

## Analysis of Unsupervised Segmentation Methods using SLIC Refinement

Presented by: Md. Saiful Islam Sajol (30042836)

Examiners: Prof. Dr.-Ing. Andreas Schwung  
Prof. Dr. Dominik Aufderheide

Supervisor: Jan Niclas Reimann, M.Sc.

# Outline

## 1. Image Segmentation

- What is image?
- Different Types of Segmentation and Applications
- Deep learning based Image Segmentation Approaches

## 2. Motivation

- Deep learning based Image Segmentation Approaches
- Architecture for Unsupervised Approach

## 3. Method Description and Problem Statement

- Base Architecture
- Results of Base Architecture
- Challenge and Limitations

## 4. Proposed Architecture 1 : FCN with Classification Approach

- Results of FCN with Classification Approach

## 5. Proposed Architecture 2 : FCN with Autoencoder Approach

- Results of FCN with Classification Approach

# Outline

5. Comparisons of Different Architecture
6. Challenges and Future Work
7. Conclusion
8. Appendix

# Image Segmentation

## What is an Image?

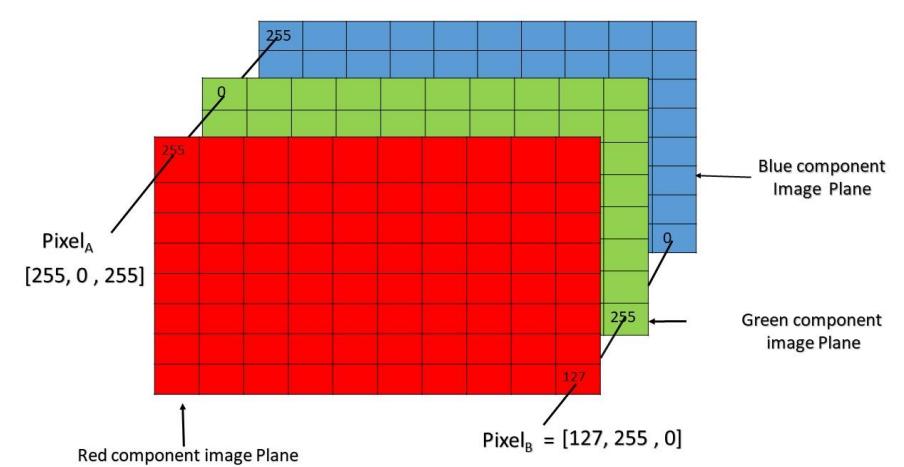
- An image is an array of numbers.



0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29	
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0	
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1	
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49	
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36	
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62	
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0	
0	87	252	250	248	215	60	0	1	121	252	255	248	144	6	0	
0	13	113	255	255	245	255	182	181	248	252	242	208	36	0	19	
1	0	5117	251	255	241	255	247	255	241	162	17	0	7	0		
0	0	0	4	58	251	255	246	254	253	259	120	11	0	1	0	
0	0	4	97	255	255	255	248	252	255	244	255	182	10	0	4	
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9	0	
0	111	255	242	255	158	24	0	0	6	39	255	232	230	56	0	
0	218	251	250	137	7	11	0	0	0	2	62	255	250	125	3	
0	173	255	255	101	9	20	0	13	3	13	182	251	245	61	0	
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4	
0	18	146	250	255	247	255	255	255	249	255	240	255	129	0	5	
0	0	23	113	215	255	250	248	255	255	248	248	118	14	12	0	
0	0	6	1	0	52	153	233	255	252	147	37	0	0	4	1	
0	0	5	5	0	0	0	0	0	0	14	1	0	6	6	0	

Grey Image

```
0 2 15 0 0 11 10 0 0 0 0 9 9 0 0 0  
0 0 0 4 60 157 236 255 255 177 95 61 32 0 0 29  
0 10 16 119 238 255 244 245 243 250 249 255 222 103 10 0  
0 14 170 255 255 244 254 255 253 245 255 249 253 251 124 1  
2 98 255 228 255 251 254 211 141 116 122 215 251 238 255 49  
13 217 243 255 155 33 226 52 2 0 10 13 232 255 255 36  
16 229 252 254 49 12 0 0 7 7 0 70 237 252 235 62  
6 141 245 255 212 25 11 9 3 0 115 236 243 255 137 0  
0 87 252 250 248 215 60 0 1 121 252 255 248 144 6 0  
0 13 113 255 255 245 255 182 181 248 252 242 208 36 0 19  
1 0 5117 251 255 241 255 247 255 241 162 17 0 7 0  
0 0 0 4 58 251 255 246 254 253 259 120 11 0 1 0  
0 0 4 97 255 255 255 248 252 255 244 255 182 10 0 4  
0 22 206 252 246 251 241 100 24 113 255 245 255 194 9 0  
0 111 255 242 255 158 24 0 0 6 39 255 232 230 56 0  
0 218 251 250 137 7 11 0 0 0 2 62 255 250 125 3  
0 173 255 255 101 9 20 0 13 3 13 182 251 245 61 0  
0 107 251 241 255 230 98 55 19 118 217 248 253 255 52 4  
0 18 146 250 255 247 255 255 255 249 255 240 255 129 0 5  
0 0 23 113 215 255 250 248 255 255 248 248 118 14 12 0  
0 0 6 1 0 52 153 233 255 252 147 37 0 0 4 1  
0 0 5 5 0 0 0 0 0 0 14 1 0 6 6 6 0 0
```



Pixel of an RGB image are formed from the corresponding pixel of the three component images

<https://www.geeksforgeeks.org/matlab-rgb-image-representation/>

Yamashita, R., Nishio, M., Do, R.K.G. et al. Convolutional neural networks: an overview and application in radiology. *Insights Imaging* 9, 611–629 (2018).  
<https://doi.org/10.1007/s13244-018-0639-9>

# Image Segmentation

## Four Main Categories of Computer Vision Tasks

### Image Classification

- Classify the object within an image



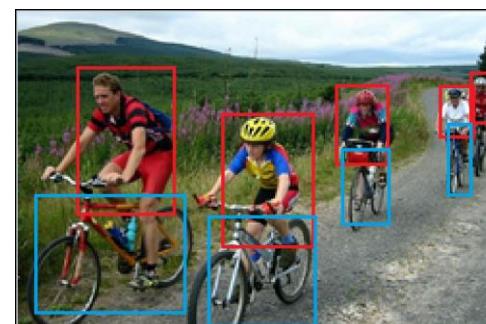
### Semantic Segmentation

- Classify the object class for each pixel (pixelwise classification)



### Object Detection

- Classify and detect the objects with bounding boxes



### Instance Segmentation

- Label each foreground pixel with object and instance



# Motivation

## Popular Deep Learning Based Image Segmentation Models

- Fully convolutional networks
- Convolutional models with graphical models
- Convolutional models with active contour models
- Dilated convolutional models and DeepLab family
- Encoder-decoder based models
- Multi-scale and pyramid network based models
- Generative models and adversarial training
- R-CNN based models (for instance segmentation)
- Attention-based models
- Recurrent neural network based models



Supervised Learning  
Methods

*But supervised labeling is expensive , time consuming and difficult to obtain!!*

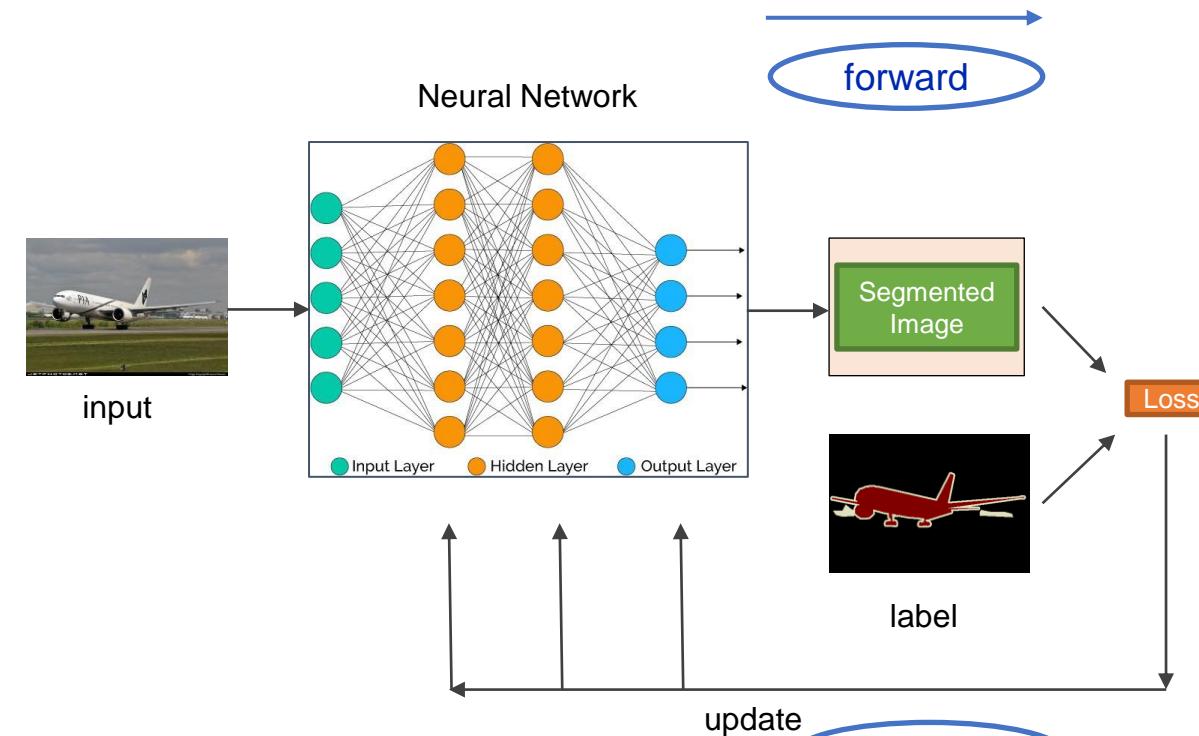


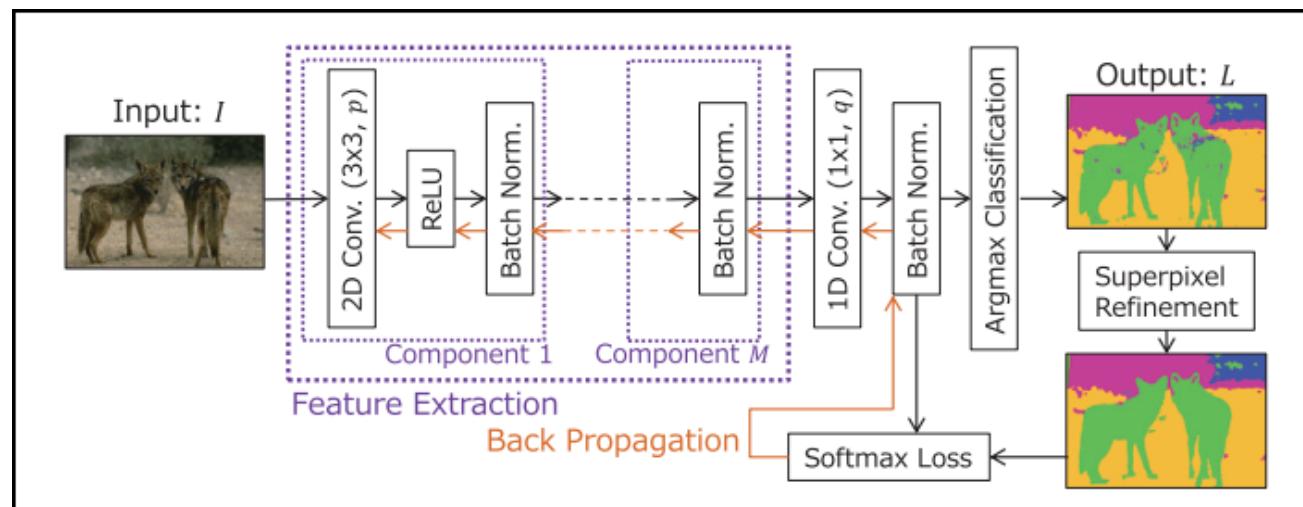
Figure: A Basic Neural Network for Supervised Segmentation

*So, we tried to find an unsupervised way*

# Method Description

## Architecture for Unsupervised Approach

We Used CNN for Unsupervised Image Segmentation with SLIC Refinement to calculate Loss for Backpropagation



A. Kanezaki, "Unsupervised Image Segmentation by Backpropagation," 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Calgary, AB, Canada, 2018, pp. 1543-1547, doi: 10.1109/ICASSP.2018.8462533.

### Main Components

- ❑ Convolutional Operation
- ❑ Argmax Classification
- ❑ Superpixel Refinement (SLIC)
- ❑ Softmax Loss/ Cross Entropy Loss

# Method Description

## Architecture for Unsupervised Approach

### Argmax Classification

$$g(x) = x^2$$

$$g(1) = 1$$

$$g(2) = 4$$

$$g(3) = 9$$

$$g(4) = 16$$

$$\text{Argmax}(g(x)) = 4$$



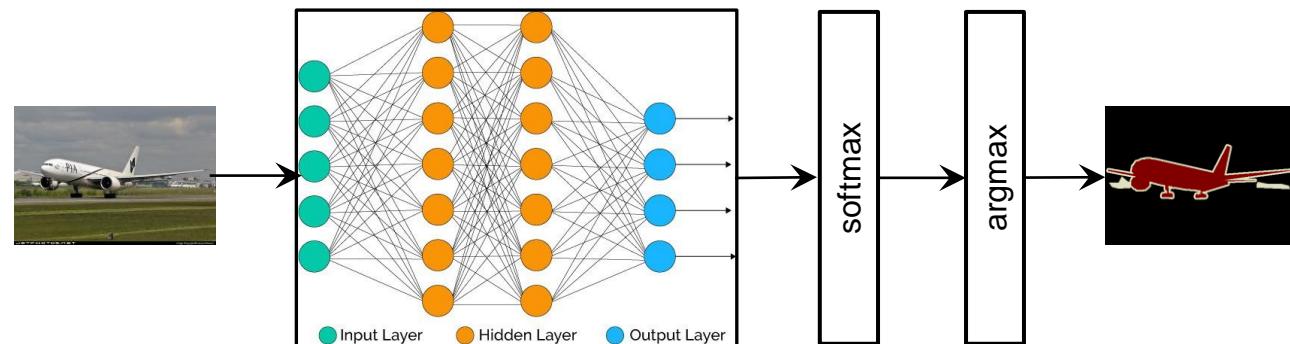
segmented →

1: Person  
2: Purse  
3: Plants/Grass  
4: Sidewalk  
5: Building/Structures

3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	5	5	5	5	5	
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	5	5	5	5	5
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	5	5	5	5	5
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	5	5	5	5	5
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	5	5	5	5	5
5	5	3	3	3	3	3	3	3	1	1	3	3	3	3	3	3	3	3	3	3	5	5	5	5	5	5
4	4	3	4	1	1	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4	5	5	5	5	5
4	4	3	4	1	1	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4	5	5	5	5	5
4	4	4	1	1	1	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
3	3	3	1	1	1	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
3	3	3	1	2	2	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
3	3	3	1	2	2	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4

Semantic Labels

Image source: <https://www.jeremyjordan.me/semantic-segmentation/>



# Method Description

## Architecture for Unsupervised Approach

### Superpixel

A superpixel is a group of pixels that shares common characteristics (like color, pixel intensity)

- Carry more information
- Similar visual color
- Convenient and compact representation of image



SLIC Superpixels Compared to State-of-the-art  
Superpixel Methods . Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi,  
Pascal Fua, and Sabine Süsstrunk

### SLIC Superpixel

Simple Linear Iterative Clustering (SLIC) adapts **kmeans** clustering to generate superpixels based on their **color similarity and proximity**

### Why SLIC Superpixel?

- Assigning same label to all pixels in a superpixel
- Very good edge detector
- Outperforms other state of the art algorithms

# Method Description

## Architecture for Unsupervised Approach

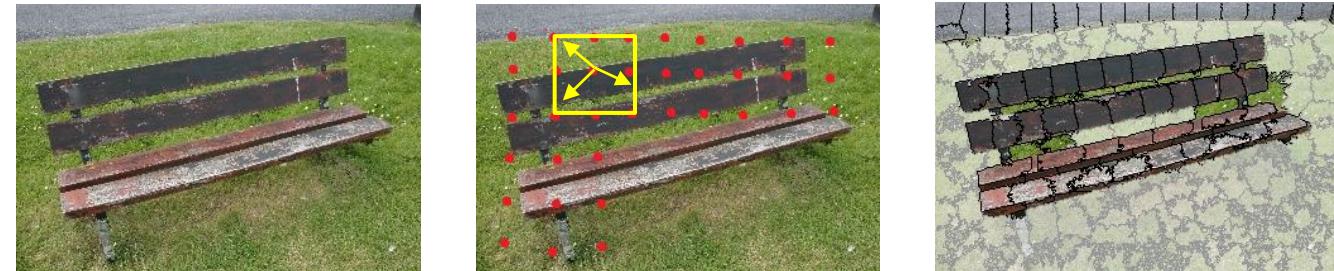
### How SLIC algorithm works?

- Convert RGB image to CIELAB Color space (5d vector)  
 $\text{RGB} \rightarrow [l \ a \ b \ x \ y]$   
 where  $[l \ a \ b]$  is the pixel color vector  
 $x \ y$  is the pixel position
- Take no. of superpixel,  $K$  and Calculate  $S = \sqrt{N/K}$
- Create a grid of initial centers with  $S$  spacing
- Calculate gradients and move the grid centres which are on edges
- Perform  $K$  means locally

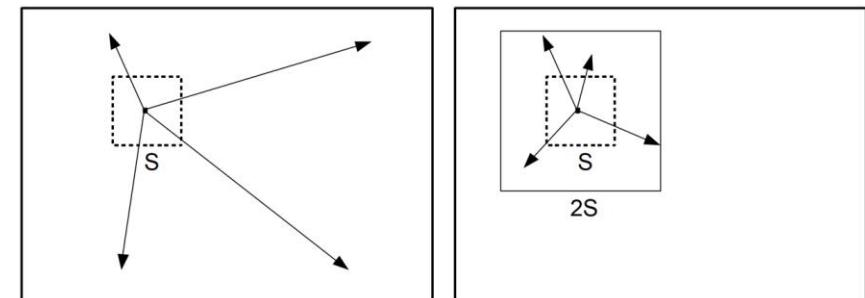
$$d_c = \sqrt{(l_j - l_i)^2 + (a_j - a_i)^2 + (b_j - b_i)^2}$$

$$d_s = \sqrt{(x_j - x_i)^2 + (y_j - y_i)^2}$$

$$D = \sqrt{d_c^2 + \left(\frac{d_s}{S}\right)^2} \ m^2$$



N	Number of pixels in the input image
K	Number of Superpixels used to segment the input image
N\K	Approximate number of pixels in each superpixel
$S=\sqrt{(N/K)}$	Approximate length of a superpixel



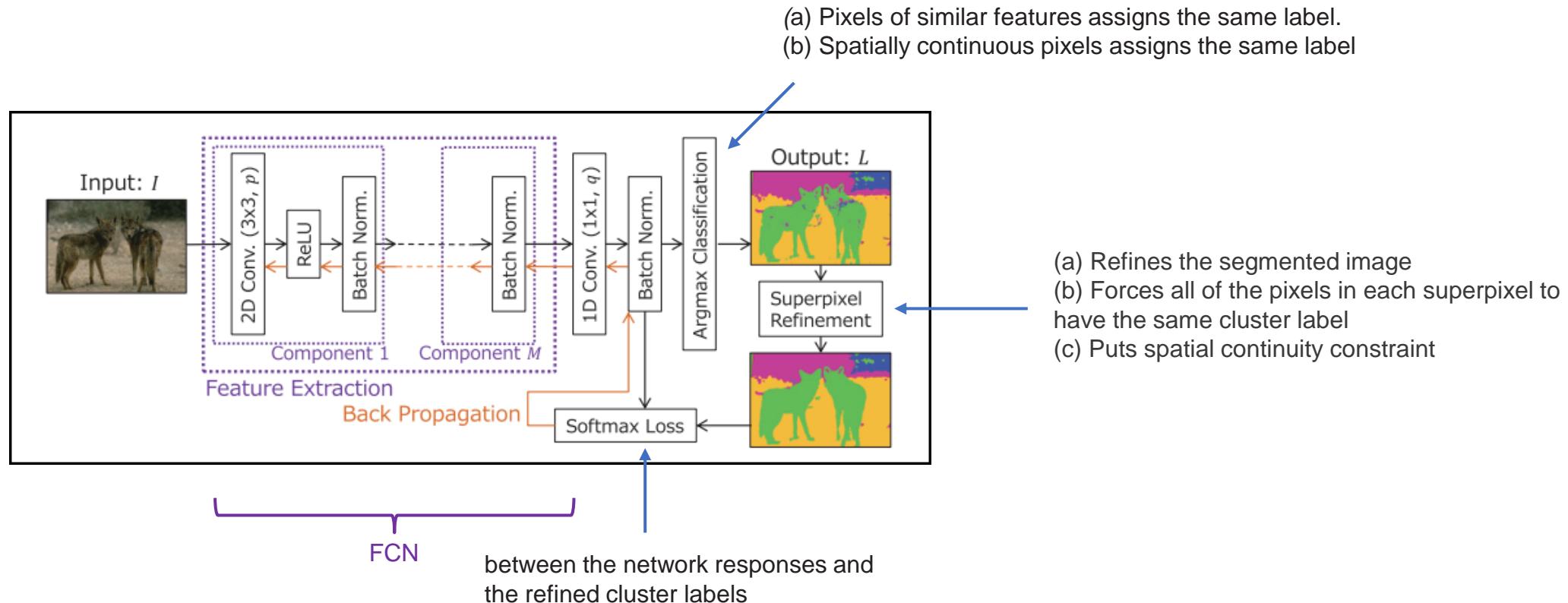
(a) standard  $k$ -means searches the entire image

(b) SLIC searches a limited region

# Method Description

## Architecture for Unsupervised Approach

### FCN for Unsupervised Image Segmentation with SLIC Refinement



The above architecture is used to reproduce the results

# Datasets

CIFAR10



- 60,000 images
- 32x32 pixels
- 50,000 train images
- 10,000 test images
- 10 classes

STL10



- 13,000 images
- 96x96 pixels
- 5,000 train images
- 8,000 test images
- 10 classes

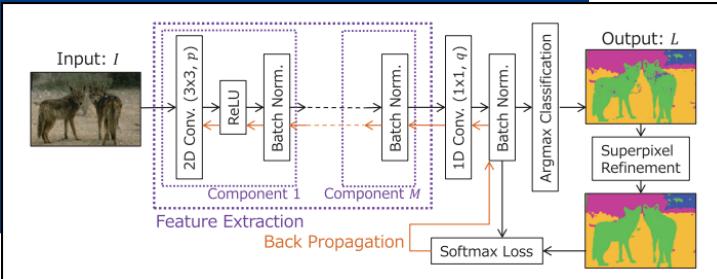
Pascal VOC 2012



- 2,913 images
- 500x350 pixels (minimum)
- 1,464 train images
- 1,449 test images
- 20 classes
- We down sampled to

224x224pixels   
University of Applied Sciences

# Base Architecture Results



Dataset: CIFAR10

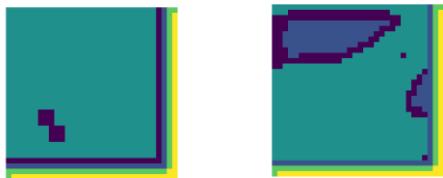
Main Image



Argmax Output

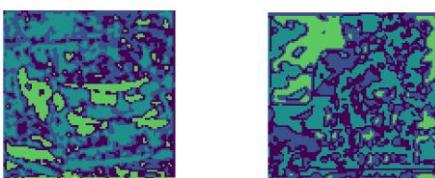


SLIC Output



Epochs=30, LR=0.001,  
BS=16, Opt=SGD

Dataset: STL10



Epochs=30, LR=0.0001,  
BS=5, Opt=SGD

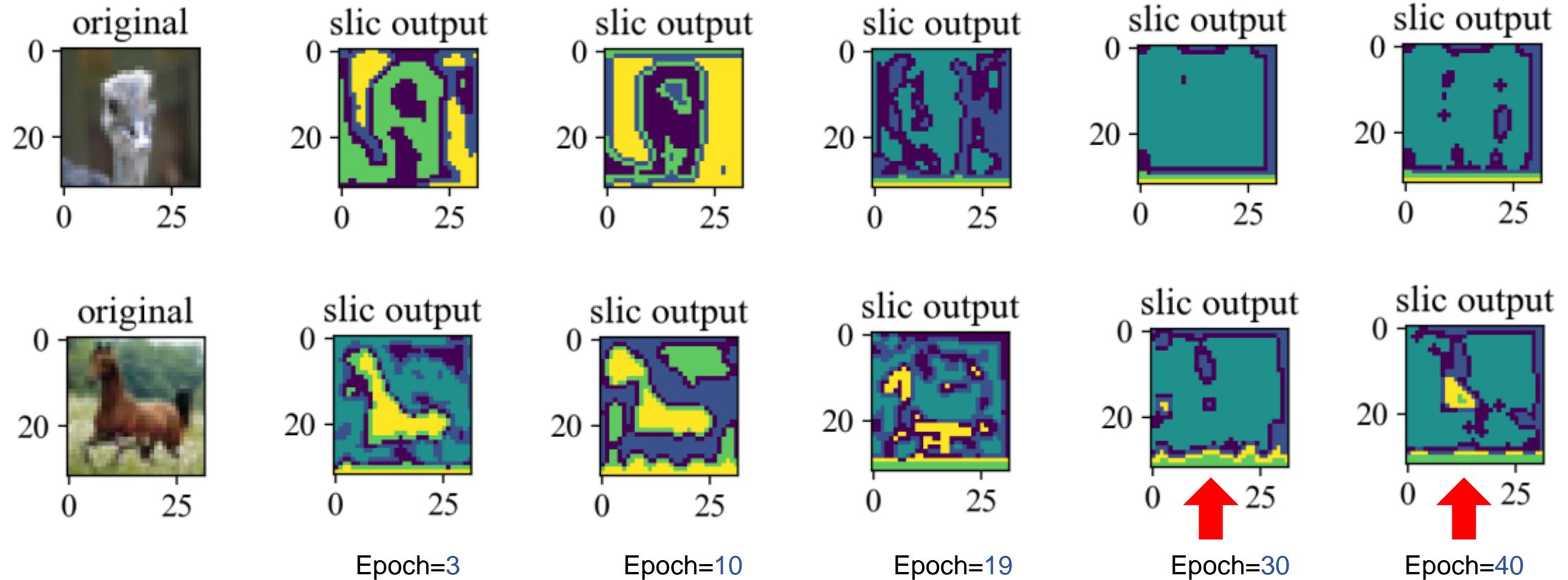
Dataset: Pascal VOC



Epochs=30, LR=0.001,  
BS=16, Opt=Adam

# Base Architecture Challenges and Limitations

Segmentation During the Training

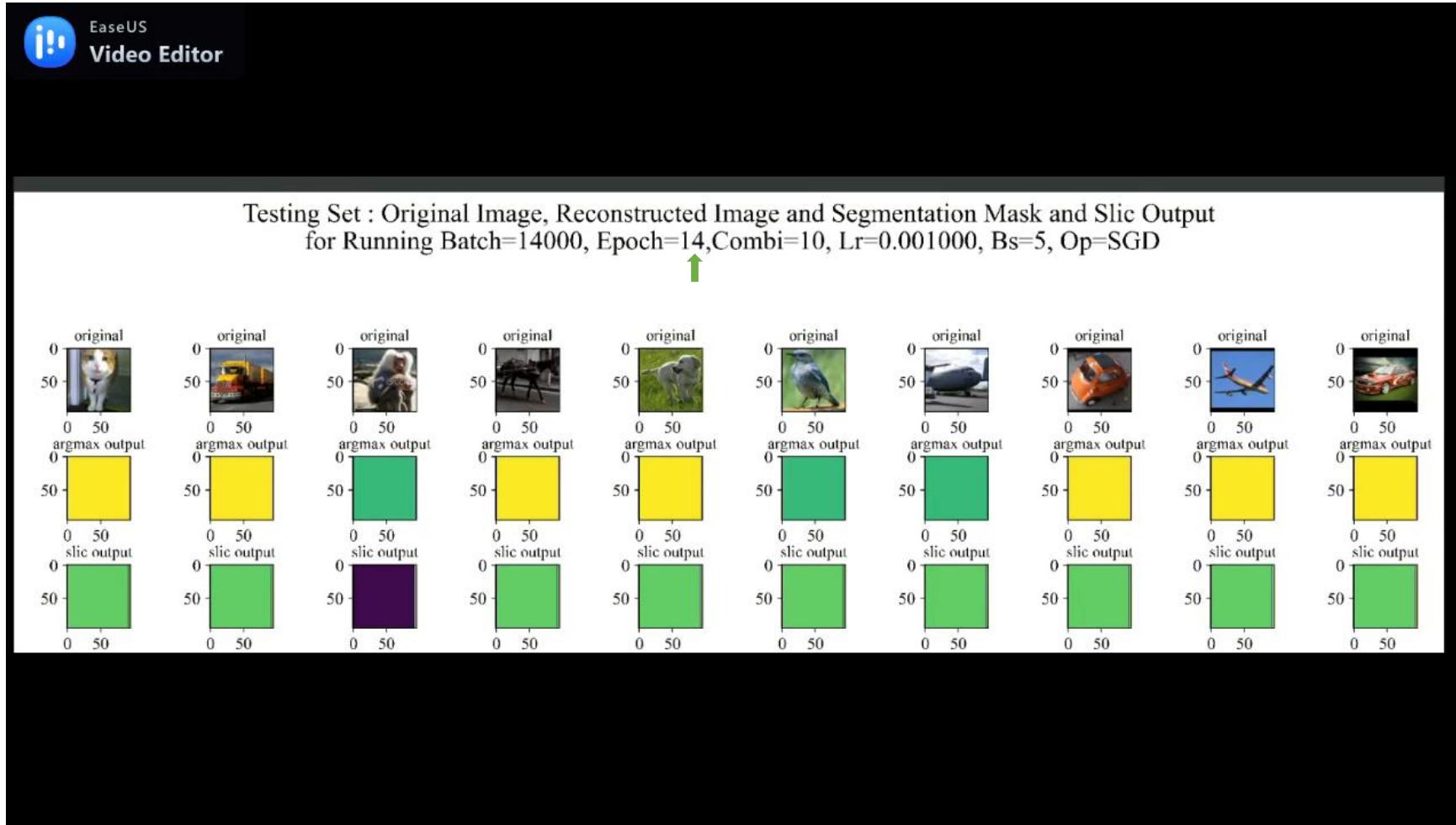


# Base Architecture

## Problem Description of the Base Architecture

# Segmentation During the Training

## Dataset: STL10



## Challenges and Limitations

- The model is not robust enough
  - Segmentations were merging over training
  - Greatly affected by parameters
  - Results are different at multiple random restarts
  - No mechanism to restricts learned experience

# Proposed Architectures for Problem Solution

We propose two possible approaches to solve these problems

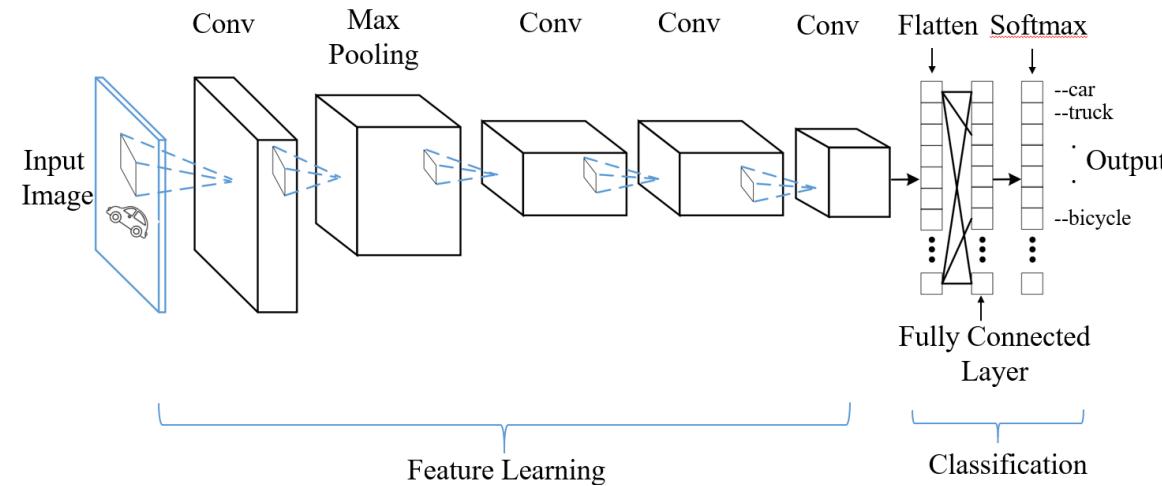
1. FCN with Classification Approach
2. FCN with Autoencoder Approach

## Image Classification

Image classification is where a network classifies images into given classes or labels

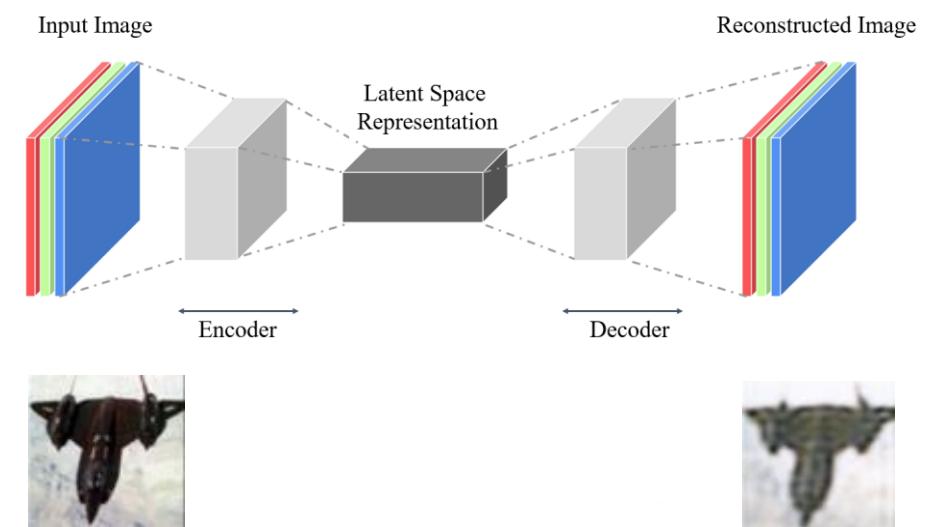
A class is essentially a label, for instance,

'car' → 0  
'truck' → 1  
'van' → 2



## Autoencoder

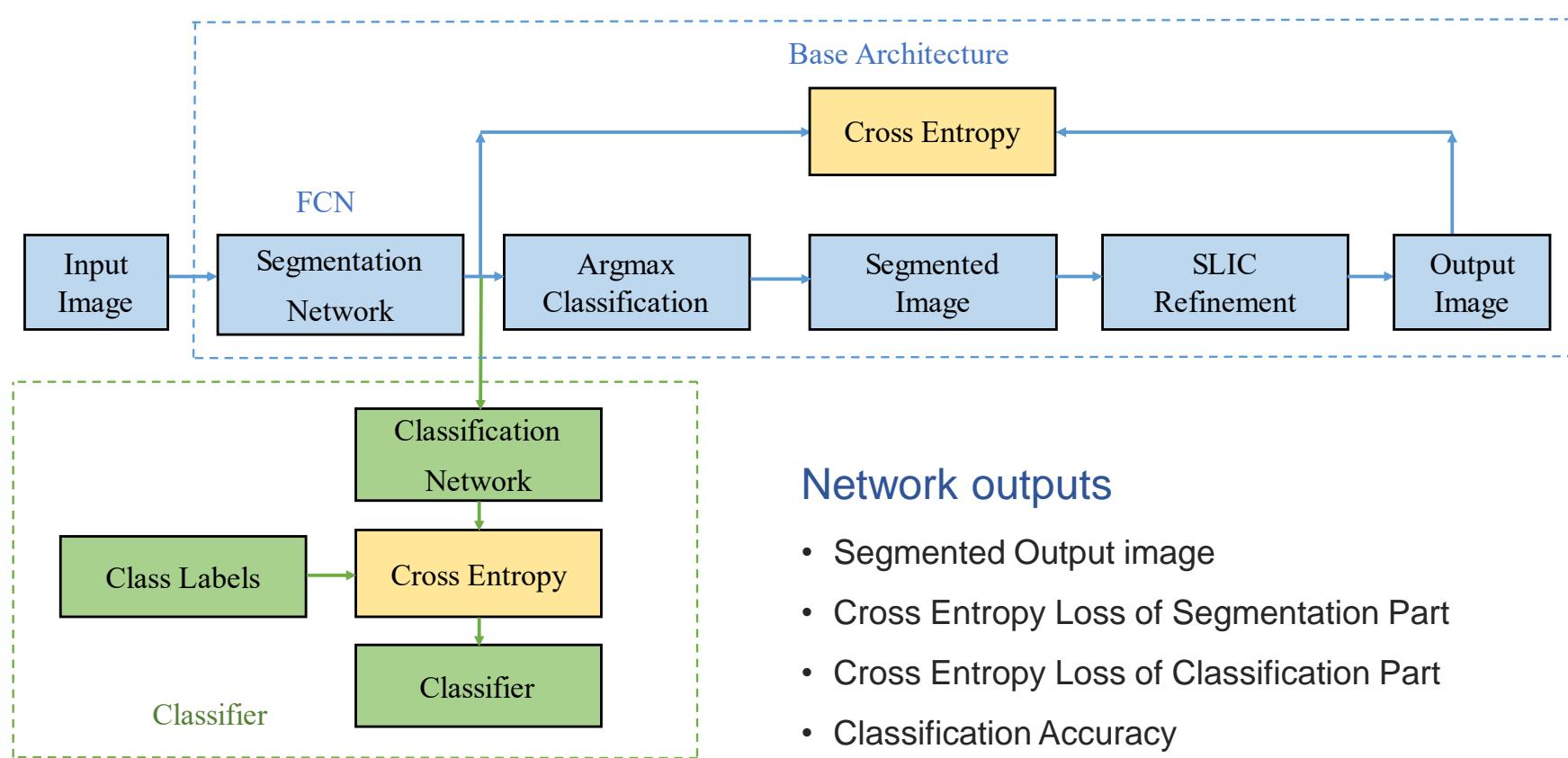
The autoencoder extracts features, compresses all the important attributes and again reconstructs the original input.



# Proposed Architecture 1 for Problem Solution

## Proposed Architecture 1: FCN with Classification Approach

Hyperparameters for Network Tuning



Hyperparameters	Range / Value
Epochs	30, 40, 50, 60, 100
Normalization	minmax (0 to 1)
Activation Function	ReLU
Batch Size	5, 16, 32, 50, 100
Learning Rate	0.1, 0.01, 0.001, 0.0001
Optimizer	Adam, SGD
Initializer	Xavier Initializer
Scheduler	StepLR

# Results

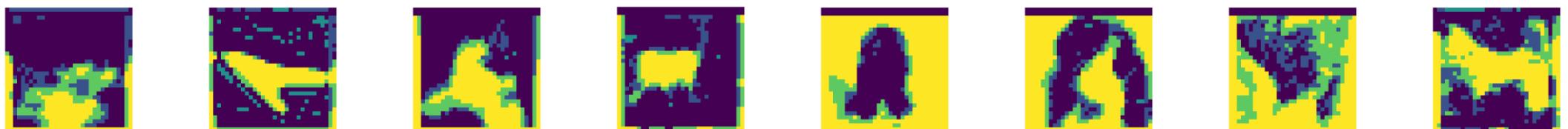
## Proposed Architecture 1: FCN with Classification Approach

Dataset: CIFAR10

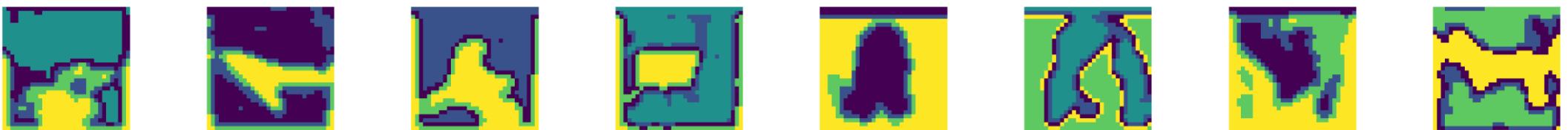
Main Image



Argmax Output



SLIC Output



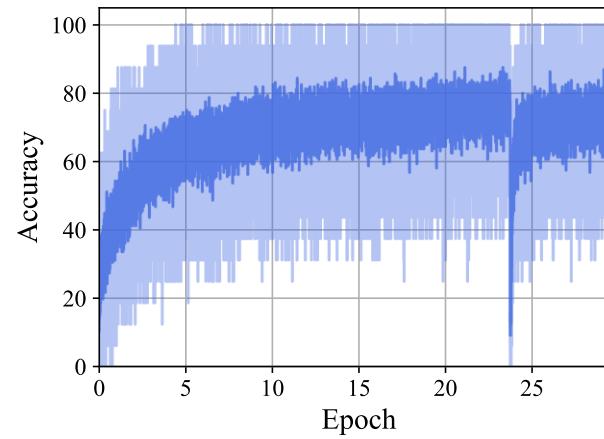
Epochs=30, LR=0.01, BS=16, Opt=SGD

# Results

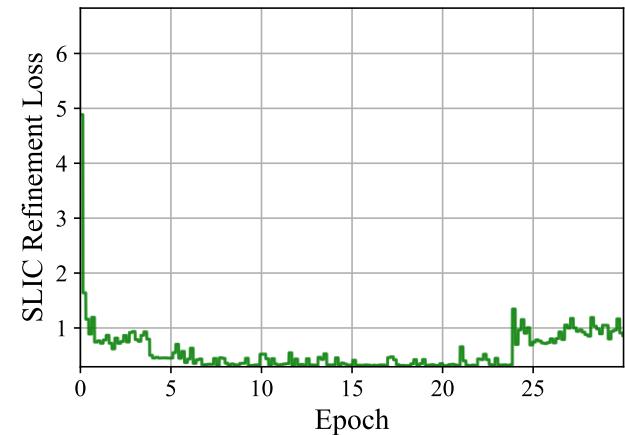
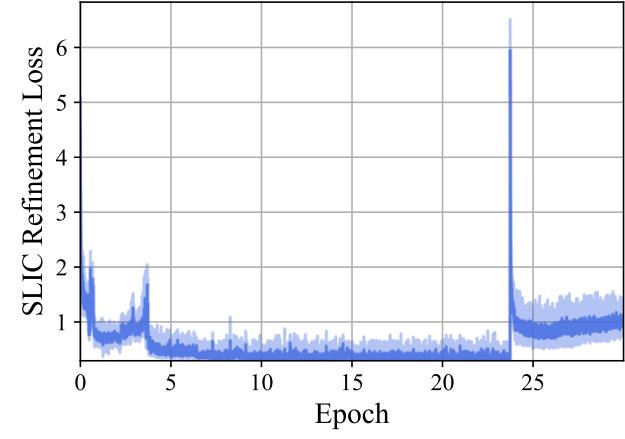
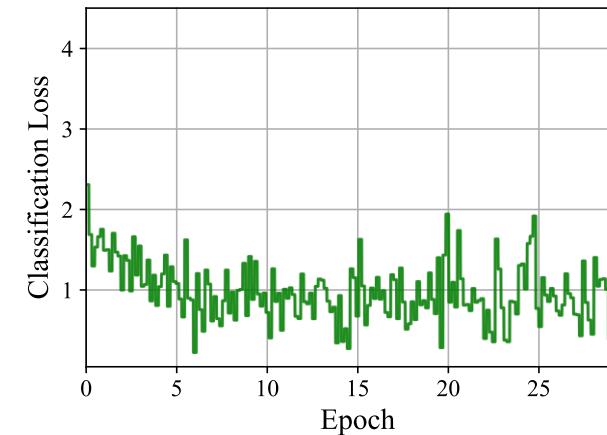
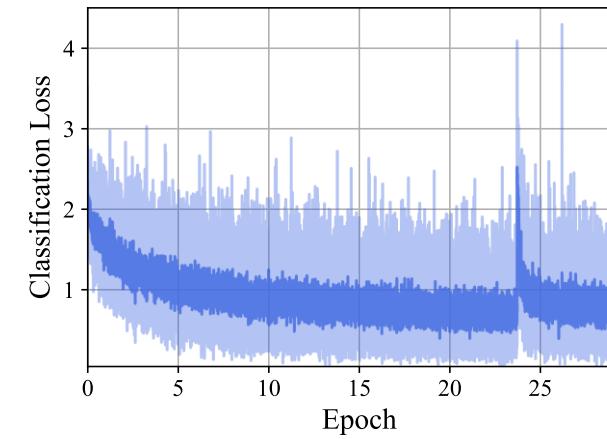
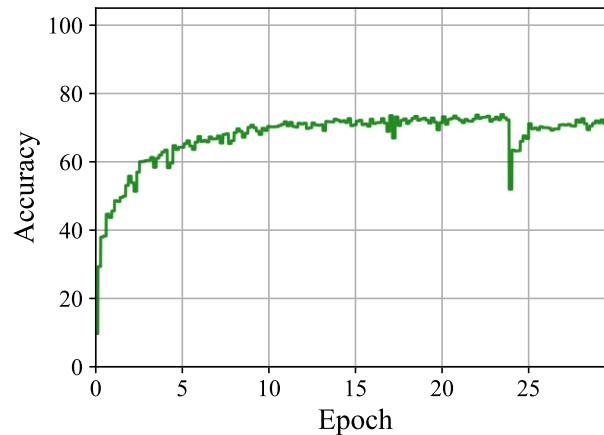
## Proposed Architecture 1: FCN with Classification Approach

Dataset: CIFAR10

Training



Testing



Epochs=30, LR=0.01, BS=16, Opt=SGD

# Results

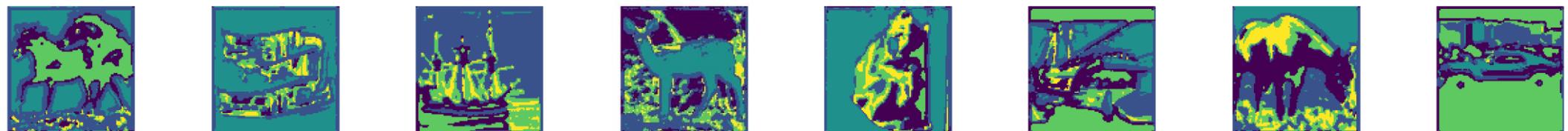
## Proposed Architecture 1: FCN with Classification Approach

Dataset: STL10

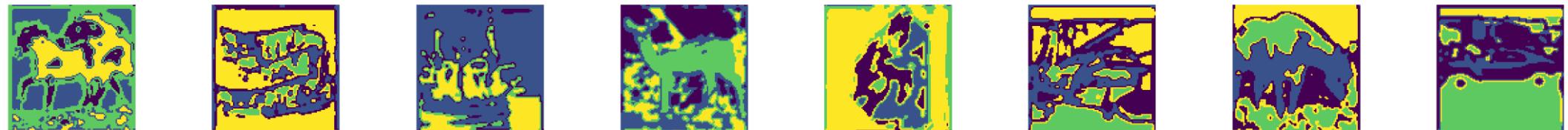
Main Image



Argmax Output



SLIC Output



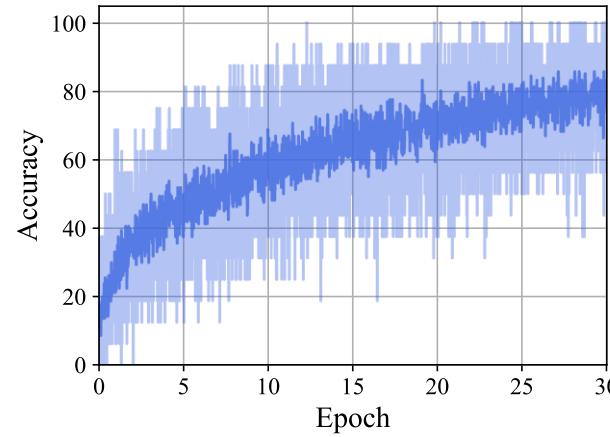
Epochs=30, LR=0.001, BS=16, Opt=SGD

# Results

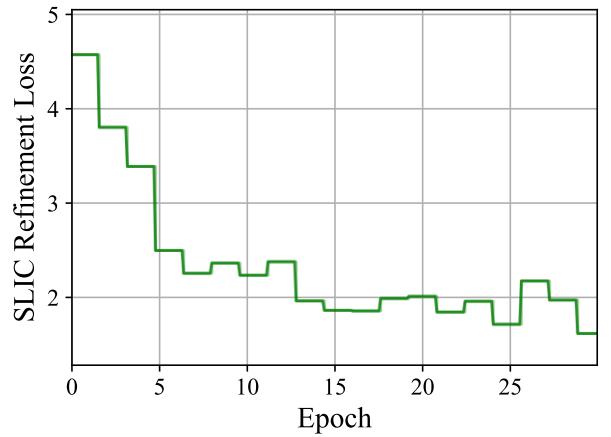
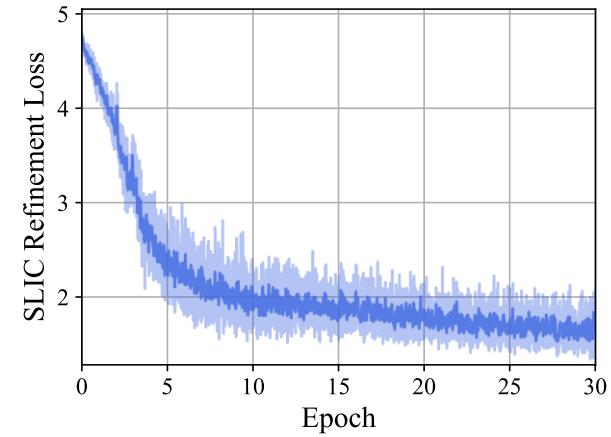
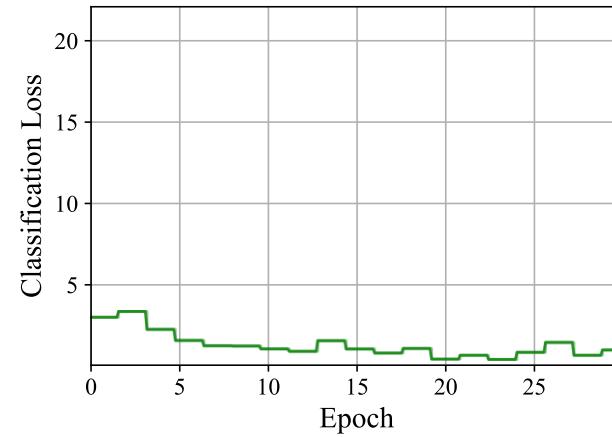
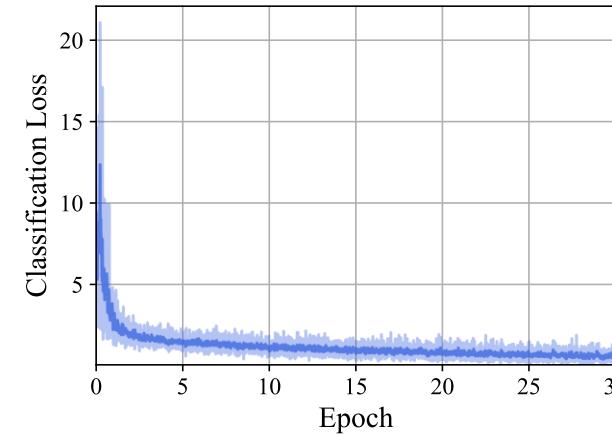
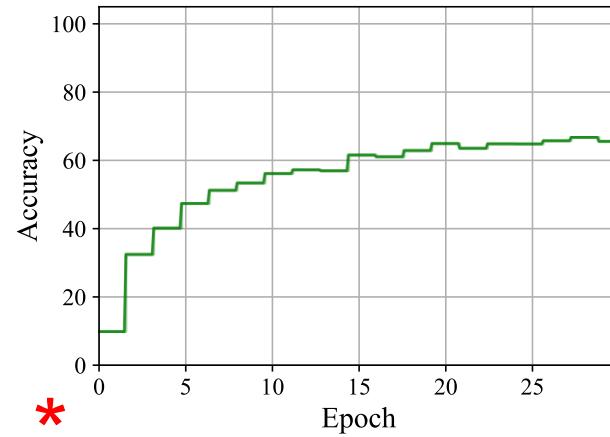
## Proposed Architecture 1: FCN with Classification Approach

Dataset: STL10

Training



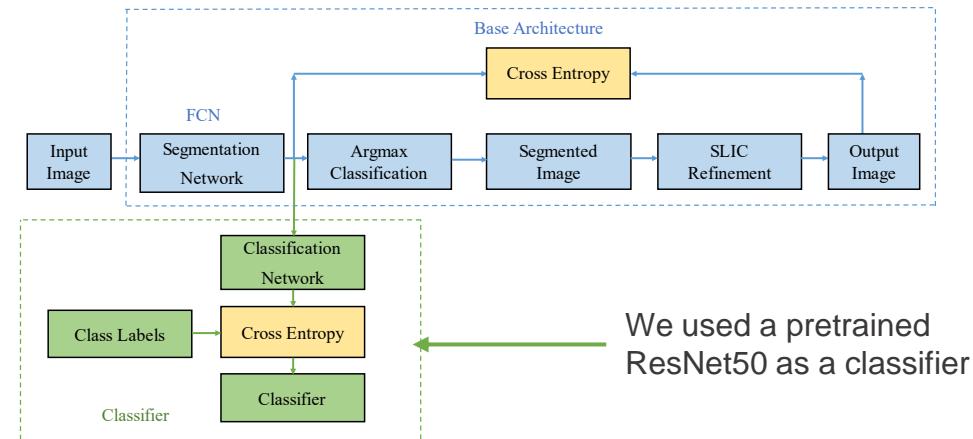
Testing



# Results

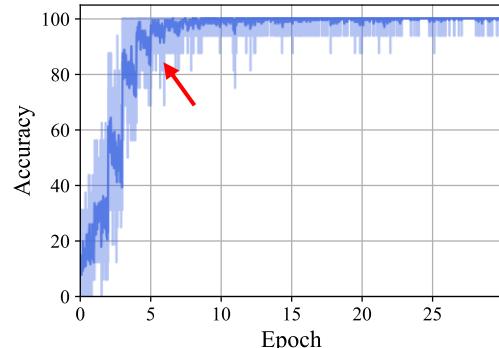
## Proposed Architecture 1: FCN with Classification Approach

Dataset: Pascal VOC

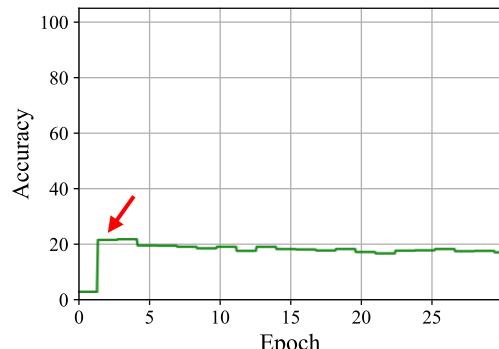


- Initially four convolutional units were used for classification part of the Pascal VOC dataset
- The model was highly overfitting
- So, we used **a pretrained ResNet50** model (pretrained on ImageNet Dataset) to avoid over fitting , and used **Transfer Learning** method to boost up the training process

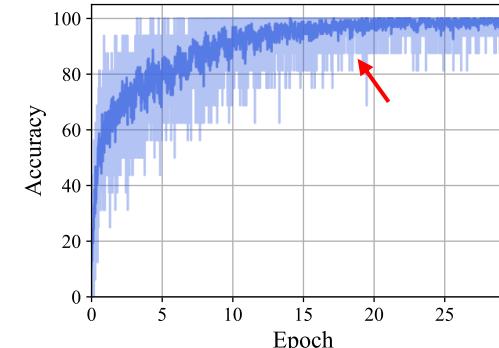
Training



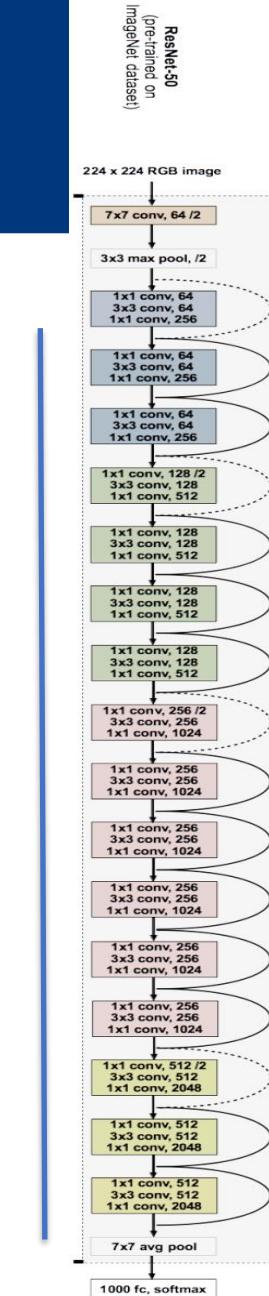
Testing



Before using ResNet50



After using ResNet50



# Results

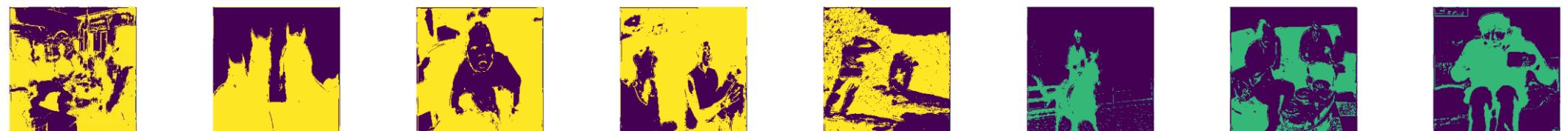
## Proposed Architecture 1: FCN with Classification Approach

Dataset: Pascal VOC

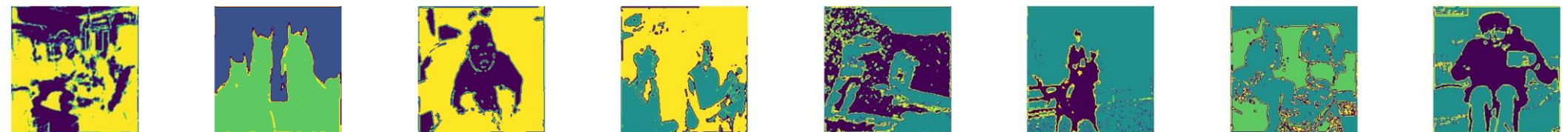
Main Image



Argmax Output



SLIC Output



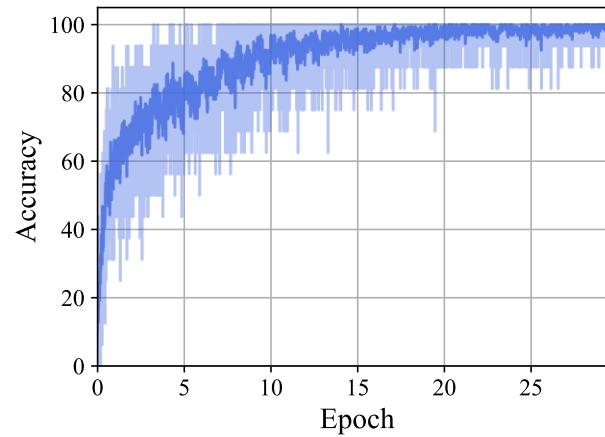
Epochs=30, LR=0.001, BS=16, Opt=SGD

# Results

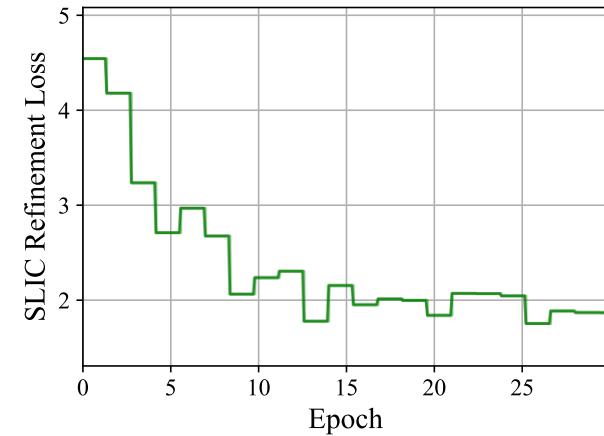
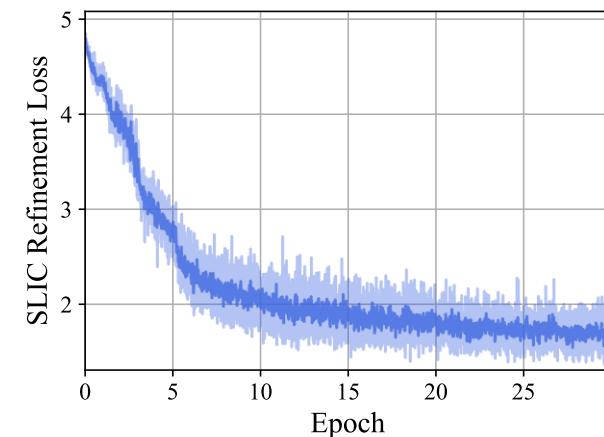
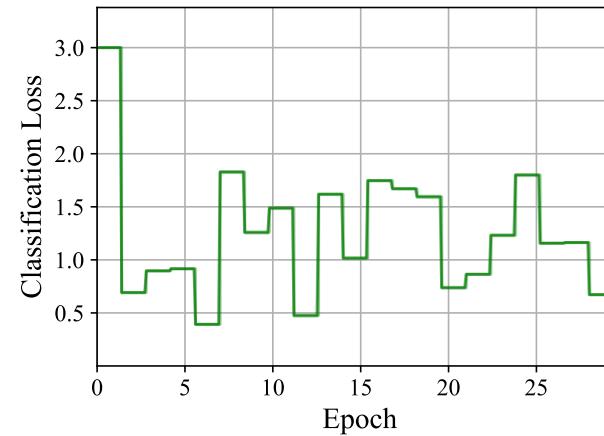
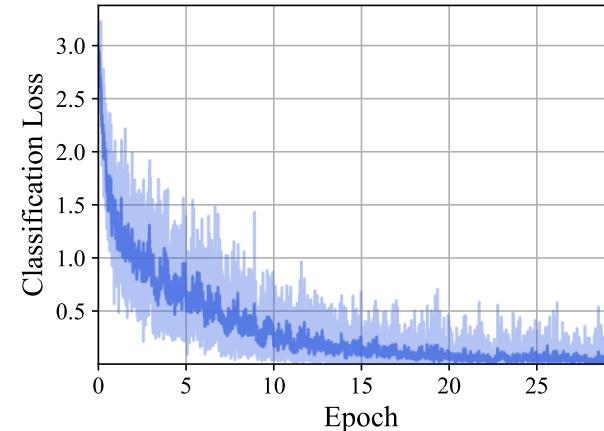
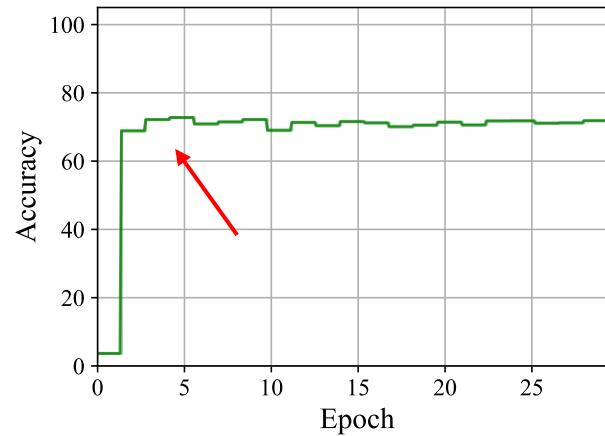
## Proposed Architecture 1: FCN with Classification Approach

Dataset: Pascal VOC

Training

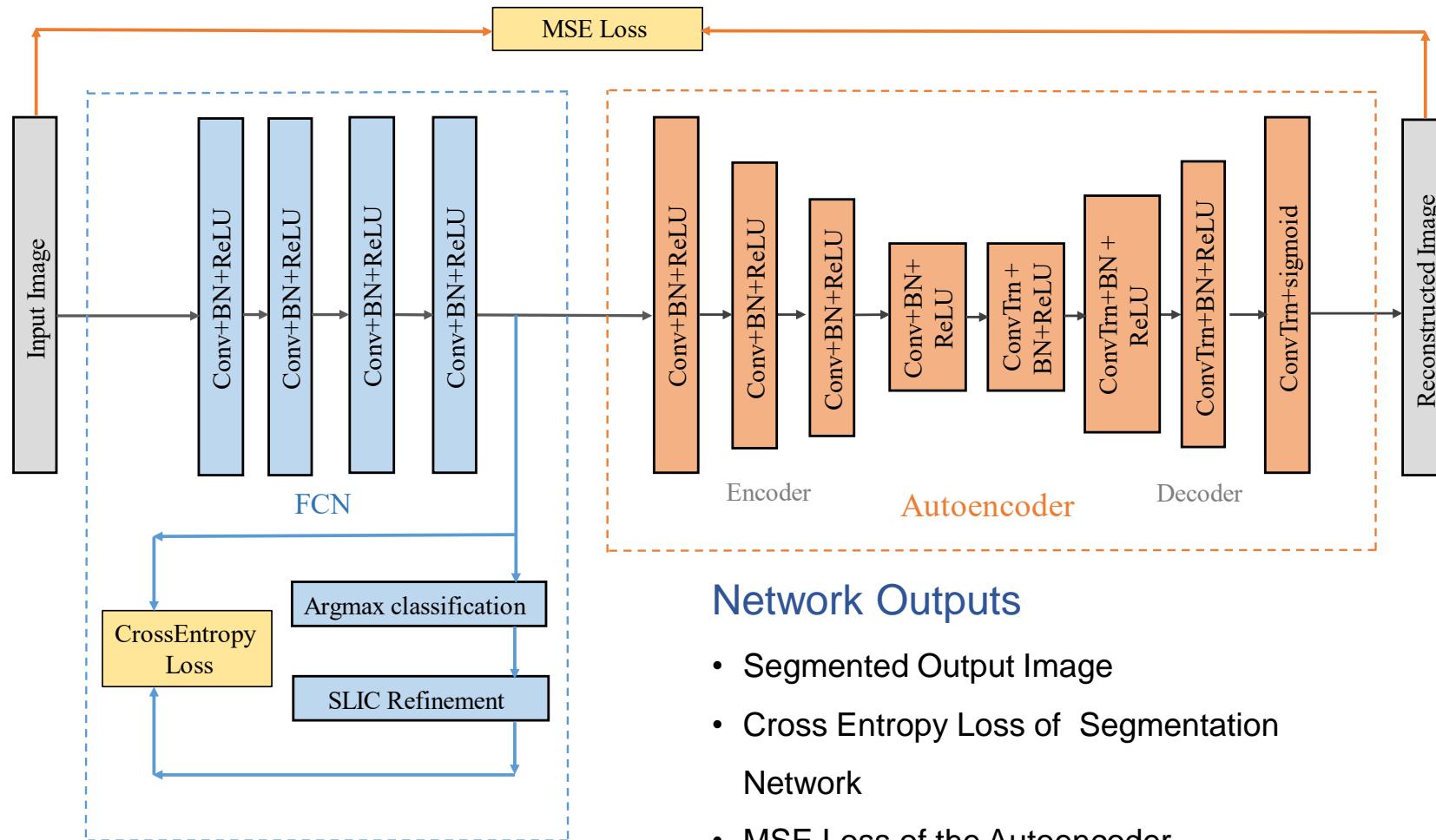


Testing



# Proposed Architecture 2 for Problem Solution

## Proposed Architecture 2: FCN with Autoencoder Approach



Hyperparameters for Network Tuning

Hyperparameters	Range / Value
Epochs	30, 40, 50, 60, 100
Normalization	minmax (0 to 1)
Activation Function	ReLU
Batch Size	5, 16, 32, 50, 100
Learning Rate	0.1, 0.01, 0.001, 0.0001
Optimizer	Adam, SGD
Initializer	Xavier Initializer
Scheduler	StepLR

## Network Outputs

- Segmented Output Image
- Cross Entropy Loss of Segmentation Network
- MSE Loss of the Autoencoder
- Reconstructed Image

# Results

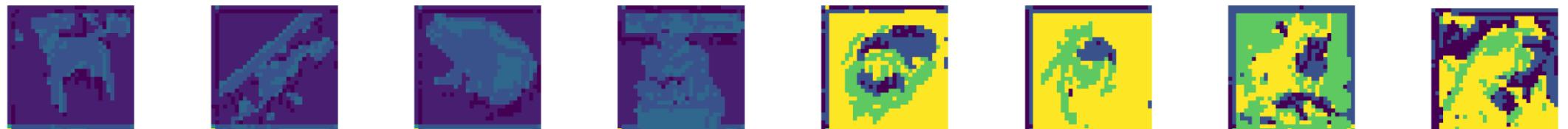
## Proposed Architecture 2: FCN with Autoencoder Approach

Dataset: CIFAR10

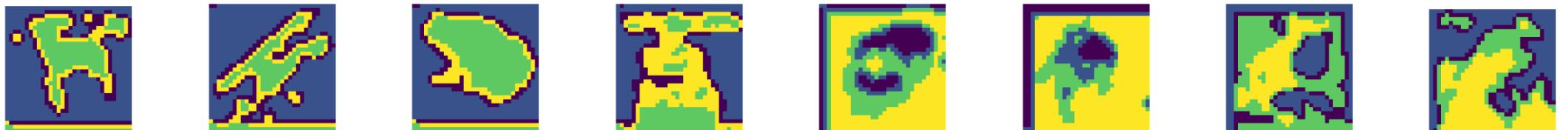
Main Image



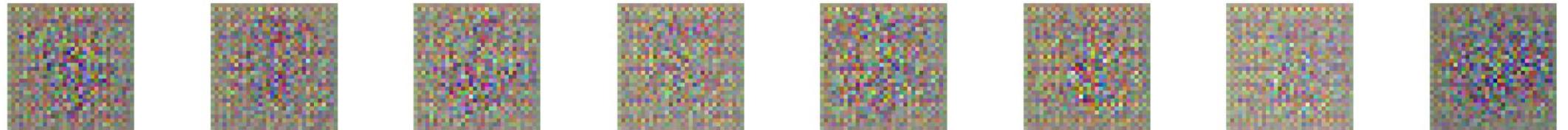
Argmax Output



SLIC Output



Reconstructed Image

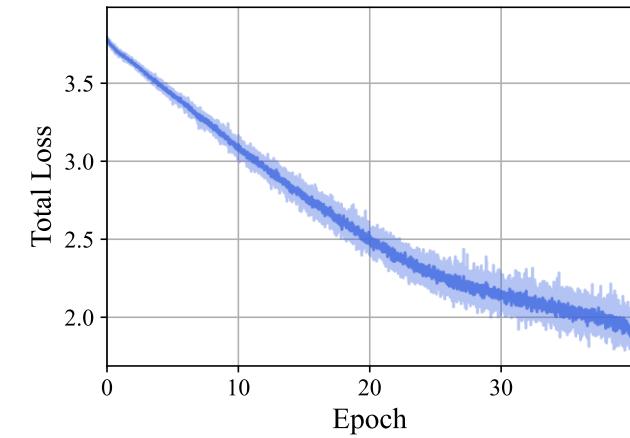
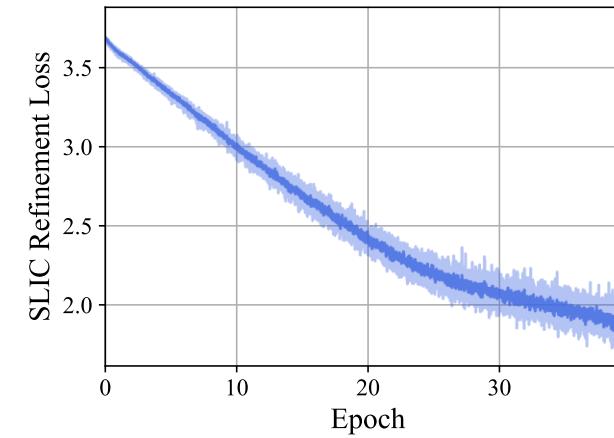
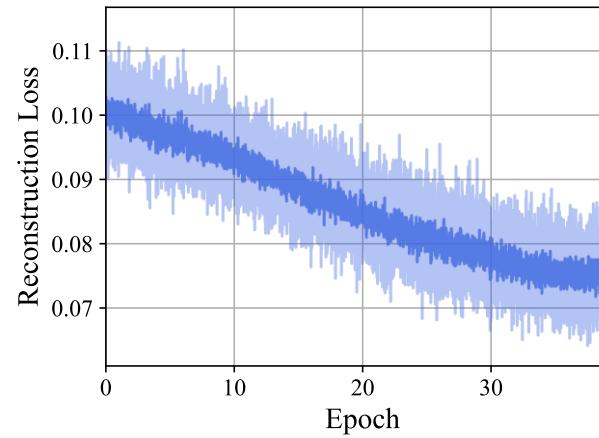


# Results

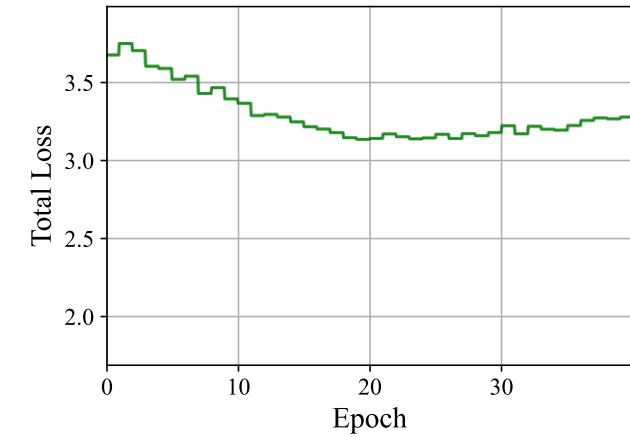
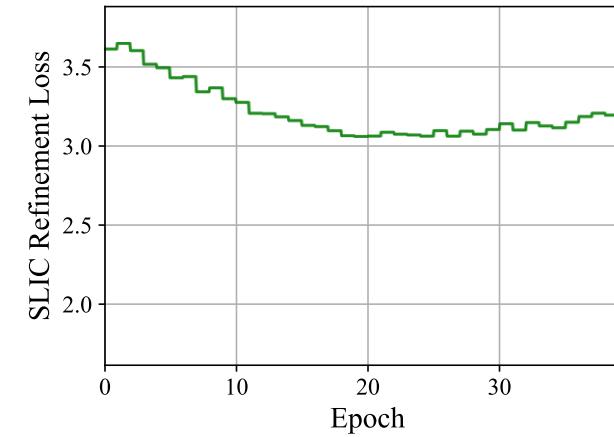
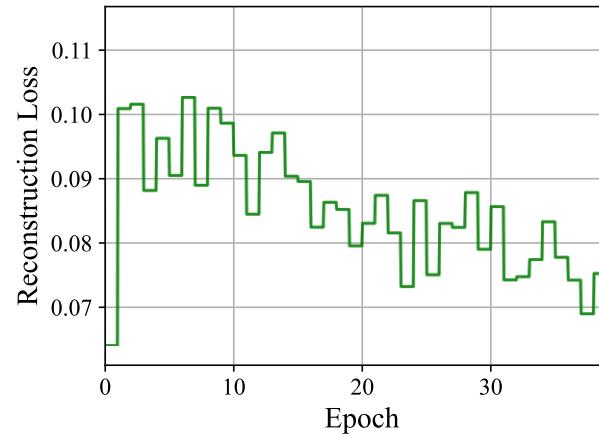
## Proposed Architecture 2: FCN with Autoencoder Approach

Dataset: CIFAR10

Training



Testing



# Results

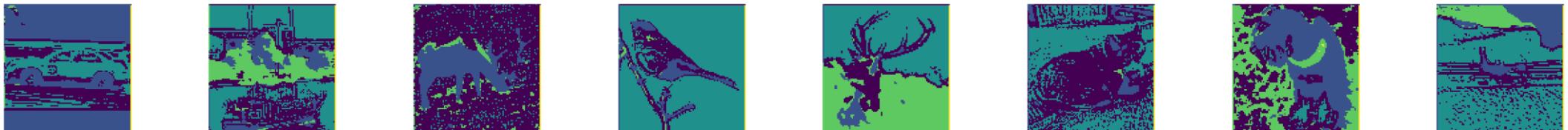
## Proposed Architecture 2: FCN with Autoencoder Approach

Dataset: STL10

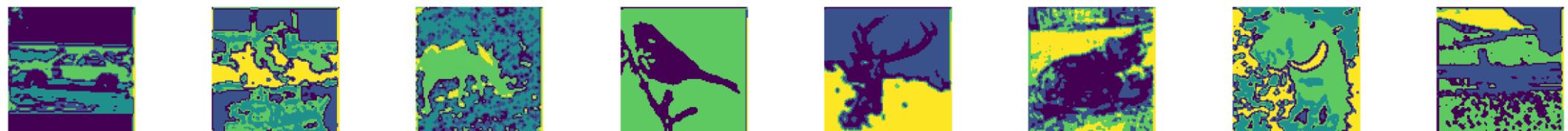
Main Image



Argmax Output



SLIC Output



Reconstructed Image

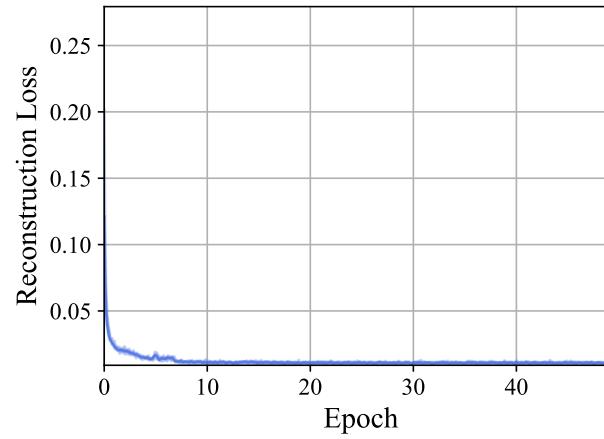


# Results

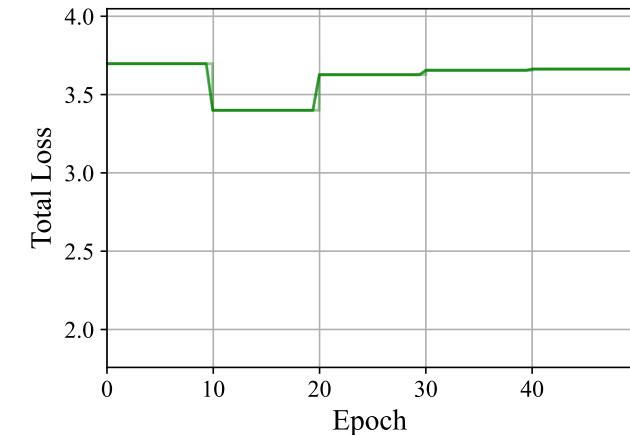
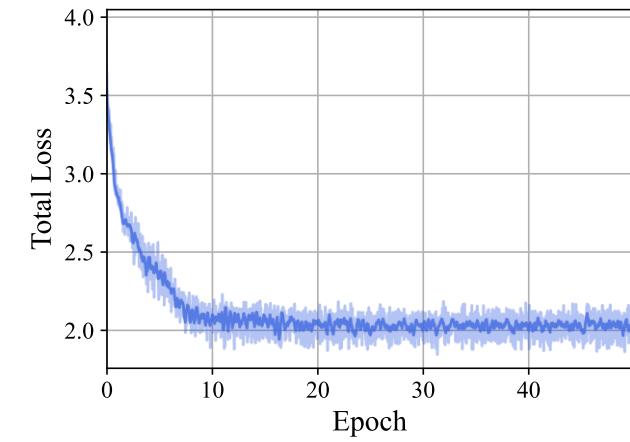
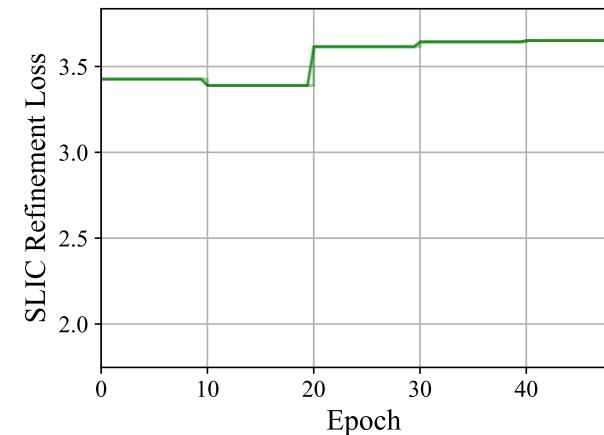
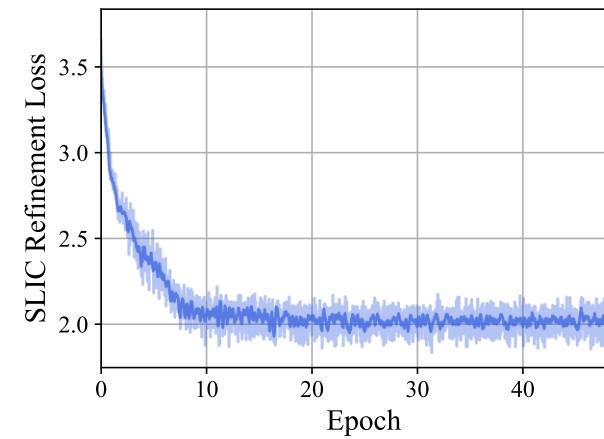
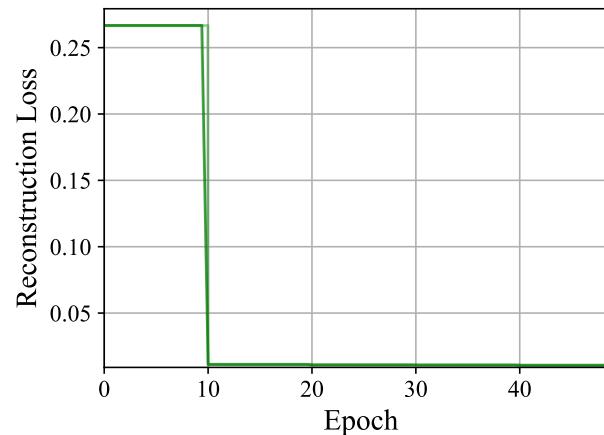
## Proposed Architecture 2: FCN with Autoencoder Approach

Dataset: STL10

Training



Testing



Epochs=50, LR=0.005, BS=100, Opt=Adam

# Results

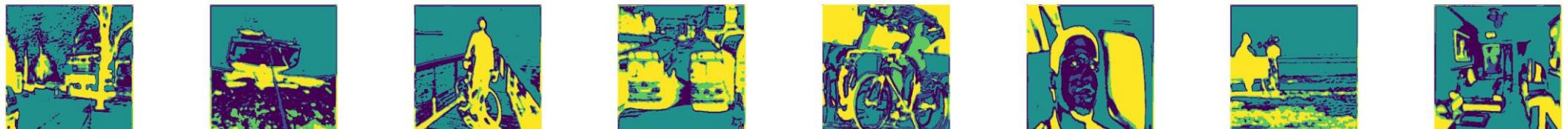
## Proposed Architecture 2: FCN with Autoencoder Approach

Dataset: Pascal VOC

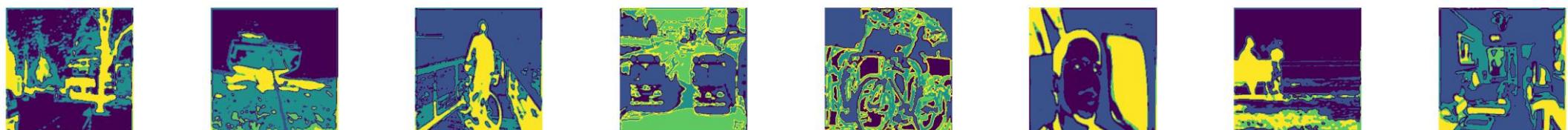
Main Image



Argmax Output



SLIC Output



Reconstructed Image



Epochs=100, LR=0.001, BS=16, Opt=Adam

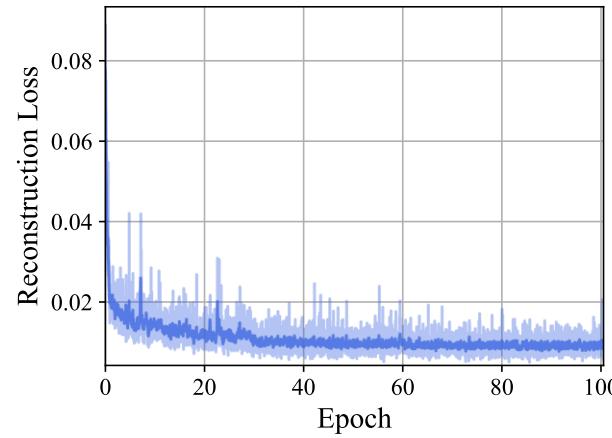


# Results

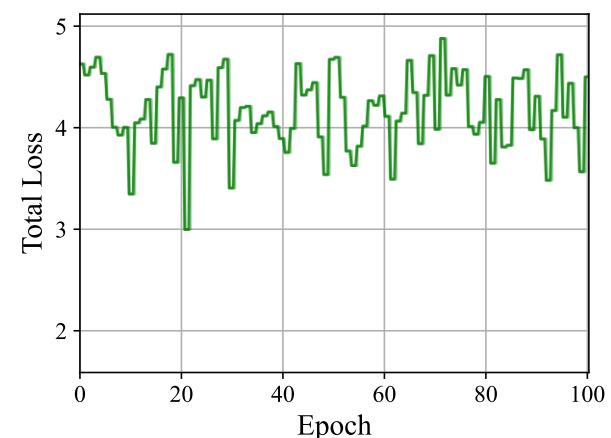
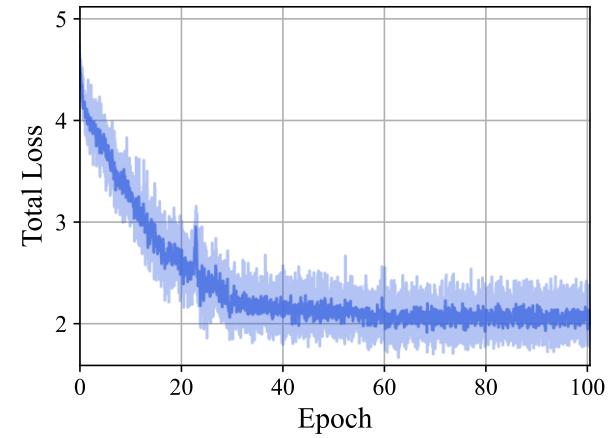
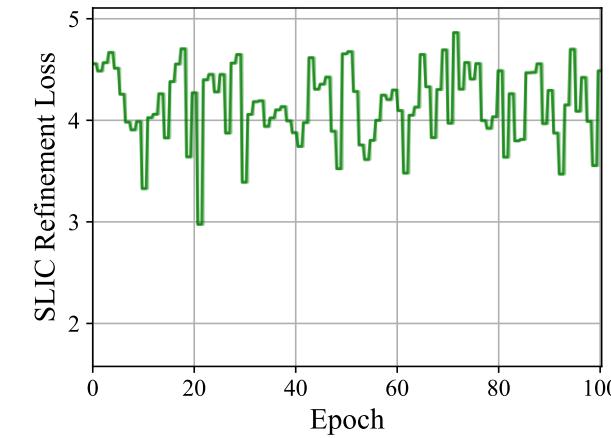
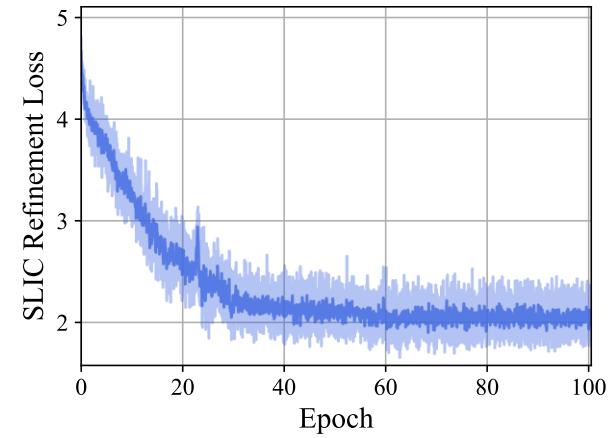
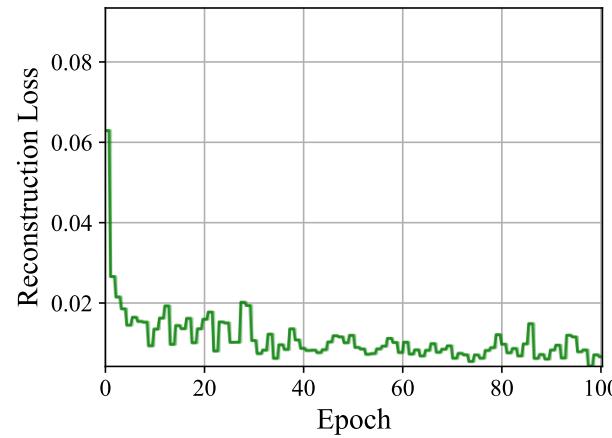
## Proposed Architecture 2: FCN with Autoencoder Approach

Dataset: Pascal VOC

Training

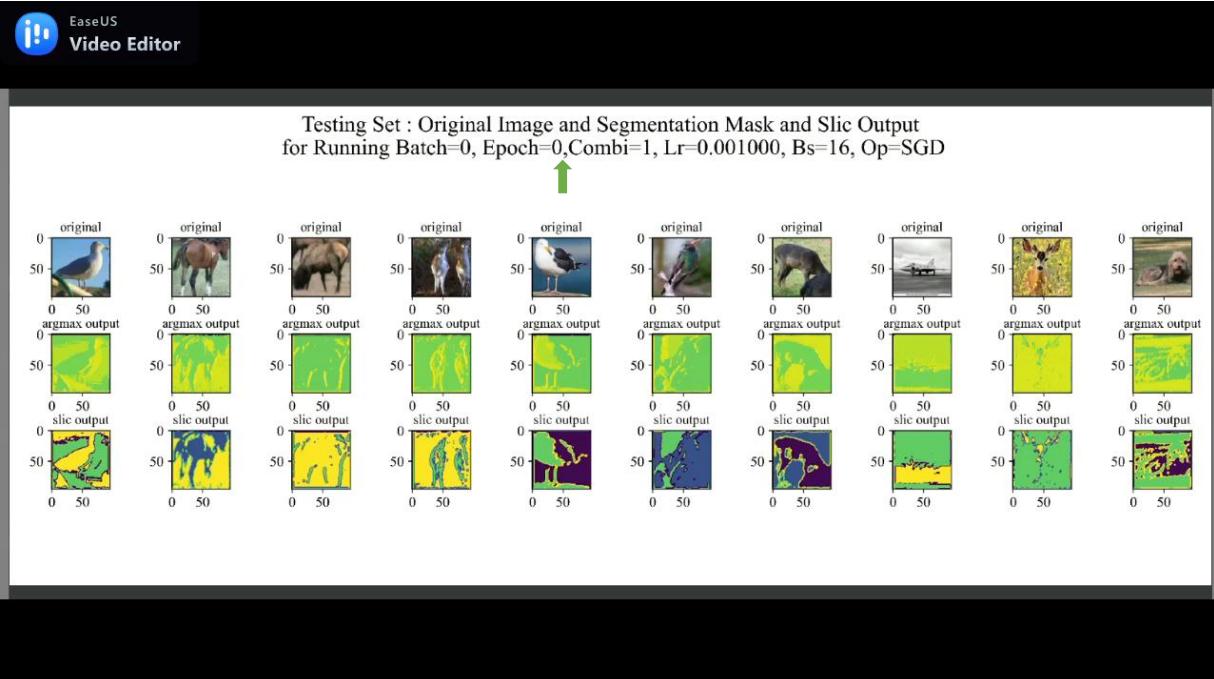


Testing



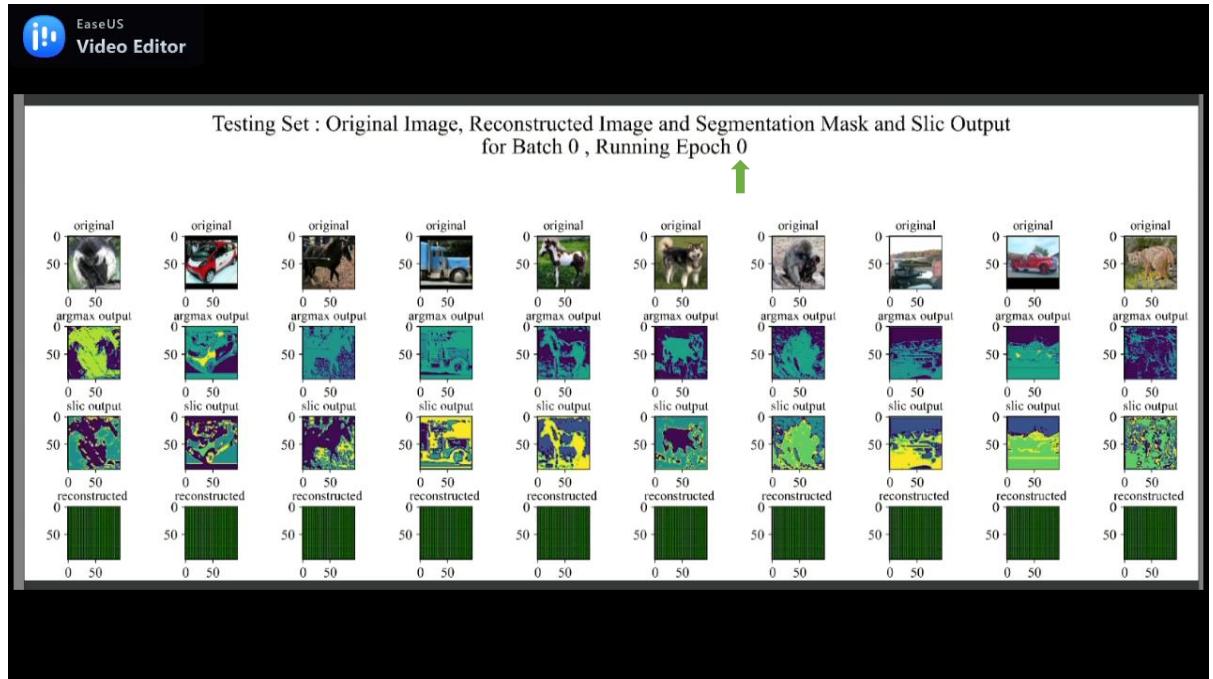
Dataset: STL10

## Effect of Using Classifier



Epochs=30, LR=0.001, BS=16, Opt=SGD

## Effect of Using Autoencoder

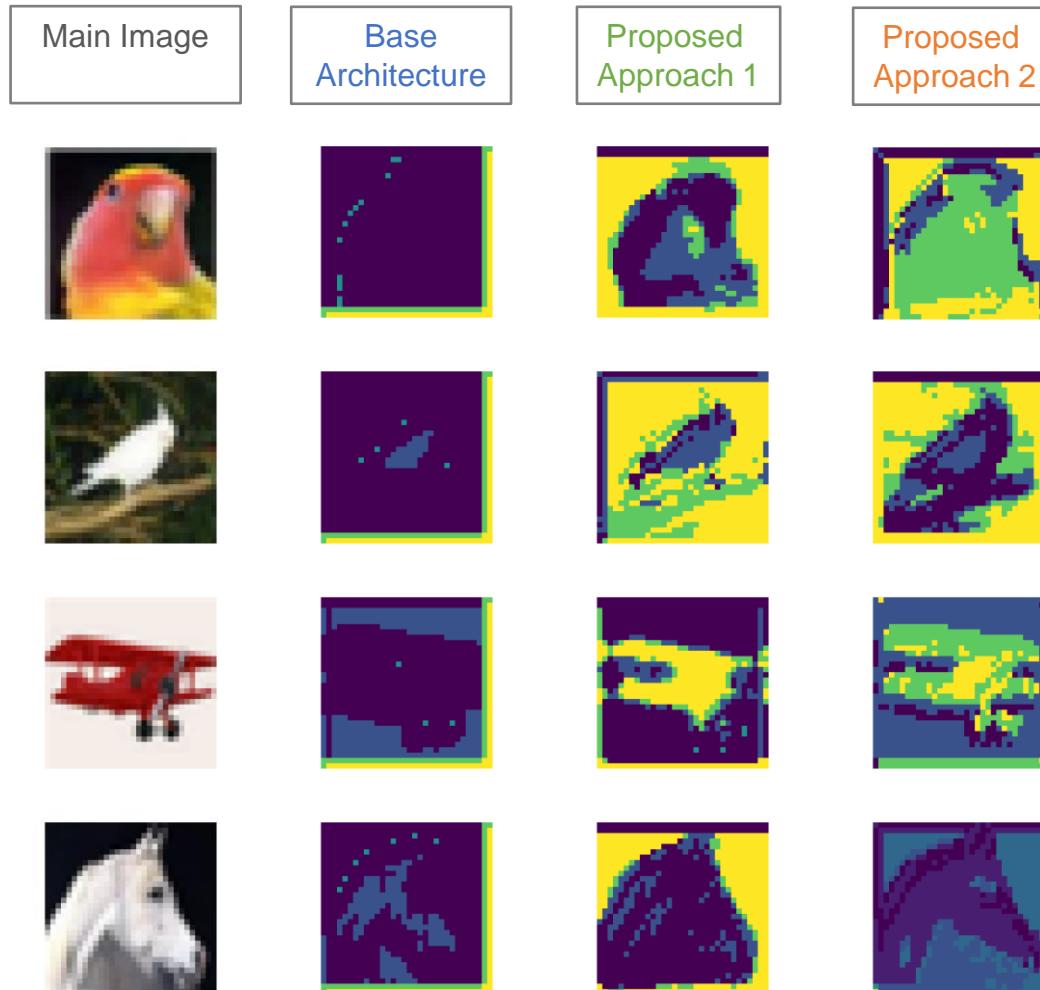


Epochs=50, LR=0.005, BS=100, Opt=Adam

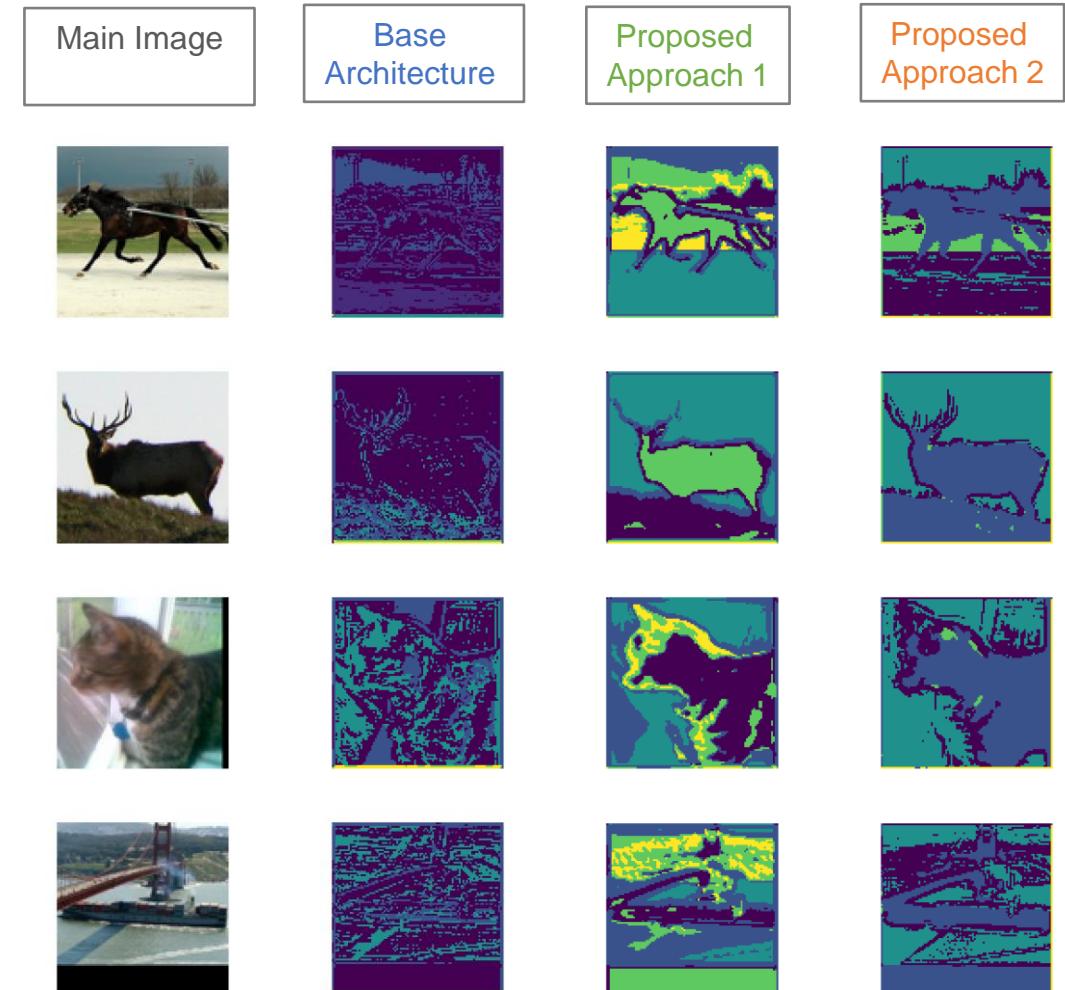
Both of the proposed approaches successfully restrict  
the network from merging away

# Comparison of Results

Dataset: CIFAR10

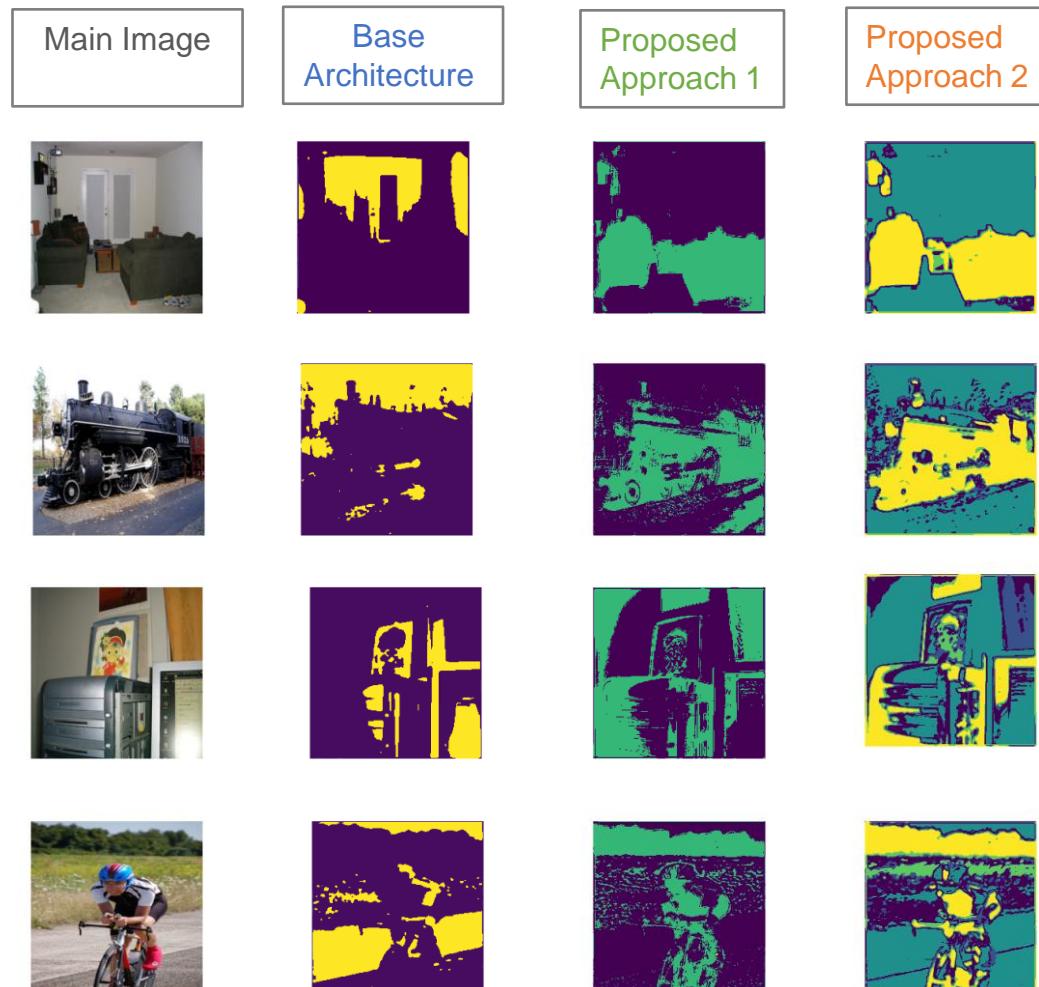


Dataset: STL10

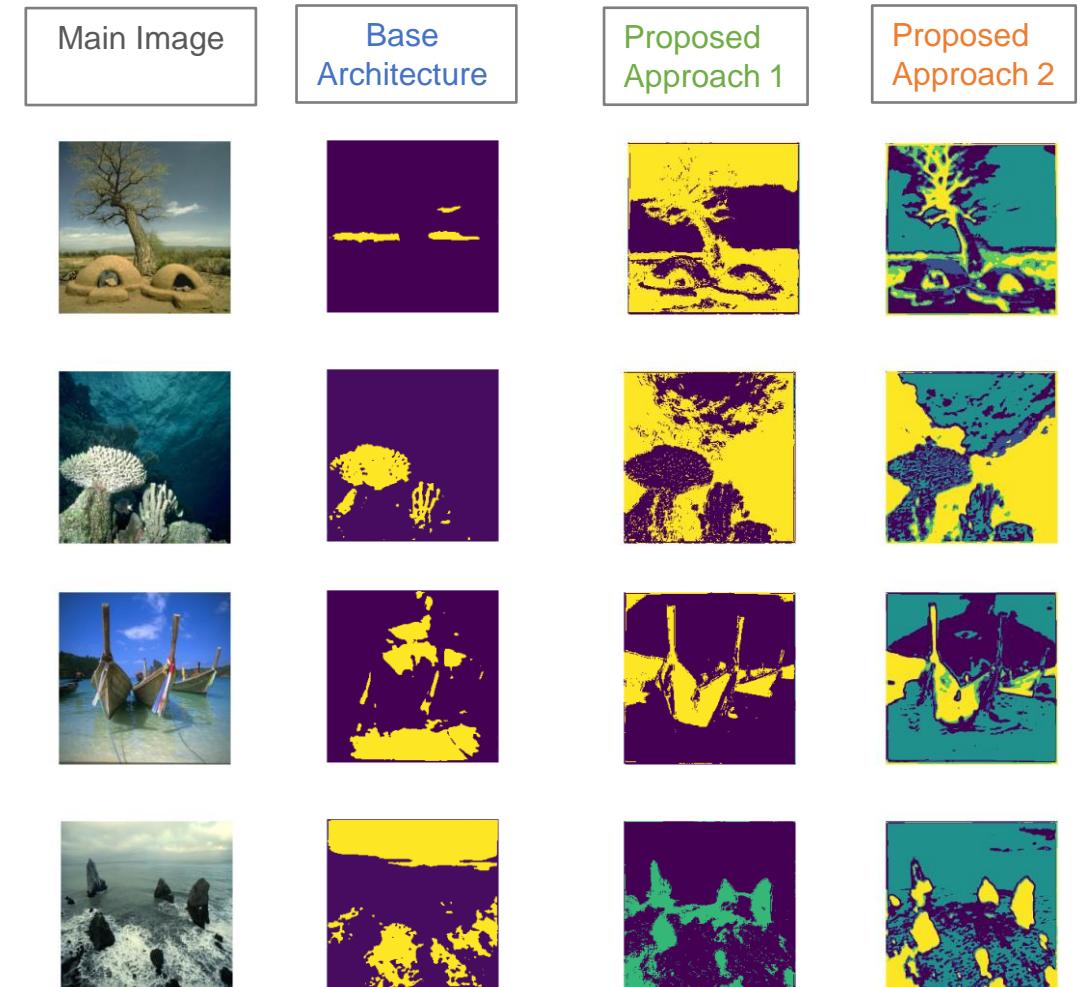


# Comparison of Results

Dataset: Pascal VOC



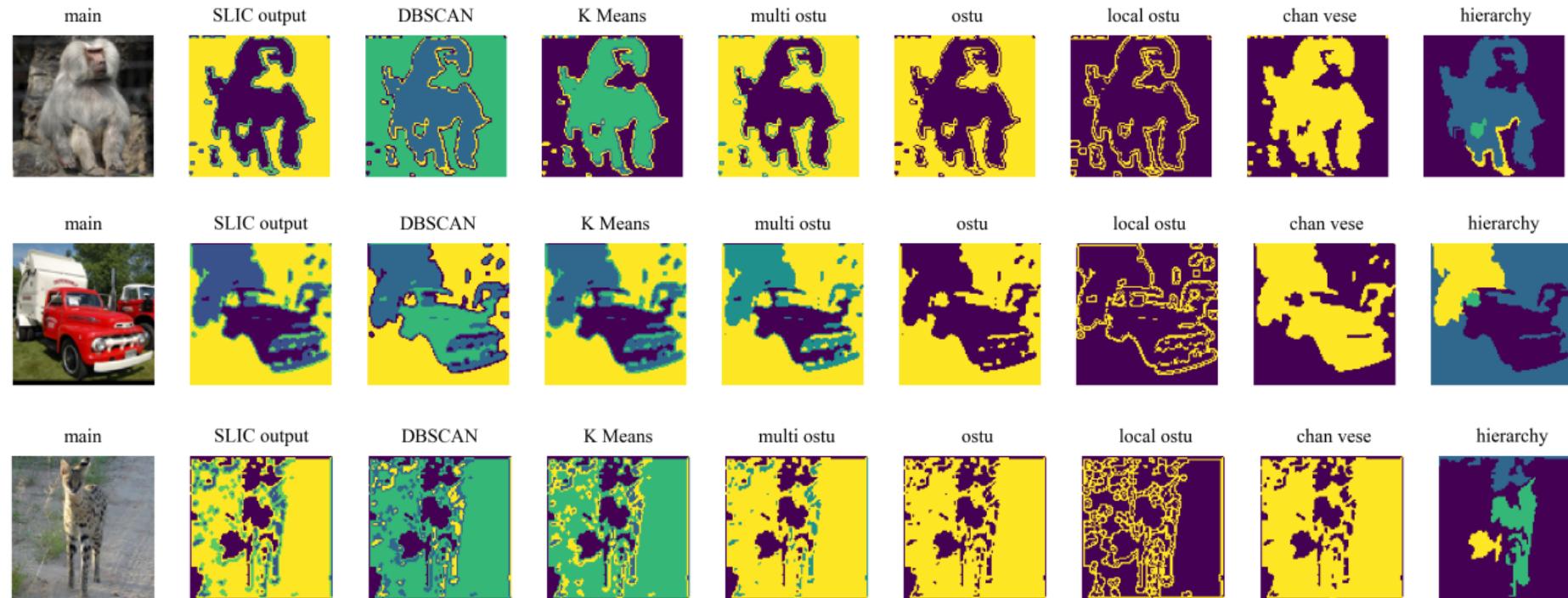
Dataset: BSDS500



# Future Work

## Post Processing of Segmented Images

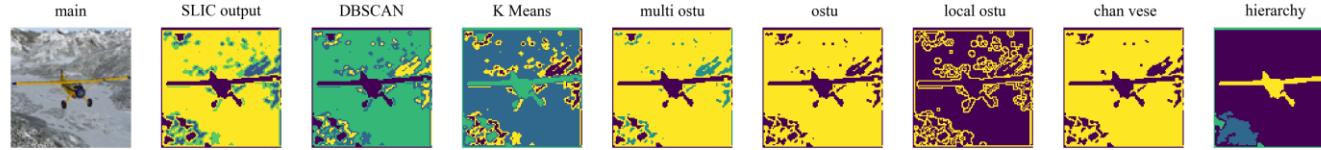
Different clustering techniques are tried to reduce noise of segmented images and to evaluate segmentation performance.



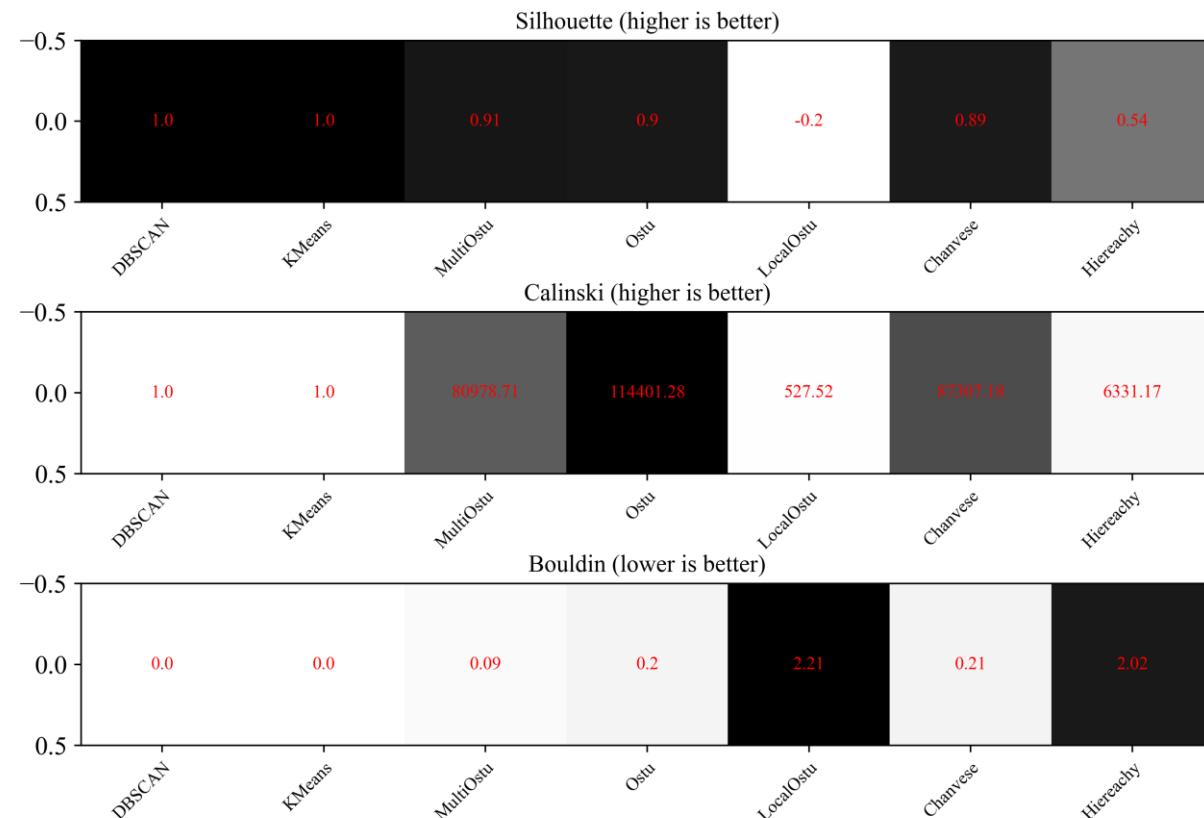
# Future Work

## Post Processing of Segmented Images

### Comparison of Different Clustering Technique



We tried different clustering techniques to reduce noise of segmented images and evaluate performance



Silhouette Score

Calinski-Harabasz Index

Davies-Bouldin Index

# Conclusion

- Very less work has been done on unsupervised segmentation deep learning techniques.
- We took the idea of combining CNN with SLIC Refinements from Asako Kanezaki .
- And develop two more stable architectures.
- We also used the new architectures on three popular image datasets and found promising results.
- The learnt models are checked on BSDS500 dataset.
- In addition, to find the best post processing methods and to design a performance metric for segmented images, 7 different cluster techniques. Further research is needed to design performance metric of segmentation in fully unsupervised way.

# Thank you!

## Questions ?

# Appendix

## Total Number of Trainable Parameters

	CIFAR10	STL10	Pascal VOC
Base Architecture	24832	24832	24832
Proposed FCN with classifier approach	177,506	2,288,468	23,954,173
Proposed FCN with autoencoder approach	2,25,883	20,5115	6,09,915

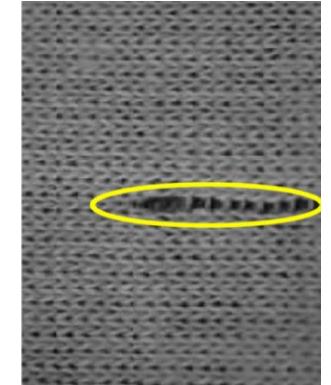
## Approximate Training Time

	CIFAR10	STL10	Pascal VOC
Base Architecture	5:39:53 h(epoch 30)	7:15:10 h(epoch 30)	8:37:34h (epoch 30)
Proposed FCN with classifier approach	3:54:18 h(epoch 30)	6 hAppx. (epoch 30)	1 day, 8:37:21h(epoch 30)
Proposed FCN with autoencoder approach	7:26:49h (epoch 100)	9h Appx. (epoch 50)	21:35:56h (epoch 100)

# Appendix

## Applications

- Road condition monitoring
- Surface defect inspection in industrial monitoring
- Fabric defect detection systems and methods
- Medical image analysis
- Anomaly detection on industrial products



Fabric defect detection systems and methods—A systematic literature review  
Author links open overlay panel Kazim Hanbay Muhammed Fatih Talu



<https://devisionx.com/industrial-applications/food-and-beverage/>

# Layer Description of Base Architecture (Appendix)

## Dataset: CIFAR10

Table 5.1: Description of convolutional neural network (CNN) architecture used in the CIFAR10 dataset

Operation Layer	Input Size (HxWxC)	Filter Num- bers	Each Filter Size	Stride Value	Padding Value	Size of Feature Map (HxWxC)
Conv1+BN+Relu	32x32x3	64	3x3x3	1x1	1x1	32x32x64
Conv2+BN+Relu	32x32x64	32	3x3x64	1x1	0x0	32x32x32
Conv3+BN+Relu	32x32x32	64	3x3x32	1x1	1x1	32x32x64
Conv4+BN+Relu	32x32x64	32	3x3x64	1x1	0x0	32x32x32

Modules	Parameters
seq1.0.weight	1728
seq1.0.bias	64
seq1.1.weight	64
seq1.1.bias	64
seq1.3.weight	2048
seq1.3.bias	32
seq1.4.weight	32
seq1.4.bias	32
seq1.6.weight	18432
seq1.6.bias	64
seq1.7.weight	64
seq1.7.bias	64
seq1.9.weight	2048
seq1.9.bias	32
seq1.10.weight	32
seq1.10.bias	32

Total Trainable Params: 24832

## Dataset: STL10

Table 5.3: Description of convolutional neural network (CNN) architecture used in STL10 dataset

Operation Layer	Input Size (HxWxC)	Filter Num- bers	Each Filter Size	Stride Value	Padding Value	Size of Feature Map (HxWxC)
Conv1+BN+Relu	96x96x3	64	3x3x3	1x1	1x1	96x96x64
Conv2+BN+Relu	96x96x64	32	3x3x64	1x1	0x0	96x96x32
Conv3+BN+Relu	96x96x32	64	3x3x32	1x1	1x1	96x96x64
Conv4+BN+Relu	96x96x64	32	3x3x64	1x1	0x0	96x96x32

Modules	Parameters
seq1.0.weight	1728
seq1.0.bias	64
seq1.1.weight	64
seq1.1.bias	64
seq1.3.weight	2048
seq1.3.bias	32
seq1.4.weight	32
seq1.4.bias	32
seq1.6.weight	18432
seq1.6.bias	64
seq1.7.weight	64
seq1.7.bias	64
seq1.9.weight	2048
seq1.9.bias	32
seq1.10.weight	32
seq1.10.bias	32

Total Trainable Params: 24832

# Layer Description of Base Architecture (Appendix)

## Dataset: Pascal VOC

Table 5.4: Description of convolutional neural network (CNN) architecture used in Pascal VOC dataset

Operation Layer	Input Size (HxWxC)	Filter Numbers	Each Filter Size	Stride Value	Padding Value	Size of Feature Map (HxWxC)
Conv1+BN+Relu	224x224x3	64	3x3x3	1x1	1x1	224x224x64
Conv2+BN+Relu	224x224x64	32	3x3x64	1x1	0x0	224x224x32
Conv3+BN+Relu	224x224x32	64	3x3x32	1x1	1x1	224x224x64
Conv4+BN+Relu	224x224x64	32	3x3x64	1x1	0x0	224x224x32

Modules	Parameters
seq1.0.weight	1728
seq1.0.bias	64
seq1.1.weight	64
seq1.1.bias	64
seq1.3.weight	2048
seq1.3.bias	32
seq1.4.weight	32
seq1.4.bias	32
seq1.6.weight	18432
seq1.6.bias	64
seq1.7.weight	64
seq1.7.bias	64
seq1.9.weight	2048
seq1.9.bias	32
seq1.10.weight	32
seq1.10.bias	32

Total Trainable Params: 24832

# Layer Description of Proposed Approach 1 (Appendix)

## Dataset: CIFAR10

Table 6.1: Segmentation network for CIFAR10 dataset

Operation Layer	Input Size (HxWxC)	Filter Numbers	Each Filter Size	Stride Value	Padding Value	Size of Feature Map (HxWxC)
Conv1+BN+Relu	32x32x3	64	3x3x3	1x1	1x1	32x32x64
Conv2+BN+Relu	32x32x64	128	3x3x64	1x1	1x1	32x32x128
Conv3+BN+Relu	32x32x128	64	3x3x128	1x1	0x0	32x32x64
Conv4+BN+Relu	32x32x64	128	3x3x64	1x1	0x0	32x32x128

Table 6.2: Classification network for CIFAR10 dataset

Operation Layer	Input Size (HxWxC)	Filter Numbers	Each Filter Size	Stride Value	Padding Value	Size of Feature Map (HxWxC)
Conv1+BN+Relu	32x32x128	12	3x3x12	1x1	1x1	32x32x12
Conv2+BN+Relu	32x32x12	12	3x3x12	1x1	1x1	32x32x12
Max Pool2d	32x32x12	1	2x2	2x2	0	16x16x12
Conv3+BN+Relu	16x16x12	24	3x3x24	1x1	1x1	16x16x24
Conv4+BN+Relu	16x16x24	24	3x3x24	1x1	1x1	16x16x24
Dropout (P=0.2)	16x16x24	-	-	-	-	16x16x24
FC Layer	6144	-	-	-	-	10

Modules	Parameters
seq.0.weight	1728
seq.0.bias	64
seq.1.weight	64
seq.1.bias	64
seq.3.weight	73728
seq.3.bias	128
seq.4.weight	128
seq.4.bias	128
seq.6.weight	8192
seq.6.bias	64
seq.7.weight	64
seq.7.bias	64
seq.9.weight	8192
seq.9.bias	128
seq.10.weight	128
seq.10.bias	128
Total Trainable Params: 92992	

Modules	Parameters
conv1.weight	13824
conv1.bias	12
batchnorm2.weight	12
batchnorm2.bias	12
conv2.weight	1296
conv2.bias	12
batchnorm3.weight	12
batchnorm3.bias	12
conv3.weight	2592
conv3.bias	24
batchnorm4.weight	24
batchnorm4.bias	24
conv4.weight	5184
conv4.bias	24
fc.weight	61440
fc.bias	10
Total Trainable Params: 84514	

# Layer Description of Proposed Approach 1 (Appendix)

## Dataset: STI10

Table 6.3: Segmentation network for STL10 dataset

Operation Layer	Input Size (HxWxC)	Filter Numbers	Each Filter Size	Stride Value	Padding Value	Size of Feature Map (HxWxC)
Conv1+BN+Relu	32x32x3	64	3x3x64	1x1	1x1	32x32x64
Conv2+BN+Relu	32x32x64	128	3x3x128	1x1	1x1	32x32x128
Conv3+BN+Relu	32x32x128	256	3x3x256	1x1	1x1	32x32x256
Conv4+BN+Relu	32x32x256	96	3x3x96	1x1	0x0	32x32x96

Table 6.4: Classification network for STL10 dataset

Operation Layer	Input Size (HxWxC)	Filter Numbers	Each Filter Size	Stride Value	Padding Value	Size of Feature Map (HxWxC)
Conv1+BN+Relu	96x96x96	32	3x3x3	1x1	1x1	96x96x32
Conv2+BN+Relu	96x96x32	64	3x3x64	1x1	1x1	96x96x64
Max Pool2d	96x96x64	1	2x2	2x2	0	48x48x64
Conv3+BN+Relu	48x48x64	128	3x3x128	1x1	1x1	48x48x128
Dropout(P= 0.05)	48x48x128	-	-	-	-	48x48x128
Max Pool2d	48x48x128	1	2x2	2x2	0	24x24x128
Conv4+BN+Relu	24x24x128	256	3x3x256	1x1	1x1	24x24x256
Conv5+BN+Relu	24x24x256	256	3x3x256	1x1	1x1	24x24x256
Max Pool2d	24x24x128	1	2x2	2x2	0	12x12x256
Dropout (P= 0.1)	12x12x256	-	-	-	-	12x12x256
FC Layer	36864	-	-	-	-	10

Modules	Parameters
seq.0.weight	1728
seq.0.bias	64
seq.1.weight	64
seq.1.bias	64
seq.3.weight	73728
seq.3.bias	128
seq.4.weight	128
seq.4.bias	128
seq.6.weight	294912
seq.6.bias	256
seq.7.weight	256
seq.7.bias	256
seq.9.weight	24576
seq.9.bias	96
seq.10.weight	96
seq.10.bias	96

Total Trainable Params: 396576

Modules	Parameters
conv_layer.0.weight	27648
conv_layer.0.bias	32
conv_layer.1.weight	32
conv_layer.1.bias	32
conv_layer.3.weight	18432
conv_layer.3.bias	64
conv_layer.4.weight	64
conv_layer.4.bias	64
conv_layer.7.weight	73728
conv_layer.7.bias	128
conv_layer.8.weight	128
conv_layer.8.bias	128
conv_layer.11.weight	147456
conv_layer.11.bias	128
conv_layer.12.weight	128
conv_layer.12.bias	128
conv_layer.15.weight	294912
conv_layer.15.bias	256
conv_layer.16.weight	256
conv_layer.16.bias	256
conv_layer.18.weight	589824
conv_layer.18.bias	256
conv_layer.19.weight	256
conv_layer.19.bias	256
fc_layer.1.weight	737280
fc_layer.1.bias	20

Total Trainable Params: 1891892

# Layer Description of Proposed Approach 1 (Appendix)

## Dataset: Pascal VOC

Table 6.5: Segmentation network for Pascal VOC dataset

Operation Layer	Input Size (HxWxC)	Filter Num- bers	Each Filter Size	Stride Value	Padding Value	Size of Fea- ture Map (HxWxC)
Conv1+BN+Relu	32x32x3	64	3x3x64	1x1	1x1	32x32x64
Conv2+BN+Relu	32x32x64	128	3x3x128	1x1	1x1	32x32x128
Conv3+BN+Relu	32x32x128	256	3x3x256	1x1	1x1	32x32x256
Conv4+BN+Relu	32x32x256	96	3x3x96	1x1	0x0	32x32x96

Table 6.6: Classification network for Pascal VOC dataset

Operation Layer	Input Size (HxWxC)	Filter Num- bers	Each Filter Size	Stride Value	Padding Value	Size of Fea- ture Map (HxWxC)
Conv1+BN+Relu	224x224x96	64	3x3x64	1x1	0x0	224x224x64
Dropout (P=0.05)	224x224x64	-	-	-	-	224x224x64
Conv2+BN+Relu	224x224x64	32	3x3x132	1x1	0x0	224x224x32
Dropout (P=0.05)	224x224x32	-	-	-	-	224x224x32
Conv3+BN+Relu	224x224x32	3	3x3x3	1x1	0x0	224x224x3
Resnet50 [74]						

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112x112					
				7x7, 64, stride 2		
					3x3 max pool, stride 2	
conv2_x	56x56	$\left[ \begin{array}{c} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$
conv3_x	28x28	$\left[ \begin{array}{c} 3 \times 3, 128 \\ 3 \times 3, 128 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 3 \times 3, 128 \\ 3 \times 3, 128 \end{array} \right] \times 4$	$\left[ \begin{array}{c} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 4$	$\left[ \begin{array}{c} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 4$	$\left[ \begin{array}{c} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 8$
conv4_x	14x14	$\left[ \begin{array}{c} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 6$	$\left[ \begin{array}{c} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 6$	$\left[ \begin{array}{c} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 23$	$\left[ \begin{array}{c} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 36$
conv5_x	7x7	$\left[ \begin{array}{c} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$
	1x1					average pool, 1000-d fc, softmax

Total Trainable Parameters: 23557597

Modules	Parameters
seq.0.weight	1728
seq.0.bias	64
seq.1.weight	64
seq.1.bias	64
seq.3.weight	73728
seq.3.bias	128
seq.4.weight	128
seq.4.bias	128
seq.6.weight	294912
seq.6.bias	256
seq.7.weight	256
seq.7.bias	256
seq.9.weight	24576
seq.9.bias	96
seq.10.weight	96
seq.10.bias	96

Total Trainable Params: 396576

# Layer Description of Proposed Approach 2 (Appendix)

## Dataset: CIFAR10

Table 7.1: Segmentation network for CIFAR10 dataset

Operation Layer	Input Size (HxWxC)	Filter Num- bers	Each Filter Size	Stride Value	Padding Value	Size of Fea- ture Map (HxWxC)
Conv1+BN+Relu	32x32x3	64	3x3x3	1x1	1x1	32x32x64
Conv2+BN+Relu	32x32x64	32	3x3x64	1x1	0x0	32x32x32
Conv3+BN+Relu	32x32x32	64	3x3x32	1x1	1x1	32x32x64
Conv4+BN+Relu	32x32x64	32	3x3x64	1x1	0x0	32x32x32

Table 7.2: Autoencoder network for CIFAR10 dataset

Operation Layer	Input Size (HxWxC)	Filter Num- bers	Each Filter Size	Stride Value	Padding Value	Size of Fea- ture Map (HxWxC)
Conv1+BN+Relu	32x32x32	12	4x4x12	2x2	1x1	16x16x12
Conv2+BN+Relu	16x16x12	24	4x4x24	2x2	1x1	8x8x24
Conv3+BN+Relu	8x8x24	48	4x4x48	2x2	1x1	4x4x48
Conv4+BN+Relu	4x4x48	96	4x4x96	2x2	1x1	2x2x96
ConvTr1+BN+Relu	2x2x96	48	4x4x48	2x2	1x1	4x4x48
ConvTr2+BN+Relu	4x4x48	24	4x4x24	2x2	1x1	8x8x24
ConvTr3+BN+Relu	8x8x24	12	4x4x12	2x2	1x1	16x16x12
ConvTr4+BN+Relu	16x16x12	3	4x4x3	2x2	1x1	32x32x3

Modules	Parameters
seq1.0.weight	1728
seq1.0.bias	64
seq1.1.weight	64
seq1.1.bias	64
seq1.3.weight	2048
seq1.3.bias	32
seq1.4.weight	32
seq1.4.bias	32
seq1.6.weight	18432
seq1.6.bias	64
seq1.7.weight	64
seq1.7.bias	64
seq1.9.weight	2048
seq1.9.bias	32
seq1.10.weight	32
seq1.10.bias	32
seq2.0.weight	6144
seq2.0.bias	12
seq2.1.weight	12
seq2.1.bias	12
seq2.3.weight	4608
seq2.3.bias	24
seq2.4.weight	24
seq2.4.bias	24
seq2.6.weight	18432
seq2.6.bias	48
seq2.7.weight	48
seq2.7.bias	48
seq2.9.weight	73728
seq2.9.bias	96

Total Trainable Params: 225883

# Layer Description of Proposed Approach 2 (Appendix)

## Dataset: STL10

Table 7.3: Segmentation network for STL10 dataset

Operation Layer	Input Size (HxWxC)	Filter Numbers	Each Filter Size	Stride Value	Padding Value	Size of Feature Map (HxWxC)
Conv1+BN+Relu	32x32x3	64	3x3x3	1x1	1x1	32x32x64
Conv2+BN+Relu	32x32x64	32	3x3x64	1x1	0x0	32x32x32

Table 7.4: Autoencoder Part for STL10 dataset

Operation Layer	Input Size (HxWxC)	Filter Numbers	Each Filter Size	Stride Value	Padding Value	Size of Feature Map (HxWxC)
Conv1+BN+Relu	96x96x32	12	4x4x12	2x2	1x1	48x48x12
Conv2+BN+Relu	48x48x12	24	4x4x24	2x2	1x1	24x34x24
Conv3+BN+Relu	24x24x24	48	4x4x48	2x2	1x1	12x12x48
Conv4+BN+Relu	12x12x48	96	4x4x96	2x2	1x1	6x6x96
ConvTr1+BN+Relu	6x6x96	48	4x4x48	2x2	1x1	12x12x48
ConvTr2+BN+Relu	12x12x48	24	4x4x24	2x2	1x1	24x24x24
ConvTr3+BN+Relu	24x24x24	12	4x4x12	2x2	1x1	48x48x12
ConvTr4+BN+Relu	48x48x12	3	4x4x3	2x2	1x1	96x96x3

Modules	Parameters		
seq1.0.weight	1728	seq2.10.weight	96
seq1.0.bias	64	seq2.10.bias	96
seq1.1.weight	64	seq2.12.weight	73728
seq1.1.bias	64	seq2.12.bias	48
seq1.3.weight	2048	seq2.13.weight	48
seq1.3.bias	32	seq2.13.bias	48
seq1.4.weight	32	seq2.15.weight	18432
seq1.4.bias	32	seq2.15.bias	24
seq2.0.weight	6144	seq2.16.weight	24
seq2.0.bias	12	seq2.16.bias	24
seq2.1.weight	12	seq2.18.weight	4608
seq2.1.bias	12	seq2.18.bias	12
seq2.3.weight	4608	seq2.19.weight	12
seq2.3.bias	24	seq2.19.bias	12
seq2.4.weight	24	seq2.21.weight	576
seq2.4.bias	24	seq2.21.bias	3
		Total Trainable Params:	205115

# Layer Description of Proposed Approach 2 (Appendix)

## Dataset: Pascal VOC

Table 7.5: Segmentation network of Pascal VOC dataset

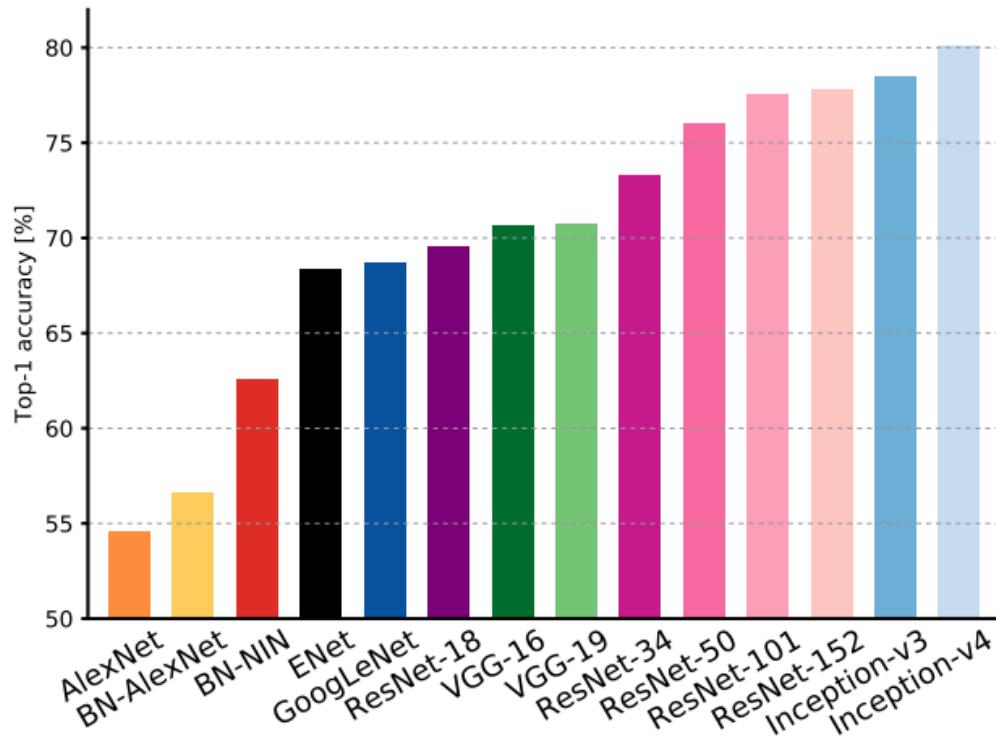
Operation Layer	Input Size (HxWxC)	Filter Numbers	Each Filter Size	Stride Value	Padding Value	Size of Feature Map (HxWxC)
Conv1+BN+Relu	32x32x3	64	3x3x64	1x1	1x1	32x32x64
Conv2+BN+Relu	32x32x64	128	3x3x128	1x1	1x1	32x32x128
Conv3+BN+Relu	32x32x128	256	3x3x256	1x1	1x1	32x32x256
Conv4+BN+Relu	32x32x256	96	3x3x96	1x1	0x0	32x32x96

Table 7.6: Autoencoder part of Pascal VOC dataset

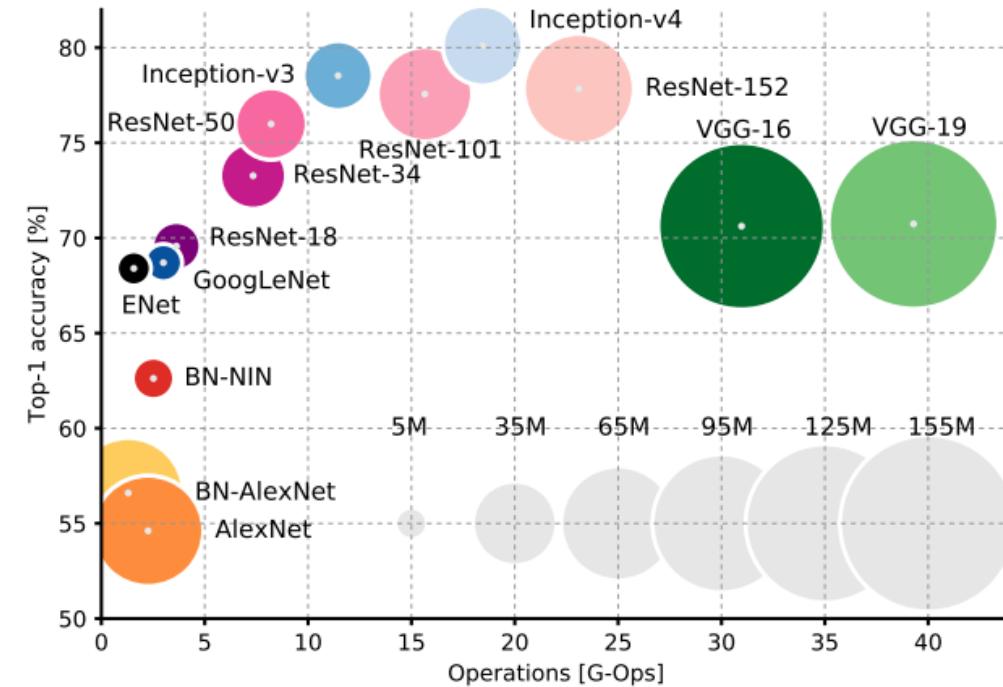
Operation Layer	Input Size (HxWxC)	Filter Numbers	Each Filter Size	Stride Value	Padding Value	Size of Feature Map (HxWxC)
Conv1+BN+Relu	224x224x32	12	4x4x12	2x2	1x1	112x112x12
Conv2+BN+Relu	112x112x12	24	4x4x24	2x2	1x1	56x56x24
Conv3+BN+Relu	56x56x24	48	4x4x48	2x2	1x1	28x28x48
Conv4+BN+Relu	28x28x48	96	4x4x96	2x2	1x1	14x14x96
ConvTr1+BN+Relu	14x14x96	48	4x4x48	2x2	1x1	28x28x48
ConvTr2+BN+Relu	28x28x48	24	4x4x24	2x2	1x1	56x56x24
ConvTr3+BN+Relu	56x56x24	12	4x4x12	2x2	1x1	112x112x12
ConvTr4+BN+Relu	112x112x12	3	4x4x3	2x2	1x1	224x224x3

Modules	Parameters	
seq1.0.weight	1728	
seq1.0.bias	64	
seq1.1.weight	64	
seq1.1.bias	64	
seq1.3.weight	73728	
seq1.3.bias	128	
seq1.4.weight	128	
seq1.4.bias	128	
seq1.6.weight	294912	
seq1.6.bias	256	
seq1.7.weight	256	
seq1.7.bias	256	
seq1.9.weight	24576	
seq1.9.bias	96	
seq1.10.weight	96	
seq1.10.bias	96	
seq2.0.weight	18432	
seq2.0.bias	12	
seq2.1.weight	12	
seq2.1.bias	12	
seq2.3.weight	4608	
seq2.3.bias	24	
seq2.4.weight	24	
seq2.4.bias	24	
seq2.6.weight	18432	
seq2.6.bias	48	
seq2.7.weight	48	
seq2.7.bias	48	
seq2.9.weight	73728	
seq2.9.bias	96	
Total Trainable Params: 609915		

# Reason for choosing ResNet50 (Appendix)



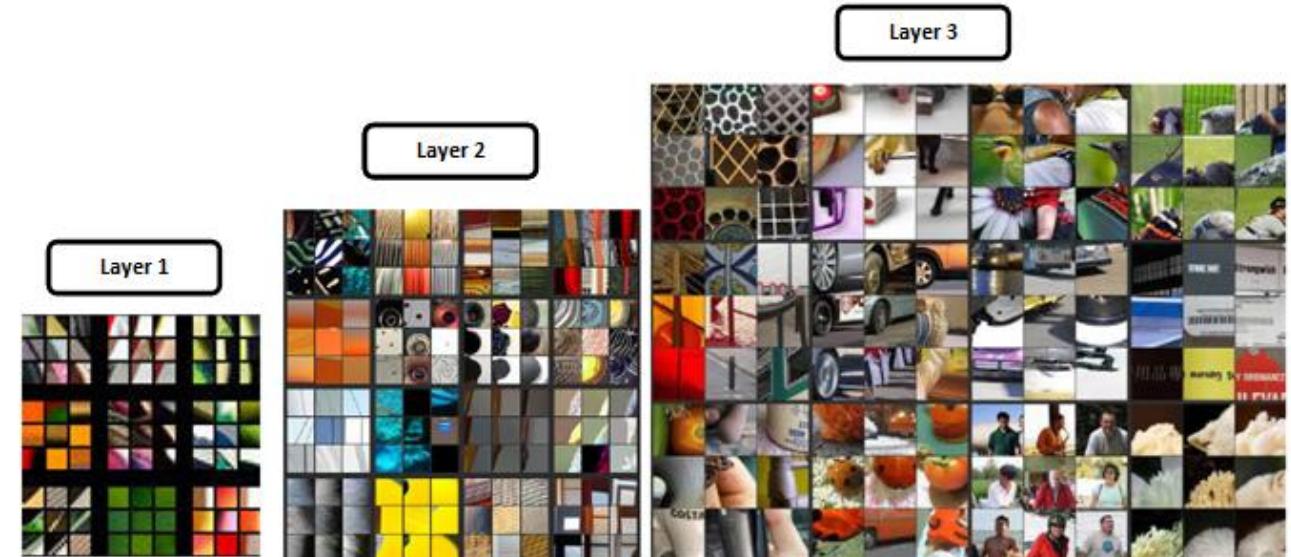
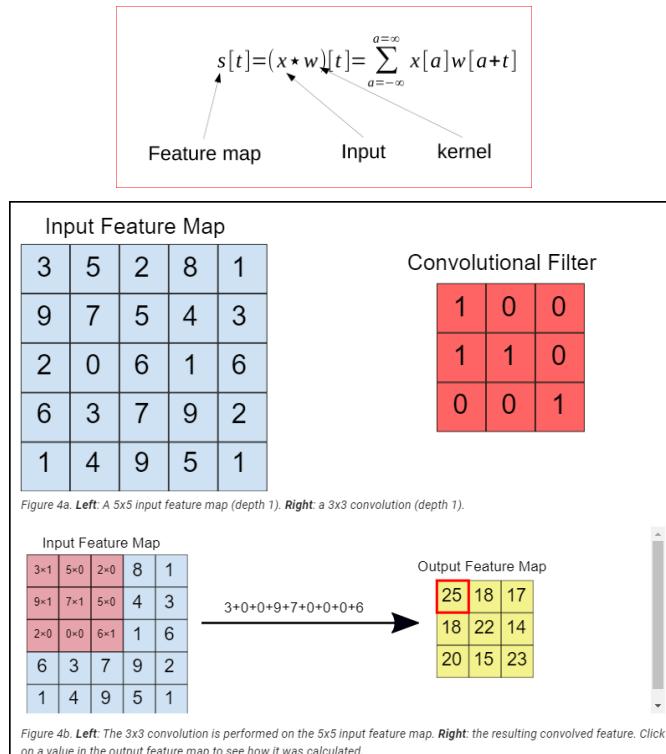
Source: AN ANALYSIS OF DEEP NEURAL NETWORK MODELS FOR PRACTICAL APPLICATIONS  
Alfredo Canziani & Eugenio Culurciello



# Appendix

## Architecture for Unsupervised Approach

### Convolutional Operation



Source:Matthew D. Zeiler and Rob Fergus), {Visualizing and Understanding Convolutional Networks},2013

- Filter extracts local features
- In practice, a number of filters are used to extract different types of local features