

Planning Paths through Occlusions in Urban Environments

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** This paper presents a novel framework for planning in unknown and
2 occluded urban spaces. We specifically focus on turns and intersections where oc-
3 clusions significantly impact navigability. Our approach uses an inpainting model
4 to fill in a sparse, occluded, semantic lidar point cloud and plans dynamically fea-
5 sible paths for a vehicle to traverse through the open and inpainted spaces. We
6 demonstrate our approach using a car’s lidar data with real-time occlusions, and
7 show that by inpainting occluded areas, we can plan longer paths, with more turn
8 options compared to without inpainting; in addition, our approach more closely
9 follows paths derived from a planner with no occlusions (called the ground truth)
10 compared to other state of the art approaches.

11 **Keywords:** Navigation, Occluded Environments, Semantic Scene Understanding

12 1 Introduction

13 Planning in environments with unknown spaces is a challenging topic in robotics, as they can limit
14 the speed of travel and decision making in real-time. Unknown spaces occur from occluding objects,
15 such as cars in the road, buildings, trees or fences, or from limitations in sensor range and resolution.
16 The number and extent of these unknown spaces increase dramatically in cluttered environments,
17 such as in an airport or in an urban city. Traditionally, path planners in these types of environments
18 typically either plan only in the known spaces (which limits speed and navigability) or assumes
19 unknown space is free (which increases the chances of varied maneuvering and collisions).

20 Consider an example of a car turning at an upcoming intersection, but there are pedestrians on the
21 corner and a truck is in the intersection, such that space beyond the truck/pedestrians is unknown. If
22 a path planner only plans in known spaces, the planner will not have the range to plan a path through
23 the intersection. If the planner assumes unknown space is free, it could potentially run into dead
24 ends. Similar examples exist regularly in urban environments. Comparatively, a human driving in
25 such an environment makes predictions about what lies beyond the occluding objects in order to
26 navigate more smoothly. If unable to predict what is behind occluding objects drivers will slow
27 down significantly before advancing further.

28 Our approach, shown in Figure 1, addresses this problem by filling in unknown spaces using image
29 inpainting techniques. We define unknown spaces as regions of previously unmapped environments
30 occluded from sensor measurements by obstacles such as buildings. Our work predicts the underly-
31 ing static structure in these unknown spaces, for example, a turn or intersection in the road for more
32 informed path planning. First, a sensor measurement is projected to the bird’s eye view (BEV) in
33 order to reveal occluded spaces. Next, a data-driven image inpainting model fills in the occluded
34 spaces. Finally, the inpainted map is used for planning a dynamically feasible path. We specifically
35 focus on turns and intersections where occlusions significantly impact navigation.

36 The main contributions of this paper are:

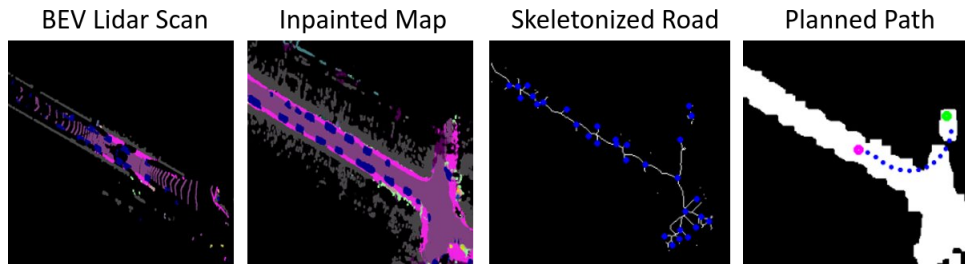


Figure 1: (L→R) (1) Initial lidar scan rendered from a BEV perspective annotated with semantic labels (purple = road). (2) Output of an inpainting model used on the initial lidar scan; unknown pixels are filled in with semantic labels. (3) The skeletonized road; blue dots represent nodes where multiple edges connect. (4) The planned path through white pixels denoting traversable road; the pink dot is the initial pose, the green dot is the goal, and blue dots are the path sequence.

- 37 • a predictive inpainting model for navigation in urban outdoor environments which fills in
- 38 unknown spaces, extending the perception range allowing for smoother and faster naviga-
- 39 tion, especially for turns and intersections in cluttered environments. The predictive model
- 40 is easily incorporated into existing path planning methods.
- 41 • a modified loss to the Pix2Pix [1] network tailored to the occlusion inpainting task.
- 42 • a dataset for training models of occluded turns and intersections for filling in unknown
- 43 spaces in urban environments.

44 The proposed framework is compared against a planner which does not use any inpainting to fill in
 45 occluded spaces. Experimental results demonstrate the novel framework is highly effective at (1)
 46 identifying turns and intersections through occluding objects, (2) generating a traversability graph
 47 which closely matches the ground truth (i.e. path assuming no occlusions), and (3) extending the
 48 range of dynamically feasible trajectory planning in occluded environments.

49 2 Related Works

50 Path planning in known environments is a well studied problem in robotics. Grid-based map repre-
 51 sentations are generally used and many optimal algorithms exist (A* [2], D* [3]). Unknown spaces
 52 (due to occlusions) are generally modeled as free or ignored [4]. Approaches that model the un-
 53 known space as free must frequently replan online due to discrepancies between what is actually in
 54 the unknown space and the optimistic assumption that unknown space is free [5].

55 Ref. [6] shows the range of a path planner in outdoor urban environments can be extended by using
 56 semantic segmentations. This approach demonstrates that by extending the trajectory length, the
 57 robot navigates faster and traverses smoother paths compared to a traditional metric grid planner.
 58 This approach however does not address planning in occluded spaces.

59 **Data-driven predictive modeling** of unknown spaces has been studied in recent work. The previous
 60 work most similar for outdoor environments [5] uses image inpainting to fill in unknown spaces from
 61 sensor measurements projected to a BEV. This work requires explicit labeling of the map pixels
 62 which must be inpainted. In our approach, we do not explicitly label which pixels must be inpainted
 63 as the pixel inpainting is completely data-driven. Our approach also considers the dynamics of the
 64 vehicle, while [5] does not. Refs. [7, 8] use predictive modeling of unknown spaces for planning,
 65 but their scope is restricted to indoor environments.

66 [9, 10, 11] use image inpainting to fill in foreground objects with background classes. However,
 67 these approaches do not directly integrate with path planning.

68 **Generative Adversarial Networks** (GANs) [12] are a powerful tool used for a variety of tasks, such
 69 as natural language processing [13], super resolution [14], and image translation [15]. Generally,
 70 GANs aim to model the underlying distribution lying in the target data domain. Traditional GANs
 71 for image inpainting tasks can be divided into two categories. The first are GANs targeted for paired

72 image-to-image translation such as Pix2Pix [1] and Pix2Pix HD [16]. The other type of GANs,
 73 considered to be state-of-the-art image inpainting models [17, 18], use free-form or rectangular
 74 masks as inpainting signals.

75 3 Technical Approach

76 Figure 1 shows the pipeline for our proposed approach. First, a semantic lidar point cloud derived
 77 from a lidar scan is created. The point cloud is transformed to a BEV and an inpainting model fills
 78 in the unknown, occluded pixels. Next, the traversable road pixels are skeletonized into a graph with
 79 waypoints. Finally, a dynamically feasible path is planned from the vehicle pose to a given goal.

80 3.1 Dataset Generation

81 Our dataset is generated based on the KITTI-360 dataset [19]. In the KITTI-360 dataset, a vehicle
 82 is driven around Karlsruhe, Germany, and sensor data is collected and annotated for 73.7 km. We
 83 specifically focus on semantic lidar data in this work. Lidar provides more reliable depth than stereo
 84 cameras, but still suffers from issues with occluding obstacles blocking out unknown space. KITTI-
 85 360 provides semantic annotations (19 classes) for the *aggregated* lidar point clouds for each route.

86 For our work, we transform point clouds to a bird’s eye view (BEV) because BEV allows for clear
 87 observation of occluded spaces. We project semantic annotations from the fully aggregated lidar
 88 point clouds to each individual lidar scan. Then, we render both the aggregated point clouds and the
 89 individual semantic lidar scans from a BEV to complete our dataset. We remove points with class
 90 label vegetation since vegetation can obscure the underlying road and sidewalks, and also unknown
 91 spaces underneath in BEV images. Since the occlusions are already present in the point cloud,
 92 removing the vegetation points does not impact our inpainting results. We have rendered 22698
 93 frames for three driving sequences.

94 3.2 Image Inpainting

95 In the lidar measurements there are unknown spaces such as shown in Figure 3 (top left). The
 96 orange oval shows semantic annotations for a region that is occluded in Figure 3 (top second). These
 97 occluded spaces are a function of obstacles such as buildings and fences, and also sensor resolution.
 98 To predict the semantics in the unknown spaces, we employ a paired generation model to recover
 99 the lost semantic information.

100 3.2.1 Baseline Paired Generation Algorithm

101 For translating an image $x^{(n)}$ in source domain \mathcal{X} (\mathcal{X} is defined as the 2D lidar space) to an image
 102 $y^{(n)}$ in target domain \mathcal{Y} (\mathcal{Y} is defined as the 2D ground truth semantic space), we first introduce a
 103 baseline paired generation algorithm, inspired by [1]. Isola et al. [1] assumes the training data are
 104 perfectly paired images, and an L1 loss in the pixel space compares $x^{(n)}$ and $y^{(n)}$:

$$\mathcal{L}_1(G, X, Y) = \frac{1}{N} \sum_n \frac{1}{K^{(n)}} \sum_{i,j} \|G(x^{(n)})_{i,j} - y_{i,j}^{(n)}\|_1. \quad (1)$$

105 where $G(x^{(n)})_{i,j}$ and $y_{i,j}^{(n)}$ are pixel values at the (i, j) location and $K^{(n)}$ is the number of pixels
 106 in image $y^{(n)}$. N refers to the number of images in the training set. Additionally, [1] further
 107 employs a standard GAN loss to restore a realistic image. We ask the reader to refer to [12] for the
 108 implementation of a standard GAN.

109 3.2.2 Implementation

110 To ensure the paired generation model restores most of the semantics lost in lidar space, we propose
 111 a modification of the Pix2Pix network. First of all, we adopt a coarse-to-fine generator G and a
 112 multi-scale discriminator D from Pix2PixHD [16] which is the state-of-the-art paired generation
 113 model.

114 3.2.3 The inpainting-targeted L1 loss

115 Although the generative model restores some unknown semantics in lidar space, the original seman-
 116 tics sometimes can get lost during the feature extraction process (see Figure 2) due to the sparsity
 117 of extractable features. To ensure the input semantics are preserved during the feature extraction
 118 process, we propose an inpainting-targeted L1 loss. Specifically, for an input image $x^n \in \mathcal{X}$ and the
 119 output image $G(x^n \in \mathcal{Y})$, the loss is defined as:

$$\mathcal{L}_1^*(G, X) = \frac{1}{N} \sum_n \frac{1}{K^{x^n}} \sum_{(i,j) \in x^{*,n}} \|G(x^n)_{i,j} - x^n_{i,j}\|_1. \quad (2)$$

120 where K^{x^n} refers to the number of non-zero pixels in the input lidar image x^n and $(i, j) \in x^{*,n}$
 121 refers to the pixel location (i, j) in the nonzero set $x^{*,n}$ of the original input lidar image. Thus, from
 122 the loss function definition, the inpainting-targeted L1 loss aims to ensure the nonzero semantics
 123 lying in the original input image x^n still exist in the generated image $G(x^n)$. Figure 2 shows an
 124 ablation study of $\mathcal{L}_1^*(G, X)$.

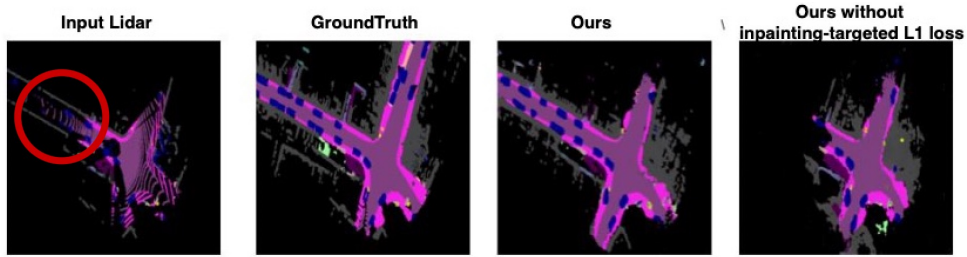


Figure 2: Comparison of the inpainting performance on our model with and without inpainting-targeted L1 loss $\mathcal{L}_1^*(G, X)$ (Sec. 3.2.3). L→R: (1) input lidar, (2) GT semantics, inference results from our model (3) with and (4) without $\mathcal{L}_1^*(G, X)$. Sparse semantics in the faraway region of the road in the original lidar map (red circle) are predicted from our model with $\mathcal{L}_1^*(G, X)$, but not without.

125 3.2.4 Ranking-Based Loss to Overcome Large Stochastic Variations

126 The inpainting-targeted L1 loss effectively compensates for the semantics lost during the feature
 127 extraction stage. However, restoration of small details, especially semantics lying in the intersection
 128 areas are not solved effectively by the loss terms described above.

129 To address this, we employ a loss function focused on local features (*i.e.*, image patches). Patches of
 130 the generated and target images at the same position should capture similar structures and content,
 131 such that their similarity is larger than patches at different positions. We realize this idea using a
 132 loss similar to the patchNCE loss [20]¹. Let H be a feature extractor and H_s be the feature at the
 133 spatial location s on the feature map. Our usage of the patchNCE loss is defined as follows

$$\sum_s \ell(H_s(G(x)), H_s(y), \{H_{s'}(y) | s' \neq s\}), \quad (3)$$

$$\text{where } \ell(v, v^+, \{v_{s'}^-\}) = -\log \left[\frac{\exp(v \cdot v^+ / \tau)}{\exp(v \cdot v^+ / \tau) + \sum_{s'} \exp(v \cdot v_{s'}^- / \tau)} \right]. \quad (4)$$

134 Here, τ is a temperature constant scalar. Through minimizing this loss with respect to G , we en-
 135 courage $\mathbf{H}_s(G(\mathbf{x}))$ and $\mathbf{H}_s(\mathbf{y})$ to have higher inner product similarity than $\mathbf{H}_s(G(\mathbf{x}))$ and other
 136 patches of \mathbf{y} . By applying this loss to all the training pairs, we have

$$\mathcal{L}_{\text{patchNCE}}(G, H, X, Y) = \frac{1}{N} \sum_s \ell(H_s(G(x^{(n)})), H_s(y^{(n)}), \{H_{s'}(y^{(n)}) | s' \neq s\}). \quad (5)$$

137 In our implementation, we follow [20] to extract features from multiple layers.

¹[20], being an unpaired translation algorithm, computes the patchNCE loss between the source and gener-
 ated images. We compute the patchNCE loss between the generated and coarsely-aligned target images.

Table 1: Dataset splits (# of images) for image inpainting.

Route	Train Set	Val Set	Test Set
Route 0	8386	1048	1049
Route 2	8982	1123	1122
Route 3	790	99	99

138 **3.2.5 Final Training Objective**

139 Combining the loss terms introduced in Sections 3.2.3 and 3.2.4, our training objective is

$$\mathcal{L}_{GAN}(G, D, X, Y) + \mathcal{L}_{\text{patchNCE}}(G, H, X, Y) + \mathcal{L}_1^*(G, X). \quad (6)$$

140 Specifically, we apply mini-batch stochastic gradient descent and learn G and H to minimize the
 141 objective while learning D to maximize the objective.

142 **3.3 Path Planning**

143 Previous work [5] uses an inpainting model to predict occluded spaces in order to plan more in-
 144 formed paths. However, [5] does not account for vehicle dynamics and uses an A* planner to plan
 145 with the pixels of the BEV image as possible states. In our paper, we use a hybrid A* planner [21],
 146 which plans optimal and dynamically feasible paths.

147 Figure 3 shows our planning pipeline. First we take as input the BEV map of either the original
 148 lidar scan (OL), inpainted image map (IM), or GT map (Figure 3 top row). Next a mask is fit to
 149 the road pixels, since the vehicle only navigates on the road. We dilate and then erode the road
 150 mask in order to connect the individual road points from the lidar scan. This allows us to visualize a
 151 fully connected road. Zhang’s method [22] is used to skeletonize the road. The intersections of the
 152 skeleton define waypoints w for the planner to navigate (Figure 3 middle row).

153 Given the final goal location, which is a point all the way around a turn at an intersection (Figure 3
 154 top row, green dot), the closest waypoint w to the goal is used as a local goal to plan to. A hybrid
 155 A* planner plans a dynamically feasible path from the vehicle pose to the goal. The planner states
 156 are $[px, py, \theta]$, where px and py are pixelwise coordinates and θ is orientation in the world frame.

157 We specifically focus on predictive modeling of the underlying static structures in the unknown
 158 spaces, such as the road. Our inpainting model can be easily combined with an existing full planning
 159 stack by using it to generate a predictive model given a sensor measurement. We present a high-
 160 level planner here, with evaluation in the experimental section, and assume an off-the-shelf low-level
 161 obstacle/collision avoidance module is used in addition to our planner.

162 **4 Experimental Evaluation**

163 **4.1 Network Training**

164 To train the network, described in section 3.2, we generate 2D lidar maps and the corresponding
 165 semantic maps using three traversals (sequences 0, 2, and 3) in the KITTI-360 dataset [19]. Table 1
 166 shows the number of images in the training and test sets. The train and test sets do not contain
 167 overlapping locations. We train our model for 200 epochs. During training, the initial learning rate
 168 $lr = 0.0002$ for the first 100 epochs. For the remaining 100 epochs, the lr linearly decreases to
 169 zero. We use the Adam optimizer [23] for both the generator and the discriminator ($\beta_1 = 0.5$ and
 170 $\beta_2 = 0.999$). The model is trained on one NVIDIA RTX3090 GPU. For the baselines in section 4.5,
 171 we use the same dataset split and network parameters (lr , Adam optimizer, and training epochs).

172 **4.2 Dataset Description**

173 We evaluate our framework by testing performance on the turns and intersections of the test set. We
 174 identify regions where the vehicle makes a turn and evaluate our method on those regions. For each

175 turn, the dataset frames corresponding to it are selected from just before the turn is seen in the ground
 176 truth (GT) to when the vehicle reaches a point where the turn is no longer feasible. This allows us to
 177 fully evaluate the effect of the inpainting model on planning for the turns and intersections, which
 178 are the most difficult regions to plan through. For example, for the first turn, fifty frames are used
 179 for evaluation. In each frame a path is planned from the current vehicle location to the selected goal
 180 point from the skeletonized nodes. The paths consist of a set of nodes as described in section 3.3.

181 4.3 Evaluation metrics

182 **Frechet distance** [24] is a popular metric for evaluating the similarity between two curves. It
 183 is defined as the shortest pairwise distance between the two curves able to traverse both curves
 184 completely. Since the planned paths are 2D curved shapes, we choose the Frechet distance [24] to
 185 measure the similarity between the planned paths using our model and the GT paths.

186 **Average Angle Difference** evaluates if the inpainted map (IM) allows the vehicle to make more
 187 informed decisions around turns by comparing the orientation (θ) in the world frame of the planned
 188 trajectories for the original lidar (OL) and IM compared to the GT map. For each trajectory \mathcal{T}
 189 (from either OL or IM) we compare (θ) of each path node to (θ) of the closest node from the GT
 190 trajectory. The average of the angle differences for each trajectory evaluates how accurately the
 191 planned trajectory follows the GT trajectory. The Average Angle Difference (AAD) calculation is
 192 introduced

$$AAD(in, gt) = \frac{1}{n} \sum_{i \in \mathcal{T}} |\theta_{in,i} - \theta_{gt,i}| \quad (7)$$

193 where $\theta_{in,i}$ refers to the i th node of the OL or the IM and $\theta_{gt,i}$ refers to the i th node of the GT map.
 194 n refers to the number of planned trajectory nodes in the IM or the OL.

195 **Accuracy of Major Branch Prediction per Skeleton** compares the skeletonized road graphs (Fig-
 196 ure 3 middle row) for OL and IM to skeletonized graphs for GT maps. Comparing the skeletonized
 197 graphs allows for evaluation of the improvement the IMs have on the range of planner and how
 198 closely the predictions from the inpainting model match the GT map from a practical planning per-
 199 spective. If the generated skeleton provides more planning hypotheses than the OL, it allows for
 200 planning with a wider range of start and end locations. The number of road branches (defined as the
 201 number of different turns that can be taken at an intersection) in the skeleton is the metric used to
 202 indicate the possible range of planning hypotheses. To compute the road branch prediction accuracy
 203 we calculate the percentage $\frac{N_{im}}{N_{gt}} \times 100$ where N_{im} and N_{gt} are the number of branches in the IM
 204 and GT skeletons respectively. We perform the same evaluation for OL.

205 **Path Length** directly reflects if the planner plans to the desired goal location in the GT map pre-
 206 cisely. We evaluate the accuracy of the planned trajectory by comparing its length with the GT
 207 trajectory. The length of the trajectories are calculated using the L2 distances between path nodes.
 208 We compare the length of the IM trajectory \mathcal{L}_{im} with the length of the GT trajectory \mathcal{L}_{gt} using the
 209 fraction $\frac{\mathcal{L}_{im}}{\mathcal{L}_{gt}}$. We perform the same evaluation for the length of the trajectory planned using the OL.

210 **Planning Ahead Frames** measures how far in advance a turn is detected. If a map is sparse and
 211 limited by occlusions around turns, the planner cannot predict the turning behaviour because it does
 212 not realize the existence of a road turn. Thus, to evaluate how far ahead a turn is detected, we count
 213 how many frames before the turn or intersection the planner plans turning behavior.

214 4.4 Qualitative Evaluation

215 We show the evaluation results qualitatively in Figure 3. We first evaluate the top row for the
 216 semantic inpainting results. On the GT map (top second), the orange oval indicates a region occluded
 217 from the original lidar (OL) scan (top left). This is because the road section in the orange oval occurs
 218 after a turn so buildings and obstacles in the scene occlude the road after the turn. Our inpainting
 219 model (top middle) is able to predict the road’s existence after the turn and accurately fill out the rest
 220 of the intersection pixels occluded in the OL scan.

221 The middle row compares the skeletonized road maps for our planner. Because the inpainting model
 222 fills in the occluded regions of the map around the turns, the skeletonized map for the inpainted map
 223 (IM) is much more similar in its road branching structure to the GT than the OL.

224 The bottom row compares the planned paths from hybrid A*. In the OL scan (bottom left), the
 225 planner is unable to see the turn highlighted by the orange oval (top second) and so the planner can
 226 only plan a path going straight forwards. In the IM (bottom middle), the planner is able to predict
 227 the upcoming turn due to the inpainting model so it plans a path around the turn that matches the
 228 GT path.

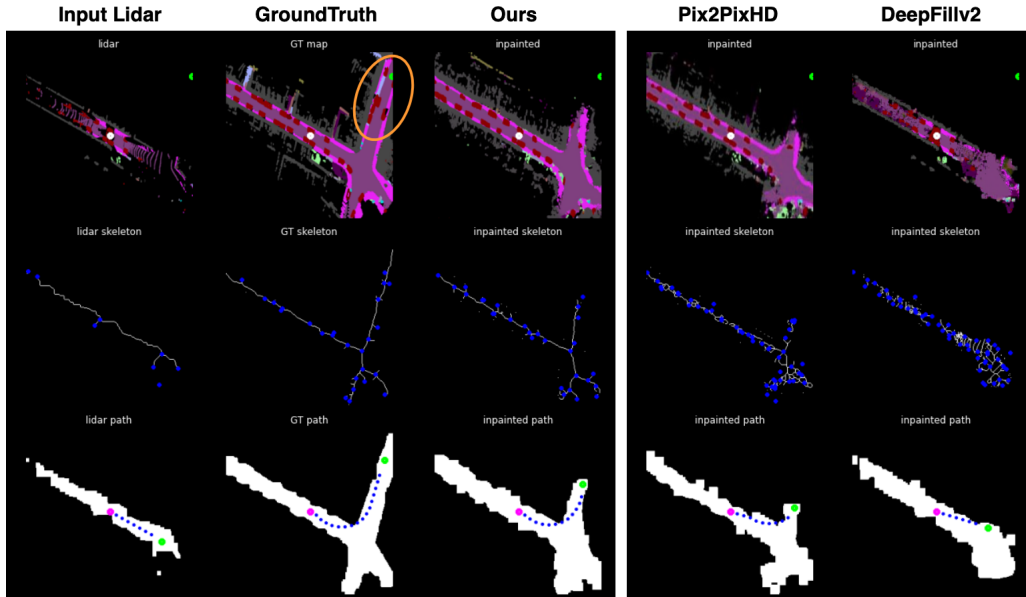


Figure 3: Comparison of inpainting and planner performance with our model and baselines: Pix2PixHD [16], DeepFillv2 [18] and ground truth (GT). Columns: (leftmost) OL, (second) GT, (third) Ours, (fourth) Pix2PixHD [16] and (rightmost) DeepFillv2 [18] Rows: (top) the semantic point cloud; the white dot is vehicle pose and the green dot is the end point. The orange oval on the the top row of the GT semantic map highlights a section of the road that is occluded from the OL. (Middle) The skeletonized road with waypoints (blue dots), (bottom) the planned path; the start is in pink, the goal is in green, and the planned nodes are in blue.

229 4.5 Quantitative Evaluation

230 To justify our model for the planning-under-occlusions task, we choose three different state-of-the-
 231 art inpainting networks (Pix2Pix [1], Pix2PixHD [16] and DeepFillv2 [18]) to conduct baseline
 232 comparisons. We also compare our model with the evaluation results we obtain from original lidar
 233 maps (OL) to show that generative network’s predictive capabilities are an improvement for path
 234 planning compared to the OL. In addition we split our data into easy and hard sets. The easy set
 235 includes straight roads which are simple to navigate even with occluding objects and the hard set
 236 includes turns and intersections which are difficult to navigate with occluding objects. This allows
 237 for evaluation for how our framework performs in environments with different levels of difficulty
 238 (Table 3).

239 The evaluation results shown in Table 2 demonstrate our inpainting model outperforms the OL
 240 scans and the baselines for the path planning task. The worst performer across all metrics is the
 241 path planner on the OL scans, which means predicting semantics in unknown and occluded spaces
 242 is useful for path planning. Our model achieves the lowest Frechet Distance (8.38) and the lowest
 243 Average Angle Difference (6.73) which means our model plans the most similar paths to the GT.
 244 Using the GT map as reference, our model also achieves the highest accuracy of major branch

Table 2: **Task planning evaluation on baselines and our model using the described metrics.** Rows are models and columns are metrics. The best result for each column is **bold**.

Model \ Metric	Frechet distance (pixel)	Average Angle Difference ($^{\circ}$)	Accuracy of Major Branch Prediction (%)	Frame Planned Ahead (#)	Path length (%)	Inference Time(FPS)
OL	20.80	12.81	62.47	32.25	65.21	N/A
DeepFillv2 [18]	18.58	12.73	63.45	35.50	71.90	14
Pix2Pix [1]	12.09	10.33	75.40	42.50	76.63	16
Pix2PixHD [16]	10.06	9.67	77.88	47.25	77.02	33
Ours	8.38	6.73	93.01	52.25	82.97	33

Table 3: **Task planning evaluation on data divided into hard and easy subsets.** Rows are models and columns are metrics. Entries are split into (hard/easy) results. The best result for each column is **bold**.

Model \ Metric	Frechet distance (hard / easy) (pixel)	Average Angle Difference (hard / easy) ($^{\circ}$)	Accuracy of Major Branch Prediction (hard / easy) (%)	Frame Planned Ahead (hard / easy) (#)	Path length (hard / easy) (%)
OL	22.67 / 17.82	14.90 / 7.74	58.44 / 70.25	19.75 / 12.5	62.67 / 93.9
DeepFillv2 [18]	18.91 / 15.17	12.89 / 6.76	61.56 / 79.40	23.00 / 12.5	69.92 / 94.81
Pix2Pix [1]	11.13 / 9.24	11.69 / 6.95	76.10 / 76.87	29.00 / 13.5	72.63 / 95.84
Pix2PixHD [16]	10.51 / 8.43	10.44 / 6.38	77.66 / 79.40	33.75 / 13.5	73.38 / 95.79
Ours	8.55 / 5.33	8.85 / 3.39	90.06 / 89.16	38.75 / 13.5	83.75 / 96.63

245 prediction (93.01%) and path length (82.97%), which shows our model generates a road network
 246 that is the most accurate compared to the ground truth.

247 The evaluation results in Table 3, where our test set splits into hard and easy subsets, demonstrate
 248 the inpainting model improves performance more for complex scenes. For the Frechet Distance, our
 249 model improves over the OL by 14.12 and 12.49 (pixels) for hard and easy data respectively. For
 250 Average Angle Difference the improvement is 6.05 and 4.35 ($^{\circ}$). For accuracy of branch predic-
 251 tion, improvement is 31.62 and 18.91 (%). We especially note that for frames planned ahead the
 252 improvement is the most noticeable at 19 and 1 frames demonstrating our framework is especially
 253 beneficial for more complex navigation scenarios such as turns and intersections. For path length,
 254 the improvement difference from hard and easy data is also large at 21.08 and 2.73 (%) respectively,
 255 demonstrating that in complex scenarios the inpainting model makes a large difference compared to
 256 the OL.

257 We list the inference time tested on one NVIDIA RTX3090 in the last column, which shows our
 258 model achieves acceptable inference speed (33 FPS) on available computation resources.

259 5 Limitations

260 Our paper assumes pose of the vehicle is known and accurate. The community has a large body of
 261 work addressing the SLAM problem which can be used for localization [25]. For our experiments
 262 we assume semantic segmentations of raw lidar point clouds can be generated in real time. Future
 263 work can use models such as [26] to generate real time semantic lidar segmentations. In addition, we
 264 assume that during online deployment, the data will be from the same domain as the training data.
 265 While this is a reasonable assumption for urban environments given the large amount of available
 266 data, in the future our framework can be extended to scenarios where online and training data are
 267 from varied domains by using domain adaptation techniques [27, 28].

268 6 Conclusion

269 This paper presents a novel framework for planning around unknown spaces in occluded, outdoor,
 270 urban environments. We use a data-driven inpainting model to fill in occluded regions and plan
 271 dynamically feasible paths given the model predictions. We introduce a modified loss to the pix2pix
 272 network and also render a dataset for the planning-through-occlusions problem. We specifically
 273 focus on turns and intersections for our evaluation as these are the regions where navigation is most
 274 heavily effected by occlusions. Experiments validate that our framework and model allow for more
 275 informed navigation in occluded spaces (especially turns), and show our planner plans paths much
 276 closer to the ground truth compared to the original sensor measurement.

277 **References**

- 278 [1] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional
279 adversarial networks. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Con-*
280 *ference on*, 2017.
- 281 [2] P. E. Hart, N. J. Nilsson, and B. Raphael. A formal basis for the heuristic determination of
282 minimum cost paths. *IEEE T Syst Sci Cyb*, 4(2):100–107, 1968.
- 283 [3] A. Stentz. Optimal and efficient path planning for partially known environments. In *Proc.*
284 *ICRA*, 1994.
- 285 [4] Y. Han, H. Lin, J. Banfi, K. Bala, and M. Campbell. DeepSemanticHPPC: Hypothesis-based
286 planning over uncertain semantic point clouds. In *Proc. ICRA*, pages 4252–4258, 2020.
- 287 [5] Y. Han, J. Banfi, and M. E. Campbell. Planning paths through unknown space by imagining
288 what lies therein. In *CoRL*, 2020.
- 289 [6] M. Ryll, J. Ware, J. Carter, and N. Roy. Semantic trajectory planning for long-distant
290 unmanned aerial vehicle navigation in urban environments. In *2020 IEEE/RSJ Interna-*
291 *tional Conference on Intelligent Robots and Systems (IROS)*, pages 1551–1558, 2020. doi:
292 [10.1109/IROS45743.2020.9341441](https://doi.org/10.1109/IROS45743.2020.9341441).
- 293 [7] G. Georgakis, B. Bucher, A. Arapin, K. Schmeckpeper, N. Matni, and K. Daniilidis.
294 Uncertainty-driven planner for exploration and navigation. In *2022 International Conference*
295 *on Robotics and Automation (ICRA)*, pages 11295–11302, 2022. doi:[10.1109/ICRA46639.](https://doi.org/10.1109/ICRA46639.2022.9812423)
296 [2022.9812423](https://doi.org/10.1109/ICRA46639.2022.9812423).
- 297 [8] M. Narasimhan, E. Wijmans, X. Chen, T. Darrell, D. Batra, D. Parikh, and A. Singh. Seeing
298 the un-scene: Learning amodal semantic maps for room navigation. In *Computer Vision –*
299 *ECCV 2020*, pages 513–529, Cham, 2020. Springer International Publishing. ISBN 978-3-
300 030-58523-5.
- 301 [9] C. Lu and G. Dubbelman. Semantic foreground inpainting from weak supervision. *IEEE RA-L*,
302 5(2):1334–1341, 2020.
- 303 [10] P. Purkait, C. Zach, and I. Reid. Seeing behind things: Extending semantic segmentation to
304 occluded regions. In *Proc. IROS*, pages 1998–2005, 2019.
- 305 [11] S. Schulter, M. Zhai, N. Jacobs, and M. K. Chandraker. Learning to look around objects for
306 top-view representations of outdoor scenes. *ArXiv*, abs/1803.10870, 2018.
- 307 [12] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville,
308 and Y. Bengio. Generative adversarial nets. *Advances in neural information processing sys-*
309 *tems*, 27, 2014.
- 310 [13] D. Croce, G. Castellucci, and R. Basili. Gan-bert: Generative adversarial learning for robust
311 text classification with a bunch of labeled examples. In *Proceedings of the 58th annual meeting*
312 *of the association for computational linguistics*, pages 2114–2119, 2020.
- 313 [14] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Te-
314 jani, J. Totz, Z. Wang, et al. Photo-realistic single image super-resolution using a generative
315 adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern*
316 *recognition*, pages 4681–4690, 2017.
- 317 [15] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-
318 consistent adversarial networkss. In *Computer Vision (ICCV), 2017 IEEE International Con-*
319 *ference on*, 2017.

- 320 [16] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-resolution im-
321 age synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE*
322 *conference on computer vision and pattern recognition*, pages 8798–8807, 2018.
- 323 [17] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang. Generative image inpainting with
324 contextual attention. *arXiv preprint arXiv:1801.07892*, 2018.
- 325 [18] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang. Free-form image inpainting with gated
326 convolution. In *Proceedings of the IEEE International Conference on Computer Vision*, pages
327 4471–4480, 2019.
- 328 [19] Y. Liao, J. Xie, and A. Geiger. KITTI-360: A novel dataset and benchmarks for urban scene
329 understanding in 2d and 3d. *arXiv preprint arXiv:2109.13410*, 2021.
- 330 [20] T. Park, A. A. Efros, R. Zhang, and J.-Y. Zhu. Contrastive learning for unpaired image-to-
331 image translation. In *European Conference on Computer Vision*, pages 319–345. Springer,
332 2020.
- 333 [21] D. Dolgov, S. Thrun, M. Montemerlo, and J. Diebel. Practical search techniques in path plan-
334 ning for autonomous driving. In *Proceedings of the First International Symposium on Search*
335 *Techniques in Artificial Intelligence and Robotics (STAIR-08, 2008)*.
- 336 [22] T. Y. Zhang and C. Y. Suen. A fast parallel algorithm for thinning digital patterns. *Commun.*
337 *ACM*, 27(3):236–239, mar 1984. ISSN 0001-0782. doi:10.1145/357994.358023. URL <https://doi.org/10.1145/357994.358023>.
- 338 [23] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *The International*
339 *Conference on Learning Representations (ICLR)*, 2015.
- 341 [24] H. Alt and M. Godau. Computing the fréchet distance between two polygonal curves. *Inter-*
342 *national Journal of Computational Geometry & Applications*, 5(01n02):75–91, 1995.
- 343 [25] H. Durrant-Whyte and T. Bailey. Simultaneous localisation and mapping (slam): Part i the
344 essential algorithms. *IEEE ROBOTICS AND AUTOMATION MAGAZINE*, 2:2006, 2006.
- 345 [26] A. Milioto, I. Vizzo, J. Behley, and C. Stachniss. RangeNet++: Fast and Accurate LiDAR
346 Semantic Segmentation. In *IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*,
347 2019.
- 348 [27] G. Csurka, R. Volpi, and B. Chidlovskii. Unsupervised domain adaptation for semantic image
349 segmentation: a comprehensive survey. *arXiv preprint arXiv:2112.03241*, 2021.
- 350 [28] Y. Liu, W. Zhang, and J. Wang. Source-free domain adaptation for semantic segmentation. In
351 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
352 1215–1224, 2021.