

---

# Fully-Convolutional Networks: Ship Detection via Semantic Segmentation

---

**Ross Greer**

Electrical And Computer Engineering  
University of California, San Diego  
regreer@eng.ucsd.edu

**Vivian Meng**

Electrical And Computer Engineering  
University of California, San Diego  
vimeng@eng.ucsd.edu

**Winson Quan**

Electrical And Computer Engineering  
University of California, San Diego  
w1quan@eng.ucsd.edu

**Yangting Sun**

Electrical And Computer Engineering  
University of California, San Diego  
yas109@eng.ucsd.edu

**Yifan Xu**

Electrical And Computer Engineering  
University of California, San Diego  
yix247@eng.ucsd.edu

## Abstract

Object detection is a familiar problem in current machine learning and computer vision research. One challenge with object detection, and learning algorithms in general, is the trade-off between generalized approaches and optimized performance on specific tasks, cited as Wolpert and Macready's "no free lunch" theorem. In this project, we investigate performance of common convolutional neural network architectures applied to the problem of detecting ships in water under a variety of image conditions. In particular, we implement Shelhamer, Long, and Darrell's "Fully Convolutional Networks for Semantic Segmentation", using both AlexNet and VGGNet as internal network architectures. While using VGGnet produced reasonable results, using AlexNet failed to learn during the training process.

## 1 Introduction

This project is derived from the Kaggle "Airbus Ship Detection Challenge," which motivates the problem as follows:

*Shipping traffic is growing fast. More ships increase the chances of infractions at sea like environmentally devastating ship accidents, piracy, illegal fishing, drug trafficking, and illegal cargo movement. This has compelled many organizations, from environmental protection agencies to insurance companies and national government authorities, to have a closer watch over the open seas.*

In our project, we worked to create a model that detects all ships in satellite images as quickly as possible.

### 1.1 Motivation

This project can help to support the maritime industry to increase knowledge, anticipate threats, trigger alerts, and improve efficiency at sea.

## 24 2 Description of Method

### 25 2.1 Algorithm

26 **Fully-Convolutional Networks** We implement a variant of the algorithm by Long, Shelhamer,  
27 and Darrell for semantic segmentation using fully-convolutional networks (FCN). Input images are  
28 passed through a series of convolutional and pooling layers, following the architectures of AlexNet  
29 and VGGNet. At different checkpoints along these architectures, the resulting tensor has a fractional  
30 size of the original image. Long et al define three such checkpoints; one at 1/8 of the original size,  
31 one at 1/16, and one at 1/32. For each of these tensors, we can immediately apply upsampling to  
32 create a representation of the image features at decreasing resolution; that is, upsampling the smallest  
33 image to restore it to a full 256x256 image will give a blurry feature representation, while the 1/8 size  
34 image can be restored to a finer (albeit still granular) resolution. This algorithm makes use of the  
35 tradeoff between resolution and parameter specificity. The image with the lowest resolution has the  
36 greatest number of extracted features, while the image with the greatest resolution has less extracted  
37 features. These three extraction sets are fused as follows: the tensor with 1/32 size is upsampled  
38 to match the image at 1/16 size, then their weighted score representations are summed. Next, this  
39 combined tensor is upsampled to match 1/8 the original image size, scored, and summed with the  
40 score 1/8 size image. Finally, the tensor is upsampled to its original size of 256x256, and a final  
41 convolutional layer is applied to classify each pixel into its expected binary class.

### 42 2.2 Architecture

#### 43 2.2.1 Convolutional Neural Network

44 Convolutional Neural Networks (CNN) are a special kind of multi-layer neural network, designed to  
45 recognize visual patterns directly from pixel images with minimal preprocessing. We chose to use  
46 established classification architectures AlexNet and VGGNet in order to compare their performance  
47 within the larger algorithm.

48 **AlexNet** In 2012, AlexNet significantly outperformed all the prior competitors and won the Im-  
49 ageNet challenge by reducing the top-5 error from 26 percent to 15.3 percent. It makes use of  
50 convolutions with filter sizes 11x11, 5x5, 3x3, max pooling, dropout, and ReLU activation layers,  
51 data augmentation, and SGD with momentum. ReLU activations are present after every convolutional  
52 and fully-connected layer.

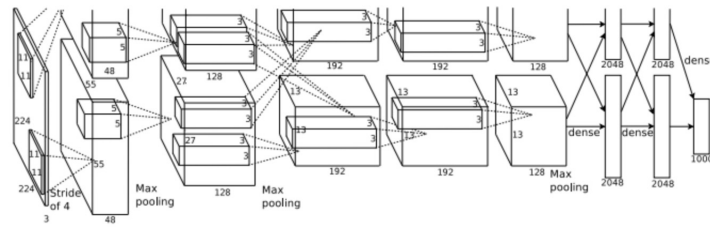


Figure 1: AlexNet Architecture.

53 **VGGNet** VGGNet was developed by Simonyan and Zisserman, achieving second place per-  
54 formance at the ILSVRC 2014 competition. VGGNet consists of 16 convolutional layers and is appealing  
55 because of its uniform architecture. The architecture is similar to AlexNet, however, convolutional  
56 layers use only 3x3 convolutions with a large bank of filters.

#### 57 2.2.2 Fully Convolutional Networks

58 Fully Convolutional Networks (FCNs) are built only from locally connected layers, such as con-  
59 volution, pooling, and upsampling. The network outputs a segmentation map. The architecture  
60 proposed by Long et al. reduces the number of parameters and computation time. Also, given that all

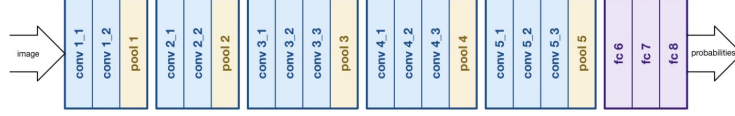


Figure 2: VGGNet Architecture.

connections are local, the network can work regardless of the original image size, without requiring any fixed number of units at any stage.

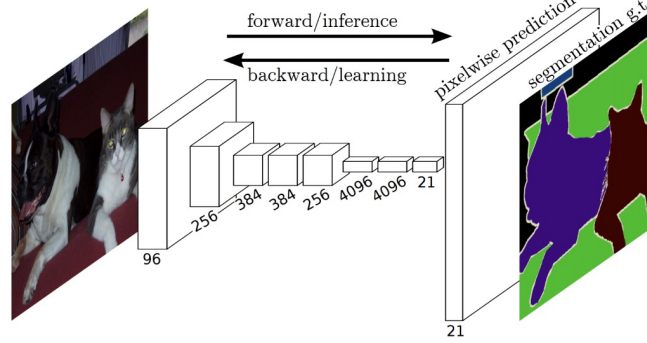


Figure 3: FCN Architecture.

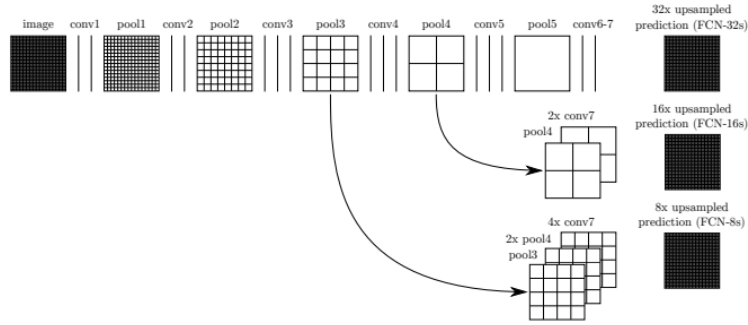


Figure 4: The above figure demonstrates the method Long et al implemented to fuse output from layers at different depths of the network. The two columns on the right represent the layer's extracted features, followed by the upsampled version which is summed to provide the final output.

### 2.3 Equations

Our metric for performance on the test set will be Intersection over Union (IoU) for known segments of ships on the labeled test data. Our test results on a per-image basis will be in a form of labeled pixels, with an active label corresponding to a prediction of 'ship' class, and an inactive label corresponding to a prediction of 'background'.

**Intersection over Union (IoU)** The equation for Intersection over Union is given in equation 1. The intersection is comprised of pixels found in both the ground truth and prediction masks. The union is comprised of all pixels found in either the ground truth or prediction masks.

$$IoU = \frac{\text{Area of Overlap}}{\text{Area of Union}} = \frac{\text{ground truth} \cap \text{prediction}}{\text{ground truth} \cup \text{prediction}} \quad (1)$$

71 The IoU score is calculated for each image in the testing set, then averaged for a mean IoU score for  
72 the segmentation predictions.

73 **Weighted Cross-Entropy Loss** We use a weighted cross-entropy loss function in our model to  
74 rectify the effects of having many more samples of oceanic background than images with ships.

$$L_n(W) = -w_{c_n} \log y_{c_n}(x_n, W) \quad (2)$$

75 In equation 2,  $w_{c_n}$  represents the weight of class  $c_n$ , and  $y_{c_n}(x_n, W)$  represents the output of the  
76 network on input  $x_n$  using weights  $W$ . In our task, we assign a greater weight to class 1 (ship  
77 present), since this class is underrepresented in the dataset.

### 78 3 Implementation Details

79 Many implementations of fully convolution networks for semantic segmentation use publicly available  
80 pre-trained weights trained on well known architectures and ImageNet. Investigation with AlexNet  
81 using pre-trained weights showed that such parameters have difficulty recognizing the ships in our  
82 dataset, possibly because of the small size of the ships in relation to the background. In addition,  
83 networks trained on ImageNet have 1000 classes whereas we only need 2.

84 To decrease the time required to training, images were cropped with the ship centered in the image.  
85 This also increases the ratio of ship pixels to background pixels. For images with multiple ships,  
86 multiple cropped images were produced with each centered on a different ship. For images without  
87 ships, a single image was cropped from the center.

88 The AlexNet classifier was trained with cross-validation and had the training and validation accuracy  
89 evaluated every epoch (as in section 4). VGGNet trains much slower so these steps were omitted.

90 After the classifiers were trained, the fully convolutional layers, including deconvolution, were  
91 attached and the networks were trained with respect to a mask indicating the pixels of the ship.

92 In defining layers and the forward method, it is worth noting that layers which require particular  
93 learning rates for the fully-connected portion of the architecture must be separated from their  
94 preceding sequential batch in order to later define their learning rate and momentum in the model's  
95 parameters, since Pytorch handles this for an entire set of sequential-layers. In this network, we add a  
96 function "remove\_classifier", which specifies if training should be performed for only the VGGNet  
97 or AlexNet architecture, or combined with the fusion, scoring, deconvolution, and convolution layers  
98 used for semantic segmentation. This abstraction allows us to take advantage of pre-trained weights  
99 for the classification networks, and to save our model and make use of transfer learning to attach  
100 learned weights to our greater context.

101 Following the classifier architecture, we define class FCN8s, which contains the layers of fusion in  
102 the architecture designed by Long et al. This is the portion of the code which contains upsampling  
103 via deconvolution, increasing the tensors to gradually larger and larger size, fusing them with tensors  
104 from earlier layers, until a resulting tensor of original image size is generated and its pixels are  
105 classified.

106 With the layers and network defined, training parameters, weighted cross-entropy loss function, and  
107 optimization method are established next. Both AlexNet and VGGNet were trained with Xavier  
108 initializations and using standard SGD with a momentum of 0.9. The fully convolutional network of  
109 VGG used the Adam optimizer. Classifiers used cross-entropy loss and the segmentation networks  
110 used weighted cross-entropy loss. The AlexNet fully convolutional network used scheduler to reduce  
111 the learning rate at a plateau, but this was unnecessary for the VGGNet.

112 The dataset was divided up to allow for different training, validation, and test sets. However, due to  
113 the lack of time for tuning parameters, the validation set was used as the test set and validation was  
114 not performed.

115 The implementation of FCN with AlexNet closely follows a implementation in CAFFE available at  
116 ([github.com/shelhamer/fcn.berkeleyvision.org](https://github.com/shelhamer/fcn.berkeleyvision.org)).

## 117 4 Experimental Setting

### 118 4.1 Dataset

119 The dataset for this project is distributed by Kaggle, and sourced from the company Airbus. The  
120 dataset is available for download at

121 <https://www.kaggle.com/c/airbus-ship-detection/data>

122 The dataset is comprised of 192,556 satellite images, each 768x768 pixels, along with the pixel  
123 locations that comprise of the ships in the image. There are an additional 15,606 unlabeled images,  
124 provided by Kaggle for online testing.

125 In the images, there may be no ships, a single ship, or multiple ships. Backgrounds vary from open  
126 water to a crowded marina. Additionally, the images may have clouds, haze, or other imperfections.

127 The first 40,000 images were used as test data and the images in the range of 40,000 to 160,000 were  
128 used as training data.

### 129 4.2 Training Parameters

130 For AlexNet, the classifier used a minibatch size of 8, a learning rate of 0.005. Training was stopped  
131 early at 10 epochs as that produced the lowest training and validation error. The VGGNet classifier  
132 was trained for 13 epochs.

### 133 4.3 Hardware Used

134 We trained our network on UCSD's Data Science and Machine Learning Platform (DSMLP). The  
135 network runs on the full allotment of 8 CPU cores, 32 GB of memory, and one GPU.

## 136 5 Results

### 137 5.1 Early Experiments

138 The robustness of the AlexNet classifier was shown with four-way cross validation, as shown in  
139 figures 5 and 6.

140 In figure 7, the error does not differ much between a minibatch size of 8 and 32 with the same  
141 number of epochs. However, the training is less stable with smaller minibatch sizes.

142 The training error of AlexNet and VGGNet are similar given the same number of epoch, but the  
143 validation error of AlexNet is much lower than that of VGGNet, as seen in figure 8. This is not  
144 unexpected because VGGNet is much deeper and hence more prone to over-fitting.

### 145 5.2 Success Cases

146 Our implementation of semantic segmentation using VGGNet in a fully-convolutional network shows  
147 success on cases involving ships on an open ocean background, as in Figure 5.1, and ships which  
148 blend in to the oceanic background, as in Figure 9. Our results show that we can successfully detect  
149 multiple ships in a scene. While the resolution is not as precise as the ground truth, the approximate  
150 centroids of each detected ship closely align with the ground truth centroids, and variance is contained  
151 by the ships width in any direction.

152 A separate success of our model is its ability to not throw false positives over hazy blocking conditions  
153 or general oceanic background, as in Figure 11. Though there is a small pattern of noise present, in  
154 general, the segmentation does not identify ambient background as ship presence.

### 155 5.3 Performance

156 The average Intersection over Union achieved is 0.22, averaged over images with a ship in both the  
157 ground truth and predicted segmentation mask. With respect to image-wise detection, there were

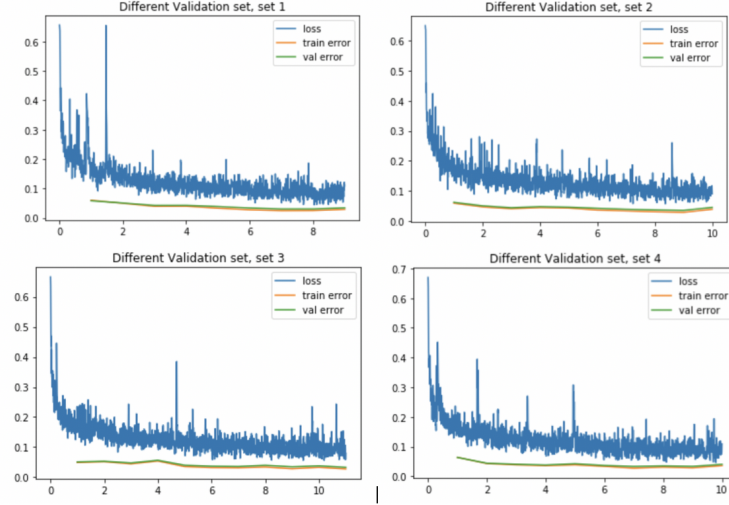


Figure 5: Validation set Error Comparisons

	training error	validation error
0	0.024482	0.030230
1	0.028355	0.035046
2	0.027317	0.032791
3	0.027177	0.033037

Figure 6: Validation set Error Comparison form

158 16630 true positives, 17108 true negatives, 14124 false positives, and 71 false negatives. IoU is only  
 159 scored for image-wise true positive cases. The distribution of these scores is in figure 13.

## 160 5.4 Failures

161 The main failure of our model is in the presence of external classes of objects, such as a harbor. Our  
 162 segmentation model is strong when the background is composed of ambient natural conditions, but in  
 163 its current training, other manmade and geographic structures tend to be classified as ships. This is  
 164 likely due to sparsity of data; while we have many training examples of open water, fog, and clouds,  
 165 examples of ships are much more rare, and ships in the presence of land are even rarer. For this  
 166 reason, differentiation between ships, land, and manmade structures is less finely tuned.

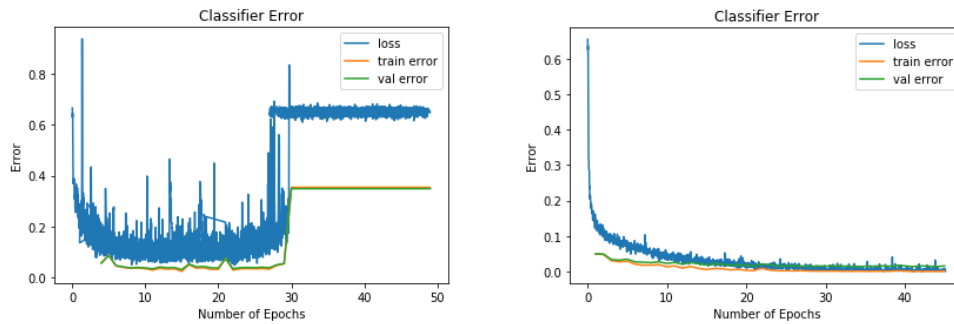


Figure 7: Classifier Error for Minibatch sizes  $B = 8$  and  $B = 32$

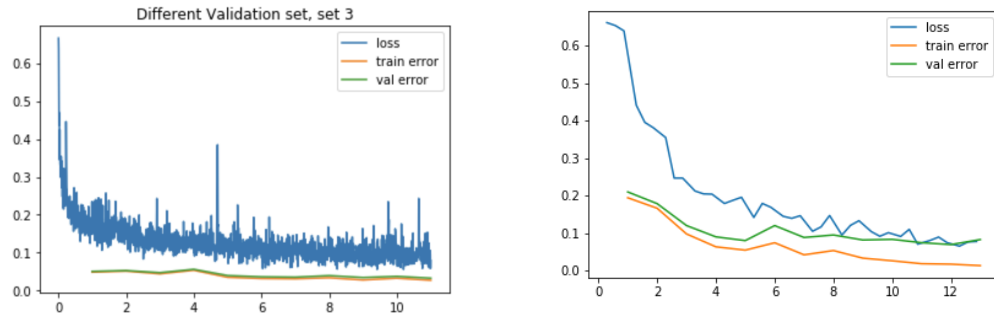


Figure 8: AlexNet and VGGNet Error

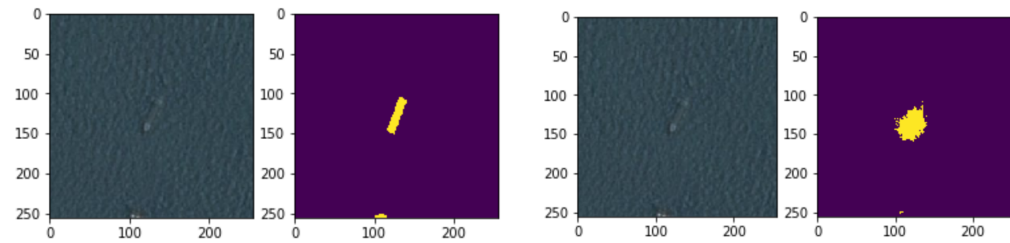


Figure 9: Ship camouflaged in water. Original images on left. Ground truth center left. Segmentation result right.

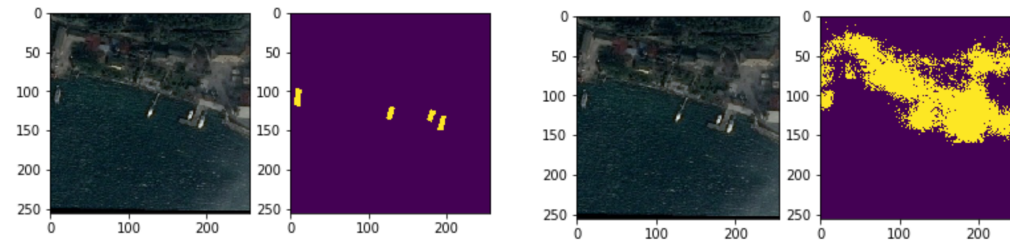


Figure 10: Ships Docked in Harbor. Original images on left. Ground truth center left. Segmentation result right.

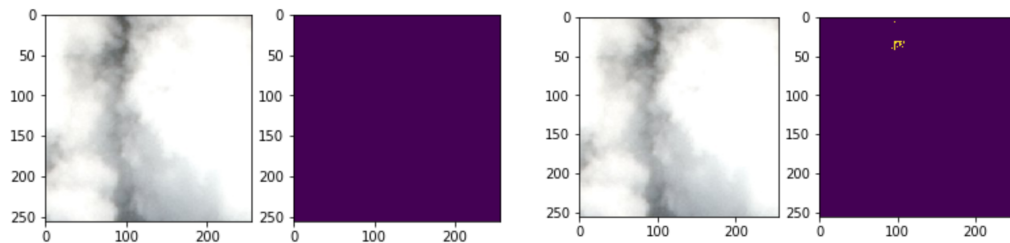


Figure 11: Scene blocked by clouds. Ground truth center left, Segmentation result right.

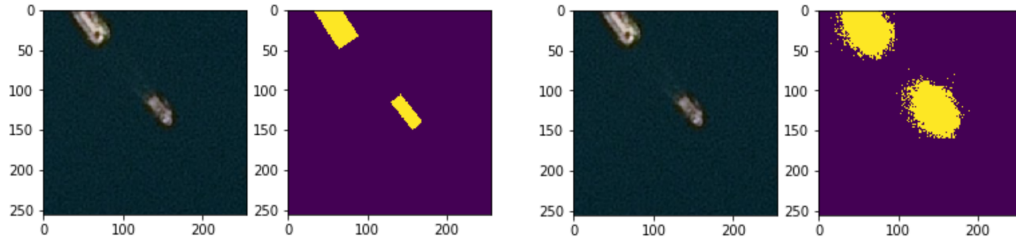


Figure 12: Two ships on open water. Ground truth center left, Segmentation result right.

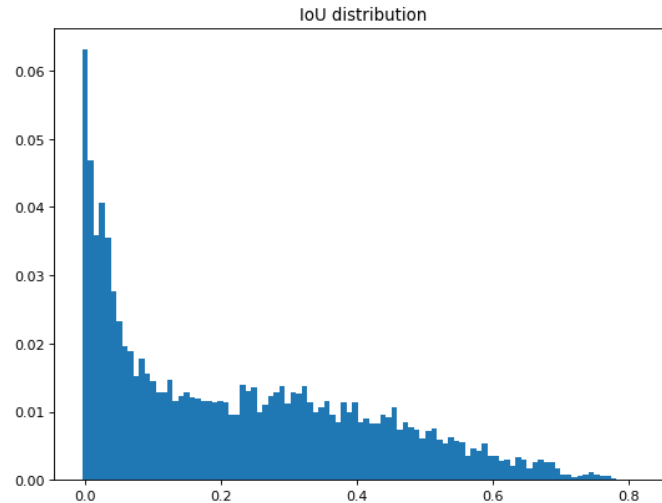


Figure 13: Distribution of intersection over union on validation set.

167 The implementation of AlexNet in a fully convolutional network for segmentation failed to learn  
 168 during training despite having good classifier performance. This may possibly be due to shallowness  
 169 of AlexNet. Unlike, VGGNet, AlexNet does not have many pooling layers from which to produce  
 170 different resolution feature maps that can be upsampled and fused together.

## 171 6 Discussion

### 172 6.1 Conclusion

173 Fully convolutional networks presents a viable solution for semantic segmentation. However, it  
 174 suffers from having to train a classifier separately in order to learn the proper feature maps, which  
 175 complicates training. Unlike a classifier, the network is not completely sequential, making the  
 176 implementation more complicated.

### 177 6.2 Learning

178 Over the course of this project, our group became familiar with two conventional CNN architectures,  
 179 AlexNet and VGGNet. We were able to observe differences in both network performance and network  
 180 training; for example, we noted that VGGNet seemed to train 10 times slower than AlexNet, perhaps  
 181 due to the relative simplicity of the AlexNet model and parameter set.

182 Our group was also able to understand the problem in the context of transfer learning. In our  
 183 architecture, we first trained the network as a binary classifier, leaving off upsampling and final  
 184 convolution layers. Then, we attached the segmentation layers, tuning the later end of the network on  
 185 the segmentation problem while transferring the learned classifier to the front.



186 A third principle learned over the course of this project was fusion of outputs at different layers of a  
187 network. This forced us to carefully consider the effects of stride and upsampling on tensor shape,  
188 and through this we were able to evaluate positions that we should sample from the network layers to  
189 later fuse and create a mixed-resolution output.

## 190 6.3 Challenges

191 **Bad data** There are some images in the dataset which do not correspond with the ground truth.

192 **Really small object** The ships in some of the images are really small. Sometimes we can only see  
193 bow waves which make the classification and segmentation tasks significantly harder.

194 **Clouds or haze** There are images that are completely covered by clouds, which confuses the  
195 classifier. Other photos have fog or haze camouflaging the ships with the water. Land A small amount  
196 of the images in the dataset includes only land and no water. Training time Doing two stages of  
197 training for each model takes more time to train and is more interactive. This is especially challenging  
198 when the trying to save and keep track of two sets of parameters so that they can be loaded separately.

## 199 6.4 Future Work

200 One approach for future experimentation would be applying progressively stronger weights to training  
201 samples containing ships in our loss function, due to the relative volume of background training  
202 samples, and creating a third class of images for non-oceanic backgrounds. By applying a greater  
203 learning weight to the ship and non-oceanic examples, this would be similar to effectively running  
204 the sample through backpropagation multiple times, compounding its learning. This is important to  
205 give balance to the sparse class presence in the given training data. While we tried using weighting in  
206 our loss function to give greater value to ship training samples, we could work to further refine this  
207 weight to have optimal effect for our sparse data.

208 Another modification we could make to our model would be image pre-processing and data aug-  
209 mentation. Providing rotations of our dataset as input would create four times the current number  
210 of training points, and this would be useful as our current training appears to show difficulty in  
211 segmenting ships in the proper orientation. The ships are segmented as elliptical figures, and perhaps  
212 if trained on ships covering a span of orientations, the network could be better tuned to a ships actual  
213 location versus a particle-like range of potential positions.

214 A more compelling result may be possible if there was enough time or compute power to run through  
215 the whole dataset without cropping and centering the ship in the images.

## 216 7 Appendix

### 217 7.1 References

- 218 [1] Garcia-Peraza-Herrera, Luis & Li, Wenqi & Gruijthuijsen, Caspar & Devreker, Alain & Attilakos, George  
219 & Deprest, Jan & Vander Poorten, Emmanuel & Stoyanov, Danail & Vercauteren, Tom & Ourselin, Sébastien.  
220 (2016). Real-Time Segmentation of Non-Rigid Surgical Tools based on Deep Learning and Tracking. Lecture  
221 Notes in Computer Science. 10170. 10.1007/978-3-319-54057-3-8.
- 222 [2] Alex Krizhevsky, Ilya Sutskever, & Geoffrey E. Hinton. 2012. ImageNet classification with deep convolu-  
223 tional neural networks. In Proceedings of the 25th International Conference on Neural Information Processing  
224 Systems - Volume 1 (NIPS'12), F. Pereira, C. J. C. Burges, L. Bottou, & K. Q. Weinberger (Eds.), Vol. 1. Curran  
225 Associates Inc., USA, 1097-1105.
- 226 [3] Evan Shelhamer, Jonathan Long, & Trevor Darrell. 2017. Fully Convolutional Networks for Se-  
227 mantic Segmentation. IEEE Trans. Pattern Anal. Mach. Intell. 39, 4 (April 2017), 640-651. DOI:  
228 <https://doi.org/10.1109/TPAMI.2016.2572683>
- 229 [4] Simonyan, Karen, & Andrew Zisserman. "Very deep convolutional networks for large-scale image recogni-  
230 tion." arXiv preprint arXiv:1409.1556 (2014).

## 231 7.2 Github Repository Link

232 The project is hosted on GitLab (<https://gitlab.com/rocketwings/ece285f>).

233

## 234 7.3 Github Repository Readme

235 # Fully-Convolutional Networks: Ship Detection via Semantic Segmentation

236

237 ## Description

238 Implementation of fully convolutional networks for semantic segmentation for the Kaggle Airbus

239

240 ## Requirements

241 Install package 'imageio' as follows:

242 \$ pip install --user imageio

243

244 ## Code organization

245 demo.ipynb — Run a demo of our code (reproduce IoU scores of our report)

246 vggnet/VggNet-classifier.ipynb — Train the VGGNet classifier

247 vggnet/FCN\_Training.ipynb — Run the fine-tuning training for the fully convolutional network

248 vggnet/fcn\_training.pth — Parameters for VGGNet FCN

249 alex/Train\_Alex\_FCN.ipynb — Failed AlexNet implementation

250 vgg2/ — Failed VGGNet implementation

251 util/ASDC\_loader.py — module for interfacing with the dataset

252 util/ModelUtils.py — misc utilities