

# See Yourself in Others: Attending Multiple Tasks for Own Failure Detection

Boyang Sun\*, Jiaxu Xing\*, Hermann Blum, Roland Siegwart, and Cesar Cadena

**Abstract**—Autonomous robots deal with unexpected scenarios in real environments. Given input images, various visual perception tasks can be performed, e.g., semantic segmentation, depth estimation and normal estimation. These different tasks provide rich information for the whole robotic perception system. All tasks have their own characteristics while sharing some latent correlations. However, some of the task predictions may suffer from the unreliability dealing with complex scenes and anomalies. We propose an attention-based failure detection approach by exploiting the correlations among multiple tasks. The proposed framework infers task failures by evaluating the individual prediction, across multiple visual perception tasks for different regions in an image. The formulation of the evaluations is based on an attention network supervised by multi-task uncertainty estimation and their corresponding prediction errors. Our proposed framework<sup>1</sup> generates more accurate estimations of the prediction error for the different task’s predictions.

## I. INTRODUCTION

Extensive research has shown that visual information is an important component in autonomous driving and many robotic perception systems [1]–[3]. Autonomous agents utilize the information from various learning-based visual perception predictions. Existing works have shown good performance on cases where the deployment environment has similar distribution to the training set [4]. However, many state-of-the-art deep learning approaches still face the lack of ability in dealing with open and unconstrained world [5]–[7], and will produce failures especially in unseen environments [8]. Thus, a method to detect prediction failures of various robotics visual perception tasks is crucial for safe robotic deployments. With higher introspection capabilities, autonomous robots will be more controllable in safety-critical scenarios.

This work focus on identifying failure predictions of various robotics perception tasks by exploiting the latent correlations among them. Those correlations have been recently used to improve tasks performance [9,10]. Our basic idea is to exploit the complementary information from multiple tasks to improve the introspection capability of the perception system on every single task based on an attention mechanism. Our failure detection model has a *unified* structure that attends the encoded multi-task feature

\* Authors contributed equally to this work. This work was partially supported by the Hilti Group and the National Center of Competence in Research (NCCR) Robotics through the Swiss National Science Foundation.

All the authors are with the Autonomous Systems Lab, ETH Zurich, 8092, Switzerland. {boysun, jixing, blumh, rsiegwart, cesarc}@ethz.ch

<sup>1</sup>Code link [https://github.com/ethz-asl/uncertainty\\_with\\_multiple\\_tasks](https://github.com/ethz-asl/uncertainty_with_multiple_tasks).

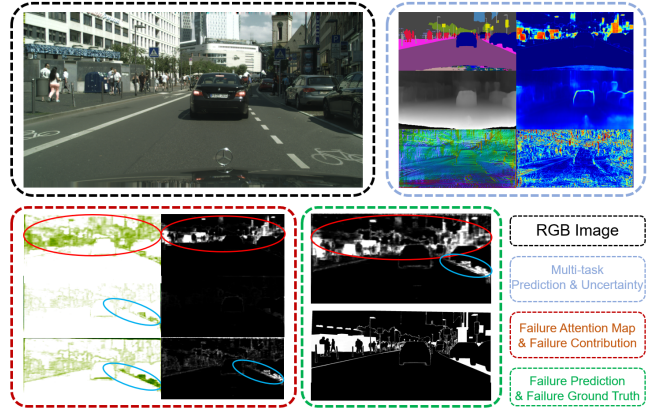


Fig. 1: **Example Result of Our Approach.** Our method captures the high prediction uncertainty regions of a single task using multiple visual tasks. The result maintains the useful uncertainty estimation from the original task (highlighted areas in red circle). Moreover, beneficial from the multi-task setup, our approach captures the relevant information from other tasks (highlighted areas in blue circle) to compensate the missed failure regions.

maps with the *expressive power* to perform failure prediction for different tasks.

Existing research recognizes uncertainty as a common measurement of the multi-task prediction’s confidence [8]. The uncertainties are ideally correlated to the corresponding task prediction errors, which measures the reliability of the predictions. General uncertainty estimation methods are based on a single task, e.g., softmax entropy from semantic segmentation [11]. However, the quality of the uncertainty estimation is limited by several factors, such as environmental conditions (e.g., clean/foggy weather), anomalies, and more [12]. In addition, to avoid the limitation brought by single task uncertainty estimation, our model takes various tasks encoded predictions as input and computes an attention map for each single task. The attention values are modeled as the regional contribution of every task uncertainty to one certain task prediction error.

To investigate the model robustness, we train our model on the Cityscapes [13] dataset and test it on several datasets with different distribution characteristics [14,15]. We evaluate the model for different tasks and compare it against several existing methods. We show that our approach outperforms all other failure prediction approaches. Moreover, our framework is flexible to the number and types of tasks with different task prediction & uncertainty estimation methods. A result example of our approach is shown in Figure 1.

In summary, the contributions of this work are:

- The first work to exploit the multiple visual tasks setup for detecting failures in their prediction in deployment.
- A novel framework with an attention mechanism over the multiple visual tasks being deployed to extract the complementary information in their uncertainty estimates for failure detection.
- A thorough evaluation of design components and their influence in open world scenarios.

## II. RELATED WORK

Works on single task failure prediction can be classified in failure detection from the estimated uncertainty and learning-based failure detection. Other works have exploited the multi-task setup to achieve better per-task prediction accuracy. However, to the best of our knowledge, failure prediction from multiple tasks has not been proposed prior to this work.

### Failure Detection via Uncertainty Estimation

Uncertainty estimation has a close relation with failure detection and introspection. The uncertainty of the prediction results reflects the level of its confidence. And the intuition follows that a low-confidence prediction is likely a failure. Therefore, uncertainty estimation could be regarded as a reference to failures prediction. A conventional way to calculate the uncertainty is to directly analyze the distribution of the model prediction such as the *softmax entropy* [16], or *softmax distance* [17] used in classification models. Besides, *image flipping* [5] investigates the model results' difference in dealing with the original and flipped image. Finally, *Bayesian estimation* [18]–[20] estimates the uncertainty by sampling multiple models, e.g. *Monte-Carlo dropout* [8,21] captures the uncertainty by randomly dropping the connection between different layers.

### Failure Detection via Learning-based methods

Given the rapid development of deep learning. Neural network has become a possible option for failure detection task. Most of the failure detection methods in the visual perception area focus on detect semantic segmentation miss-classification. These methods can be roughly divided into two categories. One group directly trains detectors with failure cases [22]–[24]. The other group uses re-synthesis methods [12,23,25,26] that rebuild image from semantic prediction, and capture the anomalies by comparing the rebuilt image and the original one.

### Learning from Multiple Tasks

Prior works have already acknowledged the relation among different perception tasks [27,28]. This latent correlated structure among visual tasks has been exposed in the work of *Taskonomy* [9]. The utilization of cross-task relations also lies in the area of domain adaptation [29,30], transfer learning [10], and multi-task learning [31]–[33]. More specifically, recent attention has focused on using cross-task supervised learning to improve the performance of a single task, such as to improve depth prediction under the supervision of semantic understanding [34,35].

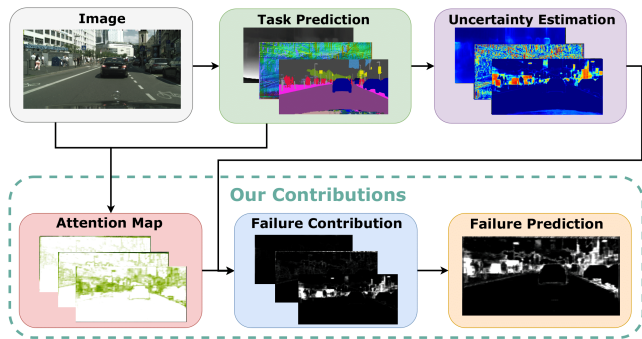


Fig. 2: **Visual Tasks Failure Detection Framework.** Given an image and its multiple tasks predictions, our approach computes the attention maps to weigh the multiple tasks uncertainty estimations. This weighted sum of attention and uncertainty maps is our failure prediction for a chosen task.

In this work, we draw inspiration from multi-task learning and single-task failure detection approaches. Our approach utilizes every single task's estimation information, exploiting their latent correlations, and infer a more robust per-task pixel-wise failure prediction.

## III. METHOD

### A. Approach Architecture

The intuition of our proposed method comes from the fact that tasks can compensate and refine each other's prediction in a multi-task setup [9,36]. Therefore, it will be easier to identify a certain task prediction failure with the existence of some other tasks. However, unlike a simple high-level task selection case, the proposed approach generates a pixel-wise attention map for each participant task. For clarification, the *Attention* used in this work is a scalar-product and it does not refer to the Transformer Networks [37] in computer vision community. These maps are used to compute a weighted sum of each task's own failure prediction. Starting with a query image, firstly, predictions of all participant tasks are generated. Secondly, these predictions and the query image are passed into our attention network. The attention network will predict a pixel-wise attention map for each task. Meanwhile, uncertainties<sup>2</sup> of each task prediction are produced. In the final step, the attention maps are applied to these uncertainties to calculate a weighted sum failure prediction for a specific task. The weighted sum result for this task is expected to have better performance, comparing with its uncertainty estimation method. Figure 2 shows the basic structure of this aforementioned procedure.

### B. Multi-task Element Generation

Our approach admits any number of visual perception tasks that provide per-pixel predictions and their per-pixel uncertainties. We rely on open-source state-of-the-art task prediction methods. As for the uncertainty estimation, many

<sup>2</sup>We use the term uncertainty here to cover all the statistical measurement of prediction reliability, which is often called uncertainty, reliability or dispersion metrics in other works.

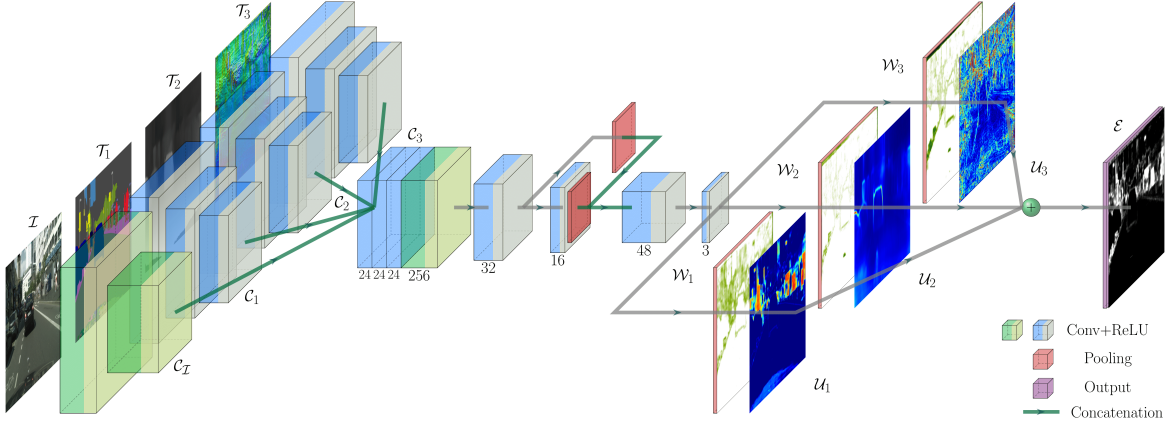


Fig. 3: **Example of the Model Architecture:** This example model uses three different tasks: semantic segmentation, depth, and normal estimation. For the uncertainties  $\mathcal{U}_1, \mathcal{U}_2, \mathcal{U}_3$  and prediction errors  $\mathcal{E}$ , we resize them to  $256 \times 256$ . Using the predefined attention patch size  $p$ , the output attention from the model  $\mathcal{W}_1, \mathcal{W}_2, \mathcal{W}_3$  will have the size  $(256/p) \times (256/p)$ . Then the output would be equally upscaled by a factor of  $p$  so that the attention maps' sizes are  $256 \times 256$ . Now the computed attention maps have the same size as the uncertainties, then element-wise multiplication can be performed.

options are available. Softmax entropy [16] and softmax distance [17] can be used for classification tasks, such as semantic segmentation. Sampling-based methods can be used for both classification and regression tasks. Even, the learning-based failure prediction approach for a single task can play the role of uncertainty estimation and be feed into our attention network. Different chosen uncertainty estimation methods will certainly influence the final output, and thus we evaluate their influence in Section IV.

### C. Attention Network Model

We denote the original image as  $\mathcal{I}$ . The predictions of all  $n$  tasks are denoted as  $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_n$ . The uncertainty estimations of them are denoted as  $\mathcal{U}_1, \mathcal{U}_2, \dots, \mathcal{U}_n$ , respectively. The architecture of our attention network is shown in Figure 3. The model first encodes the original images  $\mathcal{I}$  and its task predictions  $\{\mathcal{T}_i\}, i \in \{1, 2, \dots, n\}$ . We chose to encode the image  $\mathcal{I}$  with the first several layers of ResNet50 [38] into a 256-channel feature with a  $128 \times 128$  size. On the other hand, the predictions  $\{\mathcal{T}_i\}$  are encoded by part of MobileNetV2 [39]. Each of them is encoded into a 24-channel feature map of size  $128 \times 128$ . All tasks prediction,  $\{\mathcal{T}_i\}$ , share the same encoding structure. The encoded feature maps are denoted as  $\mathcal{C}_{\mathcal{I}}, \{\mathcal{C}_i\}$ , correspondingly.

After the encoding process, the encoded features are concatenated along the channel dimension. The resulting feature map,  $\mathcal{C}_{cat}$  is in the size of  $(256 + 24n) \times 128 \times 128$ .  $\mathcal{C}_{cat}$  is then forwarded into a neural network with four convolutional layers. In these convolutional layers, the output of the second and third layers will be concatenated together and used as the input to the last layer. In this case, the output layer has one channel for each task. Given a predefined patch size  $p$ , two pooling layers are added after the second and the third convolutional layers to resize each output channel into  $(256/p) \times (256/p)$ , which is the resolution of our attention maps. The higher the required attention resolution, the larger the resized map size of each channel. An extra nearest

neighbour rescaling layer is added here to rescale each channel to the size of the uncertainty maps. The rescaled feature maps in each channel are the final attention map generated by our attention network. We denote them as  $\mathcal{W}_1, \mathcal{W}_2, \dots, \mathcal{W}_n$ , for each task, respectively. The final failure prediction for a single task is generated by calculating the weighted sum of all tasks' uncertainty estimates, with the attention maps as the weights. The weighted sum failure prediction is denoted as  $\mathcal{E}$ .

$$\mathcal{E} = \sum_{i=1}^n \mathcal{W}_i \cdot \mathcal{U}_i \quad (1)$$

### D. Training Procedure

As introduced in the last subsection, the attention maps are predicted by our attention network:

$$\{\mathcal{W}_i\} = f_{\theta}(\mathcal{I}, \{\mathcal{T}_i\}) \quad (2)$$

We set a single error approximation loss function to learn the network parameters  $\theta$ . When training the network to learn a certain task failure, we compute the pixel-wise prediction error for this task, which is denoted as  $\epsilon_{\{\cdot\}}$ . For example,  $\epsilon_S$  is computed as the cross-entropy for semantic segmentation, and  $\epsilon_D$  is simply calculated as the L2-norm between depth prediction and ground truth depth. The general loss computation is shown as Equation 3.

$$loss_{\{\cdot\}} = \left\| \sum_{i=1}^n (\mathcal{W}_i \cdot \mathcal{U}_i) - \epsilon_{\{\cdot\}} \right\| \quad (3)$$

At last, to prevent any task's error from dominating the prediction due to the imbalanced scales, we perform an image-wise normalization process to re-scale all the  $\mathcal{C}_{cat}$  in both training and inference stages into range  $[0, 1]$ .

Task	Prediction Method	Uncertainty Estimation Methods
Semantic Segmentation	SDC Net [7]	Softmax Entropy [11] Softmax Distance [17] Synboost* [12] MC Dropout [40]
Depth Estimation	Monodepth V2 [6]	Bayesian Estimation [20] MC Dropout [40] Self Learning* [21]
Normal Estimation	VNL [41]	Flipping* [5]
Instance Segmentation	EfficientPS [42]	ROI Softmax Uncertainty* [43]

TABLE I: The selected task prediction methods and the uncertainty estimation methods for all different tasks in the experiments. \* indicates the method used by default in the experiments unless otherwise mentioned.

## IV. EXPERIMENTS

### A. Experimental Setup

1) *Model Implementation*: Inspired by task networks graph in *taskonomy* [9], we decided to choose tasks from three visual tasks in total in our different experiments: semantic segmentation, depth estimation, and normal estimation. We add later in the experiments a fourth task, instance segmentation. For each task, we implemented publicly available task prediction methods and uncertainty estimation methods. All evaluated methods are shown in Table I.

2) *Dataset*: Training was performed on Cityscapes training dataset [13], including 2975 driving scenes images with fine semantic annotations and disparity ground truth. We produced our dataset by applying the methods shown in the Table I. The training set is then composed by set in form of  $\{\mathcal{I}, \mathcal{T}_S, \mathcal{T}_D, \mathcal{T}_N, \mathcal{U}_S, \mathcal{U}_D, \mathcal{U}_N, \epsilon_{\{\cdot\}}\}$ , where  $S, D, N$  denote semantic segmentation, depth estimation, and normal estimation, tasks, respectively, and  $\epsilon_{\{\cdot\}}$  correspond to the prediction error of the chosen task to predict its failure.

To test our model’s performance, pre-processing of the various test datasets is also required. Here we performed the same pipeline as mentioned in the subsection IV-A.2 for Cityscapes validation set [13], Foggy Cityscapes validation set [44], Wilddash [15] and Dark Zurich dataset [14]. Wilddash provides a dataset and benchmark for challenging driving scenarios under real-world conditions, it contains scenarios from very diverse environments, locations, and weather conditions. Dark Zurich is a dataset designed for semantic uncertainty-aware model evaluations. It contains driving scenes images captured at night time, twilight and day time. Here we only use the night time images for our evaluation. The purpose of testing on these extra two datasets is to validate the model robustness when dealing with the *challenging unseen scenarios*.

3) *Metrics*: We choose the zero-mean normalized cross-correlation (ZNCC) in our experiments as a measurement of how close the predicted failure is to the ground-truth failure. The ZNCC of the estimation  $\mathcal{E}$  and the ground truth error  $\epsilon$

Task Entries			Semantic				Depth
$S$	$D$	$N$	ZNCC $\uparrow$	AP-Err $\uparrow$	AP-Suc $\uparrow$	FPR95 $\downarrow$	ZNCC $\uparrow$
✓	✓	✓	<b>0.649</b>	<b>0.590</b>	0.987	0.280	<b>0.646</b>
✓	✓		0.609	0.545	<b>0.990</b>	<b>0.278</b>	0.489
✓			0.494	0.413	0.978	0.570	-
	✓		-	-	-	-	0.483

TABLE II: Multiple Tasks Experiments: Comparison among multiple different task entries for both semantic segmentation’s and depth estimation’s failure prediction. The checkmark represents the task at the corresponding place is included in the model, both during training and evaluation process.

is defined as:

$$\text{ZNCC} = \frac{\sum_{(u,v)} (\mathcal{E}(u,v) - \mu_{\mathcal{E}})(\epsilon(u,v) - \mu_{\epsilon})}{\sqrt{\sum_{(u,v)} (\mathcal{E}(u,v) - \mu_{\mathcal{E}})^2} \sqrt{\sum_{(u,v)} (\epsilon(u,v) - \mu_{\epsilon})^2}} \quad (4)$$

where  $\mu_{\mathcal{E}}$  and  $\mu_{\epsilon}$  denote the mean values of  $\mathcal{E}, \epsilon$ , respectively, and  $(u, v)$  are the pixel locations.

ZNCC is invariant of affine pixel value changes [45]. Therefore, it is not be influenced by the normalization operation during the prediction and evaluation procedure, and it is also less sensitive to resizing operations than normalized cross-correlation (NCC). This metric is seamlessly applicable to classification and regression tasks.

In addition, for the classification task of semantic segmentation, previous works on failure detection have used several metrics for evaluation [12,24,26]. Thus, we also report:

- AUPR-Error: the area under the Precision-Recall (AUPR) curve, which regard incorrect prediction as the positive class.
- AUPR-Success: it computes AUPR as well, whereas treats correct prediction as the positive class.
- FPR95: the false positive rate at 95% true positive rate.

As for regression tasks, such as depth estimation, we are not aware of any previous work focusing on evaluating the failure prediction model.

### B. Comparisons

Our framework is built on a multi-task setup, and uses different single task’s uncertainty estimation methods as the input. Thus, most of the our evaluations focus on relative comparison between our model’s output and the input it uses, or among the models with different configurations. We set up several experiments to evaluate different components or factors and be able to answer a series of questions.

**Are multiple tasks beneficial for single task failure detection?** In these experiments we evaluate two aspects. The first one is the influence of increasing the number of tasks as inputs to our failure detection. And the second one, we evaluate our failure detection on different main tasks: a classification task (semantic segmentation) and a regression task (depth estimation).

We start in our network structure by only having the input of a single task (the task for which the failure is being detected), see last 2 rows of Table II, on the Cityscapes original validation set. This is equivalent to learning failures

from single task knowledge and thus can be compared to use the uncertainty input as a proxy to the failure.

Then, we continue by adding a second task (depth or semantics, as appropriate), and a third one (normal estimation), see first row in Table II.

From this experiment we have evidence that, indeed, a multi-task setup improves failure detection of one task, and, this is true for both semantic segmentation and depth estimation. This is a confirmation of our hypothesis that our framework leverages the latent correlations among tasks to improve the introspection capabilities of the every single one of them. It is also in line of the findings of [9] where is shown that the normal estimation task is more correlated to the depth estimation than to the semantic segmentation one, as seeing in the higher increase in performance of the depth estimation failure detection compared to that improvement the semantic segmentation failure detection.

**At what resolution should the attention maps be computed to capture the relevant information for failure detection?** In our network architecture, we have the choice to generate the attention maps at different resolutions by modifying the last pooling layers. Something that could be beneficial for differentiating between full image task failure or per pixel or region task failure detection.

Method	Patch	Original		Foggy	
		ZNCC↑	AP-Err↑	ZNCC↑	AP-Err↑
Ours	1	<b>0.649</b>	<b>0.590</b>	<b>0.560</b>	<b>0.518</b>
	2	0.601	0.520	0.556	0.495
	4	0.570	0.487	0.530	0.475
	8	0.529	0.439	0.530	0.470
	16	0.489	0.400	0.516	0.466
	32	0.455	0.368	0.500	0.450
	64	0.440	0.356	0.496	0.448
	128	0.437	0.357	0.501	0.454
SynBoost	-	0.450	0.387	0.506	0.480

TABLE III: Variable Attention Map Resolution: Comparison among multiple different patch sizes for semantic segmentation’s failure prediction.

Method	Patch	Original	Foggy
		ZNCC↑	ZNCC↑
Ours	1	<b>0.646</b>	0.570
	2	0.638	<b>0.573</b>
	4	0.627	0.560
	8	0.611	0.537
	16	0.569	0.529
	32	0.455	0.445
	64	0.317	0.341
	128	0.279	0.282
Self Learning	-	0.248	0.255

TABLE IV: Variable Attention Map Resolution: Comparison among multiple different patch size for depth estimation’s failure prediction.

We modify the patch size of our model from 1 to 128, while the generated attention map is up-scaled to size  $256 \times 256$ . From the evaluation results in Tables III and IV, we can conclude that with higher resolution the model is able to predict failures more accurately, for it has higher AP-Err, and higher ZNCC for both task failures. Here we included

as well the Foggy dataset to check the performance in more challenging scenarios.

**How dependant is our failure detection on the uncertainty estimate input?** The uncertainty of each task is an important component of our failure detection framework. Thus, this experiment investigates whether our conclusions have been biased to the uncertainty input used. For the semantic segmentation task we evaluate three more uncertainty inputs, and for the depth estimation another two (see Table I). For each specific uncertainty method, we select two of our models with different patch sizes (1 and 16).

Method	Patch	Cityscapes Original		Cityscapes Foggy	
		ZNCC↑	AP-Err↑	ZNCC↑	AP-Err↑
Ours with	1	<b>0.649</b>	<b>0.590</b>	<b>0.560</b>	<b>0.518</b>
Synboost	16	0.489	0.400	0.516	0.466
SynBoost	-	0.450	0.387	0.506	0.480
Ours with	1	<b>0.681</b>	<b>0.618</b>	<b>0.628</b>	<b>0.565</b>
Soft. Ent.	16	0.568	0.447	0.593	0.497
Soft. Ent.	-	0.572	0.444	0.619	0.520
Ours with	1	<b>0.576</b>	<b>0.506</b>	<b>0.430</b>	<b>0.408</b>
MC Dropout	16	0.358	0.274	0.327	0.307
MC Dropout	-	0.249	0.218	0.163	0.236
Ours with	1	<b>0.668</b>	<b>0.593</b>	<b>0.600</b>	<b>0.532</b>
Soft. Dis	16	0.540	0.409	0.574	0.479
Soft. Dis.	-	0.527	0.408	0.569	0.489

TABLE V: Effect of Changing the Semantic Uncertainty Input: Our models vs. selected uncertainty inputs for semantic estimation’s failure prediction.

Method	Patch	Original	Foggy
		ZNCC↑	ZNCC↑
Ours with	1	<b>0.646</b>	<b>0.570</b>
Self Learning	16	0.569	0.529
Self Learning	-	0.248	0.255
Ours with	1	<b>0.757</b>	<b>0.645</b>
Bayesian	16	0.684	0.584
Bayesian	-	0.091	0.074
Ours with	1	<b>0.648</b>	<b>0.516</b>
MC Dropout	16	0.575	0.446
MC Dropout	-	0.076	0.075

TABLE VI: Effect of Changing the Depth Uncertainty Inputs: Our models vs. selected uncertainty methods for depth estimation’s failure prediction.

The results of this investigation can be seen in Tables V and VI. Our framework is consistently outperforming the uncertainty input for both tasks failure detections, within both original and foggy image set from Cityscapes.

**How is the failure detection performing when adding one more task?** Finally, the question comes to our failure detection based on a multi-task setup is fixed to the already chosen tasks in the previous experiments. For that reason, we add the extra task of instance segmentation (*IS*), with its corresponding uncertainty estimate, see last row of Table I.

The results of this experiment, shown in Table VII, indicate that the addition of an extra task continues to be beneficial for the different tasks failure detections. However, we observe that the contribution is higher for the depth estimation failure detection than for the semantic segmentation. We believe this is due to less extra information provided



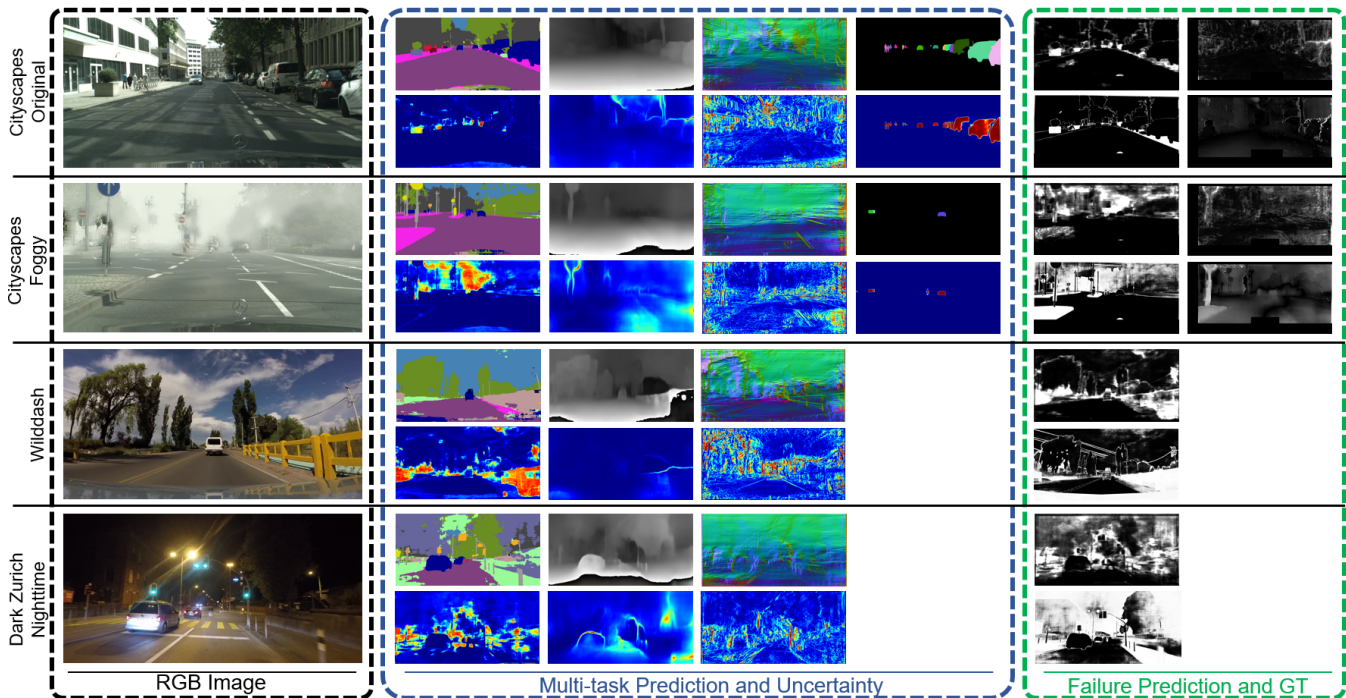


Fig. 4: **Qualitative Examples.** For an input image in the left block, the middle block contains the various visual tasks (1<sup>st</sup> row) and the corresponding uncertainties (2<sup>nd</sup> row), tasks are semantic, depth, normal, instance, from left to right. In the right block our failure predictions (1<sup>st</sup> row) and the ground truth (2<sup>nd</sup> row) are presented, semantic on the left and depth on the right. For the purpose of clear visualization, the ego-vehicle part in both Cityscapes examples are filtered out.

Patch	Task Entries				Semantic		Depth
	<i>S</i>	<i>D</i>	<i>N</i>	<i>IS</i>	ZNCC $\uparrow$	AP-Err $\uparrow$	ZNCC $\uparrow$
1	✓	✓	✓	✓	0.641	0.585	<b>0.655</b>
	✓	✓	✓	✓	<b>0.649</b>	<b>0.590</b>	0.646
16	✓	✓	✓	✓	<b>0.493</b>	<b>0.403</b>	<b>0.581</b>
	✓	✓	✓	✓	0.489	0.400	0.529

TABLE VII: Adding an Extra Task: Comparison among multiple different task entries for both semantic segmentation’s and depth estimation’s failure prediction.

by the instance segmentation with respect to the semantic segmentation task. While differentiating among different instances of the same class is highly informative for the depth estimation task.

**How is our failure detection generalizing to scenarios with larger distribution mismatch?** Here, we use our default models trained with the Cityscapes dataset, and deployed them on Wilddash and Dark Zurich (night) datasets, for the tasks of failure detection of the semantic segmentation. Results can be seen in Table VIII. Additionally, we include the evaluation of different uncertainty inputs as they are quite dependent on the distribution mismatch between the test and training set. We can conclude that the improvements brought by our framework generalize to more complex scenarios, invariant to the uncertainty estimate input.

Some qualitative results for the different datasets are visualized in figure 4.

Method	Wilddash		Dark Zurich	
	ZNCC $\uparrow$	AP-Err $\uparrow$	ZNCC $\uparrow$	AP-Err $\uparrow$
Ours with Synboost	<b>0.412</b>	<b>0.595</b>	<b>0.238</b>	<b>0.775</b>
SynBoost	0.323	0.584	0.199	0.726
Ours with Soft. Ent.	<b>0.520</b>	<b>0.630</b>	0.544	<b>0.867</b>
Soft. Ent.	0.508	0.625	<b>0.578</b>	0.830
Ours with MC Dropout	<b>0.321</b>	<b>0.537</b>	<b>0.101</b>	<b>0.734</b>
MC Dropout	0.139	0.428	-0.128	0.678
Ours with Soft. Dis	<b>0.511</b>	<b>0.628</b>	0.497	<b>0.861</b>
Soft. Dis.	0.478	0.619	<b>0.502</b>	0.797

TABLE VIII: Generalization: Our models vs. selected uncertainty inputs for semantic segmentation’s failure prediction on two other datasets: Wilddash and Dark Zurich.

## V. CONCLUSION

We propose a framework to detect visual task prediction failures. We leverage the information from multiple visual tasks simultaneously being deployed, and build a learning-based attention neural network to perform a weighted sum of task uncertainties to approximate the task prediction failure. Our approach is more accurate in detecting semantic and depth prediction errors, compared with various uncertainty estimation methods. Additionally, our thorough experimental evaluation also proves its ability to further improve the performance by increasing the attention map resolution, as well as by including in extra correlated visual tasks. Finally, we observe that the multi-task setup allows for better generalization to environments with a larger distribution mismatch to that of the training set. We believe our framework shows the potential capability to be applied on various autonomous mobile robots with multi-task visual modules, such as safety in autonomous driving.

## REFERENCES

- [1] S. Grigorescu, B. Trasnea, T. Cocias, and G. Macesanu, "A survey of deep learning techniques for autonomous driving," *Journal of Field Robotics*, vol. 37, no. 3, p. 362–386, Apr 2020. [Online]. Available: <http://dx.doi.org/10.1002/rob.21918>
- [2] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, "Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age," *IEEE Transactions on Robotics*, vol. 32, no. 6, pp. 1309–1332, 2016.
- [3] S. Garg, N. Sünderhauf, F. Dayoub, D. Morrison, A. Cosgun, G. Carneiro, Q. Wu, T.-J. Chin, I. Reid, S. Gould, and et al., "Semantics for robotic mapping, perception and interaction: A survey," *Foundations and Trends® in Robotics*, vol. 8, no. 1–2, p. 1–224, 2020. [Online]. Available: <http://dx.doi.org/10.1561/23000000059>
- [4] S. Bulusu, B. Kaikhura, B. Li, P. K. Varshney, and D. Song, "Anomalous example detection in deep learning: A survey," *IEEE Access*, vol. 8, pp. 132 330–132 347, 2020.
- [5] C. Godard, O. Mac Aodha, and G. J. Brostow, "Unsupervised monocular depth estimation with left-right consistency," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 270–279.
- [6] C. Godard, O. M. Aodha, M. Firman, and G. Brostow, "Digging into self-supervised monocular depth estimation," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 3827–3837.
- [7] Y. Zhu, K. Sapra, F. Reda, K. J. Shih, S. Newsam, A. Tao, and B. Catanzaro, "Improving semantic segmentation via video propagation and label relaxation," *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8848–8857, 2019.
- [8] Y. Gal, "Uncertainty in deep learning," Ph.D. dissertation, University of Cambridge, 2016.
- [9] A. R. Zamir, A. Sax, W. Shen, L. J. Guibas, J. Malik, and S. Savarese, "Taskonomy: Disentangling task transfer learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3712–3722.
- [10] A. R. Zamir, A. Sax, N. Cheerla, R. Suri, Z. Cao, J. Malik, and L. J. Guibas, "Robust learning through cross-task consistency," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [11] D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," *arXiv preprint arXiv:1610.02136*, 2016.
- [12] G. Di Biase, H. Blum, R. Siegwart, and C. Cadena, "Pixel-wise anomaly detection in complex driving scenes," *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [13] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [14] C. Sakaridis, D. Dai, and L. Van Gool, "Map-guided curriculum domain adaptation and uncertainty-aware evaluation for semantic nighttime image segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [15] O. Zendel, K. Honauer, M. Murschitz, D. Steininger, and G. F. Dominguez, "Wilddash - creating hazard-aware benchmarks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.
- [16] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," *Neural Information Processing Systems (NeurIPS)*, 2017.
- [17] M. Rottmann, P. Colling, T. Hack, R. Chan, F. Hüger, P. Schlicht, and H. Gottschalk, "Prediction error meta classification in semantic segmentation: Detection via aggregated dispersion measures of softmax probabilities," 07 2020, pp. 1–9.
- [18] A. Kendall and Y. Gal, "What uncertainties do we need in bayesian deep learning for computer vision?" *Neural Information Processing Systems (NeurIPS)*, 2017.
- [19] E. Ilg, O. Cicek, S. Galesso, A. Klein, O. Makansi, F. Hutter, and T. Brox, "Uncertainty estimates and multi-hypotheses networks for optical flow," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 652–667.
- [20] R. M. Neal, *Bayesian learning for neural networks*. Springer Science & Business Media, 2012, vol. 118.
- [21] M. Poggi, F. Aleotti, F. Tosi, and S. Mattoccia, "On the uncertainty of self-supervised monocular depth estimation," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [22] D. Hendrycks, M. Mazeika, and T. Dieterich, "Deep anomaly detection with outlier exposure," *arXiv preprint arXiv:1812.04606*, 2018.
- [23] C. B. Kuhn, M. Hofbauer, Z. Xu, G. Petrovic, and E. Steinbach, "Pixel-wise failure prediction for semantic video segmentation," in *2021 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2021, pp. 614–618.
- [24] Q. Marufur Rahman, N. Sünderhauf, P. Corke, and F. Dayoub, "Fsnnet: A failure detection framework for semantic segmentation," *arXiv e-prints*, pp. arXiv–2108, 2021.
- [25] K. Lis, K. Nakka, P. Fua, and M. Salzmann, "Detecting the unexpected via image resynthesis," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2152–2161.
- [26] D. Haldimann, H. Blum, R. Siegwart, and C. Cadena, "This is not what i imagined: Error detection for semantic segmentation through visual dissimilarity," *arXiv preprint arXiv:1909.00676*, 2019.
- [27] S. Ben-David and R. Borbely, "A notion of task relatedness yielding provable multiple-task learning guarantees," *Machine Learning*, vol. 73, pp. 273–287, 12 2008.
- [28] I. Kokkinos, "Ubertnet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5454–5463.
- [29] Y. Aytar and A. Zisserman, "Tabula rasa: Model transfer for object category detection," in *2011 international conference on computer vision*. IEEE, 2011, pp. 2252–2259.
- [30] V. M. Patel, R. Gopalan, R. Li, and R. Chellappa, "Visual domain adaptation: A survey of recent advances," *IEEE signal processing magazine*, vol. 32, no. 3, pp. 53–69, 2015.
- [31] X. Wang, K. He, and A. Gupta, "Transitive invariance for self-supervised visual representation learning," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1329–1338.
- [32] A. Pentina and C. H. Lampert, "Multi-task learning with labeled and unlabeled tasks," in *International Conference on Machine Learning*. PMLR, 2017, pp. 2807–2816.
- [33] I. Misra, A. Shrivastava, A. Gupta, and M. Hebert, "Cross-stitch networks for multi-task learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3994–4003.
- [34] V. Guizilini, R. Hou, J. Li, R. Ambrus, and A. Gaidon, "Semantically-guided representation learning for self-supervised monocular depth," *International Conference on Learning Representations*, 2020.
- [35] P.-Y. Chen, A. H. Liu, Y.-C. Liu, and Y.-C. F. Wang, "Towards scene understanding: Unsupervised monocular depth estimation with semantic-aware representation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [36] T. Standley, A. Zamir, D. Chen, L. Guibas, J. Malik, and S. Savarese, "Which tasks should be learned together in multi-task learning?" in *Proceedings of the 37th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, H. D. III and A. Singh, Eds., vol. 119. PMLR, 13–18 Jul 2020, pp. 9120–9132. [Online]. Available: <https://proceedings.mlr.press/v119/standley20a.html>
- [37] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Neural Information Processing Systems (NeurIPS)*, 2017.
- [38] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [39] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.
- [40] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," *Proceedings of The 33rd International Conference on Machine Learning*, 06 2015.
- [41] W. Yin, Y. Liu, C. Shen, and Y. Yan, "Enforcing geometric constraints of virtual normal for depth prediction," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 5683–5692.
- [42] R. Mohan and A. Valada, "Efficientps: Efficient panoptic segmentation," *International Journal of Computer Vision (IJCV)*, vol. 129, pp. 1551 – 1579, 2021.

- [43] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2980–2988.
- [44] C. Sakaridis, D. Dai, and L. Van Gool, "Semantic foggy scene understanding with synthetic data," *International Journal of Computer Vision*, vol. 126, no. 9, pp. 973–992, Sep 2018. [Online]. Available: <https://doi.org/10.1007/s11263-018-1072-8>
- [45] L. Di Stefano, S. Mattoccia, and F. Tombari, "Zncc-based template matching using bounded partial correlation," *Pattern Recognition Letters*, vol. 26, no. 14, pp. 2129–2134, 2005. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0167865505000905>