

RAISR: Rapid and Accurate Image Super Resolution

Yaniv Romano, John Isidoro, and Peyman Milanfar, *Fellow, IEEE*

Abstract—Given an image, we wish to produce an image of larger size with significantly more pixels and higher image quality. This is generally known as the single image super-resolution problem. The idea is that with sufficient training data (corresponding pairs of low and high resolution images) we can learn set of filters (i.e., a mapping) that when applied to given image that is not in the training set, will produce a higher resolution version of it, where the learning is preferably low complexity. In our proposed approach, the run-time is more than one to two orders of magnitude faster than the best competing methods currently available, while producing results comparable or better than state-of-the-art. A closely related topic is image sharpening and contrast enhancement, i.e., improving the visual quality of a blurry image by amplifying the underlying details (a wide range of frequencies). Our approach additionally includes an extremely efficient way to produce an image that is significantly sharper than the input blurry one, without introducing artifacts, such as halos and noise amplification. We illustrate how this effective sharpening algorithm, in addition to being of independent interest, can be used as a preprocessing step to induce the learning of more effective upscaling filters with built-in sharpening and contrast enhancement effect.

Index Terms—Filter Learning, image enhancement, image sharpening, super resolution.

I. INTRODUCTION

SINGLE Image Super Resolution (SISR) is the process of estimating a High-Resolution (HR) version of a Low-Resolution (LR) input image. This is a well-studied problem, which comes up in practice in many applications, such as zoom-in of still and text images, conversion of LR images/videos to high definition screens, and more. The linear degradation model of the SISR problem is formulated by

$$\mathbf{z} = \mathbf{D}_s \mathbf{H} \mathbf{x}, \quad (1)$$

where $\mathbf{z} \in \mathbb{R}^{M \times N}$ is the input image, $\mathbf{x} \in \mathbb{R}^{Ms \times Ns}$ is the unknown HR image, both are held in lexicographic ordering. The linear operator $\mathbf{H} \in \mathbb{R}^{MN s^2 \times MN s^2}$ blurs the image \mathbf{x} , followed by a decimation in a factor of s in each axis, which is the outcome of the multiplication by $\mathbf{D}_s \in \mathbb{R}^{MN \times MN s^2}$. In the SR

task, the goal is to recover the unknown underlying image \mathbf{x} from the known measurement \mathbf{z} . Note that, in real world scenarios, the degradation model can be non-linear (e.g. due to compression) or even unknown, and may also include noise.

The basic methods for upscaling a single image are the linear interpolators, including the nearest-neighbor, bilinear and bicubic [1], [2]. These methods are widely used due to their simplicity and low complexity, as the interpolation kernels (upsampling filters) are not adaptive to the content of the image. However, naturally, these linear methods are limited in reconstructing complex structures, often times result in pronounced aliasing artifacts and over-smoothed regions. In the last decade powerful image priors were developed, e.g., the self-similarity [3]–[6], sparsity [7]–[12], and Gaussian Mixtures [13], leading to high quality restoration with the cost of increased complexity.

In this paper we concentrate on example-based methods [8], [9], [11], [14]–[18], which have drawn a lot of attention in recent years. The core idea behind these methods is to utilize an external database of images and learn a mapping from LR patches to their HR counterparts. In the learning stage, LR-HR pairs of image patches are synthetically generated, e.g., for $2 \times$ upscaling, a typical size of the HR patch is 6×6 and the one of the synthetically downsampled LR patch is 3×3 . Then, the desired mapping is learned and regularized using various local image priors.

The sparsity model is one such prior [8], [9], where the learning mechanism results in a compact (sparse) representation of pairs of LR and HR patches over learned dictionaries. Put differently, per each LR patch, these methods construct a non-linear adaptive filter (formulated as a projection matrix), which is a combination of a few basis elements (the learned dictionary atoms) that best fit to the input patch. Applying the filter that is tailored to the LR patch leads to the desired upscaling effect.

The Anchored Neighborhood Regression (ANR) [10] keeps the high quality reconstruction of [8] and [9] while achieving a significant speed-up in runtime. This is done by replacing the sparse-coding step that computes the compact representation of each patch over the learned dictionaries, with set of pre-computed projection matrices (filters), which are the outcome of ridge regression problems. As such, at runtime, instead of applying sparse-coding, ANR suggest searching for the nearest atom to the LR patch, followed by a multiplication by the corresponding pre-computed projection matrix. A follow-up work, called A+ [11], improves the performance of ANR by learning regressors not only from the nearest dictionary atoms, but also from the locally nearest training samples, leading to state-of-the-art restoration.

SRCNN [16] is another efficient example-based approach that builds upon deep Convolutional Neural Network (CNN) [19], and learns an end-to-end mapping from LR images to

Manuscript received May 20, 2016; revised August 12, 2016 and September 30, 2016; accepted October 19, 2016. Date of publication November 15, 2016; date of current version February 8, 2017. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Sylvain Paris.

Y. Romano was with the Google Research, Mountain View, CA 94043 USA. He is now with the Department of Electrical Engineering, Technion – Israel Institute of Technology, Haifa 32000, Israel (e-mail: yromano@tx.technion.ac.il).

J. Isidoro and P. Milanfar are with Google Research, Mountain View, CA 94043 USA (e-mail: isidoro@google.com; milanfar@google.com).

More results can be found in the supplementary material: <https://goo.gl/DOETxG>.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TCI.2016.2629284

their HR counterparts. Note that, differently from sparsity-based techniques, SRCNN does not explicitly learn the dictionaries for modeling the patches. In this case, the model is implicitly learned by the hidden convolutional layers.

The above mentioned SISR methods result in impressive restoration, but with the cost of (relatively) high computational complexity. In this paper we suggest a learning-based framework, called RAISR, which produces high quality restoration while being two orders of magnitude faster than the current leading algorithms, with extremely low memory requirements.

The core idea behind RAISR is to enhance the quality of a very cheap (e.g. bilinear) interpolation method by applying a set of pre-learned filters on the image patches, chosen by an efficient hashing mechanism. Note that the filters are learned based on pairs of LR and HR training image patches, and the hashing is done by estimating the local gradients' statistics. As a final step, in order to avoid artifacts, the initial upscaled image and its filtered version are locally blended by applying a weighted average, where the weights are a function of a structure descriptor. We harness the Census Transform (CT) [20] for the blending task, as it is extremely fast and cheap descriptor of the image structure which can be utilized to detect structure deformations that occur due to the filtering step.

A closely related topic to SISR is image sharpening, aiming to amplify the structures/details of a blurry image. The basic sharpening techniques apply a linear filter on the image, as in the case of unsharp masking [21] or Difference of Gaussians (DoG) [22], [23]. These techniques are highly effective in terms of complexity, but tend to introduce artifacts such as over-sharpening, gradient reversals, noise amplification, and more. Similarly to SISR, improved results can be obtained by relying on patch priors, where the sensitivity to the content/structure of the image is the key for artifact-free enhancement [24]–[28]. For example, with the cost of increased complexity compared to the linear approach, the edge-aware bilateral filter [29], [30], Non-Local Means [3] and guided filter [25] produce impressive sharpening effect.

As a way to generate high-quality sharp images, one can learn a mapping from LR images to their *sharpened* HR versions, thus achieving a built-in sharpening/contrast-enhancement effect “for free”. Furthermore, the learning stage is not limited to a linear degradation model (as in Eq. (1)), as such, learning a mapping from *compressed* LR images to their *sharpened* HR versions can be easily done, leading to an “all in one” mechanism that not only increases the image resolution, but also reduces compression artifacts and enhances the contrast of the image.

Triggered by this observation, we develop a sharpener as well, which is of independent interest. The proposed sharpener is highly efficient and able to enhance both fine details (high frequencies) and the overall contrast of the image (mid-low frequencies). The proposed method has almost similar complexity to the linear sharpeners, while being competitive with far more complex techniques. The suggested sharpener is based on applying DoG filters [22], [23] on the image, which are capable to enhance a wide range of frequencies. Next, a CT-based structure-aware blending step is applied as a way to prevent artifacts due

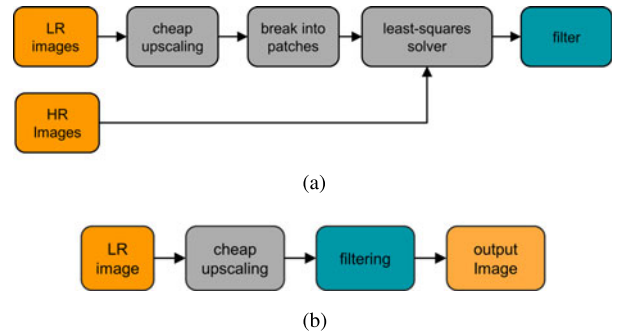


Fig. 1. The basic learning and application scheme of a global filter that maps LR images to their HR versions. (a) Learning Stage. (b) Upscaling Stage.

to the added content-aware property (similar mechanism to the one suggested in the context of SISR).

This paper is organized as follows: In Section II we describe the global learning and upscaling scheme, formulating the core engine of RAISR. In Section III we refine the global approach by integrating the initial upscaling kernel to the learning scheme. In Section IV we describe the overall learning and upscaling framework, including the hashing and blending steps. The sharpening algorithm is detailed in Section V. Experiments are brought in Section VI, comparing the proposed upscaling and sharpening algorithm with state-of-the-art methods. Conclusions and future research directions are given in Section VII.

II. FIRST STEPS: GLOBAL FILTER LEARNING

Given an initial (e.g. bilinear in our case) upscaled versions of the training database images, $\mathbf{y}_i \in \mathbb{R}^{M \times N}$, with $i = 1, \dots, L$, we aim to learn a $d \times d$ filter h that minimizes the Euclidean distance between the collection $\{\mathbf{y}_i\}$ and the desired training HR images $\{\mathbf{x}_i\}$. Formally, this is done by solving a least-squares minimization problem

$$\min_{\mathbf{h}} \sum_{i=1}^L \|\mathbf{A}_i \mathbf{h} - \mathbf{b}_i\|_2^2 \quad (2)$$

where $\mathbf{h} \in \mathbb{R}^{d^2}$ denotes the filter $h \in \mathbb{R}^{d \times d}$ in vector-notation; $\mathbf{A}_i \in \mathbb{R}^{MN \times d^2}$ is a matrix, composed of patches of size $d \times d$, extracted from the image \mathbf{y}_i , each patch forming a row in the matrix. The vector $\mathbf{b}_i \in \mathbb{R}^{MN}$ is composed of pixels from \mathbf{x}_i , corresponding to the center coordinates of \mathbf{y}_i patches. The block diagram, demonstrating the core idea of the learning process is given in Fig. 1(a).

In practice, the matrix \mathbf{A} can be very large, so we employ two separate approaches to control the computational complexity of estimating the filter. First, in general not all available patches needs to be used in order to obtain a reliable estimate. In fact, we typically construct \mathbf{A}_i and \mathbf{b}_i by sampling K patches/pixels from the images on a fixed grid, where $K \ll MN$. Second, the minimization of the least-squares problem, formulated in Eq. (2), can be recast in a way that significantly reduces both memory and computational requirements. To simplify the exposition, the following discussion is given in the context of filter learning based on just one image, but extending the idea to

several images and filters is trivial. The proposed approach results in an efficient solution for the learning stage where the memory requirements are only on the order of the size of the learned filter. The solution is based on the observation that instead of minimizing Eq. (2), we can minimize

$$\min_{\mathbf{h}} \|\mathbf{Q}\mathbf{h} - \mathbf{V}\|_2^2, \quad (3)$$

where $\mathbf{Q} = \mathbf{A}^T \mathbf{A}$ and $\mathbf{V} = \mathbf{A}^T \mathbf{b}$.

Notice that \mathbf{Q} is a small $d^2 \times d^2$ matrix, thus requiring relatively little memory. The same observation is valid for \mathbf{V} that requires less memory than holding the vector \mathbf{b} . Furthermore, based on the inherent definition of matrix-matrix and matrix-vector multiplications, we in fact avoid holding the whole matrix (and vector) in memory. More specifically, \mathbf{Q} can be computed cumulatively by summing chunks of rows (for example sub matrices $\mathbf{A}_j \in \mathbb{R}^{q \times d^2}$, $q \ll MN$), which can be multiplied independently, followed by an accumulation step; i.e.

$$\mathbf{Q} = \mathbf{A}^T \mathbf{A} = \sum_j \mathbf{A}_j^T \mathbf{A}_j \quad (4)$$

The same observation is true for matrix-vector multiplication

$$\mathbf{V} = \mathbf{A}^T \mathbf{b} = \sum_j \mathbf{A}_j^T \mathbf{b}_j, \quad (5)$$

where $\mathbf{b}_j \in \mathbb{R}^q$ is a portion of the vector \mathbf{b} , corresponding to the matrix \mathbf{A}_j . Thus, the complexity of the proposed learning scheme in terms of memory is very low – it is in the order of the filter size. Moreover, using this observation we can parallelize the computation of $\mathbf{A}_j^T \mathbf{A}_j$ and $\mathbf{A}_j^T \mathbf{b}_j$, leading to a speedup in the runtime. As for the least squares solver itself, minimizing Eq. (3) can be done efficiently since \mathbf{Q} is a positive semi-definite matrix, which perfectly suits a fast conjugate gradients solver [31].

To summarize, the learning stage is efficient both in terms of the memory requirements and ability to parallelize. As displayed in Fig. 1(b), at run-time, given a LR image (that is not in the training set), we produce its HR approximation by first interpolating it using the same cheap upscaling method (e.g. bilinear) that is used in the learning stage, followed by a filtering step with the pre-learned filter.

III. REFINING THE CHEAP UPSCALING KERNEL: DEALING WITH ALIASING

The ‘‘cheap’’ upscaling method we employ as a first step, can be any method, including a non-linear one. However, in order to keep the low complexity of the proposed approach, we use the bilinear interpolator as the initial upscaling method¹. Inspired by the work in [15], whatever the choice of the initial upscaling method, we make the observation that when aliasing is present in the input LR image, the output of the initial upscaler will generally not be shift-invariant to this aliasing.

As illustrated in Fig. 2, in the case of upscaling by a factor of 2 in each axis, the interpolation weights of the bilinear kernel

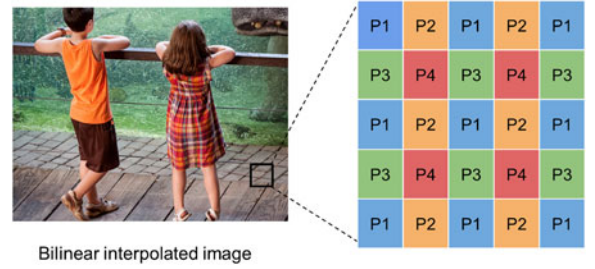


Fig. 2. Bilinear upscaling by a factor of 2 in each axis. There are four types of pixels, denoted by P1-P4, corresponding to the four kernels that are applied during the bilinear interpolation.

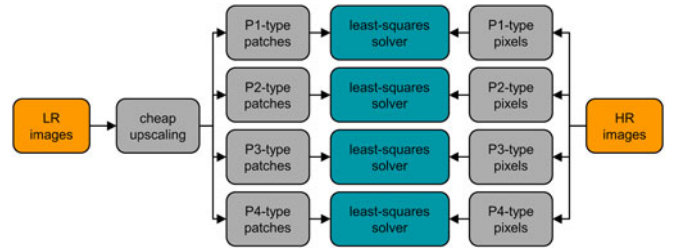


Fig. 3. Spatially varying learning scheme of four global filters, taking into consideration the internal structure of the bilinear.

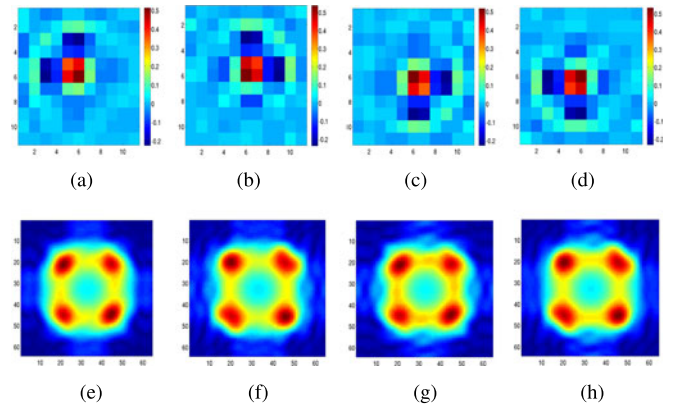


Fig. 4. Visualization of the four global filters, corresponding to P1-P4 type of pixels, in the pixel domain (a-d), along with their magnitude in the frequency domain (e-h), where the warmer the color, the larger the value. The filters are learned on Fig. 2 image: (a) P1-Filter, (b) P2-Filter, (c) P3-Filter, (d) P4-Filter, (e) P1-Spectrum, (f) P2-Spectrum, (g) P3-Spectrum, (h) P4-Spectrum.

vary according to the pixel’s location. As can be seen, there are four possible kernels that are applied on the LR image according to the type of the pixel, denoted by P1-P4. Since a convolution of two linear filters can be unified into one filter (in our case, the first is the bilinear and the second is the pre-learned one)², we should learn four different filters, corresponding to the four possible types of pixels, as demonstrated in Fig. 3.

The importance of this observation is illustrated in Fig. 4, which plots examples of actual learned filters, along with their magnitude in the frequency domain. The obtained filters act like

¹We also restrict the discussion mainly to the case of $2 \times$ upscaling to keep the discussion straightforward. Extensions will be discussed at the end of this section.

²This observation is a promising way to further speed up the algorithm and reduce the overall complexity.

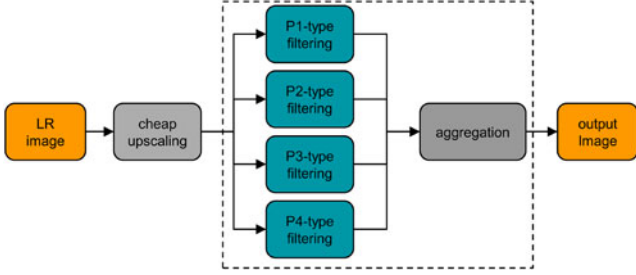


Fig. 5. Applying the four spatially varying pre-learned filters on a LR image.

bandpass filters, amplifying the mid-frequencies, and suppressing the high-frequencies (which contain aliasing components) of the interpolated image. The learned filters have similar magnitude response (Fig. 4(e)–4(h)), but different phase response (Fig. 4(a)–4(d)), standing in agreement with the four different shifted versions of the interpolation kernels.

On the application side, similarly to the core/naive upscaling idea, we first upscale the LR image using the bilinear interpolator. Then, differently from the naive approach, we apply the pre-learned filters according to the type of the pixel, followed by an aggregation step that simply combines the outcome of the filtered patches (resulting in a pixel) to an image. This process is illustrated in Fig. 5.

Notice that a similar observation holds for upscaling by any other integer factor s . For example, upscaling by a factor of 3 implies that we should learn 9 filters, one per each pixel-type. Similarly, when upscaling by a factor of 4, there are 16 types of pixels. As already mentioned, in order to keep the flow of the explanations, we will concentrate on the $2\times$ scenario since the generalization to other scaling factors is straightforward.

IV. RAISR: HASHING-BASED LEARNING AND UPSCALING

Generally speaking, the global image filtering is fast and cheap, as it implies the application of one filter per patch. Since the learning scheme reduces the Euclidean distance between the HR and the interpolated version of the LR images, the global filtering has the ability to improve the restoration performance of various linear upscaling methods. However, the global approach described so far is weaker than the state-of-the-art algorithms, e.g., sparsity-based methods [8]–[11] or the neural networks based ones [16] that build upon large amount of parameters, minimizing highly nonlinear cost functions. In contrast to these methods, the global approach is not adaptive to the content of the image, and its learning stage estimates only a small amount of parameters.

Adaptivity to the image content can be achieved by dividing the image patches into clusters, and constructing an appropriate filter per each cluster (e.g. as done in [10], [11]). However, the clustering implies the increase of the overall complexity of the algorithm, which is an outgrowth that we want to avoid. Therefore, instead of applying “expensive” clustering (e.g. K-means [32], GMM [33], [34], dictionary learning [8], [9], [11], [14]–[18]), we suggest using an efficient hashing approach, leading to adaptive filtering that keeps the low complexity of the

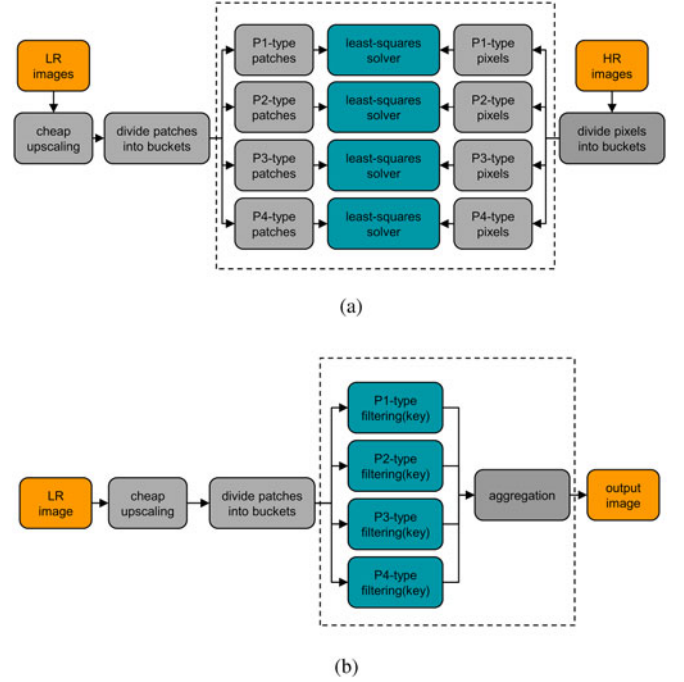


Fig. 6. Hashing based learning and upscaling schemes. We suggest dividing the patches into “buckets”, where each bucket contains patches with similar geometry (can be considered as a cheap clustering method). Then, a least squares fitting is applied per each bucket and possible shift. At run-time the hash-table key is computed per each patch, leading to the corresponding pre-learned locally adaptive filters.

linear filtering. More specifically, the local adaptivity is achieved by dividing the image patches into groups (called “buckets”) based on an informative and “cheap” geometry measures, which utilize the statistics of the gradients (a detailed description is given in Section IV-A). Then, similarly to the global approach, we also learn four filters, *but this time per each bucket*. As a consequence, the proposed learning scheme results in a hash-table of filters, where the hash-table keys are a function of the local gradients, and the hash-table entries are the corresponding pre-learned filters. An overview of the proposed hashing-based learning is shown in Fig. 6(a).

Given the hash-table, containing filters per quantized edge-statistic descriptor (more details in Section IV-A), the upscaling procedure becomes very effective. Following Fig. 6(b), we compute the hash-table key per each patch of the initial interpolated image, pointing to the relevant filters (four filters, one per patch-type), to be applied on the corresponding patch.

Similarly to the global learning process (see Section II), we utilize the matrix-matrix and matrix-vector multiplications once again. Per each bucket q , we learn a filter \mathbf{h}_q by minimizing the following cost function

$$\min_{\mathbf{h}_q} \|\mathbf{A}_q^T \mathbf{A}_q \mathbf{h}_q - \mathbf{A}_q^T \mathbf{b}_q\|_2^2, \quad (6)$$

where \mathbf{A}_q and \mathbf{b}_q are the patches and pixels that belong to the q -th bucket. In this case, the low memory requirements of the proposed learning process are crucial, especially for large hash-table that requires millions of examples to produce a

reliable estimate for the filters. As a consequence, by utilizing the observation described in Section II, we perform a sub-matrix accumulation on a sub-image block basis, leading to a learning process that can handle any desired number of examples.

A. Hash-Table Keys: Local Gradient Statistics (Angle, Strength, Coherence)

Naturally, there are many possible local geometry measures that can be used as the hash-table keys, whereas the statistics of the gradients has a major influence on the proposed approach. We suggest evaluating the local gradient characteristics via eigenanalysis [35], which yields the gradient's angle and information about the strength and coherence of the nearby gradients. Eigenanalysis also helps in cases of thin lines, stripes and other scenarios that the mean gradient might be zero, yet the neighborhood exhibits a strong directionality.

The direction, strength and coherence are computed by utilizing the $\sqrt{n} \times \sqrt{n}$ surroundings of each pixel, i.e., for the k -th pixel we consider all the pixels that are located at k_1, \dots, k_n . The basic approach starts with a computation of $2 \times n$ matrix, composed from the horizontal and vertical gradients, g_x and g_y , of the surroundings of the k -th pixel, expressed by

$$\mathbf{G}_k = \begin{bmatrix} g_{x_{k_1}} & g_{y_{k_1}} \\ \vdots & \vdots \\ g_{x_{k_n}} & g_{y_{k_n}} \end{bmatrix}. \quad (7)$$

As stated in [35], the local gradient statistics can be computed using the Singular Value Decomposition (SVD) of this matrix. The right vector corresponds to the gradient orientation, and the two singular values indicate the strength and spread of the gradients. Since the work is being performed per-pixel, we hereby focus on efficiency. We can compute those characteristics more efficiently using an eigen-decomposition of $\mathbf{G}_k^T \mathbf{G}_k$ which is a 2×2 matrix, which can be computed conveniently in a closed form. Moreover, in order to incorporate a small neighborhood of gradient samples per pixel, we employ a diagonal weighting matrix \mathbf{W}_k , constructed using a separable normalized Gaussian kernel.

Following [35], the eigenvector ϕ_1^k , corresponding to the largest eigenvalue of $\mathbf{G}_k^T \mathbf{W}_k \mathbf{G}_k$, can be used to derive the angle of the gradient θ_k , given by

$$\theta_k = \arctan(\phi_{1,y}^k, \phi_{1,x}^k). \quad (8)$$

Notice that due to the symmetry, a filter that corresponds to the angle θ_k is identical to the one corresponding to $\theta_k + 180^\circ$.

As shown in [35], the square root of the largest eigenvalue λ_1^k is analogous to the ‘‘strength’’ of the gradient. The square root of the smaller eigenvalue λ_2^k can be considered as the ‘‘spread’’ of the local gradients, or rather how much they vary in direction. Both of these can be measured in units of intensity. The two eigenvalues can be combined into a unitless measure known as ‘‘coherence’’ [35]. The coherence value μ_k ranges from 0 to 1, and formulated as

$$\mu_k = \frac{\sqrt{\lambda_1^k} - \sqrt{\lambda_2^k}}{\sqrt{\lambda_1^k} + \sqrt{\lambda_2^k}}. \quad (9)$$

Algorithm 1: Computing the hash-table keys.

Inputs

- 1: Initial interpolated version of the LR image.
- 2: Q_θ – Quantization factor for angle (e.g. 24).
- 3: Q_s – Quantization factor for strength (e.g. 3).
- 4: Q_μ – Quantization factor for coherence (e.g. 3).

Output

- 1: Hash-table keys per pixel, denoted by θ_k , λ_1^k , and μ_k .

Process

- Compute the image gradients ;
 - Construct the matrix $\mathbf{G}_k^T \mathbf{W}_k \mathbf{G}_k$, and obtain the gradients' angle θ_k , strength λ_1^k , and coherence μ_k ;
 - Quantize: $\theta_i \leftarrow \lceil \frac{\theta_i}{Q_\theta} \rceil$, $\lambda_1^i \leftarrow \lceil \frac{\lambda_1^i}{Q_s} \rceil$, $\mu_i \leftarrow \lceil \frac{\mu_i}{Q_\mu} \rceil$, where $\lceil \cdot \rceil$ is the ceiling function
-

Strength and coherence are very useful for detecting a variety of different local image properties. A low strength and low coherence often signifies a lack of image structure, and usually corresponds to noise or compression artifacts. High strength, but low coherence often indicates corners or other multi-directional structure. Having a high coherence is generally an edge, or series of stripes in the same direction, with the strength measuring the relative intensity of the stripes. Intuitively, strength and coherence allow us to detect semantically different local image properties, so by using them as part of a hash enables the filter learning process to adapt to these conditions. As such, combining angle θ_k , strength λ_1^k , and coherence μ_k into a hash function, as detailed in Algorithm 1, can produce a family of learned filters that are able to handle a variety of situations.

In Fig. 7, one can see that bucketing by angle, coherence and strength produces a wide variety of filters. The ones that correspond to low coherency and strength tend to be bandpass and directionally invariant in nature. As coherency increases, so does the directionality of the filter, smoothing orthogonal to the gradient, but strongly sharpening in the direction of the gradient.

B. Using Patch Symmetry for Nearly-Free $8 \times$ More Learning Examples

The amount of data needed for effective and stable learning of filter sets can be large. For instance, in practice, it takes at least 10^5 patches to reliably learn a given filter of size 9×9 or 11×11 . Say we decompose the patches into \mathcal{B} buckets, when using a hashed set of filters, this implies that we need 10^5 patches per bucket. However, reaching this amount using real world training data is not as simple as using $10^5 \times \mathcal{B}$ patches.

The problem emerges from the observation that certain hash values occur much more commonly than other hashes. There are often many more horizontal and vertical structures in imagery, and flat regions (such as sky, and painted surfaces) are common. Intuitively, these common structures result in the more common hashes.

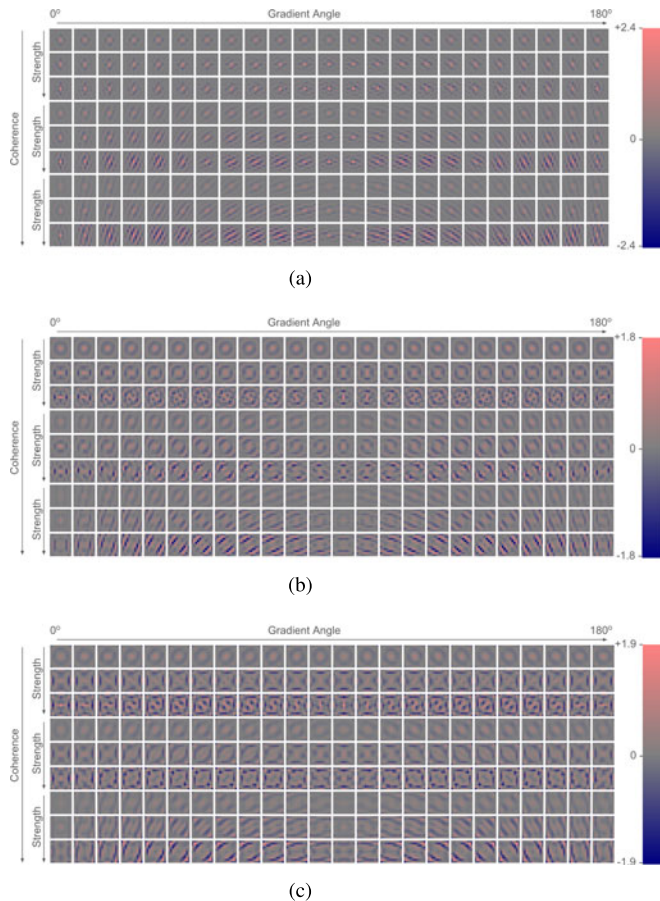


Fig. 7. Visualization of the learned filter sets for (a) $2\times$, (b) $3\times$ and (c) $4\times$ upscaling, learned from using an angle, strength, and coherence based hashing scheme. Per each subset of filters, the angle varies from left to right; the top, middle, and bottom 3 rows correspond to low, medium and high coherence. Within each set of 3 rows, gradient strength increases from top to bottom. As can be inferred, the general trend is that as coherence increases, the directionality of the filter increases. Also, as strength increases the intensity of the filter increases. Notice how the $3\times$ and $4\times$ upscaling filters are not simply scaled versions of the $2\times$ filters, but also have extracted additional information from the training data.

In order to make the patches that hash to uncommon hash values more effective, we leverage the patch symmetry and increase the learning power of each patch. More specifically, we can generate 8 different example patches; four 90° rotations, and four mirrored 90° rotations. Since 8 patches are generated for each original patch, we are effectively incorporating 8 times as much information for learning. Since each transformation is a rotation and mirroring, the transformed patches often belong to a different hash bucket and shift. For example, a patch that is rotated by 90° changes the hash angle bucket by 90° . A visual demonstration of the proposed idea is given in Fig. 8, which demonstrates a simplified hashing scheme with 4 angular buckets on a polar graph, and how they get rebucketed from x-flips, y-flips, and xy-swaps.

Moreover, the patch transformations do not need to be performed as actual image transformations of each incoming patch, which would be expensive. If the gradient angle dependent hash bucket boundaries are symmetric to flips in x, y and xy-swaps, the accumulation for transformed patches can be performed.

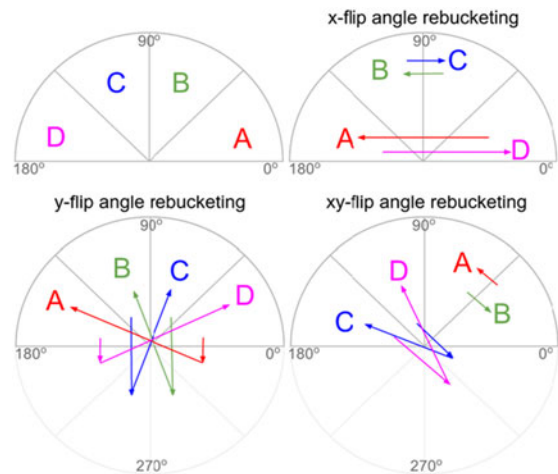


Fig. 8. Rebucketing for x-flips, y-flips, and xy-swaps. For example bucket 'B' accounts for gradient angles between 45° and 90° . When the patch is flipped on the x-axis, the bucketing also needs to be reflected, as shown in the second graph. The third graph shows y-flips, for which the subsequent modulo 180° operation maps back to the 0° to 180° range. The final diagram shows xy-swaps (flip around the $x=y$ line). The different possible combinations of x-flip, y-flip, and xy-swaps account for the 8 different patch transformations.

Achieving this symmetry is a function of the hashing function, and, in our case, having a number of angle buckets that is divisible by 4 satisfies this requirement. As such, we can accumulate the per bucket and per pixel-type matrices $\mathbf{A}_q^T \mathbf{A}_q$ and $\mathbf{A}_q^T \mathbf{b}$ across all training samples, as suggested in Eq. (6). Then, the symmetry can be applied as one last accumulation of permuted matrices, which act as a set of symmetry augmented matrices. In practice, the additional accumulation step to enable symmetry takes less than 0.1% of the learning runtime (only a few additional seconds on a 3.4 GHz 6-Core Xeon desktop computer).

C. Built-in Suppression of Compression Artifacts and Sharpening Effect

The linear degradation model that assumes blur and decimation, as expressed in Eq. (1), is very common in the literature, but less so in the real world. For example, often times the measured images are blurred with unknown kernel, compressed, post-processed (e.g. by applying gamma correction), contaminated by noise, and more.

Learning a mapping that is capable to handle highly non-linear degradation model can be done by RAISR. We found that an effective suppression of compression artifacts is achieved by learning a mapping from *compressed* LR images to their HR versions. Notice that an inherent compression parameter is the bit-rate/compression quality, affecting the outcome of the learning scheme. For example, JPEG encoders use a quality level parameter, which varies from 0 (the worst quality) up to 100 (the best quality). Our experiments show that an aggressive compression (e.g. 80) indeed suppresses the compression artifacts, but can lead to smoothed result. We also found that a moderate compression level (e.g. 95) in training helps to suppress the aliasing in addition to alleviating moderate compression artifacts.

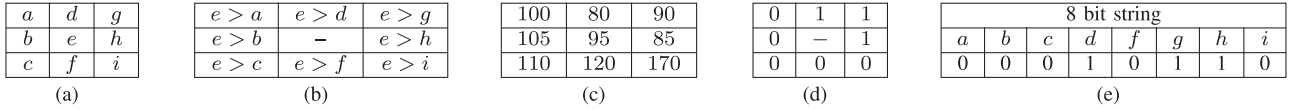


Fig. 9. Census transform. (a) 3x3 window of pixels, (b) Boolean comparisons between the center pixel and its neighbors, (c)-(e) Numerical example: (c) Intensity values, (d) The outcome of the boolean comparisons, (e) Census result, an 8 bit string that measures the local structure.

In the same spirit, gaining sharpening effect can be done by learning a mapping from LR training images to their sharpened HR versions. The stronger the sharpening during the training phase, the sharper the outcome of RAISR upscaling. We should emphasize that at runtime, we simply apply the pre-learned filters (with possibly built-in sharpening effect); we do not apply a separate sharpening step.

To conclude, by applying compression and sharpening as pre-processing steps, the learned filters are capable to map input compressed LR image to a sharpened HR output. As such, by choice, RAISR not only estimates the missing spatial information, but also suppresses the compression artifacts and amplifies the underlying signal.

D. Blending: An Efficient Structure-Preserving Solution

The proposed learning scheme results in adaptive upscaling filters, having built-in suppression of compression artifacts and sharpening effect. Concentrating on the sharpening property, two well known side-effects that exist are halos that appear along edges and the noise amplification. Put differently, applying the pre-learned filters on the initial interpolated image can lead to structure deformations due to the sharpening property.

As a way to avoid a significant modification in structure, we suggest measuring the local change in structure that occurred due to the filtering step, and blend accordingly. In areas that the structure of the initial interpolated image and the filtered one is somewhat similar – we choose the filtered version, while in areas that a major change is occurred by the filtering step – we choose the initial upscaled image. This suggestion relies on the observation that the cheap interpolated image typically does an adequate job on areas of the image that contain low spatial frequencies (e.g. flat regions). On the other hand, the higher spatial frequencies that need to be reconstructed require the careful treatment of the estimated filters. The blending therefore combines the most appropriate contributions from the cheap upscaled and the RAISR filtered image to yield the final result. One could have identified these regions ahead of time by clustering and apply different treatment, however this would have resulted in slower execution. In what follows, we present a fast alternative that works on two output images for point-wise blending.

Inspired by the CT descriptor [20], we suggest using its outcome as an engine that detect structure deformations and revert the errors of the upscaler. In order to understand the blending mechanism, a brief overview about the CT is given here. This transform maps the intensity values of the pixels within a small squared region (e.g. of size 3×3) to a bit string that captures the image structure. The CT is based on the relative ordering of local intensity values, and not on the intensity values themselves.

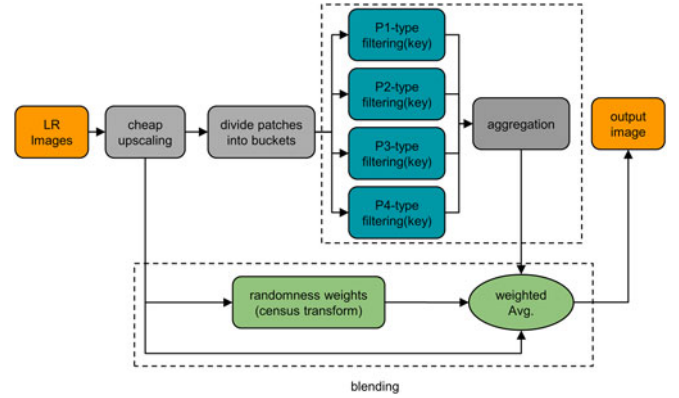


Fig. 10. RAISR upscaling that allows the amplification of high-frequencies only. Applying the pre-learned filters and avoiding halos and noise amplification. The blending select the filtered pixels in structured areas, and the cheap upscaled pixels in flat areas.

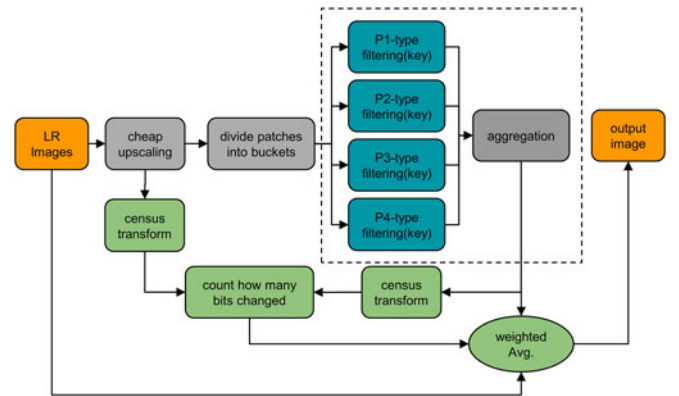


Fig. 11. RAISR upscaling that allows the amplification of wide range of frequencies, enabling a contrast enhancement effect. Applying the pre-learned filters and avoiding halos and noise amplification. The blending mechanism enables efficient structure and contrast enhancement, i.e., the amplification of a wide range of frequencies.

Following Fig. 9, the center pixel's intensity value (central cell in Fig. 9(a) and (c)) is replaced by a bit string descriptor of length 8 (Fig. 9(e)), composed of a set of boolean comparisons between the center pixel and its 3×3 neighborhood pixels (Fig. 9(b), (d)). Note that, in practice, when applying the comparisons we allow small variations in the intensity values, controlled by a threshold.

Back to RAISR, we suggest two different blending schemes, given in Fig. 10 and Fig. 11, which result in different enhancement effects (as illustrated in Sections VI-A and VI-B, respectively). The first scheme allows the pre-learned filters to enhance only the high-frequencies of the image, leading to HR images

that look natural, as required by the conventional SISR problem (a demonstration is given in Fig. 14). While the second scheme allows the enhancement of a wide range of frequencies, leading to better looking images due to the contrast enhancement effect (as shown in Fig. 17). Note that both versions aim at avoiding artifacts and structure deformations.

We start with the first blending mechanism, which enables the amplification of the high-frequencies only without modifying the low/mid frequencies, as usually done by the conventional SISR algorithms. Following Fig. 10, our suggestion is based on the observation that in flat areas (or generally low frequency areas), a linear upscaler produces good results because there are no fine details to be recovered or aliasing to be suppressed, thus there is no need to further improve the results in these areas. On the other hand, the linear interpolation fails to recover structured areas, where the proposed pre-learned filters play the key role. Moreover, within the structured areas, especially along strong edges, the pre-learned filters may introduce halos due to their size (11×11 or 9×9) and the sharpening property.

For the sake of completeness, let us explain briefly how the CT, which is indeed blind to the illumination, can be formulated as a mechanism for edge/structure-detection, thus allowing us to amplify only the high-frequencies of the image. In this case, the blending weights are the outcome of the so called “randomness” measure, which indicates how likely a pixel is in a structured area. Specifically, the size of the Least Connected Component (LCC) of the CT descriptor (in the case of Fig. 9(d) the LCC size is 3) is translated to a weight, determining the strength/amount of the structure within the descriptor window. In general, the larger the size of LCC the higher the weight. Put differently, by measuring the “randomness” of the bit string we can infer whether the pixel is a part of an edge or not, forming the blending weights map. A block diagram of the proposed upscaling scheme, which allows sharpening of high-frequencies only, is given in Fig. 10.

Although the conventional SISR algorithms result in HR images that look natural, contrast enhancement (i.e. amplification of low, mid, and high frequencies) often times can lead to better looking images (a visual demonstration is given in Fig. 17). This observation leads us to propose the second CT-based blending scheme. According to Fig. 11, in order to measure the local change in structure, we

- 1) Compute the CT of the initial upscaled image and the filtered image, and then
- 2) Per each pixel – count the number of bits that were changed, i.e., evaluate the Hamming distance, where the larger the distance, the larger the change in structure.

As a consequence, by translating the number of bits that were changed to weights, we form the desired blending map. Notice that the CT is blind to the intensity value itself. Therefore, differently from the randomness strategy (see Fig. 11), in this case, the obtained blending map allows a local change in the intensity (or contrast), while avoiding major changes in the structure.

In the learning stage, the target HR images are pre-processed by the proposed DoG sharpener (a detailed explanation about the sharpener is given in Section V), which sharpens the structures

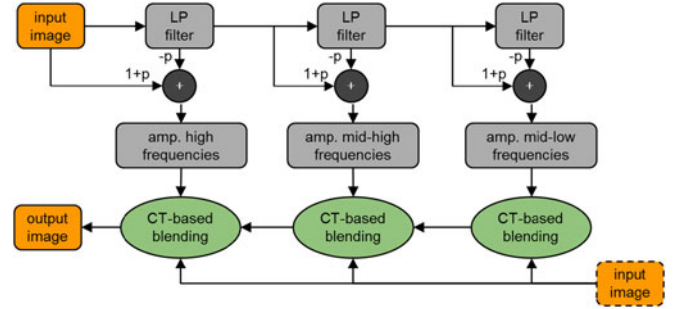


Fig. 12. CT-based DoG sharpener: Applying the DoG filters on the image, followed by blending steps that utilize the outcome of the CT.

and improves the overall contrast of the image. As a result, a built-in enhancement of both details (high-frequencies) and contrast (mid/mid-low frequencies) is achieved when applying the pre-learned filters. A block diagram of the upscaling scheme that allows contrast enhancement (i.e. amplification of a wide range of frequencies) is shown in Fig. 11.

Our experiments show that when enhancing only the high-frequencies we obtain images that look natural, having the same contrast of the LR ones (see Fig. 14). When enhancing a wider range of frequencies (i.e. allowing contrast change), RAISR produces better looking images (as illustrated in Fig. 17), but there is no guarantee that the output will have the same nature of the one of the LR. If we use the versions that clean compression artifacts, do sharpening and contrast enhancement, the effect on the PSNR or SSIM comparisons will not be very clear anymore. So in terms of this quantitative measure, we might observe deterioration, even if the images look excellent (even better than the originals!).

To conclude, we introduced the RAISR algorithm that turns a LR image into a HR image. The process is carried out in several steps:

- 1) A very cheap (e.g. bilinear) interpolation method is used to upscale the LR image.
- 2) A hash-table, containing set of filters, is learned from a training database, where the hash-table keys are a function of gradient properties. The filters are applied on the output of step (i) to improve its quality.
- 3) The outputs of steps (i) and (ii) are selectively blended (with different weights at each pixel) to produce the final result.

As a closing remark, the whole learning phase is summarized in Algorithm 2 and the upscaling stage is detailed in Algorithm 3.

V. CT-BASED DOG SHARPENER

In this section we briefly revisit the Difference of Gaussian (DoG) operator, which is widely used for edge enhancement [23]. Then, we introduce a very efficient way to sharpen an image using this operator while eliminating halos, noise amplification, and other similar common sharpening artifacts.

The DoG filter is formulated as a subtraction of two low-pass filters – Gaussians with different standard-deviation σ . In

Algorithm 2: RAISR: Hashing-based learning.**Inputs**

- 1: s – Upscaling factor.
- 2: $\{\mathbf{x}_i\}_{i=1}^L$ – The ground truth HR images (optionally sharpened/enhanced).
- 3: $\{\mathbf{z}_i\}_{i=1}^L$ – The LR (optionally compressed) versions of $\{\mathbf{x}_i\}_{i=1}^L$.

Output

- 1: $\mathcal{H}_t : (j) \mapsto \mathbf{h}_{j,t}$ – Hash-table that maps each key $j = (\theta, \lambda_1, \mu)$ to the corresponding d^2 -dimensional filter $\mathbf{h}_{j,t}$, where $1 \leq t \leq s^2$ is the pixel/patch type.

Process

- $\mathbf{Q}_{j,t} \leftarrow \mathbf{0}$, $\mathbf{V}_{j,t} \leftarrow \mathbf{0}$ for all hash-table keys denoted by j and pixel-type t
- for**
- $i = 1$
- to**
- L
- do**
- Compute \mathbf{y}_i , an initial interpolated version of \mathbf{z}_i
 - Initialize $\mathbf{A}_{j,t}$ and $\mathbf{b}_{j,t}$ to be empty matrices for all possible hash-table keys denoted by j and pixel-type t
- for**
- each pixel
- k
- in
- \mathbf{y}_i
- do**
- Denote by $j = (\theta_k, \lambda_1^k, \mu_k)$ the hash-table key of the pixel k
 - Denote by $\mathbf{p}_k \in \mathbb{R}^{d^2}$ the patch extracted from \mathbf{y}_i centered at k
 - Denote by $\mathbf{x}_i(k)$ the ground truth HR pixel of \mathbf{x}_i located at k
 - Denote by t the type of the pixel k
 - $\mathbf{A}_{j,t} \leftarrow [\mathbf{A}_{j,t}; \mathbf{p}_k^T]$, i.e., concatenate \mathbf{p}_k^T to the end of $\mathbf{A}_{j,t}$
 - $\mathbf{b}_{j,t} \leftarrow [\mathbf{b}_{j,t}; \mathbf{x}_i(k)]$, i.e., concatenate $\mathbf{x}_i(k)$ to the end of $\mathbf{b}_{j,t}$
- end**
- for**
- each key
- j
- and pixel-type
- t
- do**
- $\mathbf{Q}_{j,t} \leftarrow \mathbf{Q}_{j,t} + \mathbf{A}_{j,t}^T \mathbf{A}_{j,t}$
 - $\mathbf{V}_{j,t} \leftarrow \mathbf{V}_{j,t} + \mathbf{A}_{j,t}^T \mathbf{b}_{j,t}$
- end**
- end**
- for**
- each key
- j
- and pixel-type
- t
- do**
- $\mathbf{h}_{j,t} \leftarrow \arg \min_{\mathbf{h}} \|\mathbf{Q}_{j,t} \mathbf{h} - \mathbf{V}_{j,t}\|_2^2$
 - $\mathcal{H}_t(j) \leftarrow \mathbf{h}_{j,t}$
- end**

general, a Gaussian filter attenuates the high frequencies of the image, where the parameter σ controls the blurring effect, i.e., modifies the cut-off point of the filter. As such, when subtracting two Gaussians with different standard-deviations, we manually design a bandpass filter that reduces the amplitude of all frequencies between two cut-off points. More formally, a DoG filter $\mathcal{D}_{\sigma,\alpha,\rho}(\mathbf{z})$ can be expressed by

$$\mathcal{D}_{\sigma,\alpha,\rho}(\mathbf{z}) = (1 + \rho)\mathcal{G}_\sigma(\mathbf{z}) - \rho\mathcal{G}_{\alpha\sigma}(\mathbf{z}), \quad (10)$$

where $\mathcal{G}_\sigma(\mathbf{z})$ is a Gaussian filter with standard deviation σ , the scalar ρ controls the amplification factor, and α is a constant that controls the range of the frequencies that we wish to pass. According to the chosen parameters, the filter $\mathcal{D}_{\sigma,\alpha,\rho}(\mathbf{z})$ captures different frequency-bands, thus being an efficient mechanism to amplify a wide range of frequencies by choice [23].

Similarly to the explanation that is given in Section IV-D, when dealing with sharpening two main issues are raised: (i)

Algorithm 3: RAISR: Hashing-based upscaling.**Inputs**

- 1: s – Upscaling factor.
- 2: \mathbf{z} – LR image.
- 3: $\mathcal{H}_t : (j) \mapsto \mathbf{h}_{j,t}$ – Hash-table that maps each key $j = (\theta, \lambda_1, \mu)$ to the corresponding d^2 -dimensional filter $\mathbf{h}_{j,t}$, where $1 \leq t \leq s^2$ is the pixel/patch type.

Output

- 1: $\hat{\mathbf{x}}$ – HR estimate of \mathbf{z} .

Process

- Compute \mathbf{y} , an initial interpolated version of \mathbf{z}
- for**
- each pixel
- k
- in
- \mathbf{y}
- do**
- Denote by $j = (\theta_k, \lambda_1^k, \mu_k)$ the hash-table key of the pixel k
 - Denote by $\mathbf{p}_k \in \mathbb{R}^{d^2}$ the patch extracted from \mathbf{y} centered at k
 - Denote by t the type of the pixel k
 - $\mathbf{h}_{j,t} \leftarrow \mathcal{H}_t(j)$
 - $\hat{\mathbf{x}}(k) \leftarrow \mathbf{p}_k^T \mathbf{h}_t$
- end**
- $\hat{\mathbf{x}} \leftarrow \text{CT-Blending}(\hat{\mathbf{x}}, \mathbf{y})$

Noise amplification, and (ii) Halos artifacts. The noise may exist in the whole frequency domain, thus amplified by the naive DoG sharpener along with the underlying signal. In addition, since this sharpener is not adaptive to the content/structure of the image it tends to produce sharpening artifacts, such as halos, over-sharpening, gradient-reversals and more. Differently from the linear DoG filter, the content-aware non-linear filtering methods [24]–[28] successfully avoid these common artifacts. However, computationally, they are more complex than the linear approach.

We wish to keep the computational advantages of working with the linear DoG filter, while gaining adaptivity to the content of the image. Similarly to RAISR, we use the “blending trick” (see Section IV-D) once again. By integrating the CT blender to the naive DoG sharpener scheme we achieve the desired content-aware property, in an extremely efficient fashion. The blending is similar to the one that described in the context of RAISR, with the exception that now there are several images that are fused together. Specifically, the input image is locally blended with its different enhanced versions, obtained by applying naive DoG filters that amplifies different frequency-bands.

Moreover, in order to reduce computations, we suggest a cascade implementation that utilizes the already filtered images between the levels (see Fig. 12). For example, instead of applying a wide separable Gaussian filter in order to capture the mid-low frequency band, one can filter the already computed smoothed image, obtained in the previous cascade level. As can be easily inferred, the complexity of the proposed sharpener is very low; it is equivalent to the application of linear separable filters on the image, followed by a pixel-wise weighted average. Notice that the weights are a function of very cheap descriptor that applies extremely basic manipulations on 3×3 neighborhood of a pixel (e.g. boolean comparisons and evaluation of Hamming distance).

Following Fig. 12, notice that different blending mechanisms leads to different enhancement effect: When counting the modified bits of the CT descriptor, we are capable to enhance the contrast of the image (even in relatively low-frequencies). On the other hand, when choosing the randomness measure as a blending map, we enhance the content only along edges and structures.

VI. EXPERIMENTS

In this section we test the proposed algorithm in various scenarios. First, in the context of the conventional SISR, the effectiveness of RAISR is tested and compared to several state-of-the-art algorithms for $2\times$, $3\times$, and $4\times$ upscaling. Then, the abilities of the proposed algorithm to tackle real-world scenarios are demonstrated by applying RAISR upscaling on arbitrary compressed images. In this case, we show that RAISR is able to produce contrast-enhanced high quality images by learning filters from compressed LR images to their contrast-enhanced HR versions. The pre-processing of the training images is done by applying the proposed CT-based DoG sharpener.

A. Single Image Super-Resolution

In this subsection we compare the performance of proposed algorithm with several state-of-the-art methods for $2\times$, $3\times$, and $4\times$ upscaling factors, both quantitatively and visually. All the results are obtained by applying the scheme that is given in Fig. 10 (and its variants for $3\times$ and $4\times$ upscaling), followed by a back-projection step (see Appendix A). The filters are learned using a collection of 10,000 advertising banner images. This imagery type was chosen for its wide variety of both synthetic and real world content.

For all upscaling factors, we learn filters of size 11×11 . In the context of the hash-table, we consider a neighborhood of size 9×9 for the computation of the angle, strength, and coherence of the gradients. Also, the quantization factors of the angle Q_θ , strength Q_s , and coherence Q_μ are set to result in a total of 216 total buckets (24 angular bins by 3 strength bins by 3 coherence bins). In addition, we found that an amplification of the high-frequencies of the training-set HR images leads to an improved restoration. The proposed CT-based DoG sharpener (see Fig. 12) is utilized for this task, where we use the randomness measure as the blending map; the number of layers is set to 3, with $\sigma = 0.85$, $\alpha = \sqrt{2}$ and $\rho = 55$. Note that a fast approximation of separable Gaussian filter is suggested, obtained by applying one dimensional $[1, 2, 1]/4$ filter in a separable fashion (acts as a LP Gaussian filter with $\sigma \approx 0.85$).

As done by the leading algorithms [11], [16], [17], the proposed algorithm is tested on 2 widely used databases: (i) Set5 [36], and (ii) Set14 [9], which are composed of 5 and 14 standard images, respectively. For all scaling factors, the test LR images are generated by downscaling the original HR images using the bicubic interpolation.

The restoration performance is evaluated using the Peak Signal to Noise Ratio (PSNR), and Structural Similarity (SSIM) metrics, measured between the luminance channel of the original and the estimated images. The higher these measures the

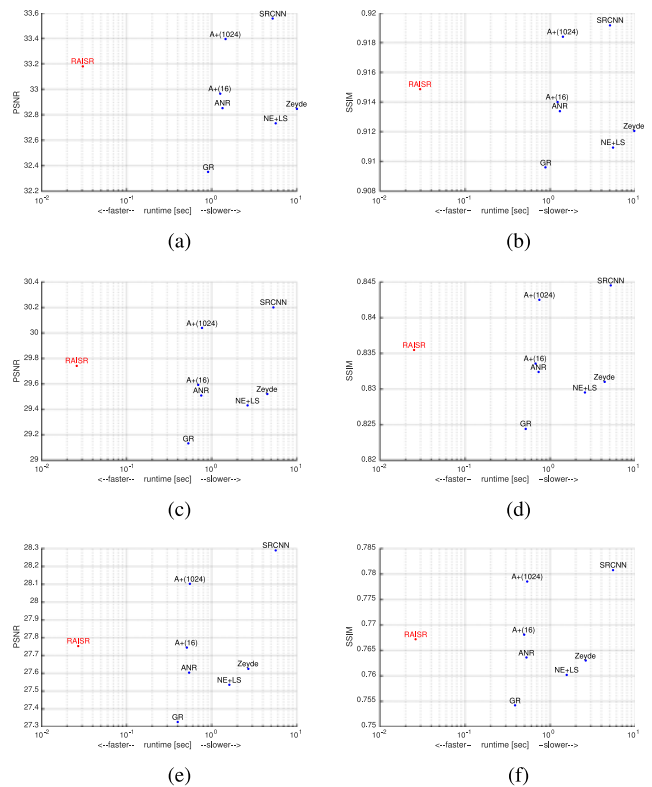


Fig. 13. Quantitative comparison between the restoration performance vs. runtime. The vertical coordinate of each point in the scatter plots corresponds to the average PSNR (a,c,e) and SSIM (b,d,f) of each method, measured on the test images of Set5 and Set14. The horizontal coordinate corresponds to the average runtime (the average size of the upscaled images is about 0.63×10^6 pixels). The size of each point reflects the standard error of the PSNR/SSIM. Detailed quantitative results can be found in the supplementary material, <https://goo.gl/D0ETxG>.

better the restoration. Note that upscaling an RGB image is done by converting it to the YCbCr color space. Then, a bicubic interpolation is applied on the luminance channel, which is used as an initialization for RAISR, while the chromatic channels are upscaled only by the bicubic interpolation. Finally, the estimated HR channels are converted back to the RGB color space.

The proposed algorithm is compared to various methods, including the sparse-coding approach of Zeyde et al. [9], along with the efficient GR and ANR [10]. Since we put emphasis on runtime, we test two versions of A+ [11], the first is a fast version that uses 16 atoms, while the second produces state-of-the-art results (with the cost of increased runtime) by utilizing 1024 atoms. In addition, we compare our algorithm to the state-of-the-art SRCNN [16], [17], which is based on a powerful convolutional neural network architecture. A neighbor embedding technique is also included, called NE+LLE, which assumes that the LR-HR patches lie on low-dimensional manifolds, having locally similar geometry [37]³. We should note that all the baseline methods significantly outperform the bicubic interpolation.

³We would like to thank the authors of [9], [11], [16], [17] for providing the software that produces the results of NE+LLE, Zeyde et al., GR, ANR, A+, and SRCNN. In all cases we use the original codes, with default parameters.

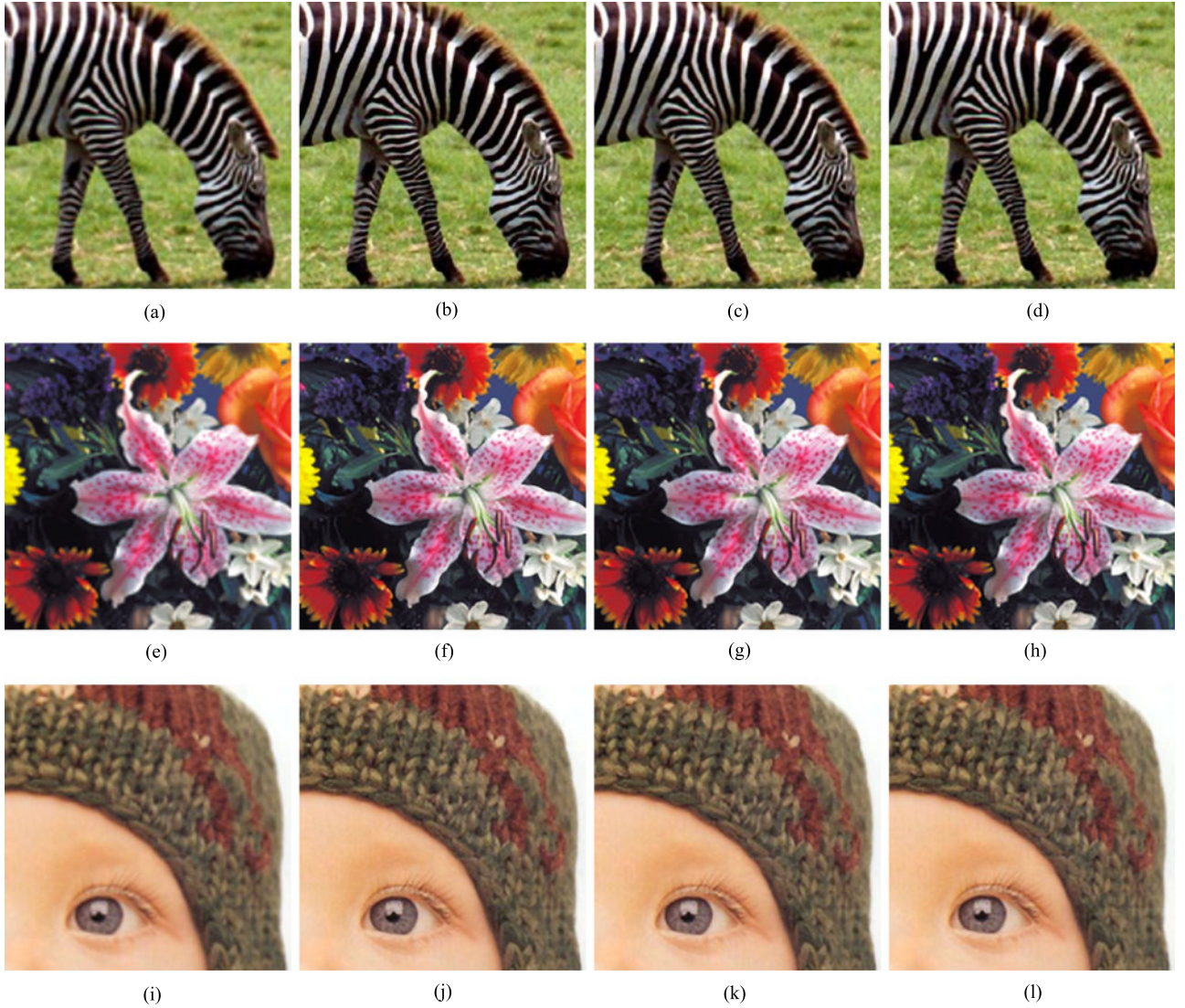


Fig. 14. Visual comparison for upscaling by a factor of 2 for the images Zebra, Flowers and Baby. (a,e,i) Bicubic interpolation, (b,f,j) SRCNN, (c,g,k) A+ with 1024 atoms, and (d,h,l) RAISR.

A quantitative comparison between RAISR and the leading SISR methods for $2\times$, $3\times$, and $4\times$ upscaling is given in Fig. 13. Per each method and upscaling factor we measure the average PSNR and SSIM over the images of Set5 and Set14 (the horizontal coordinate in Fig. 13) vs. the average runtime (the vertical coordinate in Fig. 13). The size of each point equals to $10 \cdot \sigma_{\text{PSNR}}^{\text{err}}$ and $10^3 \cdot \sigma_{\text{SSIM}}^{\text{err}}$, where σ^{err} is the standard error of the PSNR/SSIM. As can be inferred from Fig. 13, RAISR is competitive with the state-of-the-art methods in terms of these quality measures while being much faster. More specifically, for all upscaling factors, RAISR has similar restoration quality to the fast version of A+ (using 16 atoms), and it outperforms NE+LS, Zeyde et al., GR, and ANR. Notice that the more complex version of A+ (the one that uses 1024 atoms) and SRCNN performs better than RAISR, but with the cost of increased computations.

In terms of runtime, RAISR is the fastest method by far, demonstrating that high-quality results can be achieved without

sacrificing the computational complexity. Following Fig. 13, our implementation is about one to two orders of magnitude faster than the baseline methods (please refer to the supplementary material for detailed results). The runtime is evaluated on a 3.4 GHz 6-Core Xeon desktop computer⁴.

In the case of upscaling by a factor of 2, a visual comparison between RAISR, bicubic and the state-of-the-art A+ and SRCNN is provided in Fig. 14. As can be seen, the restoration quality of RAISR for both Zebra, Flowers and Baby images is competitive with much more complex algorithms. The ability of RAISR to restore continues edges is shown on the Zebra image. In addition, the effective reconstruction of fine details is demonstrated via the Flowers and Baby images.

Next, we test the effect of sharpening the HR training images. To this end, we learn a new set of filters using the same

⁴Notably, we implemented a version of RAISR running on a mobile GPU, which runs at speeds of over 200Mpix/s. The desktop GPU version of the same can be expected to perform an order of magnitude faster.

TABLE I
 QUANTIFYING THE BENEFIT OF THE BUILT-IN SHARPENING EFFECT. WE MEASURE THE AVERAGE PSNR AND SSIM (HIGHER IS BETTER) OVER SET5 AND SET14 IMAGES, ALONG WITH THE STANDARD-ERROR OF EACH QUALITY METRIC

Dataset	Scaling	Learning without sharpening				Learning with sharpening				Improvement	
		PSNR	$\sigma_{\text{PSNR}}^{\text{err}}$	SSIM	$\sigma_{\text{SSIM}}^{\text{err}}$	PSNR	$\sigma_{\text{PSNR}}^{\text{err}}$	SSIM	$\sigma_{\text{SSIM}}^{\text{err}}$	PSNR	SSIM
Set5	2×	35.913	1.374	0.947	0.015	36.153	1.337	0.951	0.015	0.241	0.004
	3×	32.061	1.392	0.895	0.019	32.211	1.363	0.901	0.019	0.150	0.006
	4×	29.689	1.479	0.839	0.021	29.837	1.481	0.848	0.021	0.147	0.009
Set14	2×	31.980	0.998	0.900	0.015	32.127	1.024	0.902	0.015	0.147	0.003
	3×	28.764	0.985	0.809	0.026	28.860	0.999	0.812	0.027	0.096	0.003
	4×	26.912	0.908	0.732	0.032	27.002	0.921	0.738	0.033	0.090	0.006



Fig. 15. Visual illustration of the evolution of Butterfly image throughout the different blocks of RAISR, along with a demonstration of the built-in sharpening effect. (a) The original HR image, (b) Bicubic interpolation by a factor of 2 of the input LR image, (c) Filtered image – the filters do not include a built-in sharpening, (d) Filtered image – the filters map the initial interpolated image to its sharpened HR version, (e) blending result of the image in (b) with the one in (c), and (f) blending result of the image in (b) with the one in (d).

collection of 10,000 images, however this time without applying the pre-sharpening step. Following Table I, this step indeed leads to an improvement both in PSNR and SSIM, emphasizing its potential. Motivated by these results, a promising future direction could be to test this concept as a way to improve various other state-of-the-art methods, e.g. A+ and SRCNN.

Fig. 15 illustrates how the Butterfly image (taken from Set5) evolves throughout the different stages of RAISR, demonstrating the effect of the sharpening and the need for the blending step. The LR image is first upsampled by a factor of 2 using the bicubic interpolation (see Fig. 15(b)). Then, the pre-learned filters are applied: Fig. 15(c) shows the result when the filters are



Fig. 16. Visual comparison for upscaling by a factor of 2 the JPEG-compressed Comic image (taken from Set14). (a) Bilinear interpolation, (b) RAISR with learning stage as in Section VI-A which does not include compression of the training LR images, (c) RAISR with learning that involves compression of the training LR images.

learned without sharpening, while Fig. 15(d) illustrates the case where the filters have a built-in sharpening effect. As can be seen, the image in Fig. 15(d) is indeed sharper than the one in Fig. 15(c), and both have noticeable halos along the edges in addition to an amplification of the noise especially in flat areas. These artifacts are then reduced by the blending step, as depicted in Fig. 15(e) and 15(f). In this case, the improvement achieved by the pre-sharpening during training is 0.44 dB and 0.009 in terms of PSNR and SSIM, respectively.

To conclude, we tested the upscaling performance of the proposed method on popular datasets and for various upscaling factors. As demonstrated, RAISR is much faster than the leading algorithms, while achieving a competitive restoration performance.

B. All in One Enhancement

The ability of the proposed algorithm to handle real-world scenarios is evaluated as well. More specifically, we test the ability of RAISR to upscale an image while reducing compression artifacts and improving the overall contrast of the image. The results are obtained using the scheme that is given in Fig. 11. Differently from the conventional SISR problem (Section VI-A), in this case we do not apply the back-projection step. In addition, in order to further reduce computations, we choose the bilinear as the initial interpolator. As a result, we manage to reduce the runtime by about a factor of 2, i.e., *this version of RAISR is more than two orders of magnitude faster than the state-of-the-art methods A+ and SRCNN*.

In the learning stage, we downsampled the training images (taken from BSDS300 database [38]) using the bicubic interpolation, followed by a compression step using JPEG, with a quality parameter that is set to 85 out of 100. In order to achieve contrast enhancement effect, the target HR images are pre-processed using the proposed CT-based DoG sharpener (see Fig. 12), where the blending is done by counting the

modified bits of the CT descriptor, which allows the amplification of a wide range of frequencies. Naturally, the wider the Gaussian filters that are used in the DoG scheme, the wider the range of frequencies that are amplified. As such, we suggest using 2 DoG layers: The first layer applies a narrow Gaussian filter, with $\sigma \approx 0.85$ (i.e., the separable one dimensional $[1, 2, 1]/4$ filter), while the second is a much larger Gaussian filter of size 64×64 , with $\sigma = 8.5$ (i.e., we set $\alpha = 10$). In both cases we use $\rho = 55$. Note that the hyper parameters of RAISR (filter-size, hashing parameters, etc.) are the same as in Section VI-A.

Before testing RAISR in a real world scenario, where the degradation model is unknown, we show the benefits of compressing the training LR images on a synthetic example. To this end, following Fig. 16, we degrade the Comic image (taken from Set14) by downscaling it by a factor of 2 in each axis and then apply a JPEG compression on the result. As can be seen in Fig. 16(b), when the filters are learned without compressing the training LR images we obtain a sharp result but with the cost of undesired amplification of compression artifacts. On the other hand, when the learning stage includes compression, the upscaled outcome has less artifacts without the loss of sharpness, as illustrated in Fig. 16(c). This result should not surprise us as we take into account the degradation model that includes both downscaling (by a factor of 2) and compression. Quantitatively, compared to a learning stage that does not involve such a compression, in terms of PSNR we achieve an average improvement of 0.43 dB and 0.34 dB for Set5 and Set14, respectively. Similarly, we obtain higher SSIM score, where the average improvement is of 0.017 and 0.012 for Set5 and Set14, respectively.

Fig. 17 shows the superiority of handling compression artifacts along with the visually pleasant outcome of our contrast-enhancement learning scheme. Upscaling results of a cropped and compressed version of the Resolution-Chart image is given in Fig. 17(a)–(c), obtained by applying SRCNN, A+, and

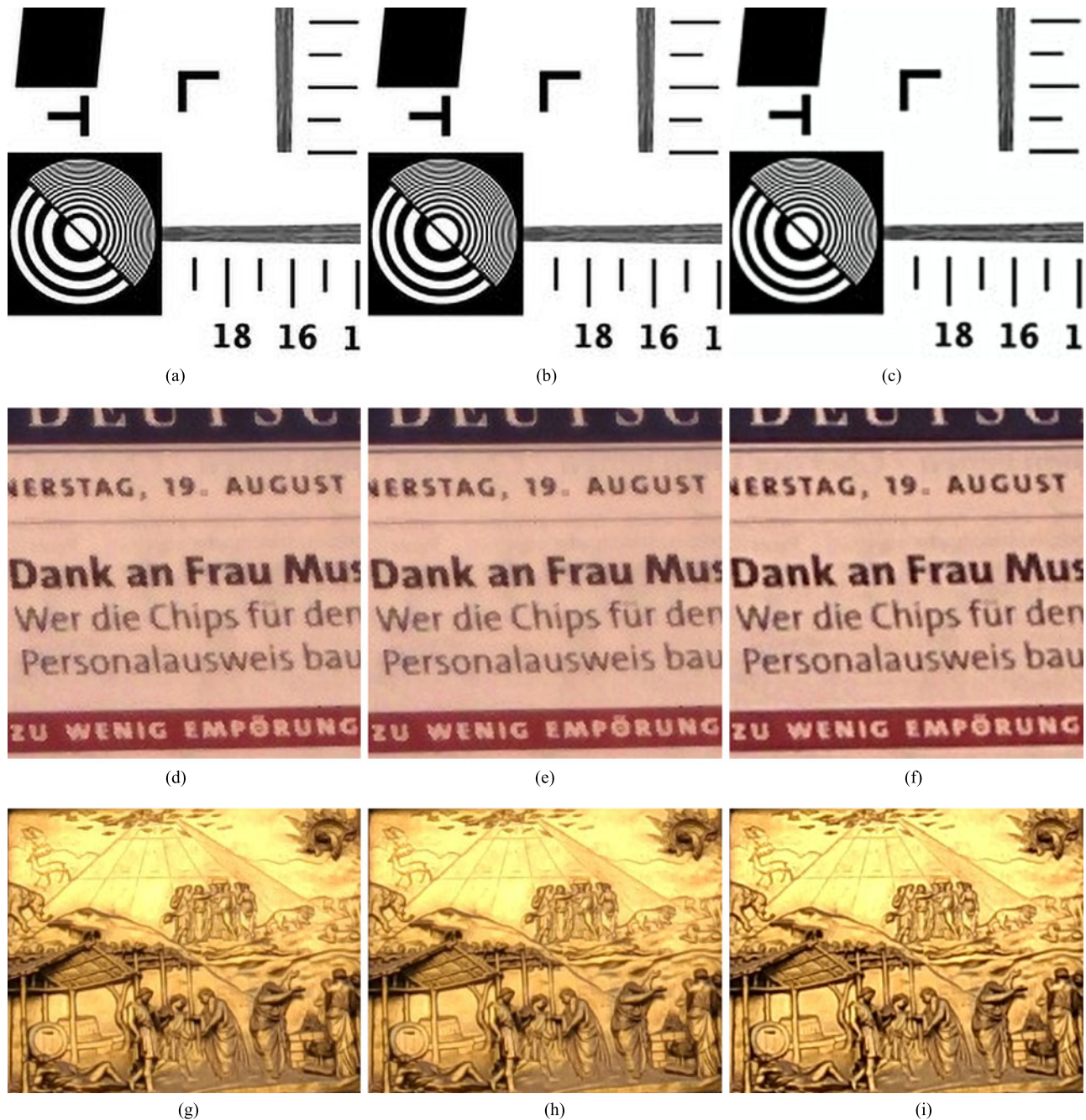


Fig. 17. Visual comparison for upscaling by a factor of 2 of Resolution-Chart, Newspaper and Painting images. (a,d,g) SRCNN, (b,e,h) A+ with 1024 atoms, and (c,f,i) RAISR. In this case, RAISR filters map compressed LR images to their contrast-enhanced HR versions.

RAISR, respectively. As can be seen, RAISR reduces some of the compression artifacts (especially around the digits and in between the circles, located in the bottom-right and left parts of the image, respectively). Despite the reduction of compression artifacts, RAISR successfully keeps the desired sharpness property (notice the effective restoration of the very fine details of the circles in the left part of the image).

The ability of RAISR to handle compressed, blurred with an unknown (possibly motion-blur) kernel, and noisy images is shown in Fig. 17(d)–(i). As can be inferred, the blending step and the built-in suppression of compression artifacts increase

the robustness of the algorithm, especially when tackling noisy images. Furthermore, the contrast-enhancement effect leads to better looking images (notice the contrast-enhanced letters in the Newspaper image and the amplification of the fine details in the Painting image). Additional examples on real images are shown in the <https://goo.gl/D0ETxG>.

To summarize, we demonstrated the ability of RAISR to achieve a high-quality restoration in various scenarios. Our experiments indicate that the proposed method stands in a line with the best algorithms that are currently available, while being about 2 orders of magnitude faster.

VII. CONCLUSION

In this paper we proposed a rapid and accurate learning-based approach for single image super-resolution (called RAISR). The suggested algorithm requires a relatively small set of training images to produce a low-complexity mechanism for increasing the resolution of any arbitrary image not seen before. The core idea behind RAISR is to learn a mapping from LR images to their HR versions. The mapping is done by filtering a “cheap” upsampled version of the LR image with a set of filters, which are designed to minimize the Euclidean distance between the input and ground-truth images.

More specifically, RAISR suggests a highly efficient locally adaptive filtering process, where the low-complexity is kept thanks to an appealing hashing scheme: In the learning stage, a database of images is divided into buckets (“cheap” clusters) of patches that share similar geometry, and a filter is learned per each bucket. In the application side, based on the geometry of the LR patch (hash-table key), the relevant pre-learned filter (hash-table entry) is chosen and applied on the patch.

Moreover, artifact-free results are obtained by blending the initial (cheap) estimation of the HR image with the filtered one. This step is based on the observation that in flat areas the reconstruction of the cheap upscaler is effective (there are no fine details or edges to be recovered in these areas). We harness the efficient CT descriptor for this task. As such, with a negligible computational cost, an accurate reconstruction is achieved. Note that two different blenders are suggested; the first allows the amplification of high frequencies only (using the “randomness” mask), while the second allows the enhancement of a wide range of frequencies. Differently from the randomness mask, the latter blending mechanism counts the number of bits of the CT descriptors that were changed due to the filtering step.

As such, the additional complexity of RAISR over a very basic interpolation method (e.g. bilinear interpolation) is roughly the application of 2-3 linear filters on the image: The first filter leads to the hash-table-key, pointing to the second filter that enhances the quality of the patch (the pre-learned filter), and the last filter (which is almost negligible) is used for the blending step. Despite the extremely low computational cost, the restoration performance of RAISR is competitive with much more complex state-of-the-art methods. For example, in the case of A+ (which is one the fastest methods that currently available), the complexity of choosing the filter is linear with the size of the dictionary, where the complexity of our approach for choosing the filter is constant due to the hashing mechanism.

Furthermore, sharpening and suppression of compression artifacts are achieved by pre-processing the training images. A sharpening/contrast-enhancement effect is obtained by amplifying the details and contrast of the HR training images. Similarly, compression artifacts are handled by compressing the LR training images. This results in filters that are unique in profile, and perform the dual purpose of de-aliasing and sharpening in a single filter. The pre-learned filters tend to have a smoothing effect in the center of the kernel to perform aliasing removal, while containing ridges near the edges of the kernel which have the effect of sharpening that edge. As such, our learning

approach is able to learn an effective set of filters that would be nearly impossible to construct by hand.

Motivated by the sharpening effect that can be achieved, we suggested a novel sharpener, which is based on applying (separable) DoG filters on the image, which are capable of flexibly enhancing a wide range of frequencies. As a way to reduce artifacts (e.g. halos and noise amplification), we suggest using once again the CT blender, leading to an extremely efficient and effective sharpener.

We should note that the quality of the hashing mechanism is crucial. Therefore, a further study of the hash function is needed, and could possibly lead to improved results. In addition, we use a simple least squares solver to learn the filters, while improved results may be achieved by regularizing the learning with efficient priors. In a wider perspective, we wish to explore the ability of RAISR to cast as a boosting mechanism [39]–[46], where in this case the learned filters can be designed to map any input images (not only the “cheap” interpolated images) to the desired outputs.

APPENDIX

A. Back Projection

As a way to reduce artifacts and estimation errors, many methods (e.g., [14], [47]–[50]) harness the popular Iterative Back Projection (IBP) [51] as a global post processing step. Following the degradation model in Eq. (1), one can demand an equivalence between the LR image and the downsampled version of the estimated HR image. More formally, given an estimation of the HR image $\tilde{\mathbf{x}}$ (i.e., the output of RAISR), we suggest minimizing the following cost

$$\min_{\mathbf{x}} \|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2 \quad \text{s.t.} \quad \mathbf{z} = \mathbf{D}_s \mathbf{H} \mathbf{x}, \quad (11)$$

leading to the desired HR estimation $\hat{\mathbf{x}}$. The solution of Eq. (11) can be obtained by applying several iterations of gradient descent. We found it helpful to apply this post-processing step in the case of the conventional SISR process (the results in Section VI-A).

REFERENCES

- [1] R. Keys, “Cubic convolution interpolation for digital image processing,” *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 29, no. 6, pp. 1153–1160, Dec. 1981.
- [2] H. Hou and H. Andrews, “Cubic splines for image interpolation and digital filtering,” *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 26, no. 6, pp. 508–517, Dec. 1978.
- [3] A. Buades, B. Coll, and J. M. Morel, “A review of image denoising algorithms, with a new one,” *Multiscale Model. Simul.*, vol. 4, no. 2, pp. 490–530, 2005.
- [4] H. Takeda, S. Farsiu, and P. Milanfar, “Kernel regression for image processing and reconstruction,” *IEEE Trans. Image Process.*, vol. 16, no. 2, pp. 349–366, Feb. 2007.
- [5] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, “Image denoising by sparse 3-D transform-domain collaborative filtering,” *IEEE Trans. Image Process.*, vol. 16, no. 8, pp. 2080–2095, Aug. 2007.
- [6] H. Takeda, P. Milanfar, M. Protter, and M. Elad, “Super-resolution without explicit subpixel motion estimation,” *IEEE Trans. Image Process.*, vol. 18, no. 9, pp. 1958–1975, Sep. 2009.

- [7] M. Aharon, M. Elad, and A. Bruckstein, "The K-SVD: An algorithm for designing of overcomplete dictionaries for sparse representation," *IEEE Trans. Signal Process.*, vol. 54, no. 11, pp. 4311–4322, Nov. 2006.
- [8] J. Yang, J. Wright, T. Huang, and Y. Ma, "Image super-resolution as sparse representation of raw image patches," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2008, pp. 1–8.
- [9] R. Zeyde, M. Elad, and M. Protter, "On single image scale-up using sparse-representations," in *Curves and Surfaces*. New York, NY, USA: Springer, 2012, pp. 711–730.
- [10] R. Timofte, V. De, and L. Van Gool, "Anchored neighborhood regression for fast example-based super-resolution," in *Proc. IEEE Int. Conf. Comput. Vis.*, Dec. 2013, pp. 1920–1927.
- [11] R. Timofte, V. De Smet, and L. Van Gool, "A+: Adjusted anchored neighborhood regression for fast super-resolution," in *Proc. Asian Conf. Comput. Vis.*, 2014, pp. 111–126.
- [12] Y. Romano, M. Protter, and M. Elad, "Single image interpolation via adaptive nonlocal sparsity-based modeling," *IEEE Trans. Image Process.*, vol. 23, no. 7, pp. 3085–3098, Jul. 2014.
- [13] D. Zoran and Y. Weiss, "From learning models of natural image patches to whole image restoration," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2011, pp. 479–486.
- [14] J. Yang, Z. Wang, Z. Lin, S. Cohen, and T. Huang, "Coupled dictionary training for image super-resolution," *IEEE Trans. Image Process.*, vol. 21, no. 8, pp. 3467–3478, Aug. 2012.
- [15] T. Peleg and M. Elad, "A statistical prediction model based on sparse representations for single image super-resolution," *IEEE Trans. Image Process.*, vol. 23, no. 6, pp. 2569–2582, Jun. 2014.
- [16] C. Dong, C. C. Loy, K. He, and X. Tang, "Learning a deep convolutional network for image super-resolution," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 184–199.
- [17] D. Chao, C. C. Loy, H. Kaiming, and T. Xiaoou, "Image super-resolution using deep convolutional networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 2, pp. 295–307, Feb. 2016.
- [18] D. Dai, R. Timofte, and L. Van Gool, "Jointly optimized regressors for image super-resolution," in *Proc. Comput. Graph. Forum*, vol. 34, no. 2, pp. 95–104, 2015.
- [19] Y. LeCun *et al.*, "Backpropagation applied to handwritten zip code recognition," *Neural Comput.*, vol. 1, no. 4, pp. 541–551, 1989.
- [20] R. Zabih and J. Woodfill, "Non-parametric local transforms for computing visual correspondence," in *Proc. Eur. Conf. Comput. Vision*, 1994, pp. 151–158.
- [21] A. Polesel, G. Ramponi, and V. J. Mathews, "Image enhancement via adaptive unsharp masking," *IEEE Trans. Image Process.*, vol. 9, no. 3, pp. 505–510, 2000.
- [22] D. Marr and E. Hildreth, "Theory of edge detection," *Proc. R. Soc. London B, Biol. Sci.*, vol. 207, no. 1167, pp. 187–217, 1980.
- [23] H. Winnemöller, J. E. Kyprianidis, and S. C. Olsen, "XDoG: An extended difference-of-Gaussians compendium including advanced image stylization," *Comput. Graph.*, vol. 36, no. 6, pp. 740–753, 2012.
- [24] B. Zhang and J. P. Allebach, "Adaptive bilateral filter for sharpness enhancement and noise removal," *IEEE Trans. Image Process.*, vol. 17, no. 5, pp. 664–678, May 2008.
- [25] K. He, J. Sun, and X. Tang, "Guided image filtering," in *Proc. Eur. Conf. Comput. Vision*, 2010, pp. 1–14.
- [26] H. Talebi and P. Milanfar, "Nonlocal image editing," *IEEE Trans. Image Process.*, vol. 23, no. 10, pp. 4460–4473, Oct. 2014.
- [27] A. Kheradmand and P. Milanfar, "Nonlinear structure-aware image sharpening with difference of smoothing operators," *Frontiers ICT*, vol. 2, 2015, Art. no. 22.
- [28] X. Liu, G. Cheung, and X. Wu, "Joint denoising and contrast enhancement of images using graph laplacian operator," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2015, pp. 2274–2278.
- [29] C. Tomasi and R. Manduchi, "Bilateral filtering for gray and color images," in *Proc. IEEE Int. Conf. Comput. Vis.*, Jan. 1998, pp. 836–846.
- [30] F. Durand and J. Dorsey, "Fast bilateral filtering for the display of high-dynamic-range images," in *Proc. ACM Trans. Graph.*, vol. 21, no. 3, 2002, pp. 257–266.
- [31] R. Barrett *et al.*, *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods*. Philadelphia, PA, USA: SIAM, 1994.
- [32] S.-C. Jeong and B. C. Song, "Training-based super-resolution algorithm using k-means clustering and detail enhancement," in *Proc. Eur. Signal Process. Conf.*, 2010, pp. 1791–1795.
- [33] G. Yu, G. Sapiro, and S. Mallat, "Solving inverse problems with piecewise linear estimators: From gaussian mixture models to structured sparsity," *IEEE Trans. Image Process.*, vol. 21, no. 5, pp. 2481–2499, May 2012.
- [34] V. Papyan and M. Elad, "Multi-scale patch-based image restoration," *IEEE Trans. Image Process.*, vol. 25, no. 1, pp. 249–261, Jan. 2016.
- [35] X. Feng and P. Milanfar, "Multiscale principal components analysis for image local orientation estimation," in *Proc. 36th Asilomar Conf. Signals, Syst. Comput.*, Nov. 2002, pp. 478–482.
- [36] M. Bevilacqua, A. Roumy, C. Guillemot, and M. L. Alberi-Morel, "Low-complexity single-image super-resolution based on nonnegative neighbor embedding," in *Proc. Brit. Mach. Vis. Conf.*, 2012, pp. 1–10.
- [37] H. Chang, D.-Y. Yeung, and Y. Xiong, "Super-resolution through neighbor embedding," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2004, <http://ieeexplore.ieee.org/document/1315043>
- [38] D. Martin, C. Fowlkes, D. Tal, and J. Malik, "A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2001, vol. 2, pp. 416–423.
- [39] H. Talebi, X. Zhu, and P. Milanfar, "How to SAIF-ly boost denoising performance," *IEEE Trans. Image Process.*, vol. 22, no. 4, pp. 1470–1485, Apr. 2013.
- [40] P. Milanfar, "A tour of modern image filtering: New insights and methods, both practical and theoretical," *IEEE Signal Process. Mag.*, vol. 30, no. 1, pp. 106–128, Jan. 2013.
- [41] Y. Romano and M. Elad, "Improving K-SVD denoising by post-processing its method-noise," in *Proc. IEEE Int. Conf. Image Process.*, 2013, pp. 435–439.
- [42] H. Talebi and P. Milanfar, "Global image denoising," *IEEE Trans. Image Process.*, vol. 23, no. 2, pp. 755–768, Feb. 2014.
- [43] A. Kheradmand and P. Milanfar, "A general framework for regularized, similarity-based image restoration," *IEEE Trans. Image Process.*, vol. 23, no. 12, pp. 5136–5151, Dec. 2014.
- [44] Y. Romano and M. Elad, "Patch-disagreement as a way to improve K-SVD denoising," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2015, pp. 1280–1284.
- [45] Y. Romano and M. Elad, "Boosting of image denoising algorithms," *SIAM J. Imag. Sci.*, vol. 8, no. 2, pp. 1187–1219, 2015.
- [46] G. Ghimpeanu, T. Batard, M. Bertalmio, and S. Levine, "A decomposition framework for image denoising algorithms," *IEEE Trans. Image Process.*, vol. 25, no. 1, pp. 388–399, Jan. 2016.
- [47] S. C. Park, M. K. Park, and M. G. Kang, "Super-resolution image reconstruction: A technical overview," *IEEE Signal Process. Mag.*, vol. 20, no. 3, pp. 21–36, May 2003.
- [48] W. Dong, L. Zhang, G. Shi, and X. Wu, "Nonlocal back-projection for adaptive image enlargement," in *Proc. IEEE Int. Conf. Image Process.*, 2009, pp. 349–352.
- [49] L. He, H. Qi, and R. Zaretzki, "Beta process joint dictionary learning for coupled feature spaces with application to single image super-resolution," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2013, pp. 345–352.
- [50] R. R. Makwana and N. D. Mehta, "Single image super-resolution via iterative back projection based Canny edge detection and a Gabor filter prior," *Int. J. Soft Comput. Eng.*, vol. 3, pp. 2231–2307, 2013.
- [51] M. Irani and S. Peleg, "Improving resolution by image registration," *CVGIP: Graph. Models Image Process.*, vol. 53, no. 3, pp. 231–239, 1991.

Authors' photographs and biographies not available at the time of publication.