

Date of publication XXXX, date of current version XXXX.

Digital Object Identifier XXXXXXXXX

# DeepNav: Joint View Learning for Direct Optimal Path Perception in Cochlear Surgical Platform Navigation

**MAJID ZAMANI, (Member, IEEE) AND ANDREAS DEMOSTHENOUS, (Fellow, IEEE)**

Department of Electronic and Electrical Engineering, University College London, Malet Place, London WC1E 7JE, United Kingdom.

Corresponding author: Andreas. Demosthenous (e-mail: a.demosthenous@ucl.ac.uk).

This work was supported by the Engineering and Physical Sciences Research Council (EPSRC) under grant EP/R511638/1.

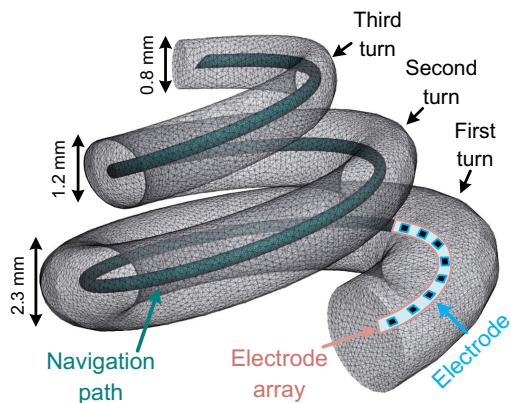
**ABSTRACT** Although much research has been conducted in the field of automated cochlear implant navigation, the problem remains challenging. Deep learning techniques have recently achieved impressive results in a variety of computer vision problems, raising expectations that they might be applied in other domains, such as identifying the optimal navigation zone (OPZ) in the cochlear. In this paper, a 2.5D joint-view convolutional neural network (2.5D CNN) is proposed and evaluated for the identification of the OPZ in the cochlear segments. The proposed network consists of 2 complementary sagittal and bird-view (or top view) networks for the 3D OPZ recognition, each utilizing a ResNet-8 architecture consisting of 5 convolutional layers with rectified nonlinearity unit (ReLU) activations, followed by average pooling with size equal to the size of the final feature maps. The last fully connected layer of each network has 4 indicators, equivalent to the classes considered: the distance to the adjacent left and right walls, collision probability and heading angle. To demonstrate this, the 2.5D CNN was trained using a parametric data generation model, and then evaluated using anatomically constructed cochlea models from the micro-CT images of different cases. Prediction of the indicators demonstrates the effectiveness of the 2.5D CNN, for example the heading angle has less than  $1^\circ$  error with computation delays of less than  $<1$  milliseconds.

**INDEX TERMS** Automated insertion, virtual surgery, cochlear implant, convolutional neural network, real-time systems, low-cost navigation, robust centerline tracing.

## I. INTRODUCTION

The cochlear implant (CI) [1] is one of the most successful implantable devices in clinical practice. It helps to restore lost hearing by delivering electrical impulses to the auditory nerves via an electrode array inserted into the cochlea in the inner ear [2]. Cochlear implant navigation involves inserting a wire containing a line of stimulating electrodes into the delicate spiral (or snail) shaped tube that varies in diameter and height along the  $Z$  plane and imposing geometrical limitations to the cochlear implant surgery as shown in Fig. 1. The quality of restored hearing sensation is strongly related to the efficacy of surgery of the cochlear device implantation, particularly the optimum positioning and the insertion depth of the electrode array inside the cochlea without further damaging the remaining hearing [3]. The present standard technique relies on the surgeon's fingertips while pushing the electrode array down the spiral-shaped cochlea. This approach requires the surgeon to identify the

optimal insertion path solely by feel. Although the tip of the electrode array is not sharp to pierce through the bony wall of cochlea, extreme pressure may increase the risk of the electrode tip crossing the auditory nerve or the modiolus. Medical-imaging techniques such as MRI, computerized tomography (CT) [5] and X-rays [6] are not practical options for guidance in implantation surgery as they cannot provide real-time imaging and they are impractical due to the very small volume of the cochlea. The systems in [7]-[13] for cochlear implant navigation derive information from impedance measurements on the electrodes at the end of the electrode array. While useful for identifying the position of the electrode tip, performance is compromised by the limited accuracy of the measured impedance values. Integration of a robotic arm [14] does not lead to better navigation performance as it similarly receives the guidance parameters from imprecise calculations. The present limited accuracy of identification of the position of the tip would be improved by



**FIGURE 1. Guidance of cochlear implant electrode array.** The mean heights at the basal, middle and apical turns are 2.3 mm, 1.2 mm and 0.8 mm respectively [4]. The quality of restored hearing sensation is strongly related to the optimum positioning and the insertion depth of the electrode array inside the cochlea without further damaging the remaining hearing.

embedding intelligence, which would require accurate navigation.

To avoid adverse consequences such as crossing anatomical wall as a result of the extreme geometrical limitations, computer-assisted surgery [15]-[16] is used to identify the extremely precise centerline trajectory required inside a three-dimensional (3D) reconstructed cochlea as priori knowledge (or post-processing) for cochlear implant electrode array insertion by automatic means.

The electrode insertion algorithm is designed based on the type of electrode: 1) lateral wall electrode [17] that slides along the spiral ligament; and 2) modular-hugging electrode [18] which tends to go closer to the inner wall (the modiolus). This paper proposes a method to significantly enhance cochlear implant navigation by identifying rapidly an interactive safe insertion zone in real-time using a novel 2.5D convolutional neural network (CNN), yielding very high insertion resolution accuracy. The proposed 2.5D algorithm navigates the tip of the electrode safely along the centerline coordinates to ensure minimal insertion risk while the rest of electrodes would slide along the cochlear wall. The electrode array model was based on a commercially available electrode [Advanced Bionics HiFocusTM SlimJ electrode (Hannover, Germany)] with 16 platinum electrodes.

The rest of the paper is organized as follows. Section II presents the prior art and the CI navigation algorithm proposed in this work. Section III describes the methods used in data generation and proposes a framework to derive the navigation indicators. It also discusses the design of the 2.5D CNN and the joint 3D operator. Section IV details the efficacy of the 2.5D CNN in different scenarios and visualizes the navigation steps for an anatomical cochlea model. Concluding remarks are drawn in Section V.

## II. RELATED WORK: CENTERLINE TRACING ALGORITHMS

There are a variety of approaches that can be utilized to identify the centerline of tubular structures. One category

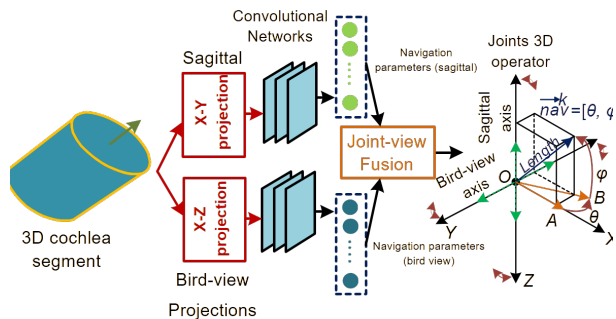
consists of skeletonization approaches [19] and those using multiscale enhancement, morphological reconstruction and segmentation methods [20]-[23]. They require the processing of full 3D volume and every image pixel with numerous operations per pixel.

A second category tracks the centerline based on a filter or an assumed model. Commonly used filters are based on eigen-structure of local Hessian [24], idealized tubular models of vessels [25] and Hough transforms [26] to locate vessel direction and its cross vectors at a reference frame. For example, Hessian of the image is interpreted as second order partial derivatives of 3D sub-images at reference nodes, which requires extensive computation time. Cylindroidal superellipsoids [27] is an advanced model of probing for 3D tubular shapes using recursive fitting methods. Although the fitting-based approaches perform well across morphological complexities, they derive model parameters using maximum likelihood which is an extremely complex and lengthy process.

A third category utilizes vectorization algorithms [28]-[30] for tubular structure boundary analysis and centerline tracing where only pixels close to the border are processed. They are well-suited to real-time and robust tracing in large image sets. The sparse exploration of the boundaries yields low computational overhead but also introduces higher sensitivity to the discontinuities and geometrical complexities. An algorithm utilizing vectorization approach to handle 3D (volumetric) data is described in [31]. It is a fully automatic 3D neuron tracing algorithm emulating a 3D cylinder model and recursively explores the neuron topology. The simulations using the 3D cylinder algorithm on constructed cochlea models illustrate that the centerline tracing does not perform reliably when it is faced with high-order tubular changes.

Machine learning offers an alternative approach to identify and trace the central coordinates [32]-[34]. Steerable features and randomized decision trees are used in [35] to perform centerline extraction by learning the structural patterns of a tubular-like object. The approach in [32] uses orientation flow field and classifier to extract blood vessel centerlines. The average computation for tracing all coronaries takes about 1 minute on an Intel Core i7 2.8 GHz processor with 32 GB RAM as reported in [34].

CNNs are a class of deep learning algorithms that have recently been utilized in 3D tubular structure tracing [35]-[37]. In [35], a 3D dilated CNN [38] was trained to predict the most likely direction and radius of an artery at any given seed point. The tracing scheme in [35] was developed based on determining a posterior probability distribution over a discrete set of possible directions as well as an estimate of the radius. The drawback with this design is that the optimal direction determination is posed as a classification problem, thus the possible directions are distributed on a sphere where each point corresponds to a class. The best classification performance was obtained for the directions {500, 1000 or 2000}. The design in [35] demands excessive computational cost in classifying directions and is not suitable for real-time applications; it requires 20 seconds for



**FIGURE 2.** Overview of the proposed joint-view navigation framework. The sagittal and bird-view views are generated by projecting the 3D points onto two orthogonal planes (i.e. X-Y and X-Z planes). Two CNNs are trained in parallel to map each view's projected image to its corresponding navigation indicators, which are then fused together to estimate 3D joint operator.

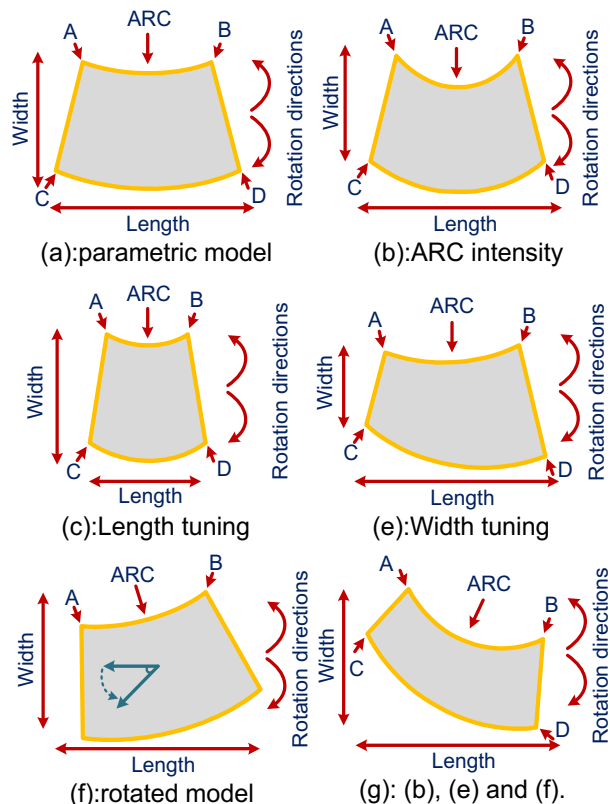
fully automatic coronary tree extraction using the Nvidia Titan Xp GPU. In [36] and [37], 3D CNNs were proposed to trace the cardiovascular tree structure. They require 58 and 25 seconds using 12 GB GPU and Tesla P40 GPU, respectively.

This paper proposes a novel low computation deep navigation method using a 2.5D multi-view CNNs that can better transform the input image to a small number of key perception indicators to recover 3D tracing information on the tubular structures, as shown in Fig. 2. The 3D cochlea segment is pre-processed and projected onto sagittal and bird-view planes and then applied to separate CNNs for mapping process. Each view decodes the relevant navigation (or tracing) information and fuses them; so contains the location distribution of the joint-view 3D tracing operator.

The proposed tracing method has the following contributions: 1) A 2.5D tracing algorithm which shows significant trade-off between the performance and processing time by removing a dimension of an image. The algorithm provides a good fit for tracing-related tasks in real-time processing images. 2) A compact residual convolutional architecture is used for each projected 2D image. It predicts the steering angle and the indicators including the collision probability and the distance to the left-right walls in real-time. 3) A direct perception approach maps an input image to a small set of indicators that are used in identifying the optimal tracing path or insertion zone for navigation of the electrodes inside, for example, a cochlea. The mapping framework performs abstraction of the images by keeping only a set of compact and yet complete descriptors which results in real-time optimal path identification. 4) A comprehensive physiological-inspired tubular dataset provides a very diverse set of virtual environments for training the 2.5D tracing algorithm. Through extensive evaluation, it is shown that the trained model is efficient and can be applied to real cochlea models. The training set-up can be completely generalized for unseen scenarios. 5) A joint 3D operator for navigation in 3D set-ups.

### III. METHODES

In this section, first the datasets used in this study are described. It is followed by the deep mapping framework for



**FIGURE 3.** Synthetic data generation. (a) Illustration of the parametric cochlea segment model. There are two arcs defined between the A and B nodes, and between the C and D nodes. Their width and the length are tunable in this proposed model. (b) shows the arc intensity change. In (c), the length of the arcs is tuned to the smaller values. (e) shows when the width is tuned based on adjusting the arc length between C and D. Cochlea segment rotation is an important factor in implant navigation and this capability is shown in (f). (g) Combining (b), (e) and (f).

extraction of the navigation indicators and the architecture of CNN. The definition of the input data and desired outputs provide a better understanding of the methods. It finally discusses the joint 3D navigation operator.

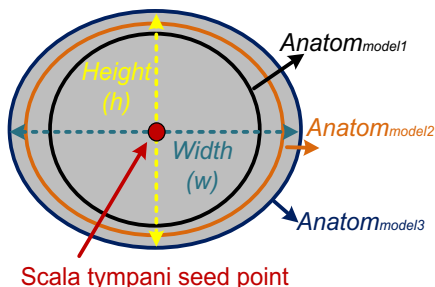
#### A. DATA

To learn the navigation indicators (or parameters) in cochlea tracing, two types of dataset were utilized. The first dataset is composed of synthetic MATLAB-generated images for training purposes. The second dataset contains anatomical cochlea models. Both are used to quantitatively analyze the navigation performance on the 2.5D multi-view CNNs.

##### 1) SYNTHETIC IMAGES

Considering the sagittal and bird views of the cochlea structure, a parametric segment model of the cochlea is proposed to accommodate all the navigation features for training the 2.5D CNNs. The model shown in Fig. 3(a) has deformation capability to emulate all the variations along the cochlea such as bend, rotation and length-width variation. For example, in the bird-view mode (i.e. looking at the cochlea from the top), the bend intensity changes constantly along the cochlea. The bend in each cochlea segment (either bird-view or sagittal) is composed of two crucial parameters; the arc intensity and the turning effect which are evident





**FIGURE 4.** The initial height ( $h$ ) and width ( $w$ ) ( $h < w$ ) of the anatomical models (Anatom<sub>model1</sub>... Anatom<sub>model3</sub>) around the scala tympani seed point. The models are designed to have  $(h/w)_{model1} < (h/w)_{model2} < (h/w)_{model3}$ . The models consider geometrical variations along the navigation paths.

when there are sharp turns. Both effects are shown in the Fig. 3(a)-(b). A closer look at Fig. 3(b) shows that the inner arc between A and B nodes is smaller compared with the outer arc between nodes C and D, which defines the turnings along cochlea. The length and the width vary radically along the cochlea path (e.g. the mean width at the basal, middle and apical turns are 2.1 mm, 1.2 mm and 0.6 mm, respectively [4]). The length and the width are, therefore, generated for various sizes to cover all the variations along the cochlea. In Fig. 3(c) the width of the cochlea segment is tuned by stacking the number of length-adjusted arcs. Orientation information is important for cochlear tracing. In the proposed parametric model it is required to obtain a rotational invariant representation for cochlea segments. In order to make the model more robust to orientation variations, the generated images are also rotated along  $z$ -axis by  $[0:360^\circ]$  to emulate the bird-view of the cochlea and along  $y$ -axis by  $[0:90^\circ]$  to generate the sagittal tracing segments. The rotation step size is  $5^\circ$ . Overall, the most practical point in data generation is to design the edges having high correlation with the cochlea projection into two orthogonal planes. Generating the right edges greatly helps to identify the navigation inferences, through the generalization capability and the noise-artefact robustness of the 2.5D CNN.

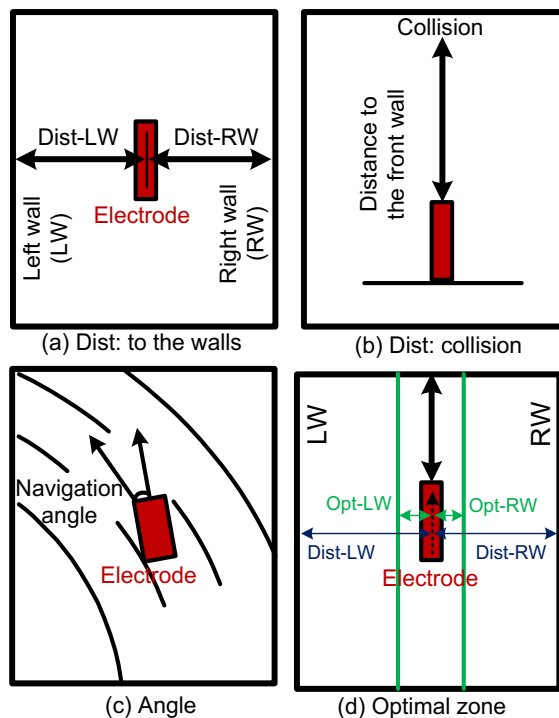
## 2) COCHLEA MODELS

The 2.5D CNN and tracing algorithms were examined with a set of three synthetic cochlea models (Synth<sub>model1</sub>... Synth<sub>model3</sub>). The purpose of utilizing synthetic data is to provide an analysis of the algorithms under controlled conditions that mimic the cochlea structure. The averaged model used for the synthetic cochlea models was generated in MATLAB 2022.b using:

$$x = \left(\frac{S}{5}\right) \sin(s), y = \left(\frac{S}{5}\right) \cos(s), z = \left(-\frac{S}{3}\right). \quad (1)$$

where  $s$  ranges from 6.5 to 21.25 to resemble the anatomical human cochlea with a mean length of 41.5 mm and diameter of 2 mm for parametric sweeping purposes [39]. The synthetic 3D cochlea models were constructed within a 10 mm  $\times$  10 mm  $\times$  10 mm volume comprising the cochlea model and the pad arrays to obtain consistent  $(x, y, z)$  dimensions for evaluation of tracing performance.

In a similar manner three anatomical cochlea models were constructed from micro-CT images (Anatom<sub>model1</sub>...



**FIGURE 5.** Illustration of navigation indicators. (a) electrode distance to the left and right walls, ( $Dist - LW$ ) and ( $Dist - RW$ ). (b) Collision probability ( $Collision$ ) which shows the distance to the front wall. (c) The navigation angle and (d) safe insertion zone for optimal navigation.

Anatom<sub>model3</sub>). These evaluate the centerline tracing algorithm against a “golden standard,” i.e., a hand-traced centerline by clinicians in realistic reconstructed cochlea models. The realistic cochlea models were derived from micro-CT images of  $512 \times 512$  pixels per slice. A manually defined ground-truth was used to quantify traversal performance. The micro-CT data was imported to Simpleware ScanIP v2016.09 (Synopsys, Mountain View, USA) for image processing and data segmentation by defining regions in the image data that belong to the same anatomical layers. Smoothing filters utilizing recursive Gaussian, median, and mean filters were used to adjust the grayscale range. Manual segmentation was used by editing the morphology or filling cavities (i.e. dilate, erode, open and closed functions) were used in ScanIP software. To obtain appropriate boundaries and remove any overlapping sections between the tissue layers, Boolean operations were applied. The volume conductor of the cochlea and the layers in its vicinity were generated based on a high-resolution ( $2.24 \mu\text{m} \times 2.24 \mu\text{m} \times 5 \mu\text{m}$ ) voxel size micro-CT image stack of a human cochlea. Due to limited computation memory, the effective operative field of the scans was rescaled to include only the cochlea and its immediate surroundings and was subsequently down sampled to an isotropic resolution of  $9.6 \mu\text{m}$  with a spatial resolution of  $930 \times 930 \times 1014$  voxels.

The constructed synthetic and anatomical models represent height ( $h$ ) and width ( $w$ ) variations ( $h < w$ ) in human cochlea anatomy. For example, the  $(h/w)$  ratio of the Anatom<sub>model1</sub>, Anatom<sub>model2</sub> and Anatom<sub>model2</sub> are

(45/62), (35/55) and (50/67). It should be noted that the reported ratios are the initial height over width as shown in Fig. 4 and are decreased along the cochlea.

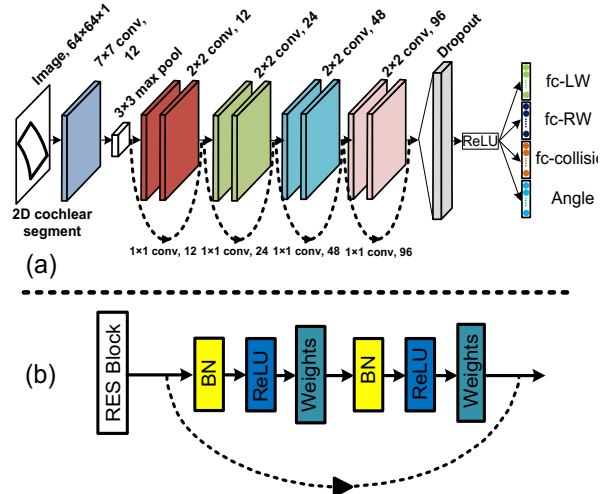
### B. DEEP MAPPING FROM AN IMAGE TO INDICATORS

A framework is laid out to map the generated image to a set of typical navigation indicators shown in Fig. 5. Three types of indicator to represent an optimal path navigation are proposed: the distance to the adjacent walls, the distance to the frontal wall (i.e. collision probability) and heading angle. The electrode array insertion is concerned with the two adjacent anatomical walls for following the centreline when the tip of the array is pushed inside the tubular structure. This is shown in Fig. 5(a) by identifying the distance of the electrode to the left and right walls indicated by  $(Dist - LW)$  and  $(Dist - RW)$  respectively. Collision probability ( $Collision$ ) is the next indicator that shows the maximum allowable navigation jump to avoid crossing the anatomical walls along insertion iterations. This is a crucial indicator as it accurately shows the stopping points specifically in the tubular turns before mapping the next image frame shown in Fig. 5(b). The navigation angles  $\gamma$  ( $\theta$ ,  $\varphi$ ) are the next indicators that direct the optimal rotation of the electrodes along the sagittal ( $\theta$ ) and bird-view ( $\varphi$ ) projection planes. In total, four affordance indicators to interpret the navigation scene are extracted from each image frame using the 2.5D CNN for each view. Considering the 2.5D view processing, a 3D safe insertion zone can be defined using all generated height and width variations of the cochlea in sagittal and bird-view projections around the predicted centerline coordinates as a hypothetical insertion cylinder [e.g. 50% of  $Dist - RW$  as shown in Fig. 5(d)].

### C. ARCHITECTURE OF THE 2.5D CNN

The 3D points are projected on two views (i.e. 2.5D view). For each view, a convolutional network having the same network architecture and architectural parameters and the outputs are constructed. Based on multi-task learning [40], a ResNet [41] architecture followed by separate outputs shown in Fig. 6 is proposed. Since residual architectures are known to help generalization on both shallow and deep networks, it is adapted to increase model performance. The architecture of the 2.5D CNN is highly compact, where the input layer has a size of  $64 \times 64$  to accept the sagittal or top views. The output of each 2D convolutional layer is activated by a rectified nonlinearity unit (ReLU) with its parameter equal to 0.1, which allows for a small, non-zero gradient when the unit is saturated and inactive.

Since most parameters in the proposed network lie in the first fully connected layer, a convolutional layer and a max-pooling layer are added to improve the degree of discrimination of the learned feature and reduce the number of parameters. Dotted lines represent skip connections defined as  $1 \times 1$  convolutional shortcuts to allow the input and output of the residual blocks to be added. After the last ReLU layer, the architecture splits into two different fully connected layers. The main branch consists of a fully connected layer and a softmax output layer to classify the



**FIGURE 6.** 2D CNN is a joint deep mapping network, from a single  $64 \times 64$  frame including  $(Dist - LW)$ ,  $(Dist - RW)$ , collision probability ( $Collision$ ) and the tracing angles along sagittal and bird views ( $\theta$ ,  $\varphi$ ). The main architecture of the CNN consists of a ResNet with 4 residual blocks. (b), followed by dropout and ReLU non-linearity. Afterwards, the network branches into 4 separated fully-connected layers. The design notation including the convolution kernel's size, the number of filters and the residual connections are shown in the figure.

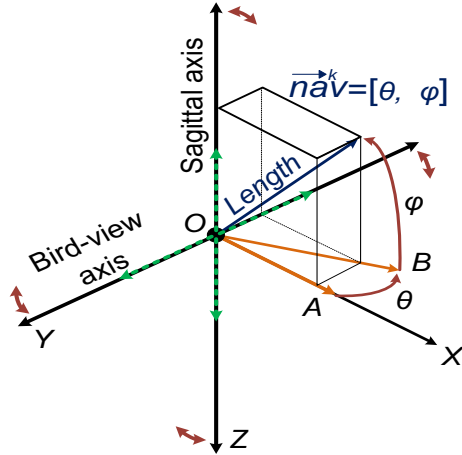
collision probability ( $Collision$ ), distance of the electrode to the left ( $Dist - LW$ ) and right ( $Dist - RW$ ) walls (see Section III.B). For the auxiliary branch, neurons are split to form a regression network for estimation of the tracing angles along sagittal or bird-view planes ( $\theta$ ,  $\varphi$ ). Mean-squared error (MSE) and cross-entropy (CEN) losses are utilized to classify the tracing angles and the affordance indicators, respectively:

$$L_{Tot} = \alpha(L_{MSE}) + \beta(L_{CEN}). \quad (2)$$

$L_{Tot}$ ,  $L_{MSE}$  and  $L_{CEN}$  represent the total loss of the model, the loss of tracing angle prediction and the loss of other indicators, and  $\alpha$ ,  $\beta$  show the loss weights. The network was designed with a compact architecture, but the joint optimization might pose a convergence problem. Specifically, imposing no weighting between the two losses during training results in convergence to a very poor solution. This is due to the fact that the MSE gradients' norms is proportional to the absolute tracing angle and initially has much higher value. Therefore,  $\alpha$  is set to 0.1 and more weight is assigned to  $L_{MSE}$  in later stages of training (i.e. 0.2-0.3). Adjusting the loss weight between the two losses would likely result in optimal performance or require much longer optimization times. The Adam optimizer [40] is used with a starting learning rate of 0.001 and an exponential per-step decay equal to  $10^{-5}$ .

### D. JOINT 3D NAVIGATION OPERATOR

A joint 3D tracing operator is proposed to flexibly position the electrode array through the complex 3D tubular structure. As illustrated in Fig. 7, the 3D navigation operator is composed of three elements: 1) the bird-view axis ( $Y$ ) to monitor the width variations in a tube, 2) the sagittal axis ( $Z$ ) to identify the height of a tube, and 3) a navigation vector  $\vec{n} \vec{a} \vec{v}^k$ . Bird view ( $Y$ ) and sagittal ( $Z$ ) axes are jointly



**FIGURE 7.** Illustrating the joint 3D navigation operator. The  $\overline{n\vec{a}v}^k = [\theta, \varphi]$  is formed by identifying the navigation indicators from the sagittal and the bird-view projections. In this example, the navigation operator is shifted by  $\theta^\circ$  to the left and  $\varphi^\circ$  upward. The length of the navigator is defined by the minimum of collision probability of sagittal and the bird-view projections. The distance to the walls in both projections also give margins for shifting the  $\overline{n\vec{a}v}^k = [\theta, \varphi]$  to left-right and up-down considering the green dotted arrows according to the optimal safe zone.

connected at node  $O$  shown in Fig. 7 and form a unified structure that is rotated based on the assigned angles to the unity vector  $\overline{n\vec{a}v}^k$ . 3D space directions are indicated by considering two angles;  $\theta$  and  $\varphi$  around a unity vector  $\overline{n\vec{a}v}^k = [\theta, \varphi]$  in Fig. 7, where  $\theta$  describes the bird-view rotations around the  $Z$  axis and  $\varphi^\circ$  describes the sagittal rotations around the  $Y$  axis after being rotated by  $\theta^\circ$  around the  $Y$  axis. The length of the navigation vector  $\overline{n\vec{a}v}^k$  also defines the maximum allowable length that the electrode array that can be pushed inside the tubular structure (cochlea in this case) in each iteration and shown in Fig. 7. The navigation vector  $\overline{n\vec{a}v}^k$  can be shifted along the identified distances  $[(Dist - LW)$  and  $(Dist - RW)]$  to the cochlea walls from the origin ( $O$ ) in both sagittal and bird-view projections. All the defined parameters in the joint 3D tracing operator introduce super-flexibility in different scenarios with highly precise tuning of the navigation of the electrodes.

#### IV. EXPERIMENTAL SETUP AND RESULTS

This section focuses on the presentation and discussion of the results, mainly using the metric-based experimental setup, CI insertion in noisy scenarios and eventually the navigation indicators prediction in a real cochlea model.

##### A. REGRESSION AND CLASSIFICATION RESULTS

In this section, the quantitative and qualitative results of the 2.5D CNN are discussed. The 2.5D CNN addresses the regression network for estimation of the tracing angles along sagittal or bird-view planes ( $\theta, \varphi$ ). To quantify the regression performance two metrics are used: root-mean-squared error (RMSE) and explained variance ratio (EVA). RMSE measures the average magnitude of the prediction

error, indicating how close the observed values  $\alpha$  are to

**TABLE 1.** Average quantitative results on cochlea models ( $Anatom_{model1} \dots Anatom_{model3}$ ): EVA and RMSE are computed on the  $(Dist - LW)$ ,  $(Dist - RW)$  and the tracing angles along sagittal or top views ( $\theta, \varphi$ ), while Avg. accuracy and F-1 score are evaluated on the collision prediction task. Despite being relatively lightweight in terms of number of parameters, 2.5D CNN maintains a very good performance on both tasks.

Model	EVA	RMSE	Avg. accuracy	F-1 score	Num. Layers	FPS*
AlexNet	0.63	0.128	88.2%	0.80	8	23
ResNet-50	0.81	0.067	97.8%	0.93	50	7
VGG-16	0.73	0.109	92.7%	0.84	16	12
2.5D CNN	0.76	0.078	95.4%	0.92	8	20

\* Processing time in frames per second (fps).

those estimated by the network  $\hat{\alpha}$ :

$$RMSE = \sqrt{\frac{1}{N} \sum_{j=1}^N (\hat{\alpha}_j - \alpha_j)^2}. \quad (3)$$

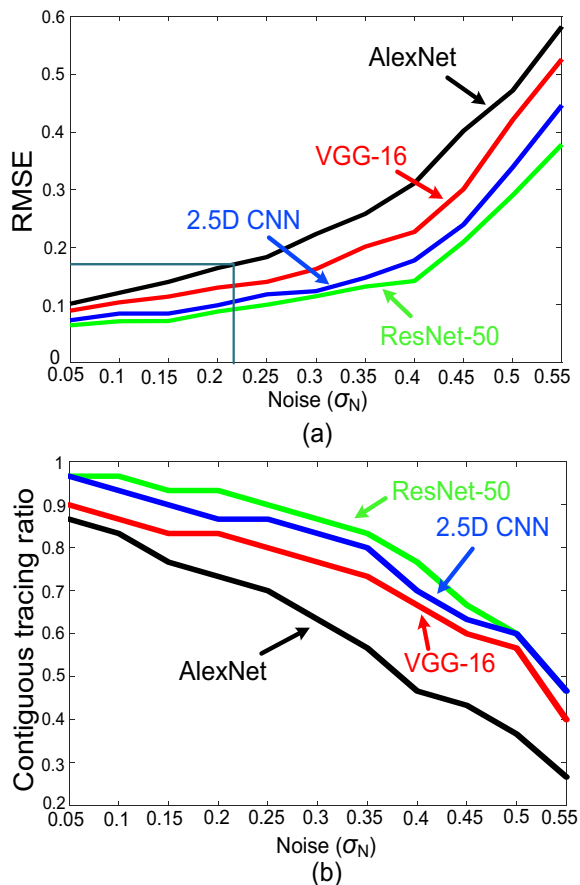
The EVA measures the proportion of variation in the predicted values with respect to those of the observed values. Such variations are given by the variance of the residuals  $Var = (\hat{\alpha} - \alpha)$  and the variance of the observed values  $Var = (\alpha)$ .

$$EVA = 1 - \frac{Var(\hat{\alpha} - \alpha)}{Var(\alpha)}. \quad (4)$$

If predicted values approximate the observed values well, the residual variance will be less than the total variance, resulting in  $EVA \lesssim 1$ . Otherwise, the residual variance will be equal or greater than the total variance, producing  $EVA = 0$  or  $EVA < 0$ , respectively. To assess the performance on collision prediction (*Collision*), the distance of the electrode to the left ( $Dist - LW$ ) and right ( $Dist - RW$ ) walls, average classification accuracy and F-1 score are used. It should be noted that training of the 2.5D CNN used the combination of the synthetic data generated by the parametric model explained in Section III-A.1 and the projection of the synthetic cochlea models ( $Synth_{model1} \dots Synth_{model3}$ ). Using the parametric synthetic data generation and synthetic cochlea models ( $Synth_{model1} \dots Synth_{model3}$ ), the sagittal and bird view networks were trained by over 1 million 2D cochlear segments with different width, length, inner and outer arcs and rotation directions. The trained networks have high generalization capability to data variation and are able to perform electrode navigation for unseen cochlea cases from different patients.

The generated data were divided into a training set containing 70% percent of the data to optimize the parameters and the hyperparameters, and the testing set consisting of the remaining 30% to evaluate the 2.5D CNN performance on the unseen data. The whole network is then examined on the anatomical models ( $Anatom_{model1} \dots Anatom_{model3}$ ) with manually defined ground-truth to quantify traversal performance. The tracing process begins by defining a sampling cube around the seed point in the scala tympani. Having sampled a segment of 3D cochlear, it





**FIGURE 8.** (a) The calculated RMSE in mapping of  $(Dist - LW)$  to the ground truth as a function of noise compiled for the anatomical models ( $Anatom_{model1} \dots Anatom_{model3}$ ). (b) Ratio of successful iterations completed by 2.5D CNN as a function of noise compared with ResNet-50, VGG-16 and AlexNet.

is projected onto the bird and sagittal views and sent to the 2.5D CNN. The sampling cube is rotated and adjusted based on the latest tracing information  $[\theta, \varphi]$  for sampling the next cochlea segment. This process continues to the last segment and sampling iterations along the cochlea and is user controlled.

Table I compares the average performance of cochlea models ( $Anatom_{model1} \dots Anatom_{model3}$ ) between the 2.5D CNN against other architectures from the literature [40], [43]-[44]. From these results, it is observed that the 2.5D CNN, even though 70 times smaller than the best architecture (ResNet-50), maintains considerable prediction performance while achieving real-time operation. Furthermore, the comparison against the VGG-16 architecture indicates the advantages in terms of generalization due to the residual learning scheme and parametric data generation model, as discussed in Section III.A.1 and Section III.B, respectively. The design succeeds at finding a good trade-off between tracing performance and the number of parameters detailed in the CNN architecture as shown in Table I. In order to enable the placement of an electrode array to promptly react to situational changes, it is necessary to reduce the network's latency as much as possible.

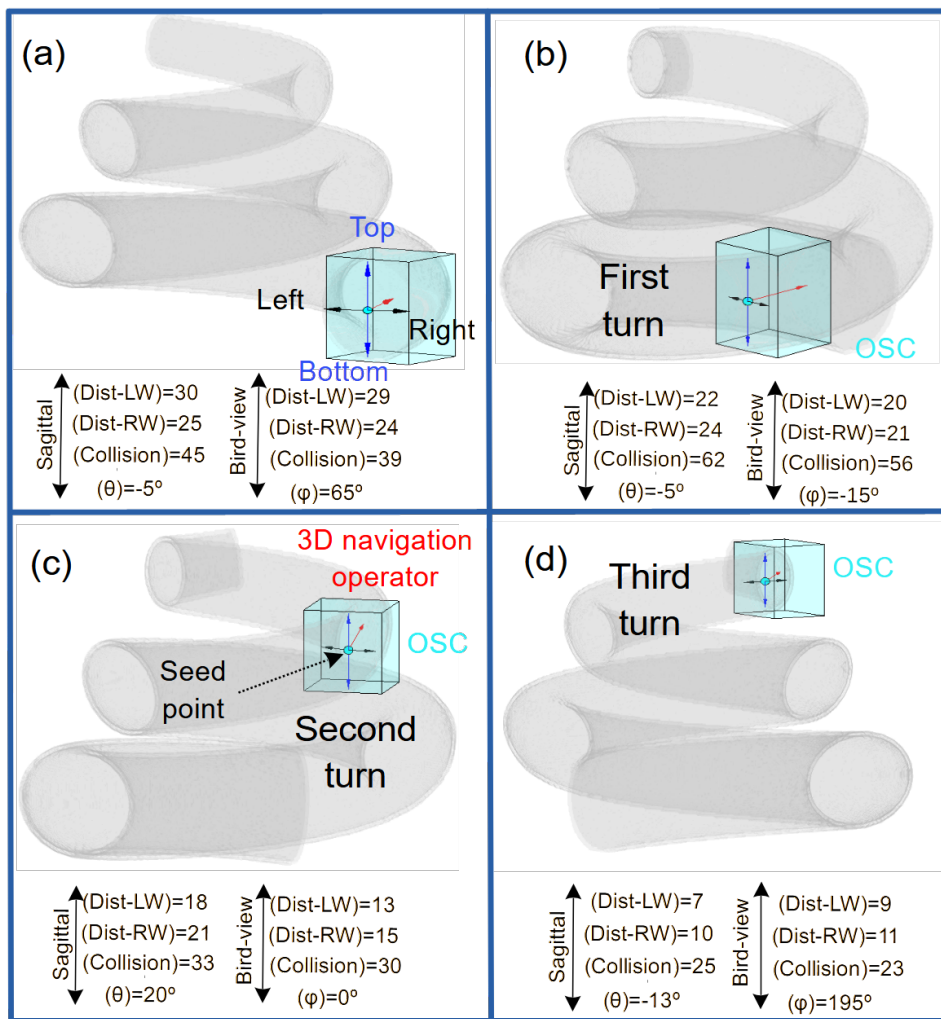
## B. DEEP MAPPING OF NAVIGATION INDICATORS IN NOISY SCENARIOS

Cochlea navigation is a difficult task, primarily because of the noise and variability associated with the real-world scenes. Computer vision has displayed a promising performance and flexibility when dealing with high degrees of noise and variability. This is because unlike most of the iterative methods where the search of true direction is determined based on a local estimate of the orientation and history information, the proposed and other CNN methods consider the whole feature map and the outline of the images (i.e. the borders). Typically, the added noise corrupts the process of mapping cochlea images to the navigation indicators including the distance to the adjacent left and right walls, collision probability and heading angle, and results in either minor or major deviations from the ground truth. The results in Fig. 8(a) show  $RMSE < 0.1$  for noise standard deviation ( $\sigma_N$ ) of  $0 < \sigma_N < 0.25$ . 2D gaussian noise was embedded to the generated images and used for deep extraction of the navigation indicators in noisy situations. Fig. 8(a) shows that for  $\sigma_N < 0.22$ , the average RMSE of  $(Dist - LW)$  (or  $(Dist - RW)$ ) in cochlea models ( $Anatom_{model1} \dots Anatom_{model3}$ ) is below 0.18. Increased noise causes more variations on the border information of the projected cochlea segments. This can be seen as a stream of images with localized amplitude variations which makes the border recognition extremely difficult. For  $\sigma_N > 0.4$ , the RMSE of all algorithms increase at a higher rate. Fig. 8(a) also shows that the ResNet-50 always shows higher noise robustness for  $0.05 < \sigma_N < 0.55$ .

Contiguous tracing which is the ratio of successful trials in tracing centerlines in all trials is calculated and shown in Fig. 8(b) for  $0.05 < \sigma_N < 0.55$ . The contiguous ratio analysis considers the randomness of the 2D noise distribution. The graphs are computed from a total number of 30 trials for the cochlea models ( $Anatom_{model1} \dots Anatom_{model3}$ ). For the 2.5D CNN, the tested cochlea models are traversed contiguously because the designed ResNet architecture helps with generalization of the border recognition in the image segments. In Fig. 8(b), the ratio of successfully traced centerline coordinates by the ResNet-50 algorithm are higher compared to 2.5D CNN but has about 3X longer execution time.

## C. JOINT-VIEW PROJECTION AND NAVIGATION: STEP-BY-STEP STUDY

Fig. 9 shows the qualitative results of progressive projection and tracing in  $Anatom_{model3}$ , its corresponding 3D operators and the identified indicators. An oriented sampling cube (OSC), which is a tight fit around 3D point in local space, is generated at four different locations of the  $Anatom_{model3}$  to show the performance of the 2.5D CNN. The considered locations capture almost all the geometrical difficulties along the navigation path (i.e. width and height variations, rotations along Z axis etc.). Fig. 9(a) is the start of the navigation and location of the OSC around the scala tympani seed point, so seed point x-y-z coordinates are set to the center of the OSC. 3D sampled points obtained from the



**FIGURE 9.** Automated tracing along a 3D cochlea using the *Anatom\_model13*. (a), (b), (c) and (d) represent the shifted OSC shown by cyan color along the *Anatom\_model13*. The superimposed OSC along cochlear samples different geometrical complexities at different turns. In (a), the OSC is placed around the scala tympani seed point, the sampled 3D cochlea segment is projected into the orthogonal sagittal and bird-view planes. The navigation indicators for both views are derived in two different columns below the projections. For example, sagittal view of 3D tracing algorithm starts from the seed point with  $\theta = -5^\circ$ ,  $Collision = 45$ ,  $Dist - LW = 30$  and  $Dist - RW = 25$ . Rotated joint 3D operators are also superimposed in each OSC for different scenarios. The derived navigation indicators {Top, Bottom, Left and Right} are shown in (a), (b), (c) and (d).

input depth image are projected onto x-y and y-z planes of the coordinate system, respectively. Notice that the projections on the three orthogonal planes may be coarse because of the resolution of the depth map [45], which can be improved by performing median filter and opening operation on the projected images. The designed CNNs for each view then process and map the input projections into the navigation's indicators. For the identified indicators including  $(Dist - LW)$ ,  $(Dist - RW)$  and  $(Collision)$  in each view, the distances from the left, right and the frontal walls are normalized between 0 and 64 (6 neurons to quantize 64 steps, with nearest points set to 0 and farthest points set to 64). The navigation angles ( $\theta$  and  $\varphi$ ) are also indicated by two numbers.

By fusing the computed navigation indicators from both sagittal and bird-view projections, a 3D joint operator is finally formed as shown in Fig. 9(a)-(d). The superimposed 3D navigators in each figure consists of blue and black arrows to quantify the height and width of the sampled

cochlear respectively. The red arrow also shows the optimal navigation path. For example, the navigation parameter for the OSC samples around the scala tympani seed point, the  $(Dist - LW)$ ,  $(Dist - RW)$  of both views are  $(Dist - LW/RW)_{Sagittal} = 30/25$  and  $(Dist - LW/RW)_{Bird-view} = 29/24$ .  $(Collision)$  which shows the length of  $\bar{n}\bar{a}\bar{v}^k$  is also set to 39, the minimum of collision in both views  $\theta$  and  $\varphi$  are also set to  $-5^\circ$  and  $115^\circ$ . This process is then repeated for four different locations by moving the OSC along cochlea as shown in Fig. 9(a)-(d). 3D depth sampling is obtained by rotating the OSC to the identified  $\theta$  and  $\varphi$  of the previous step (i.e., the  $\theta$  and  $\varphi$  history). The height, width and depth of OSC are also defined according to the derived information in the previous step [e.g.,  $(Dist - LW)$ ,  $(Dist - RW)$  and the minimum of collision probability in both views  $\theta$  and  $\varphi$ ]. This is an automated and reliable depth sampling that converts the whole cochlea to the smaller segments. The size-adjusted OSC rotates along



the cochlea; the sagittal and bird-views also rotate accordingly to capture the projections. The identified

## V. CONCLUSION

In this paper, 2.5D CNN is proposed to map the projected 2D cochlea images into accurate navigation indicators, including the distance to the adjacent left and right walls, collision probability and heading angle. A novel network architecture was designed (i.e. converting a 3D to two complementary networks) to trade off performance for processing time to enable online operation. Each network consists of 5 dense convolutional layers with  $\{(12 \times 12) \dots (96 \times 96)\}$  kernels and LeakyReLU activations, followed by just one average pooling, with size equal to the size of final feature maps and three dense layers. The training was performed by minimizing the categorical cross entropy with the Adam optimizer. Tracing of the cochlea is a laborious and dangerous task as there exist infinitesimal error margin. The proposed method learns to promptly react to the radical directional changes, geometrical variations and overall rotations along the cochlea. It was shown through extensive evaluations on processing time, navigation accuracy and noise robustness analysis that the proposed approach performs well with both synthetic MATLAB-generated images and anatomical cochlea models constructed from micro-CT images. The results confirm reliable navigation with an average of  $>98\%$  mapping accuracy. The processing time of the navigation platform which consists of 3D segment sampling, 2.5D projections, navigation indicators extraction and eventually the remapping to 3D navigators is 100 ms per insertion step. Where there are local noise and artefacts, the feature map activations clearly recognize the edges of the of the generated images by the parametric model. Future work will focus on integrating the proposed navigation method into a robotic arm with a real-time imaging module to implement a precise computer-aided system for virtual cochlear surgery.

## ACKNOWLEDGEMENT

The authors thank Dr. E. Salkim for advice on 3D cochlea model reconstruction using micro-CT images, and Advanced Bionics for providing the micro-CT images for the anatomical cochlea models.

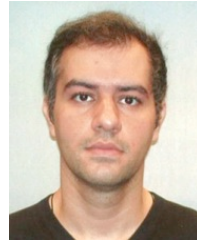
## REFERENCES

- [1] B. S. Wilson and M. F. Dorman, "Cochlear implants: A remarkable past and a brilliant future," *Hear Res.*, vol. 242, no. 1–2, pp. 3–21, Aug. 2008.
- [2] J. Wouters, H. J. McDermott, and T. Francart, "Sound coding in cochlear implants: From electric pulses to hearing," *IEEE Signal Process. Mag.*, vol. 32, no. 2, pp. 67–80, Mar. 2015.
- [3] Dorman, M.F., Loizou, P.C., & Rainey, D. (1997). *J. Acoustical Society of America*, 102, pp. 2993–2996, 1997.
- [4] E. Erixon, H. Hogstorp, K. Wadin, and H. Rask-Anderson, "Variational anatomy of the human cochlea: Implications for cochlear implantation," *Otol. Neurotol.*, vol. 30, no. 1, pp. 14–22, Jan. 2008.
- [5] G. B. Wanna, J. H. Noble, M. L. Carlson et al., "Impact of electrode design and surgical approach on scalar location and cochlear implant outcomes," *Laryngoscope*, vol. 124, Supplement 6, pp. S1–S7, 2014.
- [6] A. Hussong, T. S. Rau, T. Ortmaier, B. Heimann, T. Lenarz, and O. Majdani, "An automated insertion tool for cochlear implants: Another

coordinates and the 3D operator present the optimal navigation tool for surgical purposes.

- step towards atraumatic cochlear implant surgery," *Int. J. Comput. Assist. Radiol. Surg.*, vol. 5, no. 2, pp. 163–171, Mar. 2010.
- [7] E. Salkim, M. Zamani, D. Jiang, and A. Demosthenous, "Detection of Electrode Proximity to the Cochlea Wall Based on Impedance Variation: a Preliminary Computational Study," in *Proceedings of UKSim-AMSS 22nd International Conference on Modelling & Simulation* (Cambridge, UK:), 10.5013/IJSSST.a.21.02.13.
- [8] E. Salkim, M. Zamani, D. Jiang, S.R. Saeed, and A. Demosthenous, "Insertion Guidance Based on Impedance Measurements of a Cochlear Electrode Array," *Frontiers in Computational Neuroscience.*, vol. 16, pp. 1–14, 2022.
- [9] P. Aebischer et al., "Intraoperative impedance-based estimation of cochlear implant electrode array insertion depth," *IEEE Trans. Biomed. Eng.*, vol. 68, no. 2, pp. 545–555, Feb. 2021
- [10] C. K. Giardina, E.S. Krause, K. Koka and D.C. Fitzpatrick, "Impedance measures during in vitro cochlear implantation predict array positioning," *IEEE Trans. Biomed. Eng.*, vol.65, no. 2, pp. 327–335, Feb. 2018.
- [11] C. T. Tan et al., "Real-time measurement of electrode impedance during intracochlear electrode insertion," *Laryngoscope*, vol. 123, no. 4, pp. 1028–1032, 2013.
- [12] J. Pile et al., "Detection of modiolar proximity through bipolar impedance measurements", *Laryngoscope*, vol. 127, pp. 1413-1419, 2016.
- [13] J. Anso et al., "Electrical Impedance to Assess Facial Nerve Proximity During Robotic Cochlear Implantation," *IEEE Trans. Biomed. Eng.*, vol. 66, no. 1, pp. 237–245, Jan. 2019.
- [14] J. Pile and N. Simaan, "Modeling, Design, and Evaluation of a Parallel Robot for Cochlear Implant Surgery," *IEEE/ASME Trans. Mechatronics*, vol. 19, no. 6, pp. 1746-1755, Dec. 2014, doi: 10.1109/TMECH.2014.2308479.
- [15] N. Mangado, M. Ceresa, N. Duchateau, H. M. Kjer, S. Vera, H. Dejea Velardo, et al., "Automatic model generation framework for computational simulation of cochlear implantation", *Annals of Biomedical Engineering*, vol. 44, no. 8, pp. 2453-2463, Aug 2016.
- [16] X. Meshik, T. A. Holden, R. A. Chole, and T. E. Hullar, "Optimal Cochlear Implant Insertion Vectors.," *Otology & Neurology*, vol. 31, no. 1, pp. 58–63, Jan. 2010.
- [17] T. Bruns et al., "Real-time localization of cochlear-implant electrode arrays using bipolar impedance sensing," *IEEE Trans. Biomed. Eng.*, vol. 69, no. 2, pp. 718–724, Feb. 2022.
- [18] K. S. Min, S. B. Jun, Y. S. Lim, S. I. Park, and S. J. Kim, "Modiolus-hugging intracochlear electrode array with shape memory alloy," *Comput. Math. Methods Med.*, vol. 2013, Article ID 250915, 2013.
- [19] A. Rodriguez, D. Ehlenberger, D. Dickstein, P. Hof, and S. Wearne, "Automated three-dimensional detection and shape classification of dendritic spines from fluorescence microscopy images," *PLoS ONE*, vol. 3, no. 4, p. e1997, 2008.
- [20] B. Al-Diri, A. Hunter, and D. Steel, "An active contour model for segmenting and measuring retinal vessels," *IEEE Trans. Med. Imaging*, vol. 28, no. 9, pp. 1488–1497, Sep. 2009.
- [21] A. M. Mendonca and A. Campilho, "Segmentation of retinal blood vessels by combining the detection of centerlines and morphological reconstruction," *IEEE Trans. Med. Imaging*, vol. 25, no. 9, pp. 1200–1213, Sep. 2006.
- [22] M. Sofka and C. V. Stewart, "Retinal vessel centerline extraction using multiscale matched filters, confidence and edge measures," *IEEE Trans. Med. Imag.*, vol. 25, no. 12, pp. 1531–1546, Dec. 2006.
- [23] M. Erdt, M. Raspe, and M. Suehling, "Automatic hepatic vessel segmentation using graphics hardware," *Medical Imaging and Augmented Reality*, vol. 5128 of LNCS, pp. 403–412, Springer Berlin / Heidelberg, 2008.
- [24] A. Frangi, W. Niessen, K. L. Vincken, and M. A. Viergever. Multiscale vessel enhancement filtering. *MICCAI'98*, pages 130–137, 1998.
- [25] O. Friman, M. Hindennach, C. Kuhnel, and H.-O. Peitgen, "Multiple hypothesis template tracking of small 3D vessel structures," *Med. Imag. Anal.*, vol. 14, no. 2, pp. 160–171, Apr. 2010.
- [26] M. M. G. Macedo, C. Mekkaoui, and M. Jackowski, I. Bloch and R. M. Cesar, Eds., "Vessel centerline tracking in CTA and MRA images

- using Hough transform,” in CIARP, ser. Lecture Notes in Comput. Sci., Springer, 2010, vol. 6419, pp. 295–302.
- [27] J. Tyrrell, E. di Tomaso, D. Fuja, R. Tong, K. Kozak, R. Jain, and B. Roysam, “Robust 3-D modeling of vasculature imagery using superellipsoids,” *IEEE Trans. Med. Imag.*, vol. 26, no. 2, pp. 223–237, Feb. 2007.
- [28] H. Shen, B. Roysam, C. V. Stewart, J. N. Turner, and H. L. Tanenbaum, “Optimal scheduling of tracing computations for real-time vascular landmark extraction from retinal fundus images,” *IEEE Trans. Inform. Technol. Biomed.*, vol. 5, no. 1, pp. 77–91, Mar. 2001.
- [29] G. Lin, C. V. Stewart, B. Roysam, K. Fritzsche, G. Yang, and H. L. Tanenbaum, “Predictive scheduling algorithms for real-time feature extraction and spatial referencing: Application to retinal image sequences,” *IEEE Trans. Biomed. Imag.*, vol. 51, no. 1, pp. 115–125, Jan. 2004.
- [30] C.-L. Tsai, C. V. Stewart, H. L. Tanenbaum, and B. Roysam, “Model-based method for improving the accuracy and repeatability of estimating vascular bifurcations and crossovers from retinal fundus images,” *IEEE Trans. Inf. Technol. Biomed.*, vol. 8, no. 2, pp. 122–130, Jun. 2004.
- [31] M. Zamani, E. Salkim, S. R. Saeed and A. Demosthenous, “A Fast and Reliable Three-Dimensional Centerline Tracing: Application to Virtual Cochlear Implant Surgery,” *IEEE Access*, vol. 8, pp. 167757–167766, 2020.
- [32] M. Schneider, S. Hirsch, B. Weber, G. Székely, and B. H. Menze, “Joint 3-D vessel segmentation and centerline extraction using oblique Hough forests with steerable filters,” *Med. Image Anal.*, vol. 19, no. 1, pp. 220–249, 2015.
- [33] A. Sironi, E. Turetken, V. Lepetit, and P. Fua, “Multiscale centerline detection,” *IEEE Trans Pattern Analysis and Machine Intelligence*, vol. 38, no. 7, pp. 1327–1341, 2016.
- [34] M. A. Gulsun, G. Funka-Lea, P. Sharma, S. Rapaka, and Y. Zheng, “Coronary centerline extraction via optimal flow paths and CNN path pruning,” *Proc. MICCAI*, Athens, Greece, 2016, pp. 317–325.
- [35] J. M. Wolterink, R. W. van Hamersvelt, M. A. Viergever, T. Leiner and I. Isgum, “Coronary artery centerline extraction in cardiac CT angiography using a CNN-based orientation classifier,” *Med. Image Anal.*, vol. 51, pp. 46–60, 2019.
- [36] L. Yu, J.-Z. Cheng, Q. Dou, X. Yang, H. Chen, J. Qin, and P.-A. Heng, “Automatic 3D cardiovascular MR segmentation with densely-connected volumetric convnets,” *Proc. MICCAI*, Quebec, Canada, 2017, pp. 287–295.
- [37] K. Bin et al., “Learning tree-structured representation for 3D coronary artery segmentation,” *Comput. Med. Imaging Graph.*, vol. 80, pp. 101688, 2020.
- [38] F. Yu and V. Koltun, “Multi-scale context aggregation by dilated convolutions,” *arXiv:1511.07122*, 2015.
- [39] S. R. Stock, *Microcomputed tomography: Methodology and applications*: CRC Press, 2008.
- [40] S. Zhi, Y. Liu, X. Li, and Y. Guo, “Toward real-time 3d object recognition: a lightweight volumetric CNN framework using multitask learning,” *Comput. Graph.*, vol. 71, pp. 199–207, 2018.
- [41] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [42] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, “Curriculum learning,” *Proc. 26th Annual Int. Conf. on Machine Learning*. ACM, 2009, pp. 41–48.
- [43] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In NIPS, 2012.
- [44] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [45] W. Li, Z. Zhang, and Z. Liu. Action recognition based on a bag of 3d points. In CVPR Workshops, 2010.



**MAJID ZAMANI** (Member, IEEE) received the M.Sc. degree in microelectronics from Islamic Azad University, Science and Research Branch, Tehran, in 2011, and the Ph.D. degree from University College London (UCL), London, U.K., in 2017.



He was a Research Associate with the Bioelectronics Group, UCL. His research interests include design and fabrication of advanced and energy efficient computational systems utilizing pattern recognition, machine learning, and computer vision algorithms, especially for wearable and implantable biomedical applications. He was a recipient of the Oversea Research Scholarship and the UCL Graduate Research Scholarship to pursue his Ph.D. degree. He was also a recipient of the Best Researcher M.Sc. Student Award.

**ANDREAS DEMOSTHENOUS** (S’94–M’99–SM’05–F’18) received a B.Eng. degree in electrical and electronic engineering from the University of Leicester, Leicester, U.K., a M.Sc. degree in telecommunications technology from Aston University, Birmingham, U.K., and a Ph.D. degree in electronic and electrical engineering from University College London (UCL), London, U.K., in 1992, 1994, and 1998, respectively. He is currently a Professor with the Department of Electronic and Electrical Engineering, UCL, and leads the Bioelectronics Group. He has made outstanding contributions to improving safety and performance in integrated circuit design for active medical devices, such as spinal cord and brain stimulators. He has numerous collaborations for cross-disciplinary research, both within the U.K. and internationally. He has authored over 30 articles in journals and international conference proceedings, several book chapters, and holds several patents. His research interests include analog and mixed-signal integrated circuits for biomedical, sensor, and signal processing applications.

Dr. Demosthenous is a fellow of the Institution of Engineering and Technology and a Chartered Engineer. He was a co-recipient of a number of Best Paper Awards and has graduated many Ph.D. students. He was an Associate Editor from 2006 to 2007 and the Deputy Editor-in-Chief from 2014 to 2015 of the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS II: EXPRESS BRIEFS, and an Associate Editor from 2008 to 2009 and the Editor-in-Chief from 2016 to 2019 of the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS I: REGULAR PAPERS. He is an Associate Editor of the IEEE TRANSACTIONS ON BIOMEDICAL CIRCUITS AND SYSTEMS and serves on the International Advisory Board of Physiological Measurement. He has served on the technical committees for a number of international conferences, including the European Solid-State Circuits Conference and the International Symposium on Circuits and Systems.