# A Localisation-Segmentation Approach for Multi-label Annotation of Lumbar Vertebrae using Deep Nets

Anjany Sekuboyina[1,2,✉], Alexander Valentinitsch[2], Jan S. Kirschke[2], and Bjoern H. Menze[1]

[1]Technische Universität München, Munich, Germany
[2]Klinikum rechts der Isar, Munich, Germany
anjany.sekuboyina@tum.de

**Abstract.** Multi-class segmentation of vertebrae is a non-trivial task mainly due to the high correlation in the appearance of adjacent vertebrae. Hence, such a task calls for the consideration of both global and local context. Based on this motivation, we propose a two-staged approach that, given a computed tomography dataset of the spine, segments the five lumbar vertebrae and simultaneously labels them. The first stage employs a multi-layered perceptron performing non-linear regression for locating the lumbar region using the global context. The second stage, comprised of a fully-convolutional deep network, exploits the local context in the localised lumbar region to segment and label the lumbar vertebrae in one go. Aided with practical data augmentation for training, our approach is highly generalisable, capable of successfully segmenting both healthy and abnormal vertebrae (fractured and scoliotic spines). We consistently achieve an average Dice coefficient of over 90% on a publicly available dataset of the xVertSeg segmentation challenge of MICCAI'16. This is particularly noteworthy because the xVertSeg dataset is beset with severe deformities in the form of vertebral fractures and scoliosis.

## 1  Introduction

The identification, segmentation, and quantification of structures visible in medical images is a crucial component in the processing of medical image data. In the context of spinal images, segmentation of spine has an immediate diagnostic importance in clinical decisions around fracture detection and inter-vertebral disc pathology. Segmented spines are also used in the bio-mechanical modelling of the spine for load analysis and fracture prediction. Therefore, an automated approach attempting to segment the spine should posses two key features: (1) Highly generalisable in terms of the fields-of-view (FOV) and scanner calibrations, in addition to variability in the spine's curvature, BMD (bone mineral density) distribution, and micro-architecture and (2) Capable of segmenting images from a clinical population that consists of abnormalities such as vertebral fractures, scoliotic, and kyphotic spines.

A typical analysis pipeline for spinal images consists of three stages: spine localisation, vertebrae detection, and spine segmentation. The first two steps of localisation and detection are accomplished by basic routines such as shape matching (using generalised Hough transform [1]) and spine-curve detection (using circle detection in axial slices [2]). This is followed by a segmentation stage which may tackled using statistical mean shape models or atlases, followed by an optimisation routine that adapts the fitted model to account for local variations [3,4]. Such a pipeline has proven to be highly effective in most of the cases. However, there is a limit to the amount of generalisability such model/shape-based approach can offer in clinical cases. Its limit is determined by the robustness of the chosen model and the amount of relaxation it can withstand during the optimisation routine post fit. It is obvious that such models cannot generalise to a fractured vertebra or a deformed spine. In such cases, learning-based approaches offer respite, provided that the data that the approach can learn from is rich and diverse enough. For example, [5] and [6] solve the problem of vertebra detection in arbitrary FOVs using random forests and multi-layer regressors respectively. Chen et al. [7] make use of the omni-present convolutional neural networks to detect vertebrae using an altered cost formulation that takes into account the sequence of the vertebrae. However, there is no *end-to-end* approach that handles every problem in the analysis pipeline (localisation, detection, and segmentation) in one go, that is, takes a 3D spine scan as input and generates an annotated and segmented spine volume.

In this work, we propose an approach that segments and simultaneously labels the the lumbar vertebrae using deep neural networks. Given a CT scan volume of an arbitrary FOV, our approach performs a multi-class segmentation over five classes corresponding to the five lumbar vertebrae (L1 to L5) and a background class. This is done in a two-staged approach: (1) Localise the lumbar region, and (2) Segment the localised lumbar vertebrae in to their respective classes. Both the stages are elaborated in detail in section 2. Figure 1 gives a schematic overview of our approach. We use the dataset released as part of the xVertSeg challenge in MICCAI 2016 to test the performance of our approach. We are the first to attempt this challenge, and achieve a mean Dice coefficient upwards of 90% on both the training and the test set. Section 3 contains the implementation and experimental details.

## 2   Methodology

### 2.1   Lumbar localisation

When compared to typical binary segmentation, multi-class segmentation is inherently difficult due to a more complex representation that is to be learnt. Moreover, the appearance of adjacent vertebrae are highly correlated. Thus, instead of directly attempting the segmentation problem on the entire scan, we choose to restrict our attention to a restricted region of interest - the lumbar region.
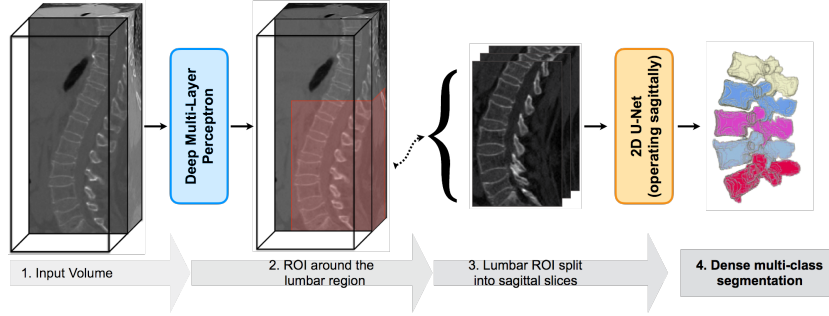
Fig. 1: A schematic diagram giving an overview of our approach.

**Non-linear Regression using deep neural network** We pose the localisation of the lumbar region as a regression problem, and employ a five-layered perceptron with ReLU (rectified linear unit) activation as a regressor. It is trained on contextual, intensity-dependant features that encode long-range spatial information, as in [5], and regresses on the location of the six planes that define a bounding box. An $n$ lenght feature, $\mathbf{f}_{i,j,k}$ can be constructed at the voxel location of $(i, j, k)$ as below:

$$\mathbf{f}_{i,j,k} = \{f^1_{i,j,k}, f^2_{i,j,k}, \ldots, f^n_{i,j,k}\}, \tag{1}$$

where $f^p_{i,j,k}, \forall p \in 1, 2, \ldots n$ is the mean intensity of the 3D image region lying inside a cuboid that is centered at a certain offset from the voxel at $(i, j, k)$. The cuboid's offset and the dimension are generated randomly for construction of a feature. Given these features, each of them corresponding to a voxel location, the regressor should predict the region-of-interest, or a bounding box around the lumbar region.

In a simple set-up, a bounding box can be defined by six planes: $x_{min}$, the smallest $x$-coordinate, $x_{max}$, the largest $x$-coordinate, and their $y$ and $z$ equivalents. These are refereed to as the *bounding planes*. Given the contextual information through the feature $\mathbf{f}_{i,j,k}$, a six-length vector encoding the voxel's location *relative* the bounding planes is learnt, as below:

$$\mathbf{y}_{i,j,k} = \{i - x_{min}, i - x_{max}, j - y_{min}, j - y_{max}, k - z_{min}, k - z_{max}\} \tag{2}$$

**Estimation of the lumbar bounding box** Each pass through the regressor using a feature corresponding to a certain voxel predicts the locations of the bounding planes with respect to that voxel. In order to speed-up the feature generation procedure without loss of useful information, only the *significant* voxels are considered for feature extraction. For this purpose, the voxels from the response of a Canny's edge detector are used for feature extraction. Thus, every significant voxel votes for a prospective bounding box of which the most representative bounding box is chosen. Figure 2 shows a few examples of the localised lumbar regions.
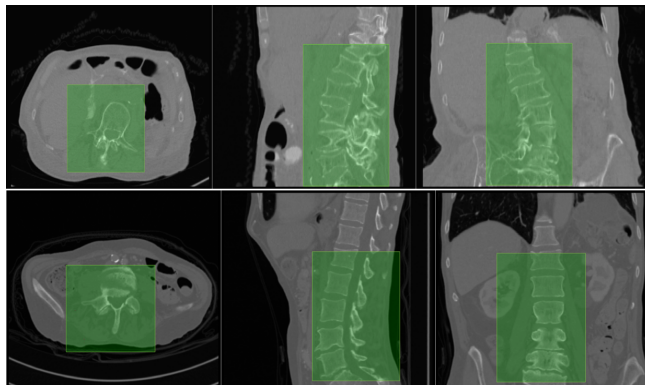
Fig. 2: Lumbar Localisation: The axial, sagittal, and coronal views of the bounding box localising the lumbar section. (Row 1) Case 18 containing severe abnormalities such as multiple fractures and scoliosis is localised perfectly. (Row 2) Case 22 shows a mild *under*-localisation, not localising a top region of L1.

### 2.2   Multi-class Segmentation

Once the lumbar region is successfully identified, the FOV is restricted, enabling a human to effectively identify the vertebrae, based on certain key points such as the sacrum, number of vertebrae in the FOV etc. We make use of a deep convolutional network to learn such key points on its own in order to segment and annotate the lumbar vertebrae.

**Fully-convolutional network for multi-class segmentation** This is the segmentation of the vertebrae is carried out by a fully convolutional network (FCN). We rely on the 2D U-net [10], but implement an architecture that one level deeper, i.e., six more convolutional layers, three each in the contracting and expanding path, joined by one additional downsampling and up-convolutional layers, and works on sagittal slices from the localised lumbar region. The motivation for a deeper network is to adapt it towards multi-label classification of vertebrae by increasing the receptive field of the coarsest level. The receptive field of our FCN ($\approx 270{\times}270$ pixel$^2$ or $27{\times}27mm^2$) covers at least two vertebrae when working at isotropic $1mm$ resolution. Such a receptive field will force the network to learn the sequence of vertebrae in pairs, L1-L2, L2-L3 etc., so that the sequence of the annotations is always in order.

**Pre-training** It is a common practice to pre-train a network for one purpose and use it as an initialisation for another network attempting a related yet more-complex task. For example, Long et al. [11] use the state-of-art recognition networks (VGG-16 etc.) as initialisation for the task of segmentation. On a similar footing, a network trained for binary segmentation of lumbar spine (spine vs. background) is employed as an initialisation for the multi-class segmentation. This alleviates the shortcoming of the limited data at our disposal to train a very deep network for the relatively complex task of multi-class segmentation.

**ROI Augmentation** Another key concept in our segmentation routine is the ROI augmentation step. However, the localisations are not *uniform* as shown in figure 2. There could either be non-lumbar vertebrae showing up in the sagittal slices (usually T11 and T12 in our experiments), or part of the lumbar region could be missing (usually L1 in our experiments). The high correlation in the appearance of the vertebrae makes this problem detrimental. Therefore, in addition to augmenting the sagittal slices using elastic and rigid transformations, we also augment based on varying bounding boxes sizes. Let $h \times w \times d$ be the dimension of the lumbar bounding box obtained from the localisation stage. We augment the sagittal slices from bounding box by randomly choosing a $\delta \in \{5, 10, 15, 20, 25\}$ so that the sagittal slice dimensions vary between $h \times w$ and $h - 2\delta \times w$. This makes the segmentation network robust to improper lumbar localisations.

**Extracting the final segmentations** Once all the sagittal slices are segmented, the final segmented volume is coronally corrected by closing the holes in every label of a coronal slice using morphological closing operation. Finally, a 3D connected-component analysis is performed on each label to discard the smaller connected components. This cleans-up a few stray segmentations in the final volume.

## 3   Experiments and Discussion

We make use of a publicly available dataset for evaluating the performance of the lumbar localisation and the multi-class segmentation stages of our approach.

**xVertSeg Dataset** The xVertSeg dataset, released as part of the xVertSeg challenge [12] in MICCAI 2016, consists of fifteen train CT volumes with ground truth segmentations of the lumbar vertebrae (into five classes, L1-L5 ) and ten test CT volumes. The participants do not have access to the ground truth segmentations of the test set. The data is very rich in terms of varying FOVs, spine curvatures, and vertebral deformities.
*Ground truth for localisation:* The first and the last slices in the three directions (sagittal, coronal, and axial) consisting of a label were considered to be the bounding planes. A tolerance of 15 slices was added on all sides of the bounding box to prevent a tight cropping of the lumbar region. This expanded bounding box was used as the ground truth for training the localisation network.
*Ground truth segmentations for test cases:* As the challenge organisers did not make the performance metrics of our approach available yet, we opted for an in-house ground truth generation. The near-perfect segmentation from our approach was given to two clinical experts (Rater-1 and Rater-2) for correction. Rater-1 was tasked to correct the entire volume, while Rater-2 was tasked to pick a random subset of sagittal, coronal, and axial slices from a volume and segment them entirely.

**Lumbar localisation** Inspired from [6], a five-layered neural network is cast as a regressor to map the features ($\mathbf{f}_{i,j,k}$) to the offsets of the bounding planes ($\mathbf{y}_{i,j,k}$). The input layer has $n$ (=500) neurons, followed by four hidden layers

with 350, 250, 150, and 50 neurons respectively. The output layer has six neurons corresponding to the offsets of the six bounding planes. All the neurons are ReLUs. The network was implemented in the Caffe [8] framework. A squared-error loss was optimised using stochastic gradient descent. The available data was augmented with rigid and elastic transformations. The network was trained for 1000 epochs over a few hours with a learning rate of 1e-3 and a momentum of 0.9. The most representative bounding box is chosen using kernel density estimation, aided by Botev bandwidth [9] selection.

To measure the performance of localisation, a measure of *sensitivity* (or true positive rate) was used, as defined: $S = 1 - \frac{|\mathcal{G} \cap \mathcal{B}^c|}{|\mathcal{G}|}$, where $\mathcal{G}$ is the of set voxels in the ground truth segmentation, and $\mathcal{B}$ is the set of voxels within the bounding box. We use the ground truth segmentation for Rater-1 for this purpose. The sensitivity measures on the test set are shown in table 1, with a few cases shown visually in figure 2. We obtain a near perfect localisation of 1.0 in all cases except one (Case025). In order to completely cover the lumbar region, a tolerance of 15 voxels is added to the bounding boxes on all sides before considering the localisation for for the next stage.

*The curious case of Case025:* Case025 of the test set is peculiar due to the presence of the entire sacrum (S1, S2, and S3) within the FOV. It is the only such case in the train and the test data. The network thus furnishes an imperfect localisation from L2 to S1. Such a behaviour can be easily rectified with additional representative training data.

| Case16 | Case17 | Case18 | Case19 | Case20 | Case21 | Case22 | Case23 | Case24 | Case25 | **Mean** |
|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|----------|
| 0.98 | 1.0 | 0.99 | 0.99 | 1.0 | 0.99 | 0.98 | 1.0 | 0.98 | 0.94 | **0.98** |

Table 1: Sensitivities of the localisation algorithms on the ten test cases. The localisation is near perfect ($\sim$1.0) for all the case except Case25.

**Multi-class segmentation** The segmentation network is implemented in the Caffe framework. A cross-entropy loss is optimised using an Adam solver with an initial learning rate of 1e-4. The binary segmentation network is run for 3000 epochs and the multi-class segmentation (with pre-training) net was run for 2000 epochs. The segmented bounding boxes are reinstated into the actual volumes to obtain the full-resolution segmentations. We report the Dice coefficient for each of the five vertebrae and for the entire lumbar region in table 2. The evaluation is carried out based on the available ground truth segmentations of the train set and those from both Rater-1 and Rater-2 in case of the test set. We also observe a mean Dice score of $\sim$92%. Since our segmentation is the starting point for Rater-1, a bias in the corresponding performance scores can be observed, with a mean Dice score of 94%. Figure 3 shows the spread of the Dice coefficients across vertebrae and among the datasets. Observe that the vertebrae in the middle (L3 and L4) are segmented well compared to the peripheral vertebrae (L1 and L5). This is expected since the uncertainty that the net has to overcome for deciding between L1 & T12 or L5 & S1 is higher compared to deciding between L2 & L3 or L3 & L4 owing to the large receptive field etc.
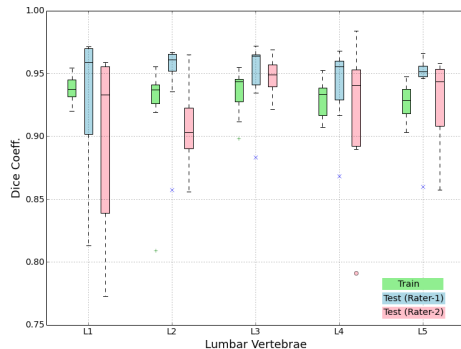
Fig. 3: Plot showing the distribution of the Dice coefficients across the vertebrae, comparing the performance among the datasets.

**Discussion** In general, both the stages in our pipeline work remarkably well as per the quantitative results in tables 1 and 2. We obtain a near perfect localisation of 1.0 for almost every case, and a mean Dice score of 92%. In addition to this, the prime motivation of our approach is to successfully segment the deformed spines where the model-based approaches fail. This can be observed visually in figure 4. Four test cases as shown highlighting the highly deformed spine and vertebrae. Observe that our algorithm successfully segments these cases in spite of the severe deformations.

| Data | L1 | L2 | L3 | L4 | L5 | Lumbar |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| Train | 92.6±6.7 | 92.7±3.0 | 93.5±1.5 | 92.9±1.3 | 92.7±1.2 | **92.7±2.3** |
| Test (Rater-1) | 92.7±5.5 | 94.9±3.2 | 95.1±2.5 | 94.2±2.9 | 94.1±2.9 | **94.3±2.8** |
| Test (Rater-2) | 89.7±6.7 | 90.8±3.2 | 94.8±1.5 | 91.8±5.6 | 92.8±3.3 | **92.0±2.3** |

Table 2: Dice coefficients (in %) of our approach on the xVertSeg dataset. Observe a consistent performance of above 90% in Dice scores. The distribution of the scores is visualised in figure 3

## 4 Conclusions

Deep-learning based algorithms are a way forward if generalisability is to be achieved. However, usage of such algorithms for dense segmentation is still in its incipient stage. The task of segmentation becomes more challenging when it involves multiple-classes over similar-looking vertebrae. We propose a two-staged approach with deep networks for localisation of the lumbar region and segmenting it into multiple classes. We are the first to present results on the xVertSeg dataset, with an incredible performance achieving a mean Dice coefficient of above 90%. We also highlight the ability of our approach to handle severe deformities in the spine, which prior approaches would struggle with. We believe that our approach can form a basis for handling more complicated tasks of multi-class segmentation of the entire spine.
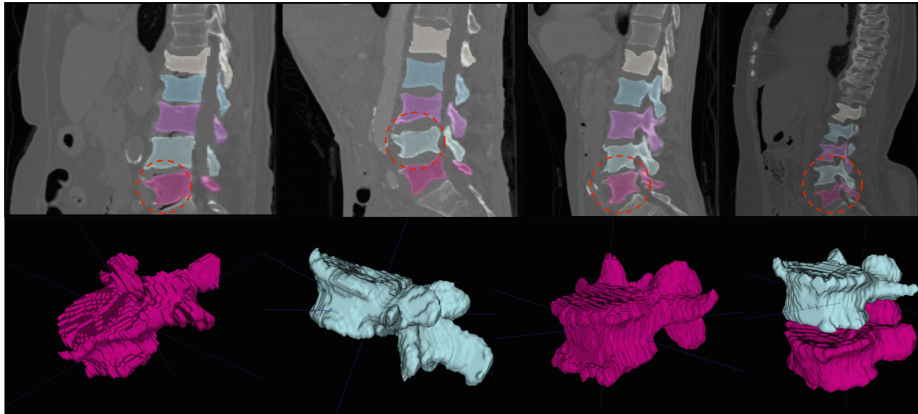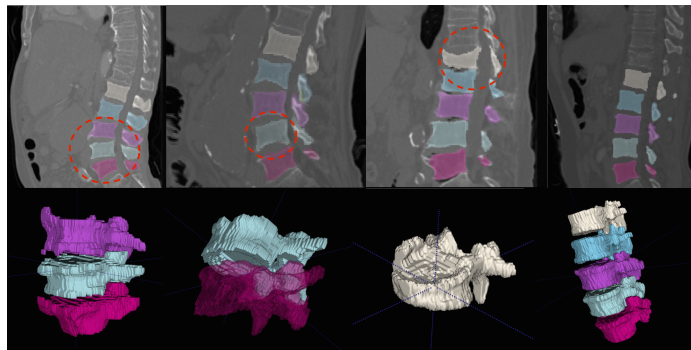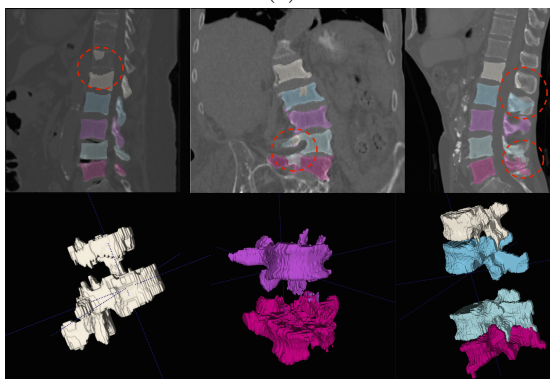
Fig. 4: Multi-class Segmentation: (Row 1) Four sagittal slices where the deformed vertebrae are highlighted, with (Row 2) the 3D rendering of the deformed vertebrae for better visualisation. More visualisations on deformed vertebrae will be made available as supplementary material.

# References

1. Klinder, T., et al.: Automated model-based vertebra detection, identification, and segmentation in CT images. In: Medical Image Analysis, 13(3):471–482 (2009)
2. Forsberg, D.: Atlas-Based Segmentation of the Thoracic and Lumbar Vertebrae. In: Proc MICCAI-CSI 2014, Boston, USA: Springer (2014)
3. Kadoury, S., Labelle, H., and Paragios, N.: Spine segmentation in medical images using manifold embeddings and higher-order MRFs. In: IEEE TMI, pages 1227–1238 (2013)
4. Korez, R., et al.: Interpolation-Based Shape Constrained Deformable Model Approach for Segmentation of Vertebrae from CT Spine Images. In: Proc MICCAI-CSI 2014, Boston, USA: Springer (2014)
5. Glocker, B., et al.: Automatic localization and identification of vertebrae in arbitrary field-of-view CT scans. In: MICCAI (2012)
6. Suzani, A., et al.: Deep learning for automatic localization, identification, and segmentation of vertebral bodies in volumetric MR images. In: SPIE Medical Imaging, International Society for Optics and Photonics, 2015, pp. 941 514–941 514 (2015)
7. Chen, H., et al.: Automatic localization and identification of vertebrae in Spine CT via a joint learning model with deep neural networks. In: MICCAI 2015. LNCS, vol. 9349, pp. 515–522. Springer, Heidelberg (2015)
8. Jia, Y., et al.: Caffe: Convolutional architecture for fast feature embedding (2014), arXiv:1408.5093 [cs.CV]
9. Botev, Z.I., Grotowski, J.F., Kroese, D.P.: Kernel density estimation via diffusion. In: Annals of Statistics. 38 (5): 2916–2957 (2010)
10. Ronneberger, O., Fischer, P., and Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: MICCAI, pp. 234–241, Springer (2015)
11. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation (2014), arXiv:1411.4038 [cs.CV]
12. The xVertSeg Challenge, http://lit.fe.uni-lj.si/xVertSeg/

(a)



(b)

Fig. 5: **(a)** Segmented lumbar region (Row 1) with a 3D rendering of the highlighted deformities (Row 2) such as osteophytes and fractures that are successfully segmented. **(b)** Slightly aberrant cases: (Left) The fracture in L4-L5 is perfectly segmented. However, notice an over-segmentation in L1. (Centre) A successful segmentation of a severely scoliotic spine and deformed vertebrae. Notice some non-homogeneity in segmented labels. Also observe in 3D, a well-captured crush in L3 and an unsegmented region in L5. (Right) The anterior regions of the vertebrae are successfully segmented. However, posterior regions of L1, L2, L4, and L5 are not fully segmented.

## 5    Appendix

### 5.1    Additional Results

We present more results of multi-class segmentation on the test set of xVertSeg (figure 5) in addition to the results in figure 4, thereby emphasising the robustness of our approach. We also present a few aberrant segmentations analysing which could further improve our approach.

### 5.2    The Outlier (Case 25 of the xVertSeg Dataset)

As mentioned in section 3 (Lumbar localisation) of the main article, the localisation in Case 25 occurs with a sensitivity of 0.94 (figure 6(c), red outline) as it is the only example in the train and test data that consists of three sacral bones (S1, S2, and S3) within the field-of-view (figure 6(a)). When the scan, with S3 manually cropped off, is used as input (figure 6(b)), the lumbar localisation is perfect (sensitivity of 1.0), as shown by the green outline in figure 6(c). It is clear that the improper localisation is a consequence of working with limited data, and can easily be averted by increasing the variability in the training dataset.
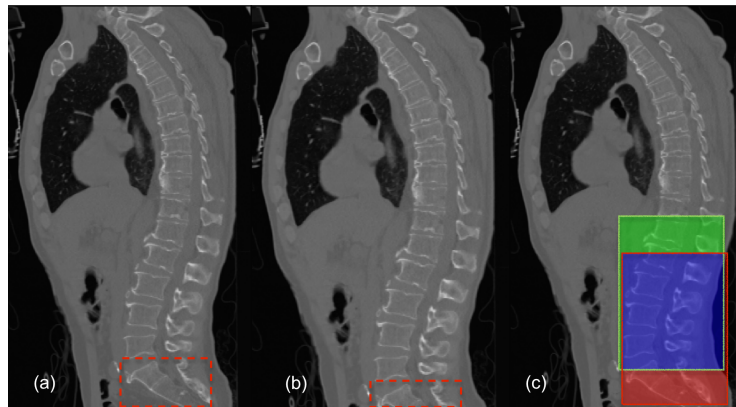


Fig. 6: **(a)** Actual Case 25. **(b)** Case 25 with S3 cropped off. **(c)** Red and green regions correspond to lumbar localisation using (a) and (b) respectively. Blue is the overlap in the regions. Both the bounding boxes are overlaid on the original image.