

Learning to segment from misaligned and partial labels

Simone Fobi*
Columbia University
sf2786@columbia.edu

Jayant Taneja
University of Massachusetts Amherst
jtaneja@umass.edu

Terence Conlon
Columbia University
tmc2180@columbia.edu

Vijay Modi
Columbia University
modi@columbia.edu

ABSTRACT

To extract information at scale, researchers increasingly apply semantic segmentation techniques to remotely-sensed imagery. While fully-supervised learning enables accurate pixel-wise segmentation, compiling the exhaustive datasets required is often prohibitively expensive. As a result, many non-urban settings lack the ground-truth needed for accurate segmentation. Existing open source infrastructure data for these regions can be inexact and non-exhaustive. Open source infrastructure annotations like OpenStreetMaps are representative of this issue: while OpenStreetMaps labels provide global insights to road and building footprints, noisy and partial annotations limit the performance of segmentation algorithms that learn from them.

In this paper, we present a novel and generalizable two-stage framework that enables improved pixel-wise image segmentation given misaligned and missing annotations. First, we introduce the Alignment Correction Network to rectify incorrectly registered open source labels. Next, we demonstrate a segmentation model – the Pointer Segmentation Network – that uses corrected labels to predict infrastructure footprints despite missing annotations. We test sequential performance on the Aerial Imagery for Roof Segmentation dataset, achieving a mean intersection-over-union score of 0.79; more importantly, model performance remains stable as we decrease the fraction of annotations present. We demonstrate the transferability of our method to lower quality data sources, by applying the Alignment Correction Network to OpenStreetMaps labels to correct building footprints; we also demonstrate the accuracy of the Pointer Segmentation Network in predicting cropland boundaries in California from medium resolution data. Overall, our methodology is robust for multiple applications with varied amounts of training data present, thus offering a method to extract reliable information from noisy, partial data.

KEYWORDS

segmentation; misaligned and missing labels; open source data

1 INTRODUCTION

Processing remotely-sensed imagery is a promising approach to evaluate ground conditions at scale for little cost. Algorithms that intake satellite imagery have accurately measured crop type [34],[21], cropped area [11], building coverage [41] [40], urbanization [1], and road networks [6] [42]. However, successful implementation of image segmentation algorithms for remote sensing applications depends on large amounts of data and high-quality annotations. Wealthy, urbanized settings can more readily apply segmentation



Figure 1: Types of label noise present in open source data. Building footprints are the class of interest.

algorithms, due to either the presence of or the ability to collect significant amounts of carefully annotated data. In contrast, more rural regions often lack the means to exhaustively collect ground truth data. Some open source datasets exist for such settings, and by successfully coupling these annotations with remotely sensed imagery, researchers can gain insights into the status of infrastructure and development where well-curated sources of these data do not exist. [20] [2].

Although these global open source ground truth datasets – e.g. OpenStreetMaps (OSM) – offer large amounts of labels for use at no cost, the annotations within suffer from multiple types of noise [28] [4]: *missing or omitted annotations*, defined as objects being present in the image and not existing in the label [28]; *misaligned annotations* occur when annotations are translated and/or rotated from its true position [38]; and *incorrect annotations* – annotations that do not directly correspond to the object of interest in the image. Figure 1 presents examples of these three types of label noise.

Noisy datasets present a training challenge when using traditional segmentation algorithms, as the model cannot learn to associate image features and target labels when the relationship is obscured by noise. To address the issues of misaligned and omitted annotations, and in order to extract information from imperfect data, we present a simple and generalizable method for pixel-wise image segmentation. First, we address annotation misalignment by proposing an Alignment Correction Network (ACN). With a small number of images and human verified ground truth annotations, the ACN learns to correct misaligned labels. Next, the corrected open source annotations are used to train the Pointer Segmentation Network (PSN), a model which takes in a point location and identifies the object containing that point. Learning associations from a representative point is a widely acknowledged method of object detection: [5] notes that an intuitive way for humans to refer to an object is through the action of pointing. By ‘pointing-out’ the

object instance of interest, our network ignores other instances that may not have corresponding annotations, therefore preventing performance degradation caused by annotation-less instances within the image. As a result, our sequential approach presents a method for handling misaligned data as well as varying levels of label completeness without explicitly changing the loss function to compensate for noise. While our approach cannot replace large amounts of carefully annotated outlines, it can complement existing open source datasets and algorithms, reduce the cost of obtaining large amounts of full annotations, and allow researchers to extract information from imperfect datasets. This paper’s key contributions are as follows:

- We introduce the Alignment Correction Network (ACN), a means to verify and correct misaligned annotations using a small amount of human verified ground truth labeled data.
- We propose the Pointer Segmentation Network (PSN), a model that can reliably predict polygon boundaries on remotely-sensed imagery despite omitted training annotations and without requiring any bespoke loss functions.
- We demonstrate the applicability of our methodology to three different segmentation problems: building footprint detection with a highly-accurate dataset, building footprint detection with noisier training data, and cropland boundary prediction.

Taken as a whole, our approach enables resource constrained actors to use large amounts of misaligned and partial labels – coupled with a very small amount of human verified ground truth annotations – to train image segmentation algorithms for a variety of tasks. The rest of the paper is organized as follows: In *Related Work*, we discuss related literature; in *Methods*, we describe our novel methodological contributions; in *Results*, we present results for the ACN and the PSN for all segmentation tasks; and in *Conclusion*, we restate our most salient findings.

2 RELATED WORK

Computer vision researchers have recently made numerous advances in semantic segmentation, in applying state-of-the art techniques to remote sensed imagery, and in learning from noisy datasets; we discuss some important contributions to the literature below.

Existing Segmentation Approaches

Primarily based on improvements to deep convolutional neural networks (DCNN) architectures, researchers have achieved record performances for a variety of different segmentation tasks. Fully convolutional encoder-decoder type architectures – one type of DCNN – take in an image and output a per-pixel prediction for the class of interest [25]. Some architectures use symmetric networks with skip connections to perform pixel-wise predictions [33] [3]. Alternatively, two-stage detection algorithms first perform region proposal – areas that have a high likelihood of containing the object of interest – and then detect objects within the identified regions [17] [16] [31]. Modifications to two-stage detection algorithms have enabled semantic segmentation of images, whereby individual pixels in an image are placed into one of a number of classes [18] [22]. Development of these segmentation architectures has been facilitated by large, comprehensive datasets which enable the implementation of these algorithms in a fully supervised approach:

here, every object in the image and its corresponding annotation are used in the learning process [12] [27] [23].

Applying Deep Learning to Remote Sensed Imagery

Multiple projects have leveraged satellite imagery to answer various questions on land use, road quality, object detection, consumption expenditure: by linking sparse ground truth with abundant imagery, researchers can extrapolate trends in existing data to areas where labeled data do not exist [35], [10], [19]. Alternatively, some works have proposed neural network architectures that sidestep training data constraints and the relative lack of labeled ground-truth in remote areas [24] [30]. Jean et al. combine Google maps daytime images (provided by DigitalGlobe), nighttime lighting, and survey data to estimate poverty for multiple African countries [29]. High resolution daytime images were used to train a model to predict nighttime lights as measured by DMSP-OLS; features extracted from the last layer of the model were then used to estimate household expenditure or wealth. Results from this paper suggest that predictions about economic development can be made from remote sensed data using features derived from imagery; this insight provides additional motivation for developing methods that extract information from noisy imagery datasets.

Learning From Noisy Annotations

The problem of poor-quality training data, especially in rural areas, for segmentation tasks is well known: [26] acknowledge the variability in coverage of open source data in Kenya and observe significant degradation of coverage as one moves away from urban settings. Coverage degradation from urban to rural areas is also seen in South Africa[36], Brazil[7] and Botswana[39]. [37] estimates the effects of multiple types of training data noise, including misalignment and missing annotations, finding that as noise levels increase, both precision and recall decrease. For applications such as measuring building or field area which are useful in downstream analysis of wealth, crop yield and more, high noise levels decrease the ability to successfully use segmentation algorithms. Several works tackle the problem of learning from imperfect labels. [28] propose new loss functions to address noisy labels in aerial images. [38] [15] both focus on the issue of misalignment: [15] uses a self-supervised approach to align cadaster maps, and while the method proposed in [38] maximizes the correlation between annotations and outputs from a building prediction CNN, it assumes buildings in small groups have the same alignment error. Our two-stage approach builds upon existing convolutional frameworks common to many noise correction approaches. However our approach relies on the well-known binary cross entropy loss function, addresses both misalignment and omitted annotation, and does not require that all misalignments are identical. Thus serving as an attractive alternative when noisy labels are present.

3 METHODS

Traditional segmentation methods take an image input x_i and aim to learn a function $f(\mathbf{x})$ that predicts a single channel label \hat{v}_i containing all building instances present in the image. Equation 1 shows the learned function given x_i , where v_i^a is the single channel label of instance a in image x_i and there are a total of A instances

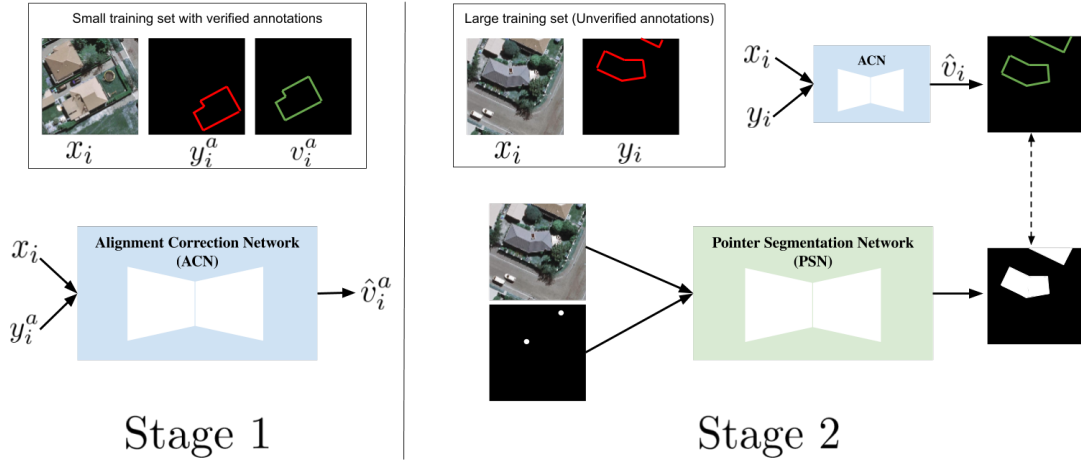


Figure 2: Summary of our two-stage approach to segment from noisy annotations. Stage 1: The ACN uses an image (x_i) and label (y_i^a) with a single misaligned annotation to predict a corrected annotation \hat{v}_i^a containing the realigned annotation. Random shifts between ± 10 pixels are applied to v_i^a to obtain y_i^a . The network is trained with a small set of images (x) and verified ground truth annotations (v). Stage 2: A large noisy training set is first realigned with the ACN. Realigned, incomplete annotations are used for supervision. The PSN uses selected points from available instances, x_i and \hat{v}_i to learn the segmentation task.

in that image:

$$f(x_i) \rightarrow \hat{v}_i \quad (1)$$

$$s.t. \quad \hat{v}_i = \hat{v}_i^1 \cup \hat{v}_i^2 \dots \cup \hat{v}_i^A$$

3.1 Alignment Correction Network

Misalignment occurs when there is a registration difference between an object in an image and its annotation. In remote sensing, misaligned annotations may occur for a number of reasons, including human error and imprecise projections of the image [15]. There are two types of annotation alignment errors: 1) translation errors, where the annotation is shifted relative to the object, and 2) rotation errors, where the annotation is rotated relative to the object. [38] suggest that translation errors are more frequent for OpenStreetMaps in rural areas. Thus in this paper, we only address translation errors present in open source data. We propose an Alignment Correction Network (ACN) that takes in an image x_i and a label y_i^a containing one misaligned instance a . The ACN outputs a label \hat{v}_i^a containing the predicted, corrected annotation. \hat{v}_i^a is compared to v_i^a to learn optimal weights for the network. During training, the misaligned label y_i^a is obtained by applying random x-y shifts, between ± 10 pixels to v_i^a . Sensitivity to the ± 10 pixels translation shift is discussed in the results.

When multiple misaligned instances are present in an image, the instances are corrected independently. This approach is chosen for two reasons: it allows instances within an image to have varying degrees of translation error and it also enables the network to be robust to incomplete labels with missing instances. Here, a small dataset of images (x) and carefully verified ground truth labels (v) are used to train the ACN as shown in Stage 1 of Figure 2.

3.2 Pointer Segmentation Network

Assuming m available annotations – $v_i^1 \dots v_i^m$, where $m < A$ – common algorithms will struggle to implement Equation 1, as some predicted object instances will not have corresponding true labels for comparison during training. To address this issue, we introduce the PSN, a network that learns to segment an image using only m available annotations. The PSN takes as inputs an image x_i and a single channel of points specifying selected instances to be segmented, and it outputs a segmentation mask only for the selected instances. We specify the fraction of instances to be used for training using a parameter α , where α is the number of selected instances divided by the number of available instances. Equation 2 shows this formulation, where $p_i(\alpha)$ specifies a point within each selected instance, and $\hat{v}_i(\alpha)$ denotes the predicted label for instances specified by $p_i(\alpha)$:

$$f(x_i, p_i(\alpha)) \rightarrow \hat{v}_i(\alpha) \quad (2)$$

By including a single channel containing points $p_i(\alpha)$, our PSN segments only instances that are associated with the points. This offers two benefits: first, we simplify the learning task to specify instances of interest, and second, the network can be trained with common binary cross entropy loss. To handle varying extents of missing annotations, the model is trained by randomly picking α for every image in each epoch; at inference time, all instances of interest are specified using points.

In the sequential training configuration, the ACN is used to correct a training dataset that is then inputted to the PSN for object segmentation; this process is shown in Stage 2 of Figure 2. Binary cross-entropy loss is used for all networks. Both ACN and PSN use the same baseline architecture (lightUNet) shown in Appendix A

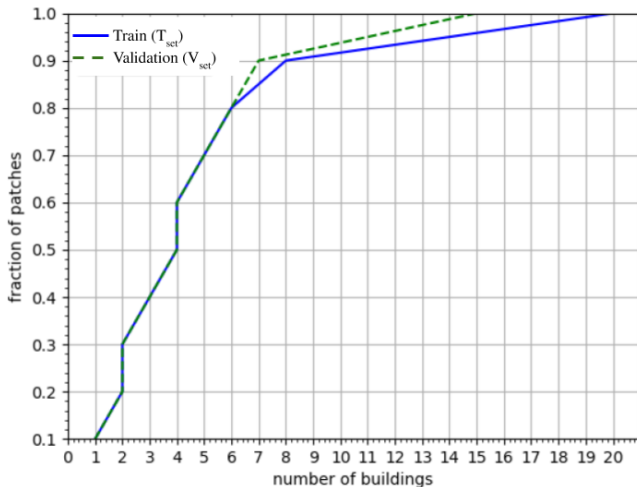


Figure 3: CDF of the number of buildings present in 128x128 patches of the 30cm-resampled AIRS dataset.

and further explained in the results, albeit modified by the number of input channels.

4 DATA

Three separate datasets are used to train and test the performance of the ACN and the PSN, all described below. During training and testing, we only use images that contain labels.

4.1 Aerial Imagery for Roof Segmentation

We use the Aerial Imagery for Roof Segmentation (AIRS) dataset to establish baseline performances for both the ACN and PSN. The AIRS dataset covers most of Christchurch (457km²), New Zealand and consists of orthorectified aerial images (RGB) at a spatial resolution of 7.5 cm with over 220,000 building annotations, split into a training set (T_{set}) and a validation set (V_{set}). The AIRS dataset provides all building footprints within the dataset coverage area. To mimic more readily-available data, we resample the imagery to 30 cm, an approach which creates imagery more similar to that provided by Google Earth. Next, we slice the resampled images into 128 by 128 pixel patches and discard all patches in which the area occupied by buildings is less than 10 % of the total area – this methodology ensures that patches with multiple buildings are selected. Other than this basic filtering, we preserve T_{set} and V_{set} .

After resampling and filtering, we obtain 99,501 and 10,108 patches from the T_{set} and V_{set} , respectively. We further split T_{set} into 80:20 fractions, where 80% is used for training and 20% for validation. V_{set} is withheld and used as a test set to evaluate performance. Figure 3 shows the fraction of patches for a given number of buildings in T_{set} and V_{set} . Note that some patches contain partial buildings.

4.2 OpenStreetMaps

Humanitarian OpenStreetMaps (OSM), through free, community-driven annotation efforts, provides building footprints by country on their Humanitarian Data Exchange (HDX) platform. While this

data provides the best (and only) ground truth for many parts of the world, label quality is highly heterogeneous, both in terms of footprint alignment and coverage. In order to test the performance of the ACN on these incomplete and misaligned building footprints, we pair OSM annotations for Kenya [13] with selected DigitalGlobe tiles from Western Kenya (a box enclosed by 0.176 S, 0.263 S, 34.365 E, and 34.453 E) and closer to Nairobi (a box enclosed by 1.230 S, 1.318 S, 36.738 E, and 36.826 E). The DigitalGlobe tiles have a 50 cm spatial resolution and were collected between 2013 and 2016. Slices measuring 128 by 128 pixels were generated from the DigitalGlobe images, which we then couple with overlapping OSM building labels. We generated human verified ground truth annotations for 500 of the image patches.

4.3 California Statewide Cropping Map

We also use crop maps and decameter imagery to demonstrate the flexibility of the PSN. The California Department of Water Resources provides a Statewide Cropping Map for 2016 [32]; we pair this shapefile with Sentinel-2 satellite imagery to learn to extract crop extents [14]. Red, blue, green, and near-infrared bands – all at 10m resolution – are acquired from a satellite pass on August 30, 2016; the bands cover the same spatial extent as Sentinel tile 11SKA (a box enclosed by 37.027 N, 36.011 N, 120.371 W, and 119.112W). Cropped polygons larger than 500m² are taken from the California cropping map and are eroded by 5m on all sides to ensure that field boundaries are distinct at a 10m spatial resolution. We split the 110km x 110km tile into images patches measuring 128 by 128 pixels and remove all slices that do not cover any cropped areas, leaving a total of 5,681 patches containing an average of 17 fields per patch; these images are split into training, validation, and test sets at a ratio of 60/20/20.

5 RESULTS

For all model testing, we report the mean intersection-over-union (mIOU), defined as the intersection of the predicted and true label footprints divided by the union of the same footprints, averaged across the testing dataset.

5.1 Baseline Model

We establish the performance of the baseline model (lightUNet) used for both the ACN and PSN by comparing the lightUNet to the UNet architecture proposed by DeepSenseAI [9]. The lightUNet¹ architecture is modified from [9] to perform segmentation with fewer parameters. We refer to the model proposed by [9] as Base-UNet; we train both the Base-UNet and lightUNet models for 30 epochs on the 30 cm resampled AIRS dataset [8], and we report their mIOU. Table 1 shows that our lightUNet model achieves comparable performance to the Base-UNet when performing routine building segmentation. Our lightUNet model has about half the number of parameters as the Base-UNet and therefore takes less time to train.

5.2 Alignment Correction Network

V_{set} in the AIRS dataset is used to evaluate the performance of the ACN. Random translations were generated between ± 10 pixels for

¹See Appendix A for details about the convolutions.

Table 1: mIOU of Base-UNet[9] and lightUNet for routine segmentation with complete and well-aligned labels. Both models are trained on 30 cm resampled AIRS imagery.

Models	mIOU
Base-UNet	0.86
lightUNet	0.85

the xy-axis and applied to ground truth AIRS annotations, resulting in unique translation shifts for each object in an image. The introduction of noise through random translation yields a baseline mIOU of 0.55 for comparison. The shifted annotations together with the images are fed into the ACN, and the corrected annotations are compared to the true annotations to drive the learning process. We report the mIOU on V_{set} when varying amounts of T_{set} data are used for training. Random translations between ± 10 pixel are applied to all objects in V_{set} . When the ACN is trained with 800, 400 and 240 images, the corresponding mIOU on all images in V_{set} are 0.81, 0.77 and 0.67 respectively, compared to the baseline of 0.55. This suggests that the ACN performs better when more images are used but can learn with only a couple hundred training images.

Table 2: mIOU before and after ACN correction.

Translation Shift (\pm pixels)	mIOU	
	Before ACN	After ACN
0 to 5	0.63	0.81
5 to 10	0.40	0.73
10 to 15	0.26	0.46
15 to 20	0.18	0.28

Using the ACN model trained with 400 images and random translation shifts between ± 10 pixels, we evaluate the robustness of the ACN to varying levels of translation shifts. Table 2 shows mIOU before and after ACN correct, when different ranges of translations shifts are applied to V_{set} . Across all translation shifts the ACN is able to perform some realignment of annotations, even for translations (>10 pixels) which the model was never trained on.

We observe two types of alignment correction as outputs from the ACN: translations and translations plus infilling. Infilling occurs when the misaligned annotation area is less than the building area. In the translation plus infilling case, the model both shifts the annotation and fills the missing portion of the annotation. Overflow is sometimes observed upon correcting the label, resulting in the corrected annotation exceeding the building outline. Figure 4 shows examples of both types of corrections when training on 800 images. This figure demonstrates how the ACN learns over time: green outlines show predictions from the ACN and blue outlines show misaligned annotations which the ACN takes as input.

5.3 Pointer Segmentation Network

As an alternative to traditional segmentation models, we propose the Pointer Segmentation Network (PSN), a network that takes in an additional channel with points of interest and returns a single channel output with annotations. The PSN was evaluated separately

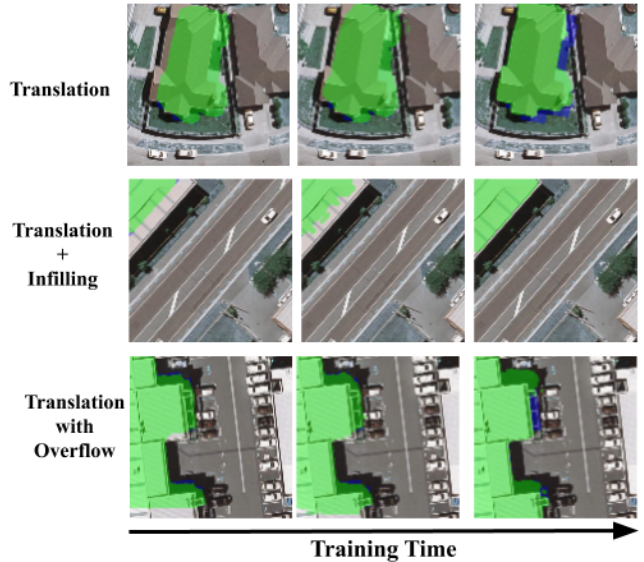


Figure 4: Types of annotation corrections performed by the ACN when trained with 800 images. Green shows corrected annotations. Blue shows misaligned annotations.

from the Alignment Correction Network (ACN); this section focuses on reporting segmentation performance on the AIRS dataset when partial – but well-aligned – labels are used. To appropriately compare the PSN with the lightUNet, we evaluate model performance using all annotations in every image of V_{set} . Here, we compare the ability of both networks to segment every building instance in the image, having learned with missing annotations. Table 3 reports the performance of the lightUNet and the PSN with varying fractions of selected annotations (α): As α decreases, performance of the PSN remains robust, indicating that the network still learns the segmentation task despite missing annotations. By specifying the points of interest, the PSN outperforms the lightUNet model.

Table 3 also presents results for two different methods of acquiring the required building points: using building centroids versus using a randomly generated point from within the corresponding annotation. By comparing the performance of the PSN using centroids with that of randomly generated points, the best annotation strategy to be used at inference can be determined. We find that the PSN performs better when centroids are used to train the model: This suggests that annotators should strive to extract points near the center of buildings to ensure better segmentation outcomes during inference. Additionally, because the extent of missing annotations may not be known *a priori* for datasets, we evaluate how the network handles heterogeneous (Het.) amounts of label completeness by sampling α from a random uniform distribution between 0 and 1. The uniform distribution ensures an equal chance for alpha to take on any value between 0 and 1. α is resampled for each image during every training epoch. Table 3 shows that the PSN remains robust at performing segmentation and works for a heterogeneous α that varies across images. Although α will likely differ across

Table 3: mIOU of PSN and lightUNet for all buildings in V_{set} images, when trained with varying α .

		mIOU
$\alpha = 1$	PSN (centroid)	0.90
	lightUNet (centroid)	0.85
$\alpha = 0.7$	PSN (centroid)	0.89
	PSN (non-centroid)	0.83
	lightUNet (centroid)	0.53
$\alpha = 0.5$	PSN (centroid)	0.87
	lightUNet (centroid)	0.18
$\alpha = \text{Het.}$	PSN (centroid)	0.87
	lightUNet (centroid)	0.71

images but remain constant for a given image at a particular time, during training we allow α to change over every training epoch for a given image, enabling our approach to be robust against images taken at different times where new construction may have occurred.

Figure 5 shows how the PSN learns – and where non-PSN type networks fail – when learning with missing annotations. The figure shows some outputs of the PSN and the lightUNet model when both are trained with $\alpha = 0.7$ and used to predict all building instances present within the image. Although both networks are trained with missing annotations, generated annotations from the PSN are more visually accurate.

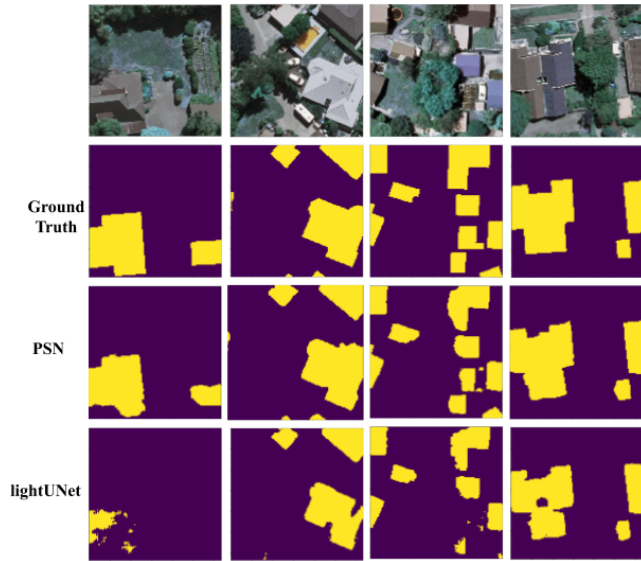


Figure 5: Annotations from PSN and lightUNet models when trained with $\alpha = 0.7$. Predictions are made for all building instances in the image and are compared to the ground truth.

5.4 Sequential Testing

The AIRS dataset is used to evaluate the sequential performance of our two-stage methodology shown in Stage 2 of Figure 2, whereby the ACN and PSN are trained and tested sequentially. Using T_{set} , we establish two training datasets for the sequential process: T1, containing misaligned labels generated from the true T_{set} ; and T2, containing ACN-corrected T1 labels. The ACN model trained with 400 training images is used to generate T2. The noise present in both training datasets is captured by the mIOU listed in Table 4. The PSN and lightUNet models are trained on T1 and T2 using $\alpha = \text{Het}$ with an identical implementation of label withholding to that described in the previous section. The trained models are used to segment V_{set} images; we compare predicted annotations to the true annotations to attain the performance metrics reported in Table 4.

Table 4 shows that, with $\alpha = \text{Het}$, the PSN performs significantly better than the lightUNet when trained on either misaligned labels (T1) or ACN-corrected labels (T2). Again, we find that with incomplete labels, regardless of alignment quality, the PSN outperforms the lightUNet. Moreover, in both training configurations, PSN mIOU performance nears that of the training dataset. As a result, we conclude that the PSN is able to predict object extents at a similar accuracy to that of the training dataset.

Table 4: Performance of the segmentation architectures. The ACN is trained with 400 images; both segmentation networks are trained with $\alpha = \text{Het}$. available annotations.

	mIOU
T1: Misaligned train dataset	0.57
PSN (trained on T1)	0.54
lightUNet (trained on T1)	0.17
T2: ACN-corrected train dataset	0.81
PSN (trained on T2)	0.79
lightUNet (trained on T2)	0.74

Figure 6 presents outputs from the PSN when trained with ACN-corrected annotations: corrected annotations from the ACN are shown in blue and predicted outputs from the PSN are shown in green. In the left half of Figure 6, we present properly corrected ACN-labels and demonstrate that the PSN is able to predict building footprints accurately when corrected annotations are accurate. The right half of the figure shows poorly corrected annotations: These corrected annotations fall on roads, grass, or across the actual building extent. In these cases, the PSN tries to predict a building footprint where there is no building. Accordingly, we conclude that improvements to the ACN can further improve PSN performance, as more accurate training labels will allow for better label prediction. Nonetheless, in the presence of misaligned annotations and partial labels, we are able to achieve better performance with our sequential architecture than with traditional segmentation approaches.

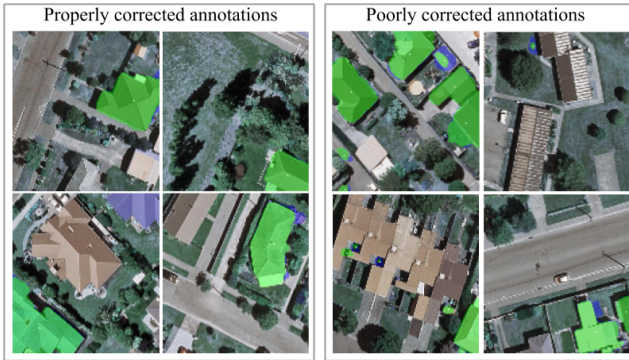


Figure 6: Sample images showing PSN performance when trained with corrected annotations. Blue footprints show ACN-corrected annotations. Green footprints show PSN-predicted annotations trained with $\alpha = Het.$ and 400 ACN-corrected labels. PSN performance is dependent on the quality of corrected annotations.

5.5 ACN Application: Realignment of OSM Annotations

In many parts of the world, ground truth is rare or nonexistent; moreover, what resources do exist often have significant accuracy issues. Despite potential shortcomings, these datasets can provide unique insight into conditions on the ground, and if their quality can be improved, they offer immense value to researchers. To confirm the performance of our realignment method on noisier images and labels, we tested the ACN on OSM building polygons in Kenya, a dataset containing considerable amounts of label misalignment. Of the 500 human-verified ground truth image labels generated for Kenya, 400 are used to train the ACN and 100 to validate. The extent of noise in OSM labels is measured by comparing the labels to the human-verified ground truth labels. mIOUs of 0.30 and 0.31 for the train and validation data respectively were recorded, when comparing OSM labels to their ground truth counterparts. OSM training labels are used to train the ACN and the trained model is ran on the 100 validation labels. A 50 % improvement in mIOU from 0.31 to 0.47 is observed on the 100 validation images. This suggests that our approach is transferable to open source labels and offers gains even with noisier images and labels, using a small dataset.

Figure 7 shows a sampling of ACN-corrected OSM annotations for images in the validation dataset: Hand-labelled annotation are shown in blue, OSM annotations are shown in red and corrected annotations are shown in green. Overall, we find that the ACN is able to correct misaligned OSM annotations both in rural and urban regions. In rural Western Kenya, where buildings tend to be smaller, the ACN shifts OSM footprints to better align with the buildings. We observe that the noisier image quality makes it more difficult for the ACN to identify extremely small buildings. In more urbanized Nairobi, the ACN also improves the alignment of OSM annotations, albeit with some failure cases.

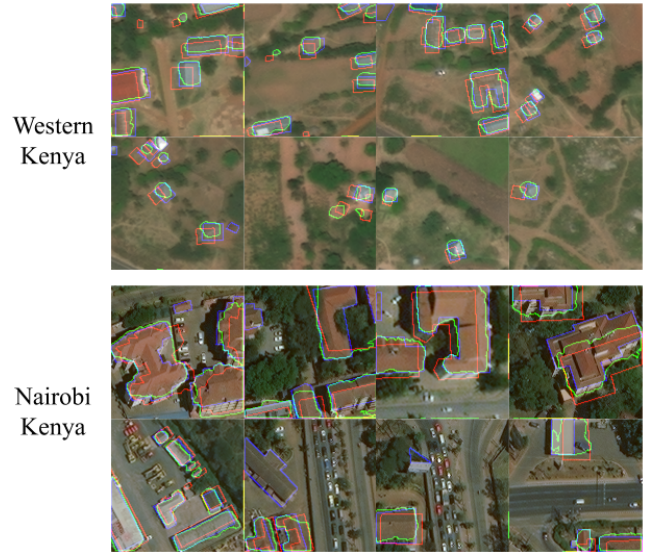


Figure 7: Hand-labelled annotations, OSM annotations and ACN-corrected annotations. The ACN is trained on 400 images from Western Kenya and Nairobi, and improves label quality despite the noisier training data.

5.6 PSN Application: Cropland Segmentation

Next, we apply the PSN to the task of cropland segmentation using Sentinel-2 imagery and a 2016 California cropping map. Knowing exact field outlines provides valuable information to farmers, planners, and governments; however, a lack of reliable, location-specific ground truth often hampers these efforts. We demonstrate the ability to accurately learn cropland extents using only a subset of fields, instead of requiring the comprehensive set of training polygons that would be necessary for traditional segmentation networks. Similar to previously described tests, we quantify the performance of the PSN in recreating these field boundaries as we select a certain fraction of the annotations, comparing results to those of the lightUNet. Table 5 presents these results.

At all fractions of available training data shown in the table, the PSN outperforms the lightUNet in segmenting croplands. After 40 training epochs, the PSN is able to predict all field boundaries for the test set across both values of α . When trained with all annotations ($\alpha = 1$), the PSN achieves a mIOU of 0.92. In contrast, the lightUNet only reaches a mIOU of 0.75 when $\alpha = 1$, and sees its performance

Table 5: mIOU for all field boundaries in test set, for varying α values.

		mIOU
$\alpha = 1$	PSN	0.92
	lightUNet	0.75
$\alpha = 0.75$	PSN	0.91
	lightUNet	0.69

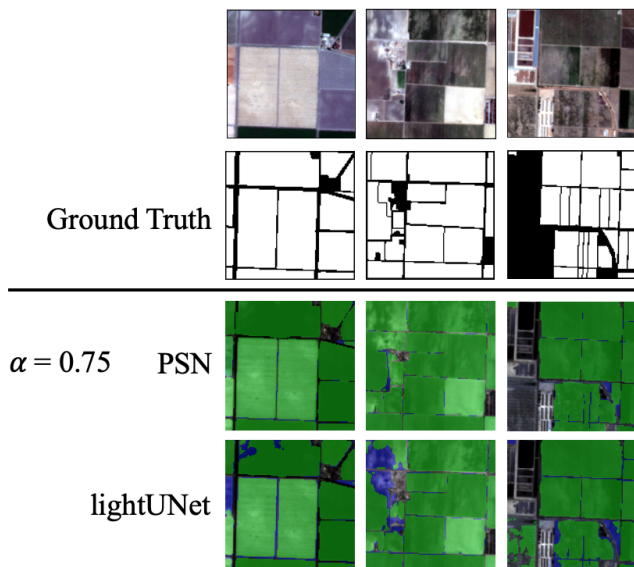


Figure 8: Sample images and ground truth labels showing cropland extent in California; also shown in green are PSN and lightUNet predicted footprints $\alpha = 0.75$, overlaid on true cropland polygons, shown in blue. PSN predictions remain highly accurate. Comparatively, the lightUNet predicts only a portion of the crop extents correctly

significantly diminish as field boundaries are withheld. Figure 8 shows the PSN- and lightUNet - recreated field polygons when the models are trained with $\alpha = 0.75$ and are asked to predict all polygons within an image. The true cropland polygons are shown in blue while the predicted polygons are shown in green; all examples shown come from the test set.

These results demonstrate the viability of the PSN in delineating field boundaries and the preferability of our method over a baseline alternative, when the acquisition of field boundaries is expensive. In locations with low data availability and smaller, non-uniform field boundaries, the PSN provides a reliable method for determining cropped area polygons.

6 CONCLUSION

As the demand for extracting information from satellite imagery increases, the value of reliable, transferable object segmentation methodologies – especially ones that compensate for noise and inaccuracies in training data – increases in parallel. In this paper, we present a novel and generalizable two-stage segmentation approach that address common issues in applying deep learning approaches to remotely-sensed imagery. First, we present the Alignment Correction Network (ACN), a model which learns to correct misaligned object annotations. We test the ACN on a set of alignment errors, including i) misalignment of the AIRS dataset, ii) existing and substantial misalignment errors within the OSM Kenyan building footprint dataset. Overall, we find that the ACN significantly improves annotation alignment accuracy.

We also introduce the Pointer Segmentation Network (PSN), a model which reliably predicts an object’s extent using only a point

from the object’s interior. The value of the PSN lies in learning to segment objects within an image despite incomplete or missing annotations, an issue which both hinders traditional segmentation efforts and is common in many ground-truth datasets. We train and test the PSN on the AIRS dataset and find that the model can accurately predict building extent regardless of the fraction of available annotations present or where the training point resides within the object. We also evaluate the performance of the PSN for cropland segmentation using Sentinel imagery and a 2016 California cropland map as inputs, demonstrating that the model can reliably learn cropland polygons regardless of the fraction of available annotations. Overall, for all testing configurations – those which vary the fraction of available training annotations and those which change the location of where the training point lies – and for both object segmentation applications presented – building footprint and cropland extent predictions – the PSN outperforms a baseline segmentation model.

Lastly, we sequentially link the ACN and PSN to demonstrate the ability of the combined networks to accurately segment objects having learnt from misaligned and incomplete training data. Taken together, we envision our proposed networks providing value to the community of researchers and scientists looking to extract information from widely-available satellite imagery and unreliable ground-truth datasets.

REFERENCES

- [1] Rasha Alshehhi, Prashanth Reddy Marpu, Wei Lee Woon, and Mauro Dalla Mura. 2017. Simultaneous extraction of roads and buildings in remote sensing imagery with convolutional neural networks. *ISPRS Journal of Photogrammetry and Remote Sensing* 130 (2017), 139–149.
- [2] Nicolas Audebert, Bertrand Le Saux, and Sebastien Lefevre. 2017. Joint Learning From Earth Observation and OpenStreetMap Data to Get Faster Better Semantic Maps. *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops* (2017), 67–75.
- [3] V. Badrinarayanan, A. Kendall, and R. Cipolla. 2017. SegNET: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation. *International Conference on Medical image computing and computer-assisted intervention*. 39, 12 (2017), 2482–2495.
- [4] Anahid Basiri, Mike Jackson, Pouria Amirian, Amir Pourabdollah, Monika Sester, Adam Winstanley, Terry Moore, and Lijuan Zhang. 2016. Quality assessment of OpenStreetMap data using trajectory mining. *Geo-spatial information science* 19 (2016), 56–68.
- [5] Amy Bearman, Vittorio Ferrari Olga Russakovsky, and Li Fei-Fei. 2016. What’s the Point: Semantic Segmentation with Point Supervision. *European conference on computer vision* (2016).
- [6] Gabriel Cadamuro, Aggrey Muhebwa, and Jay Taneja. 2019. Street smarts: measuring intercity road quality using deep learning on satellite imagery. In *Proceedings of the 2nd ACM SIGCAS Conference on Computing and Sustainable Societies (COMPASS 2019)*. 145–154.
- [7] Silvana Philippi Camboim, João Vitor Meza Bravo, and Claudia Robbi Sluter. 2015. An investigation into the completeness of, and the updates to, OpenStreetMap data in a heterogeneous area in Brazil. *ISPRS International Journal of Geo-Information* 4, 3 (2015), 1366–1388.
- [8] Qi Chen, Lei Wang, Yifan Wu, Guangming Wu, Zhiling Guo, and Steven L. Waslander. 2018. Aerial imagery for roof segmentation: A large-scale dataset towards automatic mapping of buildings. *arXiv preprint arXiv:1807.09532* (2018).
- [9] deepsense.ai. 2020. Deep learning for satellite imagery via image segmentation. Retrieved March 6, 2020 from <https://deepsense.ai/deep-learning-for-satellite-imagery-via-image-segmentation/>.
- [10] Christopher N.H. Doll, Jan-Peter Muller, and Jeremy G. Morely. 2005. Mapping regional economic activity from night-time light satellite imagery. *Ecological Economics* (2005).
- [11] Zhenrong Du, Jianyu Yang, Cong Ou, and Tingting Zhang. 2019. Smallholder crop area mapped with a semantic segmentation deep learning method. *Remote Sensing* 11, 7 (2019), 888.
- [12] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. [n.d.]. The PASCAL Visual Object Classes

- Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>.
- [13] Humanitarian Data Exchange. 2020. HOTOSM Kenya Buildings (OpenStreetMap Export). Retrieved February 27, 2020 from https://data.humdata.org/dataset/hotosm_ken_buildings.
- [14] Ferran Gascon, Catherine Bouzinac, Olivier Thépaut, Mathieu Jung, Benjamin Francesconi, Jérôme Louis, Vincent Lonjou, Bruno Lafrance, Stéphanie Massera, Angélique Gaudel-Vacaresse, Florie Languille, Bahjat Alhammoud, François Viallefont, Bringfried Pflug, Jakub Bieniarz, Sâbastien Clerc, Laëticia Pessiot, Thierry Trâlmas, Enrico Cadau, Roberto De Bonis, Claudia Isola, Philippe Martimort, and Valérie Fernandez. 2017. Copernicus Sentinel-2A Calibration and Products Validation Status. *Remote Sensing* 9, 6 (2017). <https://doi.org/10.3390/rs9060584>
- [15] Nicolas Girard, Guillaume Charpiat, and Yuliya Tarabalka. 2019. Noisy Supervision for Correcting Misaligned Cadaster Maps Without Perfect Ground Truth Data. In *IGARSS 2019-2019 IEEE International Geoscience and Remote Sensing Symposium*. IEEE, 10103–10106.
- [16] Ross Girshick. 2015. Fast R-CNN. In *The IEEE International Conference on Computer Vision (ICCV)*.
- [17] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. 2014. Rich feature hierarchies for accurate object detection and semantic segmentation.. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- [18] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*.
- [19] Bikash Joshi, Hayk Baluyan, Amer Al. Hinaï, and Wei Lee Woon. 2014. Automatic Rooftop Detection Using a Two-Stage Classification. In *Proceedings of the 2014 UKSim-AMSS 16th International Conference on Computer Modelling and Simulation (UKSim-AMSS)*. IEEE Computer Society, USA, 286–291. <https://doi.org/10.1109/UKSim.2014.89>
- [20] Pascal Kaiser, Jan Dirk Wegner, Aurélien Lucchi, Martin Jaggi, Thomas Hofmann, and Konrad Schindler. 2017. Learning Aerial Image Segmentation from Online Maps. *IEEE Transactions on Geoscience and Remote Sensing* 55, 11 (2017), 6054–6068.
- [21] Natalia Kussul, Mykola Lavreniuk, Sergii Skakun, and Andrii Shelestov. 2017. Deep Learning Classification of Land Cover and Crop Types Using Remote Sensing Data. *IEEE Geoscience and Remote Sensing Letters* 14, 5 (2017), 778–782.
- [22] Zeming Li, Chao Peng, Gang Yu, Xiangyu Zhang, and Yangdong Deng and Jian Sun. 2017. Light-Head R-CNN: In Defence of Two-Stage Object Detector. In *arXiv preprint arXiv:1711.07264*.
- [23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. *European conference on computer vision* (2014).
- [24] Jia-Qi Liu, Zhili Wang, and Kangxin Cheng. 2019. An improved algorithm for semantic segmentation of remote sensing images based on DeepLabv3+. In *ICCI'19*.
- [25] Jonathan Long, Evan Shelhamer, and Trevor Darrell. 2015. Fully Convolutional Networks for Semantic Segmentation. *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops* (2015), 3431–3440.
- [26] Ron Mahabir, Anthony Stefanidis, Arie Croitoru, Andrew T Crooks, and Peggy Agouris. 2017. Authoritative and volunteered geographical information in a developing country: A comparative case study of road datasets in Nairobi, Kenya. *ISPRS International Journal of Geo-Information* 6, 1 (2017), 24.
- [27] D. Martin, C. Fowlkes, D. Tal, and J. Malik. 2001. A Database of Human Segmented Natural Images and its Application to Evaluating Segmentation Algorithms and Measuring Ecological Statistics. In *Proc. 8th Int'l Conf. Computer Vision*, Vol. 2. 416–423.
- [28] Volodymyr Mnih and Geoffrey E. Hinton. 2012. Learning to label aerial images from noisy data. In *Proceedings of the 29th International conference on machine learning*.
- [29] Michael Xie W. Matthew Davis David B. Lobell Stefano Ermon Neal Jean, Marshall Burke. 2016. Combining Satellite Imagery and Machine Learning to Predict Poverty. *Journal of Science* 353, 6301 (2016), 790–794.
- [30] Anthony Perez, Swetava Ganguli, Stefano Ermon, George Azzari, Marshall Burke, and David B. Lobell. 2019. Semi-Supervised Multitask Learning on Multispectral Satellite Images Using Wasserstein Generative Adversarial Networks (GANs) for Predicting Poverty. *CoRR* abs/1902.11110 (2019). arXiv:1902.11110 <http://arxiv.org/abs/1902.11110>
- [31] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*.
- [32] California Department Of Water Resources. 2020. 2016 California Statewide Agricultural Land Use Map. Retrieved February 27, 2020 from <https://gis.water.ca.gov/app/CADWRLandUseViewer/>.
- [33] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation.. In *International Conference on Medical image computing and computer-assisted intervention*.
- [34] Rose M Rustowicz, Robin Cheong, Lijing Wang, Stefano Ermon, Marshall Burke, and David Lobel. 2019. Semantic Segmentation of Crop Type in Africa: A Novel Dataset and Analysis of Deep Learning Methods. *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops* (2019), 75–82.
- [35] Safyan. 2015. Overview of the planet labs constellation of earth imaging satellites. (2015).
- [36] Lindy-Anne Siebritz and George Sithole. 2014. Assessing the quality of OpenStreetMap data in South Africa in reference to national mapping standards. In *Proceedings of the Second AfricaGEO Conference, Cape Town, South Africa*. 1–3.
- [37] Benjamin Swan, Melanie Laverdiere, and H. Lexie Yang. 2018. How Good is Good Enough? Quantifying the Effects of Training Set Quality. In *Proceedings of the 2nd ACM SIGSPATIAL International Workshop on AI for Geographic Knowledge Discovery* (Seattle, WA, USA) (*GeoAI'18*). Association for Computing Machinery, New York, NY, USA, 47–51. <https://doi.org/10.1145/3281548.3281557>
- [38] John E Vargas-Munoz, Sylvain Lobry, Alexandre X. Falcao, and Devis Tuia. 2019. Correcting rural building annotations in OpenStreetMap using convolutional neural networks. *ISPRS Journal of Photogrammetry and Remote Sensing* 147 (2019), 283–293.
- [39] Alyssa Wright. 2020. Map Completeness Estimation and Experimental Analytics for Health. Retrieved March 6, 2020 from <https://www.hotosm.org/updates/experimenting-with-analytics-for-health/>.
- [40] G. Wu and Z. Guo. 2019. GeoSeg: A computer Vision Package for Automatic Building Segmentation and Outline Extraction. *IGARSS 2019-2019 IEEE International Geoscience and Remote Sensing Symposium* (2019).
- [41] Yongyang Xu, Liang Wu, Zhong Xie, and Zhanlong Chen. 2018. Building extraction in very high resolution remote sensing imagery using deep learning and guided filters. *Remote Sensing* 10, 1 (2018), 144.
- [42] Yongyang Xu, Zhong Xie, Yaxing Feng, and Zhanlong Chen. 2018. Road Extraction from High-Resolution Remote Sensing Imagery Using Deep Learning. *Remote Sensing* 10, 9 (2018), 1461.

A ARCHITECTURE

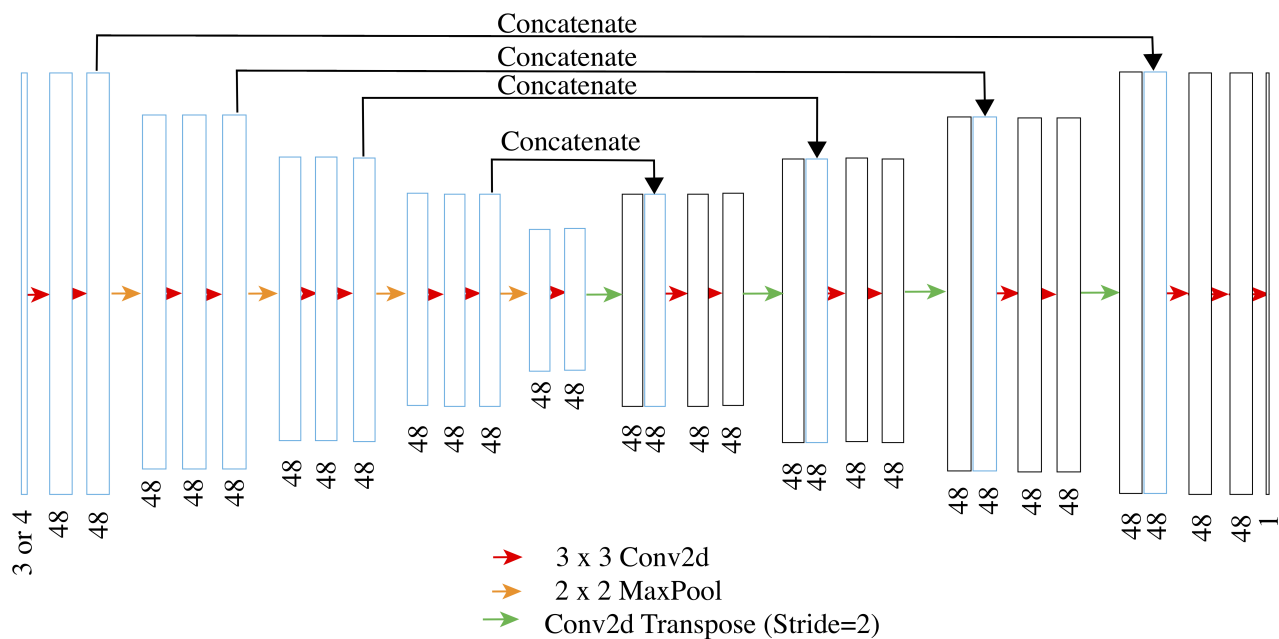


Figure 9: Architecture used for both the Alignment Correction Network (ACN) and the Pointer Segmentation Network (PSN). Four input channels are used for both ACN and PSN, while three are used for the lightUNet. This network is modified from [9] by reducing the number of filters to 48 and maintaining the same filter size through out the network. In addition, the network uses dropout in addition to batch normalization after every epoch.