# Scale Optimization Using Evolutionary Reinforcement Learning for Object Detection on Drone Imagery

**Jialu Zhang[1,2], Xiaoying Yang[1*], Wentao He[1], Jianfeng Ren[1,3†], Qian Zhang[1,3], Yitian Zhao[2], Ruibin Bai[1,3], Xiangjian He[1,3], Jiang Liu[2,4]**

[1]The Digital Port Technologies Lab, School of Computer Science, University of Nottingham Ningbo China
[2]Cixi Institute of Biomedical Engineering, Ningbo Institute of Industrial Technology, Chinese Academy of Sciences
[3]Nottingham Ningbo China Beacons of Excellence Research and Innovation Institute, University of Nottingham Ningbo China
[4]Department of Computer Science and Engineering, Southern University of Science and Technology
{sgxjz1,scxxy1,scxwh1,jianfeng.ren,qian.zhang, ruibin.bai,sean.he}@nottingham.edu.cn, yitian.zhao@nimte.ac.cn,
liuj@sustech.edu.cn

## Abstract

Object detection in aerial imagery presents a significant challenge due to large scale variations among objects. This paper proposes an evolutionary reinforcement learning agent, integrated within a coarse-to-fine object detection framework, to optimize the scale for more effective detection of objects in such images. Specifically, a set of patches potentially containing objects are first generated. A set of rewards measuring the localization accuracy, the accuracy of predicted labels, and the scale consistency among nearby patches are designed in the agent to guide the scale optimization. The proposed scale-consistency reward ensures similar scales for neighboring objects of the same category. Furthermore, a spatial-semantic attention mechanism is designed to exploit the spatial semantic relations between patches. The agent employs the proximal policy optimization strategy in conjunction with the evolutionary strategy, effectively utilizing both the current patch status and historical experience embedded in the agent. The proposed model is compared with state-of-the-art methods on two benchmark datasets for object detection on drone imagery. It significantly outperforms all the compared methods. Code is available at https://github.com/UNNC-CV/EvOD/.

## Introduction

Unmanned Aerial Vehicles have been widely used in various applications, *e.g.*, surveillance (Yun et al. 2022), autonomous detection (Ren and Jiang 2017, 2021), fleet navigation (Alami et al. 2023) and agriculture (Tokekar et al. 2016). Object detection from drone-captured images has attracted research attention recently (Xi et al. 2021, 2020; Bouguettaya et al. 2022). Although object detection on natural images has progressed significantly (Ge et al. 2021), detecting objects in aerial images remains challenging, mainly stemming from small scales and extreme scale variations (Deng et al. 2021; Xu, Li, and Wang 2022).

Objects in aerial scenes often have large-scale variations, *e.g.*, distant objects occupy few pixels while nearby objects occupy thousands. To tackle the challenges of detecting small objects and/or objects of different sizes, a common

---

*Equal Contribution.
†Corresponding author.

Figure 1: The number of objects (*y-axis*) that are optimally detected using the scaling factor (*x-axis*) for ultra-small, small, medium and large objects, respectively on the VisDrone dataset. *The optimal scales are significantly different for different objects.*

strategy is to divide an image into patches, scale the patches containing small objects to a fixed size (Deng et al. 2021; Xu, Li, and Wang 2022) or using one or more fixed scaling factors (Huang, Chen, and Huang 2022), and then feed them into an object detector. But the patch scalability is inherently limited due to the potential image artifacts caused by excessive scaling. Moreover, patches may encompass objects of different sizes. While enlarging a patch improves detecting small objects, it also enlarges large objects, potentially impeding their recognition. As shown in Fig. 1, the optimal scales for different objects vary significantly. It is hence crucial to determine the optimal scale of each patch.

However, there lacks ground-truth annotations for the optimal scales. To tackle this problem, an EVOlutionary Reinforcement Learning (EVORL) agent is designed to determine the most suitable scale for each patch, with the guidance of a carefully designed reward function. This function assesses the image patch by considering the localization accuracy, the accuracy of predicted labels, and the scale consistency among nearby patches. The first two are directly related to the performance of object detection while the last one regularizes the optimized scales. This scale consistency stems from the inherent characteristics of drone imagery, where nearby objects of the same category tend to exhibit a

similar scale. By rewarding the scale consistency, the agent is able to eliminate outliers influenced by incidental factors, thereby contributing to an improved detection performance.

Simultaneously optimizing the three rewards may result in potential conflicts, complicate the training convergence, and limit the performance. To mitigate this issue, an evolutionary strategy is integrated into the reinforcement learning framework. Specifically, the optimal scales of all patches during training are combined with sampled historical solutions to form an initial population. The proposed evolutionary algorithm refines the optimal scale determined by the current patch status by evolving the solutions using mutation and crossover, taking into account of the scale consistency among nearby patches. By incorporating both the current patch status and past experience stored in the agent's population, the proposed EVORL effectively determines the optimal scale for precise object detection.

To further boost the detection performance, a spatial-semantic attention is developed. Intuitively, spatially close objects could not only exhibit the scale consistency, but also provide the spatial-semantic attention to mutually enhance the patch features (He et al. 2023; Zhang et al. 2023). Specifically, the proposed method models the spatial and semantic attention by measuring the distances and the pairwise appearance correlations between adjacent objects, respectively, and aggregates these two to obtain the spatial-semantic attention. The proposed spatial-semantic attention could effectively model the spatial and semantic dependencies between objects, enhance the patch features and finally help to better detect objects at a most appropriate scale.

The proposed method follows a coarse-to-fine object detection pipeline (Bouguettaya et al. 2022). Specifically, a YOLOX (Ge et al. 2021) variant is utilized to coarsely generate regions of interests. These regions are expanded to include the background context and merged to form cluster regions as in (Huang, Chen, and Huang 2022). A feature extractor with the proposed spatial-semantic attention is designed to visually perceive the regions. The perceived information is transmitted to the proposed EVORL agent to determine the optimal scale for each region, with the guidance of the three carefully designed rewards. Finally, the scaled regions are fed back to the detector for fine detection.

Our contributions can be summarized as follows. 1) The proposed EVORL agent is seamlessly integrated into a coarse-to-fine object detection framework, and makes use of both the current image patch and the past experience embedded in the agent to determine the optimal scale to accurately detect objects. 2) The designed reward function well addresses the challenges of lacking ground-truth labels for optimal scales, and provides supervision signals to train the agent. The proposed scale-consistency reward considers the scales of both the current object and nearby objects, to eradicate outliers and enhance the detection performance. 3) The proposed spatial-semantic attention exploits the spatial and semantic relations between nearby patches, to enhance the discriminant power of patch features. 4) The proposed method significantly outperforms state-of-the-art methods for object detection, improving the previous best average precision from 24.6% to 28.0% on the UAVDT dataset, and

from 40.3% to 42.2% on the VisDrone dataset.

## Related Work

### Object Detection on Drone Imagery

In aerial images, there are a large number of small objects, *e.g.*, 26.5% of objects in the VisDrone dataset (Zhu et al. 2022) occupying fewer than $16^2$ pixels. Researchers have strove to improve small object detection on aerial imagery by adapting general object detectors on natural images. For example, Cheng et al. (2019) designed novel objective functions for small object detection without altering existing network architectures. Li et al. (2017) developed a super-resolution technique to enlarge the image for better detecting small objects. Bai et al. (2018) utilized a generative adversarial network to obtain fine-grained features for small blurred objects. Some researchers utilized the shallow layers of deep neural networks to alleviate the problems of low resolution and detail loss caused by down-sampling operations (Bouguettaya et al. 2022), *e.g.*, Sommer, Schuchert, and Beyerer (2017) used high-resolution feature maps from earlier layers to enhance detection performance.

Some researchers tackled the challenges of large scale variations. Wang et al. (2019) introduced a Receptive Field Expansion Block and a Spatial-Refinement Module to capture context information and refine solutions using multi-scale pyramid features. Zhang et al. (2019) developed a scale-adaptive proposal network, which consists of multi-scale region proposal networks and multi-layer feature fusion to better detect objects of different scales. The feature pyramid network is often adopted to combine low-level features from shallow layers with high-level features from deep layers for multi-scale object detection (Zhou et al. 2019).

The coarse-to-fine pipeline is often utilized for detecting objects in aerial images through extracting regions of interests using a coarse detector, scaling the image patches, and then detecting objects within them (Bouguettaya et al. 2022). Ozge Unel, Ozkalayci, and Cigla (2019) uniformly divided the high-resolution image into patches of a fixed size, and detected objects from patches. Yang et al. (2019) designed a network to crop regions of dense objects and a scale estimation network to resize the crops. Xu, Li, and Wang (2022) developed a self-adaptive region selection algorithm to focus on the dense regions, and leveraged super-resolution to enlarge the focused regions to a fixed size before fine-grained detection. Huang, Chen, and Huang (2022) first equalized the scales of all generated patches, and then fed them into a unified mosaic for inference.

Although scaling is critical to object detection, existing solutions often scale the patches to a fixed size (Xu, Li, and Wang 2022) or using fixed scaling factors (Huang, Chen, and Huang 2022). Optimal scaling has not been fully exploited.

### General Object Detection

General object detectors are often adapted for drone imagery (Cai and Vasconcelos 2018). Depending on the way of feature extraction, object detectors can be broadly divided into traditional methods and deep learning methods. Traditional methods often utilize handcrafted features such

as local binary patterns (Ren, Jiang, and Yuan 2015), scale-invariant key-points (Lowe 2004) and histograms of oriented gradients (Dalal and Triggs 2005). These handcrafted features are often task-specific, and ineffective in dealing with complex real-world problems (Bouguettaya et al. 2022).

Numerous deep learning object detectors have been developed recently (Ren et al. 2017; Ge et al. 2021), and they demonstrate superior performance thanks to the discriminative deep learning features. These models could be further categorized into two types: 1) Two-stage detectors, in which regions of interests are first extracted using a region proposal network, and objects are recognized within them. Representative models include Regions with CNN features (R-CNN) (Girshick et al. 2014), Faster R-CNN (Ren et al. 2017), and Mask R-CNN (He et al. 2017). 2) One-stage detectors, which integrate the proposal generation and object detection into one stage. YOLO-series (You Only Look Once) (Redmon et al. 2016; Redmon and Farhadi 2017; Ge et al. 2021), RetinaNet (Lin et al. 2020) and EfficientDet (Tan, Pang, and Le 2020) are the leading solutions.

General object detectors perform well on natural images, but not on aerial images. Aerial images often have higher image resolution but contain much more objects of various sizes, which imposes great challenges for detecting them.

## Proposed Method

### Overview of Proposed Method

To tackle the challenges of determining the optimal scales, an evolutionary reinforcement learning agent is proposed. The agent is integrated into a coarse-to-fine object detection framework. The proposed framework mainly consists of three modules, as shown in Fig. 2. 1) **Coarse Patch Generation**. CSPDarkNet (Wang, Bochkovskiy, and Liao 2021) is utilized as the backbone to generate the feature pyramid. In addition to the small, medium and large feature maps $\mathbf{P}^S, \mathbf{P}^M, \mathbf{P}^L$ used in YOLOX (Ge et al. 2021), an ultra-small feature map $\mathbf{P}^U$ is added, which contains low-level fine details for better detecting small objects. These features are then fed into the YOLOX decoupled heads to generate regions of interests $\mathbb{B}$. 2) **Cluster Region Generation**. The contextual information from both the background and nearby objects has shown to be helpful in recognizing objects (Zhang et al. 2022, 2023). The coarsely detected regions $\mathbb{B}$ are hence expanded by a factor of $\beta$ to include the background context as $\mathbb{B}^E = \mathcal{F}_E(\mathbb{B}; \beta)$, where $\mathcal{F}_E$ and $\mathbb{B}^E$ represent the expansion function and the expanded regions, respectively. The expanded regions are then clustered and merged into a cluster region set $\mathbb{C}$ using the Foreground Region Generation (Huang, Chen, and Huang 2022). 3) **Evolutionary Reinforcement Learning**. A visual perception network is designed to visually perceive the regions, in which a spatial-semantic attention is designed to capture the spatial and semantic relations between nearby objects. Three rewards considering localization accuracy, label accuracy and scale consistency are designed to guide training, which well addresses the problem of lacking ground-truth annotations of optimal scales. To balance these three rewards, the hybrid algorithm combining the evolutionary strategy and Proximal

Policy Optimization (PPO) strategy is designed to determine the optimal scales. The regions are then scaled accordingly, packed into mosaics as in (Huang, Chen, and Huang 2022) and fed back to the detector for fine detection. Finally, post-processing techniques such as non-maximum suppression are utilized to generate the final detection results.

### Formulation of Reinforcement Learning

The scale optimization problem is formulated as a Markov Decision Process, represented by the tuple $(\mathbb{S}, \mathbb{O}, \mathbb{A}, \mathcal{R}, p_s)$.
**State** $\mathbb{S}$ refers to the set of states of the environment, specifically, the determined scaling factors of all the generated cluster regions at a specific point in time.
**Observation** $\mathbb{O}$ encompasses the vital information about the objects, *e.g.*, spatial features, semantic features, patch attributes and the attentive information from nearby objects.
**Action** $\mathbb{A} = \{a_1, \ldots, a_N\}$ consists of a set of actions for the $N$ cluster regions, where each action $a_i$ corresponds to a specific scaling action for the cluster region $\boldsymbol{C}_i \in \mathbb{C}$.
**State transition probability** $p_s$ is defined as $p_s(s'|s, a) = \Pr\{\mathbb{S}^{t+1} = s'|\mathbb{S}^t = s, \mathbb{A}^t = a\}$, indicating the likelihood of transitioning from the current state $s$ to a new state $s'$ under the execution of action $a$.
**Reward** $\mathcal{R}$ assesses the current state based on the object detection accuracy and the scale consistency among nearby objects. More details are provided later on.

### Visual Perception with Spatial-semantic Attention

The visual perception network takes the cluster regions $\mathbb{C}$ as the input, and extracts the appearance features using a patch encoder, ResNet-18 pre-trained on ImageNet. As each region contains fewer objects than the whole image, the ResNet-18 can well extract the appearance features while keeping the network lightweight. Specifically, the appearance features are derived as $\mathbf{X} = \mathcal{F}_P(\mathbb{C}; \boldsymbol{\theta})$, where $\mathcal{F}_P$ represents the network, $\boldsymbol{\theta}$ represents the network parameters, and $\mathbf{X}$ denotes all extracted features packed together.

To capture the attentive information between nearby objects, a spatial-semantic attention is designed. Specifically, the spatial attention $\boldsymbol{S}$ is explicitly modeled by the reciprocal of the distance between the centers of two objects, where each element $\boldsymbol{S}_{ij} = 1/\mathcal{F}_D(\boldsymbol{C}_i, \boldsymbol{C}_j)$, and $\mathcal{F}_D$ calculates the spatial distance between $\boldsymbol{C}_i$ and $\boldsymbol{C}_j$. Intuitively, the smaller the spatial distance, the greater the mutual spatial attention.

To model the semantic attention, the appearance features $\mathbf{X}$ are firstly projected into three embedding spaces as the query matrix $\boldsymbol{Q} = \mathcal{F}_Q(\mathbf{X}, \boldsymbol{\theta}_Q)$, key matrix $\boldsymbol{K} = \mathcal{F}_K(\mathbf{X}, \boldsymbol{\theta}_K)$, and value matrix $\boldsymbol{V} = \mathcal{F}_V(\mathbf{X}, \boldsymbol{\theta}_V)$, where $\mathcal{F}_Q$, $\mathcal{F}_K$ and $\mathcal{F}_V$ represent the transformation networks, and $\boldsymbol{\theta}_Q$, $\boldsymbol{\theta}_K$ and $\boldsymbol{\theta}_V$ represent the learnable parameters of these three networks, respectively. The semantic attention is modeled as $\mathcal{F}_S(\boldsymbol{Q}, \boldsymbol{K}) = \frac{\boldsymbol{Q} \cdot \boldsymbol{K}^\top}{\sqrt{d}}$, where $d$ is the feature dimension, and $\sqrt{d}$ ensures a stable gradient for the attention map. The proposed semantic attention makes use of the self-attention mechanism to exploit the attentive information between nearby objects, so that correlated objects are weighted more to boost the discriminant power of the target object.

Figure 2: Overview of the proposed model. A YOLOX variant is first utilized to generate regions of interests. The regions are expanded to include the background context and merged to form cluster regions. An evolutionary reinforcement learning (EVORL) agent with three rewards is designed to determine the optimal scale for each patch. The spatial-semantic attention is designed to boost the patch features. After determining the optimal scales through the proposed EVORL, the regions are scaled and consolidated into a mosaic image, and passed back to the detector for fine detection.

The spatial-semantic attention $\boldsymbol{E}$ of all clustered regions is obtained through an aggregation network $\mathcal{F}_A(\cdot)$ by,

$$\boldsymbol{E} = \mathcal{F}_A(\mathcal{F}_S(\boldsymbol{Q}, \boldsymbol{K}) \cdot \boldsymbol{S}) \cdot \boldsymbol{V}. \tag{1}$$

The proposed spatial-semantic attention well leverages on both spatial and semantic dependencies between nearby image patches, and hence effectively boosts the discriminative power of patch features with the support of nearby objects.

### Reward Function

Three types of rewards are designed to provide feedback to the agent regarding the quality of a specific scaling action. 1) **Localization Reward** $r_l$, for accurately locating the objects. Specifically, $r_l$ calculates the average Intersection over Union (IoU) between the detected bounding boxes and the ground-truth ones, and it rewards the agent for accurately locating the objects. 2) **Labeling Reward** $r_c$, for correctly classifying the objects. Specifically, $r_c$ is defined as the average classification accuracy for objects with an IoU of at least 0.5. 3) **Scale-consistency Reward** $r_s$. In aerial images, nearby objects of the same category tend to share a similar scale. $r_s$ is designed to incentivize the scale consistency. Specifically, denote the scaling factor for $\boldsymbol{C}_i$ as $\lambda_i$. To ensure the scale consistency, for each cluster region $\boldsymbol{C}_i$, we minimize the differences between the scaling factor $\lambda_i$ and that of all its $N_i$ nearby regions of the same class, $\Delta_i = \frac{1}{N_i} \sum_{j=1}^{N_i} |\lambda_i - \lambda_j^i|$, where $\lambda_j^i$ denotes the scaling factor of the $j$-th neighboring region that has the same class label as $\boldsymbol{C}_i$. The scale-consistency reward is defined as,

$$r_s = \frac{1}{N} \sum_{i=1}^{N} e^{-\Delta_i/K}, \tag{2}$$

where $K$ is a normalization factor. $r_s$ is large if the neighboring cluster regions share similar scaling factors. Note that this reward relies on not only the optimal scaling factor of the current image patch, but also that of neighbors. Thus, the decision-making process for the optimal scaling factor of each patch becomes more complex.

The first two rewards $r_l$ and $r_c$ encourage the agent to choose a scaling factor to accurately locate and recognize the objects, and the last reward $r_s$ serves as a regularization constraint to remove the outliers in scaling factors. The reward function $\mathcal{R}$ makes use of these three rewards as,

$$\mathcal{R} = \alpha_l r_l + \alpha_c r_c + \alpha_s r_s, \tag{3}$$

where $\alpha_l$, $\alpha_c$ and $\alpha_s$ are the respective weighting factors.

### Evolutionary Reinforcement Learning Strategy

The three designed rewards may conflict with each other. Jiang et al. (2018) found that features that generated good classification scores always generated rough bounding boxes. Value-based Deep Q-Networks (Song et al. 2023) or policy-based Proximal Policy Optimization (PPO) models (Yi, Qu, and Jiao 2023) may not well address the challenges of simultaneously maximizing conflicting rewards (Bai, Cheng, and Jin 2023). Evolutionary strategies have been designed to handle conflicting rewards in multi-objective scheduling (Tu et al. 2023; Chen et al. 2022). In this paper, an evolutionary strategy is integrated with a PPO strategy, where the PPO strategy effectively makes use of the appearance features to determine a suitable scaling action under the guidance of the three rewards, and the evolutionary strategy makes use of the past experience embedded in the agent to refine the scaling action.

The PPO agent consists of an actor model to choose a proper action and a critic model to evaluate the action. Specifically, the actor model takes the spatial-semantic attended features as the input, estimates the probability distribution of feasible actions by using a squeeze-and-excitation network (Hu, Shen, and Sun 2018), and determines an appropriate scaling action for each cluster region. An action is sampled using the policy $\pi_\vartheta$, $a^t \sim \pi_\vartheta(a^t|s^t)$, and the advantage function is calculated to evaluate the action as $\mathcal{A}(s^t, a^t) = \mathcal{R}(s^t, a^t) + \gamma\mathcal{V}_\varphi(s^{t+1}) - \mathcal{V}_\varphi(s^t)$, where $\gamma$ is the discount factor and $\mathcal{V}_\varphi(\cdot)$ is the state value in a specific state estimated by the critic model. The parameters $\vartheta$ of the actor model are updated through the gradient descent as,

$$\vartheta \leftarrow \vartheta + \eta_\vartheta \nabla_\vartheta \log \pi_\vartheta(a^t|s^t)\mathcal{A}(s^t, a^t), \qquad (4)$$

where $\eta_\vartheta$ is the learning rate. The actor model performs an efficient exploration to avoid a local optimum. The critic model employs a network architecture analogous to the actor network, which takes the observations from the current state as the input and approximates the state-value function. Following the design in (Araslanov, Rothkopf, and Roth 2019), the critic loss is defined as the squared error loss of estimated state-value and discounted sum of rewards in the trajectory. The critic model is updated with a learning rate of $\eta_\varphi$ as,

$$\varphi \leftarrow \varphi - \eta_\varphi \nabla_\varphi (\mathcal{V}_\varphi(s^t) - \sum_{i=t}^{T} \gamma^{i-t}\mathcal{R}^i)^2. \qquad (5)$$

The proposed evolutionary strategy is designed to better explore and exploit the feasible action space. Specifically, denote $\boldsymbol{\lambda} = \{\lambda_i\}_{i=1}^N$ as the set of scaling factors for $N$ cluster regions. The scaling actions $\boldsymbol{\lambda}^t$ given by the actor model at Step $t$, along with the $W-1$ best solutions $\boldsymbol{H}^{W-1}$ from the history actions $\mathbb{H}$ form the initial population of size $W$. $\mathbb{H}$ contains effective solutions dominated by different rewards in different scenarios. By applying evolution operators such as crossover and mutation, the newly generated $W$ offspring could explore and exploit solutions in multiple scenarios, and balance the importance of different rewards. Specifically, the crossover of scaling factors combines historical solutions in different scenarios from more than one parent, and the mutation of scaling factors allows broader trials and escape from possible local optimums. Among $W$ parents and $W$ generated offspring, the new population is formed by $W$ individuals with the largest scale-consistency reward $r_s$, as defined in Eqn. (2). The evolution stops if $r_s \geq \delta$, where $\delta$ is a predefined threshold. The best solution after evolution is applied to scale the cluster regions, and simultaneously stored into $\mathbb{H}$. Objects are detected on the scaled regions, and the rewards are calculated to evaluate the scaling actions and update the EVORL network as in (Araslanov, Rothkopf, and Roth 2019).

The proposed evolutionary reinforcement learning for determining the optimal scales is summarized in Algo. 1.

## Experimental Results

### Experimental Settings

**Datasets** The proposed model is compared with state-of-the-art models on two benchmark drone imagery datasets.

---

**Algorithm 1: Training procedures for the proposed EVORL**

**Input**: The number of episodes $P$, the number of steps $T$, the number of evolution iterations $I$, the population size $W$
**Output**: Policy net $\pi$

1: **for** $p \leftarrow 1$ to $P$ **do**
2:     Sample a batch of $M$ images.
3:     **for** $t \leftarrow 1$ to $T$ **do**
4:         Derive the appearance features $\mathbf{X}$ from images for $N$ cluster regions as $\mathbf{X} = \mathcal{F}_P(\mathbb{C}; \boldsymbol{\theta})$.
5:         Extract the spatial-semantic features as in Eqn. (1).
6:         Obtain the scaling actions $\boldsymbol{\lambda}^t$ by using the actor.
7:         Combine $\boldsymbol{\lambda}^t$ with $\boldsymbol{H}^{W-1}$ as the initial population.
8:         **for** $i \leftarrow 1$ to $I$ **do**
9:             Yield $W$ offspring by crossover and mutation.
10:            Evaluate each offspring and parents by Eqn. (2), **break** if $r_s \geq \delta$.
11:            Choose best $W$ individuals as new population.
12:         **end for**
13:         Select the best $\boldsymbol{\lambda}^t$ from population and add to $\mathbb{H}$.
14:         Update the state using the scaling factors $\boldsymbol{\lambda}^t$.
15:         Derive the reward as $\mathcal{R}^t = \alpha_l r_l + \alpha_c r_c + \alpha_s r_s$.
16:         Estimate the state-value $\mathcal{V}_\varphi(s^t)$.
17:         Evaluate the advantage function $\mathcal{A}(s^t, a^t)$.
18:         Update the actor model by using Eqn. (4).
19:         Update the critic model by using Eqn. (5).
20:     **end for**
21: **end for**

---

**UAVDT** dataset (Du et al. 2018) is a drone imagery dataset for object detection, single-object tracking and multi-object tracking. It contains 24,143 and 16,592 images for training and testing, respectively, with an average resolution of $1,024 \times 540$ pixels. This dataset captures images in complex scenarios and is commonly utilized for detecting vehicles like cars, trucks, and buses.

**VisDrone** dataset (Zhu et al. 2022) is a large-scale benchmark collected by drone-mounted cameras, encompassing 10,209 aerial images of 10 different categories, with a size of $2,000 \times 1,500$ pixels. The dataset is officially split into training, testing and validation sets, with 6,471, 3,190 and 548 images, respectively. As ground-truth annotations of the testing set are unavailable, following (Liu et al. 2021; Ge et al. 2022), the validation set is used for evaluation.

**Compared Methods** The proposed method is compared against nine state-of-the-art models. The results of compared methods are taken directly from the original papers. Faster R-CNN (Ren et al. 2017) serves as a baseline method. HRD-Net adapts general object detectors on natural images for detecting small objects in aerial images (Liu et al. 2021). DMNet (Li et al. 2020a) adapts the Multi-Column CNN for crowd counting to estimate an object density map and crops patches for fine detection. Other models are grouped based on the way of scaling patches in the coarse-to-fine pipeline. **Resized to a fixed size:** SAIC-FPN utilizes super-resolution techniques to up-sample the input image and performs fine detection on cropped patches (Zhou et al. 2019). GLSAN (Deng et al. 2021) roughly detects patches first, and

then resizes these patches to a fixed size by super-resolution methods. AdaZoom (Xu, Li, and Wang 2022) leverages a reinforcement learning framework to determine the focused regions, and resizes them to a certain scale for fine detection. **Resized with one or a few scaling factors:** ClusDet (Yang et al. 2019) utilizes two sub-networks, one for cropping regions of dense objects and the other for adjusting the scales of crops for fine detection. UFPMP-Det (Huang, Chen, and Huang 2022) and Zoom&Reasoning Det (Ge et al. 2022) both utilize the detector with Generalized Focal Loss (Li et al. 2020b) for coarse detection. The former determines the patch scale by measuring the average object size inside the patch, and the latter incorporates a Foreground Zoom strategy to determine the patch scales.

**Implementation Details** The stochastic gradient descent strategy is employed with a weight decay rate of 0.0005, a momentum rate of 0.9, and a dropout rate of 0.5. A cosine learning rate scheduler is used with an initial learning rate of 0.01. The same $\beta = 1.5$ is used as in (Huang, Chen, and Huang 2022). For the EVORL agent, the weighting factors $\alpha_l$, $\alpha_c$ and $\alpha_s$ are set to 1, the threshold $\delta = 0.5$, the size of the population $W = 32$, the number of evolution iterations $I = 10$, the number of steps $T = 50$ for one mini-batch, and the number of episodes $P = 1000$.

## Comparison Results on UAVDT

The proposed method is compared to seven state-of-the-art methods on the UAVDT dataset, using the evaluation metrics $AP$, $AP_{50}$ and $AP_{75}$ as in (Huang, Chen, and Huang 2022; Xu, Li, and Wang 2022). As shown in Table 1, the proposed model significantly outperforms all previous solutions, specifically surpassing UFPMP-Det (Huang, Chen, and Huang 2022), the previous best performing method, by 3.4%, 5.1% and 3.5% in terms of $AP$, $AP_{50}$ and $AP_{75}$, respectively. UFPMP-Det utilizes the average object size for scaling factor selection (Huang, Chen, and Huang 2022), which struggles with large scale variations. In contrast, the proposed EVORL makes use of both the current image patch and the past experience embedded in the agent to make informed decisions, adaptively determining the optimal scale for each patch. The spatial-semantic attention mechanism exploits supportive cues between objects to enhance patch features. Moreover, the Localization Reward and Labeling Reward provide supervision signals to directly maximize detection accuracy and the Scale-consistency Reward regularizes the agent to derive a more robust solution, leading to significant performance improvements.

To further analyze the performance across objects of different sizes, $AP^S$, $AP^M$ and $AP^L$, the average precision for objects with an area smaller than $32^2$ pixels, between $32^2$ and $96^2$ pixels, and larger than $96^2$ pixels, respectively on the UAVDT dataset, are summarized in Table 2. Some methods in Table 1 did not report results for objects of different sizes. As shown in Table 2, the proposed method consistently outperforms all the compared models across three sizes, demonstrating its capability of detecting objects of various scales. Specifically, compared to Zoom&Reasoning Det (Ge et al. 2022), the performance gain is 6.5%, 7.7%

| Method | $AP$ | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|
| Faster R-CNN (TPAMI, 2017) | 12.1 | 23.5 | 10.8 |
| ClusDet (ICCV, 2019) | 13.7 | 26.5 | 12.5 |
| DMNet (CVPR Workshop, 2020) | 14.7 | 24.6 | 16.3 |
| GLSAN (TIP, 2021) | 17.0 | 28.1 | 18.8 |
| AdaZoom (TMM, 2022) | 20.1 | 34.5 | 21.5 |
| Zoom&Reasoning Det (SPL, 2022) | 21.8 | 34.9 | 24.8 |
| UFPMP-Det (AAAI, 2022) | 24.6 | 38.7 | 28.0 |
| Proposed Method | **28.0** | **43.8** | **31.5** |

Table 1: Comparison with the state-of-the-art methods on the UAVDT dataset. The proposed method consistently and significantly outperforms all the compared methods.

| Method | $AP^S$ | $AP^M$ | $AP^L$ |
|---|---|---|---|
| Faster R-CNN (TPAMI, 2017) | 8.4 | 21.5 | 14.7 |
| ClusDet (ICCV, 2019) | 9.1 | 25.1 | 31.2 |
| DMNet (CVPR Workshop, 2020) | 9.3 | 26.2 | 35.2 |
| AdaZoom (TMM, 2022) | 14.2 | 29.2 | 28.4 |
| Zoom&Reasoning Det (SPL, 2022) | 15.3 | 32.7 | 30.8 |
| Proposed Method | **21.8** | **40.4** | **35.9** |

Table 2: Comparison with state-of-the-art methods on the UAVDT dataset in terms of $AP^S$, $AP^M$ and $AP^L$.

| Method | $AP$ | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|
| Faster R-CNN (TPAMI, 2017) | 21.8 | 41.8 | 20.1 |
| SAIC-FPN (Neurocomputing, 2019) | 35.7 | 62.3 | 35.1 |
| ClusDet (ICCV, 2019) | 32.4 | 56.2 | 31.6 |
| DMNet (CVPR Workshop, 2020) | 29.4 | 49.3 | 30.6 |
| GLSAN (TIP, 2021) | 32.5 | 55.8 | 33.0 |
| HRDNet (ICME, 2021) | 35.5 | 62.0 | 35.1 |
| Zoom&Reasoning Det (SPL, 2022) | 39.0 | 66.5 | 39.7 |
| UFPMP-Det (AAAI, 2022) | 39.2 | 65.3 | 40.2 |
| UFPMP-Det+MS (AAAI, 2022) | 40.1 | 66.8 | 41.3 |
| AdaZoom (TMM, 2022) | 40.3 | **66.9** | 41.8 |
| Proposed Method | **42.2** | 66.0 | **44.5** |

Table 3: Comparisons with state-of-the-art methods on the VisDrone dataset. The proposed method significantly outperforms the compared methods in terms of $AP$ and $AP_{75}$.

and 5.1% for small, median and large objects, respectively.

## Comparison Results on VisDrone

The comparison results with nine state-of-the-art methods on the VisDrone dataset (Zhu et al. 2022) are summarized in Table 3. Key observations are summarized as follows: 1) The proposed model significantly outperforms all compared models in terms of the key evaluation metric $AP$. Specifically, it achieves an $AP$ of 42.2%, making an improvement of 1.9% over the previous best model, AdaZoom (Xu, Li, and Wang 2022). AdaZoom resizes the patches to a fixed scale, while the proposed method utilizes the current image patch, the spatial-semantic attention, the scale consistency, and the past experience embedded in the agent to adaptively determine the most appropriate scale for each patch, leading to better detection performance. 2) The most significant performance gain is observed in $AP_{75}$, with

Figure 3: Visual comparison with UFPMP-Det (Huang, Chen, and Huang 2022) on the VisDrone dataset. The proposed method correctly detects more objects than UFPMP-Det, as annotated in green.

| YOLOX | PPO | SSA | EVO | $AP$ | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|---|---|---|
| √ | | | | 37.5 | 59.3 | 39.3 |
| √ | √ | | | 39.1 | 61.9 | 41.1 |
| √ | √ | √ | | 40.7 | 64.0 | 43.0 |
| √ | √ | √ | √ | **42.2** | **66.0** | **44.5** |

Table 4: Ablation study of major components of the proposed method on the VisDrone dataset.

a 2.7% improvement over AdaZoom, thanks to the Localization Reward that enhances object localization. 3) The proposed method yields a slightly lower $AP_{50}$ than AdaZoom, because many ultra-small objects in the VisDrone dataset only contain a few pixels, while YOLOX faces challenges in detecting these objects during coarse detection (Wang et al. 2023). 4) Note that the previous best methods on the two datasets are different. Compared to the previous best method on the VisDrone dataset, AdaZoom, the proposed method achieves significant performance gains of 7.9%, 9.3%, and 10.0% in terms of $AP$, $AP_{50}$ and $AP_{75}$, respectively, on the UAVDT dataset.

## Ablation Study of Major Components

The ablation results for the proposed method on the VisDrone dataset (Zhu et al. 2022) are summarized in Table 4. 1) Compared to the YOLOX baseline, by introducing the PPO agent to determine the optimal scales based on the appearance feature extracted directly from the Patch Encoder, the $AP$, $AP_{50}$ and $AP_{75}$ are improved by 1.6%, 2.6% and 1.8%, respectively. 2) By adding the spatial-semantic attention (SSA) into the visual perception module, the $AP$, $AP_{50}$ and $AP_{75}$ are further improved by 1.6%, 2.1% and 1.9%, respectively. 3) By incorporating the evolutionary strategy into the PPO agent, the $AP$, $AP_{50}$ and $AP_{75}$ are further boosted by 1.5%, 2.0% and 1.5%, respectively. The proposed EVORL makes good use of the past experience to refine the derived scaling factors, so that it mitigates the outlier scaling factors. These ablation results show the effectiveness of the major components in the proposed method.

## Visualization of Detection Results

The proposed method is visually compared to UFPMP-Det (Huang, Chen, and Huang 2022) that yields the previous best results averaged across the two datasets. As shown in Fig. 3, the proposed model better recognizes small objects that are easily neglected, *e.g.*, the 'car' and 'person' objects at the lower left corner of the focused regions in the first two columns, and the 'tricycle' objects in the third column. The ultra-small feature map encodes more low-level but high-resolution features, partially reducing the loss of details during feature pooling. Notably, UFPMP-Det selects one of three predefined scaling factors based on the average object size in a patch, while the proposed method adaptively determines the optimal scale of each patch by utilizing both the current patch and the agent's past experience, and hence better detects small objects. Moreover, as seen from the last column of Fig. 3, UFPMP-Det wrongly classifies 'van' as 'car' whereas the proposed method can correctly classify them, thanks to the proposed scale-consistency reward and the spatial-semantic attention mechanism, which effectively utilizes supportive information from nearby objects to better distinguish challenging objects.

## Conclusion

To tackle the challenges of detecting small objects and handle the large scale variations in drone imagery, an evolutionary reinforcement learning framework has been proposed to determine the optimal scale for object detection. The designed agent combines the evolutionary strategy and the proximal policy optimization strategy to make good use of both the current patch status and the past experience embedded in the agent's population. The three designed rewards, considering the localization accuracy, the accuracy of predicted labels, and the scale consistency, address the issue of lacking ground-truth labels for optimal scales, and provide supervision signals for training the agent. Furthermore, a spatial-semantic attention has been designed to capture the mutual supportive information from nearby objects. The proposed method has been compared with nine state-of-the-art approaches on two benchmark datasets, UAVDT and VisDrone. It significantly outperforms the compared solutions.

## Acknowledgements

## References

Alami, R.; Hacid, H.; Bellone, L.; Barcis, M.; and Natalizio, E. 2023. SOREO: A System for Safe and Autonomous Drones Fleet Navigation with Reinforcement Learning. In *Proc. AAAI Conf. Artif. Intell.*, volume 37, 16398–16400.

Araslanov, N.; Rothkopf, C. A.; and Roth, S. 2019. Actor-Critic Instance Segmentation. In *Proc. Comput. Vis. Pattern Recognit.*, 8237–8246.

Bai, H.; Cheng, R.; and Jin, Y. 2023. Evolutionary Reinforcement Learning: A Survey. *Intell. Comput.*, 2: 0025.

Bai, Y.; Zhang, Y.; Ding, M.; and Ghanem, B. 2018. SOD-MTGAN: Small Object Detection via Multi-Task Generative Adversarial Network. In *Proc. Eur. Conf. Comput. Vis.*, 210–226.

Bouguettaya, A.; Zarzour, H.; Kechida, A.; and Taberkit, A. M. 2022. Vehicle Detection From UAV Imagery With Deep Learning: A Review. *IEEE Trans. Neural Netw. Learn. Syst.*, 33(11): 6047–6067.

Cai, Z.; and Vasconcelos, N. 2018. Cascade R-CNN: Delving Into High Quality Object Detection. In *Proc. Comput. Vis. Pattern Recognit.*, 6154–6162.

Chen, X.; Bai, R.; Qu, R.; and Dong, H. 2022. Cooperative Double-Layer Genetic Programming Hyper-Heuristic for Online Container Terminal Truck Dispatching. *IEEE Trans. Evol. Comput.* (Early access).

Cheng, G.; Han, J.; Zhou, P.; and Xu, D. 2019. Learning Rotation-Invariant and Fisher Discriminative Convolutional Neural Networks for Object Detection. *IEEE Trans. Image Process.*, 28(1): 265–278.

Dalal, N.; and Triggs, B. 2005. Histograms of Oriented Gradients for Human Detection. In *Proc. Comput. Vis. Pattern Recognit.*, volume 1, 886–893.

Deng, S.; Li, S.; Xie, K.; Song, W.; Liao, X.; Hao, A.; and Qin, H. 2021. A Global-Local Self-Adaptive Network for Drone-View Object Detection. *IEEE Trans. Image Process.*, 30: 1556–1569.

Du, D.; Qi, Y.; Yu, H.; Yang, Y.; Duan, K.; Li, G.; Zhang, W.; Huang, Q.; and Tian, Q. 2018. The Unmanned Aerial Vehicle Benchmark: Object Detection and Tracking. In *Proc. Eur. Conf. Comput. Vis.*, 375–391.

Ge, Z.; Liu, S.; Wang, F.; Li, Z.; and Sun, J. 2021. YOLOX: Exceeding Yolo Series in 2021. *arXiv preprint arXiv:2107.08430*.

Ge, Z.; Qi, L.; Wang, Y.; and Sun, Y. 2022. Zoom-and-Reasoning: Joint Foreground Zoom and Visual-Semantic Reasoning Detection Network for Aerial Images. *IEEE Signal Process. Lett.*, 29: 2572–2576.

Girshick, R.; Donahue, J.; Darrell, T.; and Malik, J. 2014. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. In *Proc. Comput. Vis. Pattern Recognit.*, 580–587.

He, K.; Gkioxari, G.; Dollar, P.; and Girshick, R. 2017. Mask R-CNN. In *Proc. IEEE Int. Conf. Comput. Vis.*, 2961–2969.

He, W.; Zhang, J.; Ren, J.; Bai, R.; and Jiang, X. 2023. Hierarchical ConViT with Attention-based Relational Reasoner for Visual Analogical Reasoning. In *Proc. AAAI Conf. Artif. Intell.*, volume 37, 22–30.

Hu, J.; Shen, L.; and Sun, G. 2018. Squeeze-and-Excitation Networks. In *Proc. Comput. Vis. Pattern Recognit.*, 7132–7141.

Huang, Y.; Chen, J.; and Huang, D. 2022. UFPMP-Det:Toward Accurate and Efficient Object Detection on Drone Imagery. In *Proc. AAAI Conf. Artif. Intell.*, 1, 1026–1033.

Jiang, B.; Luo, R.; Mao, J.; Xiao, T.; and Jiang, Y. 2018. Acquisition of Localization Confidence for Accurate Object Detection. In *Proc. Eur. Conf. Comput. Vis.*, 784–799.

Li, C.; Yang, T.; Zhu, S.; Chen, C.; and Guan, S. 2020a. Density Map Guided Object Detection in Aerial Images. In *Proc. Comput. Vis. Pattern Recognit. Workshops*, 737–746.

Li, J.; Liang, X.; Wei, Y.; Xu, T.; Feng, J.; and Yan, S. 2017. Perceptual Generative Adversarial Networks for Small Object Detection. In *Proc. Comput. Vis. Pattern Recognit.*, 1222–1230.

Li, X.; Wang, W.; Wu, L.; Chen, S.; Hu, X.; Li, J.; Tang, J.; and Yang, J. 2020b. Generalized Focal Loss: Learning Qualified and Distributed Bounding Boxes for Dense Object Detection. In *Proc. Adv. Neural Inf. Process. Syst.*, volume 33, 21002–21012.

Lin, T.-Y.; Goyal, P.; Girshick, R.; He, K.; and Dollár, P. 2020. Focal Loss for Dense Object Detection. *IEEE Trans. Pattern Anal. Mach. Intell.*, 42(2): 318–327.

Liu, Z.; Gao, G.; Sun, L.; and Fang, Z. 2021. HRDNet: High-Resolution Detection Network for Small Objects. In *Proc. IEEE Int. Conf. Multimedia Expo*, 1–6.

Lowe, D. G. 2004. Distinctive Image Features from Scale-Invariant Keypoints. *Int. J. Comput. Vis.*, 60: 91–110.

Ozge Unel, F.; Ozkalayci, B. O.; and Cigla, C. 2019. The Power of Tiling for Small Object Detection. In *Proc. Comput. Vis. Pattern Recognit. Workshops*.

Redmon, J.; Divvala, S.; Girshick, R.; and Farhadi, A. 2016. You Only Look Once: Unified, Real-Time Object Detection. In *Proc. IEEE Conf. Comput. Vis Pattern Recognit.*, 779–788.

Redmon, J.; and Farhadi, A. 2017. YOLO9000: Better, Faster, Stronger. In *Proc. Comput. Vis. Pattern Recognit.*, 7263–7271.

Ren, J.; and Jiang, X. 2017. Regularized 2-D Complex-log Spectral Analysis and Subspace Reliability Analysis of micro-Doppler Signature for UAV Detection. *Pattern Recognit.*, 69: 225–237.

Ren, J.; and Jiang, X. 2021. A Three-Step Classification Framework to Handle Complex Data Distribution for Radar UAV Detection. *Pattern Recognit.*, 111: 107709.

Ren, J.; Jiang, X.; and Yuan, J. 2015. A Chi-Squared-Transformed Subspace of LBP Histogram for Visual Recognition. *IEEE Trans. Image Process.*, 24(6): 1893–1904.

Ren, S.; He, K.; Girshick, R.; and Sun, J. 2017. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(6): 1137–1149.

Sommer, L. W.; Schuchert, T.; and Beyerer, J. 2017. Fast Deep Vehicle Detection in Aerial Images. In *Proc. IEEE Winter Conf. Appl. Comput. Vis.*, 311–319.

Song, X.; Jin, J.; Yao, C.; Wang, S.; Ren, J.; and Bai, R. 2023. Siamese-Discriminant Deep Reinforcement Learning for Solving Jigsaw Puzzles with Large Eroded Gaps. In *Proc. AAAI Conf. Artif. Intell.*, volume 37, 2303–2311.

Tan, M.; Pang, R.; and Le, Q. V. 2020. EfficientDet: Scalable and Efficient Object Detection. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 10781–10790.

Tokekar, P.; Hook, J. V.; Mulla, D.; and Isler, V. 2016. Sensor Planning for a Symbiotic UAV and UGV System for Precision Agriculture. *IEEE Trans. Robot.*, 32(6): 1498–1511.

Tu, C.; Bai, R.; Aickelin, U.; Zhang, Y.; and Du, H. 2023. A Deep Reinforcement Learning Hyper-heuristic with Feature Fusion for Online Packing Problems. *Expert Syst. Appl.*, 120568.

Wang, C.-Y.; Bochkovskiy, A.; and Liao, H.-Y. M. 2021. Scaled-YOLOv4: Scaling Cross Stage Partial Network. In *Proc. Comput. Vis. Pattern Recognit.*, 13029–13038.

Wang, H.; Wang, Z.; Jia, M.; Li, A.; Feng, T.; Zhang, W.; and Jiao, L. 2019. Spatial Attention for Multi-Scale Feature Refinement for Object Detection. In *Proc. IEEE Int. Conf. Comput. Vis. Workshops*.

Wang, X.; He, N.; Hong, C.; Wang, Q.; and Chen, M. 2023. Improved YOLOX-X based UAV aerial photography object detection algorithm. *Image Vis. Comput.*, 135: 104697.

Xi, Y.; Jia, W.; Zheng, J.; Fan, X.; Xie, Y.; Ren, J.; and He, X. 2021. DRL-GAN: Dual-Stream Representation Learning GAN for Low-Resolution Image Classification in UAV Applications. *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, 14: 1705–1716.

Xi, Y.; Zheng, J.; He, X.; Jia, W.; Li, H.; Xie, Y.; Feng, M.; and Li, X. 2020. Beyond context: Exploring semantic similarity for small object detection in crowded scenes. *Pattern Recognit. Lett.*, 137: 53–60.

Xu, J.; Li, Y.; and Wang, S. 2022. AdaZoom: Towards Scale-Aware Large Scene Object Detection. *IEEE Trans. Multimedia*. (Early access).

Yang, F.; Fan, H.; Chu, P.; Blasch, E.; and Ling, H. 2019. Clustered Object Detection in Aerial Images. In *Proc. IEEE Int. Conf. Comput. Vis.*, 8311–8320.

Yi, W.; Qu, R.; and Jiao, L. 2023. Automated Algorithm Design using Proximal Policy Optimisation with Identified Features. *Expert Syst. Appl.*, 216: 119461.

Yun, W. J.; Park, S.; Kim, J.; Shin, M.; Jung, S.; Mohaisen, D. A.; and Kim, J.-H. 2022. Cooperative Multiagent Deep Reinforcement Learning for Reliable Surveillance via Autonomous Multi-UAV Control. *IEEE Trans. Ind. Informat.*, 18(10): 7086–7096.

Zhang, J.; Ren, J.; Zhang, Q.; Liu, J.; and Jiang, X. 2023. Spatial Context-Aware Object-Attentional Network for Multi-Label Image Classification. *IEEE Trans. Image Process.*, 32: 3000–3012.

Zhang, J.; Zhang, Q.; Ren, J.; Zhao, Y.; and Liu, J. 2022. Spatial-Context-Aware Deep Neural Network for Multi-Class Image Classification. In *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, 1960–1964.

Zhang, S.; He, G.; Chen, H.-B.; Jing, N.; and Wang, Q. 2019. Scale Adaptive Proposal Network for Object Detection in Remote Sensing Images. *IEEE Geosci. Remote Sens. Lett.*, 16(6): 864–868.

Zhou, J.; Vong, C.-M.; Liu, Q.; and Wang, Z. 2019. Scale adaptive image cropping for UAV object detection. *Neurocomputing*, 366: 305–313.

Zhu, P.; Wen, L.; Du, D.; Bian, X.; Fan, H.; Hu, Q.; and Ling, H. 2022. Detection and Tracking Meet Drones Challenge. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(11): 7380–7399.