

Multigrid Backprojection Super-Resolution and Deep Filter Visualization

Pablo Navarrete Michelini, Hanwen Liu and Dan Zhu

BOE Technology Group Inc., Ltd.
Beijing, China

Abstract

We introduce a novel deep-learning architecture for image upscaling by large factors (e.g. $4\times$, $8\times$) based on examples of pristine high-resolution images. Our target is to reconstruct high-resolution images from their downscale versions. The proposed system performs a multi-level progressive upscaling, starting from small factors ($2\times$) and updating for higher factors ($4\times$ and $8\times$). The system is recursive as it repeats the same procedure at each level. It is also residual since we use the network to update the outputs of a classic upscaler. The network residuals are improved by Iterative Back-Projections (IBP) computed in the features of a convolutional network. To work in multiple levels we extend the standard back-projection algorithm using a recursion analogous to Multi-Grid algorithms commonly used as solvers of large systems of linear equations. We finally show how the network can be interpreted as a standard upsampling-and-filter upscaler with a space-variant filter that adapts to the geometry. This approach allows us to visualize how the network learns to up-scale. Finally, our system reaches state of the art quality for models with relatively few number of parameters.

Introduction

In this work, we focus on the problem of image upscaling using convolutional networks. Upscaling signals by integer factors (e.g. $2\times$, $3\times$) is understood in classical interpolation theory as two sequential processes: upsample (insert zeros) and filter (Proakis and Manolakis 2007; Mallat 1998). Standard upscaler algorithms, such as Bicubic or Lanczos, find high-resolution images with a narrow frequency content by using fixed low-pass filters. Similar to the classic upscaling model, the image acquisition can be modeled as low-pass filtering a high resolution image and then downsample the result (drop pixels). In test scenarios often used in benchmarks we actually know the exact downscaling model, e.g. Bicubic downscaler. The Iterative Back-Projection (IBP) algorithm (Irani and Peleg 1991) is often used to enforce the downscaling model for a given upscaler and get closer to the original image.

More advanced upscalers follow geometric principles to improve image quality. For example, *edge-directed interpolation* uses adaptive filters to improve edge smoothness (Al-gazi, Ford, and Potharlanka 1991; Li and Orchard 2001), or *bandlet* methods use both adaptive upsampling and filtering (Mallat and Peyre 2007). More recently, machine learning

has been able to use examples pairs of high and low resolution images to estimate the parameters of upscaling systems (Park, Park, and Kang 2003). In some cases, the optimization approach of machine learning hides the connection with classical interpolation theory, e.g. sparse representation with dictionaries (Yang et al. 2008; 2010). In other cases, the adaptive filter approach is explicit, e.g. RAISR (Romano, Isidoro, and Milanfar 2016).

Upscaling using convolutional networks started with SRCNN (Dong et al. 2014a; 2015) motivated by the success of deep-learning methods in image classification tasks (LeCun, Bengio, and Hinton 2015) and establishing a strong connection with sparse coding methods (Yang et al. 2008; 2010). SRCNN has later been improved, most notably by EDSR (Lim et al. 2017) and DBPN (Haris, Shakhnarovich, and Ukita 2018). Our system shares the convolutional network approach but follows a different motivation. Namely, we aim to reveal a strong connection between convolutional-networks and classical image upscaling. By doing so, we can recover the classic interpretation of upsampling and filter and visualize what is the network doing pixel by pixel. Thus, we aim to prove that convolutional networks are a natural and convenient choice for Super-Resolution (SR) tasks.

Our main contributions are:

- We extend the IBP algorithm to a **multi-level IBP**.
- We prove that our algorithm **works as well as classic IBP**, based on an unrealistic model.
- We introduce a **new network architecture** that overcomes the unrealistic model and allows us to learn both upscaling and downscaling.
- We introduce a novel **algorithm to analyze** the linear components of the network.
- We show how to **interpret the network** as a standard upscaler with adaptive filters.

Related Work

We consider the following two architectures as the most similar to our system:

- **Multi-Scale Laplacian Super-Resolution (MSLapSR)** (Lai et al. 2017b): Our system is inspired by MSLapSR to progressively upscale images using a classic upscaler and network updates, before back-projections that improve network updates in lower-resolutions. Both MSLapSR and our systems include *analysis* and *synthesis* networks to convert images to latent space and vice versa. Both systems share parameters at each scale. Our system differs mostly in the use of back-projections, that cannot be removed to recover MSLapSR because of the particular structure of our *upscaler* network module.
- **Deep Back-Projection Network (DBPN)** (Haris, Shakhnarovich, and Ukita 2018): To the extent of our knowledge, this is the first reference to use IBP in a network architecture. It is also the state-of-art in terms of image quality, surpassing EDSR(Lim et al. 2017), former winner of NTIRE 2017 SR Challenge (Timofte et al.). Their approach to use back-projections is different than ours because: first, their system is not multi-scale; and second, they iterate down and up projections. Our multi-scale architecture requires less number of parameters, because we reuse modules at every scale, and it is more flexible, because many upscaling factors can be achieved with the same modules. Also, our system is built upon an algorithm that is proven to converge, whereas, to the extent of our knowledge, mixing up and down projections has no convergence guarantees.

Our contributions add to the following lines of research:

- **Recursive Architectures:** Extensive work has been done on recursive CNN architectures inspired on iterative procedures. Our system belongs to this line of work since we propose an architecture based on IBP and multi-grid iterations. Of particular interest is the research along this line related to SR and enhancements problems. In (Kruse, Rother, and Schmidt 2017), for example, a *half quadratic splitting* iteration with a CNN denoiser prior is proposed for image denoising, deblurring and SR. Most recently, (Kokkinos and Lefkimmatis 2018) take a similar approach for image demosaicking. Next, (Gong et al. 2018) propose to learn an optimizer for image deconvolution by recurrently incorporating CNNs into a gradient descent scheme, which could be complementary to our approach. In (Zhang et al. 2017), an energy minimization iteration is performed using a CNN model for image deconvolution, which is close to the SR problem. Finally, the work of (Diamond et al. 2017) follows a similar motivation than ours by proposing a framework to incorporate knowledge of the image formation into CNNs, applied to denoising, deblurring and compressed sensing.

All these approaches differ to ours in the fact that our recursion iterates back and forth between different resolutions. In numerical methods, such iterations have been used by two types of linear equation solvers: *multi-grid* and *domain decomposition* methods (Trottenberg and Schuller 2001; Widlund and Toselli 2004). We follow the multi-grid approach, that can use inductive arguments to study convergence, and leads to specific processing workflows to move intermediate results between

scales (see for example the W-cycle in Figure 2). Systems like MS-DenseNet (Huang et al. 2017a) also move back and forth between scales with more simplified workflows but they were not designed based on classical methods. The connection to classical methods gives us a justification of these workflows that, otherwise, would be arbitrary (e.g. why traversing scales with V or W workflows in Figure 2)? which workflow is better?) We will show, for example, how different numbers of back-projections (related to depth) make outputs sharper, same as in IBP. Our main contribution here is to devise a new algorithm that introduces the multi-grid recursion into IBP. It is a different algorithm than IBP, that we prove to converge at the same rate, and then extend to a network architecture. Finally, this effort pays back since our recursion works well in experiments, with no other method reaching the same quality with the same number of parameters.

- **Network Visualization:** A major direction of research in deep-learning is how to visualize the inner processing of a given architecture (Zhang and Zhu 2018). Among these, a line of research on *feature visualization* studies what does a network detect (Olah, Mordvintsev, and Schubert 2017). Feature visualization can give example inputs that cause desired behaviors, separating image areas causing behavior from those that only relate to the causes. Another line of work on *attribution* studies how does a network assembles these individual pieces to arrive at later decisions, or why these decisions were made (Olah et al. 2018). Our visualization technique belongs to the latter because we can show how the input pixels are assembled into a particular output pixel. We target the SR problem where there is extensive knowledge of non-adaptive filters (e.g. linear, bicubic, etc.) built upon signal processing theory. The main novelty of our technique is that it provides an alternative system to replace the network for a given input. The new system generates the exact same outputs of the network from the same input images but, unlike the network, it is fully interpretable in the sense that we know what to expect from its parameters.

Multigrid Backprojections

A simple and common model for the downscaling process is

$$X = (Y * g) \downarrow s, \quad (1)$$

where Y is the high-resolution source, X is the low-resolution result, g is a blurring kernel and $\downarrow s$ is a down-sampling by factor s .

Model (1) gives additional information about the unknown high-resolution image and narrows down the search space from all possible images to images that downscaled with model (1) recover the low-resolution input image. This is the motivation behind the classic IBP algorithm (Irani and Peleg 1991). Given model (1) and an upscaled image Y , the IBP algorithm iterates:

$$e(Y_k) = X - (Y_k * g) \downarrow s \quad (2)$$

$$Y_{k+1} = Y_k + e(Y_k) \uparrow s * p. \quad (3)$$

Here, $e(Y_k)$ is the mismatch error at low-resolution, g and p are blurring and upscaling filters, respectively. The iteration

Algorithm 1 Multi-Grid Back-Projection (MGBP)

 $MGBP(X, \mu, L)$:

Input: Input image X .
Input: Integers $\mu \geq 0$ and $L \geq 1$.
Output: Images $Y_k, k = 2, \dots, L$.
 1: $Y_1 = X$
 2: **for** $k = 2, \dots, L$ **do**
 3: $u = (Y_{k-1} \uparrow s) * p$
 4: $Y_k = BP_k^\mu(u, Y_1, \dots, Y_{k-1})$
 5: **end for**

 $BP_k^\mu(u, Y_1, \dots, Y_{k-1})$:

Input: Image u , level index k , number of steps μ .
Input: Images Y_1, \dots, Y_{k-1} (only for $k > 1$).
Output: Updated image u
 1: **if** $k > 1$ **then**
 2: **for** $step = 1, \dots, \mu$ **do**
 3: $d = BP_{k-1}^\mu((u * g) \downarrow s, Y_1, \dots, Y_{k-2})$
 4: $u = u + (Y_{k-1} - d) \uparrow s * p$
 5: **end for**
 6: **end if**

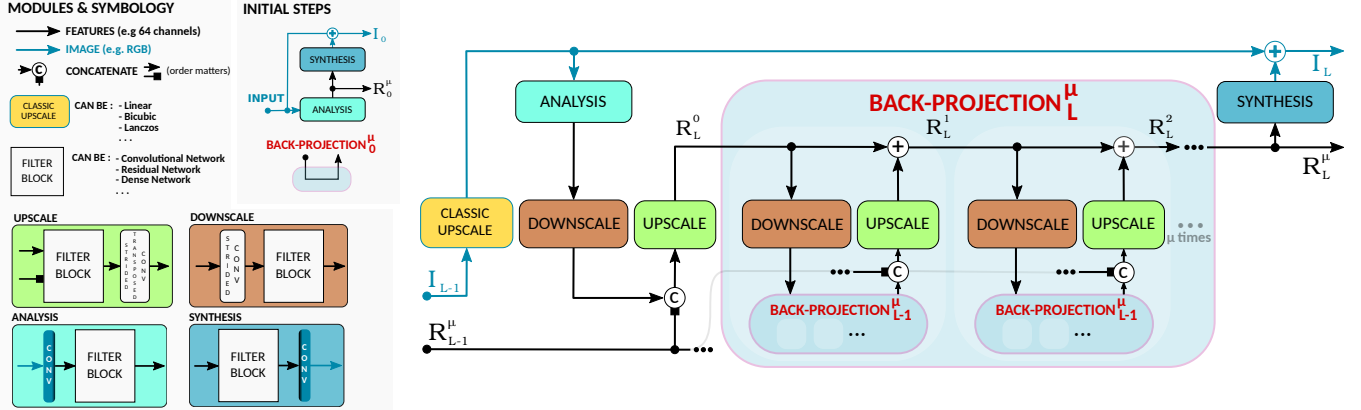


Figure 1: Multigrid Back-Projection (MGBP) system architecture. Lines indicate a number of image channels moving to different processing blocks, and line colors indicate the number of channels. Blue represents 3 channels for RGB images, and black represents the number of features managed by convolutional networks (latent space). At level k the current upscaled image is indicated as I_k and residuals in latent space are indicated as R_k . We use *Analysis* and *Synthesis* modules to transfer images into latent space and vice versa. In *Initial Steps* we obtain the first pair of output image I_0 and residual R_0 , from the input image. Then, the main diagram shows how to obtain the pair I_L and R_L from the previous pair I_{L-1} and R_{L-1} . In red we show the Back-Projection module, which repeats μ back-projection steps recursively.

is proven to converge, to enforce model (1) at exponential rate (Irani and Peleg 1991).

To make IBP work for multiple scales we change model (1) to:

$$X = (\dots \underbrace{((Y * g) \downarrow s * g) \downarrow s \dots * g) \downarrow s}_{L \text{ times}} \dots) \quad (4)$$

This is not a common downscaling procedure in practice and might be the reason why multi-scale IBP has not been considered yet. We will later replace the downscaling by a network so that model (4) becomes flexible and is able to learn a direct downscaling or even more complex models.

Upscaling images with IBP is a two-step process: first, upscale an image; and second, improve it with IBP. This is reminiscent of the way a Full-Multigrid algorithm solves linear equations (Trottenberg and Schuller 2001). This is: first, find an approximate solution; and second, improve it by solving an equation for the approximation error. Both IBP and Multigrid iterate between different scales, but IBP only uses two levels whereas Multigrid recursively move to coarser grids. We use the same strategy as in Multigrid to define a so-called *Multi-Grid Back-Projection* algorithm as shown in Algorithm 2. Here, back-projections recursively return to the lowest-resolution enforcing the downscaling

model at each scale.

For $L = 2$ the MGBP algorithm and model (4) are equivalent to the original IBP (Irani and Peleg 1991) and model (1), and thus converges at exponential rate. In section A of the supplementary material we prove convergence for $L > 2$. Basically, the algorithm inherits the exponential rate convergence from the two-level case through the recursion in Algorithm 2.

Network Architecture

We convert the Multi-Grid Back-Projection algorithm into a network structure as follows:

- Step 1: Use a classic method to upscale a low resolution image.
- Step 2: Transfer the upscale image I into latent space using a network *Analysis*(I).
- Step 3: In latent space we apply the recurrence in Algorithm 2 by changing:
 - $(u * g) \downarrow s$ into a network *Downscale*(u).
 - $(Y_{k-1} - d) \uparrow s * p$ into a network *Upscale*($[Y_{k-1}, d]$). Where $[Y_{k-1}, d]$ is the concatenation of features and replaces the subtraction. Thus, the *Upscaler* network re-

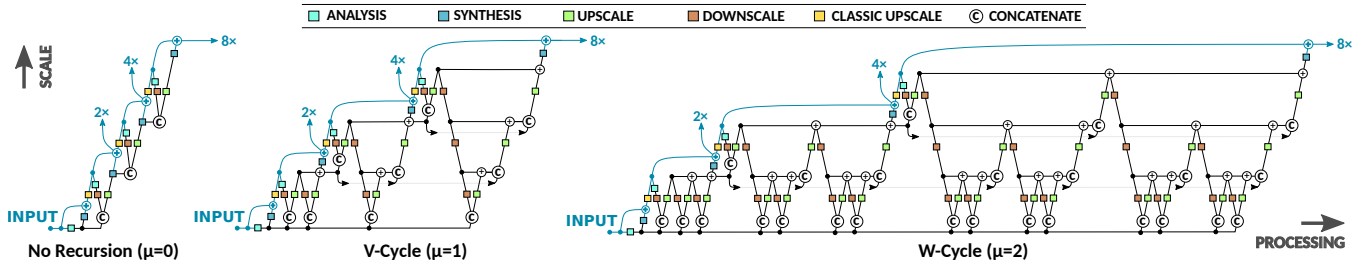


Figure 2: Multi-Grid Back-Projection (MGBP) recursion unfold from Figure 1 for different values of μ , and three levels to output $2\times$, $4\times$ and $8\times$ upscale images. Our system uses a recursion analogous to the Full-Multigrid algorithm to solve linear equations and leads to the well known workflows V-cycle for $\mu = 1$ and W-cycles for $\mu = 2$ (Trottenberg and Schuller 2001). After every upscaling step, the MGBP recursion sends the output down to every lower-resolution level in order to validate the downscaling model. The corrections are back-projected to higher resolutions.

ceives double the number of features in the input compared to the output.

Figure 1 shows the definition of our network architecture. The unfolded recursion is shown in Figure 2 for three different values of μ .

Deep Filter Visualization

Deep Learning architectures are highly non-linear, although much of their internal structure is linear (e.g. convolutions). We want to study the overall effect of the linear structure of the network. The general procedure is shown in Figure 3 and is as follows:

- An input image X passes through all layers of the network, and outputs an image Y .
- At each non-linear layer we record how much did the input change (layer gain).
- We replace (freeze) all non-linearities by the fix gain previously recorded. The overall system becomes linear, such that $Y = FX + R$.
- We obtain the *effective residual* R with an input $X = 0$ in the activation frozen system.
- We obtain the *effective filter* for a particular input pixel by using a δ input centered at the pixel location, and subtract the residual R from the output. In linear system terminology we are computing the *impulse response* of the system (Proakis and Manolakis 2007).

We can obtain explicit formulas for F and R if we consider a model of convolutional network as a sequence of linear and nonlinear layers:

$$z_n = W_n x_{n-1} + b_n \quad \text{and} \quad x_n = \sigma(z_n), \quad (5)$$

where x_n and z_n are vectorized features at layer n , after and before activations, respectively. The parameters of the network are the biases b_n and the sparse matrices W_n representing the convolutional operators (also applies for strided and transposed convolutions). For a given input image X , the input of the model is $x_0 = \text{vec}(X)$ (vectorized image). The output image Y , after n convolutional layers, is $\text{vec}^{-1}(x_n)$.

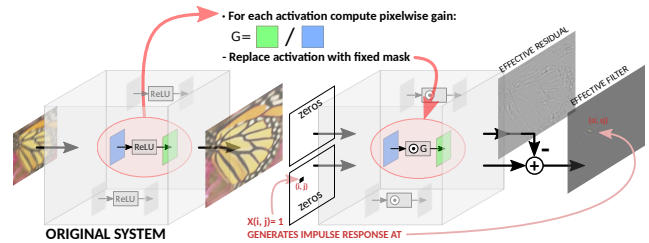


Figure 3: Activation freezing procedure to convert a network into a linear system. We can perform impulse response analysis to study the overall filter effect at each pixel location.

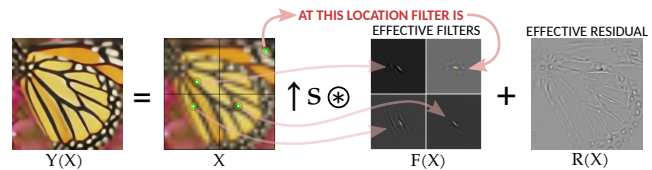


Figure 4: Interpretation of the upscaling network as a standard upscaler with adaptive filters.

Definition 1 (Activation Gain) The gain of an activation function σ is given by $G(x)[i] = \sigma(x[i])/x[i]$ and equal to 1 if $x[i] = 0$.

Theorem 1 (Activation Freeze) Let $\hat{W}_n = G(z_n)W_n$ and $\hat{b}_n = G(z_n)b_n$. Let

$$Q_0 = I, \quad Q_i = \prod_{k=n-i+1}^n \hat{W}_k, \quad \text{for } i = 1, \dots, n. \quad (6)$$

The output of the convolutional network is given by $x_n = F x_0 + R$, where $F = Q_n$ is the *effective filter* and $R = Q * \hat{b} = \sum_{k=0}^n Q_k \hat{b}_{n-k}$ is the *effective residual*.

The proof of the theorem is a direct consequence of $\sigma(x) = G(x)x$ and expansion of (5). Although this theorem only applies to sequential networks, it helps to show that the overall effective filter depends on all convolutional filters as well as biases (only through activations). Similarly, the effective

Table 1: **Quantitative evaluation** of different SR methods. In red color we show the best result, in blue the top-2 and brown the top-3 results for each column and scale. Methods are ordered by increasing number of parameters.

Algorithm	s	par [M]	Set14		BSDS100		Urban100		Manga109	
			PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
Bicubic	2	–	30.34	0.870	29.56	0.844	26.88	0.841	30.84	0.935
A+ (Timofte and Smet 2014)	2	–	32.40	0.906	31.22	0.887	29.23	0.894	35.33	0.967
FSRCNN (Dong, Loy, and Tang 2016)	2	0.01	32.73	0.909	31.51	0.891	29.87	0.901	36.62	0.971
SRCNN (Dong et al. 2014b)	2	0.06	32.29	0.903	31.36	0.888	29.52	0.895	35.72	0.968
MSLapSRN (Lai et al. 2017b)	2	0.22	33.28	0.915	32.05	0.898	31.15	0.919	37.78	0.976
Our	2	0.28	33.27	0.915	31.99	0.897	31.37	0.920	37.92	0.976
VDSR (Kim, Lee, and Lee 2016a)	2	0.67	32.97	0.913	31.90	0.896	30.77	0.914	37.16	0.974
LapSRN (Lai et al. 2017a)	2	0.81	33.08	0.913	31.80	0.895	30.41	0.910	37.27	0.974
DRCN (Kim, Lee, and Lee 2016b)	2	1.78	32.98	0.913	31.85	0.894	30.76	0.913	37.57	0.973
D-DBPN (Haris, Shakhnarovich, and Ukita 2018)	2	5.95	33.85	0.919	32.27	0.900	32.70	0.931	39.10	0.978
EDSR (Lim et al. 2017)	2	40.7	33.92	0.919	32.32	0.901	32.93	0.935	39.10	0.977
<hr/>										
Bicubic	4	–	26.10	0.704	25.96	0.669	23.15	0.659	24.92	0.789
A+ (Timofte and Smet 2014)	4	–	27.43	0.752	26.82	0.710	24.34	0.720	27.02	0.850
FSRCNN (Dong, Loy, and Tang 2016)	4	0.01	27.70	0.756	26.97	0.714	24.61	0.727	27.89	0.859
SRCNN (Dong et al. 2014b)	4	0.06	27.61	0.754	26.91	0.712	24.53	0.724	27.66	0.858
MSLapSRN (Lai et al. 2017b)	4	0.22	28.26	0.774	27.43	0.731	25.51	0.768	29.54	0.897
Our	4	0.28	28.43	0.778	27.42	0.732	25.70	0.774	30.07	0.904
VDSR (Kim, Lee, and Lee 2016a)	4	0.67	28.03	0.770	27.29	0.726	25.18	0.753	28.82	0.886
LapSRN (Lai et al. 2017a)	4	0.81	28.19	0.772	27.32	0.728	25.21	0.756	29.09	0.890
DRCN (Kim, Lee, and Lee 2016b)	4	1.78	28.04	0.770	27.24	0.724	25.14	0.752	28.97	0.886
D-DBPN (Haris, Shakhnarovich, and Ukita 2018)	4	10.4	28.82	0.786	27.72	0.740	26.54	0.795	31.18	0.914
EDSR (Lim et al. 2017)	4	43.1	28.80	0.788	27.71	0.742	26.64	0.803	31.02	0.915
<hr/>										
Bicubic	8	–	23.19	0.568	23.67	0.547	20.74	0.516	21.47	0.647
A+ (Timofte and Smet 2014)	8	–	23.98	0.597	24.20	0.568	21.37	0.545	22.39	0.680
FSRCNN (Dong, Loy, and Tang 2016)	8	0.01	23.93	0.592	24.21	0.567	21.32	0.537	22.39	0.672
SRCNN (Dong et al. 2014b)	8	0.06	23.85	0.593	24.13	0.565	21.29	0.543	22.37	0.682
MSLapSRN (Lai et al. 2017b)	8	0.22	24.57	0.629	24.65	0.592	22.06	0.598	23.90	0.759
Our	8	0.28	24.82	0.635	24.67	0.592	22.21	0.603	24.12	0.765
VDSR (Kim, Lee, and Lee 2016a)	8	0.67	24.21	0.609	24.37	0.576	21.54	0.560	22.83	0.707
LapSRN (Lai et al. 2017a)	8	0.81	24.44	0.623	24.54	0.586	21.81	0.582	23.39	0.735
D-DBPN (Haris, Shakhnarovich, and Ukita 2018)	8	23.2	25.13	0.648	24.88	0.601	22.83	0.622	25.30	0.799
EDSR (Lim et al. 2017)	8	43.1	24.94	0.640	24.80	0.596	22.47	0.620	24.58	0.778

residual depends on both convolutional filters and biases (both explicitly and through activations).

For the sake of simplicity, here we use our visualization technique from input to outputs of the network. Nevertheless, we note that we could start at any layer and stop at any posterior layer to study attributions within the network.

Experiments

We use one **single configuration** to test our system for $2\times$, $4\times$, and $8\times$ upscaling. We configure the *Analysis*, *Synthesis*, *Upscale* and *Downscale* modules in Figure 1 using 4-layer dense networks (Huang et al. 2017b) as filter-blocks. We use 48 features and growth rate 16 within dense networks. For classic upscaler we start with Bicubic and the upscaling filters are set as parameters to learn during training.

For **training** the system we fix the number of back-projections to $\mu = 2$ (W-cycle according to Figure 2). We train the system for $8\times$ upscaler with the multi-scale loss function introduced in (Lai et al. 2017b):

$$\mathcal{L}(Y, X; \theta) = \sum_{L=1,2,3} \sum_{k=1}^L \mathbb{E} [\|Y_k^L - X_k\|_1]. \quad (7)$$

We train our system with Adam optimizer and a learning rate initialized as 10^{-3} and square root decay. We use 128×128 patches (pieces of images) with batch size 16. The patches were sampled randomly from datasets DIV2K and Flickr2K, containing photographs of general natural scenes.

Comparisons

In Table 1 we compare PSNR and SSIM values for different methods. The two evaluation metrics measure the difference between an upscaler output and the original high-resolution image. Higher values are better in both cases. Roughly speaking, PSNR (range 0 to ∞) is a log-scale version of mean-square-error and SSIM (range 0 to 1) uses image statistics to better correlate with human perception. Full expressions are as follows:

$$PSNR(X, Y) = 10 \cdot \log_{10} \left(\frac{255^2}{MSE} \right), \quad (8)$$

$$SSIM(X, Y) = \frac{(2\mu_X\mu_Y + c_1)(2\sigma_{XY} + c_2)}{(\mu_X^2 + \mu_Y^2 + c_1)(\sigma_X^2 + \sigma_Y^2 + c_2)}, \quad (9)$$

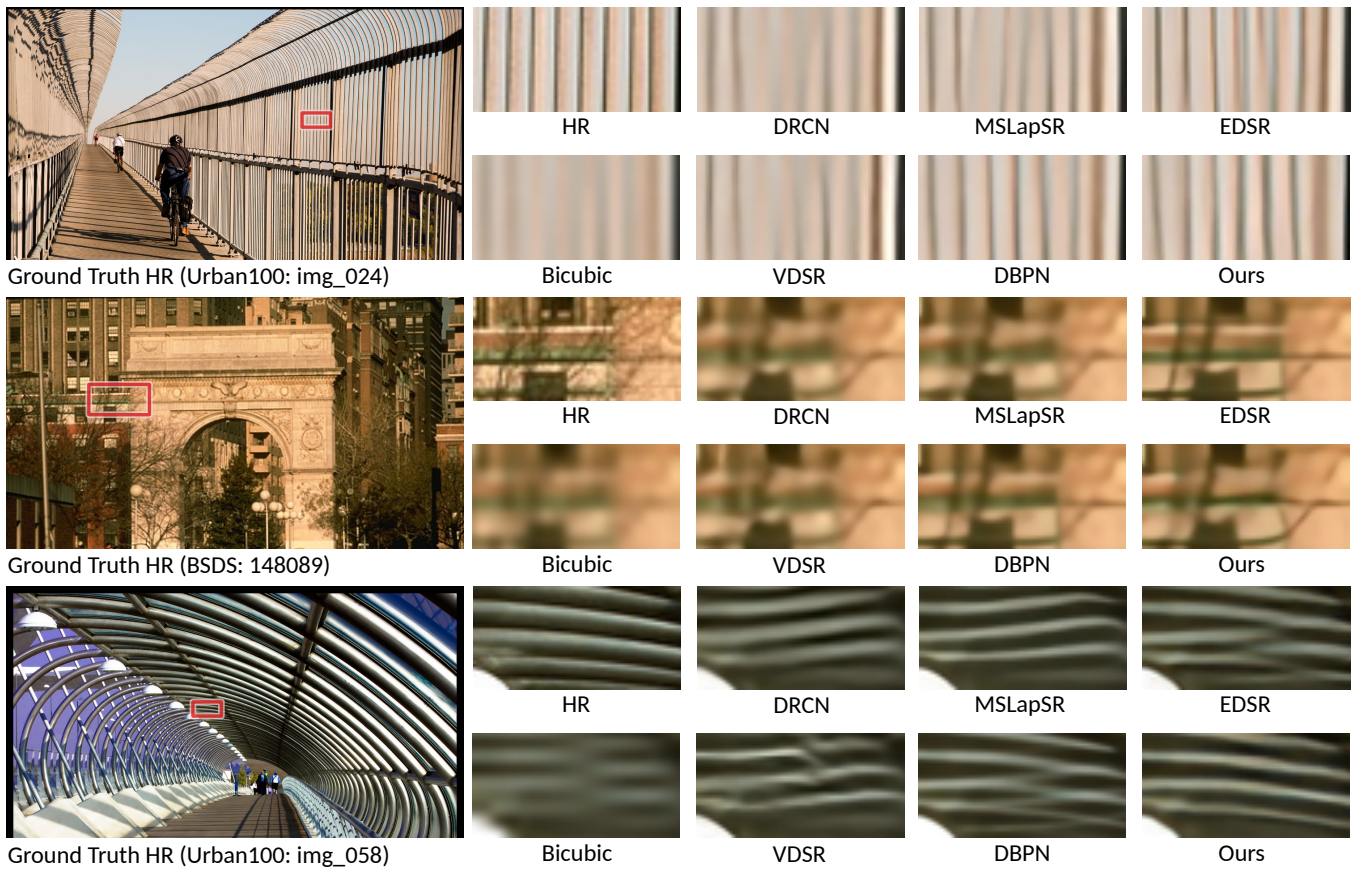


Figure 5: Perceptual evaluation of different SR methods for $4\times$ upscaling.

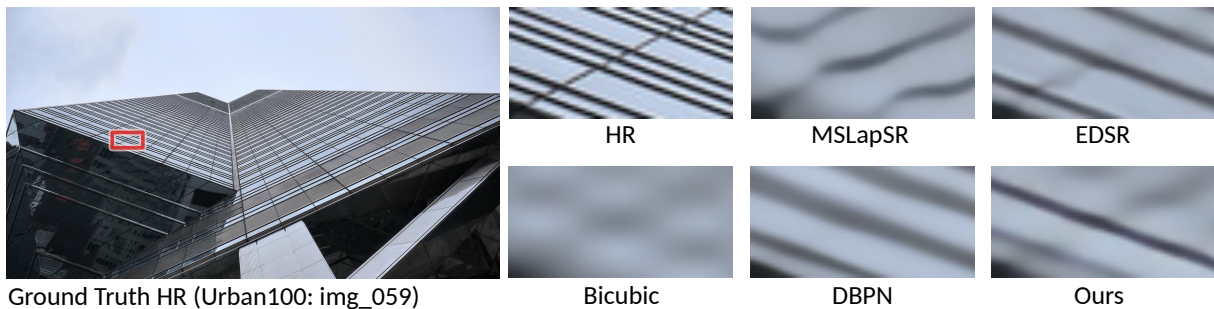


Figure 6: Perceptual evaluation of different SR methods for $8\times$ upscaling.

where $MSE = mean((X - Y)^2)$ is the mean square error of X and Y ; μ_X and μ_Y are the averages of X and Y , respectively; σ_X^2 and σ_Y^2 are the variances of X and Y , respectively; σ_{XY} is the covariance of X and Y ; $c_1 = 6.5025$ and $c_2 = 58.5225$.

Our method is outperformed only by EDSR and DBPN, that use 20 to more than 100 times the number of parameters of our system. Among systems with less than 2 million parameters, our system obtains better results, only matched by MSLapSR for $2\times$ upscaling.

In Figure 8 we better visualize the difference in complexity between different systems. It is apparent from these re-

sults that the size of the system matters to improve PSNR values, with EDSR and DBPN far from other methods, but the cost can be overwhelming in performance. Our system clearly improves the state of the art for systems with less than two million parameters, with better quality and significantly less parameters.

Figures 5 and 6 show the difference in perceptual quality for $4\times$ and $8\times$ upscaling. In general, EDSR and DBPN outputs look sharper and show consistent geometry. But often we find patches with aliasing, such as parallel lines, where the geometry becomes chaotic and perceptual quality behaves randomly. In such cases our system can take advan-



Figure 7: Output images for $8\times$ upscaling using different numbers of back-projections μ . The network was trained with $\mu = 2$ fixed. The effective filters are shown together with their frequency response (FFT).

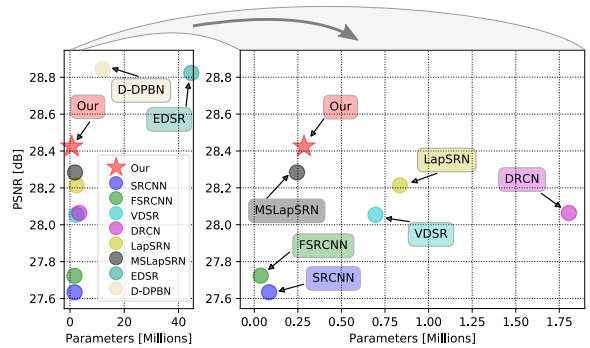
tage compared to most of other systems.

Analysis

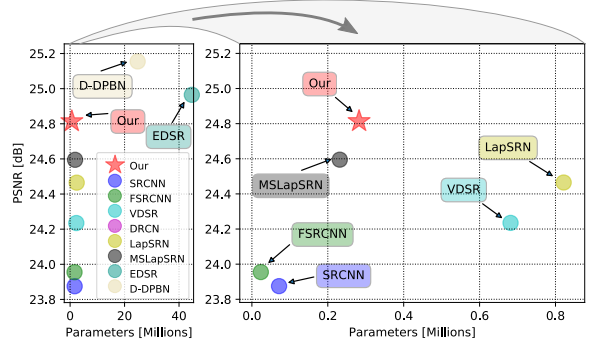
Figure 9 shows two examples of effective filters and residuals according to our novel visualization method. The shapes of filters following edges are reminiscent of the research on directional filter upscalers (Algazi, Ford, and Potharlanka 1991; Li and Orchard 2001). In general, we observe that **filters are highly adaptive**. They become symmetric and narrow in flat areas (like Linear or Bicubic upscalers), directional close to sharp lines, and increase in size close to complex features. For example, in the $4\times$ upscale of Figure 9 a red box shows an area with the eye of a baboon and the wrinkles around its eye. In the filter we can observe the shape of the eye and waves following the wrinkles. As expected, for $8\times$ upscaler the receptive field increases since the filters look bigger. We remind that we are using the same number of parameters, but passing through more layers. In the green box in 9 we highlight the face and hand of a woman. The filter at the tip of the nose follows the face features, and the filter at the hand captures features from all the edges around.

Figure 7 shows the effect of the depth on the effective filters and residuals. Here, training has tuned the network to work best with $\mu = 2$ back-projections. For $\mu > 2$ the output image looks oversharped, and too soft for $\mu < 2$. Accordingly, filters are larger and extract more high-frequencies for $\mu > 2$ and become smaller and low-pass for $\mu < 2$. The same effect is observed with classic IBP (Irani and Peleg 1991), showing the effective design of our network architecture.

We observe that residuals are in general small and help fixing textures and small details. This is an indication that



(a) $4\times$ results on Set14



(b) $8\times$ results on Set14

Figure 8: Quality vs Complexity of different SR methods.

filters are doing most of the work for upscaling. We remind that, after freezing most of the activations, the residual is a fixed component of the output that does not change with the input. Thus, residuals are very limited to estimate local details in the output as they depend only on activations for this purpose.

On the other hand, filters contain relations between neighboring pixels and we argue that because of this they are better to generalize. Effective filters also show all the details regarding the receptive field of the network. The receptive field is adaptive as the network does not use neighboring pixels if it does not need to (e.g. flat areas) and extend to large areas when the network is deep (e.g. $8\times$ upscaler) and local details are complex.

Finally, we remind that this analysis is precise. The network can be replaced by these adaptive filters plus residuals and it would give the exact same output. The strong dependency of filters and residuals on the input image shows all the non-linearities of the network. It is remarkable that convolutional networks can achieve the level of adaptivity revealed by these visualization experiments and it further justifies their success in super-resolution tasks.

Conclusions

We introduced a new architecture for single image super-resolution that reaches state of the art for methods with less than 2 million parameters and a new technique to analyze the network. The analysis shows how the network learns to upscale by capturing complex relationships between pixels.

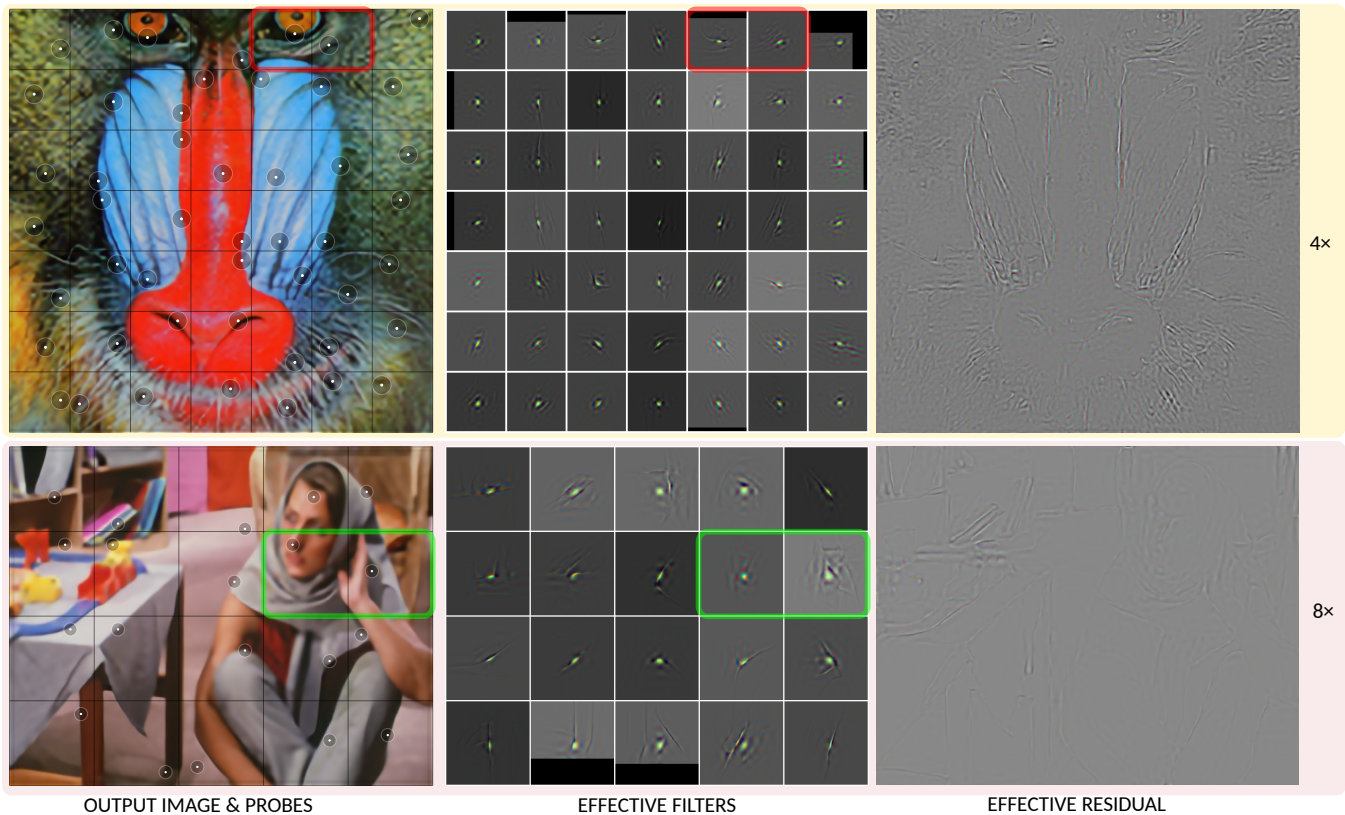


Figure 9: Effective filters and residual for 4 \times and 8 \times upscaler. . Filters are not only directional but follow several features around a pixels. Residuals are in general small and thus most of the work is done by the filters.

References

- Algazi, V. R.; Ford, G. E.; and Potharlanka, R. 1991. Directional interpolation of images based on visual properties and rank order filtering. In *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, volume 4, 30053008. Toronto, ON: IEEE Signal Processing Society.
- Dai, S.; Han, M.; Wu, Y.; and Gong, Y. 2007. Bilateral back-projection for single image super resolution. In *IEEE International Conference on Multimedia and Expo (ICME)*.
- Diamond, S.; Sitzmann, V.; Heide, F.; and Wetzstein, G. 2017. Unrolled optimization with deep priors. *CoRR* abs/1705.08041.
- Dong, C.; Loy, C.; He, K.; and Tang, X. 2014a. Learning a deep convolutional network for image super-resolution. In et al., D. F., ed., *Proceedings of European Conference on Computer Vision (ECCV)*, volume 8692 of *Lecture Notes in Computer Science*, 184–199. Zurich: Springer.
- Dong, C.; Loy, C. C.; He, K.; and Tang, X. 2014b. Learning a deep convolutional network for image super-resolution. In *in Proceedings of European Conference on Computer Vision (ECCV)*.
- Dong, C.; Loy, C.; He, K.; and Tang, X. 2015. Image super-resolution using deep convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 38(2):295–307.
- Dong, C.; Loy, C. C.; and Tang, X. 2016. Accelerating the super-resolution convolutional neural network. In *in Proceedings of European Conference on Computer Vision (ECCV)*.
- Gong, D.; Zhang, Z.; Shi, Q.; van den Hengel, A.; Shen, C.; and Zhang, Y. 2018. Learning an optimizer for image deconvolution. *CoRR* abs/1804.03368.
- Haris, M.; Shakhnarovich, G.; and Ukita, N. 2018. Deep back-projection networks for super-resolution. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Huang, G.; Chen, D.; Li, T.; Wu, F.; van der Maaten, L.; and Weinberger, K. Q. 2017a. Multi-scale dense convolutional networks for efficient prediction. *CoRR* abs/1703.09844.
- Huang, G.; Liu, Z.; van der Maaten, L.; and Weinberger, K. Q. 2017b. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Irani, M., and Peleg, S. 1991. Improving resolution by image registration. *CVGIP: Graph. Models Image Process.* 53(3):231–239.
- Kim, J.; Lee, J. K.; and Lee, K. M. 2016a. Accurate image super-resolution using very deep convolutional networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR Oral)*.

- Kim, J.; Lee, J. K.; and Lee, K. M. 2016b. Deeply-recursive convolutional network for image super-resolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR Oral)*.
- Kokkinos, F., and Lefkimmiatis, S. 2018. Deep image demosaicking using a cascade of convolutional residual denoising networks. *CoRR* abs/1803.05215.
- Kruse, J.; Rother, C.; and Schmidt, U. 2017. Learning to push the limits of efficient fft-based image deconvolution. In *The IEEE International Conference on Computer Vision (ICCV)*.
- Lai, W.-S.; Huang, J.-B.; Ahuja, N.; and Yang, M.-H. 2017a. Deep laplacian pyramid networks for fast and accurate super-resolution. In *IEEE Conference on Computer Vision and Pattern Recognition*.
- Lai, W.-S.; Huang, J.-B.; Ahuja, N.; and Yang, M.-H. 2017b. Fast and accurate image super-resolution with deep laplacian pyramid networks. *arXiv:1710.01992*.
- LeCun, Y.; Bengio, Y.; and Hinton, G. 2015. Deep learning. *Nature* 521(7553):436–444.
- Li, X., and Orchard, M. T. 2001. New edge-directed interpolation. *IEEE Transactions on Image Processing* 10(10):1521–1527.
- Lim, B.; Son, S.; Kim, H.; Nah, S.; and Lee, K. M. 2017. Enhanced deep residual networks for single image super-resolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*.
- Mallat, S., and Peyre, G. 2007. A review of bandlet methods for geometrical image representation. *Numerical Algorithms* 44(3):205–234.
- Mallat, S. 1998. *A Wavelet Tour of Signal Processing*. Academic Press.
- Olah, C.; Satyanarayan, A.; Johnson, I.; Carter, S.; Schubert, L.; Ye, K.; and Mordvintsev, A. 2018. The building blocks of interpretability. *Distill*. <https://distill.pub/2018/building-blocks>.
- Olah, C.; Mordvintsev, A.; and Schubert, L. 2017. Feature visualization. *Distill*. <https://distill.pub/2017/feature-visualization>.
- Park, S.; Park, M.; and Kang, M. 2003. Super-resolution image reconstruction: a technical overview. *Signal Processing Magazine, IEEE* 20(3):21–36.
- Proakis, J., and Manolakis, D. 2007. *Digital Signal Processing*. Prentice Hall international editions. Pearson Prentice Hall.
- Romano, Y.; Isidoro, J.; and Milanfar, P. 2016. RAISR: rapid and accurate image super resolution. *CoRR* abs/1606.01299.
- Timofte, R., and Smet, V. D. 2014. Gool, a+: Adjusted anchored neighborhood regression for fast super-resolution. In *in Proc. Asian Conf. Comput. Vis. (ACCV)*.
- Timofte, R.; Agustsson, E.; Van Gool, L.; Yang, M.-H.; Zhang, L.; et al. Ntire 2017 challenge on single image super-resolution: Methods and results. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*.
- Trottenberg, U., and Schuller, A. 2001. *Multigrid*. Orlando, FL, USA: Academic Press, Inc.
- Widlund, O., and Toselli, A. 2004. *Domain decomposition methods - algorithms and theory*, volume 34. Springer.
- Yang, J.; Wright, J.; Huang, T.; and Ma, Y. 2008. Image super-resolution as sparse representation of raw image patches. In *IEEE Conference on Computer Vision and Pattern Recognition*.
- Yang, J.; Wright, J.; Huang, T.; and Ma, Y. 2010. Image super-resolution via sparse representation. *Image Processing, IEEE Transactions on* 19(11):2861–2873.
- Zhang, Q., and Zhu, S. 2018. Visual interpretability for deep learning: a survey. *Frontiers of Information Technology and Electronic Engineering* 19(1):27–39.
- Zhang, K.; Zuo, W.; Gu, S.; and Zhang, L. 2017. Learning deep cnn denoiser prior for image restoration. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Supplementary Material

Back–Projection Convergence

In the following analysis we consider vectors instead of images. This simplifies the analysis since we can use (rectangular) matrices U and D to represent **upsampling** and **downsampling** operators, respectively. Thus:

- Dy represents **downsampling** operation $(Y * g) \downarrow s$. Where $y = \text{vec}(Y)$ for a high resolution image Y .
- Ux represents **upsampling** operation $(X \uparrow s) * p$. Where $x = \text{vec}(X)$ for a low resolution image X .

Our L -level downscaling model is now:

$$x = D^L y. \quad (10)$$

The target of back-projections is to recover this property in upscale images. Convergence is characterized by the mismatch error defined as follows:

Definition 2 (Downscale mismatch error) *The downscale mismatch error of a high-resolution image y_{L+1} at level $L + 1$ is given by:*

$$e(y_{L+1}) = x - D^L y_{L+1}. \quad (11)$$

And the classic back-projection algorithm uses the following iteration:

Definition 3 (Classic Back-projection) *The one-level ($L = 1$) back-projection iteration is given by:*

$$y_2^{t+1} = y_2^t + Ue(y_2^t). \quad (12)$$

First, we recall the classic convergence theorem (Irani and Peleg 1991), in a simplified version adapted from (Dai et al. 2007).

Theorem 2 (Classic Back-Projection Convergence) *The one-level ($L = 1$) back-projection iteration converges, $Dy \rightarrow x$, at exponential rate if the following condition holds:*

$$\|I - DU\|_1 < 1, \quad (13)$$

where $\|A\|_1 = \max_j \sum_i |a_{ij}|$.

Proof:

$$e(y_2^{t+1}) = x - Dy_2^{t+1} \quad (14)$$

$$= x - D(y_2^t + Ue(y_2^t)) \quad (15)$$

$$= \underbrace{x - Dy_2^t}_{e(y_2^t)} - DUe(y_2^t) = (I - DU)e(y_2^t), \quad (16)$$

By Hölder's inequality we get $\|e(y_2^{t+1})\|_1 \leq \|I - DU\|_1 \|e(y_2^t)\|_1$ which proves the theorem. \square

Now, we rewrite the MGBP algorithm using matrix/vector notation:

Algorithm 2 Multi-Grid Back-Projection (MGBP) algorithm in matrix/vector notation.

$MGBP(x)$:

Input: Input image x .

Input: Integers $\mu \geq 0$ and $L \geq 1$.

Output: Images y_k , $k = 2, \dots, L$.

1: $y_1 = x$

2: **for** $k = 2, \dots, L$ **do**

3: $u = Uy_{k-1}$

4: $y_k = BP_k^\mu(u, y_1, \dots, y_{k-1})$

5: **end for**

$BP_k^\mu(u, y_1, \dots, y_{k-1})$:

Input: Image u , level index k , number of steps μ .

Input: Images y_1, \dots, y_{k-1} (only for $k > 1$).

Output: Updated image u

1: **if** $k > 1$ **then**

2: **for** $step = 1, \dots, \mu$ **do**

3: $d = BP_{k-1}^\mu(Du, y_1, \dots, y_{k-2})$

4: $u = u + U(y_{k-1} - d)$

5: **end for**

6: **end if**

Definition 4 (single-level mismatch error) *Assume that the MGBP algorithm runs μ back-projection steps from level $k = 2, \dots, L$. Then, the single-level mismatch error for a $L + 1$ -resolution image y_{L+1} is given by:*

$$q(y_{L+1}) = y_L^\mu - Dy_{L+1} \quad (17)$$

This is useful since the MGBP iteration can be written as $y^{t+1} = y^t + Uq(y^t)$. It easily follows that:

$$q(y_{L+1}^{t+1}) = y_L^\mu - Dy_{L+1}^{t+1} \quad (18)$$

$$= y_L^\mu - D(y_{L+1}^t + Uq(y_{L+1}^t)) \quad (19)$$

$$= (I - DU)q(y_{L+1}^t) \quad (20)$$

Theorem 3 (MGBP Convergence) *The MGBP iteration converges, $D^L y_{L+1} \rightarrow x$, at exponential rate if the following condition holds:*

$$\|I - DU\|_1 < 1. \quad (21)$$

Proof: We follow by induction. From the classic convergence theorem we know that $\|e(y_2^{t+1})\|_1 \leq \|I - DU\|_1 \|e(y_2^t)\|_1$. For $L = 1$ the MGBP iteration is equivalent to classic back-projection, and so the theorem is verified for $L = 2$.

Next, we assume that the theorem holds for L and we study the error for $L + 1$:

$$e(y_{L+1}^\mu) = x - D^L y_{L+1}^\mu \tag{22}$$

$$= x - D^{L-1}(y_L^\mu - q(y_{L+1}^\mu)) \tag{23}$$

$$= x - D^{L-1}(y_L^\mu - (I - DU)^\mu q(y_{L+1}^0)) \tag{24}$$

$$= x - D^{L-1}(y_L^\mu - (I - DU)^\mu (y_L^\mu - DU y_L^\mu)) \tag{25}$$

$$= \underbrace{x - D^{L-1} y_L^\mu}_{e(y_L^\mu)} + D^{L-1} (I - DU)^{\mu+1} y_L^\mu \tag{26}$$

By triangle inequality and Hölder's inequality we have:

$$\|e(y_{L+1}^\mu)\|_1 \leq \|e(y_L^\mu)\|_1 + \|D^{L-1}\|_1 \|(I - DU)^{\mu+1}\|_1 \|y_L^\mu\|_1 . \tag{27}$$

The first term decreases exponentially by inductive assumption. The second term also decreases exponentially with μ if $\|I - DU\|_1 < 1$. The term $\|y_L^\mu\|_1$ is bounded since by inductive assumption we know that y_L^μ converges to a fixed L -level image. \square