

Learning from Demonstrations for Autonomous Soft-tissue Retraction*

Ameya Pore^{1,2}, Eleonora Tagliabue¹, Marco Piccinelli¹, Diego Dall’Alba¹, Alicia Casals², Paolo Fiorini¹

Abstract— The current research focus in Robot-Assisted Minimally Invasive Surgery (RAMIS) is directed towards increasing the level of robot autonomy, to place surgeons in a supervisory position. Although Learning from Demonstrations (LfD) approaches are among the preferred ways for an autonomous surgical system to learn expert gestures, they require a high number of demonstrations and show poor generalization to the variable conditions of the surgical environment. In this work, we propose an LfD methodology based on Generative Adversarial Imitation Learning (GAIL) that is built on a Deep Reinforcement Learning (DRL) setting. GAIL combines generative adversarial networks to learn the distribution of expert trajectories with a DRL setting to ensure generalisation of trajectories providing human-like behaviour. We consider automation of tissue retraction, a common RAMIS task that involves soft tissues manipulation to expose a region of interest. In our proposed methodology, a small set of expert trajectories can be acquired through the da Vinci Research Kit (dVRK) and used to train the proposed LfD method inside a simulated environment. Results indicate that our methodology can accomplish the tissue retraction task with human-like behaviour while being more sample-efficient than the baseline DRL method. Towards the end, we show that the learnt policies can be successfully transferred to the real robotic platform and deployed for soft tissue retraction on a synthetic phantom.

I. INTRODUCTION

Robot-Assisted Minimally Invasive Surgery (RAMIS) is a consolidated paradigm in the field of medical robotics. RAMIS is a viable alternative to other surgical approaches that reduce several patient complications such as excessive intraoperative blood loss, post-operative trauma, and a high amount of mortality associated with more invasive surgeries [1]. One of the widely adopted platforms for RAMIS is the da Vinci Surgical System (Intuitive Surgical Inc.), a teleoperated system that enables dexterous control with enhanced precision, stability and accuracy.

A significant portion of RAMIS procedures is spent in mobilizing and manipulating tissues to reach the region of interest. The tissue is repetitively grasped and retracted to expose the underlying anatomical area [2]. This gesture of tissue retraction occurs in multiple phases of surgery and requires extensive interaction with soft tissues having heterogeneous physical and geometric properties, such as stiffness and viscoelasticity, with high inter and intra-subject variability. During RAMIS, tissue retraction task may require

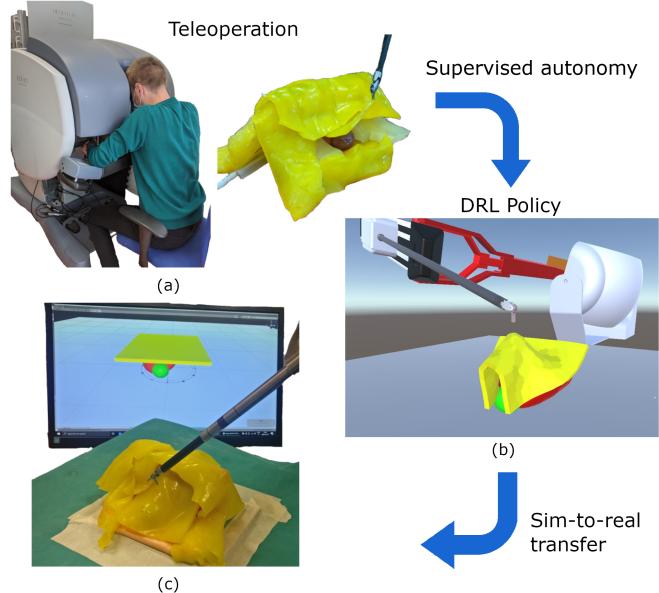


Fig. 1. The proposed methodology of LfD for the tissue retraction surgical gesture. (a) Expert demonstrations are performed and recorded using the dVRK console (b) Robotic agent is trained within a simulated environment. (c) The learnt policy is translated to the real robotic system.

the surgeon either to switch robotic arms during surgery with a different set of visuomotor feedback and limited perception or to instruct an assistant with the desired motion [3]. This involves risks such as instrument collision or tissues damage that can negatively affect the procedure. Therefore, soft tissue retraction is an ideal candidate for the automation of surgical gestures since it would help to reduce the cognitive and physical workload of the surgeons and place them in a supervisory position.

Standard motion planning techniques such as potential fields have been proposed for the automation of surgical gestures [4]. Such methods have shown to be successful in static conditions, but fail to generalize to dynamic environments like the surgical one, where human-level dexterity and adaptability are required. Hence, Learning from Demonstrations (LfD) is a preferred way to learn human gestures. However, a significant drawback of LfD methods is that they require a huge number of demonstrations to be trained properly, which is unfeasible in clinical settings considering the time, resources and ethical constraints. Moreover, LfD methods are affected by several limitations, such as errors in the acquisition and poor generalisation performance [5]. Using LfD, the robot can only become as good as the human’s demonstrations. There is no additional information

¹ Department of Computer Science, University of Verona, Italy. Email: ameya.pore@univr.it

² Automatic Control and Computer Engineering Department, Universitat Politècnica de Catalunya, Barcelona, Spain

*This project has received funding from the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie (grant agreement No. 813782 “ATLAS”) and (grant agreement No. 742671 “ARS”)

for improving the learnt behaviour.

By contrast, Deep Reinforcement Learning (DRL) allows the robot to discover new control policies through free exploration of the state-action space. DRL has shown generalisation capabilities in learning adaptable behaviours for complex and diverse scenarios such as dexterous manipulation and grasping [6]. However, DRL methods often take a long time to converge and require a well-shaped and carefully designed reward function to learn simple goal-directed behaviours. Approaches that combine LfD and DRL aim at exploiting the strength of both to overcome their respective drawbacks. Particularly, demonstrations can be used to guide the exploration done during learning, reducing the time required to find an improved control policy, which may depart from the demonstrated behaviour.

In this work, we address the limitations of LfD approaches by proposing a training methodology for surgical robots based on Generative Adversarial Imitation Learning (GAIL). GAIL relies on generative adversarial networks to generate trajectories that are similar to the acquired expert demonstrations [7]. To automate soft tissue retraction, we train GAIL in a simulated replica of the real environment using a small set of demonstrations acquired through the real robotic platform via teleoperation. This work represents an initial step for utilizing expert demonstrations to learn a control policy and to transfer the motion skills to a real robotic manipulator. Hence, our contribution is the introduction of an LfD training methodology to learn human-level surgical gestures using a small set of task demonstrations.

In Sec. II, we provide an overview of the related works. We summarise the mathematical formulation of the DRL methods used in Sec. III. In Sec. IV, we detail the training methodology and the learning setup. In Sec. V, we describe experiments conducted to validate our methodology. Further, we present the results in Sec. VI along with discussion. In the final Sec. VII, we elaborate the conclusions and highlight our future works.

II. RELATED WORKS

Previous works to automate soft-tissue retraction use motion planning algorithms such as probabilistic roadmaps for optimization based objectives [2]. This method works well for pre-operative planning in a static environment. Nagy et al. developed an approach for tissue retraction based on images, where three methods based on proportional control, Hidden Markov Models and fuzzy logic are validated [8]. Further, Attanasio et al. developed a trajectory planner based on coordinates extracted from images [3]. Both of these studies require hand-crafting control strategies and movement sequences. The execution of complex non-linear trajectories and behaviours may be challenging using these methods.

DRL for surgical tasks: Prior works have used DRL to learn the tensioning policy for tissue manipulation [9]. Pedram et al. developed a Q-learning RL algorithm based on visual features to learn a control policy for soft-tissue manipulation. Shin et al. developed an adaptive Model Predictive Controller (MPC) to learn the tissue dynamics [10].

Human demonstrations are used to pretrain and initialize the tissue dynamics model. The learnt tensioning policy is used to develop a visual model-based RL algorithm that outputs actions to manipulate the tissue points to pre-specified locations. In the former work, the tissue dynamics is learnt implicitly in the DRL setting, whereas in the latter, the dynamics is learnt via an MPC. These methods would be well suited for learning dynamics of deformable objects with similar physical properties. However, variability in soft tissues properties would require complete re-training of these methods for each subject. Hence, we base our hypothesis on learning task features of tissue manipulation from expert demonstrations in contrast to learning tissue dynamics. Task features would inherently include the surgeon's knowledge in manipulating a variety of tissues and would demonstrate better generalisation.

In the context of evaluation platform, Richter et al. proposed a framework for training DRL policies in simulation for the dVRK system and showed that the learnt policy could be transferred to a real robotic system [11]. The proposed environment does not support deformable objects, which is a major limitation when simulating the surgical scenario. Recently, Tagliabue et al. proposed a virtual framework called *UnityFlexML* for simulating deformable tissues that is well suited to train DRL methods [12]. We use this simulation framework to develop our LfD approach.

LfD for surgical tasks: Reiley et al. proposed a demonstration based framework using Gaussian Mixture Models (GMM) for motion generation [13]. Recently, a similar approach of GMM has been used to learn dynamic motion primitives from the demonstrations obtained from expert surgeons [14]. Osa et al. introduced an iterative technique to learn a single reference trajectory for knot tying [15]. However, a single demonstration does provide enough consistency to model a manipulation skill. Schulman et al. used a trajectory transfer algorithm to learn from demonstrations for the task of suturing [16]. Murali et al. devised a method to segment demonstrations into motion sequences [17]. Standard approaches for imitation learning such as Behaviour Cloning (BC) and Inverse Reinforcement Learning (IRL), tend to be successful with large amounts of data. However, they suffer from compounding error caused by covariant shift and are extremely expensive to train [6], hence not suited for the surgical paradigm. An emerging derivative of IRL is represented by GAIL [7]. Contrary to IRL, which learns a cost function, GAIL methods learn how to act by directly learning the policy. GAIL has shown an empirical improvement in reducing the number of demonstrations required to successfully learn the task, compared to other imitation learning methods. Therefore we selected GAIL as the LfD method, since it has been successfully applied in other surgical domains (e.g., endovascular manipulators [18]) but never used in the RAMIS context.

III. BACKGROUND

A general Reinforcement Learning (RL) problem can be formulated in terms of a Markov Decision Process (MDP),

where an agent learns by interacting with the environment. An MDP \mathcal{M} is defined as a tuple $(\mathcal{S}, \mathcal{A}, r, \mathcal{P}, \gamma, T)$, where \mathcal{S} is a set of possible states, \mathcal{A} is the set of actions, \mathcal{P} is the transition probability distribution, r is the reward function, $\gamma \in [0, 1]$ is the discount factor and T is the time horizon per episode. In RL, at each timestep t , the environment produces a state observation $s_t \in \mathcal{S}$. The agent then samples an action $a_t \sim \pi(s_t)$, $a_t \in \mathcal{A}$ and applies the action to the environment. As a consequence, the agent transitions to a new state s_{t+1} sampled from the transition function $p(s_{t+1}|s_t, a_t)$, $p \in \mathcal{P}$ or terminates the episode at state s_T . The agent's goal is to learn a stochastic behaviour policy π parametrised by ϕ , $\pi_\phi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ to maximise the expected future discounted reward $E[\sum_{i=0}^{T-1} \gamma^i r_i]$.

A. Proximal Policy Optimisation (PPO)

PPO is an on-policy RL algorithm capable of dealing with both continuous and discrete action spaces. PPO alternates between collecting new observations and improving the policy, while approximating the value function as well [19]. The update function for the PPO policy is the following

$$L(s_t, a_t, \theta_k, \theta) = \min \left(\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_k}(a_t|s_t)} \hat{A}^{\pi_{\theta_k}}(s_t, a_t), g(\epsilon, \hat{A}^{\pi_{\theta_k}}(s_t, a_t)) \right) \quad (1)$$

where θ_k are the parameters of the old policy, and g is defined as:

$$g(\epsilon, \hat{A}) = \begin{cases} (1 + \epsilon)\hat{A}_t, & \hat{A}_t \geq 0 \\ (1 - \epsilon)\hat{A}_t, & \hat{A}_t < 0 \end{cases} \quad (2)$$

Where \hat{A}_t is the advantage estimator function at timestep t and ϵ is a hyperparameter. The idea behind PPO is to limit the impact of the policy update using the min operator, so that the improvement in the policy is stable.

B. Generative Adversarial Imitation Learning (GAIL)

GAIL is an imitation learning algorithm [7]. The principle of GAIL is based on the generative adversarial networks, which consists of a discriminator $D_{\phi'}$ and a policy generator G_ϕ , where ϕ denotes the parameters associated with each network. The policy network generates exploration trajectories which are used by the discriminator to compute a surrogate function measuring the similarity between the generated policy and the expert policy. This similarity metric acts as reward proxy for the RL step. Unlike IRL techniques, GAIL directly generates policies instead of the reward function. The discriminator is trained to minimize the loss function:

$$L_{GAIL} = E_{\tau_\phi}[\log(D_{\phi'}(s_t, a))] + E_{\tau_E}[\log(1 - D_{\phi'}(s_t, a))] \quad (3)$$

where τ_ϕ are the trajectories generated by G_ϕ and τ_E are the expert trajectories. The policy generator G_ϕ is often adopted from methods based on stochastic policy such as PPO. There are two reasons why PPO is used for GAIL: first, PPO uses smooth policy update for stable learning and second, PPO generates diversified trajectories that act as a wide sampling range for the discriminator in GAIL.

IV. METHODS

The objective of our task is to efficiently and optimally train the robotic agent to expose the tumour region by grasping and retracting the fat tissue surrounding it. In particular, our DRL agent is represented by the End-Effector (EE) of the da Vinci Patient Side Manipulator (PSM) arm, which learns to move from an initial position \mathbf{p}_0 to a position close to the tumour \mathbf{q} , grasp the tissue and retract it to reach the desired target position \mathbf{p}_T . For the sake of simplicity, the PSM orientation is kept constant. The initial state of the anatomical environment, as well as the tumour centroid position \mathbf{q} , are assumed to be known from pre-operative data. In order to make the learnt motion primitives robust to different initial configurations, the EE starts from a different position after each episode (i.e. 2500 timesteps in our case) at training time. The considered state-space leverages solely on kinematics information defining the current robot state and the environment:

$$\mathcal{S}_t = [\mathbf{p}_t, \mathbf{q}, \mathbf{p}_T, \|\mathbf{p}_t - \mathbf{q}\|, \|\mathbf{p}_t - \mathbf{p}_T\|, g_t] \quad (4)$$

where $\|\cdot\|$ represents the Euclidean distance, \mathbf{p}_t denotes the position of the EE at time t and $g_t \in \{0, 1\}$ represents the gripper state (open/closed). At each time, the action the agent can take is defined by an increment motion of $0.5\beta mm$ in each spatial dimension, where β varies in $\{0, -1, +1\}$, corresponding respectively to the conditions of no motion, backward motion and forward motion.

The presented algorithms are trained within *UnityFlexML*¹ framework, a platform to create simulation environments supporting deformable objects and dVRK kinematics.

A. Deep Reinforcement Learning (DRL) setup

Standard DRL algorithms are trained in a simulation environment that is an exact replica of the real one. We consider PPO as the standard baseline DRL algorithm [19]. For the training phase, we rely on a reward function which varies with the gripper state, to encourage EE motion towards the tumour if the tissue has not been grasped yet (i.e. gripper is still open), and tissue retraction if the tissue has been already grasped:

$$r(s_t) = \begin{cases} -\|\mathbf{p}_t - \mathbf{q}\| * k - 0.5, & \text{before grasp} \\ -\|\mathbf{p}_t - \mathbf{p}_T\| * k, & \text{after grasp} \end{cases} \quad (2)$$

The scalar quantity of -0.5 is added to restrict the reward in the range $(-1.0, -0.5)$ before grasping and $(-0.5, 0)$ after grasping. The normalization constant k is introduced to allow re-scaling of the trajectories to different working spaces and is inversely proportional to the maximum distance the PSM can move.

B. Learning from Demonstrations (LfD) setup

Since we aim at finding a better policy than the provided demonstrations, the training is based on the linear combination of the respective DRL and GAIL losses:

$$L_{Total} = \alpha L_{DRL} + \beta L_{GAIL} \quad (5)$$

¹Project- <https://gitlab.com/altairLab/unityflexml>

where α and β represent weighting factors for the two loss functions. For a PPO agent, $\alpha = 1$ and $\beta = 0$. For the GAIL agent, our initial investigation on hyper-parameters tuning yielded best performance for $\alpha = 0.2$ and $\beta = 0.8$. Other values of α and β yielded slower convergence.

Training of a GAIL agent requires the collection of trajectory demonstrations. In this work, task demonstrations are acquired on the real dVRK and transferred to the *UnityFlexML* framework. The acquired trajectories consist of repetitive fat lifting task, performed by an expert user. Since the expert user is well aware of the final objective of tumour exposure, the grasp position is near the tumour for all the demonstrations. Moreover, the expert user is instructed to diversify the trajectories by starting each demonstration from a different initial position above the fat surface.

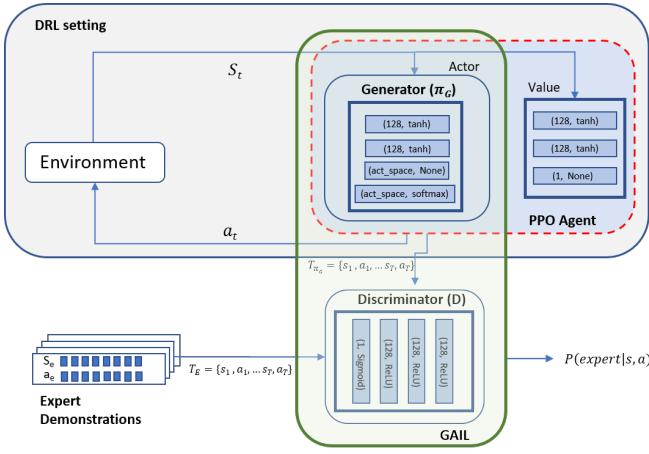


Fig. 2. Network architecture of GAIL and PPO. PPO consists of a policy (actor) and Value network. The policy network acts as Generator for GAIL. Generated trajectories and expert trajectories are passed to the Discriminator. Discriminator learns a probability function which classifies the generator trajectory as expert or non-expert. The network layer details are depicted inside each box in the format (hidden units, activation) respectively.

Acquisition of task demonstrations from the real environment leverages the communication pipeline provided by *UnityFlexML* (Fig. 1). Registration between the simulated and real environment is guaranteed following the same registration process described in [12]. The joint values are sent to *UnityFlexML* through UDP sockets and the desired configuration is reached with direct kinematics. Each recorded demonstration consists of the set of kinematic observations that define the state space (Sec. IV) and the corresponding actions at each timestep. An important aspect of this implementation is the challenge associated when we reset each episode. In the simulation, as soon as the target position is reached, the grasp is released and the episode resets. The position of the EE is then immediately teleported to the next initial point. This reset strategy has been adapted to cope with the real robotic system, during the recording of expert demonstrations. In particular, a delay of some timesteps has been added between the moment when the grasp is released and the beginning of the next episode, to allow repositioning. We make use of 35 continuous episodes recordings. Although

our simulation framework supports demonstration recordings using a keyboard or a joystick, the established communication pipeline between dVRK and *UnityFlexML* is crucial since it helps to acquire demonstrations directly with the real robotic system, thus without deviating from the surgical workflow. The network architecture used for our proposed GAIL and PPO training method is depicted in Fig. 2.

V. EXPERIMENTS

We consider a tissue retraction task in the context of a partial nephrectomy procedure. In particular, we aim at exposing a kidney tumour that is hidden by perinephric fat tissue. The real robotic setup consists of a synthetic kidney phantom covered with silicone fat tissue (Fig. 3). Methods presented in Sec. IV are trained in a simulation environment which is an exact replica of the phantom, where the deformation properties of the fat tissue have been optimized to mimic those of the real synthetic tissue used in the experiments [12]. All the simulation experiments, including DRL training and dVRK control, are executed on a workstation equipped with an AMD Ryzen 3700X processor and NVIDIA TitanX GPU.

The considered algorithms are tested both in simulation and reality based on two different criteria: sample efficiency and the optimality of the accomplished task, i.e. the ability of each method to expose the tumour. Sample efficiency is estimated by the number of time steps required by each algorithm to reach high reward values. Secondly, we estimate the optimality of the learnt behaviour by a Tumour Exposure (TE) metric. TE is assessed as the normalized percentage of tumour surface which can be seen from a camera positioned in front of the kidney, both for the simulated and the real setup. This evaluation allows us to understand if the quality of the exposure can be maximized regardless of the PSM starting position. The behaviour learnt in the simulated and real scenario is depicted in Fig. 3 top and Fig. 3 bottom, respectively.

A. Simulation experiments

The high level of realism of the simulated environment created within *UnityFlexML* does not only allow us to train the standard DRL methods with a sim-to-real approach but also provides a platform for testing the presented methods in realistic settings. In order to assess whether the behaviour learnt by the agent is robust to different starting EE positions, we perform an experiment where the trained agent has to perform the task starting from 49 different positions uniformly sampled on a 7×7 regular grid above the portion of the fat tissue. We evaluate TE each time the EE reaches p_T .

B. Real robotic experiments

Our real experimental setup (Fig. 3 bottom) is precisely aligned with the simulated scene (Fig. 3 top) by mapping the poses of the PSM in the two environments, as described in [20]. Accurate registration is essential to prevent inconsistencies between the two environments, especially considering

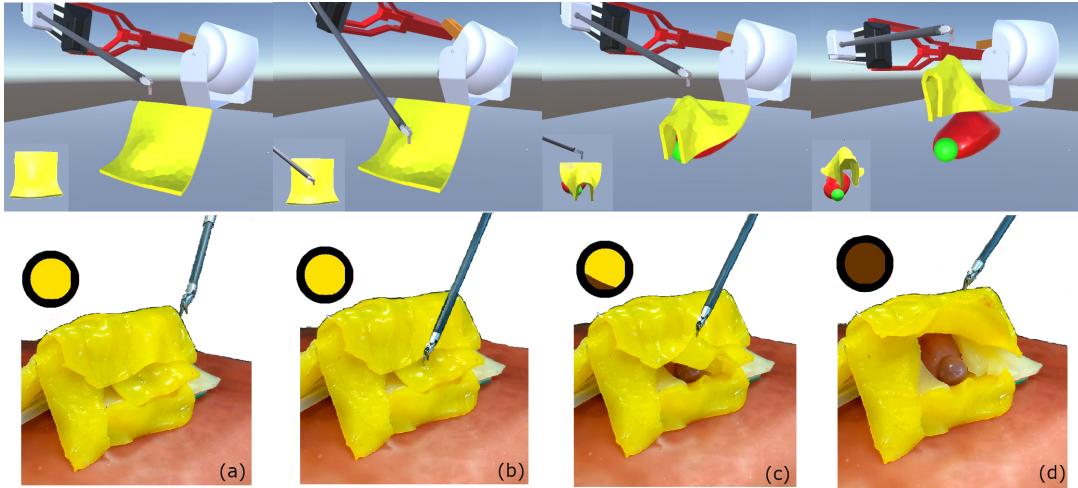


Fig. 3. Sequence of action frames for task completion in simulation (top) and reality (bottom). From left to right: approach, grasp, retract, tumour exposure. (top) Perspective of the simulated camera is overlaid on the bottom left of the simulator frames. (bottom) The real camera is placed in the same position as in the simulation (which do not correspond to the viewpoint used to take these pictures) The colour segmentation using a circular mask is illustrated for all the configurations (a) 0% exposure (b) 0% exposure (c)~15% exposure (d) 100% exposure.

that all the movements of the dVRK arm in the real system are controlled through the simulated robot, including the grasping action. To compute the TE metric, we select a circular region of interest around the tumour (exploiting the fact that its position is fixed) and we extract the visible portion by applying a mask with HSV bounds matching tumour colour (Fig. 3 bottom).

VI. RESULTS AND DISCUSSION

The results obtained will be presented in two separate sections, the first dedicated to simulation results while the second is dedicated to the results on the real setup.

A. Simulation

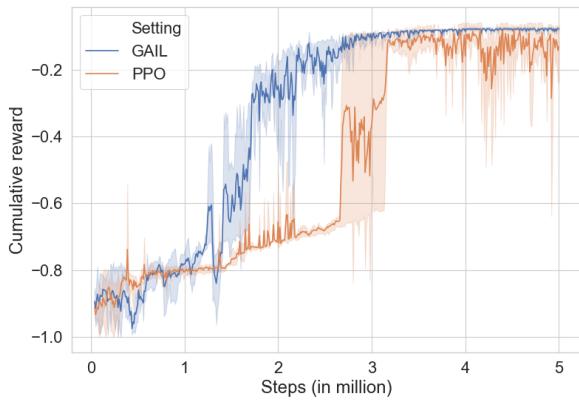


Fig. 4. The obtained learning curve for GAIL and PPO. Cumulative reward is normalised in the range $[-1, 0]$. The shaded area spans the range of values obtained when training the agent starting from three different initialization seeds.

Learning curves obtained with the considered learning configurations are showed in Fig. 4. GAIL is more efficient

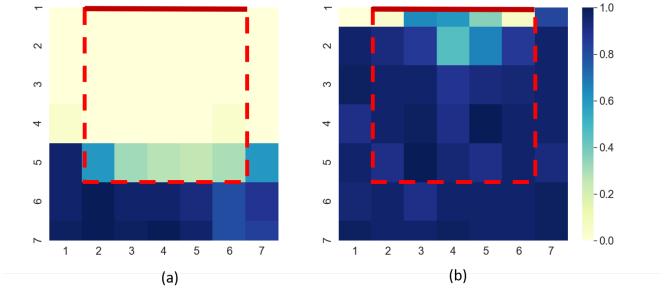


Fig. 5. Simulation experiments: TE from the camera at different initial positions, of the PSM for (a) PPO, (b) GAIL. The colour of each subregion is related to the percentage of visible tumour area when p_0 belongs to that subregion. The fat boundary from the top view is depicted in red dashed lines whereas the fat attachment is shown in the solid red line

than PPO and shows a monotonous and smooth learning pattern. GAIL learning curve begins to increase towards high-reward values and diverges from PPO around 1 million steps. PPO shows a modular reward trend: it requires 2.5 million steps to learn the approach behaviour and 1 million steps to learn the retract behaviour. This experiment shows that incorporating human demonstrations makes learning sample-efficient compared to baseline PPO. This result verifies our hypothesis that incorporating human knowledge can provide initial prior reference and takes fewer timesteps to learn the behaviour.

The plot in Fig. 5 shows the TE from the simulated camera depending on the starting position of the PSM arm above and outside the boundary of the fat tissue and points. In the case of PPO, when the starting EE position is close to the fat attachment, the agent tends to grasp near the proximity of the fat attachment, thus causing little or no tumour exposure (Fig. 5a). Note that, for PPO, the reward function is handcrafted such that before the grasp, the agent

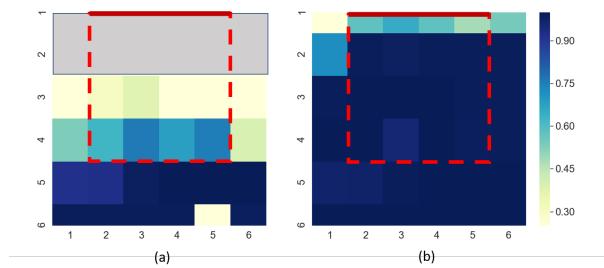


Fig. 6. Real grasp experiments: TE from the camera when starting from different initial positions of the EE, using (a) PPO (b) GAIL. The portion of fat tissue which is not considered for the experiments is coloured in grey.

is encouraged to approach the known position of the tumour, because grasping near the tumour is the best strategy to maximise tumour exposure. Low TE indicates that manually tuning the reward function to encode complex task objective such as tumour exposure can be challenging. We performed a preliminary evaluation of a scenario where we design a reward function that incorporates a TE-dependent term. However, it did not show significant improvements in the results. This might be due to the fact that before grasping the tumour exposure is always zero and represents a sparse reward scenario. We plan to further investigate this in future works.

On the contrary, when human demonstrations are incorporated, GAIL is able to grasp closer to the tumour and expose the tumour regardless of the starting position. Note that the strategy adopted by the user while acquiring demonstrations is to move and grasp towards points close to the tumour to maximize the exposure. The difference between the behaviour learnt by PPO and GAIL is in the robust selection of the grasping point when the starting position varies. When the starting position is above the attached area, PPO tends to grasp near the fat attachment, thus leading to low TE, whereas GAIL learns to grasp closer to the tumour, obtaining a higher TE.

B. Real robotic setup

We have been able to successfully replicate the learned behaviour from the simulated to the real environment without any appreciable inconsistency. The da Vinci EE successfully gets in contact with the fat tissue for all the different initial positions, and it is always able to reach the target point. The tumour exposure percentage starting from various points is illustrated in Fig. 6. For PPO, we did not initialize the EE positions near the attachment (represented as the unattempted grey region in Fig. 6a) because grasping near the attachment led to fat tissue tear. When comparing the results obtained for GAIL and PPO, it emerges that GAIL is not only able to reach higher overall exposure, but it is also more robust to changes in the initial PSM position. In particular, tumour exposure is achieved also when starting from points that were unattempted for PPO (Fig. 6b), thus suggesting an overall improvement in performances, due to grasping closer to the tumour position.

This observation indicates that the initial PSM position has a great impact on the performances in the case of PPO, whereas GAIL is able to reach optimal performance regardless of the starting position, confirming results obtained in the simulated experiments. In terms of the Average TE (ATE) computed considering all trials from different starting points, PPO obtains an ATE of 0.33 in a simulated environment and 0.38 in the real, while GAIL obtains an ATE of 0.84 in simulation and 0.90 in the real domain. We can infer that performance using demonstrations in both simulation and real-world is robust and outperforms PPO.

VII. CONCLUSIONS & FUTURE WORK

In this work, we present an LfD methodology for automation of tissue retraction based on GAIL. The proposed methodology can be trained in a simulated replica of the real scene using a small set of real robotic demonstrations and deployed in the real environment. The method builds on a consolidated DRL architecture and can learn generalised human-like trajectories in a sample-efficient way. Experiments in simulation and real environment show that, while both baseline DRL methods and GAIL can accomplish the task, the latter boosts the learning process by reducing the number of steps required and learning near-human trajectories. The learnt policy has been deployed on the dVRK and the tissue retraction task has been successfully completed.

This work has some limitations. The underlying hypothesis is based on the assumption of knowing the target positional coordinates (e.g., tumour position) pre-operatively. However, the tissue retraction surgical gesture can be carried out as an exploratory subtask without a known target. Hence, our future work will be directed in utilizing visual information to estimate the kinematic coordinates of the various image features, as described in [3], [21]. Further experimentation will be carried out to assess the impact of the quality and the number of demonstrations required to learn optimal behaviour with different surgical gestures, involving experts from various levels of expertise. Another limitation to consider is the safety issues due to free exploration of the state space that might lead to dangerous movements. Hence, in subsequent work, we will incorporate safety constraints through a Safe-Reinforcement learning technique [22].

In conclusion, surgical subtask automation is an anticipated research direction that will lead to a higher level of autonomy and put the surgeons in a supervisory role. This approach will enable a reduction in surgical workload and improved patient outcomes.

ACKNOWLEDGEMENT

The authors would like to thank Enrico Magnabosco, Sanat Ramesh for their support in the initial development of the simulator and discussion respectively. We also acknowledge the support of NVIDIA Corporation for the donation of the Titan Xp GPU used in this research.

REFERENCES

- [1] N. Simaan, R. M. Yasin, and L. Wang, "Medical technologies and challenges of robot-assisted minimally invasive intervention and diagnostics," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 1, pp. 465–490, 2018.
- [2] S. Patil and R. Alterovitz, "Toward automated tissue retraction in robot-assisted surgery," in *2010 IEEE International Conference on Robotics and Automation*. IEEE, 2010, pp. 2088–2094.
- [3] A. Attanasio, B. Scaglioni, M. Leonetti, A. F. Frangi, W. Cross, C. S. Biyani, and P. Valdastri, "Autonomous tissue retraction in robotic assisted minimally invasive surgery—a feasibility study," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 6528–6535, 2020.
- [4] M. Ginesi, D. Meli, A. Calanca, D. Dall'Alba, N. Sansonetto, and P. Fiorini, "Dynamic movement primitives: Volumetric obstacle avoidance," in *2019 19th International Conference on Advanced Robotics (ICAR)*. IEEE, 2019, pp. 234–239.
- [5] A. Hussein, M. M. Gaber, E. Elyan, and C. Jayne, "Imitation learning: A survey of learning methods," *ACM Computing Surveys (CSUR)*, vol. 50, no. 2, pp. 1–35, 2017.
- [6] J. Ibarz, J. Tan, C. Finn, M. Kalakrishnan, P. Pastor, and S. Levine, "How to train your robot with deep reinforcement learning: lessons we have learned," *The International Journal of Robotics Research*, p. 0278364920987859, 2021.
- [7] J. Ho and S. Ermon, "Generative adversarial imitation learning," in *Advances in neural information processing systems*, 2016, pp. 4565–4573.
- [8] T. D. Nagy, M. Takács, I. J. Rudas, and T. Haidegger, "Surgical subtask automation—soft tissue retraction," in *2018 IEEE 16th World Symposium on Applied Machine Intelligence and Informatics (SAMI)*. IEEE, 2018, pp. 000055–000060.
- [9] N. D. Nguyen, T. Nguyen, S. Nahavandi, A. Bhatti, and G. Guest, "Manipulating soft tissues by deep reinforcement learning for autonomous robotic surgery," in *2019 IEEE International Systems Conference (SysCon)*. IEEE, 2019, pp. 1–7.
- [10] C. Shin, P. W. Ferguson, S. A. Pedram, J. Ma, E. P. Dutson, and J. Rosen, "Autonomous tissue manipulation via surgical robot using learning based model predictive control," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 3875–3881.
- [11] F. Richter, R. K. Orosco, and M. C. Yip, "Open-sourced reinforcement learning environments for surgical robotics," *arXiv preprint arXiv:1903.02090*, 2019.
- [12] E. Tagliabue, A. Pore, D. Dall'Alba, E. Magnabosco, M. Piccinelli, and P. Fiorini, "Soft tissue simulation environment to learn manipulation tasks in autonomous robotic surgery," in *2020 IEEE International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020.
- [13] C. E. Reiley, E. Plaku, and G. D. Hager, "Motion generation of robotic surgical tasks: Learning from expert demonstrations," in *2010 Annual international conference of the IEEE engineering in medicine and biology*. IEEE, 2010, pp. 967–970.
- [14] H. Su, A. Mariani, S. E. Ovur, A. Menciassi, G. Ferrigno, and E. De Momi, "Toward teaching by demonstration for robot-assisted minimally invasive surgery," *IEEE Transactions on Automation Science and Engineering*, 2021.
- [15] T. Osa, N. Sugita, and M. Mitsuishi, "Online trajectory planning and force control for automation of surgical tasks," *IEEE Transactions on Automation Science and Engineering*, vol. 15, no. 2, pp. 675–691, 2017.
- [16] J. Schulman, A. Gupta, S. Venkatesan, M. Tayson-Frederick, and P. Abbeel, "A case study of trajectory transfer through non-rigid registration for a simplified suturing scenario," in *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2013, pp. 4111–4117.
- [17] A. Murali, S. Sen, B. Kehoe, A. Garg, S. McFarland, S. Patil, W. D. Boyd, S. Lim, P. Abbeel, and K. Goldberg, "Learning by observation for surgical subtasks: Multilateral cutting of 3d viscoelastic and 2d orthotropic tissue phantoms," in *2015 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2015, pp. 1202–1209.
- [18] W. Chi, G. Dagnino, T. M. Kwok, A. Nguyen, D. Kundrat, M. E. Abdelaziz, C. Riga, C. Bicknell, and G.-Z. Yang, "Collaborative robot-assisted endovascular catheterization with generative adversarial imitation learning," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 2414–2420.
- [19] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [20] N. Piccinelli, A. Roberti, E. Tagliabue, F. Setti, G. Kronreif, R. Muradore, and P. Fiorini, "Rigid 3d registration of pre-operative information for semi-autonomous surgery," in *2020 International Symposium on Medical Robotics (ISMR)*. IEEE, 2020.
- [21] Y. Li, F. Richter, J. Lu, E. K. Funk, R. K. Orosco, J. Zhu, and M. C. Yip, "Super: A surgical perception framework for endoscopic tissue manipulation with surgical robotics," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 2294–2301, 2020.
- [22] R. Cheng, G. Orosz, R. M. Murray, and J. W. Burdick, "End-to-end safe reinforcement learning through barrier functions for safety-critical continuous control tasks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 3387–3395.