# On Deep Recurrent Reinforcement Learning for Active Visual Tracking of Space Noncooperative Objects

Dong Zhou[1], Guanghui Sun[1,*], Zhao Zhang[1] and Ligang Wu[1]

*Abstract*—Active tracking of space noncooperative object that merely relies on vision camera is greatly significant for autonomous rendezvous and debris removal. Considering its Partial Observable Markov Decision Process (POMDP) property, this paper proposes a novel tracker based on deep recurrent reinforcement learning, named as RAMAVT which drives the chasing spacecraft to follow arbitrary space noncooperative object with high-frequency and near-optimal velocity control commands. To further improve the active tracking performance, we introduce Multi-Head Attention (MHA) module and Squeeze-and-Excitation (SE) layer into RAMAVT, which remarkably improve the representative ability of neural network with almost no extra computational cost. Extensive experiments and ablation study implemented on SNCOAT benchmark show the effectiveness and robustness of our method compared with other state-of-the-art algorithm. The source codes are available on **https://github.com/Dongzhou-1996/RAMAVT**.

*Index Terms*—Active visual tracking, Deep recurrent reinforcement learning, Space noncooperative object, Multi-head attention

## I. INTRODUCTION

With the rapid development of aerospace technology, space noncooperative object active visual tracking that drives the chasing spacecraft or space manipulator to pursue any specific noncooperative target by merely using vision camera has attracted extensive attentions. It is essential to intelligent on-orbit service such as autonomous rendezvous [1]–[3], space debris removal [4]–[6], and malfunctioning satellite maintenance [7]–[9].

Benefitting from the powerful deep reinforcement learning (DRL) that has achieved great successes in many fields like video game [10], Go [11], autonomous driving [12], and robotic manipulation [13], more and more active visual trackers [14]–[22] have been proposed in an end-to-end manner, which can learn global optimal policy after training with millions of trial-and-errors experiences.

In our preliminary work [3], we presented the SNCOAT benchmark [23] and the first active visual tracker, DRLAVT in aerospace domain. It achieves impressive tracking performance in velocity control mode by stacking multiple frames as an input. However, the stacking mechanism not only decreases the control bandwidth severely, but also makes active tracker more vulnerable to perturbations (e.g., image blur, actuator noise, computational delay).
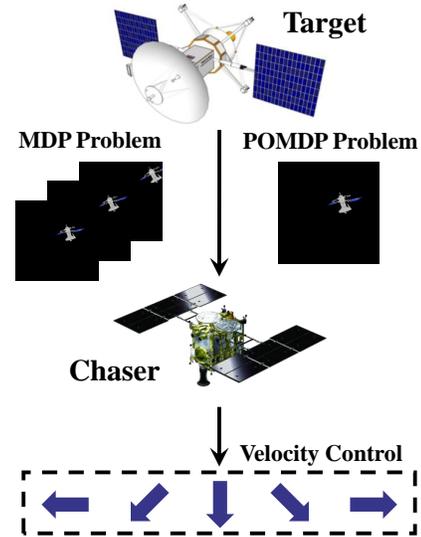
Fig. 1: Space noncooperative object active visual tracking

To this end, we take the POMDP property of space noncooperative object tracking into consideration. A novel active visual tracker based on deep recurrent reinforcement learning, RAMAVT is proposed in this paper, which can drive the chasing spacecraft to pursue arbitrary target with high-frequency and near-optimal velocity control commands. Our method features the accurate perception of target position and velocity, even though taken one image as input per time, which benefits from the recurrent neural network (RNN) in RAMAVT architecture that establishes the relationship between long-term temporal sequence.

Intuitively, it is totally enough for active visual tracker that focus attention on partial regions or partial channels of feature tensor to achieve the information of space noncooperative target. To this end, we adopt the multi-head attention module [24] and Squeeze-and-Excitation layer [25] in RAMAVT, which only increase small number of model parameters but significantly improve the representative ability of neural network. In addition, some data augmentation methods [26] are also involved to enhance the efficiency of learning process and generalization ability of active tracker.

The contributions of our work in this paper are summarized as following:

- We propose a novel active visual tracker based on deep recurrent reinforcement learning, RAMAVT which achieves excellent performance compared to other state-of-the-art methods.
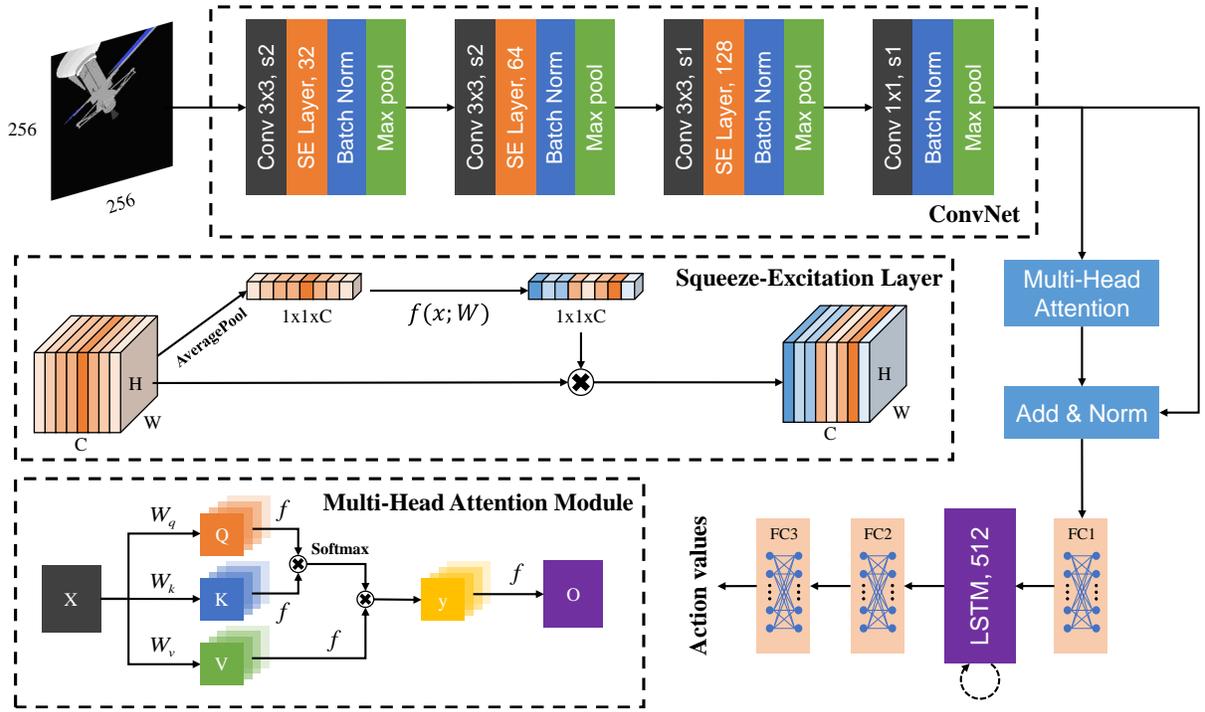
Fig. 2: The architecture of RAMAVT, which directly maps an image to optimal velocity control command attributed to the LSTM module that establish long-term relationship in temporal sequence. RAMAVT also adopts MHA module [24] and SE layer [25] to improve the representative ability of neural network.

- Multi-head attention module and SE layer are adopted into RAMAVT, combined with data augmentation, of which effectiveness has been proved by detailed ablation study.

This work proceeds as follows: Section II introduces some related works about space noncooperative object active visual tracking. Section III describes our RAMAVT method in detail. The experiments and analysis are given in Section IV. Finally, we make a conclusion in Section V.

## II. RELATED WORK

Visual object tracking is a hot topic in computer vision society, which has wide applications in civil, military, and aerospace fields. In recent years, many efforts have been devoted to studying passive methods [27]–[31], which assume that the target is always within the field of view (FOV) of vision camera. This severely limits the possibility to apply visual object tracking methods in many real-world scenarios, especially for aerospace applications where the target often maneuvers in 6-Degrees-of-Freedom (DoF) and the low-resolution camera mounted on spacecraft only has small FOV.

Therefore, active visual tracking [14]–[22], [32]–[34] has achieved more and more concerns, which not only identifies the target but also changes the pose of the chaser in real-time to keep view contact with the target. Traditional active visual trackers often adopt PBVS or IBVS framework of which modules (e.g. key-points detection, feature matching, pose estimation, and controller design) are optimized separately.

In the paper [32], a PBVS algorithm that guide space robotic manipulator to grasp noncooperative target was proposed, in which photogrammetry and adaptive extended Kalman filter are used to predict 6-DoF pose of target. However, this work is unadaptable to complex space environment as the same as other traditional active visual trackers. We also proposed a novel PBVS tracker in our preliminary work [3] that adapts state-of-the-art 2D monocular tracking method, SiamRPN [35], which achieved fairly good active tracking performance on SNCOAT benchmark but make concession in real-time capability.

Deep reinforcement learning that learns optimal action policy in an end-to-end manner with millions of trial-and-errors experiences has made great contributions in many fields, such as video game [10], Go [11], autonomous driving [12], and robotic manipulation [13], which provides a novel perspective for active visual tracking. In recent years, many DRL-based active visual trackers were proposed [14]–[22], most of which aim at terrestrial targets and only deploy on unmanned ground vehicle (UGV).

Luo et al. [14] proposed the first end-to-end active visual tracker based on A3C [36] that can only pursue two person models walking along fixed trajectory in two types of environments. In addition, The training progress takes up several days to achieve nice tracking performance, even with very low-resolution image (84×84×3). The paper [37] presented a temporal difference-based reward function adopted in PPO learning framework [38], which effectively decreases the distance error between the chasing and target

UGVs. However, this agent was only trained with single target in simulation environment. It was almost overfitted and impossible to track another target in real-world scenario. In contrast, our method is trained with 12 types of space noncooperative objects including space stations, satellites, asteroids, rockets, and return capsules, which successfully guarantees the generalization ability of RAMAVT.

Those algorithms mentioned above assume that the initial position of target is within the active tracker's FOV, which is a demanding condition for real application. To this end, Jeong et al. [18] extended the active visual tracking problem involved navigation, exploration and in-sight tracking and proposed the active tracking target network (ATTN) that learns a unified policy to track agile and anomalous object with partially known target model. This method features the incorporation that feeds egocentric maps and visit frequency to the convolutional neural network (CNN), which formulates the active visual tracking task as Markov Decision Process (MDP). In [21], Dionigi et al. also presented the DRL-based E-VAT model consisted of target-detection network and exploration-and-tracking network, which can explore the environment and track the target autonomously.

Compared with terrestrial targets, the active visual tracking tasks of space noncooperative objects are more challenging: 1) the tracked target often maneuvers agilely with complex 6-DoF trajectory; 2) less prior knowledge are available, such as geometry, texture, kinematic and dynamic parameters; 3) the images captured by vision camera on spacecraft are often low-quality, because of low resolution, small FOV, camera motion, and illumination variance.

In our preliminary work [3], we propose the first active tracker DRLAVT in aerospace domain, of which performance has a large room to be improved. The partially observable problem for active visual tracking was avoided by frame-stack mechanism, however, it severely decreases control bandwidth and tracking performance. To this end, we propose a novel active visual tracker based on deep recurrent reinforcement learning that directly maps one image to optimal velocity control command, benefitted from long-term temporal relationship established by RNN. In addition, the MHA module and SE layer are introduced to further improve network representative ability.

## III. PROPOSED METHOD

In this section, we formulate the active visual tracking problem of space noncooperative object and describe our RAMAVT algorithm thoroughly.

### A. Problem Formulation

The task of space noncooperative object active visual tracking involves the chaser mounted with vision camera and the moving target with no prior information, where the previous one should change its pose by using images to reduce the error $e_t$ which is formulated as follows:

$$e_t = \left\| r_T^B(t) - r^* \right\|_2 \tag{1}$$

in which $r_T^B(t)$ is the 3-D position of the target in the body-frame of the chaser at $t$th timestep, and $r^*$ denotes the expected distance between the chaser and target. In this work, $r^*$ is set to $\{0, 0, 5\}$.

To complete this task with DRL-based method, it can be further described as a POMDP problem. At $t$th timestep, the state of the target $s_t \in \mathcal{S}$ is observed as $o_t \in \mathcal{O}$ by agent with vision camera. Then, the agent takes action $a_t \in \mathcal{A}$ following a policy, such as the greedy policy $a_t = \max\limits_{a \in \mathcal{A}} Q(o_t, a)$ that is adopted in this article. After that, the agent would receive a reward $r_t$ from the environment which is generated by a reward function. The definition of our reward function is inherited from the paper [3], which includes a visible term $r_{vis}$ and a distance penalty term $r_{dist}$. The partially observable problem means $o_t \neq s_t$, that is, the agent can not accurately perceive actual state of the target, especially for the velocity.

### B. RAMAVT Algorithm

The POMDP problem mentioned above makes an end-to-end active visual tracker difficult to approximate optimal action value function $Q^*(o_t, a_t)$. To this end, we propose a new deep Q-network architecture based on DRQN [39], as shown in Fig. 2, which can establish the long-term relationship between temporal sequence and directly map one image to optimal velocity control command.

Meanwhile, some additional SE layers [25] and Multi-head attention module [24] are also introduced to improve the representative ability of deep Q-network and approximate better action value function $Q(o_t, a_t)$. The SE layer features the modelling of the channel-wise interdependencies of feature tensors with low computational cost, of which schematic is illustrated at the middle part of Fig. 2. In this work, We place SE layer behind every convolutional layer in ConvNet backbone.

In recent years, the self-attention mechanism derived from natural language processing filed has been widely applied to computer vision, which significantly increases the representation of neural network to images. The MHA module, shown in the bottom of Fig. 2, is one of the most famous self-attention method representing the interrelationship of different positions in one image. The basic of MHA is the scaled dot-product attention algorithm:

$$y_i = \mathbf{softmax}\left(\frac{Q_i K_i^T}{\sqrt{d_k}}\right) V_i \tag{2}$$

where $Q_i = W_q \cdot x_i$, $K_i = W_k \cdot x_i$, and $V_i = W_v \cdot x_i$ are the three feature vectors computed by different full-connected (FC) layers fed with the same input $x_i$, and $d_k$ denotes the dimension of feature vector $K_i$. Based on this, the MHA can be formulated as:

$$O = W_o \cdot \mathbf{Concat}\{y_1, y_2, \cdots, y_N\} \tag{3}$$

in which, $N$ is the number of heads that attend to information from different representation subspaces at different positions in parallel. We adopt $N = 8$ heads in this work.

TABLE I: Active Tracking Performance Comparison

| Name | Input Format | | | Episode Length | | | Episode Reward | | | Speed |
| | RGBD | Depth | Color | Avg | Min | Max | Avg | Min | Max | (Hz) |
|---|---|---|---|---|---|---|---|---|---|---|
| Random | - | - | - | 152.2 | 21 | 385 | -1545.2 | -3214.8 | -178.4 | 42703.5 |
| DRLAVT | √ | - | - | 857.9 | 6 | 1000 | -268.9 | -2341.4 | 386.4 | 63.1 |
| | - | √ | - | 901.1 | 6 | 1000 | <span style="color:green">430.3</span> | <span style="color:green">-55.69</span> | 382.2 | 66.6 |
| | - | - | √ | 841.6 | 15 | 1000 | -201.9 | -1798.7 | 320.4 | 68.9 |
| RAMAVT | √ | - | - | 952.4 | 41 | 1000 | -398.5 | -1523.2 | <span style="color:green">481.4</span> | 202.7 |
| | - | √ | - | <span style="color:green">959.1</span> | <span style="color:green">162</span> | 1000 | -59.4 | -800.5 | 445.4 | <span style="color:green">216.1</span> |
| | - | - | √ | 685.3 | 43 | 1000 | -1755.1 | -4097.5 | 434.7 | 210.9 |

(The best scores are highlighted in <span style="color:green">green</span>)

TABLE II: Training Configurations

| Params | Value | Note |
|---|---|---|
| replay buffer | 50000 | The size of replay pool |
| initial buffer | 10000 | The number of initial experiences |
| episode num | 300 | The number of episodes used to train Q-network |
| max episode len | 1000 | The max length of one episode, but if target is lost, episode will be over |
| update interval | 10 | The update interval of target network |
| gamma | 0.99 | Rewards discount factor |



Fig. 3: The training curves of different active visual trackers

Finally, we train the RAMAVT model with loss function $\mathcal{L}(\theta)$ defined as follows:

$$\mathcal{L}(\theta) = \mathbb{E}_{(o,a,r,o') \sim U(H)} \left[ (y - Q(o, a; \theta))^2 \right] \quad (4)$$

where the training data $(o, a, r, o')$ is uniformly sampled from the hierarchical memory pool $H$ proposed by us that is more suitable for deep recurrent reinforcement learning methods. $y = r + \gamma \max_{a' \in \mathcal{A}} Q(o', a'; \theta^-)$ is the Temporal-difference (TD) target estimated by the target network $\theta^-$.

## IV. EXPERIMENT

In this section, we first validate the active tracking performance of RAMAVT by using the evaluation toolkit provided in SNCOAT benchmark [23]. Then, sufficient ablation studies on RAMAVT are also implemented to show the effectiveness of our method. Finally, we further explore the working mechanism of active visual trackers. All the trackers follows the same training configurations listed in Table II. The experimental platform is HPC server equipped with Intel Xeon@E5 2650v4 CPU and Nvidia Tesla P100 GPU.

### A. RAMAVT Performance

We train the agent with 12 types of space noncooperative objects and evaluate it on other 6 different targets, including asteroids, satellite, rockets, space station, and return capsule.
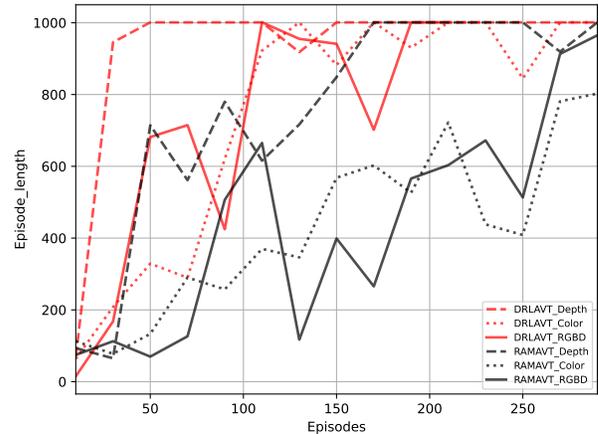
TABLE III: The Robustness Evaluation Under Different Perturbations

| Name | Perturbations | | | Metrics | |
| | Actuator Noise | Time Delay | Image Blur | AEL | AER |
|---|---|---|---|---|---|
| DRLAVT | √ | √ | √ | <span style="color:red">456.3</span> | -758.6 |
| RAMAVT | √ | - | - | 876.5 | -438.9 |
| | - | √ | - | 596.6 | -793.3 |
| | - | - | √ | 793.3 | -810.1 |
| | √ | √ | √ | 580.8 | <span style="color:red">-971.8</span> |
| | - | - | - | 952.4 | -398.5 |

(The worst scores are highlighted in <span style="color:red">red</span>)

Some data augmentations [26], such as crop, flip, cutout, and rotation are used to improve the generalization ability when trains the RAMAVT. Two metrics are adopted to measure active visual tracking performance, that is, episode length and episode reward. In this work, we utilize an agent that takes action in random as baseline and the DRLAVT algorithm as comparison.

All the training curves of two active visual trackers with

(a) active tracking trajectories



(b) active tracking errors in one trajectory



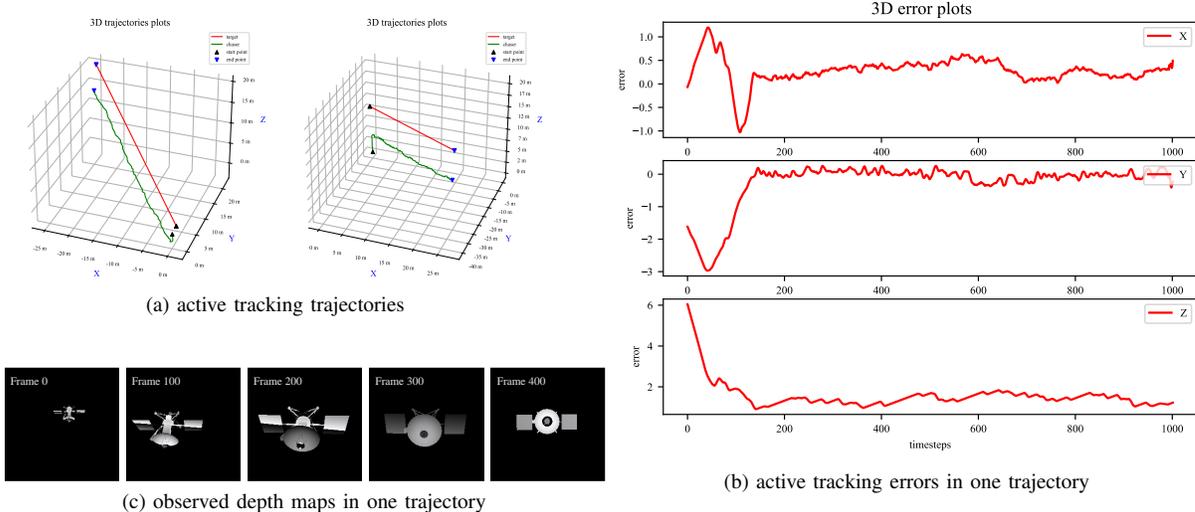(c) observed depth maps in one trajectory

Fig. 4: The results of RAMAVT with Depth image.

different inputs are depicted in Fig. 3. We find that the learning progresses of DRLAVTs are much faster than RAMAVTs whatever the input format is, because of its simple Q-network architecture and fully observable state. In addition, the depth information contained in inputs are significant for both active visual trackers to achieve higher episode length. It is worth noting that the depth map and color image adopted in this work are normalized to $[0, 1]$.

The whole evaluation results are summarized in Table I. It clearly shows that the RAMAVT taken depth map as input achieves the highest average episode length about 959.1 score with the best real-time performance, which means that our method can quickly track the target for a longer time. Meanwhile, the higher minimum episode length also proves the stability of RAMAVT. Although, its tracking accuracy (i.e. average episode reward) is slightly lower than DRLAVT. We think it results from the inaccurate target's states, such as target position and velocity, estimated by RAMAVT based on recurrent neural network. This problem gets worse when the agent is only allowed to use color images, as shown in the final row of Table I.

We visualize the tracking results of RAMAVT in Fig 4. It can be seen from Fig 4a that our method can precisely track the target in the whole episodes. In particular, the tracking errors in X, Y, and Z axes rapidly shrink to 0 and oscillate in a small range, as shown in Fig 4b. The noncooperative target is also steadily kept in the center of FOV after 200 frames (see in Fig 4c), even it moves fast with high-speed rotation.

Furthermore, we evaluate the RAMAVT tracker under three types of perturbations, involving actuator noise, time delay, and image blur, to show the robustness of our method. The experiment results are listed in Table III. It is obvious that all the three perturbations have influences on the active visual tracking performance in terms of tracking period and accuracy, especially for the time delay which decreases the
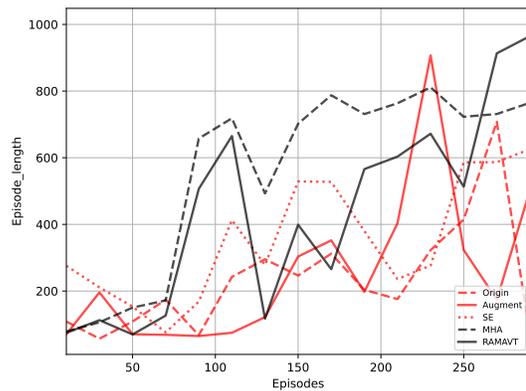


Fig. 5: The training curves of ablation models

TABLE IV: Ablation Study on RAMAVT with RGBD image

| Name | AEL | AER | Speed |
|---|---|---|---|
| Origin | 368.5 | -1876.3 | 214.3 |
| Augment | 419.3 | -1283.6 | 215.7 |
| SE | 625.8 | -956.0 | 211.3 |
| MHA | 731.1 | -845.7 | 206.6 |
| RAMAVT | 952.4 | -398.5 | 202.7 |

average episode length (AEL) about $37.4\%$. We think it is because of the inconsistency of target velocity between the training and evaluation, introduced by the random time delay. In the training stage, the target velocity is set as a random constant in one episode, which has naturally been learned by our RAMAVT model. When the 3 types of perturbations work simultaneously, our method is much robuster than the DRLAVT which has 124.5 scores less under AEL metric.

(a) Input image    (b) 1st Conv block    (c) 2nd Conv block    (d) 3rd Conv block    (e) 4th Conv block

(f) Input image    (g) 1st Conv block    (h) 2nd Conv block    (i) 4th Conv block    (j) MHA module
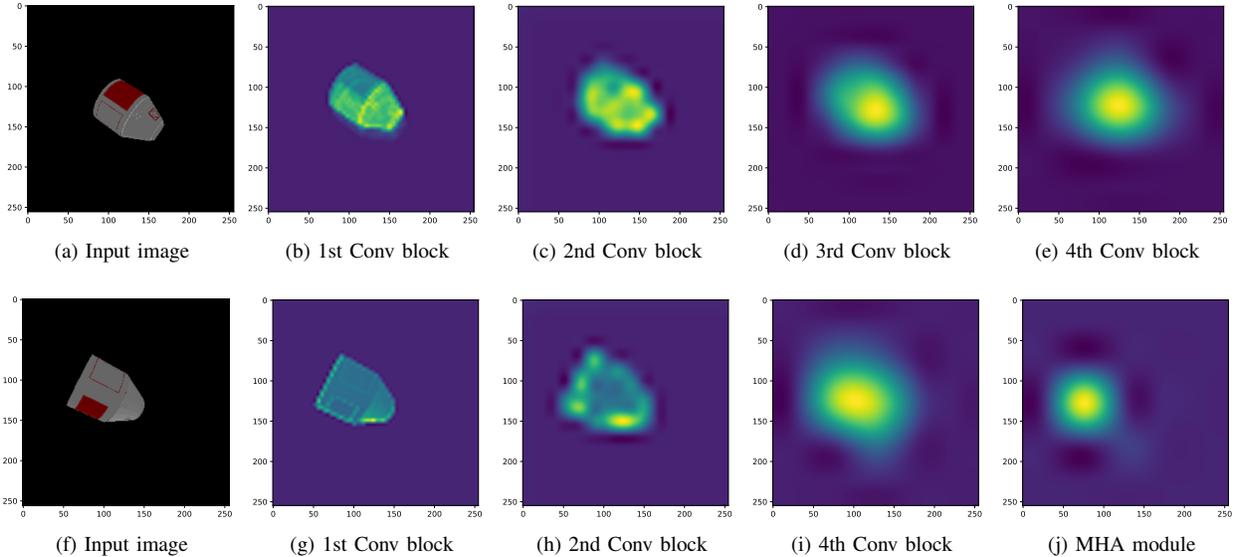
Fig. 6: The interpretability research of active visual trackers. The feature tensors of DRLAVT and RAMAVT are respectively visualized in the first and second rows. It is worthwhile noting that only the first frame of stacked input is depicted at 6a.

### B. Ablation Study

To show the effectiveness of RAMAVT model, we respectively add data augmentations, SE layer, and MHA module to the original DRQN architecture [39]. All the ablation models are trained from scratch with the same training configurations, of which training curves are illustrated in Fig 5. It can be clearly seen that the MHA module not only accelerates the learning progress of agent, but also significantly improves the episode length. This advancement attributes to the attention mainly focusing on the target, which makes the agent more sensitive and accurate to the movement of target.

The final evaluation results of ablation models are summarized in Table IV. It proves that the MHA module achieves the highest AEL and AER scores compared to the others, which only decreases $3.6\%$ running speed. The SE layer also increases the AEL measurement up to 1.7 times with almost no real-time performance loss. In addition, the data augmentation algorithms adopted in this paper including crop, cutout, flip, and rotation work unsatisfactorily, although it does not induce any computational burden during evaluation stage.

In a word, the RAMAVT model proposed in this work achieves excellent active visual tracking performance in less computational cost, mainly benefitted from spatial-wise and channel-wise attention mechanism induced by the MHA module and SE layers.

### C. Interpretability Research

We utilize the neural network interpretability method [40] that summarizes the squares of activation values along channel-wise axis and follow with 2-D Softmax operation to explore the inner working mechanism of active visual trackers. The features extracted by each layer of neural network are separately visualized in Fig 6.

The Fig. 6a-6e illustrates different levels of features extracted by DRLAVT of which architecture only contains 4 convolutional blocks. Each convolutional block involves a convolutional layer, a batch-normalization layer, and ReLU activation function. It is worthwhile noting that DRLAVT stacks 4 consecutive frames in channel-wise as one input, however, only the first frame is depicted at Fig. 6a. We clearly see that the first convolutional block extracts all the edges of target and generates higher feature value to the white body of noncooperative target. The second convolutional block further enhances the reactions to parts of edges. In the subsequent convolutional blocks, more high-level features without specific implications are extracted. The final output of ConvNet backbone looks like a point light source that follows normal Gaussian distribution.

In comparison, the visualization of RAMAVT is not the same, because of two significant differences between the backbone of RAMAVT and DRLAVT: (1) SE layer is added into the first three convolutional blocks, (2) the ConvNet backbone follows with a MHA module. Therefore, the first convolutional block no longer focus on object color (see Fig. 6b and 6g). The second convolutional block is more interested in the contour of noncooperative target. Furthermore, the distribution of final output turns more compact caused by the MHA module, which helps to estimate more accurate action value.

## V. CONCLUSION

In this paper, we formulate the active visual tracking task of space noncooperative object as POMDP problem and propose a novel active tracker based on deep recurrent reinforcement learning, RAMAVT of which architecture creatively adopts Squeeze-and-Excitation layer and Multi-Head Attention module. It can guide the chasing spacecraft to

approach arbitrary space noncooperative target with optimal and high-speed velocity control commands. The advancement of RAMAVT has been proved by sufficient experiments, compared to the state-of-the-art method DRLAVT. To show the effectiveness of our method, we implement convincing ablation study on RAMAVT architecture. In addition, we further take an interpretability research on two active visual trackers to explore their inner working mechanism.

## REFERENCES

[1] K. Hovell and S. Ulrich, "Deep Reinforcement Learning for Spacecraft Proximity Operations Guidance," *Journal of Spacecraft and Rockets*, Jan. 2021.

[2] X. Zhao, M. R. Emami, and S. Zhang, "Image-based control for rendezvous and synchronization with a tumbling space debris," *Acta Astronautica*, vol. 179, pp. 56–68, Feb. 2021.

[3] D. Zhou, G. Sun, W. Lei, and L. Wu, "Space non-cooperative object active tracking with deep reinforcement learning," *IEEE Transactions on Aerospace and Electronic Systems*, 2022.

[4] P. Huang, F. Zhang, J. Cai, D. Wang, Z. Meng, and J. Guo, "Dexterous Tethered Space Robot: Design, Measurement, Control, and Experiment," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 53, no. 3, pp. 1452–1468, June 2017.

[5] J. L. Forshaw, G. S. Aglietti, S. Fellowes, T. Salmon, I. Retat, A. Hall, T. Chabot, A. Pisseloup, D. Tye, C. Bernal, *et al.*, "The active space debris removal mission removedebris. part 1: From concept to launch," *Acta Astronautica*, vol. 168, pp. 293–309, 2020.

[6] G. S. Aglietti, B. Taylor, S. Fellowes, T. Salmon, I. Retat, A. Hall, T. Chabot, A. Pisseloup, C. Cox, A. Mafficini, *et al.*, "The active space debris removal mission removedebris. part 2: In orbit operations," *Acta Astronautica*, vol. 168, pp. 310–322, 2020.

[7] A. Flores-Abad, O. Ma, K. Pham, and S. Ulrich, "A review of space robotics technologies for on-orbit servicing," *Progress in Aerospace Sciences*, vol. 68, pp. 1–26, July 2014.

[8] D. Fourie, B. E. Tweddle, S. Ulrich, and A. Saenz-Otero, "Flight results of vision-based navigation for autonomous spacecraft inspection of unknown objects," *Journal of spacecraft and rockets*, vol. 51, no. 6, pp. 2016–2026, 2014.

[9] W.-J. Li, D.-Y. Cheng, X.-G. Liu, Y.-B. Wang, W.-H. Shi, Z.-X. Tang, F. Gao, F.-M. Zeng, H.-Y. Chai, W.-B. Luo, *et al.*, "On-orbit service (oos) of spacecraft: A review of engineering developments," *Progress in Aerospace Sciences*, vol. 108, pp. 32–120, 2019.

[10] V. Mnih, K. Kavukcuoglu, D. Silver, and et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, Feb. 2015.

[11] D. Silver, A. Huang, C. J. Maddison, A. Guez, and et al., "Mastering the game of Go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, pp. 484–489, Jan. 2016.

[12] B. R. Kiran, I. Sobh, V. Talpaert, P. Mannion, A. A. Al Sallab, S. Yogamani, and P. Pérez, "Deep reinforcement learning for autonomous driving: A survey," *IEEE Transactions on Intelligent Transportation Systems*, 2021.

[13] B. Singh, R. Kumar, and V. P. Singh, "Reinforcement learning in robotic applications: a comprehensive survey," *Artificial Intelligence Review*, pp. 1–46, 2021.

[14] W. Luo, P. Sun, F. Zhong, W. Liu, T. Zhang, and Y. Wang, "End-to-end Active Object Tracking via Reinforcement Learning," in *International Conference on Machine Learning*, July 2018, pp. 3286–3295.

[15] G. Cruciata, L. Lo Presti, and M. L. Cascia, "On the Use of Deep Reinforcement Learning for Visual Tracking: A Survey," *IEEE Access*, vol. 9, pp. 120 880–120 900, 2021.

[16] A. Devo, A. Dionigi, and G. Costante, "Enhancing continuous control of mobile robots for end-to-end visual active tracking," *Robotics and Autonomous Systems*, vol. 142, p. 103799, 2021.

[17] F. Zhong, P. Sun, W. Luo, T. Yan, and Y. Wang, "AD-VAT+: An Asymmetric Dueling Mechanism for Learning and Understanding Visual Active Tracking," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 5, pp. 1467–1482, May 2021.

[18] Heejin Jeong, H. Hassani, M. Morari, D. D. Lee, and G. J. Pappas, "Deep Reinforcement Learning for Active Target Tracking," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, 2021, pp. 1825–1831.

[19] F. Zhong, P. Sun, W. Luo, T. Yan, and Y. Wang, "Towards distraction-robust active visual tracking," in *International Conference on Machine Learning*, 2021, pp. 12 782–12 792.

[20] P. Tiritiris, N. Passalis, and A. Tefas, "Temporal Difference Rewards for End-to-end Vision-based Active Robot Tracking using Deep Reinforcement Learning," in *2021 International Conference on Emerging Techniques in Computational Intelligence (ICETCI)*, 2021, pp. 21–25.

[21] A. Dionigi, A. Devo, L. Guiducci, and G. Costante, "E-VAT: An Asymmetric End-to-End Approach to Visual Active Exploration and Tracking," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 4259–4266, 2022.

[22] M. Xi, Y. Zhou, Z. Chen, W. Zhou, and H. Li, "Anti-Distractor Active Object Tracking in 3D Environments," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 6, pp. 3697–3707.

[23] D. Zhou, "Space noncooperative object active tracking benchmark," https://github.com/Dongzhou-1996/SNCOAT.

[24] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[25] J. Hu, L. Shen, S. Albanie, G. Sun, and E. Wu, "Squeeze-and-excitation networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 8, pp. 2011–2023, 2020.

[26] M. Laskin, K. Lee, A. Stooke, L. Pinto, P. Abbeel, and A. Srinivas, "Reinforcement learning with augmented data," *Advances in neural information processing systems*, vol. 33, pp. 19 884–19 895, 2020.

[27] X. Chen, B. Yan, J. Zhu, D. Wang, X. Yang, and H. Lu, "Transformer tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 8126–8135.

[28] M. Ondrašovič and P. Tarábek, "Siamese Visual Object Tracking: A Survey," *IEEE Access*, pp. 1–1, 2021, conference Name: IEEE Access.

[29] D. Zhou, G. Sun, J. Song, and W. Yao, "2D vision-based tracking algorithm for general space non-cooperative objects," *Acta Astronautica*, vol. 188, pp. 193–202, Nov. 2021.

[30] S. Hu, X. Zhao, L. Huang, and K. Huang, "Global instance tracking: Locating target more like humans," *IEEE Transactions on Pattern Analysis & Machine Intelligence*, no. 01, pp. 1–1, 2022.

[31] M. Dunnhofer, A. Furnari, G. M. Farinella, and C. Micheloni, "Visual object tracking in first person vision," *International Journal of Computer Vision*, pp. 1–25, 2022.

[32] G. Dong and Z. H. Zhu, "Autonomous robotic capture of non-cooperative target by adaptive extended kalman filter based visual servo," *Acta Astronautica*, vol. 122, pp. 209–218, 2016.

[33] L. Sun and Z. Zheng, "Adaptive relative pose control of spacecraft with model couplings and uncertainties," *Acta Astronautica*, vol. 143, pp. 29–36, 2018.

[34] J. Liu, H. Li, Y. Luo, and J. Zhang, "Robust adaptive relative position and attitude integrated control for approaching uncontrolled tumbling spacecraft," *Proceedings of the Institution of Mechanical Engineers, Part G: Journal of Aerospace Engineering*, vol. 234, no. 2, pp. 361–374, 2020.

[35] B. Li, J. Yan, W. Wu, Z. Zhu, and X. Hu, "High performance visual tracking with siamese region proposal network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8971–8980.

[36] V. Mnih, A. P. Badia, M. Mirza, A. Graves, and et al., "Asynchronous Methods for Deep Reinforcement Learning," in *International Conference on Machine Learning*. PMLR, June 2016, pp. 1928–1937.

[37] P. Tiritiris, N. Passalis, and A. Tefas, "Temporal Difference Rewards for End-to-end Vision-based Active Robot Tracking using Deep Reinforcement Learning," in *International Conference on Emerging Techniques in Computational Intelligence, ICETCI 2021*, Virtual, India, 2021.

[38] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms," *arXiv:1707.06347 [cs]*, Aug. 2017.

[39] M. Hausknecht and P. Stone, "Deep recurrent q-learning for partially observable mdps," in *2015 aaai fall symposium series*, 2015.

[40] S. Zagoruyko and N. Komodakis, "Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer," in *International Conference on Learning Representations*, 2017.