

Flat3D: Browsing Stereo Images on a Conventional Screen

Wenjing Geng, Ran Ju, Xiangyang Xu, Tongwei Ren, and Gangshan Wu

State Key Laboratory for Novel Software Technology
Nanjing University, China
jenneng@gmail.com, {juran,xiangyang.xu}@smail.nju.edu.cn,
{rentw,gswu}@nju.edu.cn

Abstract. Expensive and cumbersome 3D equipment currently limits the popularization of emerging stereo media on the Internet. Particularly for stereo images, as a major kind of stereo media widespread on the Internet, there is not yet a good solution to show stereoscopy in conventional displays. By investigating the principles of human visual system (HVS), this paper proposes a method, called Flat3D for animating stereoscopy only through a conventional screen (2D) based on the motion parallax. The way for exhibition is dynamically transforming consistent views from left to right and then playing back reversely. The relative motion impresses viewers with strong depth perception. We investigate some factors which affect viewing experience in Flat3D and find that a reasonable fixation point and structure-preserved view transition contribute the most. Based on the above findings, we develop an adaptive fixation acquisition approach combining color and depth cues, as well as employing a probability-based view synthesis to generate the view sequences. Experiments which compared the above factors in and out of consideration show that our approach is a more convenient, effective and automatic alternative for browsing stereo images in common flat screens.

Keywords: Stereo images, fixation, view synthesis, transfer sequences.

1 Introduction

The amount of stereo images is soaring up on the Internet due to the popularization of 3D acquisition devices, for example Fujifilm 3D camera. Nevertheless, viewers cannot perceive 3D without the help of 3D equipment such as 3D glasses and 3DTV. Undoubtedly, there is a seemingly wide gap between viewers and stereo images because of extra demands for 3D equipment. In other words, stereo photographs captured by people or downloaded from the Internet, would degrade to side-by-side or red-cyan pictures in the eyes, as shown in Fig. 1 (a)-(b). Most people even do not understand what the meaning of stereo media is without glasses. This limitation is also another main reason that blocks the further widespread of 3D media. Therefore, a convenient and effective method to exhibit stereo images widespread on the Internet turns out to be an essential and

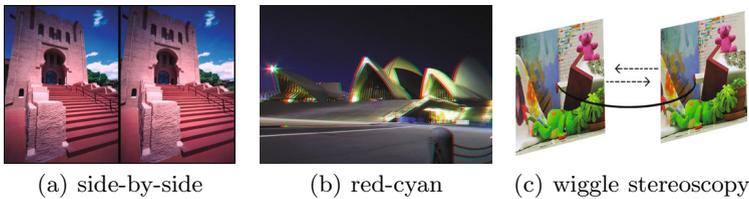


Fig. 1. Display of stereo images

significant problem. Besides, how to build a more pleasing 3D effects to cater for feelings of human vision through flat screens is another challenge.

To this problem, there are mainly two solutions separately from hardware [20] and software perspective. The former usually makes changes in the structure of screens based on optical principles, including lenticular lens, parallax barrier [14] and directional backlight [18]. Although these techniques have been applied in many devices, serving for advertising media and many other applications, it usually costs a lot, which makes it not applicable for the visualization of daily stereo images. Besides, there still exist some defects on resolution, visual angle and distance. The latter always utilizes computer graphics or visual rules to alter the content of images to create the third dimension, depth, such as wiggle stereoscopy¹, adding auxiliary lines² and perspective. Actually, these methods can indeed create stereoscopy, but there are something should not be ignored. First, the stereoscopic effect is usually achieved by human interference which is impossible to apply to the numerous stereo images on the Internet, let alone offers good and cosy user experience between left and right views based on uncertain fixation (Fig. 1 (c)). The 3D effect of this method inevitably flickers a lot. In a word, relying on photo editing and lacking complete automatic operating mechanisms make it impossible to apply to arbitrary stereo images. Second, some 3D effects generated by adding auxiliary or perspective cannot directly apply to stereo images at present and the generated results are obvious artifact.

In this paper, we propose a novel display approach named Flat3D for browsing stereo images on a conventional screen based on two findings of improving user experience for depth perception. First, considering the different influence upon region of interest (ROI) of human visual system (HVS), convergent point is calculated based on color and depth information, which is liable to find the optimal and pleasing focus point just as human do. Second, a probability-based rendering for view synthesis [2] is used to synthesize spatial-temporal consistent views for its insensitiveness to depth inaccuracy. Then the structure-preserved views are aligned based on the fixation. The format of demonstration is GIF or video sequences following visual persistence. By adaptively fixing human focus, the transfer sequences keep cyclical view transition from left to right and then from right to left. Compared with the existing methods, our approach is automatic, cheap, and efficient, and thus competent in exhibiting stereo images on a conventional screen

¹ http://en.wikipedia.org/wiki/Wiggle_stereoscopy

² <http://www.mymodernmet.com/profiles/blogs/3d-gifs>

such as PC screens and mobile screens, with good user experience. In summary, the major contributions include:

- An automatic and systematic approach for browsing stereo images is proposed by utilizing principles of motion parallax, fixation and visual persistence to improve user experience, which provides a cheap, convenient and comfortable alternative for demonstration of stereo images in daily life;
- Color and depth cues are combined to predict the fixation and structure-preserved views are synthesized by referring to the convergence of human eyes, which brings preferable 3D effects when browsing stereo images through common screens.

The rest of this paper is organized as follows. The background is briefly introduced in Sect. 2. Then a brief review of the related work is introduced in Sect. 3. The approach is detailedly described in Sect. 4 and evaluated by a few experiments shown in Sect. 5. Finally, we give a conclusion in Sect. 6.

2 Background

Main factors for HVS to produce stereoscopic impression are binocular parallax, motion parallax, accommodation and convergence [3]. The left and right eyes capture different content of the same scene. Due to accommodation and convergence, one image is composed in the brain. Existing 3D displays usually utilize horizontal parallax to create 3D mostly relying on expensive equipment. While we follow the motion parallax depending on relative movement, synthesize the intermediate views and align all the views based on the fixation marked as F_c in the Fig. 2. For stereo images are usually captured by binocular camera, we have got the first factor of binocular parallax from left and right views.

Motion parallax is a psychological phenomenon that observers view objects that are closer to them moving faster than objects that are further away from them. As to our framework, we utilize motion parallax to produce stereoscopy. We need to find the convergent point closer to HVS, or human would feel uncomfortable if the stationary point is not attractive. In other words, we need

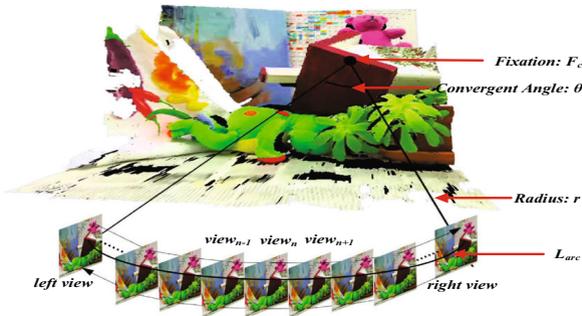


Fig. 2. The mechanism of the proposed approach

to automatically detect region of interest (ROI) first and calculate its centroid as fixation point then. For centroid is uniform distribution of quality, it is the optimal point to be fixation. We then transform this problem to saliency map calculation, because saliency determines attentional selection in psychology [11]. Furthermore, the transfer sequences are composed according to visual persistence that human can preserve images within the range of 0.1s to 0.4s. In order to build spatial-temporal consistent stereoscopy in the display, the duration of each view should not be set too long. Eq. (1) shows that it would take more time to transfer the sequences from different views with increasing number of view frames. And the velocity of the transfer sequences is inverse proportion to number of views.

$$v = \frac{L_{arc}}{T_{total}} = \frac{\theta\pi r}{Num_{view} \times t_{per}}, \quad (1)$$

where L_{arc} is the arc path that views are transforming by and T_{total} is the total time cost of transfer sequences from left to right view. θ is the angle of view and r is the radius labeled in Fig. 2. Num_{view} denotes the number of view frames and t_{per} shows the time persistence of each frame.

Based on the above methodology of psychology, physiology and HVS, we propose an automatic display approach for browsing 3D in a flat screen which caters to the requirement for comfortable 3D effects.

3 Related Work

We briefly review the similar research on wiggle stereoscopy and also the relevant techniques on saliency detection and multi-view synthesis.

Wiggle Stereoscopy. Wiggle stereoscopy describes technique giving an illusion of 3D by showing two images in rapid alternation. With the spring up of stereo images, it has gradually become a useful method when amateurs want to create stereoscopy. Therefore, there are several step-by-step tutorials distributed in the Internet. Most of these procedures rely on image editing software with human interference, for example Photoshop, which means it is impossible for thousands of stereo images to generate stereoscopic wiggle sequences. Besides, the common practice is composing dynamic switching between left and right views, for example StereoPhoto Maker, which would cause flicker to the result images. Because it lacks comparatively complete theoretical instructions, the generation GIF results always cause a little dizzy when watching it, which limits the widely use of this technique. Different from wiggle stereoscopy, the proposed approach automatically generates the results by exploring the principles of HVS, eliminating strong flicker by adaptively calculating fixation and consistent transforming.

Saliency. Saliency is a physiological and psychological phenomenon of visual attention, and has been used as a fundamental in many vision and multimedia tasks[16]. Itti et al. [4] used color, intensity and orientation to compute the

saliency map. Judd et al. [7] learnt a saliency model on a thousand images. Cheng et al. [1] introduced a global contrast based salient region detection by automatic estimating salient object regions. We choose the absorbing markov chain method to detect saliency proposed by Jiang et al. [5] for its efficiency and robustness to multi-focus images. Furthermore, it has been proved that depth perception has a strong impact on visual attention [15]. But most existing work always take a 2D color image as input which lacks depth cues. Lang et al. [9] modeled saliency as the conditional probability given depth and depth range. But the limitation is that global depth structure information is missing. Therefore, we utilize a depth saliency based on anisotropic center-surround difference proposed by Ju et al. [6] for its superiority to the state-of-the-arts.

Multi-View Synthesis. Multi-view synthesis is an essential step in many 3D display technologies. Min et al. [13] synthesized a virtual view by adapting a reverse warping instead of a forward warping. Mahajan et al. [12] interpolated a virtual view based on the idea that the given pixel to be synthesized in the virtual views traces out the path in reference images. These methods are both discrete formulation which would bring the problem of hole region. Lang et al. [10] presented a method free from hole filling, but the geometric distortion is still inevitable. Another challenge is depth inaccuracy on view interpolation. Kunita et al. [8] introduced the layered probability maps for dealing with depth ambiguities. Although the similar concept is used in the rendering process, including many geometric prior makes it not suitable for our framework. A probability-based rendering (PBR) method is proposed by Ham et al. [2] who addressed view synthesis as an image fusion. This method gives a good solution to depth inaccuracy and generates consistent images from different views.

4 Approach

Given an input of stereo images, we aim to automatically produce transfer sequences animating its stereoscopy in common displays. We formulate the problem by simulating HVS and utilizing motion parallax. The framework of our approach is shown in Fig. 3, which is made up of three main modules. The detailed procedures of the proposed approach are presented in Algorithm 1.

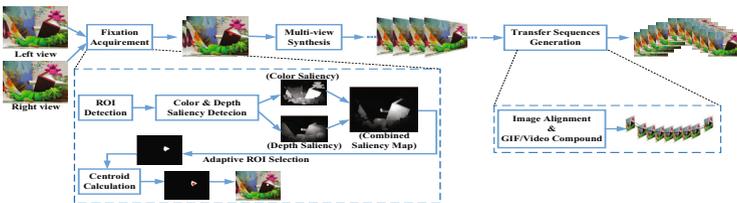


Fig. 3. The overview of the proposed approach

Algorithm 1. Outline of Flat3D

Input : A rectified stereo images**Output**: Dynamic transfer sequences from different views**begin**

1. Calculate visual fixation to acquire the fixed point as reference. (Sect. 4.1)
 - (a) Detect ROI based on saliency detection calculated by Eq. (2)-(5).
 - (b) Calculate visual fixation, centroid of ROI according to Eq. (6).
2. Synthesize multi-view images from left to right angel. A probability-based rendering method is adopted based on Eq. (7). (Sect. 4.2)
3. Generate transfer sequences. (Sect. 4.3)
 - (a) Align different views from left to right based on the fixation point according to the transformation matrix defined in Eq. (8).
 - (b) Compound transfer sequences based on visual persistence.

end

4.1 Fixation Acquisition

Visual fixation is the maintaining of visual gaze on a single location and plays an important role in human stereo vision. Saliency detection, or gaze prediction, computationally detects the fixation. Different from existing methods working on single color images, we believe many factors having effects on fixation selection, such as color contrast, spatial structure and depth. Hence, a combination of color and depth cues calculated from the stereo image pairs is adopted. This module consists of ROI detection and centroid calculation.

ROI Detection. In our approach, ROI is detected by calculating saliency map, for saliency is the most informative and interesting region in a scene. Stereo images contain more information compared to monocular images. Not only color and structure, but also depth can contribute to saliency detection. Therefore, we present a method of color saliency plus depth saliency to acquire a compound saliency map and then get ROI. Considering the appearance divergence and spatial distribution of salient objects and background, we employ the saliency detection via absorbing markov chain for color saliency [5]. The saliency map calculated by [5] with factors of color and spatial distribution is $Sal_{c,s}$, and the equation is as follows. Where i indexes the transient nodes on graph and y_w denotes the normalized weighted absorbed time vector.

$$Sal_{c,s}(i) = y_w(i), i = 1, 2, \dots, t. \quad (2)$$

For the other powerful cue, depth, we calculate the saliency based on anisotropic center-surround difference [6], denoted by Sal_d as shown in Eq. (3).

$$Sal_d(i) = \bar{D}_{acsd}(i) = \sum_{k=1}^n \sum_{t=1}^8 D_{acsd}^t(p) / \sum_{i=1}^m D_{acsd}(i), \quad (3)$$

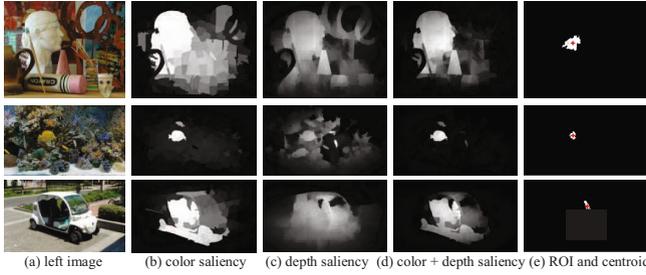


Fig. 4. Results of saliency detection and centroid calculation

where $\bar{D}_{acsd}(i)$ dedicates the normalized ACSD value of superpixel i and $D_{acsd}(i)$ is the ACSD value of superpixel i . $D_{acsd}^t(p)$ represents the ACSD value of pixel p in the t^{th} direction. The number of pixels in the superpixel is n and the number of superpixel is m .

Then the candidate saliency map Sal is calculated as Eq. (4) and normalized to $[0, 1]$. Then a threshold function is used to determine the interest region, shown in Eq. (5). $I_{bw}(\cdot)$ is the binary image labeled the connected region, τ is an adaptive threshold calculated using OTSU. Results of saliency maps for each processing step are listed in Fig. 4 (b)-(d).

$$Sal = Sal_{c,s} \times Sal_d. \quad (4)$$

$$I_{bw}(x) = \begin{cases} 1 & Sal(x) > \tau \\ 0 & \text{otherwise} \end{cases}. \quad (5)$$

Centroid Calculation. We choose centroid to represent visual fixation point, which is the point representing mass balance in physics. Let $c(x, y)$ denotes centroid of a region, we calculate its coordinate as Eq. (6), where $I(x, y)$ is the intensity of the binary image.

$$\begin{cases} x_c = \sum_{(x,y) \in T} xI(x, y) / \sum_{(x,y) \in T} I(x, y) \\ y_c = \sum_{(x,y) \in T} yI(x, y) / \sum_{(x,y) \in T} I(x, y) \end{cases}. \quad (6)$$

Besides, connected regions obeying the rule defined in Eq. (6) may surpass one. We firstly sort the connected regions and choose the largest one as the ROI and compute the centroid of it. Some examples of binary images marked with red centroid are shown in Fig. 4 (e).

4.2 Multi-view Synthesis

The aim of our framework is generating consistent view transition in a specific time according to visual persistence. In fact, fusion based image morphing (FBIM) and motion based image morphing (MBIM) could also create motion

parallax which seems to be unnecessary for multi-view synthesis. However, image wiggling only between left and right views would be too flickering. We resort to an image fusion method which is performed in a probabilistic way [2] for its insensitivity to depth inaccuracy to preserve change of view structure.

Let $I_l(m)$ and $I_r(m)$ represent left and right images separately. Their corresponding steady state matching probability defined in [2] is $P_l(m, d)$ and $P_r(m, d)$. Assuming that the baseline between left and right images is normalized to 1, then the location of virtual view $I_v(m)$ is denoted by α , where $0 < \alpha < 1$. In order to minimize the inaccuracy brought by depth map, the synthetic process is transformed to image fusion. First candidate re-sampled color images on the virtual view $I_l^v(m)$ from left images and virtual view $I_r^v(m)$ from right images are calculated. Then view synthesis turns to image fusion, shown in Eq. (7).

$$\begin{aligned} I_v(m) &= \alpha I_l^v(m) + (1 - \alpha) I_r^v(m) \\ &= \sum_d (\alpha I_l^v(m, d) P_l(m, d) + (1 - \alpha) I_r^v(m, d) P_r(m, d)), \end{aligned} \quad (7)$$

where $I_l^v(m, d)$ and $I_r^v(m, d)$ are image intensity along the different disparity hypothesis d . The detailed description of this method is described in [2]. By utilizing PBR method, we calculate the intermediate views and some of them are shown in Fig. 5. The red straight line is a fixed marker used for showing the consistent view change from right to left.

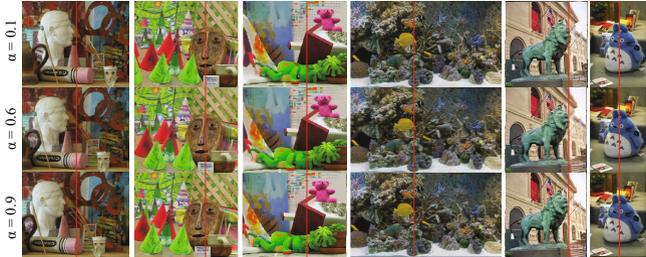


Fig. 5. Results of synthesized views

4.3 Transfer Sequences Generation

The centroid of salient region is chosen as the stationary point in which the disparity should be zero. The reason is that if the most attractive point has a frequent change, human would feel uneasy and difficult to perceive 3D. Admittedly, multi-view based image change without convergence (MBICWC) could also bring 3D effect, it is simply image switching and easily degrades stereoscopy for not obeying rules of stereo vision. In order to make the disparity of fixation to be zero, each view should be aligned based on the transformation matrix:

$$T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ disp & 0 & 1 \end{bmatrix}. \quad (8)$$

$$disp = D_{target}(centroid), \quad (9)$$

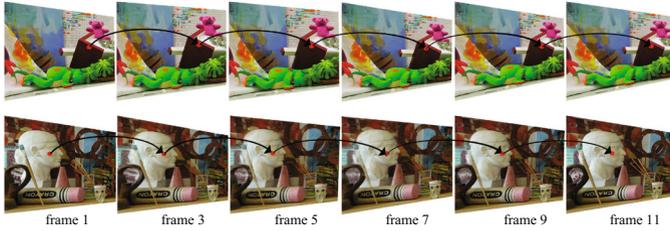


Fig. 6. Examples of image alignment based on fixation

where $disp$ is disparity of centroid from target image to reference image derived from depth map D . The direction is negative in our method because right images are used as reference images. We trace the trajectory of fixation listed in Fig. 6, align them and make their disparity to be zero.

The final processing is generating dynamic sequences of all the aligned views. Either GIF or video format is feasible based on the actual requirement. However, on the basis of visual persistence about 0.1s to 0.4s, the frame rate should obey this rule which has been discussed in Sect. 2.

5 Experiments

5.1 Datasets and Experimental Settings

Our approach is implemented using Matlab on a desktop PC with an Intel i7_4770 CPU and 16GB memory. It is noted that the generated results of our approach contain continuous transfer sequences from left to right view and then from right to left. Hence the 3D effect could be dynamically displayed in any flat screen as GIF or video format³.

To show the efficiency and robustness, we evaluate our approach on Middlebury Stereo Datasets [17] and OBSIR datasets [19]. The details of the datasets are listed in Table 1, where GT is the ground truth. We use 0.1s for each frame and 20 frames are used to create smooth view transition but only 9 views need to be interpolated. The frame rate of the generated transfer sequences is 10fps.

Table 1. Description of Middlebury and OBSIR dataset

Dataset	Num.	Num. of test	GT. of depth map	GT. of multi-view images
Middlebury	39	39	Y	Y
OBSIR	10513	1382	N	N

5.2 Results and Discussion

We make a thorough comparison by considering the influence of fixation, multi-view synthesis and a user study is conducted to evaluate the 3D effect.

³ <http://mcg.nju.edu.cn/publication/MMM15-GengWJ.html>

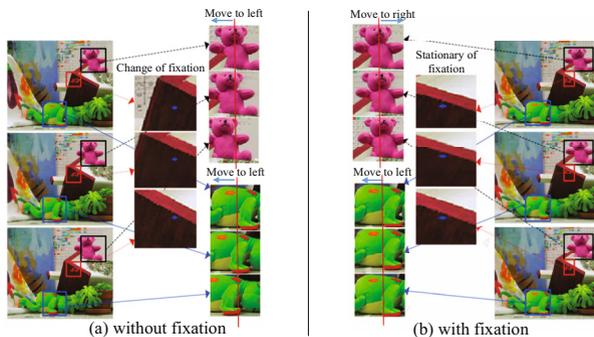


Fig. 7. Comparison with view change without image alignment based on fixation

Automatic Fixation. Fig. 7 gives a concrete comparison between MBICWC and Flat3D. We extract the same contents from three views, marked by black, red and blue windows. Red stars mark the visual fixation in each frame. It can be seen that all the pixels are moving from right to left in Fig. 7 (a) because of no fixation. On the contrary, in Fig. 7 (b) objects behind the fixation move from left to right and objects in front of the fixation move from right to left. Above all, it is the fixation that causes the relative motion before and after the fixation in 3D space. Objects behind the fixation, e.g. pink bear, move to the right while the objects in front of fixation, e.g. green legs move to the left.

View Interpolation. Not every image movement can create stereoscopy in flat screens. The proposed method is not a plain change among images. It is an imitation of human stereo vision which makes flat screens to be retina. In order to show the efficiency and reasonability, we also compare our method with FBIM and MBIM. The first line of Fig. 8 shows several frames generated by FBIM from six views separately. The middle frames become obscure due to the image fusion of linear interpolation, which looks like playing slides and hardly generates 3D effects. The second line is created by MBIM. Although frames are more distinct than the above line, there are inevitable holes caused by depth inaccuracy. The third line is the same frames built by the proposed approach, which has consistent trajectory to ensure smooth transition on views and one stationary fixation to create 3D effects.

User Study. To further evaluate the efficiency, we conducted a 3D effects user study aiming at whether the transfer sequences automatically produced by our method are preferred by users to those of sequences generated by other methods. We invite 36 subjects (20 males and 16 females with ages from 16 to 47) with normal vision and different educational background participated in the experiment. The dynamic sequences produced by 4 different approaches are compared, including FBIM, MBIM, MBICWC and the proposed approach. Each type of

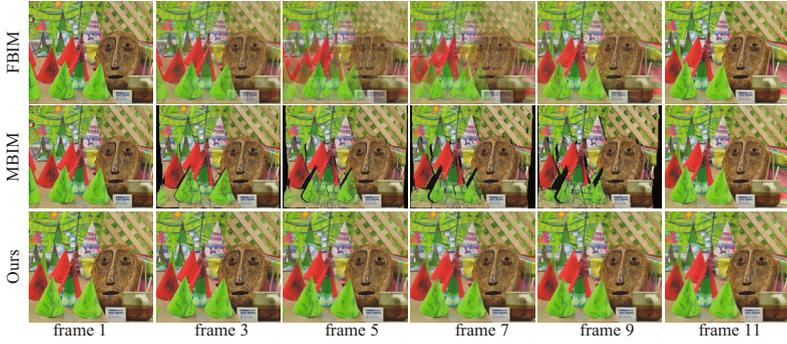


Fig. 8. Comparison results

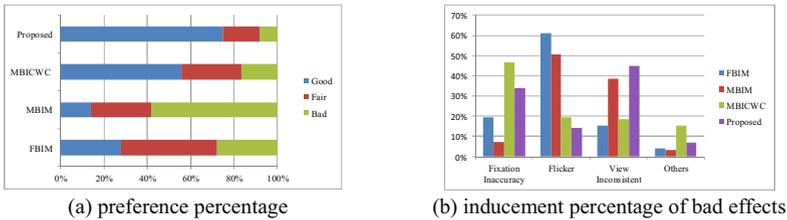


Fig. 9. The comparison of the four evaluated approaches on user ratings

dynamic sequences is numbered as 1, 2, 3, 4 with GIF format for its convenience. In the user study, the feelings for 3D effects are divided into 3 quality levels: Good, Fair, Bad, according to the observers’ own perspective. After above process, the ratings from the users on each level are accumulated and the proportions of each quality level by the different approaches are listed in Fig. 9 (a). The Good rates are separately 75.13%, 55.56%, 13.89% and 27.78%. Our approach received most preference for its good performance in 3D effects. Besides, we offer 4 main options of affecting stereoscopy for subjects to choose if bad 3D effects are labeled. The options are fixation inaccuracy, flicker, view inconsistent and something others. The statistics are the average percentage based on each subject’s option within bad levels, shown in Fig. 9 (b). 3D effects of MBICWC are close to ours because it is also a smooth view transition, but no fixation degrades its 3D effect. MBIM received the most Bad for annoying holes during view change which easily brings strong flicker. And FBIM is accurately a process of image fusion with inevitable flicker.

6 Conclusion

This paper proposes an automatic and systematic approach for browsing stereo images in flat screens by generating transfer sequences. Based on findings of human stereo vision, we believe our method is a good solution to exhibit stereo im-

ages stereoscopic in daily screens, e.g., mobile screens. The experiments and user evaluations demonstrate the reasonability and effectiveness of our method.

Acknowledgement. This work is supported by the National Science Foundation of China (No.61321491, 61202320), Research Project of Excellent State Key Laboratory (No.61223003), Natural Science Foundation of Jiangsu Province (No.BK2012304), and National Special Fund (No.2011ZX05035-004-004HZ). It was also partially supported by Collaborative Innovation Center of Novel Software Technology and Industrialization.

References

1. Cheng, M., Mitra, N., Huang, X., Torr, P., Hu, S.: Global contrast based salient region detection. PAMI (2014)
2. Ham, B., Min, D., Oh, C., Do, M., Sohn, K.: Probability-based rendering for view synthesis. TIP (2014)
3. Howard, I.P.: Binocular vision and stereopsis. Oxford University Press (1995)
4. Itti, L., Koch, C., Niebur, E.: A model of saliency-based visual attention for rapid scene analysis. PAMI (1998)
5. Jiang, B., Zhang, L., Lu, H., Yang, C., Yang, M.H.: Saliency detection via absorbing markov chain. In: ICCV. IEEE (2013)
6. Ju, R., Ge, L., Geng, W., Ren, T., Wu, G.: Depth saliency based on anisotropic center-surround difference. In: ICIP. IEEE (2014)
7. Judd, T., Ehinger, K., Durand, F., Torralba, A.: Learning to predict where humans look. In: CVPR. IEEE (2009)
8. Kunita, Y., Ueno, M., Tanaka, K.: Layered probability maps: basic framework and prototype system. In: VRST. ACM (2006)
9. Lang, C., Nguyen, T.V., Katti, H., Yadati, K., Kankanhalli, M., Yan, S.: Depth matters: Influence of depth cues on visual saliency. In: Fitzgibbon, A., Lazebnik, S., Perona, P., Sato, Y., Schmid, C. (eds.) ECCV 2012, Part II. LNCS, vol. 7573, pp. 101–115. Springer, Heidelberg (2012)
10. Lang, M., Hornung, A., Wang, O., Poulakos, S., Smolic, A., Gross, M.: Nonlinear disparity mapping for stereoscopic 3D. TOG (2010)
11. Ma, L., Xu, K., Wong, T., Jiang, B., Hu, S.: Change blindness images. TVCG (2013)
12. Mahajan, D., Huang, F.C., Matusik, W., Ramamoorthi, R., Belhumeur, P.: Moving gradients: a path-based method for plausible image interpolation. TOG (2009)
13. Min, D., Kim, D., Yun, S., Sohn, K.: 2D/3D freeview video generation for 3dtv system. SPIC (2009)
14. Neil, A.: Autostereoscopic 3d displays. Computer (2005)
15. Niu, Y., Geng, Y., Li, X., Liu, F.: Leveraging stereopsis for saliency analysis. In: CVPR. IEEE (2012)
16. Ren, T., Ju, R., Liu, Y., Wu, G.: How important is location in saliency detection. In: ICIMCS. ACM (2014)
17. Scharstein, D., Szeliski, R.: High-accuracy stereo depth maps using structured light. In: CVPR. IEEE (2003)
18. Wetzstein, G., Lanman, D., Hirsch, M., Raskar, R.: Tensor displays: compressive light field synthesis using multilayer displays with directional backlighting. TOG (2012)

19. Xu, X., Geng, W., Ju, R., Yang, Y., Ren, T., Wu, G.: Obsir: Object-based stereo image retrieval. In: ICME. IEEE (2014)
20. Zhang, Y., Ji, Q., Zhang, W.: Multi-view autostereoscopic 3D display. In: OPEE. IEEE (2010)