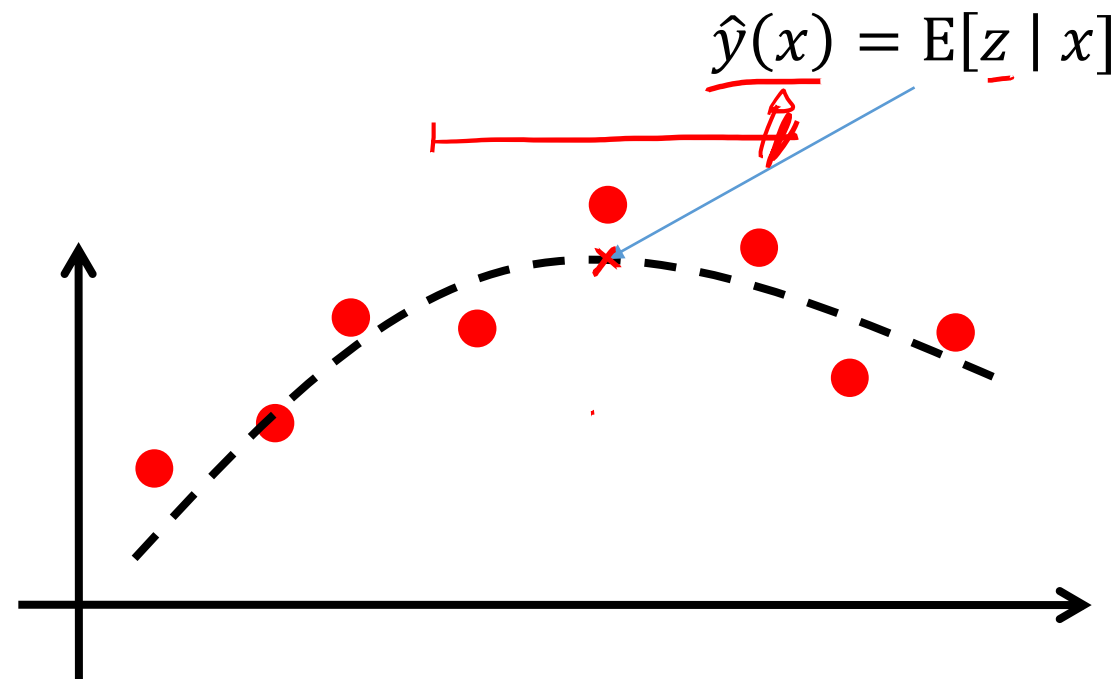


# Denoising and Regression

Observation model is

$$z(x) = y(x) + \eta(x), \quad x \in \mathcal{X}$$

Denoising can be formulated as a regression problem



# Denoising by Fitting a Constant Value

**The problem:** Estimate the constant  $C$  that minimizes the distance w.r.t. noisy observations in  $U_0$ , a neighborhood of  $x_0$

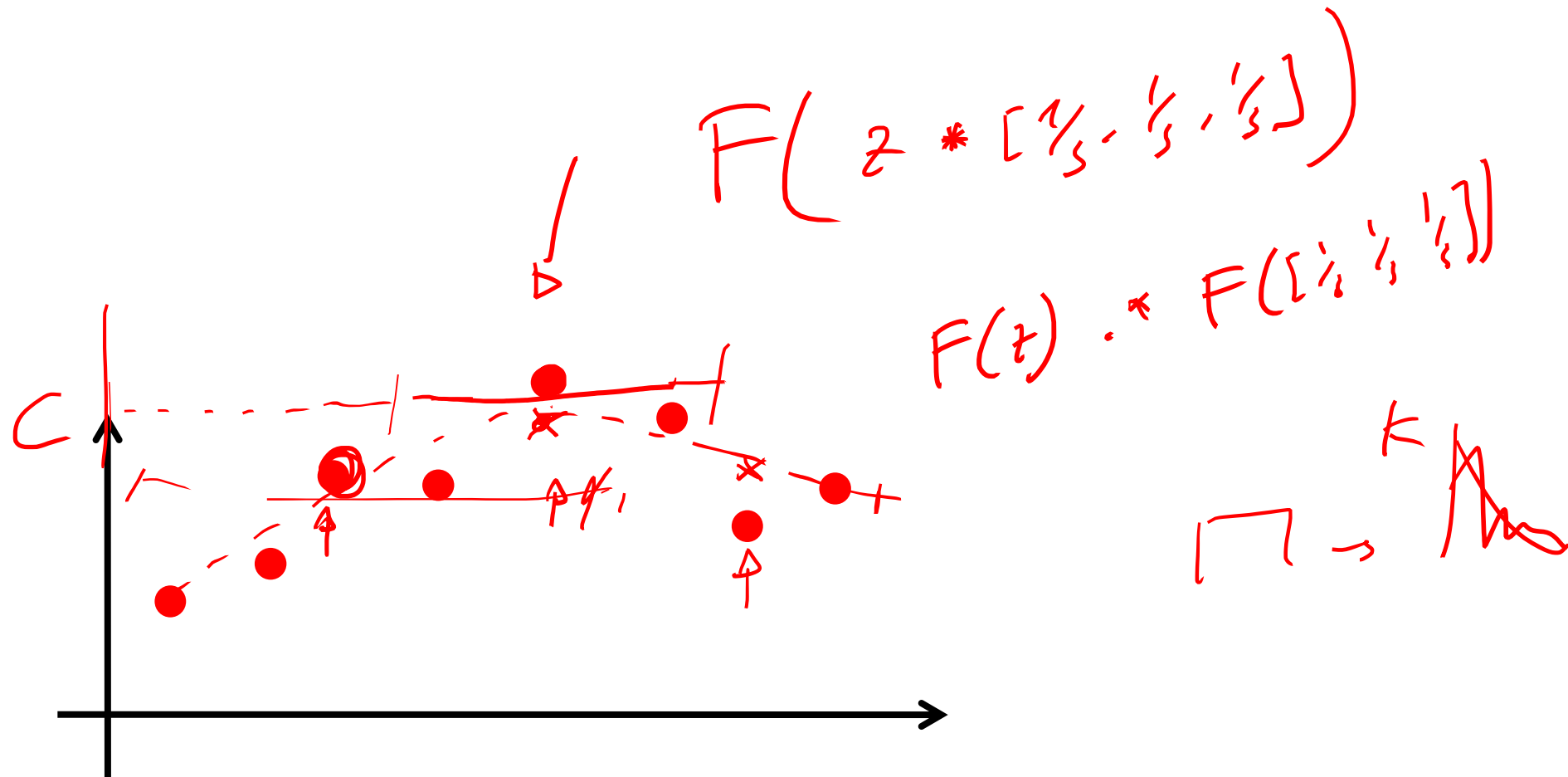
$$\hat{y}(x_0) = \underset{C}{\operatorname{argmin}} \sum_{x \in U_0} (z(x) - C)^2$$

There are many unbiased estimators (namely such that  $E_\eta[\hat{y}] = y$ ), the BLUE (Best Linear Unbiased Estimator) is the averaging noisy samples  $U_0$ , namely:

$$\hat{y}(x_0) = \frac{1}{\#U_0} \sum_{x \in U_0} z(x)$$

# Denoising by Fitting a Constant Value

The above method replaces each noisy input  $z(x)$  with the average over a fixed neighborhood  $U_x$  centered in  $x$ .



How would you implement this?

# Denoising by Weighted Fit of a Constant

The problem: Estimate the constant  $C$  that minimizes a weighted loss over noisy observations

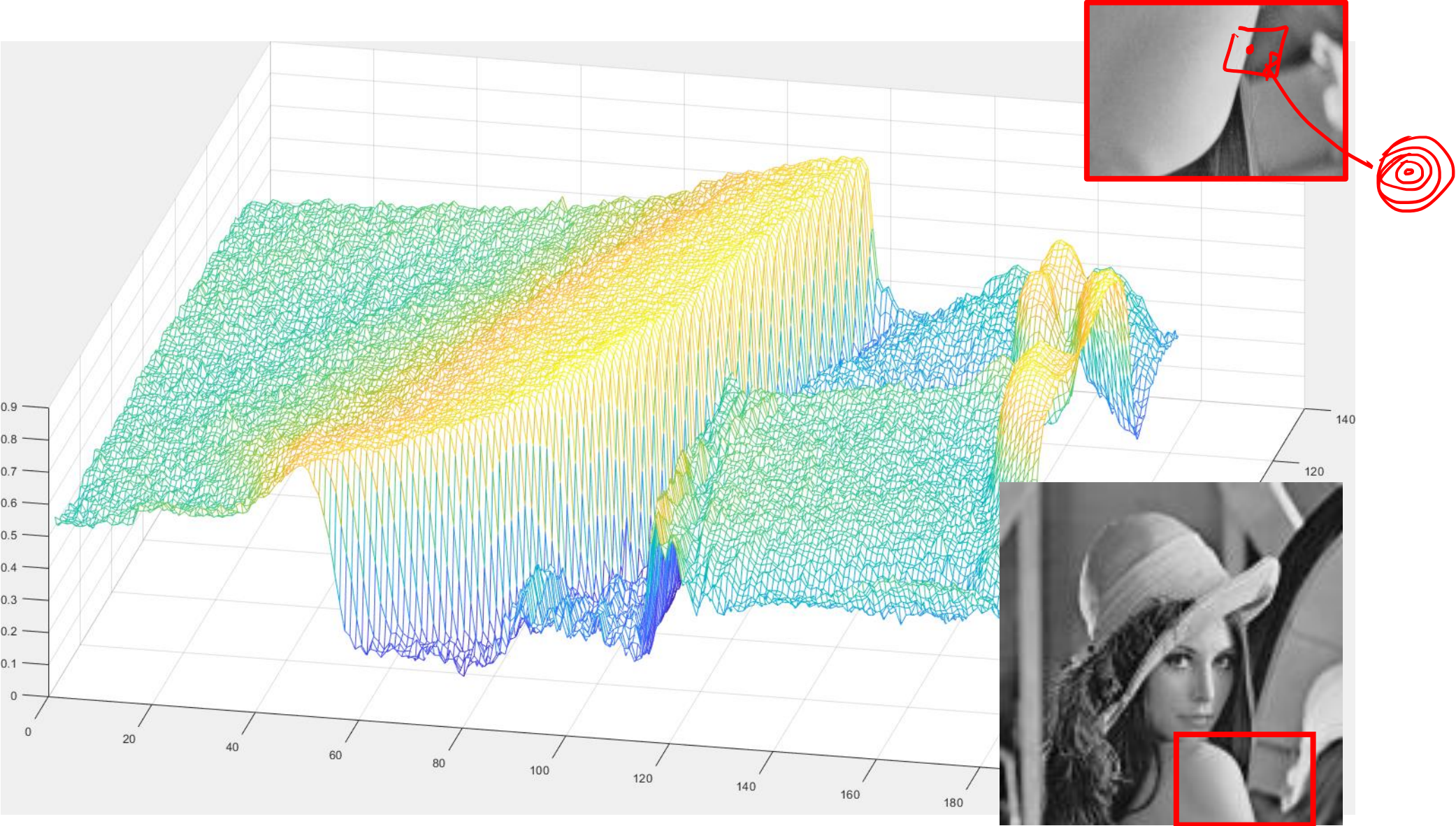
$$\hat{y}_h(x_0) = \underset{C}{\operatorname{argmin}} \sum_{x_s \in \mathcal{X}_0} w(x_0 - x_s) (z(x_s) - C)^2$$

$$w = \{w(x)\} \text{ s. t. } \sum_{x_s \in \mathcal{X}} w(x_s) = 1$$

This problem can be solved by computing the convolution of the image  $z$  against a filter whose coefficients are the error weights

$$\hat{y}(x_0) = (z \circledast w)(x_0)$$

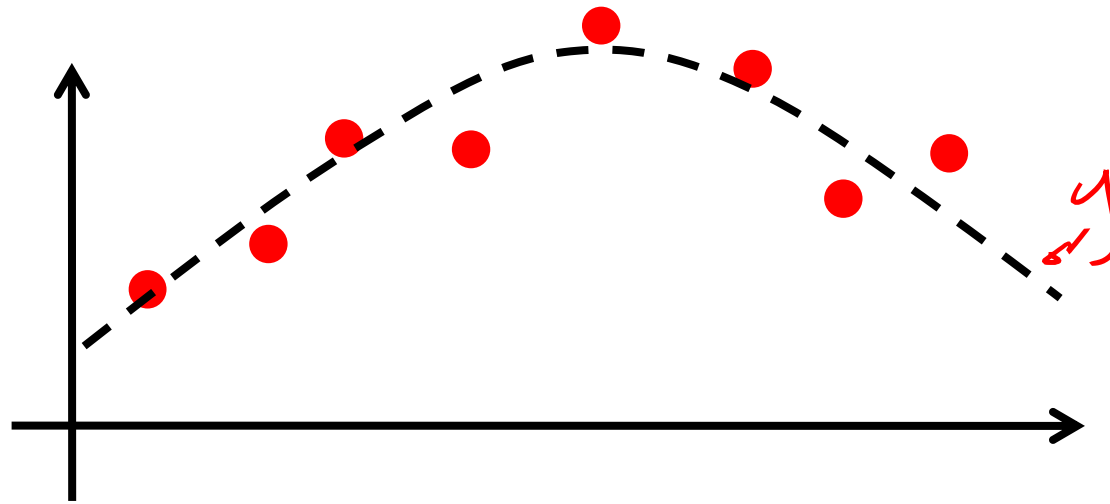
# Smoothing is Agnostic to Image Content





... of course

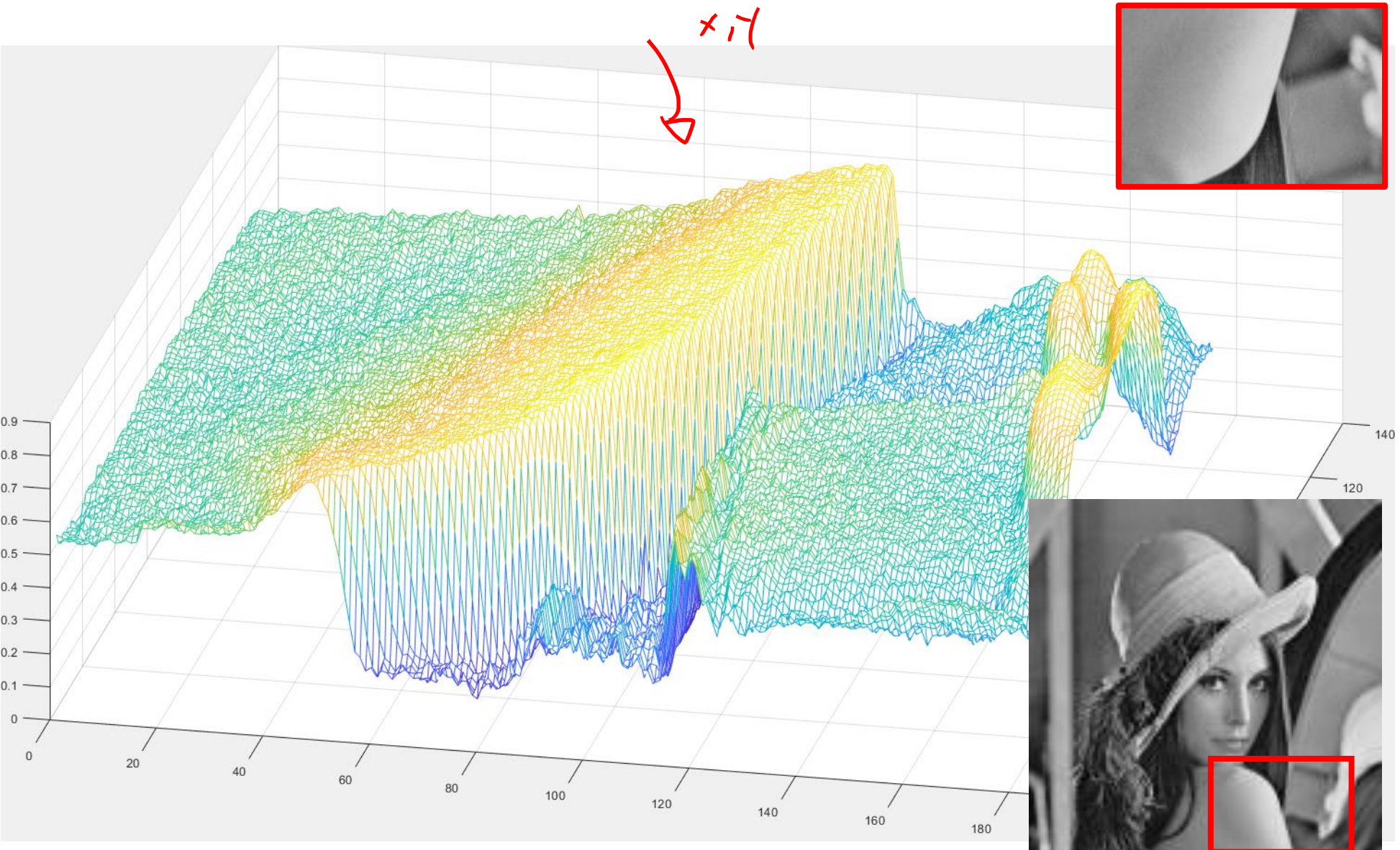
To perform regression more effectively, we can leverage a model describing the true signal  $y$



In this illustration, assuming the noise-free signal follows a polynomial trend might help!

We can pursue a “regression-approach”, but on images it is difficult to define a parametric expression for  $y$  on  $X$

# Difficult to Find a Good Model for Lena..



# Image Restoration

Giacomo Boracchi

CVPR USI, May 26 2020

[giacomo.boracchi@polimi.it](mailto:giacomo.boracchi@polimi.it)

<https://boracchi.faculty.polimi.it/>

# An Overview on Denoising Approaches

# Pointwise vs Multipoint Methods



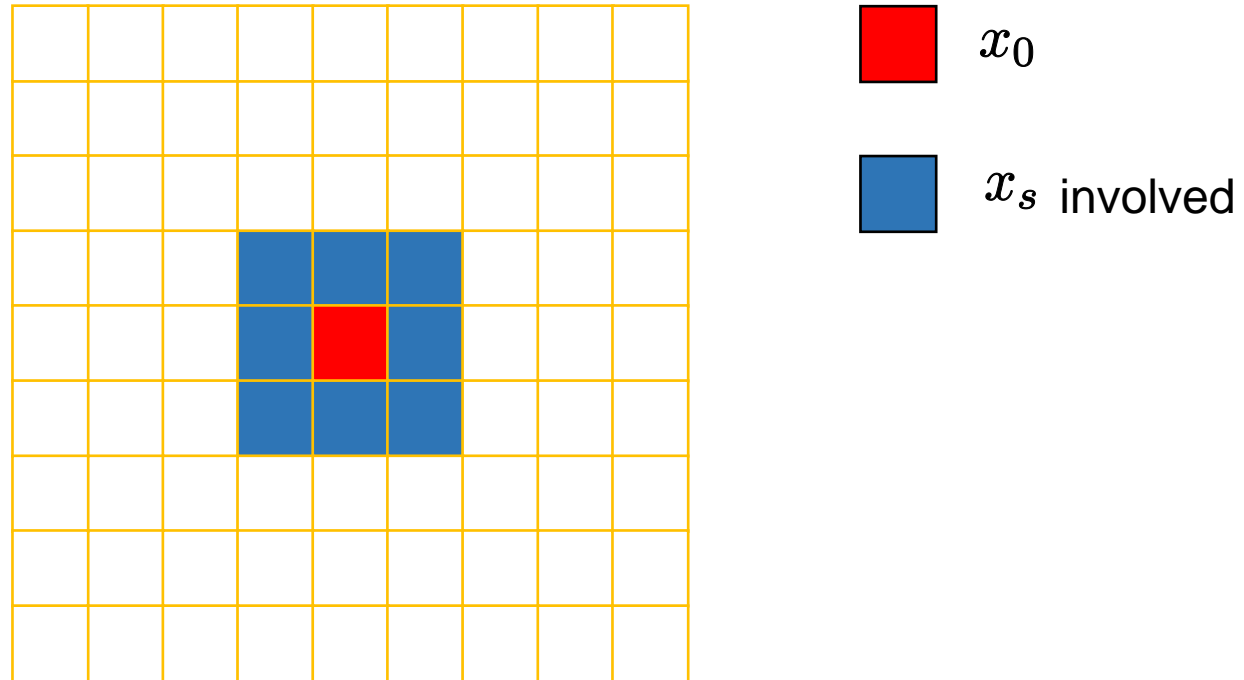
# Another view on Denoising Approaches

## Pointwise / Multipoint

- **Pointwise:** the estimation of noise-free signal is computed for the central point only  $y_0$ , and not for all the other points considered
- **Multipoint:** the estimation of the noise-free signal is computed for all the points  $y_s$  used by the estimator to estimate  $y_0$ .

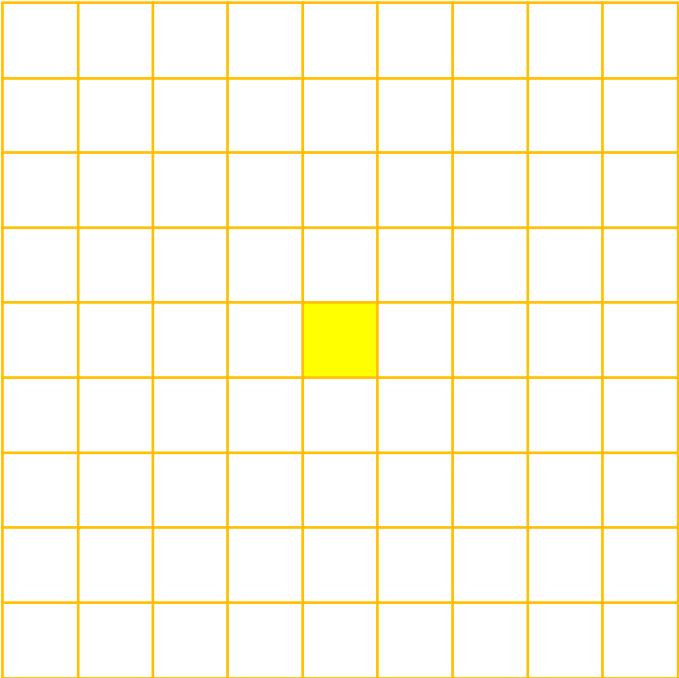
# Pointwise vs Multipoint

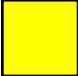
Pointwise: the estimate is given for the central point only



# Pointwise vs Multipoint

Pointwise, the estimate is given for the central point only

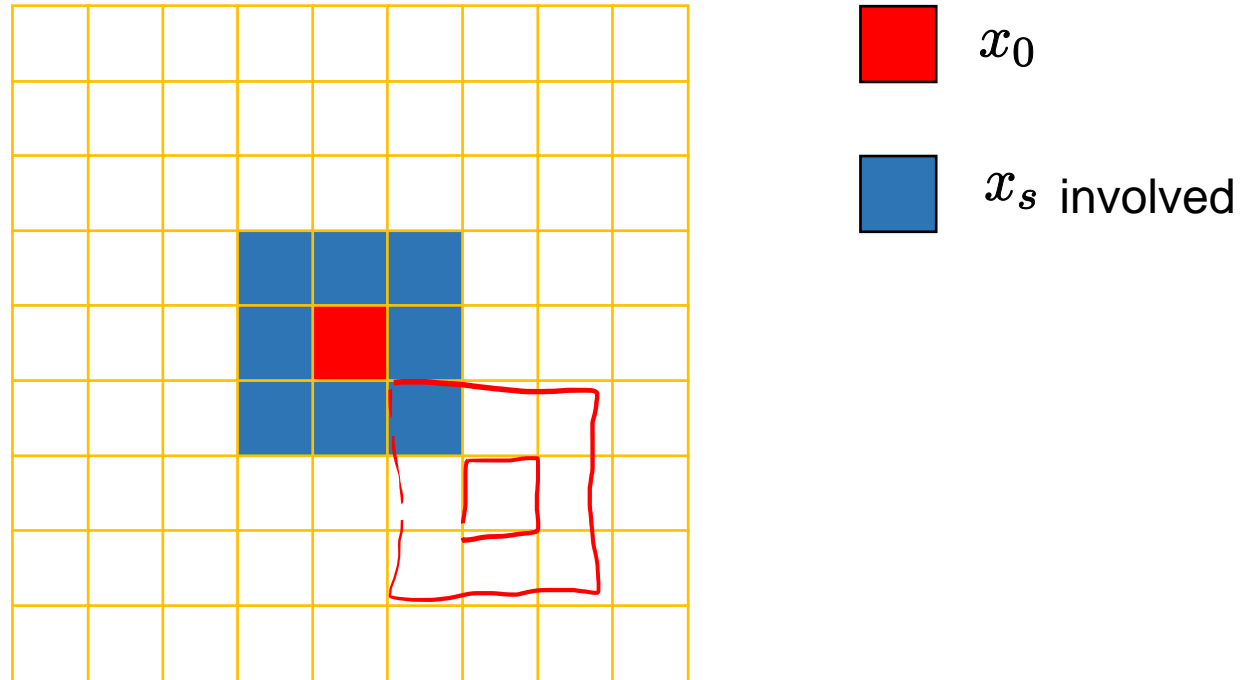


  $x_0$

Pixels where  
the true signal  
is estimated

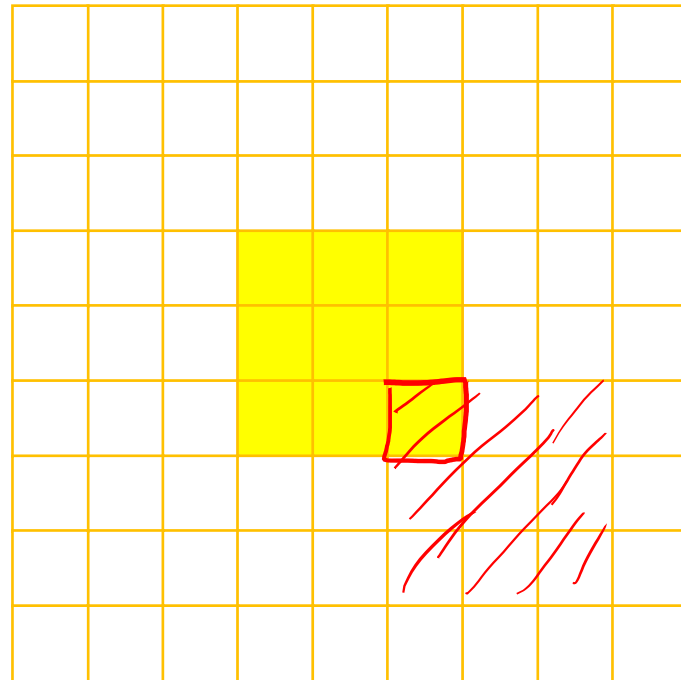
# Pointwise vs Multipoint

Multipoint, the original image is estimated in all the pixels considered in the filtering



# Pointwise vs Multipoint

Multipoint, the original image is estimated in all the pixels considered in the filtering



Pixels where  
the true signal  
is estimated



# Parametric vs Non-Parametric

# Denoising Approaches

## Parametric Approaches

- Assume the noisy-free signal  $y$  features some **sparsity property** in a suitable domain (e.g. Fourier, DCT, Wavelet) or w.r.t. some dictionary based decomposition. This implies that your image admits a **global parametric representation**

## Non Parametric Approaches

- Local Smoothing / Local Approximation
- Non Local Methods

*Bilateral filter / LPA-IC*

## Deep Learning


- Entirely data driven, these methods do not rely on an explicit prior for handling images

# Denoising Approaches

## Non Parametric Approaches

- Local Smoothing / Local Approximation
- Non Local Methods

Estimating  $y(x)$  from  $z(x)$  is treated as regression of  $z$  given  $x$

$$\hat{y}(x) = E[z \mid x]$$


# Bilateral Filter

# Bilateral Filter

The denoised image  $\hat{y}$  is a weighted average of (in principle) all the pixels

$$\hat{y}(x_1) = \frac{1}{W_1} \sum_{x_2 \in \mathcal{X}} w(x_1, x_2) z(x_2), \quad \forall x_1 \in \mathcal{X}$$

*weights depend on  $x_1$  and  $x_2$*   
 *$x_2 \in \mathcal{U}_{x_1}$*

where weights  $\{w(x_1, x_2)\}$  are adaptively defined depending on the image content and distance from pixel position

And  $W_1$  is a normalization factor for  $\hat{y}(x_1)$

$$W_1 = \sum_{x_2 \in \mathcal{X}} w(x_1, x_2)$$

# Weights in Bilateral Filters

$$f: \mathbb{R} \mapsto \mathbb{R}^+$$

In particular, weights are function of photometric and spatial distance between  $x_1$  and  $x_2$

$$w(x_1, x_2) = f_r(z(x_1) - z(x_2)) g_s(\|x_1 - x_2\|)$$



Often functions  $f_r$  and  $g_s$  are Gaussians

$$f_r(x_1, x_2) = e^{-\frac{\|z(x_1) - z(x_2)\|^2}{2\sigma_r^2}}, \quad g_s(x_1, x_2) = e^{-\frac{\|x_1 - x_2\|^2}{2\sigma_s^2}}$$

Where the parameters  $\sigma_r$  and  $\sigma_s$  regulate the weight decay factors and compensate for different ranges (image intensity / space)



# Weights in Bilateral Filters

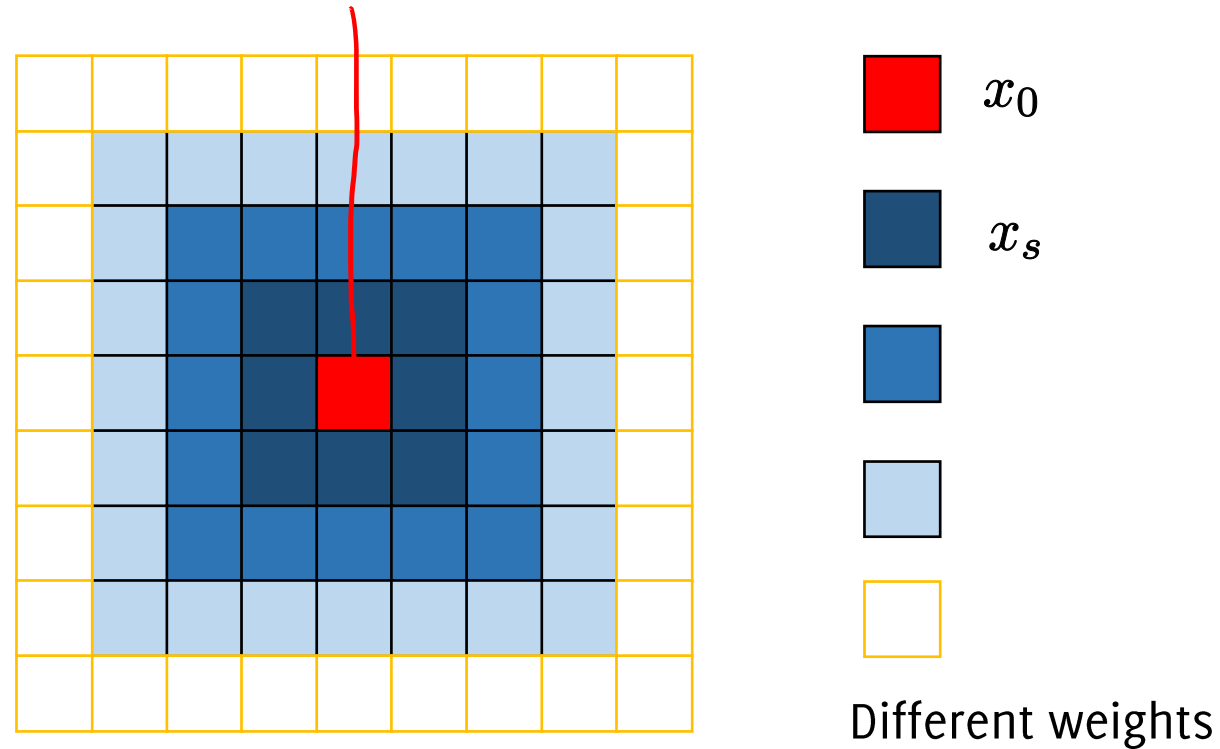
Weights take into account two different priors:

- Local similarity
- Non-local similarity

What do we obtain if we ignore the photometric contribution in the weight definition, namely we set  $f_r = 1$ ?

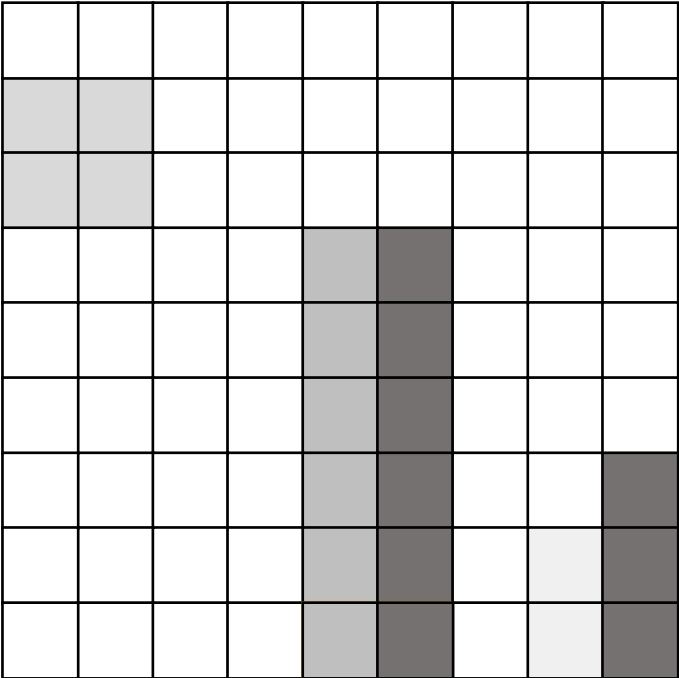
# Local vs Non-Local

Local: weights are determined by the pixel distance (regardless of the image content)



# Local vs Non-Local

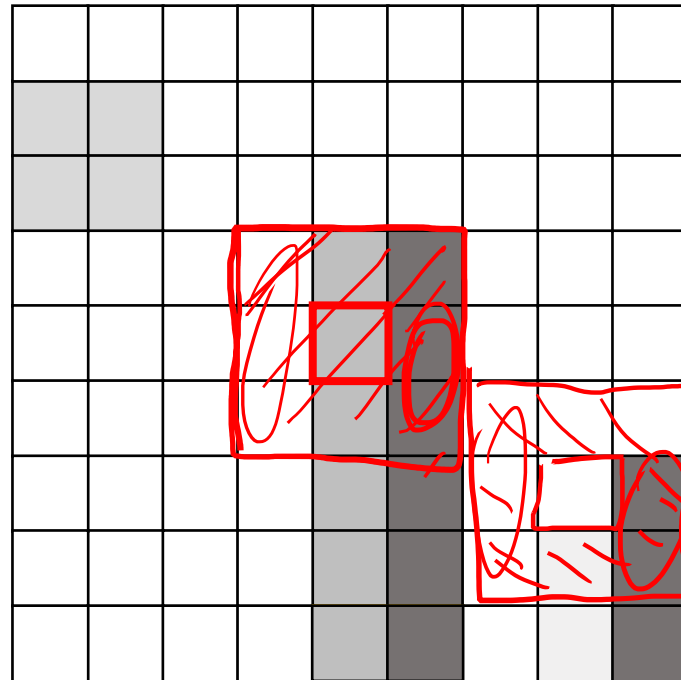
Non Local: weights are determined by the image similarity



Example of observation

# Local vs Non-Local

Non Local: weights are determined by the image similarity

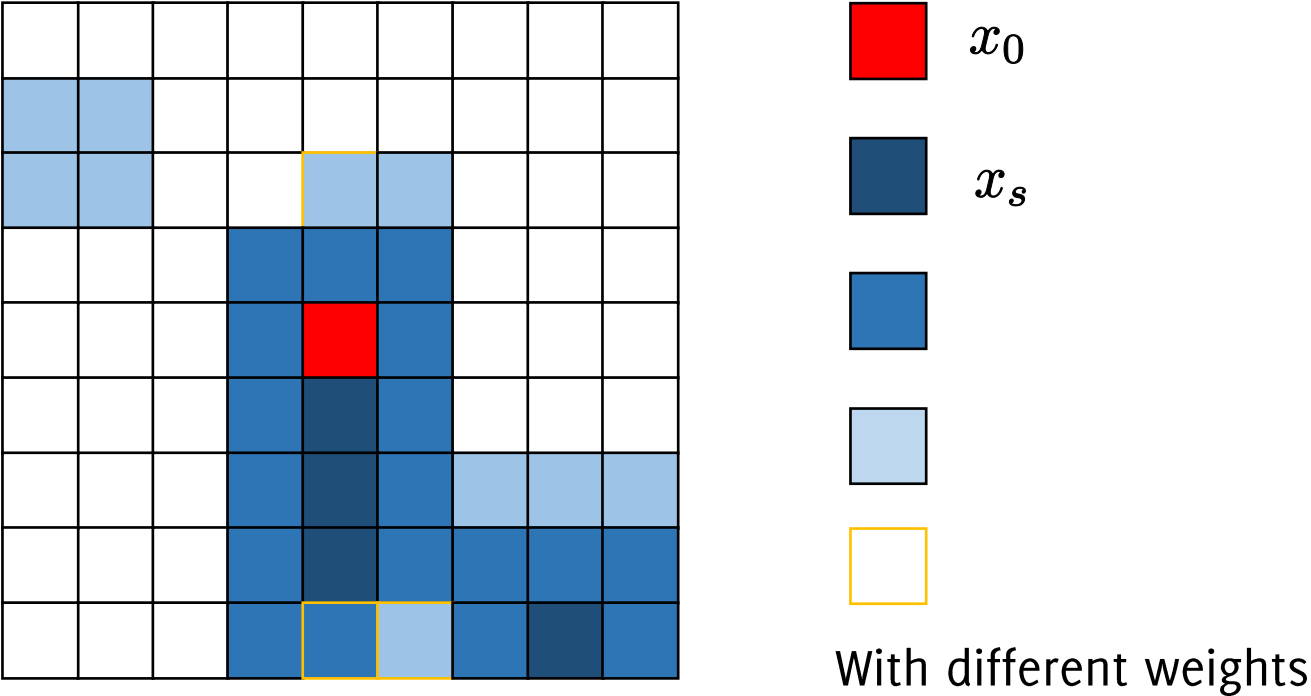


Example of observation



# Local vs Non-Local

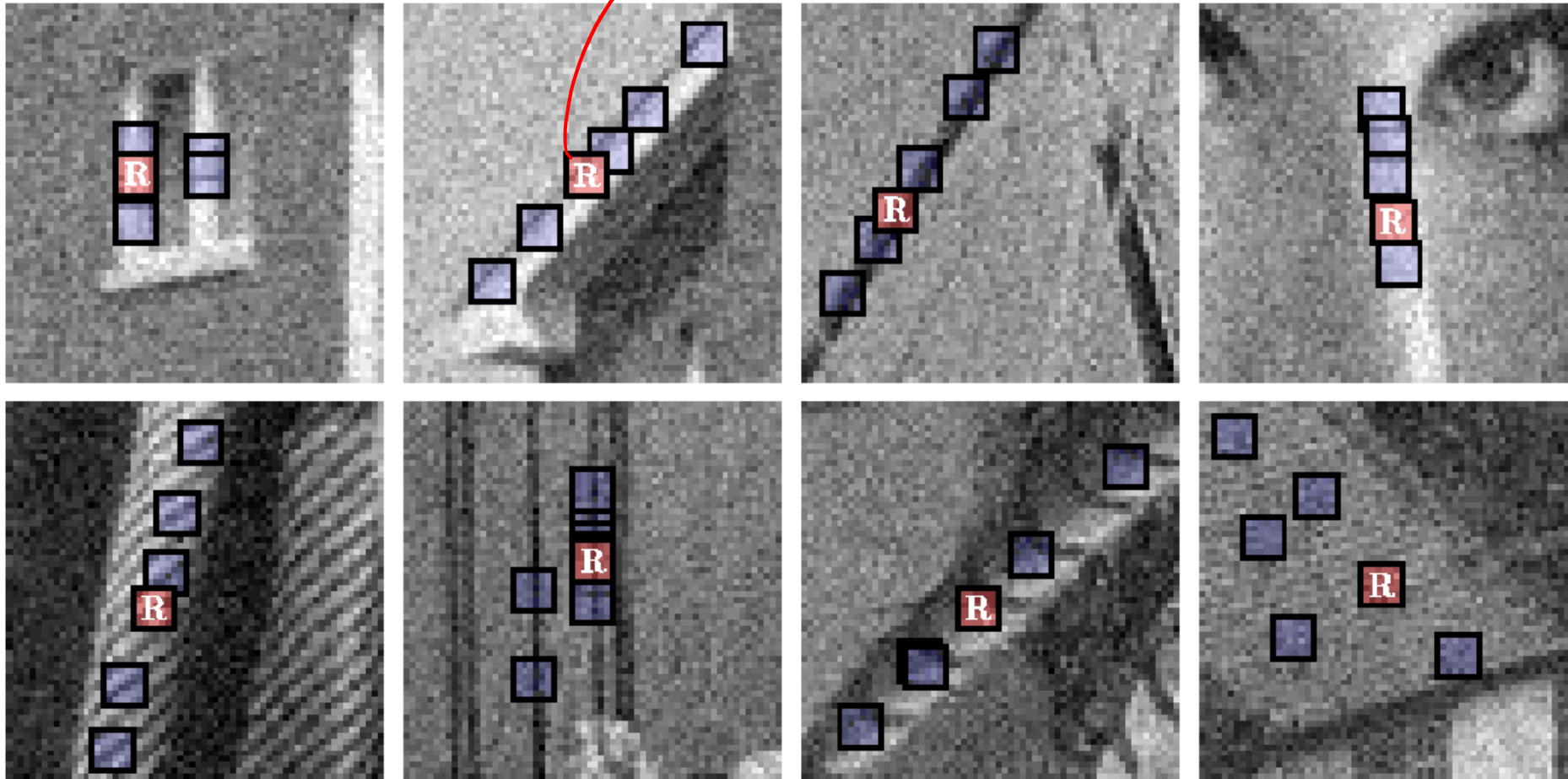
Non Local, weights are determined by the image similarity



# Nonlocal Self Similarity denoising



# NonLocal Self Similarity *Reference*



*In a natural image, for any given patch there exist **many** other **similar** looking **patches** at **different** **spatial** **locations**.*

# Nonlocal Self-similarity



# NonLocal Self Similarity in Image Processing

Traced back to **fractal models** of natural images (Barnsley, 1993) and fractal block coding (Jacquin, 1992)

*.. self-transformability on a blockwise basis...*

**Texture synthesis and completion** (Efros and Leung, 1999; Wei and Levoy, 2000).

Predicting the **central pixel** of a patch by exploiting the *long-range correlation* of natural images (Zhang and Wang, 2002)

Nonlocal self-similarity as an **effective regularity assumption** at the heart of many successful image **denoising algorithms** (NL-means, BM<sub>3D</sub>, etc.).

Nonlocal self-similarity was successfully used for **several image/video processing tasks**.

# Denoising methods based on self-similarity

These methods leverage **self similarity** of image patches as a form of **regularization** for distinguishing natural images from noise

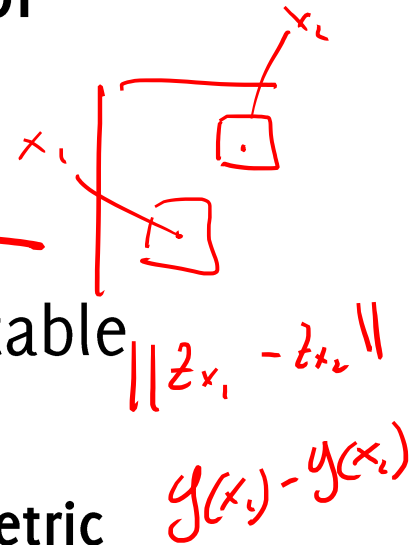
Self-similarity is assessed at patch-level and in a non-local manner.

Similar patches have to be correctly identified on the basis of a suitable patch distance measure.

- This is different from Bilateral Filter which considers only the photometric distance at pixel-wise level

Such a distance implies the assumption of a **specific descriptive model for natural images** and their self-similarity.

The denoising effectiveness actually depends on the validity of such underlying model.



# Image denoising (NL-Means)

# Patches

$$g_s(x_1, -x_2) = 1$$

Let  $U \subset \mathbb{Z}^2$  be a spatial neighborhood centered at the origin  $(0,0) \in \mathbb{Z}^2$ ,

- we define a patch centered at a pixel  $x \in X$  in the observation  $z$

$$\mathbf{z}_x(u) = \{z(x + u), \quad u \in U\}$$

- a patch centered at a pixel  $x \in X$  in the original image  $y$

$$\mathbf{y}_x(u) = y(x + u), \quad u \in U$$





# Patch Distance

## The idea:

The patch-wise distance correlates well with the distance of the two central pixels in the noisy-free patches

The idea is to assess the similarity between pixels  $y(x_1)$  and  $y(x_2)$  (not available), via the similarity of the corresponding noisy patches  $\mathbf{z}_{x_1}$  and  $\mathbf{z}_{x_2}$ .

Thus, weights based on photometric differences should be computed from  $\|\mathbf{z}_{x_1} - \mathbf{z}_{x_2}\|_2$  rather than  $z(x_1) - z(x_2)$ .



# Windowed Patch Distance in NL-means

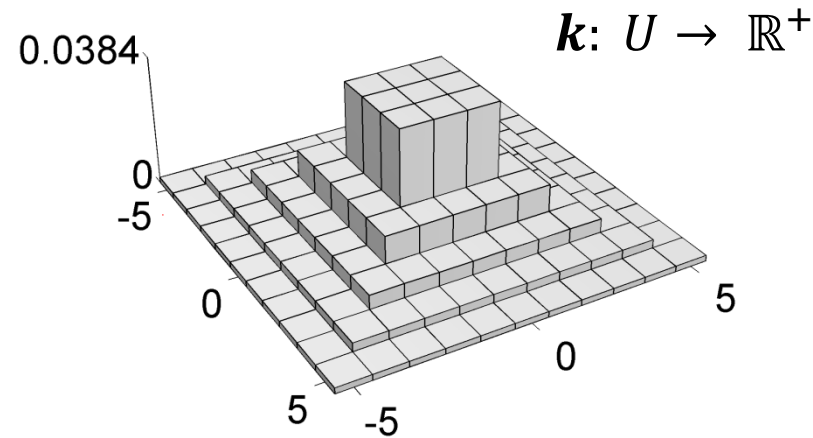
The distance operator is defined as a *windowed quadratic distance* between patches

$$d(x_1, x_2) = \sum_{u \in U} (z(x_1 + u) - z(x_2 + u))^2 \underline{\mathbf{k}(u)}$$

Where  $\mathbf{k}: U \rightarrow \mathbb{R}^+$  is a windowing kernel (we use  $\mathbf{k}$  since  $w$  is being used for the aggregation weights)

# Windowed patch distance in NL-means (cnt.)

The windowing kernel  $\mathbf{k}: U \rightarrow \mathbb{R}^+$  adjusts the contribution of each difference term depending on the position of  $u$  with respect to the patch center.



Hand-drawn diagram showing two patches,  $z_{x_1}$  and  $z_{x_2}$ , represented as red rectangles. A red line connects the centers of the two patches. Below the diagram, the following formulas are written in red:

$$\|z_{x_1} - z_{x_2}\|_2^2$$
$$\|(z_{x_1} - z_{x_2}) \sqrt{\mathbf{k}}\|_2^2$$

$d$  performs a **pixel-wise** comparison of the patches

the decay of  $\mathbf{k}$  reflects how much similarity between  $y(x_1)$  and  $y(x_2)$  may be implied from the similarity between  $y(x_1 + u)$  and  $y(x_2 + u)$  when  $u \neq 0$ .

# Non Local Means Filter (NL-means)

The denoised image  $\hat{y}$  is a weighted average of all image pixels

$$\hat{y}(x_1) = \sum_{x_2 \in X} w(x_1, x_2) z(x_2), \quad \forall x_1 \in X$$

where weights  $\{w(x_1, x_2)\}$  are adaptively defined depending on the similarity between two noisy patches  $\mathbf{z}_{x_1}$  and  $\mathbf{z}_{x_2}$

$$w(x_1, x_2) = \frac{e\left(-\frac{d(x_1, x_2)}{h^2}\right)}{\sum_X e\left(-\frac{d(x_1, x_2)}{h^2}\right)}$$

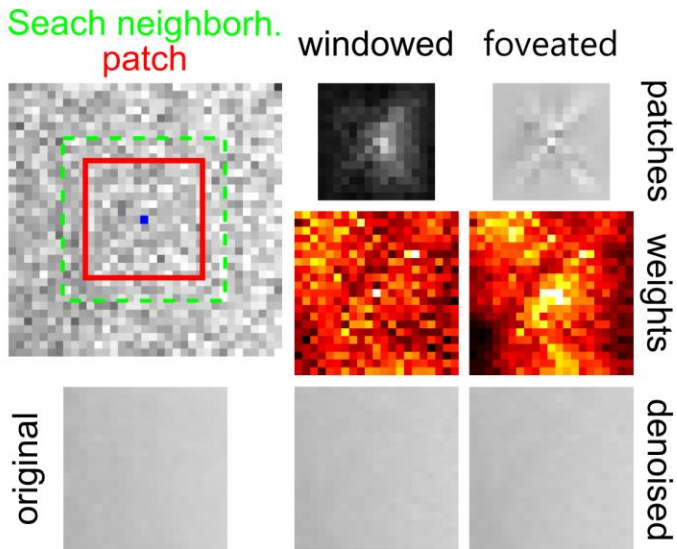
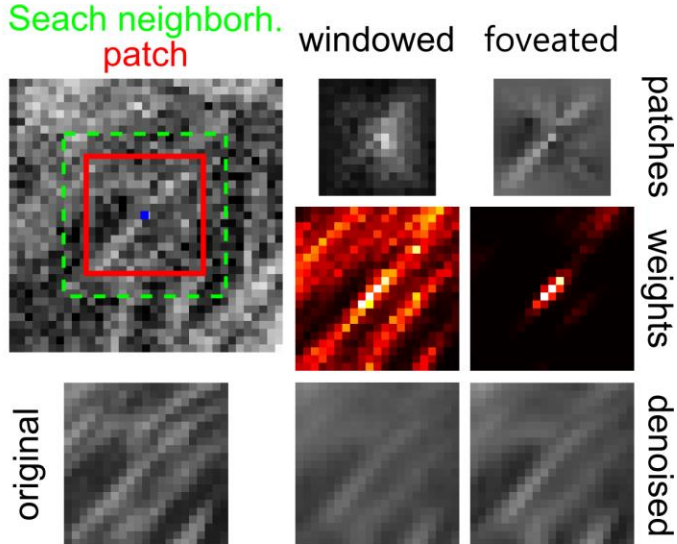
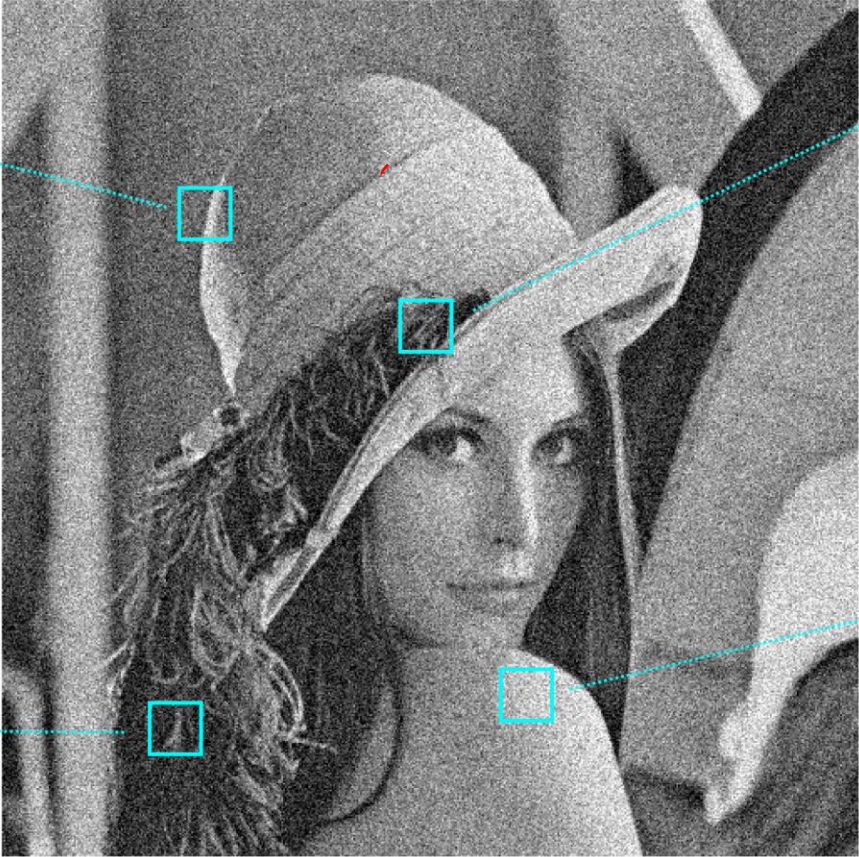
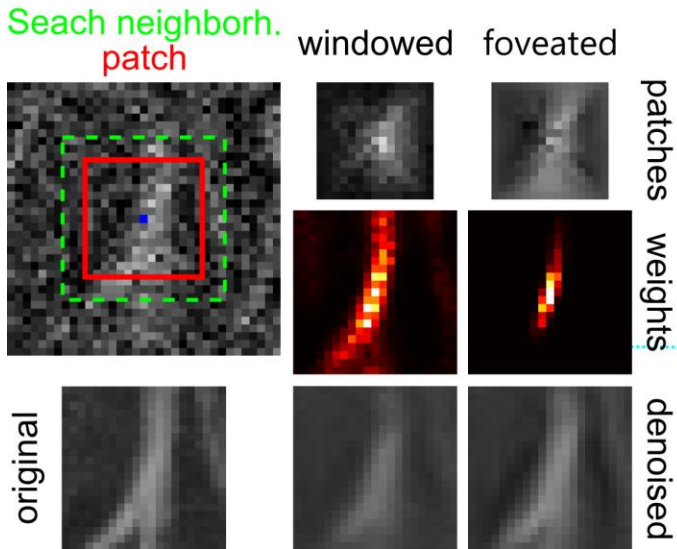
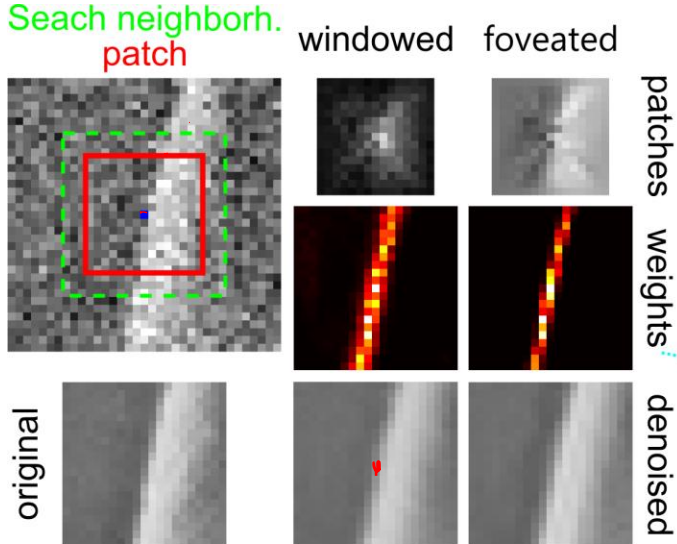
$d(x_1, x_2) = \|\mathbf{z}_{x_1} - \mathbf{z}_{x_2}\|_2$

$w_i(x_1)$

- $d(x_1, x_2)$ : distance measure between patches in  $x_1$  and  $x_2$ ,
- $h > 0$  is a smoothing parameter ( $h = \sigma$ ).
- $\mathbf{z}_{x_1}$  similar to  $\mathbf{z}_{x_2} \Rightarrow d(x_1, x_2)$  is small  $\Rightarrow w(x_1, x_2)$  large
- NL-means operates pixel-wise



# Weights from Patch-distances



# Spatially Adaptive Methods

# Denoising Approaches

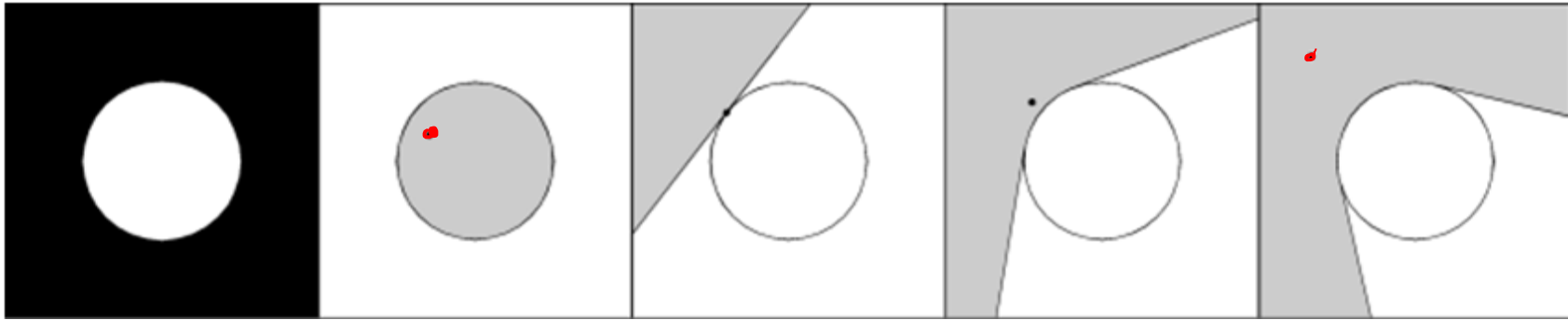
**Spatially adaptive methods**, The basic principle:

- there are no simple models able to describe the whole image  $y$ , thus perform the regression  $\hat{y}(x) = E[z | x]$
- Adopt a simple model in small image regions. For instance
$$\forall x \in X, \quad \exists \tilde{U}_x \text{ s. t. } y|_{\tilde{U}_x} \text{ is a polynomial}$$
- Define, in each image pixel, the “**best neighborhood**” where a simple parametric model can be enforced to perform regression.
- For instance, assume that on a suitable pixel-dependent neighborhood, where the image can be described by a polynomial



# Ideal neighborhood – an illustrative example

Ideal in the sense that it defines the support of a pointwise Least Square Estimator of the reference point.



Typically, even in simple images, every point has its own different ideal neighborhood.

For practical reasons, the ideal neighborhood is assumed starshaped

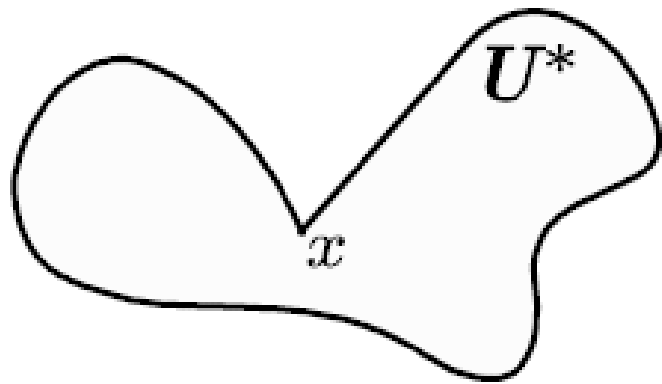
Further details at LASIP c/o Tampere University of Technology

<http://www.cs.tut.fi/~lasip/>

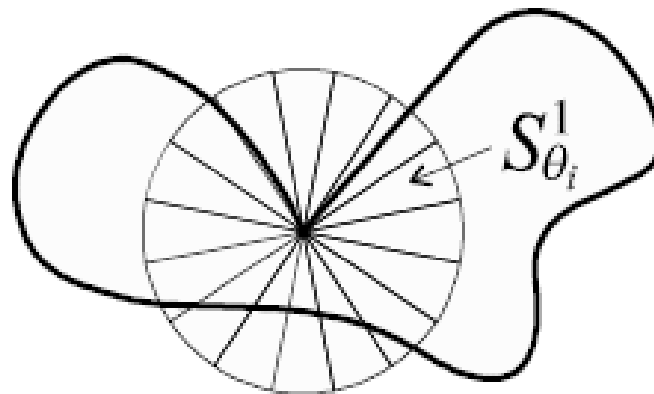
# Neighborhood Discretization

A suitable discretization of this neighborhood is obtained by using a set of directional LPA kernels  $\{g_{\theta,h}\}_{\theta,h}$ , where

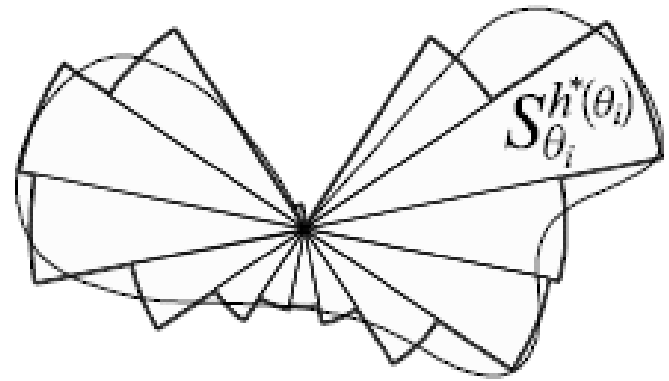
- $\theta$  determines the orientation of the kernel support,
- $h$  controls the scale of kernel support



Ideal Neighborhood



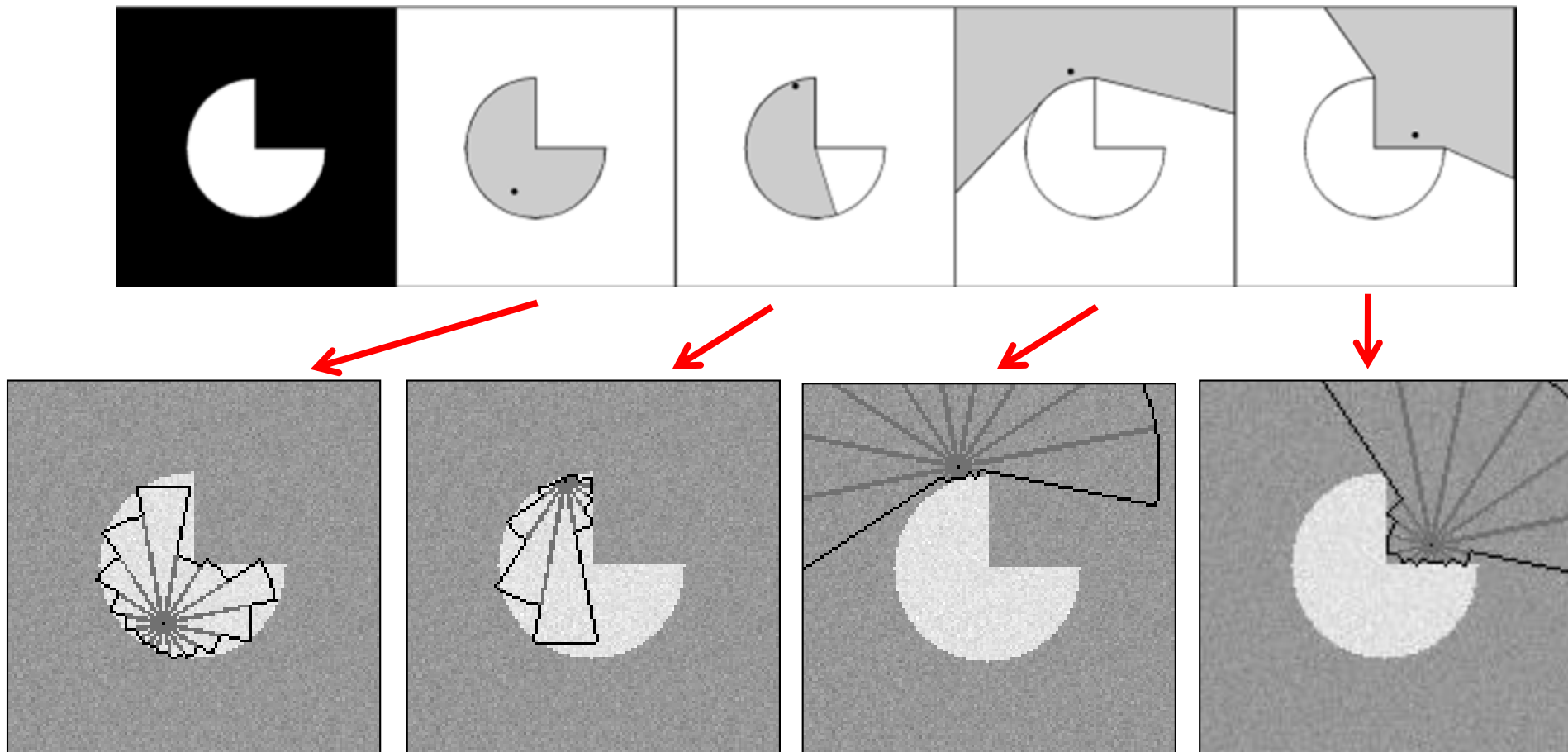
Directional kernels



Discrete Adaptive Neighborhood

# Ideal neighborhood – an illustrative example

Each point has an ideal neighborhood that defines the support of pointwise Least Square Estimator for the reference point.

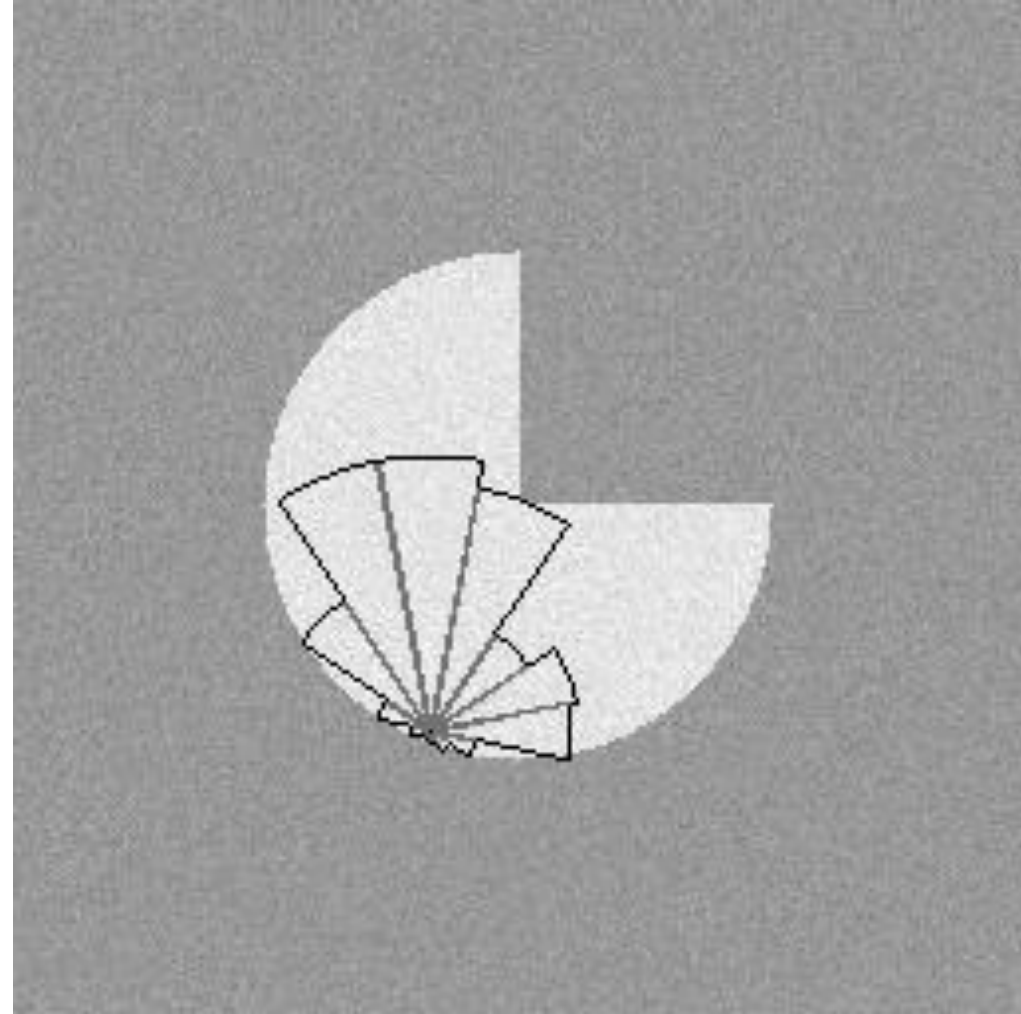


# Examples of Adaptively Selected Neighborhoods

Define,  $\forall x \in X$ , the “ideal” neighborhood  $\tilde{U}_x$  where to perform regression by fitting a given polynomial

Compute the denoised estimate at  $x$  by “using” only pixels in  $\tilde{U}_x$

$$\hat{y}(x) = E[z | x, \tilde{U}_x]$$



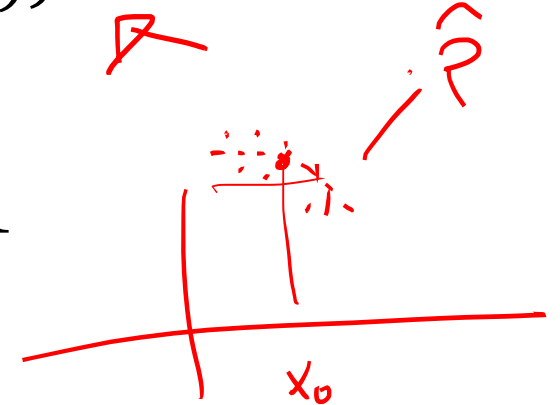
# Local Polynomial Approximation

The problem: Estimate the polynomial  $p$  of degree  $m$  that minimizes a weighted loss over noisy observations in a neighborhood of  $x_0$

$$\hat{p} = \operatorname{argmin}_{\underline{p \in \mathcal{P}_m}} \sum_{x_s \in \mathcal{X}} w(x_0 - x_s) (z(x_s) - p(x_s))^2$$

$$\text{where } w = \{w(x)\} \text{ s.t. } \sum_{x_s \in \mathcal{X}} w(x_s) = 1$$

$$\hat{y}(x_0) = \hat{p}(x_0)$$



This is directly solved by convolution of  $z$  against a filter  $g$  whose coefficients depend on the polynomial degree  $m$  and the weights  $w$

$$\hat{y}(x_0) = (z \circledast g)(x_0)$$

# Local Polynomial Approximation

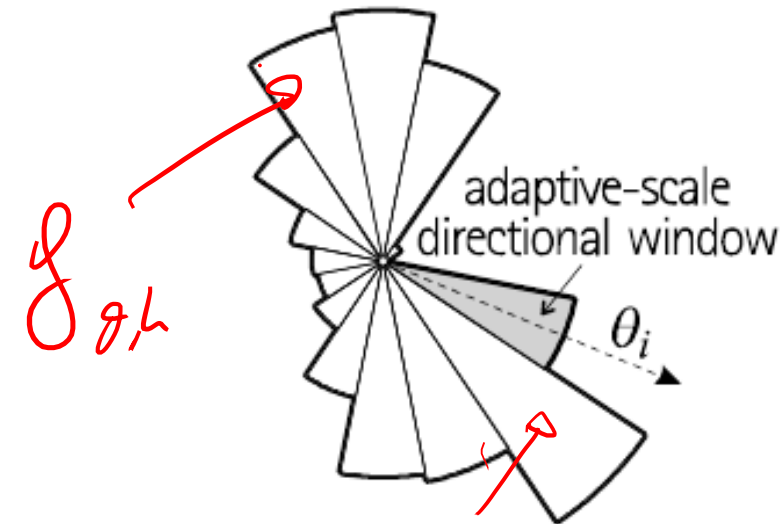
To achieve spatial adaptation we can consider multiple windows which are obtained by rotation and scaling of a basic window  $w$

Let  $w_{h,\theta} = w_\theta \left( \frac{\cdot}{h} \right)$  the scaled window over the direction  $\theta$

Consider multiple directions  $\theta \in \Theta$  and scales  $h \in H$  to define a collection of filters  $g_{\theta,h}$  yielding a batch of estimates

$$\hat{y}_{\theta,h}(x_0) = (z \circledast g_{\theta,h})(x_0)$$

Adaptation is performed by selecting, along each direction  $\theta$ , the best scale  $h^*$

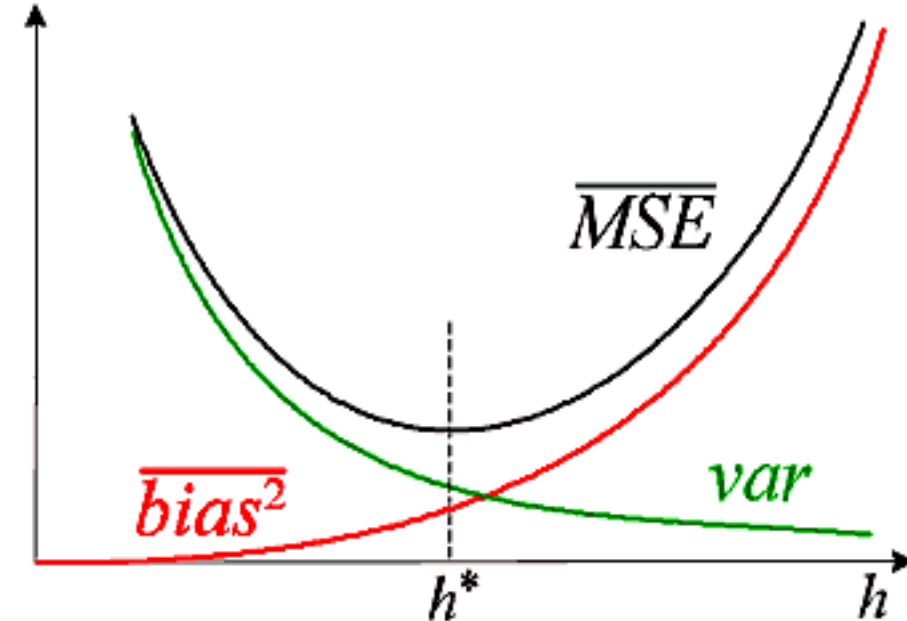


# Intersection of Confidence Interval Rule

Intersection of Confidence Interval (ICI) Rule is an adaptive scale selection criteria from statistical literature.

The scale parameter  $h$  controls the trade-off between bias and variance in the LPA estimates.

- Large  $h$  corresponds to a large window and smooth estimates, with lower variance and typically increased estimation bias.
- A small  $h$  corresponds to noisier estimates, less biased, and with higher variance.



# ICI Rule

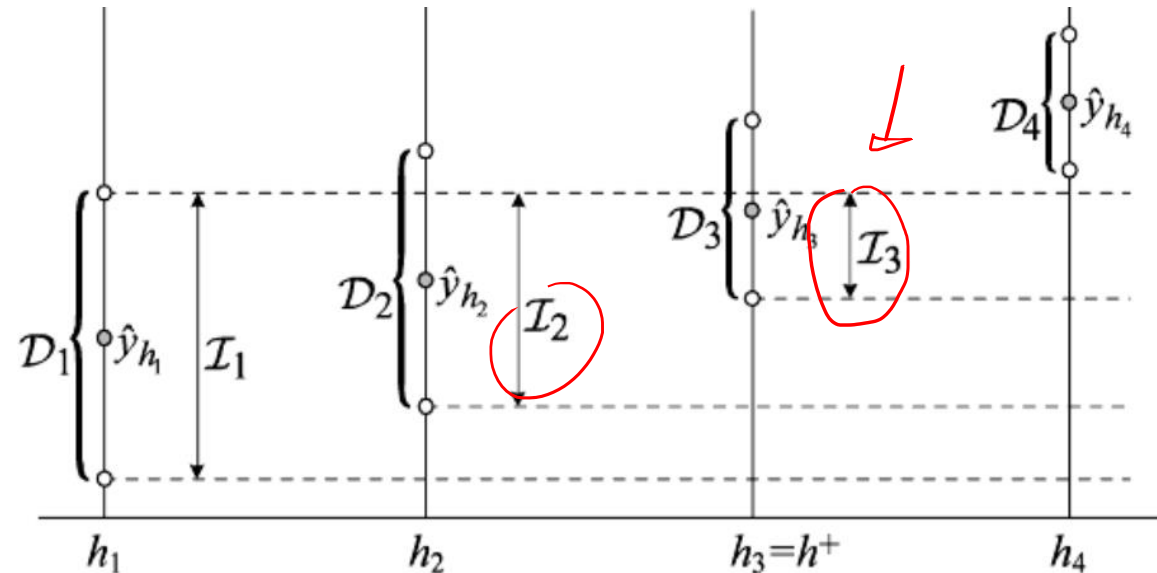
$\sigma$  is fixed

Consider the intersection  $\mathcal{I}_j$  of all the confidence intervals  $D_j$

$$\mathcal{I}_j = \bigcap_{i \leq j} D_i$$

Where  $D_i = [\hat{y}_i(x_0) - \Gamma\sigma_i, \hat{y}_i(x_0) + \Gamma\sigma_i]$  and  $\sigma_i$  is the standard deviation of the estimator associated to  $\hat{y}_i$ , and  $\Gamma > 0$  is a tuning parameter.

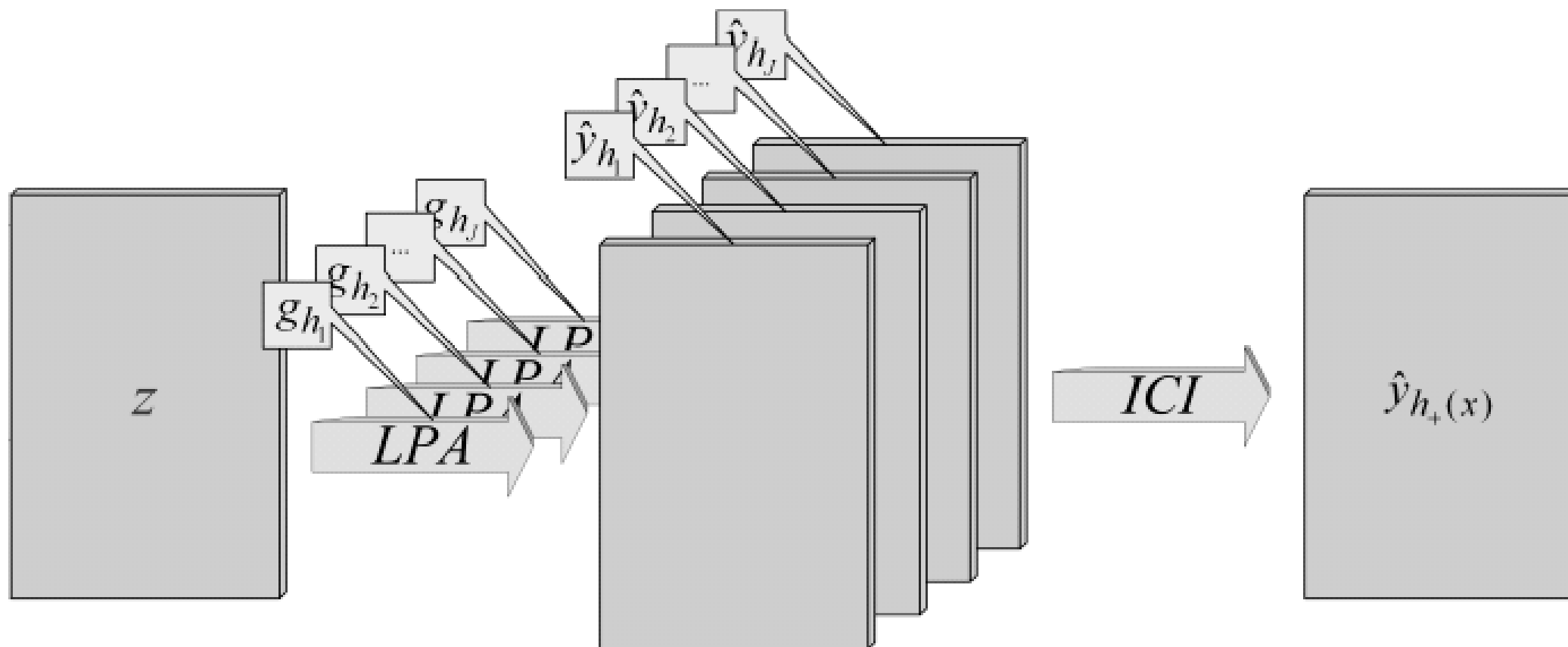
Then, the ICI rule selects the largest scale  $j^+$  such that  $\mathcal{I}_{j^+}$  is not empty.





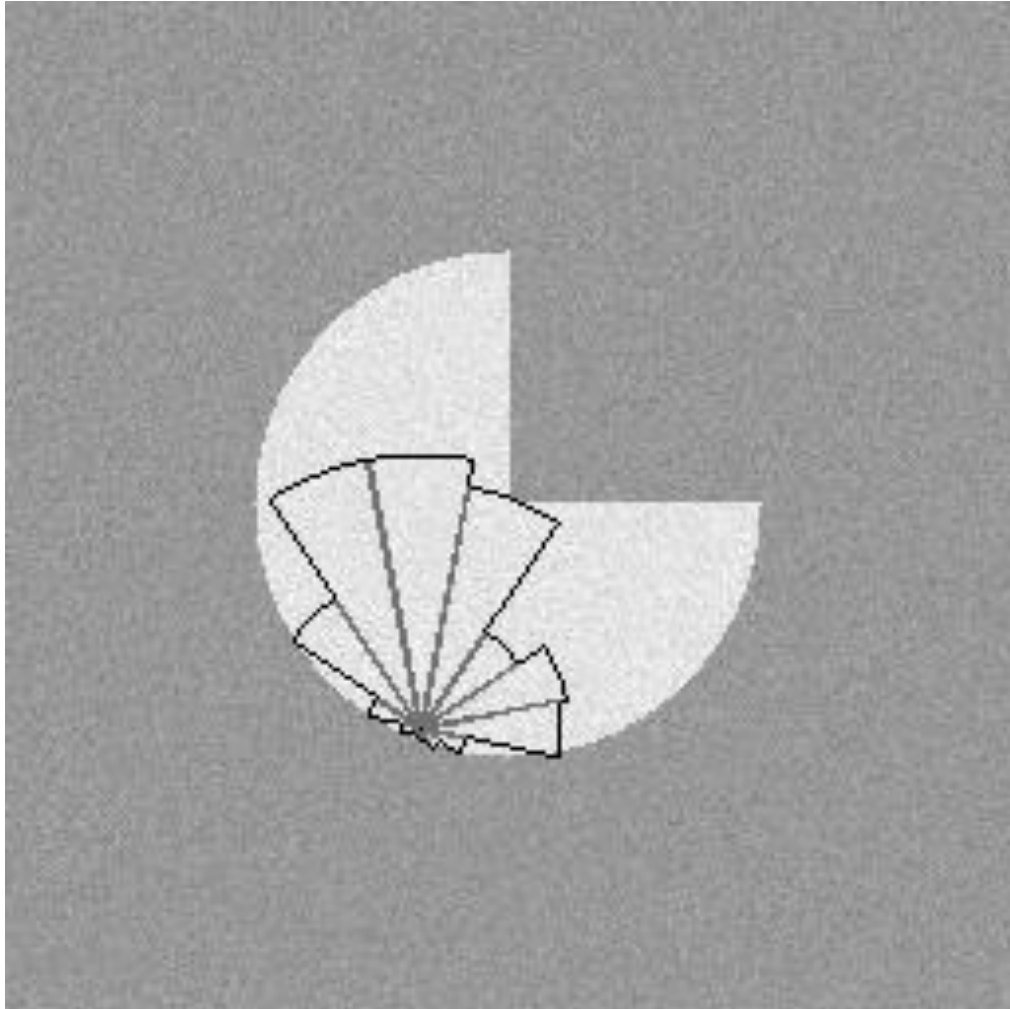
# LPA – ICI Rule

*ICI* is applied independently for each direction, yielding the adaptive scale (length) of the corresponding sector.  $\forall \theta$



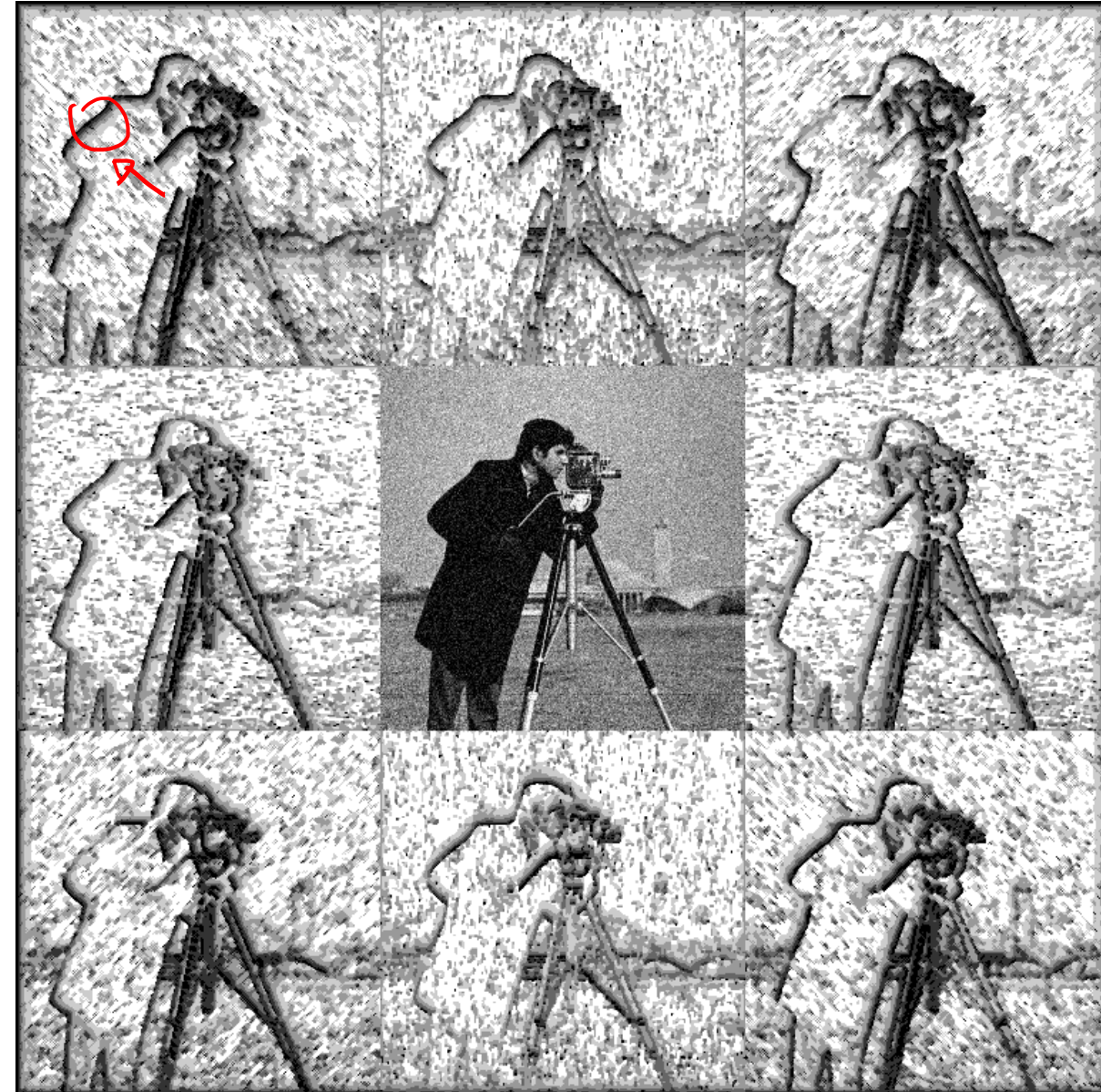
# Examples of adaptively selected neighborhoods

Neighborhoods adaptively selected using the LPA-ICI rule



# Adaptive Scales

The adaptive scales reveal the distribution of features (such as edges) across the corresponding direction.



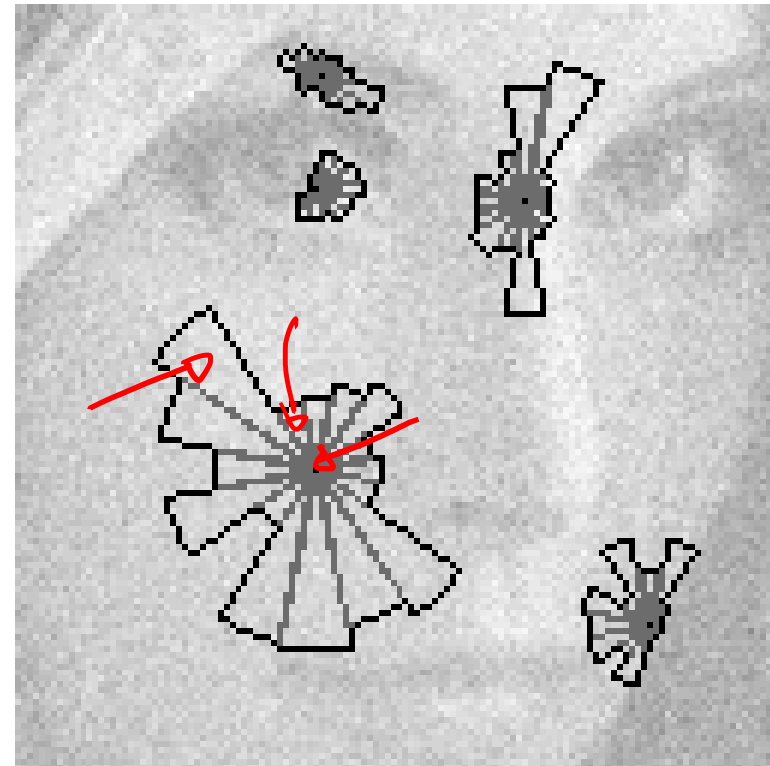
# Final Estimate

Directional adaptive-scale estimates are then *fused* together into the anisotropic estimate

$$\hat{y}(x_0) = \sum_{\theta \in \Theta} \lambda(x, \theta) \hat{y}_{h^+, \theta}(x)$$

And aggregation weights  $\lambda(x, \theta)$  are inversely proportional to the variance of the estimators

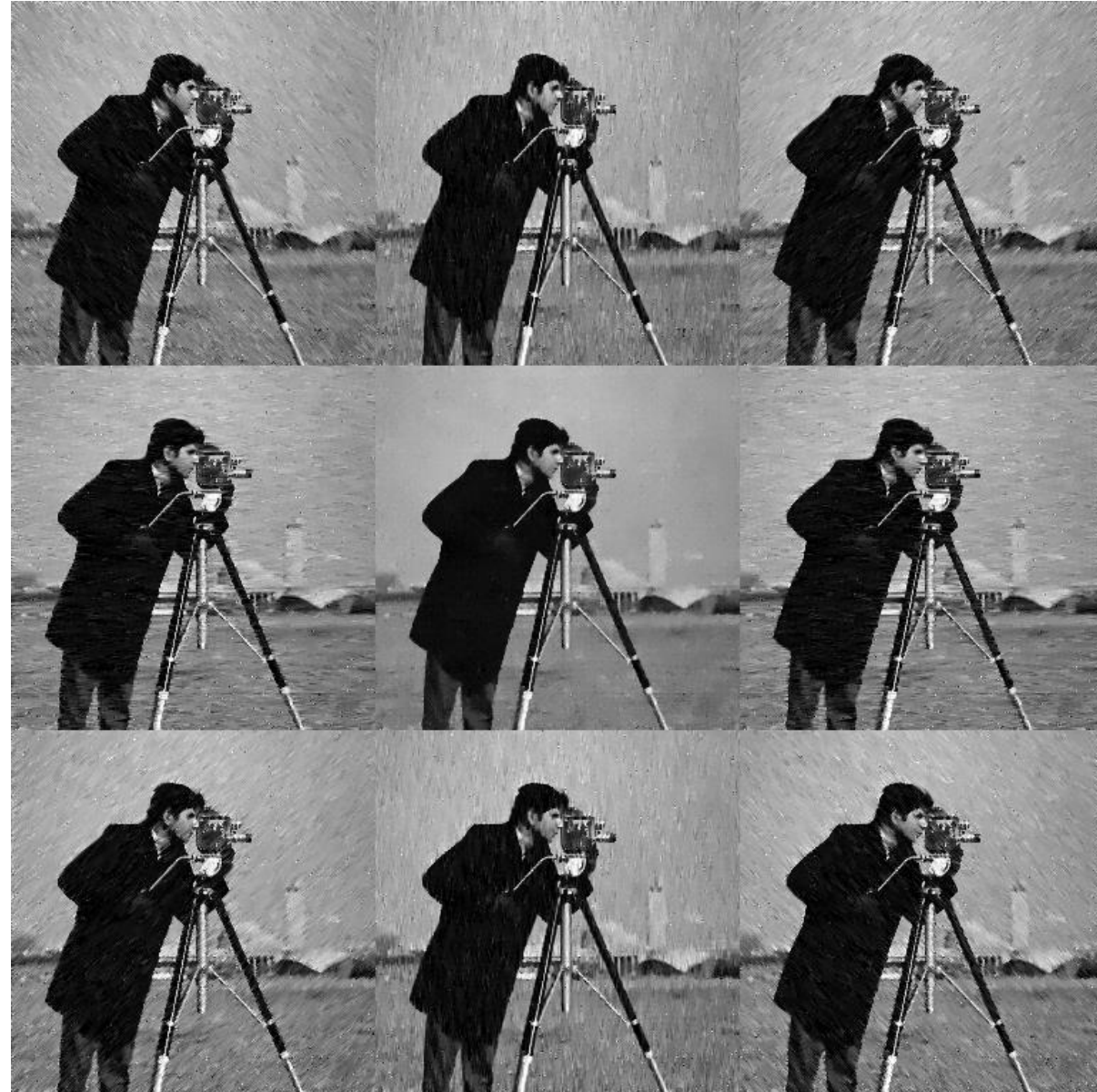
$$\lambda(x, \theta) = \frac{\sigma^{-2}(h^+, \theta, x)}{\sum_{\theta} \sigma^{-2}(h^+, \theta, x)}$$





# Final Estimate

The directional adaptive-scale estimates are fused together to obtain the final *anisotropic* estimate (shown in the center).



# Example of Performance

Original, noisy, denoised using polynomial regression on adaptively defined neighborhoods (LPA-ICI)

original



noisy



# Parametric Approaches

# Denoising Approaches

## Parametric Approaches

- Assume the noisy-free signal  $y$  features some **sparsity property** in a suitable domain (e.g Fourier, DCT, Wavelet) or w.r.t. some dictionary based decomposition. This implies that your image admits a **global parametric representation**

$$y = \sum_{i=1}^L x_i \underline{e}_i$$



# Sparsity as a Denoising Prior

# Sparsity and Parsimony

The principle of sparsity or “parsimony” consists in *representing some phenomenon with as few variables as possible*

Stretch back to philosopher William Ockham in 14<sup>th</sup> Century

Wrinch and Jeffreys [1921] relate simplicity to parsimony:

*The existence of simple laws is, then, apparently, to be regarded as a quality of nature; and accordingly we may infer that it is justifiable to prefer a simple law to a more complex one that fits our observations slightly better.*

Simplicity is the number of learning parameters

# Sparsity in Statistics

**Statistics:** simple models are preferred.

**Sparsity** is used to **prevent overfitting** and **improve interpretability** of learned models. *LASSO*

In model fitting, the number of parameters is typically used as a criterion to perform model selection.

See Bayes Information Criterion (BIC), Akaike Information Criterion (AIC), ..., Lasso.

# Sparsity in Signal Processing

**Signal Processing:** similar concepts but different terminology.

**Vectors** corresponds to signals and **data modeling** is crucial for solving inverse problems.

A very successful **prior** is to assume that **images are approximated by sparse linear combinations of prototypes** (basis elements / atoms of a dictionary), resulting in simpler and compact model.

**Noise**, having no structure, is **not expected to be sparse**.

$$z \mapsto z \in \mathbb{R}^d \quad \text{vector space}$$

# How to promote sparsity?

Sparisity does not make much sense in space domain

256

→ MOVE TO  
A DIFFERENT  
DOMAIN

256



cameraman



Sparse cameraman (only 50%  
of pixels are preserved)

# Basic Notions for Transformed-Domain Image Processing

# Preliminary Notation: Linear Independence

A set of vectors  $\{\mathbf{e}_i\}_{i=1,\dots,N}$  living in a vector space  $V$  is said to be linearly independent if

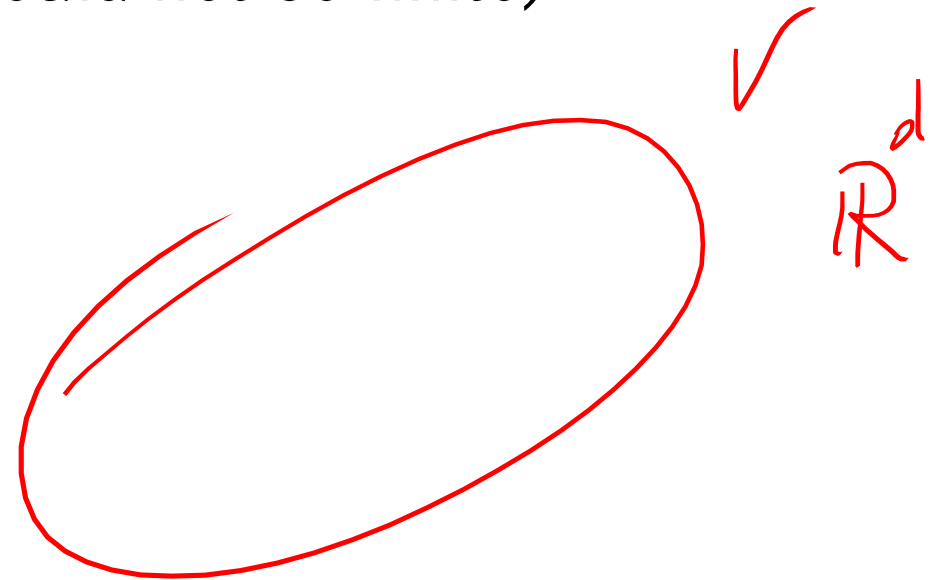
$$\sum_i x_i \mathbf{e}_i = \mathbf{0} \Leftrightarrow x_i = 0 \forall i$$

# Preliminary Notation: Basis

A basis  $\{\mathbf{e}_i\}_{i=1,\dots,d}$  of a vector space  $V$  is a set of linearly independent vectors that spans the whole vector space  $V$ , namely:

$$\forall \mathbf{v} \in V, \exists \{x_i\}_{i=1,\dots,d} \text{ s. t. } \mathbf{v} = \sum_i x_i \mathbf{e}_i$$

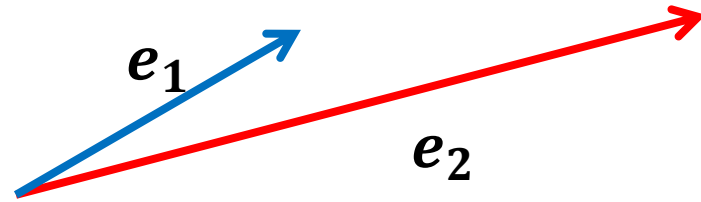
$d$  is said the dimension of  $V$  (potentially it could not be finite)





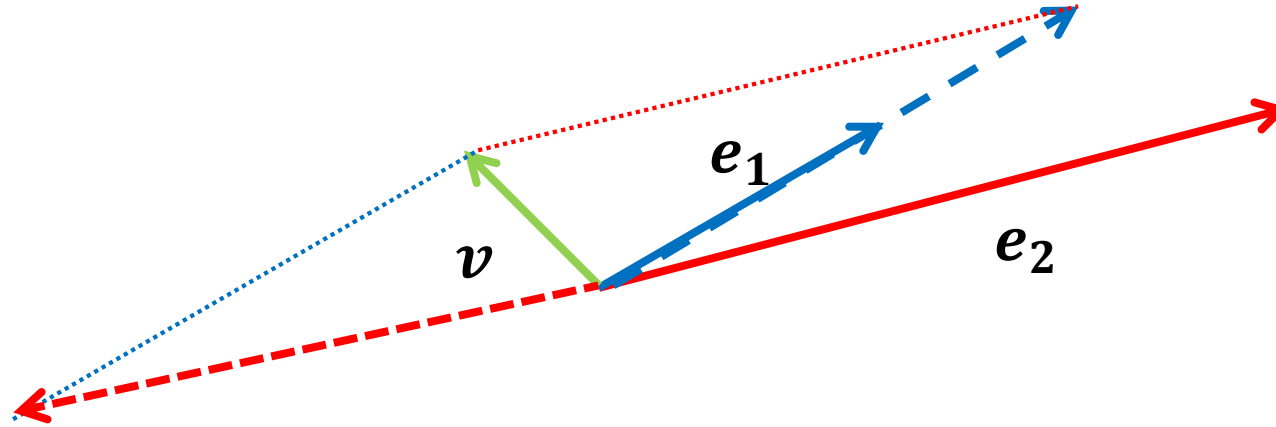
# Example

Th two vectors  $\{e_1, e_2\}$  represent a basis for the plane



# Example

Th two vectors  $\{e_1, e_2\}$  represent a basis for the plane

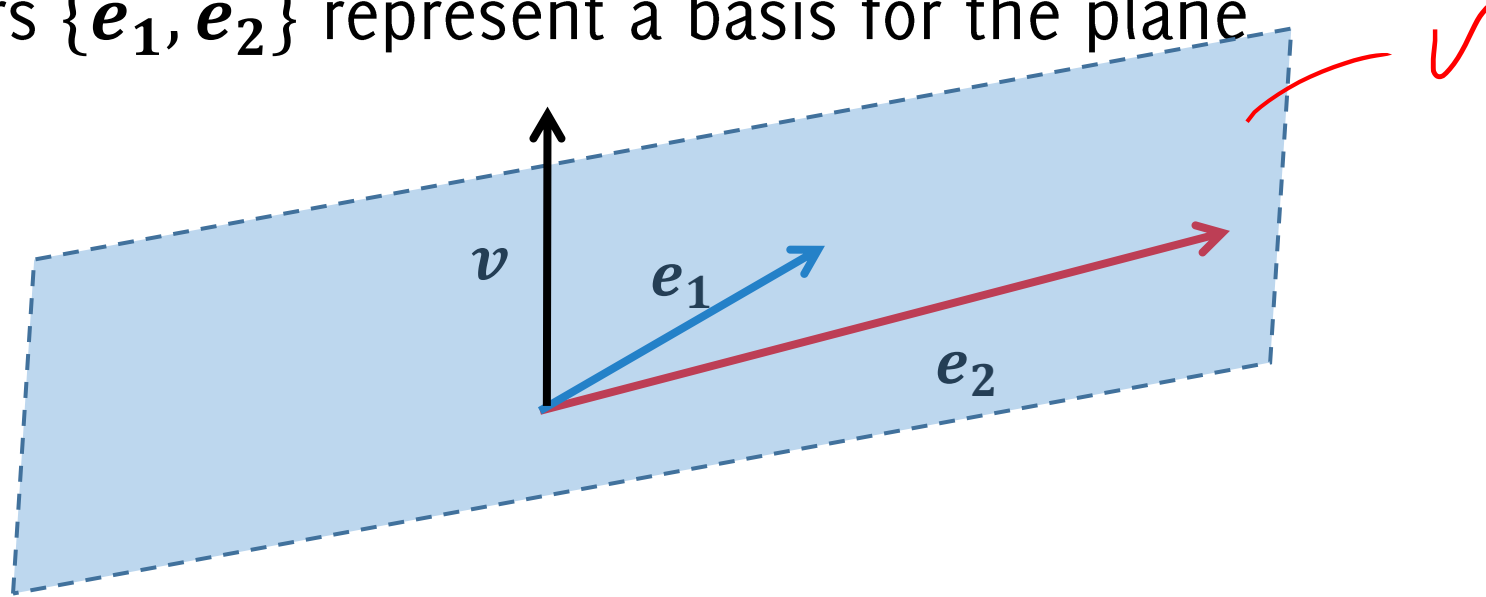


Since any vector in this plane can be written as a linear combination of these two

$$v = 1.5 * e_1 + (-1) * e_2$$

# Example

Th two vectors  $\{e_1, e_2\}$  represent a basis for the plane



This holds for any vector except the orthogonal ones

# A very important property

Let  $\{\mathbf{e}_i\}_{i=1,\dots,d}$  be a set of linearly independent generators of  $V$  then:

$$\forall \underline{\mathbf{v}} \in V, \exists \{x_i\}_{i=1,\dots,d} \text{ s. t. } \mathbf{v} = \sum_i x_i \mathbf{e}_i$$

And such representation is **unique**

$$\mathbf{v} \rightarrow \{x_i\}_{i=1,\dots,d}$$

*signs  $\mapsto$  coefficients*

**Proof:** follows trivially from the definition of linear independence

# A very important property

Let  $\{\mathbf{e}_i\}_{i=1,\dots,d}$  be a set of linearly independent generators of  $V$  then:

$$\forall \mathbf{v} \in V, \exists \{x_i\}_{i=1,\dots,d} \text{ s.t. } \mathbf{v} = \sum_i x_i \mathbf{e}_i$$

And such representation is **unique**

$$\mathbf{v} \rightarrow \{x_i\}_{i=1,\dots,d}$$

**Proof:** follow

Ok, we know that this representation exists and is unique

$$\forall \mathbf{v} \in V, \exists \{x_i\}_{i=1,\dots,d} \text{ s.t. } \mathbf{v} = \sum_i x_i \mathbf{e}_i$$

.. But given an input  $\mathbf{v}$ , how to compute  $\mathbf{v} \rightarrow \{x_i\}_{i=1,\dots,d}$

# Orthonormal Basis in Euclidean Space

An orthonormal basis is basis  $\{\mathbf{e}_i\}_{i=1,\dots,d}$  is such that

$$\mathbf{e}_i^\top \mathbf{e}_j = \delta_{i,j}$$

The advantage of using orthonormal basis is that

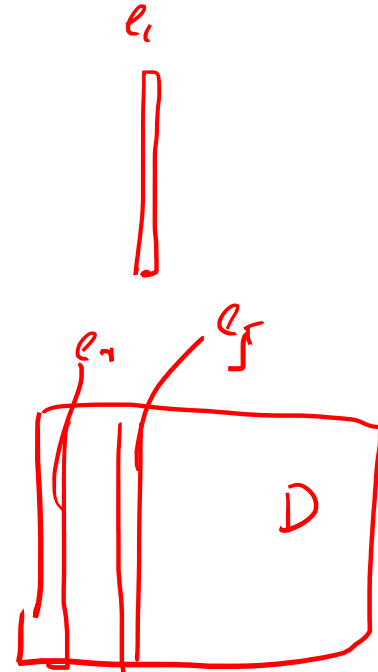
$$\forall \mathbf{v} \in V, \exists \{x_i\}_{i=1,\dots,d} \text{ s.t. } \mathbf{v} = \sum_i x_i \mathbf{e}_i$$

and

$$x_i = \mathbf{e}_i^\top \mathbf{v} = \mathbf{v}^\top \mathbf{e}_i$$

If we arrange the basis  $\{\mathbf{e}_i\}_{i=1,\dots,d}$  in the columns of a  $d \times d$  matrix  $D$

$$\mathbf{x} = \mathbf{D}^\top \mathbf{v}$$



# Ok, let's go back to images

Which is the basis we used for representing digital images so far?

$$v = \sum x_i e_i$$

*(Handwritten red text with arrows pointing to the summation and a question mark below it)*

$v =$

116	23	33
16	3	73
5	4	30

# Ok, let's go back to images

Which is the basis we used for representing digital images so far?

$$\mathbf{v} =$$

116	23	33
16	3	73
5	4	30

$$\mathbf{e}_1 =$$

1	0	0
0	0	0
0	0	0

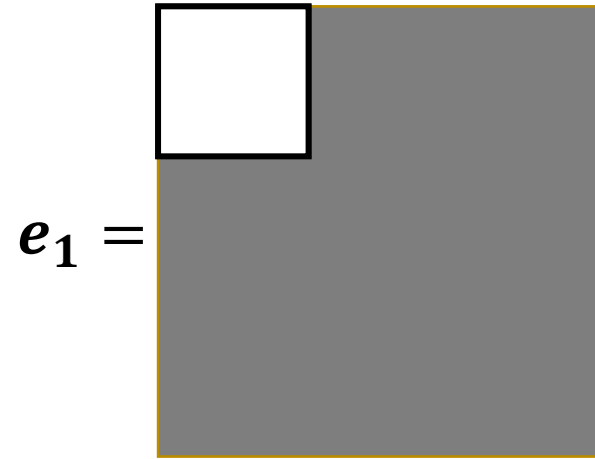


# The canonical basis of $\mathbb{R}^d$

Which is the basis we used for representing digital images so far?

$v =$

116	23	33
16	3	73
5	4	30

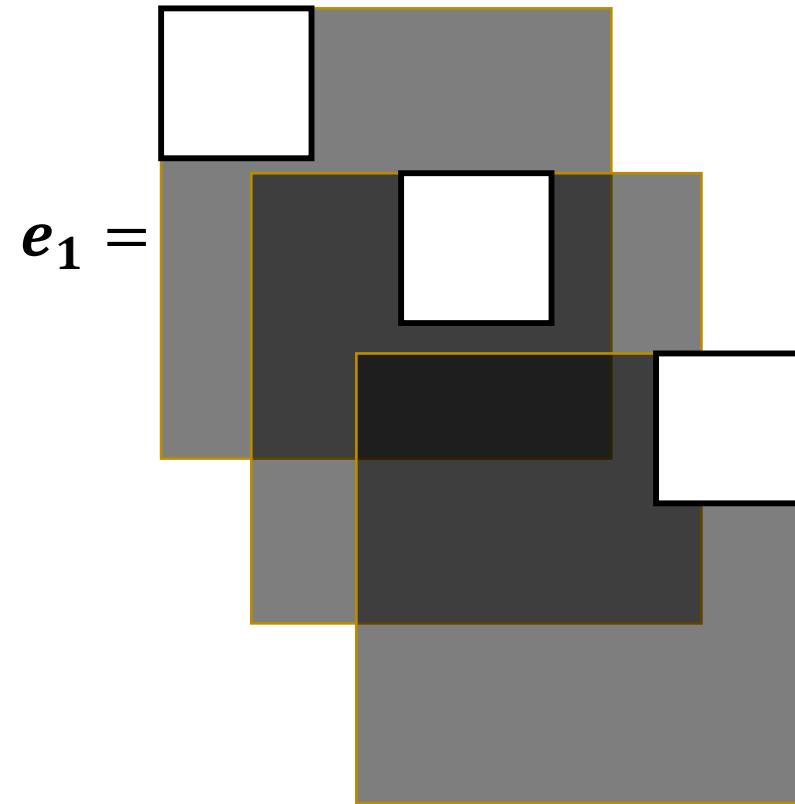


# The canonical basis of $\mathbb{R}^d$

Which is the basis we used for representing digital images so far?

$v =$

116	23	33
16	3	73
5	4	30

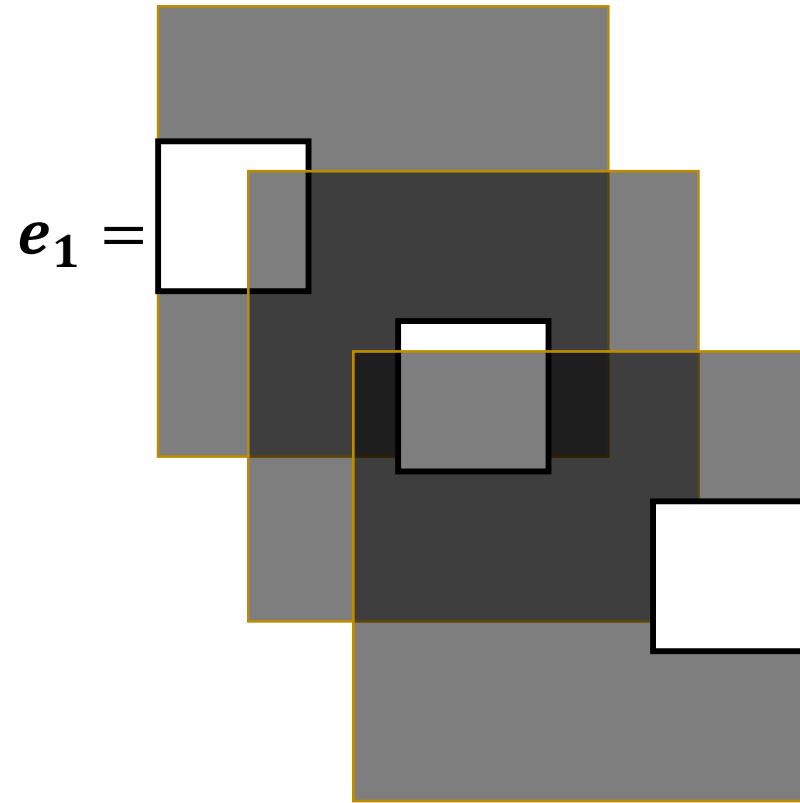


# The canonical basis of $\mathbb{R}^d$

Which is the basis we used for representing digital images so far?

$v =$

116	23	33
16	3	73
5	4	30



# The canonical basis of $\mathbb{R}^d$

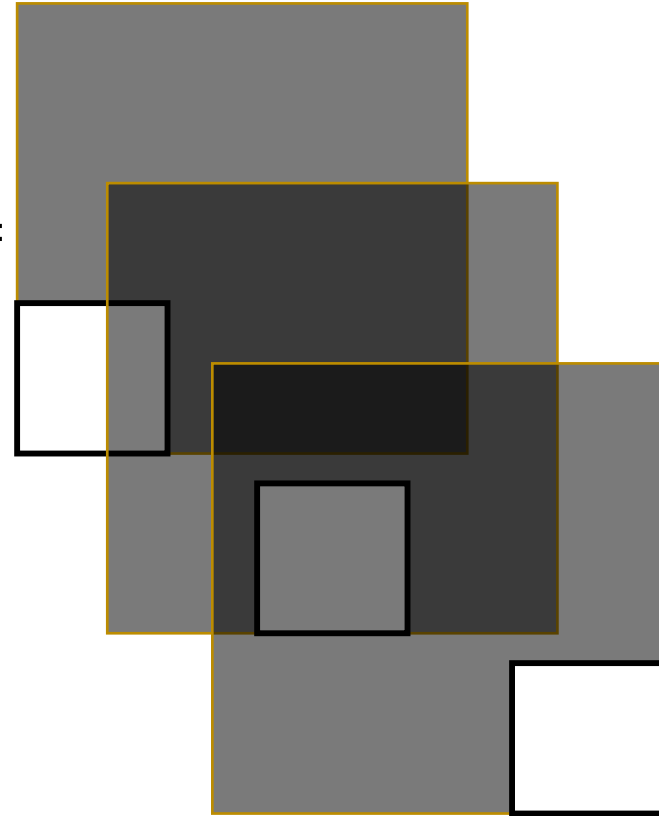
Which is the basis we used for representing digital images so far?

$$116 \cdot e_1 + 23 \cdot e_2 + 33 \cdot e_3$$

$v =$

116	23	33
16	3	73
5	4	30

$e_1 =$



# Thus the canonical basis

Canonical basis

$$\{\mathbf{e}_i\}_{i=1,\dots,d}$$

Being  $\mathbf{e}_j = \text{zeros}(1, d)$ ;  $\mathbf{e}_j(j) = 1$

Uses each coefficient to represent a pixel:

- all coefficients are equally meaningful
- thus, it is not usefull at all for compression

Are there basis that ease image processing tasks?

$$\mathbf{e}_1 = \begin{array}{|c|c|c|} \hline 1 & 0 & 0 \\ \hline 0 & 0 & 0 \\ \hline 0 & 0 & 0 \\ \hline \end{array}$$

# Fourier Transform

# 2D Fourier Transform

The  $(u,v)$ -element of the 2D Fourier basis is defined as

$$e^{-i2\pi(ux+vy)} = \cos(2\pi(ux + vy)) + i \sin(2\pi(ux + vy))$$

Each Fourier coefficient is computed with an inner product with the corresponding function.

The Fourier basis functions are constant where  $y = -\frac{ux}{c} + c$

The Fourier basis functions have unlimited support.

The Fourier Transform is invertible (it is an orthonormal transform)

Fourier domain is also called frequency domain.

# 2D Discrete Fourier Transform

The Discrete Fourier Transform (DFT) is defined as

$$F[m, n] = \sum_{k=0}^{M-1} \sum_{l=0}^{N-1} f[k, l] e^{-\pi i \left( \frac{km}{M} + \frac{ln}{N} \right)}$$

The Fourier Transform admit a fast implementation (FFT) when the signal/image sizes are powers of 2.

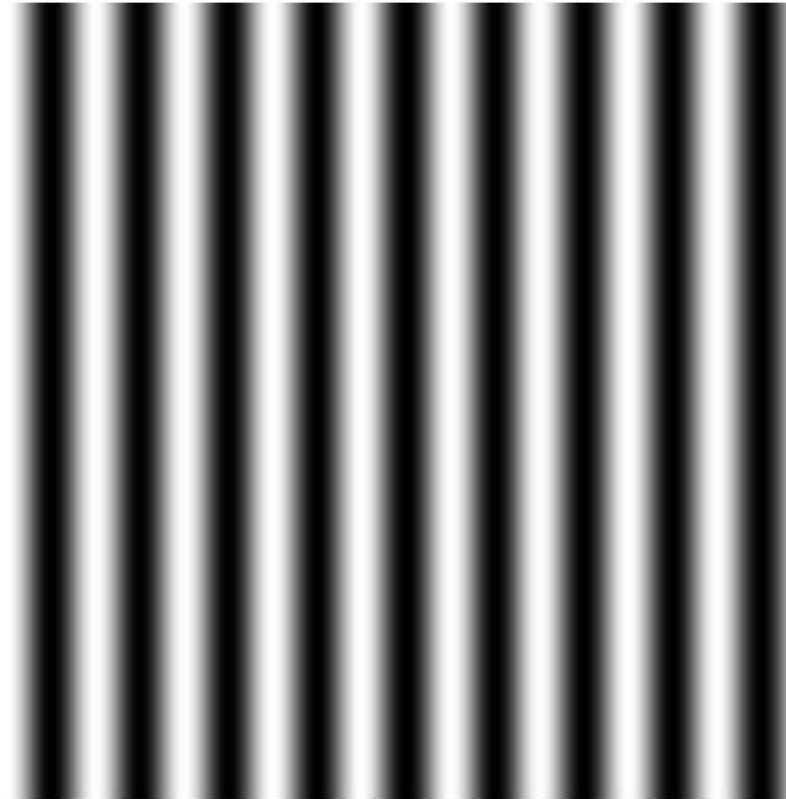


# 2D Fourier Basis Elements

## Frequency and Orientation of the 2D Fourier Basis Elements

$$b_{u,v} \quad u = 1, v = 10$$

the  $u=1, v=10$  Fourier domain basis element

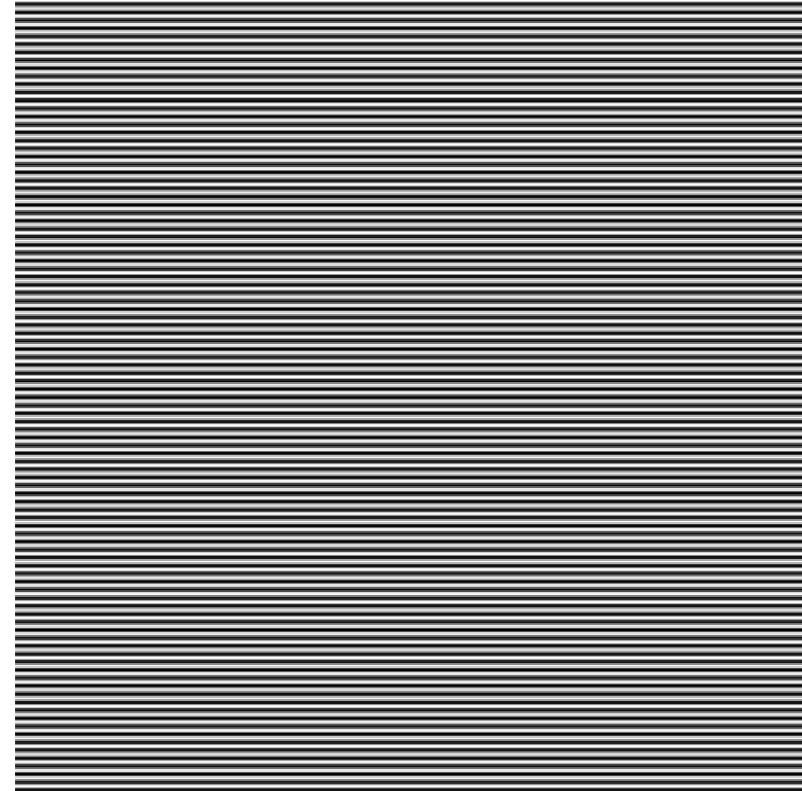


# 2D Fourier Basis Elements

## Frequency and Orientation of the 2D Fourier Basis Elements

$$b_{u,v} \quad u = 100, v = 1$$

the  $u=100, v=1$  Fourier domain basis element

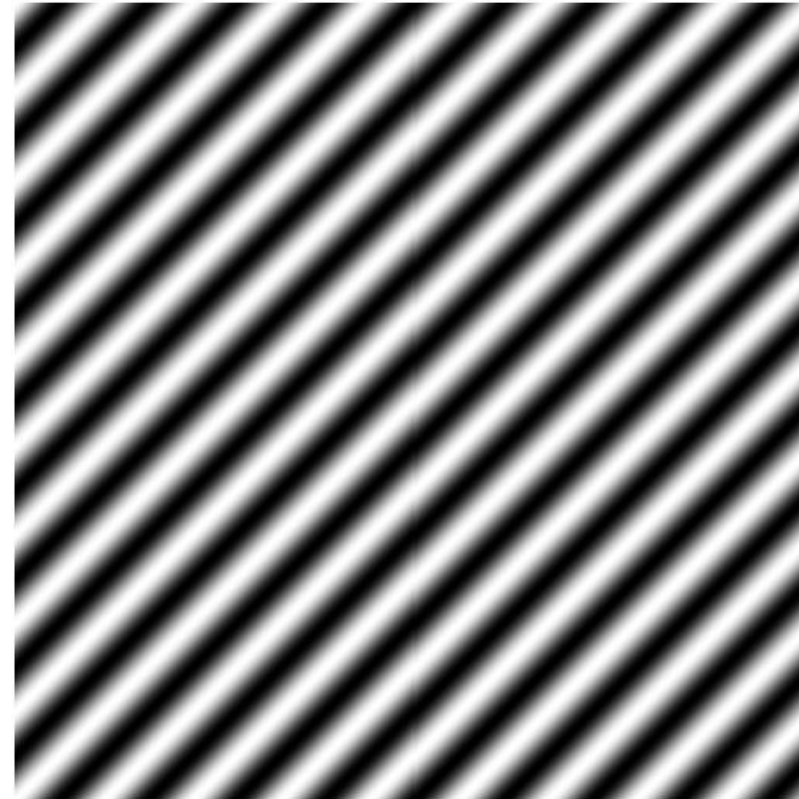


# 2D Fourier Basis Elements

## Frequency and Orientation of the 2D Fourier Basis Elements

$$b_{u,v} \quad u = 10, v = 10$$

the  $u=10, v=10$  Fourier domain basis element



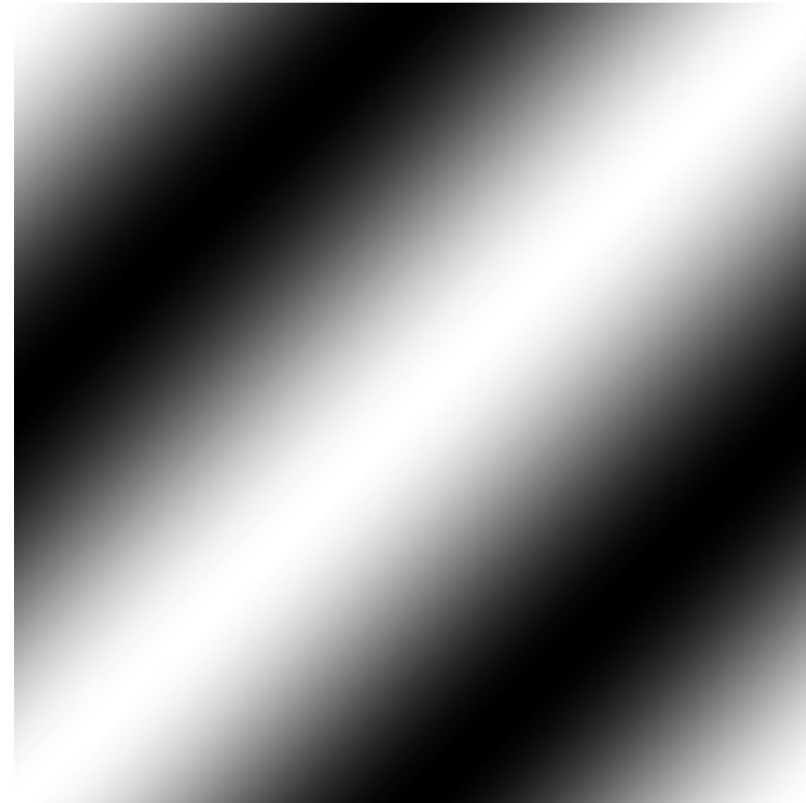
# 2D Fourier Basis Elements

## Frequency and Orientation of the 2D Fourier Basis Elements

$$b_{u,v} \quad u = 2, v = 2$$

note that in matlab, Fourier coefficients are indexed starting from 1

the  $u=2, v=2$  Fourier domain basis element



# Try it in Matlab

```
%Examples of Basis Elements
```

```
Y=zeros(512);
```

```
u=10
```

```
v=10
```

```
Y(u,v)=1;
```

```
% Y is the image in Fourier Domain having only the  $(u,v)$ -coefficient set to 1 and the other set to 0
```

```
figure(1),imshow(Y,[],),title(['the u=',num2str(u), 'v=',num2str(v), 'space domain basis element']  
)
```

```
% by inverting the Fourier transform one get the  $(u,v)$ -element of the 2D discrete Fourier basis
```

```
y=ifft2(Y);
```

```
figure(2),imshow(real(y),[],),title(['the u=',num2str(u), 'v=',num2str(v), 'Fourier domain basis  
element'] )
```

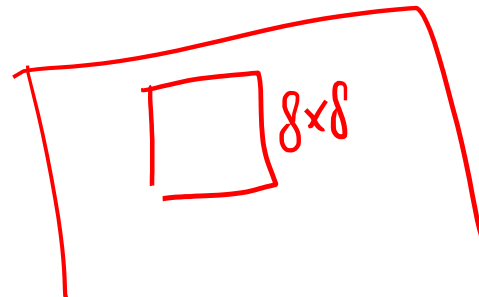
# Limitations of Fourier Transform

- Fourier Transform yield complex coefficients,
- Fourier Transformation it's a global transformation (each and every pixel affects in principle each and every coefficient)

DCT

Difficult to find a parametric model describing the whole image and enabling separation between noise and the image content

- Better operate with a real-valued basis
- Better operate patch-wise, such that the parametric model applies to small portions of the image



# Denoising by Thresholding in Transform Domain

# Denoising by Thresholding

The underlying assumption is that a clean signal admits a sparse representation w.r.t. a suitable basis  $\{\mathbf{e}_i\}_{i=1,\dots,d}$  (or set of generators)

Noise, being unstructured, spreads in all the coefficients

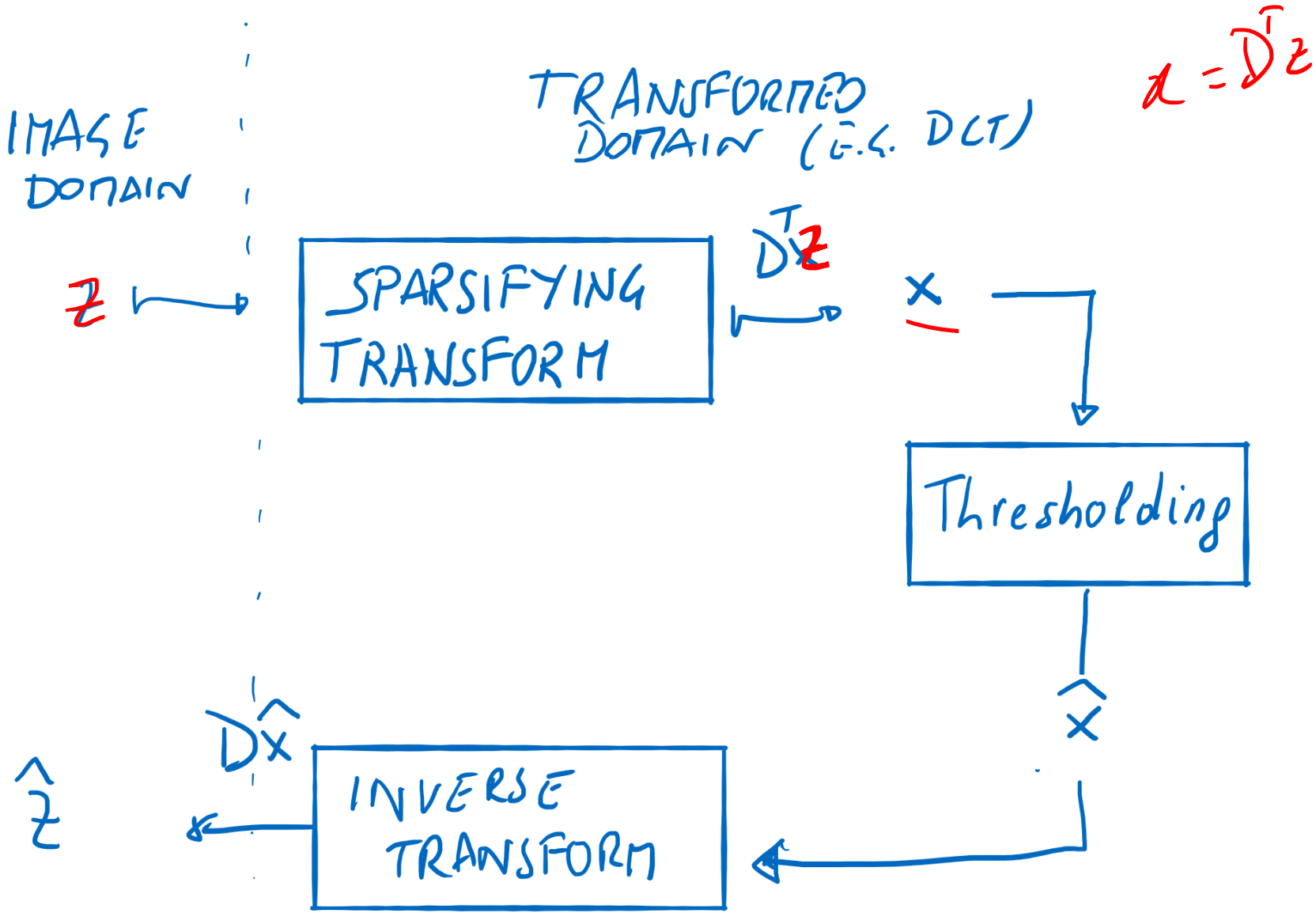
DCT

A general denoising approach consists in

- Cropping the image patch-wise
- Transform each patch according to a sparsifying transformation
- Perform Thresholding
- Invert Transformation



# Sparsity Promoting Denoising



$$x = D^T z$$

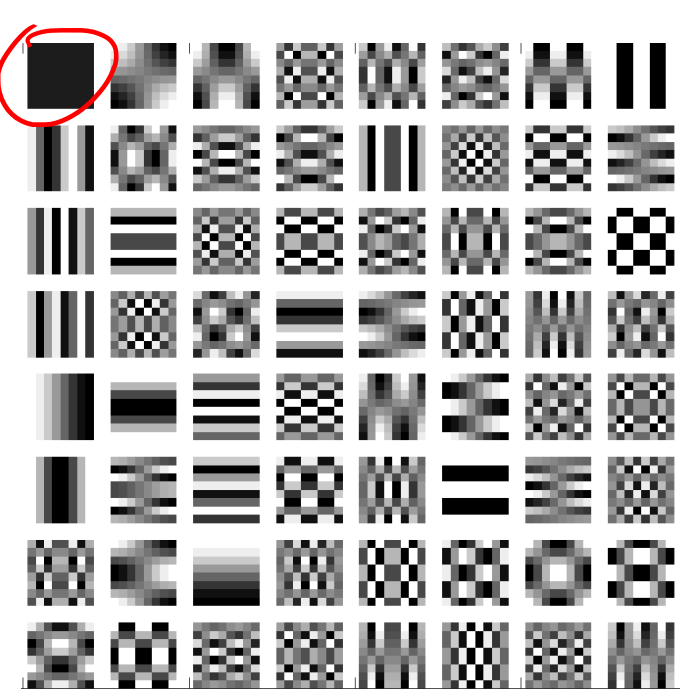
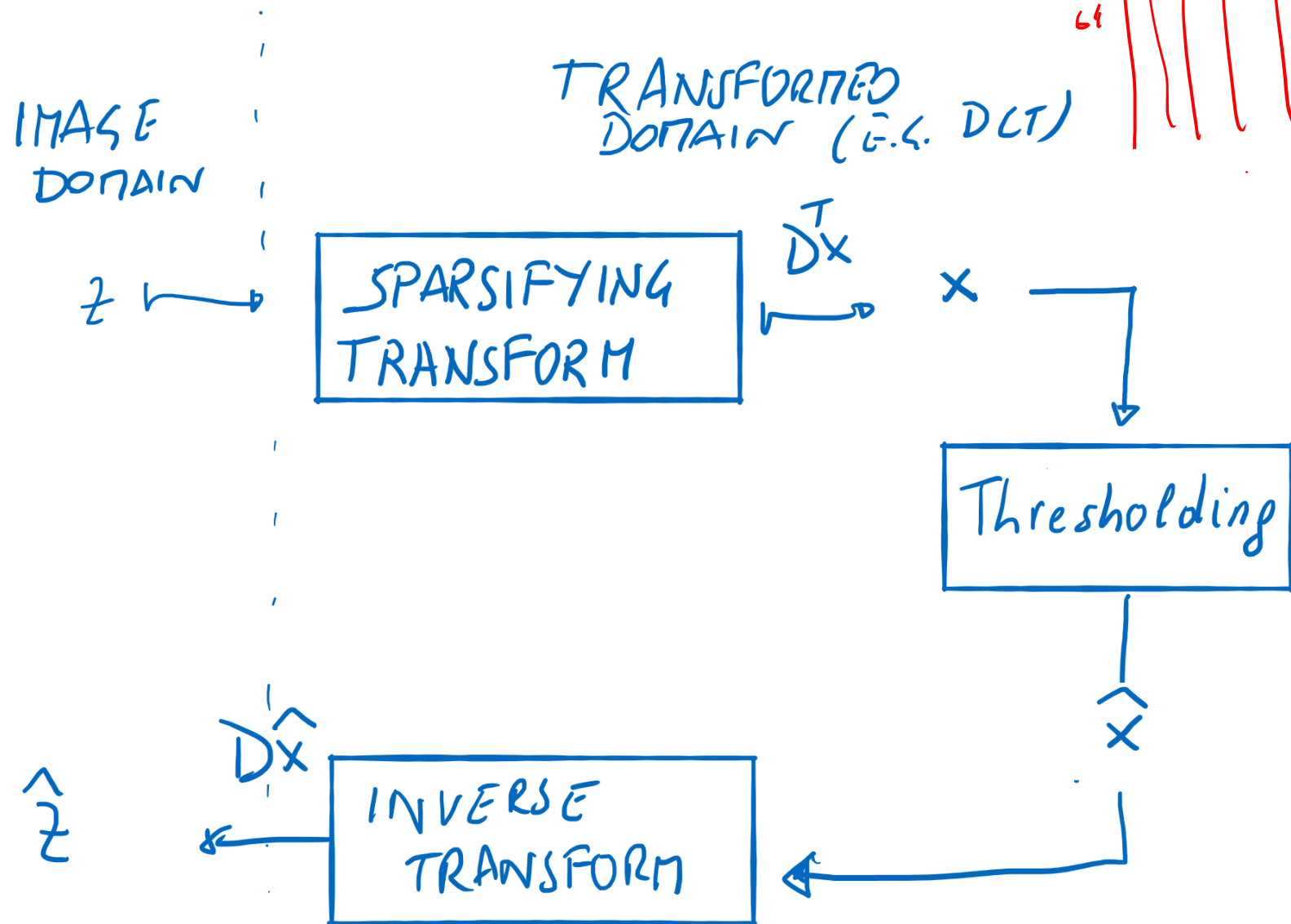
$$D^T(z) =$$

$$D^T(y+n) =$$

$$D^T y + D^T n$$

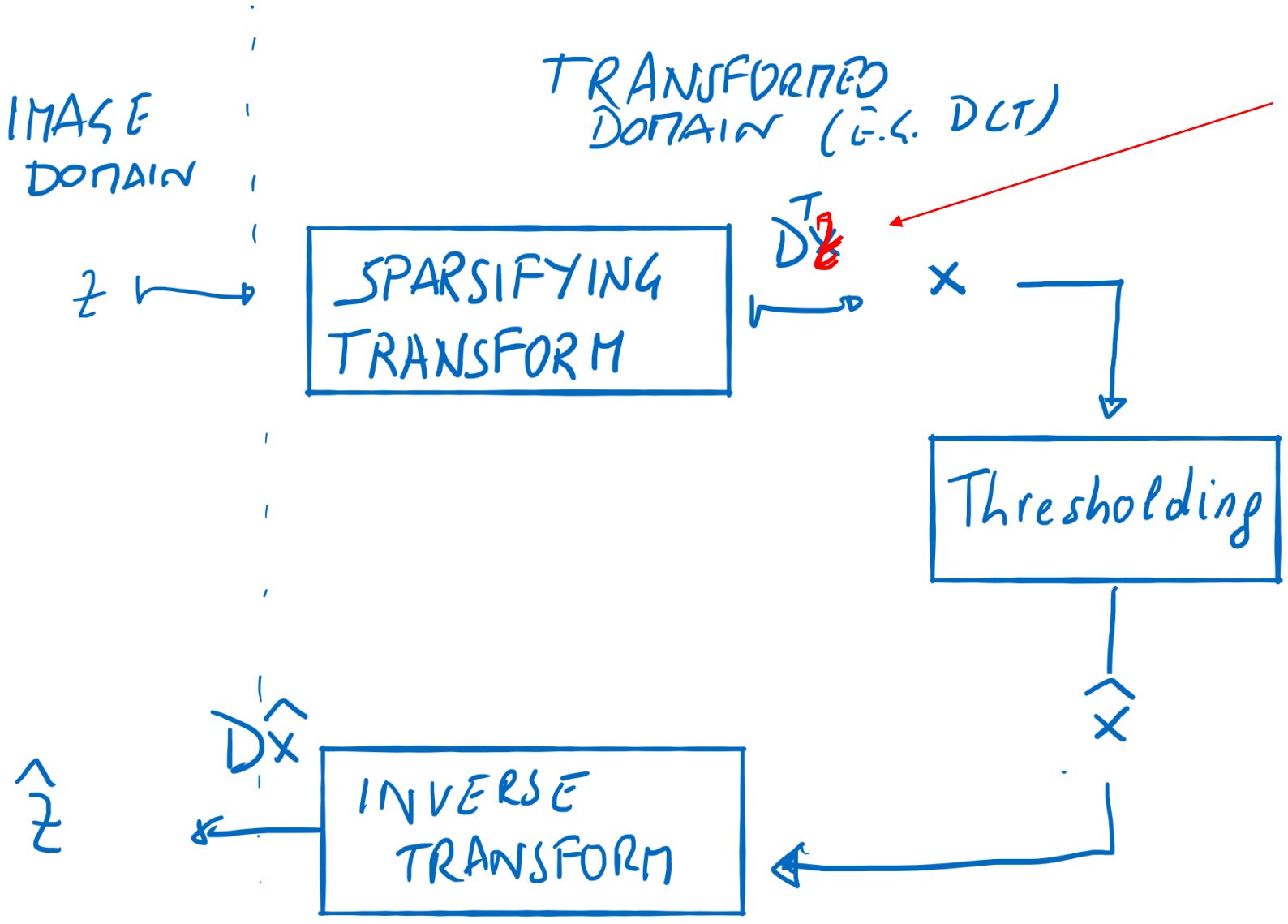
↑ sparse                      ↑ non sparse

# Sparsity Promoting Denoising



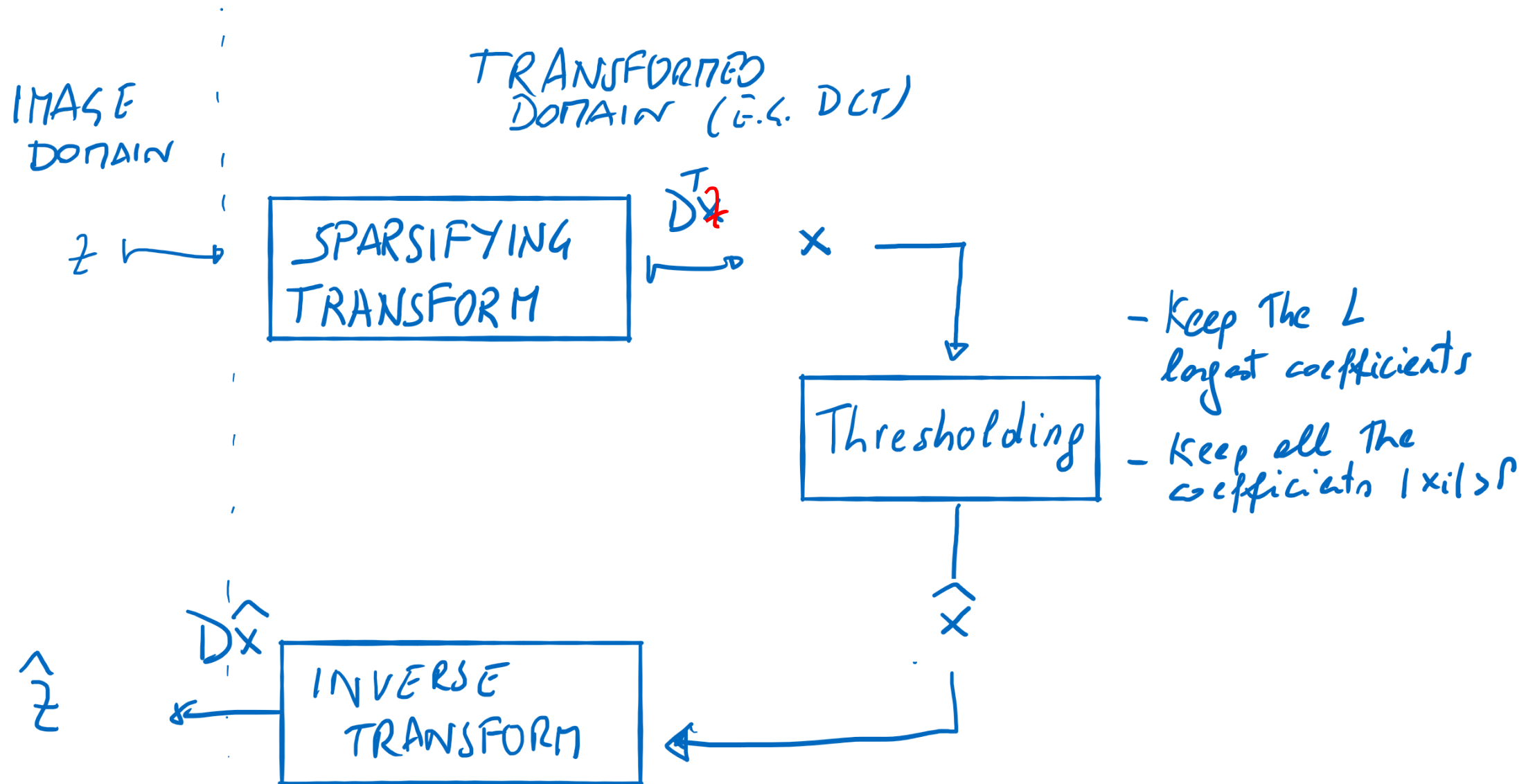
DCT: Discrete Cosine Transform is often used at patch-level. These are the elements of the basis, which are then arranged column-wise in the matrix  $D$

# Sparsity Promoting Denoising

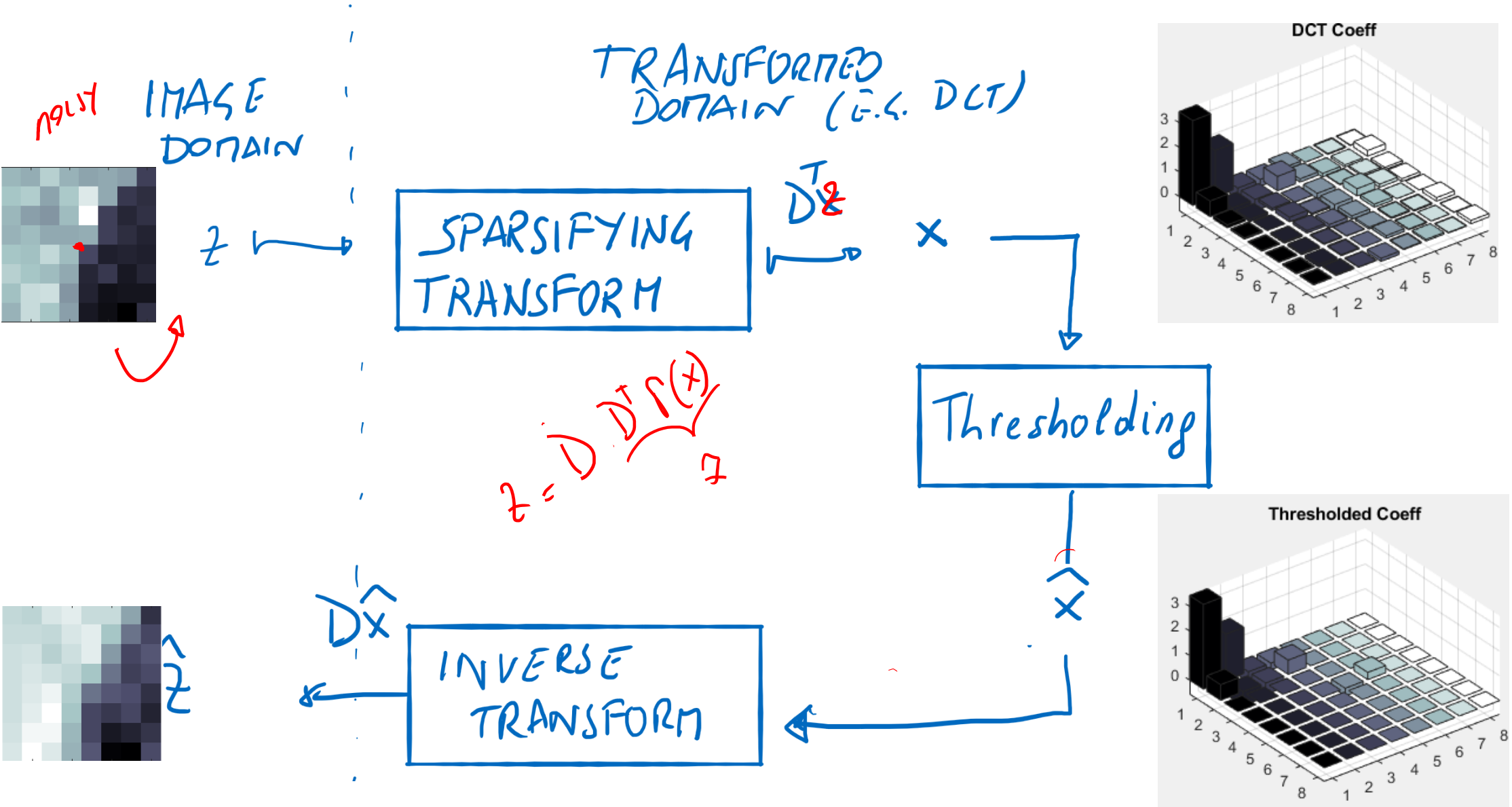


When the transformation is not w.r.t. an orthonormal basis, decomposition equation is not as simple. In case of redundant set of generators the representation is not unique and it has to «be pursued»

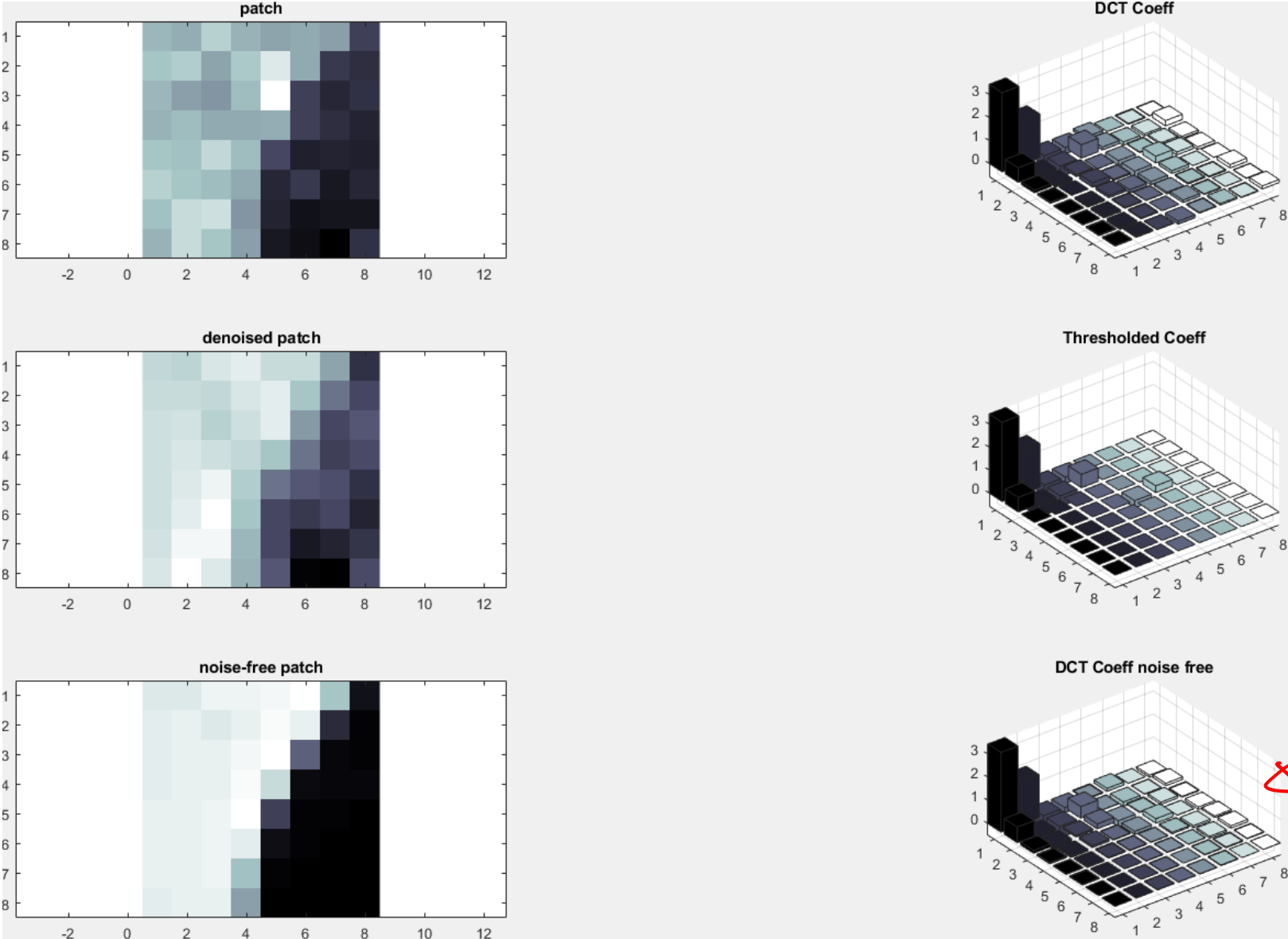
# Sparsity Promoting Denoising



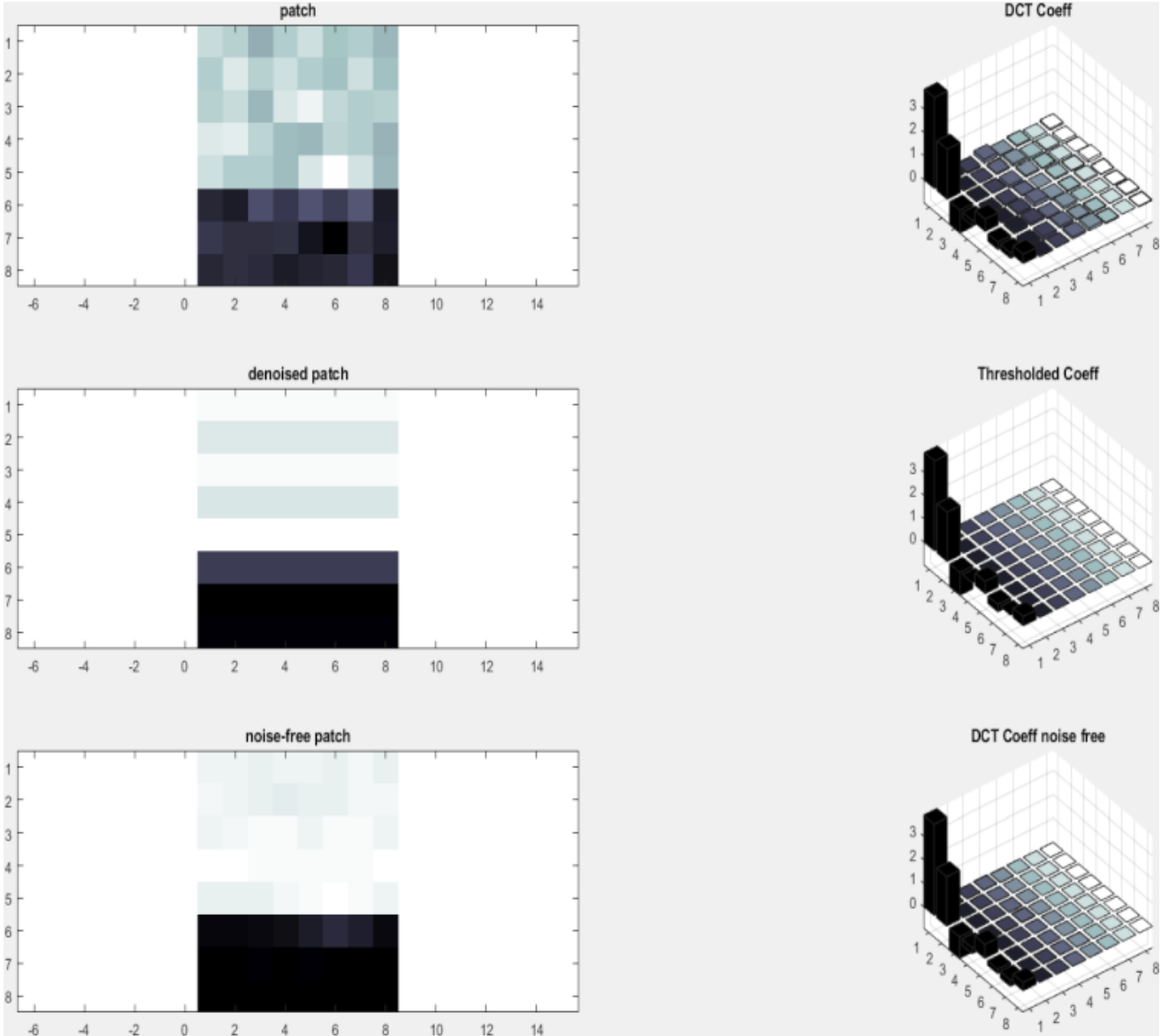
# Sparsity Promoting Denoising



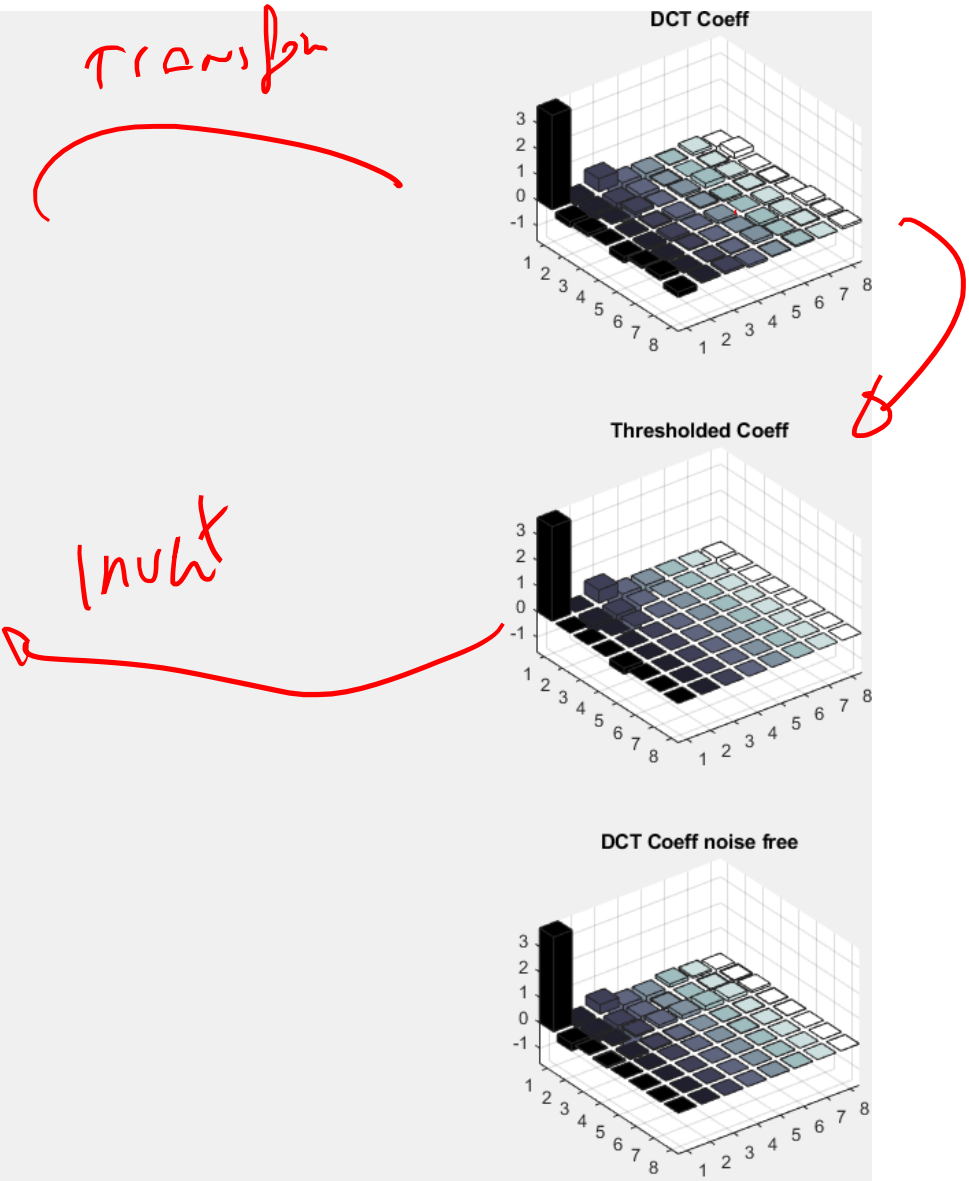
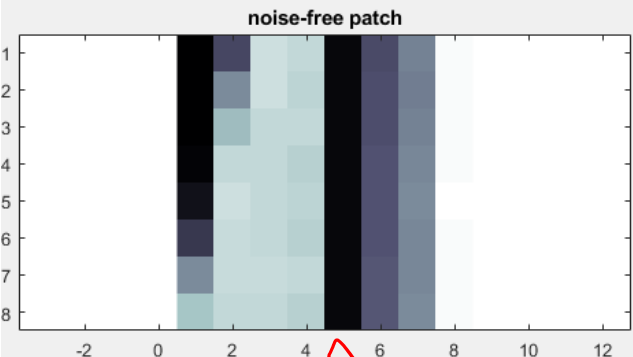
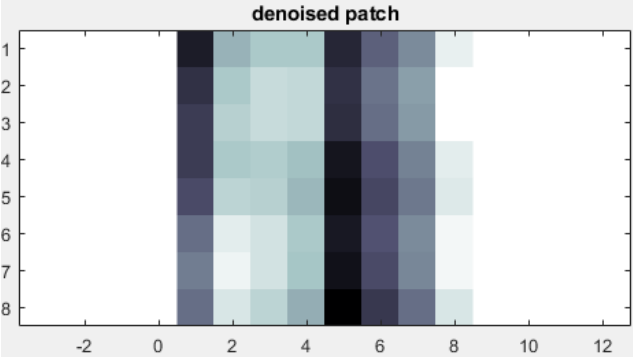
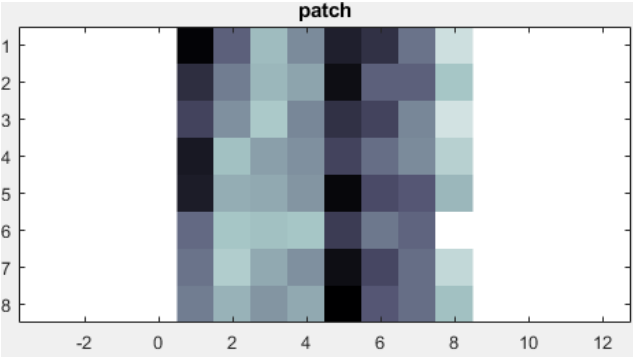
# Sparsity Promoting Denoising



# Sparsity Promoting Denoising

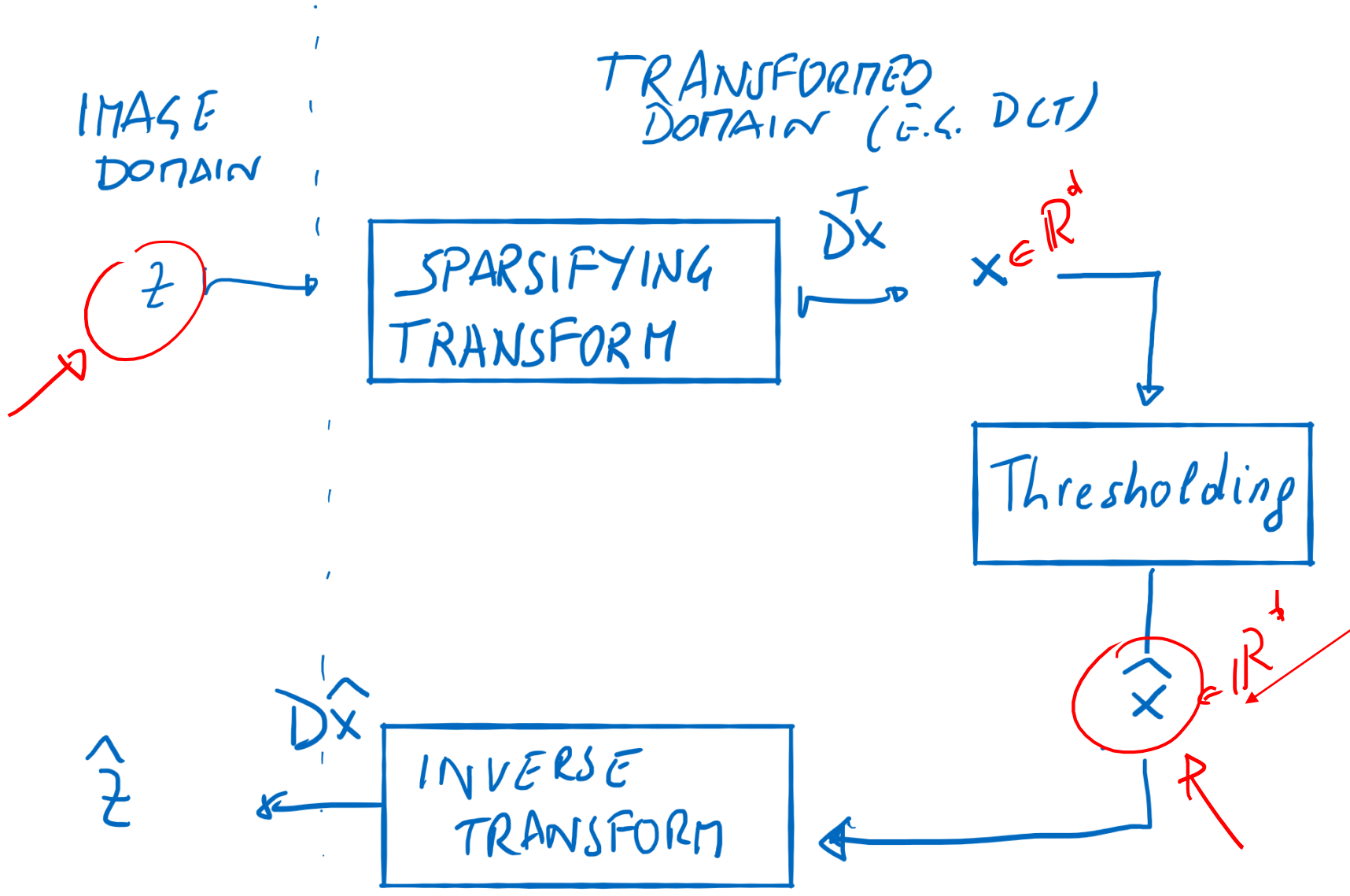


# Sparsity Promoting Denoising



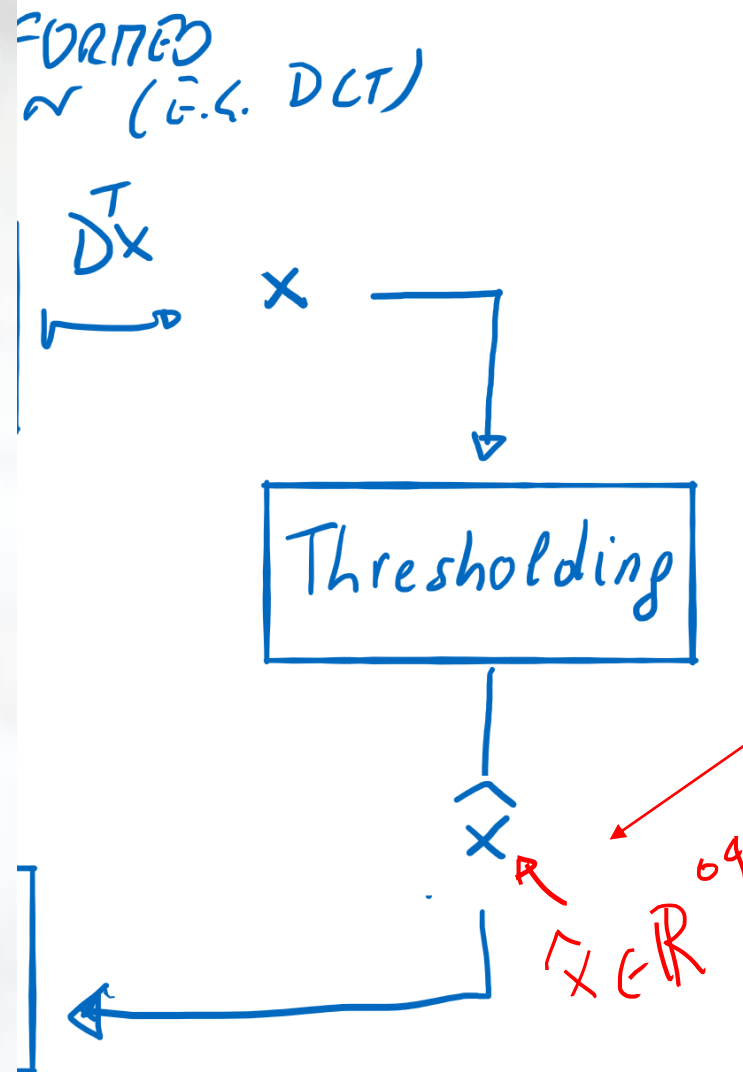


# Sparsity Promoting Denoising



This is a much more compressed representation than  $z$ .  
Encoding  $\hat{x}$  instead of  $z$  can significantly reduce the size of the image  
JPEG Compression performs encoding of  $\hat{x}$

# Sparsity Promoting Denoising



This is a much more compressed representation than  $z$ .  
Encoding  $\hat{x}$  instead of  $z$  can significantly reduce the size of the image  
JPEG Compression performs encoding of  $\hat{x}$

# Recent Denoising Algorithms

# Hand-Crafted Algorithms

Recent denoising algorithms combine principles from different approaches:

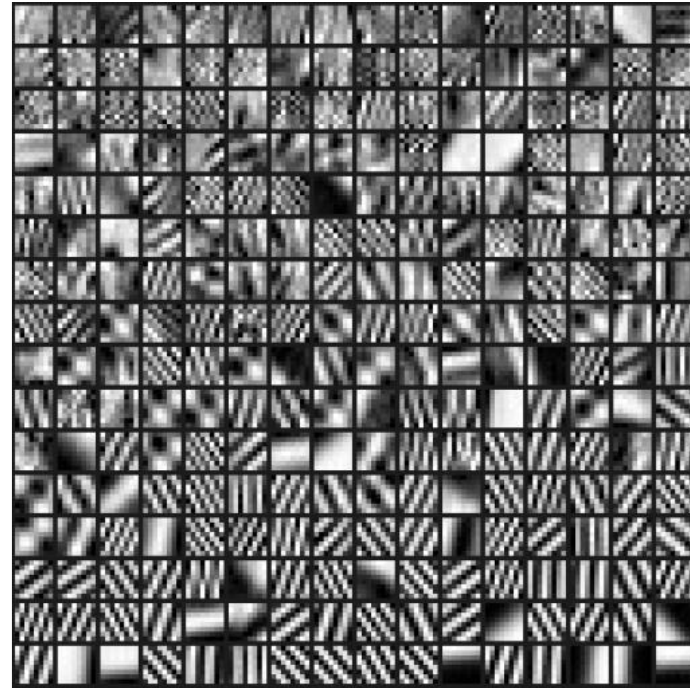
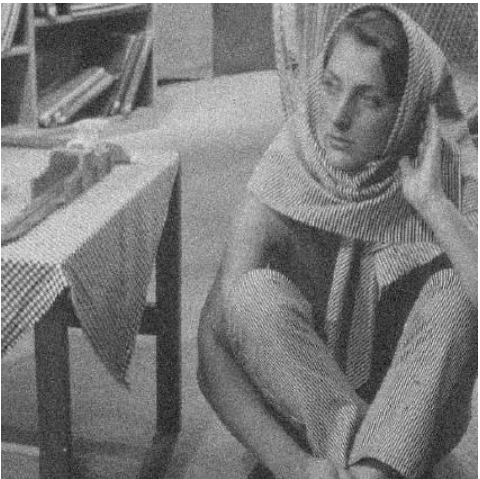
- Spatial Adaptivity
- Nonlocal self-similarity
- Sparsity w.r.t. to a basis / a learned set of generators
- ...

# K-SVD Image Denoising



original

Noisy image ( $\sigma = 20$ )



The obtained dictionary after 10 iterations

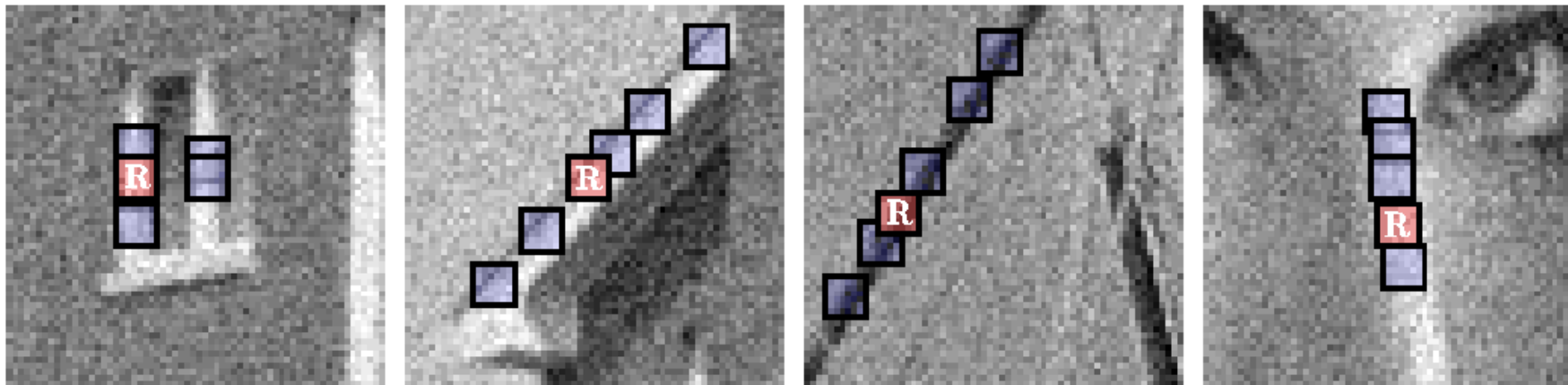
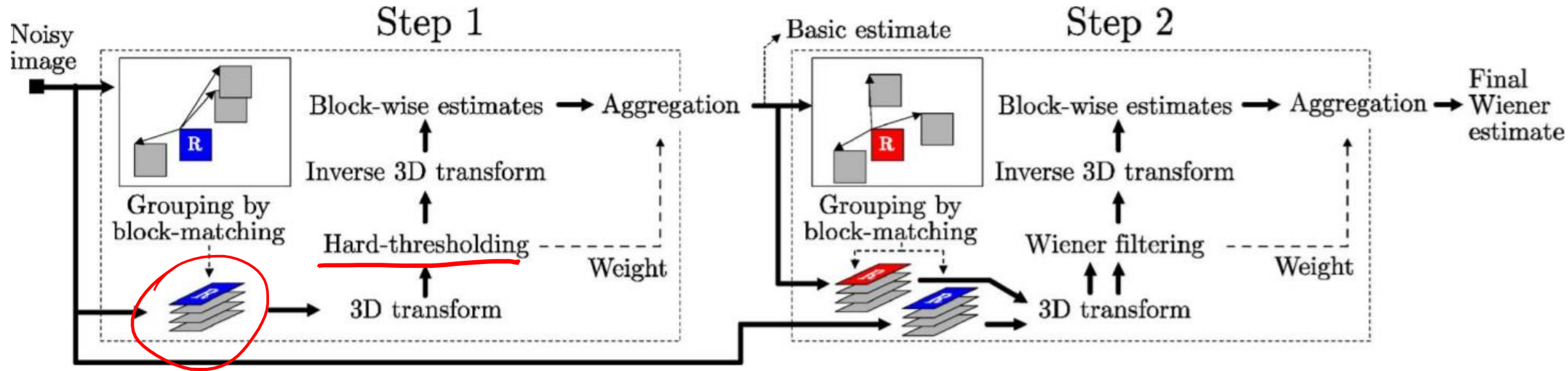


Result 30.829dB

Aharon, M., Elad, M., & Bruckstein, A. (2006). K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on signal processing*, 54(11), 4311-4322.

Sparse and Redundant Representation Modeling of Signals – Theory and Applications By: Michael Elad

# BM<sub>3</sub>D: Block Matching and 3D collaborative Filtering

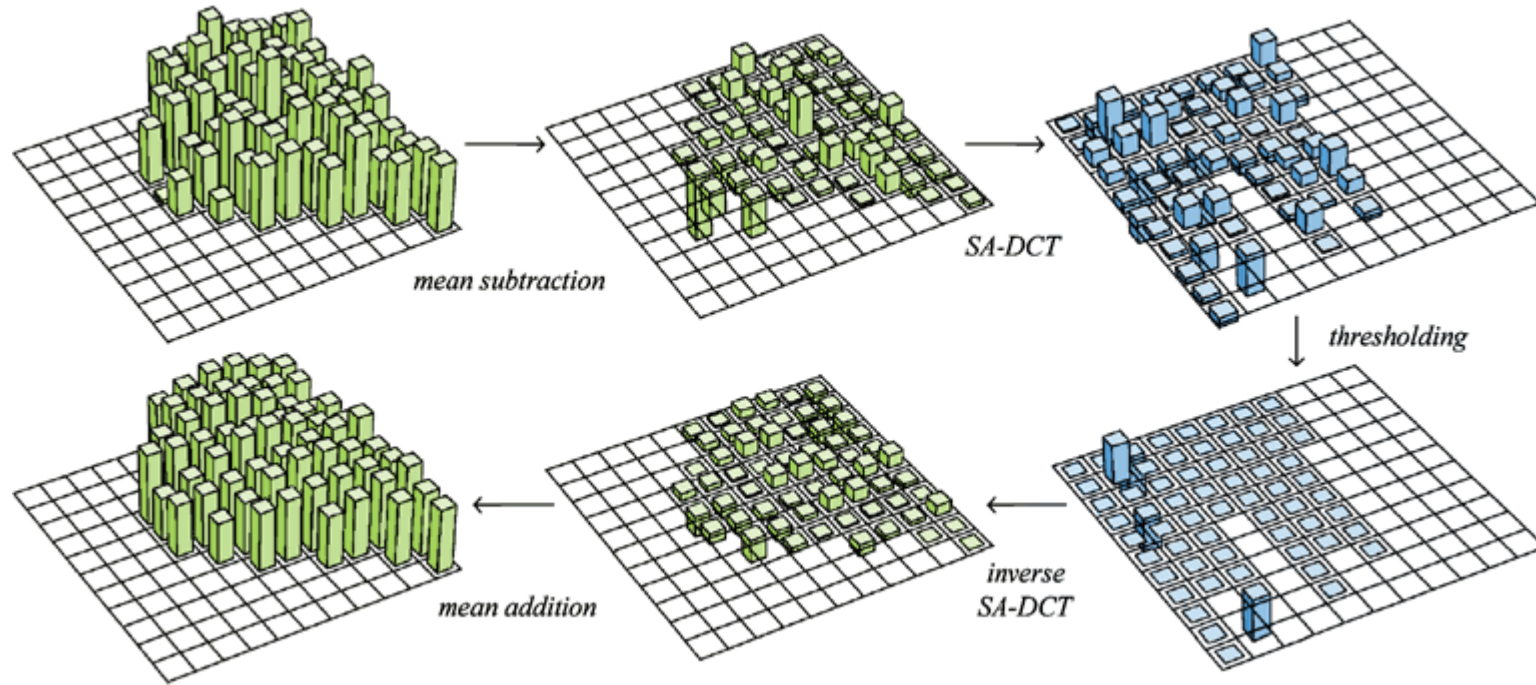




# BM<sub>3</sub>D denoising ( $\sigma = 35$ )



# Shape Adaptive DCT



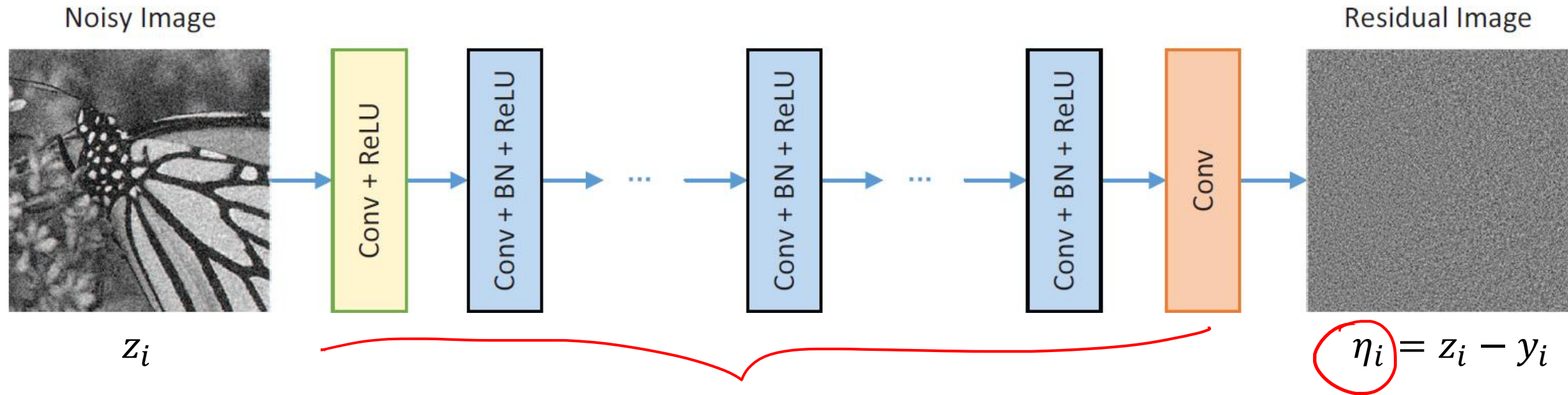


# Shape Adaptive DCT ( $\sigma = 35$ )



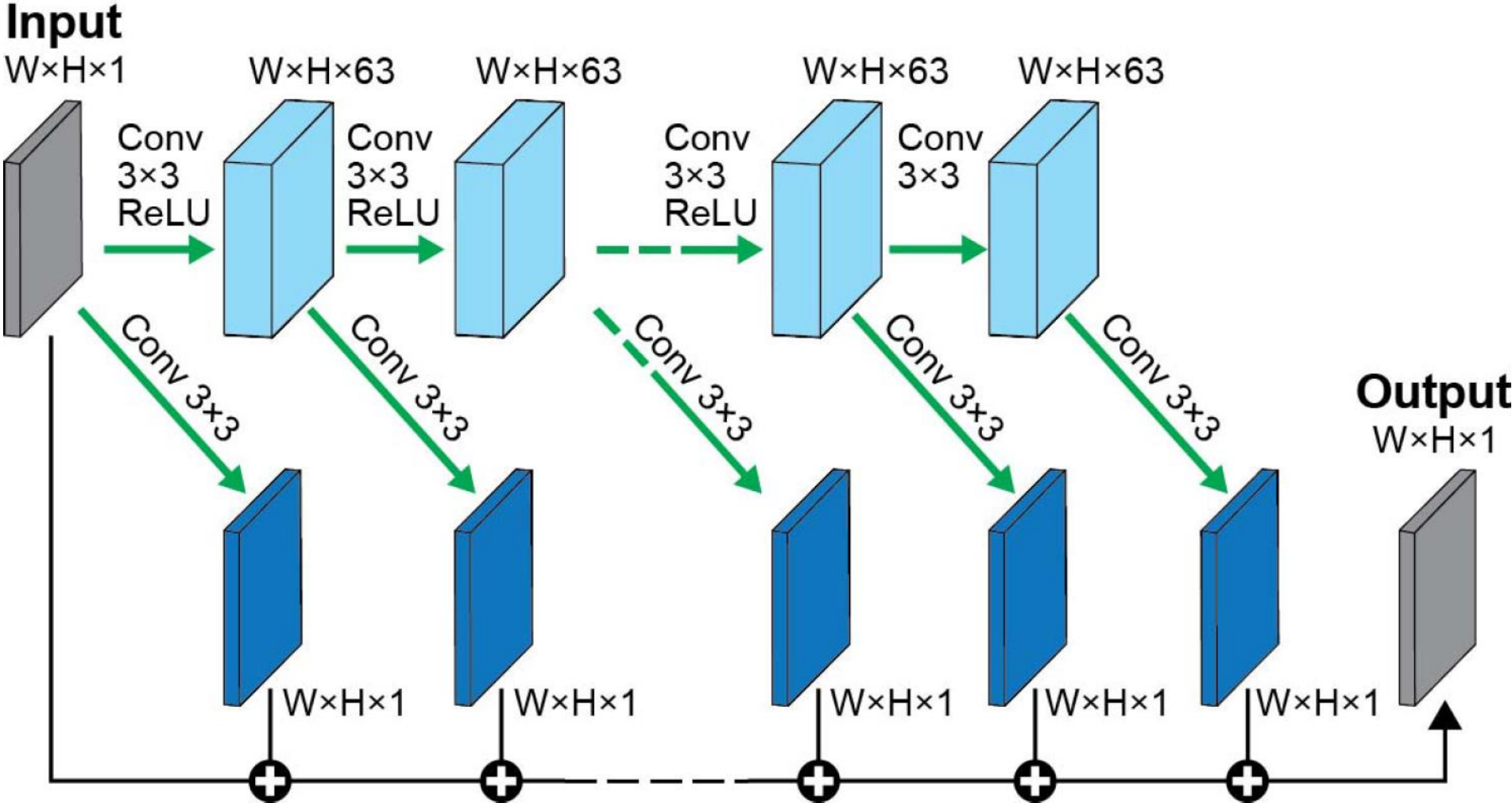
# Deep Learning Methods

# Dn-CNN





# The Network Architecture



Remez, T., Litany, O., Giryas, R., & Bronstein, A. M. (2018). Class-Aware Fully Convolutional Gaussian and Poisson Denoising. IEEE Transactions on Image Processing, 27(11), 5707-5722.

That's all, folks!

