

edge-SR: Super-Resolution For The Masses

Pablo Navarrete Michelini, Yunhua Lu, Xingqun Jiang
BOE Technology Group Co., Ltd.

Abstract

Classic image scaling (e.g. bicubic) can be seen as one convolutional layer and a single upscaling filter. Its implementation is ubiquitous in all display devices and image processing software. In the last decade deep learning systems have been introduced for the task of image super-resolution (SR), using several convolutional layers and numerous filters. These methods have taken over the benchmarks of image quality for upscaling tasks. Would it be possible to replace classic upscalers with deep learning architectures on edge devices such as display panels, tablets, laptop computers, etc.? On one hand, the current trend in Edge-AI chips shows a promising future in this direction, with rapid development of hardware that can run deep-learning tasks efficiently. On the other hand, in image SR only few architectures have pushed the limit to extreme small sizes that can actually run on edge devices at real-time. We explore possible solutions to this problem with the aim to fill the gap between classic upscalers and small deep learning configurations. As a transition from classic to deep-learning upscaling we propose edge-SR (eSR), a set of one-layer architectures that use interpretable mechanisms to upscale images. Certainly, a one-layer architecture cannot reach the quality of deep learning systems. Nevertheless, we find that for high speed requirements, eSR becomes better at trading-off image quality and runtime performance. Filling the gap between classic and deep-learning architectures for image upscaling is critical for massive adoption of this technology. It is equally important to have an interpretable system that can reveal the inner strategies to solve this problem and guide us to future improvements and better understanding of larger networks.

1. Introduction

A market is growing rapidly and steadily to provide so-called Edge-AI chips that will be able to spread the success of deep-learning systems to edge devices [13, 17, 4]. This is a massive market that includes phones, tablets and high resolution TV displays, among others. For some applica-

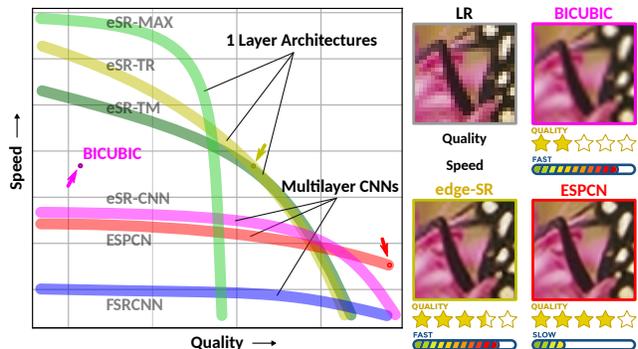


Figure 1. Trade-off pattern observed in our experiments. A standard upscaler (e.g. bicubic) is fast and can be deployed in any display hardware. Multilayer CNNs reach better quality but show sharp loss in quality when their size is reduce to make them faster. We propose edge-SR (eSR) one layer architectures that can adjust the size to reach the best possible quality for a display hardware.

tions the success is guaranteed, such as image classification or object detection, where input images are relatively small (e.g. 256×256) and the output data is low dimensional (e.g. labels or bounding boxes). For other applications such as recovering a high-resolution image from a small-resolution image, also known as image super-resolution (SR), the future is less certain since both input and output images can contain a large amount of data. Consider upscaling images from Full-HD to 4K resolution in TV displays for example. The input layer needs to handle 2 megapixels and the output layer needs to deliver 8 megapixels at a rate of at least 24 frames per second. Interestingly, upscaling with small factors (e.g. $2\times$) is both the easiest problem for networks to fix, typically requiring less number of parameters to learn, and at the same time the most difficult solution to deploy. The latter is due to the fact that display devices have a fixed output resolution. For small upscaling factors the input images are still large and demand higher input throughput compared to higher upscaling factors, where input images get smaller and smaller. Small upscaling factors are also of primary concern in applications since they are the most critical technology for transitions between current and new standards (e.g. FHD to 4K, 4K to 8K, etc). Thus, the

problem of image SR becomes both more interesting and more challenging given extreme performance constraints.

History of SR. Standard upscaler algorithms, such as linear or bicubic upscalers, apply a low-pass filter on a high resolution image created by inserting zeros between adjacent pixels in the low resolution [32, 23]. Modern tensor processing frameworks (e.g. Pytorch, Tensorflow, etc.) implement this process using a so-called *strided transposed convolutional layer* with a single filter per input channel. More advanced upscalers have followed geometric principles to improve image quality. For example, *edge-directed interpolation* uses adaptive filters to improve edge smoothness [2, 18], or *bandlet* methods use both adaptive upsampling and filtering [24]. Later on, machine learning has been able to use examples of pristine high-resolution images to learn a mapping from low-resolution [30]. The rise of deep-learning and convolutional networks in image classification tasks [15] quickly saw a series of important improvements. Many of these improvements followed the progress in network architectures for image classification, as seen for example with CNNs applied in SRCNN [5], ResNets [8] applied in EDSR [20], DenseNets [10] applied in RDN [43], attention [9] applied in RCAN [41], non-local attention [36] applied in RNAN [42], and swin transformers [22] applied in SwinIR [19].

Real-time SR. The first deep learning system proposed for image SR, namely SRCNN [5], used a relatively small number of parameters ($60k$) and became a suitable candidate for edge devices. Soon after, FSRCNN [6] realized that significant improvements in quality and performance can be achieved by performing computations at low resolution. They proposed a *short* configuration using $4k$ parameters in a sequence of 4 convolutional layers, plus a final strided transposed convolution to perform upscaling, reaching real-time performance for small resolutions. The next major progress towards real-time applications was made by ESPCN [33] that made popular the application of pixel-shuffle layers, multiplexing several network channels to form higher resolution outputs [29, 27]. They proposed a configuration using $20k$ parameters and 3 convolutional layers with all computations performed at low resolution. Both FSRCNN and ESPCN left a strong mark on future image SR research that very often performs computations at low resolution and use pixel-shuffle layers. Nevertheless, the research clearly shifted to networks of larger sizes that can achieve much better quality. But large networks that contain several million parameters, for example EDSR [20] (combining ResNets and pixel-shuffle), are currently unable to reach the throughput needed for real-time applications on edge devices. Several so-called lightweight networks have been proposed for middle ground applications [40, 16, 21, 37, 3, 12]. Typical lightweight networks use hundred of thousands parameters and are still beyond

the capabilities of real-time applications on edge devices.

The Problem. Despite the promising advances in technology, the challenge of image SR for edge devices remains largely unresolved. One might expect Edge-AI chips to get faster and cheaper but standards also evolve to make problems more difficult (e.g. BT.2020 [34]) with more pixels, higher bit depths, higher framerates, etc.. Thus, the success of AI chips to deploy image SR technologies and reach massive markets strongly depends on better algorithm solutions. The major challenge is how to simplify network structures all the way down to reach performance levels comparable to those of classic non-adaptive upscalers. A classic $2\times$ bicubic, doubling the horizontal and vertical resolution, can be implemented using a transposed convolutional layer with a single filter using 121 parameters. We can think of this as the simplest possible network configuration for image SR. A configuration that is interpretable in the sense that we understand what the interpolation filter values represent. Our main task here is to explore the landscape between classic upscaling on one hand, and small deep-learning systems on the other hand, in order to provide practical solutions for the current state of applications in edge devices.

Towards a solution. Exploring different configurations for existing networks, such as FSRCNN and ESPCN, is a straightforward and necessary task to undertake. But we propose to move a step further, introducing a minimal set of architectures, **edge-SR (eSR)**, that can perform image SR even with a single convolutional layer. We explore both a straightforward 1-layer Maxout network (eSR-MAX) as well as self-attention strategies (eSR-TM and eSR-TR) that provide a semi-classical interpretation. The latter approaches use a single layer both to detect local patterns (e.g. edges or textures) as well as to generate candidate upscale solutions. Generally speaking, the detection mechanism estimates the probability of the best upscale solution and it is used to compute a weighted average of the candidate output images that gives the final output. We will show how to implement this solution efficiently using standard deep learning modules that can run on AI chips.

Contributions. Our major contributions include:

- The **proposal** of several one-layer architectures that strive for simplicity to fill the gap between classic and deep learning upscalers.
- An **exhaustive search** among 1,185 network models, including different configurations of eSR, FSRCNN, and ESPCN. Each architecture was trained under identical conditions and tested for speed, power consumption and image quality. The results allows us to visualize the trade-off between image quality and runtime performance that is critical for our purpose. Figure 1 shows the general pattern observed in our results. We found that different architectures show very different balance in the

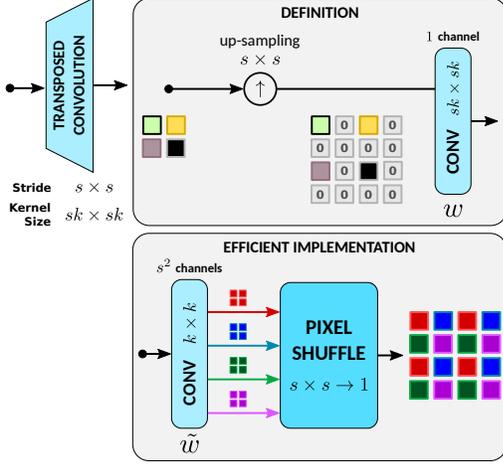


Figure 2. Classic $s \times s$ image upscaling is performed by a transposed convolutional layer. An efficient implementation splits the filter into s^2 smaller filters that work at LR. The final output is obtained by multiplexing the s^2 channels using a pixel-shuffle layer.

trade-off between speed and image quality. Multi-layer networks (deep learning) show a strong advantage at low speed and high quality, and our proposed one-layer solutions show a clear advantage at high speed requirements.

- The **interpretation and analysis** of strategies learned by self-attention in one-layer architectures. We provide a novel interpretation of the self-attention mechanism based on the simple principles of template matching and classic upscaling. Here, training results indicate that one-layer networks do not use smooth upscaling kernels and rely mostly on independent sub-pixel solutions.

These results may bring about the following **future impact**: 1) the possibility of image SR systems that can be massively deployed on edge devices, 2) a better understanding of the internal learning mechanisms of small network architectures, and 3) a better appreciation of the trade-off between image quality and runtime performance for future applications and research.

2. Super-Resolution for Edge Devices

Classical. Image *upscaling* and *downscaling* refer to the conversion of low resolution (LR) images to high resolution (HR) and vice versa. These two processes are closely related. The simplest way to downscale an image from HR to LR is known as *pooling* or *downsample*. The process of downsample uniformly drops pixels in both horizontal and vertical directions. The problem with such downscalers is that groups of high and low frequency components of the HR image can end up in the same low frequency component at LR, leading to well known *aliasing* artifacts [32, 23]. To avoid this problem a classic *linear downscaler* first removes high frequencies using an *anti-aliasing* low-pass filter and

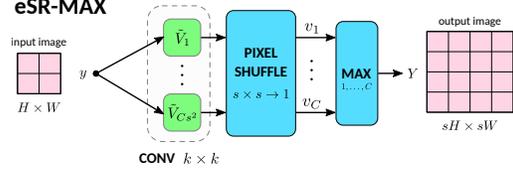


Figure 3. Diagram of edge-SR Maximum. Uses one convolutional layer followed by a pixel-shuffle multiplexer and a non-linear module that chooses the maximum pixel value among all filters.

then downsamples the image. This process is implemented in tensor processing frameworks with *strided convolutional layers* where the kernel or weight parameters correspond to the low-pass filter coefficients. The process of classic *linear upscaling* corresponds to the transposed of the downscaling linear transformation and it is illustrated in Figure 2. The transposition reverts the ‘filter-then-downsampling’ operation into an ‘upsampling-then-filter’ operation where the *upsampling* increases the resolution of an image by inserting zeros between LR pixels. The upsampling introduces high frequencies that are removed by a so-called *interpolation* filter with coefficients w . The interpolation filter is the transposed of the anti-aliasing filter, typically identical because most upscalers are symmetric. Tensor processing frameworks implement this process using *strided transposed convolutional layers*.

The upscaling definition in Figure 2 is clearly inefficient as the upsampling introduces many zeros that will waste resources when multiplied by filter coefficients. A very well know optimization, widely used in practical implementations of classic upscalers is to split or demultiplex the interpolation filter from size $sk \times sk$ in Figure 2 to s^2 so-called *efficient* filters of size $k \times k$ working at LR [32, 23]. The outputs of the s^2 filters are then multiplexed by a pixel-shuffle operation to obtain the upscaled image, as illustrated in Figure 2. Let $\tilde{w}_i \in \mathbb{R}^{k \times k}$, with $i = 1, \dots, s^2$, be the coefficients of the efficient filters. The interpolation filter can then be recovered by multiplexing the efficient coefficients back to their original place. This is,

$$w = \text{Pixel-Shift}_{s \times s}(\tilde{w}_i, i = 1 \dots, s^2). \quad (1)$$

In our experiments we will compare different architectures including a bicubic upscaler. In order to remove implementation advantages we implemented the upscaler using the efficient implementation in Figure 2. We used standard bicubic interpolation filter coefficients and verified that we obtain the same outputs as other software implementations up to floating point precision.

Maxout. Our first proposal is edge-SR Maximum (eSR-MAX). This is an attempt to obtain the fastest solution from a single convolutional layer that outputs several upscaled candidates. A quick decision is made by choosing the maximum value across all channels as shown in Figure 3. This

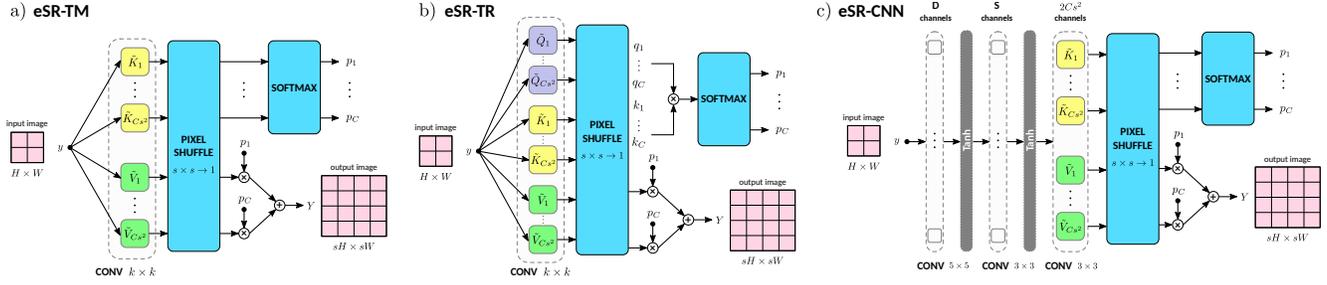


Figure 4. Diagrams of edge-SR architectures using self-attention: a) edge-SR Template Matching (eSR-TM) runs both template matching (K filters) and upscaling modules (V filters) using a single convolutional layer (see Figures 5 and 10 for details on this interpretation), b) edge-SR TRansformer (eSR-TR) uses two sets of *query* (Q) and *key* (K) filters to estimate the best upscaler model, and c) edge-SR CNN (eSR-CNN) starts with a multilayer network following ESPCN and ends with a eSR-TM module.

Algorithm 1 edge Super-Resolution (eSR). **Input** : y (1-channel). **Output** : Y (1-channel).

eSR-MAX(y, C, k, s):

Parameters: Integer $C > 1, k > 1, s > 1$.

1: $Y = \text{Max}_{1 \rightarrow C} \text{Pixel-Shuffle}_{s \times s}(\text{Conv}^{k \times k}(y))$

eSR-TR(y, C, k, s):

Parameters: Integer $C > 1, k > 1, s > 1$.

1: $f = \text{Pixel-Shuffle}_{s \times s}(\text{Conv}^{k \times k}(y))$

2: $p = \text{SoftMax}(f_{1 \rightarrow C} \otimes f_{C+1 \rightarrow 2 \cdot C})$

3: $Y = \sum_{1 \rightarrow C} (f_{2 \cdot C+1 \rightarrow 3 \cdot C} \otimes p)$

eSR-TM(y, C, k, s):

Parameters: Integer $C > 1, k > 1, s > 1$.

1: $f = \text{Pixel-Shuffle}_{s \times s}(\text{Conv}^{k \times k}(y))$

2: $Y = \sum_{1 \rightarrow C} (f_{C+1 \rightarrow 2 \cdot C} \otimes \text{SoftMax}(f_{1 \rightarrow C}))$

eSR-CNN(y, C, D, S, s):

Parameters: Integer $C > 1, D > 1, S > 1, s > 1$.

1: $f = \text{Pixel-Shuffle}_{s \times s} \circ \text{Conv}^{5 \times 5} \circ \text{Tanh} \circ \text{Conv}^{3 \times 3} \circ \text{Tanh} \circ \text{Conv}^{3 \times 3}(y)$

2: $Y = \sum_{1 \rightarrow C} (f_{C+1 \rightarrow 2 \cdot C} \otimes \text{SoftMax}(f_{1 \rightarrow C}))$

corresponds to a particular case of a Maxout network [7].

Self-Attention. Our second proposal is edge-SR Template Matching (eSR-TM) that follows a semi-classical strategy. The basic idea is explained in Figure 5. First, a template matching module detects patterns (e.g. edge directions) and gives us the probability for each pattern. This is achieved by: first, use matching filter coefficients that resemble the pattern, and second, normalize pixel values across channels to represent the probability of each template. A set of upscale images are computed at the same time for each one of the patterns. Since both the matching and the upscaling filters follow the same patterns, *we expect the filter coefficients to look similar* as displayed in Figure 5 for the case of edge patterns. Thus, we can verify if an eSR-TM configuration learned to perform template matching by checking the correlations between filter coefficients. The optimal prediction for the output image is the expected value over all templates. Thus, the probabilities are used to compute the expected value by weighing the solution of different upscales that when combined give the final output.

Figure 4.a shows the diagram of the efficient implementation of this idea using $C \in \mathbb{N}^+$ templates. In this efficient implementation of a transposed convolution the C matching filters K split into Cs^2 efficient filters \tilde{K} , before multiplexing with pixel-shuffle. We can always get the interpolation filters K from \tilde{K} using equation (1). The outputs of the

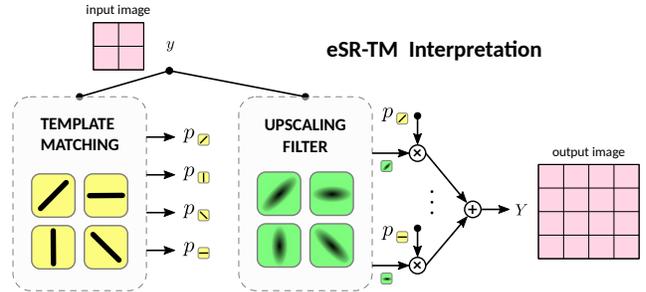


Figure 5. Interpretation of edge-SR Template Matching. A template matching module computes the probability of finding one of the template image features learned from examples. The input image is upscaled using a set of different upscales also learned by examples. The output is the expected value computed by the weighted average of upscale images and template probabilities.

filters are then normalized among all the channels using a softmax module. This gives us the pixel-wise probabilities:

$$p_i = e^{K_i \otimes (y \uparrow s)} / \sum_{j=1}^C e^{K_j \otimes (y \uparrow s)}, \quad (2)$$

where $i = 1, \dots, C$, \otimes is the convolution operator, \uparrow refers to the upsampling operation defined in Figure 2. The same convolutional layer in Figure 4.a runs Cs^2 efficient filters \tilde{V}

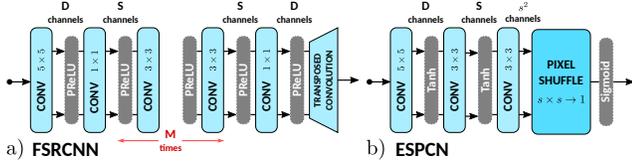


Figure 6. Deep-learning architectures selected for experiments: a) FSRCNN from [6], and b) ESPCN from [33].

to get C high resolution candidates after pixel-shuffle. The final luminance HR output image Y is given by:

$$Y = \mathbb{E}[V_i \otimes (y \uparrow s)] = \sum_{i=1}^C p_i \otimes (V_i \otimes (y \uparrow s)), \quad (3)$$

where \otimes represent a Hadamard (or pixel-wise) product.

The eSR-TM system is essentially a self-attention module, except for the pixel-shuffle layer and the sum over all channels in the last stage. These two differences are significant since: first, they embed the upscaling process within the attention module, and second, they make explicit use of probabilities to compute an expected value thus providing a clear interpretation of this module.

Our third proposal is edge-SR TRansformer (eSR-TR) that uses the popular *transformer* self-attention module from [35]. Figure 4.b shows the efficient implementation of this system. Here, the matching filters from eSR-TM are replaced by two sets of *query* (Q) and *key* (K) filters to estimate the probabilities. This changes the template matching interpretation of eSR-TM, using a rank-1 quadratic form with Q and K filters instead of a single template matching filter. The purpose of this architecture is to test any advantage that this change could bring given the increasing popularity and success of this module in recent research.

The code for all eSR systems is given in Algorithm 1.

Deep-Learning. We consider FSRCNN [6] and ESPCN [33] as candidate deep learning architectures for image SR on edge devices. Figure 6 shows the detail structure of FSRCNN and ESPCN network architectures. In comparison, FSRCNN uses more layers (at least 5) and smaller number of channels per layer than ESPCN. Another difference is the upscaling strategy, with FSRCNN using a strided transposed convolution and ESPCN using pixel-shuffle. According to classic interpolation theory these two approaches are equivalent as shown in Figure 2 (see also [32, 23]), but implementations can be different. Tensor processing frameworks typically implement transposed convolution using the gradient of a convolutional layer[31], based on the vector calculus property for gradients of linear transformations: $\nabla_x(Ax + b)y = A^T y$. This very different approach might lead to differences in performance.

Finally, we also propose the edge-SR CNN (eSR-CNN) architecture in Figure 4.c and Algorithm 1. This is sim-

Table 1. Set of hyper-parameters used to create a pool of 1, 185 models that were trained and tested in our experiments.

Bicubic	Total : 1 model per scale factor.
	C : 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16.
	k : 3, 5, 7.
eSR	Type : Maximum (MAX), Template Matching (TM), Transformer (TR).
	Total : 144 models per scale factor.
eSR-CNN	C : 2, 4, 6, 8.
	D : 1, 3, 5, 7, 9.
	S : 3, 6, 9, 12, 15.
	Total : 100 models per scale factor.
	D : 6, 19, 32, 44, 56.
	S : 1, 3, 6, 9, 12.
FSRCNN	M : 1, 4
	Total : 50 models per scale factor.
	D : 0, 4, 6, 10, 12, 16, 28, 40, 52, 64.
	S : 3, 6, 9, 12, 15, 18, 21, 24, 27, 32.
ESPCN	Total : 100 models per scale factor.
	Factors : 2 \times , 3 \times , 4 \times .
	Total : 1,185 models.

ply an extension of the single convolutional layer in eSR-TM into a multi-layer structure identical to ESPCN. Here, the purpose is to test if ESPCN, that achieves better results compared to FSRCNN in our tests, can be improved by using a self-attention module to upscale.

3. Experiments

Models. Candidate models for test evaluations include: bicubic, FSRCNN, ESPCN and eSR. From these, the bicubic classic upscaler is the only one without hyper-parameters and fixed configuration that do not require training. For other architectures we need to train a model for each set of hyper-parameters. Table 1 shows the list of hyper-parameters chosen for our experiments. These include default settings of FSRCNN and ESPCN as well as configurations with very small number of parameters. Our model pool includes a total of 1, 185 models to evaluate.

Training. We need to train a total of 1, 185 models that include different scaling factors, network architectures and model hyper-parameters. We trained all these models independently using an identical procedure. We used the General-100 dataset [6] combined with 91-image dataset [38] to extract training patches. For each image in the dataset we randomly cut a HR patch of size 78×78 for $2\times$ and $3\times$ upscaling factors, and 76×76 for $4\times$ factor. The images were converted to grayscale using BT.609 color matrix and downsampled using a standard Bicubic algorithm. We used minibatch size 16 and trained each model for 25, 000 epochs using a standard mean-square-error (MSE) loss. We started with a learning rate of 10^{-3} and reduce it to half once every 3, 000 epochs. We used Adam optimizer [14] with $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. We used seven Tesla M40 GPUs for training with the whole process completed in about two months.

Measurements. To test our final models we considered two inference devices: 1) Nvidia Jetson AGX Xavier, an embedded system-on-module (SoM) from the Nvidia

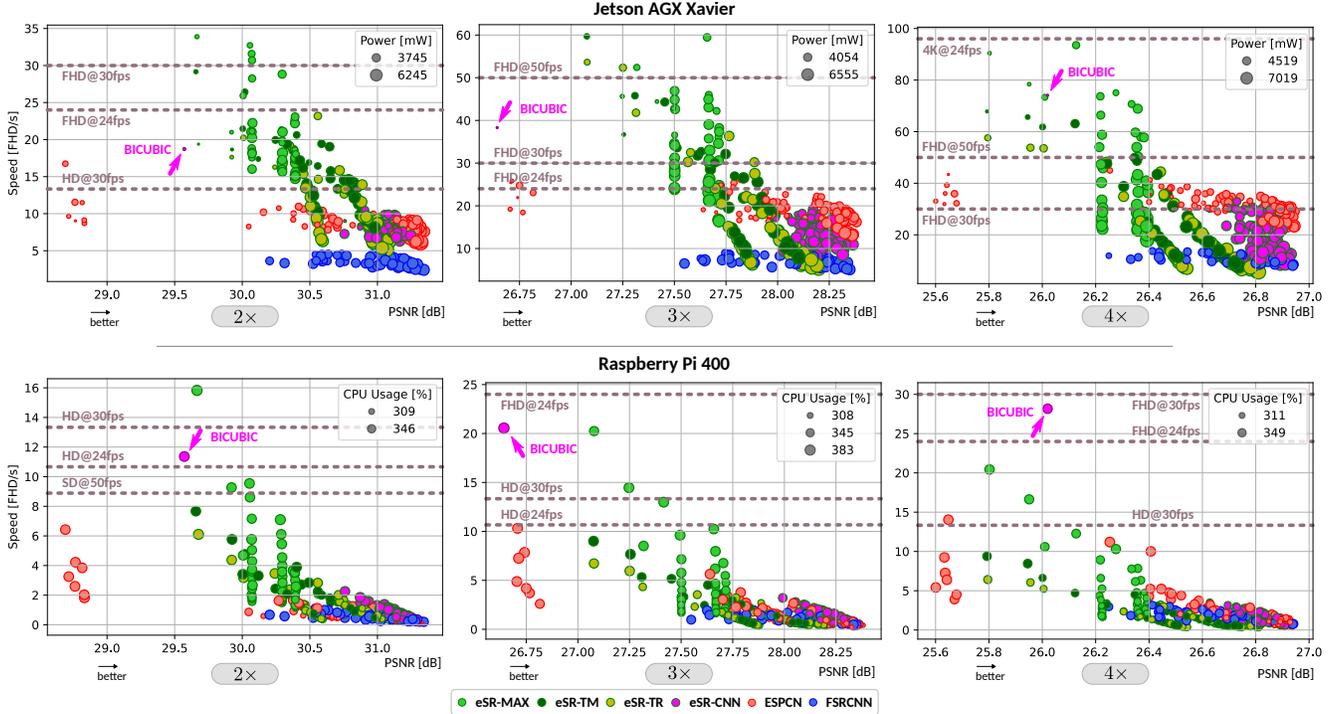


Figure 7. Scatter plot to compare speed, in number of Full-HD pixels per second, with respect to quality, measured as PSNR for the BSDS-100 dataset. A total of 1,185 models were identically trained considering different upscaling factors (2 \times , 3 \times and 4 \times) and architectures (eSR, ESPCN and FSRCNN). We run all models on edge devices: Jetson AGX Xavier (GPU with 16-bit floating point precision) and Raspberry Pi 400 (CPU with 32-bit floating point precision). Magnified plots with model annotations are provided in the Appendix.

AGX Systems family, including an integrated Volta GPU with tensor cores, and 2) a Raspberry Pi 400, an embedded device featuring a quad-core 1.8GHz, 64-bit ARM Cortex CPU processor. The power consumption of the Jetson AGX is set to a 30 Watt profile, while the Raspberry Pi 400 nominal consumption is 15 Watt.

We run each model to output a set of 14 Full-HD images, downscaling appropriately from randomly selected images of the DIV2K dataset [1]. We use 16-bit floating point precision during inference. For each image we run the model 10 times to avoid warm-up effects, measuring the minimum CPU and GPU processing time from profiler’s data. We computed the speed of a model using the total number of pixels processed (considering only one run per image) divided by the processing time (using the minimum time over each one of the 10 runs). To make the measurement of speed easier to read we use units of [FHD/s], this is, number of Full-HD pixels (1920 \times 1080) per second.

Image quality was measured separately using the standard datasets: Set-5, Set-14[39], BSDS-100[25], Urban-100[11] and Manga-109[26]. We also measured maximum power consumption for the Jetson AGX and CPU usage for the Raspberry Pi that does not include power sensors.

Results. Figure 7 shows scatter plots to compare speed with respect to image quality, measured as PSNR for the

BSDS-100 dataset. Results for other datasets, metrics (SSIM) and devices (GTX 1080 Max-Q) are shown in the Appendix with similar conclusions. The size of the circles are proportional to the power consumption and CPU usage for the AGX and Raspberry Pi devices, respectively. Finally, Table 2 shows detailed results per dataset for a subset of the models selected according to different criteria.

4. Analysis

Trade-off. The results displayed in Figure 7 allow us to fully appreciate the trade-off between image quality and runtime performance. The bicubic upscaler sets the target as we know that it can be massively deployed in display devices at large scale. Between the bicubic upscaler and deep-learning configurations using FSRCNN, ESPCN or eSR-CNN we observe a large empty region. Our proposed edge-SR (eSR) architectures succeeds to fill this gap in edge GPU devices (AGX and also GTX 1080 MaxQ available in the Appendix) and improve bicubic upscaler both in speed and image quality. In the Raspberry Pi CPU device edge-SR partially succeeds to fill this gap for 2 \times and 3 \times upscaling factor and fails at 4 \times factor where bicubic reaches a better performance. The best results of edge-SR is observed for 2 \times upscaling factor. The distribution of scatter points in Figure 7 for 2 \times upscaling shows that deep-learning meth-

Table 2. Image quality and performance metrics for selected methods among all 1, 185 models trained in our experiments. Values of speed, measured in number of Full-HD pixels per second, and power, in units of Milliwatts, are specific of a Jetson AGX Xavier GPU. Methods are selected based on best speed, PSNR in BSDS-100 dataset, and default configurations. Best results are shown in bold (ignoring bicubic).

Algorithm	s	Selection	Configuration	Speed [FHD/s]	Power [mWatts]	Set5		Set14		BSDS100		Urban100		Manga109	
						PSNR	SSIM								
Bicubic	2	–	–	19	1550	33.73	0.928	30.29	0.869	29.57	0.842	26.89	0.841	30.85	0.934
eSR	2	PSNR	CNN: $C = 6, D = 3, S = 15$	8	3868	36.58	0.953	32.38	0.905	31.25	0.885	29.26	0.891	35.33	0.965
eSR	2	speed	MAX: $k = 3, C = 1$	34	1859	33.15	0.928	30.16	0.882	29.66	0.862	26.94	0.857	30.46	0.937
ESPCN	2	default [33]	$D = 64, S = 32$	6	6800	36.64	0.953	32.46	0.907	31.32	0.887	29.37	0.893	35.76	0.967
ESPCN	2	PSNR	$D = 22, S = 32$	8	4945	36.70	0.953	32.47	0.907	31.35	0.887	29.44	0.894	35.79	0.967
ESPCN	2	speed	$D = 0, S = 3$	17	2324	29.76	0.919	28.96	0.881	28.69	0.865	26.38	0.853	27.66	0.938
FSRCNN	2	default [6]	$D = 32, S = 6, M = 1$	4	4793	36.29	0.951	32.20	0.904	31.10	0.884	28.91	0.886	35.03	0.963
FSRCNN	2	PSNR	$D = 56, S = 12, M = 4$	2	5566	36.74	0.954	32.45	0.907	31.34	0.887	29.42	0.895	35.87	0.967
FSRCNN	2	speed	$D = 6, S = 3, M = 1$	5	3560	35.36	0.943	31.52	0.898	30.64	0.878	28.01	0.870	33.13	0.951
Bicubic	3	–	–	38	1705	29.24	0.849	26.73	0.757	26.64	0.720	23.84	0.717	25.87	0.838
eSR	3	PSNR	CNN: $C = 8, D = 3, S = 15$	11	5873	32.75	0.906	29.27	0.820	28.36	0.782	26.14	0.796	30.16	0.907
eSR	3	speed	TM: $k = 3, C = 1$	60	2632	29.77	0.853	27.31	0.780	27.08	0.748	24.31	0.742	26.56	0.851
ESPCN	3	default [33]	$D = 64, S = 32$	13	6027	32.73	0.905	29.26	0.821	28.36	0.783	26.12	0.795	30.36	0.908
ESPCN	3	PSNR	$D = 22, S = 32$	16	4945	32.77	0.906	29.30	0.821	28.38	0.784	26.15	0.797	30.37	0.909
ESPCN	3	speed	$D = 0, S = 21$	26	4176	31.30	0.889	28.51	0.808	27.82	0.772	25.41	0.774	28.29	0.887
FSRCNN	3	default [6]	$D = 32, S = 6, M = 1$	8	4640	32.43	0.901	29.07	0.816	28.22	0.780	25.82	0.787	29.61	0.899
FSRCNN	3	PSNR	$D = 56, S = 12, M = 4$	5	5566	32.74	0.906	29.25	0.820	28.35	0.784	26.10	0.797	30.13	0.907
FSRCNN	3	speed	$D = 6, S = 1, M = 1$	9	3560	31.30	0.879	28.31	0.803	27.75	0.768	25.06	0.761	27.98	0.870
Bicubic	4	–	–	74	2170	28.60	0.808	26.09	0.705	26.02	0.672	23.17	0.660	24.96	0.787
eSR	4	PSNR	CNN: $C = 8, D = 9, S = 6$	13	7100	30.62	0.860	27.48	0.751	26.93	0.714	24.42	0.718	27.27	0.845
eSR	4	speed	MAX: $k = 3, C = 2$	94	3867	28.64	0.806	26.12	0.712	26.13	0.684	23.28	0.668	25.08	0.789
ESPCN	4	default [33]	$D = 64, S = 32$	23	6952	30.57	0.858	27.50	0.752	26.92	0.715	24.42	0.718	27.44	0.848
ESPCN	4	PSNR	$D = 16, S = 32$	29	4640	30.59	0.859	27.53	0.753	26.95	0.715	24.43	0.719	27.46	0.849
ESPCN	4	speed	$D = 1, S = 3$	45	3096	28.93	0.820	26.49	0.725	26.25	0.694	23.56	0.680	25.49	0.804
FSRCNN	4	default [6]	$D = 32, S = 6, M = 1$	12	4795	30.16	0.845	27.19	0.742	26.74	0.707	24.09	0.702	26.63	0.826
FSRCNN	4	PSNR	$D = 44, S = 12, M = 4$	9	5257	30.61	0.861	27.52	0.753	26.94	0.716	24.44	0.721	27.40	0.849
FSRCNN	4	speed	$D = 6, S = 1, M = 1$	14	3715	29.31	0.823	26.62	0.730	26.41	0.699	23.62	0.683	25.72	0.802

ods are better at image quality, with ESPCN achieving the best speed in the high quality range. Interestingly, eSR-CNN does not improve ESPCN at high quality and barely improves it at high speed, despite using the same multi-layer configuration. eSR-MAX shows the best performance at high speeds but it is unable to make significant improvements in image quality. eSR-TM and eSR-TR show the best performance at intermediate speed and image quality. They perform very similar with a slight but not conclusive advantage of eSR-TR on GPU devices. FSRCNN shows the worst performance at 2× factor with no improvements in speed as image quality decreases. One possible reason for this result is that the higher network depth of FSRCNN might become a disadvantage at 2× where large receptive fields are unnecessary.

The bold values in Table 2 highlight the best metrics for different columns, ignoring bicubic. edge-SR systems reach the best speed and lowest power consumption except for 4× where ESPCN gets better. They also succeed to improve bicubic’s image quality for small upscaling factors.

Filters. In Figure 10 we display the step by step processing of 2× upscaling using eSR-TM with kernel size $k = 7$ and $C = 4$ number of matching/upscaling filters. Here, we used equation (1) to reconstruct the 4 matching/upscaling filters from the efficient implementation containing $4 \cdot 2^2 = 16$ filters. In addition to the filter coefficients we also display the FFT computed using a Kaiser-Bessel window for

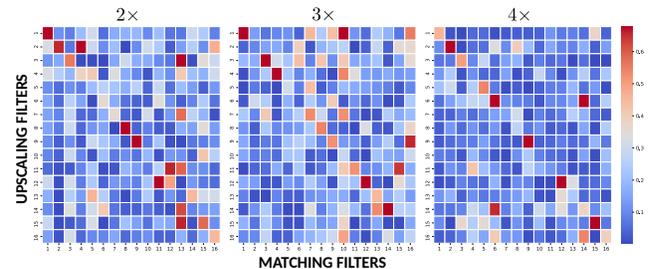


Figure 8. Correlations between *upscaling* and *matching* filters in eSR-TM $k = 7, C = 16$. Higher correlations along the diagonal mean that the model is performing template matching, with upscaling and matching filters that resemble a common template.

better frequency visualization [32]. The output for this particular image is about 1.5 dB better than the bicubic output and it is displayed next to the outputs of ESPCN and FSRCNN models with similar image quality. Here, eSR-TM achieves roughly the same speed of bicubic upscaler.

The efficient filters use kernel size $k \times k$, and after multiplexing them with a pixel-shuffle layer we can recover the original filters of size $sk \times sk$. Thus, the filter sizes of eSR models grows with the upscaling factors as seen in Figure 9. The filter coefficients in frequency domain show that each filter is processing different frequency bands. Although the filters are not smooth, they do show a level of discrimination between different directions.

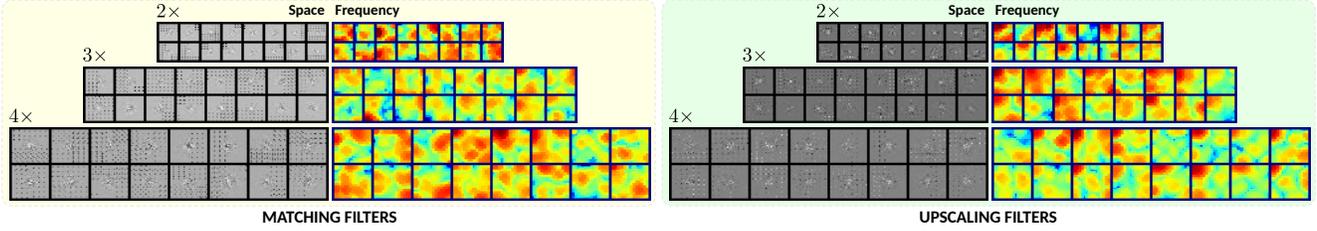


Figure 9. Matching and upscaling filters obtained after training a one-layer architecture eSR-TM with kernel size $k = 7$ and $C = 18$ number of filters for $2\times$, $3\times$ and $4\times$ upscaling factors. Filters are displayed in the original spatial format as well as in frequency domain by using FFT visualization. The filters do not change smoothly within a single filter but show diverse directionality among different filters.

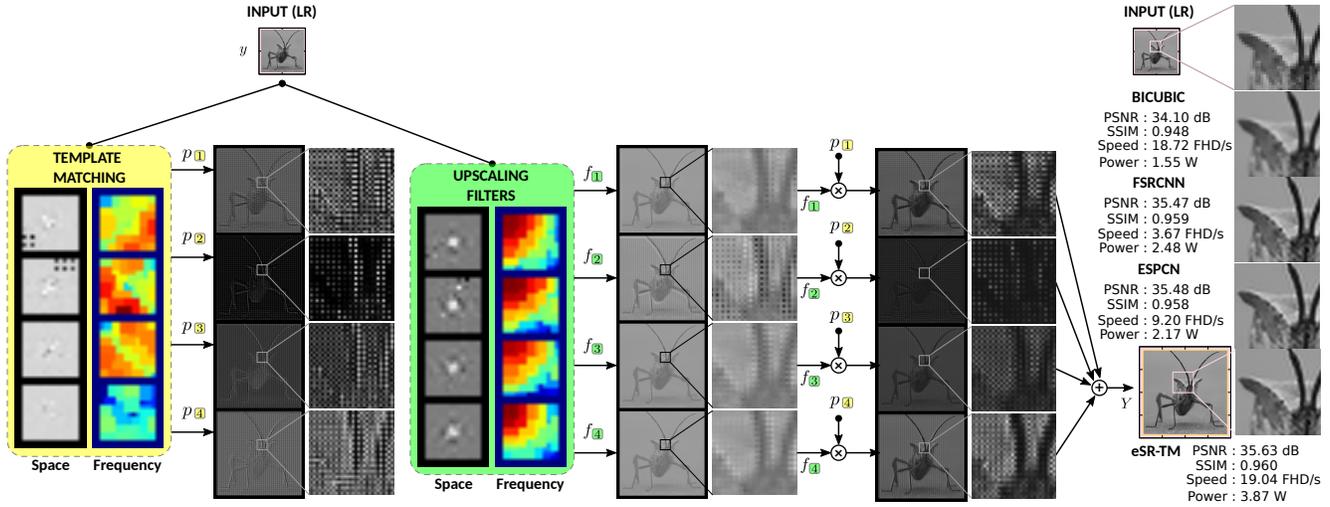


Figure 10. Inspection of all intermediate outputs and filter coefficients for the eSR-TM $2\times$ architecture with kernel size $k = 7$ and $C = 4$ number of matching/upscaling filters. The diagram follows the interpretation in Figure 5. Filters are displayed both in the original spatial format as well as in frequency domain by using FFT visualization. Each of the 4 branches is focusing on a particular sub-pixel array.

Now, moving one step inside the network from the output in Figure 10, we observe that the 4 components of the sum are clearly focusing on different sub-pixel images. This pattern is also visible in the outputs of upscaling filters and template matching modules. Both matching and upscaling filters are not smooth and also show signs of different sub-pixel processing with some degree of directionality. This indicates that the different branches of the single convolutional layer used in eSR-TM are solving the upscaling problem independently for each sub-pixel image. This is in contrast with the smooth scaling filters used in the classical edge-directed interpolation [2, 18] and also compared to smooth directional filters observed in CNNs super-resolution interpretations in [28]. Next, in Figure 8 we compute the Pearson correlation between upscaling and matching filters for eSR-TM with $k = 7$ and $C = 16$. The results show dominant correlations along the diagonal, stronger for $2\times$ factor and reducing strength towards $4\times$ factor. Strong correlations along the diagonal indicate a template matching strategy where upscaling and matching filters are simi-

lar for the same pattern and different to other patterns (see Figure 5). Thus, we confirm that the training process has a tendency to converge towards a template matching strategy that is particularly strong for small upscaling factors.

5. Conclusions

The current trend in Edge-AI chips offers the chance to deploy efficient AI solutions at massive scale. But there is a vast range of performance requirements for which these solutions are unavailable for image SR. We propose the edge-SR architectures with the aim to fill the gap between classic and deep learning upscalers. We performed an exhaustive search among more than a thousand different models identically trained, revealing the gap between classic upscalers and deep-learning solutions. Our edge-SR configurations using a single convolutional layer showed promising results to fill this gap for small upscaling factors. The simplicity of the model also makes it interpretable and allows to visualize and understand all the intermediate steps of the process.

References

- [1] Eirikur Agustsson and Radu Timofte. NTIRE 2017 challenge on single image super-resolution: Dataset and study. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, July 2017.
- [2] V. R. Algazi, G. E. Ford, and R. Potharlanka. Directional interpolation of images based on visual properties and rank order filtering. In *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, volume 4, page 3005–3008, Toronto, ON, May 1991. IEEE Signal Processing Society.
- [3] Mustafa Ayazoglu. Extremely lightweight quantization robust real-time single-image super resolution for mobile devices. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2472–2479, 2021.
- [4] Jiayi Chen and Xukan Ran. Deep learning with edge computing: A review. *Proc. IEEE*, 107(8):1655–1674, 2019.
- [5] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *in Proceedings of European Conference on Computer Vision (ECCV)*, 2014.
- [6] Chao Dong, Chen Change Loy, and Xiaoou Tang. Accelerating the super-resolution convolutional neural network. In *in Proceedings of European Conference on Computer Vision (ECCV)*, 2016.
- [7] Ian Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. Maxout networks. In *International conference on machine learning*, pages 1319–1327. PMLR, 2013.
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [9] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7132–7141, 2018.
- [10] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [11] Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. Single image super-resolution from transformed self-exemplars. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5197–5206, 2015.
- [12] Andrey Ignatov, Radu Timofte, Maurizio Denna, and Abdel Younes. Real-time quantized image super-resolution on mobile NPUs, Mobile AI 2021 challenge: Report. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2525–2534, 2021.
- [13] Alex James. The why, what and how of artificial general intelligence chip development. *IEEE Transactions on Cognitive and Developmental Systems*, 2021.
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 2015.
- [15] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553):436–444, May 2015.
- [16] Royson Lee, Stylianos I Venieris, Lukasz Dudziak, Sourav Bhattacharya, and Nicholas D Lane. Mobisr: Efficient on-device super-resolution through heterogeneous mobile processors. In *The 25th Annual International Conference on Mobile Computing and Networking*, pages 1–16, 2019.
- [17] Bingzhen Li, Jiaojiao Gu, and Wenzhi Jiang. Artificial intelligence (AI) chip technology review. In *2019 International Conference on Machine Learning, Big Data and Business Intelligence (MLBDBI)*, pages 114–117. IEEE, 2019.
- [18] Xin Li and Michael T. Orchard. New edge-directed interpolation. *IEEE Transactions on Image Processing*, 10(10):1521–1527, October 2001.
- [19] Jingyun Liang, Jie Zhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. SwinIR: Image restoration using swin transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1833–1844, 2021.
- [20] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, July 2017.
- [21] Xin Liu, Yuang Li, Josh Fromm, Yuntao Wang, Ziheng Jiang, Alex Mariakakis, and Shwetak Patel. Splitsr: An end-to-end approach to super-resolution on mobile devices. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 5(1):1–20, 2021.
- [22] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *International Conference on Computer Vision (ICCV)*, 2021.
- [23] S. Mallat. *A Wavelet Tour of Signal Processing*. Academic Press, 1998.
- [24] S. Mallat and G. Peyre. A review of bandlet methods for geometrical image representation. *Numerical Algorithms*, 44(3):205–234, March 2007.
- [25] David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, volume 2, pages 416–423. IEEE, 2001.
- [26] Yusuke Matsui, Kota Ito, Yuji Aramaki, Azuma Fujimoto, Toru Ogawa, Toshihiko Yamasaki, and Kiyoharu Aizawa. Sketch-based manga retrieval using manga109 dataset. *Multimedia Tools and Applications*, 76(20):21811–21838, 2017.
- [27] P. Navarrete and H. Liu. Upscaling beyond super-resolution using novel a deep-learning system. In *GPU Technology Conference 2017*, Talk S7231, San Jose, CA, USA, May 2017.
- [28] Pablo Navarrete Michelini, Hanwen Liu, Yunhua Lu, and Xingqun Jiang. A tour of convolutional networks guided by linear interpreters. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4753–4762, 2019.
- [29] P. Navarrete Michelini, L. Zhang, and J. He. Upscaling with deep convolutional networks and MuxOut layers. In *GPU*

- Technology Conference 2016*, Poster P6324, San Jose, CA, USA, May 2016.
- [30] S.C. Park, M.K. Park, and M.G. Kang. Super-resolution image reconstruction: a technical overview. *Signal Processing Magazine, IEEE*, 20(3):21–36, May 2003.
- [31] Terence Parr and Jeremy Howard. The matrix calculus you need for deep learning. *arXiv preprint arXiv:1802.01528*, 2018.
- [32] J.G. Proakis and D.G. Manolakis. *Digital Signal Processing*. Prentice Hall international editions. Pearson Prentice Hall, 2007.
- [33] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1874–1883, 2016.
- [34] M. Sugawara, Choi S-Y, and D. Wood. Ultra-high-definition television (Rec. ITU-R BT.2020): A generational leap in the evolution of television. *Signal Processing Magazine, IEEE*, 31(3):170–174, May 2014.
- [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [36] Xiaolong Wang, Ross B Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. *computer vision and pattern recognition*, pages 7794–7803, 2018.
- [37] Yan Wu, Zhiwu Huang, Suryansh Kumar, Rhea Sanjay Sukthanker, Radu Timofte, and Luc Van Gool. Trilevel neural architecture search for efficient single image super-resolution. *arXiv preprint arXiv:2101.06658*, 2021.
- [38] Jianchao Yang, John Wright, Thomas S Huang, and Yi Ma. Image super-resolution via sparse representation. *IEEE Transactions on Image Processing*, 19(11):2861–2873, 2010.
- [39] Roman Zeyde, Michael Elad, and Matan Protter. On single image scale-up using sparse-representations. In *International conference on curves and surfaces*, pages 711–730. Springer, 2010.
- [40] Kai Zhang, Martin Danelljan, Yawei Li, Radu Timofte, Jie Liu, Jie Tang, Gangshan Wu, Yu Zhu, Xiangyu He, Wenjie Xu, et al. Aim 2020 challenge on efficient super-resolution: Methods and results. In *European Conference on Computer Vision*, pages 5–40. Springer, 2020.
- [41] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 286–301, 2018.
- [42] Yulun Zhang, Kunpeng Li, Kai Li, Bineng Zhong, and Yun Fu. Residual non-local attention networks for image restoration. *International Conference on Learning Representations*, 2019.
- [43] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2472–2481, 2018.