# Fast Multiple Landmark Localisation Using a Patch-based Iterative Network

Yuanwei Li<sup>1</sup>, Amir Alansary<sup>1</sup>, Juan J. Cerrolaza<sup>1</sup>, Bishesh Khanal<sup>2</sup>, Matthew Sinclair<sup>1</sup>, Jacqueline Matthew<sup>2</sup>, Chandni Gupta<sup>2</sup>, Caroline Knight<sup>2</sup>, Bernhard Kainz<sup>1</sup>, and Daniel Rueckert<sup>1</sup>

 $^1\,$ Biomedical Image Analysis Group, Imperial College London, UK  $^2\,$  School of Biomedical Engineering & Imaging Sciences, King's College London, UK

**Abstract.** We propose a new Patch-based Iterative Network (PIN) for fast and accurate landmark localisation in 3D medical volumes. PIN utilises a Convolutional Neural Network (CNN) to learn the spatial relationship between an image patch and anatomical landmark positions. During inference, patches are repeatedly passed to the CNN until the estimated landmark position converges to the true landmark location. PIN is computationally efficient since the inference stage only selectively samples a small number of patches in an iterative fashion rather than a dense sampling at every location in the volume. Our approach adopts a multi-task learning framework that combines regression and classification to improve localisation accuracy. We extend PIN to localise multiple landmarks by using principal component analysis, which models the global anatomical relationships between landmarks. We have evaluated PIN using 72 3D ultrasound images from fetal screening examinations. PIN achieves quantitatively an average landmark localisation error of 5.59mm and a runtime of 0.44s to predict 10 landmarks per volume. Qualitatively, anatomical 2D standard scan planes derived from the predicted landmark locations are visually similar to the clinical ground truth. Source code is publicly available at https: //github.com/yuanwei1989/landmark-detection.

# 1 Introduction

Anatomical landmark localisation is a key challenge for many medical image analysis tasks. Accurate landmark identification can be used for (a) extracting biometric measurements of anatomical structures, (b) landmark-based registration of 3D volumes, (c) extracting 2D clinical standard planes from 3D volumes and (d) initialisation of tasks such as image segmentation. However, manual landmark detection is time-consuming and suffers from high observer variability. Thus, there is a need to develop automatic methods for fast and accurate landmark localisation. Recently, deep learning approaches have been proposed for this purpose [5,3,6,8,7,1] but there remain major challenges: (a) typically only a limited amount of annotated medical images is available, (b) model training and inference for 3D medical images is computationally intensive, making

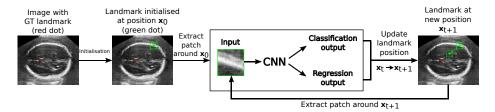


Fig. 1: Overall framework of PIN for single landmark localisation.

real-time applications challenging and (c) when multiple landmarks are detected jointly, their spatial relationships should be taken into account.

Related work: Deep learning methods for landmark localisation can be divided into two categories: The first category adopts an end-to-end learning strategy where the entire image is taken as input to a convolutional neural network (CNN) while the output is a map from which the landmark coordinates can be inferred directly. Payer et al. [5] and Laina et al. [3] output a heatmap in which Gaussians are located at the landmark positions. Xu et al. [6] train a supervised action classifier (SAC) that outputs an action map whose classification labels denote the direction towards the true landmark location. However, end-to-end learning methods are typically applied to 2D images since 3D volumetric networks require large receptive fields for landmark tasks. Such 3D networks are computationally intensive, which inhibits real-time performance, and require a large amount of memory during training, which is beyond current hardware's capabilities.

The second category uses image patches as training samples to learn a classification or regression model. Zheng et al. [8] extract a patch around each voxel in the image and use a neural network to classify if a landmark is present at the patch centre. Zhang et al. [7] and Aubert et al. [1] use a CNN-based regression model that learns the association between an image patch and its 3D displacement to the true landmark. Ghesu et al. [2] propose a deep reinforcement learning (DRL) approach that also operates on patches. Most patch-based methods require dense sampling of many image patches during prediction which is computationally intensive. Furthermore, most methods require the training of separate models to detect each landmark. This is time-consuming and neglects the spatial relationships among multiple landmarks.

Contribution: In this paper, we propose a novel landmark localisation approach that uses a patch-based CNN to predict multiple landmarks efficiently in an iterative manner. We term this approach Patch-based Iterative Network (PIN). PIN has distinct advantages that address the key challenges of landmark localisation in 3D medical images: (1) During inference, PIN guides the patch towards the true landmark location using iterative sparse sampling. This approach reduces the computational cost by avoiding dense sampling at every voxel of the volume. (2) PIN uses a 2.5D representation to approximate the 3D patch as network input. This accelerates computation as only 2D convolutions are required. (3) PIN treats landmark localisation as a combined regression and classification

problem for which a joint network is learned via multi-task learning. This prevents model overfitting, improves generalisation ability of the learned features and increases localisation accuracy. (4) PIN detects multiple landmarks jointly using a single model and takes the global anatomical spatial relationships among landmarks into account. We evaluate the landmark localisation accuracy of PIN using 3D ultrasound images of the fetal brain. In addition, clinically useful scan planes can be extracted from the predicted landmarks which visually resemble the anatomical standard planes as defined by fetal screening standards, e.g., [4].

## 2 Method

Overall Framework: Fig. 1 illustrates the overall PIN framework for single landmark localisation. We show the 2D case for clarity but the method works similarly in 3D. Given an image, the goal is to predict the true landmark coordinates (red dot in Fig. 1). A position  $x_0$  is first initialised at instant t=0 and a patch centred around  $x_0$  is extracted (solid green box in Fig. 1). The CNN takes the patch as input and predicts regression and classification outputs that are used to compute a new position  $x_{t+1}$  from the previous position  $x_t$ , bringing the patch closer to the true landmark location. The patch at  $x_{t+1}$  (dashed green box in Fig. 1) is then given as input to the CNN and the process is repeated until the patch reaches the true landmark position.

**Network Input:** For 3D data, the CNN input can be a 3D volume patch. However, 3D convolution operations on volume patches are computationally expensive. To this end, we use a 2.5D representation to approximate the full 3D patch. Specifically, given a particular position  $\mathbf{x} = (x, y, z)^T$  in a volume V, we extract three 2D image patches centred around  $\mathbf{x}$  at the three orthogonal planes (Fig. 2a). The patch extraction function is denoted as  $I(V, \mathbf{x}, s)$  where s is the length of the square patch. The three 2D patches are then concatenated together as a 3-channel 2D patch which is passed as input to the CNN. Such a representation is computationally efficient since it requires only 2D convolutions and still provides a good approximation of the full 3D volume patch.

Joint Regression and Classification: PIN jointly predicts the magnitude and direction of movement of a current point towards the true landmark by combining a regression and a classification task together in a multi-task learning framework. This joint framework shares model parameters in the convolutional layers and is experimentally shown to learn more generalisable features, which improves overall performance.

The regression task estimates how much the point at the current position should move to get to the true landmark location. The regression output  $d = (d_1, d_2, \ldots, d_{n_o})^T$  is a displacement vector that predicts the relative distance between the current and true landmark positions. In single landmark localisation, d has  $n_o = 3$  elements which give the displacement along each coordinate axis.

The classification task estimates the direction of current point movement towards the true landmark by dividing direction into 6 discrete classification categories: positive and negative direction along each coordinate axes [6]. Denoting

#### 4 Y. Li et al.

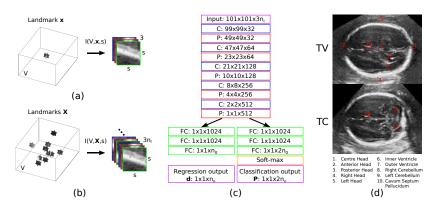


Fig. 2: (a) Patch extraction of a single landmark. (b) Patch extraction of multiple landmarks. (c) CNN architecture combining regression and classification. Output size of each layer is represented as width × height × (# feature maps). (d) Landmarks defined on the TV and TC plane for fetal sonographic examination.

c as the classification label, we have  $c \in \{c_1^+, c_1^-, c_2^+, c_2^-, c_3^+, c_3^-\}$ . For instance,  $c_1^+$  is the category representing movement along the direction of positive x-axis. The classification output  $\boldsymbol{P}$  is then a vector with  $2n_o = 6$  elements, each representing the probability/confidence of movement in that direction. Mathematically,  $\boldsymbol{P} = (P_{c_1^+}, P_{c_1^-}, \dots, P_{c_{n_o}^+}, P_{c_{n_o}^-})^T$  where  $P_{c_1^+} = \operatorname{Prob}(c = c_1^+)$ .

Given a volume V and its ground truth landmark point  $\boldsymbol{x}^{GT}$ , a training sample is represented by  $(I(V,\boldsymbol{x},s),\boldsymbol{d}^{GT},\boldsymbol{P}^{GT})$  where  $\boldsymbol{x}$  is a point randomly sampled from V and  $I(V,\boldsymbol{x},s)$  is its associated patch. The ground truth displacement vector is given by  $\boldsymbol{d}^{GT} = \boldsymbol{x}^{GT} - \boldsymbol{x}$ . To obtain  $\boldsymbol{P}^{GT}$ , we first determine the ground truth classification label  $c^{GT}$  by selecting the component of  $\boldsymbol{d}^{GT}$  with the maximum absolute value and taking into account its sign,

$$c^{GT} = \begin{cases} c_i^+, & \text{if } d_i^{GT} > 0\\ c_i^-, & \text{otherwise.} \end{cases}$$
 (1)

where  $i = \operatorname{argmax}(\operatorname{abs}(\boldsymbol{d}^{GT}))$ . For a vector  $\boldsymbol{a}$ ,  $\operatorname{argmax}(\boldsymbol{a})$  returns the index of the vector component with maximum value. During training, a hard classification label is used. As such, the probability vector  $\boldsymbol{P}^{GT}$  is obtained as a one-hot vector where component  $P_{cGT}$  is set to 1 and all others set to 0. The CNN is trained by minimising the following combined loss function:

$$L = (1 - \alpha) \frac{1}{n_0 n_{batch}} \sum_{n=1}^{n_{batch}} \left\| d_n^{GT} - d_n \right\|_2^2 - \alpha \frac{1}{n_{batch}} \sum_{n=1}^{n_{batch}} \log \left( P_{cGT, n} \right)$$
 (2)

The first term is the Euclidean loss of the regression task and the second term is the cross-entropy loss of the classification task.  $\alpha$  is the weighting between the two losses.  $n_{batch}$  is the number of training samples in a mini-batch.  $d_n$  and

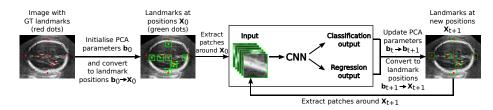


Fig. 3: Overall framework of PIN for multiple landmarks localisation.

 $P_{c^{GT},n}$  denote respectively the regression and classification outputs predicted by the CNN on the *n*th sample.

CNN Architecture: Fig. 2c shows the PIN CNN architecture combining the classification and regression tasks. The network comprises five convolution (C) layers, each followed by a max-pooling (P) layer. These layers are shared by both tasks. After the 5<sup>th</sup> pooling layer, each task has three separate fully-connected (FC) layers to learn the task-specific features. All convolution layers use 3x3 kernels with stride=1 and all pooling layers use 2x2 kernels with stride=2. ReLU activation is applied after all convolution and FC layers except for the last FC layer of each task. Drop-out is added after each FC layer.

**PIN Inference:** Given an unseen 3D volume, we initialise 19 points in the volume (one at the volume centre and 18 others at fixed distance of one-quarter image size around it). The patch extracted from each point is forward passed into the CNN and the point is moved to its new position based on the CNN outputs (d and P) and a chosen update rule. This process is repeated for T iterations until there is no significant change in the displacement of the point. The final positions of the 19 points at iteration T is averaged and taken to be the final landmark prediction. Multiple initialisations average out errors and improve the overall localisation accuracy.

**PIN update rules:** We proposed three update rules (A-C). Let  $x_t$  be the position of a point at iteration t and  $x_{t+1}$  be the new updated position. Rule A is based only on the classification output P. It updates the current landmark position by moving it one pixel in the direction category which has the highest probability as predicted by P. Rule B is based only on the regression output d and is given by:  $x_{t+1} = x_t + d$ . Rule C uses both the classification and regression outputs for the update and is given by:  $x_{t+1} = x_t + P_{max} \odot d$  where  $\odot$  is the element-wise multiplication operator and  $P_{max} = (\max(P_{c_1}^+, P_{c_1}^-), \max(P_{c_2}^+, P_{c_2}^-), \ldots, \max(P_{c_{no}}^+, P_{c_{no}}^-))^T$ . Intuitively, Rule C moves the point to its new position by an amount specified by the regression output weighted by a confidence probability specified by the classification output. This ensures smaller movement in the less confident direction and vice versa.

Multiple Landmarks Localisation: The above approach for single landmark localisation has two drawbacks: (1) Separate CNN models are required for each landmark which increase the parametrisation significantly and thus computational cost for training and inference. (2) Individual landmark prediction ignores the anatomical relationships between the different landmarks. To overcome these

problems, we extend our approach to localise multiple landmarks simultaneously using only one CNN model which also accounts for inter-landmark relationships by working in a reduced dimensional space.

Let  $\mathbf{X} = (x_1, y_1, z_1, \dots, x_{n_l}, y_{n_l}, z_{n_l})^T$  be the 3D coordinates of all  $n_l$  landmarks of one volume. Given a training set of  $\mathbf{X}$ , we use PCA to transform  $\mathbf{X}$  into a lower dimensional space. The transformations between the original and reduced dimensional spaces are given by,

$$X = \bar{X} + Wb \tag{3}$$

$$\boldsymbol{b} = \boldsymbol{W}^T (\boldsymbol{X} - \bar{\boldsymbol{X}}), \tag{4}$$

where  $\bar{X}$  is the mean of the training set, b is a  $n_b$ -element vector where  $n_b < 3n_l$  and the columns of matrix W are the  $n_b$  eigenvectors. In our case,  $n_l = 10$  and we set  $n_b = 15$  to explain 99.5% of the total variations in the training set.

We can directly apply our PIN approach to the reduced dimensional space by replacing all occurrences of  $\boldsymbol{x}$  by  $\boldsymbol{b}$ . Fig. 3 illustrates the PIN approach for multiple landmarks. Specifically, 3 orthogonal patches are extracted for every landmark and concatenated together so that a  $s \times s \times 3n_l$  block is passed as CNN input (Fig. 2b).  $\boldsymbol{d}$  becomes the displacement vector in the reduced dimensional space with  $n_o = n_b$  elements. The number of classification categories becomes  $2n_b$  which include positive and negative directions along each dimension of the reduced space. Hence,  $\boldsymbol{P}$  is a  $2n_b$ -element vector. Training can be carried out similar to Eq. 2 with the only difference being  $\boldsymbol{d}^{GT} = \boldsymbol{b}^{GT} - \boldsymbol{b}$  where  $\boldsymbol{b}^{GT}$  is transformed from  $\boldsymbol{x}^{GT}$  using Eq. 4 and  $\boldsymbol{b}$  is randomly sampled. During inference, we update  $\boldsymbol{b}$  iteratively using  $\boldsymbol{b}_{t+1} = \boldsymbol{b}_t + \boldsymbol{P}_{max} \odot \boldsymbol{d}$  (Rule C) and use Eq. 3 to convert  $\boldsymbol{b}_{t+1}$  back to  $\boldsymbol{X}_{t+1}$  for patch extraction in the next iteration. We use multiple initialisations of  $\boldsymbol{b}_0$  (one initialisation with  $\boldsymbol{b}_0 = \boldsymbol{0}$  and five random initialisations) and take their mean results as the final landmarks prediction.

#### 3 Experiments and Results

Data: PIN is evaluated on 3D ultrasound volumes of the fetal head from 72 subjects. Each volume is annotated by a clinical expert with 10 anatomical landmarks that lie on two standard planes (transventricular (TV) and transcerebellar (TC)) commonly used for fetal sonographic examination as defined in the UK FASP handbook [4] (Fig. 2d). 70% of the dataset is randomly selected for training and the remaining 30% is used for testing. All volumes are processed to be isotropic and resized to  $324 \times 207 \times 279$  voxels with voxel size  $0.5 \times 0.5 \times 0.5 \times 0.5$  mm<sup>3</sup>. Experiment Setup: PIN is implemented using Tensorflow running on a machine with Intel Xeon CPU E5-1630 at 3.70 GHz and one NVIDIA Titan Xp 12GB GPU. Patch size s is set to 101. During training, we set  $n_{batch}$ =64. Weights are initialised randomly from a distribution with zero mean and 0.1 standard deviation. Optimisation is carried out for 100,000 iterations using the Adam algorithm with learning rate=0.001,  $β_1$ =0.9 and  $β_2$ =0.999. We choose α=0.5 empirically unless otherwise stated. During inference, T=350 for Rule A and T=10 for Rule B and C.

Table 1: Localisation error (mm) and runtime (s) of different approaches for single landmark (CSP) localisation. C and R denote classification and regression training loss respectively. Results presented as (Mean  $\pm$  Standard Deviation).

	PIN1	PIN2	PIN3	PIN4	PIN5	DRL [2]
Training loss	С	R	C+R	C+R	C+R	-
Inference rule	Rule A	Rule B	Rule A	Rule B	Rule C	-
Localisation error	$7.53\pm6.48$	$6.45 \pm 3.96$	$6.34 \pm 3.62$	$6.08 \pm 3.90$	$5.47{\pm}4.23$	$7.37 \pm 5.86$
Running time (s)	3.56	0.09	3.50	0.09	0.09	6.58

Table 2: Localisation error (mm) of PIN for single and multiple landmark localisation. Results presented as (Mean  $\pm$  Standard Deviation).

Landmarks	1	2	3	4	5	6	7	8	9	10	Overall
PIN-Single	$5.62 \pm 2.85$	$11.30 \pm 7.24$	8.13±3.90	$7.23\pm3.73$	$7.11 \pm 4.73$	$4.39\pm2.07$	$5.45 \pm 2.73$	$4.04{\pm}2.22$	$5.50 \pm 3.64$	$5.47{\pm}4.23$	6.42±4.49
PIN-Multiple	$4.34{\pm}2.21$	$8.80 \pm 4.27$	$6.28 \pm 2.77$	$6.31 \pm 3.32$	$5.56 \pm 2.71$	$4.68 \pm 2.27$	$5.15 \pm 2.90$	$4.70 \pm 2.33$	$4.57{\pm}1.92$	$5.50 \pm 2.79$	$5.59 \pm 3.09$

**Results:** Table 1 compares the landmark localisation errors of a single landmark, cavum septum pellucidum (CSP), using several PIN variants which differ in the CNN model training and the inference update rule. Given the same inference rules, the model trained using both classification and regression losses ( $\alpha$ =0.5) achieves lower error than the models trained using either loss alone ( $\alpha=1$  or 0) (PIN1 vs PIN3, PIN2 vs PIN4). This illustrates the benefits of multi-task learning. Using the model trained with joint losses, we then compare the effect of different inference rules. PIN3 uses only the classification output which can result in landmarks getting stuck and oscillating between two opposing classification categories during iterative testing (e.g.,  $c_1^+$  and  $c_1^-$ ). PIN3 also takes longer during inference since the patch moves by one pixel at each test iteration and requires more iterations to converge. PIN4 uses only the regression output, which improves the localisation accuracy and runtime as the patch 'jumps' towards the true landmark position at each iteration. This requires much fewer iterations to converge. PIN5 achieves the best localisation accuracy by combining the classification and regression outputs where the regression output gives the magnitude of movement weighted by the classification output giving the probability of movement in each direction. Our proposed PIN approach also outperforms a recent state-of-the-art landmark localisation approach using DRL [2].

Table 2 shows the localisation errors for all ten landmarks. PIN-Single trains a separate model for each landmark while PIN-Multiple trains one joint model that predicts all the landmarks simultaneously. Since PIN-Multiple accounts for anatomical relationships among the landmarks, it has a lower overall localisation error than PIN-Single. PIN-Single needs a total of 0.94s to predict all ten landmarks in sequence while PIN-Multiple needs 0.44s to predict all ten landmarks simultaneously. Fig. 4 shows the TV and TC planes containing the ground truth landmarks as red dots. The landmarks predicted by PIN-Multiple are projected onto these standard planes as green dots. The supplementary materials provide visual comparison of standard planes obtained from ground truth and predicted landmarks as well as videos showing several initialisations converging towards the true landmark positions (and standard planes) after ten inference updates.

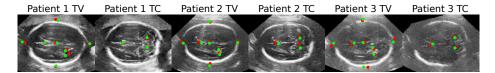


Fig. 4: Visualisation of landmarks predicted by PIN-Multiple (green dots) vs. ground truth landmarks (red dots).

## 4 Conclusion

We have presented PIN, a new approach for anatomical landmark localisation. Its patch-based and iterative nature enables training on limited data and fast prediction on large 3D volumes. A joint regression and classification model is trained by multi-task learning to improve localisation accuracy. PIN is capable of multiple landmark localisation and uses PCA to impose anatomical constraints among landmarks. PIN is generic to landmark localisation and as future work, we are extending PIN to other medical applications. It is also worthwhile to replace PCA with an autoencoder to model non-linear correlations among landmarks.

**Acknowledgments.** Supported by the Wellcome Trust IEH Award [102431]. The authors thank Nvidia Corporation for the donation of a Titan Xp GPU.

#### References

- Aubert, B., Vazquez, C., Cresson, T., Parent, S., Guise, J.D.: Automatic spine and pelvis detection in frontal x-rays using deep neural networks for patch displacement learning. In: ISBI 2016. pp. 1426–1429 (April 2016)
- Ghesu, F.C., Georgescu, B., Mansi, T., Neumann, D., Hornegger, J., Comaniciu,
  D.: An artificial agent for anatomical landmark detection in medical images. In: MICCAI 2016. pp. 229–237 (2016)
- Laina, I., Rieke, N., Rupprecht, C., Vizcaíno, J.P., Eslami, A., Tombari, F., Navab, N.: Concurrent segmentation and localization for tracking of surgical instruments. In: MICCAI 2017. pp. 664–672 (2017)
- NHS: Fetal anomaly screening programme: programme handbook June 2015. Public Health England (2015)
- Payer, C., Štern, D., Bischof, H., Urschler, M.: Regressing heatmaps for multiple landmark localization using cnns. In: MICCAI 2016. pp. 230–238 (2016)
- Xu, Z., Huang, Q., Park, J., Chen, M., Xu, D., Yang, D., Liu, D., Zhou, S.K.: Supervised action classifier: Approaching landmark detection as image partitioning. In: MICCAI 2017. pp. 338–346 (2017)
- Zhang, J., Liu, M., Shen, D.: Detecting anatomical landmarks from limited medical imaging data using two-stage task-oriented deep neural networks. IEEE Transactions on Image Processing 26(10), 4753–4764 (Oct 2017)
- 8. Zheng, Y., Liu, D., Georgescu, B., Nguyen, H., Comaniciu, D.: 3d deep learning for efficient and robust landmark detection in volumetric data. In: MICCAI 2015. pp. 565–572 (2015)