# Reference-based Restoration of Digitized Analog Videotapes

Lorenzo Agnolucci    Leonardo Galteri    Marco Bertini    Alberto Del Bimbo

University of Florence - Media Integration and Communication Center (MICC)
Florence, Italy

[name.surname]@unifi.it

## Abstract

*Analog magnetic tapes have been the main video data storage device for several decades. Videos stored on analog videotapes exhibit unique degradation patterns caused by tape aging and reader device malfunctioning that are different from those observed in film and digital video restoration tasks. In this work, we present a reference-based approach for the resToration of digitized Analog videotaPEs (TAPE). We leverage CLIP for zero-shot artifact detection to identify the cleanest frames of each video through textual prompts describing different artifacts. Then, we select the clean frames most similar to the input ones and employ them as references. We design a transformer-based Swin-UNet network that exploits both neighboring and reference frames via our Multi-Reference Spatial Feature Fusion (MRSFF) blocks. MRSFF blocks rely on cross-attention and attention pooling to take advantage of the most useful parts of each reference frame. To address the absence of ground truth in real-world videos, we create a synthetic dataset of videos exhibiting artifacts that closely resemble those commonly found in analog videotapes. Both quantitative and qualitative experiments show the effectiveness of our approach compared to other state-of-the-art methods. The code, the model, and the synthetic dataset are publicly available at* https://github.com/miccunifi/TAPE.

## 1. Introduction

Analog magnetic tapes were widely used for video data storage from the 1950s to the late 1990s, and are still found in numerous archives worldwide with inestimable cultural value. The aging of recording supports and the malfunctioning of reader devices introduce several types of artifacts that are unfortunately very common for analog videotapes, such as drop out, tape crease, scanline flickering, and tape mistracking [44, 47]. These degradations limit the distribution and consumption of content by the general public if they are not properly restored. Nowadays, the restoration is usually
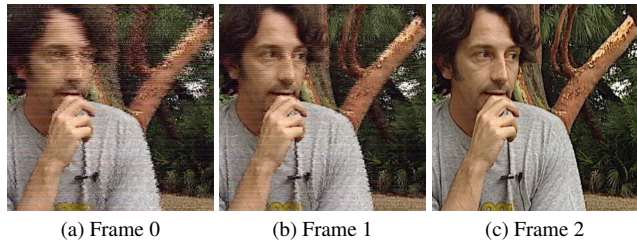


(a) Frame 0        (b) Frame 1        (c) Frame 2

Figure 1. Three consecutive frames of a real-world archival analog video. The degradations have significantly different intensities and positions and the temporal consistency is completely lost, which shows the time-varying nature of the artifacts. The third frame is essentially clean and could serve as a reference for the restoration.

performed frame-by-frame by experienced archivists using specialized commercial solutions, hence at great economic and time cost. Since restoring a few hours of content can take up to weeks of work, applying this process to entire archives is essentially unfeasible. For this reason, the development of methods for the automatic restoration of analog videos is a necessity.

Standard video restoration works [20, 25–27, 51] are designed for digital videos and do not consider the artifacts caused by media issues, which are more peculiar, more severe, and may occur simultaneously (see Fig. 1a). Additionally, due to the severity of the degradation, the temporal consistency between the video frames is completely lost, which makes optical flow-based frame alignment techniques commonly used in video restoration [10, 54] detrimental. Old video restoration methods consider the degradation related to the recording medium but focus on structured defects such as scratches [19, 48] or do not take advantage of the characteristics of the artifacts of analog videos [3, 45]. For instance, Agnolucci *et al*. [3] tackles analog video restoration with a transformer-based architecture that simply exploits the spatio-temporal information of input neighboring frames, thus without designing any strategy tailored for the task. In contrast, we propose an approach for analog video restoration that aims to remove the degradation caused by media issues – especially the most severe

1

ones, such as tape mistracking – by leveraging their temporally varying nature.

Figure 1 shows three consecutive frames that come from a real-world analog archive video. In the first two frames, the artifacts differ greatly in intensity and location, resulting in a complete loss of temporal consistency. In contrast, the last frame is essentially clean, so it contains a lot of high-quality details that could be useful for restoring temporally distant but similar degraded frames. This example shows the time-varying nature of the degradation in analog videos and motivates the main idea proposed in this work: identify the least damaged frames of a video and employ them as references for the restoration. To this end, we propose employing CLIP [38] for frame classification through zero-shot artifact detection by measuring the similarity between each frame and multiple textual prompts corresponding to different types of degradation. Then, within the frames classified as clean, we select those most similar to the input ones and employ them as references. We present a transformer-based Swin-UNet architecture that restores multiple input frames at once by exploiting the spatio-temporal information of neighboring frames while taking advantage of the references through our Multi-Reference Spatial Feature Fusion (MRSFF) block. MRSFF blocks employ multi-reference spatial cross-attention and attention pooling to restore the high-quality details that are lost in the input frames but still present in the cleaner reference ones. Figure 2 shows an overview of the proposed approach, named TAPE (resToration of digitized Analog videotaPEs). To overcome the lack of ground truth in real-world videos, we introduce a synthetic dataset of videos with artificially generated artifacts that emulate the degradation commonly found in analog ones. We hope that releasing our dataset will stimulate the research on analog video restoration to help preserve decades of video archives. Experimental results show the effectiveness of TAPE for both synthetic and real-world videos and its superior performance compared to state-of-the-art video restoration methods.

Our contributions can be summarized as follows:

- We propose TAPE, a reference-based approach for analog video restoration that exploits the time-varying nature of the artifacts typical of this media by exploiting the cleanest frames of a video to restore the degraded ones;
- We propose to leverage CLIP for zero-shot artifact detection to identify the clean frames and select the references based on the similarity to the input ones;
- We develop a Swin-UNet architecture that takes advantage of the reference frames through our Multi-Reference Spatial Feature Fusion blocks;
- Quantitative and qualitative experiments on synthetic and real-world videos prove the effectiveness of our method when compared with state-of-the-art video restoration techniques.

## 2. Related Work

**Video Restoration** Video restoration comprises various tasks that aim to enhance video quality, such as super-resolution and deblurring, each with its unique characteristics and methods [9, 15, 27, 34, 46]. For example, MANA [51] tackles video super-resolution through cross-frame non-local attention and a memory bank learned during training, while [20] proposes a multi-scale memory-based architecture for video deblurring. Other works propose approaches that achieve promising performance on multiple tasks [10, 25, 26, 54]. For example, RVRT [26] presents a recurrent transformer network with guided deformable attention, whereas BasicVSR++ [10] proposes a framework with second-order grid propagation and flow-guided deformable alignment. Despite the variety of degradations that are considered, these restoration techniques are designed for digital videos and do not address artifacts due to media issues typical of old ones.

**Old Video Restoration** Restoration of old videos must address degradation and issues related to the recording medium. Most of the existing methods focus on structured defects such as scratches and cracks typical of old films. Traditional approaches [16, 18, 22, 40] rely on a detection network based on handcrafted features followed by an inpainting pipeline. More recently, deep learning-based methods have been proposed. DeOldify [5] is an open-source tool to restore old films based on a NoGAN training approach that is commonly used to colorize grayscale videos. [19] presents a fully 3D convolutional for old video restoration and colorization. [48] develops RTN, a recurrent transformer network with a hidden state that improves the temporal consistency of the results. Although impressive results are observed for structured defects, none of these techniques are able to correct such large and strong artifacts as tape mistracking.

To the best of our knowledge, only a few works aim to remove the typical degradations of analog videotapes. [44, 45] exploit handcrafted features to detect artifacts and restore the corrupted tracking lines of interlaced videos by replacing them with those of adjacent frames. The most similar to our work is Agnolucci *et al.* [3], which restores analog videos with a Swin-UNet network that leverages the spatio-temporal information of neighboring input frames. Therefore, [3] proposes a generic method without any specific strategy developed to exploit the intrinsic characteristics of the degradation. In contrast, our method is tailored for analog video restoration and takes advantage of the time-varying nature of the artifacts. Indeed, we first identify the cleanest frames of a video with CLIP and then leverage them as references through our MRSFF blocks.
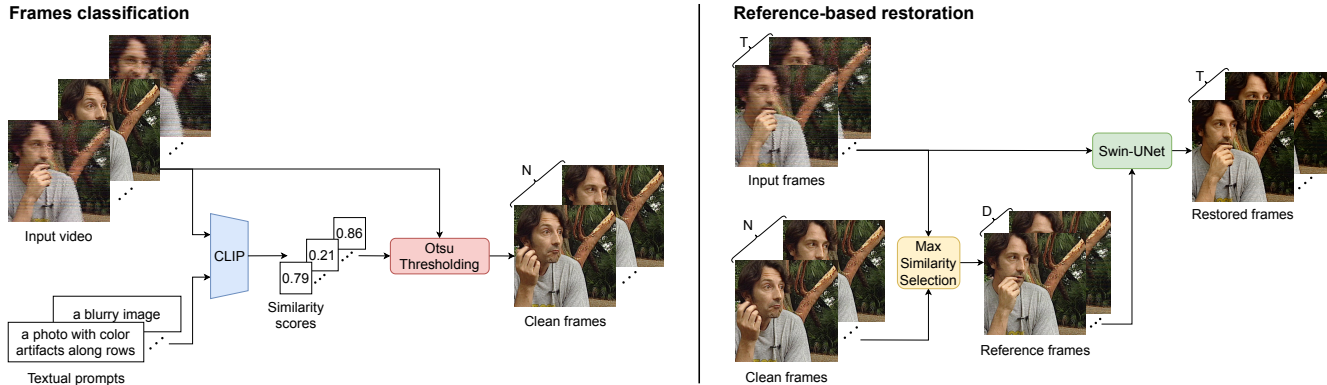
Figure 2. Overview of the proposed approach, named TAPE. *Left:* given a video, we identify the cleanest frames with CLIP. First, we measure the similarity between the frames and textual prompts that describe different artifacts. Then, we employ Otsu's method to define a threshold for classifying the frames based on their similarity scores, resulting in a set of clean frames. *Right:* given a window of $T$ degraded input frames, we select the most similar $D$ clean frames based on the CLIP image features and employ them as references. The proposed Swin-UNet then restores the input frames while effectively leveraging the references.

## 3. Proposed Approach

### 3.1. Synthetic Dataset

*Archivio Storico Luce*, renowned as the largest Italian historical video archive and one of the largest in the world, houses an extensive collection of material spanning the entirety of the 1900s, sourced from various origins. The curators provided us with some degraded analog archival videos. In the following, we refer to these videos as "real-world dataset", which consists of 4303 frames. Since these videos are obtained from the only available copy of the archive, they do not have a ground-truth counterpart and can not be used for training. Therefore, we create a synthetic dataset as similar as possible to real-world videos to train our model. Starting from high-quality videos belonging to the Harmonic dataset [2], we employ Adobe After Effects [12] to recreate different types of degradations [1, 17], which are visible in Fig. 1. In particular, we focus on: 1) tape mistracking and VHS edge waving, responsible for the horizontal displacement artifacts. These are the most complex distortions as they are the main cause of the temporal inconsistency; 2) chroma loss along the scanlines, visible from the horizontal cyan, magenta, and green lines; 3) tape noise, which is similar to Gaussian noise; 4) undersaturation, which makes the color of the frames look dull. We add these artifacts with random combinations, as well as positions and intensities, to each frame of the videos, purposely not to preserve temporal consistency. Indeed, as Fig. 1 shows, in real-world videos the degradations occur at the same time and change abruptly between consecutive frames. Since we start from the same real-world videos, our pipeline for generating synthetic artifacts resembles that of [3]. However, contrary to [3], we intend to release our synthetic dataset, in the hope of fostering further research on

this task. In the end, we have 26,392 frames corresponding to 40 clips, which we divide into training and test sets with a 75%-25% ratio. More details on the synthetic dataset and a qualitative comparison with the real-world dataset are provided in Sec. 6.

### 3.2. Frame Classification and Reference Selection

Given the time-varying nature of the artifacts in analog videos, we can identify the cleanest (*i.e.* least degraded) frames of a video to exploit them as references for restoring the other frames. In other words, we can classify the frames into fairly clean (*i.e.* the possible references) and very degraded ones. A possible approach could be training a binary classifier, but that would introduce two disadvantages: 1) we would have to choose a predefined threshold over which we consider a frame clean, which is an ill-posed problem; 2) all frames of a video could be classified as degraded, meaning that we would have no references. Hence, a continuous no-reference metric is more suitable for this task because it would allow one to define a different threshold for each video in an unsupervised manner.

We rely on CLIP [38] for our purpose. CLIP is a multimodal model trained with an image-caption alignment objective that obtained remarkable results in several downstream tasks, such as image generation [14], retrieval [6, 7], and quality assessment [49, 53]. Given a video, we propose to employ CLIP for zero-shot artifact detection by measuring the similarity between each frame and a textual prompt describing some kind of artifact (*e.g.* "*an image with color artifacts along rows*"). Therefore, a lower similarity corresponds to a less degraded frame. To improve the robustness of the predictions, we combine prompts referring to different types of artifacts with prompt ensembling [38] by averaging CLIP text features. We choose to use prompts
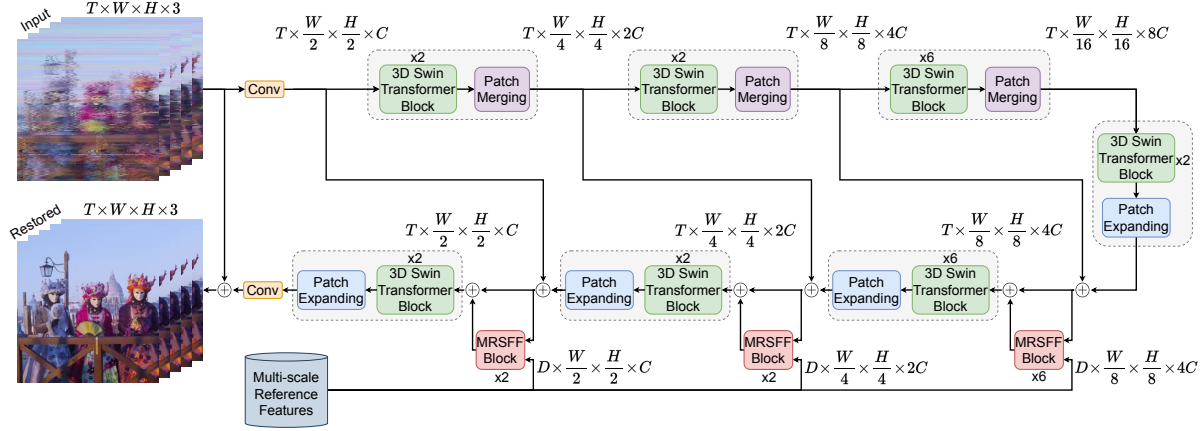
Figure 3. Overview of the proposed Swin-UNet architecture. We extract multi-scale features from the reference frames through a pre-trained Swin Transformer feature extractor. Then we exploit them through the proposed MRSFF blocks in the decoding part of the UNet.

that identify degraded frames because the visual artifacts they present are more easily described in natural language compared to the unconstrained domain of the content of clean ones. To adapt our method to address new types of degradation, it is enough to update the list of prompts, *e.g.* to specialize the approach for specific videos or types of medium. We provide more details on the list of prompts in Sec. 8.1. Then, we build a histogram of the CLIP similarities of all the frames. Intuitively, we expect this histogram to be bimodal, with the two peaks corresponding to degraded and almost clean frames. To define a threshold to classify the frames into two classes, we propose using Otsu's method [37], which is commonly adopted for automatic image thresholding [42]. The algorithm returns a single threshold determined by minimizing the intra-class variance in an unsupervised manner, thus based on the histogram of the similarities of each individual video. We classify as clean each frame with a CLIP similarity score lower than the threshold. In this way, we obtain a set of $N$ clean frames from the video, which we can leverage as references for the restoration. The number of clean frames $N$ changes for every video, as it is not a hyperparameter but rather depends on the degradation to which the video is subjected. The left-hand side of Fig. 2 shows an overview of the frame classification process.

Finally, given a window of $T$ input neighboring frames, we compute the cosine similarity of the CLIP image features between the central frame and the set of clean ones. Then, we take the $D$ most similar clean frames and use them as references. In our experience, artifacts in real-world videos never severely degrade all frames, so we always have significantly more than $D$ clean frames to choose from (see the supplementary material for more details). In other words, $N$ always exceeds $D$. Note that the arrangement of the reference frames is arbitrary as they are not nec-

essarily temporally neighboring. The right section of Fig. 2 shows the reference frames selection process and the subsequent reference-based restoration.

## 3.3. Swin-UNet Architecture

Figure 3 shows the proposed reference-based video restoration network. As in [3], our architecture is based on a Swin-UNet that restores a window of $T$ degraded frames at once. Since the complexity of attention is quadratic to the number of elements within the attention window, global attention on the entire input frames is unfeasible. Hence, we leverage the Swin Transformer [28, 29] to reduce the computational cost while maintaining the ability to learn long-range dependencies. Indeed, the Swin Transformer computes global attention only within local windows and then enables cross-window connections with a shifted-window mechanism. In this way, we can rely on transformers and their attention mechanism efficiently, in order to also take advantage of their ability to deal with frames that are not perfectly aligned. Indeed, due to the strong horizontal displacements caused by the artifacts, we cannot correctly align the input frames with the combination of explicit motion estimation and image warping. However, it has been shown that the expressiveness of the transformers alleviates this problem [25, 48].

First, we extract shallow features from the input with a convolutional layer. Then, in the encoder, we reduce the patch size and increase the number of channels through cascaded Swin 3D transformer blocks and patch merging layers [29]. In the decoder, we elaborate the processed features (*i.e.* those that go through the encoder and decoder) via Swin 3D transformer blocks and patch expanding layers, which consist of pixel shuffle layers. We use skip connections to add the encoder residual features to the processed ones. A skip connection between the input and the output
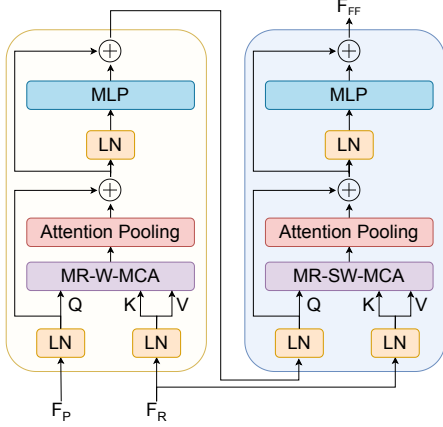
Figure 4. Two successive MRSFF blocks. LN stands for Layer Normalization. MR-(S)W-MCA represents the Multi-Reference-(Shifted)Window-Multi-head Cross Attention.

makes the network learn the residual of each frame.

Using Swin 3D transformer blocks allows taking advantage of the spatio-temporal information of the input frames thanks to self-attention. Indeed, by having the processed features attend themselves, we intuitively make each region of the input frames look at similar parts of the other frames. This is particularly useful due to the time-varying nature of the artifacts, as a highly degraded portion of one frame may be less damaged in one of the neighboring ones. However, if a given region is severely degraded in all the input frames, some details will be permanently lost. For this reason, we propose employing clean reference frames that do not belong to the window of the input frames. Contrary to [3], our model also takes as input $D$ reference frames – selected as explained in Sec. 3.2 – from which we extract multi-scale features with a pre-trained frozen Swin Transformer feature extractor. In the decoder, the proposed MRSFF blocks combine processed and reference features to recover details that would be lost otherwise.

## 3.4. Multi-Reference Spatial Feature Fusion Block

Given the reference features, our aim is to exploit them for the restoration. Since the reference frames are not contiguous in time (see Sec. 3.2), the use of temporal attention with Swin 3D transformer blocks is not suitable because it assumes some sort of temporal correlation. Furthermore, due to the variability in artifact locations across input frames, we want the processed features to independently attend to the reference ones along the spatial dimension. In this way we allow the processed features to focus on recovering the missing details in the corresponding input frames. For these reasons, we present our Multi-Reference Spatial Feature Fusion block, which makes use of Multi-Reference-(Shifted)Window-Multi-Head Cross Attention (MR-(S)W-

MCA) and attention pooling. MRSFF blocks are inspired by Swin 2D transformer blocks [28]. Indeed, both architectures rely on a shifted window-based attention mechanism along the spatial dimension. However, while Swin 2D uses self-attention, MRSFF blocks perform cross-attention between two inputs that may have different dimensions, since the number of input and reference frames T and D may not coincide. In addition, MRSFF blocks employ attention pooling to combine the information coming from the reference frames.

Similarly to [28], we partition each feature $F \in \mathbb{R}^{H \times W \times C}$ into $\frac{HW}{M^2}$ non-overlapping $M \times M$ windows. Here, $H$ and $W$ are the height and width, respectively, and $C$ is the number of channels. We flatten the $M^2$ elements of each window and compute the local attention. In the successive MRSFF block, a $\lfloor \frac{M}{2} \rfloor \times \lfloor \frac{M}{2} \rfloor$ cyclic spatial shift [28] of the features before window partitioning enables cross-window connections. MR-(S)W-MCA is based on a cross-attention mechanism, in which the processed features associated with each input frame act as independent queries, while the reference features serve as both keys and values. Intuitively, each input frame looks at similar parts of the reference frames and then exploits them to recover its missing details. Figure 4 shows two consecutive MRSFF blocks.

Formally, let $F_P \in \mathbb{R}^{T \times M^2 \times C}$ and $F_R \in \mathbb{R}^{D \times M^2 \times C}$ be, respectively, the processed and reference features of a given local window. For each head, we compute:

$$Q_P = F_P P_Q \in \mathbb{R}^{T \times M^2 \times C}$$
$$K_R = F_R P_K, \ V_R = F_R P_V \in \mathbb{R}^{D \times M^2 \times C} \quad (1)$$

where $P_Q$, $P_K$, $P_V \in \mathbb{R}^{C \times C}$ are projection matrices. Let $Q_{P_i}$, $K_{R_j}$, $V_{R_j} \in \mathbb{R}^{M^2 \times C}$ be the $i$-th and $j$-th elements of the queries, keys and values, respectively. We compute the attention map $A_{ij} = SoftMax(Q_{P_i} K_{R_j}^T / \sqrt{C} + B) \in \mathbb{R}^{M^2 \times M^2}$. As in [28], we include a learnable relative positional embedding $B \in \mathbb{R}^{M^2 \times M^2}$ in each attention head. Conceptually, $A_{ij}$ represents the correlation between the $i$-th input frame and the $j$-th reference frame and is used for a weighted sum of $V_{R_j}$, formulated as $A_{ij} V_{R_j} \in \mathbb{R}^{M^2 \times C}$. For each input frame $i \in \{1, \ldots, T\}$, we repeat this computation for all the $D$ references and concatenate the results:

$$F_G = \text{MR-(S)W-MCA}(F_P, F_R) \in \mathbb{R}^{T \times D \times M^2 \times C} \quad (2)$$

Intuitively, the features pertaining to each of the $T$ frames represent the fusion of the information from the corresponding input frame with that of the $D$ reference ones. Consequently, it is imperative to effectively combine these features on the $D$ axis to fully leverage the reference frames. Drawing inspiration from CLIP [38], we propose to use attention pooling to achieve this goal. Attention pooling relies on the self-attention mechanism to select the most pertinent

features for a pooled representation, rather than just taking the average. In other words, attention pooling allows us to take advantage of the most valuable parts of the information coming from the reference frames to obtain an aggregate representation. Note that CLIP uses attention pooling to reduce the 3D feature maps of a CNN into single vectors, thus removing the height and width dimensions. Our aim, on the contrary, is to aggregate $D$ sequences of transformer tokens into a single one, thereby pooling along a single axis. As in [38], we condition the query on the global average-pooled representation. First, we compute the average of the features along the $D$ axis:

$$F_{avg} = \text{Avg}(F_G) \in \mathbb{R}^{T \times M^2 \times C} \qquad (3)$$

Then, we concatenate the $D$ sequences of tokens:

$$F_C = \text{Concatenation}(F_G) \in \mathbb{R}^{T \times DM^2 \times C} \qquad (4)$$

Finally, for each input frame, we project the features and rely on the attention mechanism by considering the average features $F_{avg}$ as the query and the concatenated ones $F_C$ as both the keys and values:

$$F_{FF} = \text{MHA}(F_{avg}, F_C, F_C) \in \mathbb{R}^{T \times M^2 \times C} \qquad (5)$$

where MHA represents both the projection and multi-head attention operations. In this way, we obtain an aggregate representation of the features that encapsulate the information from the reference frames. Then, we further transform the features via a layer normalization layer and an MLP. In the successive MRSFF block, the fused features $F_{FF}$ serve as the queries of the transformer, assuming the role held by the processed features $F_P$ in the preceding block. Nevertheless, in this case, the features are shifted before window partitioning to enable cross-window connections and therefore learn long-range dependencies.

Ultimately, the fused $F_{FF}$ and the processed $F_P$ features of each window share the same dimensions and can consequently be added together to recover the lost details.

## 4. Experimental Results

### 4.1. Implementation Details

We train the proposed model on the synthetic dataset for 100 epochs using the ADAMW optimizer [30] with weight decay equal to 0.01 and $(\beta_1, \beta_2) = (0.9, 0.99)$. The learning rate is set to $2e - 5$. During training, we randomly crop $128 \times 128$ patches from the frames. We set the number of input frames $T$ to 5; the number of reference frames $D$ to 5; $M$ and $C$ in Sec. 3.4 to 8 and 96. We train TAPE with a weighted sum of the Charbonnier loss [11] and a perceptual loss [13, 21, 23]. More details on the training loss are reported in Sec. 7. During testing, we crop the $512 \times 512$

center patch of each frame. Quantitative results are the average over the videos of the test set. On a single A100-40GB GPU, our model takes 2 days for training and runs at 15FPS at inference time, which is more than suitable for a task with no real-time requirements.

### 4.2. Baselines

We compare TAPE with the following baselines: 1) MANA [51]: a memory-augmented architecture with cross-frame non-local attention for video super-resolution; 2) MemDeblur [20]: a multi-scale memory-augmented recurrent architecture for video deblurring; 3) BasicVSR++ [10]: a recurrent framework with second-order grid propagation and flow-guided deformable alignment for video restoration; 4) RVRT [26]: a recurrent transformer with guided deformable attention for video restoration; 5) RTN [48]: a recurrent transformer network for old film restoration; 6) Agnolucci *et al*. [3]: a Swin-UNet architecture for analog video restoration.

The baselines comprise standard and old video restoration works to show the difference from analog video restoration. For a fair comparison, we trained all the baselines from scratch on our training dataset using the official repositories.

### 4.3. Evaluation Metrics

We evaluate the performance on the synthetic dataset using four full-reference metrics: two signal-based metrics, PSNR and SSIM [50], to measure the low-level difference between restored and ground-truth frames; LPIPS [52], which better correlates with human perceived visual quality; VMAF [24], a perceptual video quality assessment model that combines multiple elementary quality metrics, including one that accounts for the temporal difference between adjacent frames, thus evaluating the presence of motion jitter and flicker.

Since the ground truth is not available for the real-world dataset, we employ three no-reference image quality assessment metrics: BRISQUE [35], NIQE [36] and CONTRIQUE [32].

### 4.4. Quantitative Results

We report the quantitative results for the synthetic and real-world datasets in Tabs. 1 and 2, respectively. Considering the synthetic dataset, TAPE outperforms all baselines on all metrics by a large margin. In particular, LPIPS shows that our method produces results that are more perceptually accurate than the other techniques, while VMAF proves that our restored videos are more temporally consistent and with less motion jitter. Furthermore, we observe that BasicVSR++ [10] and RVRT [26] perform poorly, even though they represent the state of the art in standard video restoration. We attribute this result to the use of optical flow for frame alignment, which is detrimental for analog videos, as

| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ | VMAF ↑ |
|---|---|---|---|---|
| MANA [51] | 27.81 | 0.843 | 0.206 | 40.28 |
| MemDeblur [20] | 33.22 | 0.911 | 0.106 | 71.55 |
| BasicVSR++ [10] | 31.66 | 0.916 | 0.098 | 78.91 |
| RVRT [26] | 32.47 | 0.896 | 0.117 | 72.41 |
| RTN [48] | 31.46 | 0.905 | 0.100 | 56.76 |
| Agnolucci *et al.* [3] | 34.96 | 0.940 | 0.060 | 77.83 |
| **TAPE** | **35.53** | **0.946** | **0.052** | **83.61** |

Table 1. Quantitative results for the synthetic dataset. ↑ means that higher values are better, ↓ means that lower values are better. The best results are highlighted in bold.

| Method | BRISQUE ↓ | NIQE ↓ | CONTRIQUE ↓ |
|---|---|---|---|
| MANA [51] | **41.80** | **5.90** | 48.18 |
| MemDeblur [20] | 51.20 | 8.89 | 45.82 |
| BasicVSR++ [10] | 59.19 | 8.42 | 48.44 |
| RVRT [26] | 47.61 | 8.39 | 48.64 |
| RTN [48] | 53.27 | 6.94 | 46.17 |
| Agnolucci *et al.* [3] | 59.44 | 7.90 | 45.45 |
| **TAPE** | 56.04 | 7.74 | **42.99** |

Table 2. Quantitative results for the real-world dataset. ↓ means that lower values are better. The best results are highlighted in bold.

the degradation is so severe that it completely disrupts the temporal consistency. This outcome further shows the difference between standard and analog video restoration.

Considering the real-world dataset, our approach achieves the best results for CONTRIQUE, while MANA [51] performs better for BRISQUE and NIQE. However, the qualitative results presented in Sec. 4.5 show that MANA adds high-frequency artifacts to restored frames. We argue that BRISQUE and NIQE are misled by these artifacts and mistake them for high-frequency details that are instead typical of high-quality images [4, 41]. See Sec. 9 for more details.

## 4.5. Qualitative Results

Figure 5 shows the qualitative results for the synthetic dataset. We observe that TAPE generates the most detailed and photorealistic images. MANA and BasicVSR++ yield results that lack photorealism and fine-grained details, while MemDeblur, RVRT, and RTN generate images that exhibit a higher degree of accuracy, albeit accompanied by a discernible amount of noise. The results of Agnolucci *et al.* are satisfactory, but with fewer details compared to our approach. See for example the white boat in the first row or the brick wall in the second row.

We report qualitative results of the real-world dataset in Fig. 6. As mentioned in Sec. 4.4, despite the promising quantitative results, MANA and RTN generate clearly un-

| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ | VMAF ↑ |
|---|---|---|---|---|
| w/o classification | 35.49 | 0.945 | 0.054 | 82.83 |
| w/o CLIP | 35.42 | 0.945 | 0.053 | 83.37 |
| w/o MRSFF | 35.37 | 0.944 | 0.055 | 82.94 |
| w/o att. pooling | **35.56** | 0.945 | 0.053 | 83.22 |
| $D = 1$ | 35.37 | 0.945 | 0.053 | 83.05 |
| $D = 3$ | 35.49 | 0.946 | 0.053 | 83.53 |
| **TAPE** | 35.53 | **0.946** | **0.052** | **83.61** |

Table 3. Ablation studies on the synthetic dataset. Best results are highlighted in bold.

satisfactory images with many unpleasant artifacts, showing the unreliability of BRISQUE and NIQE for this task. RVRT proves to be unable to generalize to real-world videos, since it fails to remove most of the degradation from the input frames. MemDeblur, BasicVSR++, and Agnolucci *et al.* yield acceptable results, but with visible artifacts. Regarding the synthetic dataset, TAPE generates the most detailed and photorealistic images. In our results, the eyes of the subjects in the first and second rows show fewer artifacts, and the overall images are considerably cleaner, as can be seen from the tree in the background in the first row.

We provide some video restoration examples in the supplementary material.

## 4.6. Ablation Studies

**Frame Classification and Reference Selection** We conduct ablation studies on frame classification and reference selection: 1) *w/o classification*: we select the references from all the frames of the video, without performing frame classification; 2) *w/o CLIP*: we follow the same approach described in Sec. 3.2 but employing the similarity scores and features of BRISQUE [35] instead of relying on CLIP and textual prompts. The upper section of Tab. 3 shows the results. First, we observe that limiting the selection of the references to frames classified as clean improves the results. Differently, degraded frames can be used as references if they are similar enough to the input ones, bringing little improvement to the restoration process. Second, CLIP proves to be more effective than BRISQUE for frame classification. See Sec. 8.2 for more ablation studies on frame classification.

**MRSFF block** We perform ablation studies on the MRSFF block: 1) *w/o MRSFF*: we substitute the MRSFF blocks with the standard Swin 3D [29] ones. We make use of spatio-temporal cross-attention by considering the processed features as queries and the reference features as keys and values; 2) *w/o att. pooling*: we substitute attention pooling with a simple average pooling. We report the results
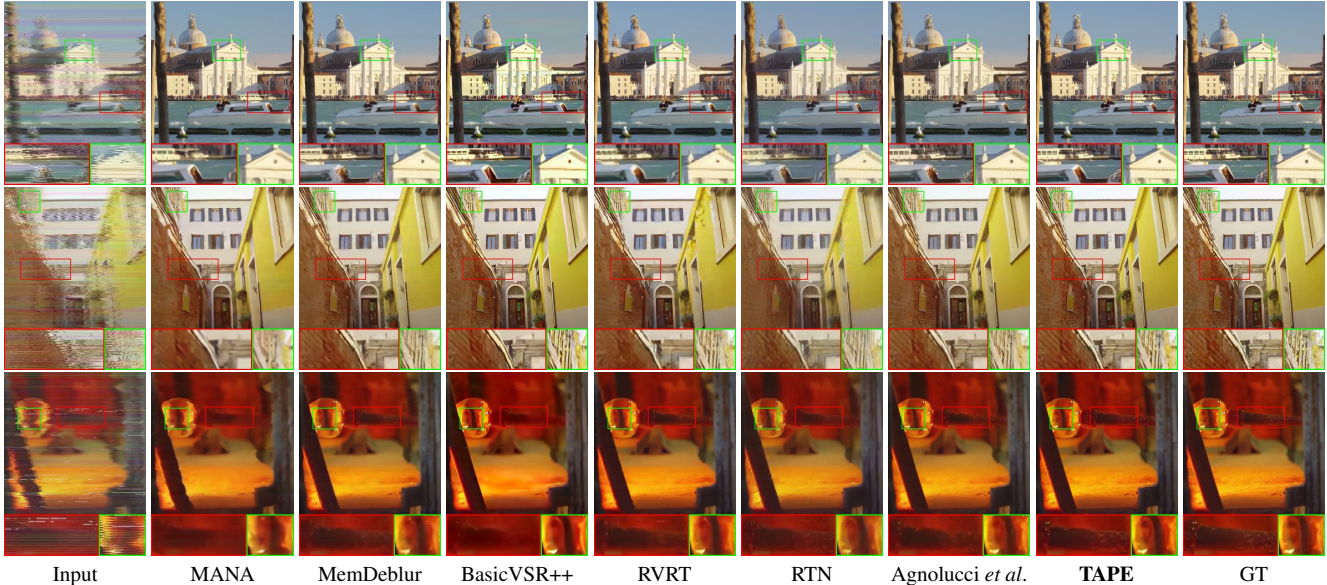
Figure 5. Qualitative results for the synthetic dataset. Best viewed in PDF.

| Input | MANA | MemDeblur | BasicVSR++ | RVRT | RTN | Agnolucci *et al*. | **TAPE** | GT |



Figure 6. Qualitative results for the real-world dataset. Best viewed in PDF.

| Input | MANA | MemDeblur | BasicVSR++ | RVRT | RTN | Agnolucci *et al*. | **TAPE** |

in the central section of Tab. 3. We notice that the spatio-temporal attention of the Swin 3D blocks performs worse than the spatial-only attention of the MRSFF ones. This is due to the fact that the reference frames are not contiguous in time – see Sec. 3.2 – so there is no temporal correlation between them. Moreover, attention pooling proves to be more effective in taking full advantage of the reference frames when compared to simple average pooling. Indeed, it allows us to select the most valuable features with self-attention to combine the information coming from the reference frames into an aggregate representation.

**Number of Reference Frames** We conduct experiments on the number of reference frames $D$, reducing it from 5

to 1 and 3. We expect that a higher number of references will lead to an improvement in the performance, as more information would be available for the processed features to leverage. The lower section of Tab. 3 shows that our expectation aligns with the experimental results.

## 5. Conclusion

In this paper, we present TAPE, a novel reference-based approach for the restoration of analog videotapes. Starting from real-world videos from an archive, we create a synthetic dataset degraded by artifacts typical of analog videotapes. We identify the cleanest frames of a video with CLIP by using textual prompts that describe different types of artifact. Then, we exploit those most similar to the de-

graded input frames as references via the proposed Swin-UNet architecture and the Multi-Reference Spatial Feature Fusion blocks. MRSFF blocks rely on cross-attention and attention pooling to recover the missing details in the input frames. Extensive experiments show the effectiveness of TAPE compared to several state-of-the-art techniques on both synthetic and real-world datasets. Our results demonstrate the differences between standard and analog video restoration, highlighting the need for approaches specifically designed for this task. In future work, we will develop a learned degradation model, similarly to [8, 31, 33], to efficiently create more accurate synthetic videos.

# References

[1] AV Artifact Atlas. http://www.avartifactatlas.com/. Accessed: 2022-02-17. 3

[2] Harmonic free 4K demo footage. https://www.harmonicinc.com/free-4k-demo-footage/, 2019. 3, 12

[3] Lorenzo Agnolucci, Leonardo Galteri, Marco Bertini, and Alberto Del Bimbo. Restoration of analog videos using Swin-UNet. In *Proceedings of the ACM International Conference on Multimedia*, pages 6985–6987, 2022. 1, 2, 3, 4, 5, 6, 7

[4] Lorenzo Agnolucci, Leonardo Galteri, Marco Bertini, and Alberto Del Bimbo. Perceptual quality improvement in videoconferencing using keyframes-based gan. *IEEE Transactions on Multimedia*, 2023. 7, 15

[5] Jason Antic. DeOldify. https://github.com/jantic/DeOldify, 2018. 2

[6] Alberto Baldrati, Marco Bertini, Tiberio Uricchio, and Alberto Del Bimbo. Conditioned and composed image retrieval combining and partially fine-tuning clip-based features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4959–4968, 2022. 3

[7] Alberto Baldrati, Lorenzo Agnolucci, Marco Bertini, and Alberto Del Bimbo. Zero-shot composed image retrieval with textual inversion. *arXiv preprint arXiv:2303.15247*, 2023. 3

[8] Adrian Bulat, Jing Yang, and Georgios Tzimiropoulos. To learn image super-resolution, use a GAN to learn how to do image degradation first. In *Proceedings of the European conference on computer vision (ECCV)*, pages 185–200, 2018. 9

[9] Jiezhang Cao, Yawei Li, Kai Zhang, and Luc Van Gool. Video super-resolution transformer. *arXiv preprint arXiv:2106.06847*, 2021. 2

[10] Kelvin CK Chan, Shangchen Zhou, Xiangyu Xu, and Chen Change Loy. BasicVSR++: Improving video super-resolution with enhanced propagation and alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5972–5981, 2022. 1, 2, 6, 7

[11] Pierre Charbonnier, Laure Blanc-Feraud, Gilles Aubert, and Michel Barlaud. Two deterministic half-quadratic regularization algorithms for computed imaging. In *Proceedings of the International Conference on Image Processing*, pages 168–172. IEEE, 1994. 6, 13

[12] Mark Christiansen. *Adobe After Effects CC Visual Effects and Compositing Studio Techniques*. Adobe Press, 2013. 3, 12

[13] Alexey Dosovitskiy and Thomas Brox. Generating images with perceptual similarity metrics based on deep networks. In *Proc. of Advances in Neural Information Processing Systems (NeurIPS)*, 2016. 6, 13

[14] Federico A Galatolo, Mario GCA Cimino, and Gigliola Vaglini. Generating images from caption and vice versa via CLIP-guided generative latent space search. *arXiv preprint arXiv:2102.01645*, 2021. 3

[15] Zhicheng Geng, Luming Liang, Tianyu Ding, and Ilya Zharkov. RSTT: Real-time spatial temporal transformer for space-time video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17441–17451, 2022. 2

[16] Ioannis Giakoumis, Nikos Nikolaidis, and Ioannis Pitas. Digital image processing techniques for the detection and removal of cracks in digitized paintings. *IEEE Transactions on Image Processing*, 15(1):178–188, 2005. 2

[17] Adam Hawkes. Identifying common tape defects for commercial video restoration. https://medium.com/pfclean-blog/identifying-common-tape-defects-for-commercial-video-restoration-677880c38daa, 2018. Accessed: 2022-10-30. 3

[18] Zhang Hongying, Wu Yadong, and Kuang Zhonglin. An efficient scratches detection and inpainting algorithm for old film restoration. In *Proc. of International Conference on Information Technology and Computer Science*, pages 75–78. IEEE, 2009. 2

[19] Satoshi Iizuka and Edgar Simo-Serra. Deepremaster: temporal source-reference attention networks for comprehensive video enhancement. *ACM Transactions on Graphics (TOG)*, 38(6):1–13, 2019. 1, 2

[20] Bo Ji and Angela Yao. Multi-scale memory-based video deblurring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1919–1928, 2022. 1, 2, 6, 7

[21] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Proc. of European Conference on Computer Vision (ECCV)*, 2016. 6, 13

[22] Seong-Whan Kim and Ki-Hong Ko. Efficient optimization of inpainting scheme and line scratch detection for old film restoration. In *Proc. of Pacific Rim International Conference on Artificial Intelligence*, pages 623–631. Springer, 2006. 2

[23] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 6, 13

[24] Zhi Li, Anne Aaron, Ioannis Katsavounidis, Anush Moorthy, and Megha Manohara. Toward a practical perceptual video quality metric. https://netflixtechblog.com/toward-a-practical-perceptual-video-quality-metric-653f208b9652, 2016. 6

[25] Jingyun Liang, Jiezhang Cao, Yuchen Fan, Kai Zhang, Rakesh Ranjan, Yawei Li, Radu Timofte, and Luc Van Gool. Vrt: A video restoration transformer. *arXiv preprint arXiv:2201.12288*, 2022. 1, 2, 4

[26] Jingyun Liang, Yuchen Fan, Xiaoyu Xiang, Rakesh Ranjan, Eddy Ilg, Simon Green, Jiezhang Cao, Kai Zhang, Radu Timofte, and Luc V Gool. Recurrent video restoration transformer with guided deformable attention. *Advances in Neural Information Processing Systems*, 35:378–393, 2022. 2, 6, 7

[27] Jing Lin, Yuanhao Cai, Xiaowan Hu, Haoqian Wang, Youliang Yan, Xueyi Zou, Henghui Ding, Yulun Zhang, Radu Timofte, and Luc Van Gool. Flow-guided sparse transformer for video deblurring. In *Proceedings of the 39th International Conference on Machine Learning*, pages 13334–13343. PMLR, 2022. 1, 2

[28] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 4, 5

[29] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3202–3211, 2022. 4, 7

[30] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6

[31] Zhengxiong Luo, Yan Huang, Shang Li, Liang Wang, and Tieniu Tan. Learning the degradation distribution for blind image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6063–6072, 2022. 9

[32] Pavan C Madhusudana, Neil Birkbeck, Yilin Wang, Balu Adsumilli, and Alan C Bovik. Image quality assessment using contrastive learning. *IEEE Transactions on Image Processing*, 31:4149–4161, 2022. 6, 14

[33] Shunta Maeda. Unpaired image super-resolution using pseudo-supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 291–300, 2020. 9

[34] Matteo Maggioni, Yibin Huang, Cheng Li, Shuai Xiao, Zhongqian Fu, and Fenglong Song. Efficient multi-stage video denoising with recurrent spatio-temporal fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3466–3475, 2021. 2

[35] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. No-reference image quality assessment in the spatial domain. *IEEE Transactions on image processing*, 21(12):4695–4708, 2012. 6, 7, 14

[36] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. Making a "completely blind" image quality analyzer. *IEEE Signal processing letters*, 20(3):209–212, 2012. 6, 14

[37] Nobuyuki Otsu. A threshold selection method from gray-level histograms. *IEEE transactions on systems, man, and cybernetics*, 9(1):62–66, 1979. 4

[38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proc. of International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 2, 3, 5, 6, 14

[39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Prompt engineering for ImageNet. https://github.com/openai/CLIP/blob/e58d49454c92986a1d2a6a48add2333bbfbeaf51/notebooks/Prompt_Engineering_for_ImageNet.ipynb, 2021. Accessed: 2022-11-16. 14

[40] Takahiro Saito, Takashi Komatsu, Tomohisa Hoshi, and Toshiaki Ohuchi. Image processing for restoration of old film sequences. In *Proceedings of the International Conference on Image Analysis and Processing*, pages 709–714. IEEE, 1999. 2

[41] Lorenzo Seidenari, Leonardo Galteri, Pietro Bongini, Marco Bertini, and Alberto Del Bimbo. Language based image quality assessment. In *ACM Multimedia Asia*, New York, NY, USA, 2022. Association for Computing Machinery. 7, 15

[42] Mehmet Sezgin and Bülent Sankur. Survey over image thresholding techniques and quantitative performance evaluation. *Journal of Electronic imaging*, 13(1):146–165, 2004. 4

[43] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 13

[44] Filippo Stanco, Dario Allegra, and Filippo Luigi Maria Milotta. Detection and correction of mistracking in digitalized analog video. In *Proc. of International Conference on Image Analysis and Processing*, pages 218–227. Springer, 2013. 1, 2

[45] Filippo Stanco, Dario Allegra, and Filippo Luigi Maria Milotta. Tracking error in digitized analog video: automatic detection and correction. *Multimedia Tools and Applications*, 75(7):3733–3746, 2016. 1, 2

[46] Maitreya Suin and AN Rajagopalan. Gated spatio-temporal attention-guided video deblurring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7802–7811, 2021. 2

[47] John WC Van Bogart. What can go wrong with magnetic media? *Publishing Research Quarterly*, 12(4):65–77, 1996. 1

[48] Ziyu Wan, Bo Zhang, Dongdong Chen, and Jing Liao. Bringing old films back to life. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17694–17703, 2022. 1, 2, 4, 6, 7

[49] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. Exploring CLIP for assessing the look and feel of images. *arXiv preprint arXiv:2207.12396*, 2022. 3

[50] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 6

[51] Jiyang Yu, Jingen Liu, Liefeng Bo, and Tao Mei. Memory-augmented non-local attention for video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17834–17843, 2022. 1, 2, 6, 7, 15

[52] 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. 6

[53] Weixia Zhang, Guangtao Zhai, Ying Wei, Xiaokang Yang, and Kede Ma. Blind image quality assessment via vision-language correspondence: A multitask learning perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14071–14081, 2023. 3

[54] Kun Zhou, Wenbo Li, Liying Lu, Xiaoguang Han, and Jiangbo Lu. Revisiting temporal alignment for video restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6053–6062, 2022. 1, 2

[55] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022. 14

# Reference-based Restoration of Digitized Analog Videotapes

## Supplementary Material

## 6. Synthetic Dataset

### 6.1. Dataset Generation

To generate our synthetic dataset, we start with *4K* videos belonging to the Harmonic dataset [2]. In particular, we employ videos related to the "Venice" scenes, which comprise 26,392 frames. We extract the $788 \times 576$ center crop from each frame. Then, we leverage Adobe After Effects [12] to add several types of degradation. We aim to make the artifacts in the synthetic frames as similar as possible to those found in real-world videos.

To reach our goal, we first reduce the saturation with the *Hue/Saturation Effect* to make the colors appear duller. Second, we reproduce the tape noise by adding Gaussian noise with the *Noise Effect*. Third, we replicate the tape dropout through an overlay with artifacts typical of VHS tapes blended into the frames in *lighten mode*. Then, we create *six* horizontal grids composed of black and white lines arranged in different ways. A *Wiggle Effect* is applied to the grids to change their vertical position over time without modifying their arrangement. In correspondence with the white lines of the grids, we add horizontal displacement artifacts with the *Displacement Map Effect* to replicate tape mistracking and VHS edge waving. Finally, we use the same approach to reproduce scanline flickering by leveraging the *CC Toner Effect* to blend horizontal cyan, magenta, and green lines into the frames.

At this point, our synthetic artifacts look similar to real ones. However, we also need to replicate the time-varying nature of real degradation. Indeed, artifacts in real-world videos change abruptly between consecutive frames and occur at the same time, leading to a disruption of temporal consistency. For this reason, we randomize all the effects we apply to the synthetic videos to make the degradations appear with different intensities, positions, and combinations for each frame. This way, we obtain a dataset of mainly temporally inconsistent videos composed of both almost clean and severely degraded frames, thus resembling real-world videos. Figure 7 shows an example of a high-quality frame and the corresponding synthetically degraded one belonging to our dataset.

### 6.2. Comparison with Real-world Videos

We provide a qualitative comparison between the synthetic and real-world videos in Figs. 8 and 9 to show their similarity. We focus on a static portion of a video, *i.e.* a patch in which both the content and the camera do not move significantly for multiple frames. In this way, we can assess how the appearance of the patch changes between clean and de-



(a) High-quality frame     (b) Synthetically degraded frame
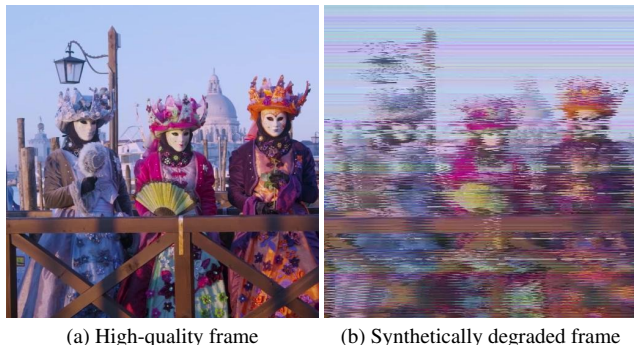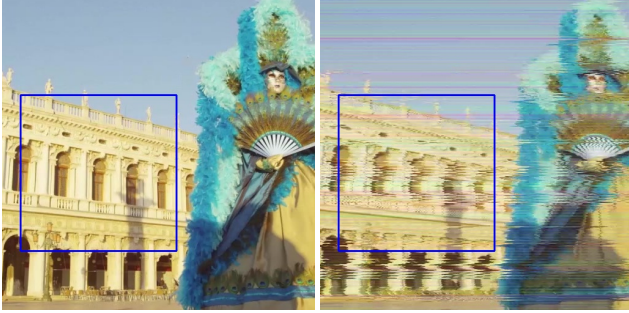
Figure 7. Example of a high-quality frame and the corresponding degraded frame belonging to the synthetic dataset.

graded frames and study the temporal consistency.

Figures 8a and 8b and Figures 9a and 9b illustrate a pair of nearly clean and severely degraded frames belonging to a synthetic and real-world video, respectively. We crop a static patch from each frame and show it in Figs. 8c and 8d and Figs. 9c and 9d. We observe that the horizontal colored lines and displacement artifacts in the synthetic frames closely resemble those found in the real-world ones. Moreover, we notice how both the balustrade in Fig. 8d and the branches in Fig. 9d are completely unrecognizable due to the severity of the degradation.
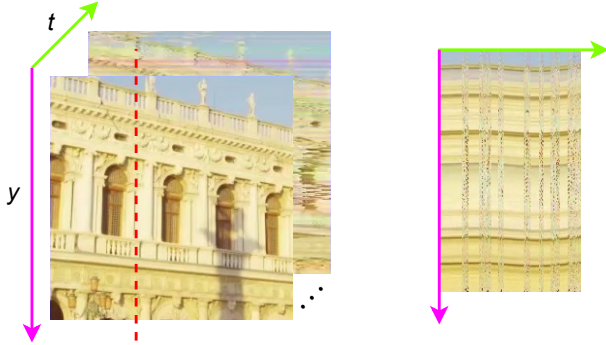
To study the temporal consistency, we select a column of pixels from the static patch of each frame and observe how it varies across time. Without artifacts, we would expect a (nearly) constant temporal profile, since the values of pixels belonging to a still patch do not change as the video progresses. Nevertheless, Figs. 8e and 9e show that in several sequences of frames in both the synthetic and real-world video the temporal consistency is completely lost. Indeed, we observe a smooth temporal transition only for frames in which the considered patches are clean. On the contrary, when the artifacts affect the crops under consideration, the temporal profile contains significant noise. However, the temporal profile also shows how the synthetic video closely resembles the real-world one, proving the quality of our dataset. Moreover, we observe that both the synthetic and real-world videos contain a relatively high quantity of clean frames. Therefore, we always have significantly more than $D$ frames to select the references from. Finally, the temporal profile justifies the main idea underlying our approach: identifying the cleanest frames of each video and exploiting them as references for restoring the sequences of severely degraded frames.

12

(a) Synthetic clean frame      (b) Synthetic degraded frame



(c) Synthetic clean patch      (d) Synthetic degraded patch



(e) Temporal profile of a column of pixels (red dashed line) belonging to a static synthetic patch. *y* and *t* represent the vertical and temporal axis, respectively.

Figure 8. Analysis of the degradation and the temporal consistency of a synthetic video.



(a) Real-world clean frame      (b) Real-world degraded frame



(c) Real-world clean patch      (d) Real-world degraded patch



(e) Temporal profile of a column of pixels (red dashed line) belonging to a static real-world patch. *y* and *t* represent the vertical and temporal axis, respectively.

Figure 9. Analysis of the degradation and the temporal consistency of a real-world video.

# 7. Training Loss

We train our model with a weighted sum of the Charbonnier loss [11] and a perceptual loss [13, 21, 23]. Let $I_R$ and $I_{GT}$ be a restored and ground-truth frame, respectively. To make the reconstructed frames faithfully approximate the ground truth ones we employ the Charbonnier loss, defined as:

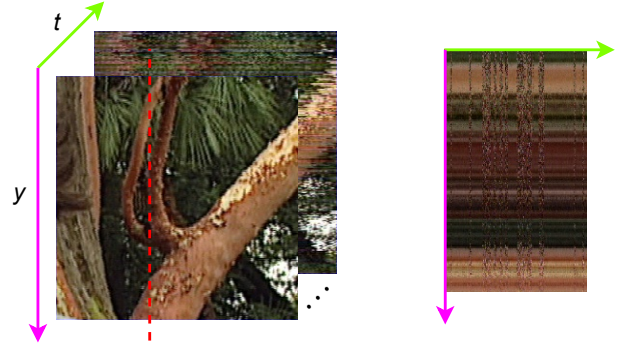$$\mathcal{L}_{char} = \sqrt{\|I_R - I_{GT}\|^2 + \epsilon^2} \qquad (6)$$

where $\epsilon$ is a constant equal to $10^{-12}$. To improve the perceptual quality and the photorealism of the results, we adopt the perceptual loss [13, 21, 23], defined in the VGG-19 [43]

feature space. The perceptual loss is formulated as follows:

$$\mathcal{L}_{perc} = \sum_{l \in L} \frac{1}{C_l H_l W_l} \|\Psi_l(I_R) - \Psi_l(I_{GT})\|^2 \qquad (7)$$

where $C_l$, $H_l$, $W_l$, and $\Psi_l$ represent the channel, height, width, and the features from the $l$-th layer of a pre-trained VGG-19 model, respectively, and $L = \{$*relu2_2*, *relu3_4*, *relu4_4*, *conv5_4*$\}$. Therefore, the overall training loss is:

$$\mathcal{L} = \lambda_{char}\mathcal{L}_{char} + \lambda_{perc}\mathcal{L}_{perc} \qquad (8)$$

where $\lambda_{char}$ and $\lambda_{perc}$ are the loss weights that we set to 200 and 1, respectively, making their values comparable during training.

13

# 8. Frame Classification

## 8.1. List of Prompts

As explained in Sec. 3.2 of the paper, we adopt prompt ensembling [38] to improve the frame classification results. The prompts we employ are the following: 1) "*an image with color artifacts along rows*"; 2) "*an image with interlacing artifacts*"; 3) "*an image of a noisy photo*"; 4) "*an image of a degraded photo*"; 5) "*a photo with distortions*"; 6) "*an image of a bad photo*"; 7) "*a jpeg corrupted image of a photo*"; 8) "*a pixelated image of a photo*"; 9) "*a blurry image of a photo*"; 10) "*a jpeg corrupted photo*"; 11) "*a pixelated photo*"; 12) "*a blurry photo*". We crafted prompts $\{1, \ldots, 5\}$ by converting the artifacts we observed in the real-world dataset into natural language. Prompts $\{6, \ldots, 12\}$ are more generic and derived from the templates employed for zero-shot classification on the ImageNet dataset by the authors of CLIP [39]. Simply updating the list of prompts makes our approach adapt to different types of degradation, *e.g.* to specialize our method for specific videos or types of medium. We find our list of prompts to be effective, but prompt learning [55] could be considered for future work.
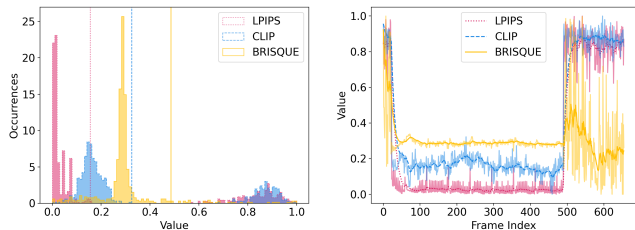
## 8.2. Analysis of Frame Classification

We evaluate the effectiveness of CLIP in the identification of the cleanest frames of a video by comparing it with several no-reference image quality assessment metrics: BRISQUE [35], NIQE [36] and CONTRIQUE [32]. Since we classify the frames into fairly clean and degraded ones, we can treat it as a binary classification task.

As explained in Sec. 3.2, we compute the histogram of the values of each metric and compute the threshold with Otsu's method for each video of the synthetic test set. Then, we classify as clean every frame with a value of the metric lower than the respective threshold. Additionally, we follow the same procedure for LPIPS, given that it is the most reliable metric to evaluate the perceptual quality of an image. Being a full-reference metric, we can leverage LPIPS only because the synthetic test set has a ground-truth counterpart. We consider the classification performed by LPIPS as the ground truth and the clean frames as positive examples. We measure the performance with standard metrics for binary classification: 1) Accuracy; 2) Precision; 3) Recall 4) F1 score. As can be seen in Tab. 5, CLIP outperforms all the no-reference metrics. NIQE tends to classify too many degraded frames as clean, achieving a high recall at the cost of lower precision. CLIP proves to be more balanced with the highest accuracy and F1 score.

In Figs. 10a and 10b we provide a visualization of the histograms and of the values of the metrics for each frame of a given video, respectively. We exclude NIQE and CONTRIQUE for clearer visualization, but the same

| Metric | Acc ↑ | P ↑ | R ↑ | F1 ↑ |
|---|---|---|---|---|
| BRISQUE [35] | 0.745 | 0.657 | 0.801 | 0.710 |
| NIQE [36] | 0.853 | 0.778 | **0.973** | 0.845 |
| CONTRIQUE [32] | 0.862 | 0.823 | 0.931 | 0.862 |
| CLIP [38] | **0.901** | **0.881** | 0.889 | **0.882** |

Table 4. Quantitative results for frame classification. ↑ means that higher values are better. Best results are highlighted in bold.



(a) Normalized histograms of the values of the metrics for a given video. The vertical lines represent the threshold values.

(b) Values of the metrics for a given video stream. The darker lines represent the smoothed values.

Figure 10. Visualization of the histograms and of the values of the metrics for a given video. Lower values are better.

considerations made for BRISQUE apply to them. Note that, for all the considered metrics, a lower value corresponds to a higher quality. Figure 10a shows how the histograms of LPIPS and CLIP are bimodal, while that of BRISQUE is unimodal. This illustrates how CLIP, contrary to BRISQUE, is capable of distinguishing the clean and degraded frames effectively. Looking at the values of LPIPS in Fig. 10b shows that the video presents degraded frames at the beginning, a long sequence of clean frames, and then other degraded frames. CLIP manages to capture the profile of the video stream and therefore correctly split the frames into the two classes. On the contrary, BRISQUE values are noisier and some of the degraded frames actually correspond to low values, thus leading to a wrong classification.

We also conduct an ablation study on the use prompt ensembling [38] for the frame classification with CLIP. We recall that prompt ensembling improves the robustness of the predictions by averaging the CLIP text features corresponding to multiple prompts. To evaluate the impact of prompt ensembling on the performance, we follow the procedure described in Sec. 3.2 but rely on a single prompt. In particular, we employ the prompt "*an image of a degraded photo*" (*i.e.* prompt 4 of the list reported in Sec. 8.1). Table 5 reports the results. Using a single prompt corresponds to the highest precision but, as expected [38], prompt ensembling achieves the best overall performance.

| Metric | Acc ↑ | P ↑ | R ↑ | F1 ↑ |
|---|---|---|---|---|
| Single prompt | 0.888 | **0.915** | 0.832 | 0.853 |
| Prompt ensembling | **0.901** | 0.881 | **0.889** | **0.882** |

Table 5. Evaluation of the effectiveness of prompt ensembling for the frame classification with CLIP. ↑ means that higher values are better. Best results are highlighted in bold.



(a) MANA: **43.77/4.90**/24.42    (b) TAPE: 51.49/6.55/**23.76**

(c) MANA: **40.70/4.72**/41.86    (d) TAPE: 41.27/6.08/**37.87**

Figure 11. Comparison between TAPE and MANA [51] on frames belonging to real-world videos. The reported values represent BRISQUE↓/NIQE↓/CONTRIQUE↓, respectively, where ↓ means that lower values are better. Best results for each pair of images are highlighted in bold.

## 9. Analysis of Quantitative Results

In Sec. 4.4 of the paper, we report the quantitative results for the real-world dataset. Our approach achieves the best results for CONTRIQUE, while MANA [51] is better for BRISQUE and NIQE. However, the qualitative results show that MANA adds high-frequency artifacts to restored frames. We argue that these artifacts deceive BRISQUE and NIQE, which mistake them for high-frequency details that are distinctive of high-quality images [4, 41]. In Fig. 11 we show two real-world examples supporting our argument. We report images obtained with TAPE and MANA and the corresponding BRISQUE, NIQE and CONTRIQUE values. Our approach generates results that are clearly more satisfying and photorealistic and achieves the best values for CONTRIQUE. However, MANA obtains the highest BRISQUE

and NIQE values, despite the presence of highly visible artifacts. We suppose that these artifacts are caused by the use of the memory bank learned during the training with our synthetic dataset which makes MANA fail to generalize to a different dataset. The lack of these artifacts in the r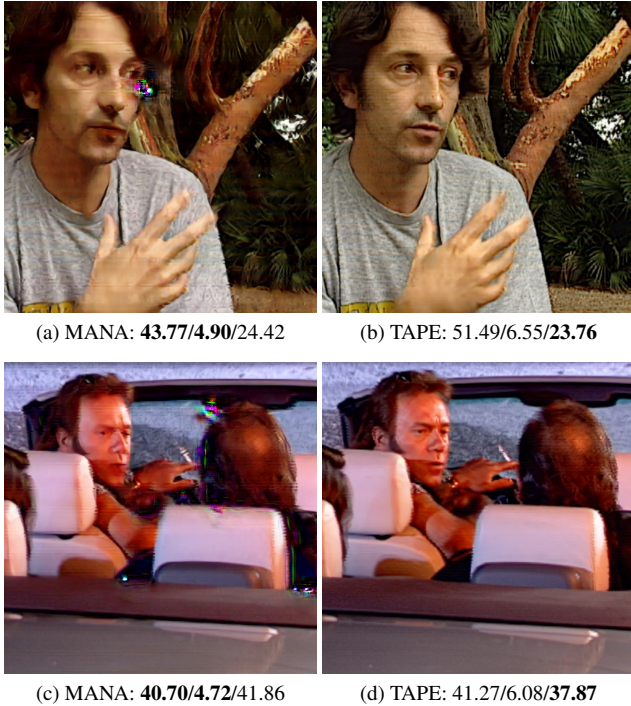esults of the synthetic test set supports our hypothesis.