
Multi-Task Learning Using Uncertainty to Weigh Losses for Scene Geometry and Semantics

Alex Kendall
University of Cambridge
agk34@cam.ac.uk

Yarin Gal
University of Cambridge
yg279@cam.ac.uk

Roberto Cipolla
University of Cambridge
rc10001@cam.ac.uk

Abstract

Numerous deep learning applications benefit from multi-task learning with multiple regression and classification objectives. In this paper we make the observation that the performance of such systems is strongly dependent on the relative weighting between each task’s loss. Tuning these weights by hand is a difficult and expensive process, making multi-task learning prohibitive in practice. We propose a principled approach to multi-task deep learning which weighs multiple loss functions by considering the homoscedastic uncertainty of each task. This allows us to simultaneously learn various quantities with different units or scales in both classification and regression settings. We demonstrate our model learning per-pixel depth regression, semantic and instance segmentation from a monocular input image. Perhaps surprisingly, we show our model can learn multi-task weightings and outperform separate models trained individually on each task.

1 Introduction

Multi-task learning aims to improve learning efficiency and prediction accuracy by learning multiple objectives from a shared representation [1]. Multi-task learning is prevalent in many applications of machine learning – from computer vision [2] to natural language processing [3] to speech recognition [4].

We explore multi-task learning within the setting of visual scene understanding in computer vision. Scene understanding algorithms must understand both the geometry and semantics of the scene at the same time. This forms an interesting multi-task learning problem because scene understanding involves joint learning of various regression and classification tasks with different units and scales. Multi-task learning of visual scene understanding is of crucial importance in systems where long computation run-time is prohibitive, such as the ones used in robotics. Combining all tasks into a single model reduces computation and allows these systems to run in real-time.

Prior approaches to simultaneously learning multiple tasks use a naïve weighted sum of losses, where the loss weights are uniform, or manually tuned [2, 5, 6]. However, we show that performance is highly dependant on an appropriate choice of weighting between each task’s loss. Searching for an optimal weighting is prohibitively expensive and difficult to resolve with manual tuning. We observe that the optimal weighting of each task is dependant on the measurement scale (e.g. meters, centimetres or millimetres) and ultimately the magnitude of the task’s noise. In this work we propose a principled way of combining multiple loss functions to simultaneously learn multiple objectives using homoscedastic uncertainty. We interpret homoscedastic uncertainty as task-dependant weighting and show how to derive a principled multi-task loss function which can learn to balance various regression and classification losses. Our method can learn to balance these weightings optimally, resulting in superior performance, compared with learning each task individually.

Specifically, we demonstrate our method in learning scene geometry and semantics with three tasks. Firstly, we learn to classify objects at a pixel level, also known as semantic segmentation [7–11].

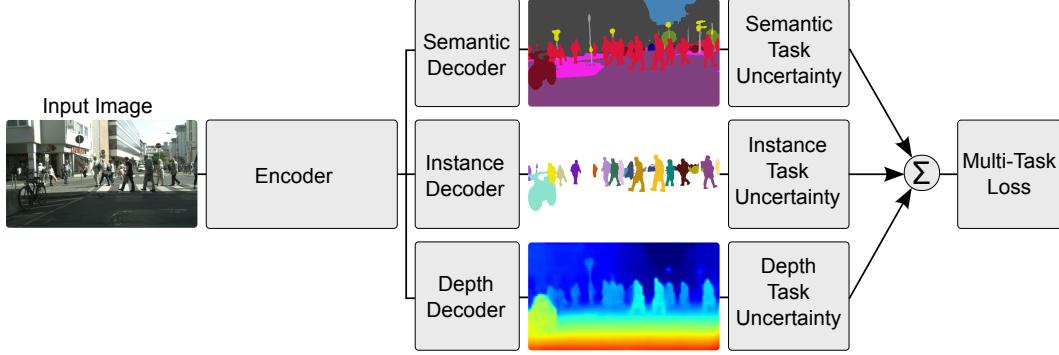


Figure 1: Multi-task deep learning. We derive a principled way of combining multiple regression and classification loss functions for multi-task learning. Our architecture takes a single monocular RGB image as input and produces a pixel-wise classification, an instance semantic segmentation and an estimate of per pixel depth. Multi-task learning can improve accuracy over separately trained models because cues from one task, such as depth, are used to regularize and improve the generalization of another domain, such as segmentation.

Secondly, our model performs instance segmentation, which is the harder task of segmenting separate masks for each individual object in an image (for example, a separate, precise mask for each individual car on the road) [12–15]. This is a more complicated task than semantic segmentation, as it requires not only an estimate of each pixel’s class, but also which object that pixel belongs to. It is also more complicated than object detection, which often predicts object bounding boxes alone [16]. Finally, our model predicts pixel-wise metric depth. Depth by recognition has been demonstrated using dense prediction networks with supervised [6] and unsupervised [17] deep learning. However it is very hard to estimate depth in a way which generalises well. We show that we can improve our estimation of geometry *and* depth by using semantic labels and multi-task deep learning.

In existing literature, separate deep learning models would be used to learn depth regression, semantic segmentation and instance segmentation to create a complete scene understanding system. Given a single monocular input image, our system is the first to produce a semantic segmentation, a dense estimate of metric depth and an instance level segmentation jointly (Figure 1). While other vision models have demonstrated multi-task learning, we show how to learn to combine semantics and geometry. Combining these tasks into a single model ensures that the model agrees between the separate task outputs while reducing computation. Finally, we show that using a shared representation with multi-task learning improves performance on various metrics, making the models more effective.

In summary, the key contributions of this paper are:

1. a novel and principled multi-task loss to simultaneously learn various classification and regression losses of varying quantities and units using homoscedastic task uncertainty,
2. a unified architecture for semantic segmentation, instance segmentation and depth regression,
3. demonstrating the importance of loss weighting in multi-task deep learning and how to obtain superior performance compared to equivalent separately trained models.

2 Related Work

Multi-task learning aims to improve learning efficiency and prediction accuracy for each task, when compared to training the models separately [18, 19]. It can be considered an approach to inductive knowledge transfer which improves generalisation by sharing the domain information between complimentary tasks. It does this by using a shared representation to learn multiple tasks – what is learned from one task can help other tasks be learned better [1].

Fine-tuning [20, 21] is a basic example of multi-task learning, where we can leverage different learning tasks by considering them as a pre-training step. Other models alternate learning between each training task, for example in natural language processing [3]. Multi-task learning can also be used in a data streaming setting [18], or to prevent forgetting previously learned tasks in reinforcement

learning [22]. It can also be used to learn unsupervised features from various data sources with an auto-encoder [23].

In computer vision there are many examples of methods for multi-task learning. Many focus on semantic tasks, such as classification and semantic segmentation [24] or classification and detection [5]. MultiNet [25] proposes an architecture for detection, classification and semantic segmentation. CrossStitch networks [26] explore methods to combine multi-task neural activations. Uhrig et al. [27] learn semantic and instance segmentations under a classification setting. Multi-task deep learning has also been used for geometry and regression tasks. Eigen and Fergus [6] show how to learn semantic segmentation, depth and surface normals. PoseNet [28] is a model which learns camera position and orientation. UberNet [2] learns a number of different regression and classification tasks under a single architecture. In this work we are the first to propose a method for jointly learning depth regression, semantic and instance segmentation. Like the model of Eigen and Fergus [6], our model learns both semantic and geometry representations, which is important for scene understanding. However, our model learns the much harder task of instance segmentation which requires knowledge of both semantics and geometry. This is because our model must determine the class and spatial relationship for each pixel in each object for instance segmentation.

More importantly, all previous methods which learn multiple tasks simultaneously use a naïve weighted sum of losses, where the loss weights are uniform, or crudely and manually tuned. In this work we propose a principled way of combining multiple loss functions to simultaneously learn multiple objectives using homoscedastic task uncertainty. We illustrate the importance of appropriately weighting each task in deep learning to achieve good performance and show that our method can learn to balance these weightings optimally.

3 Multi Task Learning with Homoscedastic Uncertainty

Multi-task learning concerns the problem of optimising a model with respect to multiple objectives. It is prevalent in many deep learning problems. The naive approach to combining multi objective losses would be to simply perform a weighted linear sum of the losses for each individual task:

$$L_{total} = \sum_i w_i L_i. \quad (1)$$

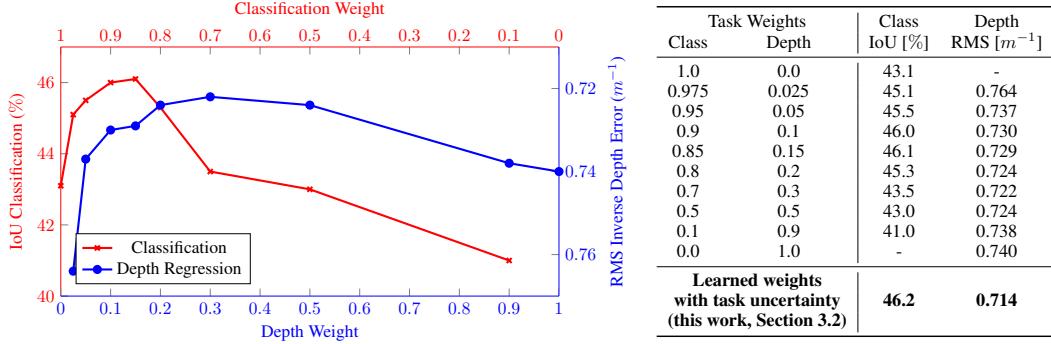
This is the dominant approach used by prior work [5, 24, 25, 27], for example for dense prediction tasks [2], for scene understanding tasks [6] and for rotation (in quaternions) and translation (in meters) for camera pose [28]. However, there are a number of issues with this method. Namely, model performance is extremely sensitive to weight selection, w_i , as illustrated in Figure 2. These weight hyper-parameters are expensive to tune, often taking many days for each trial. Therefore, it is desirable to find a more convenient approach which is able to learn the optimal weights.

More concretely, let us consider a network which learns to predict pixel-wise depth and semantic class from an input image. In Figure 2 the two boundaries of each plot show models trained on individual tasks, with the curves showing performance for varying weights w_i for each task. We observe that at some optimal weighting, the joint network performs better than separate networks trained on each task individually (performance of the model in individual tasks is seen at both edges of the plot: $w = 0$ and $w = 1$). At near-by values to the optimal weight the network performs worse on one of the tasks. However, searching for these optimal weightings is expensive and increasingly difficult with large models with numerous tasks. We next show how to learn optimal task weightings using ideas from probabilistic modelling. Additionally, we show a similar result for two regression tasks; instance segmentation and depth regression.

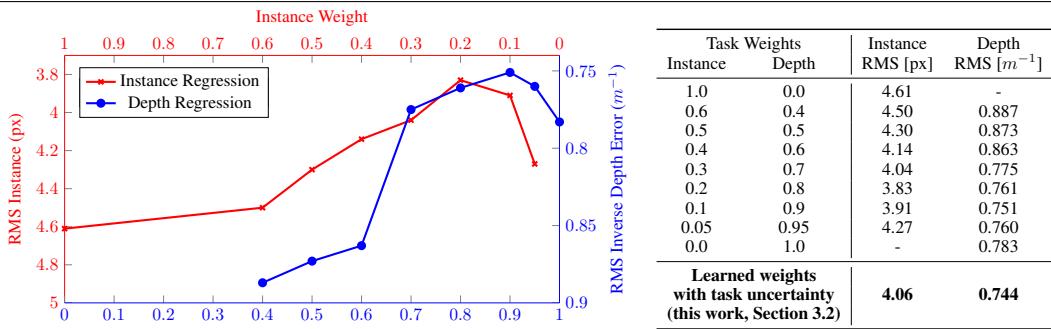
3.1 Homoscedastic uncertainty as task-dependant uncertainty

In Bayesian modelling, there are two main types of uncertainty one can model [29].

- *Epistemic uncertainty* is uncertainty in the model, which captures what our model doesn't know due to lack of training data. It can be explained away with increased training data.
- *Aleatoric uncertainty* captures our uncertainty with respect to information which our data cannot explain. Aleatoric uncertainty can be explained away with the ability to observe all explanatory variables with increasing precision.



(a) Comparing loss weightings when learning **semantic classification and depth regression**



(b) Comparing loss weightings when learning **instance regression and depth regression**

Figure 2: **Learning multiple tasks improves the model’s representation and individual task performance.** These figures and tables illustrate the advantages of multi-task learning for (a) semantic classification and depth regression and (b) instance and depth regression. Performance of the model in individual tasks is seen at both edges of the plot where $w = 0$ and $w = 1$. For some balance of weightings between each task, we observe improved performance for both tasks. All models were trained with a learning rate of 0.01 with the respective weightings applied to the losses using the loss function in (1). Results are shown using the Tiny CityScapes validation dataset using a down-sampled resolution of 128×256 .

Aleatoric uncertainty can again be divided into two sub-categories.

- *Data-dependant* or *Heteroscedastic* uncertainty is aleatoric uncertainty which depends on the input data and is predicted as a model output.
- *Task-dependant* or *Homoscedastic* uncertainty is aleatoric uncertainty which is not dependant on the input data. It is not a model output, rather it is a quantity which stays constant for all input data and varies between different tasks. It can therefore be described as task-dependant uncertainty.

In a multi-task setting, we show that the task uncertainty captures the relative confidence between tasks, reflecting the uncertainty inherent to the regression or classification task. It will also depend on the task’s representation or unit of measure. We propose that we can use homoscedastic uncertainty as a basis for weighting losses in a multi-task learning problem.

3.2 Multi-task likelihoods

In this section we derive a multi-task loss function based on maximising the Gaussian likelihood with homoscedastic uncertainty. Let $\mathbf{f}^{\mathbf{W}}(\mathbf{x})$ be the output of a neural network with weights \mathbf{W} on input \mathbf{x} . We define the following probabilistic model. For regression tasks we define our likelihood as a Gaussian with mean given by the model output:

$$p(\mathbf{y}|\mathbf{f}^{\mathbf{W}}(\mathbf{x})) = \mathcal{N}(\mathbf{f}^{\mathbf{W}}(\mathbf{x}), \sigma^2) \quad (2)$$

with an observation noise scalar σ . For classification we often squash the model output through a softmax function, and sample from the resulting probability vector:

$$p(\mathbf{y}|\mathbf{f}^{\mathbf{W}}(\mathbf{x})) = \text{Softmax}(\mathbf{f}^{\mathbf{W}}(\mathbf{x})). \quad (3)$$

In the case of multiple model outputs, we often define the likelihood to factorise over the outputs, given some sufficient statistics. We define $\mathbf{f}^{\mathbf{W}}(\mathbf{x})$ as our sufficient statistics, and obtain the following multi-task likelihood:

$$p(\mathbf{y}_1, \dots, \mathbf{y}_K | \mathbf{f}^{\mathbf{W}}(\mathbf{x})) = p(\mathbf{y}_1 | \mathbf{f}^{\mathbf{W}}(\mathbf{x})) \dots p(\mathbf{y}_K | \mathbf{f}^{\mathbf{W}}(\mathbf{x})) \quad (4)$$

with model outputs $\mathbf{y}_1, \dots, \mathbf{y}_K$ (such as semantic segmentation, depth regression, etc).

In *maximum likelihood* inference, we maximise the log likelihood of the model. In regression, for example, the log likelihood can be written as

$$\log p(\mathbf{y} | \mathbf{f}^{\mathbf{W}}(\mathbf{x})) \propto -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{f}^{\mathbf{W}}(\mathbf{x})\|^2 - \log \sigma^2 \quad (5)$$

for a Gaussian likelihood (or similarly for a Laplace likelihood) with σ the model's observation noise parameter – capturing how much noise we have in the outputs. We then maximise the log likelihood with respect to the model parameters \mathbf{W} and observation noise parameter σ .

Let us now assume that our model output is composed of two vectors \mathbf{y}_1 and \mathbf{y}_2 , each following a Gaussian distribution:

$$\begin{aligned} p(\mathbf{y}_1, \mathbf{y}_2 | \mathbf{f}^{\mathbf{W}}(\mathbf{x})) &= p(\mathbf{y}_1 | \mathbf{f}^{\mathbf{W}}(\mathbf{x})) \cdot p(\mathbf{y}_2 | \mathbf{f}^{\mathbf{W}}(\mathbf{x})) \\ &= \mathcal{N}(\mathbf{y}_1; \mathbf{f}^{\mathbf{W}}(\mathbf{x}), \sigma_1^2) \cdot \mathcal{N}(\mathbf{y}_2; \mathbf{f}^{\mathbf{W}}(\mathbf{x}), \sigma_2^2). \end{aligned} \quad (6)$$

This leads to the following *minimisation* objective (our loss) for our multi-output model:

$$\begin{aligned} \mathcal{L}(\mathbf{W}, \sigma_1, \sigma_2) &= -\log p(\mathbf{y}_1, \mathbf{y}_2 | \mathbf{f}^{\mathbf{W}}(\mathbf{x})) \\ &\propto \frac{1}{2\sigma_1^2} \|\mathbf{y}_1 - \mathbf{f}^{\mathbf{W}}(\mathbf{x})\|^2 + \frac{1}{2\sigma_2^2} \|\mathbf{y}_2 - \mathbf{f}^{\mathbf{W}}(\mathbf{x})\|^2 + \log \sigma_1^2 \sigma_2^2 \\ &= \frac{1}{2\sigma_1^2} \mathcal{L}_1(\mathbf{W}) + \frac{1}{2\sigma_2^2} \mathcal{L}_2(\mathbf{W}) + \log \sigma_1^2 \sigma_2^2 \end{aligned} \quad (7)$$

Where we wrote $\mathcal{L}_1(\mathbf{W}) = \|\mathbf{y}_1 - \mathbf{f}^{\mathbf{W}}(\mathbf{x})\|^2$ for the loss of the first output variable, and similarly for $\mathcal{L}_2(\mathbf{W})$.

We interpret minimising this last objective with respect to σ_1 and σ_2 as learning the relative weight of the losses $\mathcal{L}_1(\mathbf{W})$ and $\mathcal{L}_2(\mathbf{W})$ adaptively, based on the data. As σ_1 – the noise parameter for the variable \mathbf{y}_1 – increases, we have that the weight of $\mathcal{L}_1(\mathbf{W})$ decreases. On the other hand, as the noise decreases, we have that the weight of the respective objective increases. The noise is discouraged from increasing too much (effectively ignoring the data) by the last term in the objective, which acts as a regulariser for the noise terms.

This construction can be trivially extended to multiple regression outputs. However, the extension to classification likelihoods is more interesting. We adapt the classification likelihood to squash a *scaled* version of the model output through a softmax function:

$$p(\mathbf{y} | \mathbf{f}^{\mathbf{W}}(\mathbf{x}), \sigma) = \text{Softmax}\left(\frac{1}{\sigma^2} \mathbf{f}^{\mathbf{W}}(\mathbf{x})\right) \quad (8)$$

with a positive scalar σ . The log likelihood for this output can then be written as

$$\log p(\mathbf{y} = c | \mathbf{f}^{\mathbf{W}}(\mathbf{x}), \sigma) = \frac{1}{\sigma^2} f_c^{\mathbf{W}}(\mathbf{x}) - \log \sum_{c'} \exp\left(\frac{1}{\sigma^2} f_{c'}^{\mathbf{W}}(\mathbf{x})\right) \quad (9)$$

with $f_c^{\mathbf{W}}(\mathbf{x})$ the c 'th element of the vector $\mathbf{f}^{\mathbf{W}}(\mathbf{x})$.

Next, assume that a model's multiple outputs are composed of a continuous output \mathbf{y}_1 and a discrete output \mathbf{y}_2 , modelled with a Gaussian likelihood and a softmax likelihood, respectively. Like before,

the joint loss is given as:

$$\begin{aligned}
\mathcal{L}(\mathbf{W}, \sigma_1, \sigma_2) &= -\log p(\mathbf{y}_1, \mathbf{y}_2 = c | \mathbf{f}^{\mathbf{W}}(\mathbf{x})) \\
&= -\log \mathcal{N}(\mathbf{y}_1; \mathbf{f}^{\mathbf{W}}(\mathbf{x}), \sigma_1^2) \cdot \text{Softmax}(\mathbf{y}_2 = c; \mathbf{f}^{\mathbf{W}}(\mathbf{x}), \sigma_2) \\
&= \frac{1}{2\sigma_1^2} \|\mathbf{y}_1 - \mathbf{f}^{\mathbf{W}}(\mathbf{x})\|^2 + \log \sigma_1^2 - \log p(\mathbf{y}_2 = c | \mathbf{f}^{\mathbf{W}}(\mathbf{x}), \sigma_2) \\
&= \frac{1}{2\sigma_1^2} \mathcal{L}_1(\mathbf{W}) + \frac{1}{\sigma_2^2} \mathcal{L}_2(\mathbf{W}) + \log \sigma_1^2 + \log \frac{\sum_{c'} \exp\left(\frac{1}{\sigma_2^2} f_{c'}^{\mathbf{W}}(\mathbf{x})\right)}{\left(\sum_{c'} \exp(f_{c'}^{\mathbf{W}}(\mathbf{x}))\right)^{\frac{1}{\sigma_2^2}}} \quad (10) \\
&\approx \frac{1}{2\sigma_1^2} \mathcal{L}_1(\mathbf{W}) + \frac{1}{\sigma_2^2} \mathcal{L}_2(\mathbf{W}) + \log \sigma_1^2 + \log \sigma_2^2,
\end{aligned}$$

where again we write $\mathcal{L}_1(\mathbf{W}) = \|\mathbf{y}_1 - \mathbf{f}^{\mathbf{W}}(\mathbf{x})\|^2$ for the Euclidean loss of \mathbf{y}_1 , write $\mathcal{L}_2(\mathbf{W}) = -\log \text{Softmax}(\mathbf{y}_2, \mathbf{f}^{\mathbf{W}}(\mathbf{x}))$ for the cross entropy loss of \mathbf{y}_2 (with $\mathbf{f}^{\mathbf{W}}(\mathbf{x})$ not scaled), and optimise with respect to \mathbf{W} as well as σ_1, σ_2 . In the last transition we introduced the explicit simplifying assumption $\frac{1}{\sigma_2^2} \sum_{c'} \exp\left(\frac{1}{\sigma_2^2} f_{c'}^{\mathbf{W}}(\mathbf{x})\right) \approx \left(\sum_{c'} \exp(f_{c'}^{\mathbf{W}}(\mathbf{x}))\right)^{\frac{1}{\sigma_2^2}}$ which becomes an equality when $\sigma_2^2 \rightarrow 1$. This has the advantage of simplifying the optimisation objective, as well as empirically improving results.

This last objective can be seen as learning the relative weights of the losses for each output. Large scale values σ_2 will decrease the contribution of $\mathcal{L}_2(\mathbf{W})$, whereas small scale σ_2 will increase its contribution. The scale is regulated by the last term in the equation. The objective is penalised when setting σ_2 too large (with the last term contributing a constant value $\log C$ – with C classes – to the loss).

The multi-task objective with homoscedastic task uncertainty now becomes:

$$\mathcal{L}(\mathbf{W}, \sigma_1, \sigma_2, \dots, \sigma_i) = \sum_i \frac{1}{2\sigma_i^2} \mathcal{L}_i(\mathbf{W}) + \log \sigma_i^2 \quad (11)$$

over all tasks indexed by i . Again, we write $\mathcal{L}_i(\mathbf{W}) = \|\mathbf{y}_i - \mathbf{f}^{\mathbf{W}}(\mathbf{x})\|^2$ for regression losses \mathbf{y}_i , and $\mathcal{L}_i(\mathbf{W}) = -\log \text{Softmax}(\mathbf{y}_i, \mathbf{f}^{\mathbf{W}}(\mathbf{x}))$ for classification losses. This construction can be trivially extended to arbitrary combinations of discrete and continuous variables, allowing us to learn the relative weights of each loss in a principled and well-founded way. This loss is smoothly differentiable, and is well formed such that the task weights will not converge to zero. In contrast, directly learning the weights using a simple linear sum of losses (1) would result in weights which quickly converge to zero. In the following sections we introduce our experimental model and present empirical results.

4 Scene Understanding Model

To understand semantics and geometry we first propose an architecture which can learn regression and classification outputs, at a pixel level. Our architecture is a deep convolutional encoder decoder network [8]. Our model consists of a number of convolutional encoders which produce a shared representation, followed by a corresponding number of task-specific convolutional decoders. A high level summary is shown in Figure 1 and additional model details are explained in Appendix A. Our encoder is based on ResNet-101 [30] (without the final fully connected layer). We then split the 2048 dimensional shared feature representation into individual decoders for each task. Each decoder consists of three convolutional layers for each task.

Semantic Segmentation. We use the cross-entropy loss to learn pixel-wise class probabilities, averaging the loss over the pixels with semantic labels in each mini-batch.

Instance Segmentation. An intuitive method for defining which instance a pixel belongs to is an association to the instance’s centroid. We use a regression approach for instance segmentation [31]. This approach is inspired by [32] which identifies instances using Hough votes from object parts. In this work we extend this idea by using votes from individual pixels using deep learning. We learn an

Loss	Task Weights			Classification IoU [%]	Instance RMS Error [px]	Inverse Depth RMS Error [px]
	Cls.	Inst.	Depth			
Class only	1	0	0	43.1%	-	-
Instance only	0	1	0	-	4.61	-
Depth only	0	0	1	-	-	0.783
Unweighted sum of losses	0.333	0.333	0.333	43.6%	3.92	0.786
Approx. optimal weights	0.8	0.05	0.15	46.3%	3.92	0.799
2 task uncertainty weighting	✓	✓		46.5%	3.73	-
2 task uncertainty weighting	✓		✓	46.2%	-	0.714
2 task uncertainty weighting		✓	✓	-	4.06	0.744
3 task uncertainty weighting	✓	✓	✓	46.6%	3.91	0.702

Table 1: Quantitative improvement when learning semantic segmentation, instance segmentation and depth with our multi-task loss. Experiments were conducted on the Tiny CityScapes dataset, sub-sampled to a resolution of 256x512, results are shown from the validation set. We observe an improvement in performance when training with our multi-task loss, over both single-task models and weighted losses. Additionally, we observe an improvement when training on all three tasks ($3 \times \checkmark$) using our multi-task loss, compared with all pairs of tasks alone (denoted by $2 \times \checkmark$). This shows that our loss function can automatically learn an better performing weighting between the tasks.

instance vector, \hat{x}_n , for each pixel coordinate, c_n , which points to the centroid of the pixel’s instance, i_n , such that $i_n = \hat{x}_n + c_n$. We train this regression with an L_1 loss using ground truth labels x_n , averaged over all labelled pixels, N_I , in a mini-batch: $\mathcal{L}_{\text{Instance}} = \frac{1}{|N_I|} \sum_{N_I} \|x_n - \hat{x}_n\|_1$. An illustrated example is given in Appendix C.

To obtain segmentations for each instance, we now need to estimate the instance centres, \hat{i}_n . We propose to consider the estimated instance vectors, \hat{x}_n , as votes in a Hough parameter space and use a clustering algorithm to identify these instance centres. OPTICS [33], is an efficient density based clustering algorithm. It is able to identify an unknown number of multi-scale clusters with varying density from a given set of samples. We chose OPICS for two reasons. Crucially, it does not assume knowledge of the number of clusters like algorithms such as k-means [34]. Secondly, it does not assume a canonical instance size or density like discretised binning approaches [35]. Using OPTICS, we cluster the points $c_n + \hat{x}_n$ into a number of estimated instances, \hat{i} . We can then assign each pixel, p_n to the instance closest to its estimated instance vector, $c_n + \hat{x}_n$.

Depth Regression. We train with supervised labels using pixel-wise metric inverse depth using a L_1 loss function: $\mathcal{L}_{\text{Depth}} = \frac{1}{|N_D|} \sum_{N_D} \|d_n - \hat{d}_n\|_1$. Our architecture estimates inverse depth, \hat{d}_n , because it can represent points at infinite distance (such as sky). We can obtain inverse depth labels, d_n , from a RGBD sensor or stereo imagery. Pixels which do not have an inverse depth label are ignored in the loss.

5 Experiments

We demonstrate the efficacy of our method on CityScapes [36], a large dataset for road scene understanding. It comprises of stereo imagery, from automotive grade stereo cameras with a 22cm baseline, labelled with instance and semantic segmentations from 20 classes. Depth images are also provided, labelled using SGM [37], which we treat as ground truth. Additionally, we assign zero inverse depth to pixels labelled as sky. The dataset was collected from a number of cities in fine weather and consists of 3,250 training and 750 validation images at 2048 \times 1024 resolution. 1,000 images are used for testing on an online evaluation server.

5.1 Model Evaluation

In Table 1 we compare individual models to multi-task learning models using a naïve weighted loss or the task uncertainty weighting we propose in this paper. To reduce the computational burden, we train each model at a reduced resolution of 128 \times 256 pixels. This clearly illustrates the benefit of

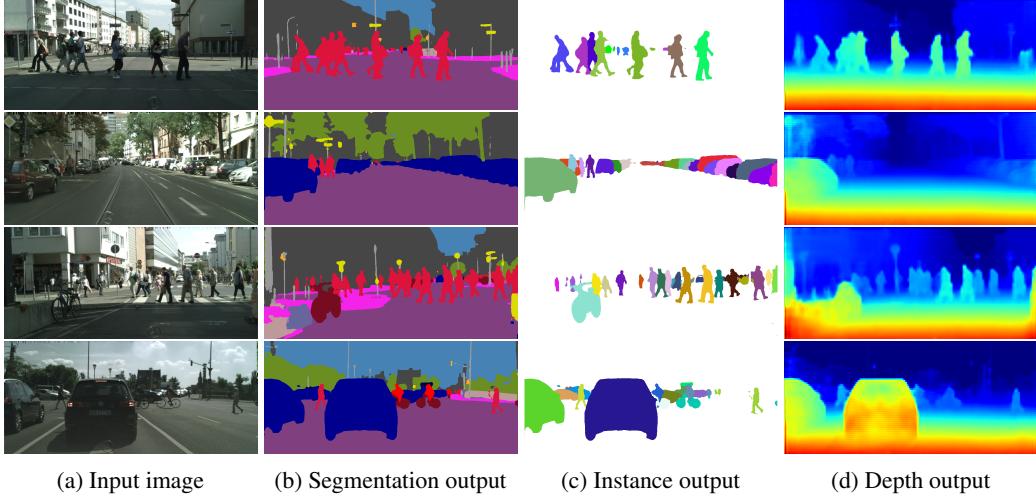


Figure 3: **Qualitative results for multi-task learning of geometry and semantics for road scene understanding.** Results are shown on test images from the CityScapes dataset using our multi-task approach with a single network trained on all tasks. We observe that multi-task learning improves the smoothness and accuracy for depth perception because it learns a representation that uses cues from other tasks, such as segmentation (and vice versa).

multi-task learning, which obtains significantly better performing results than individual task models. For example, using our method we improve classification results from 43.1% to 46.6%.

We also compare to a number of naïve multi-task losses. We compare weighting each task equally and using approximately optimal weights. Using a uniform weighting results in poor performance, in some cases not even improving on the results from the single task model. Obtaining approximately optimal weights is difficult with increasing number of tasks as it requires an expensive grid search over parameters. However, even these weights perform worse compared with our proposed method. Figure 2 shows that using task uncertainty weights can even perform better compared to optimal weights found through fine-grained grid search. We believe that this is due to two reasons. First, grid search is restricted in accuracy by the resolution of the search. Second, optimising the task weights using a homoscedastic noise term allows for the weights to be dynamic during training. In general, we observe that the uncertainty term decreases during training which improves the optimisation process.

This loss is also robust to the value we use to initialise the weights. In Appendix B we show with any reasonable initialisation of $\log \sigma^2$ from -2.0 to 5.0 , the homoscedastic uncertainty terms converge to the same value after 100 training iterations. This is significantly less than the 30,000 training iterations for the network. Therefore our model is robust to the choice of initial value for the weighting terms. Interestingly, we observe that this loss allows the network to dynamically tune the weighting. Typically, the homoscedastic noise terms decrease in magnitude as training progresses.

6 Conclusions

We have shown that correctly weighting loss terms is of paramount importance for multi-task learning problems. We demonstrated that homoscedastic (task) uncertainty is an effective way to weight losses. We derived a principled loss function which can learn a relative weighting automatically from the data and is robust to the weight initialization. We showed that this can improve performance for scene understanding tasks with a unified architecture for semantic segmentation, instance segmentation and per-pixel depth regression. We demonstrated modelling task-dependant homoscedastic uncertainty improves the model’s representation and each task’s performance when compared to separate models trained on each task individually.

An interesting question left unanswered is where the optimal location is for splitting the shared encoder network into separate decoders for each task? And, what network depth is best for the shared multi-task representation?

References

- [1] Rich Caruana. Multitask learning. In *Learning to learn*, pages 95–133. Springer, 1998.
- [2] Iasonas Kokkinos. Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory. *arXiv preprint arXiv:1609.02132*, 2016.
- [3] Ronan Collobert and Jason Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167. ACM, 2008.
- [4] Jui-Ting Huang, Jinyu Li, Dong Yu, Li Deng, and Yifan Gong. Cross-language knowledge transfer using multilingual deep neural network with shared hidden layers. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 7304–7308. IEEE, 2013.
- [5] Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *International Conference on Learning Representations (ICLR)*, 2014.
- [6] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2650–2658, 2015.
- [7] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2015.
- [8] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for scene segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [9] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. In *ICLR*, 2016.
- [10] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*, 2015.
- [11] Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip Torr. Conditional random fields as recurrent neural networks. In *International Conference on Computer Vision (ICCV)*, 2015.
- [12] Pedro O Pinheiro, Ronan Collobert, and Piotr Dollar. Learning to segment object candidates. In *Advances in Neural Information Processing Systems*, pages 1990–1998, 2015.
- [13] Bharath Hariharan, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Hypercolumns for object segmentation and fine-grained localization. In *In Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, pages 447–456. IEEE, 2014.
- [14] Jifeng Dai, Kaiming He, and Jian Sun. Instance-aware semantic segmentation via multi-task network cascades. In *In Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2016.
- [15] Min Bai and Raquel Urtasun. Deep watershed transform for instance segmentation. *arXiv preprint arXiv:1611.08303*, 2016.
- [16] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *In Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, pages 580–587, 2014.
- [17] Ravi Garg and Ian Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. *Computer Vision-ECCV 2016*, pages 740–756, 2016.
- [18] Sebastian Thrun. Is learning the n-th thing any easier than learning the first? In *Advances in neural information processing systems*, pages 640–646. MORGAN KAUFMANN PUBLISHERS, 1996.
- [19] Jonathan Baxter et al. A model of inductive bias learning. *J. Artif. Intell. Res.(JAIR)*, 12(149-198):3, 2000.
- [20] Pulkit Agrawal, Joao Carreira, and Jitendra Malik. Learning to see by moving. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 37–45, 2015.
- [21] Maxime Oquab, Leon Bottou, Ivan Laptev, and Josef Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *In Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, pages 1717–1724. IEEE, 2014.
- [22] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, page 201611835, 2017.
- [23] Jiquan Ngiam, Aditya Khosla, Mingyu Kim, Juhan Nam, Honglak Lee, and Andrew Y Ng. Multimodal deep learning. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 689–696, 2011.

- [24] Yiyi Liao, Sarath Kodagoda, Yue Wang, Lei Shi, and Yong Liu. Understand scene categories by objects: A semantic regularized scene classifier using convolutional neural networks. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2318–2325. IEEE, 2016.
- [25] Marvin Teichmann, Michael Weber, Marius Zoellner, Roberto Cipolla, and Raquel Urtasun. Multinet: Real-time joint semantic reasoning for autonomous driving. *arXiv preprint arXiv:1612.07695*, 2016.
- [26] Ishan Misra, Abhinav Shrivastava, Abhinav Gupta, and Martial Hebert. Cross-stitch networks for multi-task learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3994–4003, 2016.
- [27] Jonas Uhrig, Marius Cordts, Uwe Franke, and Thomas Brox. Pixel-level encoding and depth layering for instance-level semantic labeling. *arXiv preprint arXiv:1604.05096*, 2016.
- [28] Alex Kendall, Matthew Grimes, and Roberto Cipolla. Convolutional networks for real-time 6-dof camera relocalization. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2015.
- [29] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *arXiv preprint arXiv:1703.04977*, 2017.
- [30] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *In Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2016.
- [31] Xiaodan Liang, Yunchao Wei, Xiaohui Shen, Jianchao Yang, Liang Lin, and Shuicheng Yan. Proposal-free network for instance-level object segmentation. *arXiv preprint arXiv:1509.02636*, 2015.
- [32] Bastian Leibe, Aleš Leonardis, and Bernt Schiele. Robust object detection with interleaved categorization and segmentation. *International Journal of Computer Vision (IJCV)*, 77(1-3):259–289, 2008.
- [33] Mihael Ankerst, Markus M Breunig, Hans-Peter Kriegel, and Jörg Sander. Optics: ordering points to identify the clustering structure. In *ACM Sigmod Record*, volume 28, pages 49–60. ACM, 1999.
- [34] James MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA., 1967.
- [35] Dorin Comaniciu and Peter Meer. Mean shift: A robust approach toward feature space analysis. *IEEE Transactions on pattern analysis and machine intelligence*, 24(5):603–619, 2002.
- [36] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *In Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2016.
- [37] Heiko Hirschmuller. Stereo processing by semiglobal matching and mutual information. *IEEE Transactions on pattern analysis and machine intelligence*, 30(2):328–341, 2008.
- [38] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *arXiv preprint arXiv:1606.00915*, 2016.

A Scene Understanding Model Details

The purpose of the encoder is to learn a deep mapping to produce rich, contextual features, using domain knowledge from a number of related tasks. Our encoder is based on ResNet-101 [30] (without the final fully connected layer). We apply this encoder in a convolutional manner over the input image, which results in a 2048 dimensional shared feature representation. Inspired by the dilated convolutional approach of [38], this encoder feature map is sub-sampled by a factor of 8 compared to the input image dimensions.

We then split the network into separate decoders (with separate weights) for each task. The purpose of the decoder is to learn a mapping from the shared features to an output. Each decoder consists of three convolutional layers with kernel size 3×3 , 1×1 and 1×1 respectively, and feature size 512, 512 and the number of output dimensions respectively.

B Training Convergence Results

One of the attractive properties of our approach to weighting multi-task losses is that it is robust to the initialisation choice for the homoscedastic noise parameters. Figure 4 shows that for an array of initial choices of $\log \sigma^2$ from -2.0 to 5.0 the homoscedastic noise and task loss is able to converge to the same minima. Additionally, the homoscedastic noise terms converges after only 100 iterations, while the network requires 30,000+ iterations to train.

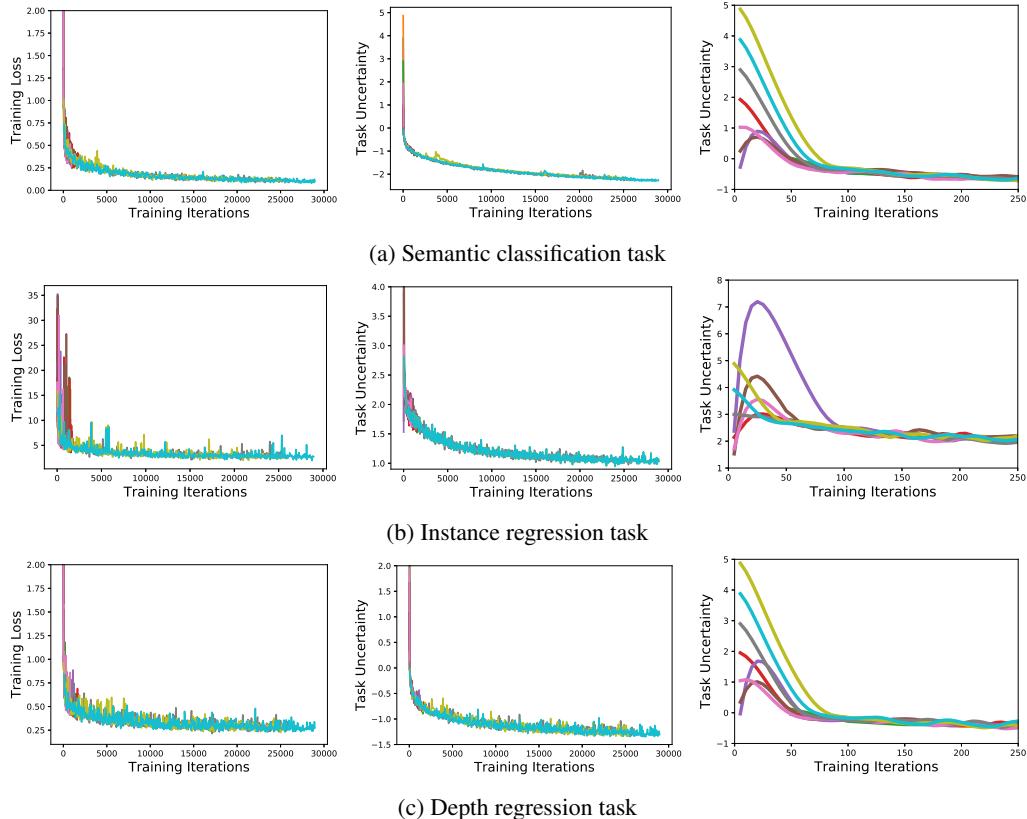


Figure 4: Training plots showing convergence of homoscedastic noise and task loss for an array of initialisation choices for the homoscedastic uncertainty terms for all three tasks. The left plot shows that the loss converges to the same minimum from varying initialisation choices. The centre plot shows the the homoscedastic noise value optimises to the same solution from a variety of initialisations. The plots on the right show a zoomed in view of the homoscedastic noise plot, showing the initialisation and convergence over a few hundred training iterations. Despite the network taking 10,000+ iterations for the training loss to converge, the task uncertainty converges very rapidly after only 100 iterations.

C Instance Segmentation Parametrisation with Centroid Vectors

Figure 5 details the representation we use for instance segmentation. Figure 5(a) shows the input image and a mask of the pixels which are of an instance class (at test time inferred from the predicted semantic segmentation). Figure 5(b) and Figure 5(c) show the ground truth and predicted instance vectors for both x and y coordinates. We then cluster these votes using OPTICS [33], resulting in the predicted instance segmentation output in Figure 5(d).

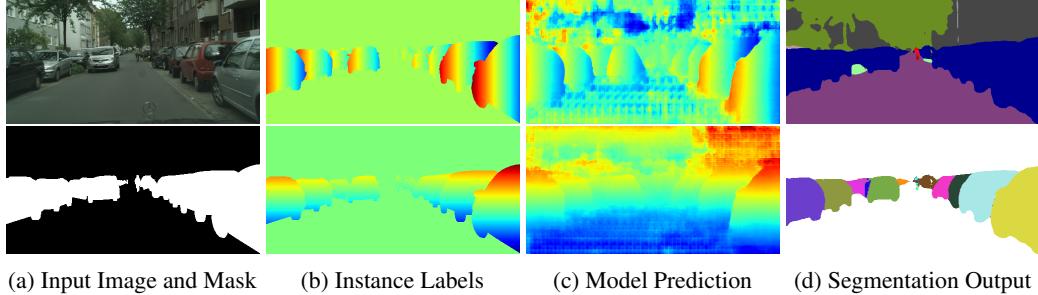


Figure 5: **Instance centroid regression training data.** For each pixel, we regress a vector pointing to the instance’s centroid. The loss is computed with the mask, and x and y instance centroid vectors.

One of the most difficult cases for instance segmentation algorithms to handle is when the instance mask is split due to occlusion. Figure 6 shows that our method can handle these situations, by allowing pixels to vote for their instance centroid with geometry. Methods which rely on watershed approaches [15], or instance edge identification approaches fail in these scenarios.

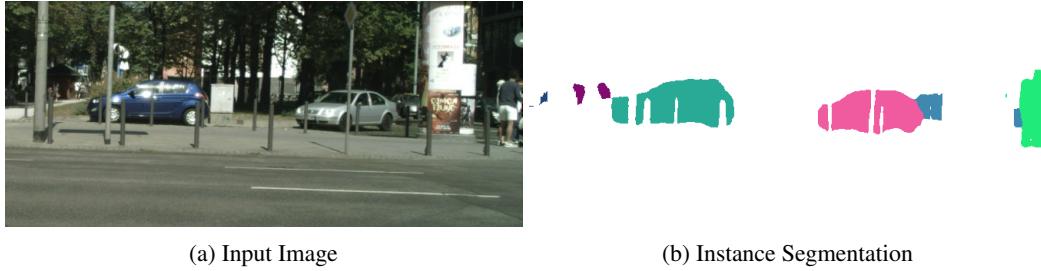


Figure 6: This example shows two cars which are occluded by trees and lampposts, making the instance segmentation challenging. Our instance segmentation method can handle occlusions effectively. We can correctly handle segmentation masks which are split by occlusion, yet part of the same instance. By incorporating semantics and geometry.

D Further Qualitative Results

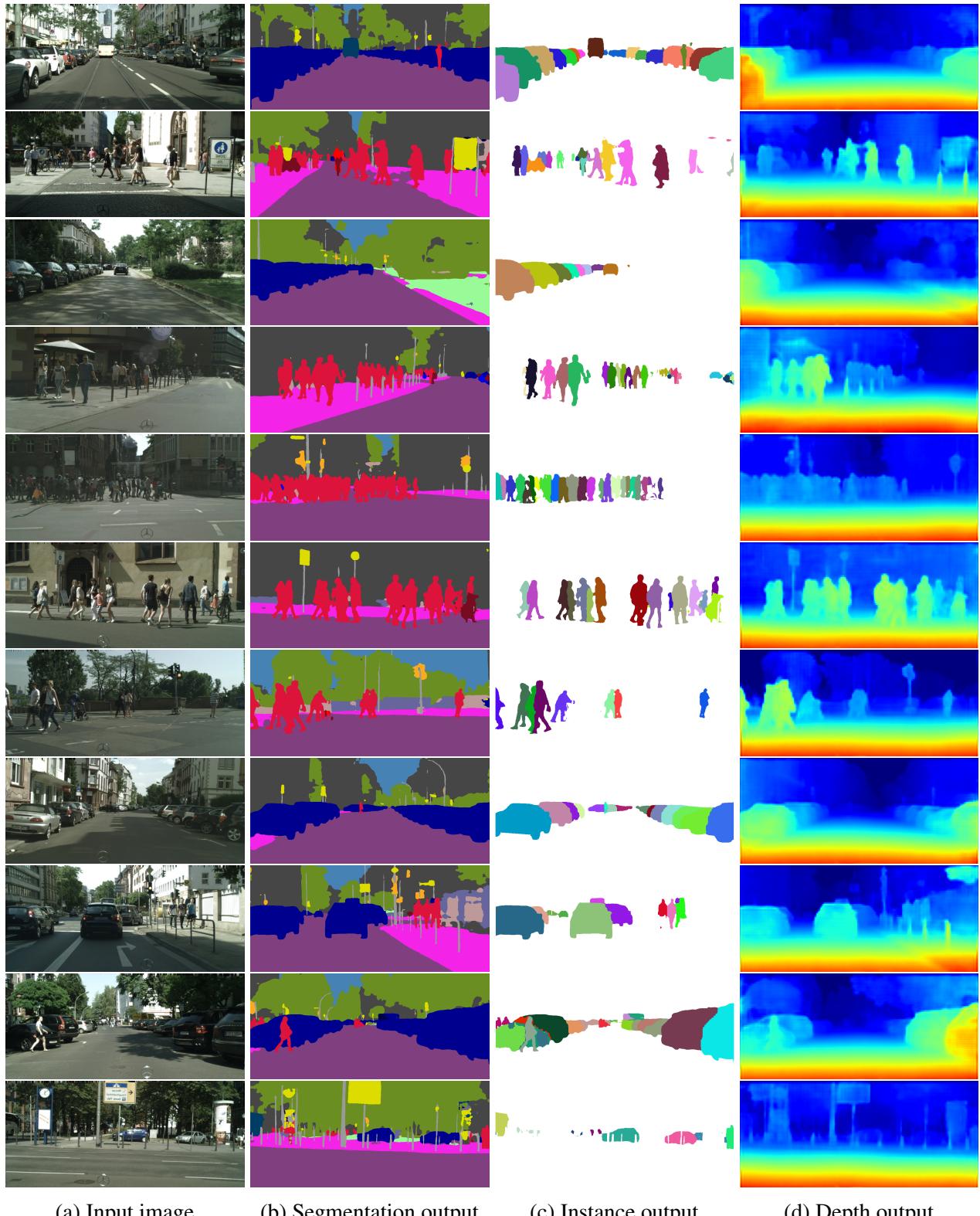


Figure 7: More qualitative results on test images from the CityScapes dataset.