



Contents lists available at ScienceDirect

Signal Processing: *Image Communication*journal homepage: www.elsevier.com/locate/image

Objective quality assessment in free-viewpoint video production

J. Kilner*, J. Starck, J.Y. Guillemaut, A. Hilton

Centre for Vision, Speech and Signal Processing, Faculty of Engineering and Physical Sciences, University of Surrey, Guildford GU2 7XH, UK

ARTICLE INFO

Article history:

Received 10 October 2008

Accepted 19 October 2008

Keywords:

Free-viewpoint video

Image-based reconstruction

Image-based rendering

Quality assessment

ABSTRACT

This paper addresses the problem of objectively quantifying accuracy in free-viewpoint video production. Free-viewpoint video makes use of geometric scene reconstruction and renders novel views using the appearance sampled in multiple camera images. Previous work typically adopts an objective evaluation of geometric accuracy against ground-truth data or a subjective evaluation of visual quality in view synthesis. We consider two production scenarios, human performance capture in a highly constrained studio environment and sports production in a large-scale external environment. The accuracy of scene reconstruction is typically limited and absolute geometric accuracy does not necessarily reflect the quality of free-viewpoint rendering. A framework is introduced to quantify error at the point of view synthesis. The approach can be applied as a full-reference metric to measure fidelity to a ground-truth image or as a no-reference metric to measure the error in rendering. The framework is applied to a data set with known geometric accuracy and a comparison is presented for studio based and sports production scenarios.

© 2008 Elsevier B.V. All rights reserved.

1. Introduction

Over the past decade multiple-view capture of events has gained increasing interest as a means to create three-dimensional (3D) video content. Application areas range from on-line visualisation for mixed reality environments [2], communications [22], and production or pre-visualisation in television [21], games [47] and 3DTV [50]. In 3DTV applications, cameras are typically arranged with a relatively short baseline to synthesise virtual views directly from the camera images [38]. Free-viewpoint video on the other hand is typically based on a relatively sparse set of cameras that surround a scene and makes use of 3D geometry to synthesise arbitrary viewpoints. A novel view is synthesised as shown in Fig. 1 by rendering a geometric proxy to a virtual viewpoint using the appearance sampled from the adjacent camera images.

Research in free-viewpoint video currently lacks a consistent framework for quality assessment. Previous work in image-based reconstruction for static scenes evaluates geometric accuracy using ground-truth 3D shape [43]. In image-based rendering relatively little work has addressed the accuracy or quality of view synthesis, relying instead on a subjective visual assessment of performance [49] or pixel-wise error metrics with respect to ground-truth images [60]. We find that geometric accuracy does not necessarily reflect the quality of view synthesis and ground-truth 3D shape is often unavailable for the dynamic scenes used in free-viewpoint video. Furthermore, conventional pixel-wise error metrics do not necessarily reflect perceived visual quality and provide no intuition as to the underlying geometric errors that cause artefacts in rendering.

This paper presents a quality assessment framework to quantify the accuracy of free-viewpoint video production. An objective measure of quality is required as a means to optimise the parameters of different algorithms and to benchmark the performance of different production

* Corresponding author. Tel.: +44 1483 689842; fax: +44 1483 686031.
E-mail address: J.Kilner@surrey.ac.uk (J. Kilner).



Fig. 1. Free-viewpoint video synthesis: a geometric proxy is used to reproject the appearance sampled in camera images to a new viewpoint.

techniques. The key contributions in the paper are as follows:

- (i) Two different production scenarios are considered: human performance capture in a highly constrained studio environment and sports production in a large-scale external environment. A taxonomy of errors in free-viewpoint video synthesis is presented.
- (ii) A framework is proposed to objectively quantify errors. Simple objective measures of fidelity in view synthesis are presented that are designed to reflect perceived visual artefacts and to provide a well-understood measure of accuracy.
- (iii) The framework is applied in two cases *full-reference* (FR) against ground truth from a novel viewpoint where available; and *no-reference* (NR) to measure artefacts in free-viewpoint rendering for the general case where ground truth is not available. Results of FR evaluation are presented for both the studio based and sports production scenarios and a NR evaluation is presented for the studio scenario.

Objective quality assessment is restricted here to the problem of evaluating geometric production for free-viewpoint video synthesis. Production techniques should provide high-quality view synthesis independent of coding, transmission and display. This paper is based on the objective evaluation of production for human performance capture [48] and sporting events [27]. In Section 2 previous work on free-viewpoint video production is presented along with approaches to geometric and image-based quality assessment. Section 3 provides an overview of the framework adopted for quality assessment. Section 4 introduces evaluation in the constrained studio environment and Section 5 large-scale external environments for sporting events. An evaluation of different production techniques is then presented in Section 6. Finally, Section 7 draws conclusions and considers future directions.

2. Background

2.1. Free-viewpoint video

In traditional video and film production, an event is recorded from a specific location using a camera. When the video is rendered it provides the fixed viewing experience dictated by the director. Free-viewpoint video attempts to break this restriction by allowing the specification of the camera location at the point of rendering using multiple-video streams recorded from different locations. Research to-date has focused on the multiple camera acquisition systems and the computer vision algorithms required to achieve robust reconstruction and high-quality view synthesis either in real-time or as an off-line post-process [49]. Recent advances have exploited image-based reconstruction and image-based rendering to produce free-viewpoint video at a quality comparable to captured video [62].

Image-based reconstruction deals with the problem of deriving scene geometry from the appearance sampled in camera images. Conventional stereo-based techniques [26] reconstruct a 2.5D depth image representation from two or more cameras through a regularised search for image correspondence. Volumetric techniques allow inference of visibility and integration of appearance across all camera views without image correspondence. Shape-from-silhouette (SFS) techniques [32] derive the *visual-hull*, the maximal volume consistent with foreground silhouettes. This is refined in space-carving techniques [31] which provide the *photo-hull*, the maximal volume that has a consistent appearance across all visible cameras. Multiple shape cues and iterative refinement techniques have been adopted to search for optimal surfaces in a scene [18,24,55].

Image-based rendering is the process of synthesising novel views from camera images. With no geometric scene information, synthesis is performed directly by treating multiple-view images as a set of samples from the light-field in a scene. The light-field is then interpolated and resampled to generate a novel view [33]. In this way, highly realistic view synthesis can be achieved at the cost of a requirement for dense image sampling to avoid interpolation artefacts. Image-based reconstruction and image-based rendering have been combined [12,6] to synthesise novel views from a sparse set of cameras by using scene geometry to provide image correspondence. The advantage of view-dependent rendering is that it can overcome inaccuracies in geometric scene reconstruction by reproducing the change in surface appearance that is sampled in the original camera images.

Free-viewpoint video production within controlled studio environments, which provide fixed illumination, static backgrounds, stationary cameras and a relatively small capture volume, has received considerable interest. The *virtualised reality* system [26] first used 51 cameras distributed over a 5 m dome to capture the performance of an actor in a studio. Fast and robust reconstruction from silhouettes [2,8,21,37] has allowed real-time systems to be developed for mixed reality environments. Off-line systems have the potential for more accurate geometric

scene reconstruction. Model-based techniques make use of a prior humanoid model to track the changes in shape, appearance and pose over time [1,7,46,54]. Data-driven approaches use image-based reconstruction without restriction to a prior model [19,36,47,53]. Video-based rendering blends the appearance sampled in camera images using view-dependent rendering [14,61,57].

In the sports environment, free-viewpoint video has been used to manipulate the viewpoint for playback of sequences within a match. Sports environments such as a soccer stadium require capture of a large volume (50 m × 100 m × 2 m) with severely limited control over the environment. As a result reconstruction must cope with a relatively wide-baseline between cameras, the use of moving broadcast cameras, less accurate calibration and difficulties in player segmentation due to the uncontrolled illumination and backgrounds. Image morphing has been applied for view interpolation with weak camera calibration [11,25], a simplified geometric representation using planar billboards has been proposed for real-time view synthesis [29] and graph cuts [23] and a deformable model [28] have been proposed for high-quality off-line view synthesis.

This field has also been exploited commercially. Eye-vision [16] used a bank of robotic cameras and specialist capture equipment to provide free rotation around a fixed instant in time for use during coverage of the Superbowl. Piero [4] use either hardware or software to calibrate broadcast cameras and generate a 3D billboard representation of the players which then allows limited viewpoint mobility. Most recently, LiberoVision [34] uses a 3D representation of a soccer match to allow a variety of novel camera viewpoints to be rendered. Most techniques either require a high degree of user input or are limited in the range of viewpoints allowed. High-quality view synthesis from wide-baseline camera positions in an uncontrolled external environment remains an open and challenging problem.

2.2. Objective quality assessment

In image-based reconstruction, geometric accuracy has been evaluated using ground-truth 3D shape. Seitz et al. [43] present a comprehensive framework to compare reconstruction techniques against 3D geometry acquired from a laser stripe scanner. Recent work [49] in free-viewpoint production has demonstrated that with current camera hardware, geometric accuracy is inherently insufficient to represent the detailed geometry of a scene and that where display resolution reflects camera resolution, image-based rendering is required to achieve sub-pixel accuracy to minimise visual artefacts in view synthesis. An evaluation of free-viewpoint video should therefore target the accuracy, or quality of view synthesis rather than ground-truth accuracy in geometric reconstruction.

The problem of defining video quality metrics has received significant interest in the image processing community to assess degradations introduced by video

acquisition, processing, coding, transmission and display [41]. There is also active research studying the degradation introduced by watermarking schemes [17]. An overview of the field can be found in Eckert and Bradley [13]. Research into image quality assessment can be broken down into two broad categories, those attempting to model the human visual system (HVS) and those using more direct pixel fidelity criteria.

There has been much work focusing on HVS-based measures of the fidelity of an image. Examples include measuring mutual information in the wavelet domain [45], contrast perception modeling [42] and modeling the contrast gain control of the HVS [58]. However, HVS techniques do not necessarily reflect the true complexity of the visual system and objective measurement of perception remains an open research problem [5,44,59].

Pixel-wise fidelity metrics such as mean square error (MSE) and peak signal to noise ratio (PSNR) remain widely adopted as simple, well-understood measures of fidelity despite a poor correlation with visual quality [56]. An overview of several measures and their performance can be found in Eskicioglu and Fisher [15] while a statistical analysis of various techniques encompassing both pixel metrics and HVS-based metrics can be found in Avcibas et al. [3].

Objective evaluation should ideally provide simple, repeatable quality measures that afford a clear physical interpretation tailored to perceived visual quality. However, the lack of effective standard measures are testament to the difficulties of achieving this in the general case, and so a more fruitful approach is to use domain-specific measures to target the quality assessment of particular types of images.

3. Quality assessment framework

In this section a framework is proposed to quantify accuracy in multiple camera geometric production for free-viewpoint video synthesis. The following definition is provided as the basis for quality assessment.

Free-viewpoint video production should recover a sufficiently accurate 3D scene representation for view synthesis free from visual artefacts. View synthesis should in turn target the resolution of the input camera images such that 3D video provides an acceptable alternative to conventional video.

The HVS is highly adapted to perceive structural and temporal detail in a scene [56] and errors in the visual assessment of free-viewpoint video become apparent where prominent features or geometric structure is incorrectly reproduced. The accuracy of free-viewpoint video is therefore quantified as the structural error in reproducing the appearance of a scene at the point of view synthesis.

Typical image quality assessment metrics either directly measure image fidelity (i.e. the pixel-wise distance between two images) or image quality (i.e. HVS-based systems). A similar classification can be made in the domain of reconstruction quality assessment

between techniques which measures the piece-wise accuracy of the reconstruction against some ground-truth surface (fidelity measures) versus those that measure the overall suitability of a reconstruction for the task of representing the original surface (quality measures). Typical reconstruction quality assessment measures only reconstruction fidelity [43]. However, in some cases reconstruction fidelity has little or no effect on output image quality, which is itself independent of output image accuracy (for example an image shifted exactly one pixel to the right is preferred to an image which has had shot noise added). Our approach is an attempt to avoid simply measuring reconstruction fidelity and instead measure reconstruction quality by measuring the artefacts in the rendered image.

Seitz et al. [43] present a framework to evaluate geometric error in image-based reconstruction with respect to ground-truth 3D geometry. An error metric δ is defined as the distance such that 90% of the reconstructed geometry lies within the distance δ of the ground-truth surface. Fig. 2 illustrates the geometry reconstructed using a free-viewpoint video production technique [47] with an accuracy of $\delta = 1.01$ mm.

We extend this framework to the image domain to evaluate free-viewpoint video. The accuracy of a synthesised image I is quantified as the registration error with respect to a reference image I' . An image I is represented by a function on a set of pixels where each pixel value (typically a three component RGB vector) that makes up the image is obtained by evaluating $I(p)$. We denote those pixels which are in the domain of a specific image with the notation $p \in I$. The error at each pixel is computed as a distance $d(p, I')$ to the reference image. Here for illustration we compute distance using image similarity $S(\cdot)$.

$$d(p, I') = \|p - p'\|_2, \quad \operatorname{argmax}_{p' \in I'} S(I(p), I'(p')) \quad (1)$$

The maximum error in the distribution can be used to characterise the maximum visual artefact that will be apparent, defined by the *Hausdorff distance* from image I to the reference I' :

$$d(I, I') = \max_{p \in I} d(p, I') \quad (2)$$

In practice the Hausdorff metric is sensitive to outliers in the data and the generalised Hausdorff distance is taken

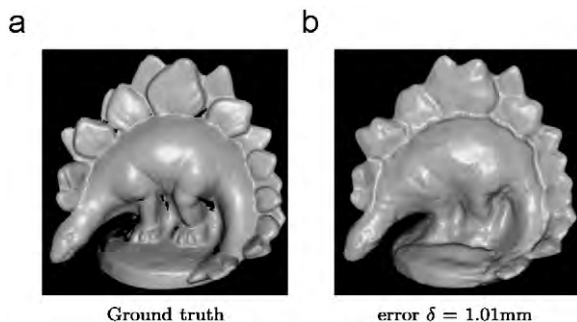


Fig. 2. Geometric evaluation of multiple-view reconstruction courtesy of the Multi-View Stereo Evaluation Homepage (<http://vision.middlebury.edu/mview/>): (a) ground-truth; (b) error $\delta = 1.01$ mm.

as the k th ranked distance in the distribution, where $Q_{x \in X}^k f(x)$ is the quantile of rank k for $f(x)$ over the set X :

$$d^k(I, I') = Q_{p \in I}^k d(p, I') \quad (3)$$

For consistency with [43] we adopt the 90th percentile measure d^{90} .

The error for a synthesised view I is now defined by a single metric $d^k(I, I')$ that quantifies mis-registration between two images. Intuitively the distance measure is related to the geometric error in the underlying geometry of the scene. With a larger error in the 3D geometry of the scene, there will be a shift in reprojected 2D appearance and greater mis-registration.

This framework provides a single intuitive value in terms of pixel accuracy in view synthesis that can be applied at the resolution of the input video images. The approach can be applied as an FR measure of fidelity in aligning structural detail with respect to a ground-truth image, or as a NR metric where mis-registration of structural detail can be evaluated between the appearance sampled from different camera viewpoints in view-dependent rendering.

4. Studio production error metric

Studio-based production provides a highly constrained environment for multiple-view video capture. Studio systems typically consist of between four and 51 cameras [49] configured to surround a specific volume of interest. Fig. 3 illustrates one such production studio. In this environment the camera parameters defining the projective transformation to the image plane can be pre-calibrated to subpixel accuracy [40,51]. The scene can also be recorded using a single ambient lighting environment with a fixed backdrop such as a blue-screen for accurate foreground matting.

In this section the different visual artefacts are categorised and a metric is presented to evaluate the accuracy of view synthesis for the studio environment. The metric is demonstrated against ground-truth geometric accuracy using the 16-view data set courtesy of the multi-view evaluation project [47]. In Section 6 the application to multiple-view video of people is presented where no ground-truth geometry is available.



Fig. 3. A typical image produced in a studio environment.

4.1. A taxonomy of errors

With accurate camera parameters, visual artefacts in view synthesis arise principally from an inexact geometric representation of the scene. Geometric error results in an incorrect projected shape or an incorrect sampled appearance for the scene surface. Reconstruction errors can be categorised as global errors in gross shape or local errors in exact surface position. Examples of typical reconstruction errors are shown in Fig. 4. Large-scale errors are apparent as extraneous or phantom volumes or protrusions that do not correspond to the true surface of the scene. Local errors in surface position give inexact sampling of appearance from the multiple-view video images leading to blurring or ghosting. Errors are summarised in Table 1.

Real-time systems [2,21] typically make use of fast and robust reconstruction from image silhouettes. Accurate foreground mattes can be derived in a studio environment; however, the resulting visual hull representation provides the *maximum* volume that is consistent with the silhouettes. These techniques therefore suffer from gross geometric errors where false positive volumes are derived. Such errors often appear as protrusions particularly when only a few cameras are used in reconstruction or there are multiple occlusions in the scene.

Off-line systems [1,47] on the other hand optimise the scene geometry and often adopt a minimal surface as a regularisation constraint in reconstruction. These techniques tend to suffer less from gross errors in geometry and visual artefacts arise from local errors in surface extraction.

4.2. Registration error metric

Here we focus on evaluating off-line scene reconstruction as the basis for high-quality free-viewpoint video production. Visual artefacts arise where local errors in surface geometry result in incorrectly sampled surface appearance. These artefacts become apparent at the prominent visual features in the scene. In areas of uniform surface appearance inexact geometry will correctly sample a similar appearance across camera viewpoints.

Structural registration error is computed between images using a public domain optic flow algorithm [9] that accounts for the expected image variance in uniform areas of appearance. Registration is derived as the displacement at each pixel $\underline{d}(p, l)$, providing the motion between two images l, l' .

The 90th percentile Hausdorff distance d^{90} now provides a single measure of the maximum error at distinct image regions, where the effect of geometric error is apparent in view synthesis.

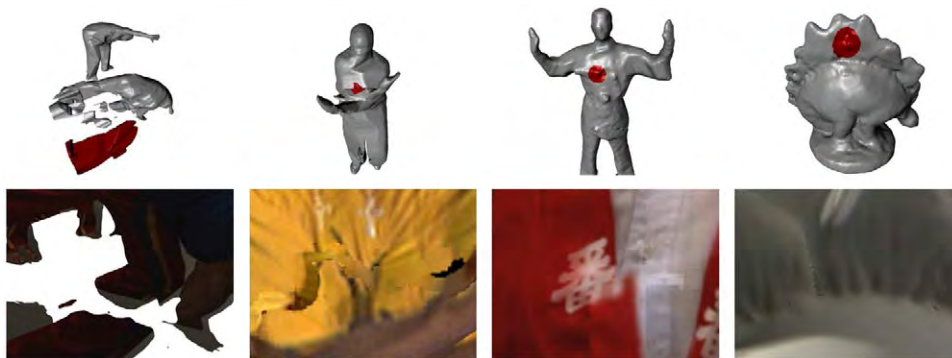


Fig. 4. Gross and local errors. The top row shows the error in reconstruction highlighted in red while the bottom row shows the corresponding artefact in the rendered view. From left to right, phantom volume, phantom protrusion, a raised surface area and blurring. Capoeira and Kimono data sets (columns 1 and 3) courtesy of Matsuyama lab., Kyoto University, Japan [36]. Dancer data set (column 2) courtesy of the public domain multiple camera data set [47]. Dinosaur model (column 4) courtesy of the Multi-View Stereo Evaluation Homepage (<http://vision.middlebury.edu/mview>).

Table 1

Classification of errors in foreground synthesis.

Type	Error	Cause	Description
Gross errors	Phantom volume	Ambiguity due to occlusions	The visual hull contains disconnected components due to the ambiguity in shape caused by occlusions
	Phantom protrusion	Ambiguity due to occlusions	The visual hull is extended into a region where no genuine surface exists This is often caused by self-occlusion e.g. in the region between the legs
	Holes	Incorrect matting	Areas incorrectly marked as background will cause the space carving algorithm to generate a visual hull with holes in it where none existed in the original shape
Local errors	Sunken or raised surface	Stereo ambiguity	Lack of structure or repeating patterns can cause an ambiguity in the stereo scores for a region. As a result stereo based optimisation cannot determine the true depth of the surface
	Blurring	Quantisation	Fundamental limits to precision of surface placement and resolution of input images lead to blurring when multiple images are combined

4.3. FR evaluation

An FR comparison makes use of a ground-truth reference for a frame-by-frame evaluation of the accuracy in view synthesis. This is illustrated using a leave-one-out test where the reconstructed geometry shown in Fig. 2(b) is used to synthesise a view that is excluded in synthesis. A comparison is made between the registration error d^{90} , the RMSE registration error as well as PSNR. The comparison is made for varying degrees of geometric error introduced by inflating the surface by 1 and 2 mm beyond the known geometric error $\delta = 1.01$ mm [47].

Fig. 5 shows the synthesised views compared to the reference image. The corresponding error metrics for the leave-one-out tests against the ground-truth image are presented in Table 2. As the geometric accuracy is reduced, double exposure effects can be observed in structured image regions such as shadow boundaries. The error metrics follow this subjective decrease in image quality with a reduced PSNR, an increased RMSE and an

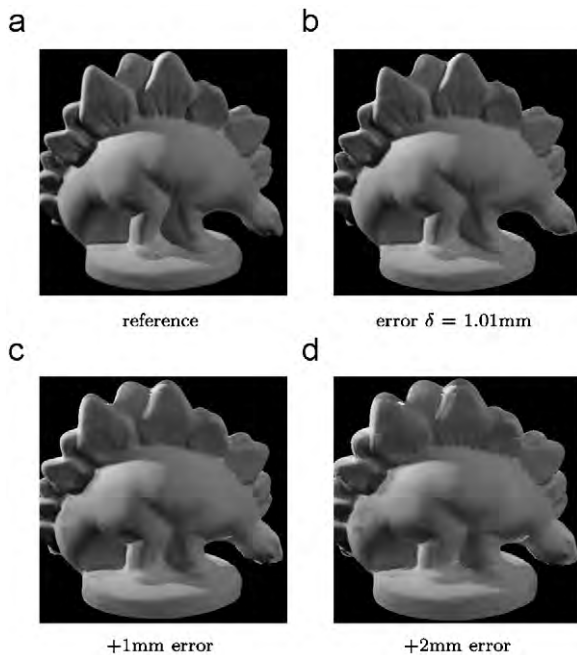


Fig. 5. Synthesised views in a leave-one-out test in comparison to (a) the reference image, with (b) baseline error $\delta = 1.01$ mm and additional geometric error (c) +1 mm (d) +2mm introduced in order to generate typical mis-registration artefacts such as blurring and double images.

Table 2

Error metrics for FR comparison in a leave-one-out test with varying degrees of additional geometric error.

Geometric error (mm)	d^{90}	RMSE	PSNR (dB)
+0	0.70	0.65	31.50
+1	2.12	1.44	24.40
+2	3.81	2.36	21.90

Note that the RMSE presented in these tables is derived from the registration error whilst PSNR is derived from pixel value differences.

Table 3

Error metrics for NR comparison where the reprojected appearance from two camera images is compared.

Geometric error (mm)	d^{90}	RMSE	PSNR (dB)
+0	1.09	0.76	31.60
+1	2.22	1.24	24.90
+2	3.74	2.13	21.70

increase in the generalised maximum error d^{90} . Both the RMSE and d^{90} metrics reflect the change in geometric scene accuracy. The d^{90} measure also provides intuition as to the maximum error that is apparent in Fig. 5 as the numerical value indicates the expected scale of mis-registration errors.

4.4. No-reference (NR)

An NR comparison requires no explicit ground truth. In view synthesis the appearance of a 3D scene is sampled in two or more camera images and reprojected to a new view. In the absence of a ground-truth reference the reprojected appearance from different cameras can be compared directly. This is illustrated using a virtual viewpoint placed at the mid-point between two cameras in the 16-view data set. From this viewpoint, two images are rendered, each using one of the original cameras. These images are then compared against each other to produce the results shown in Table 3. These results demonstrate that the error metric values obtained in the NR case are strongly correlated with those in the FR case with known ground truth.

The NR comparison provides a measure of the potential artefacts in view synthesis without the requirement for a ground-truth image. Note that visual artefacts are observed where the reprojected appearance is mis-registered in Fig. 5(c) and (d). The generalised maximum error d^{90} provides a metric for the maximum apparent error which mirrors the FR metrics shown in Table 2.

4.5. Summary

A simple metric has been presented to quantify structural errors in view synthesis. Conventional error metrics such as the *root mean square error* (RMSE) can be adopted. However, a simple mean across an entire image can mask visually distinct errors in highly structured image regions. For example where there are extended areas of uniform appearance the RMSE will naturally tend to zero as the images I, I' are similar. Here the metric targets distinct image regions and features where geometric errors become visually apparent. The technique is relatively simple to implement using public domain software and can be used to compare a synthesised view to a ground-truth image for an FR evaluation, or to compare the appearance sampled from different camera images in a virtual viewpoint as an NR evaluation.

5. Sports production error metric

The sports production environment typically consists of a set of cameras (up to 30) arranged at various mounting points around a stadium. These cameras are either static “locked off” cameras or operator controlled cameras. The cameras are connected to an editing suite in a mobile control unit. The environment is unconstrained in terms of lighting (which may be natural or floodlit), camera calibration (cameras are often hired and cameras are moving constantly during the game) and backgrounds (backgrounds are natural, noisy and temporally varying). In addition, moving broadcast and fixed cameras capture the scene at different resolutions and levels of motion blur. A typical image captured from a sporting event is shown in Fig. 6.

As such, attempting to use free-viewpoint video techniques developed for studio use results in very inaccurate reconstructions [27]. Errors in calibration and matting quickly swamp small ambiguities in shape and stereo-based optimisation becomes unfeasible as the baseline between cameras is widened.

5.1. A taxonomy of errors

When reconstructing an outdoor sporting event, the magnitude of errors increases and the output resolution decreases dramatically such that the distinction between gross and local errors becomes redundant. All errors, whether they are caused by ambiguity in the input data or by inaccuracies in the calibration and matting, cause significant deviation of the reconstructed shape from the true shape of the scene. As such a different and more fundamental taxonomy of errors is employed to characterise the performance of a reconstruction in this environment.

Free-viewpoint video in the sports environment suffers from input images with a greatly reduced resolution compared to the studio environment. Similarly the output images generated are at a much lower resolution and the foreground elements account for a much lower percentage



Fig. 6. A typical image captured during a soccer match broadcast.

of the image pixels. This renders the techniques developed for studio use inappropriate for use in the sports environment. Also the lack of overlap between the wide-baseline images coupled with their varying surface sampling rates (some images may zoom in on one or two players while others may cover half the field) and differences in motion blur, render an NR metric impractical. However, similar principles may be used to derive an FR metric for evaluation in this environment.

The errors inherent in free-viewpoint video synthesis in the sports environment can be classified by considering the ways in which a synthetic video sequence generated for a given viewpoint can differ from the video sequence that would have been captured by a real camera placed at that viewpoint.

5.2. Errors in shape and appearance

Errors in shape are errors where I is missing a foreground element that is present in I' , or I contains an extraneous foreground element that was not present in I' . Examples are missing limbs or double images as shown in Fig. 7(a) and (b). In both cases pixels in I have been incorrectly synthesised as either foreground or background.

Errors in appearance occur when a region $R \in I$ contains different pixel values to $R' \in I'$. This can occur through the rendering of the incorrect surface or due to incorrect sampling of the input image. Examples are the rendering of surfaces in incorrect locations and blurred rendering of surfaces, as shown in Fig. 7(c) and (d).

Table 4 summarises these classifications. Through the distinction between foreground and background (and hence between shape and appearance), a measurement of image fidelity can be composed that conveys more information about the errors in the synthesis pipeline. If the entirety of I' is treated as foreground then this analysis reduces to a comparison of pixel values across the image.

5.3. Registration error metric

We now derive metrics to measure errors in appearance and shape as described above. We also derive a completeness metric to measure only missing foreground, as this allows us to perform meaningful analysis in the presence of matting errors which may introduce extraneous elements of the scene into the reconstruction.

An r -shuffle is a perturbation of an image such that if I is an r -shuffle of I' then every pixel $p' \in I'$ will be transformed to a pixel $p \in I$ such that $\|p' - p\|_2 < r$ [30]. Due to the accumulation of errors in the reconstruction process the generated image I will be a distortion of I' . By modelling this distortion as an r -shuffle, the accumulation of errors in the view synthesis pipeline can be accounted for, allowing an assessment of the fidelity of the underlying reconstruction technique. This measure is similar to the measure d used in the studio evaluation; however, computation of d requires some image registration technique such as optic flow to be performed reliably on

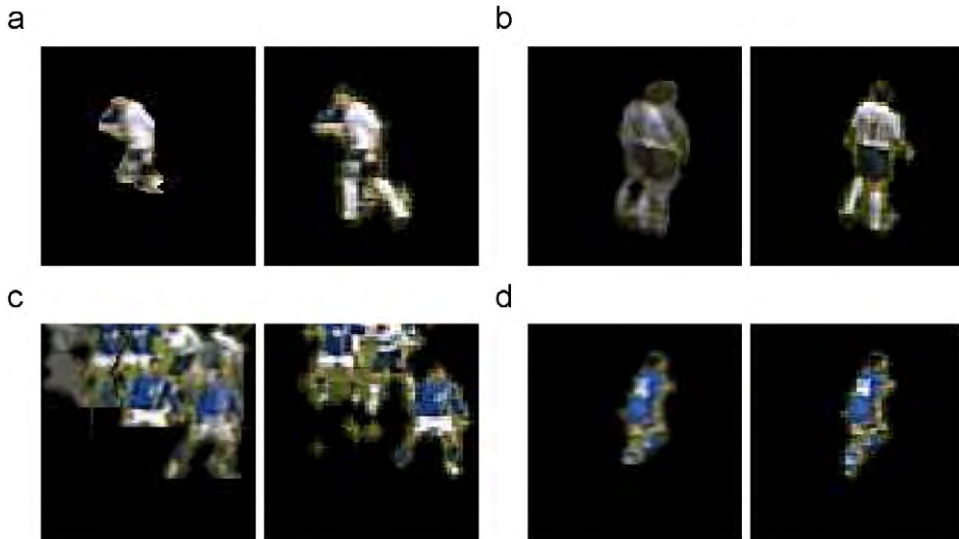


Fig. 7. A comparison of synthetic images to their corresponding ground truths. In each pair the synthetic image is on the left and the ground truth on the right (a) shows an incomplete synthetic image, (b) shows a synthetic image where the player is incorrectly rendered twice, (c) shows a player incorrectly rendered to a foreground region and (d) shows a blurred player.

Table 4

Classification of errors in foreground synthesis comparing a region of the ground-truth image R' to a corresponding region of the synthetic image R .

Error	Image in R'	Image in R	Classification
Missing foreground	Present	Absent	Error in shape
Extraneous foreground	Absent	Present	Error in shape
None	Present	Present	Correct shape
Incorrect image	Image of β	Image of α	Error in appearance
Distorted image	Distorted image of α	Image of α	Error in appearance
None	Image of α	Image of α	Correct view synthesis

α and β denote different elements within the scene.

the images, whereas r simply requires an estimate of the mean pixel error in the system which can be inferred directly from known camera calibration and matting errors.

The r -neighbourhood N_r of a pixel p on the image I is defined such that for some other pixel q

$$q \in N_r(p) \iff \|q - p\|_2 < r \quad (4)$$

We then define the pixel-wise shape matching function s in terms of the function $F(v)$ which returns 1 if pixel value v is foreground and 0 otherwise:

$$s(p, I, I') = F(I(p)) \max_{q \in N_r(p)} (F(I'(q))) \quad (5)$$

Summing over the entire image and normalising by the number of foreground pixels gives the shape score. When

r is taken as 0 this is simply the area of the intersection of both foreground regions divided by the area of their union:

$$\text{shape}(I, I') = \frac{\sum_p s(p, I, I')}{\sum_p \max(F(I(p)), F(I'(p)))} \quad (6)$$

The pixel-wise completeness function is defined similarly to s except that it does not penalise for extraneous foreground:

$$c(p, I, I') = \max_{q \in N_r(p)} (\max(F(I(p))F(I'(q)), 1 - F(I'(q)))) \quad (7)$$

Summing over the entire image and normalising by the number of foreground pixels gives the completeness score:

$$\text{comp}(I, I') = \frac{\sum_p c(p, I, I')}{\sum_p \max(F(I(p)), F(I'(p)))} \quad (8)$$

The pixel-wise appearance matching function a can be defined:

$$a(p, I, I') = \max_{q \in N_r(p)} (T(\|I(p) - I'(q)\|_2)) \quad (9)$$

where $T(x)$ returns 1 if $x \geq \tau$ and 0 otherwise, and τ is some chosen small threshold. Summing over the entire image and normalising by the common foreground region gives the appearance score:

$$\text{app}(I, I') = \frac{\sum_p a(p, I, I')s(p, I, I')}{\sum_p F(I(p))F(I'(p))} \quad (10)$$

These measures are compared against the PSNR which is given by

$$\text{PSNR}(I, I') = 20 \log_{10} \left(\frac{K\sqrt{n}}{\sqrt{\sum_p \|I(p) - I'(p)\|_2^2}} \right) \quad (11)$$

where K is the maximum value that can be given by $\|I(p) - I'(p)\|_2^2$ and n is the number of pixels in the image.

These relationship between r and the d^{90} metric can now be seen. Finding a value of r which yields an appearance score of 0.9 is analogous to the d^{90} score used in the studio evaluation as both identify the distance within which 90% of pixels can be matched to a corresponding pixel in the comparison image. However, the size and distribution of the errors in images from the sports outdoor broadcast environment render this measure less useful than it is in the studio environment. Also the wide baseline and small size of the players makes the optic flow algorithm, upon which the d^{90} measure is based, unstable. Plotting the shape, appearance and completeness scores against the size of the r -shuffle yields more insight as to the nature of the errors present in the reconstruction.

5.4. FR evaluation

Table 5 shows a comparison of these derived measures against the visual information fidelity (VIF) measure [45] of visual quality in an image which was chosen as a baseline FR quality metric. The comparison was carried out on several test images, some consisting of filtered versions of an original image and others on reconstructions of the scene using billboards [29], visual hull [52] and the view-dependent visual hull (VDVH) [39] (a fuller description of the experimental setup for reconstruction is provided in Section 6.2). It can be seen that all measures are in broad agreement, justifying the use of the proposed measures in the comparison. However, the shape, completeness and appearance measures give more information than VIF or PSNR as to the nature of the reconstruction. For example all scores agree that the blur transformation degrades the image more than the median filter, but the shape and completeness scores correctly indicate that the median filter preserves the shape of the image while the blur does not, and that the distortion caused by blurring is an expansion (shape < 1 is but comp \approx 1).

Table 5

Comparison of evaluation techniques of a single reconstructed or processed image vs. a ground-truth image.

Image	Shape	Comp	App	PSNR	VIF
Original	1	1	1	Inf	1
Median	0.99	1	0.98	17.35	0.36
Blur	0.75	0.99	0.96	16.64	0.32
Visual hull	0.86	0.87	0.99	14.02	0.22
Visual hull-1	0.84	0.88	0.95	11.91	0.15
Billboards	0.81	0.89	0.95	11.03	0.17
Billboards-1	0.70	0.89	0.86	8.94	0.08
VDVH	0.56	0.58	0.98	8.88	0.07
VDVH-1	0.56	0.58	0.95	8.71	0.06
Blank	0	0	0	6.00	0

The suffix "-1" refers to a reconstruction generated from the "leave one out" data set.

5.5. Summary

A methodology has been presented to quantitatively evaluate free-viewpoint video production in the sports environment. A set of metrics are introduced that measure errors in shape, completeness and appearance of view synthesis. The measures have been compared to PSNR and VIF as benchmark NR error metrics. The techniques can be used to gain more information as to the nature of errors in foreground view synthesis and are suitable for the low resolution images generated in the sports environment. The techniques are simple to implement and provide an intuitive metric for quality assessment.

6. Results

This section will now evaluate the framework and metrics previously introduced by applying them to two production environments. Firstly the studio-based metrics will be applied to capture from a studio environment and the sports-based metrics are applied to video from a soccer match.

6.1. Studio production

In free-viewpoint video production sparse camera sets are typical and additional camera views are not necessarily available for an FR quality assessment. An NR comparison is now presented to evaluate two free-viewpoint video production techniques.

6.1.1. Capture environment

Two data sets are considered, the first courtesy of [47] consists of a street-dancer performing fast acrobatic motions wearing everyday clothing recorded from 8, 1920 \times 1080 resolution cameras, the second courtesy of [36] consists of a Maiko wearing a brightly coloured Kimono performing a slow dance recorded from 16, 1024 \times 768 resolution cameras. Figs. 8 and 9 illustrate the 3D geometry recovered using a surface optimisation technique [36] with a computational cost of 1 min/frame on an Intel(R) Xeon(TM) 3.6GHz CPU and a global optimisation technique [47] with a cost of 38 min/frame on an Intel(R) Xeon(TM) 3 GHz CPU.

6.1.2. Reconstruction techniques

Free-viewpoint video in the studio environment requires high quality reconstruction techniques that work with well-calibrated high-quality cameras arranged along a wide baseline. The techniques compared in the studio environment are a deformable 3D mesh model and a graph cut optimised surface. Both approaches use image segmentation techniques to generate silhouette images and then use SFS to calculate the visual hull.

The deformable 3D mesh model described in Matsuyama et al. [36] then uses a temporo-spatial deformable model to optimise the surface for photo-consistency, silhouette fitting, smoothness, 3D motion flow and inertia.

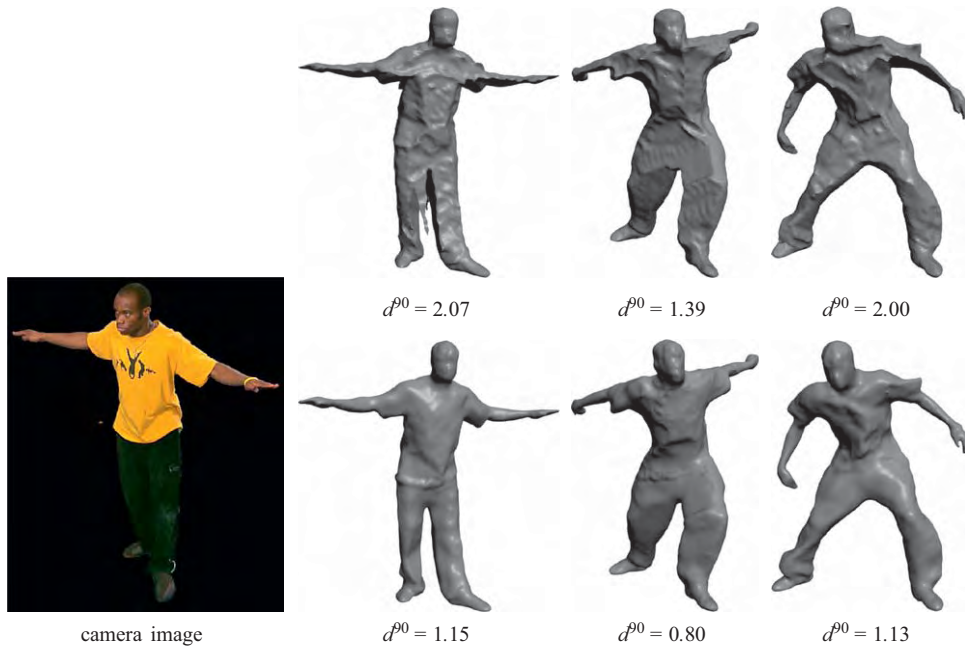


Fig. 8. NR evaluation of two techniques (top) [36], (bottom) [47] for the street dancer sequence courtesy of [47].

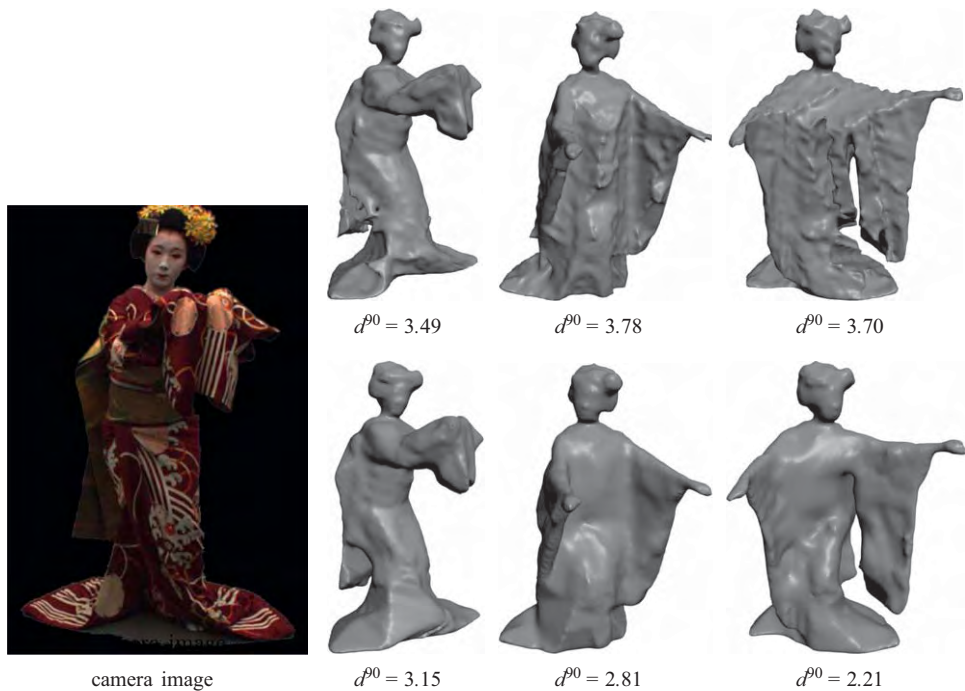


Fig. 9. NR evaluation of two techniques (top) [36], (bottom) [47] for the kimono sequence courtesy of [36].

The graph cut optimised surface described in Starck et al. [47] uses a max-flow min-cut algorithm to find a global optimum surface that maximises stereo correspondence while simultaneously fulfilling silhouette constraints and minimising intra-frame distances.

6.1.3. Evaluation

The NR comparison is performed using a virtual viewpoint placed at the mid-point between a set of cameras in the studio, a pair of cameras for the planar 8 camera setup and a camera triplet for the non-planar 16 camera setup.

6.1.4. Discussion

The d^{90} metric given in Figs. 8 and 9 provides an objective comparison of the quality of view synthesis in the virtual viewpoint. The relatively computationally expensive global optimisation approach [47] gives lower d^{90} rendering error. This is expected as the global optimisation, which includes a stereo correspondence term, is expected to increase the geometric accuracy of reconstruction. Note that the errors for the kimono sequence are higher as the appearance in the scene is more highly structured. The metric will be both camera configuration and data set dependent.

6.2. Sports production

6.2.1. Capture environment

The data set chosen for this evaluation was a recording of a soccer match. The recording was made with 15 static cameras arranged around 120° of the stadium [20]. The cameras were interlaced PAL broadcast cameras captured at a resolution of 702 × 288 (one deinterlaced widescreen field). The camera configuration is illustrated in Fig. 10. This data set was chosen as stable calibration data is available for the cameras, and the arrangement of the cameras allows “leave one out” and “leave three out” comparisons without excessively reducing the quality of the reconstruction. This allows consideration of the behaviour of each technique as information becomes more limited.

6.2.2. Reconstruction techniques

Different techniques for scene reconstruction have different overheads and trade-offs in terms of quality and fidelity. For the proposed application of sporting event reconstruction, real-time playback capabilities along with the ability to work on data from sparse viewpoints are important properties. This has led to a number of alternatives considered here: billboards [29], visual hull [52] and the VDVH [39].



Fig. 10. Arrangement of cameras in soccer stadium. Camera highlighted in red provided ground-truth images and was not used for “leave one out” experiment and cameras highlighted in blue were additionally removed for “leave three out” experiment.

In billboarding a single polygon is rotated around the Y-axis so that it is constantly facing the virtual camera. An image of the original object is then applied to the polygon as a texture map.

A volumetric SFS approach is adopted [52] that divides space into a voxel grid and back-projects every voxel to each image. By comparing the overlap of the voxel to the silhouettes for each camera you can determine the voxel's occupancy. Overlap can be tested up to a given reprojection error to account for calibration errors in the camera system. The surfaces generated by the voxels are then triangulated using the marching cubes algorithm [35] to produce a mesh.

These techniques are compared to the VDVH [39] derives a view-dependent 2.5D depth representation for the visual hull with respect to a given viewpoint. Surface geometry is derived in the image domain by reprojecting a ray from this given viewpoint and deriving the exact intersection with each image silhouette to provide a depth-per-pixel.

6.2.3. Evaluation

The techniques of billboarding, visual hull, and VDVH were then compared. A single renderer capable of rendering all the scene reconstructions was created. This avoided differences in camera representation or lighting creating significant variation in global error between the techniques. A sequence of 100 frames of video from the 15-camera recording was matted using Bayesian matting [10] and 100 frames of synthetic video produced from the viewpoint of camera 5. In the first experiment, all available data sets were used, in the second, camera 5 was omitted from the data set and in the third, cameras 4–6 were omitted (as shown in Fig. 10). The synthetic video streams were then compared with the ground-truth video stream using the technique described in Section 5. It should be noted that the ground truth is also segmented using Bayesian matting and as such is not perfect. However, any gross errors in the matting were corrected and the remaining errors are small compared to the reconstruction errors (Figs. 11 and 12).

6.2.4. Discussion

The left hand column of Fig. 13 shows the behaviour of the shape scores for each of the techniques as the magnitude of the estimated system error is increased.

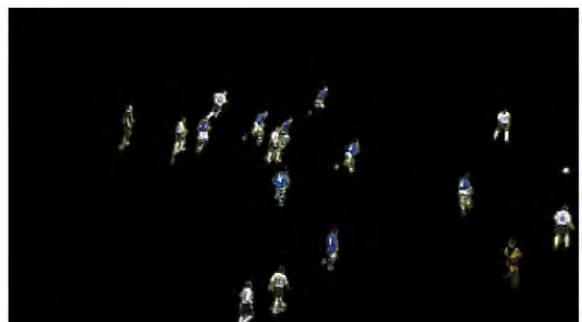


Fig. 11. Ground-truth image from camera 5.

The shape graph for the “leave none out” test shows that none of the techniques achieve a score of 1, even when compensating for large system errors. Some of this is due to errors in the original matting.

The VDVH suffers quite extensively from missing and clipped players. These problems are not exacerbated by the removal of input cameras from the system, hence it does not degrade as much as the billboard or the visual hull. Billboards, being the simplest geometrical representation of the scene, degrade most as cameras are removed

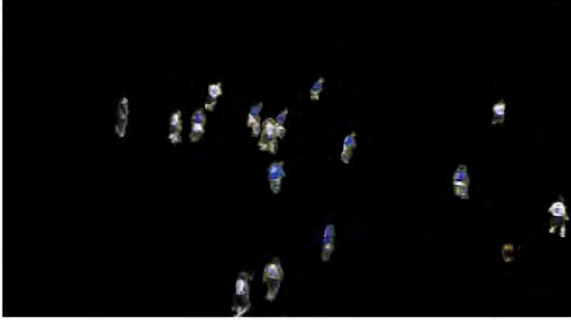


Fig. 12. Visual hull reconstruction not using camera 5 with compensation for a silhouette error of 3 pixels.

from the system as they provide the worst interpolation between views.

The central column of Fig. 13 shows the completeness scores for the techniques. This graph clearly shows the difference between both the billboard and visual hull techniques which are generally overcomplete, and the VDVH which is generally incomplete. This can be interpreted in terms of the requirement for agreement between cameras. In the VDVH agreement is required between silhouette projections from all cameras to generate a surface and so errors in calibration and matting erode the shape. The visual hull and billboard techniques used for these experiments both have some tolerance for disagreement between cameras and so errors do not erode the shape as badly. However, the cost for this is an increase in false positives which can adversely affect the shape through the generation of phantom volumes, as seen in the shape scores in the “leave three out” experiment.

The right hand column of Fig. 13 shows the behaviour of the appearance scores for each of the techniques as the magnitude of the estimated system error is increased. When a system error of 0 pixels is estimated, the appearance score for all techniques is significantly decreased. This is due to resampling errors when the camera images are converted for use in the renderer.

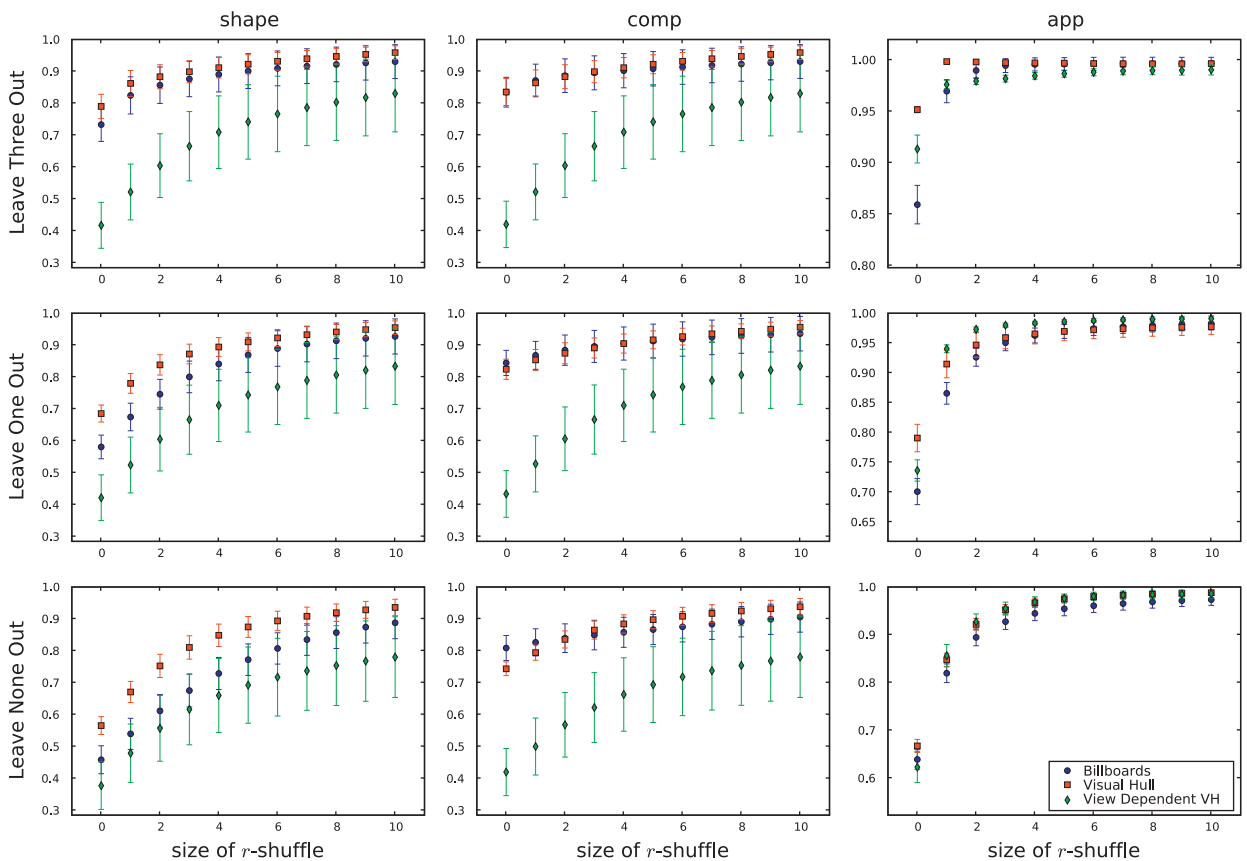


Fig. 13. Plots showing the mean and standard deviation over a sequence of 100 frames for shape, completeness and appearance scores vs. the size of the r -shuffle.

It should be noted that even in this worst case “leave three out” test, appearance is generally synthesised with high fidelity compared to shape. The performance of the billboarding technique, which ignores small-scale surface shape, indicates that the pressing problem with all current techniques is one of generating a scene reconstruction that is accurate and complete at the large scale; small-scale parallax having little effect on the appearance.

This conclusion is further supported by the shape of the graphs. The “dog-leg” shape of the appearance graph implies that a large number of pixels are rendered to within a small distance of their correct locations. This is consistent with small displacements of areas of correct view synthesis. However, the shape and completeness graphs show a much smoother gradient as the size of the r -shuffle is increased, implying that errors in shape cannot be accounted for by a simple displacement and that more serious distortions such as truncations have occurred.

7. Conclusion

An objective error measure for evaluation of free-viewpoint rendering has been presented based on the registration error with respect to either ground-truth reference images or with respect to appearance sampled from multiple input images. This allows the calculation of either FR or NR image synthesis quality assessment. Results demonstrate that the proposed error measures correlate strongly with the error in scene reconstruction and give an improved estimate of image quality over previously used RMSE and PSNR. Results show that the NR metric gives error values in agreement with those obtained from the FR metric.

Objective quality evaluation is performed for two free-viewpoint video scenarios: studio capture in a controlled environment; and sports production in an uncontrolled stadium environment with unconstrained illumination, natural backgrounds and moving cameras. Results demonstrate that the proposed error measures are able to quantify perceived visual artefacts and give a good indication of relative algorithm performance.

Acknowledgements

This work was supported by the DTI Technology programme under Free-viewpoint video for interactive entertainment production TP/3/DSM/6/1/15515 and EPSRC Grant EP/D033926, lead by BBC Research and Development. The authors gratefully acknowledge the project partners for providing the sports footage and discussion of the research reported in this paper. For further details visit the iview project (<http://www.bbc.co.uk/rd/iview>).

References

- [1] E. Aguiar, C. Stoll, C. Theobalt, N. Ahmed, H.-P. Seidel, S. Thrun, Performance capture from sparse multi-view video, *ACM Trans. Graphics (SIGGRAPH)* (2008), in press.
- [2] J. Allard, J.-S. Franco, C. Menier, E. Boyer, B. Raffin, The grimage platform: a mixed reality environment for interactions, in: *IEEE International Conference on Computer Vision Systems (ICVS)*, 2006, p. 46.
- [3] I. Avciabas, B. Sankur, K. Sayood, Statistical evaluation of image quality measures, *J. Electron. Imaging* 11 (2) (2002) 206–223.
- [4] BBC, The Piero System, BBC Production Magic (<http://www.bbc.co.uk/rd/projects/virtual/piero>).
- [5] T. Brandao, P. Queluz, Towards objective metrics for blind assessment of images quality, in: 2006 IEEE International Conference on Image Processing, 2006, pp. 2933–2936.
- [6] C. Buehler, M. Bosse, L. McMillan, S. Gortler, M. Cohen, Unstructured lumigraph rendering, *ACM Trans. Graphics (SIGGRAPH)* (2001) 425–432.
- [7] J. Carranza, C.M. Theobalt, M. Magnor, H. Seidel, Free-viewpoint video of human actors, *ACM Trans. Graphics (SIGGRAPH)* 22 (3) (2003) 569–577.
- [8] G. Cheung, T. Kanade, J. Bouguet, M. Holler, A real time system for robust 3D voxel reconstruction of human motions, in: *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2000, pp. 714–720.
- [9] W. Christmas, Filtering requirements for gradient-based optical flow measurement, *IEEE Trans. Image Process.* 9 (10) (2000) 1817–1820.
- [10] Y. Chuang, B. Curless, D. Salesin, R. Szeliski, A bayesian approach to digital matting, in: *Proceedings of IEEE CVPR 2001*, vol. 2, 2001, pp. 264–271.
- [11] K. Connor, I. Reid, A multiple view layered representation for dynamic novel view synthesis, in: *Proceedings of the 14th British Machine Vision Conference (BMVC)*, 2000.
- [12] P. Debevec, Y. Yu, G. Borshukov, Efficient view-dependent image-based rendering with projective texture-mapping, in: *Proceedings of Eurographics Workshop on Rendering*, 1998, pp. 105–116.
- [13] M.P. Eckert, A.P. Bradley, Perceptual quality metrics applied to still image compression, *Signal Process.* 70 (3) (1998) 177–200.
- [14] M. Eisemann, B. Decker, M. Magnor, P. Bekaert, E. Aguiar, N. Ahmed, C. Theobalt, A. Sellent, Floating textures, in: *Computer Graphics Forum*, *Proceedings of the Eurographics EG'08*, vol. 27(2), 2008, pp. 409–418.
- [15] A. Eskicioglu, P. Fisher, Image quality measures and their performance, *IEEE Trans. Comm.* 43 (12) (1995) 2959–2965.
- [16] Eye-Vision, Carnegie Mellon goes to the Super Bowl (<http://www.ri.cmu.edu/events/sb35/tksuperbowl.html>).
- [17] M.C.Q. Farias, S.K. Mitra, M. Carli, A. Neri, A comparison between an objective quality measure and the mean annoyance values of watermarked videos, in: *Proceedings of the IEEE International Conference on Image Processing*, 2002, pp. 469–472.
- [18] Y. Furukawa, J. Ponce, Carved visual hulls for image-based modeling, *Internat. J. Comput. Vision* (2008), doi: 10.1007/s11263-008-0134-8.
- [19] B. Goldluecke, M. Magnor, Space-time isosurface evolution for temporally coherent 3D reconstruction, in: *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 1, 2004, pp. 350–355.
- [20] O. Grau, A. Hilton, J. Kilner, G. Miller, T. Sargeant, J. Starck, A free-viewpoint video system for visualisation of sport scenes, *SMPTE Motion Imaging* (2007) 213–219.
- [21] O. Grau, T. Pullen, G. Thomas, A combined studio production system for 3D capturing of live action and immersive actor feedback, *IEEE Trans. Circuits Syst. Video Technol.* 14 (3) (2003) 370–380.
- [22] M. Gross, S. Würmlin, M. Naef, E. Lamboray, C. Spagno, A. Kunz, E. Koller-Meier, T. Svoboda, L. Van Gool, S. Lang, K. Strehlke, A. Vande Moere, O. Staadt, blue-c: a spatially immersive display and 3d video portal for telepresence, *ACM Trans. Graphics (SIGGRAPH)* 22 (3) (2003) 819–827.
- [23] J.-Y. Guillemaut, A. Hilton, J. Starck, J. Kilner, O. Grau, A bayesian framework for simultaneous matting and 3d reconstruction, in: *3DIM '07: Proceedings of the Sixth International Conference on 3-D Digital Imaging and Modeling*, 2007, pp. 167–176.
- [24] C. Hernandez, F. Schmitt, Silhouette and stereo fusion for 3D object modeling, *Comput. Vision Image Understanding* 96 (3) (2004) 367–392.
- [25] N. Inamoto, H. Saito, Arbitrary viewpoint observation for soccer match video, in: *The 1st European Conference on Visual Media Production (CVMP)*, 2004, pp. 21–30.
- [26] T. Kanade, P. Rander, P. Narayanan, Virtualized reality: constructing virtual worlds from real scenes, *IEEE Multimedia* 4 (1) (1997) 34–47.
- [27] J. Kilner, J. Starck, A. Hilton, A comparative study of free-viewpoint video techniques for sports events, in: *European Conference on Visual Media Production*, 2006, pp. 87–96.

- [28] J. Kilner, J. Starck, A. Hilton, O. Grau, Dual-mode deformable models for free-viewpoint video of sports events, in: 3DIM '07: Proceedings of the Sixth International Conference on 3-D Digital Imaging and Modeling, 2007, pp. 177–184.
- [29] T. Koyama, I. Kitahara, Y. Ohta, Live mixed-reality 3d video in soccer stadium, in: The 2nd IEEE and ACM International Symposium on Mixed and Augmented Reality, 2003, pp. 178–186.
- [30] K. Kutulakos, Approximate n-view stereo, *ECCV* 1 (2000) 67–83.
- [31] K. Kutulakos, S. Seitz, A theory of shape by space carving, *Internat. J. Comput. Vision* 38 (3) (2000) 199–218.
- [32] A. Laurentini, The visual hull concept for silhouette based image understanding, *IEEE Trans. Pattern Anal. Mach. Intell.* 16 (2) (1994) 150–162.
- [33] M. Levoy, P. Hanrahan, Light field rendering, *ACM Trans. Graphics (SIGGRAPH)* 30 (1996) 31–42.
- [34] LiberoVision, LiberoVision GmbH website (<http://www.liberovision.com/>).
- [35] W. Lorenson, H. Cline, Marching cubes: a high resolution 3d surface construction algorithm, *Comput. Graphics* 21 (4) (1987) 163–169.
- [36] T. Matsuyama, X. Wu, T. Takai, S. Nobuhara, Real-time 3D shape reconstruction, dynamic 3D mesh deformation, and high fidelity visualization for 3D video, *Comput. Vision Image Understanding* 96 (3) (2004) 393–434.
- [37] W. Matusik, C. Buehler, L. Mcmillan, Polyhedral visual hulls for real-time rendering, in: Proceedings of Eurographics Workshop on Rendering, 2001, pp. 115–126.
- [38] W. Matusik, H. Pfister, 3D TV: a scalable system for real-time acquisition, transmission, and autostereoscopic display of dynamic scenes, *ACM Trans. Graphics (SIGGRAPH)* (2004) 814–824.
- [39] G. Miller, A. Hilton, Exact view-dependent visual hulls, in: Proceedings of the 18th International Conference on Pattern Recognition (ICPR), 2006.
- [40] J. Mitchelson, A. Hilton, Wand-based multiple camera studio calibration, CVSSP Technical Report VSSP-TR-2/2003.
- [41] T. Pappas, R. Sfrank, Perceptual criteria for image quality evaluation, *Handbook of Image and Video Processing*, 2000, pp. 669–684.
- [42] A. Pons, J. Malo, J. Artigas, P. Capilla, Image quality metric based on multidimensional contrast perception models, *Displays* 20 (25 August 1999) 93–110, (18).
- [43] S. Seitz, B. Curless, J. Diebel, D. Scharstein, R. Szeliski, A comparison and evaluation of multi-view stereo reconstruction algorithms, in: IEEE International Conference on Computer Vision and Pattern Recognition, vol. 1, 2006, pp. 519–526.
- [44] K. Seshadrinathan, A. Bovik, A structural similarity metric for video based on motion models, in: IEEE Conference on Acoustics, Speech and Signal Processing (ICASSP), vol. 1, 2007, pp. 869–872.
- [45] H. Sheikh, A. Bovik, Image information and visual quality, *IEEE Trans. Image Process.* 15 (2) (2006) 430–444.
- [46] J. Starck, A. Hilton, Model-based multiple view reconstruction of people, in: IEEE International Conference on Computer Vision (ICCV), 2003, pp. 915–922.
- [47] J. Starck, A. Hilton, Surface capture for performance based animation, *IEEE Comput. Graphics Appl.* 27 (3) (2007) 21–31.
- [48] J. Starck, J. Kilner, A. Hilton, Objective quality assessment in free-viewpoint video production, in: IEEE Conference on 3DTV, 2008, pp. 225–228.
- [49] J. Starck, A. Maki, S. Nobuhara, A. Hilton, T. Matsuyama, The 3D production studio, Technical Report VSSP-TR-4/2007.
- [50] E. Stoykova, A. Alatan, P. Benzie, N. Grammalidis, S. Malassiotis, J. Ostermann, S. Piekh, V. Sainov, C. Theobalt, T. Thevar, X. Zabulis, 3-d time-varying scene capture technologies—a survey, *IEEE Trans. Circuits Systems Video Technol.* 17 (11) (2007) 1568–1586.
- [51] T. Svoboda, D. Martinec, T. Pajdla, A convenient multicamera self-calibration for virtual environments, *Presence: Teleoper. Virtual Environ.* 14 (4) (2005) 407–422.
- [52] R. Szeliski, Rapid octree construction from image sequences, *CVGIP: Image Understanding* 58 (1) (1993) 23–32.
- [53] S. Vedula, S. Baker, T. Kanade, Image-based spatio-temporal modeling and view interpolation of dynamic events, *ACM Trans. Graphics* 24 (2) (2005) 240–261.
- [54] D. Vlasic, I. Baran, W. Matusik, J. Popovic, Articulated mesh animation from multi-view silhouettes, *ACM Trans. Graphics (SIGGRAPH)* (2008), in press.
- [55] G. Vogiatzis, C. Hernandez, P. Torr, R. Cipolla, Multiview stereo via volumetric graph-cuts and occlusion robust photo-consistency, *IEEE Trans. Pattern Anal. Mach. Intell.* 29 (12) (2007) 2241–2246.
- [56] Z. Wang, A. Bovik, H. Sheikh, E. Simoncelli, Image quality assessment: from error visibility to structural similarity, *IEEE Trans. Image Process.* 13 (4) (2004) 600–612.
- [57] M. Waschbusch, S. Wurmlin, M. Gross, 3d video billboard clouds, *Comput. Graphics Forum* 26 (September 2007) 561–569, (9).
- [58] S. Winkler, A perceptual distortion metric for digital color video, in: *SPIE*, 1999, pp. 175–184.
- [59] S. Winkler, Video quality and beyond, in: Proceedings of European Signal Processing Conference, 2007.
- [60] O. Woodford, I. Reid, A. Fitzgibbon, Efficient new view synthesis using pairwise dictionary priors, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- [61] S. Yamazaki, R. Sagawa, H. Kawasaki, K. Ikeuchi, M. Sakauchi, Microfacet billboarding, in: Eurographics Workshop on Rendering (EGWR), vol. 97(2), 2002, pp. 169–180.
- [62] C. Zitnick, S.B. Kang, M. Uyttendaele, S.A.J. Winder, R. Szeliski, High-quality video view interpolation using a layered representation, *ACM Trans. Graphics (SIGGRAPH)* 23 (3) (2004) 600–608.