# Automatic View Planning with Multi-scale Deep Reinforcement Learning Agents

Amir Alansary, Loic Le Folgoc, Ghislain Vaillant, Ozan Oktay,
Yuanwei Li, Wenjia Bai, Jonathan Passerat-Palmbach, Ricardo Guerrero,
Konstantinos Kamnitsas, Benjamin Hou, Steven McDonagh, Ben Glocker,
Bernhard Kainz, and Daniel Rueckert

Imperial College London,UK
{a.alansary14}@imperial.ac.uk

**Abstract.** We propose a fully automatic method to find standardized view planes in 3D image acquisitions. Standard view images are important in clinical practice as they provide a means to perform biometric measurements from similar anatomical regions. These views are often constrained to the native orientation of a 3D image acquisition. Navigating through target anatomy to find the required view plane is tedious and operator-dependent. For this task, we employ a multi-scale reinforcement learning (RL) agent framework and extensively evaluate several Deep $Q$-Network (DQN) based strategies. RL enables a natural learning paradigm by interaction with the environment, which can be used to mimic experienced operators. We evaluate our results using the distance between the anatomical landmarks and detected planes, and the angles between their normal vector and target. The proposed algorithm is assessed on the mid-sagittal and anterior-posterior commissure planes of brain MRI, and the 4-chamber long-axis plane commonly used in cardiac MRI, achieving accuracy of 1.53mm, 1.98mm and 4.84mm, respectively.

## 1 Introduction

In medical imaging, obtaining accurate biometric measurements that are comparable across populations is essential for diagnosis and supporting critical decision making. For this purpose, standard view planes through a defined anatomy are commonly used in clinical practice to establish comparable metrics. Finding these planes in an imaging examination through a 3D volume is slow and suffers from inter-observer variability. The neuro-imaging community defines a standard (axial) image plane by adopting the anterior-posterior commissure (ACPC) line. Transforming an image to the ACPC coordinate system includes a number of steps: (i) marking the AC point, (ii) obtaining the optimal view of the ACPC and the mid-sagittal plane, and (iii) marking the PC point. Accurate detection of the mid-sagittal plane is useful for the initial step in image registration [1]. It is also used in evaluation of pathological brains by estimating the departures from bilateral symmetry in the cerebrum [12]. Similarly, in cardiac MRI standard views are used to assess anomalies. Because of the complexity of cardiac

anatomy, the appearance of relevant structures can exhibit large variance according to the positioning of the imaging plane. During conventional cardiac MRI acquisition, the localization of short and long-axis of the heart requires a multi-step approach that involves double-oblique slices, exhibiting both inter and intra-observer variance [6]. These steps include: (i) whole 3D pilot image acquisition, (ii) left ventricle (LV) localization, (iii) short axis orientation, (iv) 3-chamber view calculation, (v) landmark detection in mid-ventricular slices, and (vi) 4- and 2-chamber view calculation.

In this work, we aim to automate the view planning process by using reinforcement learning (RL) where an agent learns to make comprehensive and sensible decisions by mimicking navigation processes as outlined above, in a manner that allows medical experts to gain confidence in fully automatic methods. RL constitutes a sub-field of machine learning concerned with how agents take actions in an environment. In contrast to supervised learning, RL involves learning by interacting with an environment instead of using a set of labeled examples that is typically provided by a knowledgeable supervisor. This learning paradigm allows RL agents to learn complex tasks that may need several steps to find a solution [13]. Mnih et al. [9] adopted a deep convolution network for RL function approximation, known as the Deep $Q$-Network (DQN), achieving human-level performance in a suite of Atari games. Recently, DQN has shown promising results when employed in related applications in the medical imaging domain. Ghesu et al. [3] introduced an automatic landmark detection approach using a DQN-agent to navigate in 3D images with fixed step actions. Maicas et al. [7] proposed a similar method for breast lesion detection using actions to control the location and size of the bounding box. Liao et al. [5] presented an image registration approach using actions to explore transformation parameters. We adopt different DQN-based architectures as a solution for the proposed RL formulation of the view planning task.

**Related Work:** Ardekani et al. [1] proposed a method to automatically detect the mid-sagittal plane in 3D brain images by maximizing the cross-correlation between the two image sections on either side of the sought plane. Stegmann et al. [12] proposed to use a sparse set of profiles in the plane normal direction and maximize the local symmetry around them. In [4,6], they proposed an automatic view planning algorithm for cardiac MRI acquisition. Their methods are based on learning the anatomy segmentation and detecting anchor landmarks in order to calculate standard cardiac views. These methods require prior knowledge of the whole 3D image for the purpose of plane detection. This involves manual annotation of anatomical landmarks, which is a tedious and time-consuming task. In our method, we use the acquired standardized views for cardiac scans in training without any manual labeling.

**Contribution:** We propose a novel RL-based approach for fully automatic standard view plane detection from volumetric MRI data. The proposed model follows a multi-scale search strategy with hierarchical action steps in a coarse-to-fine fashion. By sequentially updating plane parameters, our algorithm is able to reach the target plane. We run extensive experiments for evaluating different

DQN baselines on detecting 3 different planes. Applications of our method to brain and cardiac MRI data show a target plane detection in real time with accuracy around 2 and 5 mm, respectively.

## 2   Background

An RL agent learns by interacting with an environment, $E$. At every state, $s$, a single decision is made to choose an action, $a$, from a set of multiple discrete actions, $A$. Each valid action choice results in an associated scalar reward, defining the reward signal, $R$. The agent attempts to learn a policy to maximize both immediate and subsequent future rewards (optimal policy).

$Q$**-Learning:** The optimal action-selection policy can be identified by learning a state-action value function, $Q(s, a)$ [16]. The $Q$-function is defined as the expected value of the accumulated discounted future rewards $E[r_{t+1} + \gamma r_{t+2} + \cdots + \gamma r_{t+n}|s, a]$. $\gamma \in [0, 1]$ is a discount factor that represents the uncertainty in the agent's environment and is used to weight future rewards accordingly. This value function can be unrolled recursively (using the Bellman Equation [2]) and can thus be solved iteratively: $Q_{i+1}(s, a) = E\left[r + \gamma \max_{a'} Q_i(s', a')\right]$.

**Deep $Q$-Learning:**  Mnih et al. [9] proposed the Deep $Q$-Network (DQN) and implemented a standard $Q$-learning algorithm with the addition of approximating the $Q$-function using a ConvNet, $Q(s, a) \approx Q(s, a; \omega)$, where $\omega$ represents the network's parameters. The DQN loss function is defined as:

$$L(\omega) = E\left[\left(r + \gamma \max_{a'} Q_{target}(s', a'; \omega^-) - Q_{net}(s, a; \omega)\right)^2\right],$$

Approximating the $Q$-function in this manner allows to learn from larger data sets using mini-batches. The DQN uses $Q_{target}(\omega^-)$, a fixed version of $Q_{net}(\omega)$, that is periodically updated. This is used to stabilize rapid policy changes, due to the quick variations in $Q$-values and the distribution of the data. Another problem that may cause divergence is successive data sampling. To avoid this, an experience replay memory that stores transitions of $(s_t, a_t, r_{t+1}, s_{t+1})$ is randomly sampled to create the mini-batches used for training. We outline below two recent state-of-the-art improvements to the standard DQN.

**Double DQN (DDQN):** It has been shown that DQN is susceptible to bias in noisy environments, where the target network may cause upward bias due to delayed updates. Van Hasselt et al. [14] proposed a solution that replaces the maximum approximated action from $Q_{target}(s', a'; \omega^-)$ with an action selected from the $Q_{target}(s', Q_{net}(s', a', \omega); \omega^-)$. This strategy is able to mitigate bias by decoupling the selected action from $Q_{target}$. DDQN improves the stability of learning, which can translate to the ability to learn more complicated tasks but may not necessarily improve the performance [14].

**Duel DQN:** Wang et al. [15] showed improved performance over the original DQN by defining two separate channels: *(i)* an action-independent value function $V(s)$ to provide an estimate of the value of each state, and *(ii)* an

action-dependent advantage function $A(s, a)$ to calculate potential benefits of each action. Both functions are then combined into a single action-advantage $Q$-function, $Q(s, a) = A(s, a) + V(s)$. Duel DQN may achieve more robust estimates of state value by decoupling it from specific actions, so $s$ could be more explicitly modelled, which yields higher performance in general.
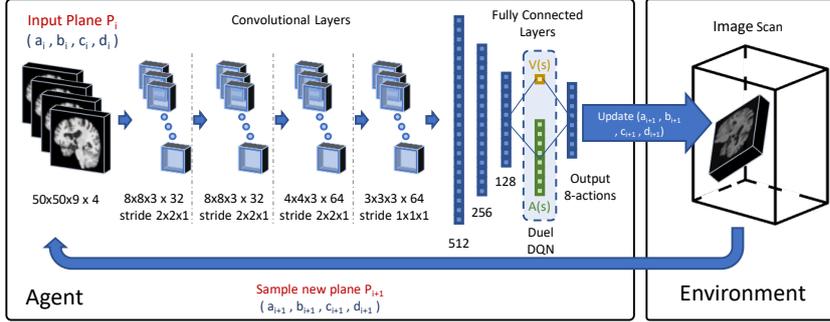
## 3    Method



Fig. 1: The pipeline of the proposed multi-scale RL agent. Initially, the environment samples a plane: $ax + by + cz + d = 0$, from the 3D image scan. The agent selects an action to update a single parameter for sampling the next plane. This process is repeated till the agent reaches a terminal state (oscillation).

A plane $P$, in the Cartesian coordinate frame of the 3D image, is defined as: $ax + by + cz + d = 0$. Where $(a, b, c)$ represent the normal direction (cosine) to this plane and $d$ is the distance of the plane from the origin. To automate standard view planning, we aim to find the appropriate parameterization of the target plane. We formulate our RL framework by defining the following elements:

- **States:** Our Environment $E$ is represented by a 3D scan and $s$ is a 3D region of interest that contains $P$. A frame history buffer is used for storing the last planes from previous steps to stabilize search trajectories and prevent getting stuck in repeated cycles. We choose a history size of 4 frames similar to [9].
- **Actions:** The agent interacts with $E$ by taking action steps $a \in A$ to modify the position parameters of the plane. The action space consists of eight actions, $\{\pm a_{\theta_x}, \pm a_{\theta_y}, \pm a_{\theta_z}, \pm a_d\}$, which update the plane parameters $a = cos(\theta_x + a_{\theta_x})$, $b = cos(\theta_y + a_{\theta_y})$, $c = cos(\theta_z + a_{\theta_z})$ and $d = d + a_d$.
- **Reward:** The RL reward function forms a proxy for the true task goal and care must be taken to capture exactly what this goal entails. In our problem instance, the difficulty comes from designing a reward that encourages the agent to move towards the target plane while still being learnable. With these considerations, we define the reward $R = \text{sgn}(D(P_{i-1}, P_t) - D(P_i, P_t))$, where

$D$ is a function to take the Euclidean distance between plane parameters. We further denote $P_i$ as the current predicted plane at step $i$, with $P_t$ the target ground truth plane. The difference of the parameter distances, between the previous and current steps, signifies whether the agent is moving closer to or further away from the desired plane parameters. $R \in \{+1, 0, -1\}$ provide the agent with a per step (non-sparse) reward signal, with zero-valued $R$ presents plane oscillations around the correct solution.

– **Terminal State:** The final state is defined as the state in which the agent finds the target plane $P_t$. A trigger action can be used to signal when the target state is reached [7]. However, adding extra actions increases the action space size, which may in turn increase the complexity of the task to be learned. The maximum number of interactions should also be defined in such a setting. We found that terminating the episode when oscillation is detected heuristically works in practice without the need to expand the action space. However, in contrast to [3], we choose the terminating action with the lower $Q$-value. We find that $Q$-values are lower when the target plane is closer. Intuitively, the DQN encourages awarding higher $Q$-values to actions when the current plane is far from the target.

**Multi-scale Agent:** In order to provide more structural information, we introduce a novel multi-resolution approach in a coarse-to-fine fashion with hierarchical action steps. In this scenario, $E$ samples a grid of a fixed plane size $(P_x, P_y, P_z)$ of voxels around the plane origin $P_o$ and initial spacing $(S_x, S_y, S_z)$ mm. Initially, the agent searches for the plane with higher action steps. Once the target plane is found, $E$ samples the new planes with smaller spacing and the agent uses smaller action steps. Coarser levels in the hierarchy provide additional guidance to the optimization process by enabling the agent to see larger context of the image. Whereas, finer scales provide sharper adjustments for the final estimation of the plane. Similarly, larger step actions speed up the solution towards the target plane, while smaller steps fine tune the final estimation of plane parameters. The same DQN is shared between all levels in the hierarchy, see Fig. 1. The next section exhibits results of utilizing this multi-scale approach.

## 4  Experiments and Results

The proposed algorithm is assessed using 12 different experiments; a combination of four different DQN-based methods with three target planes. We evaluate our results using the distance between anatomical landmarks and the detected planes. We also measure the orientation error by calculating the angle between normal vectors of the detected and target planes.

**Datasets:** A set of 832 isotropic 1mm MR scans were obtained from the ADNI database [10] to evaluate the proposed method. While, a subset of 728 and 104 images are used for training and testing. All brain images were skull stripped and affinely aligned to the MNI space, thus allowing ground truth planes to be extracted in the standard directions. For cardiac images, we use 455 short-axis cardiac MR of resolution $(1.25 \times 1.25 \times 2)$ mm obtained from the UK Digital

Heart Project [8]. A subset of 364 and 91 images are used for training and testing. ACPC planes are evaluated using the AC and PC landmarks for the distance error calculation. Similarly, we use the outer aspect, inferior tip and inner aspect points of splenium of corpus callosum for mid-sagittal planes. For cardiac MRI, we use six landmarks projected on the 4-chamber plane; the two right ventricle (RV) insertion points, right and left ventricles (LV) lateral wall turning points, apex, and the center of the mitral valve, See Fig. 2.



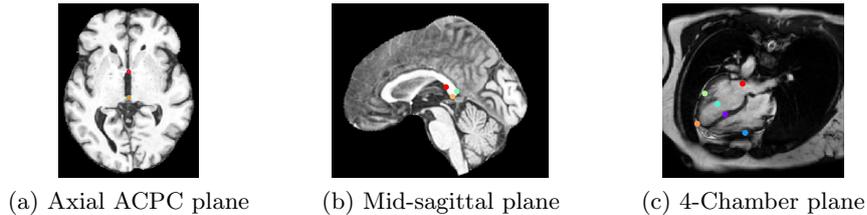| (a) Axial ACPC plane | (b) Mid-sagittal plane | (c) 4-Chamber plane |

Fig. 2: Ground truth planes from brain and cardiac MRI scans. (a) ACPC axial plane marking AC (red) and PC (yellow) points. (b) Mid-sagittal plane with outer aspect (green), inferior tip (yellow) and inner aspect (red) points of splenium of corpus callosum. (c) 4-Chamber view with the projected two RV insertion points(violet, green), RV and LV lateral wall turning points (blue, lime), apex (orange), and the center of the mitral valve (red).

**Experiments:** During training, a random point is sampled from the 3D input image. The initial random plane is then defined using the normal vector between the center of the image and the random point. The origin of this plane is the projected point of the center of the input image. Finally, a plane of size $(50, 50, 9)$ voxels is sampled around the plane origin with initial spacing $3 \times 3 \times 3$ mm. Initial $a_{\theta_x}, a_{\theta_y}, a_{\theta_z}$ equal 8 and $a_d$ equals 4. With every new scale $a_{\theta_x}, a_{\theta_y}, a_{\theta_z}$ decrease by a factor of 2 and $a_d$ decrease 1 unit. 3-levels of scale with spacing from 3 to 1 mm are used for the brain experiments, and 4-levels of scale from 5 to 2 mm for the cardiac experiment. For experiments on cardiac images, initial planes are sampled randomly from the 3D input image within 20% around the center of the image, to avoid sampling outside the field of view.

**Results:** During inference, the environment samples a plane and the agent updates sequentially new plane's parameters until reaching the terminal state. In order to have a fair comparison between different variants of the proposed method, we fix the initial plane for all models during evaluation. Table 1 shows the results from these comparative experiments. All methods share similar performance including speed and accuracy, and there is no unique winner for the best overall method. Best performing agents for detecting the mid-sagittal and ACPC planes achieve accuracy of $1.53\pm2.2$ mm and $2.44\pm5.04°$, and $1.98\pm2.23$ mm and $4.48 \pm 14.0°$, respectively. Where in cardiac, the task is more complex due to the lower quality and higher variability between different scans. The agent

has to navigate in a bigger field of view compared to brain images. Thus Duel DQN-based architectures achieve the best results for detecting the 4-chamber plane with $4.84 \pm 3.03$ mm and $8.86 \pm 12.42°$ accuracy, as a result from learning a better state value function by decoupling it from action-value function. These results are better than the state-of-the-art [6], which achieves an accuracy of $5.7 \pm 8.5$mm and $17.6 \pm 19.2°$. Unlike [6], our method does not require manual annotation of landmarks. More visualization results are published on our github.

Table 1: Results of our multi-scale RL agent detecting 3 different MRI planes.

| Model | Mid-sagittal brain | | ACPC brain | | 4-Chamber cardiac | |
|---|---|---|---|---|---|---|
| | $e_d(mm)$ | $e_\theta(°)$ | $e_d(mm)$ | $e_\theta(°)$ | $e_d(mm)$ | $e_\theta(°)$ |
| DQN | $1.65 \pm 1.99$ | $2.42 \pm 5.27$ | $2.61 \pm 5.44$ | $\mathbf{3.23 \pm 6.03}$ | $5.61 \pm 4.09$ | $10.16 \pm 10.62$ |
| DDQN | $2.08 \pm 2.58$ | $3.44 \pm 7.46$ | $\mathbf{1.98 \pm 2.23}$ | $4.48 \pm 14.00$ | $5.79 \pm 4.58$ | $11.20 \pm 14.86$ |
| Duel DQN | $1.69 \pm 1.98$ | $3.82 \pm 7.15$ | $2.13 \pm 1.99$ | $5.24 \pm 13.75$ | $\mathbf{4.84 \pm 3.03}$ | $8.86 \pm 12.42$ |
| Duel DDQN | $\mathbf{1.53 \pm 2.20}$ | $\mathbf{2.44 \pm 5.04}$ | $5.30 \pm 11.19$ | $5.25 \pm 12.64$ | $5.07 \pm 3.33$ | $\mathbf{8.72 \pm 7.44}$ |

**Implementation** Training times are around $12 - 24$ hours for the brain experiments and $2 - 4$ days for the cardiac experiments using an NVIDIA GTX 1080Ti GPU. During inference, the agent finds the target plane using iterative steps, where each step takes ~0.02s. The details of the our proposed network for DQN are in Figure 1. The source code of our implementation is publicly available on github https://git.io/vhuMZ.

## 5 Discussion and Conclusion

We proposed a novel approach based on multi-scale reinforcement learning agents for automatic standard view extraction. Our approach is capable of finding standardized planes in real time, which in turn enables accelerated image acquisition. Consequently, it can alleviate the comparison between different imaging examinations using anatomically standardized biometric measurements. We extensively evaluated several DQN based strategies for the detection of three different planes. Our approach achieved good results for the automatic detection of the ACPC and mid-sagittal planes from brain MRI with distance error less than 2 mm, and for the detection of the 4-chamber plane from cardiac MRI with distance error around 5 mm.

**Limitations:** Our results show that the optimal algorithm for achieving the best performance is environment-dependant. In general, reinforcement learning is a difficult problem that needs a careful formulation of its elements such as states, rewards and actions. For example, RL tends to overfit to the reward signals, which may cause unexpected behaviours. Therefore the design of the reward function has to capture exactly the desired task, and still be learnable.

**Future Work:** we will investigate using a continuous action space to improve the performance through reduction of quantization errors introduced by

fixed action steps. We will also explore the use of either competitive or collaborative multi-agents to detect the same or different anatomical planes. Another future direction is inspired by AlphaGo [11], where an RL agent could mimic the moves of a human expert and accumulate this experience, thus learning from experienced operators during real time observation.

## References

1. Ardekani, B.A., Kershaw, J., Braun, M., Kanuo, I.: Automatic detection of the mid-sagittal plane in 3-D brain images. TMI 16(6), 947–952 (1997)
2. Bellman, R.: Dynamic programming. Courier Corporation (2013)
3. Ghesu, F.C., Georgescu, B., Zheng, Y., Grbic, S., Maier, A., Hornegger, J., Comaniciu, D.: Multi-Scale Deep Reinforcement Learning for Real-Time 3D-Landmark Detection in CT Scans. PAMI (2017)
4. Le, M., Lieman-Sifry, J., Lau, F., Sall, S., Hsiao, A., Golden, D.: Computationally efficient cardiac views projection using 3D Convolutional Neural Networks. In: Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support, pp. 109–116. Springer (2017)
5. Liao, R., Miao, S., de Tournemire, P., Grbic, S., Kamen, A., Mansi, T., Comaniciu, D.: An Artificial Agent for Robust Image Registration. In: AAAI. pp. 4168–4175 (2017)
6. Lu, X., Jolly, M.P., Georgescu, B., Hayes, C., Speier, P., Schmidt, M., Bi, X., Kroeker, R., Comaniciu, D., Kellman, P., et al.: Automatic view planning for cardiac MRI acquisition. In: MICCAI. pp. 479–486. Springer (2011)
7. Maicas, G., Carneiro, G., Bradley, A.P., Nascimento, J.C., Reid, I.: Deep Reinforcement Learning for Active Breast Lesion Detection from DCE-MRI. In: MICCAI. pp. 665–673. Springer (2017)
8. de Marvao, A., Dawes, T.J., Shi, W., Minas, C., Keenan, N.G., Diamond, T., Durighel, G., et al.: Population-based studies of myocardial hypertrophy: high resolution cardiovascular magnetic resonance atlases improve statistical power. Journal of Cardiovascular Magnetic Resonance 16(1), 16 (2014)
9. Mnih, V., Kavukcuoglu, K., Silver, D., et al.: Human-level control through deep reinforcement learning. Nature 518(7540), 529 (2015)
10. Mueller, S.G., Weiner, M.W., Thal, L.J., Petersen, R.C., Jack, C., Jagust, W., Trojanowski, J.Q., Toga, A.W., Beckett, L.: The Alzheimer's disease neuroimaging initiative. Neuroimaging Clinics 15(4), 869–877 (2005)
11. Silver, D., Huang, A., Maddison, C.J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., et al.: Mastering the game of go with deep neural networks and tree search. nature 529(7587), 484–489 (2016)
12. Stegmann, M.B., Skoglund, K., Ryberg, C.: Mid-sagittal plane and mid-sagittal surface optimization in brain MRI using a local symmetry measure. In: Medical Imaging: Image Processing. vol. 5747, pp. 568–580 (2005)
13. Sutton, R.S., Barto, A.G.: Reinforcement learning: An introduction, vol. 1. MIT press Cambridge (1998)
14. Van Hasselt, H., Guez, A., Silver, D.: Deep Reinforcement Learning with Double Q-Learning. In: AAAI. vol. 16, pp. 2094–2100 (2016)
15. Wang, Z., Schaul, T., Hessel, M., Van Hasselt, H., Lanctot, M., De Freitas, N.: Dueling network architectures for deep reinforcement learning. arXiv preprint arXiv:1511.06581 (2015)
16. Watkins, C.J., Dayan, P.: Q-learning. Machine learning 8(3-4), 279–292 (1992)