

# APEX: A Decoupled Memory-based Explorer for Asynchronous Aerial Object Goal Navigation

## Supplementary Material

### 6. Algorithm Flow

Algorithm 1 and Algorithm 2 illustrates the overall workflow of our agent. We employ a multi-process architecture to enable parallel execution of the three modules.

---

**Algorithm 1** Overall Workflow

---

**Input:** Observations  $O$

**Input:** States  $S$

**Input:** Target Description  $D$

- $MEM$  denotes the shared information, including  $M_{attr}$ ,  $M_{expl}$ ,  $M_{obst}$ , and others.

```
1: function MAP UPDATING( $MEM, O, S, D$ )
2:   while not IsDetected do
3:      $c, s \leftarrow CAP(O, D)$ 
4:      $\mathcal{M} \leftarrow SEG(O, c)$ 
5:      $MEM \leftarrow Update(MEM, \mathcal{M}, s, O, S)$ 
6:   end while
7:   process.join()
8: end function

9: function ACTION DECISION( $MEM, S, O, p$ )
10:  while True do
11:    if IsDetected then
12:      Navigate( $p$ )
13:    else
14:      break
15:    end if
16:     $MEM \leftarrow Update(MEM, O, S)$ 
17:     $a \leftarrow Policy(MEM, S)$ 
18:    Execute( $a$ )
19:  end while
20:  process.join()
21: end function

22: function TARGET GROUNDING
23:  while not IsDetected do
24:     $b \leftarrow GD(O, D)$ 
25:     $p \leftarrow RP(b, O, S)$ 
26:  end while
27:  process.join()
28: end function
```

---

### 7. More Visualization Results

#### 7.1. Cases

We provide additional visualizations to further illustrate APEX’s performance. Fig. 8 shows four navigation episodes: three successful searches for ”Car”, ”Playground”, and ”Umbrella”, and one failure case targeting a ”Sandbox”. The failure was caused by a collision with a tree, which suggests a potential limitation in our Obstacle Map’s ability to detect sparse obstacles.

#### 7.2. Trajectories

Fig. 6 demonstrates the complete trajectories for these four cases (labeled as Case 1-4), in addition to the main paper’s example (Case 0).

### 8. Details of Heuristic Algorithm

The heuristic algorithm mentioned in Ablation Study is detailed in Algorithm 3. The core idea is to simulate each possible action, evaluate the potential reward based on attraction and exploration criteria, and select the action with the highest predicted reward.

### 9. More Training Details

#### 9.1. Offline Map Generation

Directly incorporating the VLM for attraction map generation within the reinforcement learning loop is infeasible due to its high inference latency. To address this, we pre-generate a complete attraction map for each of the 36 tasks before training begins, which are shown in Fig. 7.

For each task, we designed a flight path to cover the entire relevant search area. An agent was simulated along this path to generate and store a full attraction map. During RL training, instead of calling the VLM, the agent identifies its currently visible grid cells via the depth map and populates its local attraction map by looking up the pre-computed values. This method bypasses the VLM’s latency, making the RL training process computationally available.

#### 9.2. Policy Network Architecture

Fig. 5 illustrates the architecture of the policy network in our RL-based Action Decision Module. It utilizes three CNNs to extract features from the Attraction, Exploration, and Obstacle maps, respectively. To process the 3D map data, the z-axis is treated as the channel dimension for the

---

**Algorithm 2** Memory Updating Workflow

---

**Input:** Masks  $\mathcal{M}$ **Input:** Observations  $O$ **Input:** States  $S$ **Input:** Attraction Scores  $s$ 

```
1: function ATTRUPDATE( $M_{\text{attr}}, s, \mathcal{M}, O, S$ )
2:    $\mathcal{G}_{\text{vis}}$ : the set of visible grid cells in the current view.
3:   for  $g \in \mathcal{G}_{\text{vis}}$  do
4:      $P_g$ : the set of pixels from  $O$  that project to  $g$ .
5:      $Mask^* \leftarrow \arg \max_{Mask_i \in \mathcal{M}} |P_g \cap Mask_i|$ 
6:     if  $Mask^*$  is found then
7:        $s^* \leftarrow$  score in  $s$  corresponding to  $Mask^*$ 
8:        $d_{\text{new}} \leftarrow \text{Dis}(S.\text{pos}, g.\text{center})$ 
9:       if  $d_{\text{new}} < M_{\text{attr}}[g].\text{dis}$  then
10:         $M_{\text{attr}}[g].\text{score} \leftarrow s^*$ 
11:         $M_{\text{attr}}[g].\text{dis} \leftarrow d_{\text{new}}$ 
12:       end if
13:     end if
14:   end for
15:   return  $M_{\text{attr}}$ 
16: end function
```

```
17: function EXPLUPDATE( $M_{\text{expl}}, O, S$ )
18:    $\mathcal{P}_{\text{end}}$ : the set of 3D world points from  $O$ .
19:    $\mathcal{G}_{\text{traversed}} \leftarrow \text{TraceRay}(S.\text{pos}, \mathcal{P}_{\text{end}})$ 
20:   for  $g \in \mathcal{G}_{\text{traversed}}$  do
21:      $d \leftarrow \text{distance}(S.\text{pos}, g.\text{center})$ 
22:      $v_{\text{current}} \leftarrow M_{\text{expl}}[g]$ 
23:      $v_{\text{new}} \leftarrow \text{ExplRate}(d, v_{\text{current}})$ 
24:      $M_{\text{expl}}[g] \leftarrow v_{\text{current}} + v_{\text{new}}$ 
25:   end for
26:   return  $M_{\text{expl}}$ 
27: end function
```

```
28: function OBSTUPDATE( $M_{\text{obst}}, O, S$ )
29:    $\mathcal{G}_{\text{depth}}$ : the set of grid cells hit by depth rays from  $O$ .
30:    $M_{\text{obst}}[g] \leftarrow 1, \forall g \in \mathcal{G}_{\text{depth}}$ 
31:   return  $M_{\text{obst}}$ 
32: end function
```

```
33: function TRACERAY( $S_{\text{start}}, \mathcal{P}_{\text{end}}$ )
34:    $\mathcal{G}_{\text{traversed}}$ : the set of grid cells
35:   for  $p_{\text{end}} \in \mathcal{P}_{\text{end}}$  do
36:      $\vec{r} = p_{\text{end}} - S_{\text{start}}$ 
37:     Sample points  $\mathcal{P}_s$  along  $\vec{r}$  at regular intervals.
38:      $\mathcal{G}_{\text{traversed}} \leftarrow \text{RP}(\mathcal{P}_s)$ 
39:   end for
40:   return  $\mathcal{G}_{\text{traversed}}$ 
41: end function
```

---

---

**Algorithm 3** Heuristic Action Decision

---

**Input:** Attraction Map  $M_{\text{attr}}$ **Input:** Exploration Map  $M_{\text{expl}}$ **Input:** Obstacle Map  $M_{\text{obst}}$ **Input:** Current State  $s$ 

- Let  $a$  be an action from the action set  $\mathcal{A}$ .
- Let  $s'$  be the hypothetical state after action  $a$ .
- The **attraction reward**  $R_{\text{attr}}$  and **exploration reward**  $R_{\text{expl}}$  is defined in main paper.

```
1: Initialize  $R$  to store the predicted reward.
2: for  $a \in \mathcal{A}$  do
3:    $s' \leftarrow \text{Simulate}(a, s, M_{\text{attr}}, M_{\text{expl}}, M_{\text{obst}})$ 
4:   if IsUnsafe( $s'$ ) then
5:      $R[a] \leftarrow -\infty$ 
6:     continue
7:   else
8:      $R[a] \leftarrow W_{\text{attr}}R_{\text{attr}}(s') + W_{\text{expl}}R_{\text{expl}}(s')$ 
9:   end if
10: end for
11: return  $\arg \max_{a \in \mathcal{A}} R[a]$ 
```

---

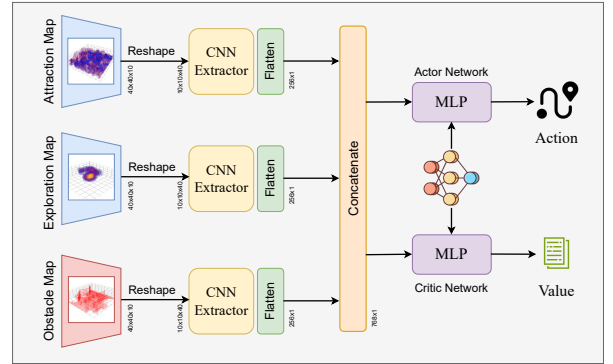


Figure 5. The architecture of policy network.

convolutional layers. The resulting feature vectors are concatenated and shared between an Actor (MLP) and a Critic (MLP) head. The Actor network outputs the final action policy. This implementation is based on the MultiInputPolicy architecture provided by the Stable-Baselines3 [29].

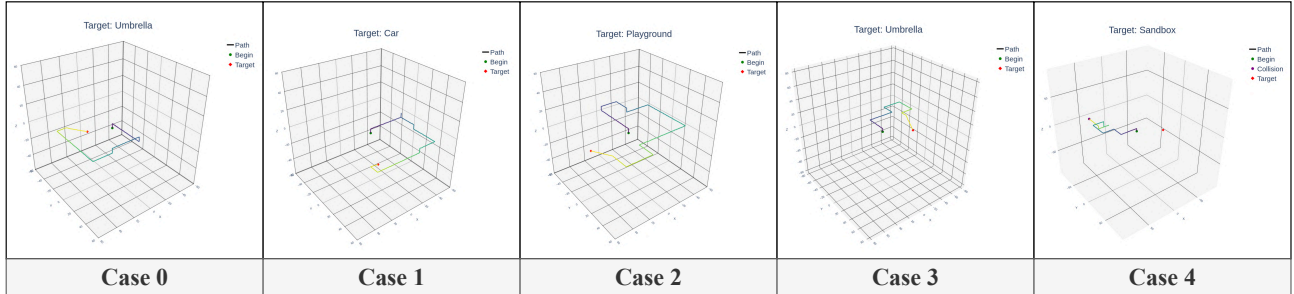


Figure 6. Complete trajectories of visualization examples.

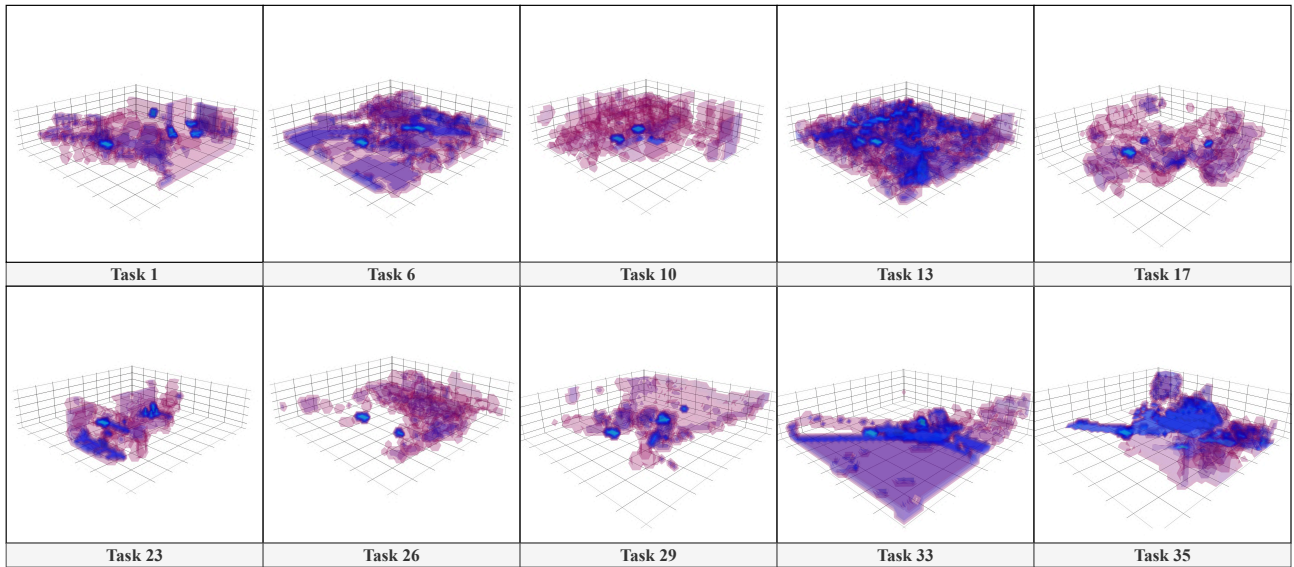


Figure 7. Examples of collected attraction maps.



Figure 8. More visualized results. The first three examples are successful cases, with the targets "Car", "Playground", and "Umbrella" respectively. The last one is a failure case, with the target being "Sandbox".