

A Fast Approximation of the Bilateral Filter Using a Signal Processing Approach

Sylvain Paris · Frédo Durand

Received: 15 October 2006 / Accepted: 5 November 2007 / Published online: 18 December 2007
© Springer Science+Business Media, LLC 2007

Abstract The bilateral filter is a nonlinear filter that smoothes a signal while preserving strong edges. It has demonstrated great effectiveness for a variety of problems in computer vision and computer graphics, and fast versions have been proposed. Unfortunately, little is known about the accuracy of such accelerations. In this paper, we propose a new signal-processing analysis of the bilateral filter which complements the recent studies that analyzed it as a PDE or as a robust statistical estimator. The key to our analysis is to express the filter in a higher-dimensional space where the signal intensity is added to the original domain dimensions. Importantly, this signal-processing perspective allows us to develop a novel bilateral filtering acceleration using downsampling in space and intensity. This affords a principled expression of accuracy in terms of bandwidth and sampling. The bilateral filter can be expressed as linear convolutions in this augmented space followed by two simple nonlinearities. This allows us to derive criteria for downsampling the key operations and achieving important acceleration of the bilateral filter. We show that, for the same running time, our method is more accurate than previous acceleration techniques. Typically, we are able to process a 2 megapixel image using our acceleration technique in less than a second, and have the result be visually similar to the exact computation that takes several tens of minutes. The acceleration is most effective with large spatial kernels. Furthermore, this approach extends naturally to color images and cross bilateral filtering.

Keywords Image processing · Bilateral filter · Cross bilateral filter · Color image filtering

1 Introduction

The bilateral filter is a nonlinear filter proposed by Aurich and Weule (1995), Smith and Brady (1997), and Tomasi and Manduchi (1998) to smooth images. It has been adopted for several applications such as image denoising (Tomasi and Manduchi 1998; Liu et al. 2006), relighting and texture manipulation (Oh et al. 2001), dynamic range compression (Durand and Dorsey 2002), illumination correction (Elad 2005), and photograph enhancement (Eisemann and Durand 2004; Petschnigg et al. 2004; Bae et al. 2006). It has also been adapted to other domains such as mesh fairing (Jones et al. 2003; Fleishman et al. 2003), volumetric denoising (Wong et al. 2004), optical flow and motion estimation (Xiao et al. 2006; Sand and Teller 2006), and video processing (Bennett and McMillan 2005; Winnemöller et al. 2006). This large success stems from several origins. First, its formulation and implementation are simple: a pixel is simply replaced by a weighted mean of its neighbors. And it is easy to adapt to a given context as long as a distance can be computed between two pixel values, *e.g.* distance between hair orientations (Paris et al. 2004). The bilateral filter is also non-iterative, thereby achieving satisfying results with only a single pass. This makes the filter's parameters relatively intuitive since their effects are not cumulated over several iterations.

The bilateral filter has proven to be very useful, however it is slow. It is nonlinear and its evaluation is computationally expensive since traditional accelerations, such as performing convolution after an FFT, are not applicable. Brute-force computation is on the order of tens of minutes.

S. Paris (✉) · F. Durand
Computer Science and Artificial Intelligence Laboratory,
Massachusetts Institute of Technology, Cambridge, MA 02139,
USA
e-mail: sparis@csail.mit.edu

Nonetheless, solutions have been proposed to speed up the evaluation of the bilateral filter (Durand and Dorsey 2002; Elad 2002; Pham and van Vliet 2005; Weiss 2006). Unfortunately, most of these methods rely on approximations that are not grounded on firm theoretical foundations, and it is difficult to evaluate the accuracy that is sacrificed.

Overview In this paper, we build on this body of work but we interpret the bilateral filter in terms of signal processing in a higher-dimensional space. This allows us to derive an improved acceleration scheme that yields equivalent running times but dramatically improves accuracy. The key idea of our technique is to analyze the frequency content of this higher-dimensional space. We demonstrate that in this new representation, the signal of interest is mostly low-frequency and can be accurately approximated using coarse sampling, thereby reducing the amount of data to be processed. The quality of the results is evaluated both numerically and visually to characterize the strengths and limitations of the proposed method. Our study shows that our technique is especially fast with large kernels and achieves a close match to the exact computation. As a consequence, our approach is well-suited for applications in the field of computational photography, as illustrated by a few examples. Furthermore, we believe that our new high-dimensional interpretation of images and its low-sampling-rate encoding provide novel and powerful means for edge-preserving image manipulation in general.

This article extends our conference paper (Paris and Durand 2006). We provide more detailed description and discussion, including the algorithm pseudo-code. We conducted new conceptual and quantitative comparisons with existing approximations of the bilateral filter. We describe a new and faster implementation based on direct convolution with a small kernel rather than FFT, and demonstrate and discuss the extension of our acceleration scheme to color-image filtering and cross bilateral filtering. We also refer readers to our recent work that extends the image representation described in this article to other edge-aware applications beyond bilateral filtering such as scribble interpolation, painting, and local histogram equalization (Chen et al. 2007).

2 Related Work

The bilateral filter was first introduced by Aurich and Weule (1995) under the name “nonlinear Gaussian filter”, then by Smith and Brady (1997) as part of the “SUSAN” framework. It was rediscovered later by Tomasi and Manduchi (1998) who called it the “bilateral filter” which is now the most commonly used name. The filter output at each pixel is a weighted average of its neighbors. The weight assigned to

each neighbor decreases with both the distance in the image plane (the *spatial domain* \mathcal{S}) and the distance on the intensity axis (the *range domain* \mathcal{R}). Using a Gaussian G_σ as a decreasing function, and considering a gray-level image I , the result I^b of the bilateral filter is defined by:

$$I_{\mathbf{p}}^b = \frac{1}{W_{\mathbf{p}}^b} \sum_{\mathbf{q} \in \mathcal{S}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) G_{\sigma_r}(|I_{\mathbf{p}} - I_{\mathbf{q}}|) I_{\mathbf{q}} \quad (1a)$$

$$\text{with } W_{\mathbf{p}}^b = \sum_{\mathbf{q} \in \mathcal{S}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) G_{\sigma_r}(|I_{\mathbf{p}} - I_{\mathbf{q}}|). \quad (1b)$$

The parameter σ_s defines the size of the spatial neighborhood used to filter a pixel, and σ_r controls how much an adjacent pixel is downweighted because of the intensity difference. W^b normalizes the sum of the weights.

2.1 Choice of the Weighting Functions

The bilateral filter can be defined with various weighting functions. The choice of the spatial function is mainly driven by its frequency spectrum. Two options have been explored in the literature: a box function (Yaroslavsky 1985; Weiss 2006) that leads to simple computation due to its binary nature but introduces Mach bands and thus requires several iterations of the bilateral filter (see Sect. 6.3 for details), and a Gaussian kernel which is computationally more expensive but does not require iterations. In practice, the latter option is mostly used.

Durand and Dorsey (2002) showed that the choice of the range function can be interpreted in terms of robust statistics (Huber 1981; Hampel et al. 1986; Black et al. 1998). Different functions yield to different behaviors with respect to outliers. For instance, compactly supported functions are insensitive to gross outliers but are not able to filter out extreme defects such as salt-and-pepper noise.

In this article, we focus on the Gaussian bilateral filter which uses Gaussian kernels for both the spatial and range weights because all the practical applications use this version. Furthermore, we will see that these two Gaussian kernels can be elegantly combined into a single higher-dimensional Gaussian kernel leading to a well-defined notion of bandwidth for the bilateral filter. We keep the study of other weighting functions as future work.

2.2 Link with Other Filters

Barash (2002) showed that the two weight functions are actually equivalent to a single weight function based on a distance defined on $\mathcal{S} \times \mathcal{R}$. Using this approach, he related the bilateral filter to adaptive smoothing. Our work follows a similar idea and also uses $\mathcal{S} \times \mathcal{R}$ to describe bilateral filtering. Our formulation is nonetheless significantly different because we not only use the higher-dimensional space for

the definition of a distance, but we also use convolution in this space. Elad (2002) demonstrated that the bilateral filter is similar to using Jacobi iterations to minimize an energy function defined over a large neighborhood. Buades et al. (2005) exposed an asymptotic analysis of the Yaroslavsky filter which is a special case of the bilateral filter with a step function as spatial weight (Yaroslavsky 1985). They proved that asymptotically, the Yaroslavsky filter behaves as the Perona-Malik filter, *i.e.* it alternates between smoothing and shock formation depending on the gradient intensity. Aurich and Weule (1995) pointed out the link with robust statistics (Huber 1981; Hampel et al. 1986; Black et al. 1998). Durand and Dorsey (2002) also cast their study into this framework and showed that the bilateral filter is a w -estimator (Hampel et al. 1986, p. 116). This explains the role of the range weight in terms of sensitivity to outliers. They also pointed out that the bilateral filter can be seen as an extension of the Perona-Malik filter using a larger neighborhood and a single iteration. Weickert et al. (1998) and Barash et al. (2003) described acceleration techniques for PDE filters. They split the multi-dimensional differential operators into combinations of one-dimensional operators that can be efficiently integrated. The obtained speed-up stems from the small spatial footprint of the 1D operators, and the extension to bilateral filtering is unclear. van de Weijer and van den Boomgaard (2001) demonstrated that the bilateral filter is the first iteration of a process that seeks the local mode of the intensity histogram of the adjacent pixels. Mrázek et al. (2006) related bilateral filtering to a large family of nonlinear filters. From a single equation, they expressed filters such as anisotropic diffusion and statistical estimators by varying the neighborhood size and the involved functions.

The main difference between our study and existing work is that the previous approaches link bilateral filtering to another nonlinear filter based on PDEs or statistics whereas we cast our study into a signal processing framework. We demonstrate that the bilateral filter can be mainly computed with linear operations, leaving the nonlinearities to be grouped in a final step.

2.3 Variants of the Bilateral Filter

Higher-Order Filters The bilateral filter implicitly assumes that the desired output should be piecewise constant: such an image is unchanged by the filter when the step discontinuities between constant parts are high enough. Several articles (Elad 2002; Choudhury and Tumblin 2003; Buades et al. 2005) extended the bilateral filter to a *piecewise-linear* assumption. They share the same idea and characterize the local “slope” of the image intensity to better represent the local shape of the signal. Thus, they define a modified filter that better preserves the image characteristics. In particular, they avoid the formation of shocks. We have not explored this direction but it is an interesting avenue for future work.

Cross Bilateral Filter In computational photography applications, it is often useful to decouple the data I to be smoothed from the data E defining the edges to be preserved. For instance, in a “flash no-flash” scenario (Eisemann and Durand 2004; Petschnigg et al. 2004), a picture P^{nf} is taken in a dark environment without flash and another picture P^{f} is taken with flash. Directly smoothing P^{nf} is hard because of the high noise level typical of low-light images. To address this problem, Eisemann and Durand (2004) and Petschnigg et al. (2004) introduced the cross bilateral filter (*a.k.a.* joint bilateral filter) as a variant of the classical bilateral filter. This filter smooths the no-flash picture $P^{\text{nf}} = I$ while relying on the flash version $P^{\text{f}} = E$ to locate the edges to preserve. The definition is similar to (1) except that E replaces I in the range weight G_{σ_r} :

$$I_{\mathbf{p}}^c = \frac{1}{W_{\mathbf{p}}^c} \sum_{\mathbf{q} \in \mathcal{S}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) G_{\sigma_r}(|E_{\mathbf{p}} - E_{\mathbf{q}}|) I_{\mathbf{q}}$$

$$\text{with } W_{\mathbf{p}}^c = \sum_{\mathbf{q} \in \mathcal{S}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) G_{\sigma_r}(|E_{\mathbf{p}} - E_{\mathbf{q}}|).$$

Aurich and Weule (1995) introduced ideas related to the cross bilateral filter, but for a single input image when the filter is iterated. After a number of iterations of bilateral filtering, they filter the original images using range weights derived from the last iteration.

The method described in this article also applies to cross bilateral filtering.

Channel Smoothing Felsberg et al. (2006) described an efficient smoothing method based on a careful design of the intensity weighting function. They showed that B-splines enable the discretization of the intensity range into a small set of *channels*. Filtering these channels yields smooth images with preserved edges akin to the output of the bilateral filter. B-splines allowed for a precise theoretical characterization of their filter using robust statistics. The downside of B-splines is the higher computational effort required to handle them. This approach and ours are complementary and closely related on a number of points that we discuss in Sect. 9.2.

2.4 Fast Methods for Bilateral Filtering

The work most related to ours are the techniques that speed up the evaluation of the bilateral filter. There are two categories of acceleration schemes: specialized filters that perform an exact computation but are restricted to a specific scenario, and approximated filters that are more general but do not produce exact results.

Exact Computation Elad (2002) used Gauss-Seidel iterations to accelerate the convergence of iterative filtering. This technique is only useful when the filter is iterated to reach the stable point, which is not the standard use of the bilateral filter (one iteration or only a few). Weiss (2006) describes an efficient algorithm to incrementally compute the intensity histogram of the square windows surrounding each pixel. This technique is primarily used for median filtering. As an extension, Weiss showed that integrating these histograms weighted by a range function G_{σ_r} (cf. (1)) is equivalent to calculating a bilateral filter where a step function is used on a square window instead of the isotropic Gaussian G_{σ_s} . This filter actually corresponds to a Yaroslavsky filter computed on square neighborhoods. The achieved computation times are on the order of a few seconds for 8 megapixel images. The downside is that the spatial weighting is restricted to a step function that incurs defects such as ripples near strong edges because its Fourier transform contains numerous zeros. In addition, this technique can only handle color images channel-per-channel, which can introduce color-bleeding artifacts. Although these two techniques are fast and produce an exact result, they are too specialized for many applications such as image editing as shown later.

Approximate Computation At the cost of an approximate result, several authors propose fast methods that address more general scenarios. For instance, Weiss (2006) iterated his filter based on square windows to obtain a smoother profile, thereby removing the ripple defects. To avoid shocks that sharpen edges and result in a cartoon look, the range weights G_{σ_r} are kept constant through the iterations *i.e.* they are always evaluated according to the original input picture. van de Weijer and van den Boomgaard (2001) showed that it corresponds to a search for the closest local mode in the neighborhood histogram.

Durand and Dorsey (2002) linearized the bilateral filter which makes possible the use of fast Fourier transforms. They also downsample the data to accelerate the computation to a second or less for one-megapixel images. Although their article mentions FFT computation for the linearized bilateral filter, once the data is downsampled, a direct convolution is more efficient because the kernel is small enough. While their article does not emphasize it, their final results are obtained with direct convolution, without FFT. Our technique is related to their work in that we also express the bilateral filter with linear operations and draw much of our speedup from downsampling. However, our formulation relies on a more principled expression based on a new higher dimensional interpretation of images. This affords a solid signal processing perspective on the bilateral filter and improved accuracy.

Pham and van Vliet (2005) applied a 1D bilateral filter independently on each image row and then on each column. It

produces smooth results that still preserve edges. The convolutions are performed in the spatial domain (without FFT). Compared to brute-force computation, all these approximations yield a better computational complexity and shorter running times suitable for interactive applications, and even real-time processing using modern graphics cards (Winemöller et al. 2006). However, no theoretical study is proposed, and the accuracy of these approximations is unclear. In contrast, we base our technique on signal processing grounds which helps us to define a new and meaningful numerical scheme. Our algorithm performs low-pass filtering in a higher-dimensional space. We show that our approach extends naturally to color images and can also be used to speed up cross-bilateral filtering (Eisemann and Durand 2004; Petschnigg et al. 2004). The cost of a higher-dimensional convolution is offset by downsampling the data without significant accuracy loss, thereby yielding a better precision for running times equivalent to existing methods.

2.5 Contributions

This paper introduces the following contributions:

- An interpretation of the bilateral filter in a signal processing framework. Using a higher dimensional space, we formulate the bilateral filter as a convolution followed by simple nonlinearities.
- Using this higher dimensional space, we demonstrate that the convolution computation can be downsampled without significant impact on the resulting accuracy. This approximation technique enables a speed-up of several orders of magnitude while controlling the induced error.
- We evaluate the accuracy and performance of the proposed acceleration over several scenarios. The obtained results are compared to existing techniques, thereby characterizing the strengths and limitations of our approach.
- We show that this method naturally handles color images and can be easily adapted to cross bilateral filtering.

2.6 Notation

Table 1 summarizes the main notation we use throughout this article. All the vectors in this paper are column vectors. We sometimes use a row notation and omit the transpose sign to prevent clutter in the equations.

3 Signal Processing Approach

We propose a new interpretation of the bilateral filter as a higher-dimensional convolution followed by two nonlinearities. For this, we propose two important re-interpretations of the filter that respectively deal with its two nonlinear components: the normalization division and the data-dependent

Table 1 Notation used in the paper

\mathcal{S}	Spatial domain
$I, W \dots$	2D functions defined on \mathcal{S}
$\mathbf{p}, \mathbf{C} \dots$	Vectors
$\ \mathbf{x}\ $	L_2 norm of vector \mathbf{x}
\otimes	Convolution operator
G_σ	1D Gaussian: $x \mapsto \exp(-\frac{x^2}{2\sigma^2})$
s_s, s_r	Sampling rates (space and range)
I^b	Result of the bilateral filter
\mathcal{R}	Range domain
$i, w \dots$	3D functions defined on $\mathcal{S} \times \mathcal{R}$
$\mathbf{p} \in \mathcal{S}$	Pixel position (2D vector)
$I_{\mathbf{p}} \in \mathcal{R}$	Image intensity at \mathbf{p}
$\delta(x)$	Kronecker symbol (1 if $x = 0$, 0 otherwise)
g_{σ_s, σ_r}	3D Gaussian: $(\mathbf{x}, \zeta) \in \mathcal{S} \times \mathcal{R} \mapsto \exp(-\frac{\mathbf{x}\cdot\mathbf{x}}{2\sigma_s^2} - \frac{\zeta^2}{2\sigma_r^2})$
σ_s, σ_r	Gaussian parameters (space and range)
W^b	Normalization factor

weights (1). First, we define a homogeneous intensity that will allow us to obtain the normalization term $W_{\mathbf{p}}^b$ as a homogeneous component. Second, we introduce an additional dimension to the 2D image domain, corresponding to the image intensity *a.k.a.* range. While the visualization of images as height fields in such 3D space is not new, we actually go further and interpret filtering using functions over this full 3D domain, which allows us to express the bilateral filter as a linear shift-invariant convolution in 3D. This convolution is followed by simple pixel-wise nonlinearities to extract the relevant output and perform the normalization by our homogeneous component.

3.1 Homogeneous Intensity

Similar to any weighted average, the two lines of (1) that define the bilateral filter output and its normalization factor are almost the same. The main difference is the normalizing factor $W_{\mathbf{p}}^b$ and the image intensity $I_{\mathbf{q}}$ in the first equation. We emphasize this similarity by multiplying both sides of (1a) by $W_{\mathbf{p}}^b$. We then rewrite the equations using two-dimensional vectors:

$$\begin{pmatrix} W_{\mathbf{p}}^b I_{\mathbf{p}}^b \\ W_{\mathbf{p}}^b \end{pmatrix} = \sum_{\mathbf{q} \in \mathcal{S}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) G_{\sigma_r}(|I_{\mathbf{p}} - I_{\mathbf{q}}|) \begin{pmatrix} I_{\mathbf{q}} \\ 1 \end{pmatrix}, \quad (2)$$

where G_{σ_s} and G_{σ_r} are Gaussian functions, \mathcal{S} is the spatial domain, I the input image, and I^b the result of the bilateral filter. To maintain the property that the bilateral filter is a weighted mean, we assign a weight $W = 1$ to the input values:

$$\begin{pmatrix} W_{\mathbf{p}}^b I_{\mathbf{p}}^b \\ W_{\mathbf{p}}^b \end{pmatrix} = \sum_{\mathbf{q} \in \mathcal{S}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) G_{\sigma_r}(|I_{\mathbf{p}} - I_{\mathbf{q}}|) \begin{pmatrix} W_{\mathbf{q}} I_{\mathbf{q}} \\ W_{\mathbf{q}} \end{pmatrix}. \quad (3)$$

By assigning a couple $(W_{\mathbf{q}} I_{\mathbf{q}}, W_{\mathbf{q}})$ to each pixel \mathbf{q} , we express the filtered pixels as linear combinations of their adjacent pixels. Of course, we have not “removed” the division since to access the actual value of the intensity, the first coordinate (WI) still has to be divided by the second one (W). This is similar to homogeneous coordinates used in projective geometry. Adding an extra coordinate to our data makes most of the computation pipeline computable with linear operations; a division is made only at the final stage. Inspired by this parallel, we call the two-dimensional vector (WI, W) the *homogeneous intensity*. It is also related to the use of pre-multiplied alpha in image algebra (Porter and Duff 1984; Blinn 1996; Willis 2006). We discuss this aspect in Sect. 7.1.1.

An important aspect of this homogeneous formulation is its *exact equivalence* with the original formulation (1). Using a homogeneous representation enables a simpler formulation of the filter. However, although (3) is a linear combination, this does not define a linear filter yet since the weights depend on the actual values of the pixels. The next section addresses this issue.

3.2 The Bilateral Filter as a Convolution

If we ignore the term $G_{\sigma_r}(|I_{\mathbf{p}} - I_{\mathbf{q}}|)$, (3) is a classical convolution by a Gaussian kernel: $(W^b I^b, W^b) = G_{\sigma_s} \otimes (WI, W)$. Furthermore, it has been pointed out that the product of the spatial and range Gaussian defines a higher dimensional Gaussian in the 3D product space of the domain and range of the image. However, this 3D Gaussian interpretation has so far only been used to illustrate the weights of the filter, not to linearize computation. Since these weights are in 3D but the summation in (1) is only over the 2D spatial domain, it does not define a convolution. To overcome this, we push the 3D interpretation further and define an intensity value for each point of the product space so that we can define a summation over this full 3D space.

Formally, we introduce an additional dimension ζ and define the intensity I for each point (x, y, ζ) . With the Kronecker symbol $\delta(\zeta)$ ($\delta(0) = 1$, $\delta(\zeta) = 0$ otherwise) and \mathcal{R} the interval on which the intensity is defined, we rewrite (3) using $\delta(\zeta - I_{\mathbf{q}})$ such that the terms are cancelled when $\zeta \neq I_{\mathbf{q}}$:

$$\begin{pmatrix} W_{\mathbf{p}}^b I_{\mathbf{p}}^b \\ W_{\mathbf{p}}^b \end{pmatrix} = \sum_{\mathbf{q} \in \mathcal{S}} \sum_{\zeta \in \mathcal{R}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) \times G_{\sigma_r}(|I_{\mathbf{p}} - \zeta|) \delta(\zeta - I_{\mathbf{q}}) \begin{pmatrix} W_{\mathbf{q}} I_{\mathbf{q}} \\ W_{\mathbf{q}} \end{pmatrix}. \quad (4)$$

Equation (4) is a sum over the product space $\mathcal{S} \times \mathcal{R}$. We now focus on this space. We use lowercase names for the

functions defined on $\mathcal{S} \times \mathcal{R}$. The product $G_{\sigma_s}G_{\sigma_r}$ defines a separable Gaussian kernel g_{σ_s, σ_r} on $\mathcal{S} \times \mathcal{R}$:

$$g_{\sigma_s, \sigma_r}: (\mathbf{x} \in \mathcal{S}, \zeta \in \mathcal{R}) \mapsto G_{\sigma_s}(\|\mathbf{x}\|)G_{\sigma_r}(|\zeta|). \tag{5}$$

From the remaining part of (4), we build two functions i and w :

$$i: (\mathbf{x} \in \mathcal{S}, \zeta \in \mathcal{R}) \mapsto I_{\mathbf{x}} \tag{6a}$$

$$w: (\mathbf{x} \in \mathcal{S}, \zeta \in \mathcal{R}) \mapsto \delta(\zeta - I_{\mathbf{x}}) = \delta(\zeta - I_{\mathbf{x}})W_{\mathbf{x}} \text{ since } W_{\mathbf{x}} = 1. \tag{6b}$$

With definitions (6), we rewrite the right side of (4):

$$\begin{aligned} \delta(\zeta - I_{\mathbf{q}}) \begin{pmatrix} W_{\mathbf{q}}I_{\mathbf{q}} \\ W_{\mathbf{q}} \end{pmatrix} &= \begin{pmatrix} \delta(\zeta - I_{\mathbf{q}})W_{\mathbf{q}}I_{\mathbf{q}} \\ \delta(\zeta - I_{\mathbf{q}})W_{\mathbf{q}} \end{pmatrix} \\ &= \begin{pmatrix} w(\mathbf{q}, \zeta)i(\mathbf{q}, \zeta) \\ w(\mathbf{q}, \zeta) \end{pmatrix}. \end{aligned} \tag{7}$$

Then with definition (5), we get:

$$\begin{aligned} \begin{pmatrix} W_{\mathbf{p}}^b I_{\mathbf{p}}^b \\ W_{\mathbf{p}}^b \end{pmatrix} &= \sum_{(\mathbf{q}, \zeta) \in \mathcal{S} \times \mathcal{R}} g_{\sigma_s, \sigma_r}(\mathbf{p} - \mathbf{q}, I_p - \zeta) \\ &\quad \times \begin{pmatrix} w(\mathbf{q}, \zeta)i(\mathbf{q}, \zeta) \\ w(\mathbf{q}, \zeta) \end{pmatrix}. \end{aligned} \tag{8}$$

The above formula corresponds to the value at point (\mathbf{p}, I_p) of a convolution between g_{σ_s, σ_r} and the two-dimensional function (wi, w) :

$$\begin{pmatrix} W_{\mathbf{p}}^b I_{\mathbf{p}}^b \\ W_{\mathbf{p}}^b \end{pmatrix} = \left[g_{\sigma_s, \sigma_r} \otimes \begin{pmatrix} wi \\ w \end{pmatrix} \right] (\mathbf{p}, I_p). \tag{9}$$

According to the above equation, we introduce the functions i^b and w^b :

$$(w^b i^b, w^b) = g_{\sigma_s, \sigma_r} \otimes (wi, w). \tag{10}$$

Thus, we have reached our goal. The bilateral filter is expressed as a convolution followed by nonlinear operations:

linear: 3D convolution

$$(w^b i^b, w^b) = g_{\sigma_s, \sigma_r} \otimes (wi, w), \tag{11a}$$

nonlinear: slicing + division

$$I_{\mathbf{p}}^b = \frac{w^b(\mathbf{p}, I_p) i^b(\mathbf{p}, I_p)}{w^b(\mathbf{p}, I_p)}. \tag{11b}$$

The nonlinear section is actually composed of two operations. The functions $w^b i^b$ and w^b are evaluated at point

(\mathbf{p}, I_p) . We name this operation *slicing*. The second nonlinear operation is the division that retrieves the intensity value from the homogeneous vector. This division corresponds to applying the normalization that we delayed earlier. In our case, slicing and division commute *i.e.* the result is independent of their order because g_{σ_s, σ_r} is positive and w values are 0 and 1, which ensures that w^b is positive.

3.3 Intuition

To gain more intuition about our formulation of the bilateral filter, we propose an informal description of the process before further discussing its consequences.

The spatial domain \mathcal{S} is a classical xy image plane and the range domain \mathcal{R} is a simple axis labelled ζ . The w function can be interpreted as “the plot in the $xy\zeta$ space of $\zeta = I(x, y)$ ” *i.e.* w is null everywhere except on the points $(x, y, I(x, y))$ where it is equal to 1. The wi product is similar to w . Instead of using binary values 0 or 1 to “plot I ”, we use 0 or $I(x, y)$ *i.e.* it is a plot with a pen whose brightness equals the plotted value. An example is shown in Fig. 1.

Using these two functions wi and w , the bilateral filter is computed as follows. First, we “blur” wi and w , that is we convolve wi and w with a Gaussian defined on $xy\zeta$. This results in the functions $w^b i^b$ and w^b . For each point of the $xy\zeta$ space, we compute $i^b(x, y, \zeta)$ by dividing $w^b i^b(x, y, \zeta)$ by $w^b(x, y, \zeta)$. The final step is to get the value of the pixel (x, y) of the filtered image I^b . This corresponds directly to the value of i^b at $(x, y, I(x, y))$ which is the point where the input image I was “plotted”. Figure 1 illustrates this process on a simple 1D image.

Note that although the 3D functions get smoothed aggressively, the slicing nonlinearity makes the output of the filter piecewise-smooth and the strong edges of the input are preserved.

4 Fast Approximation

The key idea which speeds up the computation is computing the 3D convolution at a coarse resolution. For this, we demonstrate that the wi and w functions can be downsampled without introducing significant errors. In fact, we never construct the full-resolution product space. This ensures the good memory and speed performance of our method. We discuss the practical implementation of this strategy and analyze the accuracy and performance of the proposed technique.

4.1 Downsampling the Data

We have shown that the bilateral filter can be interpreted as a Gaussian filter in a product space. Our acceleration scheme

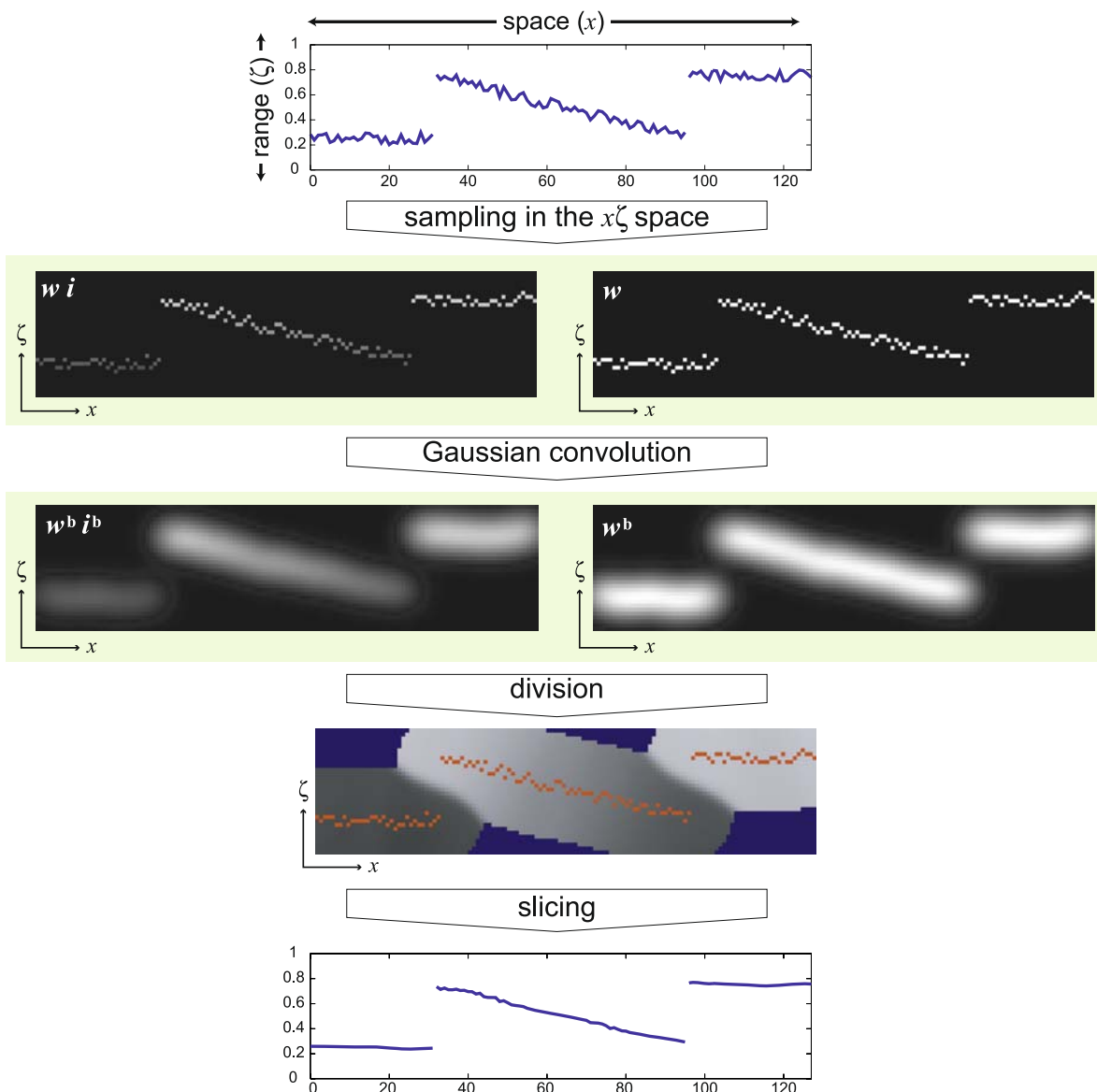


Fig. 1 (Color online) Our computation pipeline applied to a 1D signal. The original data (top row) are represented by a two-dimensional function (w_i, w) (second row). This function is convolved with a Gaussian kernel to form $(w^b i^b, w^b)$ (third row). The first component is then di-

vided by the second (fourth row, blue area is undefined because of numerical limitation, $w^b \approx 0$). Then the final result (last row) is extracted by sampling the former result at the location of the original data (shown in red on the fourth row)

directly follows from the fact that this operation is a low-pass filter. $(w^b i^b, w^b)$ is a band-limited function which is well approximated by its low frequencies. According to the sampling theorem (Shannon 1949; Smith 2002, p. 35), it is sufficient to sample with a rate at least half of the smallest wavelength considered. We exploit this idea by performing the convolution at a lower resolution. Formally, we down-sample (w_i, w) , perform the convolution, and up-sample the result as indicated by the following equations. Note, however, that our implementation never stores full-resolution data: the high-resolution data is built at each pixel and down-sampled on the fly, and we up-sample only at the slicing

location. The notion of using a high-resolution 3D space which we then downsample is used only for formal exposition (cf. Sect. 4.2 for implementation details):

$$(w_{\downarrow} i_{\downarrow}, w_{\downarrow}) = \text{downsample}(w_i, w) \quad [\text{computed on the fly}], \tag{12a}$$

$$(w_{\downarrow}^b i_{\downarrow}^b, w_{\downarrow}^b) = g_{\sigma_s, \sigma_r} \otimes (w_{\downarrow} i_{\downarrow}, w_{\downarrow}), \tag{12b}$$

$$(w_{\downarrow \uparrow}^b i_{\downarrow \uparrow}^b, w_{\downarrow \uparrow}^b) = \text{upsample}(w_{\downarrow}^b i_{\downarrow}^b, w_{\downarrow}^b) \quad [\text{evaluated only at slicing location}]. \tag{12c}$$

The rest of the computation remains the same except that we slice and divide $(w_{\downarrow\uparrow}^b i_{\downarrow\uparrow}^b, w_{\downarrow\uparrow}^b)$ instead of $(w^b i^b, w^b)$, using the same $(\mathbf{p}, I_{\mathbf{p}})$ points. Since slicing occurs at points where $w = 1$, it guarantees $w^b \geq g_{\sigma_s, \sigma_r}(0)$, which ensures that we do not divide by small numbers that would degrade our approximation.

We use box-filtering for the prefilter of the downsampling (*a.k.a.* average downsampling), that is, we first convolve the data with a box profile before downsampling. For upsampling, we use linear interpolation. While these filters do not have perfect frequency responses, they offer much better performances than schemes such as tri-cubic filters. We name s_s and s_r the sampling rates of \mathcal{S} and \mathcal{R} , *i.e.* we use a box function that measures $s_s \times s_s$ pixels and s_r intensity units.

4.2 Implementation

This section details the actual implementation of our algorithm. We have not performed low-level optimization. Although some optimization may be introduced by compilers (we used GCC 3.3.5), the code used in these tests does not explicitly rely on vector instructions of modern CPU nor on the streaming capacities of recent GPU. In a dedicated article (Chen et al. 2007), we describe a graphics-hardware implementation that achieves performances two orders of magnitude faster than the CPU implementation tested here. Our code is publicly available on our webpage: <http://people.csail.mit.edu/sparis/bf#code>. The software is open-source and under the MIT license.

4.2.1 Design Overview

In order to achieve high performance, we never build the $\mathcal{S} \times \mathcal{R}$ space at the fine resolution. We only deal with the downsampled version. In practice, this means that we directly construct the downsampled $\mathcal{S} \times \mathcal{R}$ domain from the image, *i.e.* we store each pixel directly into the corresponding coarse bin. At the slicing stage, we evaluate the upsampled values only at the points $(\mathbf{p}, I_{\mathbf{p}})$, *i.e.* we do not upsample the entire $\mathcal{S} \times \mathcal{R}$ space, only the points corresponding to the final result.

4.2.2 Pseudo-Code

Our algorithm is summarized in Fig. 2. Step 3 directly computes the downsampled versions of wi and w from the input image I . The high-resolution value (wi, w) computed in Step 3a is a temporary variable. Downsampling is performed on the fly: We compute (wi, w) (Step 3a) and directly add it into the low-resolution array $(w_{\downarrow} i_{\downarrow})$ (Step 3c). In Step 3b, we offset the \mathcal{R} coordinate to ensure that the array indices

start at 0. Note that in Step 3c, we do not need to normalize the values by the size of the downsampling box. This results in a uniform scaling that propagates through the convolution and upsampling operations. It is cancelled by the final normalization (Step 5b) and thus does not affect the result. Step 4 performs the convolution. We discuss two different options in the following section. Steps 5a and 5b correspond respectively to the slicing and division nonlinearities described previously.

For completeness, Fig. 3 gives the pseudo-code of a direct implementation of (1).

4.2.3 Efficient Convolution

Convolution is at the core of our technique and its efficiency is important to the overall performance of our algorithm. We describe two options that will be evaluated in Sect. 5.

Full Kernel One can use the fast Fourier transform to efficiently compute the convolution. This has the advantage that the computation does not depend on the Gaussian size, but on only the domain size. Thus, we can use an exact Gaussian kernel. In this case, the approximation comes only from the downsampling and upsampling applied to the data before and after the convolution and from cross-boundary artifacts inherent in FFT. To minimize these artifacts, the domain is padded with zeros over 2σ .

Truncated Kernel When the convolution kernel is small, an explicit computation in the spatial domain is an effective alternative since only a limited number of samples are involved. Although the Gaussian kernel has no compact support, its tail falls off quickly. We therefore use the classical approximation by truncating the kernel beyond 2σ . We only apply this technique for the case where the sampling rate equals the Gaussian standard deviation (*i.e.* $s = \sigma$) since, in this case, the downsampled kernel is isotropic and has a variance equal to 1. The truncated kernel then covers a $5 \times 5 \times 5$ region which is compact enough to ensure fast computation. Since the kernel is shift-invariant and separable, we further shorten the running times by replacing the 3D kernel by three 1D kernels. This well-known technique reduces the number of points to an average of 15 ($= 3 \times 5$) instead of 125 ($= 5^3$). Unlike the separable approximation (Pham and van Vliet 2005), this separation is exact since the original 3D kernel is Gaussian and therefore separable.

We shall see that this convolution in the spatial domain with a truncated kernel yields a better ratio of numerical accuracy over running time than the frequency-space convolution with the full kernel. This latter option remains nonetheless useful to achieve high numerical accuracy, at the cost of slower performances.

FAST BILATERAL FILTER

Input: Image I Gaussian parameters σ_s and σ_r Sampling rates s_s and s_r Output: Filtered image I^b

1. Initialize all $w_{\downarrow}i_{\downarrow}$ and w_{\downarrow} values to 0.
2. Compute the minimum intensity value:

$$I_{\min} \leftarrow \min_{(X,Y) \in \mathcal{S}} I(X, Y)$$

3. For each pixel $(X, Y) \in \mathcal{S}$ with an intensity $I(X, Y) \in \mathcal{R}$
 - (a) Compute the homogeneous vector (wi, w) :

$$(wi, w) \leftarrow (I(X, Y), 1)$$

- (b) Compute the downsampled coordinates (with $[\cdot]$ the rounding operator)

$$(x, y, \zeta) \leftarrow \left(\left[\frac{X}{s_s} \right], \left[\frac{Y}{s_s} \right], \left[\frac{I(X, Y) - I_{\min}}{s_r} \right] \right)$$

- (c) Update the downsampled $\mathcal{S} \times \mathcal{R}$ space

$$\begin{pmatrix} w_{\downarrow}i_{\downarrow}(x, y, \zeta) \\ w_{\downarrow}(x, y, \zeta) \end{pmatrix} \leftarrow \begin{pmatrix} w_{\downarrow}i_{\downarrow}(x, y, \zeta) \\ w_{\downarrow}(x, y, \zeta) \end{pmatrix} + \begin{pmatrix} wi \\ w \end{pmatrix}$$

4. Convolve $(w_{\downarrow}i_{\downarrow}, w_{\downarrow})$ with a 3D Gaussian g whose parameters are σ_s/s_s and σ_r/s_r

$$(w_{\downarrow}^b i_{\downarrow}^b, w_{\downarrow}^b) \leftarrow (w_{\downarrow}i_{\downarrow}, w_{\downarrow}) \otimes g$$

5. For each pixel $(X, Y) \in \mathcal{S}$ with an intensity $I(X, Y) \in \mathcal{R}$
 - (a) Tri-linearly interpolate the functions $w_{\downarrow}^b i_{\downarrow}^b$ and w_{\downarrow}^b to obtain $W^b I^b$ and W^b :

$$W^b I^b(X, Y) \leftarrow \text{interpolate} \left(w_{\downarrow}^b i_{\downarrow}^b, \frac{X}{s_s}, \frac{Y}{s_s}, \frac{I(X, Y)}{s_r} \right)$$

$$W^b(X, Y) \leftarrow \text{interpolate} \left(w_{\downarrow}^b, \frac{X}{s_s}, \frac{Y}{s_s}, \frac{I(X, Y)}{s_r} \right)$$

- (b) Normalize the result

$$I^b(X, Y) \leftarrow \frac{W^b I^b(X, Y)}{W^b(X, Y)}$$

Fig. 2 Pseudo-code of our algorithm. The algorithm is designed such that we never build the full-resolution $\mathcal{S} \times \mathcal{R}$ space

5 Evaluation of Our Approximation

This section investigates the accuracy and performance of our technique compared to the exact computation. The timings are measured on an Intel Xeon 2.8 GHz with 1 MB cache using double-precision floating-point numbers.

On Ground Truth In practical applications such as photograph enhancement (Bae et al. 2006) or user-driven image editing (Weiss 2006), the notions of numerical accuracy and

ground truth are not always well defined. Running times can be objectively measured, but other aspects are elusive and difficult to quantify. Furthermore, whether the bilateral filter is the ideal filter for an application is a separate question.

Since our method achieves acceleration through approximation, we have chosen to measure the numerical accuracy by comparing the outputs of our technique and of other existing approximations with the result of the original bilat-

BRUTE FORCE BILATERAL FILTER

Input: Image I
 Gaussian parameters σ_s and σ_r

Output: Filtered image I^b

1. Initialize all I^b and W^b values to 0.
2. For each pixel $(X, Y) \in \mathcal{S}$ with an intensity $I(X, Y) \in \mathcal{R}$
 - (a) For each pixel $(X', Y') \in \mathcal{S}$ with an intensity $I(X', Y') \in \mathcal{R}$
 - (i) Compute the associated weight:

$$\text{weight} \leftarrow \exp\left(-\frac{(I(X', Y') - I(X, Y))^2}{2\sigma_r^2} - \frac{(X' - X)^2 + (Y' - Y)^2}{2\sigma_s^2}\right)$$

- (ii) Update the weight sum $W^b(X, Y)$:

$$W^b(X, Y) \leftarrow W^b(X, Y) + \text{weight}$$

- (iii) Update $I^b(X, Y)$:

$$I^b(X, Y) \leftarrow I^b(X, Y) + \text{weight} \times I^b(X', Y')$$

- (b) Normalize the result:

$$I^b(X, Y) \leftarrow \frac{I^b(X, Y)}{W^b(X, Y)}$$

Fig. 3 Pseudo-code of the brute force algorithm

eral filter. This comparison pertains only to the “numerical quality” of the approximation, and the readers should keep in mind that numerical differences do not necessarily produce unsatisfying outputs. To balance this numerical aspect, we also provide visual results to let the readers examine by themselves the outputs. Important criteria are then the regularity of the achieved smoothing, artifacts that add visual features which do not exist in the input picture, tone (or color) faithfulness, and so on.

5.1 Numerical Accuracy

To evaluate the error induced by our approximation, we compare the result $I_{\downarrow\uparrow}^b$ from our fast algorithm to the exact result I^b obtained from (1). We have chosen three images as different as possible to cover a broad spectrum of content (Fig. 12):

- An artificial image with various edges, frequencies, and white noise.
- An architectural picture structured along two main directions.
- And a photograph of a natural scene with more stochastic structure.

To express numerical accuracy, we compute the peak signal-to-noise ratio (PSNR) considering $\mathcal{R} = [0; 1]$:

$\text{PSNR}(I_{\downarrow\uparrow}^b) = -10 \log_{10}(\frac{1}{|\mathcal{S}|} \sum_{\mathbf{p} \in \mathcal{S}} |I_{\downarrow\uparrow}^b(\mathbf{p}) - I^b(\mathbf{p})|^2)$. For instance, considering intensity values encoded on 8 bits, if two images differ from one gray level at each pixel, the resulting PSNR is 48 dB. As a guideline, PSNR values above 40 dB often corresponds to limited differences, almost invisible. This should be confirmed by a visual inspection since a high PSNR can “hide” a few large errors because it is a mean over the whole image.

The box downsampling and linear upsampling schemes yield very satisfying results while being computationally efficient. We tried other techniques such as tri-cubic downsampling and upsampling. Our tests showed that their computational cost exceeds the accuracy gain, *i.e.* a similar gain can be obtained in a shorter time using a finer sampling of the $\mathcal{S} \times \mathcal{R}$ space. The results presented in this paper use box downsampling and linear upsampling. We experimented with several sampling rates (s_s, s_r) for $\mathcal{S} \times \mathcal{R}$. The meaningful quantities to consider are the ratios $(\frac{s_s}{\sigma_s}, \frac{s_r}{\sigma_r})$ that indicate the relative position of the frequency cutoff due to downsampling with respect to the bandwidth of the filter we apply. Small ratios correspond to limited approximations and high ratios to more aggressive downsamplings. A consistent approximation is a sampling rate proportional to the Gaussian bandwidth (*i.e.* $\frac{s_s}{\sigma_s} \approx \frac{s_r}{\sigma_r}$) to achieve similar accuracy on the whole $\mathcal{S} \times \mathcal{R}$ domain. The results plotted in Fig. 4 show that

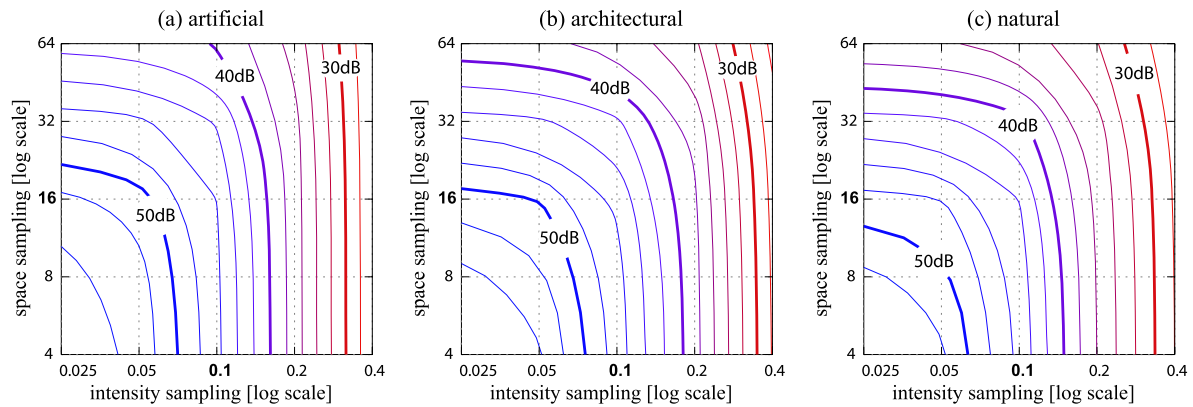


Fig. 4 Accuracy evaluation. All the images are filtered with ($\sigma_s = 16$, $\sigma_r = 0.1$). The PSNR in dB is evaluated at various sampling rates of \mathcal{S} and \mathcal{R} (greater is better). Our approximation scheme is more robust to space downsampling than range downsampling. It is also slightly

more accurate on structured scenes (a, b) than on stochastic ones (c). This test uses the full-kernel implementation of the 3D convolution based on FFT

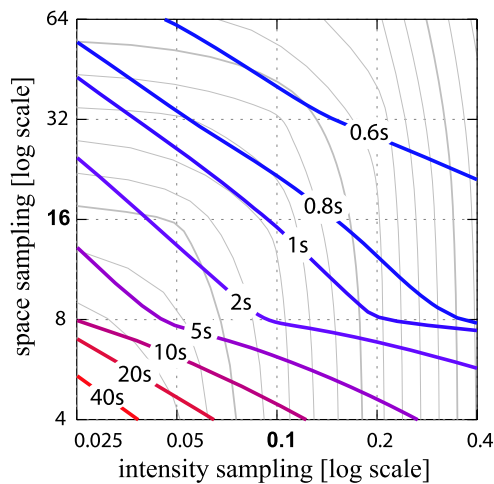


Fig. 5 Running times on the architectural picture with ($\sigma_s = 16$, $\sigma_r = 0.1$). The PSNR isolines are plotted in gray. Exact computation takes several tens of minutes (varying with the actual implementation). This test uses the full-kernel implementation of the 3D convolution based on FFT

this remark is globally valid in practice. A closer look at the plots reveals that \mathcal{S} can be slightly more downsampled than \mathcal{R} . This is probably due to the nonlinearities and the anisotropy of the signal.

5.2 Running Times

Figure 5 shows the running times for the architectural picture with the same settings. In theory, the gain from space downsampling should be twice the one from range downsampling since \mathcal{S} is two-dimensional and \mathcal{R} one-dimensional. In practice, the nonlinearities and caching issues induce minor deviations. Combining this plot with the PSNR plot (in gray under the running times) allows for selecting the best sampling parameters for a given error tol-

erance or a given time budget. As a simple guideline, using sampling steps equal to σ_s and σ_r produces results without visual difference with the exact computation (see Fig. 12). Our scheme achieves a dramatic speed-up since direct computation of (1) lasts several tens of minutes (varying with the actual implementation). Our approximation requires one second.

Effect of the Kernel Size An important aspect of our approximation is visible on Fig. 6 (right). Our technique runs faster with larger kernels. Indeed, when σ increases, keeping the level of approximation constant (the ratio $\frac{s}{\sigma}$ is fixed) allows for a more important downsampling, that is, larger sampling steps s . The rationale of this behavior is that the more the image is smoothed, the more the high frequencies of $\mathcal{S} \times \mathcal{R}$ are attenuated, and the more we can afford to discard them without incurring visible differences.

Truncated Kernel versus Full Kernel Figure 6 also shows that, for a limited loss of accuracy, the direct convolution with a truncated kernel yields running times significantly shorter than the Fourier-domain convolution with a full kernel. We believe that the direct convolution with the 5^3 kernel is a better choice for most applications.

High and Low Resolution Filtering To evaluate the usability of our technique in a professional environment, we processed a 8 megapixel image and obtained a running time of 2.5 s. This is in the same order of magnitude as the iterated-box filter and orders of magnitude faster than currently available professional software (see Weiss 2006 for details). We also experimented with an image at DVD resolution, *i.e.* $576 \times 320 \approx 0.2$ megapixel, and we obtained a running time of about 57 ms near the 40 ms required

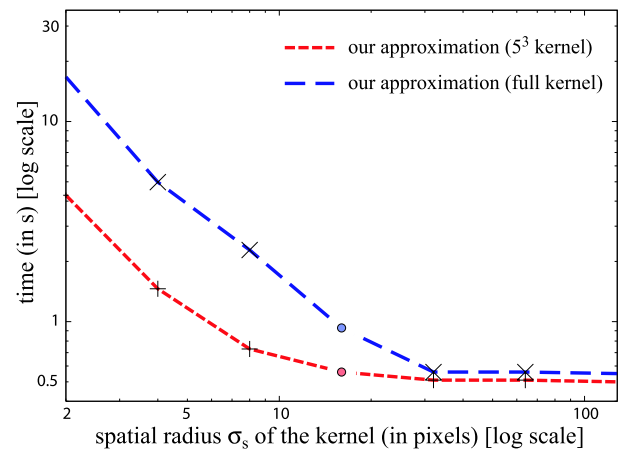
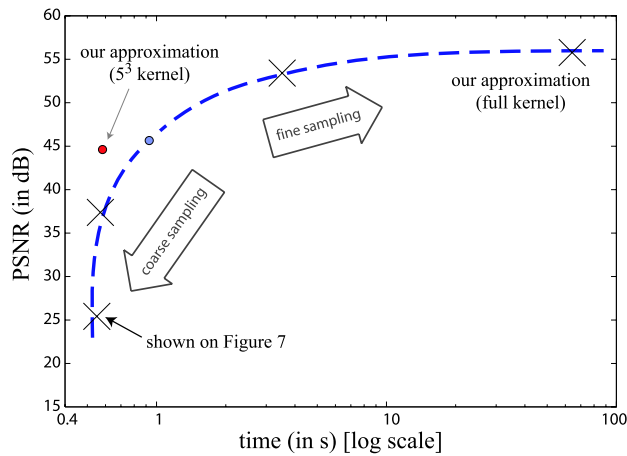


Fig. 6 (Color online) *Left*: Accuracy-versus-time comparison. Our approximations are tested on the architectural picture (1600 × 1200) using $\sigma_s = 16$ and $\sigma_r = 0.1$ as parameters. The full-kernel can achieve variable degrees of accuracy and speed-up by varying the sampling rates of $S \times R$. We tested it with the following sampling rates of $S \times R$ (from left to right): (4, 0.025), (8, 0.05), (16, 0.1), (32, 0.2), (64, 0.4). The result with the coarsest sampling

($s_s = 64, s_r = 0.4$) is shown on Fig. 7. *Right*: Time-versus-radius comparison. Our approximations are tested on the architectural image using $\sigma_r = 0.1$. The full-kernel technique uses a sampling rate equal to (σ_s, σ_r) . Remember that an exact computation lasts at least several tens of minutes (varying with the implementation). The *color dots* in the *left* and *right* plots correspond to the same measurements

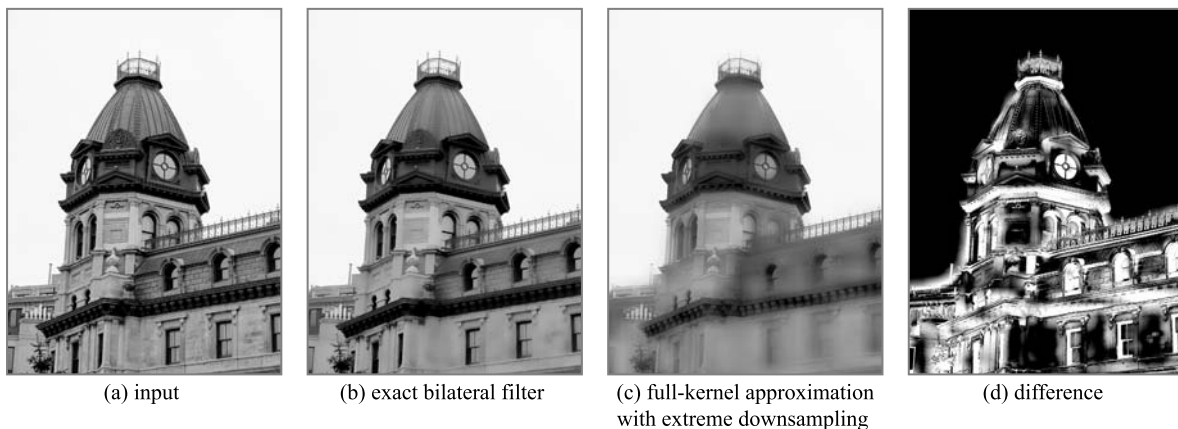


Fig. 7 We filtered the architectural picture (a) and applied an extreme downsampling to the $S \times R$ domain, four times the filter bandwidth, *i.e.* $(\sigma_s, \sigma_r) = (16, 0.1)$ and $(s_s, s_r) = (64, 0.4)$. The result (c) exhibits numerous visual artifacts: the whole image is overly smoothed and the smoothing strength varies spatially (observe the roof railing for instance). Compared to the exact computation (b), the achieved

PSNR is 25 dB. The difference image (with a 10× multiplier) shows that large errors cover most of the image (d). We do not recommend such extreme downsampling since the speed gain is limited while the accuracy loss is important as shown on Fig. 6 (left) (this result corresponds to the lower-left measured point on the dashed curve)

to achieve real-time performances, *i.e.* a 25 Hz frame rate. This motivated us to run our algorithm on our most recent machine, an AMD Opteron 252 at 2.6 GHz with 1 MB of cache, in order to simulate a professional setup. With this machine, the running time is 38 ms. This means that our algorithm achieves real-time video processing in software. Although this does not leave any time to apply other effects as demonstrated by Winnemöller et al. (2006) using the GPU, we believe that our technique paves the way toward interactive video applications that are purely software-based.

Bottlenecks Figures 6 and 7 show that when the sampling rate becomes too coarse, the accuracy suffers dramatically but the running time does not improve. This is because slicing becomes the bottleneck of the algorithm: The trilinear interpolation necessary at each output pixel becomes significantly more costly than the operations performed on the 3D data (see Table 2). On the other extreme, when the sampling rate becomes very fine, the running time grows significantly but the plot shows little improvement because the errors become too small to be reliably measured after the quantization on 8 bits.

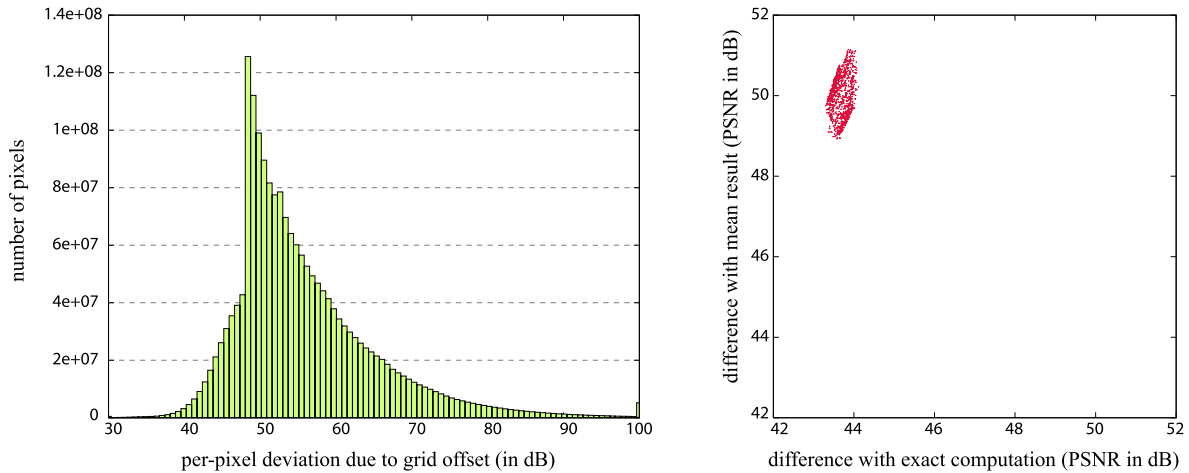


Fig. 8 *Left:* Per-pixel distance to the mean. We filter 1000 times the architectural picture using the 5^3 kernel with a different offset of the downsampling grid. For each pixel \mathbf{p} , we computed the mean value $M_{\mathbf{p}}$. Then for each run, we measured the distance between the result value $I_{\downarrow\uparrow}^b(\mathbf{p})$ and the mean $M(\mathbf{p})$. The histogram shows the distribution of these distances across the 1000 runs. *Right:* Comparison between the distance to the mean and the distance to the exact result. For each

run, we measured the PSNR between the approximated result $I_{\downarrow\uparrow}^b$ and the mean result M . We also measured the PSNR between $I_{\downarrow\uparrow}^b$ and the exact result I^b . The plot shows these two measures for each run. The approximated results are one order of magnitude closer to the mean M . This demonstrates that the position of the downsampling grid a limited influence on the produced results. These tests use the 5^3 -kernel implementation

5.3 Effect of the Downsampling Grid

Our method uses a coarse 3D grid to downsample the data. The position of this grid is arbitrary and can affect the produced results as pointed out by Weiss (2006). To evaluate this effect, we filter 1000 times the same image and applied each time a different random offset to the grid. We average all the outputs to get the mean result. For each grid position, we measured at each pixel the difference between the intensity obtained for this offset and the mean value over all offsets. Figure 8 (left) shows the distribution of these variations caused by the grid offsets. We can see that there are a few large variations (under 40 dB) but that most of them are small (≈ 50 dB and more). These variations have to be compared with the accuracy of the approximation: Fig. 8 (right) shows that they are significantly smaller than the error incurred by the approximation. In other words, the potential variations caused by the position of the downsampling grid are mostly negligible in regards of the error stemming from the approximation itself.

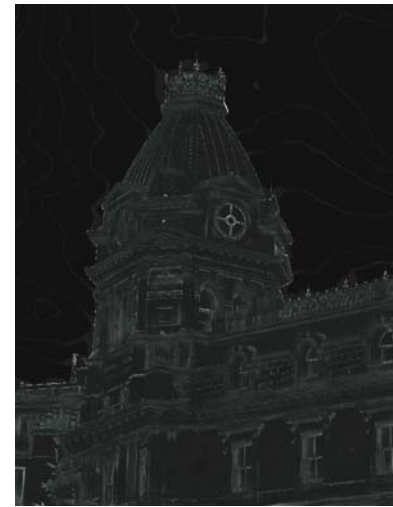
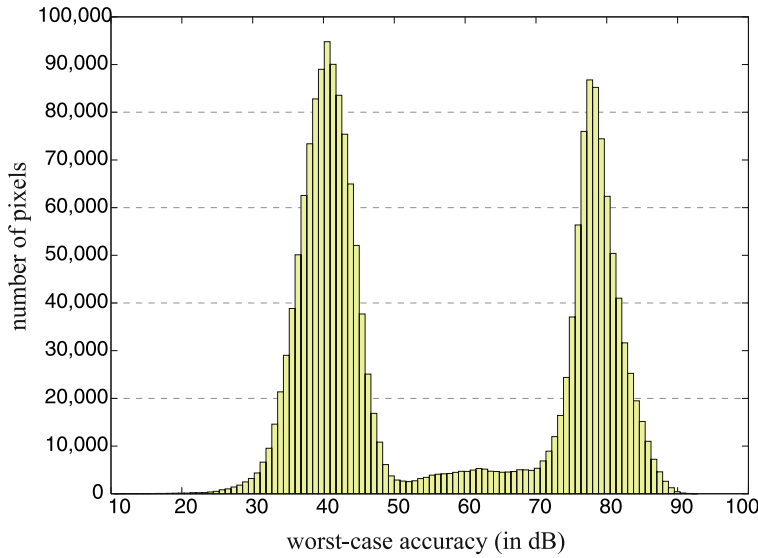
Although these variations have never been a problem in our applications (cf. Sect. 8), they might be an issue in other contexts. To better appreciate the potential defects that may appear, we built a worst-case scenario with the statistics of the 1000 runs. We computed the minimum and maximum results:

$$I_{\min}^b(\mathbf{p}) = \min_{\text{all runs}} I_{\downarrow\uparrow}^b(\mathbf{p}) \quad \text{and} \quad I_{\max}^b(\mathbf{p}) = \max_{\text{all runs}} I_{\downarrow\uparrow}^b(\mathbf{p}). \quad (14)$$

These two images correspond to the worst-case hypotheses where all the negative (resp. positive) deviations happened at the same time which is unlikely. Figure 9 (left) plots the histogram of the differences between I_{\max}^b and I_{\min}^b : Even in this worst-case scenario, most variations remain above 40 dB. Globally, the maximum and minimum results still achieve satisfying accuracies of 42 dB and 43 dB compared to the exact computation. Figure 9 (right) reveals that the errors are larger on discontinuities (edges, corners) and textured regions. As a consequence, the variation distribution (Fig. 9 (left)) is bimodal with a high-accuracy mode corresponding to the smooth sky and a lower-accuracy mode stemming from the building. Nonetheless, if higher accuracy is required, one can refine the downsampling grid at the cost of longer running times as shown earlier (Fig. 5).

6 Comparisons with Other Acceleration Techniques

We now compare our method to other fast approximations of the bilateral filter. To better understand the common points and differences with our approach, we look into the details of these methods. We selected the methods that explicitly aim for bilateral filtering (Durand and Dorsey 2002; Pham and van Vliet 2005; Weiss 2006). We did not test the Gauss-Seidel scheme proposed by Elad (2002) since it requires a large number of iterations whereas practical applications use non-iterative bilateral filtering. As we shall see, our technique performs especially well on large kernels and has the advantage to be general and simple to implement.



worst-case variations

Fig. 9 *Left:* Influence of the grid position in a worst-case scenario. The histogram shows the distances between the maximum and minimum of all the runs (14). We used the architectural picture. The two modes of the distribution correspond to the building area that contains many contrasted features likely to be affected by the grid position, and

to the smooth sky region with almost no such features. *Right:* Variation map. The *thin features* are the most affected by the downsampling grid position whereas the smooth areas are mostly unaltered. The variation amplitudes are mapped to gray levels after a 10× multiplication

6.1 Comparison with the Piecewise-Linear Approximation

Durand and Dorsey (2002) describe a piecewise-linear approximation of the bilateral filter that is closely related to our approach. Our framework provides better understanding their method: As we will see, the main difference with our technique lies in the downsampling strategy.

Using evenly spaced intensity values $\zeta_1.. \zeta_n$ that cover \mathcal{R} , the piecewise-linear scheme can be summarized as (for convenience, we also name G_{σ_s} the 2D Gaussian kernel):

$$\iota_{\downarrow} = \text{downsample}(I) \quad [\text{image downsampling}], \quad (15a)$$

$$\forall k \in \{1..n\}, \quad \omega_{\downarrow k}(\mathbf{p}) = G_{\sigma_r}(|\iota_{\downarrow}(\mathbf{p}) - \zeta_k|) \quad [\text{range weight evaluation}], \quad (15b)$$

$$\forall k \in \{1..n\}, \quad \omega_{\downarrow k}(\mathbf{p}) = \omega_{\downarrow k}(\mathbf{p})\iota_{\downarrow}(\mathbf{p}) \quad [\text{intensity multiplication}], \quad (15c)$$

$$\forall k \in \{1..n\}, \quad (\omega_{\downarrow k}^b, \omega_{\downarrow k}^b) = G_{\sigma_s} \otimes_{\mathcal{S}} (\omega_{\downarrow k}, \omega_{\downarrow k}) \quad [\text{spatial convolution on } \mathcal{S}], \quad (15d)$$

$$\forall k \in \{1..n\}, \quad \iota_{\downarrow k}^b = \omega_{\downarrow k}^b / \omega_{\downarrow k}^b \quad [\text{normalization}], \quad (15e)$$

$$\forall k \in \{1..n\}, \quad \iota_{\downarrow \uparrow k}^b = \text{upsample}(\iota_{\downarrow k}^b) \quad [\text{layer upsampling}], \quad (15f)$$

$$I_{pl}^b(\mathbf{p}) = \text{interpolation}(\iota_{\downarrow \uparrow k}^b)(\mathbf{p}) \quad [\text{linear layer interpolation}]. \quad (15g)$$

Without downsampling (*i.e.* $\{\zeta_k\} = \mathcal{R}$ and Steps 15a–f ignored), the piecewise-linear scheme is equivalent to ours because Steps 15b–d correspond to a convolution on $\mathcal{S} \times \mathcal{R}$. Indeed, $\iota_{\downarrow}(\mathbf{p}) = I_{\mathbf{p}}$, Steps 15b and 15c can be rewritten in a vectorial form, and the value of the Gaussian can be expressed as a convolution on \mathcal{R} with a Kronecker symbol:

$$\begin{aligned} [\text{Step 15c}] \quad \begin{pmatrix} \omega_{\downarrow k}(\mathbf{p}) \\ \omega_{\downarrow k}(\mathbf{p}) \end{pmatrix} &= \begin{pmatrix} I_{\mathbf{p}} \\ 1 \end{pmatrix} G_{\sigma_r}(|I_{\mathbf{p}} - \zeta_k|) \\ [\text{Step 15b}] \quad &= \begin{pmatrix} I_{\mathbf{p}} \\ 1 \end{pmatrix} [\delta_{I_{\mathbf{p}}} \otimes_{\mathcal{R}} G_{\sigma_r}](\zeta_k) \quad (16a) \end{aligned}$$

$$\text{with } \delta_{I_{\mathbf{p}}}(\zeta \in \mathcal{R}) = \delta(I_{\mathbf{p}} - \zeta). \quad (16b)$$

With Step 15d, the convolution on \mathcal{S} , these three steps perform a 3D convolution using a separation between \mathcal{R} and \mathcal{S} .

The main differences comes from the downsampling approach, where Durand and Dorsey downsample in 2D while we downsample in 3D. They also interleave linear and non-linear operations differently from us: Their division is done after the convolution 15d but before the upsampling 15f. There is no simple theoretical ground to estimate the error. More importantly, the piecewise-linear strategy is such that the intensity ι and the weight ω are functions defined on \mathcal{S} only. A given spatial pixel in the downsampled image has only one intensity and one weight. After downsampling, both sides of a discontinuity may be represented by the same

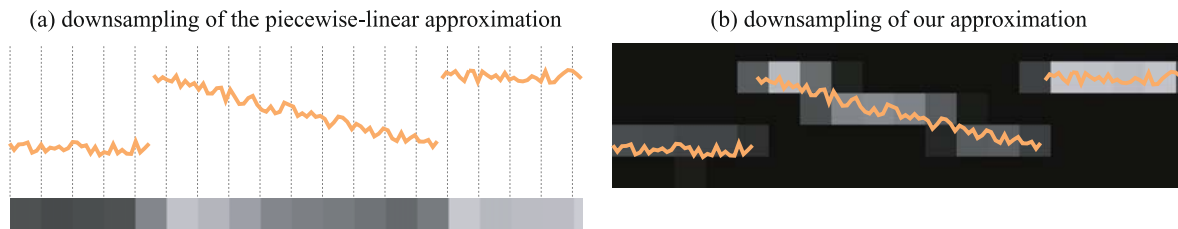


Fig. 10 (Color online) **a** The piecewise-linear approximation is such that only a single value can represent at each position. After downsampling discontinuities are represented by only one intensity value which poorly approximates the discontinuity. **b** With our scheme, the

discontinuities are represented by two distinct values in the downsampled $\mathcal{S} \times \mathcal{R}$ domain, even after downsampling. The original function (in red) is the same as in Fig. 1. The corresponding downsampled representation of the intensity is shown below

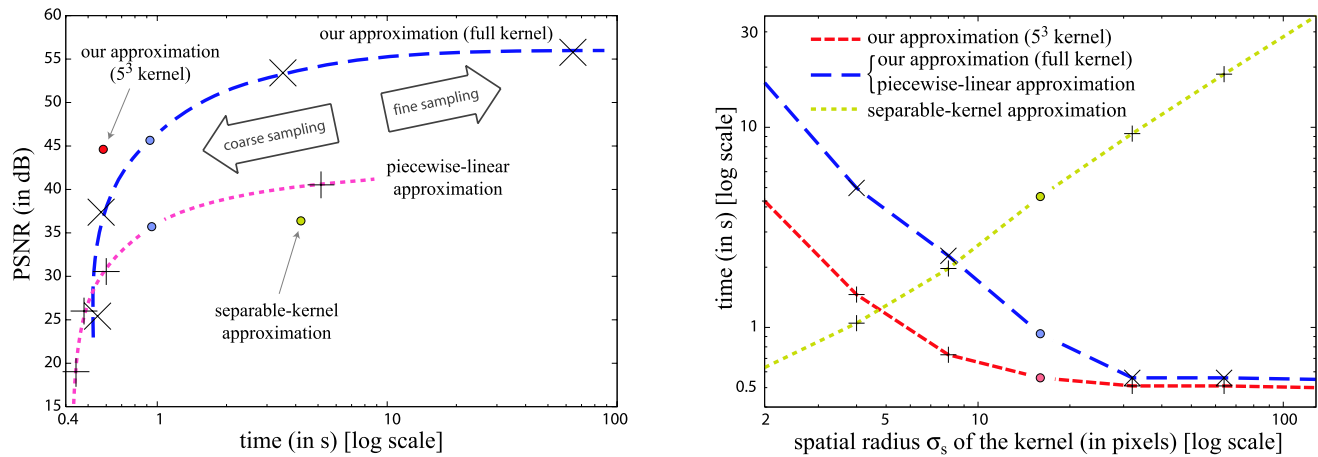


Fig. 11 (Color online) *Left*: Accuracy-versus-time comparison. The methods are tested on the architectural picture (1600×1200) using $\sigma_s = 16$ and $\sigma_r = 0.1$ as parameters. The piecewise-linear approximation and our method using an untruncated kernel can achieve variable degrees of accuracy and speed-up by varying the sampling rates of $\mathcal{S} \times \mathcal{R}$. We tested them with the same sampling rates of $\mathcal{S} \times \mathcal{R}$ (from left to right): (4, 0.025), (8, 0.05), (16, 0.1), (32, 0.2), (64, 0.4).

Right: Time-versus-radius comparison. The methods are tested on the architectural image using $\sigma_r = 0.1$. Our method using the full kernel uses a sampling rate equal to (σ_s, σ_r) . Its curve and the piecewise-linear curve are identical because we allow the same computation budget to both methods. Remember that an exact computation lasts at least several tens of minutes (varying with the implementation). The *color spots* in the left and right plots correspond to the same measurements

values of ι and ω . This is a poor representation of discontinuities since they inherently involve several values. In comparison, we define functions on $\mathcal{S} \times \mathcal{R}$. For a given image point in \mathcal{S} , we can handle several values on the \mathcal{R} domain. The advantage of working in $\mathcal{S} \times \mathcal{R}$ is that this characteristic is not altered by downsampling (cf. Fig. 10). It is the major reason why our scheme is more accurate than the piecewise-linear technique, especially around discontinuities.

Numerical and Visual Comparisons We have implemented the piecewise-linear technique with the same code base as our technique. Figure 11 shows the precision and running time achieved by both approaches. Both techniques exhibit similar profiles but our approach achieves significantly higher accuracies for the same running times (except for extreme downsampling, but these results are not satisfying as shown on Fig. 7). Figure 12 confirms visually this measure: Our method approximates better the exact result. As a consequence, our technique can advantageously replace the

piecewise-linear approximation since it is both simpler and more precise.

6.2 Comparison with the Separable-Kernel Approximation

Pham and van Vliet (2005) approximate the 2D bilateral filter with a separable kernel by first filtering the image rows and then the columns. This dramatically shortens the processing times since the number of pixels to process is proportional to the kernel radius instead of its area in the 2D case. Pham (2006) also describes an extension to his work where the separation axes are oriented to follow the image gradient. However, we did not learn of this extension in time to include it in our tests. Our experiments focus on the axis-aligned version that has been shown to be extremely efficient by Winnemöller et al. (2006). We have implemented this technique in C++ using the same code base as for our method. Our C++ code on 2.8 GHz Intel Xeon achieves a

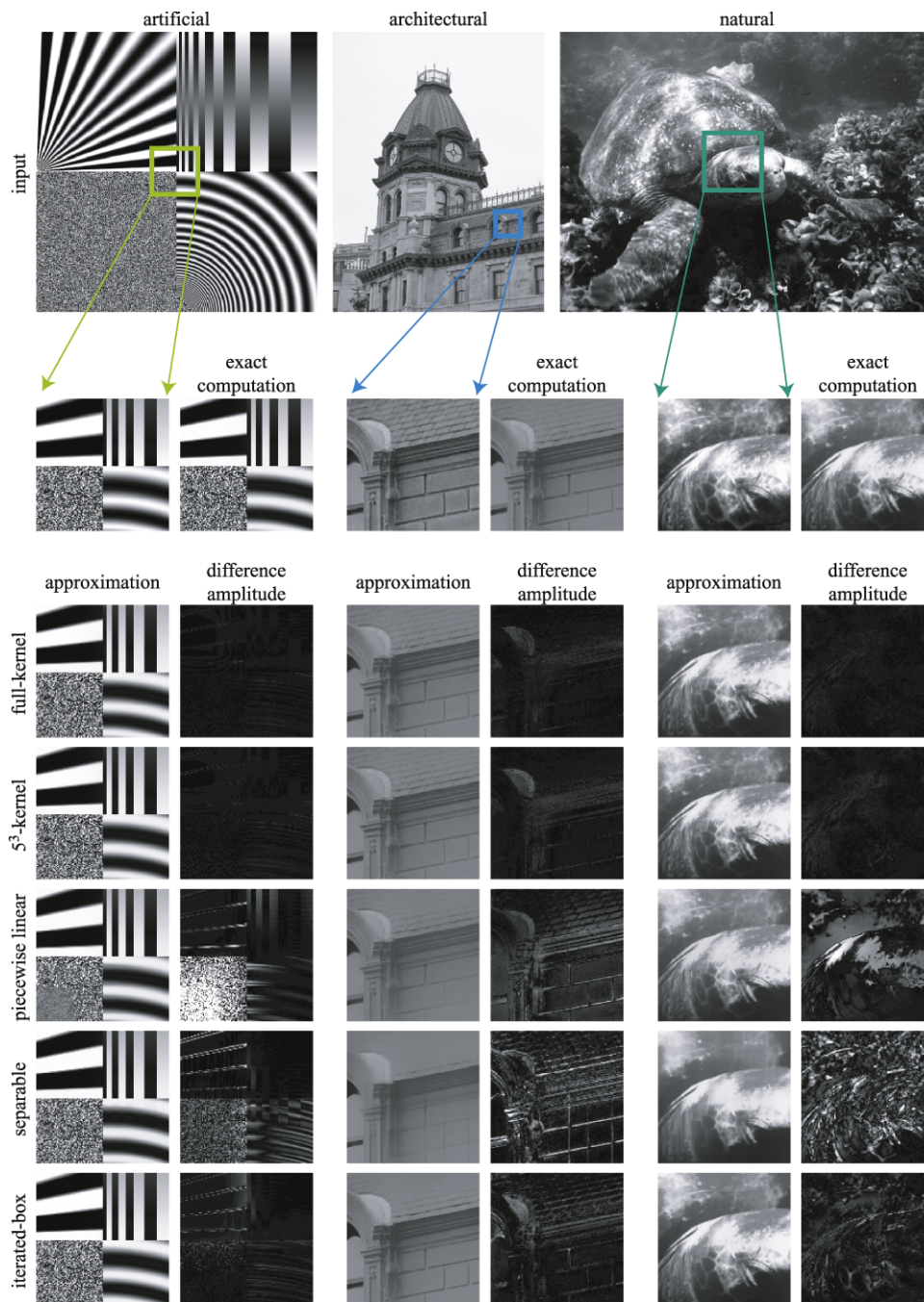


Fig. 12 We have tested our approximated scheme on three images (*first row*): an artificial image (512×512) with different types of edges and a white noise region, an architectural picture (1600×1200) with strong and oriented features, and a natural photograph (800×600) with more stochastic textures. For clarity, we present representative close-ups. Full resolution images are provided as supplemental material. Our approximations produces results visually similar to the exact computation. Although we applied a $10\times$ multiplier to the difference amplitudes, the image differences do not show significant errors. In comparison, the piecewise-linear approximation introduces large visual discrepancies: small yet contrasted features are washed out. Furthermore this defect depends the neighborhood. For instance, on the artificial image, the noise region is less smoothed near the borders because of

the adjacent white and black bands. Similar artifacts are visible on the brick shadows of the architectural image. The separable kernel introduces axis-aligned streaks which are visible on the bricks of the architectural picture. In comparison, the iterated-box kernel does not incur artifacts but the filter does not actually approximate a Gaussian kernel (Sect. 6.3). This results in large differences. All the filters are computed for $\sigma_s = 16$ and $\sigma_r = 0.1$. Our filter uses a sampling rate of $(16, 0.1)$. The piecewise-linear filter is allocated the same time budget as our full-kernel approximation: its sampling rate is chosen in order to achieve the same (or slightly superior) running time. The iterated-box method is set with 2 additional iterations as described by Weiss (2006) and $\sigma'_r = \frac{0.1}{\sqrt{3}} \approx 0.0577$ (Sect. 6.3)

5× speed-up compared to the experiment in Matlab on an 1.47 GHz AMD reported in the original article.

As in the original paper, our tests confirm that the separable-kernel strategy is useful for small kernels. Figure 11 (right) shows that this advantage applies to radii that do not exceed a few pixels ($\sigma_s \approx 5$ pixels with our implementation). Winnemöller et al. (2006) have shown that it is a suitable choice to simplify the content of videos by applying several iterations of the bilateral filter using a small kernel. This approximation is less suitable for larger kernels as shown by the running times on Fig. 11 (right) and because artifacts due to the approximation become visible (Figs. 12 and 15f). As the kernel becomes bigger, the pixel neighborhoods contain more and more complex features, e.g. several edges and corners. These features are poorly handled by the separable kernel because it considers the rows and columns separately. This results in axis-aligned artifacts which may be undesirable in a number of cases. In comparison, our approach handles the complete neighborhoods. The incurred error is significantly lower and follows the image structure, yielding mostly imperceptible differences (Figs. 12 and 15g).

6.3 Comparison with the Iterated-Box Method

Weiss (2006) describes an algorithm that computes efficiently the exact bilateral filter using a square box function to define the spatial influence. It is based on the maintenance of local intensity histograms, and Sect. 9.5 relates this concept with our higher-dimensional interpretation. In this section, we discuss performances and accuracy. The complexity of Weiss's algorithm evolves as the logarithm of the kernel radius, thereby achieving running times in the order of a few seconds even for large kernels (100 pixels and above for an 8 megapixel image).

However, the box function is known to introduce visual artifacts because its Fourier transform is not band-limited and it introduces Mach bands. To address this point, Weiss iterates the filter three times while keeping the range weights constant, i.e. the range values are always computed according to the input image. Weiss motivates this approach by the fact that, on regions with smooth intensity, the range weights have little influence and the iteration corresponds to spatially convolving a box function β^0 with itself, i.e. $\beta^0 \otimes_s \beta^0 \otimes_s \beta^0$. This results in a smooth quadratic B-spline β^2 that approximates a Gaussian function. This interpretation does not hold on non-uniform areas because the spatial component interacts with the range function. van de Weijer and van den Boomgaard (2001) have shown that this iteration scheme that keeps the weights constant eventually leads to the nearest local maximum of the local histogram. A downside of this sophisticated scheme is that the algorithm is more difficult to adapt. For instance, extension to color images seems nontrivial.

6.3.1 Running Times

A performance comparison with our method is difficult because Weiss's implementation relies on the vector instructions of the CPU. Implementing this technique as efficiently as the original article could not be done in a reasonable amount of time. And since the demonstration plug-in¹ does not run on our platform, we cannot make side-by-side timing comparisons. Nonetheless, we performed some tests on a 8 megapixel image as Weiss uses in his article. For a small kernel ($\sigma_s = 2$ pixels), our software runs in 18 s whereas the iterated-box technique runs in ≈ 2 s. Although the CPU and the implementation are not the same, the difference is significant enough to conclude that the iterated-box method is faster on small kernels. We also experimented with bigger kernels (recall that our running times decrease): our running times are under 3 s for radii over 10 pixels and stabilize at 2.5 s for larger radii ($\sigma_s \geq 30$ pixels) because of the fixed time required by downsampling, upsampling, and nonlinearities. These performances are in the same order as the iterated-box technique. In such situations, the running times depend more on the implementation than on the actual algorithm. Our code can be clearly improved to exploit the vector and multi-core capacities of modern architectures. Weiss also describes avenues to speed up his implementation. We believe that both techniques are approximately equivalent in terms of performance for intermediate and large kernels. The major difference is that our running times decrease with big kernels whereas the ones of the iterated-box slowly increase.

6.3.2 Numerical and Visual Evaluation

Comparing the accuracy is not straightforward either since the iterated version of this filter approximates a Gaussian profile only on constant regions while the kernel has not been studied in other areas. To better understand this point, we conducted a series of tests to characterize the produced results. Weiss (2006) points out that convolving a box function n times with itself produces a B-spline β^n , and as n increases, β^n approximates a Gaussian profile. The iterated-box method is however more complex since at each iteration, the spatial box function β^0 is multiplied by the Gaussian range weight G_{σ_r} and then normalized. This motivated us to investigate the link between the iterated-box scheme and the non-iterative bilateral filter with a B-spline as spatial weight.

Test Setup We have coded a brute-force implementation of the iterated-box filter using classical “for loops”. Although our piece of software is slow, it uses exactly the same space

¹<http://www.shellandslate.com/fastmedian.html>

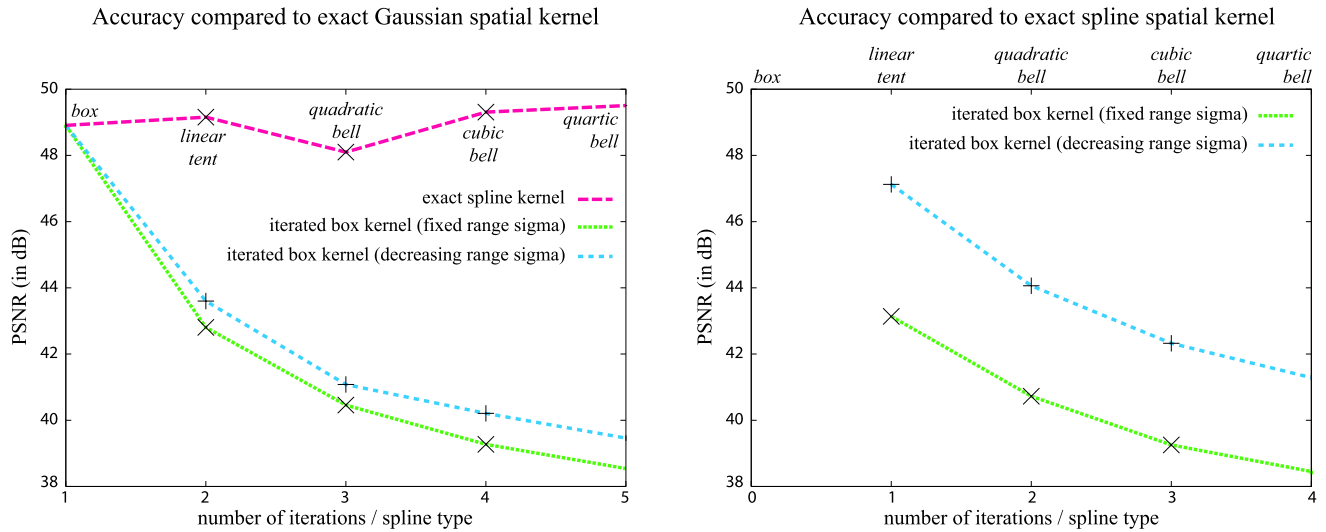


Fig. 13 Comparisons with exact bilateral filtering using a Gaussian spatial kernel (*left*) and using a B-spline spatial kernel (*right*). As the number of iterations increases, the iterated-box filter becomes more different from the non-iterated bilateral filters, for both spatial profiles (Gaussian and spline) and both range settings. On the *left*, the spline

profile converges toward a Gaussian kernel with higher spline order although numerical precision issues limit the convergence speed. The spline shape is indicated on the curve. On the *right*, the iterated-box kernel does not converge towards a spline kernel. The shape of the spline used for the comparison is indicated on the *top axis*

and range functions as the iterated-box method. Thus we can perform visual and numerical comparisons safely as long as we do not consider the running times.

We computed the iterated-box filter with increasing numbers of iterations and compared the results with the exact computation of the bilateral filter using a Gaussian as spatial weight as well as with the exact computation of the bilateral filter using a B-spline β^n with n matching the number of iterations. In both cases, the range weight remains a Gaussian as described by Weiss. We set the radius r_{box} of the box function so as to achieve the same standard deviation σ_s for the all the tests. For a square box function (*i.e.* no iteration), a classical variance computation gives $\sigma_s^2 = \frac{2}{3}r_{\text{box}}^2$. For n iterations, we use $r_{\text{box}} = \sigma_s \sqrt{\frac{3}{2(n+1)}}$.

Setting the range influence is equally important. A first solution is to use the same value σ_r independently of the number of iterations. As a consequence, the smoothing strength varies with the number of iterations (Fig. 14). Our interpretation of the bilateral filter shows that the range domain also undergoes a convolution. This suggests to set $\sigma'_r = \sigma_r \sqrt{\frac{1}{n+1}}$ when the filter is applied a total of $n + 1$ times. It would result in a range Gaussian of size σ_r if the filter was purely a convolution. Let's use our framework to examine the effects of the iterations with fixed weights. Because the weights are kept constant, the same 3D kernel and normalization factors are applied at each iteration. However, the normalization factors vary spatially and slicing occurs after each pass. This corresponds to setting to 0 all the $\mathcal{S} \times \mathcal{R}$ points (\mathbf{x}, ζ) such that $\zeta \neq I_{\mathbf{x}}$. Hence, the process

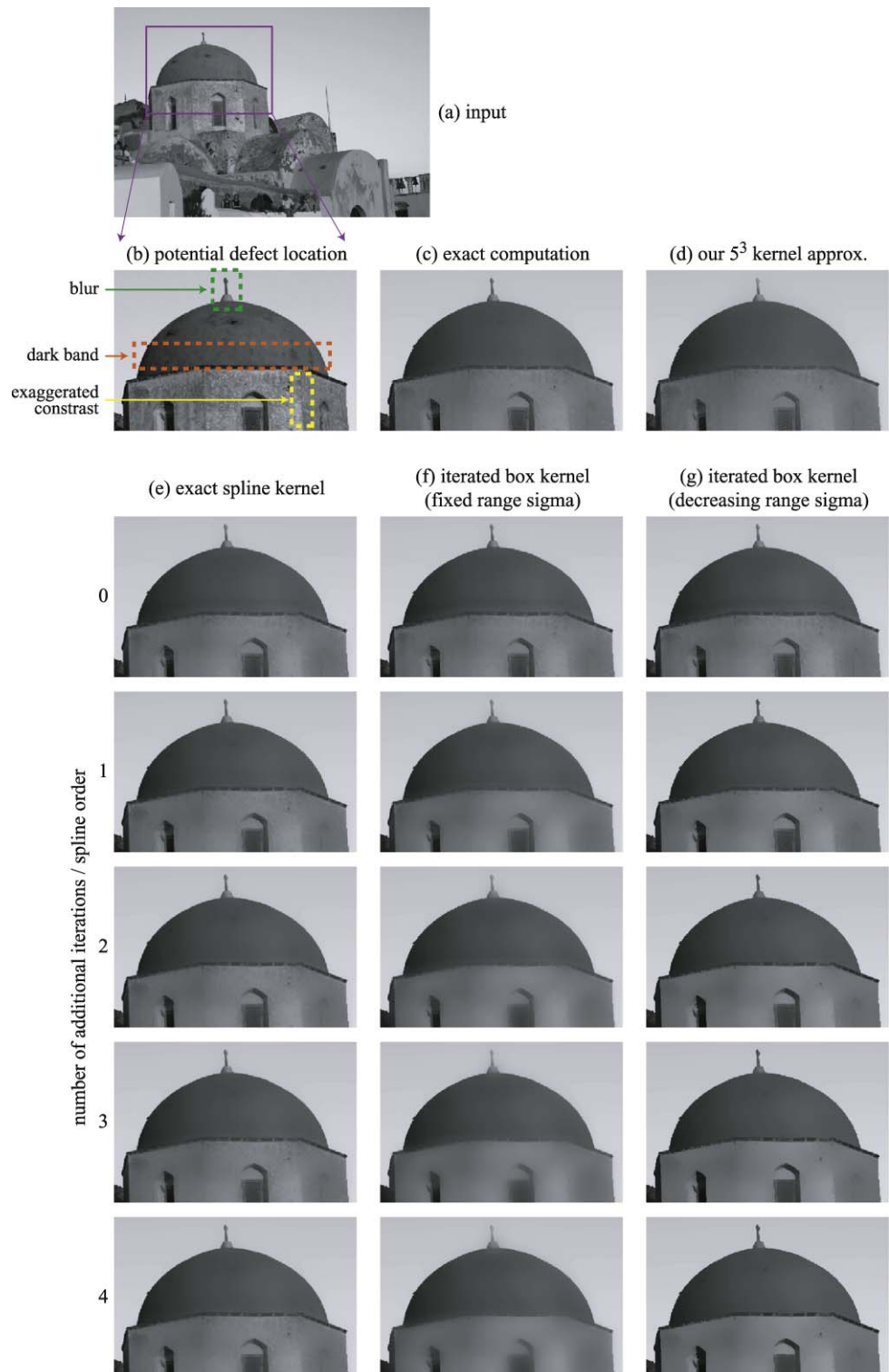
is not purely a succession of convolutions. Nevertheless, we shall see that adapting σ_r to the total number of iterations has some visual advantages over the fixed value.

Note that Weiss experimented with adapting both the space and range sigmas to the total number of iterations although his article did not focus on approximating any kernel and showed only results from three iterations.

Numerical Evaluation Figure 13 shows that the iterated-box scheme does not converge either toward Gaussian bilateral filtering or toward spline bilateral filtering, at least with the parameters that we used. It might not be a problem depending on the application, since it does not affect the visual quality of the result. For instance, this scheme is well-suited for applications such as denoising or smoothing (Weiss 2006). A formal characterization of the achieved kernel would nevertheless be a valuable contribution since it would enable a better comparison with the other existing methods and allows for motivated choices.

Visual Results Figure 14 illustrates the typical defect appearing with a square box function as spatial weight and a single iteration: A faint band appears above the dark cornice. It is worth mentioning that we tested several images and this defect was less common than we expected. As soon as the contrast is large enough or when the picture is textured, no defect is visible. As an example, no banding is visible around the dome contour and near the windows. However, when this defect appears, it impacts the overall image quality. The solution is to apply a smoother kernel, either by

Fig. 14 Results from the spline and iterated-box bilateral filters. Several defects may appear depending on the chosen scenario; their location is shown on the close-up (b). Using a square box function without iteration produces a faint band above the cornice (first row (e, f, g)). Iterating the filter makes this band disappear. If the range Gaussian is not adapted according to the number of iterations, the results become blurry (see the cross and the cornice in (f)). Decreasing σ_r as the square root of the number of iterations prevents blurriness and produces more stable results (g). But iterating the filter can introduce exaggerated contrast variations (see the edge of the wall). The contrast of the close-ups has been increased for clarity purpose. Original images are available in supplemental material



using a higher-order spline, or by iterating the filter. Our experiment shows that using a spline of order 1 (a linear tent) solves the problem while a single additional iteration only attenuates the band without completely removing it. As suggested by Weiss (2006), three iterations yield results without

visible artifacts. Furthermore, this figure shows that the results are progressively washed out if the range setting is not adapted according to the number of iterations. Adapting σ_r as the square root of the number of iterations prevents this effect and produces more stable results. The downside of it-

erating can be seen on the wall corner at the bottom right end of the close-up: Stacking up the nonlinearities results in an exaggerated contrast variation because the top part of the corner has been smoothed at each iteration while the bottom part has been preserved. The non-iterated versions of the filter (either with a spline or a Gaussian kernel) do not present such contrast variations. This effect may be a concern depending on the application.

7 Extensions

An advantage of our approach is that it can be straightforwardly adapted to handle color images and cross bilateral filtering. In the following sections, we describe the details of these extensions.

7.1 Color Images

For color images, the range domain \mathcal{R} is typically a three-dimensional color space such as RGB or CIE-Lab. We name $\mathbf{C} = (C_1, C_2, C_3)$ the vector in \mathcal{R} describing the color of a pixel. The bilateral filter is then defined as:

$$\mathbf{C}_p^b = \frac{1}{W_p^b} \sum_{\mathbf{q} \in \mathcal{S}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) G_{\sigma_r}(\|\mathbf{C}_p - \mathbf{C}_q\|) \mathbf{C}_q \quad (17a)$$

$$\text{with } W_p^b = \sum_{\mathbf{q} \in \mathcal{S}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) G_{\sigma_r}(\|\mathbf{C}_p - \mathbf{C}_q\|). \quad (17b)$$

The definition is similar to (1), except that we handle color vectors \mathbf{C} instead of scalar intensities I . We derive formulae that isolate the nonlinearities similarly to (11) by propagating this change through all the equations. We give the main steps in the following paragraphs and the details are identical to the gray-level case.

Akin to the homogeneous intensity defined in Sect. 3.1, we name *homogeneous color* the 4D homogeneous vector $(W\mathbf{C}, W)$. This let us write (17) in a vector form:

$$\begin{pmatrix} W_p^b \mathbf{C}_p^b \\ W_p^b \end{pmatrix} = \sum_{\mathbf{q} \in \mathcal{S}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) G_{\sigma_r}(\|\mathbf{C}_p - \mathbf{C}_q\|) \times \begin{pmatrix} W_q \mathbf{C}_q \\ W_q \end{pmatrix}. \quad (18)$$

Then, we define the functions \mathbf{c} and w on the joint domain $\mathcal{S} \times \mathcal{R}$:

$$\mathbf{c}: (\mathbf{x} \in \mathcal{S}, \zeta \in \mathcal{R}) \mapsto \mathbf{C}_x, \quad (19a)$$

$$w: (\mathbf{x} \in \mathcal{S}, \zeta \in \mathcal{R}) \mapsto \delta(\|\zeta - \mathbf{C}_x\|) W_x. \quad (19b)$$

We obtain formulae similar to (11):

$$\text{linear: } (w^b \mathbf{c}^b, w^b) = g_{\sigma_s, \sigma_r} \otimes (w\mathbf{c}, w), \quad (20a)$$

$$\text{nonlinear: } \mathbf{C}_p^b = \frac{w^b(\mathbf{p}, \mathbf{C}_p) \mathbf{c}^b(\mathbf{p}, \mathbf{C}_p)}{w^b(\mathbf{p}, \mathbf{C}_p)}. \quad (20b)$$

7.1.1 Discussion

Dimensionality For color images, the $\mathcal{S} \times \mathcal{R}$ domain is 5D: 2 dimensions for the spatial position and 3 dimensions to describe the color. Each point of $\mathcal{S} \times \mathcal{R}$ is a 4D vector since \mathbf{C} is 3D, thus the homogeneous color $(W\mathbf{C}, W)$ is four-dimensional. This dimensionality increase becomes a problem when the $\mathcal{S} \times \mathcal{R}$ domain is finely sampled because of the memory required to store the data. Nonetheless, in many cases, the computation is still tractable. The validation section details this aspect.

Pre-Multiplied Alpha Transparency is classically represented by an additional image channel named *alpha* (Porter and Duff 1984; Blinn 1996; Willis 2006). Pixel values are then 4D vectors (\mathbf{C}, α) and it is well-known that most filtering operations such as blurring must be done on the *pre-multiplied alphas* $(\alpha\mathbf{C}, \alpha)$ in order to obtain correct color and transparency values. Intuitively, transparent pixels contribute less to the output. Although this representation “looks” similar to the homogeneous colors, there are a few important differences. An α value has a physical meaning: the transparency of a pixel. As a consequence, it lies in the $[0; 1]$ interval. In contrast, a W value carries statistical information since it corresponds to a pixel weight. It is non-negative and has no upper bound. A major difference between both representations is that W values can be scaled uniformly across the image without changing their meaning. This is a known property of homogeneous quantities: they are defined up to a scale factor. On the other side, a global scaling on α values alters the result since objects become more or less transparent.

7.1.2 Validation

We tested several alternative strategies on a color image (Fig. 15). First, we filtered an RGB image as three independent channels and compared the output to the result of filtering colors as vectors. We found that vector filtering yields better results but color inconsistencies remained. Thus, we experimented with the CIE-Lab color space which is known to be perceptually meaningful.

Per-channel Filtering versus Vector Filtering Processing the RGB channels independently induces inconsistencies when an edge has different amplitudes depending on the channel. In that case, an edge may be smoothed in a channel while it is preserved in another, inducing incoherent results between channels. Figure 15d shows our truncated-kernel approximation on an example where the blue sky

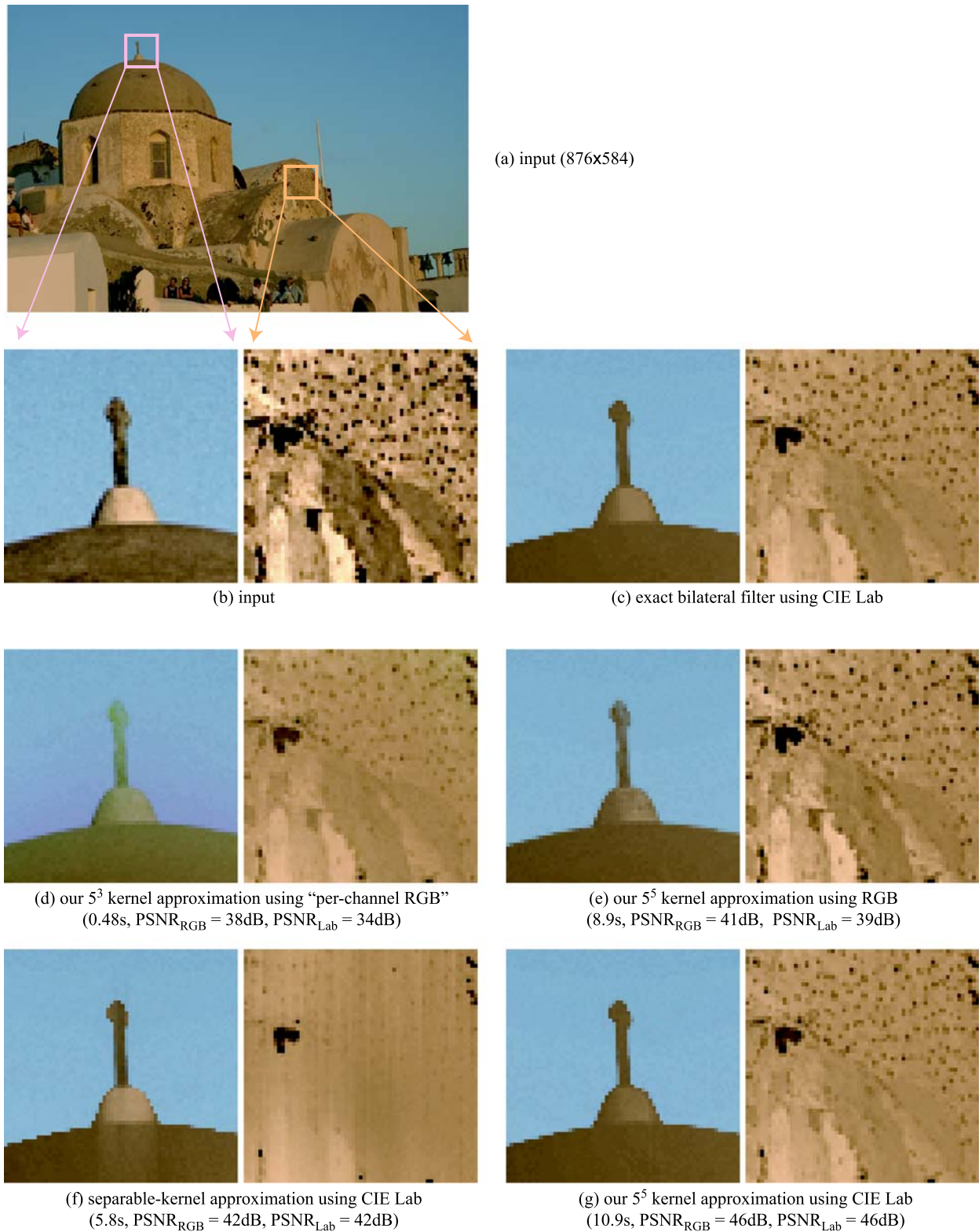
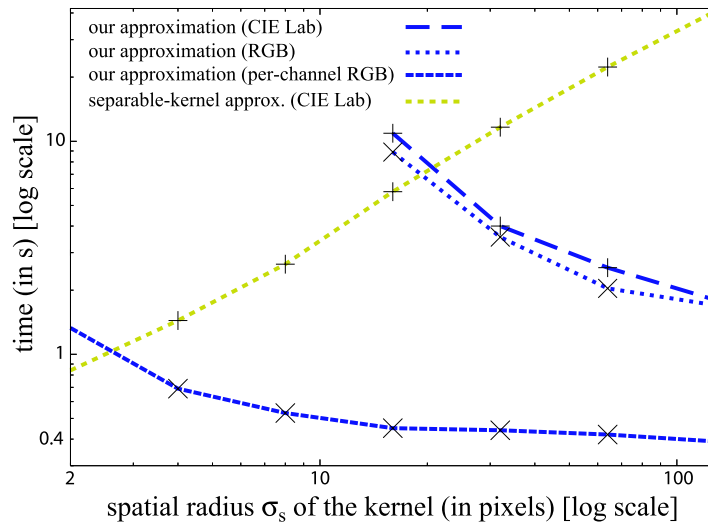


Fig. 15 (Color online) Comparison on a color image. We tested various strategies to filter a color image (a, b). Processing the red, green, and blue channels independently results in color bleeding that makes the cross disappear in the sky (d). Dealing with the RGB vector as described in (20) improves this aspect but some bleeding still occurs (e). In contrast, working in the CIE-Lab space achieves satisfying results (c, g). Comparing our method (g) to the separable-kernel tech-

nique (f) shows our technique is slower but produces a result closer to the exact computation (c). Especially, the separable kernel incurs axis-aligned streaks (f) that may be undesirable in a number of applications. These remarks are confirmed by the numerical precision evaluated with the PSNR computed in the RGB and CIE-Lab color spaces. The contrast of the close-ups has been increased for clarity purposes. Original images are available in supplemental material

Fig. 16 Time-versus-radius comparison. We processed the image shown in Fig. 15 with $\sigma_r = 0.1$ in RGB color space, and $\sigma_r = 10$ in the CIE-Lab color space. Both settings correspond to 10% of the intensity range. The plots show the running times of various options depending on the spatial extent of the kernel. Our approximation is not able to process small kernels because of memory limitations (see the text for details). We used a truncated kernel for our method



“leaks” over the brown cross. Considering the RGB channels altogether using (20) significantly reduces this bleeding defect (Fig. 15e). The downside is the longer computation times required to process the 5D space (cf. Sect. 7.1.1 and Fig. 16). It precluded our algorithm from handling small kernels because of the fine sampling required. In that case, the memory usage was over 1 GB for an image of 0.5 megapixel with sampling rates $(s_s, s_r) = (8, 0.1)$. We did not measure the running time because it was perturbed by disk swapping. This limitation makes our technique not suitable for small kernels on color images. For these cases, one should prefer per-channel filtering and since the kernel is small, the iterated-box method (Weiss 2006) seems an appropriate choice. However, even the RGB-vector filtering result exhibits some color leaking (observe the cross in Fig. 15e). This motivated us to experiment with the CIE-Lab color space which is known to provide better color management (Margulis 2005).

RGB versus CIE-Lab The CIE-Lab space solves the color-bleeding problem. Visual and numerical comparisons show that our technique produces results close to the exact computation (Fig. 15c, g). In contrast, the separable-kernel technique (Pham and van Vliet 2005) applied in the CIE-Lab space yields results degraded by axis-aligned streaks (Fig. 15f). The iterated-box method is described only for per-channel processing, and the extension to vector filtering is unclear.

7.2 Cross Bilateral Filter

Another advantage of our approach is that it can be extended to cross bilateral filtering. The cross bilateral filter has been simultaneously discovered by Eisemann and Durand (2004) and Petschnigg et al. (2004) (who named it “joint bilateral

filter”). It smoothes an image I while respecting the edges of another image E :

$$I_{\mathbf{p}}^c = \frac{1}{W_{\mathbf{p}}^c} \sum_{\mathbf{q} \in \mathcal{S}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) G_{\sigma_r}(|E_{\mathbf{p}} - E_{\mathbf{q}}|) I_{\mathbf{q}} \quad (21a)$$

$$\text{with } W_{\mathbf{p}}^c = \sum_{\mathbf{q} \in \mathcal{S}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) G_{\sigma_r}(|E_{\mathbf{p}} - E_{\mathbf{q}}|). \quad (21b)$$

The only difference with a classical bilateral filter is the range weight G_{σ_r} which uses E instead of I . Applying this change to all the equations results in a new definition for the w function while the function i remains unmodified (6):

$$i: (\mathbf{x} \in \mathcal{S}, \zeta \in \mathcal{R}) \mapsto I_{\mathbf{x}}, \quad (22a)$$

$$w: (\mathbf{x} \in \mathcal{S}, \zeta \in \mathcal{R}) \mapsto \delta(\zeta - E_{\mathbf{x}}) W_{\mathbf{x}}. \quad (22b)$$

These new functions are used to obtain a set of equations that isolate the nonlinearities. Note that the slicing has to be done at the points $(\mathbf{p}, E_{\mathbf{p}})$:

$$\text{linear: } (w^c i^c, w^b) = g_{\sigma_s, \sigma_r} \otimes (wi, w), \quad (23a)$$

$$\text{nonlinear: } I_{\mathbf{p}}^c = \frac{w^c(\mathbf{p}, E_{\mathbf{p}}) i^c(\mathbf{p}, E_{\mathbf{p}})}{w^c(\mathbf{p}, E_{\mathbf{p}})}. \quad (23b)$$

Intuition Following the example of a pen whose ink varies (cf. Sect. 3.3), the function wi can be interpreted as the plot of E using a pen with a brightness equal to I (cf. (22)). The brightness of the plot depends on I and its shape is controlled by E .

Validation Figure 17 illustrates the effects of the bilateral filter on a flash/no-flash pair. The flash picture has a low level of noise but an unpleasant illumination. The no-flash

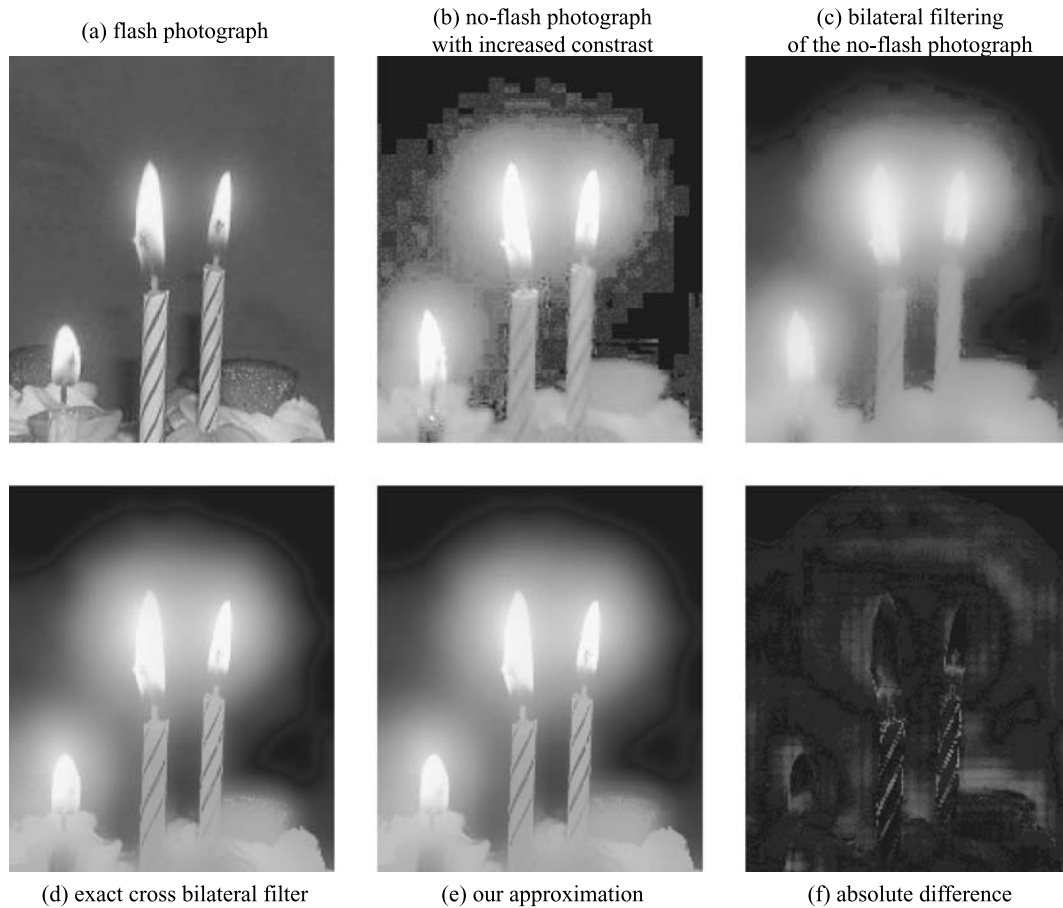


Fig. 17 Example of cross bilateral filtering with a flash/no-flash pair of photographs (**a**, **b**). Directly denoising the no-flash picture is hard because of the high level of noise (**c**). Relying on the flash picture to define the edges to preserve greatly improves the result quality (**d**, **e**). Our approximation scheme produces a result (**e**) visually similar to the exact computation (**d**). The difference image (**f**) reveals subtle devi-

ations (a $10\times$ multiplication has been applied). The PSNR is 41 dB. We used $\sigma_s = 6$ for all the experiments, and $\sigma_r = 0.25$ for direct bilateral filtering and $\sigma_r = 0.1$ for cross bilateral filtering. These values produce the most pleasing results from all our tests. The input images are courtesy of Elmar Eisemann

image has a better illumination but is very noisy. We denoised it with the cross bilateral filter by using the flash image to define the range influence. The achieved result exhibits crisper details than a direct bilateral filter based on the no-flash image alone. Our approximation of the cross bilateral filter achieves the same accuracy and performances as demonstrated for bilateral filtering. We refer the reader to the articles by Petschnigg et al. (2004) and Eisemann and Durand (2004) for more advanced approaches to processing flash/no-flash pairs. Bae et al. (2006) propose another use of the cross bilateral filter to estimate the local amount of texture in a picture. All the results of this article have been computed with our technique.

8 Applications

In this section, we illustrate the capacities of our technique to achieve quality results and meet high standards. We reproduce several of results published previously by the authors.

Tone Mapping We used our technique to manipulate the tone distribution of pictures. First, we implemented a simplified version of the tone-mapping operator described by Durand and Dorsey (2002). Given a high dynamic range image H whose intensity values span a range too broad to be displayed properly on a screen, the goal is to produce a low dynamic range image L that fits the display capacities. The technique can be summarized as follows (see Durand and Dorsey 2002 for details):

1. Compute the logarithmic image $\log(H)$ and apply the bilateral filter to split it into a large-scale component (*a.k.a. base*) $B = bf(\log(H))$ and a small-scale component (*a.k.a. detail*) $D = \log(H) - B$.
2. Compress the base: $B' = \gamma B$ where $\gamma = \text{contrast} / (\max(B) - \min(B))$. ‘contrast’ is set by the user to achieve the desired rendition.
3. Compute the result $L = \exp(B' + D)$.



Fig. 18 Example of tone mapping. The input has a dynamic range too large to be displayed correctly, either the sky is over-exposed or the city is under-exposed. Tone-mapping computes a new image that contains all the details in the sky and in the city and that can viewed on a stan-

dard display. We implemented a simplified version of the tone-mapping operator described by Durand and Dorsey (2002). Our approximation does not incur any artifact in the result. The input image is courtesy of Paul Debevec

Figure 18 shows a sample result that confirms the ability of our technique to achieve high-quality outputs.

Tone Management We also used our method within a more complex pipeline to manipulate the “look” of digital photographs. This technique starts from the same idea as the tone-mapping operator by splitting the image into two layers but then the two layers undergo much stronger transformations in order to adjust both the global contrast and the local amount of the texture of the picture. The bilateral filter is used early in the pipeline to separate the picture into a large-scale layer and a small-scale layer. Cross bilateral filtering is applied later to compute a map that quantifies the local amount of texture at each pixel. Our fast scheme allows this application to be interactive since the result is computed in a few seconds. The details of the algorithm are given in a dedicated paper (Bae et al. 2006). A sample result is shown in Fig. 19. Notice how the result is free from artifacts although an extreme increase of contrast has been applied.

9 Discussion

9.1 Dimensionality

Our approach may seem counterintuitive at first since to speed up the computation, we increase the dimensionality of the problem. Our separation into linear and nonlinear parts comes at the cost of additional range dimensions (one

for gray-level images, three for color images). One has to be careful before increasing the dimensionality of a problem since the incurred performance overhead may exceed the gains, restricting our study to a theoretical discussion. The key of our performances is the possibility to perform an important downsampling of the $\mathcal{S} \times \mathcal{R}$ domain without incurring significant errors. Our tests have shown that the dimensionality is a limiting factor only with color images when using small kernels. In that case, the 5D $\mathcal{S} \times \mathcal{R}$ space is finely sampled and the required amount of memory becomes prohibitive (Sect. 7.1). Adapting the *narrow band* technique (Adalsteinsson and Sethian 1995) to bilateral filtering would certainly help on this aspect. Nevertheless, in all the other scenarios (color images with bigger kernels, and single-channel images with kernel of any size), we have demonstrated that our formalism allows for a computation scheme that is several orders of magnitude faster than a straightforward application of the bilateral filter. This advocates performing the computation in the $\mathcal{S} \times \mathcal{R}$ space instead of the image plane. This strategy is reminiscent of level sets (Osher and Sethian 1988) which alleviate topology management by representing surfaces in a higher-dimensional space. In comparison, we introduce a higher-dimensional image representation that enables dramatic speed-ups through signal downsampling.

Note that using the homogeneous intensities and colors does not increase the dimensionality since (1) and (17) compute the W^b function in addition to I^b or C^b .



Fig. 19 Example of tone management. Our approximation of the bilateral filter and of the cross bilateral filter has been used to enhance digital photographs. See the original article for details (Bae et al. 2006)

9.2 Comparison with *Channel Smoothing*

Felsberg et al. (2006) describe an edge-preserving filter that represents an image as a set of *channels* corresponding to regularly spaced intensity values. For each pixel, the three closest channels are assigned a value corresponding to a second-order B-spline centered on the pixel intensity. Then, the channels are smoothed independently. The output values are computed by reconstructing a B-spline from three channel values. In our bilateral filter framework, this technique can be interpreted as the replacement of the Gaussian by a B-spline to define the range influence. The channels are similar to our downsampling strategy applied only to the range domain except that we use a box function and a Gaussian instead of a B-spline. This suggests that further speed-up can be obtained by downsampling the channels as well, akin to our space downsampling. The strength of this approach is that aspects such as the *influence function* (Huber 1981; Hampel et al. 1986; Black et al. 1998; Durand and Dorsey 2002) can be analytically derived and studied. The downside is the computational complexity introduced to deal with splines. In particular, a relatively complex process is run at each pixel during the final reconstruction to avoid aliasing. In comparison, we cast our approach as an approximation problem and characterize the achieved accuracy with signal processing arguments. Thus, we do not define a new filter and focus on bilateral filtering to improve dramatically its computational efficiency. For instance, we prevent aliasing with a simple linear interpolation and downsample both the range and space domains. Further study of both approaches would be a valuable research contribution.

9.3 Comparison with Image Manifolds

Sochen et al. (1998) describe the *geometric framework* that handles images as manifolds in the $\mathcal{S} \times \mathcal{R}$ space. For in-

stance, a gray-level image I is seen as a 2D surface embedded in a 3D space, *i.e.* $z = I(x, y)$. This representation leads to techniques that define functions on the manifold itself instead of the xy plane, and that can be interpreted as deformations of the image manifold. In this context, the bilateral filter is shown to be related to the *short-time kernel* of the heat equation defined directly on the image manifold: Sochen et al. (2001) demonstrate that bilateral filtering using a small window and blurring using a Gaussian kernel embedded in the image manifold yield results close to bilateral filtering. This interpretation based on small neighborhoods is related to other results linking bilateral filtering to anisotropic diffusion using partial differential equations (Durand and Dorsey 2002; Elad 2002; Buades et al. 2005) since these filters involve only the eight neighbors of a given pixel. In a similar spirit, Barash (2002) uses points in the $\mathcal{S} \times \mathcal{R}$ domain to interpret the bilateral filter. He handles $\mathcal{S} \times \mathcal{R}$ to compute distances and express the difference between adaptive smoothing and bilateral filtering as a difference of distance definitions.

The main difference between these techniques and our interpretation stems from the image representation: in image manifolds, images remain fundamentally two-dimensional—but embedded in a 3D space, while our representation stores values in the whole 3D space. In the geometric framework, each pixel is mapped to a point in $\mathcal{S} \times \mathcal{R}$. Given a point in $\mathcal{S} \times \mathcal{R}$, either it belongs to the manifold and its \mathcal{S} and \mathcal{R} coordinates directly indicate its position and intensity, or it is not on the manifold and is ignored by the algorithm. These methods also use the intrinsic metric of the image manifold. In contrast, we deal with the entire $\mathcal{S} \times \mathcal{R}$ domain, we use its Euclidean metric, we define functions on it, resample it, perform convolutions, and so on. Another major difference is that we define the intensity through a function (i or i^b), and that in general, the intensity of a point (\mathbf{x}, ζ) is *not* its range coordinate ζ , *e.g.* $i^b(\mathbf{x}, \zeta) \neq \zeta$. To our knowledge, this

use of the $\mathcal{S} \times \mathcal{R}$ domain has not been described before and opens new avenues to deal with images. We have shown that it enables the use signal-processing techniques and theory to better compute and understand the bilateral filter. Our approach is complementary to existing frameworks and has potential to inspire new image processing techniques.

9.4 Frequency Content

Our interpretation shows that the bilateral filter output is obtained from smooth low-frequency data which may seem incompatible with the feature-preserving aspect of bilateral filtering. The convolution step of our algorithm indeed smoothes the functions defined on the $\mathcal{S} \times \mathcal{R}$ domain. Nevertheless, the final result exhibits edges and corners because during the slicing nonlinearity, two adjacent pixels, *i.e.* two adjacent points in the spatial domain, can sample this smooth signal at points distant in the range domain, thereby generating discontinuities in the output.

9.5 Comparison with Local-Histogram Approaches

Local histograms are classical intensity histograms where each pixel contributes only a fraction defined by a spatial influence function. Koenderink and van Doorn (1999) have shown that local histograms are a useful tool to study image structure and content. van de Weijer and van den Boomgaard (2001) and Weiss (2006) demonstrated that the result of the bilateral filter at a given pixel is the average intensity of its local histogram with each bin weighted by the range function. This average is normalized by the sum of the weights. Our technique can be interpreted in this framework by remarking that the ζ axis represents histograms. When we first build the $\mathcal{S} \times \mathcal{R}$ domain, the homogeneous coordinate of the w function is a trivial histogram at each pixel \mathbf{p} : $w(\mathbf{p}, \zeta) = 1$ at $\zeta = I_{\mathbf{p}}$ and 0 at $\zeta \neq I_{\mathbf{p}}$. The bins are point-wise, *i.e.* there is a single intensity value per bin. Then when we downsample the \mathcal{R} domain by a factor s_r , the bins become wider and cover an intensity interval of size s_r . At this stage, each pixel has still its own histogram. After downsampling the \mathcal{S} domain with a box function of size s_s , each group of $s_s \times s_s$ pixels shares the same histogram, and in general, several bins are occupied, *i.e.* $w_{\downarrow}(\mathbf{p}, \zeta) > 0$ for several ζ values. Then applying the spatial Gaussian approximates a Gaussian window for the histograms, and the range Gaussian is equivalent to weighting the histogram bins.

Let's have a look at a few specifics of our approach. First, although groups of pixels share the same histogram, there are no blocks in the results since the linear interpolation applied at the end of the algorithm assigns to each pixel its own histogram. In addition, we compute the w_i function that stores the mean intensity of each bin. Using this value, we adapt the intensity associated to each bin to the data it contains instead of using a generic value. For instance, if we



Fig. 20 A simple edge picture. Our approximations produce the exact result. Although we downsample both the space and range domains, there are no spatial blocking artifacts nor intensity aliasing defects. See the text for details

put a single sample into a bin, this bin is represented by the sample value and not by its midpoint, thereby preventing any error. If several samples are stored together, the bin is assigned their mean. Coupled with the linear interpolation on the range domain, this ensures that our technique does not suffer from intensity aliasing. Figure 20 shows a simple example that confirms that our approximated schemes do not introduce blocks (the diagonal edge is preserved) and that there is no intensity aliasing (constant intensity regions are unaltered).

9.6 Complexity

Our algorithm operates on two types of data: the original 2D full-resolution image of size $|\mathcal{S}|$ and the low-resolution 3D representation of size $\frac{|\mathcal{S}|}{s_s^2} \times \frac{|\mathcal{R}|}{s_r}$ where $|\cdot|$ indicates the cardinality of a set, and s_s and s_r are the sampling rates of the space and range domains. The complexity of the method depends on operations performed on both types. The complexity of the convolution is:

- $\mathcal{O}\left(\frac{|\mathcal{S}|}{s_s^2} \frac{|\mathcal{R}|}{s_r} \log\left(\frac{|\mathcal{S}|}{s_s^2} \frac{|\mathcal{R}|}{s_r}\right)\right)$ for the full-kernel option computed with fast Fourier transform and multiplication in the frequency domain.
- $\mathcal{O}\left(\frac{|\mathcal{S}|}{s_s^2} \frac{|\mathcal{R}|}{s_r}\right)$ for the 5^3 -kernel option computed explicitly in the spatial domain. Since in this case, we have $(s_s, s_r) = (\sigma_s, \sigma_r)$, the complexity can be expressed as $\mathcal{O}\left(\frac{|\mathcal{S}|}{\sigma_s^2} \frac{|\mathcal{R}|}{\sigma_r}\right)$.

Downsampling, upsampling, slicing, and dividing are done pixel by pixel and are linear in the image size *i.e.* $\mathcal{O}(|\mathcal{S}|)$. The total algorithm complexity is thus $\mathcal{O}\left(|\mathcal{S}| + \frac{|\mathcal{S}|}{s_s^2} \frac{|\mathcal{R}|}{s_r} \log\left(\frac{|\mathcal{S}|}{s_s^2} \frac{|\mathcal{R}|}{s_r}\right)\right)$ with the full-kernel convolution and $\mathcal{O}\left(|\mathcal{S}| + \frac{|\mathcal{S}|}{s_s^2} \frac{|\mathcal{R}|}{s_r}\right)$ using a truncated kernel. Hence, depending on the sampling rates, the algorithm's order of growth is dominated either by the convolution or by the downsampling, upsampling, and nonlinearities. Table 2 illustrates the impact of the sampling rates on the running times: Convolution takes most of the time with small sampling values whereas it becomes negligible for large values. This complexity also characterizes the effects of large kernels when the sampling rates are equal to the Gaussian sigmas: The

Table 2 Time used by each step at different sampling rates of the architectural image. Upsampling is reported with the nonlinearities because our implementation computes $i_{\downarrow\uparrow}^b$ only at the

$(\mathbf{x}, I_{\mathbf{x}})$ points rather than upsampling the whole $\mathcal{S} \times \mathcal{R}$ space (cf. Sect. 4.2). We use the full-kernel implementation with $\sigma_s = 16$ and $\sigma_r = 0.1$

Sampling (s_s, s_r)	(4, 0.025)	(8, 0.05)	(16, 0.1)	(32, 0.2)	(64, 0.4)
Downsampling	1.3 s	0.23 s	0.09 s	0.07 s	0.06 s
Convolution	63 s	2.8 s	0.38 s	0.02 s	0.01 s
Alicing, upsampling, division	0.48 s	0.47 s	0.46 s	0.47 s	0.46 s

per-pixel operations are not affected but the convolution becomes “cheaper”. This is confirmed in practice on the running times shown in Figs. 16 and 11. Finally, photography applications set the kernel spatial radius as a portion of the image size (Durand and Dorsey 2002; Eisemann and Durand 2004) and the range parameter can also be proportional to the intensity range (Petschnigg et al. 2004; Bae et al. 2006). In the latter case, the ratios $\frac{|\mathcal{S}|}{\sigma_s^2}$ and $\frac{|\mathcal{R}|}{\sigma_r}$ are fixed. This leads to a constant-cost convolution and a global complexity linear in the image size $\mathcal{O}(|\mathcal{S}|)$. For these applications, the running time variations are purely due to the per-pixel operations (downsampling, upsampling, and division).

9.7 Practical Use

Our experiments have highlighted the differences between a number of accelerations of the bilateral filter. In practice, the choice of a method will depend on the specific application, and in particular on the spatial kernel size, the need for accuracy vs. performance, the need for extensions such as cross-bilateral filtering and color filtering, as well as ease of implementation constraints.

Brute-Force Bilateral Filter The brute-force bilateral filter is a practical choice only for small kernels (in the 5×5 range) and when accuracy is paramount. It is trivial to implement but can be extremely slow (minutes per megapixel).

Separable Kernel The separable kernel approximation should be used only with small spatial kernels (no more than 10 pixels). It is simple to implement and can be adapted to graphics hardware. However, processing time and accuracy suffer dramatically with bigger kernels. This technique is thus well adapted to noise-removal tasks that typically use small spatial and range sigmas.

Iterated Box Weiss’s iterated box method is very efficient for kernels of all sizes. It is only an approximation of the Gaussian bilateral filter, but for most applications such as computational photography, this is probably not an issue. The speed of Weiss’s implementation draws heavily from the vector instruction set of modern CPUs and the algorithm implementation might be more complex than the one in this

paper. The main drawbacks of Weiss’s method are the restriction to box spatial weights, to single-channel images, and the difficulty in extending it to cross bilateral filtering.

Our Method Our method is most efficient for large kernels such as the ones used in computational photography to extract a large-scale component of an image, e.g. (Durand and Dorsey 2002; Eisemann and Durand 2004; Petschnigg et al. 2004; Bae et al. 2006). It is very easy to implement and to extend to color images and cross bilateral filtering. The main weakness of our technique is for small kernels where the size of the higher-dimensional representation becomes prohibitive. Our method subsumes Durand and Dorsey’s fast bilateral filter (Durand and Dorsey 2002) and, in our opinion, there is no compelling reason to use that older method since the new version is conceptually simpler and implementation-wise, and has better accuracy-performance tradeoffs.

The above discussion is summarized in Table 3.

10 Conclusions

We have presented a fast approximation of the bilateral filter based on a signal processing interpretation. From a theoretical point of view, we have introduced the notion of homogeneous intensity and demonstrated a new approach of the space-intensity domain: We define intensities through functions that are resampled and convolved in this space whereas existing frameworks use it to represent images as manifolds. Although smooth functions are at the core of our approach, the results exhibit sharp features because of the slicing non-linearity. We believe that these concepts can be applied beyond bilateral filtering, and we hope that these contributions will inspire new studies. From a practical point of view, our approximation technique yields results visually similar to the exact computation with interactive running times. We have demonstrated that this technique enables interactive applications relying on quality image smoothing. Our experiments characterize the strengths and limitations of our technique compared to existing approaches. Our study casts a new light on these other methods, leading for instance to a consistent strategy to set the parameters of the iterated-box

Table 3 Summary of the properties of the various implementations of the bilateral filter

	Complexity	Pros	Cons
Brute force	$ \mathcal{S} ^2$	Accurate	Very slow with kernels larger than 5×5
Separable kernel	$ \mathcal{S} \sigma_s$	Fast with small kernels, easy to implement and adapt	Axis-aligned artefacts
Iterated box	$ \mathcal{S} \log(\sigma_s)$	Always fast	Gray-level images only, hard to implement and adapt
Our method	$ \mathcal{S} + \frac{ \mathcal{S} }{s_s^2} \frac{ \mathcal{R} }{s_r} \log\left(\frac{ \mathcal{S} }{s_s^2} \frac{ \mathcal{R} }{s_r}\right)$ or $ \mathcal{S} + \frac{ \mathcal{S} }{s_s^2} \frac{ \mathcal{R} }{s_r}$	Visually similar to brute force, fast with large kernels, easy to implement and adapt	Slow with small kernels, large memory requirement for color images

technique. It also points out its lack of convergence toward Gaussian and B-spline kernels. We have listed a few guidelines stemming from these tests to select an appropriate bilateral filter implementation depending on the targeted application. Our technique is best at dealing with big kernels. Furthermore, our method is extremely simple to implement and can be easily extended to cross bilateral filtering and color images.

Acknowledgements We are grateful to Ben Weiss, Tuan Pham, and Lucas van Vliet for their remarks and feedback about our experiments, and to Michael Cohen, Todor Georgiev, Pierre Kornprobst, Bruno Lévy, Sing Bing Kang, Richard Szeliski, and Matthew Uyttendaele for their insightful discussions and suggestions to extend our conference paper. We thank Tilke Judd, Paul Green, and Sara Su for their help with the article.

This work was supported by a National Science Foundation CAREER award 0447561 “Transient Signal Processing for Realistic Imagery,” an NSF Grant No. 0429739 “Parametric Analysis and Transfer of Pictorial Style,” a grant from Royal Dutch/Shell Group, and the Oxygen consortium. Frédo Durand acknowledges a Microsoft Research New Faculty Fellowship and a Sloan Fellowship. Sylvain Paris was partially supported by a Lavoisier Fellowship from the French “Ministère des Affaires Étrangères.”

References

- Adalsteinsson, D., & Sethian, J. A. (1995). A fast level set method for propagating interfaces. *Journal of Computational Physics*, 118, 269–277.
- Aurich, V., & Weule, J. (1995). Non-linear Gaussian filters performing edge preserving diffusion. In *Proceedings of the DAGM symposium*.
- Bae, S., Paris, S., & Durand, F. (2006). Two-scale tone management for photographic look. *ACM Transactions on Graphics*, 25(3), 637–645. Proceedings of the ACM SIGGRAPH conference.
- Barash, D. (2002). A fundamental relationship between bilateral filtering, adaptive smoothing and the nonlinear diffusion equation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(6), 844.
- Barash, D., Schlick, T., Israeli, M., & Kimmel, R. (2003). Multiplicative operator splittings in non-linear diffusion: from spatial splitting to multiplicative timesteps. *Journal of Mathematical Imaging and Vision*, 19, 33–48.
- Bennett, E. P., & McMillan, L. (2005). Video enhancement using per-pixel virtual exposures. *ACM Transactions on Graphics*, 24, 845–852. Proceedings of the ACM SIGGRAPH conference.
- Black, M. J., Sapiro, G., Marimont, D. H., & Heeger, D. (1998). Robust anisotropic diffusion. *IEEE Transactions on Image Processing*, 7(3), 421–432.
- Blinn, J. F. (1996). Fun with premultiplied alpha. *IEEE Computer Graphics and Applications*, 16(5), 86–89.
- Buades, A., Coll, B., & Morel, J.-M. (2005). *Neighborhood filters and PDE’s* (Technical Report 2005-04). CMLA.
- Chen, J., Paris, S., & Durand, F. (2007). Real-time edge-aware image processing with the bilateral grid. *ACM Transactions on Graphics* 26(3). Proceedings of the ACM SIGGRAPH conference.
- Choudhury, P., & Tumblin, J. E. (2003). The trilateral filter for high contrast images and meshes. In *Proceedings of the Eurographics symposium on rendering*.
- Durand, F., & Dorsey, J. (2002). Fast bilateral filtering for the display of high-dynamic-range images. *ACM Transactions on Graphics*, 21(3), 257–266. Proceedings of the ACM SIGGRAPH conference.
- Eisemann, E., & Durand, F. (2004). Flash photography enhancement via intrinsic relighting. *ACM Transactions on Graphics*, 23(3), 673–678. Proceedings of the ACM SIGGRAPH conference.
- Elad, M. (2002). On the bilateral filter and ways to improve it. *IEEE Transactions on Image Processing*, 11(10), 1141–1151.
- Elad, M. (2005). Retinex by two bilateral filters. In *Proceedings of the scale-space conference*.
- Felsberg, M., Forssén, P.-E., & Schar, H. (2006). Channel smoothing: Efficient robust smoothing of low-level signal features. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(2), 209–222.
- Fleishman, S., Drori, I., & Cohen-Or, D. (2003). Bilateral mesh denoising. *ACM Transactions on Graphics*, 22(3), 950–953. Proceedings of the ACM SIGGRAPH conference.
- Hampel, F. R., Ronchetti, E. M., Rousseeuw, P. M., & Stahel, W. A. (1986). *Robust statistics—the approach based on influence functions*. New York: Wiley–Interscience. ISBN 0-471-73577-9.
- Huber, P. J. (1981). *Robust statistics. Probability and statistics*. New York: Wiley–Interscience.
- Jones, T. R., Durand, F., & Desbrun, M. (2003). Non-iterative, feature-preserving mesh smoothing. *ACM Transactions on Graphics*, 22(3), 943–949. Proceedings of the ACM SIGGRAPH conference.
- Koenderink, J. J., & van Doorn, A. J. (1999). The structure of locally orderless images. *International Journal of Computer Vision*, 31(2–3), 159–168.
- Liu, C., Freeman, T., Szeliski, R., & Kang, S. (2006). Noise estimation from a single image. In *Proceedings of the computer vision and pattern recognition conference*. New York: IEEE Press.
- Margulis, D. (2005). *Photoshop LAB color: The canyon conundrum and other adventures in the most powerful colorspace*. Berkeley: Peachpit. ISBN: 0321356780.

- Mrázek, P., Weickert, J., & Bruhn, A. (2006). Geometric properties from incomplete data. In *On robust estimation and smoothing with spatial and tonal kernels*. Berlin: Springer.
- Oh, B. M., Chen, M., Dorsey, J., & Durand, F. (2001). Image-based modeling and photo editing. In *Proceedings of the ACM SIGGRAPH conference*. New York: ACM.
- Osher, S., & Sethian, J. A. (1988). Fronts propagating with curvature-dependent speed: Algorithms based on Hamilton-Jacobi formulations. *Journal of Computational Physics*, 79, 12–49.
- Paris, S., & Durand, F. (2006). A fast approximation of the bilateral filter using a signal processing approach. In *Proceedings of the European conference on computer vision*.
- Paris, S., Briceño, H., & Sillion, F. (2004). Capture of hair geometry from multiple images. *ACM Transactions on Graphics*, 23(3), 712–719. Proceedings of the ACM SIGGRAPH conference.
- Petschnigg, G., Agrawala, M., Hoppe, H., Szeliski, R., Cohen, M., & Toyama, K. (2004). Digital photography with flash and no-flash image pairs. *ACM Transactions on Graphics*, 23(3), 664–672. Proceedings of the ACM SIGGRAPH conference.
- Pham, T. Q. (2006). *Spatiotonal adaptivity in super-resolution of undersampled image sequences*. Ph.D. thesis, Delft University of Technology.
- Pham, T. Q., & van Vliet, L. J. (2005). Separable bilateral filtering for fast video preprocessing. In *International conference on multimedia and expo*. New York: IEEE Press.
- Porter, T., & Duff, T. (1984). Compositing digital images. *Computer Graphics*, 18(3), 253–259.
- Sand, P., & Teller, S. (2006). Particle video: Long-range motion estimation using point trajectories. In *Proceedings of the computer vision and pattern recognition conference*.
- Shannon, C. E. (1949). Communication in the presence of noise. *Proceedings of the Institute of Radio Engineers*, 37(1), 10–21.
- Smith, S. (2002). *Digital signal processing*. London: Newnes. ISBN: 075067444X.
- Smith, S. M., & Brady, J. M. (1997). SUSAN—a new approach to low level image processing. *International Journal of Computer Vision*, 23(1), 45–78.
- Sochen, N., Kimmel, R., & Malladi, R. (1998). A general framework for low level vision. *IEEE Transactions in Image Processing*, 7, 310–318.
- Sochen, N., Kimmel, R., & Bruckstein, A. M. (2001). Diffusions and confusions in signal and image processing. *Journal of Mathematical Imaging and Vision*, 14(3), 237–244.
- Tomasi, C., & Manduchi, R. (1998). Bilateral filtering for gray and color images. In *Proceedings of the international conference on computer vision* (pp. 839–846). New York: IEEE Press.
- van de Weijer, J., & van den Boomgaard, R. (2001). Local mode filtering. In *Proceedings of the conference on computer vision and pattern recognition*.
- Weickert, J., ter Haar Romeny, B. M., & Viergever, M. A. (1998). Efficient and reliable schemes for nonlinear diffusion filtering. *IEEE Transactions on Image Processing*, 7, 398–410.
- Weiss, B. (2006). Fast median and bilateral filtering. *ACM Transactions on Graphics*, 25(3), 519–526. Proceedings of the ACM SIGGRAPH conference.
- Willis, P. J. (2006). Projective alpha colour. *Computer Graphics Forum*, 25(3), 557–566. Proceedings of the Eurographics conference.
- Winnemöller, H., Olsen, S. C., & Gooch, B. (2006). Real-time video abstraction. *ACM Transactions on Graphics*, 25(3), 1221–1226. Proceedings of the ACM SIGGRAPH conference.
- Wong, W. C. K., Chung, A. C. S., & Yu, S. C. H. (2004). Trilateral filtering for biomedical images. In *Proceedings of the international symposium on biomedical imaging*. New York: IEEE Press.
- Xiao, J., Cheng, H., Sawhney, H., Rao, C., & Isnardi, M. (2006). Bilateral filtering-based optical flow estimation with occlusion detection. In *Proceedings of the European conference on computer vision*.
- Yaroslavsky, L. P. (1985). *Digital picture processing. An introduction*. Berlin: Springer.