

# VUNet: Dynamic Scene View Synthesis for Traversability Estimation using an RGB Camera

Noriaki Hirose, Amir Sadeghian, Fei Xia, Roberto Martín-Martín, and Silvio Savarese

**Abstract**—We present VUNet, a novel view(VU) synthesis method for mobile robots in dynamic environments, and its application to the estimation of future traversability. Our method predicts future images for given virtual robot velocity commands using only RGB images at previous and current time steps. The future images result from applying two types of image changes to the previous and current images: 1) changes caused by different camera pose, and 2) changes due to the motion of the dynamic obstacles. We learn to predict these two types of changes disjointly using two novel network architectures, SNet and DNet. We combine SNet and DNet to synthesize future images that we pass to our previously presented method GONet [1] to estimate the traversable areas around the robot. Our quantitative and qualitative evaluation indicate that our approach for view synthesis predicts accurate future images in both static and dynamic environments. We also show that these virtual images can be used to estimate future traversability correctly. We apply our view synthesis-based traversability estimation method to two applications for assisted teleoperation.

**Index Terms**—Robot safety, computer vision for other robotic applications, collision avoidance.

## I. INTRODUCTION

**A**UTONOMOUS robots can benefit from the ability to predict how their actions affect their input sensor signals. The ability to predict future states provides an opportunity for taking better actions. This ability can be applied to a variety of tasks from perception and planning to safe navigation. In robot visual navigation the actions are the velocity commands given to a robot, the input sensor signals are the images captured from the robot’s RGB camera. And the predicted future images used to better understand the consequence of actions, can be predicted using scene view synthesis methods. In this context, a view synthesis model can determine which actions bring the desired sensor outcomes or can cause future hazards.

Previous approaches have addressed the scene view synthesis problem assuming that a 3D model of the environment is available to virtually move the camera [2, 3]. Recently, several approaches have relaxed this assumption and synthesized images using only a small set of previous images and a virtual action [4, 5]. However, none of these approaches can be applied to predict images for navigation in unknown environments with dynamic obstacles. In this scenario, scene

This paper was recommended for publication by Editor Paolo Rocco upon evaluation of the Associate Editor and Reviewers’ comments. The Toyota Research Institute (“TRI”) provided funds to assist with this research, but this article solely reflects the opinions and conclusions of its authors and not TRI or any other Toyota entity. The TOYOTA Central R&D Labs., INC. supported N. Hirose at Stanford University.

The authors are with Computer Science Department, Stanford University, USA [hirose@mosk.tytlabs.co.jp](mailto:hirose@mosk.tytlabs.co.jp)

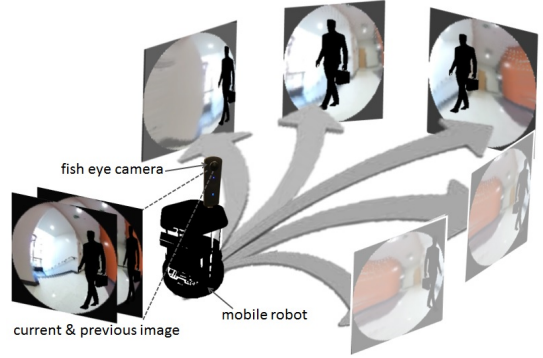


Fig. 1: Scene view synthesis for mobile robots. Our method, VUNet predicts multiple future images assuming different navigation commands and integrating the image changes caused by changes in robot pose and by dynamic objects. The input to our method are current and previous images from an on-board fish-eye camera. The robot image is cited from [11].

view synthesis is extremely challenging because it needs to account for both the changes in the camera pose and the motion of the dynamic obstacles.

Being able to predict the future state of both the static and the dynamic parts of an environment has multiple direct applications towards safe navigation. One of these applications is traversability estimation: to identify traversable and non-traversable spaces in the surroundings of the robot. Traditionally, traversability estimation methods have relied on depth sensors or on LIDARs [6, 7, 8, 9]. However, depth sensors can fail in outdoor conditions or when the surface of the static or dynamic obstacles are reflective, and LIDARs are very expensive compared to more affordable RGB cameras, or might fail to detect glasses. A new family of methods uses only RGB images to estimate traversability [1, 10]. However, these RGB-based methods do not have the predictive power to estimate the traversability of the locations the robot will need to navigate in the future.

In this work we propose a novel deep neural network-based method for dynamic-scene view synthesis, VUNet in the context of robot navigation and its application for future traversability estimation. Our synthesis method can predict the appearance of both static (e.g., walls, windows, stairs) and dynamic (e.g., humans) elements of the environment from different camera poses in future time steps. To do that, our method requires only as input the last two acquired images and a virtual navigation command, i.e., a linear and angular velocity (Fig. 1). We combine this method to predict future images with our previously presented RGB-based traversability estimation algorithm, GONet [1], into a system that identifies the traversable areas around the robot as well as the various

velocity commands that can be executed safely by the robot.

The main contributions of this work are thus twofold: First, we propose a novel dynamic scene view synthesis method, VUNet. The technical novelty of our method is the combination of two different networks that can separately model static and dynamic transformations conditioned on robot’s actions. The proposed view synthesis method outperforms state-of-the-art methods in both static and dynamic scenes. And second, we propose a system to estimate traversability in future steps based on the synthesized images. We also propose two applications of the system in assisted teleoperation: early obstacle detection and multi-path traversability estimation.

## II. RELATED WORK

We will cover in this section two main research areas: scene view synthesis, and traversability estimation.

**Scene View Synthesis** is the problem of generating images of the environment from virtual camera poses. For unknown environments, a variant of the problem assumes the only input are real images taken at a certain pose. This problem has been widely studied both in computer vision [12, 13, 14] and in computer graphics [15, 16, 17] using two main types of methods. The first type of methods synthesizes pixels from an input image and a pose change with an Encoder-Decoder structure [18, 19, 20]. The second type reuses pixels from an input image with a sampling mechanism [12]. Instead of generating pixels, this type of method generates a flow field to morph the input image. If information from multiple views is available, a smart selection mechanism needs to be used to choose which image to sample pixels from [5]. Previous methods focus on predicting either changes due to camera motion or due to dynamic objects [21, 22], but not both. Our method is able to deal with changes both in camera view and dynamic objects, making it suitable for dynamic scenes.

**Traversability Estimation:** Estimating which areas around the robot are safe to traverse has been traditionally done using Lidar or other depth sensors [6, 7, 8, 9, 23, 24]. These methods estimate the geometry of the surroundings of the robot and use it to infer the traversable areas. However, lidar sensors are expensive and depth measurements can be affected by surface textures and materials, e.g. highly reflective surfaces and transparent objects such as mirrors and glass doors. These issues have motivated the use of RGB images for traversability estimation [25, 26, 27]. Some RGB-based methods try to first estimate depth from RGB and then apply a method based on depth images [28, 29]. Other methods learn a generative deep neural network and formulate it as anomaly detection problems [1, 10]. For example, GONet[1], which we use in our system, contains a Generative Adversarial Network (GAN) trained in a semi-supervised manner from traversable images of a fisheye camera. Since the GAN only learns to generate traversable images, GONet uses the similarity between the input image and its GAN regenerated image to estimate traversability. GONet and other RGB-based methods estimate traversability only in the the space right next to the robot; our proposed approach predicts traversability for longer horizon trajectories.

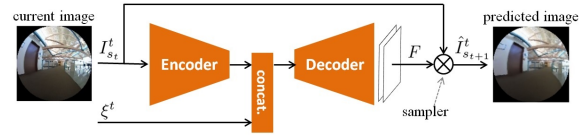


Fig. 2: Static Transformation Network (SNet) structure. The input is an RGB image,  $I_{s_t}^t$  from the robot’s fish-eye camera at time  $t$  and a camera pose  $s_t$ , and a virtual future twist velocity,  $\xi^t$ . The output is a predicted image from the future robot pose,  $I_{s_{t+1}}^t$ . The network decides how to change (sample from) the current image to generate the virtual image based on a flow field,  $F$ .

## III. METHODS

### A. Dynamic Scene View Synthesis

In this section, we introduce VUNet for view synthesis in dynamic environments. Our goal is to generate predicted images by altering both the spatial and the temporal domains (see Fig. 4a). Formally, let  $\{I_{s_t}^t\}$  be a set of consecutive images, each of them captured at different time steps  $t$  and (possibly) different poses  $s_t$ . Given this set, we aim to predict an image  $I_{s_{t'}}^{t'}$  at a new robot pose  $s_{t'}$  and a new time-step  $t'$ . Usually  $t'$  will be the next time step  $t + 1$ . Since we are working with mobile robots the pose is parameterized by robot’s position and orientation,  $s_t = (x, y, \theta)$ ,  $s_t \in \mathbb{R}^3$ . The robot command at time  $t$  is a velocity in robot frame, expressed as a twist  $\xi^t = (v^t, \omega^t)$ , where  $v^t$  and  $\omega^t$  are the linear and angular components. We assume the mobile robot is nonholonomic and the velocity is two dimensional,  $v^t, \omega^t \in \mathbb{R}$ .

Our general approach is to apply changes to the last acquired real images to generate virtual images. Changes in the image are caused by two factors: changes in the viewpoint of the camera (spatial domain) and changes due to dynamically moving objects (temporal domain). For the static parts of the environment, the image changes are only caused by the change of viewpoint, while for dynamic objects both factors contribute to the appearance change. It is difficult to learn both factors simultaneously, so we propose to learn them in a disentangled manner: we completely separate the models to predict image changes in the static and moving parts of the environment, and individually train each model. In the following subsections we will first present our model to predict changes in the static parts of the environment due to robot motion, SNet. Then we will present the model to predict appearance changes due to motion of the dynamic parts of the environment, DNet. Finally we will explain VUNet, the combination of SNet and DNet to synthesize complete images in future time steps from different viewpoints in dynamic environments.

**Static Transformation Network (SNet):** Figure 2 shows the network structure of SNet. SNet uses an image from a camera pose  $s_t$  and a virtual velocity  $\xi^t$  to predict an image from a different camera pose  $s_{t+1}$  (changes in the spatial domain). The architecture is based on the encoder-decoder architecture (ED). Our encoder-decoder has two special characteristics: 1) the virtual velocity input is concatenated to the low dimensional image representation to realize the spatial transformation before the decoding phase, and 2) the output of the decoder is a 2D flow field image ( $F$ ) that is used to sample the original input images and generate the

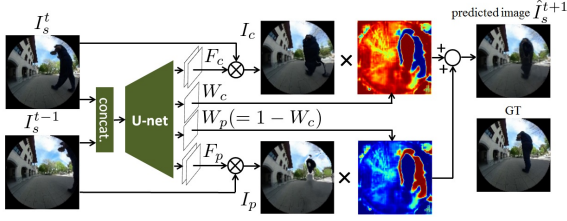


Fig. 3: Dynamic Transformation Network (DNet) structure. The input are two RGB images at time  $t$  and  $t-1$  at location  $s$ . The output is a predicted image in the next time step,  $t+1$ , from the same location considering dynamic moving objects. Our network decides how to sample from the images at previous and current time steps based on the generated flow fields  $F_p$  and  $F_c$ , and how to alpha-blend the samples based on the probabilistic selection masks  $W_p$  and  $W_c$ . Both sampled images and the probabilistic masks are depicted (red indicates high weight for the merge, blue color indicates low weight)

predicted future image. SNet generates sharper images than the classical ED architectures because the sampling procedure reuses original pixels of the input image using the internal flow field representation (see Fig. 6).

**Dynamic Transformation Network (DNet):** Figure 3 depicts the architecture of the DNet. DNet takes as input two images (real or virtual) acquired from the same camera pose ( $s$ ) in consecutive time steps and synthesizes a virtual image in the next time step. The synthesized image accounts for the changes due to the motion of the dynamic objects in the scene (changes in the temporal domain).

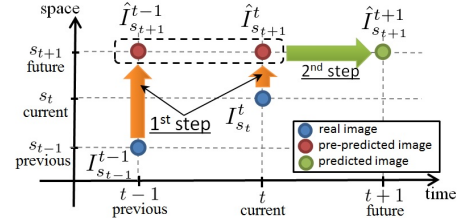
To synthesize the image, DNet generates four intermediate representations: two 2D flow field images,  $F_c$  and  $F_p$ , to sample pixels from current and previous images respectively, and two 1D probabilistic selection masks,  $W_c$  and  $W_p$ , to weight the contribution of the samples from the current and previous images in a final alpha-blend merge. We use a softmax function to generate  $W_p$  and  $W_c$  that satisfies  $W_p(u, v) + W_c(u, v) = 1$  for any same image coordinates  $(u, v)$ . The intermediate representation is generated by a U-net architecture [30] that has been successfully applied before to other virtual image synthesis tasks by Isola et al. [2].

We use two consecutive images in DNet for two reasons: 1) a single image does not have the pixel information of the parts of the environment occluded by the dynamic object, and 2) a single image does not contain motion information of the dynamic object. Using two images we can acquire pixel information behind the dynamic obstacles and also understand their motion behavior as illustrated in Fig. 3.

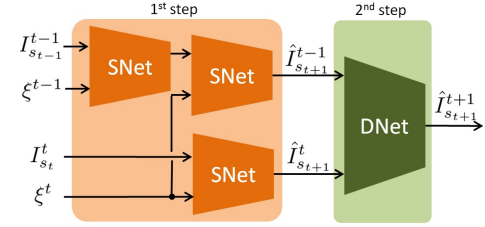
#### Dynamic-Scene View Synthesis Architecture (VUNet):

Figure 4 shows the overall structure of proposed approach, VUNet for view synthesis in dynamic environments by combining SNet and DNet. VUNet is composed of two steps. In the first step, the method applies SNet on the previous and current images to predict the virtual images of current and previous time step as they would be seen from robot's next pose,  $s_{t+1}$ . This step is shown as an orange arrows in Fig. 4a and as a light orange block in Fig. 4b. Note that, after this step the difference between the two predicted virtual images is expected to be caused only by the motion of the dynamic objects.

In the second step, our method feeds the two virtual images



(a) Generation of virtual images by altering both the spatial and temporal dimensions. First we generate virtual images from a different pose based on the real images (orange arrows). Then we generate a virtual image from that pose in the next time step for predicting the future position of the moving objects (green arrow). Blue circles indicate images at previous ( $I_{s_{t-1}}^{t-1}$ ) and current ( $I_{s_t}^t$ ) time steps.



(b) Overall network structure of VUNet for dynamic-scene view synthesis. We use SNet (twice on the previous image  $I_{s_{t-1}}^{t-1}$ , once on the current image  $I_{s_t}^t$ ) to generate virtual images from the future robot pose. Then we use DNet to generate a virtual image from the future robot pose in the next time step combining changes in the static and dynamic parts of the environment

Fig. 4: Overview of our system, VUNet for view synthesis in dynamic environments. First step: we generate virtual images as seen from the location the robot will move to ( $s_{t+1}$ ) at previous ( $t-1$ ) and current ( $t$ ) time steps using SNet. Second step: we generate a virtual image as seen from the location the robot will move to ( $s_{t+1}$ ) at the next time step ( $t+1$ ) using DNet over the previously generated spatially altered virtual images.

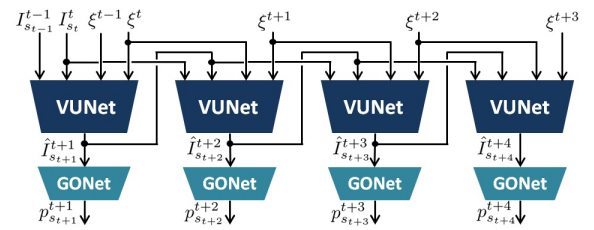


Fig. 5: System for multi-step future image prediction and traversability estimation. The input to our system are the current and previous images, the last velocity twist command  $\xi^{t-1}$  and a series of virtual velocity twists  $\xi^{t+i}$ . Each block of VUNet represent the system of Fig. 4b. They generate predicted virtual images at next time steps. These images are 1) passed to our previously presented method GONet [1] to estimate the traversable probability  $p_{s_{t+i}}^{t+i}$ , and 2) passed as input to the next VUNet block to predict the next image.

predicted by the previous step into DNet. DNet predicts the pixel changes caused by the motion of the dynamic objects and generates the final synthesized image: an image at the new pose  $s_{t+1}$  at the future time step  $t+1$ . This step is shown as a green arrow in Fig. 4a and a light green block in Fig. 4b)

By combining SNet and DNet, VUNet can predict future images that satisfy both static and dynamic changes in an environment caused by robot and dynamic objects movements.

## B. Future Traversability Estimation

We now show how we apply VUNet to the estimation of traversability in future steps. This process is composed of two steps: First, we use the current and previous acquired images as well as the last robot’s velocity ( $\xi^{t-1}$ ) and a virtual velocity ( $\xi^t$ ) to generate a first predicted future image (most left VUNet block in Fig. 5). The generated virtual image is passed to our previously proposed GONet architecture [1] – an RGB-based method that estimates the traversable probability of a scene image. The virtual image is also passed to the next VUNet block to generate a virtual image at the next time step. By repeating this process for multiple time steps, our approach is able to estimate the traversable probability at  $n$  consecutive future time steps, assumed  $n$  future virtual velocities. The general approach to estimate traversability in future steps is depicted in Fig. 5.

## IV. EXPERIMENTAL SETUP

### A. Implementation

**SNet:** In SNet we use as network structure a regular encoder-decoder (ED) similar to the network used by Isola et al. [2] but without the skip connections and the dropout. The input to the encoder is a three channel (RGB)  $128 \times 128$  image. The output of the encoder is a 512 dimensional feature vector that we concatenate with the two dimensional velocity vector  $\xi^t$  to feed it to the decoder. The decoder generates a flow field image  $F$  of size  $2 \times 128 \times 128$  that we use for bilinear sampling. The synthesized image resulting from SNet is a three channel (RGB)  $128 \times 128$  image.

**DNet:** The DNet is based on the network architecture of U-net [30] as used in [2]. Different to [2], in DNet the inputs are two 3 channel (RGB)  $128 \times 128$  images and the outputs (intermediate representation) are two  $2 \times 128 \times 128$  flow field images and two  $128 \times 128$  probabilistic masks. The final output of DNet is a three channel (RGB)  $128 \times 128$  image resulting from an alpha-blend process.

The GONet network used for future traversability estimation is a pretrained network as explained in [1]. All networks are implemented in Chainer [31] and our sampling period (time between consecutive steps) is 0.33 s.

### B. Training

To train the different components of our model we will need two different types of data: data where the robot moves in a static environment to train SNet, and data that includes dynamic objects without the robot motion to train DNet.

To train SNet we use the *GO Stanford 2* (GS2) dataset presented in [1]. GS2 contains 16 hours and 42 minutes of videos from 27 campus buildings acquired from a Ricoh THETA S fisheye camera on a teleoperated Turtlebot2 robot, and the velocity commands to this robot. Even though some few sequences in GS2 include dynamic objects, their number is very small and they do not affect the training process of SNet. We randomly flip the image and invert the angular velocity for the data augmentation to avoid overfitting in the training process.

To train DNet we record new data from a constant robot position observing dynamically moving objects (humans, vehicles, ...). We maintain the robot at a fixed position to have only image changes caused by the motion of the dynamic objects. We use the same robot and camera to record 4 hours and 34 minutes (47730 images) of videos at 46 different points in 23 different indoor and outdoor environments. We also make this new dataset “GO Stanford 3” (GS3) available to the community<sup>1</sup>.

We could train directly DNet on pairs of current and previous images from GS3. However, as explained in Sec. III-A, the input to DNet in our method is the output of SNet. SNet often generates some small disturbances. To train DNet on data with these disturbances we preprocess GS3 images passing them through a trained SNet with a small random velocity perturbation  $\epsilon_\xi$ , uniformly distributed between  $\pm 0.05$  in all dimensions.

To train both SNet and DNet we use data from separate locations (i.e. different buildings or campus areas) for training, test, and validation. This way, the evaluation on test and validation assesses how well our method generalizes to completely new environments. This location-based splits the data to 70% training, 15% test, and 15% validation.

We iteratively train all networks with a batch size of 80 using Adam optimizer [32] and with a learning rate of 0.0001. Our networks are trained by minimizing the L1 norm. For real-world experiments, our proposed system for view synthesis is implemented on a robot with a laptop equipped with a Nvidia Geforce GTX 1070 GPU that allows us to maintain 3 fps of constant computation time.

Additionally, we collected videos of teleoperated robot trajectories in dynamic environments with humans to use for the evaluation of our view synthesis and future traversability estimation methods. We recorded 26 minutes of videos in six different environments and include it as part of GS3 and more than one hour of highly dynamic environments.

## V. EXPERIMENTS

We conducted two sets of experiments. In the first set we evaluate quantitatively the performance of our view synthesis method and our system for future traversability estimation for mobile robots in dynamic environments. In the second set of experiments we evaluate all methods qualitatively.

### A. Quantitative Analysis

**DNet:** First we evaluate the performance of only DNet on the test data of GS3 where the robot is static. In this scenario it is not necessary to use SNet, because the camera viewpoint does not change. Hence, we can evaluate DNet separately. We compare DNet to several baselines: DNet variants using a regular encoder-decoder (ED) without path skips (instead of the U-net architecture), without multi-image merge, and a variant using an extrapolation based on optical flow (OF) between previous and current image using FlowNet [33]. We report mean L1 norm (pixel difference) and structural

<sup>1</sup><http://svl.stanford.edu/projects/vunet/>

TABLE I: Evaluation of DNet

	OF [33]	ED+S	U-net+S	ED+S+M	U-net+S+M
L1	0.146	0.135	0.119	0.113	<b>0.104</b>
SSIM	0.649	0.698	0.710	0.703	<b>0.727</b>

similarity (SSIM [34]) between the generated images and the ground truth future images of the test data.

Table I depicts the result of our quantitative analysis on DNet and the comparing baselines. In this table “+S” indicates sampling, and “+M” indicates multi-image merge (last step of DNet). The U-net architecture used in DNet outperforms the methods using a regular encoder-decoder(ED). Also, our multi-image merge (M) approach leads to better results (lower L1 and higher SSIM) than the single image approaches. Moreover, our method also outperforms the optical flow (OF) based baseline.

**Dynamic-Scene View Synthesis:** We evaluate the performance of our complete view synthesis method, VUNet for dynamic environments in test data with i) static environments in GS2, ii) dynamic environment with fixed robot position from GS3, and iii) dynamic environment with moving robot from GS3. This last group of sequences is especially challenging because this type of data has not been seen during SNet and DNet training. We compare multiple structures for SNet and DNet using regular encoder-decoder versus U-net architecture, training with and without GAN, and optical flow. We report mean L1 norm (pixel difference) and structural similarity (SSIM) between the generated images and the ground truth future images in the test data.

Table II shows the quantitative results of the baseline methods and the ablation study on the proposed method to generate virtual images (SNet+DNet). For each model and type of data the first value indicates the mean pixel L1 difference (smaller is better) and the second value indicates the structural similarity, SSIM (larger is better). The six rows from (b) to (g) depict the results of synthesizing images using **only** different variations of the SNet architecture. Without sampling, U-net (row (d)) outperforms the regular encoder-decoder (ED, row (b)). However, with sampling the encoder-decoder architecture (row (f)) improves the performance, especially in our application scenario, dynamic environments with robot motion (third column). We also observed that using a GAN training does not improve the results (rows (c) and (e)).

The last six rows (from (h) to (m)) depict the result of synthesizing images combining different variants of SNet and DNet, including our proposed method (row (m)). Our proposed method, VUNet is one of the best three methods in all scenarios. Specifically, in our application domain, dynamic scenes with robot motion, our method outperforms all baselines. We observe that the usage of DNet does not improve the results on static environments (first column). In these scenarios it is simpler to use only SNet. However, as expected, SNet alone fails in scenarios with dynamic objects.

**Future Traversability Estimation:** We evaluate the accuracy of the estimation of traversability based on the images generated by our view synthesis and the previously proposed baselines. We randomly sample images from GS2 and GS3 and hand-label them until we collect 200 traversable and 200

untraversable images. The untraversable images are images just before the robot collides or falls. We take the two previous images to each selected image and use them together with the ground truth commanded velocity to predict the selected image. We feed then the generated image to GONet [1] to calculate the future traversable probability. If the probability is over a threshold  $p_{min} = 0.5$  we label the image as traversable, otherwise we label it as untraversable. We estimate the accuracy of the traversability estimation by comparing the predicted and the manually assigned (ground truth) traversability labels.

The left side of the last column of Table II depicts the results of this quantitative evaluation on the future traversability estimation in GS2 and GS3. Feeding the images generated by VUNet for view synthesis (row (m)) to GONet yields the highest accuracy for the traversability estimation. The higher accuracy is the result of a higher quality in the predicted images. We note that the accuracy in the traversability estimation of some variants without DNet is also high (e.g. ED+S, row (f), and U-net+S, row (g)). This is an artifact caused by the distribution of our evaluation data: even though we sample 100 evaluation images from GS3 depicting dynamic objects, they are usually far and failing to predict their changes in the image do not usually affect the traversability estimation.

In order to evaluate more clearly the benefits of using our DNet component, we collected an additional dataset of about one hour (*Ped. DS* in Table II) where the dynamic obstacles (pedestrians) and the robot often cross their trajectories. We compare our future traversability estimation method to the baselines and list their accuracy in the right side of the last column in Table II. Our method achieves the highest accuracy and shows a clear quantitative advantage against baseline methods without DNet, i.e. not accounting for the motion of the dynamic obstacles.

We also compared VUNet to a baseline using depth images from a Kinect sensor included also in GS2, GS3 and the pedestrian dataset. We turn the Kinect sensor pointing forward into a proximity sensor and develop a baseline that indicates that an area is untraversable if there are obstacles closer than a distance threshold. We determined the optimal threshold as the threshold that leads to the maximum accuracy in the validation set. The accuracy of the Kinect-based baseline is listed in the first row of Table II. The baseline using depth images performs worse than our RGB-based method because the Kinect images contain noise due to reflections in mirrors, glass walls, and missing points in dark objects. Also, the traversability estimation using Kinect do not consider the motion of the pedestrian in the future step, which leads to the poor performance in the pedestrian dataset.

Additionally, we evaluated our proposed approach, VUNet for future traversability to predict two, three and four steps obtaining 91.3%, 89.3% and 88.0%, respectively. The decreased accuracy is caused by the more difficult predictions in longer horizons. Considering safety applications we evaluated the accuracy using a more conservative traversability threshold of  $p_{min} = 0.7$  decreasing the amount of non-predicted untraversable future steps (the riskiest case) to less than 6% in one to four future steps.

TABLE II: Evaluation of View Synthesis and Traversability Estimation

Models:			GS2	GS3	GS3	Trav. Accuracy	
SNet Variants	+	DNet Variants	Static Environment	Dynamic env. w/ robot motion	Dynamic env. w/ robot motion	GS2 & 3	Ped. DS
(a) Kinect			-	-	-	0.818	0.735
(b) ED[20]	+	-	0.117 / 0.556	0.225 / 0.395	0.151 / 0.501	0.690	0.620
(c) ED+GAN[35]	+	-	0.147 / 0.468	0.253 / 0.333	0.188 / 0.400	0.660	0.540
(d) U-net[30]	+	-	<b>0.064</b> / <b>0.779</b>	0.148 / <b>0.698</b>	0.115 / 0.644	0.920	0.735
(e) U-net+GAN	+	-	0.069 / 0.752	0.148 / <b>0.698</b>	0.124 / 0.602	0.920	0.675
(f) ED+S[12]	+	-	<b>0.065</b> / <b>0.777</b>	0.155 / 0.672	0.116 / 0.647	0.947	0.777
(g) U-net+S	+	-	0.067 / 0.765	0.158 / 0.663	0.117 / 0.642	0.945	0.770
(h) U-net+S	+	OF	0.086 / 0.706	0.155 / 0.607	0.143 / 0.548	0.905	0.822
(i) U-net+S	+	ED+S+M	0.068 / 0.761	0.123 / 0.668	0.108 / 0.647	0.945	0.797
(j) U-net+S	+	U-net+S+M	0.068 / 0.765	0.116 / 0.686	0.105 / 0.653	0.937	0.810
(k) ED+S	+	OF	0.092 / 0.680	0.158 / 0.594	0.594 / 0.529	0.905	0.817
(l) ED+S	+	ED+S+M	0.068 / 0.766	0.123 / 0.686	0.110 / 0.644	0.937	0.800
(m) ED+S	+	U-net+S+M (VUNet)	<b>0.065</b> / <b>0.776</b>	<b>0.113</b> / <b>0.698</b>	<b>0.104</b> / <b>0.657</b>	<b>0.950</b>	<b>0.830</b>

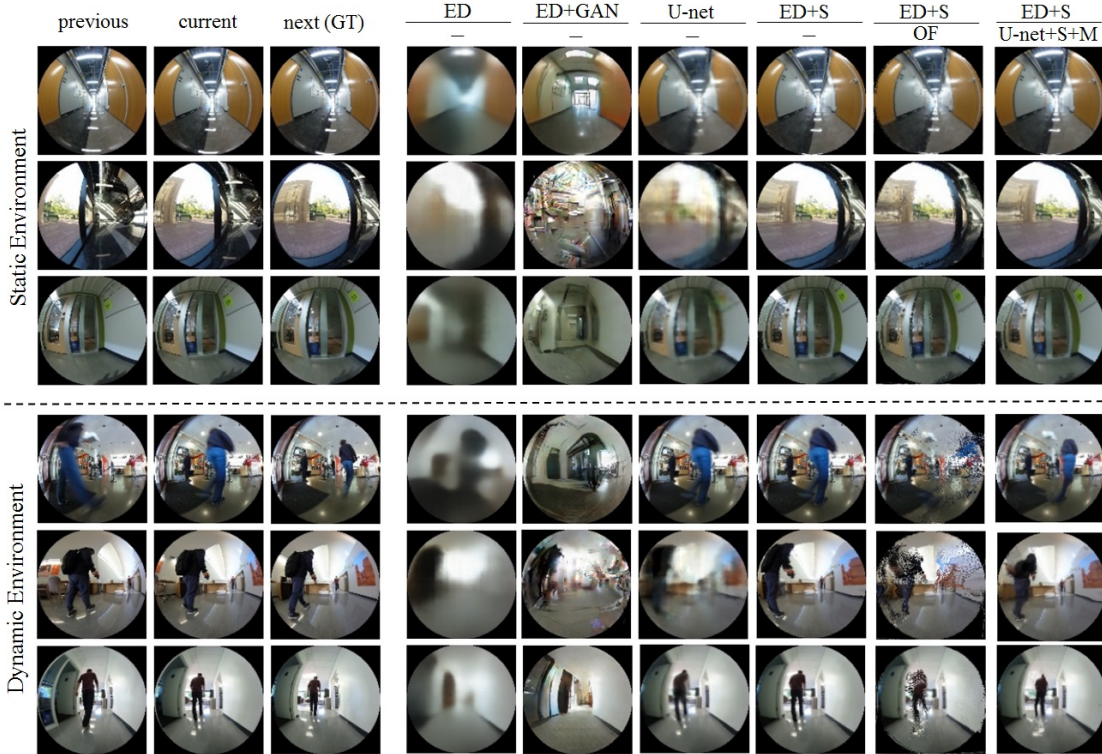


Fig. 6: Predicted images in static (first three rows) and dynamic (last three rows) environments. From left to right: previous, current and future (ground truth for view synthesis, GT) images as viewed by the robot, predicted images from baselines (SNet and DNet variants), and predicted images from VUNet (most right)

### B. Qualitative Analysis

In the second set of experiments we evaluate qualitatively the results of our dynamic-scene view synthesis, VUNet and traversability estimation approaches.

**Dynamic-Scene View Synthesis:** First, we compare the generated images from our method and the baselines methods side by side (see Fig. 6). The first three columns show the previous, current, and future (ground truth for the synthesis, GT) RGB images as viewed by the robot. We observe that the GAN training improves the sharpness of the blurred predicted image from the encoder-decoder (ED and ED+GAN, 4th and 5th columns). However, while being sharper, some of the generated images by ED+GAN do not resemble much the real image (e.g. 2nd and 5th row of 5th column). U-net (6th column) can generate very clear images when the current image is similar to the future image, but it does not perform

as well when it has to predict dynamic obstacles. Similarly, all baseline methods without DNet are not able to predict the appearance changes due to moving objects (e.g. humans). We observe that the location of humans in the predicted images without DNet is same as the future predicted image (three last rows). We can also see that, ED+S+OF can not predict accurately the human movement: there are speckle patterns in the predicted images. This is because the errors in SNet cause wrong extrapolations with OF.

In contrast, our method, VUNet (SNet+DNet) is able to predict the image changes due to both robot pose changes and motion of the dynamic objects (i.e. humans). For example, in the scene shown in the 5th row, the robot is turning to the right side while a human is crossing by. Surprisingly, even the unseen part of the picture on the right side wall in the future image can be correctly constructed in the predicted image.

Additionally, the human is correctly moved towards the right in the predicted image (a failure in ED+S).

In the last scene (last row), both the robot and the human are moving forward in the corridor. While several methods can correctly predict that the static parts of the environment (e.g. the door) will appear closer to the robot, only our method predicts that the human, which is faster than the robot (as can be observed from previous and current images), should be farther away in the predicted image.

**Future Traversability Estimation:** To evaluate qualitatively our method for future traversability estimation we propose two applications based on it for assisted teleoperation: early obstacle detection or multi-path future traversability estimation. These methods are implementations of the system depicted in Fig. 5 with different ways of generating virtual future robot velocity commands.

*Early obstacle detection:* In this application, the teleoperator uses a joystick to control the robot and gets audio input (warning) or emergency stops from the proposed system. To predict the images and the traversability in the future our method assumes that the upcoming robot commands will be the robot’s maximum linear velocity and a constant angular velocity  $\xi = (v_{max}, v_{max}/r_c)$ . With this assumption our method assumes the riskiest possible future. The robot’s maximum linear velocity is  $0.5 \text{ m s}^{-1}$ , and  $r_c$  is the turning radius last used by the teleoperator calculated as  $r_c = v_c/\omega_c$ , where  $v_c$  and  $\omega_c$  are the teleoperator’s last commands. A safety alarm is fired when the traversable probability for the third ( $p_{s_{t+3}}^{t+3}$ ) or fourth ( $p_{s_{t+4}}^{t+4}$ ) time steps in the future are less than 0.5. Additionally, an emergency stop interrupts the current teleoperation of the robot if the traversable probability of the current state ( $p_{s_t}^t$ ), next ( $p_{s_{t+1}}^{t+1}$ ), or second next steps ( $p_{s_{t+2}}^{t+2}$ ) are less than 0.5.

Fig. 7 shows three examples of our application for early obstacle detection of cases in dynamic environment with moving obstacles (pedestrians). The traversable probability of each image is shown under each image by applying GONet to each image generated with VUNet. We compare the results of future traversability estimation based on images using only SNet or our proposed method, VUNet (SNet+DNet). In the predicted images without DNet, the changes due to the motion of the human cannot be predicted. Therefore, the model without DNet wrongly estimates the traversability assuming a non-moving pedestrian. Our proposed approach using DNet predicts the motion of the human in the image, which leads to a more accurate prediction of the future traversability. Additional qualitative results can be seen in our supplementary video. Our proposed application for early obstacle detection correctly estimates the traversable probability in the future and indicates this to the teleoperator with warning signals and emergency stop commands.

*Multi-path traversability estimation:* For the future traversability estimation for multiple paths, we propose to apply our system of Section III-B in the way depicted in Fig.5, top. The system generates virtual velocities for five different paths around the robot, predicts the images using our scene view method and calculate the traversability for each of the paths. To generate the virtual robot velocities

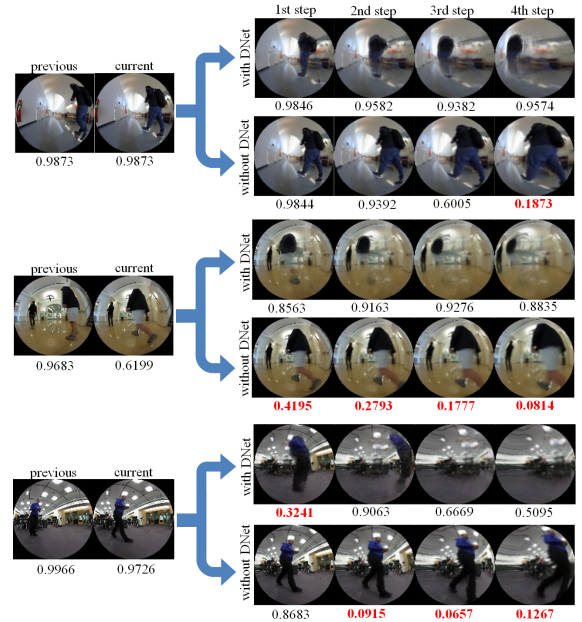


Fig. 7: Application of our future traversability estimation for early obstacle detection in three dynamic environment examples. The inputs to our application system are the previous and current images, and the last teleoperation command. The system predicts the images at four consecutive time steps and estimates the traversable probability, depicted under each image. Red probabilities indicate values under 0.5. We show the different predictions using DNet(VUNet) or without DNet (only SNet.)

$\xi^{t+i}, i \in 1 \dots 4$  our system assumes a constant maximum linear velocity,  $v_{max} = 0.5 \text{ m s}^{-1}$  and five different angular velocities multiple of  $\omega_0 = 0.5 \text{ rad s}^{-1}$ :  $\omega_{LL} = 2\omega_0$ ,  $\omega_L = \omega_0$ ,  $\omega_C = 0$ ,  $\omega_R = -\omega_0$ , and  $\omega_{RR} = -2\omega_0$ . We ask a teleoperator to navigate the robot in different scenarios and collect the multiple path traversability predictions.

Figure 8 shows an example of our multi-path traversability estimation. The figure shows previous and current images (left side) as well as the predicted images on each path on the test set (bottom). The traversable probabilities are shown under each image. In this example, the robot is moving in a narrow corridor with tables and chairs on both sides. Our method can correctly predict the safe path in front of the robot based on the synthesized future images. Additional qualitative examples of this application are included in our supplementary video.

## VI. LIMITATIONS AND FUTURE WORK

The quantitative and qualitative analysis pointed out some limitations but we don’t deem them severe for our applications. The quality of the synthesized images decreases for longer time horizon predictions. This affects the accuracy of the future traversability estimation. The degradation of future image predictions is caused by two main factors: 1) the part of the environment that is now visible was not visible in the images we used to synthesize the view (e.g. due to large occlusions or abrupt rotations), and 2) the dynamic objects present a complex motion pattern (e.g. different parts of the human body like legs and arms). However, even in these scenarios, the quality of the generated images is good enough to predict with high accuracy the future traversable probability for the spaces around the mobile robot. To alleviate further these effects we will explore in future work methods to reflect

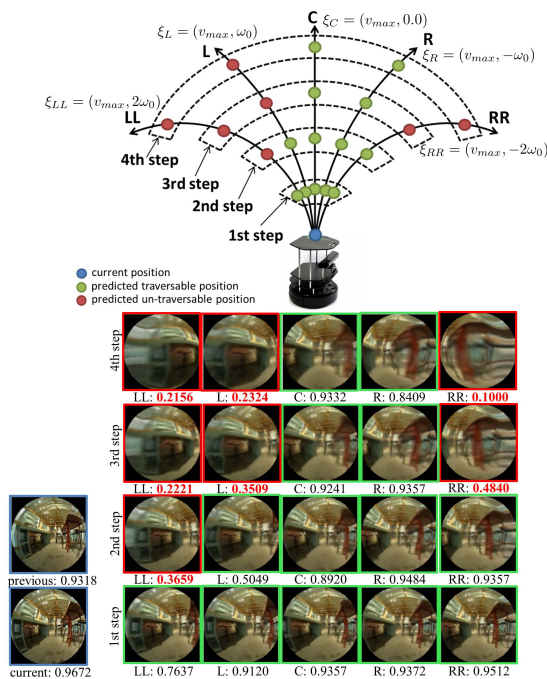


Fig. 8: Application of our multi-path traversability estimation system for assisted teleoperation. Top: Spatial diagram of the five paths where we predict traversability around the robot. We estimate traversability in a most left (LL), left (L), central (C), right (R), and most right (RR) paths. Bottom: the input previous and current images (left) and the generated images (right) with associated traversable probability. The robot image is cited from [11].

the uncertainty on the predictions, both due to odometry errors and due to non-deterministic dynamic obstacle motion[36]. We also plan to combine our approach with a vision-based target navigation into a full autonomous navigation system that avoids obstacles and reaches a target destination.

## VII. CONCLUSION

In this paper we propose a novel dynamic-scene view synthesis method for robot visual navigation using an RGB camera, VUNet. We show the proposed method is applicable for future traversability estimation. Our view synthesis method predicts accurate future images given virtual robot velocity commands. Our method is able to predict the changes caused both from the moving camera viewpoint and the dynamically moving objects. Our synthesized images outperform both quantitatively and qualitatively the images generated by state of the art baseline methods. We use the synthesized images to predict traversability in future steps for multiple paths and show its application to assisted teleoperation scenarios.

## REFERENCES

- [1] N. Hirose, A. Sadeghian, M. Vázquez, P. Goebel, and S. Savarese, "Gonet: A semi-supervised deep learning approach for traversability estimation," in *IROS*, 2018, pp. 3044–3051.
- [2] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," *CVPR*, 2017.
- [3] F. Xia, A. R. Zamir, Z. He, A. Sax, J. Malik, and S. Savarese, "Gibson env: Real-world perception for embodied agents," in *CVPR*, 2018.
- [4] M. Babaeizadeh, C. Finn, D. Erhan, R. H. Campbell, and S. Levine, "Stochastic variational video prediction," *arXiv preprint arXiv:1710.11252*, 2017.

- [5] J. Flynn, I. Neulander, J. Philbin, and N. Snavely, "Deepstereo: Learning to predict new views from the world's imagery," in *CVPR*, 2016, pp. 5515–5524.
- [6] D. Fox, W. Burgard, and S. Thrun, "The dynamic window approach to collision avoidance," *RAM*, vol. 4, no. 1, pp. 23–33, 1997.
- [7] F. Flacco, T. Kröger, A. De Luca, and O. Khatib, "A depth space approach to human-robot collision avoidance," in *ICRA*. IEEE, 2012, pp. 338–345.
- [8] M. Pfeiffer, M. Schaeuble, J. Nieto, R. Siegwart, and C. Cadena, "From perception to decision: A data-driven approach to end-to-end motion planning for autonomous ground robots," in *ICRA*. IEEE, 2017, pp. 1527–1533.
- [9] M. Cinietal, "Real-time obstacle avoidance by an autonomous mobile robot using an active vision sensor and a vertically emitted laser slit," in *Intelligent Autonomous Systems*, vol. 7, no. 1, 2002, pp. 301–309.
- [10] C. Richter and N. Roy, "Safe visual navigation via deep learning and novelty detection," in *RSS*, 2017.
- [11] <https://www.robot-advance.com/EN/cat-turtlebot-150.htm>.
- [12] T. Zhou *et al.*, "View synthesis by appearance flow," in *ECCV*. Springer, 2016, pp. 286–301.
- [13] D. Ji, J. Kwon, M. McFarland, and S. Savarese, "Deep view morphing," Technical report, Tech. Rep., 2017.
- [14] E. Park *et al.*, "Transformation-grounded image generation network for novel 3d view synthesis," in *CVPR*. IEEE, 2017, pp. 702–711.
- [15] S. M. Seitz and C. R. Dyer, "View morphing," in *SIGGRAPH*. ACM, 1996, pp. 21–30.
- [16] P. Hedman, T. Ritschel, G. Drettakis, and G. Brostow, "Scalable inside-out image-based rendering," *TOG*, vol. 35, no. 6, p. 231, 2016.
- [17] J. Shade *et al.*, "Layered depth images," in *SIGGRAPH*. ACM, 1998, pp. 231–242.
- [18] A. Dosovitskiy, J. T. Springenberg, and T. Brox, "Learning to generate chairs with convolutional neural networks," in *CVPR*, 2015.
- [19] T. D. Kulkarni, W. F. Whitney, P. Kohli, and J. Tenenbaum, "Deep convolutional inverse graphics network," in *NIPS*, 2015.
- [20] M. Tatarchenko, A. Dosovitskiy, and T. Brox, "Multi-view 3d models from single images with a convolutional network," in *ECCV*. Springer, 2016, pp. 322–337.
- [21] C. Vondrick, H. Pirsiavash, and A. Torralba, "Generating videos with scene dynamics," in *NIPS*, 2016, pp. 613–621.
- [22] W. Liu, D. L. W. Luo, and S. Gao, "Future frame prediction for anomaly detection – a new baseline," in *CVPR*, 2018.
- [23] I. Bogoslavski *et al.*, "Efficient traversability analysis for mobile robots using the kinect sensor," in *ECMR*. IEEE, 2013, pp. 158–163.
- [24] A. Cherubini, B. Grechanichenko, F. Spindler, and F. Chaumette, "Avoiding moving obstacles during visual navigation," in *ICRA*. IEEE, 2013, pp. 3069–3074.
- [25] B. Suger, B. Steder, and W. Burgard, "Traversability analysis for mobile robots in outdoor environments: A semi-supervised learning approach based on 3d-lidar data," in *ICRA*, 2015, pp. 3941–3946.
- [26] D. Kim, S. M. Oh, and J. M. Rehg, "Traversability classification for ugv navigation: A comparison of patch and superpixel representations," in *IROS*. IEEE, 2007, pp. 3166–3173.
- [27] I. Ulrich and I. Nourbakhsh, "Appearance-based obstacle detection with monocular color vision," in *AAAI/IAAI*, 2000, pp. 866–871.
- [28] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," in *NIPS*, 2014.
- [29] D. Eigen and R. Fergus, "Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture," in *ICCV*, 2015, pp. 2650–2658.
- [30] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *MICCAI*. Springer, 2015, pp. 234–241.
- [31] S. Tokui *et al.*, "Chainer: a next-generation open source framework for deep learning," in *NIPS: Workshop LearningSys*, 2015.
- [32] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [33] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. Van Der Smagt, D. Cremers, and T. Brox, "Flownet: Learning optical flow with convolutional networks," in *CVPR*, 2015, pp. 2758–2766.
- [34] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [35] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther, "Autoencoding beyond pixels using a learned similarity metric," *arXiv preprint arXiv:1512.09300*, 2015.
- [36] A. Sadeghian, V. Kosaraju, A. Sadeghian, N. Hirose, and S. Savarese, "Sophie: An attentive gan for predicting paths compliant to social and physical constraints," *arXiv preprint arXiv:1806.01482*, 2018.