

Siamese Object Tracking for Vision-Based UAM Approaching with Pairwise Scale-Channel Attention

Guangze Zheng¹, Changhong Fu^{1,*}, Junjie Ye¹, Bowen Li¹, Geng Lu², and Jia Pan³

Abstract—Although the manipulating of the unmanned aerial manipulator (UAM) has been widely studied, vision-based UAM approaching, which is crucial to the subsequent manipulating, generally lacks effective design. The key to the visual UAM approaching lies in object tracking, while current UAM tracking typically relies on costly model-based methods. Besides, UAM approaching often confronts more severe object scale variation issues, which makes it inappropriate to directly employ state-of-the-art model-free Siamese-based methods from the object tracking field. To address the above problems, this work proposes a novel Siamese network with pairwise scale-channel attention (SiamSA) for vision-based UAM approaching. Specifically, SiamSA consists of a pairwise scale-channel attention network (PSAN) and a scale-aware anchor proposal network (SA-APN). PSAN acquires valuable scale information for feature processing, while SA-APN mainly attaches scale awareness to anchor proposing. Moreover, a new tracking benchmark for UAM approaching, namely UAMT100, is recorded with 35K frames on a flying UAM platform for evaluation. Exhaustive experiments on the benchmarks and real-world tests validate the efficiency and practicality of SiamSA with a promising speed. Both the code and UAMT100 benchmark are now available at <https://github.com/vision4robotics/SiamSA>.

I. INTRODUCTION

For multiple practical scenarios [1]–[4], *e.g.*, autonomous grasping [1], aerial pick-and-place [4], and wristband placement [2], the applications of unmanned aerial manipulator (UAM) mainly consist of two stages, *i.e.*, visual approaching and precisely manipulating the object. The key to visual UAM approaching lies in the object tracking method, which provides the continuous visual perception of the object. However, current tracking methods for vision-based UAM approaching are generally model-based [2] [4] [5]. These methods commonly rely on predefined class labels and require specific large-scale training datasets, which are inadequate to track arbitrary objects and increase heavy workload for source-limited UAM platforms. Besides, collecting specific datasets for training model-based methods also causes huge development costs. An intuitive solution is to introduce the model-free object tracking methods for UAM approaching. In the visual object tracking field, Siamese network-based methods [6]–[10] have surpassed correlation filter-based methods [11]–[15] and reached state-of-the-art (SOTA)

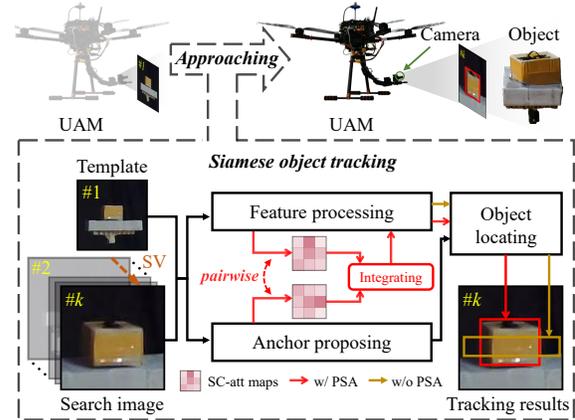


Fig. 1. A UAM approaching process and qualitative comparison. The results demonstrate the effects of the proposed pairwise scale-channel attention (PSA). The result in red box shows that tracking with PSA is more robust to deal with severe scale variation (SV) in UAM approaching than without PSA (brown box).

performance. The model-free Siamese trackers can perform on any online-assigned object to meet the requirements of practical application with an affordable camera. Therefore, this work introduces Siamese network as the object tracking method for vision-based UAM approaching.

Another issue of vision-based UAM approaching is the extreme object scale variation (SV), which is rarely taken into account. From the perspective of the UAM onboard camera, the scale of objects will get larger as the UAM approaches it. Because the distance between the UAM and the object is usually limited in practical applications, even modest changes in relative distance can result in significant scale variation, posing a major barrier for general object tracking methods. To address the SV problem, this work uses scale-equivariant convolution as the first step. In addition, attention-based strategies are elaborately designed to further solve SV issues. Because of the success in feature refinement [16] [17], these strategies have sparked considerable concern. In general, attention-based methods infer attention maps among spatial and channel dimensions, while this work attempts to uncover the scale attention for UAM approaching tasks. Specifically, a novel pairwise scale-channel attention (PSA) is proposed to extract important scale information across channels. As shown in Fig. 1, in both feature processing and anchor proposing, scale-channel attention (SC-att) maps are inferred in a pairwise structure and integrated to locate the object with stronger scale awareness.

Two branches are built based on the proposed PSA: a pairwise scale-channel attention network (PSAN) and a scale-aware anchor proposal network (SA-APN). PSAN seeks

*Corresponding Author

¹G. Zheng, C. Fu, J. Ye, and B. Li are with the School of Mechanical Engineering, Tongji University, Shanghai 201804, China. changhongfu@tongji.edu.cn

²G. Lu is with the Department of Automation, Tsinghua University, Beijing 100084, China.

³J. Pan is with the Department of Computer Science, the University of Hong Kong, Hong Kong, China

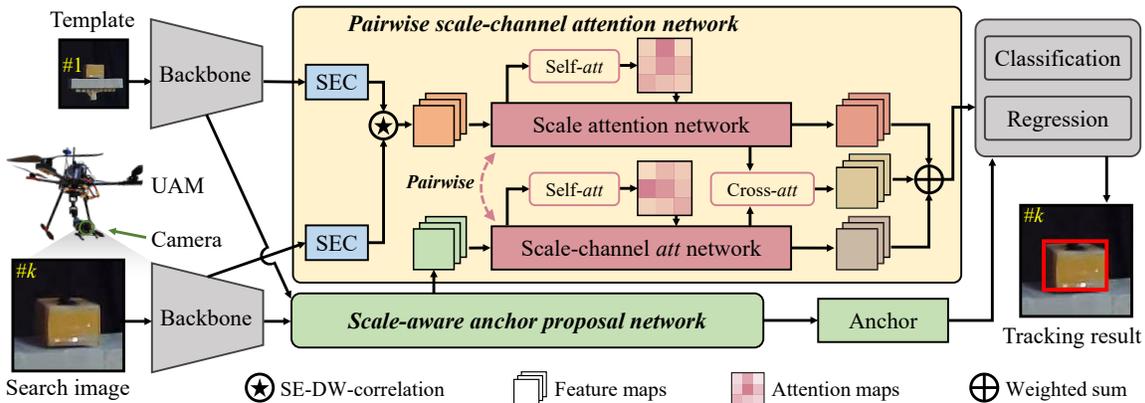


Fig. 2. An overview of the proposed Siamese tracking with pairwise scale-channel attention (SiamSA) for UAM approaching. SEC refers to scale-equivariant convolution, and *att* denotes attention.

to express more resilient features against extreme object scale variation, whereas SA-APN links scale information to anchors in the anchor proposal network. Self-attention and cross-attention tactics are used.

A fair evaluation of object tracking methods is required for vision-based UAM approaching. Despite vision-based UAM approaching being critical for practical applications, there are no publicly available benchmarks to assess UAM tracking methods for approaching yet. Therefore, this work records UAMT100 on a flying UAM platform with 100 image sequences. The UAMT100 benchmark covers common object tracking challenges and also introduces the unique challenge of UAM approaching, *i.e.*, UAM attack. The videos are taken in an indoor environment with a motion-capture system. SiamSA is also evaluated on the challenging authoritative UAV tracking benchmark to verify the generality of the proposed pairwise scale awareness. Finally, the proposed SiamSA has been deployed on a UAM platform and has proven to be highly efficient. The main contributions of this work can be summarized as follows:

- A Siamese object tracking method with novel pairwise scale-channel attention (SiamSA) is proposed for vision-based UAM approaching.
- A new pairwise scale-channel attention network and a scale-aware anchor proposal network are proposed to solve the scale variation issues in UAM approaching.
- A novel UAMT100 benchmark with precise annotations is built for evaluating vision-based UAM tracking methods for approaching.
- Exhaustive experiments on the new UAMT100 benchmark, the authoritative aerial tracking benchmark, and real-world tests have verified the practicality and effectiveness of SiamSA for efficient UAM approaching.

II. RELATED WORKS

A. Object tracking for UAM approaching

In numerous practical vision-based UAM approaching scenarios, continuous perception of the object location has been critical, resulting in urgent demand for tracking methods' efficiency and practicability. G. Garimella *et al.* [4] use a detection method for UAM tracking when dealing with aerial pick-and-place. However, the detection method relies

on LED markers, which has constrained the automation of UAM. During the research on picking and delivery of magnetic objects, A. Gawel *et al.* [5] perform object detection instead of tracking methods to track the object's center of gravity. They also find the deficiency of detection methods for UAM approaching and prepare to implement object tracking. Similarly, J. M. Gómez-de Gabriel *et al.* [2] adopt an object detection method on UAM in the wristband placement task. Such kind of detection methods depends on predefined class labels, which cannot track various kinds of object and brings inflexibility. Besides, collecting specific datasets for training model-based methods also causes enormous extra development costs. To better accomplish UAM tracking requirements and promote automation to a higher degree, this work introduces the model-free Siamese tracking methods for UAM approaching.

B. Siamese tracking for aerial systems

Model-free Siamese tracking has been prevalent in the object tracking field due to its SOTA efficiency and robustness. SiamRPN [18] introduces region proposal network (RPN) into Siamese tracking to increase precision and robustness with anchors. SiamRPN++ [6] employs deeper features and further improves effectiveness. SiamFC++ [7] abandons anchors by regressing offsets to reduce hyper-parameters and classification errors. Afterward, SiamAPN [8] improves Siamese network for UAV tracking, which streamlines anchor generation strategy by regressing offsets, thereby reducing hyper-parameters for UAV tracking with high efficiency. However, the scale variation issue in UAM approaching has posed a formidable challenge for directly employing Siamese tracking and requires an urgent solution. On the other hand, attention-based approaches have been widely used to capture specialized information in object tracking. These methods can tell where to focus inside the region of interest while also improving the tracking object's representation. RASNet [17] designs residual, channel, and general attention modules to learn corresponding information. SiamAttn [16] introduces deformable self-attention and cross-attention to learn context information and contextual inter-dependencies. These methods mainly infer spatial or channel attention maps, while this study focuses on scale variation issues and proposes scale

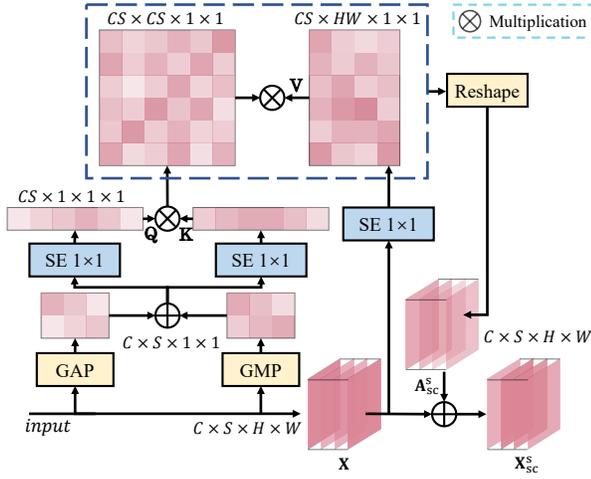


Fig. 3. Main workflow of scale-channel self-attention strategy. The refined features are attached with scale information to better estimate object scale. attention for UAM approaching.

III. PROPOSED METHOD

The detailed introduction of SiamSA is mainly divided into three parts. First, scale-equivariant convolution is briefly reviewed. Afterward, the pairwise scale-channel attention network (PSAN) and the scale-aware anchor proposal network (SA-APN) are discussed in detail respectively. Figure 2 demonstrates an overview of the proposed SiamSA method for UAM approaching.

A. Scale-equivariant convolution

Scale-equivariant (SE) convolution [19] in SiamSA enables UAM object tracking with preliminary object scale awareness for approaching. Common convolutional neural networks (CNNs) lack corresponding techniques for SV difficulties. However, SE convolution can tackle the problem with high computational efficiency, as described below:

$$[f \star_H \kappa_\sigma](s, t) = \sum_{s'} [f(s', \cdot) \star \kappa_{s, \sigma}(s^{-1}s', \cdot)](t) \quad , \quad (1)$$

where $f \star_H(s, t)$ is a function of scale s and translation t . $\kappa_\sigma(s, t)$ stands for a kernel. Besides, \star represents common convolution operation, and s' represents SV degree, where $s' > 1$ means upscale while $s' < 1$ denotes downscale. In this step, a preliminary perception of object scale is obtained by adding an additional scale dimension to image features. **Remark 1:** Based on depthwise correlation in [6], correlation is improved as SE-DW-correlation for UAM tracking.

B. Pairwise scale-channel attention network

PSAN aims to represent scale-aware feature maps during severe scale variation in UAM approaching. PSA utilizes correlation results and feature maps from SA-APN to excavate object scale information in PSAN, as shown in Fig. 2.

Given input scaled feature maps $\mathbf{X} \in \mathbb{R}^{C \times S \times H \times W}$ from the SE-DW-correlation, where superscript S represents the additional dimension related to scale, scale-channel self-attention is first employed. B. Li *et al.* [6] demonstrate that a single channel in feature maps may focus on a specific

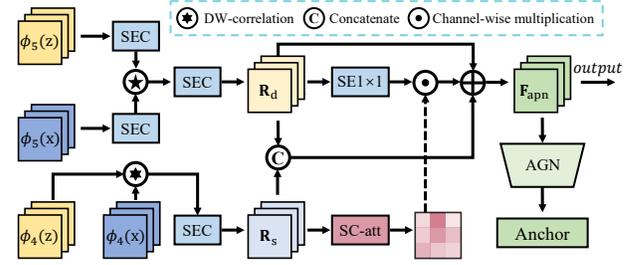


Fig. 4. Main framework of scale-aware anchor proposal network (SA-APN). The output feature maps are sent to form pairwise scale-channel attention, and the generated anchors are equipped with scale awareness.

class of objects, while the proposed SC self-attention takes advantage of this observation, and exploits scale information among different channels for stronger scale awareness. Figure 3 demonstrates the general pipeline of SC self-attention. First, through the operation of global average pooling (GAP) and global max pooling (GMP) respectively, the input 4-dimensional features $\mathbf{X} \in \mathbb{R}^{C \times S \times H \times W}$ are turned into 2-dimension. Afterward, 2-dimensional features pass a fast 1×1 scale convolution layer and are flattened to 1-dimension to be query \mathbf{Q} and key \mathbf{K} . \mathbf{Q} is generated as:

$$\mathbf{Q} = \text{Flat}([f \star_H \kappa'_\sigma](s, \text{GAP}(\mathbf{X}); \text{GMP}(\mathbf{X}))) \quad , \quad (2)$$

where $\mathbf{Q} \in \mathbb{R}^{CS}$, and $[f \star_H \kappa'_\sigma]$ denotes fast 1×1 scale convolution. Value \mathbf{V} is acquired directly from the feature maps with fast 1×1 scale convolution, where value is reshaped to 2-dimension $\mathbf{V} \in \mathbb{R}^{CS \times HW}$. Therefore, the scale-channel self-attention \mathbf{A}_{sc}^s is represented as:

$$\mathbf{A}_{sc}^s(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}(\mathbf{Q}\mathbf{K}^T)\mathbf{V} \quad , \quad (3)$$

where the superscript s denotes self-attention.

The generated scale-channel self-attention refines the feature maps following the formula below:

$$\mathbf{X}_{sc}^s = \mathbf{X} + \gamma^s \mathbf{A}_{sc}^s \quad , \quad (4)$$

where $\mathbf{X}_{sc}^s \in \mathbb{R}^{C \times S \times H \times W}$ represents feature maps with scale-channel self-attention and γ denotes a weight.

Remark 2: The same operations are performed on feature maps from both correlation results and SA-APN in a pairwise structure, exploiting scale information from both aspects.

For further learning significant scale clues, the internal relationship between feature maps from SE-DW-correlation and SA-APN is worth learning by the cross-attention strategy. As scale-channel self-attention are generated, query $\mathbf{Q}_c \in \mathbb{R}^{CS}$, key $\mathbf{K}_c \in \mathbb{R}^{CS}$ from the correlation and value $\mathbf{V}_a \in \mathbb{R}^{CS \times hw}$ from SA-APN are ready for cross-attention:

$$\mathbf{A}_{sc}^c(\mathbf{Q}_c, \mathbf{K}_c, \mathbf{V}_a) = \text{Softmax}(\mathbf{Q}_c\mathbf{K}_c^T)\mathbf{V}_a \quad , \quad (5)$$

where superscript c means cross-attention. Similar to Eq. 4, the scale-channel cross-attention $\mathbf{A}_{sc}^c \in \mathbb{R}^{CS \times hw}$ assists in the following generation of the refined feature maps \mathbf{X}_{sc}^c :

$$\mathbf{X}_{sc}^c = \mathbf{X} + \gamma^c \mathbf{A}_{sc}^c \quad , \quad (6)$$

where $\mathbf{X}_{sc}^c \in \mathbb{R}^{C \times S \times H \times W}$ refers to feature maps from the SA-APN with scale-channel cross-attention related to PSAN. **Remark 3:** Cross-attention integrates the useful scale information, achieving better estimation of the object scale for UAM approaching.

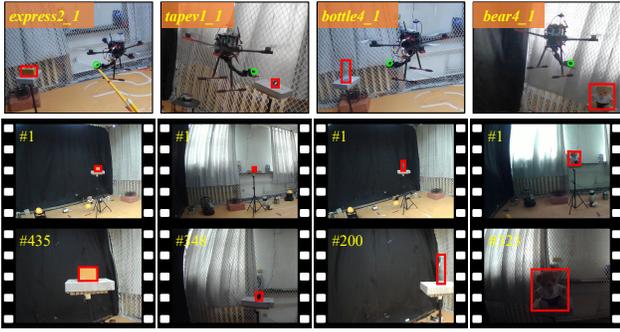


Fig. 5. Examples of scenes from UAMT100. The first row is from a fixed camera while the others are from the UAM onboard camera. The scale of the objects suffers drastic variation during UAM approaching. The red box and green circle denote the object and the onboard camera, respectively.

C. Scale-aware anchor proposal network

To attach meaningful scale information to anchor proposing, the scale-channel cross-attention strategy is also employed to form a scale-aware anchor proposal network (SA-APN). As shown in Fig. 4, ϕ_i denotes the backbone layer corresponding to the subscript $i \in \{4, 5\}$, x refers to the search region, z points to the template, and AGN denotes the anchor generation network, which consists of two convolution layers. Since the feature maps from shallower layers generally represent meaningful spatial information while the ones from deeper layers contain richer semantic clues, the scale-channel cross-attention is adopted to aggregate the information for stronger scale awareness. First, the SE-DW-correlation results \mathbf{R}_d from deeper layers are acquired and expanded with a scale dimension by scale-equivariant convolution, so as the results of depthwise correlation \mathbf{R}_s from the shallower layers. Second, scale-channel cross-attention \mathbf{A}_{apn} and concatenation \mathbf{C}_{apn} are calculated respectively. By adding weights to each item, refined features \mathbf{F}_{apn} are acquired by the following formula:

$$\mathbf{F}_{apn} = \mathbf{R}_d + \lambda_1 \mathbf{A}_{apn} + \lambda_2 \mathbf{C}_{apn}, \quad (7)$$

where \mathbf{C}_{apn} is the feature concatenation of both shallow and deep correlation results, and λ refers to a weight. Third, as the refined features are obtained, the features will pass to PSA, as shown in Fig. 2. On the other hand, scale-aware anchors are also acquired by AGN.

Remark 4: In Eq. 7, the design of SA-APN is concerned with scale information of features from the last two backbone layers. SC cross-attention focuses on scale information particularly, while concatenation directly combines the information between two layers in general as a supplement.

IV. EXPERIMENT

In this section, the proposed benchmark to evaluate object tracking methods for UAM approaching, *i.e.*, UAMT100, is first introduced. A comparison between object tracking in UAM approaching and general aerial scenes is described. Details about the evaluation experiments on UAMT100 and the authoritative aerial tracking benchmark, *i.e.*, UAV123@10fps [20], are also given. Finally, real-world tests of SiamSA are demonstrated.

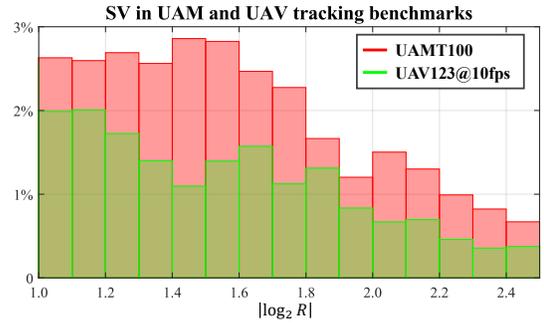


Fig. 6. SV comparison plots of UAMT100 and UAV123@10fps. The larger area of the histogram represents more severe SV issues.

A. UAMT100 benchmark

UAMT100 is recorded to evaluate object tracking methods for UAM approaching. The benchmark consists of 100 image sequences, which are captured from a flying UAM platform. It represents various scenarios of UAM’s tracking an object for approaching. 16 kinds of objects are involved, and 11 attributes are annotated for each sequence. The attributes include common challenges of object tracking, *i.e.*, aspect ratio change (ARC), out-of-view (OV), background clutter (BC), fast motion (FM), low illumination (LI), object blur (OB), partial occlusion (POC), scale variation (SV), similar object (SOB), and viewpoint change (VC). Besides, some challenging scenarios with the UAM attack (UAM-A) by a stick are also considered especially for UAM tracking. The videos are recorded at 10 frames per second (FPS), with the JPG image resolution of 800×600 pixels.

Remark 5: The system to record UAMT100 benchmark is described as follows. The OptiTrack¹ Flex 13 camera from Quanser² acquires the pose information of UAM and reports it to the NVIDIA Xavier NX through ROS³ client nodes for the VRPN library. Besides, the communication between the onboard computer and Pixhawk relies on serial. Moreover, QGroundControl acts as the ground control station.

Figure 6 quantitatively demonstrates the difference in SV issue between UAM object tracking for approaching and common UAV tracking. R denotes the degree of SV, which is measured by the ratio of the current object’s ground truth bounding box area to the initial one. SV is measured when R is outside the range $[0.5, 2]$, *i.e.*, $|\log_2 R| > 1$. The percentage of frames whose $|\log_2 R|$ is with a certain section is drawn as the SV histogram, with an interval length of 0.1 over the range of 1 to 2.5. The proportion of sections where $|\log_2 R| > 2.5$ is less than 0.5% and not of reference significance. Therefore, the larger area of the histogram means the higher frequency of object SV, and during UAM tracking for approaching, SV is more common and more severe than in UAV tracking.

B. Implementation details

In SiamSA, The last three layers of the backbone are fine-tuned, and the input size of the template and search image

¹<https://optitrack.com/>

²<https://www.quanser.com/>

³<https://www.ros.org/>

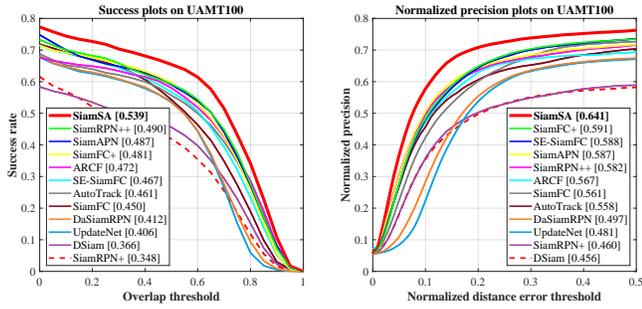


Fig. 7. Overall performance evaluation on UAMT100 benchmark. SiamSA surpasses all other trackers in AUC score and NP.

are set to 127×127 and 287×287 pixels respectively. With stochastic gradient descent (SGD), SiamSA is trained on common object tracking training sets, including COCO [21], ImageNet VID [22], GOT-10k [23], and Youtube-BB [24], without extra specific training dataset. All training and evaluating process is implemented with an Intel i9-9920X CPU, a 32GB RAM, and two NVIDIA TITAN RTX GPUs.

C. Evaluation metrics

The evaluation complies with the classic standard of one-pass evaluation (OPE) [25], including success rate and normalized precision (NP) [26]. Success rate reflects intersection over union (IoU) score, while NP is concerned with the percentage of frames where the normalized distance between the estimated and ground truth positions is within a threshold. The area under the curve (AUC) of the plot is used to rank the trackers in the experiment. Notably, NP is less vulnerable to object scale and can better demonstrate tracking robustness compared with the traditional precision metric. Therefore, NP is adopted for UAM approaching.

D. Evaluation on UAMT100 benchmark

The evaluation of UAMT100 consists of overall performance, attribute-based performance, comparison with deeper backbones, and ablation study.

Remark 6: To be fair, [6], [8], and [27] are equipped with AlexNet [28]. All parameters of trackers are consistent with the paper, and all trackers are evaluated on the same platform.

1) *Overall performance:* As shown in Fig. 7, SiamSA outperforms other 11 SOTA trackers [6] [8] [13] [15] [19] [27] [29]–[32], including Siamese-based and correlation filter-based ones. Compared with the second-best performance, gains on AUC score and NP are **10.0%** and **8.5%** respectively, which represents the effectiveness of the proposed PSA for UAM tracking.

TABLE I

COMPARISON WITH DEEPER BACKBONES ON UAMT100. THE BEST TWO PERFORMANCES ARE HIGHLIGHTED WITH RED AND BLUE COLORS.

Trackers	Source	Backbone	AUC	NP	FPS
Ocean	ECCV 2020	ResNet50	0.418	0.525	95
SiamBAN	CVPR 2020	ResNet50	0.494	0.582	67
SiamMask	CVPR 2019	ResNet50	0.506	0.605	72
SiamFC++	AAAI 2020	GoogleNet	0.521	0.614	106
LightTrack	CVPR 2021	SuperNet	0.536	0.624	56
SiamSA	Ours	AlexNet	0.539	0.641	122

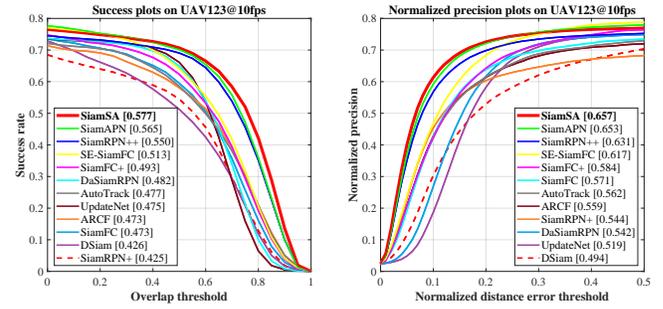


Fig. 8. Overall performance evaluation on UAV123@10fps. SiamSA yields competitive performance in terms of AUC score and NP.

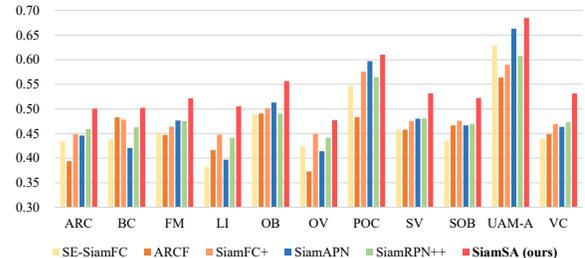


Fig. 9. Attribute-based evaluation on UAMT100. SiamSA performs best in all challenges.

2) *Attribute-based performance:* Figure 9 demonstrates the attribute-based comparison of 6 top trackers in UAMT100 benchmark. SiamSA ranks first in all the attributes. In challenges that are important to UAM approaching, *e.g.*, SV, ARC, and OV, SiamSA is particularly outstanding, surpassing the second place by **10.6%**, **8.9%**, and **6.2%**. Besides, on UAM attack, SiamSA also outperforms the second-best performance by **3.3%**. The comparison validates that the proposed PSA is competent for scale variation issues, and SiamSA can handle the challenges that are faced in object tracking for UAM approaching.

3) *Comparison with deeper backbones:* To validate the efficiency of UAM approaching, SiamSA is also compared with 5 SOTA trackers with deeper backbones. Notably, deeper backbones generally bring better performance but heavier computational burdens, which can not meet the requirement of real-time tracking. In TABLE I, despite SiamSA is equipped with the lightweight AlexNet, it outperforms other 5 trackers [7] [33]–[36] with SuperNet, ResNet50, or GoogLeNet. Specifically, SiamSA is **2.18** \times faster than the second-best tracker, while exceeding the second-fastest tracker by **3.5%**. The comparison shows that SiamSA targets well at UAM tracking efficiently.

4) *Ablation study:* TABLE II validates the effectiveness of the proposed pairwise scale-channel attention. With only the SA-APN, anchors that adapt to the object scale are generated for object locating, which promotes the success

TABLE II

ABLATION STUDY OF SIAMSA ON UAMT100. Δ DENOTES GAINS.

Trackers	AUC	$\Delta_{auc}(\%)$	NP	$\Delta_{np}(\%)$
Baseline	0.487	-	0.587	-
Baseline+SA-APN	0.504	3.49%	0.610	3.91%
Baseline+PSAN	0.512	5.13%	0.625	6.47%
Baseline+SA-APN+PSAN (SiamSA)	0.539	10.68%	0.641	9.20%

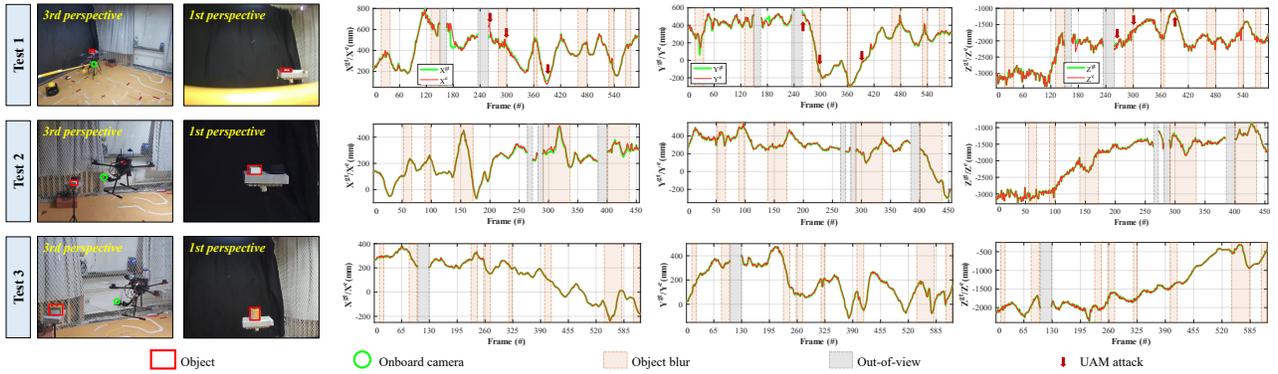


Fig. 10. Real-world vision-based UAM approaching tests with SiamSA for object tracking. The estimated positions X^e , Y^e and scale S^e (red line) is compared with the ground truth X^{gt} , Y^{gt} , and S^{gt} (green line) according to the pinhole camera model. Demo videos of the real-world tests are also available at <https://github.com/vision4robotics/SiamSA>.

rate by 3.49%. On the other hand, with PSAN, the feature maps are equipped with scale awareness, which improves the tracker by 5.13%. Combining both modules, scale attention works in a pairwise structure, which further excavates the scale information to track the object robustly and is more competent for vision-based UAM approaching.

E. Evaluation on the UAV benchmark

To validate the generality of the proposed PSA, SiamSA is also evaluated on UAV123@10fps with 11 SOTA trackers.

1) *Overall performance*: The overall performance on UAV123@10fps [20] is shown in Fig. 8. The benchmark contains SV and other challenges from aerial perspectives, which provides a reference for the generality of the proposed PSA. Compared with the second-best performance, SiamSA improves **2.1%** on AUC score attributed to PSA.

2) *Attribute-based performance*: Various challenges in UAV benchmarks have reference significance for analyzing UAM tracking. Especially, SV, ARC, FM, and OV, which also saliently influence UAM tracking, are discussed. Performance on these attributes is displayed in TABLE III. PSA contributes to the promotion of SiamSA tracker on the SV issue, with an increase of **1.7%** on AUC score. Compared with baseline, AUC score increase on ARC, CM, and OV are **1.5%**, **3.6%**, and **6.2%**, which validate the effectiveness and generality of SiamSA against these challenges.

V. REAL-WORLD UAM APPROACHING TESTS

Real-world tests of SiamSA on three UAM approaching scenes are shown in Fig. 10. In these tests, SiamSA tracks the object robustly and accurately with over 10 FPS on a UAM, which is equipped with an NVIDIA Jetson Xavier NX and an

onboard camera. TensorRT acceleration is performed locally. In test 1, a yellow stick is employed to attack the UAM three times during approaching a battery. The attacks generally damage the tracking stability by causing sudden deviations in the tracking results. Test 2 confronts periods of UAM fast motion, which bring severe object blur and negatively impact the object perception. In test 3, the object out-of-view also poses a formidable challenge for stable tracking. Notably, all the tests suffer from severe object scale variation during the UAM approaching the object. But SiamSA manages to handle these issues and represents strong robustness.

Remark 7: Three real-world vision-based UAM approaching tests of SiamSA further prove the practicality and effectiveness of SiamSA in practical UAM tracking.

VI. CONCLUSION

In this work, a Siamese network with pairwise scale-channel attention (SiamSA) is proposed especially for object tracking in vision-based UAM approaching. Specifically, to address the severe scale variation issues in UAM approaching, the novel pairwise scale-channel attention (PSA) is proposed. Based on PSA, a pairwise scale-channel attention network and a scale-aware anchor proposal network are designed. The former mainly refines image features with scale information, while the latter proposes anchors that adapt to the object scale. To fairly evaluate object tracking methods for vision-based UAM approaching, UAMT100 benchmark is recorded with a flying UAM platform. Comprehensive experiments on the proposed UAM tracking benchmark and the authoritative UAV tracking benchmark have validated the effectiveness, efficiency, and generality of SiamSA. Furthermore, the real-world tests with a flying UAM and an onboard camera also prove the practicality of SiamSA for various applications. It is convinced that both SiamSA and the UAMT100 benchmark will facilitate the improvement of vision-based UAM-approaching-related applications.

ACKNOWLEDGMENT

This work is supported by the National Natural Science Foundation of China (No. 62173249), the Natural Science Foundation of Shanghai (No. 20ZR1460100), and the Key R&D Program of Sichuan Province (No. 2020YFSY0004).

TABLE III

ATTRIBUTE-BASED EVALUATION ON UAV123@10FPS. THE BEST TWO RESULTS ARE HIGHLIGHTED WITH RED AND BLUE COLORS.

Attribute	SV		OV		ARC		FM	
	AUC	NP	AUC	NP	AUC	NP	AUC	NP
DaSiamRPN	0.463	0.529	0.423	0.519	0.451	0.530	0.380	0.467
SiamFC+	0.460	0.554	0.425	0.554	0.429	0.534	0.377	0.490
SE-SiamFC	0.483	0.590	0.453	0.597	0.461	0.583	0.401	0.522
SiamRPN++	0.522	0.604	0.474	0.589	0.498	0.596	0.431	0.532
SiamAPN	0.542	0.633	0.504	0.629	0.519	0.623	0.506	0.628
SiamSA (ours)	0.551	0.634	0.535	0.652	0.527	0.620	0.524	0.623

REFERENCES

- [1] K. M. Popek, M. S. Johannes, K. C. Wolfe, R. A. Hegeman, J. M. Hatch, J. L. Moore, K. D. Katyal, B. Y. Yeh, and R. J. Bamberger, "Autonomous Grasping Robotic Aerial System for Perching (AGRASP)," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2018, pp. 1–9.
- [2] J. M. Gómez-de Gabriel, J. M. Gandarias, F. J. Pérez-Maldonado, F. J. García-Núñez, E. J. Fernández-García, and A. J. García-Cerezo, "Methods for Autonomous Wristband Placement with a Search-and-Rescue Aerial Manipulator," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2018, pp. 7838–7844.
- [3] G. Heredia, A. Jimenez-Cano, I. Sanchez, D. Llorente, V. Vega, J. Braga, J. Acosta, and A. Ollero, "Control of a Multirotor Outdoor Aerial Manipulator," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2014, pp. 3417–3422.
- [4] G. Garimella and M. Kobilarov, "Towards Model-Predictive Control for Aerial Pick-and-Place," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2015, pp. 4692–4697.
- [5] A. Gawel, M. Kamel, T. Novkovic, J. Widauer, D. Schindler, B. P. von Altshofen, R. Siegart, and J. Nieto, "Aerial Picking and Delivery of Magnetic Objects with MAVs," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2017, pp. 5746–5752.
- [6] B. Li, W. Wu, Q. Wang, F. Zhang, J. Xing, and J. Yan, "SiamRPN++: Evolution of Siamese Visual Tracking with Very Deep Networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4282–4291.
- [7] Y. Xu, Z. Wang, Z. Li, Y. Yuan, and G. Yu, "SiamFC++: Towards Robust and Accurate Visual Tracking with Target Estimation Guidelines," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2020, pp. 12 549–12 556.
- [8] C. Fu, Z. Cao, Y. Li, J. Ye, and C. Feng, "Siamese Anchor Proposal Network for High-Speed Aerial Tracking," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2021, pp. 510–516.
- [9] J. Ye, C. Fu, G. Zheng, D. P. Paudel, and G. Chen, "Unsupervised Domain Adaptation for Nighttime Aerial Tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 8896–8905.
- [10] Z. Cao, Z. Huang, L. Pan, S. Zhang, Z. Liu, and C. Fu, "TC-Track: Temporal Contexts for Aerial Tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 14 798–14 808.
- [11] J. Ye, C. Fu, F. Lin, F. Ding, S. An, and G. Lu, "Multi-Regularized Correlation Filter for UAV Tracking and Self-Localization," *IEEE Transactions on Industrial Electronics*, vol. 69, no. 6, pp. 6004–6014, 2022.
- [12] G. Zheng, C. Fu, J. Ye, F. Lin, and F. Ding, "Mutation Sensitive Correlation Filter for Real-Time UAV Tracking with Adaptive Hybrid Label," in *Proceedings of the International Conference on Robotics and Automation (ICRA)*, 2021, pp. 503–509.
- [13] Y. Li, C. Fu, F. Ding, Z. Huang, and G. Lu, "AutoTrack: Towards High-Performance Visual Tracking for UAV with Automatic Spatio-Temporal Regularization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 11 920–11 929.
- [14] C. Fu, B. Li, F. Ding, F. Lin, and G. Lu, "Correlation Filters for Unmanned Aerial Vehicle-Based Aerial Tracking: A Review and Experimental Evaluation," *IEEE Geoscience and Remote Sensing Magazine*, vol. 10, no. 1, pp. 125–160, 2022.
- [15] Z. Huang, C. Fu, Y. Li, F. Lin, and P. Lu, "Learning Aberrance Repressed Correlation Filters for Real-Time UAV Tracking," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 2891–2900.
- [16] Y. Yu, Y. Xiong, W. Huang, and M. R. Scott, "Deformable Siamese Attention Networks for Visual Object Tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 6728–6737.
- [17] Q. Wang, Z. Teng, J. Xing, J. Gao, W. Hu, and S. Maybank, "Learning Attention: Residual Attentional Siamese Network for High Performance Online Visual Tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 4854–4863.
- [18] 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/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 8971–8980.
- [19] I. Sosnovik, A. Moskalev, and A. W. Smeulders, "Scale Equivariance Improves Siamese Tracking," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2021, pp. 2765–2774.
- [20] M. Mueller, N. Smith, and B. Ghanem, "A Benchmark and Simulator for UAV Tracking," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016, pp. 445–461.
- [21] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common Objects in Context," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014, pp. 740–755.
- [22] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein et al., "Imagenet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [23] L. Huang, X. Zhao, and K. Huang, "GOT-10k: A Large High-Diversity Benchmark for Generic Object Tracking in the Wild," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 5, pp. 1562–1577, 2021.
- [24] E. Real, J. Shlens, S. Mazzocchi, X. Pan, and V. Vanhoucke, "Youtube-Boundingboxes: A Large High-Precision Human-Annotated Data Set for Object Detection in Video," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5296–5305.
- [25] Y. Wu, J. Lim, and M.-H. Yang, "Online Object Tracking: A Benchmark," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013, pp. 2411–2418.
- [26] M. Muller, A. Bibi, S. Giancola, S. Alsubaihi, and B. Ghanem, "TrackingNet: A Large-Scale Dataset and Benchmark for Object Tracking in the Wild," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 300–317.
- [27] Z. Zhu, Q. Wang, B. Li, W. Wu, J. Yan, and W. Hu, "Distractor-Aware Siamese Networks for Visual Object Tracking," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 101–117.
- [28] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2012, pp. 1097–1105.
- [29] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. Torr, "Fully-Convolutional Siamese Networks for Object Tracking," in *Proceedings of the European Conference on Computer Vision Workshops (ECCVW)*, 2016, pp. 850–865.
- [30] Z. Zhang and H. Peng, "Deeper and Wider Siamese Networks for Real-Time Visual Tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4591–4600.
- [31] L. Zhang, A. Gonzalez-Garcia, J. v. d. Weijer, M. Danelljan, and F. S. Khan, "Learning the Model Update for Siamese Trackers," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 4010–4019.
- [32] Q. Guo, W. Feng, C. Zhou, R. Huang, L. Wan, and S. Wang, "Learning Dynamic Siamese Network for Visual Object Tracking," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 1763–1771.
- [33] B. Yan, H. Peng, K. Wu, D. Wang, J. Fu, and H. Lu, "LightTrack: Finding Lightweight Neural Networks for Object Tracking via One-Shot Architecture Search," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 15 180–15 189.
- [34] Q. Wang, L. Zhang, L. Bertinetto, W. Hu, and P. H. Torr, "Fast Online Object Tracking and Segmentation: A Unifying Approach," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1328–1338.
- [35] Z. Chen, B. Zhong, G. Li, S. Zhang, and R. Ji, "Siamese Box Adaptive Network for Visual Tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 6667–6676.
- [36] Z. Zhang, H. Peng, J. Fu, B. Li, and W. Hu, "Ocean: Object-Aware Anchor-Free Tracking," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 771–787.