

MonoLayout: Amodal scene layout from a single image

Kaustubh Mani^{*1}, Swapnil Daga¹, Shubhika Garg², N. Sai Shankar¹, Krishna Murthy Jatavallabhula^{†3,4}, and K. Madhava Krishna^{†1}

¹Robotics Research Center, IIIT Hyderabad, ²IIT Kharagpur, ³Mila - Quebec AI Institute, Montreal, ⁴Université de Montréal



Figure 1: **MonoLayout**: Given only a single image of a road scene, we propose a neural network architecture that reasons about the *amodal* scene layout in bird’s eye view in *real-time* (30 fps). Our approach, dubbed *MonoLayout* can *hallucinate* regions of the static scene (road, sidewalks)—and traffic participants—that do not even project to the visible regime of the image plane. Shown above are example images from the KITTI [10] (left) and Argoverse [5] (right) datasets. *MonoLayout* outperforms prior art (by more than a 20% margin) on hallucinating occluded regions.

Abstract

In this paper, we address the novel, highly challenging problem of estimating the layout of a complex urban driving scenario. Given a single color image captured from a driving platform, we aim to predict the bird’s eye view layout of the road and other traffic participants. The estimated layout should reason beyond what is visible in the image, and compensate for the loss of 3D information due to projection. We dub this problem amodal scene layout estimation, which involves hallucinating scene layout for even parts of the world that are occluded in the image. To this end, we present MonoLayout, a deep neural network for real-time amodal scene layout estimation from a single image. We represent scene layout as a multi-channel semantic occupancy grid, and leverage adversarial feature learning to “hallucinate” plausible completions for occluded image parts. We extend several state-of-the-art approaches for road-layout estimation and vehicle occupancy estimation in bird’s eye view to the amodal setup for rigorous evaluation. By leveraging temporal sensor fusion to generate training labels, we significantly outperform current art over a number of datasets. A video abstract

of this paper is available [here](#).

1. Introduction

The advent of autonomous driving platforms has led to several interesting, new avenues in perception and scene understanding. While most of the industrially-led solutions leverage powerful sensors (eg. lidar, precision GPS, etc.), an interesting research question is to push the capabilities of monocular vision sensors. To this end, we consider the novel and highly challenging task of estimating *scene layout* in bird’s eye view, given only a single color image.

Humans have a remarkable cognitive capability of perceiving *amodal* attributes of objects in an image. For example, upon looking at an image of a vehicle, humans can nominally *infer* the occluded parts, and also the potential geometry of the surroundings of the vehicle. While modern neural networks outperform humans in image recognition and object detection [8, 11, 15, 16, 29, 32, 41], they still lack this innate cognitive capability of reasoning beyond image evidence. With this motivation, we propose *MonoLayout*, a neural architecture that takes as input a color image of a road scene, and outputs the *amodal* scene layout in bird’s eye view. *MonoLayout* maps road regions, sidewalks, as well as regions occupied by other

*Corresponding author: kaustubh3095@gmail.com

Project page: <https://hbutsuak95.github.io/monolayout>

[†]Both (un)supervisors procrastinated equally

traffic participants such as cars, to bird’s eye view, in a single pass, leveraging adversarial feature learning.

To the best of our knowledge, *MonoLayout* is the first approach to *amodally* reason about the static and dynamic objects in a scene. We show that, by using a shared context to reason about scene entities, *MonoLayout* achieves better performance on each task when compared to approaches that train only for a particular task. On the task of amodal scene layout estimation, *MonoLayout* outperforms all evaluated baselines by a significant margin on several subsets of the KITTI [10] and the Argoverse [5] datasets. Further, *MonoLayout* achieves state-of-the-art object detection performance in bird’s eye view, without using any form of thresholding / postprocessing. In summary, our contributions are the following:

1. We propose *MonoLayout*, a practically motivated deep architecture to estimate the *amodal* scene layout from just a single image (*c.f.* Fig. 1).
2. We demonstrate that adversarial learning can be used to further enhance the quality of the estimated layouts, specifically when hallucinating large missing chunks of a scene (*c.f.* Fig. 1, Sec. 3).
3. We evaluate against several state-of-the-art approaches, and outperform all of them by a significant margin on a number of established benchmarks (KITTI-Raw, KITTI-Object, KITTI-Odometry [10], Argoverse [5], *c.f.* Sec. 4, Table 1).
4. Further, we show that *MonoLayout* can also be efficiently trained on datasets that do not contain lidar scans by leveraging recent successes in monocular depth estimation. [12] (*c.f.* Table 2).

Please refer to the appendix for more results, where we demonstrate that the extracted amodal layouts can suit several higher level tasks, such as (but not limited to) multi-object tracking, trajectory forecasting, etc.

2. Related Work

To the best of our knowledge, no published approach has tackled the task of simultaneous road layout (static scene) and traffic participant (dynamic scene) estimation from a single image. However, several recent approaches have addressed the problem of estimating the layout of a road scene, and several other independent approaches have tackled 3D object detection. We summarize the most closely related approaches in this section.

Road layout estimation

Schulter et al. [34] proposed one of the first approaches to estimate an occlusion-reasoned bird’s eye view road layout from a single color image. They use monocular

depth estimation [12] as well as semantic segmentation to bootstrap a CNN that predicts occluded road layout. They use priors from OpenStreetMap [26] to adversarially regularize the estimates. More recently, Wang *et al.* [38] builds on top of [34] to infer parameterized road layouts. Our approach does not require to be bootstrapped using either semantic segmentation or monocular depth estimation, and can be trained end-to-end from color images.

Perhaps the closest approach to ours is MonoOccupancy [24], which builds a variational autoencoder (VAE) to predict road layout from a given image. They also present results for extracting regions close to roads (eg. sidewalk, terrain, non-free space). However, they reason only about the pixels present in the image, and not beyond occluding obstacles. Also, the bottleneck enforced by that leads to non-sharp, blob-like layouts. On the other hand, *MonoLayout* estimates *amodal* scene layouts, reasoning beyond occlusion boundaries. Our approach produces crisp road edges as well as vehicle boundaries, by leveraging adversarial feature learning and sensor fusion to reduce noise in the labeled ground-truth training data.

Object detection in bird’s eye view

There exist several approaches to 3D object detection that exclusively use lidar [2, 35, 39], or a combination of camera and lidar sensors [7, 20, 22]. However, there are only a handful of approaches that purely use monocular vision for object detection [6, 21, 25]. Most of these are two stage approaches, comprising a region-proposal stage, and a classification stage.

Another category of approaches map a monocular image to a bird’s eye view representation [30], thereby reducing the task of 3D object detection to that of 2D image segmentation. Recently, BirdGAN [36] leveraged adversarial learning for mapping images to bird’s eye view, where lidar object detectors such as [2] were repurposed for object detection.

Such techniques usually require a pre-processing stage (usually a neural network that maps an image to a bird’s eye view) after which further processing is applied. On the other hand, we demonstrate that we can achieve significantly higher accuracy by directly mapping from the image space to objects in bird’s eye view, bypassing the need for a pre-processing stage altogether.

More notably, all the above approaches require a post-processing step that usually involves non-maximum suppression / thresholding to output object detections. *MonoLayout* neither requires pre-processing nor post-processing and it directly estimates

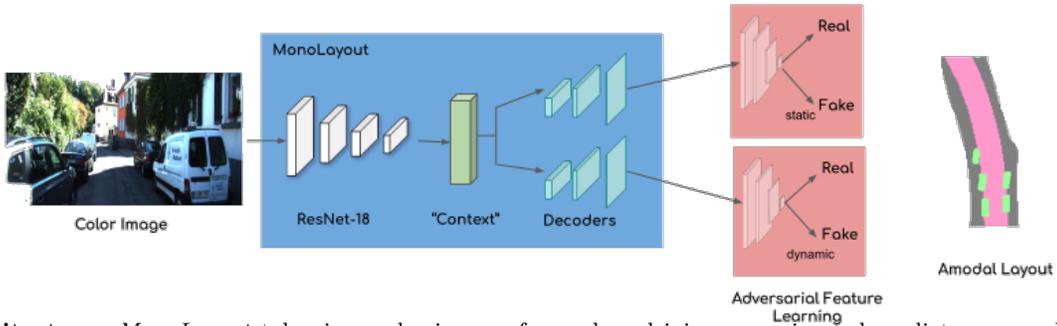


Figure 2: **Architecture:** *MonoLayout* takes in a color image of an urban driving scenario, and predicts an amodal scene layout in bird’s eye view. The architecture comprises a *context encoder*, *amodal layout decoders*, and *two discriminators*.

scene layouts that can be evaluated (or plugged into other task pipelines).

3. MonoLayout: Monocular Layout Estimation

3.1. Problem Formulation

In this paper, we address the problem of *amodal* scene layout estimation from a single color image. Formally, given a color image \mathcal{I} captured from an autonomous driving platform, we aim to predict a bird’s eye view layout of the static and dynamic elements of the scene. Concretely, we wish to estimate the following three quantities.¹

1. The set of all static scene points \mathcal{S} (typically the road and the sidewalk) on the ground plane (within a rectangular range of length L and width W , in front of the camera), *regardless of whether or not they are imaged by the camera*.
2. The set of all dynamic scene points \mathcal{D} on the ground plane (within the same rectangular range as above) occupied by vehicles, *regardless of whether or not they are imaged by the camera*.
3. For each point discerned in the above step as being occupied by a vehicle, an instance-specific labeling of which vehicle the point belongs to.

3.2. MonoLayout

The problem of amodal scene layout estimation is challenging from a neural networks standpoint in several interesting ways. First, it necessitates that we learn *good* visual representations from images that help in estimating 3D properties of a scene. Second, it requires these representations to reason beyond classic 3D reconstruction; these representations must enable us to *hallucinate* 3D geometries of image regions that

¹*Flat-earth assumption:* For the scope of this paper, we assume that the autonomous vehicle is operating within a bounded geographic area of the size of a typical city, and that all roads in consideration are *somewhat* planar, i.e., no steep/graded roads on mountains.

are occupied by occluders. Furthermore, the learned representations must implicitly disentangle the static parts of the scene (occupied by road points) from the dynamic objects (eg. parked/moving cars). With these requirements in mind, we design *MonoLayout* with the following components.

Maximum a posteriori estimation

We formulate the amodal road layout estimation problem as that of recovering the Maximum a posteriori (MAP) estimate of the distribution of scene statics and dynamics. Given the image \mathcal{I} , we wish to recover the posterior $P(\mathcal{S}, \mathcal{D}|\mathcal{I})$, over the domain $\Omega \triangleq \{(x, H, z) \mid \|x - x_0\|_1 \leq L; \|z - z_0\|_1 \leq W; (z - z_0) > 0\}$ ². Note that the static (road) and dynamic (vehicle) marginals are not independent. They are not independent - they exhibit high correlation (vehicles ply on roads). Hence, we introduce an additional conditioning context variable \mathcal{C} that can be purely derived only using the image information \mathcal{I} , such that, \mathcal{S} and \mathcal{D} are conditionally independent given \mathcal{C} . We term this conditioning variable as the "shared context" as it necessarily encompasses the information needed to estimate the static and dynamic layout marginals. This allows the posterior to be factorized in the following form.

$$\begin{aligned}
 P(\mathcal{S}, \mathcal{D}|\mathcal{I}) &\propto P(\mathcal{S}, \mathcal{D}, \mathcal{C}|\mathcal{I}) \\
 &= \underbrace{P(\mathcal{S}|\mathcal{C}, \mathcal{I})}_{\text{static decoder}} \underbrace{P(\mathcal{D}|\mathcal{C}, \mathcal{I})}_{\text{dynamic decoder}} \underbrace{P(\mathcal{C}|\mathcal{I})}_{\text{context encoder}}
 \end{aligned} \tag{1}$$

In accordance with the above factorization of the posterior, the architecture of *MonoLayout* comprises three subnetworks (*c.f.* Fig. 2).

1. A *context encoder* which extracts multi-scale feature representations from the input monocular image. This provides a shared context that captures

²This domain is a rectangular region in bird’s eye view. H is the height of the camera above the ground.

static as well as dynamic scene components for subsequent processing.

2. An *amodal static scene decoder* which decodes the shared context to produce an amodal layout of the static scene. This model consists of a series of deconvolution and upsampling layers that map the shared context to a static scene bird’s eye view.
3. A *dynamic scene decoder* which is architecturally similar to the road decoder and predicts the vehicle occupancies in bird’s eye view.
4. Two *discriminators* [17, 33] which regularize the predicted static/dynamic layouts by regularizing their distributions to be *similar* to the *true* distribution of plausible road geometries (which can be easily extracted from huge unpaired databases such as OpenStreetMap [26]) and ground-truth vehicle occupancies.

Feature Extractor

From the input image, we first extract meaningful image features at multiple scales, using a ResNet-18 encoder (pre-trained on ImageNet [9]). We finetune this feature extractor in order for it to learn low-level features that help reason about static and dynamic aspects of the scene.

Static and dynamic layout decoders

The static and dynamic layout decoders share an identical architecture. They decode the shared context from the feature extractor by a series of upsampling layers to output a $D \times D$ grid each³.

Adversarial Feature Learning

To better ground the likelihoods $P(S|C, \mathcal{I})$, $P(D|C, \mathcal{I})$ (c.f. Eq. 1), we introduce adversarial regularizers (discriminators) parameterized by θ_S and θ_D respectively. The layouts estimated by the static and dynamic decoders are input to these patch-based discriminators [17]. The discriminators regularize the distribution of the output layouts (*fake data* distribution, in GAN [13] parlance) to match a prior data distribution of *conceivable* scene layouts (*true data* distribution). This prior data distribution is a collection of road snippets from OpenStreetMap [26], and rasterized images of vehicles in bird’s eye view. Instead of training with a paired

³We tried training a single decoder for both the tasks, but found convergence to be hard. We attribute this to the extreme change in output spaces: while roads are large, continuous chunks of space, cars are tiny, sparse chunks spread over the entire gamut of pixels. Instead, we chose to train two decoders over a shared context, which bypasses this discrepancy in output spaces, and results in sharp layout estimates.

OSM for each image, we choose to collect a set of diverse OSM maps representing the true data distribution of road layouts in bird’s eye view and train our discriminators in an unpaired fashion. This mitigates the need to have perfectly aligned OSM views to the current image, making *MonoLayout* favourable compared to approaches like [34] that perform an explicit alignment of the OSM before beginning processing.

Loss function

The parameters ϕ, ν, ψ of the context encoder, the amodal static scene decoder, and the dynamic scene decoder respectively are obtained by minimizing the following objective using minibatch stochastic gradient descent.

$$\begin{aligned} & \min_{\phi, \nu, \psi, \theta_S, \theta_D} \mathcal{L}_{sup}(\phi, \nu, \psi) + \mathcal{L}_{adv}(\phi, \theta, \psi) + \mathcal{L}_{discr}(\phi, \nu) \\ \mathcal{L}_{sup} &= \sum_{i=1}^N \|\mathcal{S}_{\phi, \nu}(\mathcal{I}^i) - \mathcal{S}_{gt}^i\|^2 + \|\mathcal{D}_{\phi, \psi}(\mathcal{I}^i) - \mathcal{D}_{gt}^i\|^2 \\ \mathcal{L}_{adv}(S, D; \phi, \theta, \psi) &= \mathbb{E}_{\theta \sim p_{fake}} [(D(\theta_S) - 1)^2] \\ &+ \mathbb{E}_{\theta \sim p_{fake}} [(D(\theta_D) - 1)^2] \\ \mathcal{L}_{discr}(D; \theta) &= \sum_{\theta \in \{\theta_D, \theta_S\}} \mathbb{E}_{\theta \sim p_{true}} [(D(\theta) - 1)^2] \\ &+ \mathbb{E}_{\theta \sim p_{fake}} [(D(\theta) - 1)^2] \end{aligned}$$

Here, \mathcal{L}_{sup} is a supervised ($L2$) error term that penalizes the deviation of the predicted static and dynamic layouts ($\mathcal{S}_{\phi, \nu}(\mathcal{I})$, $\mathcal{D}_{\phi, \nu}(\mathcal{I})$) with their corresponding ground-truth values (\mathcal{S}_{gt}^i , \mathcal{D}_{gt}^i). The adversarial loss \mathcal{L}_{adv} encourages the distribution of layout estimates from the static/dynamic scene decoders (p_{fake}) to be close to their true counterparts (p_{true}). The discriminator loss \mathcal{L}_{discr} is the discriminator update objective [13].

3.3. Generating training data: sensor fusion

Since we aim to recover the amodal scene layout, we are faced with the problem of extracting training labels for even those parts of the scene that are occluded from view. While recent autonomous driving benchmarks provide synchronized lidar scans as well as semantic information for each point in the scan, we propose a sensor fusion approach to generate more robust training labels, as well as to handle scenarios in which direct 3D information (eg. lidar) may not be available.

As such, we use either monocular depth estimation networks (Monodepth2 [12]) or raw lidar data and initialize a pointcloud in the camera coordinate

Dataset	Method	Static Layout Estimation			Vehicle Layout	
		Road mIoU	Sidewalk mIoU	Road + Sidewalk occl mIoU	mIoU	mAP
KITTI Raw	MonoOccupancy [24]	56.16	18.18	28.24	-	-
	Schulter <i>et al.</i> [34]	68.89	30.35	61.06	-	-
	MonoLayout-static (Ours)	73.86	32.86	67.42	-	-
KITTI Odometry	MonoOccupancy [24]	64.72	12.08	34.87	-	-
	MonoLayout-static (Ours)	80.08	42.66	72.46	-	-
KITTI Object	MonoOccupancy-ext	-	-	-	20.45	22.59
	Mono3D [6]	-	-	-	17.11	29.62
	OFT [30]	-	-	-	25.24	34.69
	MonoLayout-dynamic (Ours)	-	-	-	26.08	40.79
KITTI Tracking	MonoLayout (Ours)	53.94	-	-	24.16	36.83
Argoverse	MonoOccupancy-ext	32.70	-	-	16.22	38.66
	MonoLayout (Ours)	58.33	-	-	32.05	48.31

Table 1: **Quantitative results:** We evaluate the performance of *MonoLayout* on several datasets, on amodal scene layout estimation. As there is no existing baseline that simultaneously estimates static (road/sidewalk) as well as dynamic (vehicle) layout, we evaluate under multiple settings. On the KITTI Raw and KITTI Odometry [10] datasets, we evaluate *MonoLayout-static*. On the KITTI Object [10] dataset, we evaluate *MonoLayout-dynamic*. On the KITTI Tracking [10] and Argoverse [5] datasets, we evaluate *MonoLayout*, the full architecture that estimates both static and dynamic layouts. We outperform existing approaches by a significant margin, on all metrics.

frame. Using odometry information over a window of W frames, we aggregate/register the sensor observations over time, to generate a more dense, noise-free pointcloud. Note that, when using monocular depth estimation, we discard depths of points that are more than 5 meters away from the car, as they are noisy. To compensate for this narrow field of view, we aggregate depth values over a much larger window size (usually 40 – 50) frames.

The dense pointcloud is then projected to an occupancy grid in bird’s eye view. If ground-truth semantic segmentation labels are available, each occupancy grid cell is assigned the label based on a simple majority over the labels of the corresponding points. For the case where ground-truth semantic labels are unavailable, we use a state-of-the-art semantic segmentation network [32] to segment each frame and aggregate these predictions into the occupancy grid.

For vehicle occupancy estimation though, we rely on ground-truth labels in bird’s eye view, and train only on datasets that contain such labels⁴.

4. Experiments

To evaluate *MonoLayout*, we conduct experiments over a variety of challenging scenarios and against multiple baselines, for the task of *amodal* scene layout estimation.

⁴For a detailed description of the architecture and the training process, we refer the reader to the appendix

4.1. Datasets

We present our results on two different datasets - KITTI [10] and Argoverse [5]. The latter contains a *high-resolution* semantic occupancy grid in bird’s eye view, which facilitates the evaluation of amodal scene layout estimation. The KITTI dataset, however, has no such provision. For a semblance of *ground-truth*, we register depth and semantic segmentation of lidar scans in bird’s eye view.

To the best of our knowledge, there is no published prior work that reasons jointly about road and vehicle occupancies. However, there exist approaches for road layout estimation [24, 34], and a separate set of approaches for vehicle detection [6, 30]. Furthermore, each of these approaches evaluate over different datasets (and in cases, different train/validation splits). To ensure fair comparison with all such approaches, we organize our results into the following categories.

- Baseline Comparison:** For a fair comparison with state-of-the-art road layout estimation techniques, we evaluate performance on the KITTI RAW split used in [34] (10156 training images, 5074 validation images). For a fair comparison with state-of-the-art 3D vehicle detection approaches we evaluate performance on the KITTI 3D object detection split of Chen et al. [6] (3712 training images, 3769 validation images).
- Amodal Layout Estimation:** To evaluate layout estimation on both static and dynamic scene attributes (road, vehicles), we use the KITTI Tracking [10] and Argoverse [5] datasets. We an-

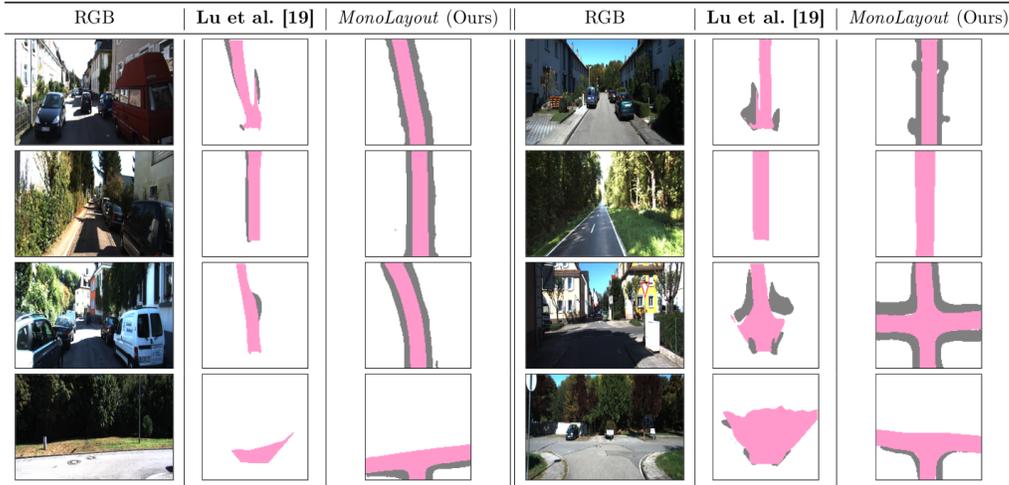


Figure 3: **Static layout estimation:** Observe how *MonoLayout* performs amodal completion of the static scene (road shown in pink, sidewalk shown in gray). MonoOccupancy [24] fails to reason beyond occluding objects (top row), and does not hallucinate large missing patches (bottom row), while *MonoLayout*(Ours) is accurately able to do so. Furthermore, even in cases where there is no occlusion (row 2), *MonoLayout*(Ours) generates road layouts of much sharper quality. Row 3 show extremely challenging scenarios where most of the view is blocked by vehicles, and the scenes exhibit high-dynamic range (HDR) and shadows.

notate sequences from the KITTI Tracking dataset for evaluation (5773 training images, 2235 validation images). Argoverse provides HD maps as well as vehicle detections in bird’s eye view (6723 training images, 2418 validation images).

3. **Temporal sensor fusion for supervision:** We then present results using our data generation approach (*c.f.* Sec. 3.3) on the KITTI Odometry dataset. This also uses the dense semantic segmentation labels from the Semantic KITTI dataset [1].

4.2. Approaches evaluated

We evaluate the performance of the following approaches.

- *Schulter et al.*: The static scene layout estimation approach proposed in [34].
- *MonoOccupancy*: The static scene layout estimation approach proposed in [24].
- *Mono3D*: The monocular 3D object detection approach from [6].
- *OFT*: A recent, state-of-the-art monocular bird’s eye view detector [30].
- *MonoOccupancy-ext*: We extend MonoOccupancy [24] to predict vehicle occupancies.
- *MonoLayout-static*: A version of *MonoLayout* that only predicts static scene layouts.

- *MonoLayout-dynamic*: A version of *MonoLayout* that only predicts vehicle occupancies.
- *MonoLayout*: The full architecture, that predicts both static and dynamic scene layouts.

4.3. Static layout estimation (Road)

We evaluate *MonoLayout-static* against Schuster *et al.* [34] and MonoOccupancy [24] on the task of static scene (road/sidewalk) layout estimation. Note that Schuster *et al.* assume that the input image is first passed through monocular depth estimation and semantic segmentation networks, while we operate directly on raw pixel intensities. Table 1 summarizes the performance of existing road layout estimation approaches (Schuster *et al.* [34], MonoOccupancy [24]) on the KITTI Raw and KITTI Odometry benchmarks. For KITTI Raw, we follow the exact split used in Schuster *et al.* and retrain MonoOccupancy and *MonoLayout-static* on this train split. Since the manual annotations for semantic classes for KITTI Raw aren’t publicly available, we manually annotated sequences with semantic labels (and will make them publicly available).

From Table 1 (*c.f.* “KITTI Raw”), we see that *MonoLayout-static* outperforms both MonoOccupancy [24] and Schuster *et al.* [34] by a significant margin. We attribute this to the strong hallucination capabilities of *MonoLayout-static* due to adversarial feature learning (*c.f.* Fig. 3). Although Schuster *et al.* [34] use a discriminator to regularize layout predictions, they

seem to suffer from cascading errors due to sequential processing blocks (eg. depth, semantics extraction). MonoOccupancy [24] does not output sharp estimates of road boundaries by virtue of being a variational autoencoder (VAE), as mean-squared error objectives and Gaussian prior assumptions often result in *blurry* generation of samples [23]. The hallucination capability is much more evident in the occluded region evaluation, where we see that *MonoLayout*-static improves by more than 10% on prior art.

4.4. Dynamic (vehicle) layout estimation

To evaluate vehicle layout estimation in bird’s eye view, we first present a comparative analysis on the KITTI Object [10] dataset, for a fair evaluation with respect to prior art. Specifically, we compare against *Orthographic Feature Transform* (OFT [30]), the current best monocular object detector in bird’s eye view. We also evaluate against Mono3D [6], as a baseline. And, we extend MonoOccupancy [24] to perform vehicle layout estimation, to demonstrate that VAE-style architectures are ill-suited to this purpose. This comparison is presented in Table 1 (“KITTI Object”).

We see again that *MonoLayout*-dynamic outperforms prior art on the task of vehicle layout estimation. Note that Mono3D [6] is a two-stage method and requires strictly additional information (semantic and instance segmentation), and OFT [30] performs explicit orthographic transforms and is parameter-heavy (23.5M parameters) which slows it down considerably (5 fps). We make no such assumptions and operate on raw image intensities, yet obtain better performance, and about a 6x speedup (32 fps). Also, MonoOccupancy [24] does not perform well on the vehicle occupancy estimation task, as the variational autoencoder-style architecture usually *merges* most of the vehicles into large *blob-like* structures (*c.f.* Fig. 4).

4.5. Amodal scene layout estimation

In the previous sections, we presented results individually for the tasks of static (road) and dynamic (vehicle) layout estimation, to facilitate comparison with prior art *on equal footing*. We now present results for *amodal* scene layout estimation (i.e., both static and dynamic scene components) on the Argoverse [5] dataset. We chose Argoverse [5] as it readily provides ground truth HD bird’s eye view maps for both road and vehicle occupancies. We follow the train-validation splits provided by [5], and summarize our results in Table 1 (“Argoverse”).

We show substantial improvements at nearly 20% on mIoU vis-a-vis the next closest baseline [24]. This is a demonstration of the fact that, even when perfect

ground-truth is available, approaches such as MonoOccupancy [24] fail to reach the levels of performance as that of *MonoLayout*. We attribute this to the *shared context* that encapsulates a rich feature collection to facilitate both static and dynamic layout estimation. Note that, for the Argoverse [5] dataset we *do not* train our methods on the ground-truth HD maps, because such fine maps aren’t typically available for all autonomous driving solutions. Instead, we train our methods using a semblance of ground-truth (generated by the process described in Sec 3.3), and use the HD maps only for evaluation. This validates our claim that our model, despite being trained using noisy ground estimates by leveraging sensor fusion, is still able to hallucinate and complete the occluded parts of scenes correctly as shown in Fig 6.

4.6. Ablation Studies

Using monocular depth as opposed to lidar

While most of the earlier results focused on the scenario where explicit lidar annotations were available, we turn to the more interesting case where the dataset only comprises monocular images. As described in Sec 3.3, we use monocular depth estimation (MonoDepth2 [12]) and aggregate/register depthmaps over time, to provide training signal. In Table 2, we analyze the impact of the availability of lidar data on the performance of amodal scene layout estimation.

We train MonoOccupancy [24] and *MonoLayout*-static on the KITTI Raw dataset, using monocular depth estimation-based ground-truth, as well as lidar-based ground-truth. While lidar based variants perform better (as is to be expected), we see that self-supervised monocular depth estimation results in reasonable performance too. Surprisingly, for the per-frame case (i.e., no sensor fusion), monocular depth based supervision seems to fare better. Under similar conditions of supervision, we find that *MonoLayout*-static outperforms MonoOccupancy [24].

Impact of sensor fusion

The sensor fusion technique described in Sec 3.3 greatly enhances the accuracy of static layout estimation. Aggregating sensor observations over time equips us with more comprehensive, and noise-free maps. Table 2 presents results for an analysis of the performance benefits obtained due to temporal sensor fusion.

Impact of adversarial learning

With the discriminators, we not only improve qualitatively (sharper/realistic samples) (*c.f.* Fig 7), but

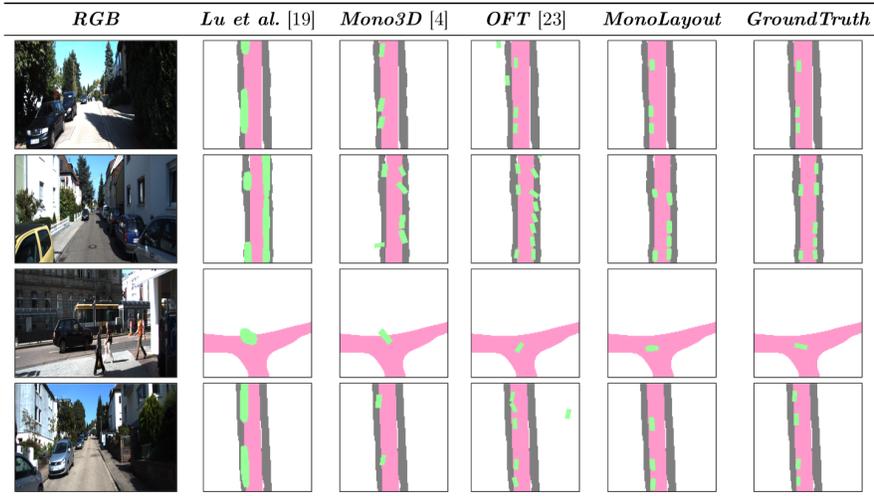


Figure 4: **Dynamic layout estimation:** We show vehicle occupancy estimation results on the KITTI [10] 3D Object detection benchmark. From left to right, the column corresponds to the input image, MonoOccupancy [24], Mono3D [6], OFT [30], *MonoLayout* (Ours), and ground-truth respectively. While the other approaches miss out on detecting cars (top row), or split a vehicle detection into two (second row), or stray detections off road (third row), *MonoLayout* (Ours) produces crisp object boundaries while respecting vehicle and road geometries.

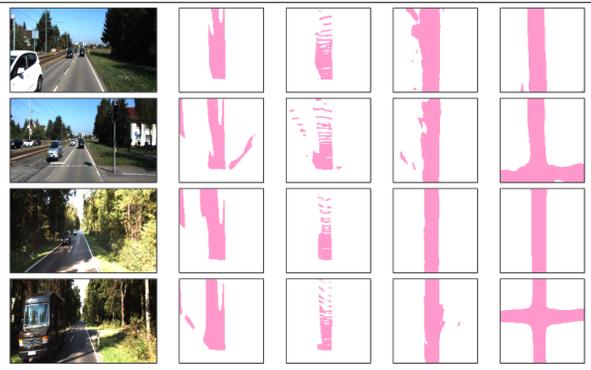


Figure 5: **Impact of sensor fusion:** *Col 1:* Input images. *Col 2:* Using per-frame monocular depth, *Col 3:* Using per-frame lidar depth, *Col 4:* Sensor-fused monocular depth, *Col 5:* Sensor-fused lidar depth. (Refer to Table 2 and Sec. 4.6 for more details)

also gain significant performance boosts. This is even more pronounced on Argoverse [5], as shown in Table 3. In case of vehicle occupancy estimation, while using a discriminator does not translate to quantitative performance gains, it often results in qualitatively sharper, aesthetic estimates as seen in Fig 7.

Timing analysis

We also show the computation test time of our method as compared to other baselines in Table 6. Unlike Schuster *et al.* [34], our network does not require discrimi-

Supervision	<i>MonoOccupancy-ext</i>	<i>MonoLayout-static</i> (Ours)
Per-frame monocular depth	56.16	58.87
Sensor-fused monocular depth	64.81	66.71
Per-frame lidar	44.29	48.29
Sensor-fused lidar	71.67	73.86

Table 2: **Monocular depth:** If lidar data is unavailable, we leverage self-supervised monocular depth estimation to generate training data for *MonoLayout-static* and achieve reasonable static layout estimation (rows 1-2). Although performance is inferior to the case when lidar is available (rows 3-4), this is not unexpected. **Sensor fusion:** Regardless of the modality of depth information, sensor-fusing depth estimates over a window of 40 frames dramatically improves performance (row 2, row 4).

Dataset	<i>MonoLayout-no-disc</i>			<i>MonoLayout</i>		
	Road	Vehicle (mIoU)	Vehicle (mAP)	Road	Vehicle (mIoU)	Vehicle (mAP)
KITTI Raw	70.95	-	-	73.86	-	-
KITTI Object	-	26.25	37.66	-	25.47	41.52
Argoverse	51.66	32.84	44.07	58.33	32.06	48.31

Table 3: **Effect of discriminator:** Adding a discriminator clearly translates to an accuracy boost in static (road) layout estimation. For vehicle occupancy estimation, while a quantitative boost is not perceived, the generated layouts are sharper, and aesthetic as opposed to when not using the discriminator (*c.f.* Fig. 7)

nator to be used during inference time. It achieves real time inference rate of approx. 32 Hz for an input size $3 \times 512 \times 512$ and an output size $2 \times 128 \times 128$



Figure 6: **Amodal scene layout estimation** on the Argoverse [5] dataset. The dataset comprises multiple challenging scenarios, with low illumination, large number of vehicles. *MonoLayout* is accurately able to produce sharp estimates of vehicles and road layouts. (Sidewalks are not predicted here, as they aren’t annotated in Argoverse).

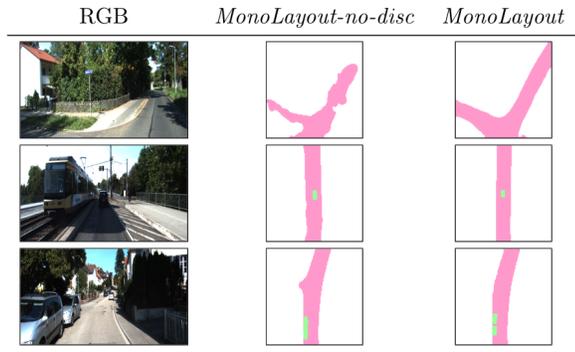


Figure 7: **Effect of adversarial learning**: As can be clearly seen here, the discriminators help enhance both the static (road) layout estimation (top and middle rows), as well as produce sharper vehicle boundaries (bottom row). While this translates to performance gains in static layout estimation (*c.f.* Table 3), the gains in dynamic layout estimation are more cosmetic in nature.

on an NVIDIA GeForce GTX 1080Ti GPU. Note that in *MonoLayout* the static and dynamic decoders are executed in parallel, maintaining comparable runtime. *MonoLayout* is an order of magnitude faster than previous methods, making it more attractive for on-road implementations.

5. Discussion and conclusions

This paper proposed *MonoLayout*, a versatile deep network architecture capable of estimating the amodal layout of a complex urban driving scene in real-time.

Method	Parameters	Computation Time
Mono3D [6]	>> 20 M	0.24 fps
OFT [30]	23.5 M	< 5 fps
MonoOccupancy [24]	27.5 M	15 fps
Schulter et al. [34]	>> 20 M	< 3 fps
<i>MonoLayout</i> (Ours)	19.6M	32 fps

Table 4: A comparative study of inference time of various methods. *MonoLayout* is about 2x faster and significantly more accurate compared to prior art. (*c.f.* Table 1).

In the appendix, we show several additional results, including extended ablations, and applications to multi-object tracking and trajectory forecasting. A promising avenue for future research is the generalization of *MonoLayout* to unseen scenarios, as well as incorporating temporal information to improve performance.

Acknowledgements

We thank Nitesh Gupta, Aryan Sakaria, Chinmay Shirore, Ritwik Agarwal and Sachin Sharma for manually annotating several sequences from the KITTI RAW [10] dataset.

References

- [1] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall. Semantickitti: A dataset for semantic scene understanding of lidar sequences. In *ICCV*, 2019. 6
- [2] J. Beltrán, C. Guindel, F. M. Moreno, D. Cruzado, F. Garcia, and A. De La Escalera. Birdnet: a 3d object detection framework from lidar information. In *ITSC*, 2018. 2
- [3] E. Bochinski, V. Eiselein, and T. Sikora. High-speed tracking-by-detection without using image information. In *International Workshop on Traffic and Street Surveillance for Safety and Security at IEEE AVSS 2017*, 2017. 13, 14
- [4] J.-R. Chang and Y.-S. Chen. Pyramid stereo matching network. In *CVPR*, 2018. 11, 12
- [5] M.-F. Chang, J. Lambert, P. Sangkloy, J. Singh, S. Bak, A. Hartnett, D. Wang, P. Carr, S. Lucey, D. Ramanan, et al. Argoverse: 3d tracking and forecasting with rich maps. In *CVPR*, 2019. 1, 2, 5, 7, 8, 9, 16
- [6] X. Chen, K. Kundu, Z. Zhang, H. Ma, S. Fidler, and R. Urtasun. Monocular 3d object detection for autonomous driving. In *CVPR*, 2016. 2, 5, 6, 7, 8, 9, 11, 15
- [7] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia. Multi-view 3d object detection network for autonomous driving. In *CVPR*, 2017. 2
- [8] D. Cireřan, U. Meier, and J. Schmidhuber. Multi-column deep neural networks for image classification. *arXiv preprint*, 2012. 1

- [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 4
- [10] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012. 1, 2, 5, 7, 8, 9, 11, 12, 15
- [11] R. Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 1
- [12] C. Godard, O. Mac Aodha, M. Firman, and G. Brostow. Digging into self-supervised monocular depth estimation. *arXiv preprint*, 2018. 2, 5, 7, 12
- [13] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems 27*. 2014. 4
- [14] D. Hall, F. Dayoub, J. Skinner, H. Zhang, D. Miller, P. Corke, G. Carneiro, A. Angelova, and N. Sünderhauf. Probabilistic object detection: Definition and evaluation. *arXiv preprint arXiv:1811.10800*, 2018. 11
- [15] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 1
- [16] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*, 2015. 1
- [17] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 4
- [18] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 11, 13
- [19] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 11
- [20] J. Ku, M. Mozifian, J. Lee, A. Harakeh, and S. L. Waslander. Joint 3d proposal generation and object detection from view aggregation. In *IROS*, 2018. 2, 12
- [21] B. Li, W. Ouyang, L. Sheng, X. Zeng, and X. Wang. Gs3d: An efficient 3d object detection framework for autonomous driving. In *CVPR*, 2019. 2
- [22] M. Liang, B. Yang, S. Wang, and R. Urtasun. Deep continuous fusion for multi-sensor 3d object detection. In *ECCV*, 2018. 2
- [23] S. Lin, R. Clark, R. Birke, N. Trigoni, and S. Roberts. Wise-ale: Wide sample estimator for aggregate latent embedding. 2019. 7
- [24] C. Lu, M. J. G. van de Molengraft, and G. Dubbelman. Monocular semantic occupancy grid mapping with convolutional variational encoder-decoder networks. *IEEE Robotics and Automation Letters*, 2019. 2, 5, 6, 7, 8, 9, 12, 15
- [25] A. Mousavian, D. Anguelov, J. Flynn, and J. Kosecka. 3d bounding box estimation using deep learning and geometry. In *CVPR*, 2017. 2
- [26] OpenStreetMap contributors. Planet dump retrieved from <https://planet.osm.org>. <https://www.openstreetmap.org>, 2017. 2, 4, 13
- [27] A. Paszke, A. Chaurasia, S. Kim, and E. Cukurciello. Enet: A deep neural network architecture for real-time semantic segmentation. *arXiv preprint arXiv:1606.02147*, 2016. 12, 13
- [28] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 11, 13
- [29] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 1
- [30] T. Roddick, A. Kendall, and R. Cipolla. Orthographic feature transform for monocular 3d object detection. *arXiv preprint*, 2018. 2, 5, 6, 7, 8, 9, 11, 12, 15
- [31] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, 2015. 13
- [32] S. Rota Bulò, L. Porzi, and P. Kotschieder. In-place activated batchnorm for memory-optimized training of dnn. In *CVPR*, 2018. 1, 5
- [33] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. In *Advances in neural information processing systems*, 2016. 4
- [34] S. Schuster, M. Zhai, N. Jacobs, and M. Chandraker. Learning to look around objects for top-view representations of outdoor scenes. In *ECCV*, 2018. 2, 4, 5, 6, 8, 9, 12
- [35] S. Shi, X. Wang, and H. Li. Pointcnn: 3d object proposal generation and detection from point cloud. In *CVPR*, 2019. 2, 11, 12
- [36] S. Srivastava, F. Jurie, and G. Sharma. Learning 2d to 3d lifting for object detection in 3d for autonomous vehicles. *arXiv preprint*, 2019. 2
- [37] Y. Wang, W.-L. Chao, D. Garg, B. Hariharan, M. Campbell, and K. Q. Weinberger. Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In *CVPR*, 2019. 12
- [38] Z. Wang, B. Liu, S. Schuster, and M. Chandraker. A parametric top-view representation of complex road scenes. In *CVPR*, 2019. 2
- [39] B. Yang, W. Luo, and R. Urtasun. Pixor: Real-time 3d object detection from point clouds. In *CVPR*, 2018. 2
- [40] Y. You, Y. Wang, W.-L. Chao, D. Garg, G. Pleiss, B. Hariharan, M. Campbell, and K. Q. Weinberger. Pseudo-lidar++: Accurate depth for 3d object detection in autonomous driving. *arXiv preprint*, 2019. 11
- [41] Q. Yu, Y. Yang, F. Liu, Y.-Z. Song, T. Xiang, and T. M. Hospedales. Sketch-a-net: A deep neural network that beats humans. *International journal of computer vision*, 2017. 1

A. Implementation Details

In this section, we describe the network architecture and training procedure in greater detail.

A.1. Network Architecture

MonoLayout comprises the following four blocks: a feature extractor, a static layout decoder, dynamic layout decoder, and two discriminators.

Feature extractor

Our feature extractor is built on top of a ResNet-18 encoder⁵. The network usually takes in RGB images of size $3 \times 512 \times 512$ as input, and produces a $512 \times 32 \times 32$ feature map as output. In particular, we use the ResNet-18 architecture without bottleneck layers. (Bottleneck layers are absent, to ensure fair comparison with OFT [30]). This extracted feature map is what we refer to as the *shared context*.

Layout Decoders

We use two parallel decoders with identical architectures to estimate the static and dynamic layouts. The decoders consists 2 convolution layers and take as input the $512 \times 32 \times 32$ shared context. The first convolution block maps this shared context to a $128 \times 16 \times 16$ feature map, and the second convolution block maps the output of the first block to another $128 \times 8 \times 8$ feature map.

At this point, 4 deconvolution (transposed convolution) blocks are applied on top of this feature map. Each block increases spatial resolution by a factor of 2, and decreases the number of channels to 64, 32, 16, and O respectively, where O is the number of channels in the output feature map ($O \in \{1, 2\}$ for the static layout decoder, and $O = 1$ for the dynamic layout decoder). This results in an output feature map of size $O \times 128 \times 128$. We also apply a spatial dropout (ratio of 0.4) to the penultimate layer, to impose stochastic regularization. The output 128×128 grid corresponds to a rectangular region of area $40m \times 40m$ on the ground plane.

Discriminators

The discriminator architecture is inspired by Pix2Pix [18]. We found the patch based regularization in Pix2Pix to be much better than a standard DC-GAN [28]. So, we use patch-level discriminators

⁵We also tried other feature extraction / image-image architectures such as UNet and ENet, but found them to be far inferior in practice.

that contain four convolution layers (kernel size 3×3 , stride 2), that output an 8×8 feature map. This feature map is passed through a *tanh* nonlinearity and used for patch discrimination.

A.2. Training details

We train *MonoLayout* with a batch size of 16 for 200 epochs using the Adam [19] optimizer with initial learning rate $5e - 5$. The input image is reshaped to $3 \times 512 \times 512$ and further augmented to make our model more robust. Some of the augmentation techniques we use include random horizontal flipping and color jittering.

A.3. Metrics used

Road layout estimation

To evaluate estimated road layouts, we use intersection-over-union (IoU) as our primary metric. We split IoU evaluation into two parts and measure IoU for the entire static scene, as well as IoU for occluded regions (i.e., regions of the road that are occluded in the image and were hallucinated by *MonoLayout*).

Vehicle occupancy estimation

While most approaches to vehicle detection evaluate only mean Average Precision (mAP), it has been shown to be a grossly inaccurate measure of how tight a bounding box is [14]. We hence adopt mean Intersection-over-Union (mIoU) as our primary basis of evaluation. To ensure a fair comparison with prior art, we also report mAP⁶. (Note that, since we are evaluating object detection in bird’s eye view, and not in 3D, we use the mIoU and mAP, as is common practice [6, 10]). This choice of metrics is based on the fact that *MonoLayout* is not an object “*detection*” approach; it is rather an object occupancy estimation approach, which calls for mIoU and mAP evaluation. We extend the other object “*detection*” approaches (such as pseudo-lidar based approaches [30, 35, 40]) to the occupancy estimation setting, for a fair comparison.

B. Timing analysis

We also show the computation test time of our method as compared to similar methods in Table 6. Our network does not require discriminator to be used during inference time. It achieves real time inference rate of approx. 30 Hz for an input image with a resolution of $512 * 512$ pixels and an output map with $128 * 128$ occupancy grid cells using a Nvidia GeForce GTX

⁶We outperform existing methods under both these metrics

Method	Vehicle (mAP)	Vehicle (mIoU)	Road (mIoU)	Frame Rate
ENet + Pseudo lidar input(Monodepth2)	0.37	0.24	0.62	12.34 fps
PointRCNN + Pseudo lidar input(Monodepth2)	0.43	0.26	-	5.76 fps
<i>MonoLayout</i> (Ours)	0.41	0.26	0.80	32 fps
AVOD + Pseudo lidar input(PSMNet) (Stereo)	0.59	0.43	-	< 1.85 fps

Table 5: **Comparison with Pseudo-lidar [37]**: We also evaluate *MonoLayout* against several variants of pseudo-lidar [37] approaches. While the usage of increasingly heavy processing blocks for pseudo-lidar variants improves accuracy, it drastically increases computation time. On the other hand, *MonoLayout* offers—by far—the best mix of accuracy and runtime for real-time applications. Also note that the comparison with AVOD + Pseudo-lidar (PSMNet) is unfair, since it uses stereo disparities.

1080Ti GPU. The code for [34] is not publicly available, and the computation time is based on the PSMNet [4] backbone they use. Here again the proposed method is almost an order faster than previous methods making it more attractive for on-road implementations.

Method	Parameters	Computation Time
OFT [30]	24.5 M	2 fps
MonoOccupancy [24]	27.5 M	15 fps
Schulter et al. [34]	>> 20 M	< 3 fps
<i>MonoLayout</i> (Ours)	19.6M	32 fps

Table 6: A comparative study of test computation time on NVIDIA GeForce GTX 1080Ti GPU for different methods on the images of KITTI [10] RAW dataset.

C. Comparison with pseudo-lidar

There is another recent set of approaches to object detection in bird’s eye view—*pseudo-lidar* approaches [37]. At the core of these approaches lies the idea that, since lidar object detection works exceedingly well, monocular images can be mapped to (*pseudo*) lidar-like maps in bird’s eye view, and object detection networks tailored to lidar bird’s eye view maps can readily be applied to this setting. Such approaches are primarily geared towards detecting objects in lidar-like maps.

MonoLayout, on the other hand, intends to estimate an *amodal* scene layout, and to do so, it must reason not only about vehicles, but also about the static scene layout. Table 5 compares *MonoLayout* with a set of *pseudo-lidar* approaches, in terms of vehicle occupancy estimation and road layout estimation. Specifically, we evaluate the following pseudo-lidar based methods.

1. *ENet + Pseudo-lidar input (Monodepth2)*: Uses an ENet [27]-style encoder-decoder architecture that uses Monodepth2 [12] to get monocular depth estimates.

2. *PointRCNN + Pseudo-lidar input (Monodepth2)*: Uses a PointRCNN [35] architecture (a two-stage object detector comprising a region proposal network, and a classification network) to detect vehicles in bird’s eye view.
3. *AVOD + Pseudo-lidar input (PSMNet)*: A **stereo, supervised** method. Uses the aggregated view object detector AVOD [?] and pseudo-lidar input computed from a disparity estimation network (PSMNet [4]).

The comparison is shown in Table 5. The PointRCNN [35] and AVOD [20] are tailored specifically for object detection, and hence *cannot be repurposed to estimate road layouts*. However, the ENet [27] architecture can, and we trained it for the task of road layout estimation. We observe that, among all approaches, *MonoLayout* is the fastest (about an order of magnitude speedup over pseudo-lidar methods). Furthermore, the accuracy is competitive, if not greater, compared to pseudo-lidar based approaches.

We also evaluate against a stereo pseudo-lidar baseline (AVOD + pseudolidar PSMNet [4]). By virtue of using stereo images, and being supervised on the KITTI Object dataset ⁷, achieves superior performance. However, the comparison is unfair, and is provided only for a reference, to enable progress in amodal layout estimation from a monocular camera.

Another shortcoming of pseudolidar-style approaches is that, it is not possible to learn *visual* (i.e., image intensity based) features that are extremely useful in road layout estimation).

D. Application: Trajectory forecasting

One of the use-cases of *MonoLayout* is to forecast future trajectories from the estimated *amodal* scene layout.

⁷Monodepth2 [12] is unsupervised, and has not been finetuned on the KITTI Object dataset

We demonstrate accurate trajectory forecasting performance by training a Convolutional LSTM that operates over the estimated layouts from *MonoLayout*. Specifically, we adopt an encoder-decoder structure similar to ENet [27], but add a convolutional LSTM between the encoder and the decoder. We also add a convolutional LSTM over each of the skip connections present in ENet.

We *pre-condition* the trajectory forecasting network for 1 second, by feeding it images, and then predict future trajectories for the next 3 seconds. Note that, when predicting future trajectories, no images are fed to the forecasting network. Rather, the network operates in an autoregressive manner, by producing an output trajectory estimate, and using this estimate as the subsequent input. We also tried predicting static scene layouts, to forecast future static scenes, but without any success, owing to the high variability in static scene layouts. The resultant model, called *MonoLayout-forecast*, works in real-time, and accurately forecasts the future trajectory of a moving vehicle, as shown in Fig. 8. As with the layout estimation task, the forecast trajectories are confined to a 128×128 grid, equivalent to a $40m \times 40m$ square area in front of the ego-vehicle.

E. Application: Multi-object tracking

Further, we propose an extension to generate accurate *tracks* of multiple moving vehicles by leveraging the vehicle occupancy estimation capabilities of *MonoLayout*. We also construct a strong baseline multi-object tracker using the open-source implementation from IoU-Tracker [3] as a reference. We term this baseline *BEVTracker*. Specifically, we use disparity estimates from stereo cameras, semantic and instance segmentation labels, to segment and identify unique cars in bird’s eye view. We then run IoU-Tracker [3] on these estimates.

We demonstrate in Table 7 that *MonoLayout* outperforms BEVTracker [3], without access to any such specialized information (disparity, semantics, instances). Instead, we run a traditional OpenCV blob detector on *MonoLayout* vehicle occupancy estimates, and use the estimated instances to obtain the coordinates of the center of the vehicle. We then use a *maximum-overlap* data association strategy across time, using intersection-over-union to measure overlap. We run our approach on all 24 sequences of the KITTI Tracking (train) benchmark, and present the results in Table 7.

F. More qualitative results

Additional qualitative results of joint static and dynamic scene layout estimation are presented in Fig. 9 and Fig. 10.

G. Shortcomings

Despite outperforming several state-of-the-art approaches and achieving real-time performance, *MonoLayout* also suffers a few shortcomings. In this section, we discuss some failure cases of *MonoLayout* and also some negative results. Please watch the supplementary video for additional results.

G.1. Failure cases

Fig. 11 shows a few scenarios in which *MonoLayout* fails to produce accurate layout estimates. Recall that *MonoLayout* uses adversarial feature learning to estimate plausible road layouts. We leverage OpenStreetMap [26] to randomly extract road patches to use as the true data distribution. However, no such data is available for *sidewalks*, and this results in a few artifacts.

As shown in the bottom row of Fig. 11, high-dynamic range and shadows coerce the network into predicting *protrusions* along these directions. Also, as shown in the top row, sometimes, sidewalk predictions are not coherent with road predictions. In the bottom row, we show failure cases for vehicle occupancy estimation.

G.2. Negative results

Produced below is a list of experiments that the authors tried, but were unsuccessful. We hope this will expedite progress in this nascent field, by saving fellow researchers significant amounts of time.

These did not work!

- Using a single encoder-decoder network to estimate both static and dynamic scene layouts.
- Using an ENet [27] or a UNet [31] architecture as opposed to the ResNet-18 encoder and customized decoder we employed.
- Using a DCGAN [28] as opposed to patch-discriminators [18].
- Employing a variational autoencoder-style latent code between the encoder and decoder (to allow for sampling)

Method	Mean Error in Z (m)	Mean Error in X (m)	Mean L2 error (m)
BEVTracker (Open-source [3], enhanced with stereo, semantics, etc.)	1.08	0.51	1.27
<i>MonoLayout</i> (Ours)	0.23	0.47	0.58

Table 7: **Multi-object tracking performance:** We show the performance of *MonoLayout* on a multi-object tracking task, in bird’s eye view. The comparison is unfair as *BEVTracker* uses strictly more information (disparities, semantic segmentation, instance segmentation). However, *MonoLayout* does not employ any such privileged information, and as such uses OpenCV blob detection to identify instances, and a *maximum-IoU-overlap* data association framework. These estimates are obtained over 24 sequences of the KITTI tracking benchmark. (Caveat: We only evaluate tracking accuracies for cars within a $40m \times 40m$ square region around the ego-vehicle, as *MonoLayout* estimates are confined to this region).

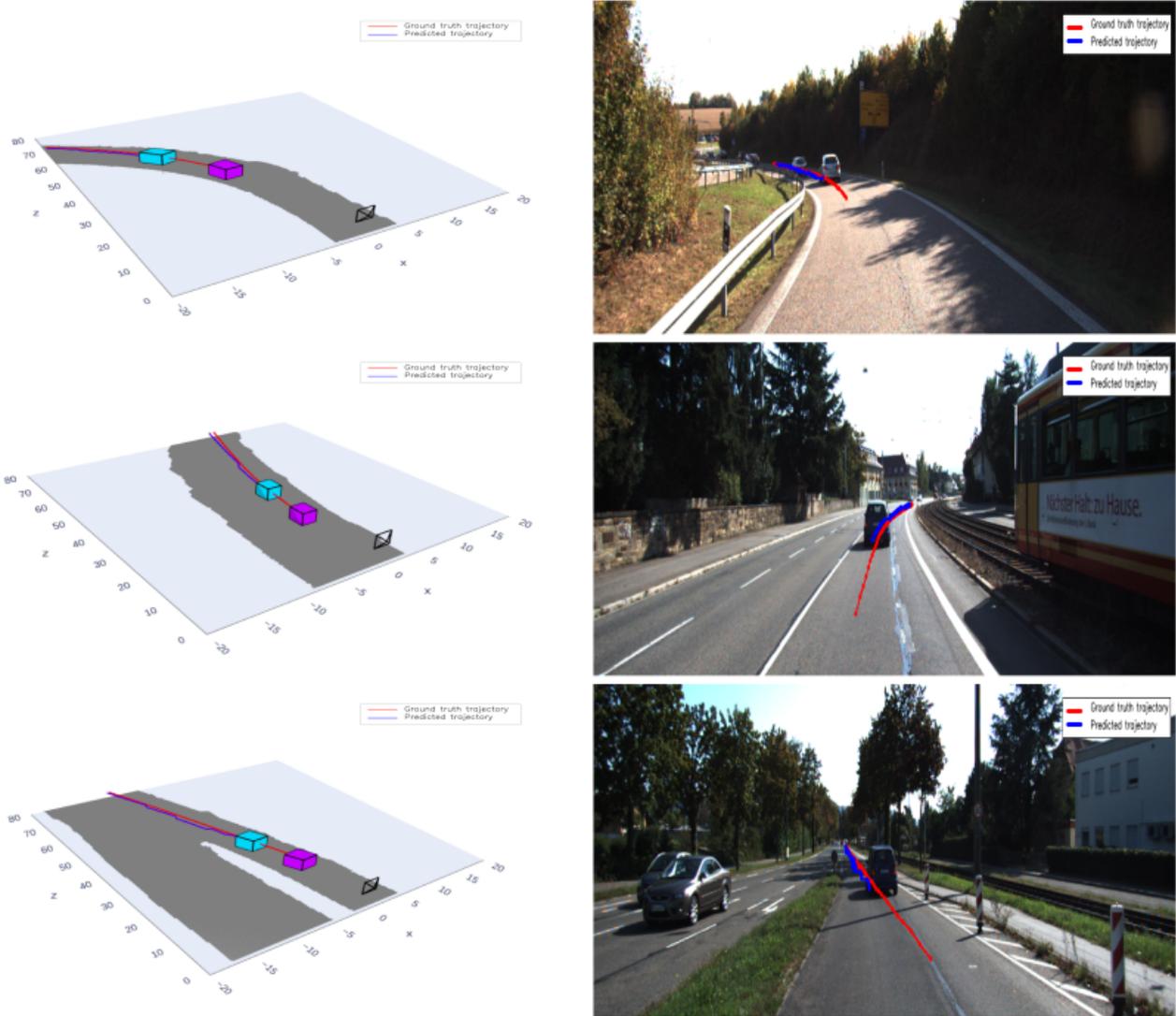


Figure 8: **Trajectory forecasting:** *MonoLayout-forecast* accurately estimates future trajectories of moving vehicles. (Left): In each figure, the magenta cuboid shows the initial position of the vehicle. *MonoLayout-forecast* is pre-conditioned for 1 seconds, by observing the vehicle, at which point (cyan cuboid) it starts forecasting future trajectories (shown in blue). The ground-truth trajectory is shown in red, for comparison. (Right): Trajectories visualized in image space. Notice how *MonoLayout-forecast* is able to forecast trajectories accurately despite the presence of moving obstacles (top row), turns (middle row), and merging traffic (bottom row).

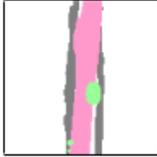
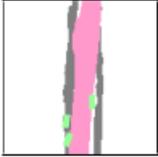
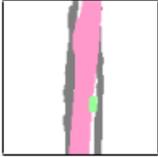
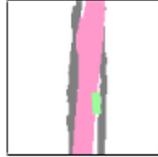
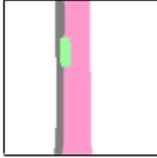
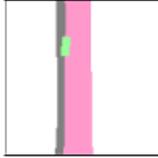
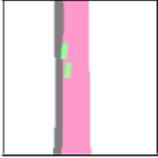
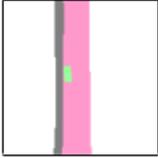
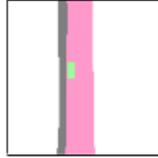
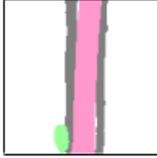
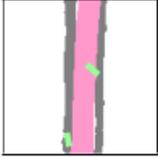
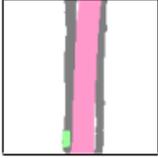
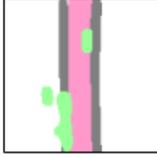
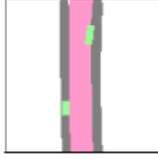
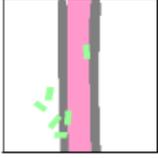
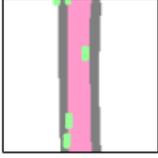
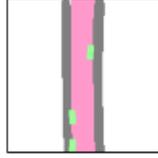
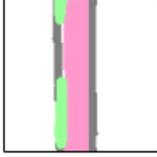
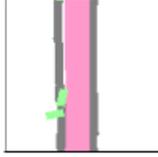
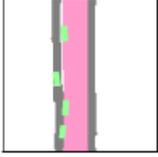
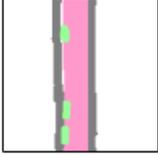
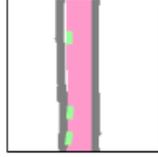
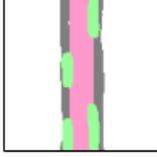
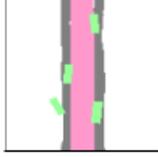
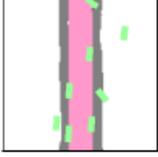
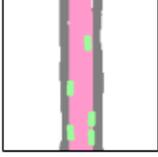
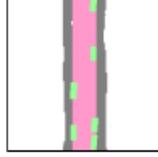
<i>RGB</i>	<i>Lu et al. [11]</i>	<i>Mono3D [4]</i>	<i>OFT [15]</i>	<i>MonoLayout</i>	<i>Ground Truth</i>
					
					
					
					
					
					

Figure 9: **Dynamic layout estimation on KITTI [10]**: Additional qualitative results for dynamic scene layout estimation on the KITTI [10] dataset. From left to right, the column corresponds to the input image, MonoOccupancy [24], Mono3D [6], OFT [30], *MonoLayout* (Ours), and ground-truth respectively. *MonoLayout* (Ours) produces crisp object boundaries while respecting vehicle and road geometries.

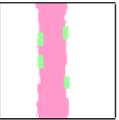
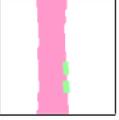
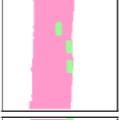
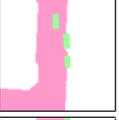
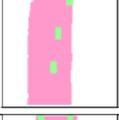
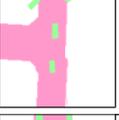
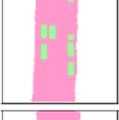
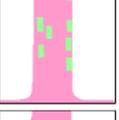
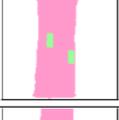
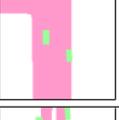
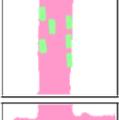
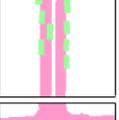
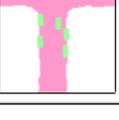
	RGB	<i>MonoLayout</i>	Ground-truth
Night			
			
Sunny			
			
Cloudy			
			
Overcast			
			

Figure 10: **Qualitative results on Argoverse:** Additional qualitative results on the Argoverse [5] dataset (road shown in pink, vehicles shown in green). *MonoLayout* (center column) uses both static and dynamic layout discriminators and produces sharp estimates, and is robust to varying weather conditions, high dynamic range (HDR), and shadows.

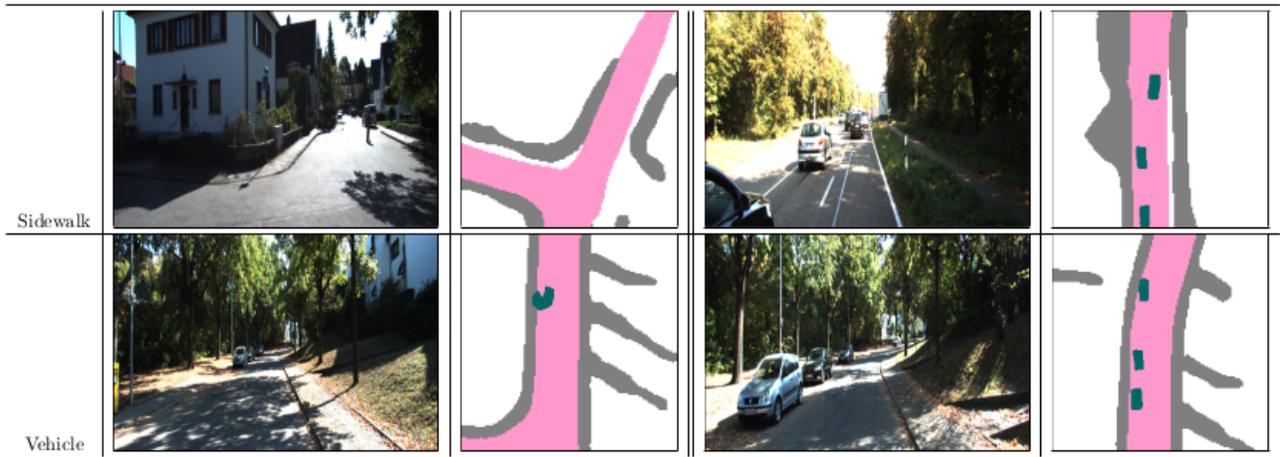


Figure 11: **Failure cases:** This figure highlights a few cases in which *MonoLayout* produces incorrect layout estimates. Adverse lighting conditions and sharp turns, in some cases effect sidewalk estimation accuracy(Row 1). Also, multiple near-by vehicles in an image get merged into a single estimate, at times. (As shown in Row 2, only when the ego-vehicle gets close to the two cars parked in close vicinity to each other, the model is able to distinguish the two cars.)