

---

# Reproducibility of ”Pixel-wise Anomaly Detection in Complex Driving Scenes ” for ML Reproducibility Challenge 2021

---

Anonymous Author(s)

Affiliation

Address

email

## Reproducibility Summary

1

2 The following paper is a reproducibility report for **Pixel-wise Anomaly Detection in Complex Driving Scenes**[4]  
3 published in CVPR 2021 as part of the ML Reproducibility Challenge 2021. We reproduced the results quantitatively,  
4 performed ablation studies, re-implemented the model in PyTorch Lightning and integrated WandB[1].  
5 **Links:** [PyTorch Repository](#), [PyTorch Lightning Repository](#), & [Our WandB Report](#).

### 6 Scope of Reproducibility

7 Our efforts are focused on validating the authors’ proposed anomaly segmentation framework, which employs the latest  
8 re-synthesis approaches[8] and extends them to incorporate the advantages of uncertainty estimation methods[6] [7][11].  
9 This proposed model outperforms existing re-synthesis techniques by a significant margin on the task of anomaly  
10 segmentation on the Fishyscapes dataset.

### 11 Methodology

12 We initially re-implemented the dissimilarity module in PyTorch Lightning using the authors’ publicly available  
13 source code. PyTorch Lightning increases the readability, reproducibility, and robustness of the code. It also provides  
14 distributed training. We used pre-trained weights for image segmentation [14], image reconstruction[9] [13] and trained  
15 the dissimilarity model on the Cityscapes dataset[3]. We trained all the models on a single P-100 GPU offered by  
16 Kaggle for over 850 training hours.

### 17 Results

18 Overall, our results back the original paper’s claims. As shown in Table 3, our model outperforms the original study on  
19 a few metrics but slightly falls behind on others on the benchmarked Fishyscapes dataset; discussed in section 5.1.

### 20 What was easy

21 The paper was well-written and easy to understand. The provided open-source code is well-structured and modular.  
22 Having pre-trained weights available for standard segmentation and reconstruction models reduced computational load.

### 23 What was difficult

24 Even with modular code available, re-implementing the code in PyTorch Lightning proved more challenging than  
25 expected. Our experiments were limited by the model’s computational constraints, with an average training duration of  
26 25 hours per model and kaggle only providing 9 hours of continuous training time. The data generation method is not  
27 specified in the original repository; We have created a single script in our repository for the same.

### 28 Communication with original authors

29 Authors were contacted via email to help clarify queries about the code, dataset generation, and discrepancies in the  
30 results. Our final report received a positive review from the authors.

31 **1 Introduction**

32 Deep neural networks are becoming increasingly prominent in various applications, including computer vision, speech  
 33 recognition, and language modeling. One such application is Semantic Segmentation, the per-pixel categorization  
 34 of objects in an image. Recent advancements in Deep learning have brought significant improvement in this task,  
 35 providing very accurate results, yet these networks fail to detect objects which are not part of the training dataset. Since  
 36 anomalies are a component of many critical real-world applications, such as autonomous driving, implementing these  
 37 networks for real-time situations necessitates overcoming issues like recognizing anomalous objects (any impediment  
 38 on the road) and misclassification.

39 The existing methodologies either rely on predicted segmentation maps and their confidence scores or compare the  
 40 reconstructed image to the predicted segmentation map to detect anomalies. However, these methods may fail as the  
 41 segmentation network might make noisy anomaly predictions. Building over the existing resynthesis methods, the  
 42 authors use uncertainty measures [6] [7] [11] to assist the dissimilarity network in differentiating the input and generated  
 43 images that successfully generalizes to all anomalies.

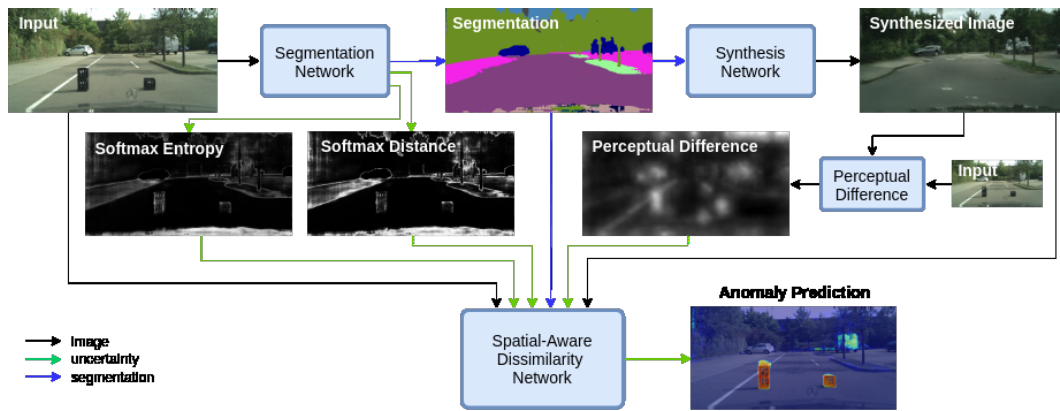


Figure 1: Anomaly Segmentation Framework described in Section 3.2

44 **2 Scope of reproducibility**

45 We sought to answer the following questions by reproducing the results of the paper, as well as conducting additional  
 46 experiments to corroborate the paper’s primary argument:

- 47 1. What is the importance of the uncertainty map’s softmax entropy  $H(1)$ , softmax distance  $D(2)$ , and perceptual  
 48 difference  $V(3)$  (section 4.2.1)?
- 49 2. Can the network generalize to other segmentation and resynthesis networks (section 4.1.3)?
- 50 3. How does the performance of segmentation and resynthesis networks affect the overall performance of the  
 51 model (section 4.2.2)?
- 52 4. Does the additional data generator added in this method improve the results (section 4.1.1)?
- 53 5. Is it better to choose weights for uncertainty maps for prediction during training or at the end through grid  
 54 search (section 4.1.2)?
- 55 6. We further evaluate the sensitivity of our model to changes in the feature extractor module (encoder), activation  
 56 function, and learning rate (section 4.2.3).

57 **3 Methodology**

58 **3.1 Model descriptions**

59 The model mainly has three modules namely segmentation, synthesis, dissimilarity and an ensemble:

60 **Segmentation Module :** We employ the pre-trained weights of the model as trained in [14] on Cityscapes dataset. In  
 61 addition to generating a segmented image, we generate two dispersion maps, softmax entropy  $H$  and softmax  
 62 distance  $D$ , which prove beneficial in understanding anomalies within the generated segmentation map ( $p(c)$  is the softmax  
 63 probability for class  $c$ ). For each pixel  $x$ ,  $H$  and  $D$  are calculated as follows:

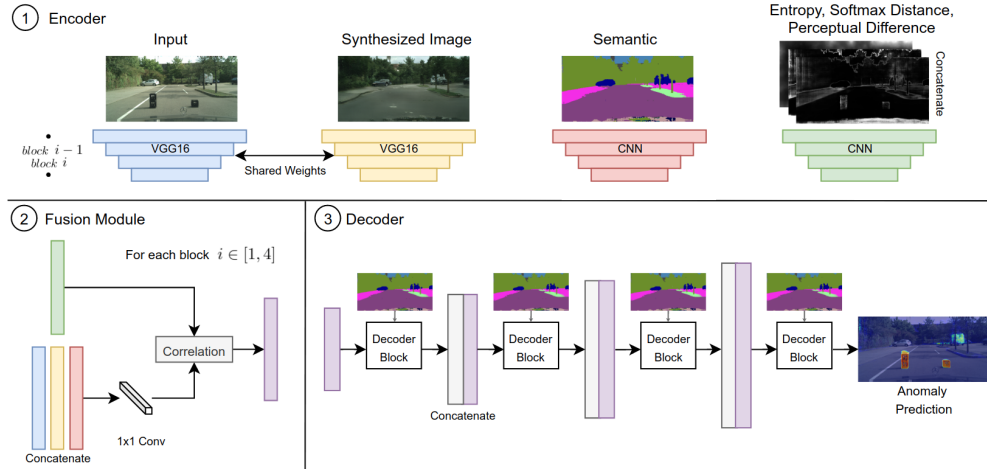
$$H_x = - \sum_{c \in \text{classes}} p(c) \log_2 p(c) \quad (1)$$

$$D_x = 1 - \max_{c \in \text{classes}} p(c) + \max_{c^1 \in \text{classes} \setminus (\arg \max_c p(c))} p(c^1) \quad (2)$$

64 **Synthesis Module** : To build a realistic image out of the segmentation map, we employ pre-trained weights from  
 65 the model trained on Cityscapes dataset as a conditional generative adversarial network (c-GAN) [9] [13]. However,  
 66 because the semantic map lacks information such as color appearance, per-pixel value comparison between the original  
 67 input and the synthesized image is not possible. As a result, we use perceptual difference, which employs a pre-trained  
 68 VGG16 model as a feature extractor to compare overall spatial structure rather than features such as color and texture,  
 69 allowing us to better classify anomalies. For every pixel  $x$  of the input image and corresponding pixel  $r$  from the  
 70 synthesized image  $V$  is defined as follows :

$$V(x, r) = \sum_1^N \frac{1}{M_i} \|F^i(x) - F^i(r)\|_1 \quad (3)$$

71  $F(i)$  denotes the  $i$ -th layer with  $M_i$  elements of the VGG network and  $N$  layers, normalized between  $[0, 1]$ .



**Figure 2: Spatial-Aware Dissimilarity Module** described in **Section 3.1**

72 **Spatial-Aware Dissimilarity Module** : This module takes as input the original image, generated image, semantic map,  
 73 and uncertainty maps (softmax entropy, softmax distance, perceptual difference) calculated in the previous steps to  
 74 predict the anomaly segmentation map. It is mainly divided into three modules namely encoder, fusion and decoder:

- 75 1. **Encoder** : Encoder uses a pretrained VGG 16[12] to extract features of resynthesis and input image. Whereas  
 76 a simple CNN is used to extract features from the uncertainty maps and semantic map.
- 77 2. **Fusion Module** : Concatenates and passes features extracted from resynthesis, input, segmented maps through  
 78 a  $1 \times 1$  convolution which is then passed into correlation block along with encoded uncertainty map where  
 79 pointwise correlation is performed outputting four feature map resolutions corresponding to each of the four  
 80 layers of the decoder.
- 81 3. **Decoder** : There are four decoder blocks used in the dissimilarity network. The first decoder block takes the  
 82 lowest resolution feature map. The concatenation of the feature map from the fusion module and the output of  
 83 the preceding decoder block is used as the input for all subsequent decoder blocks.

84 **Ensemble** : Dissimilarity module predictions are usually overconfident, which is reduced by combining final predictions  
 85 with uncertainty maps (1),(2) ,(3) using appropriate weights obtained by grid search.

### 86 3.2 Training Procedure

87 We collect semantic maps and two uncertainty maps( $H$  and  $D$ ) by passing input images through the segmentation  
 88 network. The synthesis network subsequently processes the predicted semantic map, which results in a photo-realistic

89 image. The perceptual difference is then calculated by comparing features between the input and generated images. The  
 90 spatial-aware dissimilarity module receives uncertainty maps (1),(2) ,(3), input image, semantic map and resynthesized  
 91 image to generate the anomaly prediction, which is then combined with uncertainty maps using optimal weights found  
 92 via grid search.

### 93 3.3 Datasets

94 We employed the data generator technique provided in [8], which was expanded by assuming void labels in the ground  
 95 truth semantic map as anomalies.

96 Evaluation is done on two sets of Fishyscapes Benchmark[2]:

- 97 1. **FS Lost & Found[10]**: A collection of 275 images with small objects (e.g. toys, boxes) on the roadway.
- 98 2. **FS Static[5]**: A collection of 1000 images from Cityscapes blended with anomalous Pascal VOC objects.

### 99 3.4 Hyperparameters

Hyperparameter	Value
Number of epochs	25 or 50
Learning rate	1E-4
Learning rate Policy	ReduceLROnPlateau
Weight decay	0.0000
Power	0.9
Patience	10
Batch Size	4
$\beta_1$	0.5
$\beta_2$	0.999

**Table 1:** We observed that the optimal model is obtained between 5 and 15 epochs. As a result, for additional experiments, we just train for 25 epochs.

### 100 3.5 Experimental setup

101 The training code was run on KAGGLE with Tesla P100-PCIE-16GB GPU (NVIDIA-SMI 450.119.04, Driver Version:  
 102 450.119.04, CUDA Version: 11.0).

### 103 3.6 Computational requirements

104 Using the pre-trained weights, the generation of semantic maps (softmax entropy, softmax distance) and synthesised  
 images (along with perceptual difference) took 3 hours and 2 hours, respectively.

Training time per epoch (Dissimilarity Module)	30 minutes
GPU Requirement	9 - 10 Gb
CPU memory Requirement	2-3 Gb
Inference time (Dissimilarity Module)	63 ms

**Table 2:** Computational Requirement

## 106 4 Results

### 107 4.1 Results reproducing original paper

#### 108 4.1.1 Main Results

Sample	FS L&F		FS Static	
	$\uparrow$ AP	$\downarrow$ FPR95	$\uparrow$ AP	$\downarrow$ FPR95
Original Study	43.22	<b>15.79</b>	<b>72.59</b>	18.75
Reproduced Results	<b>44.47</b>	18.7	71	<b>17.17</b>

**Table 3:** Benchmarked on Fishyscapes dataset[2] through model submission on official website

109 Table 3 shows that our model outperforms the original study in FPR 95 (False Positive Rate at 95% True Positive Rate)  
 110 of FS STATIC and AP (Average Precision) of FS L&F on the benchmarked fishyscapes dataset, but slightly falls behind  
 111 on AP of FS STATIC and FPR 95 of FS L&F; discussed in section 5.1

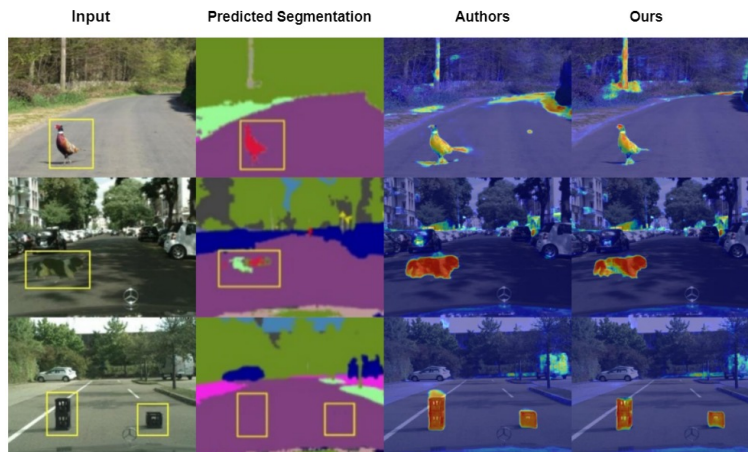
Method	FS L&F (OURS)		FS L&F		FS Static (OURS)		FS Static	
	↑AP	↓FPR95	↑AP	↓FPR95	↑AP	↓FPR95	↑AP	↓FPR95
Full Framework	51 ± 7	41 ± 1	55 ± 5	40 ± 5	57 ± 4	28 ± 2	62 ± 5	26 ± 1
w/o ensemble	59 ± 5	61 ± 2	58 ± 9	66 ± 8	55 ± 6	32 ± 4	57 ± 6	41 ± 12
w/o unc. maps	28 ± 6	60 ± 8	39 ± 9	64 ± 10	28 ± 6	64 ± 4	38 ± 8	51 ± 4
w/o data generator								
& w/o unc. maps	14 ± 5	55 ± 9	15 ± 5	63 ± 12	11 ± 3	58 ± 9	14 ± 4	57 ± 11

**Table 4:** The following results are calculated on publicly available Fishyscapes validation datasets

112 The above results are reported as an average and standard deviation across five random weight initializations, as done in  
 113 the original paper. Below are the training approaches whose results have been provided in Table 4:

- 114 1. **w/o unc maps** : Training performed without the use of uncertainty maps (1),(2) ,(3).  
 115 2. **w/o data generator & w/o unc. maps** : Training performed without using uncertainty maps (1), (2), (3) and  
 116 without additional data generator as explained in section 3.3.  
 117 3. **Full Framework**: Full training as instructed in section 3.2 along with data generator.  
 118 4. **w/o ensemble**: Trained as instructed in section 3.2 except the ensemble.

119 Table 4 shows that most of the results are within the authors’ given range. The minor differences can be explained  
 120 since the model adjusts for larger variances shown in the results, and the five runs performed may not be adequate to  
 121 generalize the model results. We can also see that uncertainty maps, data generator help to improve overall performance.



**Figure 3:** Our predictions are comparable with the authors’ in detecting the main anomaly \*

122

#### 123 4.1.2 End to End ensemble Vs Ensemble Grid Search

Method	FS L&F		FS Static	
	↑AP	↓FPR95	↑AP	↓FPR95
ensemble grid search	51.54	<b>42.36</b>	57.36	<b>28.22</b>
end to end ensemble (OURS)	<b>61.01</b>	63.38	47.02	46.31
end to end ensemble	59.6	58.6	<b>61.1</b>	37.3

**Table 5:** Comparison between end to end ensemble training and ensemble through grid search.

124 Below are the training approaches whose results have been provided in the above table:

- 125 1. **ensemble through grid search**: Training is carried out as stated in section 3.2, with weights for uncertainty  
 126 determined at the end of training using grid search.  
 127 2. **end to end ensemble**: Here, weights for uncertainty maps are learnable parameters optimized during training.

\*More images comparing author’s outputs and ours are provided in the Supplementary material

128 We attribute higher deviation from results again to the fact that the experiments have a higher standard deviation, and  
 129 even five runs may not be sufficient to generalize. However, in FS Static, our results show that an ensemble with  
 130 empirical weights is more likely to identify intentionally blended objects, whereas end-to-end training is less likely  
 131 since the network improves the weights before the final prediction. The increase in AP of FS L& F in end to end  
 132 ensemble is expected as the network optimizes the weights efficiently. Our results supports the authors’ contention  
 133 that utilising an ensemble through grid search yields lower FPR 95 values. In safety-critical contexts, this makes the  
 134 ensemble with empirical weights for final prediction more practical (e.g., autonomous driving).

### 135 4.1.3 Generalizability of the model

Method	FS L&F		FS Static	
	↑AP	↓FPR95	↑AP	↓FPR95
Original	<b>51.54</b>	<b>42.36</b>	<b>57.36</b>	<b>28.22</b>
Light ( <b>OURS</b> )	32.88	50.96	28.68	31.84
Light	36	46.4	33.4	36.1
Im. Resyn.++	5.7	47.7	8	62.7

**Table 6:** Testing dissimilarity module on lighter segmentation and resynthesis frameworks.

136 Below are the combination of segmentation and resynthesis networks used, whose results have been shown in Table 6:

- 137 1. **Original** : SDCNet Segmentation + C-GAN Resynthesis Networks.
- 138 2. **LIGHT** : ICNet Segmentation + Spade Resynthesis Networks.
- 139 3. **Im. Resyn.++** : Image Resynthesis++ is used for comparison as our method builds upon it.

140 According to Table 6, the results for the light framework are significantly lower than the original model but significantly  
 141 better than the original image resynthesis method(Im. Resyn.++), proving the authors’ claim that this method generalises  
 142 well enough to be used as a wrapper to already existing segmentation and resynthesis networks.

## 143 4.2 Results beyond original paper

### 144 4.2.1 Importance of Uncertainty Maps

Method	FS L&F		FS Static	
	↑AP	↓FPR95	↑AP	↓FPR95
Full Framework	51.54	42.36	<b>57.36</b>	28.22
W/O Softmax Entropy	37.03	42.8	40.17	35.9
W/O Softmax Distance	45.93	41.47	53.47	<b>27.47</b>
W/O Perceptual Difference	<b>51.7</b>	<b>39.57</b>	51.8	27.83

**Table 7:** Assessing the significance of each uncertainty map

145 From Table 7, we can observe that omitting Softmax Entropy lowers all outcomes. Even Softmax Distance has a  
 146 noticeable impact on the model’s performance. However, removing Perceptual Difference has almost no effect (except  
 147 on AP of FS STATIC), which we believe is because the resynthesized image is already passed to the dissimilarity  
 148 module, which does the task of differentiating input and resynthesized image.

### 149 4.2.2 Importance of performance of Segmentation and Resynthesis networks

Method	FS L&F		FS Static	
	↑AP	↓FPR95	↑AP	↓FPR95
Moderate1	44.767	<b>38.88</b>	31.6	39.167
Moderate2	49.87	43.83	50.93	29.2

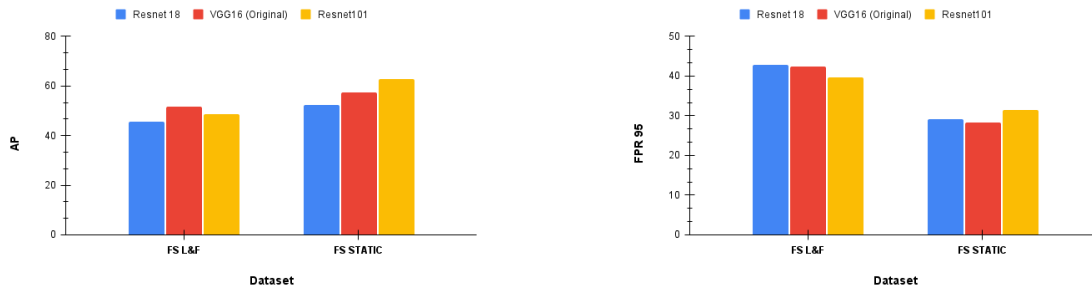
**Table 8:** Effect of segmentation and resynthesis networks.

150 Below are the combination of segmentation and resynthesis networks used, whose results have been shown in Table 8:

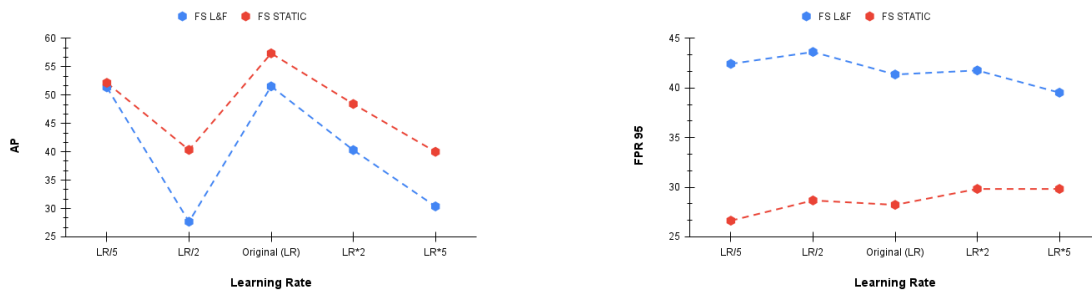
- 151 1. **Moderate1** : SDCNet Segmentation + Spade Resynthesis Networks.
- 152 2. **Moderate2** : ICNet Segmentation + C-GAN Resynthesis Networks.

153 From Table 8, the quality of the segmentation and resynthesis networks has an impact on the outcomes, and it is clear that  
 154 the original performs better than light because it employs superior segmentation and synthesis methods. Results of Light  
 155 version and original can be referenced from Table 6. Original findings outperform Moderate1 (lighter segmentation,  
 156 same resynthesis network) and Moderate2 (same segmentation, lighter resynthesis network), demonstrating that both  
 157 synthesis and segmentation are critical components of the model and are directly connected to its overall performance.

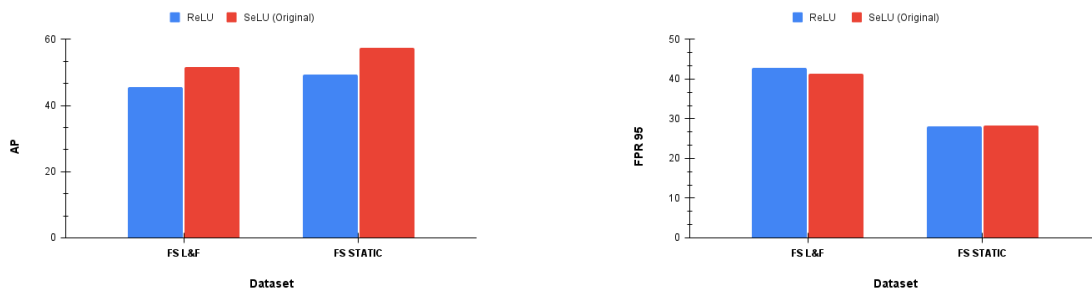
### 158 4.2.3 Network Tuning



(a) AP Vs Encoder (b) FPR 95 Vs Encoder  
**Figure 4:** VGG16 dominates ResNet18 & ResNet101 dominates VGG16 in FPR 95 of FS L&F and AP of FS STATIC



(a) AP Vs Learning Rate (b) FPR 95 Vs Learning Rate  
**Figure 5:** Overall the Original Learning Rate(1e-4) seems to perform better.



(a) AP Vs Activation Function (b) FPR 95 Vs Activation Function  
**Figure 6:** We can observe the supremacy of SeLU over ReLU when used in dissimilarity module in all the outcomes.

## 159 5 Discussion

160 Our results are on par with those of the authors' on the private test data but do not exactly match the results on the  
 161 validation datasets despite utilizing the exact parameters described in the paper. Nonetheless, they are practically on par  
 162 with them, proving the authors' claim. Section 4.1.3 demonstrates that the model can generalize to even lower-level  
 163 segmentation and synthesis networks, allowing it to be used as a wrapper to any pre-trained segmentation and resynthesis  
 164 models. Section 4.2.2 demonstrates the importance of segmentation and resynthesis network performance. In Section  
 165 4.1.2, we can see that ensemble using grid search outperforms end-to-end ensemble.

166 We provide two arguments to explain why our results differ from the authors. Firstly with five random initializations,  
167 we can see very high standard deviation; so, other standard deviation remains to be accounted for, and these runs are  
168 insufficient to generalize conclusions. Secondly, the weights for FS L&F and FS STATIC for uncertainty maps after  
169 ensemble are different, as mentioned in section 5.1.

170 We quantitatively confirmed the relevance of uncertainty maps in section 4.2.1, where we observed that removing  
171 the softmax entropy map or the softmax distance map has a considerable effect on outcomes, although removing  
172 the perceptual difference had little to no effect. To improve anomaly segmentation, all uncertainty maps are heavily  
173 employed in conjunction with their complimentary information.

174 We also attempted to enhance the findings in section 4.2.3; (4a, 4b) we replaced the VGG16 encoder with ResNet18 and  
175 ResNet101; (6a, 6b) We tried commonly used ReLU instead of SeLU activation function in the dissimilarity module;  
176 (5a, 5b) we modified learning rates. Despite improvements in a few outcomes, the parameters used by authors proved to  
177 be the best, with a strong balance among all metrics.

## 178 5.1 Further discussion on discrepancies of the results

179 Table 4 shows that ensemble weights are an important aspect of the model, as they reduce the model’s overconfidence.  
180 The challenge, however, is determining the optimal ensemble weights for the model. The sensitivity to these ensemble  
181 weights is one of the key causes for our results’ divergence.

182 For most of our runs, the ensemble weights for FS L&F differed from those for FS STATIC. This makes keeping the  
183 final ensemble weights a little more challenging. In most circumstances, [0.75,0.25,0,0] (the four weights being for  
184 original output, entropy, perceptual difference, and softmax distance, respectively) yields a better outcome. However,  
185 employing these weights degrades the results of FS L&F; which can be seen from its FPR 95 in Table 3.

186 Table 4 shows that the model accounts for more significant variations, and the five runs done may be insufficient to  
187 generalise the results. This large deviation in model results might be the reason for lower AP of FS Static in Table 3.

188 Another thing to note is that all of the tuning and validation is done on the Fishyscapes validation datasets, even though  
189 the Fishyscapes paper expressly warns that this may result in overfitting to the validation set and will most likely not  
190 result in better performance on the unexpected test data. The fact that our model doesn’t perform on par with the  
191 original model on the validation datasets, but performs well on private test data proves this claim. But owing to a  
192 shortage of suitable anomaly detection datasets we are forced to use these validation datasets for tuning. Furthermore,  
193 choosing tuning/validation datasets for anomaly segmentation is always difficult since we are constantly in danger of  
194 overfitting the parameters to specific abnormalities.

### 195 What was easy

196 The paper was well-written and easy to understand. The provided open-source code is well-structured and modular.  
197 Having pre-trained weights available for standard segmentation and reconstruction models reduced computational load.

### 198 What was difficult

199 Even with modular code available, re-implementing the code in PyTorch Lightning proved more challenging than  
200 expected. Our experiments were limited by the model’s computational constraints, with an average training duration of  
201 25 hours per model and kaggle only providing 9 hours of continuous training time. The data generation method is not  
202 specified in the original repository; We have created a single script in our repository for the same.

### 203 Communication with original authors

204 Authors were contacted via email to help clarify queries about the code, dataset generation, and discrepancies in the  
205 results. Our final report received a positive review from the authors.

### 206 Acknowledgement

207 We thank the original authors for their quick and thorough response, facilitating reproducibility. We also thank Boyang  
208 SUN for assisting us with benchmarking results on the Fishyscapes website, making reproducibility complete.

### 209 Future Scope

210 From Table 4, we see that the dissimilarity module alone takes 63ms. Further, considering inference time of segmentation  
211 and resynthesis networks limits the model from applying it for real-time tasks. Hence improving the model for real-time  
212 performance would be a considerable task for future. One approach to get around this could be to use distributed  
213 training. Because PyTorch Lightning is hardware agnostic and easy to scale, it streamlines this task.



214 **References**

- 215 [1] Lukas Biewald. *Experiment Tracking with Weights and Biases*. Software available from wandb.com. 2020. URL:  
216 <https://www.wandb.com/>.
- 217 [2] Hermann Blum et al. “The Fishyscapes Benchmark: Measuring Blind Spots in Semantic Segmentation”. In:  
218 *CoRR* abs/1904.03215 (2019). arXiv: [1904.03215](https://arxiv.org/abs/1904.03215). URL: <http://arxiv.org/abs/1904.03215>.
- 219 [3] Marius Cordts et al. “The Cityscapes Dataset for Semantic Urban Scene Understanding”. In: *Proc. of the IEEE*  
220 *Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- 221 [4] Giancarlo Di Biase et al. “Pixel-Wise Anomaly Detection in Complex Driving Scenes”. In: *Proceedings of the*  
222 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2021, pp. 16918–16927.
- 223 [5] M. Everingham et al. “The Pascal Visual Object Classes Challenge: A Retrospective”. In: *International Journal*  
224 *of Computer Vision* 111.1 (Jan. 2015), pp. 98–136.
- 225 [6] Yarin Gal and Zoubin Ghahramani. “Dropout as a Bayesian Approximation: Representing Model Uncertainty  
226 in Deep Learning”. In: *Proceedings of The 33rd International Conference on Machine Learning*. Ed. by Maria  
227 Florina Balcan and Kilian Q. Weinberger. Vol. 48. Proceedings of Machine Learning Research. New York, New  
228 York, USA: PMLR, 20–22 Jun 2016, pp. 1050–1059. URL: <https://proceedings.mlr.press/v48/gal16.html>.  
229 [html](https://proceedings.mlr.press/v48/gal16.html).
- 230 [7] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. “Perceptual losses for real-time style transfer and super-  
231 resolution”. In: *European conference on computer vision*. Springer. 2016, pp. 694–711.
- 232 [8] Krzysztof Lis et al. *Detecting the Unexpected via Image Resynthesis*. 2019. arXiv: [1904.07595](https://arxiv.org/abs/1904.07595) [cs.CV].
- 233 [9] Xihui Liu et al. “Learning to Predict Layout-to-image Conditional Convolutions for Semantic Image Synthesis”.  
234 In: *CoRR* abs/1910.06809 (2019). arXiv: [1910.06809](https://arxiv.org/abs/1910.06809). URL: <http://arxiv.org/abs/1910.06809>.
- 235 [10] Peter Pinggera et al. “Lost and Found: detecting small road hazards for self-driving vehicles”. In: *2016 IEEE/RSJ*  
236 *International Conference on Intelligent Robots and Systems (IROS)* (2016), pp. 1099–1106.
- 237 [11] Matthias Rottmann et al. “Prediction Error Meta Classification in Semantic Segmentation: Detection via Ag-  
238 gregated Dispersion Measures of Softmax Probabilities”. In: *2020 International Joint Conference on Neural*  
239 *Networks (IJCNN)*. 2020, pp. 1–9. DOI: [10.1109/IJCNN48605.2020.9206659](https://doi.org/10.1109/IJCNN48605.2020.9206659).
- 240 [12] Karen Simonyan and Andrew Zisserman. “Very Deep Convolutional Networks for Large-Scale Image Recogni-  
241 tion”. In: *arXiv 1409.1556* (Sept. 2014).
- 242 [13] Ting-Chun Wang et al. “High-Resolution Image Synthesis and Semantic Manipulation with Conditional GANs”.  
243 In: *CoRR* abs/1711.11585 (2017). arXiv: [1711.11585](https://arxiv.org/abs/1711.11585). URL: <http://arxiv.org/abs/1711.11585>.
- 244 [14] Yi Zhu et al. “Improving Semantic Segmentation via Video Propagation and Label Relaxation”. In: *Proceedings*  
245 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2019.