

Efficient Meta-Tuning for Content-aware Neural Video Delivery

Xiaoqi Li^{1*}, Jiaming Liu^{1,2*}, Shizun Wang^{3*}, Cheng Lyu³, Ming Lu^{4†},
Yurong Chen⁴, Anbang Yao⁴, Yandong Guo², and Shanghang Zhang^{1†}

¹ Peking University

² OPPO Research Institute

³ Beijing University of Posts and Telecommunications

⁴ Intel Labs China

clorisleef0313@gmail.com, shzhang.pku@gmail.com, yandong.guo@live.com

Abstract. Recently, Deep Neural Networks (DNNs) are utilized to reduce the bandwidth and improve the quality of Internet video delivery. Existing methods train corresponding content-aware super-resolution (SR) model for each video chunk on the server, and stream low-resolution (LR) video chunks along with SR models to the client. Although they achieve promising results, the huge computational cost of network training limits their practical applications. In this paper, we present a method named Efficient Meta-Tuning (EMT) to reduce the computational cost. Instead of training from scratch, EMT adapts a meta-learned model to the first chunk of the input video. As for the following chunks, it fine-tunes the partial parameters selected by gradient masking of previous adapted model. In order to achieve further speedup for EMT, we propose a novel sampling strategy to extract the most challenging patches from video frames. The proposed strategy is highly efficient and brings negligible additional cost. Our method significantly reduces the computational cost and achieves even better performance, paving the way for applying neural video delivery techniques to practical applications. We conduct extensive experiments based on various efficient SR architectures, including ESPCN, SRCNN, FSRCNN and EDSR-1, demonstrating the generalization ability of our work. The code is released at <https://github.com/Neural-video-delivery/EMT-Pytorch-ECCV2022>.

Keywords: Neural Video Delivery, Super-Resolution, Meta Learning

1 Introduction

With the popularity of High-Definition (HD) display devices, high-resolution videos are strongly demanded by end users. This brings a huge burden to the video delivery infrastructure. As the development of deep learning, several recent

*Equal contribution

†Corresponding Author

works are proposed to reduce the bandwidth of video delivery [14,29,19,13]. The motivation of these works is to stream both the low-resolution videos and content-aware SR models from servers to clients. The clients run the inference of SR models to super-resolve the LR videos. In this manner, high-resolution videos can be delivered under limited Internet bandwidth.

In contrast to existing approaches on Single Image Super-Resolution (SISR) [23,8,18,31,15] and Video Super-Resolution (VSR) [2,25,3,4], content-aware models utilize the overfitting property of DNNs to achieve higher SR performance. To be more specific, a video is divided into several video chunks, and a corresponding SR model is trained for each chunk. This type of DNN-based video delivery system can achieve better performance even compared with commercial techniques like WebRTC [14].

Although neural video delivery is a promising technique, the huge computational cost of training content-aware SR models limits its practical applications. For example, existing methods [19,29] uniformly divide a 45s/1080P/30FPS video into 5-second chunks, and train the SR models for all chunks. However, even with efficient SR architectures like ESPCN [23], it still takes about 10.2 hours to train the content-aware SR models on a high-end NVIDIA V100 GPU. Therefore, reducing the computational cost of network training is crucial for neural video delivery.

In order to pave the way for practical applications, we propose Efficient Meta-Tuning (EMT) in this paper. Instead of training from scratch [19,29], EMT sequentially adapts a meta-learned model to the video chunks, delivering all the content-aware SR models. Compared with random initialization or pre-trained initialization, a meta-learned model can transfer better to different video chunks. We collect a large-scale dataset of diverse video chunks and take each chunk as one specific task. MAML [9] is adopted to train the meta-learned model, whose parameters are shared by all content-aware SR models. For the chunks of the input video, EMT adapts the meta-learned model to the first chunk. As for the following chunks, it can fine-tune the partial parameters of the previous adapted model due to the temporal consistency between neighboring chunks. The partial parameters are selected by gradient masking, which masks a fraction of most significant parameters after a few gradient updates. Since EMT sequentially adapts the meta-learned model, each chunk simply needs to store the selected partial parameters. The current content-aware SR model can be constructed by updating the partial parameters of the previous model. This is important to compress all the models into one shared model and a few private parameters. Compared with CaFM [19], our method is more compact since the meta-learned model is shared by all chunks, while CaFM can only share one model for chunks within the input video.

To further reduce the computational cost, we propose a novel sampling strategy for EMT, which selects the most challenging patches from video frames. Our motivation is that previous adapted SR model already possesses the ability to super-resolve current chunk due to temporal consistency. Therefore, the training efforts of EMT should focus on challenging patches, which cannot be well

handled by the previous model. However, performing the evaluation of previous model on all patches of current chunk is time-consuming and brings additional cost. Inspired by video codec, we extract the I-frames from the input video and only perform the evaluation on I-frames. The positions of challenging patches are extracted based on I-frames and propagated to other frames. Since I-frame is very sparse within a video, the computational cost of the evaluation is negligible. On the other side, as the frames between two I-frames are temporally consistent, the propagated positions can extract reasonable patches on the in-between frames. Our sampling strategy is simple yet effective and can further reduce the computational cost of EMT.

Our contributions can be concluded as follows:

- We propose Efficient Meta-Tuning (EMT) for neural video delivery, significantly reducing the cost of training content-aware SR models and achieving even better performance.
- We present a novel challenging patch sampling strategy, which further reduces the cost of EMT. Our strategy improves the convergence of EMT with negligible additional cost.
- We conduct detailed experiments based on various efficient SR architectures to evaluate the advantage and generalization of our method.

2 Related work

DNN-based Image Super-Resolution SRCNN [8] is the first work that introduces DNNs to SR task. Their method consists of three stages, namely feature extraction, non-linear mapping and image reconstruction. With the rapid advance of DNN, plenty of methods are proposed to improve the performance of SISR following the pipeline of SRCNN. For example, VDSR [15] adopts a very deep DNN to predict the image residual instead of HR image. Motivated by ResNet [10], SRResNet [16] introduces Residual Block to the network and improves the SR performance. EDSR [18] modifies the structure of SRResNet by removing the Batch Normalization layer [11], further boosting the SR performance. RCAN [31] introduces the attention mechanism to the networks and presents deeper DNNs for SR. However, RCAN is computationally complicated, which limits its practical usage. To reduce the computational cost, many efficient methods are proposed for SR. ESPCN [23] uses LR image as input and up-samples the feature map by the pixel-shuffle layer to obtain the HR output. LAPAR [17] proposes a method based on linearly-assembled adaptive regression network. All of those methods are external methods, which train one model on large-scale image databases like DIV2K [1] and test on given input images. However, external methods fail to explore the overfitting property of DNNs, which can significantly boost the performance for practical video delivery system.

DNN-based Video Super-Resolution Different from image super-resolution, video super-resolution can additionally exploit the neighboring frames for SR. Therefore, temporal alignment plays an essential role and should be thoroughly considered. VESPCN [2] first predicts the motions between neighboring frames,

and then performs image warping before feeding neighboring frames into the SR network. However, it is difficult to accurately estimate the optical flow. TOFlow [27] proposes a task-oriented flow designed for specific video processing tasks. They jointly train the motion estimation component and video processing component in a self-supervised manner. DUF [12] solves the problem of accurate explicit motion compensation by training a network to generate dynamic up-sampling filters and a residual image. In order to reduce the computational cost of VSR, FRVSR [22] presents a recurrent framework that uses the previous SR result to super-resolve the following frame. Their recurrent framework naturally ensures temporally consistency and reduces the computational cost by warping only one image in each step. All these VSR approaches also belong to external methods that fail to explore the overfitting property of DNN. Apart from this, handling temporal alignment brings huge additional computational and storage costs, which limits their practical applications in resource-limited devices like mobile phone.

Neural Video Delivery NAS [29] is a promising Internet video delivery framework that integrates DNN for quality enhancement. It can solve the video quality degradation problem under limited Internet bandwidth. NAS can enhance the average Quality of Experience (QoE) by 43.08% using the same bandwidth budget, or saving 17.13% of bandwidth while providing the same user QoE. The main idea is to leverage DNN’s overfitting property and use the training accuracy to deliver high SR performance. Many following works are proposed to apply the idea of NAS to different scenarios, like UAV video streaming [26], live streaming [14], 360 video streaming [7,5], volumetric video streaming [30], and mobile video streaming [28], etc. Recent methods [19,13] propose to further reduce the bandwidth budget by sharing most of the parameters over video chunks. Therefore, only a small portion of private parameters are streamed for each video chunk. However, they still need huge computational cost for network training and fail to study the scene conversion for constructing optimal video chunks.

3 Method

3.1 Overview

In this section, we present our method to significantly accelerate the training of content-aware SR models. Our method adapts the meta-learned model to the first chunk of the video, and sequentially adapts partial parameters of the previous model to the following chunks. The partial parameters are selected by gradient masking and the challenging patches are extracted for adaption. The pipeline of our method is illustrated in Fig. 1. We first introduce Efficient Meta-Tuning (EMT) to sequentially deliver the models from a meta-learned model in Sec. 3.2. Then we propose a novel challenging patch sampling strategy to further accelerate EMT in Sec. 3.3.

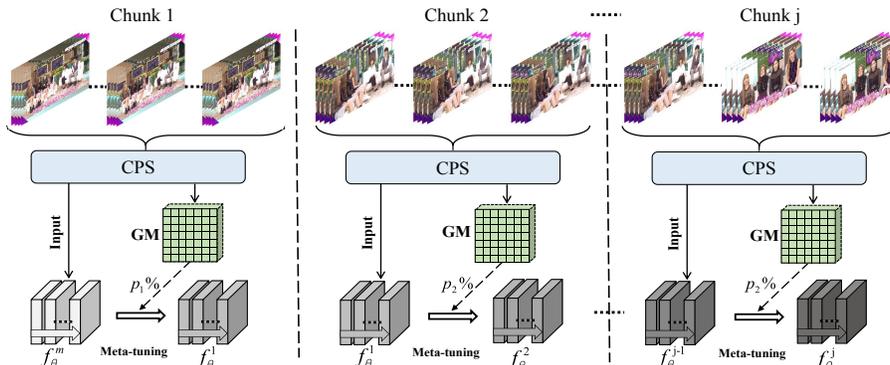


Fig. 1: The pipeline of our method. CPS and GM indicate challenging patch sampling and gradient masking respectively. The selected challenging patches are used to mask the parameters and fine-tune the models.

3.2 Efficient Meta-Tuning

Following former works [19,29], we uniformly divide the input video into chunks, and train the corresponding content-aware SR model for each chunk. [29] proposes to apply deep super-resolution networks to video delivery by training one model for each chunk from scratch. [19] presents a method to compress all the models by one shared model and a few private parameters. However, both methods train the content-aware SR models from scratch, resulting in huge computational cost. Since neighboring chunks are temporally consistent, fine-tuning is much more reasonable compared with training from scratch. Precisely, finding a generic initial model that can not only generalize over diverse video chunks but also adapt rapidly to any specific video chunk, plays a key role in fine-tuning. In order to obtain a better initialization, we adopt MAML [9] to train a meta-learned model. Although MAML has been applied to Zero-Shot SR [21,24] and video frame interpolation [6], it has never been studied in neural video delivery to the best of our knowledge. Compared with random initialization or pretrained initialization, a meta-learned model has better transferability. To obtain the content-aware SR models, we sequentially fine-tune partial parameters of the previous model. In contrast to fine-tuning the whole model, our method can compress the parameters of all models into one shared meta-learned model and a few partial parameters.

Meta-Learned Initialization We take one chunk as a specific task and aim to learn a SR model that can adapt to various chunks. Specifically, we first pretrain the SR model f_θ on DIV2K [1], then we utilize meta learning to optimize the pretrained parameters as illustrated by Alg. 1. This step enables the SR model to converge to a transferable point, which can be rapidly fine-tuned. In order to build a variety of tasks, we collect several video sequences and uniformly divide them into video chunks. Totally, we obtain N chunks for

Algorithm 1 Meta-Learned Initialization

Input: Initialized SR model f_θ , meta-learning dataset D_N
Output: Meta-learned model f_θ^m

- 1: **while** not done **do**
- 2: Sample n tasks D_n from D_N
- 3: **for** $t_i \in D_n$ **do**
- 4: Sample pairs (I_{HR}^i, I_{LR}^i) from t_i
- 5: Copy f_{θ_i} from the latest f_θ
- 6: Evaluate training loss according to Eq. 1
- 7: Update parameters according to Eq. 2
- 8: **end for**
- 9: Calculate $\mathcal{L}_{f_{\theta_i}}$ with respect to t_i
- 10: Update f_θ according to Eq. 3
- 11: **end while**
- 12: **return** Meta-learned model f_θ^m

meta-learning, and each chunk is set as the task t_i . The collected dataset is denoted as D_N . We apply bicubic downsampling to the frames and generate LR-HR pairs (I_{LR}^i, I_{HR}^i) . Our goal is to optimize f_θ according to each LR-HR pair by minimizing the L1 loss as shown in Eq. 1.

$$\mathcal{L}_f = |f_\theta(I_{LR}) - I_{HR}|_1 \quad (1)$$

During the inner loop (Line 4-7), we conduct one or more gradient updates for the task t_i in each iteration. The temporary model for task t_i is denoted as f_{θ_i} . During each inner gradient update, the task-specific parameters are updated according to Eq. 2, where α is the inner learning rate.

$$f_{\theta_i} \leftarrow f_{\theta_i} - \alpha \nabla_{f_{\theta_i}} \mathcal{L}_{f_{\theta_i}} \quad (2)$$

As for the outer loop (Line 9-10), we evaluate the loss of f_{θ_i} on each t_i and sum up the losses of all tasks to update the SR model f_θ . For one outer gradient update, it considers the gradients from all tasks. The outer update can be formulated as Eq. 3, where β is the outer learning rate.

$$f_\theta \leftarrow f_\theta - \beta \nabla_{f_\theta} \sum_i \mathcal{L}_{f_{\theta_i}} \quad (3)$$

Partial Model Adaption The meta-learned model f_θ^m is shared by all the video chunks. To obtain the content-aware SR models for the input video, we adapt the meta-learned model to the first chunk, and sequentially fine-tune the partial parameters of previous adapted model. Formally, we denote the content-aware SR model of j^{th} chunk as f_θ^j . For the first chunk, we rapidly fine-tune the meta-learned model f_θ^m to obtain the adapted model f_θ^1 . We use the gradient masking to select $p_1\%$ most significant parameters before adapting f_θ^m . As for other chunks, we fine-tune previous model f_θ^{j-1} on j^{th} chunk, delivering the adapted model f_θ^j . We also adopt the gradient masking to select $p_2\%$ most

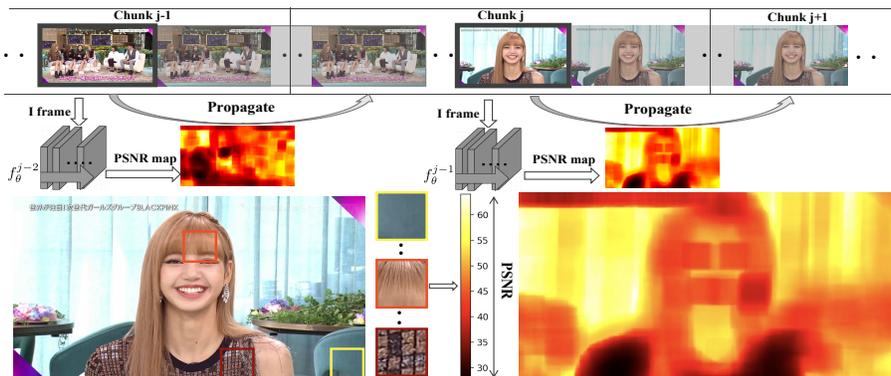


Fig. 2: Illustration of challenging patch sampling. Different colors in PSNR map represent patch’s difficult levels for previous chunk’s model. ‘propagate’ arrow indicates the propagation of PSNR map from I frame to its in-between frames.

significant parameters before fine-tuning f_{θ}^{j-1} . It has to be noted that we adopt different percentages for the first chunk and other chunks in this work. Our partial model adaption requires much fewer epochs for both the first chunk and other chunks. Therefore, the computational cost can be greatly reduced compared with training from scratch under the same performance.

Gradient Masking In order to compress the parameters of content-aware models, we design a simple yet effective strategy to find the $p\%$ most significant parameters. Given a reference model f , we need to find a fraction of parameters before fine-tuning f . Specifically, we adopt a few iterations to update f , obtaining a temporal model \hat{f} . Afterwards, we calculate $|\theta_{\hat{f}} - \theta_f|$ and choose the $p\%$ parameters that vary most, delivering the parameter mask $M(\theta_f)$.

Once we collect the $p\%$ most significant parameters, we can fine-tune the reference model f and simply update the significant parameters. In this manner, our method only needs to store $p\%$ private parameters of the reference model. When fine-tuning for the first chunk, we choose the meta-learned model f_{θ}^m as the reference model. For j^{th} chunk of the input video, we choose the previous adapted model f_{θ}^{j-1} as the reference model.

3.3 Challenging Patch Sampling

Previous adapted SR model already possesses the ability to super-resolve current chunk due to the temporal consistency. Thus it can achieve satisfying results on the majority of regions. However, some hard regions are still challenging for the previous adapted model. Since SR networks are fine-tuned on sampled LR-HR patch pairs, we further present a strategy to sample the $r\%$ most challenging patches for EMT. Our strategy is highly efficient and brings negligible additional cost. Inspired by the video codec, it first locates the positions of challenging

patches in I-frames, then propagates the positions to in-between frames as shown in Fig. 2.

Formally, we denote the frames of input video as I_1, \dots, I_T , where T is the number of frames. The I-frames are denoted as I_{k_1}, \dots, I_{k_M} , where k_1, \dots, k_M are the indices of M I-frames. We also denote the chunk indices of M I-frames as c_1, \dots, c_M . In order to localize the challenging patches for I_{k_m} , we run the inference of previous adapted model $f_\theta^{c_m-1}$ to super-resolve the downsampled I-frame $I_{k_m}^{LR}$:

$$I_{k_m}^{SR} = f_\theta^{c_m-1}(I_{k_m}^{LR}) \quad (4)$$

We can calculate the PSNR between $I_{k_m}^{SR}$ and I_{k_m} in terms of all possible patches as illustrated by Fig. 2. The time of calculating the PSNR map for 1080P frames is 0.1 seconds. Therefore, it is time-consuming to produce PSNR maps for all frames. For instance, when dealing with a 45-second video, it usually takes around 135 seconds to generate PSNR maps for all frames. Instead, we localize the positions of $r\%$ ($r=20$) most challenging patches in I_{k_m} , and then extract the patches from frames between I_{k_m} and $I_{k_{m+1}}$ according to the coordinates. As the frames between two I-frames are temporally consistent, the localized positions at I-frame can also choose reasonable patches on the in-between frames. Since I-frame is very sparse within a video, the computational cost is negligible. For a 45-second video, the total time of position localization and patch extraction is 0.7 seconds. However, the results of EMT using challenging patches are the same as results using all frames to some extent. In this way, the training efforts of EMT focus on challenging patches, resulting in faster convergence.

4 Experiments

In this section, we conduct extensive experiments to show the advantages of our method. The experimental details are given in Sec. 4.1. We first present the comparison with baseline and codec standards in Sec. 4.2, and then compare EMT with other neural video delivery methods in Sec. 4.3. We also conduct comprehensive ablation study to evaluate the contribution of each component in Sec. 4.4. In order to show the generalization ability, we report results across different scaling factors and architectures in Sec. 4.5.

4.1 Experimental Details

For meta-learning, we conduct two gradient updates for each individual task in the inner loop. After updating on all sampled tasks, we conduct one outer gradient update on the SR model. We randomly select training patches with a resolution of 144×144 and set mini-batch size as $16 * n$, where n is the number of sampled tasks. Particularly, we set n to 15 and each task consists of 50 frames. We set inner learning rate $\alpha = 0.5e - 5$ and outer learning rate $\beta = 1e - 3$. We adopt Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$. ESPCN serves as

the default architecture, x2 is utilized as the default scaling factor and PSNR is the default metric.

For fine-tuning, we conduct experiments on two video lengths, including 45 seconds and 2 minutes. The batch size is 16 and learning rate is $1e-4$. We set $p_1\%$ as 20% and $p_2\%$ as 1% to compress the parameters as default. We design three settings of our method, including S, M, and L. S and M settings adopt 0.1 epoch and 3 epochs for fine-tuning respectively. As for L setting, we alter $p_1\%$ to 100%. We conduct all the experiments on NVIDIA V100 GPUs.

4.2 Comparison with Baseline and Codec Standards

In this section, we compare our method against the baseline [29] and two codec standards. The baseline uniformly divides a video into chunks and trains one SR model for each chunk from scratch with 300 epochs. We denote the baseline as C_{1-n} . For the two commercial codec standards H.264 and H.265, we use ffmpeg with libx264 codec and libx265 codec to compress the HR videos to lower bit-rate while maintaining the resolution. The compressed videos are of the same storage size as our method (LR videos and SR models). We report three variants of our method under S, M, and L settings. Under the L setting, we aim to show the potential of meta-tuning by updating all the parameters for the first chunk. As shown in Tab. 1, our method achieves better performance with less time and parameters compared with baseline [29]. Our results with 0.1 epoch already outperform the baseline with 300 epochs. In terms of parameter compression, given a video with n chunks, we compress all models $n * P$ to $1 * p_1\%P + (n - 1) * p_2\%P$. In comparison with H.264 and H.265, our results outperform H.264 and H.265 in most cases as shown in Tab. 2. We also show the qualitative comparison in Fig. 3. As can be seen, our method can restore better details compared with codec standards.

4.3 Comparison with Neural Video Delivery Methods

In this section, we compare our work with other neural video delivery methods. CaFM [19] uses content-aware models to compress parameters and achieve competitive performance. However, its joint training strategy takes huge computational cost. Therefore, CaFM is not practical for delivering long videos. SRVC [13] also sequentially delivers the content-aware SR models by fine-tuning previous model. However, they fail to generalize on various architectures and have to train from scratch for a new video. Deep Video Compression (DVC) [20] is an end-to-end DNN-based video compression method. We also compare our method with (DVC) at four different bitrate-distortion trade-off operating points $\lambda \in \{256, 512, 1024, 2048\}$ (DVC1, DVC2, DVC3, DVC4). In Fig. 4 and Tab. 3, we demonstrate the advantages of our method in terms of accuracy, training time, and storage. In Fig. 4, we calculate the average PSNR and storage cost on 45s videos from VSD4K [19]. Under the same storage, our method outperforms other methods at most circumstances. Though CaFM achieves promising results, it takes a huge computational cost. As shown in Tab. 3, we demonstrate

Table 1: Comparisons with baseline [29]. We show the results of our method under S, M, and L settings. Paras indicates the model parameters and P denotes the parameters of ESPCN. m and h in Time column represent minutes and hours respectively.

	inter-45s			sport-45s			game-45s		
	PSNR	Paras	Time	PSNR	Paras	Time	PSNR	Paras	Time
C_{1-n}	38.95	9P	11.2h	46.03	9P	11.2h	35.61	9P	11.2h
Ours(S)	39.08	0.28P	1.2m	46.11	0.28P	1.2m	36.32	0.28P	1.2m
Ours(M)	39.18	0.28P	7.6m	46.25	0.28P	7.6m	36.51	0.28P	7.6m
Ours(L)	39.56	1.08P	55.5m	46.41	1.08P	1.76h	37.09	1.08P	1.64h
	dance-45s			vlog-45s			game-2min		
	PSNR	Paras	Time	PSNR	Paras	Time	PSNR	Paras	Time
C_{1-n}	43.47	9P	11.2h	46.20	9P	11.2h	34.46	24P	11.2h
Ours(S)	44.13	0.28P	1.2m	46.48	0.28P	1.2m	34.35	0.43P	1.2m
Ours(M)	44.24	0.28P	7.6m	46.71	0.28P	7.6m	34.50	0.43P	7.6m
Ours(L)	44.59	1.08P	1.53h	46.97	1.08P	48.6m	35.33	1.23P	56.3m

Table 2: Comparisons with H.264 and H.265. Under the same storage size, our PSNR results outperform H.264 and H.265 in most cases.

	inter		sport		game		dance		vlog	
	45s	2min								
H.264	36.32	44.82	38.29	36.75	37.29	35.13	28.86	29.36	42.37	43.02
H.265	36.78	45.47	39.98	38.44	38.17	36.87	30.83	31.06	43.35	43.97
Ours(M)	39.18	46.73	46.25	40.44	36.51	34.50	44.24	42.74	46.71	48.66
Storage(mB)	13.19	35.21	13.65	37.43	13.97	35.17	13.81	36.80	13.83	36.37

the trade-off between accuracy and training time. Our method shows competitive performance while maintaining low computational cost.

4.4 Ablation Study

Variants of EMT We intend to evaluate the contribution of each component in EMT. Firstly, we compare our method with pretrained initialization, which is denoted as P_{1-n} . Similar to our meta-learned initialization, P_{1-n} first initializes the SR model on DIV2K, and is normally finetuned on the meta-learning dataset D_N . We replace the meta-learned model of EMT by the normally finetuned model to obtain the results of P_{1-n} . To be mentioned, P_{1-n} still utilizes gradient masking and challenging patch sampling for a fair comparison. Therefore, we are able to evaluate the effectiveness of meta-learned initialization. To evaluate the effectiveness of Challenging Patch Sampling (CPS), we remove the step of patch sampling in our full pipeline and the results are denoted as MT_{1-n} . As shown in Tab. 4, in order to achieve the same PSNR, P_{1-n} takes extra cost compared with

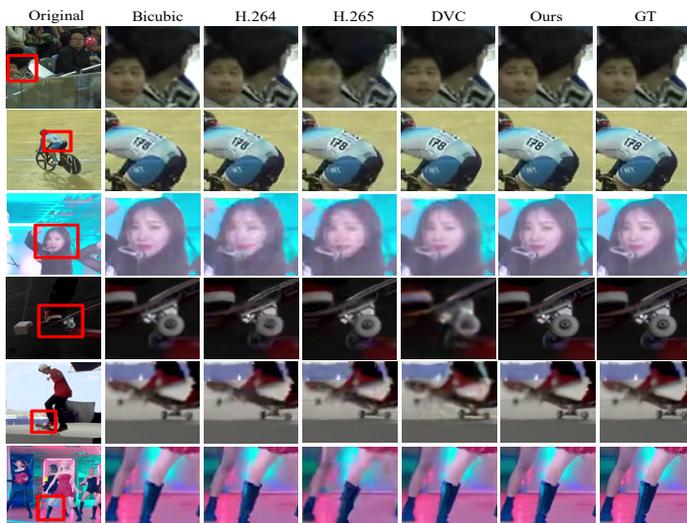


Fig. 3: Qualitative results. This figure shows the qualitative comparison against H.264, H.265 and DVC. Our method can restore better details compared with other methods. Best viewed by zooming x4.

our meta-learned initialization. Meanwhile, our method outperforms MT_{1-n} in regard to efficiency, demonstrating the effectiveness of CPS.

Variants of Meta Learning During meta-learning, we randomly sample 15 tasks and each task contains 50 frames. We also study the effect of different numbers of tasks and frames. For the number of tasks, we compare the results of 10, 15 and 20. For the number of frames, we evaluate the results of 10, 50 and 150. As shown in Tab. 5, all the variants achieve similar results, and the setting of 15 tasks with 50 frames already achieves competitive performance.

Variants of Gradient Masking We conduct extensive experiments to explore the performance under different variants of gradient masking. Since our method uses different portions of parameters for fine-tuning the first chunk and other chunks. We report the results of $p_1 \in \{10, 20, 30, 100\}$ and $p_2 \in \{0.5, 1, 5, 10\}$. As can be seen from Tab.6, we empirically set $p_1 = 20$ and $p_2 = 1$ as our default setting since it can already achieve satisfying results.

4.5 The Generalization of Our Method

Generalization of Various Scaling Factors We evaluate the generalization ability of EMT using ESPCN as the backbone across scaling factors x2, x3 and x4. As shown in Tab. 7, EMT outperforms C_{1-n} under various scaling factors, demonstrating the generalization ability of EMT across various scaling factors.

Generalization of Various Efficient Backbones In this part, we present the results of our method using different efficient SR backbones. We evaluate

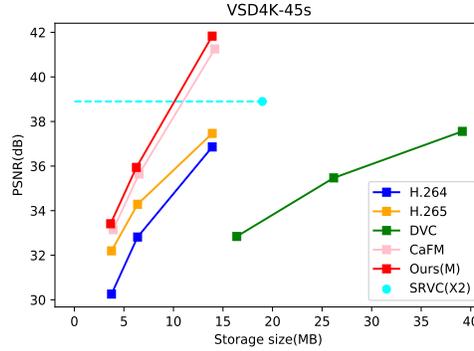


Fig. 4: Comparisons with neural video delivery methods in terms of PSNR and storage.

Table 3: Comparisons with neural video delivery in terms of PSNR and training time. Red and blue indicate the best and the second best results among all methods.

	inter-45s		sport-45s		game-45s		dance-45s		vlog-45s	
	Acc	Time								
C1-n	38.95	11.2h	46.03	11.2h	35.61	11.2h	43.47	11.2h	46.20	11.2h
CaFM	38.90	10.2h	46.12	10.2h	35.96	10.2h	43.63	10.2h	46.45	10.2h
SRVC	37.26	12.1m	41.38	12.1m	33.34	12.1m	40.87	12.1m	45.59	12.1m
DVC1	31.98	35.6m	35.52	33.6m	31.76	35.8m	27.67	34.8m	37.86	35.3m
DVC2	34.44	36.1m	37.45	34.3m	33.93	36.0m	32.46	35.2m	39.92	35.8m
DVC3	36.60	37.1m	39.58	34.8m	36.10	36.5m	34.40	35.4m	41.67	36.2m
DVC4	38.70	38.0m	41.28	34.8m	38.10	36.2m	36.33	35.8m	43.22	36.4m
Ours(M)	39.18	7.6m	46.25	7.6m	36.51	7.6m	44.24	7.6m	46.71	7.6m

Table 4: Variants of EMT. We show the results of our method under S setting. MT stands for meta-tuning strategy. m and h in Time column represent minutes and hours respectively. P denotes the parameters of SR architecture (ESPCN).

Method	MT	CPS	inter-45s			sport-45s		
			PSNR	Paras	Time	PSNR	Paras	Time
C_{1-n}	-	-	38.95	9P	11.2h	46.03	9P	11.2h
P_{1-n}	-	✓	39.08	0.28P	18.4m	46.11	0.28P	4.2m
MT_{1-n}	✓	-	39.08	0.28P	9.7m	46.11	0.28P	2.7m
Ours(S)	✓	✓	39.08	0.28P	1.2m	46.11	0.28P	1.2m

our method on 45 seconds videos and adopt three additional efficient backbones, including SRCNN [8], FSRCNN [8], and EDSR with one residual block in body

Table 5: Variants of meta-learning. Ours(S) adopt S setting.

	task	frame	inter-45s	sport-45s	game-45s	vlog-45s
C_{1-n}	-	-	38.95	46.03	35.61	46.20
Ours(S)	10	50	39.06	46.06	36.25	46.51
Ours(S)	15	50	39.08	46.11	36.32	46.48
Ours(S)	20	50	39.07	46.09	36.22	46.45
Ours(S)	15	10	39.09	45.94	36.34	46.57
Ours(S)	15	50	39.08	46.11	36.32	46.48
Ours(S)	15	150	39.03	46.20	36.39	46.58

Table 6: Variants of gradient masking. Ours(S) adopt S setting.

	p1%	p2 %	inter-45s	sport-45s	game-45s	vlog-45s
C_{1-n}	-	-	38.95	46.03	35.61	46.20
Ours(S)	10	1	39.04	46.00	36.28	46.42
Ours(S)	20	1	39.08	46.11	36.32	46.48
Ours(S)	30	1	39.11	46.13	36.33	46.49
Ours(S)	100	1	39.19	46.13	36.35	46.56
Ours(S)	20	0.5	39.08	46.10	36.31	46.47
Ours(S)	20	1	39.08	46.11	36.32	46.48
Ours(S)	20	5	39.09	46.18	36.34	46.53
Ours(S)	20	10	39.11	46.23	36.37	46.57

Table 7: Results of ESPCN on various scaling factors. We show the results of our method under M setting for fine-tuning.

	inter-45s			sport-45s			vlog-45s		
	x2	x3	x4	x2	x3	x4	x2	x3	x4
C_{1-n}	38.95	32.19	28.73	46.03	40.43	37.21	46.20	41.68	39.52
Ours(M)	39.18	32.55	29.05	46.25	40.51	37.22	46.71	42.25	40.08
	dance-45s			game-45s			city-45s		
	x2	x3	x4	x2	x3	x4	x2	x3	x4
C_{1-n}	43.47	36.86	35.22	35.61	30.67	28.80	36.44	31.60	29.23
Ours(M)	44.24	37.42	35.88	36.51	31.13	29.01	36.42	31.56	29.16

(EDSR-1). As shown in Tab. 8, our method also generalizes well to different efficient backbones, validating the generalization ability of our method.

5 Extension to Long Videos

In this section, we extend EMT to long videos beyond 2 minutes. Since former works [19,29] take too much time to train the content-aware SR models, we

Table 8: Results of various SR architectures. We show the results of our method under M setting for fine-tuning.

Backbone	inter-45s		sport-45s		game-45s		dance-45s		vlog-45s	
	C_{1-n}	Ours								
SRCNN	39.02	39.06	46.21	46.19	35.37	35.64	43.28	43.92	46.25	46.42
FSRCNN	39.06	39.25	46.29	46.32	35.84	35.90	43.61	44.08	46.09	46.47
EDSR-1	39.17	39.13	45.99	46.02	35.51	35.58	44.04	44.05	46.28	46.28

Table 9: Comparisons with H.264 and H.265 on long videos. We show the results of our method using 3 epochs for fine-tuning. The storage is measured in megabytes.

	vlog-5min		vlog-10min		vlog-20min		vlog-30min	
	PSNR	Storage	PSNR	Storage	PSNR	Storage	PSNR	Storage
Ours(M)	37.67	18.62	38.33	35.64	38.31	71.19	38.41	145.12
H.264	34.68	18.62	37.15	35.64	35.29	71.19	35.02	145.12
H.265	36.75	18.62	35.78	35.64	37.18	71.19	37.07	145.12

only compare with commercial codec standards. Directly applying EMT to long videos may achieve degraded results since the temporal consistency between neighboring chunks is not always true for long videos. Therefore, in order to extend EMT to long videos, we first divide the video frames into groups and apply EMT to each group. To be specific, we extract all the I-frames from the input video and make each group contain 30 I-frames. As shown in Tab. 9, our method outperforms commercial codec standards even on long videos, showing the great potential of our method.

6 Conclusion

To pave the way for practical applications, we propose Efficient Meta-Tuning (EMT) to significantly reduce the computational cost of neural video delivery. Instead of training from scratch, EMT sequentially fine-tunes a meta-learned model to deliver the content-aware SR models of the input video. Gradient masking is introduced to select partial parameters for fine-tuning, compressing all the models into one shared model and a few private parameters. In addition, we also present a sampling strategy to extract the challenging patches for fine-tuning, further reducing the cost of EMT. We conduct detailed comparisons with the commercial codec standards and other neural video delivery methods to demonstrate the advantages of our approach. We hope this paper can inspire future work on neural video delivery.

References

1. Agustsson, E., Timofte, R.: 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. Caballero, J., Ledig, C., Aitken, A., Acosta, A., Totz, J., Wang, Z., Shi, W.: Real-time video super-resolution with spatio-temporal networks and motion compensation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4778–4787 (2017)
3. Chan, K.C., Wang, X., Yu, K., Dong, C., Loy, C.C.: Basicvsr: The search for essential components in video super-resolution and beyond. arXiv preprint arXiv:2012.02181 (2020)
4. Chan, K.C., Zhou, S., Xu, X., Loy, C.C.: Basicvsr++: Improving video super-resolution with enhanced propagation and alignment. arXiv preprint arXiv:2104.13371 (2021)
5. Chen, J., Hu, M., Luo, Z., Wang, Z., Wu, D.: Sr360: boosting 360-degree video streaming with super-resolution. In: Proceedings of the 30th ACM Workshop on Network and Operating Systems Support for Digital Audio and Video. pp. 1–6 (2020)
6. Choi, M., Choi, J., Baik, S., Kim, T.H., Lee, K.M.: Scene-adaptive video frame interpolation via meta-learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9444–9453 (2020)
7. Dasari, M., Bhattacharya, A., Vargas, S., Sahu, P., Balasubramanian, A., Das, S.R.: Streaming 360-degree videos using super-resolution. In: IEEE INFOCOM 2020-IEEE Conference on Computer Communications. pp. 1977–1986. IEEE (2020)
8. Dong, C., Loy, C.C., He, K., Tang, X.: Learning a deep convolutional network for image super-resolution. In: European conference on computer vision. pp. 184–199. Springer (2014)
9. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: International Conference on Machine Learning. pp. 1126–1135. PMLR (2017)
10. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
11. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: International conference on machine learning. pp. 448–456. PMLR (2015)
12. Jo, Y., Oh, S.W., Kang, J., Kim, S.J.: Deep video super-resolution network using dynamic upsampling filters without explicit motion compensation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3224–3232 (2018)
13. Khani, M., Sivaraman, V., Alizadeh, M.: Efficient video compression via content-adaptive super-resolution. arXiv preprint arXiv:2104.02322 (2021)
14. Kim, J., Jung, Y., Yeo, H., Ye, J., Han, D.: Neural-enhanced live streaming: Improving live video ingest via online learning. In: Proceedings of the Annual conference of the ACM Special Interest Group on Data Communication on the applications, technologies, architectures, and protocols for computer communication. pp. 107–125 (2020)
15. Kim, J., Kwon Lee, J., Mu Lee, K.: Accurate image super-resolution using very deep convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1646–1654 (2016)

16. Ledig, C., Theis, L., Huszár, F., Caballero, J., Cunningham, A., Acosta, A., Aitken, A., Tejani, A., Totz, J., Wang, Z., et al.: Photo-realistic single image super-resolution using a generative adversarial network. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4681–4690 (2017)
17. Li, W., Zhou, K., Qi, L., Jiang, N., Lu, J., Jia, J.: Lapar: Linearly-assembled pixel-adaptive regression network for single image super-resolution and beyond. *Advances in Neural Information Processing Systems* **33** (2020)
18. Lim, B., Son, S., Kim, H., Nah, S., Mu Lee, K.: Enhanced deep residual networks for single image super-resolution. In: Proceedings of the IEEE conference on computer vision and pattern recognition workshops. pp. 136–144 (2017)
19. Liu, J., Lu, M., Chen, K., Li, X., Wang, S., Wang, Z., Wu, E., Chen, Y., Zhang, C., Wu, M.: Overfitting the data: Compact neural video delivery via content-aware feature modulation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4631–4640 (2021)
20. Lu, G., Ouyang, W., Xu, D., Zhang, X., Cai, C., Gao, Z.: Dvc: An end-to-end deep video compression framework. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11006–11015 (2019)
21. Park, S., Yoo, J., Cho, D., Kim, J., Kim, T.H.: Fast adaptation to super-resolution networks via meta-learning. *arXiv preprint arXiv:2001.02905* **5** (2020)
22. Sajjadi, M.S., Vemulapalli, R., Brown, M.: Frame-recurrent video super-resolution. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 6626–6634 (2018)
23. Shi, W., Caballero, J., Huszár, F., Totz, J., Aitken, A.P., Bishop, R., Rueckert, D., Wang, Z.: 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. pp. 1874–1883 (2016)
24. Soh, J.W., Cho, S., Cho, N.I.: Meta-transfer learning for zero-shot super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3516–3525 (2020)
25. Wang, X., Chan, K.C., Yu, K., Dong, C., Change Loy, C.: Edvr: Video restoration with enhanced deformable convolutional networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. pp. 0–0 (2019)
26. Xiao, X., Wang, W., Chen, T., Cao, Y., Jiang, T., Zhang, Q.: Sensor-augmented neural adaptive bitrate video streaming on uavs. *IEEE Transactions on Multimedia* **22**(6), 1567–1576 (2019)
27. Xue, T., Chen, B., Wu, J., Wei, D., Freeman, W.T.: Video enhancement with task-oriented flow. *International Journal of Computer Vision* **127**(8), 1106–1125 (2019)
28. Yeo, H., Chong, C.J., Jung, Y., Ye, J., Han, D.: Nemo: enabling neural-enhanced video streaming on commodity mobile devices. In: Proceedings of the 26th Annual International Conference on Mobile Computing and Networking. pp. 1–14 (2020)
29. Yeo, H., Jung, Y., Kim, J., Shin, J., Han, D.: Neural adaptive content-aware internet video delivery. In: 13th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 18). pp. 645–661 (2018)
30. Zhang, A., Wang, C., Liu, X., Han, B., Qian, F.: Mobile volumetric video streaming enhanced by super resolution. In: Proceedings of the 18th International Conference on Mobile Systems, Applications, and Services. pp. 462–463 (2020)
31. Zhang, Y., Li, K., Li, K., Wang, L., Zhong, B., Fu, Y.: Image super-resolution using very deep residual channel attention networks. In: Proceedings of the European conference on computer vision (ECCV). pp. 286–301 (2018)