# Label Refinement Network for Coarse-to-Fine Semantic Segmentation

Md Amirul Islam[1], Shujon Naha[2], Mrigank Rochan[1], Neil Bruce[1], and Yang Wang[1]

[1]University of Manitoba, Winnipeg, MB
[2]Indiana University, Bloomington, IN

## Abstract

*We consider the problem of semantic image segmentation using deep convolutional neural networks. We propose a novel network architecture called the* label refinement network *that predicts segmentation labels in a coarse-to-fine fashion at several resolutions. The segmentation labels at a coarse resolution are used together with convolutional features to obtain finer resolution segmentation labels. We define loss functions at several stages in the network to provide supervisions at different stages. Our experimental results on several standard datasets demonstrate that the proposed model provides an effective way of producing pixel-wise dense image labeling.*

## 1. Introduction

We consider the problem of semantic image segmentation, where the goal is to densely label each pixel in an image according to the object class that it belongs to. We propose a convolutional neural network architecture called the *label refinement network* (LRN) that performs semantic segmentation in a coarse-to-fine fashion.

Deep convolutional neural networks (CNNs) have been successfully applied to a wide variety of visual recognition problems, such as image classification [13], object detection [19], action recognition [24], etc. CNNs extract deep feature hierarchies using alternating layers of operations, such as convolution, pooling, etc. Features from the top layers of CNNs tend to be invariant to "nuisance factors" such as pose, illumination, small translations, etc. The invariance properties of these features make them particularly useful for vision tasks such as whole image classification. However, these features are not well suited for other vision tasks (e.g. semantic segmentation) that require precise pixel-wise information.

There have been some recent efforts [1, 4, 8, 9, 16, 17, 18] on adapting CNNs for semantic segmentation. Some of these approaches (e.g. [9, 16, 17]) are based on combining the con-volutional features from multiple layers in a CNN. Another popular approach (e.g. [1, 18] is to use upsampling (also known as deconvolution) to enlarge the spatial dimensions of the feature map at the top layer of a CNN, e.g. to the same spatial dimensions as the original image, then predict the pixel-wise labels from the enlarged feature map. In both the cases, the final full-sized semantic segmentation result is obtained in a "single shot" at the very end of the network architecture.

In this paper, we introduce a novel CNN-based architecture for semantic segmentation. Different from previous approaches, our proposed model predicts semantic labels at several different resolutions in a coarse-to-fine fashion. We use resized ground-truth segmentation labels as the supervision at each resolution level. The segmentation labels at a coarse scale are combined with convolutional features to produce segmentation labels at a finer scale. See Fig. 1 for an illustration of our model architecture.

We make threefold contributions in this paper which are as follows:

- We introduce a new perspective on the semantic segmentation (or more generally, pixel-wise labeling) problem. Instead of predicting the final segmentation result in a single shot, we propose to solve the problem in a coarse-to-fine fashion by first predicting a coarse labeling, then progressively refine the result to get the finer scale results.

- We propose an end-to-end CNN architecture to learn to predict the segmentation labels at multiple resolutions. Unlike most of the previous methods that only have supervision at the end of their network, our model has supervision at multiple resolutions in the network. Although we focus on semantic image segmentation in this paper, our network architecture is general enough to be used for any pixel-wise labeling task.

- We perform extensive experiments on several standard datasets to demonstrate the effectiveness of our proposed model.
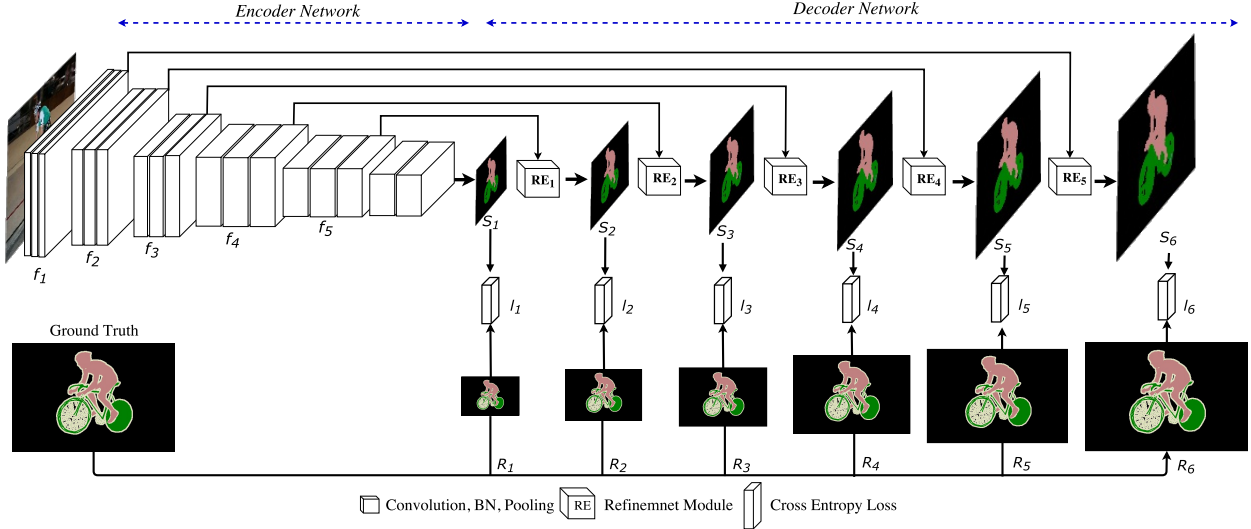
Figure 1. Overall architecture of the Label Refinement Network (LRN). LRN is based on encoder-decoder framework. The encoder network produces a sequence of feature maps with decreasing spatial dimensions. The decoder network produces label maps with increasing spatial dimensions. A larger label map is obtained by combining the previous (smaller) label map and the corresponding convolutional features from a layer in the encoder network indicated by the solid line. We use down-sampled ground-truth label maps to provide the supervision at each stage of the decoder network.

The rest of the paper is organized as follows. Section 2 presents related work. Section 3 discusses the essential background of encoder-decoder architecture. Section 4 introduces our label refinement network. Section 5 describes the experiment details and also presents some analysis. Finally, we conclude the paper in Section 6.

## 2. Related Work

CNNs have shown tremendous success in various visual recognition tasks, e.g. image classification [13], object detection [19], action recognition [24]. Recently, there has been work [1, 4, 8, 9, 16, 17, 18] on adapting CNNs for pixel-wise image labeling problems such as semantic segmentation.

Directly related to our work is the idea of multi-scale processing in computer vision. Early work on Laplacian pyramids [3] is an example of capturing image structures at multiple scales. Eigen et al. [6] use a multi-scale CNN for predicting depths, surface normals and semantic labels from a single image. Denton et al. [5] use coarse-to-fine idea in the context of image generation. Honari et al. [10] combine coarse-to-fine features in CNNs for facial keypoint localization.

Similar to our proposed model, some papers have also used the idea of defining loss functions at multiple stages in a CNN to provide more supervisions in learning. The Inception model [22] use auxiliary classifiers at the lower stages of the network to encourage the network to learn discriminative features. Lee et al. [15] propose a deeply-supervised network for image classification. Xie et al. [25]

use a similar idea for edge detection.

## 3. Background: Encoder-Decoder Network for Semantic Segmentation

In recent years, convolutional neural networks have been driving a lot of success in computer vision. Following the early success of CNNs on image classification [13], there has been a line of work [1, 16, 18] on adapting CNNs for pixel-wise dense image labeling problems such as semantic image segmentation. Most of the CNN-based semantic segmentation algorithms can be characterized in terms of a generic encoder-decoder architecture [1, 18]. The encoder network is usually a convolutional network that extracts features from an input image. Similar to the CNNs used for image classification, the encoder network typically consists of alternating layers of convolution, subsampling (e.g. via max-pooling), non-linear rectification, etc.

The subsampling layers in CNNs allow the networks to extract high-level features that are translation invariant, which are crucial for image classification. However, they reduce the spatial extent of the feature map at each layer in the network. Let $I \in \mathbb{R}^{h \times w \times d}$ be the input image (where $h$, $w$ are spatial dimensions $d$ is the color channel dimension) and $f(I) \in \mathbb{R}^{h' \times w' \times d'}$ be the feature map volume at the end of the encoder network. The feature map $f(I)$ has smaller spatial dimensions than the original image, i.e. $h' < h$ and $w' < w$. In order to produce a full-sized segmentation result as the output, an associated decoder network is applied
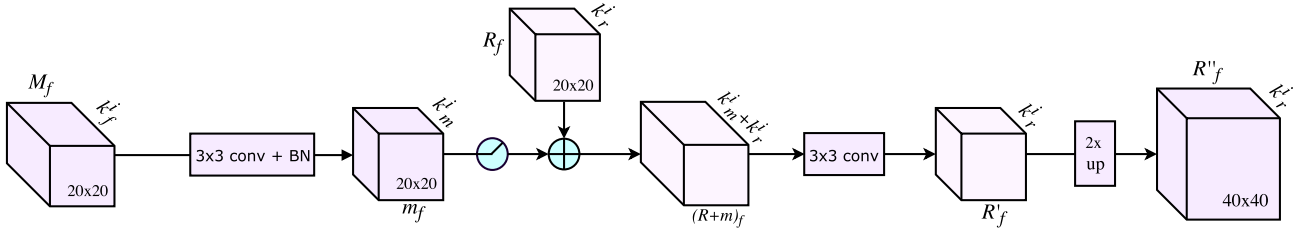
Figure 2. Detailed illustration of the refinement module (RE). The refinement module (see RE in Fig. 1) is used after each segmentation map in the decoder network. Here we unfold and illustrate RE for $i^{th}$ ($i = 1, 2, .., 5$) stage of the decoder. RE takes two inputs. The first input is the coarse label map $R_f$ with channel dimension $k_r^i$ generated at a decoding stage, and the second input is the skip feature map $M_f$ with channel dimension $k_f^i$ from the corresponding convolution layer (we use the last convolution layer of each stage ($f_c$ in Fig. 1) where $c = 1, 2, .., 5$)) of the encoder. $R_f$ and $M_f$ have the same spatial dimension but different channel dimensions (typically $k_m^i \gg k_r^i$). In order to combine $R_f$ and $M_f$, we apply a 3x3 convolution on $M_f$ followed by a batch normalization and ReLU operation. These operations results in a skip feature map $m_f$ (with channel dimension $k_m^i$) that matches in terms of the spatial dimension and channel dimension of the coarse label map $R_f$. We then apply a $3 \times 3$ convolution to the combined feature map $(R+m)_f$ to produce a feature map $R_f'$ whose channel dimension is same as of coarse label feature map $R_f$. $R_f'$ is the $i^{th}$ stage prediction map. Then we upsample the prediction map $R_f'$ by a factor of 2 ($R_f''$ in the figure which is the segmentation map $s_k(I)$ in Fig. 1) and feed it to the next stage $(i + 1)^{th}$ of the decoder. Note that convolution and upsampling layers are shown in rectangles, whereas feature maps and concatenation operations are represented by 3D boxes and circles respectively. In summary, the refinement module is composed of convolution, batch normalization, ReLU, concatenation, and bilinear upsampling operations.

to the output of the encoder network to produce the result that has spatial size of the original image. The fundamental differences between our work and various research contributions in prior work mainly lie in choices of the structure and composition of the decoder network. The decoder in SegNet [1] progressively enlarges the feature map using an upsampling technique without learnable parameters. The decoder in DeconvNet [18] is similar, but has learnable parameters. FCN [16] uses a single layer interpolation for deconvolution.

At the end of the decoder network, the full-sized feature map is mapped to predict the class label for each pixel, and a loss function is defined to compare the predicted labels with the ground-truth. It is important to note that the network makes the prediction only at the end of the decoder network. As a result, the supervision only happens at the last layer of the entire architecture, with which the loss function is associated.

## 4. Label Refinement Network

In this section, we propose a novel network architecture called the *Label Refinement Network (LRN)* for dense image labeling problems. An overview of the LRN architecture is shown in Fig. 1. Similar to prior work [1, 16, 18], LRN also uses an encoder-decoder framework. The encoder network of LRN is similar to that of SegNet [1] which is based on the VGG16 network [20]. The novelty of LRN lies in the decoder network. Instead of making the prediction at the end of the entire architecture, the decoder network in LRN makes predictions in a coarse-to-fine fashion at several stages. In addition, LRN introduces supervisions early in the network

by adding loss functions at multiple stages (not just at the last layer) in the decoder network.

The LRN architecture is motivated by the following observations: Due to the convolution and subsampling operations, the feature map $f(I)$ obtained at the end of the encoder network mainly contains high-level information about the image (e.g. objects). Spatially precise information is lost in the encoder network, and therefore $f(I)$ cannot be used directly to recover a full-sized semantic segmentation which requires pixel-precise information. However, $f(I)$ contains enough information to produce a *coarse* semantic segmentation. In particular, we can use $f(I)$ to produce a segmentation map $s(I)$ of spatial dimensions $h' \times w'$, which is smaller than the original image dimensions $h \times w$. Our decoder network then progressively refines the segmentation map $s(I)$. Let $s_k(I)$ be the $k$-th segmentation map (left to right in Fig. 1). If $k > k'$, the segmentation map $s_k(I)$ will have a larger spatial dimensions than $s_{k'}(I)$. Note that one can interpret $s_k(I)$ as the feature map in the decoder network in most work (e.g. [1, 16, 18]). However, our model enforces the channel dimension of $s_k(I)$ to be the same as the number of class labels, so $s_k(I)$ can be considered as a (soft) label map.

The network architecture in Fig. 1 has 5 convolution layers in the encoder. We use $f_k(I) \in \mathbb{R}^{h_k \times w_k \times d_k}$ ($k = 1, 2, ..., 5$) to denote the feature map after the $k$-th convolution layer. The decoder network generates 6 label maps $s_k$ ($k = 1, 2, ..., 6$). After the last convolution layer of the encoder network, we use a $3 \times 3$ convolution layer to convert the convolution feature map $f_5(I) \in \mathbb{R}^{h_5 \times w_5 \times d_5}$ to $s_1(I) \in \mathbb{R}^{h_5 \times w_5 \times C}$, where $C$ is the number of class

labels. We then define a loss function on $s_1(I)$ as follows. Let $Y \in \mathbb{R}^{h \times w \times C}$ be the ground-truth segmentation map, where the label on each pixel is represented as a $C$-dimensional vector using the one-shot representation. We use $R_1(Y) \in \mathbb{R}^{h_5 \times w_5 \times C_5}$ to denote the segmentation map obtained by resizing $Y$ to have the same spatial dimensions as $s_1(I)$. We can then define a loss function $\ell_1$ (cross entropy loss is used) to measure the difference between the resized ground-truth $R_1(Y)$ and the coarse prediction $s_1(I)$ (after softmax). In other words, these operations can be written as:

$$s_1(I) = \text{conv}_{3\times3}\big(f_5(I)\big)$$
$$\ell_1 = \text{Loss}\Big(R_1(Y), \text{softmax}\big(s_1(I)\big)\Big) \quad (1)$$

where $\text{conv}_{3\times3}(\cdot)$ denotes the $3 \times 3$ convolution and $\text{Loss}(\cdot)$ denotes the cross entropy loss function.

Now we explain how to get the subsequent segmentation map $s_k(I)$ ($k > 1$) with larger spatial dimensions. One simple solution is to upsample the previous segmentation map $s_{k-1}(I)$. However, this upsampled segmentation map will be very coarse. In order to derive a more precise segmentation, we use the idea of skip-connections [9, 16]. The basic idea is to make use of the outputs from one of the convolutional feature layers in the encoder network. Since the convolutional feature layer contains more precise spatial information, we can combine it with the upsampled segmentation map to obtain a refined segmentation map (see Fig 2 for detailed illustration of our refinement module). In our model, we simply concatenate the outputs from the skip layer with the upsampled decoder layer to create a larger feature map, then use $3 \times 3$ convolution to convert the channel dimension to $C$. For example, the label map $s_2(I)$ is obtained from upsampled $s_1(I)$ and the convolutional feature map $f_5(I)$ (see Fig. 1 for details on these skip connections). We can then define a loss function on this (larger) segmentation map by comparing $s_k(I)$ with the resized ground-truth segmentation map $R_k(Y)$ of the corresponding size. These operations can be summarized as follows:

$$s_k(I) = \text{conv}_{3\times3}\Big(\text{concat}\Big(\text{upsample}\big(s_k(I)\big), f_{7-k}(I)\Big)\Big) \quad (2)$$
$$\ell_k = \text{Loss}\Big(R_k(Y), \text{softmax}\big(s_k(I)\big)\Big), \quad \text{where } k = 2, .., 6 \quad (3)$$

Our model is trained using back-propagation to optimize the summation of all the losses $\sum_{k=1}^{6} \ell_k$.

Fig. 3 shows the advantage of progressively refinement of the segmentation map $s(I)$ in the decoder network. The finer feature maps from the convolution layers help to recover image details, and therefore the segmentation map produced at each stage of the decoder network becomes more accurate.

# 5. Experiments

Section 5.1 explains some implementation details and several baseline methods that we compare with. We
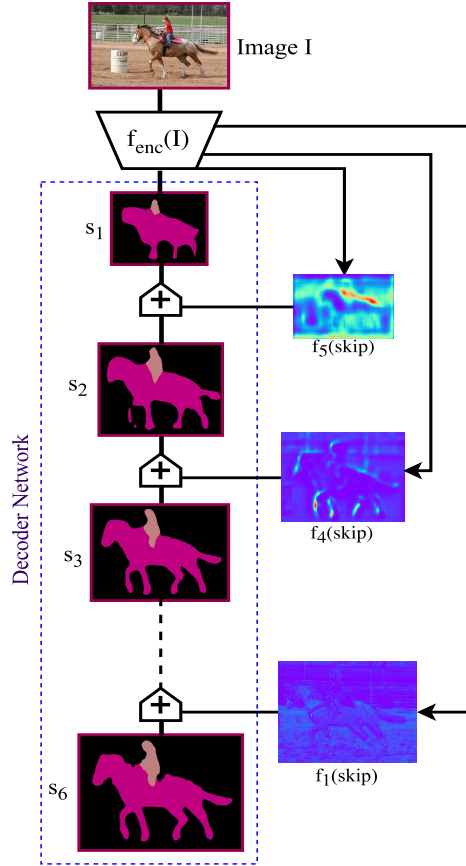


Figure 3. Visualization of the refinement process in the network. Coarse segmentation map ($s_1$) produced by the encoder network $f_{enc}$ is fed to the decoder network (shown in dotted rectangular box). The decoder network then progressively refines the obtained segmentation map using the refinement module (RE) (see Fig. 2) to recover the image details step by step (see $s_1, s_2, s_3, ...s_6$ in the figure). Basically, the refinement process is integrating the finer feature maps (activations shown as heat maps on the right) from the corresponding convolution layers with the segmentation maps. As a result, the segmentation map gets progressively refined throughout the decoder network.

then demonstrate experimental results on three challenging benchmark datasets: the PASCAL VOC 2012 dataset (Sec. 5.2), the Cambridge Driving Labeled Video (CamVid) dataset (Sec. 5.3), and the SUN RGB-D dataset (Sec. 5.4).

## 5.1. Implementation Details and Baselines

Following [1], we train our network on the CamVid and SUN RGB-D datasets using the class balancing technique in [6] as there is a large variation in the number of pixels for each category. For example, most of the pixels on the CamVid dataset are labeled as roads and cars. Each class is assigned a weight in the loss function. The weights are calculated by taking the ratio of the median of the class fre-

| method | aero | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbike | person | plant | sheep | sofa | train | tv | mean IoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| FCN-8s [16] | 76.8 | 34.2 | 68.9 | 49.4 | **60.3** | 75.3 | 74.7 | 77.6 | 21.4 | 62.5 | 46.8 | **71.8** | 63.9 | 76.5 | 73.9 | 45.2 | **72.4** | 37.4 | 70.9 | 55.1 | 62.2 |
| SegNet [1] | 74.5 | 30.6 | 61.4 | **50.8** | 49.8 | 76.2 | 64.3 | 69.7 | **23.8** | 60.8 | **54.7** | 62.0 | 66.4 | 70.2 | 74.1 | 37.5 | 63.7 | 40.6 | 67.8 | 53.0 | 59.1 |
| SegNet+ DS | 70.0 | 35.7 | 72.8 | 42.7 | 53.0 | 75.6 | 71.6 | 72.8 | 23.2 | **69.5** | 46.8 | 64.6 | 68.8 | 78.9 | 76.0 | 41.9 | 65.7 | **44.1** | 65.7 | 53.7 | 61.1 |
| **LRN** | **79.3** | **37.5** | **79.7** | 47.7 | 58.3 | **76.5** | **76.1** | **78.5** | 21.9 | 67.7 | 47.6 | 71.2 | **69.1** | **82.1** | **77.5** | **46.8** | 70.1 | 40.3 | **71.5** | **57.4** | **64.2** |

Table 1. Quantitative results on the PASCAL VOC 2012 test set. We report the IoU score for each class and the mean IoU score.



Figure 4. Qualitative comparisons between the FCN[16] model and our LRN model on the PASCAL VOC 2012 dataset.

quencies to the individual class frequencies from the training dataset. For the CamVid and SUN RGB-D datasets, training images are resized to a common size of $360 \times 480$ pixels. For the PASCAL VOC 2012 dataset, when fed into LRN, we randomly crop a $320 \times 320$ region from each image during training and subtracted a mean pixel provided by VGG net at each position. During test, we generate a $512 \times 512$ segmentation map and crop the final label out from it similar to [4].

We implemented LRN using Caffe [11]. A GTX Titan X GPU was used both in training and testing for acceleration. The whole network was trained end-to-end by using back propagation algorithm with the following hyper-parameters: mini-batch size (10), learning rate (0.001), momentum (0.9), weight decay (0.0005), number of iterations ($\sim 80,000$). After each $50,000$ iterations we reduce the learning rate by a factor of $0.1$. The encoder network is fine-tuned from a pre-trained VGG-16 network.

To demonstrate the merit of individual component of our model, we perform an ablation study by comparing with the following baselines.

**SegNet**: This is the SegNet model in [1]. Whenever possible, we use the reported numbers in [1] in our comparison.

**SegNet+Deep Supervision (SegNet+DS)**: This is the SegNet model with multiple loss functions to provide deep supervisions.

## 5.2. PASCAL VOC 2012

PASCAL VOC 2012 is a challenging segmentation benchmark [7]. It consists of 20 object categories and the background. Following previous works such as [4, 16, 1], we employ 10,582 images for training, 1,449 images for validation, and 1,456 images for testing. We train our model on the trainval set and evaluate our model on the test set. Quantitative results are shown in Table 1. We can see that SegNet+DS performs better than SegNet. This shows that deep supervisions (i.e. multiple losses) help boosting the performance. In addition, Table 1 shows that our model
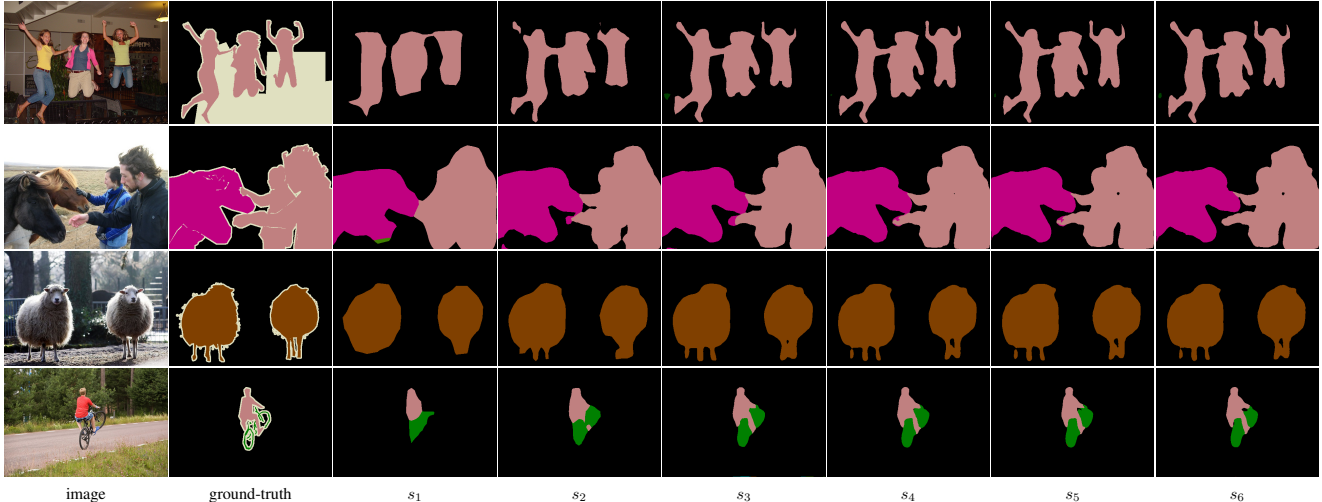
Figure 5. Stage-wise visualization of semantic segmentation results on PASCAL VOC 2012. We show the segmentation result obtained from each of the 6 label maps $s_k$ ($k = 1, 2, ..., 6$) produced by the decoder of our model.

| Method | Building | Tree | Sky | Car | Sign | Road | Pedestrian | Fence | Pole | Sidewalk | Bicyclist | Class avg | Mean IoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| [26] | 85.3 | 57.3 | 95.4 | 69.2 | 46.5 | **98.5** | 23.8 | 44.3 | 22.0 | 38.1 | 28.7 | 55.4 | - |
| [23] | 87.0 | 67.1 | **96.9** | 62.7 | 30.1 | 95.9 | 14.7 | 17.9 | 1.7 | 70.0 | 19.4 | 51.2 | - |
| [14] | 81.5 | 76.6 | 96.2 | 78.7 | 40.2 | 93.9 | 43.0 | 47.6 | 14.3 | 81.5 | 33.9 | 62.5 | - |
| SegNet [1] | 88.0 | 87.3 | 92.3 | 80.0 | 29.5 | 97.6 | 57.2 | 49.4 | 27.8 | 84.8 | 30.7 | 65.9 | 50.2 |
| SegNet+DS | 87.0 | 85.1 | 77.1 | 79.8 | 56.3 | 94.8 | 71.1 | 50.5 | 39.7 | 88.6 | 50.6 | 70.1 | 53.7 |
| **LRN** | **89.8** | **88.1** | 78.5 | **86.3** | **61.2** | 96.8 | **82.1** | **59.0** | **45.4** | **92.6** | **69.7** | **77.2** | **61.7** |

Table 2. Quantitative results on the CamVid [2] dataset. For each method, we report the accuracy for each class and the average class accuracy. We also report the mean IoU of all classes if it is available. Our model (LRN) outperforms all the other methods.

(LRN) outperforms SegNet+DS. This demonstrates that the refinement modules give additional performance improvement. Note that we also perform better than the popular FCN [16] model for semantic segmentation. We present qualitative results in Fig. 4. The results show that LRN can preserve subtle details of objects. To further illustrate the effectiveness of the coarse-to-fine segmentation in our model, Fig. 5 shows each of the label maps produced by the refinement network in our model. We can see that the coarse-to-fine refinement scheme can progressively improve the details of segmentation maps by recovering the missing parts (e.g., the leg of the person and sheep in the top and $3^{rd}$ row respectively).

## 5.3. CamVid

The CamVid dataset [2] consists of a ten minutes video footage with 701 images. The video footage was recorded from driving around various urban settings and is particu-

larly challenging given variable lighting settings (e.g. dusk, dawn). The dataset provides ground-truth labels that associate each pixel with one of 32 semantic class. Following [1], we consider 11 larger semantic classes (road, building, sky, tree, sidewalk, car, column-pole, fence, pedestrian, bicyclist, and sign-symbol) among the 32 semantic classes in our experiment. Table 2 shows the comparison of our model with other state-of-the-art approaches on this dataset. For each method, we report the accuracy for each class and the average class accuracy. We also report the mean IoU score if it is available. Again, our LRN model outperforms SegNet+DS, which in turn outperforms SegNet. We show some qualitative results on this dataset in Fig. 6.

## 5.4. SUN RGB-D

We also evaluate our model on the SUN RGB-D dataset [21] which contains 5,285 training and 5,050 test images. The images consist of indoor scenes of varying res-
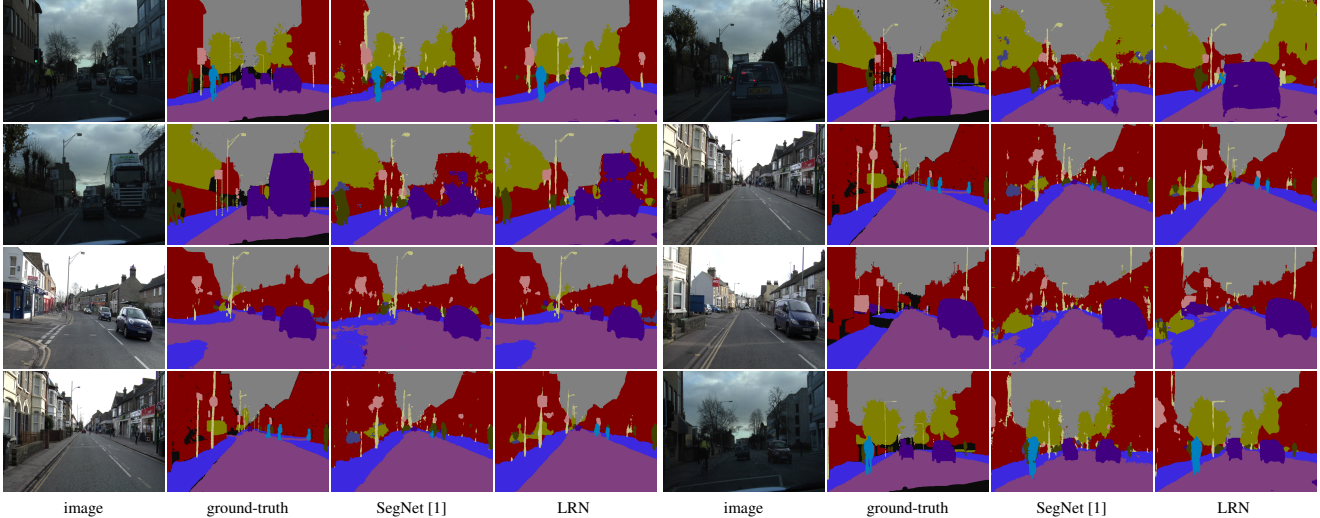
Figure 6. Qualitative results on the CamVid dataset. LRN predictions capable of retaining the shape of smaller object categories (e.g. column-pole, side-walk, bicyclist, and sign-symbols) accurately compared to SegNet [1].

olution and aspect ratio. There are 37 indoor scene classes with corresponding segmentation labels (background is not considered as a class and is ignored during training and testing). The segmentation labels are instance-wise, i.e. multiple instances of same class in an image have different labels. We convert the ground-truth labels into class-wise segmentation labels so that all instances of the same class have the same corresponding label. Although the dataset also contains depth information, we only use the RGB images to train and test our model. Quantitative results on this dataset are shown in Table 3. Our LRN model achieves better mean IoU than other baselines on this dataset. We show some qualitative results in Fig. 7.

### 5.5. Ablation Analysis

We perform further ablation analysis to demonstrate the benefit of our coarse-to-fine approach. Our LRN model produces 6 label maps $s_k$ ($k = 1, 2, ..., 6$) of different spatial dimensions (see Fig. 1). It is possible to obtain the final full-sized label map from any of these $s_k$ by upsamling $s_k$ (using bilinear interpolation) to the spatial dimensions of the input image. Table 4 shows the results of this stage-wise analysis on the three datasets. Note that on the PASCAL 2012 dataset, we perform this analysis by learning the model on the training set, and test on the validation set. The reason is that we do not have direct access to the ground-truth labels on the PASCAL 2012 test set. The evaluation on the test set has to be done on the PASCAL evaluation server. The results in Table 4 show that if we simply take a label map at one of the early stages of the refinement network and upsample it as the final prediction, the result is not as good as using the label map produced at the last stage. This is reasonable

since the label map at an early stage only provides labels corresponding to smaller spatial dimensions.

### 5.6. Discussion

The main focus of our paper is the coarse-to-fine label refinement network architecture. Smaller objects are hard to classify due to the limitation of having a fixed size kernel window. For instance, the poles and sign-symbols in the CamVid dataset [2] often carry significant importance and are important to scene understanding for many applications. This is where our model performs especially well considering the baseline [1]. Table 2 shows our improvement for classifying column-pole or sign symbols. Figure 8 shows visually that the prediction of finer details such as the column pole are improved. This improvement is mainly due to the recurrent connection of the weights in the encoder and decoder networks respectively. In addition, the weighted summation of pixels through skip connections help significantly to recover details of smaller objects even after lowering the resolution through several pooling layers.

### 6. Conclusion

We have proposed a novel CNN architecture for semantic segmentation. Our model produces segmentation labels in a coarse-to-fine manner. The segmentation labels at coarse levels are used to progressively refine the labeling produced at finer levels. We introduce multiple loss in the network to provide deep supervisions at multiple stages of the network. Our experimental results demonstrate that our proposed network improves the performance of several baseline approaches. Although we focus on semantic segmentation, the proposed architecture is quite general and can be applied for other
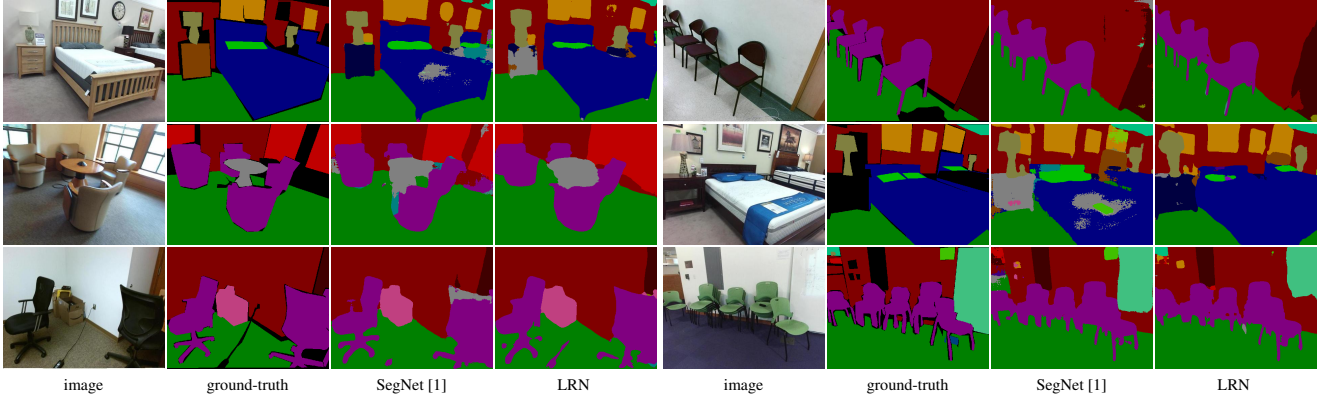
7

Figure 7. Qualitative results on the SUN RGB-D dataset.

| Method | Wall | Floor | Cabinet | Bed | Chair | Sofa | Table | Door | Window | Bookshelf | Picture | Counter | Blinds | Desk | Shelves | Curtain | Dresser | Pillow | Mirror | Floor Mat | Clothes | Ceiling | Books | Fridge | TV | Paper | Towel | Curtain | Box | Whiteboard | Person | Night Stand | Toilet | Sink | Lamp | Bathtub | Bag | Class avg | mean IoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| [12] | 80.2 | 90.9 | **58.9** | 64.8 | 76.0 | 58.6 | 62.6 | 47.7 | **66.4** | 31.2 | 63.6 | 33.8 | **46.7** | 19.7 | 16.2 | **67.0** | 42.3 | **57.1** | 39.1 | 0.1 | 24.4 | **84.0** | 48.7 | 21.3 | 49.5 | 30.6 | 18.8 | 0.1 | 24.1 | 56.8 | 17.9 | 42.9 | 73.0 | 66.2 | 48.8 | 45.1 | 24.1 | 45.9 | 30.7 |
| SegNet [1] | **86.6** | 92.0 | 52.4 | **68.4** | 76.0 | 54.3 | 59.3 | 37.4 | 53.8 | 29.2 | 49.7 | 32.5 | 31.2 | 17.8 | 5.3 | 53.2 | 28.8 | 36.5 | 29.6 | 0.0 | 14.4 | 67.7 | 32.4 | 10.2 | 18.3 | 19.2 | 11.5 | 0.0 | 8.9 | 38.7 | 4.9 | 22.6 | 55.6 | 52.7 | 27.9 | 29.9 | 8.1 | 35.6 | 26.3 |
| SegNet+DS | 78.9 | **92.4** | 52.5 | 58.1 | 74.1 | **61.9** | **62.6** | **64.1** | 60.6 | 12.1 | **70.7** | **38.7** | 37.5 | 13.1 | **21.7** | 62.3 | 41.2 | 55.1 | **40.1** | **25.0** | **42.4** | **84.0** | 58.1 | **44.6** | **51.4** | **44.6** | **27.6** | 4.9 | 24.8 | **60.9** | 14.4 | 41.1 | **77.3** | **68.7** | **60.3** | 63.4 | **28.6** | **49.2** | 31.2 |
| **LRN** | 84.6 | 90.8 | 55.6 | 65.2 | **79.4** | 46.2 | 59.8 | 44.6 | 63.6 | **32.6** | 66.1 | 37.7 | 39.8 | **21.0** | 15.7 | 60.2 | **42.9** | 46.5 | 30.4 | 20.5 | 23.0 | 82.8 | **59.9** | 26.4 | 44.4 | 32.5 | 23.4 | **10.8** | **30.9** | 51.6 | **21.7** | **43.0** | 75.7 | 61.8 | 51.3 | **63.9** | 28.4 | 46.8 | **33.1** |

Table 3. Quantitative results on the SUN RGB-D dataset. We report the accuracy for each class, average class accuracy, and the mean IoU of all classes.

| Stage | Mean IoU (%) | | Stage | Mean IoU (%) | | Stage | Mean IoU (%) |
|---|---|---|---|---|---|---|---|
| $s_1$ | 58.4 | | $s_1$ | 50.9 | | $s_1$ | 31.0 |
| $s_2$ | 61.6 | | $s_2$ | 55.5 | | $s_2$ | 31.6 |
| $s_3$ | 61.9 | | $s_3$ | 59.1 | | $s_3$ | 32.4 |
| $s_4$ | 62.6 | | $s_4$ | 60.9 | | $s_4$ | 32.7 |
| $s_5$ | 62.5 | | $s_5$ | 61.4 | | $s_5$ | 32.8 |
| $s_6$ | **62.8** | | $s_6$ | **61.7** | | $s_6$ | **33.1** |
| (a) PASCAL VOC 2012 | | | (b) CamVid | | | (c) SUN RGB-D | |

Table 4. Stage-wise mean IoU on (a) PASCAL VOC 2012 validation set; (b) CamVid dataset; (c) SUN RGB-D dataset. We can see that, from $s_1$ to $s_6$, mean IoU is progressively enhanced in all the dataset. Note that $s_1..s_6$ are not results of separate training. They are obtained from different stage of the decoder in our model. The result from $s_1$ is implicitly affected by the supervisions provided at $s_2..s_6$, so in Camvid and SUN-RGBD dataset it is reasonable for $s_1$ to produce better performance than SegNet [1].
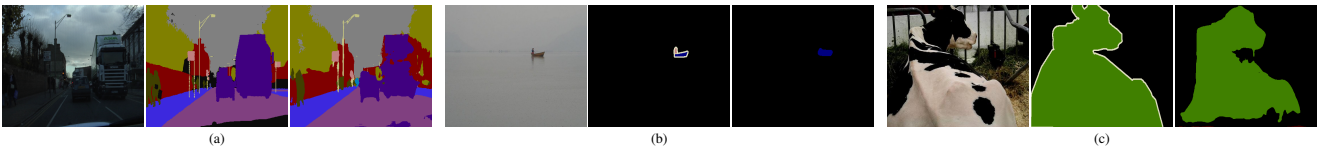


Figure 8. Qualitative analysis of dense predictions on a few problematic images of CamVid [2] and Pascal VOC 2012 [7] datasets. Our predictions are more refined and sensitive to smaller and thinner objects (see (a) and (b)). Performance is also better for larger objects (see (c)) because of the refinement modules and guidance driven by early constraints provided by deep supervision.

pixel-wise labeling problems. In addition, experimental results based on ablation analysis suggest generality in the value of coarse-to-fine predictions, deep supervisions and skip connections respectively, with the implication that a wide range of CNNs may benefit from these architectural modifications.

# References

[1] V. Badrinarayanan, A. Kendall, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *arXiv preprint arXiv:1511.00561*, 2015.

[2] G. J. Brostow, J. Fauqueur, and R. Cipolla. Semantic object classes in video: A high-definition ground-truth database. *Pattern Recognition Letters*, 30(2):88–97, 2009.

[3] P. J. Burt and E. H. Adelson. The laplacian pyramid as a compact image code. *IEEE Transactions on Communications*, 1983.

[4] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected CRFs. In *International Confernece on Learning Representations*, 2015.

[5] E. Denton, S. Chintala, A. Szlam, and R. Fergus. Deep generative image models using a laplacian pyramid of adversarial networks. In *Advances in Neural Information Processing Systems*, 2015.

[6] D. Eigen and R. Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *IEEE International Conference on Computer Vision*, 2015.

[7] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL visual object classes (VOC) challenge. *International Journal of Computer Vision*, 88(2):303–338, 2010.

[8] C. Farabet, C. Couprie, L. Najman, and Y. LeCun. Learning hierarchical features for scence labeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2013.

[9] B. Hariharan, P. Arbelaez, R. Girshick, and J. Malik. Hypercolumns for object segmentation and fine-grained localization. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015.

[10] S. Honari, J. Yosinski, P. Vincent, and C. Pal. Recombinator networks: Learning coarse-to-fine feature aggregation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.

[11] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv:1408.5093*, 2014.

[12] A. Kendall, V. Badrinarayanan, and R. Cipolla. Bayesian SegNet: Model uncertainty in deep convolutional encoder-decoder architectures for scene understanding. *arXiv preprint arXiv:1511.02680*, 2015.

[13] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, 2012.

[14] L. Ladický, P. Sturgess, K. Alahari, C. Russell, and P. H. S. Torr. What, where & how many? combining object detectors and CRFs. In *European Conference on Computer Vision*, 2010.

[15] C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu. Deeply-supervised nets. In *Proceedings of AISTATS*, 2015.

[16] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015.

[17] M. Mostajabi, P. Yadollahpour, and G. Shakhnarovich. Feed-forward semantic segmentation with zoom-out features. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015.

[18] H. Noh, S. Hong, and B. Han. Learning deconvolution network for semantic segmentation. In *IEEE International Conference on Computer Vision*, 2015.

[19] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems*, 2015.

[20] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.

[21] S. Song, S. P. Lichtenberg, and J. Xiao. SUN RGB-D: A RGB-D scene understanding benchmark suite. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015.

[22] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015.

[23] J. Tighe and S. Lazebnik. SuperParsing: Scalable nonparametric image parsing with superpixels. *International Journal of Computer Vision*, 101(2):329–349, 2013.

[24] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatialtemporal features with 3d convolutional networks. In *IEEE International Conference on Computer Vision*, 2015.

[25] S. Xie and Z. Tu. Holistically-nested edge detection. In *IEEE International Conference on Computer Vision*, 2015.

[26] C. Zhang, L. Wang, and R. Yang. Semantic segmentation of urban scenes using dense depth maps. In *European Conference on Computer Vision*, 2010.