

Racking focus and tracking focus on live video streams: a stereo solution

Zhan Yu · Xuan Yu · Christopher Thorpe ·
Scott Grauer-Gray · Feng Li · Jingyi Yu

© Springer-Verlag Berlin Heidelberg 2013

Abstract The ability to produce dynamic Depth of Field effects in live video streams was until recently a quality unique to movie cameras. In this paper, we present a computational camera solution coupled with real-time GPU processing to produce *runtime* dynamic Depth of Field effects. We first construct a hybrid-resolution stereo camera with a high-res/low-res camera pair. We recover a low-res disparity map of the scene using GPU-based Belief Propagation, and subsequently upsample it via fast Cross/Joint Bilateral Upsampling. With the recovered high-resolution disparity map, we warp the high-resolution video stream to nearby viewpoints to synthesize a light field toward the scene. We exploit parallel processing and atomic operations on the GPU to resolve visibility when multiple pixels warp to the same image location. Finally, we generate racking focus and tracking focus effects from the synthesized light field rendering. All processing stages are mapped onto NVIDIA's CUDA architecture. Our system can produce racking and tracking focus effects for the resolution of 640×480 at 15 fps.

Keywords Dynamic depth of field · Racking focus · Tracking focus · Belief propagation · Cross bilateral Filtering · Light field · CUDA

1 Introduction

Depth of field (DoF) effects are a useful tool in photography and cinematography because of their aesthetic value. In photography, they have been used to emphasize objects by creating a shallow plane of focus around the subject while blurring the rest of the scene. In cinematography, movie cameras use *dynamic* DoF effects to shift the viewer's attention from the foreground to the background or vice versa. Producing a high quality DoF, however, requires dedicated and often expensive lens systems. Commodity lenses such as the low cost Canon EF 50 mm f/1.4 use a small number of aperture blades and produce blur artifacts caused by the polygon shaped apertures.

Varying the DoF effect and displaying the results in real-time was until recently a quality unique to movie cameras. For example, classical digital SLRs, while capable of modifying the DoF effect from shot to shot, can only produce static images. Of late, high-end DSLRs can stream video, but this capability comes at increased cost to the consumer and still yields a nonideal video platform. Additionally, current DSLRs adjust focus using a ring on the lenses and DoF adjustments tend to produce momentary tilting of the captured video. Movie cameras, while capable of runtime dynamic DoF effects, tend to be bulky and unwieldy in order to accommodate complex corrective lenses arrays and storage media. Further, fully served studio lenses such as the Canon DIGI SUPER 25 XS can cost tens of thousands of dollars.

Z. Yu (✉) · X. Yu · C. Thorpe · S. Grauer-Gray · F. Li · J. Yu
University of Delaware, Newark, DE, 19716, USA
e-mail: zyu@cis.udel.edu

X. Yu
e-mail: xyu@cis.udel.edu

C. Thorpe
e-mail: thorpe@cis.udel.edu

S. Grauer-Gray
e-mail: grauerg@cis.udel.edu

F. Li
e-mail: feli@cis.udel.edu

J. Yu
e-mail: yu@cis.udel.edu

Our work is inspired by the recent light field imaging systems [9, 11, 21, 22, 27]. The Stanford light field camera array [19, 20, 23, 24] is a two-dimensional grid composed of 128 1.3 megapixel firewire cameras, which stream live video to a stripped disk array. The large volume of data generated by this array forces the DoF effect to be rendered in post processing rather than in real-time. Furthermore, the system infrastructure such as the camera grid, interconnects, and workstations are bulky, making it less suitable for on-site tasks. The MIT light field camera array [25] uses a smaller grid of 64 1.3 megapixel usb webcams instead of firewire cameras and is capable of synthesizing real-time dynamic DoF effects. Both systems, however, still suffer from spatial aliasing because of the baseline between neighboring cameras. The camera spacing creates appreciable differences between the pixel locations of the same scene point in neighboring cameras producing an aliasing effect at the DoF boundary when their images are fused.

In an attempt to reduce aliasing artifacts in light-field based solutions, Ng [13] designed a light field camera that combines a single DSLR with a microlenslet array. Each lenslet captures the scene from a different viewpoint and the lenslet array effectively emulates a camera array. By using a large number of densely packed lenslets, one can significantly reduce the spatial aliasing artifacts. It also, however, comes at the cost of reduced resolution for each light field view. A light field camera is typically paired with an ultra-high-resolution static DSLR and is therefore not applicable to video streaming.

Instead of using either a camera array or a microlenslet array, we develop a novel hybrid stereo-lightfield solution. Our goal is to first recover a high-resolution disparity map of the scene and then synthesize a virtual light field for producing dynamic DoF effects. Despite recent advances in stereo matching, recovering high-resolution depth/disparity maps from images is still too expensive to perform in real-time. We, therefore, construct a hybrid-resolution stereo camera system by coupling a high-res/low-res camera pair. We recover a low-res disparity map, and subsequently upsample it via fast cross bilateral filters. We then use the recovered high-resolution disparity map and its corresponding video frame to synthesize a light field. We implement a GPU-based disparity warping scheme and exploit atomic operations to resolve visibility. To reduce aliasing, we present an image-space filtering technique that compensates for spatial undersampling using MIPMAPPING. Finally, we generate racking focus and tracking focus effects using light field rendering. The complete processing pipeline is shown in Fig. 2.

We map all processing stages onto NVIDIA's CUDA architecture. Our system can produce racking focus and tracking focus effects with arbitrary aperture sizes and focal depths for the resolution of 640×480 at 15 fps, as shown in the supplementary video. This indicates that if we capture



Fig. 1 Depth of Field effect on a parking car scene using our system

the video streams at the same frame rate, we can display the refocused stream simultaneously. Our system thus provides a low-cost, computational imaging solution for runtime refocusing, an effect that is usually the domain of expensive movie cameras with servo-controlled lenses. Experiments on both indoor and outdoor scenes show that our framework can robustly handle complex, dynamic scenes and produce high quality results. Figure 1 shows the result of our system on a parking lot scene.

2 Hybrid-resolution stereo camera

We first construct a hybrid stereo camera for recovering high-resolution disparity map in real-time. Our system uses the Pointgrey Flea2 camera pair to produce one *high-resolution color* video stream and one *low-resolution gray-scale* video stream. We synchronize frame capture to within $125 \mu\text{s}$ by using the Pointgrey camera synchronization package. A unique feature of our approach is coupling our Hybrid-Resolution Stereo Camera with a CUDA processing pipeline for real-time DoF synthesis. Sawhney et al. proposed a hybrid stereo camera for synthesis of very high resolution stereoscopic image sequences [15]. Li et al. also proposed a hybrid camera for motion deblurring and depth map super-resolution [10]. Our configurations, however, have many more advantages. First and foremost, it provides a multi-resolution stereo matching solution that can achieve real-time performance (Sect. 3). Second, the lower bandwidth requirement also allows our system to be implemented for less expense on a greater number of platforms. Stereo systems that stream two videos at 15 fps and 640×480 resolution can produce up to 27.6 MB of data per second. By comparison, our hybrid-resolution stereo camera only produces slightly more than half that rate of data. Although our current implementation uses Firewire cameras,

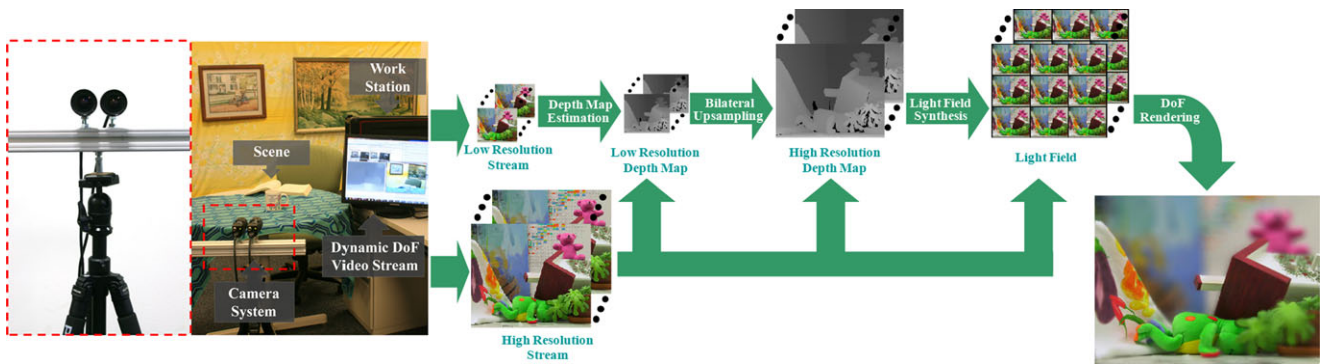


Fig. 2 The imaging hardware and the processing pipeline of our dynamic DoF video acquisition system. All processing modules are implemented on NVIDIA’s CUDA to achieve real-time performance

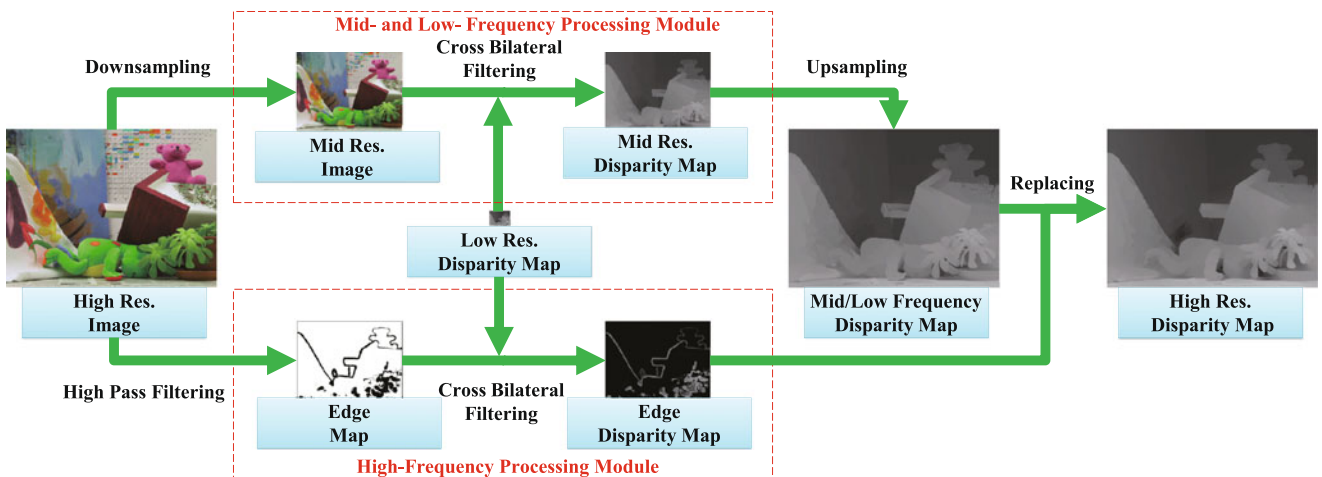


Fig. 3 Our fast cross bilateral upsampling scheme synthesizes a high-resolution disparity map from the low-resolution BP stereo matching result on CUDA

the low bandwidth demands of our solution make it possible to use a less expensive and more common alternative like USB 2.0, even for streaming higher resolutions such as 1024×768 . Finally, compared to off-the-shelf stereo cameras such as Pointgrey’s Bumblebee, our system has several advantages in terms of image quality, cost, and flexibility. For example, the form factor of the Bumblebee forces its lenses to be small and it produces image with severe radial distortion. Our system is also less expensive (\$1500 vs. \$4000), and our setup allows us to dynamically adjust the camera baseline to best fit different types of scenes unlike the Bumblebee. We calibrate the stereo pair using a planar checker board pattern the algorithm outlined by Zhang [28]. It is not necessary, however, that the calibration be absolutely accurate as the disparity map is recovered from a severely downsampled image pair. Our experiments have shown that disparity map recovery using belief propagation on the low-resolution image pair is not affected by slight changes in the camera pair geometry. The intensity calibration on the camera pair is performed prior to capture via

histogram equalization. The mappings for these processes are retained and applied to each incoming frame prior to stereo matching.

3 Real-time stereo matching

In order to efficiently generate a high-resolution disparity map from the input low-res/high-res image pairs, we implement a GPU-based stereo matching algorithm on CUDA.

3.1 CUDA belief propagation

Stereo matching is a long standing problem in computer vision [16]. Global methods based on belief propagation (BP) [18] and graph-cut [1, 5] have been known to produce highly reliable and accurate results. These methods, however, are more expensive when compared to local optimization methods such as dynamic programming. Fortunately, BP lends itself well to parallelism on the GPU [2, 4], where the core

Table 1 Performance of our CUDA stereo matching at different resolutions. Note that the number of disparity levels is proportionally scaled to the resolution. The levels of belief propagation are all set to 5 and iterations per level are all set to 10

Data sets	Resolutions		
	128 × 96	320 × 240	640 × 480
Teddy	13 ms	78 ms	446 ms
Tsukuba	8 ms	55 ms	357 ms
Cones	11 ms	69 ms	424 ms

computations can be performed at every image pixel in parallel on the device.

We utilize the methods presented by Felzenwalb [3] to speed up our implementation without affecting the accuracy: we use a hierarchical implementation to decrease the number of iterations needed for message value convergence; we apply a checkerboard scheme to split the pixels when passing messages in order to reduce the number of necessary operations and halve the memory requirements; and we utilize a two-pass algorithm to reduce the running time to generate each message from $O(n^2)$ to $O(n)$ using the truncated linear model for data/smoothness costs.

Our CUDA BP implementation uses five separate kernels, whereas the CPU only calls the appropriate kernels and adjusts the current parameters/variables. A kernel is used to perform each of the following steps in parallel, with each thread mapping to computations at a distinct pixel:

1. Compute the data costs for each pixel at each possible disparity at the bottom level.
2. Iteratively compute the data cost for each pixel at each succeeding level by aggregating the appropriate data costs at the proceeding level.
3. For each level of the implementation:
 - (a) Compute the message values at the current ‘checkerboard’ set of pixels and pass the values to the alternative set. Repeat for i iterations, alternating between the two sets.
 - (b) If not at the bottom level, copy message values at each pixel to corresponding pixels of the succeeding level.
4. Compute the disparity estimate at each pixel using the data cost and current message values corresponding to each disparity.

Table 1 shows the performance of our algorithm on some of the Middlebury datasets at different resolutions. Despite the acceleration on the GPU, we find that it is necessary to use the lower resolution images (320 × 240 or lower) as inputs to our stereo algorithm in order to achieve real-time performance.

In our experiments described in the rest of the paper, we first smooth these low-resolution image pairs using a Gaussian filter where σ equals 1.0, then process them using our

implementation with a disparity range from 0 to 35, maximum data cost and smoothness costs of 15.0 and 1.7, respectively, a data cost weight of 0.7 in relation to the smoothness cost, with 5 levels of belief propagation and 10 iterations per level. Each kernel is processed on the GPU using thread block dimensions of 32×4 .

3.2 Fast cross bilateral upsampling

Given a low-resolution disparity map D' and a high-resolution image I , we intend to recover a high-resolution disparity map D using cross bilateral filters [26], where we apply a spatial Gaussian filter to D' and a color-space Gaussian filter to I . Assuming p and q are two pixels in I ; W is the filter window size; I_p and I_q are the color of p and q in I ; and q' is the corresponding pixel coordinate of q in D' . We also use σ_c and σ_d as constants to threshold the color difference and filter size. We compute the disparity of pixel D_p as

$$D_p = \frac{\sum_{q \in W} G_d(p, q) G_c(p, q) D'_q}{K_p}, \quad (1)$$

where $K_p = \sum_{q \in W} G_d(p, q) G_c(p, q)$, $G_d(p, q) = \exp(-\frac{\|p-q\|}{\sigma_d})$, and $G_c(p, q) = \exp(-\frac{\|I_p - I_q\|}{\sigma_c})$.

The complexity of cross bilateral upsampling (CBU) is $O(NW)$ where N is the output image size and W is the filter window size. Therefore, the dominating factor to the processing time is the number of pixels that need to be upsampled, i.e., the resolution of the high-res image in the brute-force implementation.

To accelerate our algorithm, we implement a fast CBU scheme that effectively reduces the pixels to be upsampled. Paris et al. [14] have shown that the mid and low frequency components of an image remain approximately the same when downsampled. We therefore treat the high-frequency and the mid- and low-frequency components separately. Our method first applies a Gaussian high-pass filter to identify the pixels of high frequency in I and then uses a standard cross bilateral filter to estimate the disparity values at only these pixels. We store the resulting disparity map as D_{high} . We call this step the high-frequency processing module. In parallel, we downsample the color image to mid-resolution I_{mid} , apply CBU between D' and I_{mid} to obtain the mid-res disparity map D_{mid} ; and subsequently upsample D_{mid} to D_{high} using standard bilinear upsampling. We call this step the mid- and low-frequency processing module. Finally, we perform high frequency compensation by replacing the disparity value at the identified high frequency pixels \tilde{I} with D_{high} . Figure 3 shows the complete processing pipeline of our algorithm. Compared with standard CBU, our scheme only needs to upsample a small portion of the pixels and hence is much faster.

We also added a refining stage for sharpening the boundary regions and smoothing the surface regions with a cross bilateral filter, after the unsampling is done. The stage is basically the same as the upsampling stage except the input disparity map is the same size as the color image. Since the output disparity map could be treated as the input of another stage, this refining stage can be performed iteratively.

As shown in Fig. 6, if the refining stage is not performed, the edges and surfaces of the disparity map looks noisy due to the imperfection of the low resolution disparity map and the textures in the color image. However, if the refining stage contains too many iterations, then the disparities of one side of edges starts to bleed into the other side, which is the effect of over-smoothing. Therefore, a compromise number of iterations must be chosen at run time, using our interactive parameter interface (Sect. 6).

Note that the upsampled depth edges may not be consistent with the depth edges compute using the high frequency map. Here, we experimented with the following solutions: (1) Use the unsampled depth edges. (2) Use the high frequency depth edges. (3) Blend the two results. We found out that we can preserve more accurate edges and render better results using the second way.

3.3 CUDA implementation

We developed a GPU implementation of our algorithm on the CUDA architecture to tightly integrate with our CUDA BP stereo mapping algorithm. In our experiments, we found that it is most efficient to assign one thread to upsample each pixel in the disparity map. To further evaluate the throughput of our implementation, we upsampled 128×128 disparity maps with 1280×1280 color images. Our implementation achieves a processing speed of 22 ms per frame or 14 ms per megapixel with a 5×5 filter window, a significant speedup to the CPU-based scheme [6] (which was 2 seconds per megapixel).

To measure the accuracy of our scheme, we performed experiments using various stereo data sets. In Fig. 4, we show using the Teddy data set that reintroducing high frequency compensation produces sharper edges and smoother surfaces. Figure 5 illustrates our results in three regions on the Teddy data set. They are upsampled from 30×25 to 450×375 . Compared with standard bicubic or Gaussian upsampling, our method preserves fine details near the edges. It is important to note that preserving edges while removing noise in the disparity map is crucial to our DoF synthesis as DoF effects are most apparent near the occlusion boundaries. Figure 7 gives the results on an indoor scene using bicubic upsampling and our method. To further measure the accuracy, we compared our estimation with the ground truth by computing the mean squared errors over all pixels. Table 2 compares the error incurred by our method under different upsampling scales on a variety of Middlebury stereo

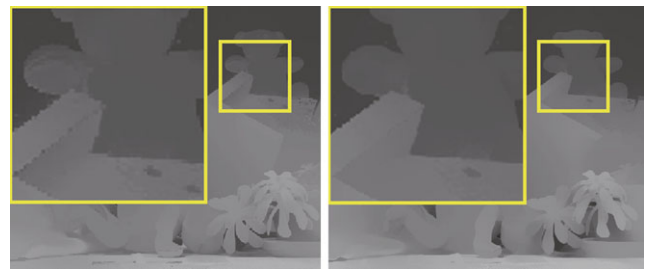


Fig. 4 Comparison of the result with (*right*) and without (*left*) high frequency compensation

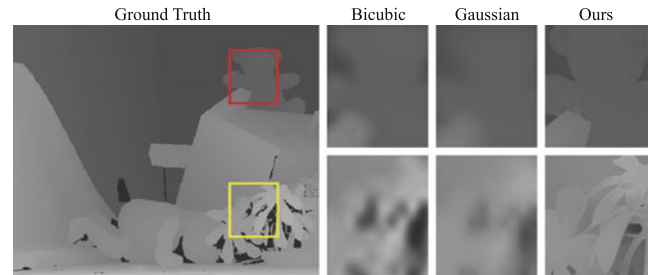


Fig. 5 Comparison of our method and other upsampling schemes on synthesize data. Both patches in the disparity map are upsampled from a resolution of 30×25 to 450×375

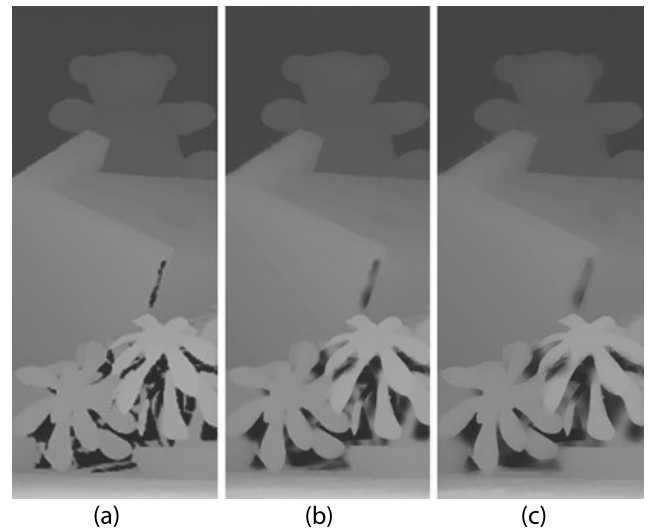


Fig. 6 Comparison of three results using different number of refining iterations. Result (a), (b), (c) are using 0, 3, and 10 iterations, respectively

data sets, and the results show that our method is reliable and accurate even with very high upsampling scales.

In our indoor and outdoor experiments, good results of disparity maps can be achieved when there are 10 iterations in refining stage. Since the system runs at interactive speed, it is impossible to use standard CBU to upsample the low resolution disparity because it would take 0.1 second (10 frame per second) to compute a single frame of resolution



Fig. 7 Comparison between our method and bicubic upsampling on real scenes. The disparity map is upsampled from 320×240 to 640×480 . Our method preserves sharp edges and maintains smoothness, which is critical to reliable DoF synthesis

Table 2 Pixels with disparity error larger than 1 under different upsampling factors on the Middlebury data sets

Data sets	Upsampling scales			
	20×20	10×10	5×5	2×2
Teddy	10.41 %	3.56 %	1.71 %	0.52 %
Plastic	8.36 %	4.23 %	2.05 %	0.91 %
Monopoly	11.76 %	5.35 %	2.96 %	1.14 %
Books	9.28 %	6.12 %	2.63 %	1.02 %
Baby2	5.76 %	2.38 %	1.61 %	0.69 %
Aloe	15.12 %	7.83 %	3.40 %	1.17 %
Cones	11.51 %	5.87 %	3.25 %	1.28 %
Art	13.47 %	7.15 %	3.43 %	1.41 %

640×480 with CUDA implementation. However, with our fast CBU framework, the speed quickly goes up to 40 frame per second with the downsampling factor 2×2 .

4 Real time DoF synthesis

Once we obtain the high-resolution disparity map, we set out to synthesize dynamic DoF effects. Previous single image based DoF synthesis algorithms attempt to estimate the circle of confusion at every pixel and then apply the spatially varying blurs on the image. These methods produce strong bleeding artifacts at the occlusion boundaries, as shown in Fig. 8. In computer graphics, the distributed ray tracing and the accumulation buffer techniques have long served as the rendering method for dynamic DoF. Both approaches are computationally expensive as they either require tracing out a large number of rays or repeated rasterization of the scene. Furthermore, to apply ray-tracing or accumulation buffer in our application requires constructing a triangulation of the scene from the depth map, which would incur additional computational cost.

In this paper, we adopt a similar approach to [27] by dynamically generating a light field from the high-resolution video stream and its depth stream, as shown in Fig. 9. Our

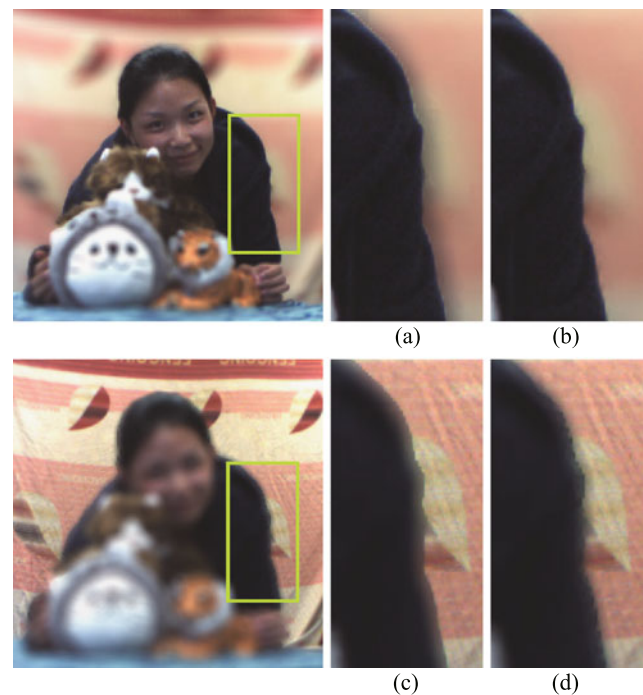


Fig. 8 Comparing results generated by image space blurring (a, c) and our light field synthesis method (b, d). Our approach effectively reduces both the intensity leakage (a) and boundary discontinuity (c) artifacts

technique, however, differs in that we directly use the disparity map for warping and filtering whereas [27] builds upon the depth map. As follows, we briefly reiterate the main steps of this light-field based DoF rendering technique.

4.1 The lens light field

The light field is a well-known image based rendering technique. It uses a set of rays commonly stored in a 2D array of images to represent a scene. Each ray in the light field can be indexed by an integer 4-tuple (s, t, u, v) , where (s, t) is the image index and (u, v) is the pixel index within a image.

Our first step generates a light field from the stereo pair. The high resolution camera in our stereo pair is used as the reference camera R_{00} .

To synthesize the light field, we use the high-resolution camera in our stereo pair as the reference camera R_{00} (i.e., $(s, t) = (0, 0)$). We can then easily find all rays that pass through a 3D point A in terms of its disparity γ from the reference view. Assuming A 's image is at pixel (u_0, v_0) in the reference camera, we can compute its image (pixel coordinate) in any light field camera R_{st} as:

$$(u, v) = (u_0, v_0) + (s, t) \cdot \gamma. \quad (2)$$

We use $L_{\text{out}}(s, t, u, v)$ to represent the out-of-lens light field and $L_{\text{in}}(x, y, s, t)$ to represent the in-camera light field.

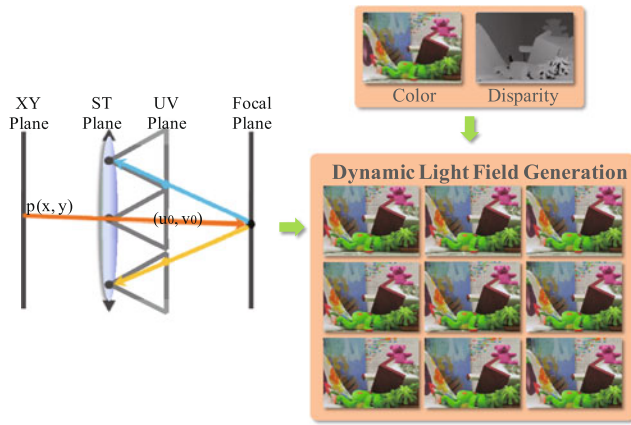


Fig. 9 We synthesize an in-lens light field (*left*) from the recovered high-resolution color image and disparity map (*right*)

The image formed by a thin lens is proportional to the irradiance at a pixel a [17], which can be computed as a weighted integral of the incoming radiance through the lens:

$$a(x, y) \approx \sum_{(s,t)} L_{in}(x, y, s, t) \cdot \cos^4 \phi \quad (3)$$

To map the in-lens light field to the out-of-lens light field, it is easy to verify that pixel $a(x, y)$ on the sensor maps to pixel $(u_0, v_0) = (w - x, h - y)$ in R_{00} . Therefore, if we want to focus at the scene depth whose corresponding disparity is γ_f , we can find the pixel index in camera R_{st} using Eq. (2). The irradiance at a can be approximated as

$$a(x, y) = \sum_{(s,t)} L_{out}(s, t, u_0 + s \cdot \gamma_f, v_0 + t \cdot \gamma_f) \cdot \cos^4 \phi.$$

To estimate the attenuation $\cos^4 \phi$ term, we can directly compute $\cos^4 \phi$ for each ray (s, t, u, v) . Notice that the ray has direction $(s, t, 1)$. Therefore, we can compute $\cos^4 \phi = \frac{1}{(s^2 + t^2 + 1)^2}$.

4.2 CUDA implementation

To synthesize the light field from the reference camera R_{00} and its disparity map, we warp it onto the rest light field cameras using Eq. (2). Note that inverse warping is impractical here because the disparity maps of target light field cameras are unknown. Therefore, we choose to forwardly constructing those cameras. A naive approach would be to directly warp the RGB color of each pixel $a(u_0, v_0)$ in R_{00} onto other light field cameras. Specifically, using a 's disparity value, we can directly compute its target pixel coordinate in camera R_{st} using Eq. (2). Since the CUDA architecture supports parallel write, we can simultaneously warp all pixels in R_{00} onto other light field cameras.

Although the warping process is straight forward, attention needs to be paid to the correctness of the light field. Since multiple pixels in R_{00} may warp to the same pixel a in

the light field camera R_{st} , a depth comparison is necessary to ensure the correct visibility. Thus, each light field camera requires an additional depth buffer. To avoid write-write conflicts in the warping process, we use atomic operations. However, current graphics hardware cannot handle atomic operations on both color and depth values at the same time. To resolve this issue, we only choose to warp the disparity value. We can easily index the RGB value for each light field ray using the stored disparity value and the camera parameters. This solution requires less video memory as the RGB value does not need to be stored in the light field.

Due to speed requirements, we can only render a small light field with 36 to 48 cameras at a 640×480 image resolution. The low spatial resolution leads to strong aliasing artifacts due to undersampling. Since our reference view does not contain information from the occluded regions, the warped light field camera images will contain holes.

To reduce the image artifacts caused by undersampling and occlusions, we develop a simple technique similar to the cone tracing method to prefilter the reference view [7]. Our method is based on the observation that out-of-focus regions exhibit most severe aliasing artifacts and occlusion artifacts since they blend rays corresponding to different 3D points. Our method compensates for undersampling by first blurring the out-of-focus rays and then blending them. A similar concept has been used in the Fourier slicing photography technique for generating a band-limited light field [12].

To simulate low-pass filtering in light field rendering, we first generate a Mipmap from the reference image using a 3×3 Gaussian kernel [8].

Gaussian Mipmaps eliminate the ringing artifacts and produce smoother filtering results than the regular box-filters. We then integrate the Gaussian Mipmap into the light field ray querying process.

Assume the scene is focused at depth d_f . For a ray (u, v) in camera R_{st} that has depth value d_r , we have a similitude relationship:

$$C_{lens}/C_{blurdisk} = d_f/(d_r - d_f) = (\gamma_r - \gamma_f)/\gamma_f \quad (4)$$

where C_{lens} is the diameter of the aperture and $C_{blurdisk}$ is the size of the blur disk in world space. The MipMap level for the ray can be calculates as

$$\begin{aligned} l &= \log_2(\gamma_r \cdot C_{blurdisk}/N) \\ &= \log_2(\gamma_r \cdot (C_{lens} \cdot \gamma_f / (\gamma_r L - \gamma_f)) / N) \\ &= \log_2(C_{lens} \cdot (\gamma_f - \gamma_r) / (B \cdot N)), \end{aligned} \quad (5)$$

where N is the number of samples, γ_r gives the pixel per length ratio which transform the size of the ray cone $C_{blurdisk}/N$ into number of pixels on the image.

4.3 Our technique vs. single-image blurring

Compared with single-image methods that apply spatially varying blurs, our light field based DoF synthesis tech-

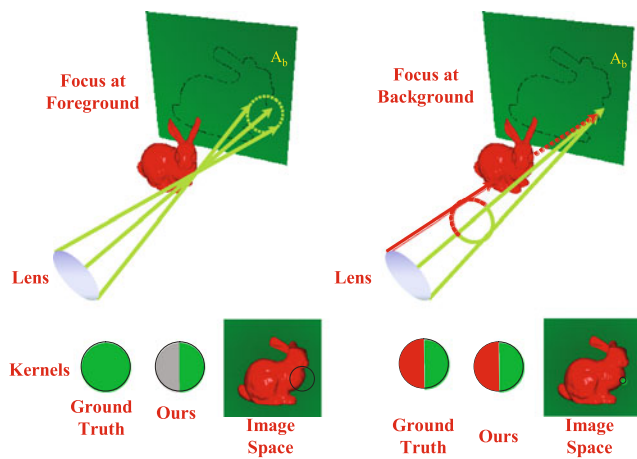


Fig. 10 Illustrations of two types of boundary artifacts. See Sect. 4.3 for details

nique significantly reduces two types of boundary artifacts. In instances where the camera focuses at the foreground, the ground truth result should blend points on the background. Conversely, single-image filtering techniques use a large kernel to blend the foreground and background pixels and hence, produce the *intensity leakage* artifact. Consider a point A_b lying on the background near the boundary, as shown in Fig. 10. Our method attempts to blend rays originating from the background. Although our technique can only access a portion of them due to occlusions, it still produces reasonable approximations.

In instances where the camera focuses on the background, the ground truth result should blend both the foreground and background points. Single-image filtering techniques, however, would consider A_b in focus and hence directly use its color as the pixel's color. In this case, the transition from the foreground to the background appears abrupt, causing the *boundary discontinuity* artifacts. Consider a point A_b on the background near the occlusion boundary in the image as shown in Fig. 10. Since rays originating from both the foreground and background are captured by our synthesized light field, our technique will produce the correct result.

Figure 8 compares the rendering results using our method and the single-image filtering approach on an indoor scene. Our technique exhibits fewer visual artifacts compared to the single-image filtering method, especially near the boundary of the girl. When examining the boundary of the sweater, the single-image method blurs the black sweater regions into the background and thus causes color bleeding, whereas our technique prevents such leakage. When focusing at the background, the single-image method exhibits discontinuous transitions from the girl to the background while our method preserves the smooth transition.

Our method also correctly preserves the boundaries between the in-focus and out-of-focus regions when synthe-

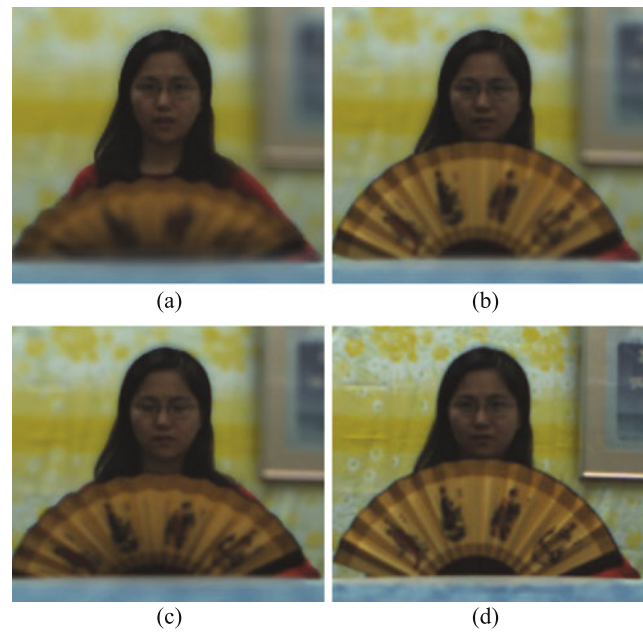


Fig. 11 Results of synthesizing changing aperture sizes. The aperture size gradually decreases from (a) to (d)

sizing changing aperture sizes. As shown in Fig. 11, we fix the focus at the woman. With the aperture fully open in (a), the blur level decreases as we decrease the aperture size.

5 Tracking and auto-refocusing with a stereo pair

One challenging task for shooting a dynamic scene with movie cameras is that it is hard to focus exactly on moving objects. Since the resolution of movie camera's viewfinder is relatively small, it is hard to tell whether the object of interest is sharp or blurry until postprocessing stage. To resolve this issue, we implemented a tracking and auto-refocusing functionality in our system.

5.1 Tracking

Like all the other classic tracking algorithms, we model our problem by reasoning probabilistically about the world based on Bayes's theorem. Since we have two images as the input, the posterior probability can be represented as

$$p(W|I_1, I_2) = \frac{p(I_1, I_2|W)p(W)}{p(I_1, I_2)},$$

where W is the latent scene, $p(I_1, I_2)$ is treated as normalizing constant, I_1 and I_2 are the images seen, and

$$p(I_1, I_2|W) = p(I_1|I_2, W)p(I_2|W)p(W).$$

Here, we use the *maximum a posteriori* estimate to find the result. Since the underline scene W does not change during one shot, $p(I_1|I_2, W)$ could be interpreted as the warping result from one of the images using the disparity map.

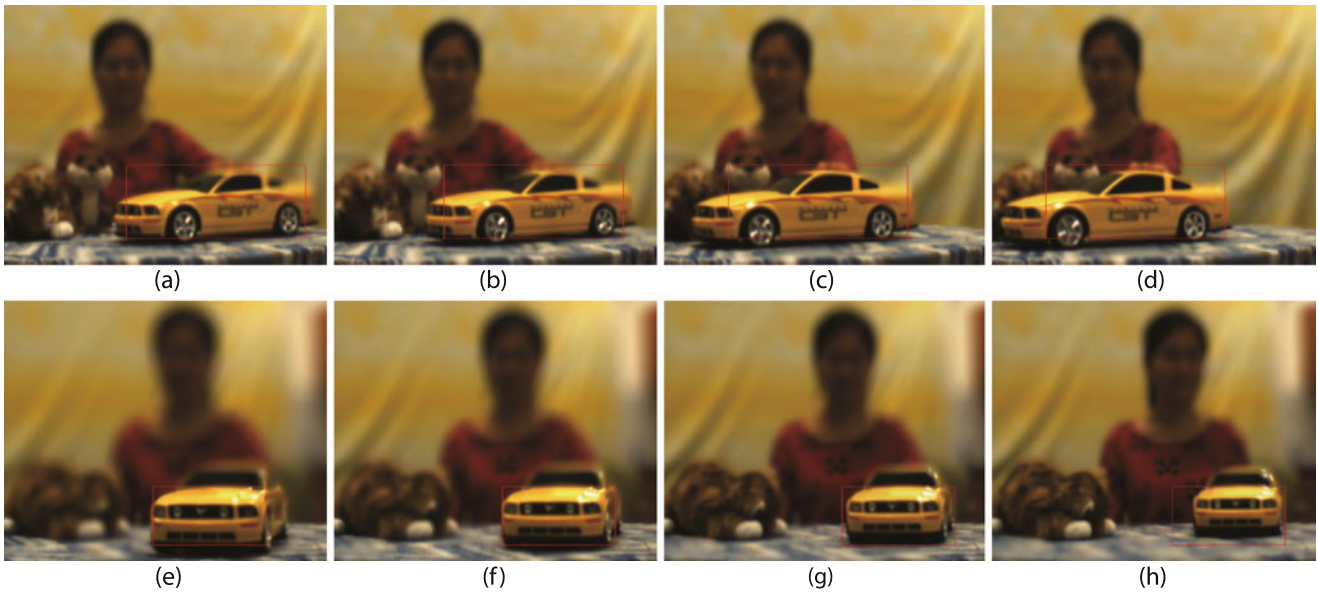


Fig. 12 Results using our tracking algorithm. Notice that with the auto-refocusing functionality, the cat on the right-hand side of the girl is becoming sharper as the toy car moves closer to its plane

Therefore, instead of dealing with multiple images, we are use both images and a disparity map as inputs. The result is estimated by $\operatorname{argmax}_W p(W|I_1, I_2)$.

We use Sum of Squared Differences (SSD) as the error function in our calculation. The estimated location of object in frame i is computed using Algorithm 1.

Algorithm 1 Compute current tracking position

```

if  $i = 0$  then
     $pos[i] \leftarrow pos[i]$ 
else
     $MinError \leftarrow INFINITE$ 
    while  $p \leftarrow next\ position$  do
         $n \leftarrow 0$ 
        for  $j = i \rightarrow \max(i - MaxLength, 0)$  do
             $e \leftarrow e + DisparitySSD(p, pos[i]) \times$ 
                 $ColorSSD(p, pos[i])$ 
             $n \leftarrow n + 1$ 
        end for
         $e \leftarrow e/n$ 
        if  $e < MinError$  then
             $MinError \leftarrow e$ 
             $pos[i] \leftarrow p$ 
        end if
    end while
end if
    
```

With the additional disparity information, the tracking result becomes very stable even though the object of interest and the background have similar colors, as shown in Row 1

of Fig. 12. The search of tracking position is also parallelized with CUDA, so the computation overhead of this step is negligible.

5.2 Auto-refocusing

With the object of interest being estimated on certain frame i , we assume that all pixels p_{ij} around pixel j inside the object should have the same disparity. A straight forward approach will be calculating the focusing disparity value $Disp$ by averaging all disparity values in this region. The result, however, is subject to noise and not robust to pixels, which are incorrectly marked as the object. To overcome these problems, we first assign different weights for pixels. Therefore, $Disp$ is computed by

$$Disp = \frac{\sum_j D_{ij} w_{ij}}{\sum_j w_{ij}},$$

where D_{ij} is the disparity of pixel p_{ij} on frame i and w_{ij} is the weight for pixel p_{ij} . The straightforward approach is assigning constant weights for all pixels. Note that user is defining the object of interest by a rectangle, pixels with different disparities or even occlusions may appear on the boundary when shooting the video. To make our disparity computation robust, here we use Gaussian weight. For each pixel in the object, we keep track of the previous assigned disparity. Since our system runs at interactive speed, we can safely assume that large disparity jumps do not occur on any pixel. If the difference between the previous and current disparities is larger than a certain threshold, we claim that this pixel is noisy and do not use it in the computation of current

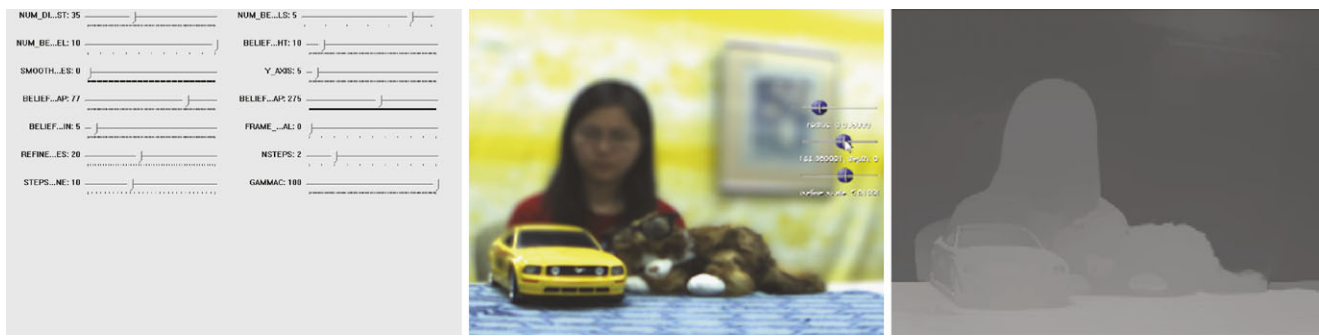


Fig. 13 The real time DoF effects (*middle*) and disparity map (*right*) given by our system after fine tuning the parameters using our interface (*left*)

Table 3 Speed up of each component in the system

Component	CPU	GPU
Depth estimation	200 ms	30–50 ms
Bilateral upsampling	100 ms	5 ms
Light field rendering	200 ms	15 ms

focused disparity $Disp$. Row 2 of Fig. 12 shows results of our auto-refocusing algorithm. The moving car stays in focus while the out of focus regions are getting sharper, such as the girl and the cat, or more blurry, such as the tablecloth in the front, as the car moves closer or further away to their planes, respectively.

6 Results and discussions

Our hybrid-resolution stereo system is connected to a workstation through a single PCI-E Firewire card. The workstation is equipped with a 3.2 GHz Intel Core i7 970 CPU, 4 GB memory and an NVIDIA Geforce GTX 480 Graphic Card with 1.5 GB memory. We implement all three processing modules (the disparity map estimation, fast CBU, dynamic DoF rendering) using NVIDIA's CUDA 3.1 with compute capability 2.0. Our system runs at the resolution of 640×480 with 15 fps. Compared with a equivalent CPU implementation at 0.2 fps, the overall speed up is over $\times 30$. Table 3 gives detailed speed up of each component in our system.

A crucial step in our real-time stereo matching module is choosing the proper parameters (e.g., the weight for the smooth/compatible terms of the energy function) to fit different types of scenes (indoor vs. outdoor). We have developed an interface to dynamically change the parameters, as shown in Fig. 13.

We have conducted extensive experiments on both indoor and outdoor scenes. We first demonstrate our system on indoor scenes with controlled lighting. Figure 14 row 1 shows four frames captured by our system of a girl drinking coffee

while reading. The coffee cup in the scene is textureless and very smooth. Our fast CBU scheme, however, still preserves the disparity edges, as shown in Fig. 14 row 1. We then dynamically change the depth of the focal plane: column (a) and column (d) focus on the front of the table, column (b) focuses on the girl, and column (c) focuses on the background. Notice how the blur varies and the in-focus regions fade into the out-of-focus regions.

Figure 14 row 2 displays a scene of a girl moving a toy car on a table. The surface of the car is specular, and the background and car have similar colors, making it challenging to prevent the disparity of the background from merging with the disparity of the car. Moreover, the motion of the car is toward the camera, causing the body of the car to have several different disparities. This makes labeling each pixel using stereo correspondence methods even more difficult. Nevertheless, our algorithm preserves the edges of the car when it is in focus and correctly blurs portions of the scene outside of the focal plane. Our system performs well indoors because the background distance is often limited, therefore, allowing one baseline to produce accurate disparity labels for the entire scene. In addition, artificially lit indoor scenes with diffuse walls and surfaces tend to have moderate dynamic range and have few poorly lit or saturated regions.

Indoor scenes undoubtedly aid the performance of our system. Our experiments on outdoor scenes, however, show promising results as well. Row 3 of Fig. 14 shows an outdoor sequence with a distant background under dynamically varying lighting conditions. Notice that in column (a), the image is brighter than the rest of the frames in the sequence and the background contains noticeable shadows. In addition to incoherent illumination, large portions of the scene such as the sky and the ground are textureless, making it difficult to achieve robust stereo matching. Since our system allows us to dynamically change the camera baseline, we use its real-time feedback to tune the parameters and increase the camera baseline to obtain a satisfactory disparity map, as shown in the supplementary video. The use of large baseline may lead to holes near the occlusion boundaries on full-resolution images. These holes are, however, less sig-



Fig. 14 Screen captures of live video streams produced by our system on both indoor (*top two rows*) and outdoor (*bottom row*) scenes

nificant in low-resolution stereo pairs and our upsampling scheme is able to borrow information from the color image to fill in the holes. The extracted frames show that we are able to correctly change the focus between the moving targets in both foreground and background.

7 Limitations

Since our system is dealing with a wide variety of indoor and outdoor scenes, it is reasonable to analyze the limitations of our frameworks and algorithms.

7.1 Scene colors

It is well known that regions without any texture cause problems when computing disparity maps. Since our system starts with a disparity matching component using Belief Propagation, visual artifacts usually show up in these regions, especially for outdoor scenes when sky is visible in the background.

Our proposed bilateral upsampling algorithm assumes that adjacent boundary regions of different depth should have different colors, in order to filter out sharp edges. If a big filter kernel is used and some parts of two or more adjacent depths happen to have similar colors, boundary bleeding artifacts will occur.

7.2 Occlusions

Occlusions on the boundaries of objects not only cause problems at the stereo matching stage, but also give unreal visual effects on the boundary regions when the Depth of Field is synthesized. This is because it is impossible to correctly recover the information of the occluded regions when the different views in the light field are warped from a single view with its disparity map, as mentioned in Sect. 4.3. Therefore, the light field rendering results on these regions is computed simply by blending the visible rays in the light field. These regions will look different if compared with the result from an optical lens.

7.3 Specular highlights

Specular highlights also introduce artifacts. As shown in Fig. 15. First, the disparity of this region is different from the neighboring regions because the high lighted portion on each object depends on the both position of the object and the light source in the scene. Second, the color of the high-light region is lost.

7.4 Translucent regions

It is impossible to compute disparity for a fuzzy region such as the hair of a person using classical stereo matching algorithms. The reason is because the model of these



Fig. 15 Observed artifacts (*high lighted with red rectangle*) at specular regions on a computed disparity map

algorithms assume that each pixel on the image has only a single disparity, which is not the case when the region is translucent, meaning rays of multiple objects are captured by the same pixel, causing the translucent effect, for example, a single hair fiber or the boundary of a human face, as shown in Fig. 16. To resolve this issue, matting needs to be performed first to separate each layer out. Then the disparity of each layer could be estimated correctly. But this method is computationally too expensive for our system.

8 Conclusion and future work

We have presented an affordable stereo solution for producing high quality live DoF effects. Our system shows promising results on indoor and outdoor scenes although it still has several limitations. First, 15 fps is a low frame rate and our resolution of 640×480 precludes our system from immediately being used in high quality HD video applications. Using multiple GPUs may address this problem as they allow greater exploitation of inherent parallelism in our computational pipeline. Second, although the high quality sensor and lens system on our camera pair significantly reduces image noise and optical distortions, this comes with a higher price. While less expensive than existing commercial movie cameras, our system is still twice the cost of most base level video cameras. Integrating existing real-time techniques to correct optical distortions and sensor noise into our pipeline would make it feasible to use lower cost webcams instead of the firewire Flea cameras.

Our future efforts include adapting our system to functional applications such as privacy protected surveillance. We plan to demonstrate the usefulness of our system in urban spaces to limit the focal plane to public areas, e.g., the sidewalks, while blurring more distant private areas like the interior of homes. Current urban surveillance networks are augmented with real-time recognition algorithms to detect illegal activity. When illegal activity is detected, our system could provide more information to law enforcement by removing the DoF effect using the stored disparity map stream for subsequent scene reconstruction. We can also

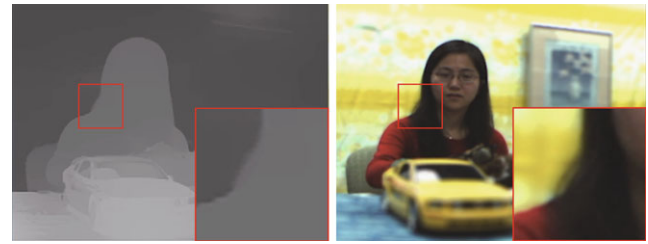


Fig. 16 Observed artifacts at translucent regions

leverage future gains in ubiquitous computing to produce a truly mobile platform which utilizes, for example, two camera phones for producing DSLR quality imagery. On the algorithm side, instead of performing a straight forward high-pass Gaussian filter to acquire high frequency information in the image, we can also perform other sophisticated frequency decomposition method such as wavelet transform.

Acknowledgements This project was partially supported by the National Science Foundation under Grants IIS-CAREER-0845268 and IIS-RI-1016395, and by the Air Force Office of Scientific Research under the YIP Award.

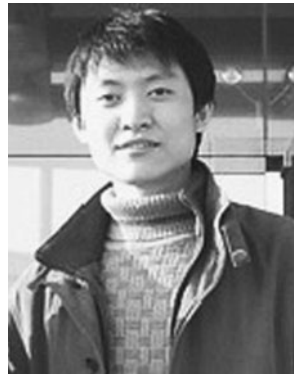
References

1. Boykov, Y., Veksler, O., Zabih, R.: Fast approximate energy minimization via graph cuts. *IEEE Trans. Pattern Anal. Mach. Intell.* **23**(11), 1222–1239 (2001)
2. Brunton, A., Shu, C., Roth, G.: Belief propagation on the GPU for stereo vision. In: *The 3rd Canadian Conference on Computer and Robot Vision* (2006)
3. Felzenszwalb, P., Huttenlocher, D.: Efficient belief propagation for early vision. In: *CVPR* (2004)
4. Grauer-Gray, S., Kambhampettu, C., Palaniappan, K.: GPU implementation of belief propagation using CUDA for cloud tracking and reconstruction. In: *2008 IAPR Workshop on Pattern Recognition in Remote Sensing (PRRS 2008)*, pp. 1–4 (2008)
5. Kolmogorov, V., Zabih, R.: What energy functions can be minimized via graph cuts? *IEEE Trans. Pattern Anal. Mach. Intell.* **26**(2), 147–159 (2004)
6. Kopf, J., Cohen, M.F., Lischinski, D., Uyttendaele, M.: Joint bilateral upsampling. In: *SIGGRAPH* (2007)
7. Lee, S., Eisemann, E., Seidel, H.P.: Depth-of-field rendering with multiview synthesis. In: *SIGGRAPH Asia* (2009)
8. Lee, S., Kim, G.J., Choi, S.: Real-time depth-of-field rendering using anisotropically filtered mipmap interpolation. *IEEE Trans. Vis. Comput. Graph.* **15**(3), 453–464 (2009)
9. Levin, A., Hasinoff, S.W., Green, P., Durand, F., Freeman, W.T.: 4d frequency analysis of computational cameras for depth of field extension. *ACM Trans. Graph.* **28** (2009)
10. Li, F., Yu, J., Chai, J.: A hybrid camera for motion deblurring and depth map super-resolution. In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR* (2008)
11. Mcmillan, L., Yang, J.C., Yang, J.C.: A light field camera for image based rendering (2000)
12. Ng, R.: Fourier slice photography. In: *SIGGRAPH* (2005)
13. Ng, R., Levoy, M., Brédif, M., Duval, G., Horowitz, M., Hanrahan, P.: Stanford tech report ctrs 2005-02 light field photography with a hand-held plenoptic camera

14. Paris, S., Durand, F.: A fast approximation of the bilateral filter using a signal processing approach. *Int. J. Comput. Vis.* **81**(1), 24–52 (2009)
15. Sawhney, H.S., Guo, Y., Hanna, K., Kumar, R.: Hybrid stereo camera: an IBR approach for synthesis of very high resolution stereoscopic image sequences. In: *SIGGRAPH*, pp. 451–460 (2001)
16. Scharstein, D., Szeliski, R.: A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *Int. J. Comput. Vis.* **47**, 7–42 (2002)
17. Stroebel, L., Compton, J., Current, I., Zakia, R.: *Photographic Materials and Processes* (1986)
18. Sun, J., Zheng, N.N., Shum, H.Y.: Stereo matching using belief propagation. *IEEE Trans. Pattern Anal. Mach. Intell.* **25**(7), 787–800 (2003)
19. Vaish, V., Levoy, M., Szeliski, R., Zitnick, C., Kang, S.B.: Reconstructing occluded surfaces using synthetic apertures: stereo, focus and robust measures. In: *CVPR* (2006)
20. Vaish, V., Wilburn, B., Joshi, N., Levoy, M.: Using plane + parallax for calibrating dense camera arrays. In: *CVPR* (2004)
21. Wang, H., Sun, M., Yang, R.: Space-time light field rendering. *IEEE Trans. Vis. Comput. Graph.* **13**, 697–710 (2007)
22. Wang, H., Yang, R.: Towards space: time light field rendering. In: *I3D* (2005)
23. Wilburn, B., Joshi, N., Vaish, V., Levoy, M., Horowitz, M.: High-speed videography using a dense camera array. In: *CVPR* (2004)
24. Wilburn, B., Joshi, N., Vaish, V., Talvala, E.V., Antunez, E., Barth, A., Adams, A., Horowitz, M., Levoy, M.: High performance imaging using large camera arrays. *ACM Trans. Graph.* **24**, 765–776 (2005)
25. Yang, J.C., Everett, M., Buehler, C., McMillan, L.: A real-time distributed light field camera. In: *Proceedings of the 13th Eurographics Workshop on Rendering, EGRW '02*, pp. 77–86 (2002)
26. Yang, Q., Yang, R., Davis, J., Nister, D.: Spatial-depth super resolution for range images. In: *CVPR* (2007)
27. Yu, X., Wang, R., Yu, J.: Real-time depth of field rendering via dynamic light field generation and filtering. *Comput. Graph. Forum* **29**(7), 2099–2107 (2010)
28. Zhang, Z.: A flexible new technique for camera calibration. *IEEE Trans. Pattern Anal. Mach. Intell.* **22**(11), 1330–1334 (2000)



Zhan Yu received his B.S. degree from the Department of Software Engineering, Xiamen University in 2008. He is now a Ph.D. student at the Department of Computer and Information Sciences, University of Delaware. His research interests include computational photography, computer vision, and computer graphics.



Xuan Yu received the B.S. degree in software engineering from Shanghai Jiao Tong University, and the Ph.D. degree in computer science from the University of Delaware in 2012. He is joining visual effect industry in Rhythm and Hue studio, Los Angeles. His research mainly focuses on real time and realistic rendering.



Christopher Thorpe received the M.S. degree in computer science from the University of Delaware in 2011. He has joined EF Technologies, Inc.



Scott Grauer-Gray is graduate student in the Department of Computer & Information Sciences at the University of Delaware. He graduated with a B.S. in Computer Science and a minor in math from the University of Massachusetts-Amherst in 2006. His current areas of interest are GPU computing, optimization, auto-tuning, and computer vision. He previously did research in database security and interoperability.



Feng Li received his B.E. from the Department of Electrical Engineering, Fuzhou University, in 2003, M.S. from the Institute of Pattern Recognition, Shanghai Jiao-tong University, in 2006, and his Ph.D. from the Department of Computer and Information Sciences, University of Delaware, in 2011. He joined Mitsubishi Electric Research Labs as a visiting research scientist in November 2011. His research interests include multi-camera system design and applications, fluid surface reconstruction, and medical imaging.



Jingyi Yu is an Associate Professor in the Department of Computer & Information Sciences and the Department of Electrical & Computer Engineering at the University of Delaware. He received his B.S. from Caltech in 2000 and Ph.D. from MIT in 2005. His research interests span a range of topics in computer vision and computer graphics, especially on computational cameras and displays. He has served as a Program Chair of OMNIVIS '11, a General Chair of

Projector-Camera Systems '08, and an Area and Session Chair of ICCV '11. He is a recipient of both the NSF CAREER Award and the AFOSR YIP Award.