# Removing Dynamic Objects for Static Scene Reconstruction using Light Fields

Pushyami Kaveti[1], Sammie Katt[1] and Hanumant Singh[2]

*Abstract*— There is a general expectation that robots should operate in environments that consist of static and dynamic entities including people, furniture and automobiles. These dynamic environments pose challenges to visual simultaneous localization and mapping (SLAM) algorithms by introducing errors into the front-end. Light fields provide one possible method for addressing such problems by capturing a more complete visual information of a scene. In contrast to a single ray from a perspective camera, Light Fields capture a bundle of light rays emerging from a single point in space, allowing us to see through dynamic objects by refocusing past them.

In this paper we present a method to synthesize a refocused image of the static background in the presence of dynamic objects that uses a light-field acquired with a linear camera array. We simultaneously estimate both the depth and the refocused image of the static scene using semantic segmentation for detecting dynamic objects in a single time step. This eliminates the need for initializing a static map . The algorithm is parallelizable and is implemented on GPU allowing us execute it at close to real time speeds. We demonstrate the effectiveness of our method on real-world data acquired using a small robot with a five camera array.

## I. INTRODUCTION

SLAM facilitates robot navigation and mapping and is one of the primary tasks in mobile robotics. Most of the research in SLAM assumes static environments, but the real world is complex and dynamic, making practical applications (autonomous motion on a crowded road or in a corridor) difficult. Algorithms tend to fail in the presence of dynamic objects due to errors in feature matching, loop closure and pose estimation. These problems have resulted in considerable research and development of techniques targeting dynamic environments [1], [2].

Dynamic scenes are handled by dividing the scene content into static and dynamic components in a variety of ways. Dynamic objects can be detected using temporal methods like dense scene flows, clustering 3D motion, epipolar constraints, moving consistency checks, by tracking changes in a static map or by using instantaneous methods such as semantic segmentation. Once detected, most approaches proceed by explicitly discarding the inputs associated with the dynamic objects as outliers in pose estimation. However, simply discarding information may fail if the dynamic portions of the image are significant occluders or tend to dominate the image in terms of feature space. Hence, extracting static features is imperative to estimating the pose of the robot accurately. Most algorithms currently require an initialization phase where they map the static landmarks and world before being able to detect and track dynamic features [3], [4], [5].

In this paper, we propose a solution to reconstruct the static scene occluded by the dynamic objects in a single time step using just semantic information of the scene and light fields as sensing modality. While light fields have seen considerable use, primarily in the computer graphics community, the applicability and usefulness of light field imaging systems for mobile robotics is far more recent [6] and [7]. Arrays of cameras can be used as a light field imaging system [8]. Cameras stand out for being robust and inexpensive with a rich history in terms of their geometric and radiometric characterization. They are also relatively easy to incorporate on most robotics systems.

A light field array captures spatial and angular radiance information at a point in space. This corresponds to a pencil of rays originating at that point and propagating in all directions as opposed to a single ray in a monocular camera. The redundancy in the pencil of rays provides information that can help us to extract the rays emerging from partially occluded portions of the scene and render synthetic aperture and digitally refocused images. In our method the static background of the scene is reconstructed not just by discarding the dynamic objects, but seeing through them via synthetic aperture refocusing. The dynamic objects are detected via deep learning based semantic segmentation. Since people are the most commonly seen dynamic class we show the results of our algorithm by detecting people in a scene. We pose the problem as a probabilistic graphical model where we jointly estimate the depth map of the static scene a well as get a refocused image of the static background using a semantic segmentation prior. We use expectation maximization (EM) [9] to refine the segmentation prior for consistent labeling of dynamic objects instead of just using the predictions from deep learning. In addition, our algorithm is completely parallelizable where each pixel can be processed independently. Our method is implemented on a GPU and is capable of running close to real-time.

## II. RELATED WORK

Our method lies at the intersection of 3D reconstruction and light field rendering. In this section we discuss some of the recent advances in these areas that are closely related.

*Static Background Reconstruction:* Most of the dense mapping algorithms perform static background reconstruction. [5] estimates a background model by accumulating

[1]Khoury College of Computer Science, Northeastern University, Boston, MA (email: kaveti.p@husky.neu.edu, katt.s@husky.neu.edu)

[2]Department of Electrical and Computer Engineering, Northeastern University,Boston, MA (email: ha.singh@northeastern.edu)
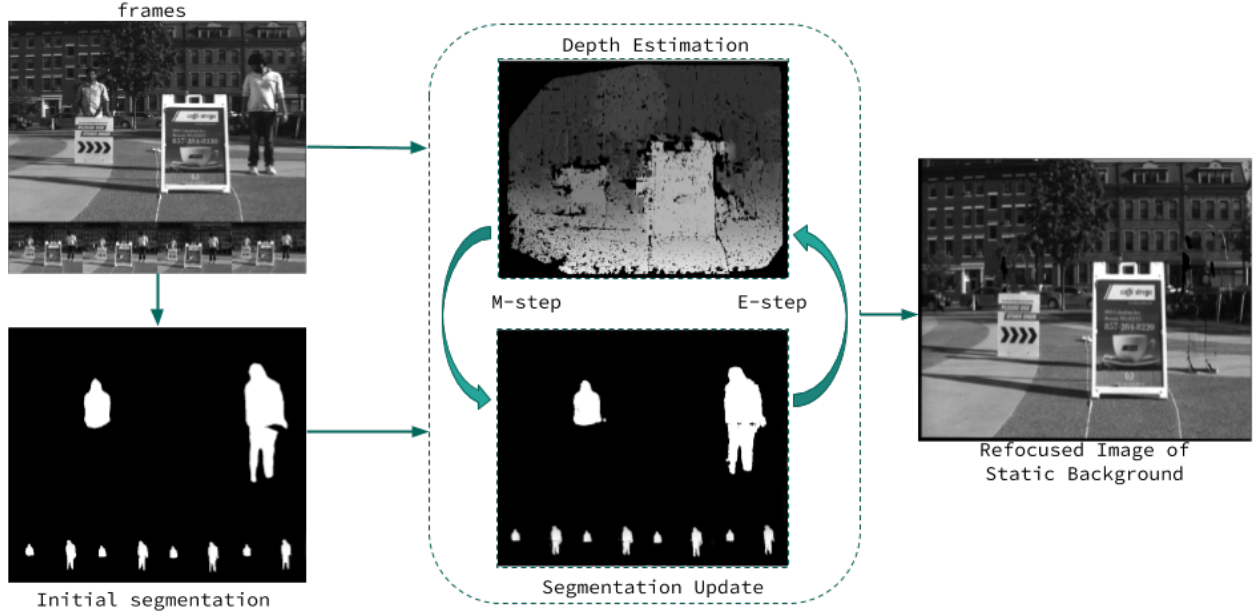
Fig. 1. Block diagram of the static background reconstruction pipeline. Frames from the linear camera array and the segmentation probabilities are given as input into the EM framework where we iteratively estimate the depth and segmentation maps. The optimized depth and segmentation is used to compute the refocused image.

warped depth between consecutive RGB-D images and applies energy-based dense visual odometry for motion estimation. StaticFusion [10] is another remarkable algorithm which not only detects moving objects, but also fuses temporally consistent data instead of discarding them. A recent approach to detect the dynamic content is by semantic segmentation using deep neural networks. Including semantic information along with geometric constraints enables detecting not just dynamic objects but also potentially movable objects [11], [3], [12], [13].

All these methods require an initialization phase to build a static map from temporal data over multiple frames and then reconstruct static background in the subsequent time steps by tracking changes. All these methods use RGB-D sensor where they already have a depth map per frame. Our goal is to compute a single dense depth map of the static background within a single processing step.

***Light-Fields for Robotics:*** Inspired by the two plane parameterization proposed by [8], light fields are a popular topic in computer vision and graphics for refocusing and rendering [14], [15], [16], super-resolution [17] and depth estimation [18], [19], [20], [21], [22]. Application of light fields in robotics is not as developed but has been discussed in [6], [23], [7]. Most of the robotics related work aims to solve the visual odometry problem [23], and to compensate for challenging light and weather conditions [24], [7], [25], [26]. None of this work addresses the problem of dynamic environments. We perform semantic guided refocusing and reconstruction to deal with dynamic objects.

Refocusing [14] is performed by blending rays emerging from a focal surface and passing through a synthetic aperture without taking into account the semantics of the rays. The light field depth reconstruction techniques typically are focused on increasing the accuracy of reconstruction compromising on speed and also do not incorporate semantics. Lastly, most of the light-field research is targeted towards micro lenslet cameras [27] which suffers from small baseline separation, but for the purpose of seeing through the dynamic objects wide-baseline arrays are more suitable. Hence, we use a linear array of five cameras similar to [7], a geometry that is most suitable for deploying on real robots.

***Multi-view Reconstruction:*** Stereo reconstruction is a well studied research area with recent methods that demonstrate real-time performance. In ELAS [28], piece-wise planar prior formed from a sparse set of matches are used to efficiently sample disparities to achieve fast computation. Inspired by this approach we also use the planar prior, but adapt the reconstruction problem to a multi-view setup incorporating scene semantics. Another important aspect of our approach is to enhance the segmentation labels produced by deep learning models so that the labels are consistent and align with 3D structure. We believe that semantic segmentation and depth reconstruction benefit from each other as shown by [29] and [30] where they jointly solve the segmentation and 3D reconstruction problems. Although these methods provide promising solutions, they do not use semantic information to mask out specific objects and are not meant for practical applications with real-time constraints. Next, We propose a solution addressing these issues.
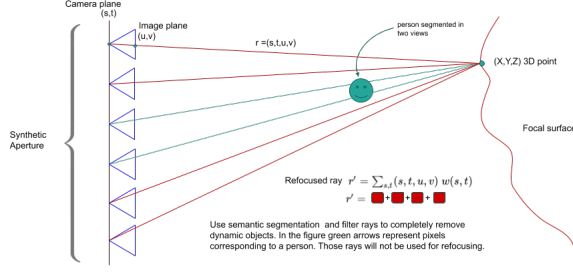
Fig. 2. Semantic refocusing using a light field array. Given a 3D point on the focal surface the rays projected into other cameras that correspond to static objects based on segmentation are combined to create a single refocused image.
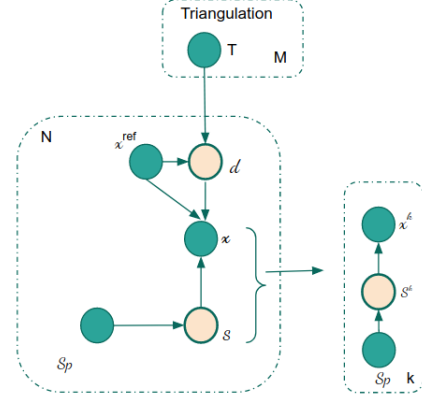


Fig. 3. The probabilistic graphical model. The known variables are colored solid green ($x^ref, x, S_P, T$) and the variables to be estimated are colored white ($d, S$). The Triangulation $T$ and segmentation prior $S_P$ are computed apriori. $d$ is disparity of a pixel in reference view, $x$ is a set of K reprojections into other cameras.

## III. LIGHT FIELD RECONSTRUCTION

This section describes our approach to light field based static scene reconstruction for dynamic environments.

### A. Problem Setup

Suppose we have a linear array of K cameras forming images $I_k$, calibrated for both intrinsic and extrinsic parameters. This setup can be seen as a single system with synthetic aperture where each camera contributes to a ray passing through the aperture. Following the two-plane parameterization of light field [14], the camera positions on the array as (s,t) coordinates lie on the entrance/camera plane and the pixels on the image plane are (u,v) coordinates, forming a unique 4D ray (s,t,u,v) as shown in fig. 2. Thus, given an arbitrary focal surface F and a 3D point (X,Y,Z) on it, we can get all the rays emitting from that point and intersecting the cameras by applying the calibrated extrinsic parameters and a projective mapping combined with camera intrinsics. These rays form our synthetic aperture and can be combined via the weighted sum of the sample weights w(s,t) to compute a single refocused ray. The size of synthetic aperture represents the angular spread of the rays and depends on the separation between the cameras on the array.

Synthetic aperture reconstruction helps to simulate varying focus and depth of field. In free-space, all the camera rays correspond to the same point on the focal surface. But, in case of occlusions, some rays will be obstructed. A large synthetic aperture provides angular spread and helps us to see-through foreground objects. Choosing a focal surface on behind the foreground occluder brings the background into focus and causes the occluder to blur. This was applied by [7] to see through rain and snow by manually selecting a focal plane on the road signs for autonomous navigation.

Given the depth map of the static background $d^{ref}$ and pixel coordinates $x^{ref}$ in the reference view, and the calibration parameters of the cameras array, then the rays $x^k$ in the other cameras $k$ can be sampled using:

$$x^k = \pi_k[R_k|t_k]\pi_{ref}^{-1}(x^{ref}, d^{ref}) \qquad (1)$$

where $\pi_k$ is a projective mapping between the a 3D point in space and 2D pixel coordinates on the image plane, and $R_k$ and $t_k$ are the rotation and translation of the $k^{th}$ camera.

These rays are typically combined by applying an average filter [15], [14] giving equal weight to all the rays, causing a foreground blur. Instead of using equal weights to all the rays we design a filter which selects only the rays that map to static pixels in the respective views obtained from semantic segmentation. This way, only the rays corresponding to the static background are considered and the dynamic objects in the foreground are completely eliminated. Denoting whether a pixel $x^k$ belongs to the static scene or the dynamic object with $s^k$, the refocused image $I^*$ can be calculated as follows:

$$I^s = \frac{\sum I(x^k) * (s^k = static)}{\sum(s^k = static)} \qquad (2)$$

As we can observe in eq. (1), the depth map of the static background is required to compute the refocused image. Previous work in this regard either required manual selection [7] or pre-computation [3] of the static background, both of which are not suitable for real time applications. In this paper we calculate the depth map of the static scene and the refocused image simultaneously without any prior knowledge of the static background. To do so, we frame our problem as a probabilistic graphical model and perform MAP estimation. We further show that we can improve the semantic segmentation maps by modeling them as hidden variables within our probabilistic framework and using EM as solution method.

### B. Probabilistic Graphical Model

The main challenge in the depth estimation of the static background is that some parts of the scene are occluded by dynamic objects (people) in some views but may be visible in others. Given an array of cameras acting as the source of image data, we need to estimate the depth map $D^*$ and refocused image $I^*$ of static background for a reference camera view $X^{ref}$. The key to finding correct depth is to choose the subset of camera views which exclude the pixels

corresponding to dynamic objects while computing the image correspondences. We utilise scene semantics to determine static and dynamic pixels. Each camera image is used to compute per-pixel semantic labels based on a deep learning model. These semantic labels might not always be perfect and consistent across camera views due to illumination factors, photo-metric properties of the scene and occlusions. The CNN model also assigns a probability to pixels being segmented as dynamic or static: $\{S_P \in \mathbb{R}^k | 0 \leq S_P^k \leq 1\}$, where $S_P^k = 0$ means *dynamic* in $k^th$ view. Instead of using the segmentation labels as ground truth the probabilities assigned to pixels are used as a prior. We introduce a set of binary random variables $S^k$ assigned to each pixel $i$ in camera $k$ which take values either 0 or 1 representing dynamic or static pixel. These variables are conditioned on the segmentation prior $S_P$ obtained from CNN and will be inferred from the graphical model.

The probabilistic graphical model is shown in the fig. 3. In this model the disparity $d$ can be obtained by maximizing the posterior probability $p(d|\mathbf{x}, x^{ref}, T, Sp)$ (MAP estimate). The joint probability from the graphical model is computed as follows (where the normalization constant can be safely ignored during maximization):

$$d^* = argmax_d \ p(d|\mathbf{x}, T, S_P, x^{ref}),$$

$$p(d|\mathbf{x}, T, S_P, x^{ref}) \propto p(d, \mathbf{x}|T, S_P, x^{ref})$$
$$= \sum_s p(d, S, \mathbf{x}|T, S_P, x^{ref}$$

In the above equation latent variables S are introduced into the joint distribution. Expectation-Maximization (EM [9]) tackles the issue of optimization with hidden (latent) variables with an iterative approach of an 'E-step' and 'M-step'. In general, assuming data $X$, the 'E-step' computes the distribution over the latent variables $Z$ according to an estimate of the parameters $\theta^{old}$: $p(Z|X, \theta^{old})$. During the 'M-step', the estimates for Z are used to update the parameters $\theta$ by optimizing $Q(\theta) = \sum_z p(Z|X, \theta^{old}) \log p(X, Z|\theta))$. The trivial extension to compute the MAP estimate is to add a prior term $\log p(\theta)$.

In our model the pixel disparities are the parameters to be optimized, and the segmentation labels are latent variables. This results in the following specification for the distributions:

$$p(Z|X, \theta^{old}) \qquad \rightarrow p(S|\mathbf{x}, d^{old}, T, S_P, x^{ref}) \qquad (3)$$
$$\log p(X, Z|\theta) \qquad \rightarrow \log p(S, \mathbf{x}|d, T, S_P, x^{ref}) \qquad (4)$$
$$\log p(\theta) \qquad \rightarrow \log p(d|T, x^{ref}) \qquad (5)$$

Which are applied in EM as follows:
1) $S \leftarrow \arg\max$ eq. (3)
2) $d \leftarrow \arg\max \sum_S$ eq. (4) $\times$ eq. (5)
3) go to (1) with $d^{old} \leftarrow d$

Since it is more natural to pick an initial segmentation assignment (according to the prior) rather than a depth-estimate, we discuss (and apply) the M-step first.

## C. Disparity Estimation (M-step)

Here we are interested in optimizing item 2, given a current estimate for the hidden variables $S$ given by the E-step. We assume a hard assignment for the hidden variables in E-step, which means that there is only one configuration of $S$ with non-zero probability. Consequently, the summation over $S$ in $Q(d)$ item 2 collapses into a single term $\log p(S, \mathbf{x}|d, T, S_P, x^{ref})$, which corresponds to the complete log likelihood of the data and latent variables. According to the graphical model fig. 3, this factorizes into the segmentation prior and pixel likelihood:

$$\log p(S, \mathbf{x}|d, T, S_P, x^{ref}) \propto \log \left[ p(\mathbf{x}|d, x^{ref}, S) + p(S|S_P) \right]$$

$S_P$ is irrelevant during the optimization (thanks to the collapse to a single $S$ assignment). As a result, the M-step reduces to picking the disparity that optimizes:

$$\arg\max_d \ \log \underbrace{p(\mathbf{x}|d, x^{ref}, S)}_{\text{pixel likelihood}} + \log \underbrace{p(d|T, x^{ref})}_{\text{disparity prior}} \qquad (6)$$

We discuss these two factors in the following paragraphs.

*1) Disparity Prior:* We use a piece-wise planar prior over the disparity space $p(d|T, x^{ref})$ by forming triangulation on a sparse set of points similar to ELAS [28]. The advantage of this prior is that it helps with poorly textured regions and gives a coarse disparity map, thus reducing the search space during optimization.

In ELAS, first a sparse set of unique points are detected and matched along the full range of disparities on the epipolar lines. These support points, together with their disparities, are used to form delaunay triangulation. Unfortunately, some of the triangulation would consider the dynamic objects, so in our algorithm we exploit the knowledge from the segmentation to filter out the support points that lie on the dynamic objects. Some portions of the static background that is occluded by foreground dynamic objects might be observed from other views. So we detect support points in other camera views, compute the disparity with respect to their neighbouring camera and choose the points that re-project on to the occluded portions of the reference image. We then filter out duplicates and inconsistent support points based on their disparity values. The final set of support points are used as vertices for the delaunay triangulation. Specifically, the prior on disparity $p(d|T, x^{ref})$ is a combination of a uniform distribution and a sampled Gaussian centered about the interpolated disparity from the triangulation, for support points in the neighbourhood.

*2) Pixel Likelihood:* The triangulation prior provides a coarse map of the interpolated disparities and a set of candidate disparities which are here used for accurate estimates based on the pixel likelihood. Given the coordinate $x^{ref}$ in the reference frame and a candidate disparity $d$, the corresponding coordinates of the source images $x^k$ can be determined using a warping function $\mathcal{W}_k(x^{ref}, d)$. This
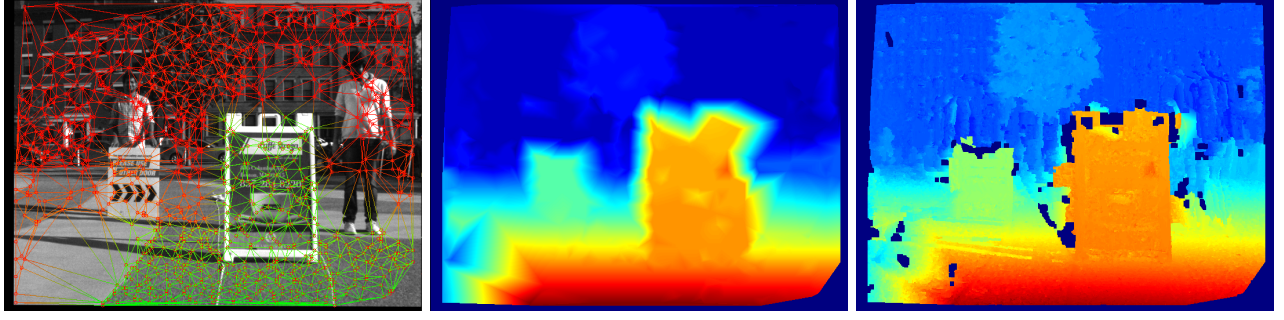
Fig. 4. Various steps involved in depth estimation (M-step) of the algorithm. *left:* Triangulation of the sparse set of unique support points colored based on their disparities (white: close, black : Far). *center:* The coarse disparity map formed by the piece-wise planar prior. *right:* The final refined disparity map of the static background.

warping function represents a homography which can be computed from the reference and $k^{th}$ camera matrices [31]. This allows us to use generalized disparity space suitable for multi-view configuration as opposed to a classic multi-baseline setup where cameras perfectly placed in a plane perpendicular to their optical axes. We work in the disparity space as opposed to depth space as it facilitates discrete optimization resulting in fast computation. The pixel likelihood relies on the fact that for correct disparity value there is a high probability that warped static pixels $x^k$ in the source images will have photo metric consistency. So, we design our likelihood function as a Laplace distribution restricted to only the static pixels such that the variance between them is minimum. The dynamic pixels don't contribute to the likelihood as they can have arbitrary intensities. As a result, the likelihood is modelled as follows:

$$p(x^1...x^K \mid d, x^{ref}, S) \propto \begin{cases} exp\big(-\beta\,Var(f(x^1),...,f(x^K))\big) \\ \quad for\ \ x^k = \mathcal{W}_k(x^{ref}, d) \\ 0 \quad otherwise \end{cases}$$

Where the variance is computed on the feature descriptors $f(x^K)$ of the pixels that are classified as static:

$$Var(...) = \frac{\sum_{k=1}^{K}(S^k=1)[f(x^k)-\hat{f}]^2}{\sum(S^k=1)}$$

with $\hat{f} = \frac{\sum_k(S^k=1)f(x^k)}{\sum_k(S^k=1)}$ as the mean of the descriptors. We use the descriptors created from 3x3 sobel filter responses similar to ELAS[28].

Plugging in the disparity and pixel likelihood formulae into eqs. (4) and (5) reveals the optimization problem of the M-step:

$$E(d) = \beta\,Var(f(x^1),...,f(x^k))$$
$$- log\left[\gamma + exp\Big(-\frac{[d-\mu(T,x^{ref})]^2}{2\sigma^2}\Big)\right]$$

In practice this is solved by considering all $2^k$ configurations of $S$.

### D. Segmentation update (E-step)

Assuming the segmentation algorithms is (near-) perfect, the solution to the graphical model described above is relatively straight forward. Unfortunately this proves difficult for many real-world scenarios as mentioned earlier. Modelling the segmentation as a hidden variable and the segmentation algorithm as a prior induces robustness, but comes at a cost of increased complexity of the E-step: the computation of the new segmentation distribution given the current estimate of disparity $d^{old}$:

$$p(S|d^{old}, \mathbf{x}, T, S_P, x^{ref}) \propto p(S, \mathbf{x}|d^{old}, T, S_P, x^{ref})$$
$$= \underbrace{p(x|d^{old}, S, x^{ref})}_{\text{pixel likelihood}}\underbrace{p(S|S_P)}_{S\ \text{prior}} \quad (7)$$

The first term term, the pixel likelihood, has been discussed prior (in the M-step). The segmentation posterior is assumed to factorize into its independent pixels, $p(S) = \prod_k p(S^K)$, where $p(S^K)$ is modelled as a delta function (e.g. 'hard assignment'). It is feasible to enumerate over these, and hence the E-step results in finding the most likely assignment of $S$: the one that maximizes (where we have substituted eq. (7) with their distributions):

$$E(S) = \beta\,Var(f(x^1),...,f(x^k))p(S|S_P)$$

This results in a more accurate estimation of segmentation, which in turn helps produce a better refocused image. The number of EM iterations depends on the initial quality of the segmentation algorithm.

### E. Refocused Image Synthesis

Once the EM optimization detailed above converges the estimated depth map and updated segmentation maps are used to compute the refocused image of the static background using eq. (2). We can observe that for every EM iteration where the depth map is calculated the refocused image is also computed. It usually takes 2-3 iterations for the EM to converge since we have a decent prior on segmentation, but doing EM helps us to enhance the depth and refocused image on the borders of the segmentation maps. Since the refocused image is calculated pixelwise, it tends to have specular noise

Fig. 5. Segmentation map update. *left*: The initial segmentation prior from Bodypix. *center*: Segmentation labels obtained by simple thresholding. *right*: Final updated segmentation map.

due to noise in the depth image. A median filter is applied to the refocused image to get rid of these artifacts.

### F. Implementation details

As mentioned earlier we restrict the classes of dynamic objects to just people and use Bodypix [32] to detect and generate pixel wise probabilities indicating the dynamic and static portions of the scene. However, our method is applicable to any moving or potentially movable objects given an appropriate segmentation. During first EM iteration $S_i^k$ values are set by thresholding the segmentation prior $S_P >= 0.7$ as suggested [32]. The depth map and refocused image of the static background are computed in a reference view. Without loss of generality we consider the left most image $X^1$ as the reference image. Thus we will have an identity transformation between $X^{ref}$ and source image $X^1$ mapping the pixel back to itself.

During the energy minimization step we work in generalized disparity space to facilitate discrete optimization. we have an array of cameras where a unit shift in disparity between reference view and a camera may result fractional shifts with respect to another camera. Thus, while considering disparities for optimizing the energy function we make sure to include non-integral disparity values to account for multi baseline properties. The algorithm is implemented using cuda on a GPU where each pixel is independently processed for depth estimation, segmentation update to finally compute the refocused image. This results in a highly parallel real-time algorithm. Refocusing can also be performed only on the dynamic pixels in the reference frame as the other pixels already have the static background image. This further increases the speed and the computational complexity scales with respect to the amount of dynamic content in the scene.

### IV. EXPERIMENTAL EVALUATION

In this section we show the results of our static background reconstruction algorithm on real-world sequences collected with a custom-built light field array.

### A. System Setup

Light field acquisition is done via two methods - using a large array of cameras [8] or by using a micro lenslet array in front of image sensor [27]. Large two dimensional arrays



Fig. 6. The data collection system with the custom-built linear camera array mounted on husky.The cameras are hardware synced with a master-slave architecture.

are impractical to mount on most robots while the lenslet cameras suffer from limited parallalax due to small baseline separation. Our approach, based on constraints associated with real mobile robots uses a custom built five camera linear array that we show is sufficient for practical robotic SLAM applications. All cameras are hardware synced for synchronized image capture. The array uses a 1.6MP Pointgrey BlackflyS global shutter cameras operating at 20 fps. The camera is calibrated for camera intrinsics and extrinsics with a 9x16 checkerboard pattern using the Kalibr multi-camera calibration package[33]. The light field data is collected by mounting the array on a Clearpath robotics' Husky UGV running ROS. All our experiments were conducted on an Microsoft Surface Pro laptop with a 6GB GEForce GTX 1060 GPU. The algorithm runs at near real-time at 2-3 frames per second. We can achieve faster frame rates by running the algorithm such that we are reconstructing only the image portions segmented as persons which brings the performance to 10 fps.
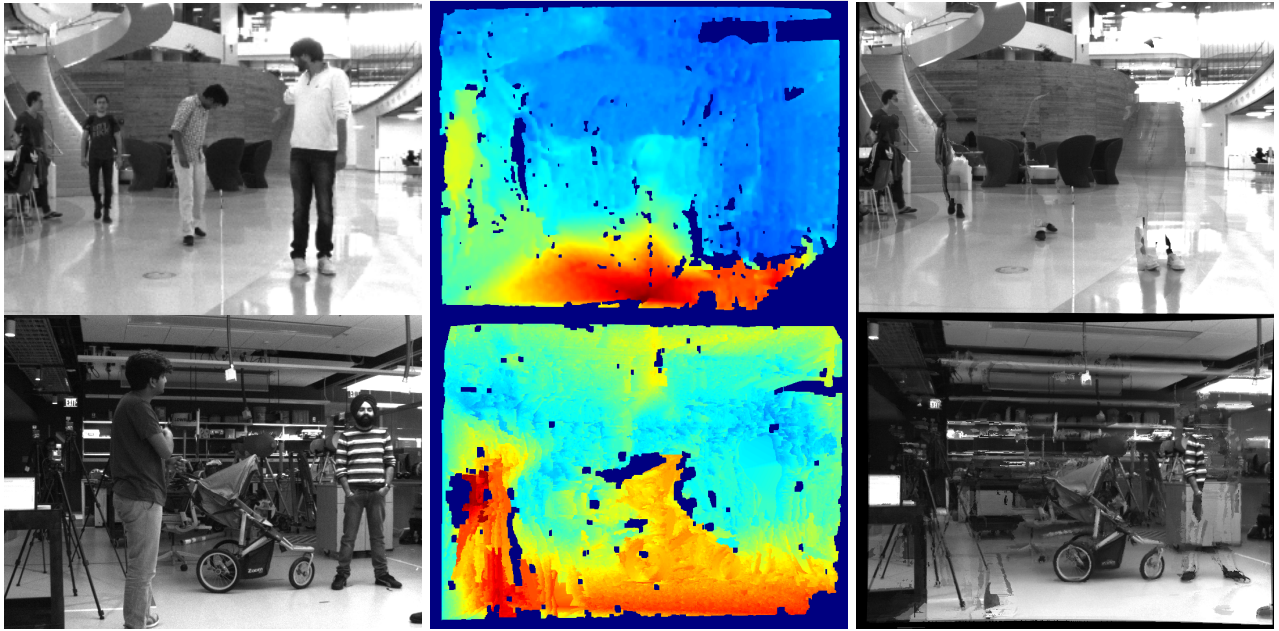
Fig. 7. *Left:* Original images from the reference view. *Center:* Disparity map of the image without people. *Right:* Refocused image of the static scene estimated with our algorithm. Top row shows the result in low texture regions with plain walls. The bottom row shows the reconstruction result in present of cluttered scene.

## B. Indoor and Outdoor datasets

We tested our approach on various datasets collected both indoors and outdoors in real-world environments with people moving in random directions. All the datasets consist of images of 720x540 pixels. We used the left most camera as the reference view for reconstruction but this choice is arbitrary and any camera can be used as the reference. We have shown the results of our static view synthesis on the outdoor dataset in fig. 4. We also tested in scenes which has relatively more people occluding majority of the background and the algorithm does a good job of reconstructing such a complex scene. The indoor data shown in top row of fig. 7 presents a case of low textured regions and shows that the algorithm can handle such a situation quite well due to the piece-wise planar triangulation prior. In the second indoor scene we can observe the quality of the reconstruction in the presence of significant detail. In the reconstructions we can see that there are some portions of the image which still retain the original image data even though it is segmented as a dynamic object eg: feet of the person, some portions of upper body in bottom row image of fig. 7. This happens when there is not enough paralallax and none of the rays can reach the background. This is a limitation of our camera array and limitations of the baseline separation. However, as we can see in the imagery we can recover the majority of the static scene.

In fig. 5 the segmentation update results are shown between the original segmentation from Bodypix and the enhanced segmentation after EM. We can observe that the updated segmentation is more consistent with the 3D structure as well across the multiple camera views which in turn helps provide an improved refocused image.

## C. Discussion

There are some advantages and shortcomings that fall out of our system that are wroth pointing out. We note that the refocused image is obtained by combining only the rays that reach the static background. If the person in the foreground is occluding the static background in such a way that no rays can reach the background, for example if the person is too close to the wall or too big for our choice of baseline separation, we will not be able to recover the background. In these cases, having temporal information could help in the reconstruction of the occluded portions of a static scene.

Note also that shadows are not usually picked up by the segmentation algorithms but are strong candidates for key points. Even though one would detect the dynamic objects using semantic segmentation and just discard the key point associated with them, moving shadows will cause significant errors in stable localization. Our approach, however, is indifferent to shadows as we are reconstructing the static 3D scene at every time step and using this directly so that shadows do not affect us at all.

Note that we have also arranged our cameras in a horizontal linear array to form the light field array. This naturally allows us to focus on dynamic classes that usually appear distributed vertically in the scene and for which we require horizontal parallalax to see through them.

## V. CONCLUSION AND FUTURE WORK

In this paper we presented a method for reconstructing the depth and 2D image of the static background from a reference view using a 4D linear light field array. We formulate this problem in a probabilistic framework where we perform an EM based optimization to estimate the depth

and refocused image of the static scene and at the same time improve the semantic segmentation masks obtained from a deep learning model so that they are consistent with the 3D structure. We show promising results by evaluating our algorithm on real-world data sets collected both indoor and outdoor with people. This can be a potential front-end for SLAM to deal with dynamic environments. Through this work we have shown that light field proves to be a good candidate for robot sensing and navigation. The main advantages of our approach are 1) We do not need an initialization phase or an initial static map for localization as we are not tracking changes in the scene. 2) Our algorithm is parallelizable and capable of running at close to real-time speed and we have demonstrated it working at 2-3 frames per second on a laptop GPU. An obvious next step would be to incorporate temporal information to get a complete dense map of the static background and We would like to further develop and demonstrate a full end-to-end light field based SLAM solution for dynamic environments.

## ACKNOWLEDGMENT

## REFERENCES

[1] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, "Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age," *IEEE Transactions on robotics*, vol. 32, no. 6, pp. 1309–1332, 2016.

[2] M. R. U. Saputra, A. Markham, and N. Trigoni, "Visual slam and structure from motion in dynamic environments: A survey," *ACM Computing Surveys (CSUR)*, vol. 51, no. 2, pp. 1–36, 2018.

[3] B. Bescos, J. M. Fácil, J. Civera, and J. Neira, "Dynaslam: Tracking, mapping, and inpainting in dynamic scenes," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 4076–4083, 2018.

[4] W. Tan, H. Liu, Z. Dong, G. Zhang, and H. Bao, "Robust monocular slam in dynamic environments," in *2013 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, 2013, pp. 209–218.

[5] D.-H. Kim and J.-H. Kim, "Effective background model-based rgb-d dense visual odometry in a dynamic environment," *IEEE Transactions on Robotics*, vol. 32, no. 6, pp. 1565–1573, 2016.

[6] F. Dong, S.-H. Ieng, X. Savatier, R. Etienne-Cummings, and R. Benosman, "Plenoptic cameras in real-time robotics," *The International Journal of Robotics Research*, vol. 32, no. 2, pp. 206–217, 2013.

[7] A. Bajpayee, A. H. Techet, and H. Singh, "Real-time light field processing for autonomous robotics," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 4218–4225.

[8] M. Levoy and P. Hanrahan, "Light field rendering," in *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996, pp. 31–42.

[9] C. M. Bishop, *Pattern recognition and machine learning*. springer, 2006.

[10] R. Scona, M. Jaimez, Y. R. Petillot, M. Fallon, and D. Cremers, "Staticfusion: Background reconstruction for dense rgb-d slam in dynamic environments," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 1–9.

[11] C. Yu, Z. Liu, X.-J. Liu, F. Xie, Y. Yang, Q. Wei, and Q. Fei, "Ds-slam: A semantic visual slam towards dynamic environments," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 1168–1174.

[12] I. A. Bârsan, P. Liu, M. Pollefeys, and A. Geiger, "Robust dense mapping for large-scale dynamic environments," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 7510–7517.

[13] E. Palazzolo, J. Behley, P. Lottes, P. Giguere, and C. Stachniss, "Refusion: 3d reconstruction in dynamic environments for rgb-d cameras exploiting residuals," *arXiv preprint arXiv:1905.02082*, 2019.

[14] A. Isaksen, L. McMillan, and S. J. Gortler, "Dynamically reparameterized light fields," in *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, 2000, pp. 297–306.

[15] T. G. Georgiev and A. Lumsdaine, "Focused plenoptic camera and rendering," *Journal of electronic imaging*, vol. 19, no. 2, p. 021106, 2010.

[16] A. Davis, M. Levoy, and F. Durand, "Unstructured light fields," in *Computer Graphics Forum*, vol. 31, no. 2pt1. Wiley Online Library, 2012, pp. 305–314.

[17] T. E. Bishop and P. Favaro, "The light field camera: Extended depth of field, aliasing, and superresolution," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 5, pp. 972–986, 2011.

[18] V. Vaish, M. Levoy, R. Szeliski, C. L. Zitnick, and S. B. Kang, "Reconstructing occluded surfaces using synthetic apertures: Stereo, focus and robust measures," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, vol. 2. IEEE, 2006, pp. 2331–2338.

[19] S. Wanner and B. Goldluecke, "Globally consistent depth labeling of 4d light fields," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 41–48.

[20] T.-C. Wang, A. A. Efros, and R. Ramamoorthi, "Occlusion-aware depth estimation using light-field cameras," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 3487–3495.

[21] I. K. Park, K. M. Lee *et al.*, "Robust light field depth estimation using occlusion-noise aware data costs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 10, pp. 2484–2497, 2017.

[22] H. Sheng, P. Zhao, S. Zhang, J. Zhang, and D. Yang, "Occlusion-aware depth estimation for light field using multi-orientation epis," *Pattern Recognition*, vol. 74, pp. 587–599, 2018.

[23] D. G. Dansereau, I. Mahon, O. Pizarro, and S. B. Williams, "Plenoptic flow: Closed-form visual odometry for light field cameras," in *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2011, pp. 4455–4462.

[24] K. A. Skinner and M. Johnson-Roberson, "Towards real-time underwater 3d reconstruction with plenoptic cameras," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 2014–2021.

[25] N. Zeller, F. Quint, and U. Stilla, "Scale-awareness of light field camera based visual odometry," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 715–730.

[26] J. Eisele, Z. Song, K. Nelson, and K. Mohseni, "Visual-inertial guidance with a plenoptic camera for autonomous underwater vehicles," *IEEE Robotics and Automation Letters*, vol. 4, no. 3, pp. 2777–2784, 2019.

[27] R. Ng, M. Levoy, M. Brédif, G. Duval, M. Horowitz, P. Hanrahan *et al.*, "Light field photography with a hand-held plenoptic camera," *Computer Science Technical Report CSTR*, vol. 2, no. 11, pp. 1–11, 2005.

[28] A. Geiger, M. Roser, and R. Urtasun, "Efficient large-scale stereo matching," in *Asian conference on computer vision*. Springer, 2010, pp. 25–38.

[29] A. Kundu, Y. Li, F. Dellaert, F. Li, and J. M. Rehg, "Joint semantic segmentation and 3d reconstruction from monocular video," in *European Conference on Computer Vision*. Springer, 2014, pp. 703–718.

[30] C. Hane, C. Zach, A. Cohen, R. Angst, and M. Pollefeys, "Joint 3d scene reconstruction and class segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 97–104.

[31] R. Szeliski and P. Golland, "Stereo matching with transparency and matting," *International Journal of Computer Vision*, vol. 32, no. 1, pp. 45–61, 1999.

[32] T. L. Zhu and D. Oved, "BodyPix - Person Segmentation in the Browser," 2019.

[33] P. Furgale, J. Rehder, and R. Siegwart, "Unified temporal and spatial calibration for multi-sensor systems," in *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2013, pp. 1280–1286.