

edge-SR: Super-Resolution For The Masses

– Appendix –

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

Abstract

We provide the following additional information:

- *Details of evaluation metrics :*
 - *image quality,*
 - *speed,*
 - *power and cpu usage.*
- *Extended analysis :*
 - *speed vs parameters,*
 - *additional edge device (GTX 1080 MaxQ),*
 - *qualitative evaluations,*
 - *effect of datasets and metrics,*
 - *magnified scatter plots.*
- *Reproducibility :*
 - *test results,*
 - *pytorch models and code.*

A. Evaluation Metrics

Image quality. Quantitative evaluations in our experiments include objective metrics PSNR and SSIM. These are reference-based metrics that measure the difference between an impaired image and ground truth. Higher values are better in both cases. The 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:

$$PSNR(X, Y) = 10 \cdot \log_{10} \left(\frac{255^2}{MSE} \right), \quad (1)$$

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

where $MSE = \mathbb{E}[(X - Y)^2]$ is the mean square error of the difference between 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 $c_2 = 58.5225$.

Benchmark in the literature can show big differences in PSNR and SSIM metrics due to different ways to evaluate color. For real-time applications it is common to process color images in YUV space, and the super-resolution task applies only to the luminance channel Y . Color in U and V

channels can use a faster bicubic upscaler with small impact in perceptual quality. We follow the implementation in [4] using a conversion of RGB to YUV color-spaces following the BT.709 standard, including offsets that are often avoided in other implementations.

In Appendix B and C we also provide measurements of the perceptual quality metric LPIPS defined in [3] and implemented using the PIQA library [2].

Speed. 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 (1920×1080) per second. Figure 1 shows the code used to parse Pytorch's profiler output to obtain the CPU and GPU processing time.

Power and CPU usage. 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. During this process we monitor the maximum power consumption using `nvidia-smi` for GTX 1080 Max-Q GPU, and `tegrastats` for Jetson AGX Xavier. We register the maximum power measured in this process. The power data provided by Max-Q driver is in units of watts, whereas the AGX device uses units of milliwatts. The AGX device allows different profiles for power consumption and in our experiments we used 30 Watts.

The Raspberry Pi 400 device does not include power sensors and in this case we replace the power measurement by CPU usage, monitor by parsing the `top` command with delay-time of $0.05s$ and registering the average CPU reading during the inference process.

```

1  import numpy as np
2  import torch.autograd.profiler as profiler
3
4  def str_to_time(s):
5      if s.endswith('ms'):
6          return float(s[:-2])*1e-3
7      if s.endswith('us'):
8          return float(s[:-2])*1e-6
9      return float(s[:-1])
10
11 def speed(model, input):
12     dt = np.inf
13     for _ in range(10):
14         with profiler.profile(
15             record_shapes=True, use_cuda=True
16         ) as prof:
17             with profiler.record_function(
18                 'model_inference'
19             ):
20                 output = model(input)
21             dt1 = str_to_time(
22                 prof.key_averages().table(
23                     sort_by='cpu_time_total', row_limit=10
24                 ).split(
25                     'CPU time total: '
26                 )[1].split('\n')[0]
27             )
28             dt2 = str_to_time(
29                 prof.key_averages().table(
30                     sort_by='cpu_time_total', row_limit=10
31                 ).split(
32                     'CUDA time total: '
33                 )[1][:-1]
34             )
35             dt = min(dt1+dt2, dt)
36
37     pix = np.asarray(output.shape).prod()
38
39     return pix / (dt * 1920*1080)

```

Figure 1. Python function used to measure the speed of models in Pytorch v1.8. It runs a model 10 times to avoid warm-up effects. Then, it parses the output of Pytorch profiler to get both CPU and GPU runtime. The speed is the number of output pixels divided by the minimum runtime in the 10 runs.

B. Extended Analysis

Speed vs parameters. In Figure 2 we show the relationship between the size of the model, given by the number of parameters, and the execution speed when running the models on GPU devices. We observe that smaller models run faster, and the speed increases exponentially. The non-linear relation between speed and number of parameters becomes critical under: 5,000 parameters for 2 \times upscaling factor, 10,000 parameters for 3 \times upscaling factor, and 15,000 parameters for 4 \times upscaling factor. This shows the importance of focusing on speed compared to number of parameters in our study. Research on lightweight SR architectures often focuses on number of parameters and typically use several hundred thousand parameters where the non-linearity is still not critical.

Additional device. In Figure 3 we show scatter plots in-

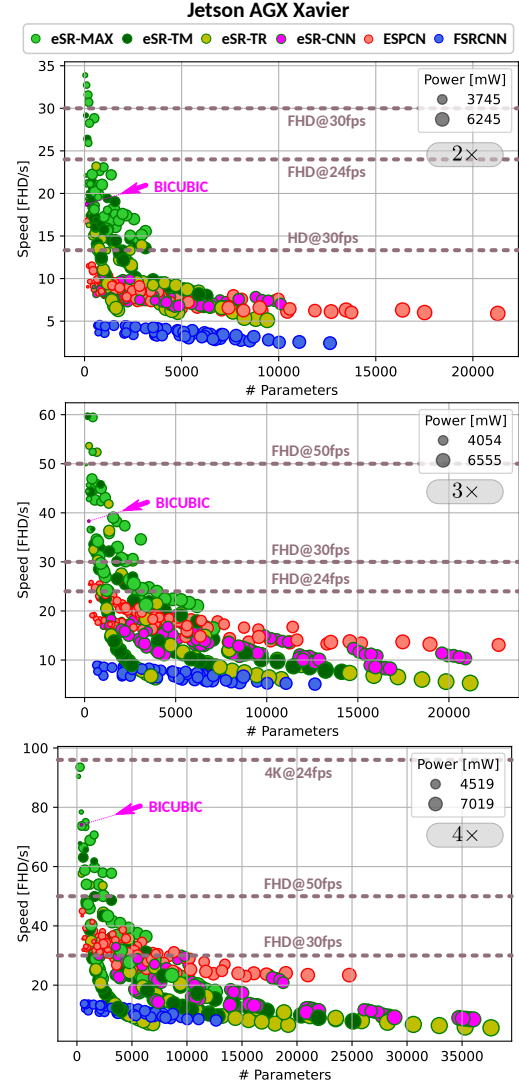


Figure 2. Example of the change in speed of the models with respect to their number of parameters for experiments on the Jetson AGX Xavier device. Models become faster at an exponential rate as the number of parameter reduces.

cluding Nvidia GeForce GTX-1080 Max-Q GPU. This is a mobile high-end GPU from the Pascal series typically used for laptop computers. The power consumption of the Max-Q design ranges between 90 and 110 Watt, compared to 30 Watt used in the Jetson AGX Xavier. Although much more powerful than a Raspberry Pi and Jetson AGX devices, the GTX 1080 Max-Q device can fit in high-end display TV panels and thus serve a different range of applications for edge devices. Consequently, Figure 3 shows a performance that can deliver 8K videos in real-time (even with 16-bit floating point precision). We also observe that ESPCN performance improves for 3 \times and 4 \times upscaling factors, indicating the important effect of the increased number of cores in GPU architectures as well as a significant increase in

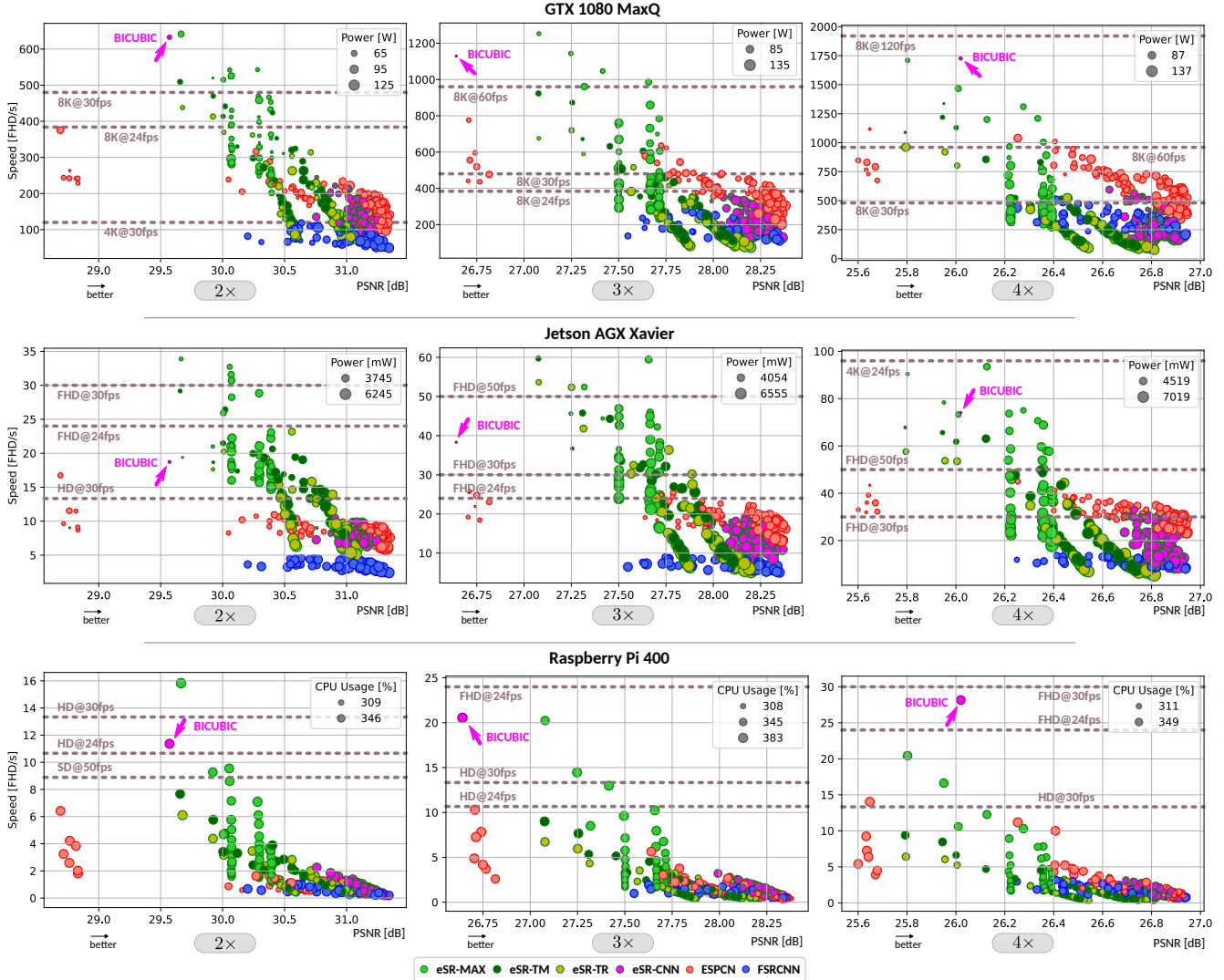


Figure 3. Scatter plot to compare speed, in number of Full-HD pixels per second, with respect to quality, measured as PSNR for the BSDS-100 dataset. A total of 1,185 models were identically trained considering different upscaling factors (2 \times , 3 \times and 4 \times) and architectures (eSR, ESPCN and FSRCNN). We run all models on edge devices: GTX-1080 Max-Q (GPU with 16-bit floating point precision), 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 Figures 20, 21, 22, 14, 15 and 16.

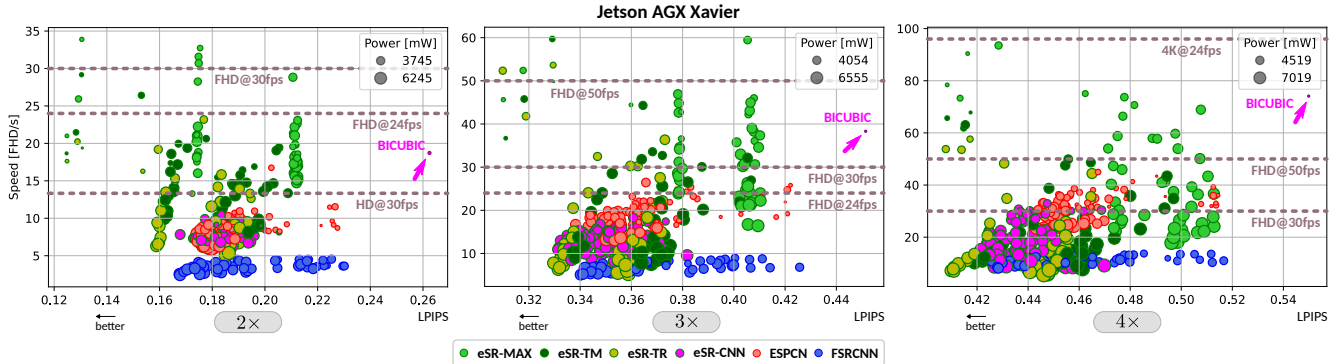


Figure 4. Scatter plot to compare speed, in number of Full-HD pixels per second, with respect to quality, measured as LPIPS for the BSDS-100 dataset. Lower values of LPIPS mean better quality as opposed to PSNR where larger values are better. Compared with the same case but using PSNR quality measure in Figure 3, eSR-TM performs much better and Bicubic shows the worst quality.

```

1 >>> import pickle
2 >>> test = pickle.load(open('tests.pkl', 'rb'))
3 >>> test['Bicubic_s2']
4 {'psnr_Set5': 33.72849620514912,
5  'ssim_Set5': 0.9283912810369976,
6  'lpips_Set5': 0.14221979230642318,
7  'psnr_Set14': 30.286027790636204,
8  'ssim_Set14': 0.8694934108301432,
9  'lpips_Set14': 0.19383049915943826,
10 'psnr_BSDS100': 29.571233006609656,
11 'ssim_BSDS100': 0.8418117904964167,
12 'lpips_BSDS100': 0.26246454380452633,
13 'psnr_Urban100': 26.89378248655882,
14 'ssim_Urban100': 0.8407461069831571,
15 'lpips_Urban100': 0.21186692919582129,
16 'psnr_Mangal09': 30.850672809780587,
17 'ssim_Mangal09': 0.9340133711400112,
18 'lpips_Mangal09': 0.102985977955641,
19 'parameters': 104,
20 'speed_AGX': 18.72132628065749,
21 'power_AGX': 1550,
22 'speed_MaxQ': 632.5429857814075,
23 'power_MaxQ': 50,
24 'temperature_MaxQ': 76,
25 'memory_MaxQ': 2961,
26 'speed_RPI': 11.361346064182795,
27 'usage_RPI': 372.8714285714285}

```

Figure 5. Example of how to read our test data from the Python dictionary file `tests.pkl`.

power consumption.

Qualitative evaluation. In Figures 11, 12 and 13 we show example output images for different models. These models were selected by moving in an approximately optimal trajectory in the Speed vs Quality (SQ) plane for the Jetson AGX Xavier device. As mentioned before, the best advantage of edge-SR models is observed for $2\times$ upscaling factor. This is both the most difficult and important factor for applications due to the high input throughput in high resolution displays (e.g. HD to FHD). The classic bicubic upscalers offers over-smooth outputs that are somehow effective to reduce *jaggy* artifacts. edge-MAX models can significantly improve sharpness with limited control over jaggies. Next, edge-TM and edge-TR models show the best trade-off with very similar performance. These models are more effective at reducing jaggies before multi-layer networks like edge-CNN and ESPCN become better with both sharp and smooth edges.

In Figure 4 we show a few scatter plots of the trade-off between image quality and runtime performance using the perceptual quality metric LPIPS [3] for the BSDS-100 dataset. Here, lower values of LPIPS mean better quality as opposed to PSNR where larger values are better. Compared with the same case but using PSNR quality measure in Figure 3, eSR-TM performs much better and Bicubic shows the worst quality.

At $3\times$ and $4\times$ upscaling factor the pattern is similar but it becomes more difficult for single-layer models to effec-

```

1 import torch
2 from torch import nn
3
4 class edgeSR_MAX(nn.Module):
5     def __init__(self, C, k, s):
6         super().__init__()
7
8         self.pixel_shuffle = nn.PixelShuffle(s)
9         self.filter = nn.Conv2d(
10             in_channels=1,
11             out_channels=s*s*C,
12             kernel_size=k,
13             stride=1,
14             padding=(k-1)//2,
15             bias=False,
16         )
17
18     def forward(self, input):
19         return self.pixel_shuffle(
20             self.filter(input)
21         ).max(dim=1, keepdim=True)[0]

```

Figure 6. Pytorch v1.8 implementation of edge-SR Maximum (eSR-MAX) single-layer architecture.

```

1 import torch
2 from torch import nn
3
4 class edgeSR_TM(nn.Module):
5     def __init__(self, C, k, s):
6         super().__init__()
7
8         self.pixel_shuffle = nn.PixelShuffle(s)
9         self.softmax = nn.Softmax(dim=1)
10        self.filter = nn.Conv2d(
11            in_channels=1,
12            out_channels=2*s*s*C,
13            kernel_size=k,
14            stride=1,
15            padding=(k-1)//2,
16            bias=False,
17        )
18
19    def forward(self, input):
20        filtered = self.pixel_shuffle(
21            self.filter(input)
22        )
23        B, C, H, W = filtered.shape
24
25        filtered = filtered.view(B, 2, C, H, W)
26        upscaling = filtered[:, 0]
27        matching = filtered[:, 1]
28        return torch.sum(
29            upscaling * self.softmax(matching),
30            dim=1, keepdim=True
31        )

```

Figure 7. Pytorch v1.8 implementation of edge-SR Template Matching (eSR-TM) single-layer architecture.

tively reduce *jaggy* artifacts. Results get worsen and we observe an increased gap between bicubic and other architectures. ESPCN becomes better at high quality ranges and overcomes edge-SR models for the most part. Finally we note that trade-off evaluations at $3\times$ and $4\times$ upscaling factors could be misleading, as flickering video artifacts are


```

1 import torch
2 from torch import nn
3
4 class edgeSR_TR(nn.Module):
5     def __init__(self, C, k, s):
6         super().__init__()
7
8         self.pixel_shuffle = nn.PixelShuffle(s)
9         self.softmax = nn.Softmax(dim=1)
10        self.filter = nn.Conv2d(
11            in_channels=1,
12            out_channels=3*s*s*C,
13            kernel_size=k,
14            stride=1,
15            padding=(k-1)//2,
16            bias=False,
17        )
18
19    def forward(self, input):
20        filtered = self.pixel_shuffle(
21            self.filter(input)
22        )
23        B, C, H, W = filtered.shape
24
25        filtered = filtered.view(B, 3, C, H, W)
26        value = filtered[:, 0]
27        query = filtered[:, 1]
28        key = filtered[:, 2]
29        return torch.sum(
30            value * self.softmax(query*key),
31            dim=1, keepdim=True
32        )

```

Figure 8. Pytorch v1.8 implementation of edge-SR TRansformer (eSR-TR) single-layer architecture.

likely to become visible at this point and video SR solutions might be needed. For this reason, the results at $2\times$ are arguably the most important for practical applications.

Effect of datasets and metrics. Scatter plots in the main text use only PSNR metric measured in the BSDS-100 dataset. In Figure 10 we show the effect of changing both the metric, from PSNR to SSIM, and dataset, among Set5, Set14, BSDS-100, Urban-100 and Manga-109. The range of values and relative position of scatter points changes. For example, SSIM shows a larger margin in image quality between bicubic and other models. PSNR values show significant changes depending on the dataset. Nevertheless, the trade-off trajectory and the advantage shown by different architectures remains the same.

Magnified scatter plots. It is useful to identify the hyper-parameters of each specific model in scatter plots. Figures 20, 21, 22, 14, 15 and 16 show magnified plots using all the page width and include annotations with model hyper-parameters. The annotations are useful at middle and high speed ranges where data is more sparse. In the high image quality range the performance of different models become clustered and the annotations are not readable. Here, we recommend to look at Figures 11, 12 and 13 to identify the best models. Finally, we make all test data available in a Python dictionary file (see Appendix C).

```

1 import torch
2 from torch import nn
3
4 class edgeSR_CNN(nn.Module):
5     def __init__(self, C, D, S, s):
6         super().__init__()
7
8         self.softmax = nn.Softmax(dim=1)
9         if D == 0:
10            self.filter = nn.Sequential(
11                nn.Conv2d(D, S, 3, 1, 1),
12                nn.Tanh(),
13                nn.Conv2d(
14                    in_channels=S,
15                    out_channels=2*s*s*C,
16                    kernel_size=3,
17                    stride=1,
18                    padding=1,
19                    bias=False,
20                ),
21                nn.PixelShuffle(s),
22            )
23        else:
24            self.filter = nn.Sequential(
25                nn.Conv2d(1, D, 5, 1, 2),
26                nn.Tanh(),
27                nn.Conv2d(D, S, 3, 1, 1),
28                nn.Tanh(),
29                nn.Conv2d(
30                    in_channels=S,
31                    out_channels=2*s*s*C,
32                    kernel_size=3,
33                    stride=1,
34                    padding=1,
35                    bias=False,
36                ),
37                nn.PixelShuffle(s),
38            )
39
40    def forward(self, input):
41        filtered = self.filter(input)
42        B, C, H, W = filtered.shape
43
44        filtered = filtered.view(B, 2, C, H, W)
45        upscaling = filtered[:, 0]
46        matching = filtered[:, 1]
47        return torch.sum(
48            upscaling * self.softmax(matching),
49            dim=1, keepdim=True
50        )

```

Figure 9. Pytorch v1.8 implementation of edge-SR CNN (eSR-CNN) multi-layer architecture.

C. Reproducibility

Test results. Test results are provided in the Python dictionary file `tests.pkl` using the native *pickle* module. A sample code to read the file is provided in Figure 5. The keys of the dictionary identify the name of each model and its hyper-parameters using the following format:

- 'Bicubic.s#',
- 'eSR-MAX.s#_K#_C#',
- 'eSR-TM.s#_K#_C#',
- 'eSR-TR.s#_K#_C#',
- 'eSR-CNN.s#_C#_D#_S#',

- 'ESPCN_s#_D#_S#', or
- 'FSRCNN_s#_D#_S#_M#',

where # represents an integer number with the value of the correspondent hyper-parameter. For each model the data of the dictionary contains a second dictionary with the information displayed in Figure 5. This includes: number of model parameters; image quality metrics PSNR, SSIM and LPIPS measured in 5 different datasets; as well as power, speed, CPU usage, temperature and memory usage for devices AGX (Jetson AGX Xavier), MaxQ (GTX 1080 MaxQ) and RPI (Raspberry Pi 400).

Pytorch implementations. Figures 6, 7, 8 and 9 show the Python code to implement all proposed edge-SR models using the Pytorch tensor processing framework version 1.8. Model files and sample code are also available in <https://github.com/pnavarre/eSR>.

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. 1
- [2] François Rozet. Pytorch image quality assessment (PIQA). <https://github.com/francois-rozet/piqa>, 2021. 1
- [3] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 586–595, 2018. 1, 4
- [4] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Evaluation code for residual dense networks. https://github.com/yulunzhang/RDN/blob/master/RDN_TestCode/Evaluate_PSNR_SSIM.m, 2018. [Online; accessed 20-May-2020]. 1

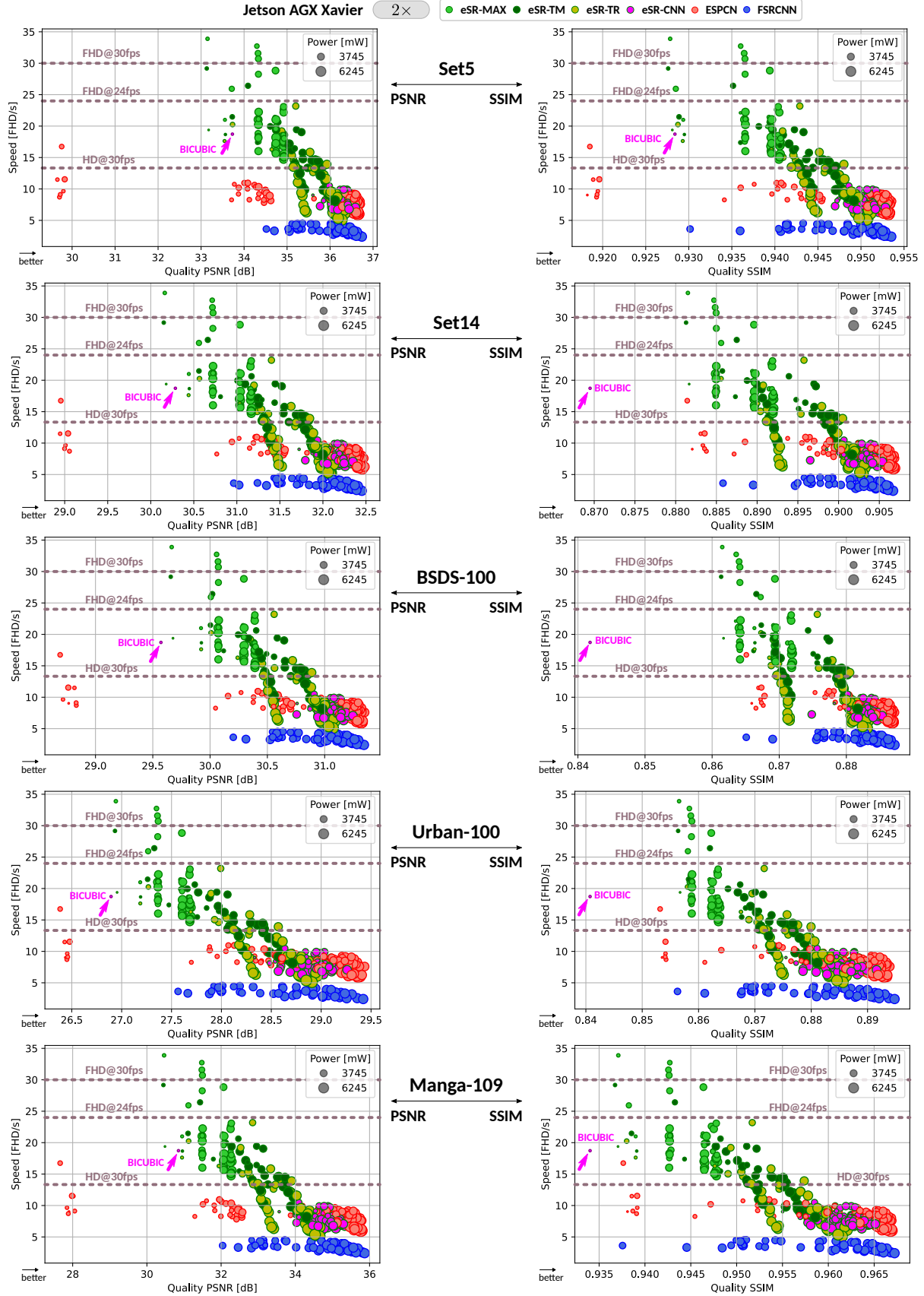


Figure 10. Example of the effect of different metric (PSNR/ SSIM) and datasets in the trade-off between image quality and runtime performance. We observe that even though the range of values and relative positions change, edge-SR models remain with better performance at high speed, and multi-layer networks remain better at low speed.

Jetson AGX Xavier

2×

- eSR-MAX
- eSR-TM
- eSR-TR
- eSR-CNN
- ESPCN
- FSRCNN

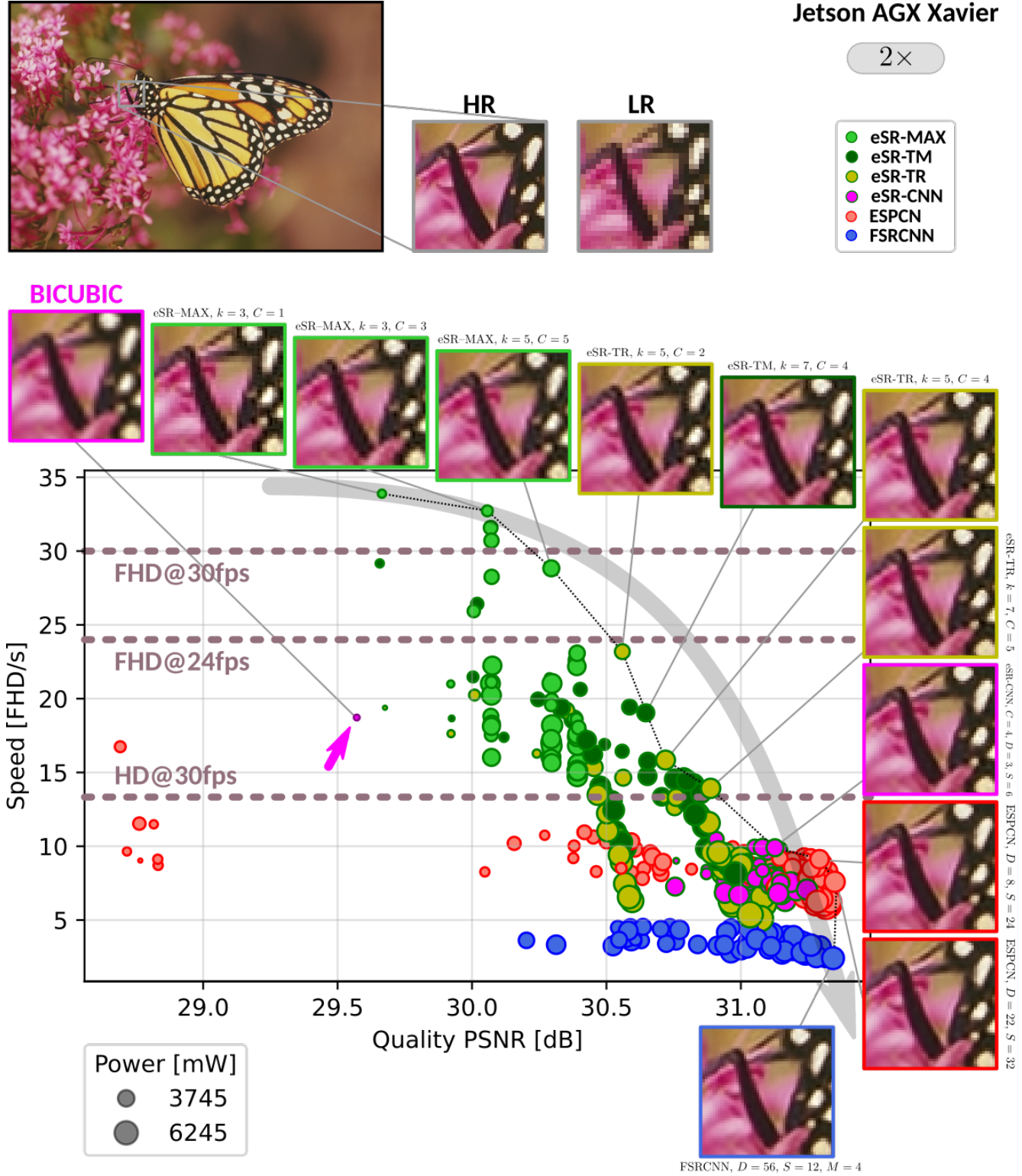


Figure 11. Example output images in the trajectory through the trade-off between image quality and runtime performance for 2× upscaling on a Jetson AGX Xavier edge device.

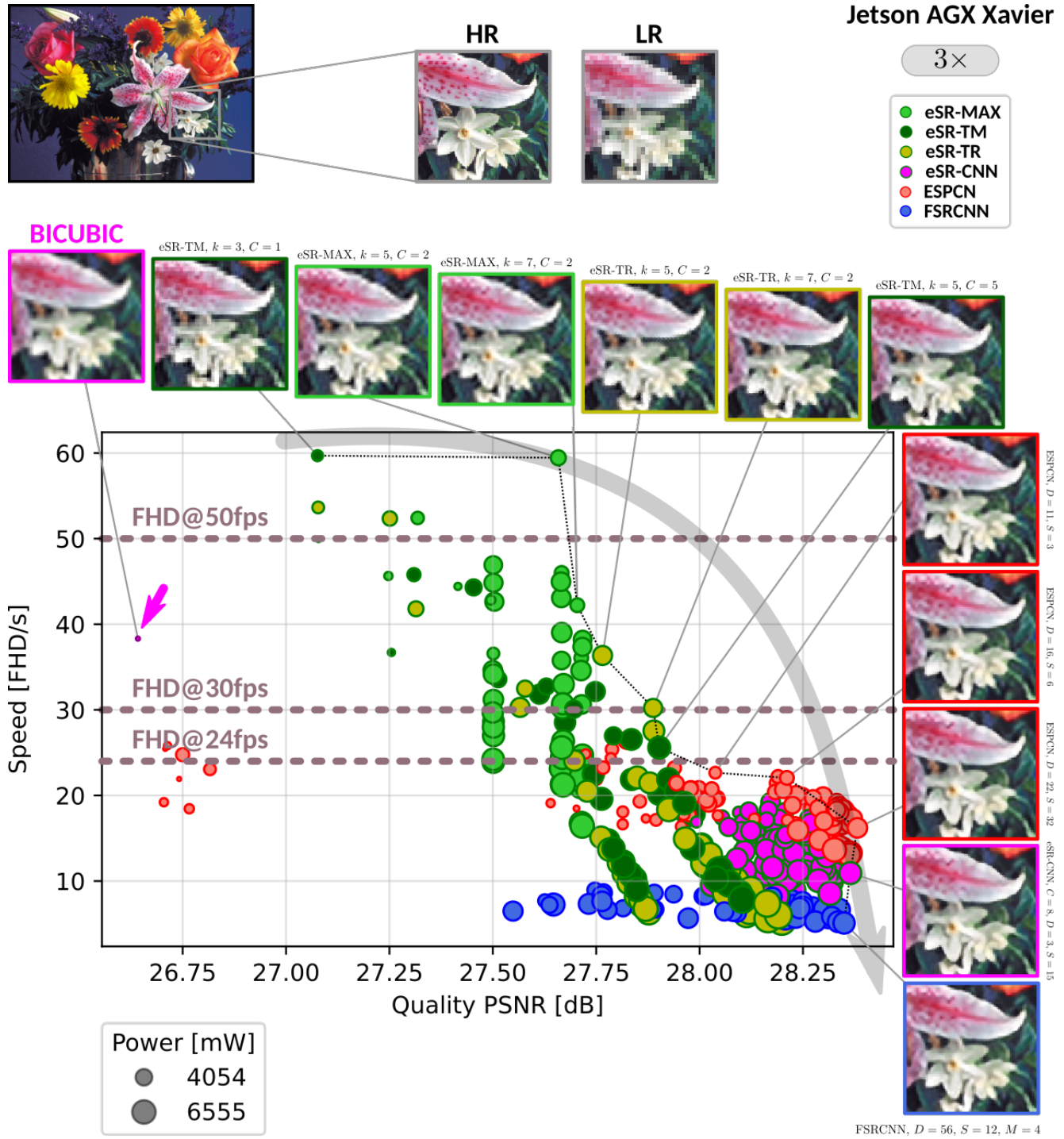


Figure 12. Example output images in the trajectory through the trade-off between image quality and runtime performance for 3× upscaling on a Jetson AGX Xavier edge device.

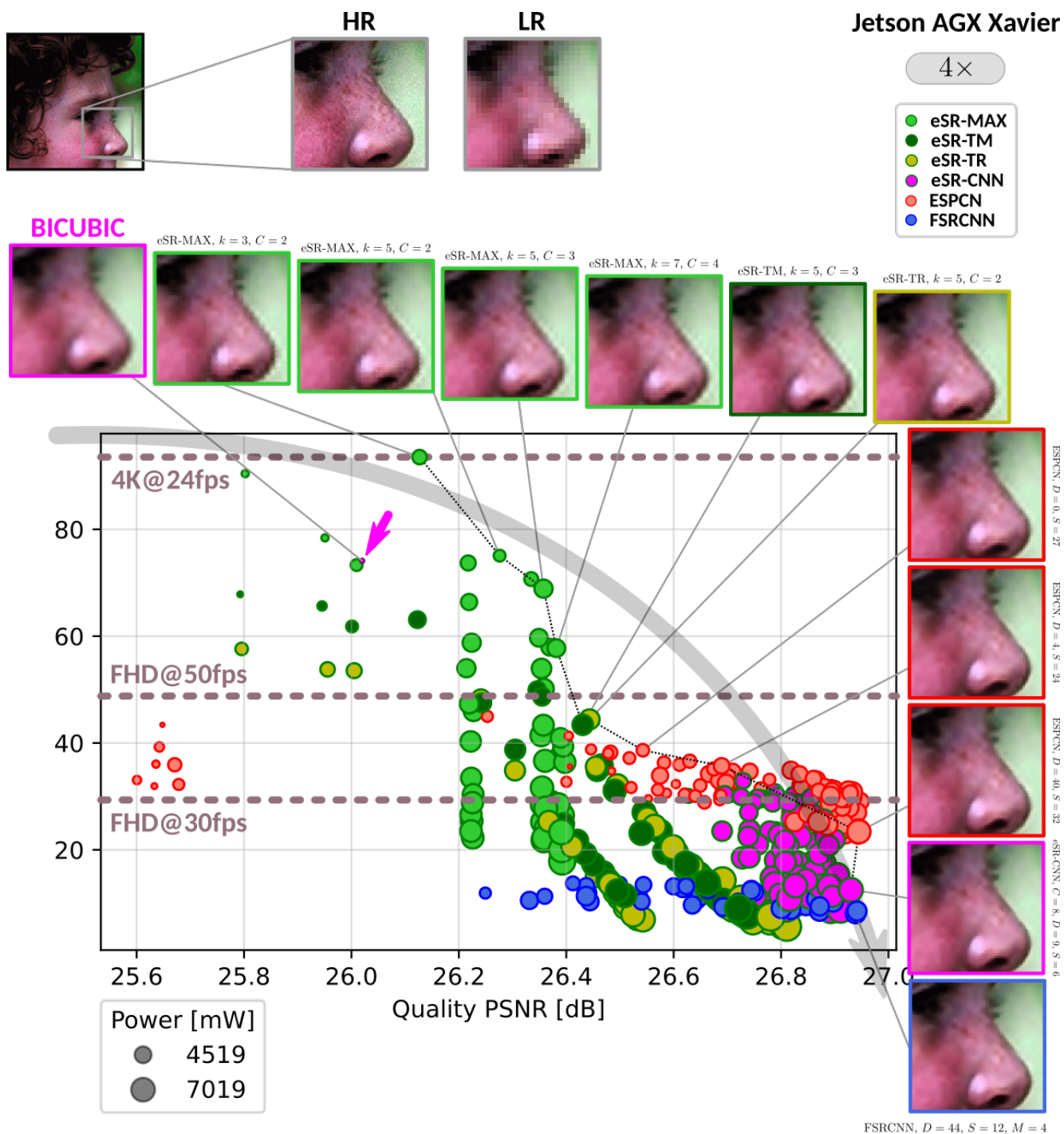


Figure 13. Example output images in the trajectory through the trade-off between image quality and runtime performance for 4× upscaling on a Jetson AGX Xavier edge device.

Jetson AGX Xavier

2×

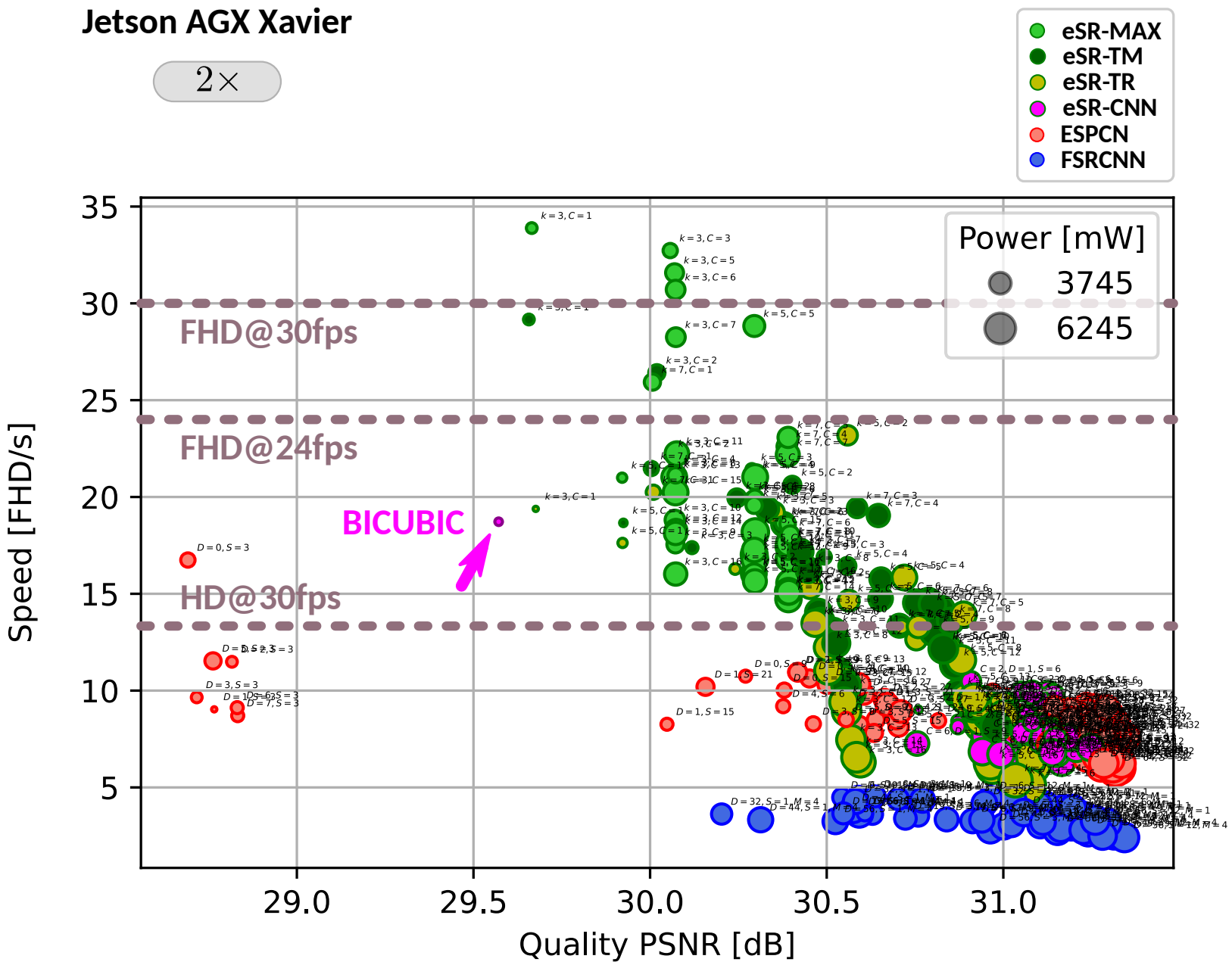
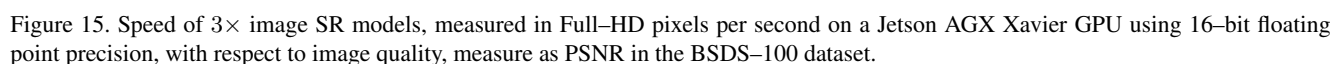


Figure 14. Speed of 2× image SR models, measured in Full-HD pixels per second on a Jetson AGX Xavier GPU using 16-bit floating point precision, with respect to image quality, measure as PSNR in the BSDS-100 dataset.

3x



Jetson AGX Xavier

4×

- eSR-MAX
- eSR-TM
- eSR-TR
- eSR-CNN
- ESPCNN
- FSRCNN

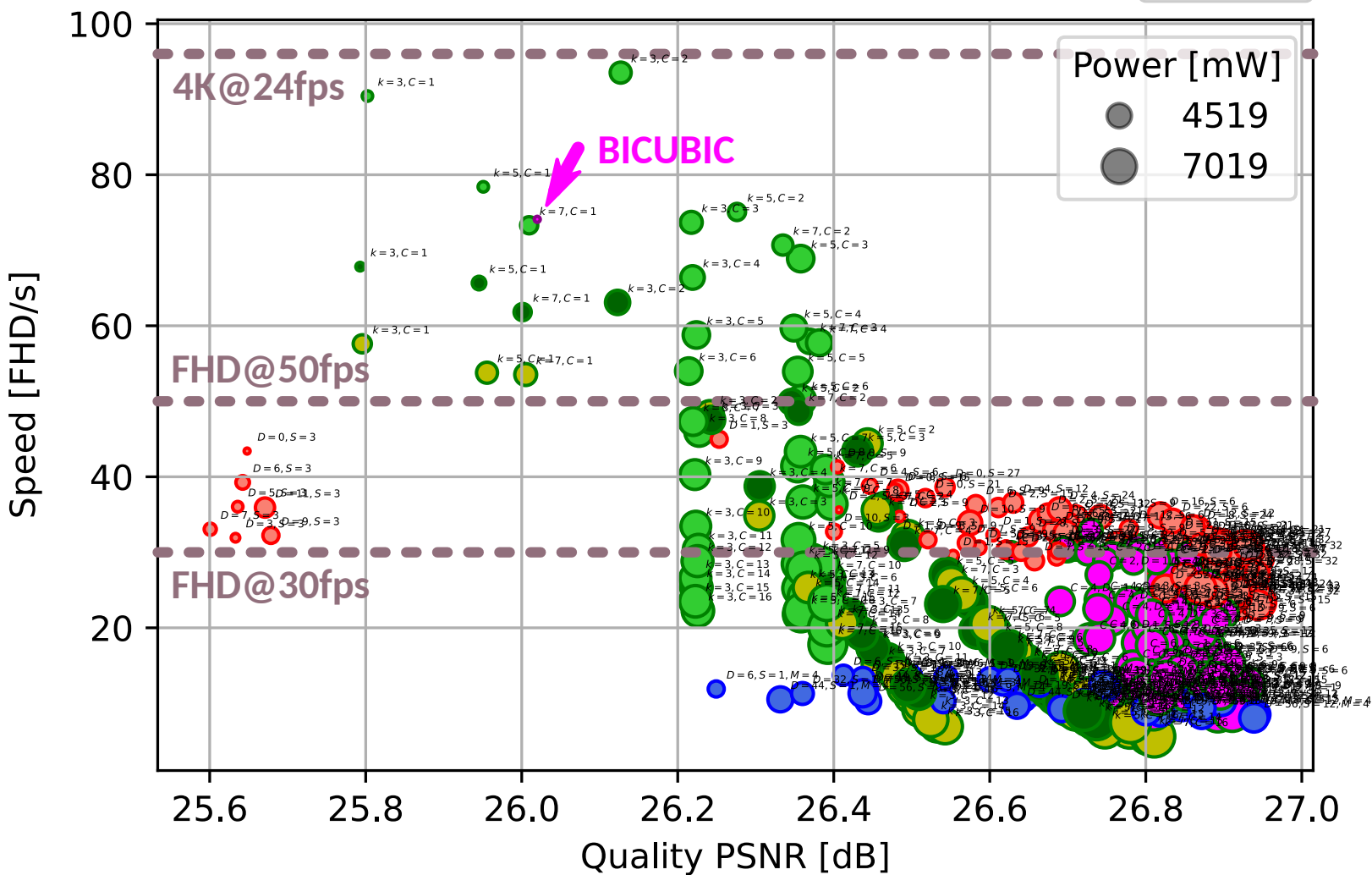
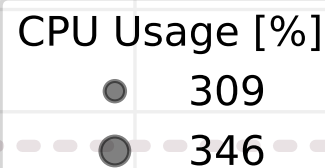


Figure 16. Speed of 4× image SR models, measured in Full-HD pixels per second on a Jetson AGX Xavier GPU using 16-bit floating point precision, with respect to image quality, measure as PSNR in the BSDS-100 dataset.

2x



14

Raspberry Pi 400

3×

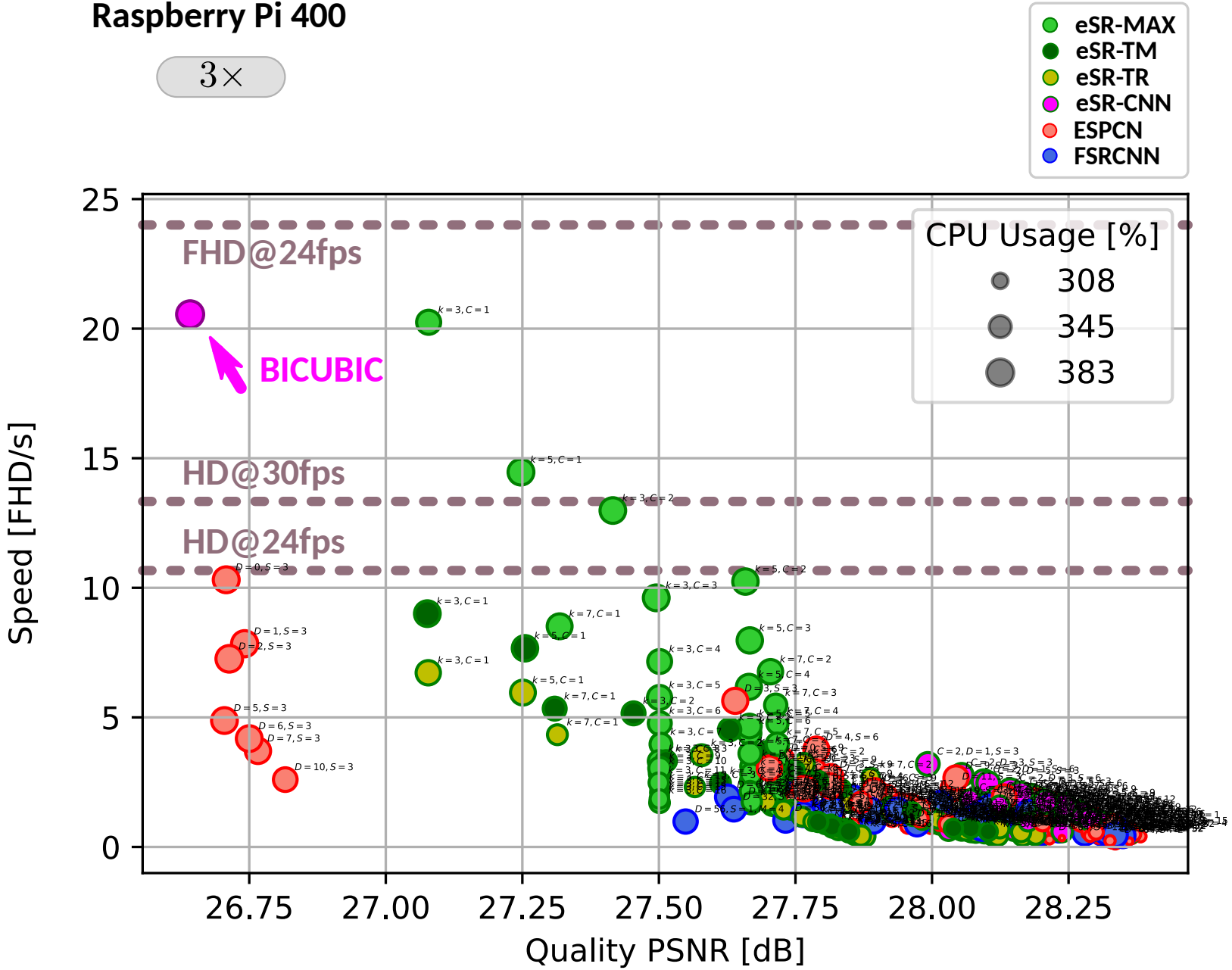


Figure 18. Speed of 3× image SR models, measured in Full-HD pixels per second on a Raspberry Pi 400 CPU using 32-bit floating point precision, with respect to image quality, measure as PSNR in the BSDS-100 dataset.

Raspberry Pi 400

4×

- eSR-MAX
- eSR-TM
- eSR-TR
- eSR-CNN
- ESPCN
- FSRCNN

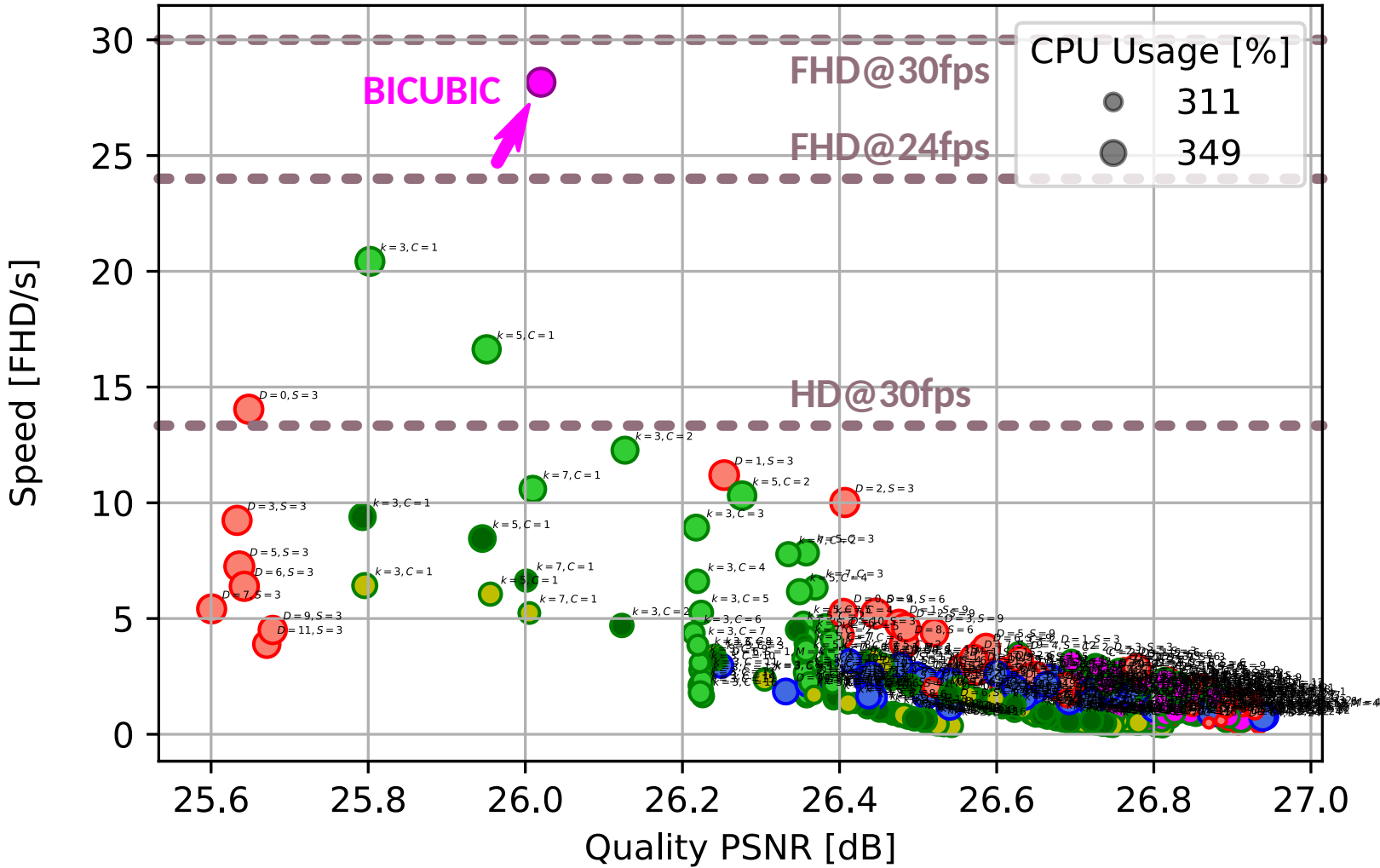
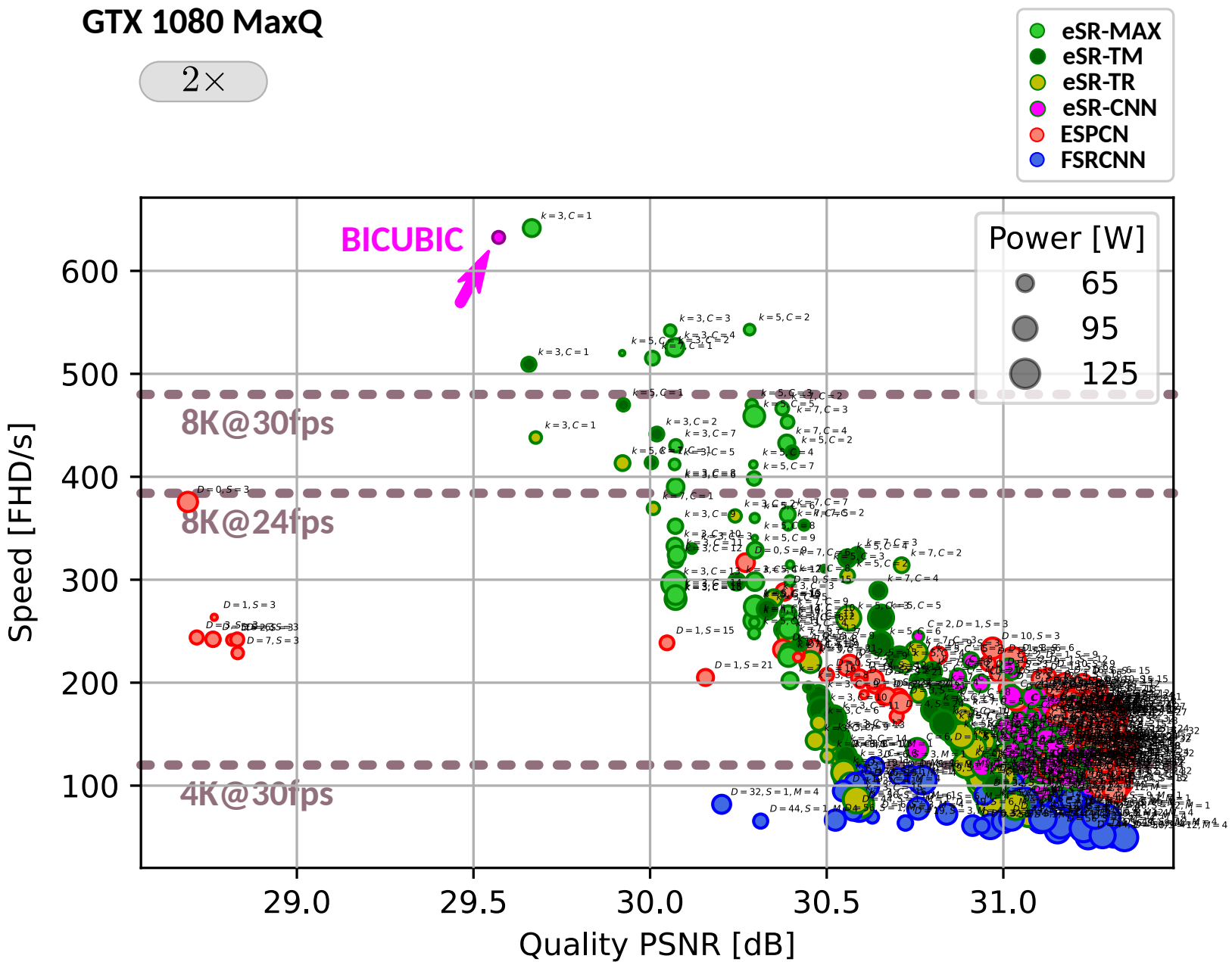


Figure 19. Speed of 4× image SR models, measured in Full-HD pixels per second on a Raspberry Pi 400 CPU using 32-bit floating point precision, with respect to image quality, measure as PSNR in the BSDS-100 dataset.

GTX 1080 MaxQ

2×



GTX 1080 MaxQ

3×

- eSR-MAX
- eSR-TM
- eSR-TR
- eSR-CNN
- ESPCN
- FSRCNN

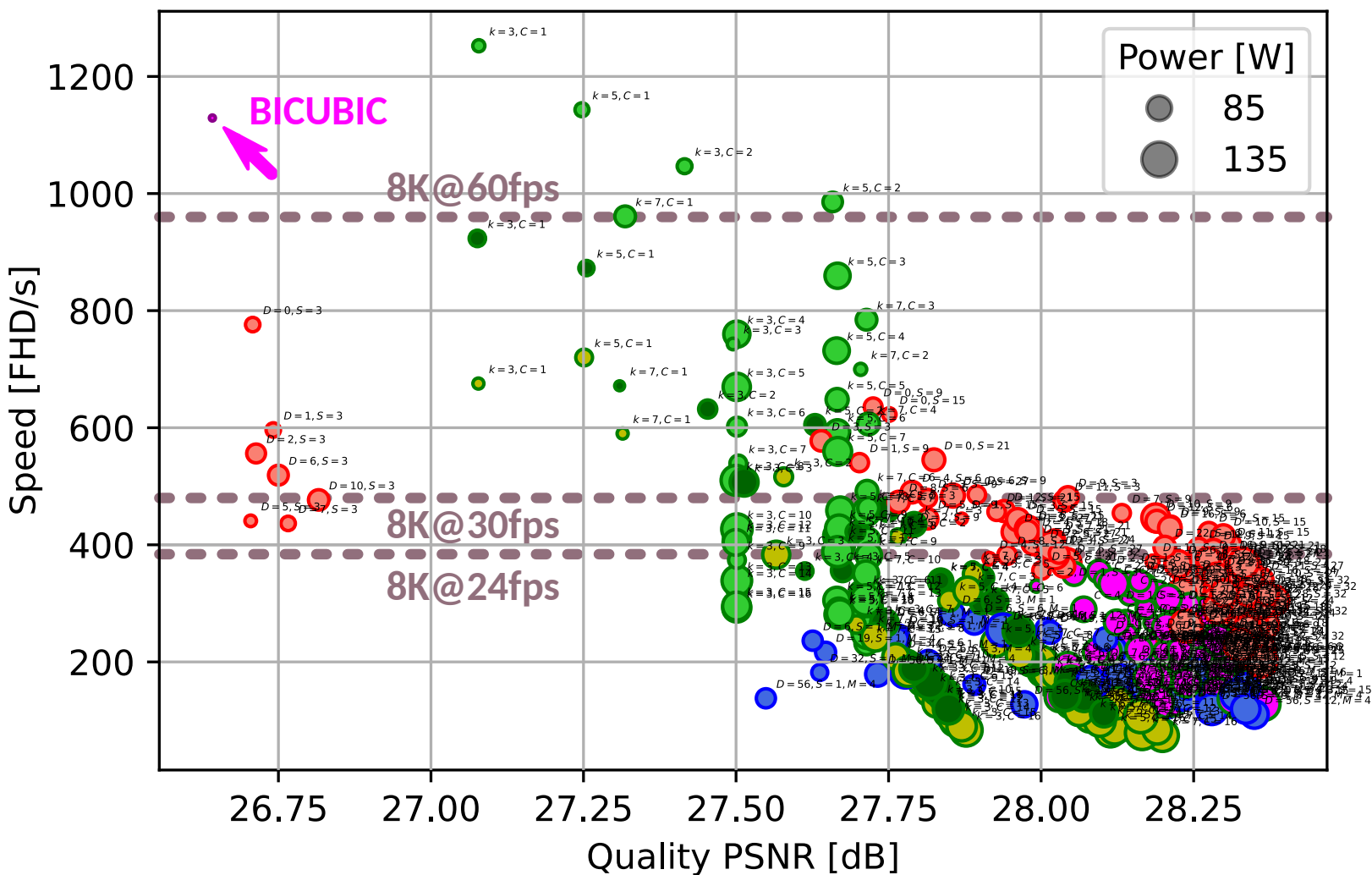


Figure 21. Speed of 3× image SR models, measured in Full-HD pixels per second on a GTX 1080 Max-Q GPU using 16-bit floating point precision, with respect to image quality, measure as PSNR in the BSDS-100 dataset.

GTX 1080 MaxQ

4×

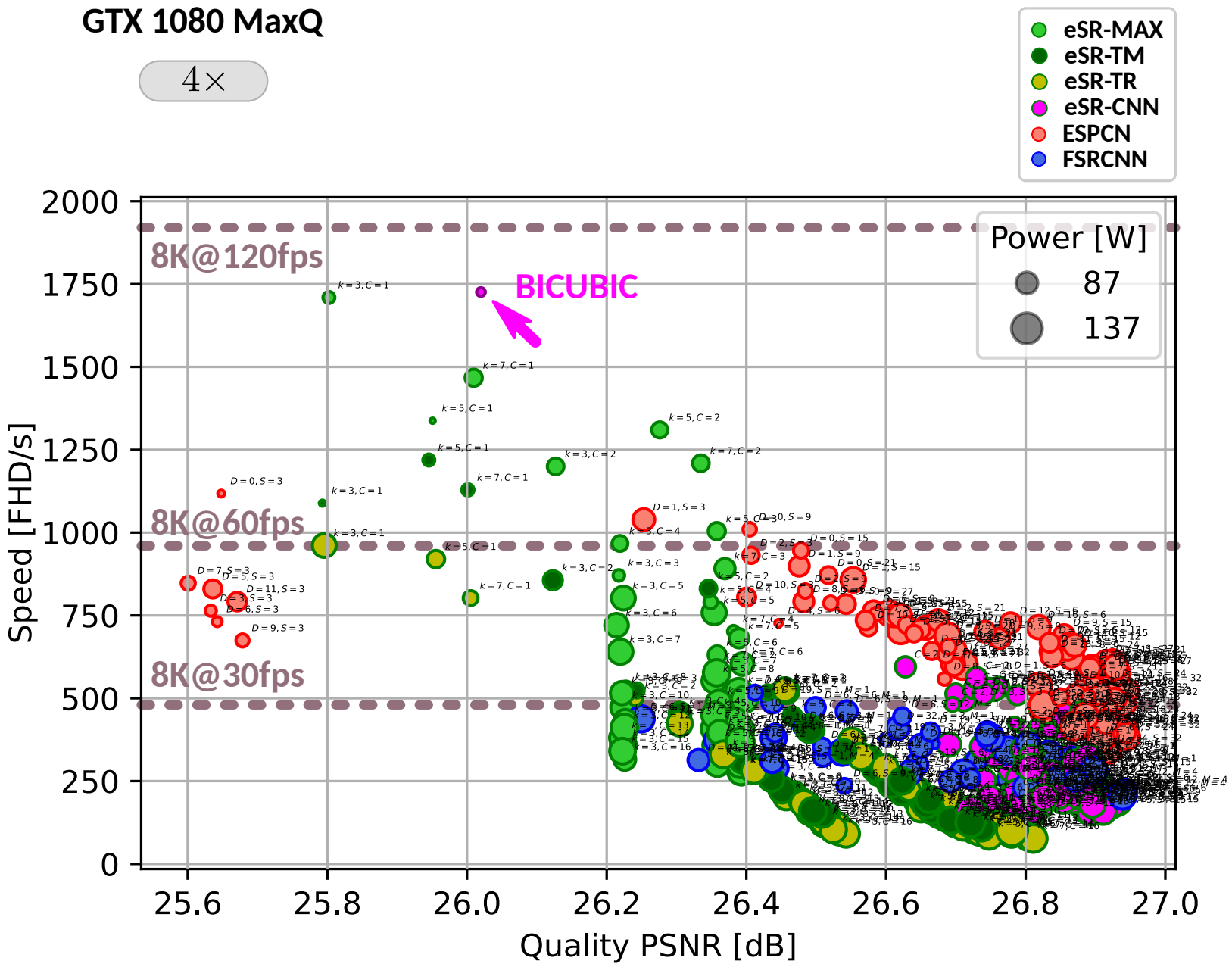


Figure 22. Speed of 4× image SR models, measured in Full-HD pixels per second on a GTX 1080 Max-Q GPU using 16-bit floating point precision, with respect to image quality, measure as PSNR in the BSDS-100 dataset.