

Joint Prediction of Depths, Normals and Surface Curvature from RGB Images using CNNs

Thanuja Dharmasiri* Andrew Spek* Tom Drummond

Abstract—Understanding the 3D structure of a scene is of vital importance, when it comes to developing fully autonomous robots. To this end, we present a novel deep learning based framework that estimates depth, surface normals and surface curvature by only using a single RGB image. To the best of our knowledge this is the first work to estimate surface curvature from colour using a machine learning approach. Additionally, we demonstrate that by tuning the network to infer well designed features, such as surface curvature, we can achieve improved performance at estimating depth and normals. This indicates that network guidance is still a useful aspect of designing and training a neural network. We run extensive experiments where the network is trained to infer different tasks while the model capacity is kept constant resulting in different feature maps based on the tasks at hand. We outperform the previous state-of-the-art benchmarks which jointly estimate depths and surface normals while predicting surface curvature in parallel.

I. INTRODUCTION

Extracting information from raw data is a well studied problem in robotics. A visual image is one such form of raw data and has been widely used in the community to tackle a range of problems including image segmentation [1], localization and mapping [2], visual servoing [3] etc. and there exist a continuous stream of research which look at maximizing the amount of information extracted. In this paper we show that we can estimate geometric quantities such as surface curvature using only RGB images as input. To our knowledge this is the first work to demonstrate such a capability.

Surface Curvature is an important geometric surface feature, that indicates the rate at which the direction of the normals change on the surface at any particular point. It has been shown to be particularly useful for the task of segmentation on range image and 3D data [4, 5, 6, 7]. A key challenge in accurately estimating surface curvature is its sensitivity to noise in the input data, as it is a second order surface derivative, it is affected quadratically by noise. Previous works have shown that neural networks can be used to provide accurate geometric estimates from just single RGB images [8, 9, 10, 11], including estimating depth and normals. In this work we extend our network to estimate principal surface curvatures as well as depth and normals and demonstrate that we can accurately perform this task from a single RGB image.

Contrary to the popular belief that hand-engineered features are inferior compared to learnt features, we argue that well designed features combined with machine

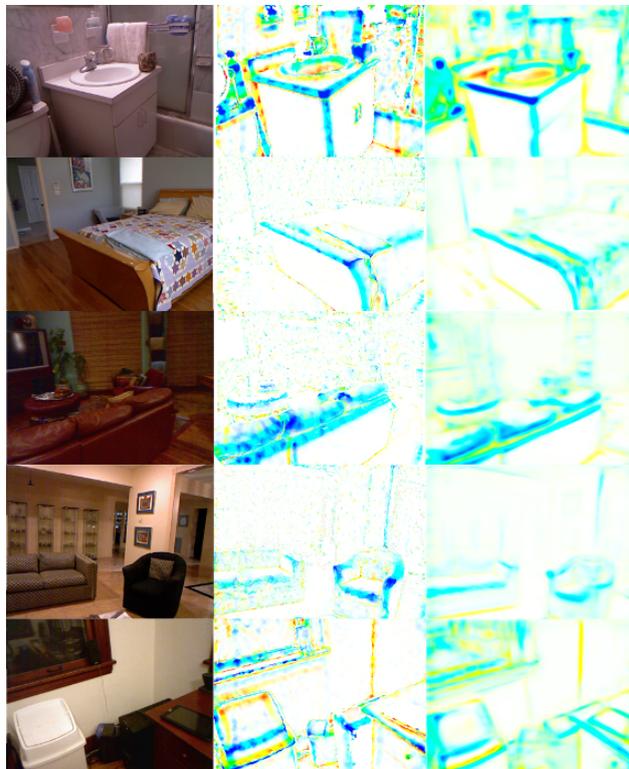


Fig. 1: A selection of curvature predictions made by our system. The left column shows the corresponding RGB image from the NYUv2 test dataset which was used as the only source of information to estimate curvature. The middle column shows the ground truth curvature computed using the depth data and the right image shows the prediction of our network. The Positive curvatures are shown in blue, Negatives in red, Saddles in green and Planes in white.

learnt representations provide improved performance. It should be stressed that the features designed are not hand calculated by us, but rather predicted by the network itself as part of the inference pipeline. More concretely, we *inform* the network in order to accurately estimate a single quantity such as depth, normals or curvature the network should learn an internal representation of the other two quantities. We demonstrate this by estimating surface curvature, surface normals and depth in a multi task learning framework which gives us superior results compared to training them as individual tasks. We employ a two-stage learning process where coarse level predictions of all three quantities are used as feature maps for the finer layers. Our work is similar to [8] in that sense, as Eigen et al. also estimated three quantities (depth, surface normals and semantic labels) using a single network. The fundamental difference between ours and their approach

*The authors contributed equally

This work was supported by the Australian Research Council Centre of Excellence for Robotic Vision (project number CE14010006).

The authors are with the Faculty of Electrical and Computer Systems Engineering, Monash University, Australia. [firstname].[lastname]@monash.edu

is that the three quantities we try to estimate are more tightly coupled at a primitive level where as the semantic labels, although clearly related, should be considered a higher order quantity compared to its counterparts depths and normals. We show both quantitatively and qualitatively that we are able to achieve better results on depth and surface normals on the NYUv2 dataset [12] by estimating a view point invariant quantity (surface curvature) jointly with depth and normals. We believe that robotic applications that revolve around segmentation tasks such as the Amazon Picking Challenge could benefit from our approach. Our contributions are as follows :

- A novel technique to estimate surface curvature of objects using purely RGB images.(Method: Section IV-C, Results: Section VI-C)
- A framework which predicts depth,surface normals and curvature jointly.(Method: Section IV, Results: Section VI)
- Demonstrate that joint training can improve the accuracy of all three tasks while keeping the model capacity fixed (Method: Section V-C, Results: Tables I, II, III.

II. RELATED WORK

In this section we review existing work in the literature that is related to this paper and in turn inspired the ideas presented. We take a look at traditional approaches used to compute surface curvature from raw depth data, then we summarize how deep learning has been used to predict information from images and finally, we discuss how the problem of learning multiple tasks in a single platform was performed using deep learning.

Surface Curvature Estimation Surface Curvature estimation is a well explored topic in robotics and computer vision. It has been shown to be useful for object segmentation [4, 5, 7, 6] in depth scans and RGB-D imagery. There are several popular approaches to estimating the surface curvature. One technique is to simply twice-differentiate the surface [4, 6], but this can lead to a high sensitivity to noise in the data and generally requires removal or rejection of surface outliers. Another technique is to estimate the surface curvature from a locally connected surface mesh based on the change in adjacent facet normal angles [7, 13, 14]. This method is predominantly used for computer graphics and low-noise data as it operates on a small neighbourhood of facets. Yet another technique is to use locally fit surface quadrics and directly extract the principal curvatures from their parameterization [5, 15], which has been shown to be robust to noisy data and fast enough to be computed in real-time [16]. In this work we use the approach in [16], to compute surface curvature and surface normals from the training data sourced from the NYUv2 dataset [12] as they have shown it performs well on range image data of the type present in the dataset.

Predicting Information using Deep Neural Networks Convolutional Neural Networks (CNNs) have been very effectively applied to a range of robotic and vision tasks including grasp pose detection [17, 18], image classification [19, 20], semantic segmentation [1], depth estimation [8, 9, 10, 11], surface normal estimation [21, 22, 8]. Our

work is more closely related to the latter two tasks as we demonstrate surface curvature can be predicted using RGB images as the only input. We began this work by using the VGG architecture [23] as a starting point to predict surface curvature in a standalone network and extended it to estimate depth and surface normals as well, in the one network.

Prior to the resurgence of neural networks depth was either computed using a Simultaneous Localisation and Mapping (SLAM) system [24, 2] or directly obtained from a range sensor such as expensive LIDAR, stereo rigs, Time of Flight (ToF) sensors or structured light sensors[25] (Microsoft Kinect). In a robotic context going from data to information as efficiently as possible is vital, and predicting quantities from a single image is a step in that direction. Saxena et al. in [10] used a supervised learning approach that combines local and global image features by using a Markov Random Field (MRF). The idea of using both global and local features was further investigated by Eigen et al. [26] using the AlexNet [19] architecture in a multi-scale scheme. Liu et al. [9] proposed to combine graphical models in the form of a Conditional Random Field (CRF) with a CNN to improve the accuracy of monocular depth estimation. More recently, Laina et al. [11] proposed to use a far superior fully convolutional residual architecture and obtain state-of-the-art results in single image depth estimation.

Data driven single image surface normal estimation was first tackled by Fouhey et al. in [27]. They used a SVM based detector followed by an iterative optimization scheme to extract geometrically informative primitives. Ladicky et al. proposed to use image cues of pixel-wise and segment based methods to generate a feature representation that can estimate surface normals in a boosting framework [28]. A ConvNet approach to estimating surface normals in global and local scales while incorporating numerous constraints such as room layout and edge labels was taken by Wang et al. [21]. Recently, Bansal et al. [22] showed that by combining hierarchy of features from different levels of activations in a skip-network architecture that you could generate much finer predictions for surface normals achieving state of the art results.

Learning Multiple Tasks In one of the earliest works in this area Caruana et al. showed in [29] that by learning related tasks in parallel, the performance of all tasks could be improved, which is consistent with our findings. Multiple tasks were learned in the form of material classification and defect detection in railway fasteners in [30] where they used Deep CNN based multi task learning for railway track inspection. They were able to show the adaptability of the multi task learning platform by using different training batch sizes (due to availability of data). In our case, all three tasks were trained with the same batch size as training data for the derived quantities (normals and curvature) were computed from depth. Multi task learning algorithms were also used to perform head pose estimation [31], web search ranking [32], face verification [33] etc. Li et al. in their work *Learning Without Forgetting* [34] demonstrated that in the presence of a model trained on one task, it can be fine-tuned to perform better on a new

task while not hindering the performance of the previous task by only using training data of the new task. However, as we have access to training data for all three tasks we train the prediction stacks jointly in order to achieve superior performance compared to fine-tuning.

III. MODEL ARCHITECTURE

The functionalities of the model can be divided into 3 main sections. Firstly, there is a set of convolutional layers based on the VGG16[23] architecture corresponding to *feature extraction* which is followed by 2 fully connected layers which can see the whole image in their field views. Secondly, we have a stack of convolutional layers corresponding to *coarse level predictions* of depths, normals and curvatures and finally, a set of convolutional layers which predicts the three quantities at a *finer resolution*.

It is worth mentioning that all the convolutional layers in the coarse and fine level prediction stacks perform 5x5 convolutions with a stride of 1 and a pad of 2. Therefore, the input resolution is preserved at the output. There is an explicit up-sampling layer which up samples the coarse level prediction from 74x55 resolution to 147x109 and this is maintained throughout, by the final convolutional stack as shown in Figure 2. At the end of each scale, there are individual solvers for each of the training tasks which essentially compute the loss and initiate the backward propagation.

Although the overall architecture is as explained above, in order to make sure the model capacity is kept constant and the contribution of each new task is indeed improving the previous tasks we make several changes during training which is explained in the Section V-C .

IV. TASKS

A. Depth

We use the raw depth data distribution given by the NYUv2 dataset [12] for training based on the official train and test scene split (that is 249 training scenes and 219 test scenes). Similar to previous approaches we train our network to estimate depth at several scales. The loss function used for calculating the error in the depth estimation includes an Euclidean loss term, a scale invariant term and a gradient term which compares the local rate of change of the predicted and ground truth depth values spatially. However, we do not the fix the coarse level feature stacks while training the fine level features (depth, normals and curvature) rather jointly train both stacks together as opposed to [8]. For the benefit of the reader we include the loss function in the following equation which is the same loss criterion employed by [8].

$$L(D, D^*) = \sum_{i=1}^n d_i^2 - \frac{1}{2n^2} \left(\sum_{i=1}^n d_i \right)^2 + \frac{1}{n} \sum_{i=1}^n ((\nabla_x d_i)^2 + (\nabla_y d_i)^2) \quad (1)$$

where d_i is the difference in predicted log depth and ground truth log depth for the valid pixels n (pixels that contain non-zero depth values in the raw depth data), $\nabla_x d_i$ is the horizontal image gradient of the difference and $\nabla_y d_i$ is it's vertical counterpart.

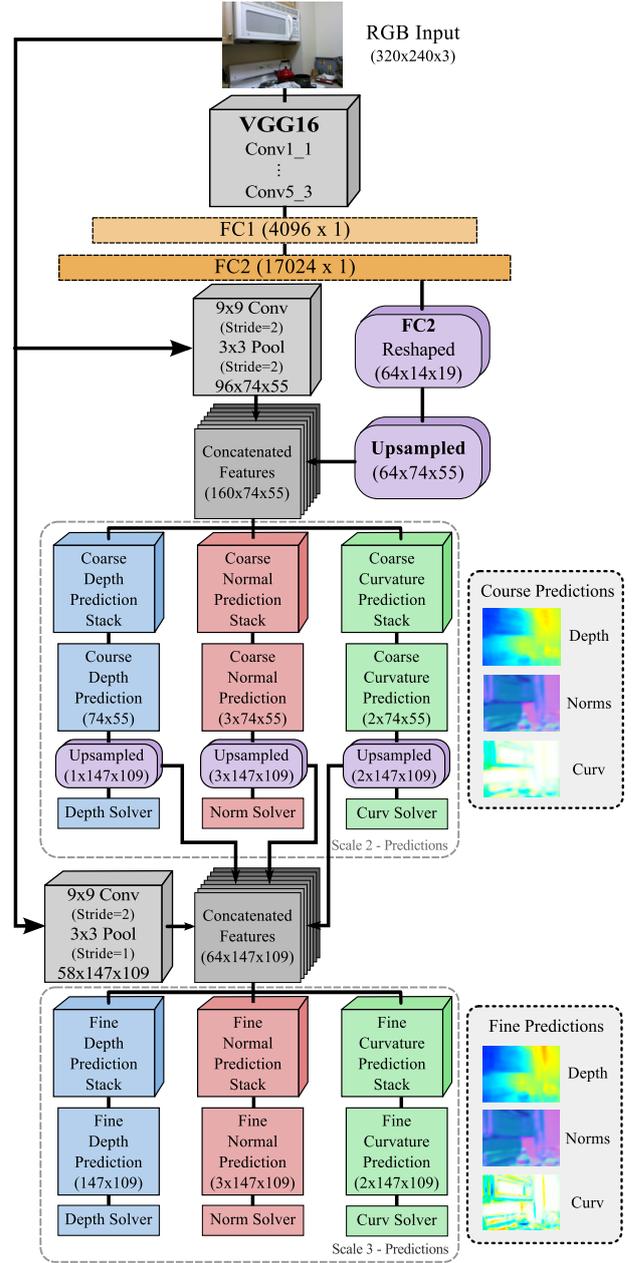


Fig. 2: Visual Representation of Model Architecture

B. Surface Normals

The ground truth normals are computed using different techniques in the literature. We estimate the normals by fitting a quadric patch to a set of nearby points in the point cloud. This gives a more accurate representation of the surface compared to just fitting planar regions, while not adding an extra time complexity as the normals are computed as part of the curvature computation pipeline. We use a combination of pixel wise Euclidean loss along the three channels corresponding to the three unit vectors i, j, k and the difference in angle between the predicted normals and the ground truth as the loss criterion when the normals are trained. This is expressed as

$$L(N, N^*) = -N \cdot N^* + \sum_{i=1}^3 (n_i - n_i^*)^2, \quad (2)$$

where N and N^* are the predicted and ground truth normal respectively, $n_i \in N$ and $n_i^* \in N^*$ are the three components (i, j, k) of each of the normals. Inclusion of the the Euclidean terms improves both the convergence rate and the final accuracy of the system, compared to using the dot product term alone.

C. Surface Curvature

We use the method from [16] to compute an estimate of principal surface curvatures, which is computed from a locally fit parabolic quadric. We use a sparsely sampled circular patch of radius 18 pixels, to fit a quadric at each point and extract the local principal curvature values. We limit the principal curvature κ_1, κ_2 to the range $\{-100, 100\}$ in order to avoid the estimation of implausible curvatures, effectively limiting the minimum detectable radius of curvature to be 1cm. This aligns with the precision of the system[35] at the distances present in the training data. This provides a dense estimate of curvature for every point (640x480), which we then bicubically downsample to 120x160 to generate the LMDBs that can be used in the training of our network. We attempt to estimate principal curvatures directly as opposed to Gaussian or mean curvature, as we found principal curvatures to provide improved performance during training.

We employed a Euclidean loss criterion with depth based weighting to predict surface curvature. Due to the inherent sensor noise the computed principal curvatures which are used as the ground truth tend to have a large uncertainty beyond a certain distance threshold. To prevent the network from learning these rather uncertain values we use the following loss function

$$L(C, C^*) = \sum_{i=1}^n \frac{(\kappa_{1i} - \kappa_{1i}^*)^2 + (\kappa_{2i} - \kappa_{2i}^*)^2}{(1 + D_i)^{-2}}, \quad (3)$$

where κ_{1i} and κ_{2i} are the predicted principal curvatures and κ_{1i}^* and κ_{2i}^* are their corresponding ground truth values while D represents the depth in meters for the i^{th} pixel.

V. TRAINING

A. Data Generation

We randomly augment the training data by performing flips, translations, rotations and variations on the color channels. The same transformation is applied to the RGB input, ground truth depth, surface curvature and surface normals in order to obtain consistent training data. Unlike some notable previous approaches [8], we use the raw depth directly from the dataset provided without any post processing to fill holes or smooth surfaces. We also use the raw depths to calculate the surface normals and surface curvature, which provides a stronger link between our three ground truth sources.

As described in Section IV-C we use the method described in [16] to produce training data for surface normals and surface curvature. Their approach to curvature and normal estimation is specifically targeted for noisy data such as that from a Microsoft Kinect, and they show that it produces good estimates for both surface normals and principal curvatures. We found that by scaling the ground truth curvature values by a factor of 0.1, to produce a

similar range of values to the input depth, improved both qualitative and quantitative results of curvature estimation. We reverse this scaling when we provide our final prediction for both principal curvatures κ_1 and κ_2 by multiplying each value by 10.

B. Hyperparameters and Weight Initialisation

We use Nesterovs accelerated gradient [36] as the optimizer with a base learning rate of 0.1 and a momentum of 0.95 and train for 50 epochs using a NVIDIA GeForce GTX 1080, which took approximately 4 days. Weights of the convolutional layers corresponding to feature extraction were initialized using VGG pretrained on ILSVRC [37] image data. We also experimented with initializing the feature extraction layers with the VGG weights of [8] and found that it did not give a qualitative or quantitative improvement, although it converged faster. All the convolutional layers corresponding to depth, normals and curvature estimation and the fully connected layers were randomly initialized using MSRA weight initialization scheme [38] which converged much faster compared to initializing the filters from a Gaussian distribution with zero mean and 0.01 standard deviation. Everytime when the training loss plateaued (approximately every 10 epochs) we halved the learning rate and continued training. Caffe [39] was used as the learning framework and all the experiments were carried out using a mini batch size of 16.

C. Training Separate Models With Equal Model Capacity

We train several models with equivalent model capacity to estimate quantities both separately and jointly. We do this to demonstrate that the improved estimates for normals and depths are not the result of increased model capacity, but more likely the result of including derived features as tasks to the network. Explicitly we train 4 models, depth only, normals only, depth+normals, depth+normals+curvature, all while maintaining a constant model capacity for each task. More concretely, when we train a single quantity model (depths only or normals only) we leave the coarse level convolutional layers corresponding to the other tasks in place. However, there is only a single solver attached at the end of scale 2 based on the training task. We would like to point the reader to Figure 3, in which we are looking at the scale 2 prediction section of our model. When we are training all three tasks jointly, there is a solver attached at the end of each prediction stack. Simultaneously, we pass the coarse level predictions to scale 3 to be refined further. In a scenario where there is only one training task, the solver corresponding to the training task is kept intact while the other solvers are removed. However, the feature maps of the other stacks are still present and now act as additional weights which are trained using the scale 3 solver. To recapitulate, we preserve the model capacity by keeping the number of feature maps a constant regardless of the task/tasks that is being trained while greatly influencing what is being learnt by the feature maps through the use of additional tasks.

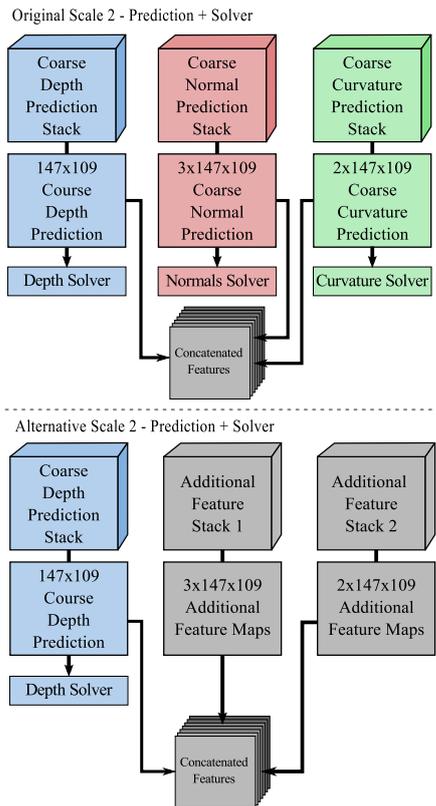


Fig. 3: A closer look at Scale 2 of the architecture for different tasks. **Top:** When all three tasks are trained jointly, there is a solver at the end of Scale 2 for all three tasks and the coarse feature maps are passed on to Scale 3 after being concatenated together. **Bottom:** When only a single task is trained (in this case depth) there is a single solver at the end of Scale 2 and the other two stacks now provide additional feature maps which can be trained by the Scale 3 solver (not shown in the figure).

VI. EXPERIMENTS AND RESULTS

In this section we evaluate the performance of our system across the three tasks. We begin with a quantitative analysis for each of the tasks using the established benchmarks. Then we present qualitative results of our system and conclude this section with a segmentation example to showcase how this work could be applied in a real life scenario.

| Depth Prediction | | | | | | | |
|------------------|----------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Type | Method | Rel_{abs} | RMS_{lin} | RMS_{log} | δ | δ^2 | δ^3 |
| single | Liu[9] | 0.230 | 0.824 | - | 0.614 | 0.883 | 0.972 |
| | Eigen[26] | 0.214 | 0.877 | 0.283 | 0.614 | 0.888 | 0.972 |
| | Ours(Depth) | 0.156 | 0.646 | 0.216 | 0.765 | 0.949 | 0.987 |
| | Laina[11] | 0.127 | 0.573 | 0.195 | 0.811 | 0.953 | 0.988 |
| joint | Eigen(Alex)[8] | 0.198 | 0.753 | 0.255 | 0.697 | 0.912 | 0.977 |
| | Ours(D+N) | 0.156 | 0.642 | 0.215 | 0.766 | 0.949 | 0.988 |
| | Eigen(VGG)[8] | 0.158 | 0.641 | 0.214 | 0.769 | 0.950 | 0.988 |
| | Ours(Full) | 0.156 | 0.624 | 0.212 | 0.776 | 0.953 | 0.989 |

TABLE I: Depth prediction Metrics: the middle three columns indicate errors (lower better) from ground truth, the final three columns indicate the percentage of points within δ^n (higher better) of the ground truth ($\delta = 1.25$). The bold values indicate the best performing method of each type (single, joint).

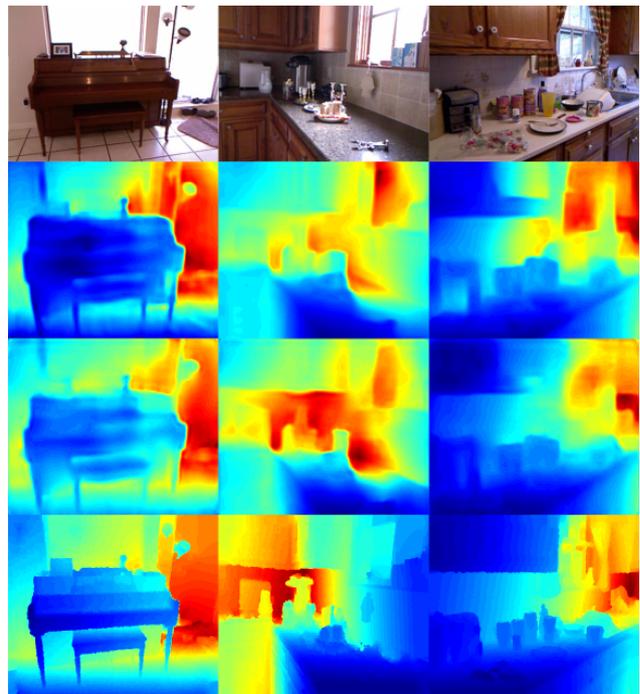


Fig. 4: Demonstrates the qualitative improvement of our approach for depth estimation. **Top:** RGB image **1st row:** Eigen's Prediction **2nd row:** Our Prediction **Bottom:** Ground Truth

A. Depth

We evaluate our depth predictions in the same manner as outlined in previous work [11],[8] and the results are tabulated in Table I. The predicted depth maps are upsampled by a factor of 4 to match the image resolution of 640x480 and are evaluated against the official ground truth depth maps including the filled in areas but limited to the region where there is a valid depth map projection. In terms of relative performance we improve mostly in terms of RMS_{lin} and have similar performance for Rel_{abs} and RMS_{log} , which are more related to the ratio of predicted and ground truth depths. We have included the methods that estimate depth alone as a single task for completeness, although we outperform all the methods except [11] which uses a much more powerful ResNet[20] architecture. Based on the results of the joint task training scheme we strongly believe that the performance of [11] could still be improved had it been trained simultaneously with normals and surface curvature. As we keep adding more tasks that are based on related quantities we can see gains in performance. Also the contribution from curvature is much more significant (reduction of RMSE by 0.02) compared to the contribution of normals (reduction of RMSE by 0.004). As it can be seen in Table I the contribution of semantic labels (Eigen VGG [8]), although very small, helps to increase the performance. But, curvature being a more tightly connected quantity to depth gives the largest improvement.

B. Surface Normals

We compare our normals in a similar way to [8, 21, 27]. As we don't have access to ground truth normal data, we compare our approach against two different methods of

| Type | Method | Angular Error | | Within t° | | |
|--------|-------------------|---------------|-------------|--------------------|-------------------|-----------------|
| | | Mean | Median | $\leq 11.25^\circ$ | $\leq 22.5^\circ$ | $\leq 30^\circ$ |
| single | Ladicky [28] | 35.3 | 31.2 | 16.4 % | 36.6% | 48.2% |
| | Wang [21] | 26.9 | 14.8 | 42.0% | 61.2% | 68.2% |
| | Ours (Norms) | 21.1 | 13.5 | 43.6% | 66.6% | 75.4% |
| | Bansal et al [22] | 19.8 | 12.0 | 47.9 % | 70.0 % | 77.8 % |
| joint | Eigen(Alex)[8] | 23.7 | 15.5 | 39.2 % | 62.0 % | 71.1% |
| | Ours (D+N) | 21.1 | 13.6 | 43.6% | 66.5% | 75.4% |
| | Eigen(VGG)[8] | 20.9 | 13.2 | 44.4% | 67.2% | 75.9% |
| | Ours (Full) | 20.6 | 13.0 | 44.9% | 67.7% | 76.3% |

| Type | Method | Angular Error | | Within t° | | |
|--------|-------------------|---------------|-------------|--------------------|-------------------|-----------------|
| | | Mean | Median | $\leq 11.25^\circ$ | $\leq 22.5^\circ$ | $\leq 30^\circ$ |
| single | Wang [21] | 36.4 | 26.2 | 27.2% | 45.6% | 53.9% |
| | Ours (Norms) | 27.7 | 20.2 | 31.8% | 53.7% | 63.8% |
| | Bansal[22] | 27.1 | 19.0 | 32.8% | 55.8% | 65.7% |
| joint | Eigen(Alexnet)[8] | 29.7 | 21.8 | 30.0% | 51.0% | 61.0% |
| | Ours (D+N) | 27.7 | 20.2 | 31.7% | 53.6% | 63.7% |
| | Eigen(VGG)[8] | 27.3 | 19.6 | 32.3% | 54.7% | 64.6% |
| | Ours (Full) | 27.2 | 19.6 | 32.9% | 54.7% | 64.7% |

TABLE II: The mean, median angular error and the percentage of points with an angular error less than a threshold t° for several normal estimation approaches evaluated against two different methods [28, 16].

estimating normals from the raw depth data, including the ground truth normals as shown in [28] and the method we use to compute our input data from [16]. Qualitatively [28] takes a more aggressive approach to noise and produces very smoothed out estimates, while the method in [16] produces smoothed normals but still provides sharp edges. During evaluation the regions corresponding to the missing depth values are masked out since the ground truth normals can not be accurately computed on those areas. We summarise these results in Table II and demonstrate improved results for each normal estimation method over previous methods. Quantitatively we approach the performance metrics of Bansal et al. [22] who used a skip architecture with a larger model capacity compared to ours, although arguably qualitatively both [8] and our approach outperform their predictions as shown in Figure 5.

Similar to depths, predicted normals also gained an increase in accuracy when the network was trained in a multi task platform. Although, having merely depths in parallel did not make a noticeable change extending the network to learn all three tasks resulted in a significant improvement.

C. Surface Curvature

In order to evaluate the accuracy of estimating surface curvature without access to true ground truth data we evaluate the performance of our approach against the method of [16]. We evaluate the predictions from our network which attempts to explicitly predict curvature, to the estimated curvature values computed from the predicted depths produced by our own network and the network from [8]. We compare the RMS error of each of the principal curvatures (κ_1, κ_2) against the computed ground truth and present the median error of the *mean curvature* $0.5 * (\kappa_1 + \kappa_2)$ across two categories, planar and non-planar regions. We define planar surfaces to be those with a radius of curvature greater than 1 meter. As expected predicted curvatures clearly outperform the computed curvatures



Fig. 5: **Demonstrates the qualitative improvement of our approach for normal estimation. Top: RGB image 1st row: Bansal[22], 2nd row: Eigen[8], 3rd row: Our Prediction Bottom: Ground Truth [16].** The missing areas in the ground truth normals coincide with those in the raw depth images.

from depths. Furthermore, the predicted curvatures using the joint model which learned surface normals and depths in parallel provide better performance compared to the model which only learnt surface curvature.

Figure 6 is included as a reference to show how the metrics in Table III translate into visual appearance.

| Method [16] | RMS (m^{-1}) | | Median (m^{-1}) | | Within σ_i | | |
|-------------------|------------------|-------------|---------------------|--------------|-------------------|--------------|--------------|
| | κ_1 | κ_2 | planar | non-planar | σ_1 | σ_2 | σ_3 |
| Eigen(Depth) [8] | 5.56 | 7.50 | 3.86 | 1.44 | 25.7% | 33.9% | 43.5% |
| Ours (Depth) | 6.03 | 6.50 | 4.23 | 1.38 | 26.9% | 34.9% | 44.2% |
| Ours (Curvatures) | 3.41 | 5.17 | 1.984 | 0.184 | 52.6% | 63.2% | 73.2% |
| Ours (joint) | 2.81 | 4.47 | 1.634 | 0.085 | 63.1% | 72.7% | 80.3% |

TABLE III: The table shows the RMS error of estimating the principal surface curvatures (κ_1, κ_2), the median error for planar and non-planar regions and the percentage of curvatures values that are within a threshold $\sigma_1 = 0.25m^{-1}$, $\sigma_2 = 0.5m^{-1}$, $\sigma_3 = 1m^{-1}$. The first two approaches do not explicitly predict curvature and are computed from the predicted depths.

D. Possible Applications For This Work

As a purely qualitative demonstration of our approach, we show a simple example of scene segmentation that combines information from the colour, depth and curvature of selected scenes. We generate a simple segmentation by

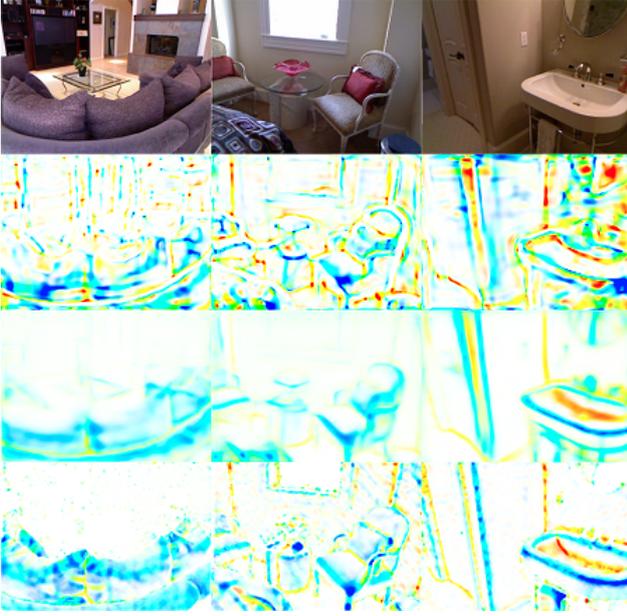


Fig. 6: Demonstrates the qualitative improvement of our approach for surface curvature estimation. **Top:** RGB image **1st row:** Computed surface curvature based on Eigen’s[8] depth prediction **2nd row:** Prediction of our system **Bottom:** Ground Truth computed from raw depth data

combining the gradients of colour and depth, and curvature values. This border function $b(u, v)$ can be expressed as

$$b(u, v) = w_I \cdot \nabla I(u, v) + w_D \cdot \nabla D(u, v) + w_C \cdot C(u, v), \quad (4)$$

where $\nabla I(u, v)$ is the the magnitude of the image intensity gradient, $\nabla D(u, v)$ is the magnitude of the depth gradient and $C(u, v)$ is the curvature value at the point u, v . That is

$$\nabla I(u, v) = \sqrt{\frac{\partial I(u, v)^2}{\partial u} + \frac{\partial I(u, v)^2}{\partial v}}, \quad (5)$$

and

$$\nabla D(u, v) = \sqrt{\frac{\partial D(u, v)^2}{\partial u} + \frac{\partial D(u, v)^2}{\partial v}}. \quad (6)$$

The segmentation is then generated by a simple single threshold on this border function. That is a pixel is considered a border ($\mathbf{B}(u, v)$) if it satisfies the condition

$$\mathbf{B}(u, v) = \begin{cases} 1 & \text{if } b(u, v) \geq \delta_{thresh} \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

We compare the performance of this segmentation method using the ground truth quantities and the predictions (depths and curvature) generated by our network. We show the results of this in Figure 7. The results are not intended to be treated as state of the art segmentations, but are included to demonstrate a possible future extension of this work and also to illustrate that the information from the network can be used to perform similar tasks.

VII. CONCLUSIONS

In this work we present a unified multi task learning platform which is capable of predicting depths, surface

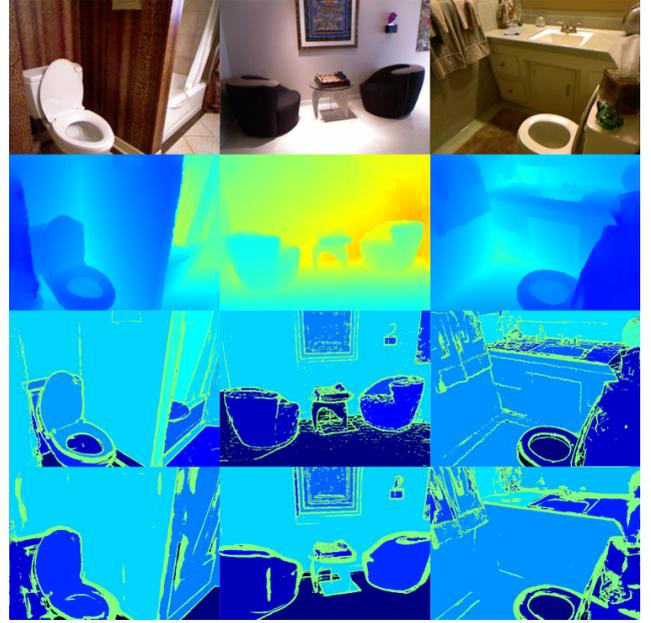


Fig. 7: Demonstrates a basic segmentation algorithm, that uses colour, depth and curvature to generate a border function. The rows of the figure are, top to bottom: Input colour image, Input ground truth depth, Segmentation From GT data, Segmentation from Predicted Data. The key contribution of the depth and curvature to the segmentations, are on the depth boundaries and wall edges that are difficult to differentiate from colour alone.

normals and surface curvatures using a single RGB image. We show that carefully chosen hand crafted feature representations can outperform the machine learnt features provided they are closely related to the prediction task. This shows that network guidance is a useful aspect and should not be completely ignored when training neural networks. We run extensive experiments where we keep the model capacity of the architecture fixed while gradually increasing the number of prediction tasks to verify the effectiveness of our hypothesis. We provide a potential application for our work in a robotic context in the form of a segmentation example as a qualitative demonstration.

REFERENCES

- [1] J. Long, E. Shelhamer, and T. Darrell, “Fully Convolutional Networks for Semantic Segmentation.”
- [2] R. A. Newcombe, S. J. Lovegrove, and A. J. Davison, “DTAM : Dense Tracking and Mapping in Real-Time,” pp. 2320–2327, 2011.
- [3] B. Espiau, F. Chaumette, and P. Rives, “A new approach to visual servoing in robotics,” *IEEE Transactions on Robotics and Automation*, vol. 8, no. 3, pp. 313–326, 2002.
- [4] P. Besl and R. Jain, “Segmentation through variable-order surface fitting,” *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, vol. 10, no. 2, pp. 167–192, mar 1988.
- [5] I. Douros and B. Buxton, “Three-Dimensional Surface Curvature Estimation using Quadric Surface Patches,” *Scanning*, vol. 44, no. 0, 2002.
- [6] Y. Alshabkeh, N. Haala, and D. Fritsch, “Range Image Segmentation Using the Numerical Description of Mean Curvature Values,” in *The International Archives of Photogrammetry, Remote Sensing and Spatial Information Sciences (ISPRS)*, 2008.

- [7] J. Lee, S. Kim, and S.-J. Kim, "Mesh segmentation based on curvatures using the GPU," *Multimedia Tools and Applications*, no. April, Jun 2014.
- [8] D. Eigen and R. Fergus, "Predicting Depth, Surface Normals and Semantic Labels with a Common Multi-Scale Convolutional Architecture."
- [9] F. Liu, C. Shen, and G. Lin, "Deep Convolutional Neural Fields for Depth Estimation from a Single Image."
- [10] A. Saxena, S. H. Chung, and A. Y. Ng, "Learning Depth from Single Monocular Images."
- [11] I. Laina, C. Rupprecht, and F. Tombari, "Deeper Depth Prediction with Fully Convolutional Residual Networks," 2016.
- [12] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor Segmentation and Support Inference from RGBD Images," pp. 1–14.
- [13] S. Rusinkiewicz, "Estimating curvatures and their derivatives on triangle meshes," in *Symposium on 3D Data Processing, Visualization and Transmission (3DPVT)*, no. 2. IEEE, 2010, pp. 486–493.
- [14] W. Griffin, Y. Wang, D. Berrios, and M. Olano, "Real-time GPU surface curvature estimation on deforming meshes and volumetric data sets." *IEEE Transactions on Visualization and Computer Graphics*, vol. 18, no. 10, pp. 1603–13, Oct 2012.
- [15] N. J. Mitra, N. Gelfand, H. Pottmann, and L. Guibas, "Registration of point cloud data from a geometric optimization perspective," in *Eurographics/ACM SIGGRAPH symposium on Geometry processing*, no. January, 2004, p. 22.
- [16] A. Spek and T. Drummond, "A Fast Method For Computing Principal Curvatures From Range Images," in *Australian Conference on Robotics and Automation (ACRA)*. ARAA, 2015.
- [17] J. Redmon and A. Angelova, "Real-Time Grasp Detection Using Convolutional Neural Networks," pp. 1316–1322, 2015.
- [18] M. Gualtieri, K. Saenko, R. Platt, and I. Science, "High precision grasp pose detection in dense clutter *," 2016.
- [19] A. Krizhevsky and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," pp. 1–9.
- [20] K. He and J. Sun, "Deep Residual Learning for Image Recognition," pp. 1–9.
- [21] X. Wang, D. Fouhey, and A. Gupta, "Designing deep networks for surface normal estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 539–547.
- [22] A. Bansal, B. Russell, and A. Gupta, "Marr revisited: 2d-3d alignment via surface normal prediction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 5965–5974.
- [23] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," in *International Conference on Learning Representations*, 2015.
- [24] G. Klein and D. Murray, "Parallel Tracking and Mapping for Small AR Workspaces," *IEEE and ACM International Symposium on Mixed and Augmented Reality*, pp. 1–10, 2007.
- [25] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon, "Kinectfusion: Real-time dense surface mapping and tracking," in *Mixed and augmented reality (ISMAR), 2011 10th IEEE international symposium on*. IEEE, 2011, pp. 127–136.
- [26] D. Eigen, C. Puhrsch, and R. Fergus, "Depth Map Prediction from a Single Image using a Multi-Scale Deep Network," pp. 1–9.
- [27] D. F. Fouhey, A. Gupta, and M. Hebert, "Unfolding an Indoor Origami World."
- [28] L. Ladicky, B. Zeisl, and M. Pollefeys, "Discriminatively trained dense surface normal estimation," in *European Conference on Computer Vision*, 2014.
- [29] R. Caruana, "Multitask learning," in *Learning to learn*. Springer, 1998, pp. 95–133.
- [30] X. Gibert, V. M. Patel, and R. Chellappa, "Deep multitask learning for railway track inspection," *IEEE Transactions on Intelligent Transportation Systems*, vol. 18, no. 1, pp. 153–164, 2017.
- [31] Y. Yan, E. Ricci, R. Subramanian, G. Liu, O. Lanz, and N. Sebe, "A multi-task learning framework for head pose estimation under target motion," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 6, pp. 1070–1083, 2016.
- [32] O. Chapelle, P. Shivaswamy, S. Vadrevu, K. Weinberger, Y. Zhang, and B. Tseng, "Boosted multi-task learning," *Machine learning*, vol. 85, no. 1-2, pp. 149–173, 2011.
- [33] X. Wang, C. Zhang, and Z. Zhang, "Boosted multi-task learning for face verification with applications to web image and video search," in *computer vision and pattern recognition, 2009. CVPR 2009. IEEE conference on*. IEEE, 2009, pp. 142–149.
- [34] Z. Li and D. Hoiem, "Learning without forgetting," in *European Conference on Computer Vision*. Springer, 2016, pp. 614–629.
- [35] K. Khoshelham and S. O. Elberink, "Accuracy and resolution of Kinect depth data for indoor mapping applications." *Sensors*, vol. 12, no. 2, pp. 1437–54, Jan 2012.
- [36] Y. Nesterov, "A method of solving a convex programming problem with convergence rate $o(1/k^2)$," in *Soviet Mathematics Doklady*, vol. 27, no. 2, 1983, pp. 372–376.
- [37] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.
- [38] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1026–1034.
- [39] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," *arXiv preprint arXiv:1408.5093*, 2014.