

Bicubic++: Slim, Slimmer, Slimmest Designing an Industry-Grade Super-Resolution Network

Bahri Batuhan Bilecen Mustafa Ayazoglu
Aselsan Research
Ankara, Türkiye

{batuhanb,mayazoglu}@aselsan.com.tr

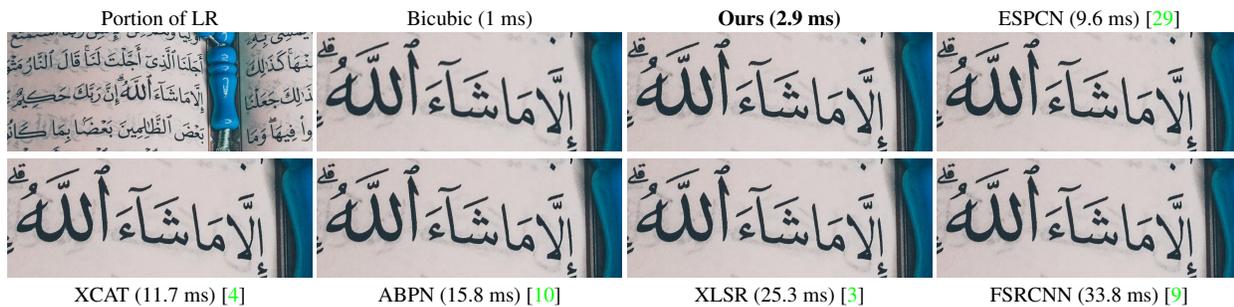


Figure 1. Image 70 from the NTIRE 2023 RTSR Challenge [7] Track 2 $\times 3$ test set [36], where the data is JPEG Q=90 degraded. Our practical method runs much faster (< 3 ms) and yields nearly identical results to the other comparable (~ 10 -30ms) models on RTX3070.

Abstract

We propose a real-time and lightweight single-image super-resolution (SR) network named Bicubic++. Despite using spatial dimensions of the input image across the whole network, Bicubic++ first learns quick reversible downgraded and lower resolution features of the image in order to decrease the number of computations. We also construct a training pipeline, where we apply an end-to-end global structured pruning of convolutional layers without using metrics like magnitude and gradient norms, and focus on optimizing the pruned network’s PSNR on the validation set. Furthermore, we have experimentally shown that the bias terms take considerable amount of the runtime while increasing PSNR marginally, hence we have also applied bias removal to the convolutional layers. Our method adds ~ 1 dB on Bicubic upscaling PSNR for all tested SR datasets and runs with ~ 1.17 ms on RTX3090 and ~ 2.9 ms on RTX3070, for 720p inputs and 4K outputs, both in FP16 precision. Bicubic++ won NTIRE 2023 RTSR Track 2 $\times 3$ SR competition and is the fastest among all competitive methods. Being almost as fast as the standard Bicubic up-sampling method, we believe that Bicubic++ can set a new industry standard.

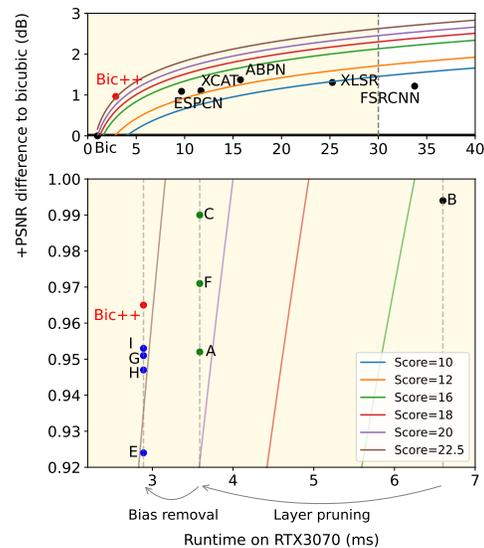


Figure 2. Runtime versus Δ PSNR wrt. Bicubic, and the effect of the proposed training pipeline. Due to Eq. (1), Bicubic score is 0. Bic++ scores the highest with < 3 ms runtime, achieving $\sim +1$ dB on top of Bicubic on our test set. A..I denote the models experimented. We follow $\mathbf{B} \rightarrow \mathbf{B}^* \rightarrow \mathbf{C} \rightarrow \mathbf{Bic++}$. \mathbf{B}^* is not in plot’s range.

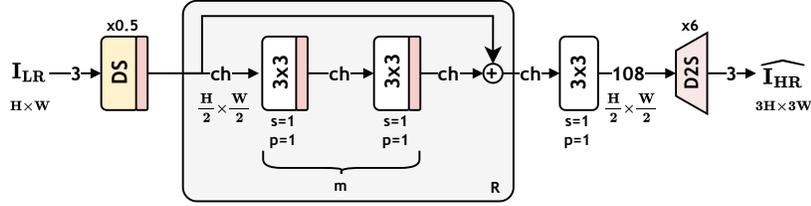


Figure 3. General network structure of Bicubic++. DS denotes the downscaling sub-network (Fig. 5), and D2S is the depth-to-space layer [29]. Residual addition blocks are shown by \mathbf{R} , and the number of convolutional layers inside \mathbf{R} is m . s and p denote the stride and padding of the convolutional layers, respectively. ch is the global channel number parameter. In the final proposed model, $\mathbf{R}=1$, $m=2$, $ch=32$, and DS is a 3×3 convolutional layer with $s=2$, $p=1$.

1. Introduction

Single-image super-resolution (SR) is the task of upscaling a low-resolution (LR) image to a high-resolution (HR) one. SR is an highly ill-posed inverse problem, and researchers are still coming up with different approaches which would yield reasonable and high-quality HR images under varying LR conditions (noise, blur, camera motion, compression, *etc.*). Besides providing high-quality HR's, many SR networks also focus on runtime performance [15, 16] since SR is also extensively utilized in real-life scenarios such as on the cameras of unmanned vehicles and surveillance systems. Nonetheless, a lot of industry-grade applications still consider Bicubic as the tried-and-true upscaling method because of its speed and ease-of-use. Despite all, we believe that a carefully constructed and trained neural network can be just as useful as the traditional Bicubic upsampling, run in real-time, and even replace Bicubic upsampling with visually more pleasing outputs.

Hence, to contribute to the ongoing development of real-time single-image super-resolution, we present our work in the *NTIRE 2023 Real-Time Super-Resolution Challenge Track 2 ($\times 3$ upscaling)* [7, 36] and propose a lightweight Bicubic upscaling alternative, **Bicubic++**. Bicubic++ first downscales the image features by $\times 2$ to reduce the number of operations significantly, and applies $\times 6$ upscaling at the end. We also train the networks with our proposed three-stage training pipeline, where we first train a network with convolutional layer channels greater than the "hardware's sweet spot". Then, we apply global structured layer pruning without conditioning on weight or gradient norms while focusing on optimizing the PSNR, convolutional layer bias removal, and fine-tuning operations to further decrease its runtime speed without sacrificing much from its visual output quality. Bicubic++ ranked 1st on the *NTIRE 2023 Real-Time Super-Resolution Challenge Track 2* and is the fastest method among all competitive methods.

Throughout the paper, we emphasise on optimizing the

score function of the challenge given in Eq. (1):

$$S(P, t) = \begin{cases} 0 & \text{if } P \leq P_{bic} \\ \frac{2^{P-P_{bic} \times 2}}{0.1 \times \sqrt{t}} & \text{else} \end{cases} \quad (1)$$

where P and P_{bic} are the PSNR values of the network and the Bicubic upsampling on the test set, respectively. t is the runtime of the network when a 720p LR image is given as an input and a 4K HR is obtained. In addition, due to the real-time requirements, t must be smaller than 30ms. This score formulation is provided by the competition holders and is calculated with their given code¹.

2. Related Works

Deep learning-based single-image super-resolution.

Dong *et al.* proposed SRCNN [8] as the first deep learning based SR algorithm. Later on, FSRCNN [9] was proposed, as a faster version for SRCNN. It replaced ReLU with a PReLU activation layer, reformed SRCNN with smaller filter sizes with more mapping layers, and proposed postponing the upscaling layer at the end of the network which resulted in a great speed-up. ESPCN [29] introduced a sub-pixel convolutional layer, often known as depth to space, which is currently used in several efficient SR networks [3, 4, 10]. The development of deep learning-based SR was carried on by VDSR [18], EDSR [23], and WDSR [35] by expanding the number of parameters in exchange for accuracy and sacrificing speed. Concepts like generative adversarial networks [20, 30, 34], recursive & residual networks [2, 5, 14], and attention mechanism [22, 28] also rapidly took place in SR network topologies with the recent advancements [27].

Efficient super-resolution. Just like the regular SR, efficient SR did not only emerge with deep learning either, as there are notable methods using sparse representations like Zeyde *et al.*'s [37], ANR [32], and A+ [33]. However, thanks to the challenges and workshops held recently [7, 15, 16, 21, 36], there is an abundance of deep

¹<https://github.com/eduardzamfir/NTIRE23-RTSR>

learning-based SR methods now, each more efficient than the previous ones most of the time. XLSR [3] utilized clipped ReLU at the end of the network to minimize the PSNR decrease in post-training INT8 quantization for the first time. ABPN [10] proposed an anchor-based residual network and applied quantization-aware training. SC-SRN [11] applied a similar idea of residual adding to ABPN but in the feature level, and performed reparametrization of the convolutional layers to train with larger number of parameters and reduce them in the inference time. RLFN [19] used three convolutional layers for residual learning to simplify feature aggregation, and stated to have achieved a balance between visual quality and runtime. A CPU-inference based method named SR-LUT [17] proposed to train a deep SR network to construct a lookup table, and matched the input LR image patches to the output HR image patches utilizing the said table.

3. Methodology

Our proposed network structure is given in Fig. 3 in its most general form. The architecture and the design choices are inspired from the efficient SR methods [15, 16, 21], and the experiments conducted. Based on this architecture, we tune its parameters and try different training strategies to maximize the scoring function given in Eq. (1).

3.1. Design Procedure and Network Architecture

One initial observation of ours during the design procedure was that the speed of the network does not directly depend on the number of parameters, but on the amount of activations as stated in [21], which is a value proportional to the volume of data passing through the network. Hence, processing the feature tensors with lower width & height than those of the input image dimensions can decrease the number of activations significantly. Furthermore, it has been shown in [31] that the downscaling operation has positive effects when learned together with the upscaling. Therefore, although it seems to complicate the problem at first, we chose to "downscale" the features ($\times 0.5$) by a strided convolution and then apply super-resolution operation ($\times 6$) in order to speed up the network. We also experiment with other downscaling approaches (Fig. 5), but go for a strided convolution at the end. In a sense, this lower dimensional feature extraction approach can be thought as a way of compressing the data spatially; however, this compression is reversible since it is trained all together with the upscaling.

In our network, we wanted to use the "optimum" number of channels for the convolutional layers. We observed that decreasing the number of channels of consecutive layers does not always reflect to a decrease in the runtime. Hence, we decided to keep the same number of channels across the network (unless necessary), even though some

SR methods do the quite opposite by squeeze and expand blocks [12] to extract features better. However, the increase in runtime caused by the blocks with varying channels outweigh its benefits considerably. To get a feeling of the optimum channel numbers, we conducted a simple experiment and observed that runtime versus the number of channels is not positively correlated all the time, and there exists some "sweet spots" (Fig. 4). For our setup and under the real-time constraints, 56 channels used by RFDN [24] and 28 channels used by ABPN [10] in the literature came out to be not the optimal points. Thus, among the options shown in Fig. 4, we opted for a channel number of 32 in our final model. However, as detailed in Sec. 4.1, to get the most out of the model, we will start from a non-optimal channel (34) and globally prune the model to achieve 32 channels.

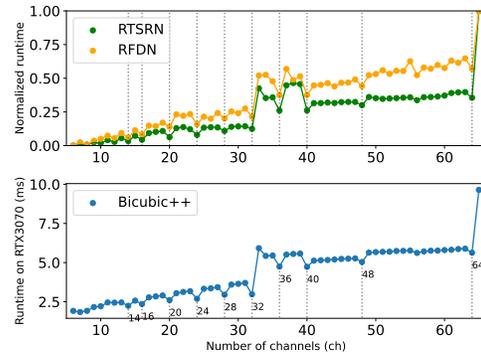


Figure 4. The effect of the channel number on the runtime and normalized runtime (runtime of $ch=6$ is mapped to the origin, and $ch=65$ to 1) of our model, RTSRN [15], and RFDN [24]. Note the "sweet spots" are on $ch \in \{16, 20, 24, 28, 32, \dots\}$, and are observed in common among different models.

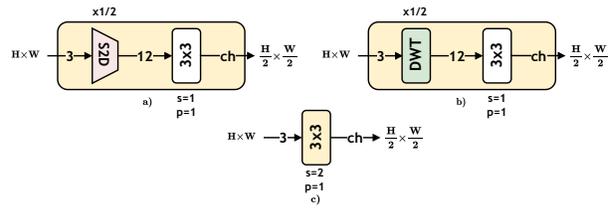


Figure 5. Different DS structures in Fig. 3. S2D (a) denotes the space-to-depth layer [29], and DWT (b) is the discrete wavelet transform. We use a strided convolution (c) for DS in the end.

In addition to all these, we also noticed that the innocent-looking bias terms of the convolutional layers are usually overlooked, and actually take a considerable portion of the overall runtime. However, the drastic decrease in the runtime caused by removing the bias terms actually overpower the marginal decrease in PSNR, which is also clearly visible from the obtained score values (Fig. 2). Hence, we wanted

to eliminate the bias terms from the layers in our final network structure as well.

To sum up those aforementioned observations and ideas:

- Data volume processed is correlated with runtime rather than number of parameters
- Varying the number of channels among blocks decreases runtime performance
- There exists some "sweet spots" for the channel size for the hardware, and reaching those "sweet spots" can be done by training a larger network and pruning afterwards
- The bias terms may take a considerable amount of time

In the end, to apply these points to our proposed network, we constructed a training pipeline. For each choice, we empirically proved its benefits. We also detail the pipeline and empirical justifications in Sec. 3.2 and Sec. 4.1, respectively.

3.2. Three-Stage Training Pipeline

We construct a three-stage training pipeline: training with a larger number of channels, global structured pruning of channels, and bias removal of convolutional layers. For the last two stages, we fine-tune the network after the modifications.

Training hyperparameters & computational issues:

For each stage, we use Adam optimizer with parameters $\beta_{1,2}$ 0.99 and 0.999, respectively. Models in all stages are trained for 1000 epochs, each epoch consuming 800 randomly cropped and rotated patches of dimension (108,108,3) (for LR) from JPEG Q=90 degraded DIV2K training [1] dataset. For the learning rate (LR), we initialize with $5e-4$ and use a decaying learning rate scheduler for all stages, where for the first 500 epochs the LR stays constant and for the last 500 epochs the LR decays linearly until $1e-8$.

To be able to train such a tiny model by utilizing the training hardware to its fullest, all the elements of the training pipeline need to be adjusted properly. This is due to the model forward and backward passes being done in almost no time, and memory accesses along with the validation steps becoming the actual bottleneck. To speed-up the validation, we use 48 images with same dimensions (680,452,3) (for LR) from JPEG Q=90 degraded DIV2K validation dataset, and hence be able to pick a batch size of 8 instead of 1. To take care of the memory access bottleneck, we pre-loaded both the training and the validation images to RAM and avoided fetching the data from disk continuously for each training and validation step. This way, we were able to complete a training cycle of 1000 epochs in only 1 hour using a single Tesla V100. Note that since the model's

capacity is small, opting for a comparably smaller training dataset (like DIV2K) instead of favoring for a larger dataset is usually viable, and does not hurt the validation score in the end generally.

Stage 1 - Slim: Training with a larger number of channels. We train the network in Fig. 3, where ch is 34, m is 2, R is 1, and DS is the strided convolutional layer in Fig. 5. All convolutional layer biases are enabled in this stage. Respective network parameter choices are justified and explained thoroughly in Sec. 4.1.

Stage 2 - Slimmer: Global structured pruning. We take the resulting checkpoint from Stage 1 and perform the proposed pruning procedure given in Fig. 6. The proposed method is a structured pruning (prune channel-wise instead of kernel element-wise), in a global sense, since we prune the entire network considering the entire network's structure through sharing and transferring pruning masks and using a global fitness measure (*i.e.* using validation PSNR of the pruned model), instead of focusing on individual elements of the network and heuristic measures for fitness.

In detail; for the convolutional layers to be pruned, we first construct all possible Mask #1's in Fig. 6 (\mathcal{M}_1), each only masking a single channel. After obtaining different \mathcal{M}_1 's, we calculate the PSNR scores on the validation set for each mask applied individually, and unify two \mathcal{M}_1 's (to reach the closest "sweet spot", $ch=32$) which cause the most marginal PSNR drop as our best & final \mathcal{M}_1 .

After constructing \mathcal{M}_1 , we move onto the Mask #2 (\mathcal{M}_2). We follow the same procedure as in the creation of \mathcal{M}_1 ; however, now we do not start from the beginning and keep the \mathcal{M}_1 's modifications (though starting from scratch resulted the same set of masks). We do not perform any additional training while \mathcal{M}_1 and \mathcal{M}_2 are being selected, but only after \mathcal{M}_1 and \mathcal{M}_2 are all applied. We do not prune the bias terms here, either. One observation we made here was the network's overall performance was a lot more sensitive to the channels in \mathcal{M}_1 than those in \mathcal{M}_2 .

Stage 3 - Slimmest: Bias removal of the convolutional layers. After applying \mathcal{M}_1 & \mathcal{M}_2 and reducing ch down from 34 to 32, we take the checkpoint from Stage 2, disable all bias from the convolutional layers, and fine-tune the overall network again.

4. Experiments

We have performed extensive experiments trying to find the optimal parameters for our general network structure in Fig. 3, with the given scoring function in Eq. (1). In addition to network structure and the proposed training pipeline in Sec. 3, we experiment with partial loading of model components to gradually decrease the size and runtime of the network. All PSNR scores in the ablation study are evaluated on our 48-image JPEG Q=90 degraded DIV2K validation set, unless stated otherwise.

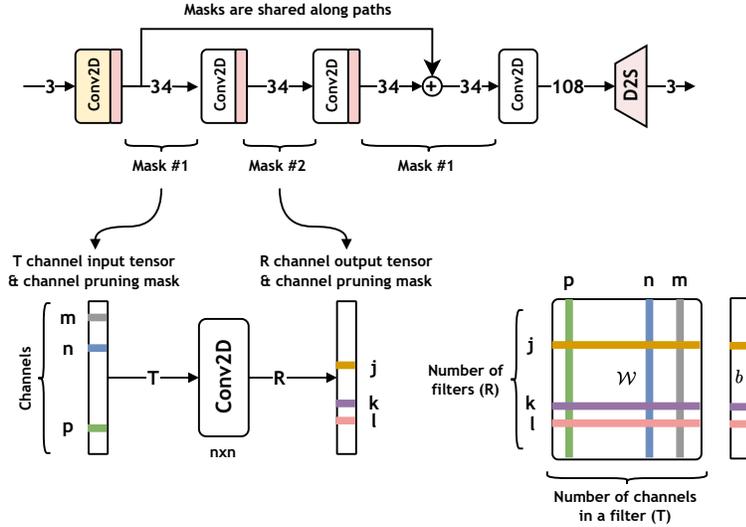


Figure 6. Global structured pruning mask generation structure. $\mathcal{W}_{R \times T}$ is the matrix consisting of the convolutional filters for a specific Conv2D layer with kernel size $n \times n$, and b is the corresponding bias term of the layer. Each entry of \mathcal{W} is an $n \times n$ kernel.

Exp	Models	DS	Act	ch	m	R	b	PSNR (Y) (dB)	Runtime (ms)	Score	Load	Operation
#1: DS, b, and activation	Q	SC	ReLU	24	2	1	✓	30.157	3.19	19.86	✗	✗
	X	DWT	ReLU	24	2	1	✓	30.219	3.59	19.55	✗	✗
	Y	S2D	ReLU	24	2	1	✓	30.172	3.33	19.64	✗	✗
	Z	S2D	ReLU	32	2	1	✓	30.269	3.72	19.88	✗	✗
	W	DWT	ReLU	32	2	1	✓	30.326	4.14	19.60	✗	✗
	T	SC	ReLU	32	2	1	✓	30.271	3.57	20.32	✗	✗
	A	SC	LReLU	32	2	1	✓	30.282	3.58	20.43	✗	✗
	E	SC	LReLU	32	2	1	✗	30.254	2.89	22.32	✗	✗
#2: m and R	E	SC	LReLU	32	2	1	✗	30.254	2.89	22.32	✗	✗
	1	SC	LReLU	32	2	2	✓	30.410	4.95	19.00	✗	✗
	2	SC	LReLU	32	4	1	✓	30.433	4.84	19.53	✗	✗
	3	SC	LReLU	32	4	1	✗	30.424	4.15	20.96	✗	✗
	4	SC	LReLU	32	2	2	✗	30.390	4.30	20.11	✗	✗
#3: ch selection for pruning	P	SC	LReLU	35	2	1	✓	30.293	6.50	15.29	✗	✗
	P*	SC	LReLU	32	2	1	✓	29.712	3.58	13.77	P	Global structured pruning
	R	SC	LReLU	33	2	1	✓	30.286	6.16	15.63	✗	✗
	R*	SC	LReLU	32	2	1	✓	30.113	3.58	18.44	R	Global structured pruning
	B	SC	LReLU	34	2	1	✓	30.324	6.60	15.50	✗	✗
	B*	SC	LReLU	32	2	1	✓	30.155	3.58	18.72	B	Global structured pruning
#4: Fine tuning and bias removal	C	SC	LReLU	32	2	1	✓	30.320	3.58	20.99	B*	Fine tuning (FT)
	(Bic++)	SC	LReLU	32	2	1	✗	30.295	2.89	22.96	C	Bias removal + FT
	F	SC	LReLU	32	2	1	✓	30.301	3.58	20.72	A	Fine tuning
	G	SC	LReLU	32	2	1	✗	30.281	2.89	22.74	A	Bias removal + FT
	H	SC	LReLU	32	2	1	✗	30.277	2.89	22.68	E	Fine tuning
	I	SC	LReLU	32	2	1	✗	30.283	2.89	22.77	F	Bias removal + FT
	J	SC	LReLU	32	2	1	✗	30.290	2.89	22.88	H	Fine tuning
	#5	5	SC	LReLU	32	2	1	✓	30.210	3.58	19.45	2
6		SC	LReLU	32	2	1	✓	30.281	3.58	20.43	1	Partial load \mathbf{R} 's + FT
Other methods	Bicubic	-	-	-	-	-	-	29.334	1.0	0	-	-
	ESPCN [29]	-	-	-	-	-	-	30.419	9.68	13.67	-	-
	XCAT [4]	-	-	-	-	-	-	30.435	11.68	12.58	-	-
	ABPN [10]	-	-	-	-	-	-	30.703	15.76	13.05	-	-
	XLSR [3]	-	-	-	-	-	-	30.637	25.25	9.84	-	-
	FSRCNN [9]	-	-	-	-	-	-	30.547	33.75	8.00	-	-
	RFDN [24]	-	-	-	-	-	-	30.921	159.3	4.77	-	-

Table 1. Quantitative ablation study to pick the best scoring model. \mathbf{b} denotes the bias of convolutional layers, and LReLU is the leaky ReLU activation. \mathbf{DS} configurations are given in Fig. 5. We follow $\mathbf{B} \rightarrow \mathbf{B}^* \rightarrow \mathbf{C} \rightarrow \mathbf{Bic++}$ for the final model. PSNR scores are evaluated on the 48-image DIV2K validation dataset described in Sec. 3.2. Runtimes are measured on RTX3070.

4.1. Ablation Studies

We provide an extensive ablation study in Tab. 1. We also divide the experiments in Tab. 1 into 5 different sections for easy readability and comparability:

Downscaling (DS), bias (b), and activation selection. From **Q-T**, we observed that using a *strided convolution (SC)* is superior compared to discrete wavelet transform (DWT) and space-to-depth (S2D), even when the number of channels are smaller (notice that **Q&Z** nearly got the same score despite **Q** having less number of channels). **T&A** reveal that using leaky ReLU instead of ReLU can provide a significant benefit in PSNR without nearly any increase in the runtime. We can also see from **A&E** that disabling bias terms can be advantageous to get higher scores.

Number of residual blocks (R) and convolutional layers (m). Results of **1-4** and **E** suggest that setting $m=2$ and $R=1$ yields the best score among tried model architectures. Hence, after the experiments done up to this stage and the observation in Fig. 4, we decided that for our final model, ch is 32, DS is SC, activations are leaky ReLU, and the bias terms should be disabled at some point.

Selection of number of channels for pruning. This part includes the experiments done for searching the ideal number of channels for reducing to 32. We mainly focus on $ch \in \{33, 34, 35\}$, and apply our proposed pruning method (Sec. 3.2) for **P, R, B**, and obtain **P*, R*, and B***, respectively. The results reveal that choosing $ch=34$ and applying the pruning obtains the highest score, hence we move on from **B***.

Fine tuning and bias removal. After training **B** and pruning to obtain **B***, we fine-tune the network by training it again, and obtain **C**. The positive effect of pruning can be observed from the comparison of **C&F**.

Afterwards, we remove the bias terms from **C**, fine-tune it once more, and obtain our proposed model, **Bicubic++**. **I**, cannot reach the final network’s score, indicating the advantages of the proposed training pipeline. We also show other possible training paths (**A** \rightarrow **G**, **E** \rightarrow **H**, **A** \rightarrow **F** \rightarrow **I**, **A** \rightarrow **E** \rightarrow **H** \rightarrow **J**) to further point out that our training pipeline yields the best score. A visual comparison is also provided for this experimental setup in Fig. 2.

Extra experiments regarding partial loading of R&m. We also experiment with the residual blocks (**R**), the number of convolutional layers (**m**), and partial loading models with different **R&m** and provide the results.

It is worthy to note that one can always increase the number of channels & blocks and construct a model obtaining higher PSNR scores with slower runtime; however, we aimed to stay under $<3ms$ on RTX3070, hence decided to stick with the choices made in this paper.

4.2. Comparative Results

We provide quantitative and qualitative comparative results in Tab. 2, Tab. 3, Fig. 7, and Fig. 8. Bicubic++ achieves better results than the traditional Bicubic upsampling and runs much faster compared to the other relevant methods, making it a good candidate for real-time SR applications.

Although relevant, we did not include the results of sparse coding methods like A+ [33] and SR-LUT [17] since they are not compliant with the challenge requirements, either due to their slower runtime or lack of implementation in GPU.

	Track 2 Score	PSNR (RGB)	SSIM	PSNR (Y)	Runtime (ms)
Aselsan Research	31.26	32.06	0.8344	34.56	1.17
Team OV	29.63	32.17	0.8376	34.72	1.51
ALONG	28.57	32.18	0.8367	34.66	1.66
RTVSR	26.89	32.22	0.8372	34.77	1.96
Noah.TerminalVision	26.68	32.65	0.8455	35.10	3.64
NJUST-RTSR	23.51	32.25	0.8384	34.90	2.68
Antins.cv	23.44	32.63	0.8457	35.21	4.60
DFCDN Team	22.64	32.07	0.8371	34.63	2.25
Multimedia	21.55	32.33	0.8398	34.83	3.56
z6	20.90	32.59	0.8446	35.05	5.47
R.I.P. ShopeeVideo	15.67	32.84	0.8469	35.30	13.79
ECNU_SR	15.39	32.64	0.8458	35.17	10.75
Touch.Fish	11.55	32.67	0.8468	35.31	19.86
P.A.I.R	8.66	32.55	0.8441	35.04	30.03
SEU.CNII	6.68	31.85	0.8326	34.52	19.05
diSRupt	6.34	31.64	0.8292	34.25	16.00
Bicubic	0.00	31.30	0.8246	33.82	0.5
Organizer’s baseline [7]	14.01	31.74	0.8299	34.25	3.74

Table 2. All methods scoring above Bicubic in Track 2 [7] and the suggested baselines. PSNR is evaluated on the competition’s test set, and runtime is measured on RTX3090/3060 with FP16 precision.

5. Conclusion and Discussion

In this paper, we proposed our real-time SR network, Bicubic++. Our model learns downsampled features of the input image to increase efficiency, and is trained with our proposed three-stage pipeline where we apply global structured pruning and bias removal. We believe that our experimental results, observations, proposed architecture, and the training pipeline may help the development of real-time SR methods. We believe that Bicubic++ sets a new practical industry standard for upscaling tasks and it can be the contemporary alternative of Bicubic upscaling.

For further research, a fusion of deep learning and lookup table along with sparse representation methods can be investigated on GPUs. Furthermore, we believe that the proposed training pipeline would benefit from including reparametrization along with INT8 quantization.

	Set5 [6]		Set14 [38]		BSD100 [25]		Urban100 [13]		Manga109 [26]		DIV2K Val [1]		Runtime (ms)
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	
ESPCN [29]	31.56	0.8737	28.53	0.7833	27.80	0.7454	25.66	0.7680	29.29	0.8768	30.37	0.8351	9.6
XCAT [4]	31.50	0.8728	28.52	0.7843	27.81	0.7464	25.67	0.7697	29.33	0.8781	30.38	0.8360	11.7
ABPN [10]	31.93	0.8812	28.76	0.7910	27.99	0.7523	26.13	0.7862	30.14	0.8928	30.65	0.8421	15.8
XLSR [3]	31.84	0.8800	28.74	0.7900	27.95	0.7509	25.98	0.7812	29.94	0.8899	30.59	0.8406	25.3
FSRCNN [9]	31.65	0.8762	28.61	0.7863	27.86	0.7477	25.81	0.7736	29.52	0.8814	30.48	0.8377	33.8
Bicubic	29.98	0.8434	27.30	0.7529	26.99	0.7180	24.31	0.7196	26.54	0.8300	29.33	0.8127	1.0
Bicubic++	31.19	0.8656	28.35	0.7799	27.68	0.7431	25.49	0.7626	28.72	0.8653	30.24	0.8324	2.9

Table 3. Quantitative comparative results of our method and other comparable models for $\times 3$ SR. All low-resolution images in the datasets are JPEG Q90 degraded, and all PSNR values are calculated for the Y channel. Runtime is measured on RTX3070 with FP16 precision.



Figure 7. Qualitative comparative results of Bicubic++ and other relevant methods on the datasets [1, 13, 25, 26].

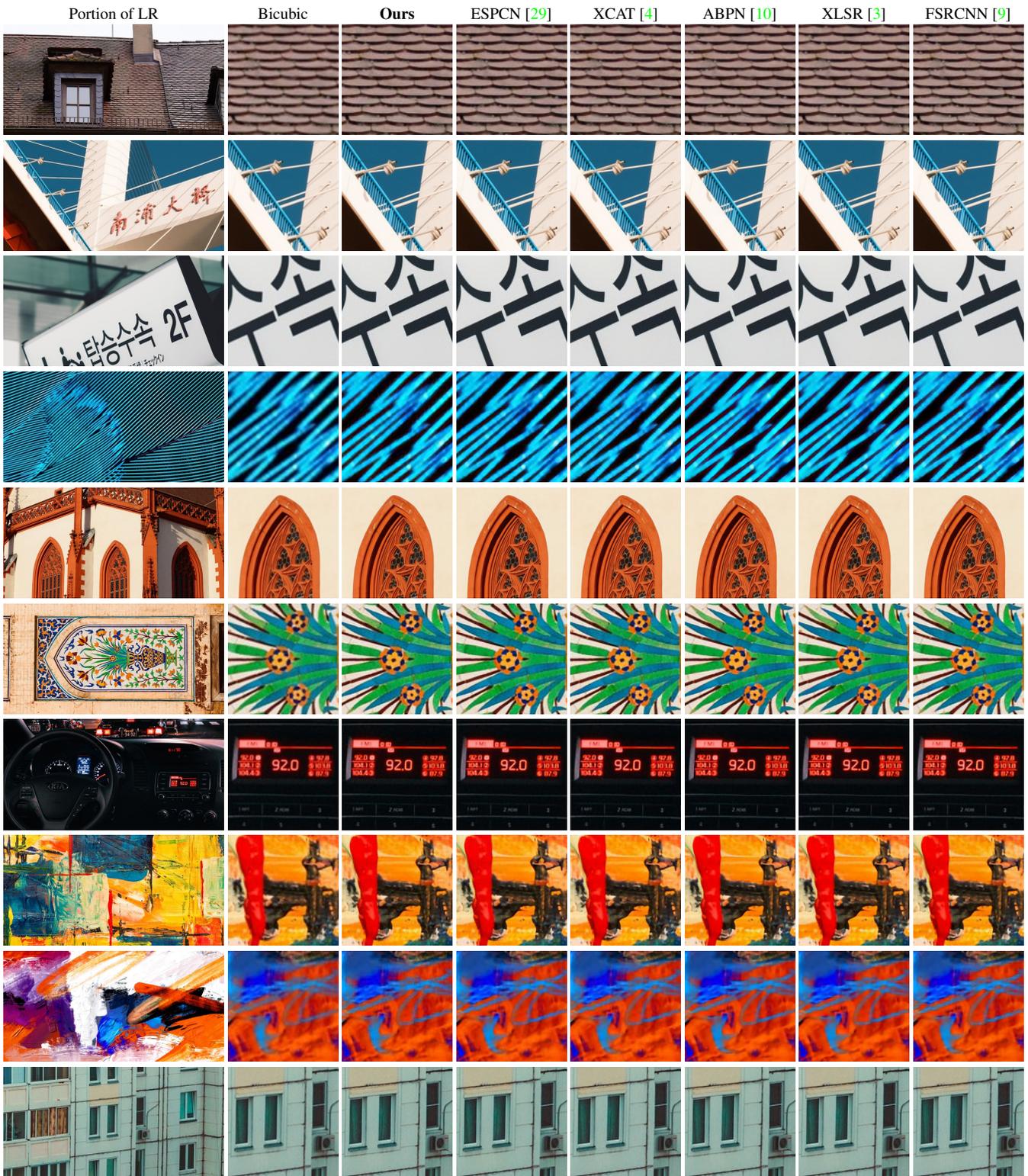


Figure 8. Qualitative comparative results of Bicubic++ and other relevant methods on the challenge test set [36].

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. 4, 7
- [2] Saeed Anwar and Nick Barnes. Densely residual laplacian super-resolution. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020. 2
- [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 (CVPR) Workshops*, pages 2472–2479, June 2021. 1, 2, 3, 5, 7, 8
- [4] Mustafa Ayazoglu and Bahri Batuhan Bilecen. Xcat - lightweight quantized single image super-resolution using heterogeneous group convolutions and cross concatenation. In *Computer Vision – ECCV 2022 Workshops*, pages 475–488, Cham, 2023. Springer Nature Switzerland. 1, 2, 5, 7, 8
- [5] Mustafa Ayazoğlu. Imdeception: Grouped information distilling super-resolution network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 756–765, June 2022. 2
- [6] Marco Bevilacqua, Aline Roumy, Christine Guillemot, and Marie line Alberi Morel. Low-complexity single-image super-resolution based on nonnegative neighbor embedding. In *Proceedings of the British Machine Vision Conference*, pages 135.1–135.10. BMVA Press, 2012. 7
- [7] Marcos V Conde, Eduard Zamfir, Radu Timofte, et al. Efficient deep models for real-time 4k image super-resolution. ntire 2023 benchmark and report. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 1, 2, 6
- [8] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38, 12 2014. 2
- [9] Chao Dong, Chen Change Loy, and Xiaoou Tang. Accelerating the super-resolution convolutional neural network. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, pages 391–407, Cham, 2016. Springer International Publishing. 1, 2, 5, 7, 8
- [10] Zongcai Du, Jie Liu, Jie Tang, and Gangshan Wu. Anchor-based plain net for mobile image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 2494–2502, June 2021. 1, 2, 3, 5, 7, 8
- [11] Ganzorig Gankhuyag, Jingang Huh, Myeongkyun Kim, Kihwan Yoon, Hyeoncheol Moon, Seungho Lee, Jinwoo Jeong, Sungjei Kim, and Yoonsik Choe. Skip-concatenated image super-resolution network for mobile devices. *IEEE Access*, 11:4972–4982, 2023. 3
- [12] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 3
- [13] Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. Single image super-resolution from transformed self-exemplars. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5197–5206, 2015. 7
- [14] Zheng Hui, Xinbo Gao, Yunchu Yang, and Xiumei Wang. Lightweight image super-resolution with information multi-distillation network. In *Proceedings of the 27th ACM International Conference on Multimedia (ACM MM)*, pages 2024–2032, 2019. 2
- [15] Andrey Ignatov and Radu Timofte. Efficient and accurate quantized image super-resolution on mobile npus, mobile ai & aim 2022 challenge: Report. In *Computer Vision – ECCV 2022 Workshops*, pages 92–129, Cham, 2023. Springer Nature Switzerland. 2, 3
- [16] 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 (CVPR) Workshops*, pages 2525–2534, June 2021. 2, 3
- [17] Younghyun Jo and Seon Joo Kim. Practical single-image super-resolution using look-up table. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021. 3, 6
- [18] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 2
- [19] Fangyuan Kong, Mingxi Li, Songwei Liu, Ding Liu, Jingwen He, Yang Bai, Fangmin Chen, and Lean Fu. Residual local feature network for efficient super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 766–776, June 2022. 3
- [20] Christian Ledig, Lucas Theis, Ferenc Huszar, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, and Wenzhe Shi. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 2
- [21] Yawei Li, Kai Zhang, Radu Timofte, and Van Gool. Ntire 2022 challenge on efficient super-resolution: Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 1062–1102, June 2022. 3
- [22] Jingyun Liang, Jiezhong Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. *arXiv preprint arXiv:2108.10257*, 2021. 2
- [23] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, July 2017. 2
- [24] Jie Liu, Jie Tang, and Gangshan Wu. Residual feature distillation network for lightweight image super-resolution, 2020. 3, 5

- [25] 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 *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, volume 2, pages 416–423 vol.2, 2001. 7
- [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, November 2016. 7
- [27] Brian B. Moser, Federico Raue, Stanislav Frolov, Sebastian Palacio, Jörn Hees, and Andreas Dengel. Hitchhiker’s guide to super-resolution: Introduction and recent advances. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–21, 2023. 2
- [28] Ben Niu, Weilei Wen, Wenqi Ren, Xiangde Zhang, Lianping Yang, Shuzhen Wang, Kaihao Zhang, Xiaochun Cao, and Haifeng Shen. Single image super-resolution via a holistic attention network. In *Computer Vision – ECCV 2020*, pages 191–207, Cham, 2020. Springer International Publishing. 2
- [29] Wenzhe Shi, Jose Caballero, Ferenc Huszar, 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 (CVPR)*, June 2016. 1, 2, 3, 5, 7, 8
- [30] Khushboo Singla, Rajoo Pandey, and Umesh Ghanekar. A review on single image super resolution techniques using generative adversarial network. *Optik*, 266:169607, 2022. 2
- [31] Wanjie Sun and Zhenzhong Chen. Learned image downscaling for upscaling using content adaptive resampler. *IEEE Transactions on Image Processing*, PP:1–1, 02 2020. 3
- [32] Radu Timofte, Vincent De, and Luc Van Gool. Anchored neighborhood regression for fast example-based super-resolution. In *2013 IEEE International Conference on Computer Vision*, pages 1920–1927, 2013. 2
- [33] Radu Timofte, Vincent De Smet, and Luc Van Gool. A+: Adjusted anchored neighborhood regression for fast super-resolution. In Daniel Cremers, Ian Reid, Hideo Saito, and Ming-Hsuan Yang, editors, *Computer Vision – ACCV 2014*, pages 111–126, Cham, 2015. Springer International Publishing. 2, 6
- [34] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In Laura Leal-Taixé and Stefan Roth, editors, *Computer Vision – ECCV 2018 Workshops*, pages 63–79, Cham, 2019. Springer International Publishing. 2
- [35] Jiahui Yu and Yuchen Fan. Wide activation for efficient image and video super-resolution. *Machine Vision and Applications*, 12 2020. 2
- [36] Eduard Zamfir, Marcos V Conde, and Radu Timofte. Towards real-time 4k image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 1, 2, 8
- [37] Roman Zeyde, Michael Elad, and Matan Protter. On single image scale-up using sparse-representations. In *Curves and Surfaces*, pages 711–730, Berlin, Heidelberg, 2012. Springer Berlin Heidelberg. 2
- [38] Roman Zeyde, Michael Elad, and Matan Protter. On single image scale-up using sparse-representations. In Jean-Daniel Boissonnat, Patrick Chenin, Albert Cohen, Christian Gout, Tom Lyche, Marie-Laurence Mazure, and Larry Schumaker, editors, *Curves and Surfaces*, pages 711–730, Berlin, Heidelberg, 2012. Springer Berlin Heidelberg. 7