



# A review of digital video tampering: From simple editing to full synthesis

Pamela Johnston<sup>\*</sup>, Eyad Elyan

Robert Gordon University, Garthdee House, Garthdee Road, Aberdeen, AB10 7QB, Scotland, UK

## ARTICLE INFO

### Article history:

Received 16 November 2018

Received in revised form

24 January 2019

Accepted 17 March 2019

Available online 22 March 2019

### Keywords:

Video tampering

Video synthesis

Deep learning

Video forgery

## ABSTRACT

Video tampering methods have witnessed considerable progress in recent years. This is partly due to the rapid development of advanced deep learning methods, and also due to the large volume of video footage that is now in the public domain. Historically, convincing video tampering has been too labour intensive to achieve on a large scale. However, recent developments in deep learning-based methods have made it possible not only to produce convincing forged video but also to fully synthesize video content. Such advancements provide new means to improve visual content itself, but at the same time, they raise new challenges for state-of-the-art tampering detection methods. Video tampering detection has been an active field of research for some time, with periodic reviews of the subject. However, little attention has been paid to video tampering techniques themselves. This paper provides an objective and in-depth examination of current techniques related to digital video manipulation. We thoroughly examine their development, and show how current evaluation techniques provide opportunities for the advancement of video tampering detection. A critical and extensive review of photo-realistic video synthesis is provided with emphasis on deep learning-based methods. Existing tampered video datasets are also qualitatively reviewed and critically discussed. Finally, conclusions are drawn upon an exhaustive and thorough review of tampering methods with discussions of future research directions aimed at improving detection methods.

© 2019 Elsevier Ltd. All rights reserved.

## 1. Introduction

area

The synthesis of convincing fake video content has increased recently due to the development of intelligent models (Le et al., 2017; Wang et al., 2018a; Bansal et al., 2018). Selective modification of image content has been possible for some years, but the application of similar techniques to video has been too labour intensive to see mass use. If each frame in a video is treated as an independent image, there are simply too many images to process efficiently. This has changed with increased computing power and the advent of deep neural networks. Deep learning techniques have seen great success in many applications recently. Generative Adversarial Networks (GANs) in particular have been used to alter source video: to re-enact human facial expressions (Suwajanakorn et al., 2017), change the weather (Liu et al., 2017) and to apply face-swapping (Dong et al., 1701). Human facial re-enactment is a relatively new but common area of research where a simple, talking

head is visually altered to mimic the facial expressions of a second actor (Suwajanakorn et al., 2017; Thies et al., 2016; Rössler et al., 2018) or to match a different audio track (Karras et al., 2017; Chen et al., 2018). This may have innocent applications such as re-dubbing a film in a different language or creating new movie scenes using old footage of an iconic actor, but it can also be used to produce convincing fake content. In some circumstances, fake content is convincing enough to reliably fool human eyes. The authors of (Rössler et al., 2018) even found that human viewers performed little better than random guessing when trying to ascertain whether facial re-enactment footage was authentic or synthesised. A deep neural network, however, could distinguish between the authentic and forged footage with ease.

Research into data-driven machine learning has also prompted the gathering of large image and video datasets such as ImageNet (Russakovsky et al., 2015), Youtube-8m (Abu-El-Haija et al., 2016) and CelebA (Liu et al., 2015). These datasets are a valuable resource for further research into convincing image and video forgery and in some cases (Rössler et al., 2018), a library of resources available for use in tampered datasets. The influence of these datasets has led to

relevance

<sup>\*</sup> Corresponding author.

E-mail address: [p.a.johnston3@rgu.ac.uk](mailto:p.a.johnston3@rgu.ac.uk) (P. Johnston).

an increase in the application of deep neural networks to tampering. Spatially localised changes in video footage, such as face swapping, can change the entire context of a news story or film and can have repercussions for the people portrayed. Already, videos which have been cleverly edited to change the context of what was said by influential people have gone viral<sup>1</sup>. If that can be done with unsophisticated editing techniques, it is worth considering what more could be achieved with modern techniques.

There are already a number of recent surveys which review tampering detection methods (Sitara and Mehtre, 2016; Singh and Aggarwal, 2017; Milani et al., 2012; Pandey et al., 2016; Al-Sanjary and Sulong, 2015). Tampering detection methods are broadly categorised as active or passive, with more focus on passive tampering detection methods. A review of passive tampering detection in video is provided in (Milani et al., 2012; Al-Sanjary and Sulong, 2015) and, more recently (Sitara and Mehtre, 2016). Inter-frame tampering detection is covered in (Singh and Aggarwal, 2017). Pandey et al. (2016) cover tampering detection through noise in images as well as video. There are, however, far fewer reviews on tampering itself. This paper aims to help redress the balance. The work in (Khodabakhsh et al., 2018) provides an overview of personation, specifically how a person's likeness in appearance and voice can be forged in videos either physically or digitally. However it is important to objectively catalogue current known techniques that can be used for video tampering in order that they can be identified and, ultimately, detected or countered. Many detection techniques are explicitly tailored to specific tampering methods. For example, the authors of (Rössler et al., 2018) trained a deep neural network to detect their own video content changes in order to assess the quality of their content-altering techniques (Wu et al., 2014; Sitara and Mehtre, 2017); focused on inter-frame tampering (Lin and Tsay, 2014); created tampered sequences using established inpainting techniques (Patwardhan et al., 2007; Criminisi et al., 2004) to assess their detector. All of these techniques worked well, but all of them required prior knowledge of the type of tampering. As tampering methods multiply, it becomes important to fully assess new detection methods, and to appreciate which types of tampering techniques they can feasibly detect and which they are blind to. This review exposes new research directions by cataloguing known tampering algorithms to aid development of automated, universal detection techniques.

Wang and Farid (2007) noted that, at the time of their 2007 publication, there were very few video tampering detection techniques. This is no longer the case, however, many published techniques, specific to particular types of tampering or source authentication, were assessed on proprietary datasets which remain unreleased (Sitara and Mehtre, 2017; Kobayashi et al., 2009; Chuang et al., 2011; Subramanyam and Emmanuel, 2012; Kaur et al., 2017). In some cases (Aghamaleki and Behrad, 2017), datasets are detailed sufficiently in the literature so that they can be replicated, provided the sequences used for dataset synthesis are available. This serves to evidence the fragmentation of the tampering detection field. In reality, a tampered video may be subject to a variety of techniques, including combinations. For tampering detection to be effective, individual detectors must be analysed and matched with an appropriate type of tampering. In order for that to happen, we must review and differentiate types of tampering.

Here, we catalogue and analyse the current trends in digital video manipulation techniques from simple edits to fully synthetic

video. This allows for work towards a method of universal video tampering detection. We list available tampered datasets and identify their potential challenges. We thoroughly examine problems in dataset gathering and dissemination, including challenges created by compression.

This review opens up a new field of research in the form of tampering classification. If video tampering is, by its nature, designed to be invisible to human eyes, then tampering classification will necessarily be done by algorithms and machines. As a first step in this process, we analyse the current classes of video tampering. Future tampering detection methods may maximise their impact by identifying and targeting tampering classes instead of individual algorithms. The purpose of this paper is to provide an in-depth analysis and review of existing methods of digital video tampering, to understand the current state-of-the-art, how it may progress in future and how this can be used to inform future development of tampering detection techniques. Our contributions are as follows:

- A thorough examination of how video tampering techniques have been categorised in the past and how these have influenced development of tampering detection.
- An in depth, original review of how video tampering has evolved in recent times with discussion of the latest deep learning techniques.
- A qualitative evaluation of existing tampered video datasets. This includes critical analysis of large, tampered image datasets, which is used to assess the challenges associated with creating and distributing video datasets and propose effective methods for overcoming these.

The paper is structured as follows: Section 2 critically analyses traditional types of video tampering and how these have been augmented by the latest developments to form a spectrum of video tampering from fully authentic video through to fully synthesised video; Section 3 examines some state-of-the-art tampering/anti-forensic detection methods and discusses how these reveal underlying information about datasets used for training; Section 4 examines existing tampered video datasets and highlights good lessons that can inform future dataset compilation; Section 5 concludes the paper and gives new research directions.

## Overview

In (Sitara and Mehtre, 2016) video tampering was defined as “a process of malicious alteration of video content, so as to conceal an object, an event or change the meaning conveyed by the imagery in the video”. Similarly (Qureshi and Deriche, 2015), described image forgery as “the digital manipulation of pictures with the aim of distorting some information in these images”. In this paper, video tampering is regarded as any technique which is intended to produce manipulated, photo-realistic content using authentic sources. There is no defined limit as to when authentic video becomes tampered video. Similarly, malicious intent is difficult to quantify, and so tampering detection and video authentication techniques must focus on forensic analysis, providing objective localisation of inconsistencies within digital footage that may imply content alteration.

It is important to note that this paper examines only digital video tampering: video content can also be “staged” whereby the video file is an authentic record of events, but the events themselves were contrived or unnatural during filming. Detection of staged video involving natural ballistic trajectories is examined in (Conotter et al., 2012), and (Khodabakhsh et al., 2018) details how plausible audiovisual personation is achieved in front of the

<sup>1</sup> “Video of Barack Obama speech circulating on the Internet was edited to change his meaning”: <https://www.politifact.com/truth-o-meter/statements/2014/jun/23/chain-email/video-barack-obama-speech-circulating-internet-was/> Accessed 2019-1-24.

practical value

theoretical value

contribution

Plan of the paper

camera. Digital video forgery can take a number of forms (Sitara and Mehtre, 2016) and Fig. 1 gives an overview of the classical interpretation.

In the past, video tampering methods have been simply classified as inter- or intra-frame (Sitara and Mehtre, 2016; Singh and Aggarwal, 2017) (Fig. 1). The terms inter- and intra-frame primarily distinguish temporal tampering from spatial tampering. Inter-frame tampering is performed on a sequence-level: the pixels of individual frames are unaltered, but the sequence as a whole is changed. Intra-frame tampering is performed on a pixel-level: some spatial regions are altered, but alterations temporally correlated to form a convincing forged region. The term “inter-video tampering” can also be used to describe the merging of content from two different videos (Ardizzone and Mazzola, 2015). Traditionally, this has been a form of splicing, where chroma-keyed objects taken from one sequence are inserted into another, as in (D’Avino et al., 2017). Recent developments, however, mean that convincing synthetic regions (Suwajanakorn et al., 2017; Dong et al., 1701; Chan et al., 1808) or even whole videos (Wang et al., 2018a) can be synthesised automatically from authentic content. This development means that we must now consider different levels and categories of video tampering. It is important to be aware of the different categories because video tampering is designed to be invisible to human eyes, and detection techniques often address only one type of tampering.

#### The spectrum of video content

The current field of video tampering may be viewed as a spectrum, as in Fig. 2, where different types of video tampering are ordered according to potential to deviate from authentic source. Whereas the traditional view in Fig. 1 provides only two categories of video tampering, the spectrum in Fig. 2 demonstrates that there are now multiple ways to produce convincing, falsified content. This distinction is important because detection methods often address one particular type of video tampering such as object forgery or inter-frame tampering. In (Sitara and Mehtre, 2016) tampering detection methods are categorised as recompression, inter-frame forgery or region tampering detection. With current tampering techniques, the distinction is less clear cut. Moreover, multiple tampering techniques can be applied to the same sequence.

Fig. 2 summarises the current categories of video tampering.

Video editing compiles single camera shots into full films complete with scene cuts. Although clever editing may change the context of a video, scene cuts are not usually deliberately concealed. Video clips of maliciously edited content exist in mainstream media and are surprisingly effective at disseminating misinformation through social media. Traditional inter-frame tampering, where edits are concealed, may reorder events or even remove or insert events into the timeline, but its content-altering effects are self-limiting. Retouching temporally or spatially upscaled content, or applying global filters to improve perceptual quality may affect every pixel in a video sequence and can cosmetically alter content. Retouching can also be applied to specific regions. Intra-frame tampering and other object forgeries such as inpainting can alter content and context, as can motion transfer. Finally, fully synthetic video or synthetic regions can be produced. Unlike historical animations, the synthetic content of today looks convincingly realistic. The following subsections 2.2 to 2.6 detail examples from each of these types of video manipulation.

Table 1 shows how motion transfer and video synthesis techniques have become common in recent years and demonstrates how methods of evaluation remain relatively underdeveloped. Evaluation techniques are difficult to define since there is no pre-defined ground truth for tampered video data. Every new method can be assessed qualitatively. Methods which seek to imitate authentic video, such as frame interpolation, can use full reference quality metrics such as SSIM and PSNR. As can be seen in Table 1, video manipulation methods use user studies to evaluate their output or simply publish examples of their methods for future evaluation. However, even user studies can vary. Some ask users to classify frames as tampered or authentic. Some request a user preference between the published method and other, similar methods. In a related field, image inpainting evaluation techniques are reviewed in (Qureshi et al., 2017) and these can all be used to assess the spatial features of inpainted video or indeed, any form of tampering which affects individual frames. No-reference video quality assessment is a large and open field and although we do not cover this here, we point to this field to at least partially inform on tampered video evaluation.

#### Editing and inter-frame tampering

Editing and inter-frame tampering both change the order of the frames in the video without changing the contents of each frame. In

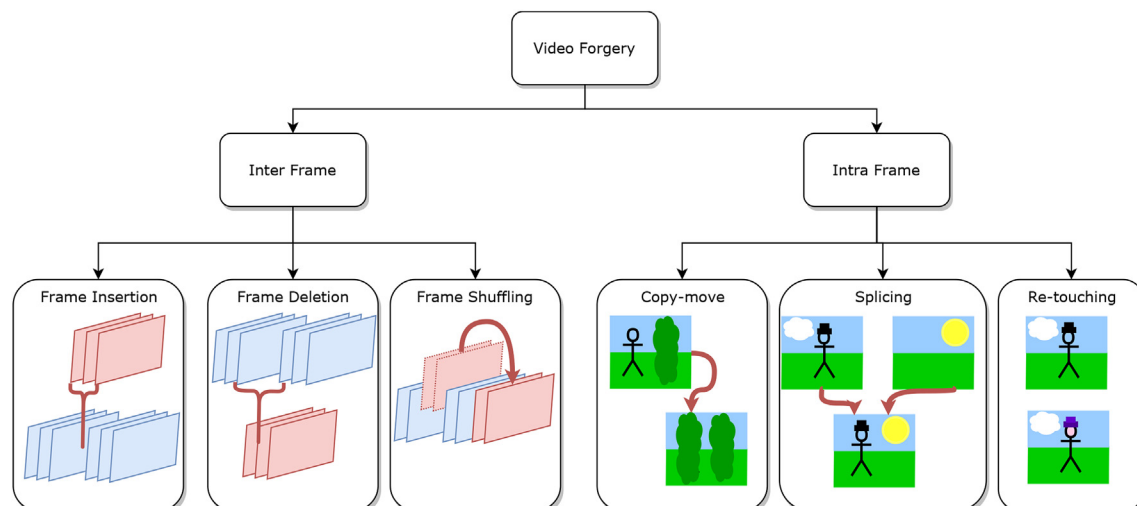


Fig. 1. Traditional video forgery categories.

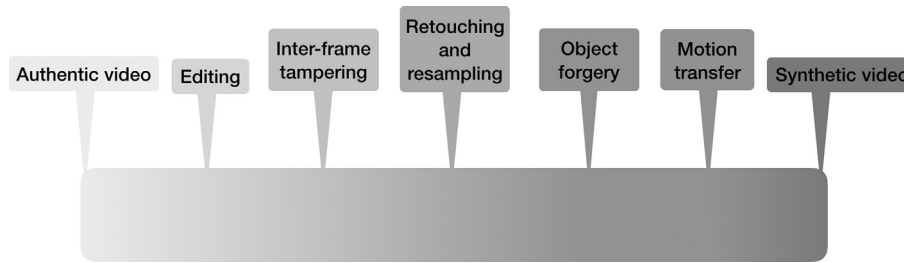


Fig. 2. Video tampering spectrum.

**Table 1**  
Video Tampering and Evaluation Methods: Qual = qualitative analysis; PSNR=Peak Signal to Noise Ratio; SSIM=Structural SIMilarity; UP=User preference to previous methods; UR=User comparison with real video; Rel = Released Sequences.

Reference	Year	Type of Tampering	Qual	PSNR/SSIM	UP	UR	Rel	Other
ETS (Criminisi et al., 2004)	2004	Inpainting	✓					
Ha et al. (Ha et al., 2004)	2004	Frame interpolation	✓	PSNR				
Patwardhan et al. (Patwardhan et al., 2007)	2007	Inpainting via temporal copy-move	✓					
Wexler et al. (Wexler et al., 2007)	2007	Inpainting, frame interpolation	✓				✓	
Shih et al. (Shih et al., 2011)	2011	Object forgery	✓				✓	
SULFA forged (Qadir et al., 2012)	2012	Object forgery	✓				✓	detection
SULFA supplemental (Bestagini et al., 2013)	2013	Object forgery	✓				✓	detection
Newson et al. (Newson et al., 2014)	2014	Inpainting	✓				✓	
Ardizzone and Mazzola (Ardizzone and Mazzola, 2015)	2015	Copy-move	✓				✓	
Ebdelli et al. (Ebdelli et al., 2015)	2015	Inpainting	✓	PSNR			✓	
Lotter et al. (Lotter et al., 2015)	2015	Frame prediction	✓					error
Dar and Bruckstein (Dar and Bruckstein, 2015)	2015	Frame interpolation	✓	PSNR				
Face2Face (Thies et al., 2016)	2016	Motion transfer	✓				✓	
Le et al. (Le et al., 2017)	2017	Inpainting	✓				✓	
Suwajanakorn et al. (Suwajanakorn et al., 2017)	2017	Motion transfer	✓					
Liu et al. (Liu et al., 2017)	2017	Style transfer	✓				✓	
Niklaus et al. (Niklaus et al., 2017)	2017	Frame interpolation	✓	PSNR				
MoCoGAN (Tulyakov et al., 2018)	2017	Motion transfer	✓			✓		ACD
Walker et al. (Walker et al., 2017)	2017	Frame prediction	✓					Inception
FaceForensics (Rössler et al., 2018)	2018	Motion transfer	✓			✓	✓	detection
Recycle-GAN (Bansal et al., 2018)	2018	Video synthesis	✓		✓			✓
Wang et al. (Wang et al., 2018a)	2018	Video synthesis (sketch)	✓		✓			
Chan et al. (Chan et al., 2018)	2018	Motion transfer	✓	SSIM				LPIPS
Jiang et al. (Jin et al., 2018)	2018	Video synthesis (blurred image)	✓	PSNR				
Wang et al. (Wang et al., 2018b)	2018	Video synthesis (smile)	✓		✓			
Xiong et al. (Xiong et al., 2018)	2018	Video synthesis (time-lapse)	✓		✓	✓		
Babaeizadeh et al. (Babaeizadeh et al., 2018)	2018	Frame prediction	✓	✓				
Zhao et al. (Zhao et al., 2018)	2018	Frame prediction	✓	PSNR		✓		ACD
SCGAN (Yang et al., 2018a)	2018	Video synthesis (human pose)	✓		✓			pose eval.
SDC-Net (Reda et al., 2018)	2018	Frame prediction	✓	✓				
Cai et al. (Cai et al., 2018)	2018	Frame prediction/interpolation	✓	✓				Inception

the case of editing, the goal is to turn a series of single camera shots into a coherent story. Clever edits can be used to turn innocent footage into propaganda,<sup>2</sup> but scene cuts are not hidden and such videos are not above suspicion. In inter-frame tampering, the goal is to *invisibly* remove, re-order or alter events.

Edits in inter-frame tampering are deliberately concealed so as to be invisible to the human eye. Detection of visible scene cuts in video has been studied extensively so that key frames can be identified for efficient compression and or used to condense/index the sequence (Cotsaces et al., 2006). Invisible scene cuts are studied in the context of inter-frame tampering detection (Huang et al., 2017; Zheng et al., 2014; Smith et al., 2017).

Such is the theoretical simplicity of generating an inter-frame tampered sequence, that many tampering detection methods,

such as (Wu et al., 2014; Sitara and Mehtre, 2017; Tralic et al., 2014; Stamm et al., 2012) generate their own datasets from single-camera video sequences such as SULFA (Qadir et al., 2012) or Derf's media collection (Xiph.org video test media) or even film their own sequences as in (Kaur et al., 2017). SULFA (Qadir et al., 2012) replicates single camera sequences as obtained from CCTV footage, and, therefore may be representative of the most likely application of inter-frame tampering: altering CCTV evidence. Derf's media collection (Xiph.org video test media), on the other hand, provides publicly available uncompressed sequences and allows researchers complete control over the forensic history of synthesised tampered sequences (Stamm et al., 2012).

It remains unclear how widespread inter-frame tampering is in the wild because, if it is done correctly, it will be undetectable by human eyes and remain above suspicion. Meanwhile, it is important that synthesised datasets are as high quality as possible. In creation of inter-frame tampered datasets (Kaur et al., 2017; Aghamaleki and Behrad, 2017; Thakur et al., 2016), simply removed pre-determined frame numbers from each sequence, and it is unclear if this caused visible effects. In (Stamm et al., 2012) frame

<sup>2</sup> "Israeli army edits video of Palestinian medic its troops shot dead to misleadingly show she was 'human shield' for Hamas", The Independent, <https://www.independent.co.uk/news/world/middle-east/gaza-protests-latest-idf-condemned-edited-video-angel-of-mercy-medic-razan-al-najjar-a8389611.html>.



addition and removal was limited to the beginning of each sequence, but again it is unclear if the additions were visible: simply reversing the sequence from the point of tampering may effectively locally conceal the edit. Recent developments in video quality assessment mean that temporal glitches in video can be objectively quantified (He et al., 2017) and also smoothed (Lai et al., 2018) to achieve temporal consistency. Future datasets for inter-frame tampering can use this to improve such that inter-frame tampering techniques can be deployed in the wild.

As noted in the review in (Sitara and Mehtre, 2016), many inter-frame tampering detection methods suffer from limitations which are often related to consistencies within the dataset which may not translate to other video data. These consistencies are often related to video compression. The authors of (Joshi and Jain, 2015) note that some tampering detection techniques are tied into the fixed Group of Pictures (GOP) size, commonly used in MPEG-2 (ITU-T, 2012) to minimise error accumulation due to non-integer frequency domain transforms. Later video compression standards, such as H.264/AVC (ITU-T, 2016a), use integer-based transforms so error accumulation drift between encoder and decoder is no longer an issue, and therefore key frames are used only as access points into the stream or efficient compression of cut scenes. Moreover, sequences compressed using H.264/AVC no longer exhibit visible evidence of key frames.

Although inter-frame tampering detection is widely studied in the literature, effects similar to inter-frame tampering can be achieved using a spatio-temporal copy-move. Rather than replacing complete frames in the sequence, only partial frames containing motion or objects to be concealed are replaced. With a static camera and consistent lighting, this is visually effective, and video edits prove near invisible to the naked eye. Indeed, some sequences which initially look like inter-frame tampering (Bestagini et al., 2013) are actually spatio-temporal copy-moves, as can be revealed by examining pixel-by-pixel difference between the authentic and tampered sequences (see Fig. 3d).

### Retouching and resampling

Retouching involves adjusting pixels within an image using transforms or filters applied to the pixels themselves which may only have a low-level interpretation of video content. As the name suggests, retouching is less invasive to content than other types of forgery but may still change the context of a video. Moreover, retouching can be used on tampered video specifically as an anti-forensic device.

A retouching function  $R$  can affect specific pixels according to a binary mask,  $M$ :

$$V_{retouched} = R(M \odot V_{original}) + ((I - M) \odot V_{original}) \quad (1)$$

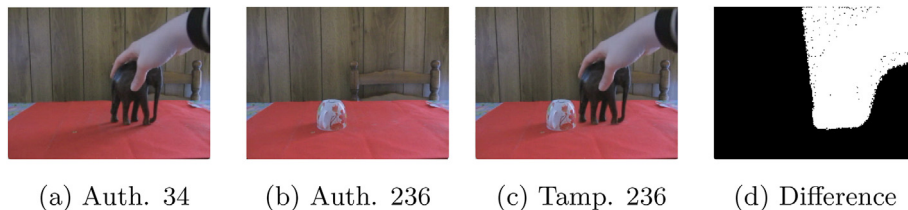
Here,  $I$  represents a matrix of ones and all matrices have equal size. Retouching can also be applied globally as in:

$$V_{retouched} = R(V_{original}) \quad (2)$$

Colour correction methods such as those available in Adobe After Effects may normalise lighting in a sequence of shots taken on different days under different weather conditions to create a convincing narrative. Similarly, colour grading can also be used to add effects or make video filmed during daylight hours appear to have been filmed during twilight. A typical colour correction model works by adjusting the histogram of colour over a specified region, however the authors of (Choi et al., 2017) found that gamma correction (a form of colour correction) was particularly difficult to detect using a deep neural network. Where median filtering and Gaussian blurring could be detected reliably with over 91% accuracy, detection of gamma correction was only 57.6%.

Compression is often a necessary part of video processing but it can also be used as an anti-forensic method. Video compression standards such as (ITU-T, 2012, 2016a) reduce video file size with no explicit understanding of video content. Compression has been found to reduce the efficacy of tampering detectors (Rössler et al., 2018; Bestagini et al., 2013; Thakur et al., 2016; Chen et al., 2016) and has also been shown to reduce the classification accuracy of convolutional neural network (CNN) based classifiers (Johnston et al., 2018). The authors of (Rössler et al., 2018) found that video compression (ITU-T, 2016a) reduced the accuracy of deep neural networks trained to detect human facial re-enactment. Of seven different forgery detectors tested, the Xception network (Chollet, 2017) was the most robust against compression, achieving 87.81% accuracy, compared with 99.93% accuracy on uncompressed sequences. Other methods (Bayar and Stamm, 2016) performed less well for forgery detectors on the compressed dataset with performance for some (Cozzolino et al., 2017) dropping as low as 55.77%. This may be attributed to the depth of the Xception network. The authors of (Chen et al., 2016) also account for compression in their SYSU-OBJFORG dataset. In benchmarking it using seven common steganographic features, they, too, found a drop in accuracy on the reduced bitrate video. While their ensemble-based detector achieved precisions in the range 78.9–93.15%, this reduced to 61.85–79.34% when bitrate of the video data was halved. Halving height and width of the video sequences also reduced precision to the range 73.02–90.28%, which was not as significant as bitrate reduction. In detection of frame deletion, it was found in (Thakur et al., 2016) that an SVM conditioned on uncompressed data to detect dropped frames did not perform well on compressed data taken from YouTube, with accuracy dropping below 37%. It is clear from this that standard video compression can reduce some of the features associated with tampering. This is most likely down to the way video compression quantises data in the frequency domain.

Compression artifact removal is another example of retouching and methods such as (Dong et al., 2015; Cavigelli et al., 2017; Guo and Chao, 2017) have been applied to JPEG images. The authors of (Kirmemis et al., 2018) used a deep residual network to reduce artifacts in a BPG (ITU-T, 2016b) compressed frame. More recently



**Fig. 3.** An intra-frame tampering example from (Bestagini et al., 2013). 3a and 3b show authentic content. 3c shows the spatio-temporal copy-move and 3d shows the difference between 3b and 3c.

compression artifacts have been removed in the video domain (Yang et al., 2018b), where videos compressed using HEVC and H.264/AVC were retouched to increase Peak Signal to Noise Ratio (PSNR). PSNR is defined as:

$$PSNR = 20 \log_{10} \left( \frac{255}{\sqrt{MSE}} \right) \quad (3)$$

where MSE is Mean Squared Error, and it is a very common full reference quality metric where processed pixels are compared directly to unprocessed pixels. It is useful for measuring pixel fidelity and is often used alongside the full reference quality metric Structural Similarity (SSIM) (Wang et al., 2004), but may not directly reflect perceptual quality. Although (Yang et al., 2018b) achieved overall improvement in PSNR, it was unclear if this resulted in a gain in perceptual quality at specific bit rates. Given that some tampering detection methods, such as (Gironi et al., 2014), actively utilise compression structures, methods which alter the underlying patterns of compression in video frames could be used as anti-forensics in the future.

Artificially upscaling video frame size (Sajjadi et al., 2018), frame rate (Jiang et al., 2018; Xia et al., 2017; Li et al., 2018) or bitrate can be a form of video tampering. High quality video content is more desirable to consumers, and larger file sizes for the same film/footage are often indicative of higher quality, with bitrate often taking the place of quality in common parlance. Compression encoders utilise a specified bitrate, even if this means compressing existing compression artifacts. Bitrate upscaling can be done innocently as researchers seek to provide high quality, “uncompressed” datasets and either overlook or deliberately replicate compression artifacts in the pixels of mined data.

Artificially increasing the spatial dimensions of video has been commonly done historically as Standard Definition (SD) content is displayed on High Definition (HD) screens. More recently, super-resolution has evolved from an image enhancement technique to use within videos (Sajjadi et al., 2018; Caballero et al., 2017; Liao et al., 2015), and the metrics commonly used for evaluation are, again, PSNR and SSIM: full reference quality metrics. It is important for spatially upscaled video to demonstrate temporal coherence. The authors of (Sajjadi et al., 2018; Caballero et al., 2017), also assessed the temporal coherence of their super-resolution sequences using a technique called “temporal profile”. This is where single rows of pixels are viewed along with their temporal neighbours from different frames, and temporal inconsistencies or video “flicker” shows up as hard edges in the resultant image. The work in (He et al., 2017) is also a method of assessing temporal consistency.

Video deblurring (Kupyn et al., 2018) is another example of retouching and a dataset exists to facilitate the development of this (Nah et al., 2017). The dataset was filmed using a GoPro camera at 240fps and then downsampled and blurred so that a non-blurred ground truth can be supplied for each blurred frame, thus enabling deblurrers to be assessed using full reference quality measures such as PSNR. Super slow motion has recently become a strong field of research with many new techniques for temporally upsampling video (Niklaus et al., 2017; Jiang et al., 2018; Meyer et al., 2018) to create a slow motion effect in the absence of a high speed camera. Previously, upsampled video would simply involve frame repetition or averaging. The field of motion compensated interpolation improved upon this (Ha et al., 2004; Dar and Bruckstein, 2015; Niklaus et al., 2017) so that interpolated frames were less obvious. The work in (Ha et al., 2004) is an early example of motion compensated interpolation, and the authors used block-based motion estimation similar to that used in video compression (ITU-T, 2012) to inform interpolation and create

sophisticated intermediate frames. The frame rate upconversion algorithm was objectively assessed by downscaling some publicly available uncompressed sequences and then comparing the original sequence with the computed upscaled version using PSNR. In (Dar and Bruckstein, 2015), the authors showed how their work in frame rate upscaling could be used to improve low bitrate video compression. In (Niklaus et al., 2017), a CNN was used to interpolate between frames. The authors obtained their training data from high quality YouTube channels and downsampled from 1080p to 720p in order to reduce the effects of compression. A user study confirmed that their interpolated frames were better than previous state-of-the-art. Objective assessment of results used sampled alternate frames from a popular YouTube video and used PSNR to compare interpolated frames with actual frames. In (Jiang et al., 2018), multiple frames were synthesised between two authentic frames using a CNN trained on high frame-rate (720p, 240fps) video from YouTube and a high frame-rate dataset (Janai et al., 2017). The synthesised frames were assessed using high frame-rate video and it was found that PSNR between interpolated frames and ground truth was improved upon previous state-of-the-art. As noted in (Ilan and Shamir, 2015), inpainting or video completion (Section 2.4) can be used to resample a video, and entire frames inpainted.

Retouching might be one step of many in tampering, and although it does not necessarily alter context, it can be used to make tampering detection much more difficult. Countermeasures for anti-forensics are well studied in the literature (Choi et al., 2017; Bayar and Stamm, 2016; Chen et al., 2015), and datasets can be generated easily. In (Bayar and Stamm, 2016), a CNN was used to classify an image in terms of its anti-forensic processing. The labels used were: original (no processing), Gaussian blurring, additive white noise, median filtering and resampling. The CNN accurately detected the presence of each process over with over 98% accuracy using only the green colour channel. A new type of convolutional layer was designed to prevent the network from learning typical image features. The authors of (Chen et al., 2015) showed how their median filter detector could be used to localise median filtering within a spliced image and hence localise image tampering. Although retouching does not always correlate with tampering, localised retouching can be a strong indicator of splicing or other object forgery.

### Intra-frame tampering

Intra-frame tampering is where spatial content of individual frames is changed, that is, individual objects are added or concealed/removed. Intra-frame tampering is also known as “region tampering” (Al-Sanjary and Sulong, 2015) and applies equally to video and still images, although the video application is more complex. Care must be taken to ensure that spatial tampering across individual frames is coherent and does not cause visual jarring in the video. Intra-frame tampering methods in images were classified as spatio-temporal copy-move, splicing and retouching in (Chauhan et al., 2016), but here we discuss retouching separately (Section 2.3).

A spatio-temporal copy move can be defined by:

$$V_t^{Lj} = \left( (I - M) \odot V_o^{Lj} \right) + \left( M \odot V_o^{Lk} \right) \quad (4)$$

where  $I$  is the matrix of ones,  $M$  is a binary mask to localise tampering,  $V_o^{Lj}$  is an authentic sequence of  $L$  frames starting on frame  $j$ ,  $V_o^{Lk}$  is the same sequence but starting on frame  $k$  where  $j \neq k$ . The frames/mask can be re-aligned or cropped so that any object or region of pixels from any spatial or temporal location can be copied to any location.  $V_t^{Lj}$  is a tampered video sequence:

$$V_t^{Lj} = [v^j, \dots, v^{j+L-1}] \quad (5)$$

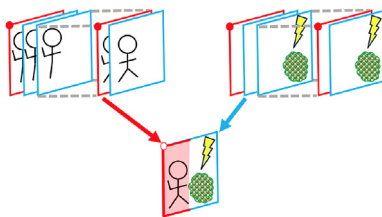
In spatio-temporal copy-move attacks, all the data used in the video forgery  $V_t$  comes from within the same video sequence  $V_o$ . For example, complete objects from frame  $k$  in the sequence are inserted into frame  $j$  using mask  $M$ . Fig. 3 shows an example. This is similar to image-based copy-move where the pixels involved in the tampered region come from within the image itself. Although this reduces the range of potential content, it helps to minimise differences between legitimate and tampered regions. There is less need to alter the colour histogram or adjust the frame rate to make tampered content consistent with authentic content if both share the same source. Using a copy-move attack, objects can be added to a sequence by adding foreground objects, or concealed/removed from a sequence by duplicating background regions from within the same frame or from within a different frame in the same sequence.

Some versions of copy-move attacks simply duplicate a still background region (Qadir et al., 2012), and these can be detected with relative ease by high coherence or abnormally low motion within the tampered region (Bestagini et al., 2013). Other methods (Bestagini et al., 2013) duplicate an entire spatio-temporal region, and this is more difficult to detect. Although duplicates can be detected by matching copied region to original data, this becomes more difficult in the presence of compression (Bestagini et al., 2013). A copy-move attack can be detected in images by identifying and locating duplicated regions, and this has been done using search based on brute-force pixel matching, region matching or key point matching (Chauhan et al., 2016; Li et al., 2015). While this type of copy-move detection is feasible in images, video adds another dimension and searches become an order of magnitude more complex. Previous video inpainting attempts such as Temporal Copy-Paste (TCP), where identical pixels are used frame after frame to conceal an object within a video (Qadir et al., 2012), or Exemplar-based Texture Synthesis (ETS) (Criminisi et al., 2004) were detected by (Lin and Tsay, 2014) where detection was based on correlation between adjacent frames which was either too strong or not strong enough to be authentic video.

Splicing is an extension of spatio-temporal copy-move. In a splicing attack, two sets of pixels from different sources are combined as in Fig. 4. The sources can be videos or even still images (as shown in Fig. 5). Eq. (6) defines splicing:

$$V_t^{Lj} = ((I - M) \odot V_{s1}^{Lj}) + P(M \odot V_{s2}^{Lk}) \quad (6)$$

Where video sequences are defined as in Eqs. (4) and (5),  $s1$  means sequence 1,  $s2$  means sequence 2 and  $P$  is an optional processing step which can be applied to aid blending between different source videos. The frames/masks of the two sources can be re-



**Fig. 4.** Forging a video (best viewed in colour), Red borders/dots indicate key frames in the sequence. Blue borders (no dots) indicate predicted frames. The hybrid frame is shaded red where the pixels have come from a key frame, and unshaded where the donor frame was a predicted frame. (For interpretation of the references to color in this figure legend, the reader is referred to the Web version of this article.)

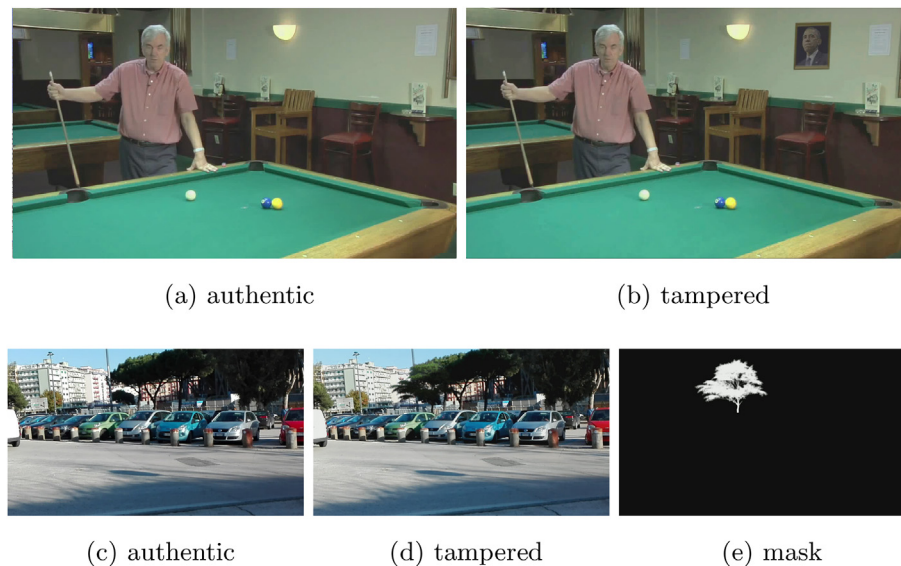
aligned or cropped so that any object or region of pixels from any location from source 2 can be pasted into any location in source 1.

Splicing has large potential for context-changing edits because two entirely different subjects can be spliced together. Any source sequences involved can be retouched (see Section 2.3) using colour correction or temporal synchronisation before or after a video splice in order to visually camouflage spliced content, or even to launder the splicing operation to make it undetectable to existing forensic tools.

Copy-move and splicing are also known as “object forgery” (Chen et al., 2016) because they often involve removing or adding complete objects to videos. Introduction of an object to a video can be done using chroma-keying techniques, as in the field of video special effects (Aksoy et al., 2016). Chroma-keying requires filming against a single colour background under specific lighting conditions to facilitate segmentation of foreground objects. An example of object forgery using chroma-keyed sources is shown in Fig. 5d. Other segmentation methods such as (Bideau and Learned-Miller, 2016; Xu and CorsoLibsvx, 2016) may be used in place of chroma-keying so that foreground objects can come from any sequence without the need for special green-screen filming. The authors of (Bideau and Learned-Miller, 2016) applied segmentation to the optical flow of videos in order to distinguish foreground and background objects in sequences with moving cameras. In (Xu and CorsoLibsvx, 2016), segmentation was achieved using supervoxels, using spatiotemporal uniformity in pixels to group them into voxels and supervoxels to represent different objects in the sequence. Masks produced by (Bideau and Learned-Miller, 2016; Xu and CorsoLibsvx, 2016) could be used in place of specially filmed green-screen sequences, thus rendering any video susceptible to use in object forgery.

Inpainting or video completion (Patwardhan et al., 2007; Wexler et al., 2007; Newson et al., 2014; Ebdelli et al., 2015; Ilan and Shamir, 2015) allows removal of objects from an image by interpolating remaining pixels to conceal a “hole” left by a removed object or corrupt section of video. It is useful for error concealment when streaming video over an unreliable channel and can be used to restore old film, but it can also be used to deliberately remove objects or even frames from a sequence. Video inpainting techniques were surveyed recently in (Ilan and Shamir, 2015), where it was noted that many methods of video inpainting rely on patch completion where the missing spatio-temporal volume is filled using small patches from within the same video sequence. This is evident in early video inpainting methods such as (Patwardhan et al., 2007; Wexler et al., 2007). In (Patwardhan et al., 2007), static background and dynamic foreground were assumed, thus highlighting one of the challenges associated with video completion: motion. This was handled by first registering or aligning the frames of the sequence. Background mosaics of the video sequence were then constructed by removing all non-static objects, and foreground mosaics contained all the moving objects. Missing data was then inpainted by finding close matches in the mosaics and interpolation with texture synthesis between the matched segments. The authors of (Wexler et al., 2007) used space-time volumes of  $5 \times 5 \times 5$  pixels taken from other areas of video sequences to fill in the space left behind by a removed object. Motion was accounted for by representing each pixel not only in terms of its RGB components but also two components based on the derivation along the x-, y- and t-dimensions. This method also allowed temporal and spatial upsizing as spatio-temporal holes had varying dimensions in the spatial and temporal axes. More recently, the work in (Newson et al., 2014) realigned source patches to create closer matches and less warping of synthesised video content. Initial values for missing pixel data were also explicitly defined in (Newson et al., 2014).





**Fig. 5.** An example of spliced content. 5a and 5b are from VTD (Al-Sanjary et al., 2016) and the spliced content (a picture on the wall) comes from a static image. 5c, 5d and 5e come from D'Avino et al. (D'Avino et al., 2017) where the spliced content comes from a chroma-keyed video.

The assessment of inpainting quality in the image domain was critically reviewed in (Qureshi et al., 2017), and user survey to assess the visibility of inpainted regions remains the gold standard. The authors also noted that Video Inpainting Quality Assessment remains an important, open area of research. Indeed, in the absence of an accepted video completion quality assessment, authors (Le et al., 2017; Newson et al., 2014; Ebdelli et al., 2015) simply publish videos of their inpainting techniques applied to standard sequences online and (Ilan and Shamir, 2015) notes this as a trend. The authors of (Ebdelli et al., 2015) also provided their original sequences along with defined masks so that future inpainting techniques can be applied to precisely the same data for comparison. The existence of these inpainted sequences provides a good source of data for video tampering detection research.

Inpainting can be used in conjunction with spatio-temporal copy move to create complex forgeries. An early example of this can be found in (Shih et al., 2011) where the authors changed the winner of a 100 m race. The authors considered the video as a series of layers. They applied inpainting using unoccluded areas of background and interpolated/sampled the motion of forged runners to make them move slower/faster relative to other objects in the video. While individual frames taken from the forged sequences looked visually convincing, full video sequences are not currently available for full analysis. Indeed, assessment of tampered video remains an open problem, one which the authors of (Shih et al., 2011) suggest is best tackled by forgery detection methods. Although the subject matter of this video was somewhat ambitious for its time, and the authors explicitly target the field of video special effects, it gives a good idea of how tampering can be used to court controversy.

#### Style and motion transfer

Style transfer is a new method of image and video manipulation which has been facilitated by the advent of Generative Adversarial Nets (GANs), which were first established in (Goodfellow et al., 2014) and extended to conditional GANs in (Mirza and Osindero, 2014). Style transfer can completely change the context of an image or the subject of a video. It is strongly related to motion transfer because the resultant video is a combination of motion from one

source video and content or subjects from another. Combining the two can be viewed as a style transfer when the style of the content source is mapped to the motion source or it can be considered motion transfer when motion is mapped to the content source.

Examples of style transfer in the image domain include (Isola et al., 2017) where features from one object are mapped to a similar object; a scene can be changed from a summer scene to a winter scene; a horse can be exchanged for a zebra (Zhu et al., 2017); Google Street View House Numbers can be translated into MNIST-style digits (Liu et al., 2017) and evaluated using accuracy on a CNN trained to classify MNIST. Examples in the video domain include motion transfer (Suwajanakorn et al., 2017; Thies et al., 2016) as well as style transfer.

An example of a conditional GAN used to perform style transfer can be found in the seminal Pix2Pix (Isola et al., 2017), which performs image to image translation. A GAN consists of a generator network and a discriminator network. In Pix2Pix, the generator network maps an observed image and a random noise vector to a generated image. The discriminator network then uses both the mapped image and the observed image to classify the mapped image as an example from the authentic dataset or one from the generator. Authentic examples given to the discriminator dictate the “style”. This architecture is distinct from a non-conditional GAN where the discriminator network sees only the mapped image. The authors of Pix2Pix noted that they could achieve very good results based on small datasets of only 400 authentic images and so the GAN can be trained for a multitude of applications. For example, the input image can be a sketch or a semantic segmentation mask and the mapped image can be photorealistic, or vice versa; daytime scenes can be mapped to night; the mapping process can even perform inpainting or background removal. The versatility of (Isola et al., 2017) has also spawned further applications in the video domain including (Bansal et al., 2018; Tulyakov et al., 2018).

Motion transfer is similar to style transfer where the motion of one object is passed on to another object. Early applications were mostly specific to human facial re-enactment such as lip synchronisations and expression translation between talking heads (Suwajanakorn et al., 2017; Thies et al., 2016). Thies et al. (2016) presented the first real-time facial re-enactment system that used only RGB as input. The method used authentic frames from a target



video and transformed them to match the facial expressions and mouth motions from a source video. In (Suwajanakorn et al., 2017), the authors added video re-timing for realistic head motion to fit the context of the spoken word and used a recurrent neural network (RNN) trained on many hours of footage of the particular subject to transform an audio track into mouth shapes. While (Suwajanakorn et al., 2017) was not real time, and required many hours of video footage to train the RNN, it was capable of producing a representative video from audio and stock footage, whereas (Thies et al., 2016) required video for both source and target. More recently, motion transfer has been achieved using models based on style transfer.

MoCoGAN (Tulyakov et al., 2018), used GANs in a similar way to Pix2Pix (Isola et al., 2017). Content and motion were treated independently in MoCoGAN and video sequences expressed as:

$$Z_i = Z_c \times Z_m \quad (7)$$

Every frame in  $Z_i$  has a content vector,  $Z_c$ , and a motion vector,  $Z_m$ , associated with it. In order to perform motion transfer, the content of one sequence was substituted with the content of another. The architecture consisted of two distinct discriminator networks: one to classify real and generated images (or frames), and one to distinguish real and generated video. The video discriminator was responsible for smooth video generation. Similarly, there were two connected generator networks: one to generate motion, the output of which was used to condition the content generator which produced video frames. The motion generator network was a recurrent neural network (RNN) which modelled motion through time. Motion content could also be extracted from a different sequence and hence motion can be transferred between two similar videos. The authors of (Bansal et al., 2018) also applied motion transfer to videos, successfully replicating lip motions. Because (Bansal et al., 2018; Tulyakov et al., 2018) are both based on style transfer, they can also be used to create photo-realistic synthetic video from semantic segmentation masks (Section 2.6).

The assessment of GAN-produced images and videos remains an open problem. In (Liu et al., 2017; Isola et al., 2017), translated images were objectively assessed using the accuracy of CNNs pre-trained on authentic images in the output-style classifying the translated images. It was found that the CNNs classified the translated image of (Liu et al., 2017) with more than 90% accuracy. Image translation methods from (Liu et al., 2017) were also applied to some street driving video sequences, and qualitative analysis of the results showed a convincing, low frame rate video where the weather had been translated from sunny to snowy or the lighting mapped from day to night. The authors of (Luan et al., 2004) applied neural style transfer to photographed objects spliced into images of paintings, thus reducing the visibility of the tampered object. A user study found that their edited image set achieved similar user scores to an unedited image set meaning people could not reliably localise such processed image edits. Although (Liu et al., 2017; Isola et al., 2017; Zhu et al., 2017; Luan et al., 2004) show a method to alter image content, they do not assess whether there is a counter method which can detect these alterations. Since all of these methods employ the use of GANs, it is implicit that there already exists a network which has been trained to discriminate between authentic examples of the style and synthesised content, but due to GAN convergence, this network may not be optimal for detection. In the video domain, the authors of (Suwajanakorn et al., 2017; Thies et al., 2016) have released examples of their work to the public. In (Karras et al., 2017), the authors used a user study to compare their audio-to-video speech synthesiser to both a previous model and a motion capture solution. The study showed that their

work advanced the state-of-the-art as their examples were preferable to human eyes when compared to previous speech synthesis, but not when compared to motion capture generated video. The authors of (Rössler et al., 2018) performed a user study on their tampered dataset and found that, when asked to differentiate between tampered and authentic videos, humans achieved no better than random guessing. Generic motion re-enactment and video generation has also been studied in (Tulyakov et al., 2018) however state-of-the-art is not yet of a standard where such tampered videos are high quality content.

Image to image style translation can be applied to video frames to universally change the overall context of the video. Complimentary work can be found in (Lai et al., 2018) where the authors examined the removal of flicker from a sequence of frames. They specifically aimed to allow the use of image style transfer on individual frames to produce a temporally coherent video sequence, independent of the style transfer method. In (Chan et al., 1808) the authors used style translation on videos to synthesise video content of people performing dance moves they had never done. Pose estimation was used as an intermediate step. A conditional GAN was trained to map a stick-man pose estimation to a photo-realistic video frame using the previous frame to condition the GAN. They then applied a spatio-temporal smoothing to generate convincing videos that showed a target actor dancing in a manner defined by a source actor from a different video. The video sequences were assessed by extracting a pose from the mapped sequence and comparing it to the pose used to generate the sequence, and manual qualitative analysis of temporal qualities including some publicly released sequences. The authors conceded that there were still a number of challenges to overcome in this field, such as loose clothing and cluttered background, but it is easy to see that motion transfer can already be convincingly applied to human faces and bodies.

#### Photo-realistic synthetic video

Although purely synthetic video in the form of animation has been around for a long time, more recently synthetic video has been generated which is so photo-realistic that it could be mistaken for authentic, filmed content. In this section, we examine the most recent techniques in photo-realistic video synthesis and discuss their evaluation. Although video synthesis is not explicitly tampering an existing video, full convincing, photo-realistic video synthesis has the potential to be just as damaging as motion transfer or inpainting. It is important to examine it with a view to detecting it as a future research direction. Current trends in top international conferences on computer vision show that video frame prediction is a strong trend.

A short video sequence was extrapolated from the motion blur of a single image in (Jin et al., 2018). The authors noted that the main challenge of this is temporal ordering. While the central frame of the synthesised sequence corresponds to the de-blurred image, the motion of individual objects in frames before and after is ambiguous. The authors proposed a pair-wise ordering invariant loss to aid convergence of their CNN, which was based on pairs of frames at an equal temporal distance from the middle frame. Although the de-blurring aspect of the technique improved on the previous method (Nah et al., 2017) for moderate blur, evaluation of the short synthetic sequences proved difficult. The ambiguity in temporal ordering could be resolved when the process is constrained to temporal super-resolution. Video generation from a single image was also covered in (Wang et al., 2018b) where Wang et al. detailed a method to produce a short photo-realistic video of a smile from a single aligned face image. The method use a series of conditional Long Short Term Memories (LSTMs) to produce a

sequence of facial landmarks moving from a neutral expression to a smile. A network similar to (Isola et al., 2017) was then used to translate the facial landmarks into a realistic video. A comparative user study found that the resulting sequences looked more realistic than a previous method, but the authors noted that it was difficult to evaluate such a method as there were no directly comparable existing methods. Both (Jin et al., 2018) and (Wang et al., 2018b) extrapolated short synthetic video sequences based on a single image, and both noted challenges in evaluation. Xiong et al. (2018) produced short, realistic time-lapse videos of skyscapes, up to 32 frames from a single image using a two-stage GAN architecture. The first stage produced a sequence of frames and the second stage refined it to produce a coherent video. They also gathered a large dataset of real time-lapse videos from YouTube for the purposes of training. Again, there was no previous work available for direct comparison, but the authors were able to repurpose other network architectures to synthesise time-lapse videos. Evaluation was by user study where users were asked to identify the more realistic of two sequences. Although the proposed method outperformed all other synthetic videos, when comparing synthetic video with real video, only 16% of synthetic video tests were preferred.

In (Lotter et al., 2011), a CNN, LSTM and deconvolutional neural net were used together to predict the next frame in a video sequence. The authors noted that natural images were much more challenging than simple moving circle animations, and that, in predicting the next frame in a face rotation sequence, the network altered sufficient features so as to change the perceived identity of the face. The work in (Babaeizadeh et al., 2010) followed this, using a variational autoencoder to predict the next 10 frames from a 10 frame sequence. The authors tested their sequences on (Ionescu et al., 2014) among others where they compared their predicted frames with ground truth frames using PSNR and SSIM. They conceded that assessing the quality of the predicted frames was difficult, and that the prediction yielding the worst PSNR was sometimes qualitatively the best. They also publicly released many examples of predicted sequences. In (Xu et al., 2018), the authors used a two stream structure and RNN to perform frame prediction. The authors of (Cai et al., 2018) viewed frame prediction as analogous to frame interpolation and fully synthetic video generation. They successfully produced short video sequences which interpolated between two frames as well as predicting short sequences given only the first frame of each sequence. Evaluation used PSNR and SSIM for interpolated sequences and Inception scores for generated sequences with no ground truth. Frame prediction was also covered in (Zhao et al., 2018; Reda et al., 2018; He et al., 2018), and evaluation also involved full reference quality metrics. Although methods were evaluated using full reference quality metrics, there is no guarantee that frame *prediction* as opposed to interpolation will predict a frame that matches the original sequence, but it may yet produce a valid, realistic frame.

Some methods create synthetic video data specifically for machine learning datasets. In (Varol et al., 2017), a dataset of synthetic videos was created from motion capture data. The motion capture data was used to generate 3D models of human bodies which were combined with a texture map to add clothes and skin, and a static background image. The synthetic videos, rendered using Blender, were found to improve body part and foreground background segmentation. The Human 3.6 M dataset (Ionescu et al., 2014), includes some mixed reality videos which consist of a moving synthetic human model combined with a real video sequence. The real backgrounds included annotated occluding items so that synthetic human models could realistically interact with authentically filmed objects. Neither (Varol et al., 2017) nor (Ionescu et al., 2014) are specifically designed to fool human eyes, but instead intended to aid development of human pose estimation and body

segmentation. Rather than annotate thousands of frames of authentic video, the synthetic human model is already annotated. This is an example of a non-malicious application of synthetic video, although the detection of the synthetic parts is often trivial to human eyes.

Wang et al. (2018a) have already synthesised coherent, photo-realistic video sequences of up to 30 s from semantic segmentation mask or pose model sequences. A GAN was used, and a discriminator part of the GAN used to classify the content as an authentic video or not, similar to MoCoGAN (Tulyakov et al., 2018). Using a discriminator in this way ensured temporally coherent video. The authors conceded that significant changes in an object's appearance is still a substantial challenge and that their model was also prone to colour drift over time, however a user study showed that (Wang et al., 2018a) produced video that was preferable to human eyes than that produced by MoCoGAN (Tulyakov et al., 2018). SCGAN (Yang et al., 2018a) also performed a user study which put their synthetic video content at a higher level of realism than MoCoGAN. Recycle-GAN (Bansal et al., 2018) was also used to generate photo-realistic video from semantic segmentation mask sequences and assessed their method's accuracy by asking users to classify videos as synthetic or real as well as comparison with existing state-of-the-art. Users were fooled into thinking synthetic video was real 28.3% of the time. Quantitative results were also obtained using the Viper dataset (Richter et al., 2017) which supplies pixel-level segmentation masks for computer game scenes with a high level of realism.

Methods of evaluation for synthetic video remains an open field. It can be seen in Table 1 that the main methods include full reference quality metrics PSNR and SSIM as well as a variety of others such as Average Content Distance (ACD) (Tulyakov et al., 2018; Zhao et al., 2018), Learned Perceptual Image Patch Similarity (LPIPS) (Chan et al., 2018) and tampering detection methods. Many of the techniques for synthetic video generation also utilise user surveys to assess the quality of synthetic video, and in many cases, evaluation is relative to previous related work (Table 1). It can be inferred from this that although current methods do not yet reliably generate video that is photo-realistic enough to fool human eyes, improvements are continuous and incremental. It is simply a matter of time before photo-realistic synthetic video becomes mainstream. As (Rössler et al., 2018) showed, some techniques are already indistinguishable from authentic video for human viewers. This raises the problem that, in future, not only will humans be unable to detect tampered or even photo-realistic synthetic video, they will also be blind to whatever tampering technique has been applied. When this is the case, universal tampering detection systems will be required to fill the gap in human perception, and these must be developed urgently if detection systems are to keep pace with tampering methods. For this, datasets are required.

#### Image tampering detection

To advance the field of universal video tampering detection, it is vital to gather datasets of independent examples of video tampering techniques. In this section, we look at the lessons relating to tampered image datasets that can be learned from the application of deep neural networks to the problem of tampering detection.

As machine learning techniques come to the fore in tampering detection, the collation of large datasets to train and test networks becomes desirable. However it is important to realise that *any* consistencies within labelled classes may be exploited as features by deep learning techniques, including any features arising during dataset generation that are unrelated to actual tampering. In 2011 Torralba and Efros (2011) discussed how bias is ubiquitous within

computer vision datasets. Images from the same dataset exhibit characteristics specific to that dataset, so much so that a basic support vector machine (SVM) classifier trained to label a given image with its associated dataset achieved reasonable accuracy of 39% over 12 datasets. Each dataset has its own inherent distribution which may be irrelevant to the real world situation, and may be overlooked by human eyes. This problem is subtly highlighted by the advance of deep learning, particularly in the field of image forensics.

A good example of unintentional features comes in the CASIA2 TIDE dataset ([. This large dataset consists of 7491 authentic and 5123 tampered images which use splicing or copy-move techniques. The size of this dataset makes it an attractive option for deep learning and over 97% classification accuracy has been achieved by \(Rao and Ni, 2016\). However, as noted in \(Sutthiwan et al., 2011\), compression applied to tampered images of the dataset differs from that applied to authentic images. Put simply, during dataset generation, tampered images were compressed twice, authentic images were compressed only once. There were also patterns in the colour space resolutions with tampered images more likely to have lower colour channel resolution. This means that classifying a CASIA2 image as tampered or authentic can be accurately achieved using features of compression, recompression and colour resolution. The recompression step may have arisen from the tools used to tamper the images, but it is independent of the tampering task itself. There is no reason that an authentic image cannot be innocently recompressed.](#)

In ([Rota et al., 2016](#)), dataset weaknesses such as those in CASIA2 were used as an explanation for the sharp drop off in CNN classification accuracy whenever the test images were compressed. Classification accuracy dropped from 97.44% on unprocessed CASIA2 image patches to 68.11% when the images were compressed with JPEG quality factor 90, a fairly light compression. The authors proposed a means to circumvent this dataset flaw by extracting authentic patches from tampered images, however they did not report whether this reduced the drop-off in accuracy when the source dataset was compressed, nor did they report on a CNN trained using compressed image patches. Maliciously tampered images in the wild are not necessarily recompressed, and authentic images are not necessarily compressed only once. Tampered and authentic patches may be extracted from only the tampered data but only if reliable localisation masks exist to differentiate tampered and authentic pixels.

High levels of classification accuracy were also achieved by a deep neural network on the large rebroadcast dataset presented in ([Agarwal et al., 2018](#)). This dataset comprises over 29,000 images, half authentic and half rebroadcast in some way. Rebroadcast techniques included printing out and rescanning/photographing the images, screen grabs and screen photography. While some traditional techniques ([Cao and Kot, 2010](#)) demonstrated poor accuracy on this new dataset, a CNN trained on 60% of the images and tested on the remainder achieved over 97% accuracy. In this case, recompression is very likely a necessary feature of retransmission,

so features that emerged during CNN training are a true reflection of the real process of retransmission. One way to objectively assess this is to check the performance on a rebroadcast test set gathered independently. If a deep neural network exploits unintentional weaknesses inherent in a particular dataset, then the learning will not transfer well to other, similar datasets unless they exhibit the same features.

[Table 2](#) shows how CNNs excel in detection of image anti-forensics. A detection or classification accuracy of over 95% is a common occurrence. Anti-forensics are methods designed to “launder” tampering and thus fool tampering detectors. Laundering techniques include general filtering methods such as compression, median filter and Gaussian blur. This field is emerging rapidly because large datasets can be synthesised with relative ease, and this makes it particularly appropriate for machine learning. Datasets such as BOSSBase ([Bas et al., 2011](#)), UCID ([Schaefer and Stich, 2003](#)) and Dresden image database ([Gloe and Böhme, 2010](#)) provide a large variety of unprocessed images to which known anti-forensic techniques can be universally applied and subsequently detected.

In ([Boroumand and Fridrich, 2018](#)), a CNN was trained to identify laundering techniques applied to an image. The laundering techniques, applied singly, were: low-pass-, high-pass- and denoising-filters and tonal sharpening. The authors first compressed the images of the dataset, then applied a single laundering technique and then rescaled, cropped and recompressed the resulting images. Compression was JPEG with a Quality Factor (QF) ranging from 75 to 95. They achieved over 95% accuracy in identification of the laundering technique used regardless of the image compression level, provided the CNN was trained on images with a QF similar to that of the test dataset. The idea that a training dataset must be well matched to the test data in terms of compression is also supported in ([Johnston et al., 2018](#)). All of these high accuracies show that machines are adept at detecting patterns in visual data which are invisible to humans. This makes designing a dataset which is representative of the problem of video tampering but immune to unintentional side effects especially important.

#### Tampered video datasets

With so many different methods of tampering already available, and the field progressing at an unprecedented rate, it is important for tampering detection techniques to keep pace. Unfortunately this is challenging because there are few large, diverse tampered video datasets. [Pandey et al. \(2016\)](#) noted that tampered video datasets lag far behind tampered image datasets in terms of maturity. In this section, we examine existing video datasets and give recommendations for the design of new datasets. [Table 3](#) provides a list of video tampering datasets, specifying the type of tampering applied and the size of the datasets.

A number of tampered video datasets already exist, but these vary both in terms of processing and parameters. In ([D’Avino et al., 2017](#)) the authors supply tampered video along with an explicit

**Table 2**  
CNNs for image anti-forensics detection.

Reference	Detection of:	Dataset	Accuracy
Bayar and Stamm ( <a href="#">Bayar and Stamm, 2016</a> )	Gaussian blurring, additive white noise, median filter, resampling	proprietary	99%
Choi et al. ( <a href="#">Choi et al., 2017</a> )	All combinations of Gaussian blurring, Median filtering	( <a href="#">Bas et al., 2011</a> ; <a href="#">Gloe and Böhme, 2010</a> )	> 91%
Choi et al. ( <a href="#">Choi et al., 2017</a> )	Gamma correction	( <a href="#">Bas et al., 2011</a> ; <a href="#">Gloe and Böhme, 2010</a> )	57.6%
Amerini et al. ( <a href="#">Amerini et al., 2017</a> )	Double compression level	<a href="#">Schaefer and Stich (2003)</a>	83.5%–99.9%
Boroumand and Fridrich ( <a href="#">Boroumand and Fridrich, 2018</a> )	Low-pass-, high-pass-, denoising- filters and tonal sharpening	<a href="#">Bas et al. (2011)</a>	> 95%
Agarwal ( <a href="#">Agarwal et al., 2018</a> )	Rebroadcast	Public ( <a href="#">Agarwal et al., 2018</a> )	> 97%



**Table 3**  
Tampered video datasets.

Name, Date, Ref.	Type of Tampering	Size	Details
VTD 2016 (Al-Sanjary et al., 2016)	splicing; copy-move; frame-shuffling	26 tampered + related authentic	Distribution on YouTube means all videos affected by varying compression
FaceForensics 2018 (Rössler et al., 2018)	motion transfer ((Thies et al., 2016))	1004 tampered $\times 2$	Taken from Youtube-8m (Abu-El-Haija et al., 2016), one set self-re-enactment, one set source-target translation
SYSU-OBJFORG 2016 (Chen et al., 2016)	object forgery	100 authentic, 100 tampered	No source for public download
D'Avino et al., 2017 (D'Avino et al., 2017)	splicing	10 tampered	Binary masks provided, tampering is easily seen.
SULFA forged 2012 (Qadir et al., 2012)	spatio-temporal copy-move	5 tampered	Static camera, background duplicated to conceal objects. Part of a larger database including untampered video for camera identification
SULFA supplemental 2013 (Bestagini et al., 2013)	spatio-temporal copy-move	10 tampered	Duplicate spatio-temporal regions to conceal/introduce objects
Lin et al., 2014 (Lin and Tsay, 2014)	inpainting (TCP and ETS)	18 tampered $\times 2$	Only 4 sequences available for direct download
VISION 2017 (Shullani et al., 2017a)	source/social media platform identification	1914 sequences	648 straight-from-device videos, 622 YouTube and 644 WhatsApp. Future dataset extension planned via "MOSES" application (Shullani et al., 2017b)
Ardizzone and Mazzola, 2015 (Ardizzone and Mazzola, 2015)	copy-move	160 sequences	Sequences synthesised from (Qadir et al., 2012) and CANTATA datasets
Newson et al., 2014 (Newson et al., 2014)	inpainting	3 tampered	Masks supplied for two sequences, demonstration of inpainting rather than explicit tampering detection dataset
Le et al., 2017 (Le et al., 2017)	inpainting	53 sequences	Both object removal and object reconstruction, demonstration of inpainting rather than explicit tampering detection dataset

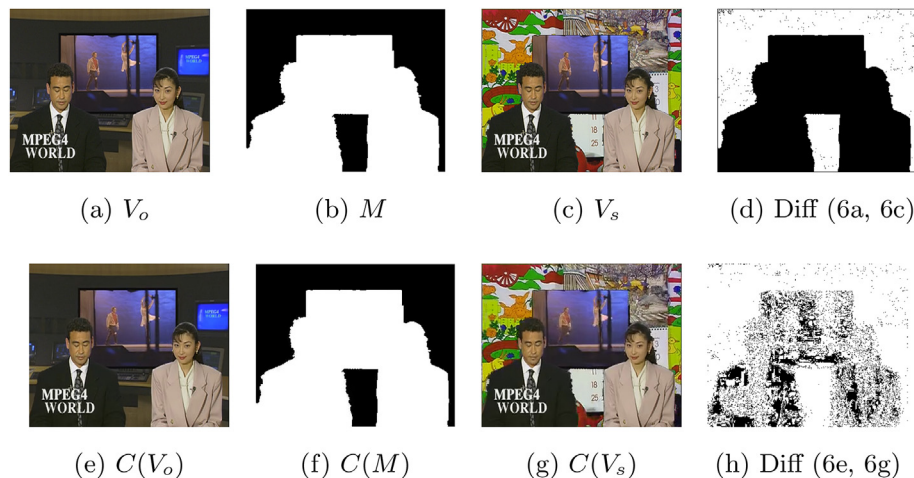
pixel level binary mask detailing the chroma-keyed addition. Some video tampering datasets come complete with original and tampered videos, thus providing a means to calculate all masks and labels associated with tampering (Al-Sanjary et al., 2016; Bestagini et al., 2013). This allows for tampering detection and localisation in spatial and temporal domains. It also allows for any differences in distribution between tampered and authentic sequences to be overcome by extracting authentic patches from tampered sequences. However, an accurate mask can only be extracted where videos can be synchronised and are identically processed post-tampering. Any recompression during distribution allows compression errors to creep in and increases the difficulty of extracting a bit-accurate tamper mask (as in Fig. 6). Moreover, some pixels may be part of the tampered region but remain unchanged in value, and this makes for noisy masks in need of post processing.

Using differences between original and tampered videos may be inappropriate for temporally tampered videos (Al-Sanjary et al., 2016), where a frame-by-frame label might provide more information. This can be achieved when unprocessed original and

tampered sequences are provided. Indeed, public inter-frame tampered video datasets are in short supply, with many inter-frame tampering detectors simply building their own datasets from available sequences (see Section 2.2).

As can be seen in Table 3, most tampered video datasets focus on a single tampering method, such as splicing or object forgery or inpainting. Only VTD (Al-Sanjary et al., 2016) demonstrates a variety of types. Variety is vital to accurately assess performance of video tampering detectors and support work towards universal video tampering detection. As discussed in (Torralba and Efros, 2011), an approach using a combination of datasets will ensure more generalisable results with little need for specific domain adaptation. It is also important that tampered sequences are independent of tampering detection, so techniques such as (Le et al., 2017; Ardizzone and Mazzola, 2015; Newson et al., 2014) which publicly release their results are important to move forward both tampering AND tampering detection.

A number of datasets are produced and benchmarked with an existing detection technique (Rössler et al., 2018; Lin and Tsay,



**Fig. 6.** The problems with recompression in the distribution of tampered datasets: Left column shows uncompressed data, right column has been lightly compressed. Fig. 6e–g are the compressed versions of 6a, 6b, 6c respectively. Fig. 6b and h are both uncompressed and show binarised differences.

2014; Shullani et al., 2017a) and many achieve high levels of precision on their selected dataset, often over 90% accuracy. This tends to portray tampering detection as a solved problem on that particular dataset, which discourages researchers from publishing lower results. A tampering detection method based on motion residue was presented in (Kancherla and Mukkamala, 2012), and the experimental dataset was gathered from several previous object forgery works (Patwardhan et al., 2007; Shih et al., 2011; Hsu et al., 2008). Accuracy was lower than 90%. This contrasts with over 90% in (Hsu et al., 2008) and over 99% benchmarked in uncompressed FaceForensics (Rössler et al., 2018). However, in a real world situation, the original method of tampering will be unknown, therefore, it is worth collating results on several different datasets such that work towards a universal tampering detection can be realised.

Methods of dataset dissemination are an important consideration as this can cause unintentional post-processing of video data. Although video sharing websites such as YouTube may seem like an attractive distribution option (Al-Sanjary et al., 2016), any processing applied during publishing must be taken into account. It is possible to apply social media platform processing by uploading/downloading a video to/from a social media website, however the effects on the video are then irreversible. While researchers may add processing to an unprocessed video, they cannot remove it. Indeed, the effects of video processing by social media platforms on video are represented in isolation in a dataset provided by (Shullani et al., 2017a) who found that sensor noise pattern used for camera source identification was adversely affected by processing on Facebook, YouTube and WhatsApp, even when using high quality settings. These results are important in themselves as they show how tampering detection methods which rely on sensor noise, such as (Kobayashi et al., 2009), can be defeated by virtue of the distribution platform alone. They also emphasise how post-processing can be easily overlooked.

Fig. 6 illustrates some of the complications associated with recompression. Starting with uncompressed data, a binary mask was created based on segmentation of static and non-static content, and two uncompressed sequences simply spliced together. Fig. 6d shows which pixels differ between Fig. 6a and c and it can be seen that it is almost the perfect inverse of the mask (Fig. 6b), with a few pixels that are identical between the original and spliced content. Fig. 6e–g show the visual effects of compression on Fig. 6a–c respectively. Fig. 6h shows how compression has introduced tiny inaccuracies between the pixels of the original and spliced sequences so that the difference between them no longer provides a mostly accurate inverse of the mask. With some thresholding and morphological processing, the difference sequence could still be used to infer a mask, but the degree of accuracy suffers even under slight compression. A compressed mask, as shown in Fig. 6f provides a more accurate ground truth than deriving the mask from the compressed tampered/untampered pair. Moreover, official mask provision rather than frame difference inference removes the philosophical debate over whether a pixel, fully within the tampered region but by chance unchanged by the tampering process, is labelled tampered or not.

## Conclusion

Many modern techniques of video tampering simply do not fit neatly into the traditional categories of inter- and intra-frame tampering. In particular, there is significant overlap between the recent categories of motion/style transfer and synthetic video generation. Changing the style of a video sequence from semantic segmentation masks to photo-realistic generates a purely synthetic video, but the same techniques can be used to perform digital

puppetry and transfer motion from one mouth to another. This means that detection of synthetic video should be viewed as an extension of tampering detection. Given the current trend of using full reference quality measures in the evaluation of retouching, frame interpolation and in video frame prediction, it is clear that one of the current goals is to replicate authentic video. What remains unclear, however, is whether these methods will deviate from authentic content as evaluation methods emerge or even help to launder video tampering evidence in the same way as video compression.

One important new research direction in digital video manipulation is an accepted method of evaluation. Many existing methods rely on only qualitative evaluation and while this is an important first step, adoption of existing video quality techniques, including no reference quality metrics will speed up development. Until then, user studies and public release of manipulated video clips remains the gold standard. In the absence of elegant quality measures, altered video and the associated methods are often publicly released for analysis, and video tampering detectors should look to utilise this provision where possible to create realistic detection methods. To facilitate this, video tamperers should release either sufficient data to simplify the creation of accurate tampering masks or release the masks themselves. Furthermore, video data should be distributed in such a way as to minimise further processing. Video processing, such as compression and retouching can effectively conceal tampering. While detection of such anti-forensics is an important research direction, processing can be applied to video independently after dataset publication, but only if the original dataset is published in such a way as to avoid unnecessary processing. With the increasing application of deep learning methods to tampering detection, any future dataset gatherers must take care to avoid potential pitfalls which cause datasets to reflect their own specific features relating to publishing platform or tool use, rather than those legitimately tied to the tampering technique.

As the variety of video manipulation techniques expands and advances, tampered and synthetic video will become indistinguishable from authentic video to human eyes. Therefore, new techniques are required which either classify tampered video according to its tampering type or perform tampering detection irrespective of the type of tampering. To maintain confidence in the authenticity of video content in future, it is crucial to develop techniques which can identify and localise video processing and manipulation. Universal video manipulation detection and localisation is essential if tampering detection is to keep pace with tampering methods.

This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

## References

- Abu-El-Haija, S., Kothari, N., Lee, J., Natsev, P., Toderici, G., Varadarajan, B., Vijayanarasimhan, S., 2016. Youtube-8m: A Large-Scale Video Classification Benchmark arXiv preprint arXiv:1609.08675.
- Agarwal, S., Fan, W., Farid, H., 2018. A diverse large-scale dataset for evaluating rebroadcast attacks. In: IEEE International Conference on Acoustics, Speech, and Signal Processing.
- Aghamaleki, J.A., Behrad, A., 2017. Malicious inter-frame video tampering detection in mpeg videos using time and spatial domain analysis of quantization effects. *Multimed. Tool. Appl.* 76 (20), 20691–20717.
- Aksoy, Y., Aydin, T.O., Pollefeys, M., Smolić, A., 2016. Interactive high-quality green-screen keying via color unmixing. *ACM Trans. Graph.* 35 (5), 152.
- Al-Sanjary, O.I., Sulong, G., 2015. Detection of video forgery: a review of literature. *J. Theor. Appl. Inf. Technol.* 74 (2).
- Al-Sanjary, O.I., Ahmed, A.A., Sulong, G., 2016. Development of a video tampering dataset for forensic investigation. *Forensic Sci. Int.* 266, 565–572.
- Amerini, I., Uricchio, T., Ballan, L., Caldelli, R., 2017. Localization of jpeg double compression through multi-domain convolutional neural networks. In: *Proc. Of IEEE CVPR Workshop on Media Forensics*, vol. 3.

- Ardizzone, E., Mazzola, G., 2015. A tool to support the creation of datasets of tampered videos. In: *International Conference on Image Analysis and Processing*. Springer, pp. 665–675.
- M. Babaeizadeh, C. Finn, D. Erhan, R. H. Campbell, S. Levine, Stochastic Variational Video Prediction, arXiv preprint arXiv:1710.11252.
- Bansal, A., Ma, S., Ramanan, D., Sheikh, Y., 2018. Recycle-gan: unsupervised video retargeting. In: *European Conference on Computer Vision*. Springer, pp. 122–138.
- Bas, P., Filler, T., Pevný, T., 2011. Break our steganographic system: the ins and outs of organizing boss. In: *International Workshop on Information Hiding*. Springer, pp. 59–70.
- Bayar, B., Stamm, M.C., 2016. A deep learning approach to universal image manipulation detection using a new convolutional layer. In: *Proceedings of the 4th ACM Workshop on Information Hiding and Multimedia Security*. ACM, pp. 5–10.
- Bestagini, P., Milani, S., Tagliasacchi, M., Tubaro, S., 2013. Local tampering detection in video sequences. In: *Multimedia Signal Processing (MMSp)*, 2013 IEEE 15th International Workshop on, IEEE, pp. 488–493.
- Bideau, P., Learned-Miller, E., 2016. Its moving! a probabilistic model for causal motion segmentation in moving camera videos. In: *European Conference on Computer Vision*. Springer, pp. 433–449.
- Boroumand, M., Fridrich, J., 2018. Deep learning for detecting processing history of images. *Electron. Imag.* 2018 (7), 1–9.
- Caballero, J., Ledig, C., Aitken, A., Acosta, A., Totz, J., Wang, Z., Shi, W., 2017. Real-time video super-resolution with spatio-temporal networks and motion compensation. In: *Computer Vision and Pattern Recognition (CVPR)*, 2017 IEEE Conference on, IEEE, pp. 2848–2857.
- Cai, H., Bai, C., Tai, Y.-W., Tang, C.-K., 2018. Deep video generation, prediction and completion of human action sequences. In: *The European Conference on Computer Vision (ECCV)*.
- Cao, H., Kot, A.C., 2010. Identification of recaptured photographs on lcd screens. In: *Acoustics Speech and Signal Processing (ICASSP)*, 2010 IEEE International Conference on, IEEE, pp. 1790–1793.
- Cavigelli, L., Hager, P., Benini, L., 2017. Cas-cnn: a deep convolutional neural network for image compression artifact suppression. In: *Neural Networks (IJCNN)*, 2017 International Joint Conference on, IEEE, pp. 752–759.
- Chan, C., Ginosar, S., Zhou, T., Efros, A.A., Everybody Dance Now. 1808.07371.
- Chauhan, D., Kasat, D., Jain, S., Thakare, V., 2016. Survey on keypoint based copy-move forgery detection methods on image. *Procedia Computer Science* 85, 206–212.
- Chen, J., Kang, X., Liu, Y., Wang, Z.J., 2015. Median filtering forensics based on convolutional neural networks. *IEEE Signal Process. Lett.* 22 (11), 1849–1853.
- Chen, S., Tan, S., Li, B., Huang, J., 2016. Automatic detection of object-based forgery in advanced video. *IEEE Trans. Circuits Syst. Video Technol.* 26 (11), 2138–2151.
- Chen, L., Li, Z., K Maddox, R., Duan, Z., Xu, C., 2018. Lip movements generation at a glance. In: *The European Conference on Computer Vision (ECCV)*.
- Choi, H.-Y., Jang, H.-U., Kim, D., Son, J., Mun, S.-M., Choi, S., Lee, H.-K., 2017. Detecting composite image manipulation based on deep neural networks. In: *Systems, Signals and Image Processing (IWSSIP)*, 2017 International Conference on, IEEE, pp. 1–5.
- Chollet, F., 2017. Xception: deep learning with depthwise separable convolutions. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1251–1258.
- Chuang, W.-H., Su, H., Wu, M., 2011. Exploring compression effects for improved source camera identification using strongly compressed video. In: *Image Processing (ICIP)*, 2011 18th IEEE International Conference on, IEEE, pp. 1953–1956.
- Conotter, V., O'Brien, J.F., Farid, H., 2012. Exposing digital forgeries in ballistic motion. *IEEE Trans. Inf. Forensics Secur.* 7 (1), 283–296.
- Cotsaces, C., Nikolaidis, N., Pitas, I., 2006. Video shot boundary detection and condensed representation: a review. *IEEE Signal Process. Mag.* 23 (2), 28–37.
- Cozzolino, D., Poggi, G., Verdoliva, L., 2017. Recasting residual-based local descriptors as convolutional neural networks: an application to image forgery detection. In: *Proceedings of the 5th ACM Workshop on Information Hiding and Multimedia Security*. ACM, pp. 159–164.
- Credits for the use of the casia image tempering detection evaluation database (casia tide) v2.0 are given to the national laboratory of pattern recognition, institute of automation, chinese academy of science, corel image database and the photographers. <http://forensics.idealtest.org>.
- Criminisi, A., Pérez, P., Toyama, K., 2004. Region filling and object removal by exemplar-based image inpainting. *IEEE Trans. Image Process.* 13 (9), 1200–1212.
- Dar, Y., Bruckstein, A.M., 2015. Motion-compensated coding and frame rate up conversion: models and analysis. *IEEE Trans. Image Process.* 24 (7), 2051–2066.
- Dong, H., Neekhar, P., Wu, C., Guo, Y., 2017. Unsupervised Image-To-Image Translation with Generative Adversarial Networks. 1701.02676.
- Dong, C., Deng, Y., Change Loy, C., Tang, X., 2015. Compression artifacts reduction by a deep convolutional network. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 576–584.
- D'Avino, D., Cozzolino, D., Poggi, G., Verdoliva, L., 2017. Autoencoder with recurrent neural networks for video forgery detection. *Electron. Imag.* 2017 (7), 92–99.
- Ebelli, M., Le Meur, O., Guillemot, C., 2015. Video inpainting with short-term windows: application to object removal and error concealment. *IEEE Trans. Image Process.* 24 (10), 3034–3047.
- Gironi, A., Fontani, M., Bianchi, T., Piva, A., Barni, M., 2014. In: *A Video Forensic Technique for Detecting Frame Deletion and Insertion*. ICASSP, pp. 6226–6230.
- Gloe, T., Böhme, R., 2010. The dresden image database for benchmarking digital image forensics. *J. Digit. Forensic Pract.* 3 (2–4), 150–159.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014. Generative adversarial nets. In: *Advances in Neural Information Processing Systems*, pp. 2672–2680.
- Guo, J., Chao, H., 2017. One-to-many network for visually pleasing compression artifacts reduction. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3038–3047.
- Ha, T., Lee, S., Kim, J., 2004. Motion compensated frame interpolation by new block-based motion estimation algorithm. *IEEE Trans. Consum. Electron.* 50 (2), 752–759.
- He, L., Lu, W., Jia, C., Hao, L., 2017. Video quality assessment by compact representation of energy in 3d-dct domain. *Neurocomputing* 269, 108–116.
- He, J., Lehmann, A., Marino, J., Mori, G., Sigal, L., 2018. Probabilistic video generation using holistic attribute control. In: *The European Conference on Computer Vision (ECCV)*.
- Hsu, C.-C., Hung, T.-Y., Lin, C.-W., Hsu, C.-T., 2008. Video forgery detection using correlation of noise residue. In: *Multimedia Signal Processing*, 2008 IEEE 10th Workshop on, IEEE, pp. 170–174.
- Huang, C.C., Zhang, Y., Thing, V.L., 2017. Inter-frame video forgery detection based on multi-level subtraction approach for realistic video forensic applications. In: *Signal and Image Processing (ICSIP)*, 2017 IEEE 2nd International Conference on, IEEE, pp. 20–24.
- Ilan, S., Shamir, A., 2015. A survey on data-driven video completion. In: *Computer Graphics Forum*, vol. 34. Wiley Online Library, pp. 60–85.
- Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C., 6m, Human3., 2014. Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Trans. Pattern Anal. Mach. Intell.* 36 (7), 1325–1339.
- Isola, P., Zhu, J., Zhou, T., Efros, A.A., 2017. Image-to-image translation with conditional adversarial networks. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5967–5976. <https://doi.org/10.1109/CVPR.2017.632>.
- ITU-T, H., 2012. 262 Information Technology - Generic Coding of Moving Pictures and Associated Audio Information: Video. ITU-T.
- ITU-T, 2016a. H.264 Advanced Video Coding for Generic Audiovisual Services. ITU-T.
- ITU-T, 2016b. H.265 High Efficiency Video Coding. ITU-T, 12.
- Janai, J., Güney, F., Wulff, J., Black, M.J., Geiger, A., 2017. Slow flow: exploiting high-speed cameras for accurate and diverse optical flow reference data. In: *Computer Vision and Pattern Recognition (CVPR)*, 2017 IEEE Conference on, IEEE, pp. 1406–1416.
- Jiang, H., Sun, D., Jampani, V., Yang, M.-H., Learned-Miller, E., Kautz, J., 2018. Super slo-mo: high quality estimation of multiple intermediate frames for video interpolation. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Jin, M., Meishvili, G., Favaro, P., 2018. Learning to extract a video sequence from a single motion-blurred image. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Johnston, P., Elyan, E., Jayne, C., 2018. Spatial effects of video compression on classification in convolutional neural networks. In: *Neural Networks (IJCNN)*, 2018 International Joint Conference on, IEEE, pp. 1370–1377.
- Joshi, V., Jain, S., 2015. Tampering detection in digital video: a review of temporal fingerprints based techniques. In: *Computing for Sustainable Global Development (INDIACom)*, 2015 2nd International Conference on, IEEE, pp. 1121–1124.
- Kancherla, K., Mulkamala, S., 2012. Novel blind video forgery detection using markov models on motion residue. In: *Asian Conference on Intelligent Information and Database Systems*. Springer, pp. 308–315.
- Karras, T., Aila, T., Laine, S., Herva, A., Lehtinen, J., 2017. Audio-driven facial animation by joint end-to-end learning of pose and emotion. *ACM Trans. Graph.* 36 (4), 94.
- Kaur, J., Upadhyay, S., Sharma, A., 2017. A video database for intelligent video authentication. In: *Computing, Communication and Automation (ICCCA)*, 2017 International Conference on, IEEE, pp. 1081–1085.
- Khodabakhsh, A., Busch, C., Ramachandra, R., 2018. A taxonomy of audiovisual fake multimedia content creation technology. In: *2018 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)*. IEEE.
- Kirmemis, O., Bakar, G., Murat Tekalp, A., 2018. Learned compression artifact removal by deep residual networks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 2602–2605.
- Kobayashi, M., Okabe, T., Sato, Y., 2009. Detecting video forgeries based on noise characteristics. In: *Pacific-rim Symposium on Image and Video Technology*. Springer, pp. 306–317.
- Kupyn, O., Budzan, V., Mykhailych, M., Mishkin, D., Matas, J., Deblurgan, 2018. Blind motion deblurring using conditional adversarial networks. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Lai, W.-S., Huang, J.-B., Wang, O., Shechtman, E., Yumer, E., Yang, M.-H., 2018. Learning blind video temporal consistency. In: *The European Conference on Computer Vision (ECCV)*.
- Le, T., Almansa, A., Gousseau, Y., Masnou, S., 2017. Motion-consistent video inpainting. In: *ICIP 2017: IEEE International Conference on Image Processing*.
- Li, J., Li, X., Yang, B., Sun, X., 2015. Segmentation-based image copy-move forgery detection scheme. *IEEE Trans. Inf. Forensics Secur.* 10 (3), 507–518.
- Li, R., Liu, Z., Zhang, Y., Li, Y., Fu, Z., 2018. Noise-level estimation based detection of motion-compensated frame interpolation in video sequences. *Multimed. Tool. Appl.* 77 (1), 663–688.
- Liao, R., Tao, X., Li, R., Ma, Z., Jia, J., 2015. Video super-resolution via deep draft-ensemble learning. In: *The IEEE International Conference on Computer Vision*



- (ICCV).
- Lin, C.-S., Tsay, J.-J., 2014. A passive approach for effective detection and localization of region-level video forgery with spatio-temporal coherence analysis. *Digit. Invest.* 11 (2), 120–140.
- Liu, Z., Luo, P., Wang, X., Tang, X., 2015. Deep learning face attributes in the wild. In: *Proceedings of International Conference on Computer Vision (ICCV)*.
- Liu, M.-Y., Breuel, T., Kautz, J., 2017. Unsupervised image-to-image translation networks. In: *Advances in Neural Information Processing Systems*, pp. 700–708.
- W. Lotter, G. Kreiman, D. Cox, Unsupervised Learning of Visual Structure Using Predictive Generative Networks, arXiv preprint arXiv:1511.06380.
- F. Luan, S. Paris, E. Shechtman, K. Bala, Deep Painterly Harmonization, arXiv preprint arXiv:1804.03189.
- Meyer, S., Djelouah, A., McWilliams, B., Sorkine-Hornung, A., Gross, M., Schroers, C., 2018. Phasenet for video frame interpolation. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Milani, S., Fontani, M., Bestagini, P., Barni, M., Piva, A., Tagliasacchi, M., Tubaro, S., 2012. An overview on video forensics. *APSIPA Trans Signal Inf Process* 1.
- M. Mirza, S. Osindero, Conditional Generative Adversarial Nets, arXiv preprint arXiv:1411.1784.
- Nah, S., Hyun Kim, T., Mu Lee, K., 2017. Deep multi-scale convolutional neural network for dynamic scene deblurring. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3883–3891.
- Newson, A., Almansa, A., Fradet, M., Gousseau, Y., Pérez, P., 2014. Video inpainting of complex scenes. *SIAM J. Imag. Sci.* 7 (4), 1993–2019.
- Niklaus, S., Mai, L., Liu, F., 2017. Video frame interpolation via adaptive separable convolution. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 261–270.
- Pandey, R.C., Singh, S.K., Shukla, K.K., 2016. Passive forensics in image and video using noise features: a review. *Digit. Invest.* 19, 1–28.
- Patwardhan, K.A., Sapiro, G., Bertalmio, M., 2007. Video inpainting under constrained camera motion. *IEEE Trans. Image Process.* 16 (2), 545–553.
- Qadir, G., Yahaya, S., Ho, A.T., 2012. Surrey University Library for Forensic Analysis (Sulfa) of Video Content.
- Qureshi, M.A., Deriche, M., 2015. A bibliography of pixel-based blind image forgery detection techniques. *Signal Process. Image Commun.* 39, 46–74.
- Qureshi, M.A., Deriche, M., Beghdadi, A., Amin, A., 2017. A critical survey of state-of-the-art image inpainting quality assessment metrics. *J. Vis. Commun. Image Represent.* 49, 177–191.
- Rao, Y., Ni, J., 2016. A deep learning approach to detection of splicing and copy-move forgeries in images. In: *Information Forensics and Security (WIFS)*, 2016 IEEE International Workshop on, IEEE, pp. 1–6.
- Reda, F.A., Liu, G., Shih, K.J., Kirby, R., Barker, J., Tarjan, D., Tao, A., Catanzaro, B., 2018. Sdc-net: video prediction using spatially-displaced convolution. In: *The European Conference on Computer Vision (ECCV)*.
- Richter, S.R., Hayder, Z., Koltun, V., 2017. Playing for benchmarks. In: *International Conference on Computer Vision (ICCV)*, vol. 2.
- Rössler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., Nießner, M., 2018. Face-forensics: A Large-Scale Video Dataset for Forgery Detection in Human Faces arXiv preprint arXiv:1803.09179.
- Rota, P., Sangineto, E., Conotter, V., Pramerdorfer, C., 2016. Bad teacher or unruly student: can deep learning say something in image forensics analysis?. In: *Pattern Recognition (ICPR)*, 2016 23rd International Conference on, IEEE, pp. 2503–2508.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al., 2015. Imagenet large scale visual recognition challenge. *Int. J. Comput. Vis.* 115 (3), 211–252.
- Sajjadi, M.S.M., Vemulapalli, R., Brown, M., 2018. Frame-recurrent video super-resolution. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Schaefer, G., Stich, M., 2003. Ucid: an uncompressed color image database. In: *Storage and Retrieval Methods and Applications for Multimedia 2004*, vol. 5307. International Society for Optics and Photonics, pp. 472–481.
- Shih, T.K., Tang, N.C., Tsai, J.C., Hwang, J.-N., 2011. Video motion interpolation for special effect applications. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* 41 (5), 720–732.
- Shullani, D., Fontani, M., Iuliani, M., Al Shaya, O., Piva, A., 2017. Vision: a video and image dataset for source identification. *EURASIP J. Inf. Secur.* 2017 (1), 15.
- Shullani, D., Al Shaya, O., Iuliani, M., Fontani, M., Piva, A., 2017. A dataset for forensic analysis of videos in the wild. In: *International Tyrrhenian Workshop on Digital Communication*. Springer, pp. 84–94.
- Singh, R.D., Aggarwal, N., 2017. Video content authentication techniques: a comprehensive survey. *Multimed. Syst.* 1–30.
- Sitara, K., Mehtre, B.M., 2016. Digital video tampering detection: an overview of passive techniques. *Digit. Invest.* 18, 8–22.
- Sitara, K., Mehtre, B., 2017. A comprehensive approach for exposing inter-frame video forgeries. In: *Signal Processing & its Applications (CSPA)*, 2017 IEEE 13th International Colloquium on, IEEE, pp. 73–78.
- Smith, E., Basharat, A., Anthony Hoogs, C., et al., 2017. A c3d-based convolutional neural network for frame dropping detection in a single video shot. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 86–94.
- Stamm, M.C., Lin, W.S., Liu, K.R., 2012. Temporal forensics and anti-forensics for motion compensated video. *IEEE Trans. Inf. Forensics Secur.* 7 (4), 1315–1329.
- Subramanyam, A., Emmanuel, S., 2012. Video forgery detection using hog features and compression properties. In: *Multimedia Signal Processing (MMSp)*, 2012 IEEE 14th International Workshop on, IEEE, pp. 89–94.
- Sutthiwan, P., Shi, Y.Q., Zhao, H., Ng, T.-T., Su, W., 2011. Markovian rake transform for digital image tampering detection. In: *Transactions on Data Hiding and Multimedia Security VI*. Springer, pp. 1–17.
- Suwajanakorn, S., Seitz, S.M., Kemelmacher-Shlizerman, I., 2017. Synthesizing obama: learning lip sync from audio. *ACM Trans. Graph.* 36 (4), 95.
- Thakur, M.K., Saxena, V., Gupta, J., 2016. Learning based no reference algorithm for dropped frame identification in uncompressed video. In: *Information Systems Design and Intelligent Applications*. Springer, pp. 451–459.
- Thies, J., Zollhöfer, M., Stamminger, M., Theobalt, C., Nießner, M., 2016. Face2face: real-time face capture and reenactment of rgb videos. In: *Computer Vision and Pattern Recognition (CVPR)*, 2016 IEEE Conference on, IEEE, pp. 2387–2395.
- Torralba, A., Efros, A.A., 2011. Unbiased look at dataset bias. In: *Computer Vision and Pattern Recognition (CVPR)*, 2011 IEEE Conference on, IEEE, pp. 1521–1528.
- Tralic, D., Grgic, S., Zovko-Cihlar, B., 2014. Video frame copy-move forgery detection based on cellular automata and local binary patterns. In: *Telecommunications (BIHTEL)*, 2014 X International Symposium on, IEEE, pp. 1–4.
- Tulyakov, S., Liu, M.-Y., Yang, X., Kautz, J., 2018. Mocogan: decomposing motion and content for video generation. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Varol, G., Romero, J., Martin, X., Mahmood, N., Black, M.J., Laptev, I., Schmid, C., 2017. Learning from synthetic humans. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2017). IEEE, pp. 4627–4635.
- Walker, J., Marino, K., Gupta, A., Hebert, M., 2017. The pose knows: video forecasting by generating pose futures. In: *Computer Vision (ICCV)*, 2017 IEEE International Conference on, IEEE, pp. 3352–3361.
- Wang, W., Farid, H., 2007. Exposing digital forgeries in interlaced and deinterlaced video. *IEEE Trans. Inf. Forensics Secur.* 2 (3), 438–449.
- Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P., 2004. Image quality assessment: from error visibility to structural similarity. *IEEE Trans. Image Process.* 13 (4), 600–612.
- Wang, T.-C., Liu, M.-Y., Zhu, J.-Y., Liu, G., Tao, A., Kautz, J., Catanzaro, B., 2018. Video-to-video synthesis. In: *Advances in Neural Information Processing Systems*.
- Wang, W., Alameda-Pineda, X., Xu, D., Fua, P., Ricci, E., Sebe, N., 2018. Every smile is unique: landmark-guided diverse smile generation. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Wexler, Y., Shechtman, E., Irani, M., 2007. Space-time completion of video. *IEEE Trans. Pattern Anal. Mach. Intell.* (3), 463–476.
- Wu, Y., Jiang, X., Sun, T., Wang, W., 2014. Exposing video inter-frame forgery based on velocity field consistency. In: *Acoustics, Speech and Signal Processing (ICASSP)*, 2014 IEEE International Conference on, IEEE, pp. 2674–2678.
- Xia, M., Yang, G., Li, L., Li, R., Sun, X., 2017. Detecting video frame rate up-conversion based on frame-level analysis of average texture variation. *Multimed. Tool. Appl.* 76 (6), 8399–8421.
- Xiong, W., Luo, W., Ma, L., Liu, W., Luo, J., 2018. Learning to generate time-lapse videos using multi-stage dynamic generative adversarial networks. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Xiph.org Video Test Media [derf's Collection]. URL <https://media.xiph.org/video/derf/>.
- Xu, C., Corso, J.J., Libsvox, 2016. A supervoxel library and benchmark for early video processing. *Int. J. Comput. Vis.* 119 (3), 272–290.
- Xu, J., Ni, B., Li, Z., Cheng, S., Yang, X., 2018. Structure preserving video prediction. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Yang, C., Wang, Z., Zhu, X., Huang, C., Shi, J., Lin, D., 2018. Pose guided human video generation. In: *The European Conference on Computer Vision (ECCV)*.
- Yang, R., Xu, M., Wang, Z., Li, T., 2018. Multi-frame quality enhancement for compressed video. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhao, L., Peng, X., Tian, Y., Kapadia, M., Metaxas, D., 2018. Learning to forecast and refine residual motion for image-to-video generation. In: *The European Conference on Computer Vision (ECCV)*.
- Zheng, L., Sun, T., Shi, Y.-Q., 2014. Inter-frame video forgery detection based on block-wise brightness variance descriptor. In: *International Workshop on Digital Watermarking*. Springer, pp. 18–30.
- Zhu, J., Park, T., Isola, P., Efros, A.A., 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In: 2017 IEEE International Conference on Computer Vision (ICCV), pp. 2242–2251. <https://doi.org/10.1109/ICCV.2017.244>.