# Enhancing Exploratory Capability of Visual Navigation Using Uncertainty of Implicit Scene Representation

Yichen Wang, Qiming Liu, Zhe Liu, and Hesheng Wang*

*Abstract*— In the context of visual navigation in unknown scenes, both "exploration" and "exploitation" are equally crucial. Robots must first establish environmental cognition through exploration and then utilize the cognitive information to accomplish target searches. However, most existing methods for image-goal navigation prioritize target search over the generation of exploratory behavior. To address this, we propose the Navigation with Uncertainty-driven Exploration (NUE) pipeline, which uses an implicit and compact scene representation, NeRF, as a cognitive structure. We estimate the uncertainty of NeRF and augment the exploratory ability by the uncertainty to in turn facilitate the construction of implicit representation. Simultaneously, we extract memory information from NeRF to enhance the robot's reasoning ability for determining the location of the target. Ultimately, we seamlessly combine the two generated abilities to produce navigational actions. Our pipeline is end-to-end, with the environmental cognitive structure being constructed online. Extensive experimental results on image-goal navigation demonstrate the capability of our pipeline to enhance exploratory behaviors, while also enabling a natural transition from the exploration to exploitation phase. This enables our model to outperform existing memory-based cognitive navigation structures in terms of navigation performance. Project page: https://github.com/IRMVLab/NUE-NeRF-nav

## I. INTRODUCTION

When searching for a target in an unfamiliar environment, our subconscious instinctively prioritizes exploring to establish cognitive. Upon locating the target, we shift to exploitation, navigating towards it. When a robot engages in visual navigation tasks in an unknown environment, as shown in Fig. 1, it also benefits from both exploratory and exploitative thinking. However, existing cognitive navigation frameworks primarily focus on the robot's performance in the exploitation phase, neglecting the design of its exploratory behavior. We intend to use implicit scene representation as the memory structure of our navigation pipeline, specifically emphasizing the enhancement of the robot's exploratory capability. This enhancement enables the robot to rapidly establish environmental awareness, discover target-related cues, and transition into the exploitation phase.

Y. Wang and Q. Liu are with the Department of Automation, Shanghai Jiao Tong University, Shanghai 200240, China. Z. Liu is with MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, Shanghai 200240, China. H. Wang is with the Department of Automation, Shanghai Jiao Tong University, Shanghai 200240, China.

*Corresponding author: Hesheng Wang (e-mail: wanghesheng@ sjtu.edu.cn).
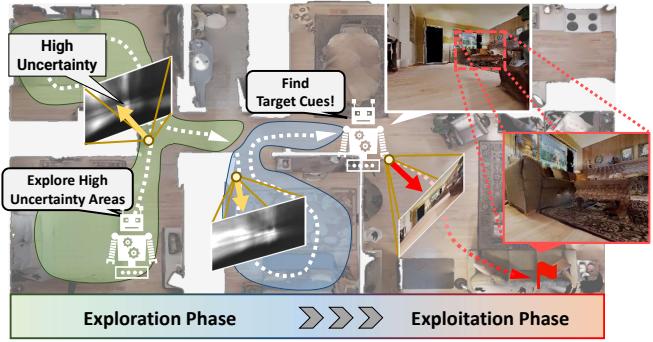
Fig. 1: **Navigation of robots in an unknown environment.** In visual navigation of robots in unknown environments, the process involves two phases: exploration and exploitation. Initially, the robot explores based on uncertainty to refine its cognitive structure, transitioning to navigation toward detected target-related cues in the environment.
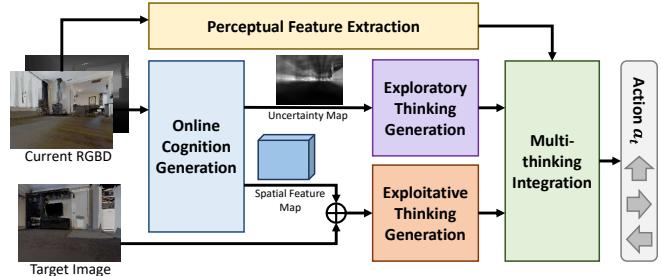


Fig. 2: **The overall architecture of NUE.** Firstly, real-time image input is used for online cognitive generation and perceptual feature extraction. Secondly, cognitive information is extracted to generate exploratory thinking and exploitative thinking. Eventually, multiple thinking is integrated, and navigational actions are generated.

In this paper, we adopt the Neural Radiance Fields (NeRF) [1] as our memory structure. In recent years, NeRF has shown excellent performance in reconstructing 3D scenes and synthesizing novel views [2]. Due to its compact implicit scene representation capability and network-based structure, NeRF has shown great potential in downstream navigation tasks that rely on long-term information representation. However, since our optimization of NeRF is performed online, its high signal-to-noise output early in navigation struggles to provide valuable environmental memory information. To overcome this challenge, we use an estimate of NeRF's uncertainty and leverage it to foster exploratory thinking within the robot, considering that exploration can speed up the establishment of implicit representations.

We propose the Navigation with Uncertainty-driven Exploration (NUE) pipeline, as illustrated in Fig. 2, which is end-to-end and fully differentiable. Firstly, we conduct online training on NeRF to generate cognition of the scenes. To enhance the exploration capability of the robot and reduce the noise of NeRF in the initial navigation phase, we include an estimate of NeRF's uncertainty. Secondly, we extract distinct components, namely uncertainty and spatial information, from NeRF. These components are subsequently compressed into feature representations. The uncertainty feature is leveraged to augment exploration capabilities, while the spatial feature is harnessed to optimize exploitation performance. Finally, we adaptively fuse the two features and output navigation actions through an action generator. Our main innovations can be summarized as follows:

- We propose NUE, an end-to-end visuomotor navigation pipeline integrating NeRF as a cognitive structure. By leveraging the compact scene representation capabilities of NeRF, we extend its application from the perception domain to the control domain.
- We utilize the estimation of NeRF's uncertainty to enable the robot to exhibit exploratory behavior. Additionally, our model successfully balances exploratory and exploitative thinking, achieving seamless integration between the exploration and exploitation stages.
- Experimental results demonstrate that NUE significantly improves navigation performance compared to existing cognitive memory structures. Interpretability experiments validate that NUE effectively generates and balances exploration and exploitation behaviors.

## II. RELATED WORKS

### A. Neural Radiance Fields in Perception

Neural Radiance Fields (NeRF) [1] is an innovative approach for synthesizing views. It combines multi-layer perceptrons (MLPs) with volume rendering techniques to generate new views. What makes NeRF impressive is its ability to achieve high-quality scene representation with a compact and concise structure. Currently, many efforts are focused on improving the training and inference efficiency of NeRF and addressing the issue of sparse view sensitivity.

Due to the computationally intensive of NeRF, some methods focus on improving the structure of NeRF [3] or exploring the use of additional depth supervision [4] to accelerate the inference process. In addition, the high-quality scene reconstruction in NeRF relies on a large number of densely sampled visual observations. However, obtaining densely sampled views can be challenging in many tasks. Therefore, some works [5] have focused on improving NeRF's sensitivity to sparse views. This is crucial for navigation tasks as well; since robots often have limited and unevenly distributed local observations during navigation.

In addition, several studies have applied NeRF to SLAM (simultaneous localization and mapping) systems. iNeRF [6] first models pose estimation as the inverse process of NeRF inference. Based on this, some studies have improved the

sampling method [7], further optimizing the performance and efficiency of pose estimation. While these studies have achieved impressive results in the perception domain, they have not yet extended the powerful characterization capabilities of NeRF to the control domain.

### B. Neural Radiance Fields in Robotics

In the field of robotics, NeRF is mainly applied to robotic arm control [8, 9] and navigation [10]–[13]. In the field of navigation, certain studies strive to establish the transformational relationship between NeRF and the geometric representation of the occupancy space. This enables precise estimation of the scene's geometric structure, and, by utilizing planning methods, facilitates the generation of collision-free and smooth trajectories [10, 11]. Although these works have made great progress in navigation, they still rely on using NeRF as an auxiliary tool for trajectory planning and have not fully integrated the implicit structure of NeRF with neural controllers.

Some recent studies [12, 13] have attempted to connect the implicit representation directly to neural controllers and output navigation actions. Additionally, researchers are focusing on the problem of active representation learning, which explores how to use navigation robots to better construct NeRF representations through exploration in the environment [14, 15]. While these studies have achieved remarkable results, they often lack a natural integration of the exploration and exploitation phases. Our focus is on enabling the robots to learn a multi-modal thinking approach, allowing for the generation of a navigation strategy that closely aligns with human thinking patterns.

## III. METHODOLOGY

### A. Overall System Architecture

We aim to enhance the spatial cognition capability of robots to achieve better exploratory behavior. To achieve this goal, we introduce NUE, as shown in Fig. 3, which consists of three key processes. The first part involves online cognition generation, where the robot stores real-time perceptual information of the environment in NeRF, providing the robot with spatial cognition. The second part is online cognitive extraction, which utilizes resnet and CBAM [16] to extract features from the uncertainty and spatial information produced by NeRF, and generate corresponding exploration and exploitation strategies. The third part encompasses multi-thinking integration, extracting perceptual features from real-time image input and then fusing perceptual features, uncertainty features, and spatial features to generate navigational actions through a neural controller.

### B. Online Cognition Generation

The spatial cognition of robots is generated online, independent of any prior knowledge of the scenes. However, due to the online generation nature, the model's cognitive awareness has significant uncertainty in the early stages of the navigation process, which is not conducive to direct
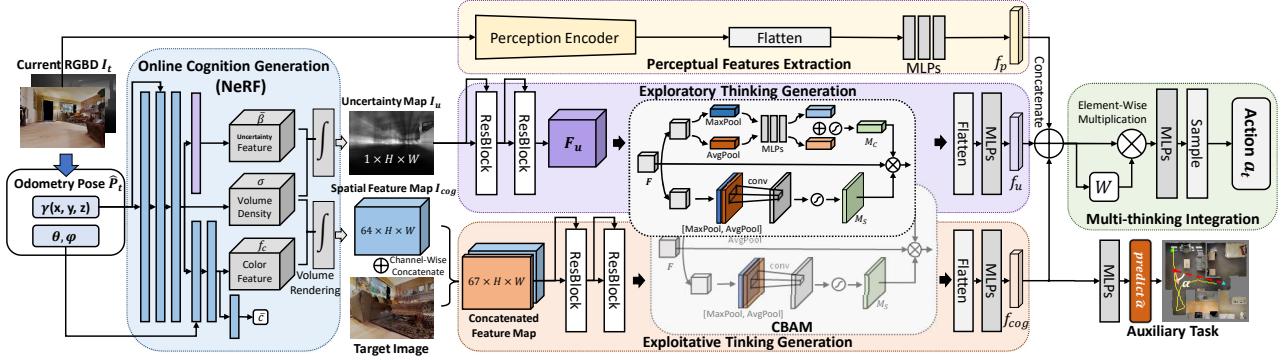
Fig. 3: **The network structure of NUE.** The overall framework first extracts real-time perceptual features and generates cognitive signals in NeRF. Subsequently, we compress the uncertainty map to generate uncertainty features and concatenate the spatial feature map with the target image in the channel dimension for spatial feature extraction. Finally, the features are concatenated and fed into an adaptive neural controller to generate the final navigation actions.

information utilization. Therefore, we enhance exploration by estimating the uncertainty of the cognitive structure.

We use NeRF to generate the spatial cognition of the robot. To incorporate uncertainty estimation into the model, we combine previous research [5] and model the emitted radiance values at each position in space as a Gaussian distribution, parameterized by mean $\bar{c}$ and variance $\bar{\beta}^2$. The formulation is as follows:

$$\left[\sigma, f, \bar{\beta}^2(\mathrm{r}(t))\right] = \mathrm{MLP}_{\theta_1, \theta_3}(\gamma_x(\mathrm{r}(t))), \quad (1)$$

$$\bar{c}(\mathrm{r}(t)) = \mathrm{MLP}_{\theta_2}(f, \gamma_d(\mathrm{r}(t))), \quad (2)$$

where $\gamma_{x,d}(\cdot)$ is the position encoding functions, $f$ represents the intermediate features, and $\mathrm{r}(t) = o + td$ represents the entire ray, with $o$ as the origin and $d$ as the direction of the ray from the camera center to the sampled point.

Due to the incorporation of uncertainty, we adopt the loss function used in ActiveNeRF [5] to optimize the NeRF network. The main loss function is as follows:

$$\mathcal{L}_i^u = \frac{\left\| C(\mathrm{r}_i) - \bar{C}(\mathrm{r}_i) \right\|_2^2}{2\bar{\beta}^2(\mathrm{r}_i)} + \frac{\log \bar{\beta}^2(\mathrm{r}_i)}{2}, \quad (3)$$

where $\mathrm{r}_i$ is sampled ray. $\bar{C}(\mathrm{r}_i)$ and $\bar{\beta}^2(\mathrm{r}_i)$ denote the mean and variance of rendering colors through sampled ray $\mathrm{r}_i$, respectively. $C(\mathrm{r}_i)$ denotes the ground truth.

*C. Online Cognition Extraction*

In this section, we aim to imbue the robot with an integration of both exploratory and exploitative cognition. This involves prioritizing unexplored areas before observing the target, collecting target-related cues, and showing a preference for the target once it is observed. We achieve this by extracting uncertainty features for exploration and spatial features from the cognitive structure for exploitation.

*1) Exploratory Thinking Generation:* To generate exploratory thinking, we extract the uncertainty features from NeRF. We render the uncertainty $\bar{\beta}^2(\mathrm{r}(t_i))$ outputted by NeRF as an uncertainty map $I_u \in \mathbb{R}^{1 \times W \times H}$, the rendering of uncertainty can be described as follows:

$$I_u(\mathrm{r}) = \sum_{i=1}^{N_s} \alpha_i \bar{\beta}^2(\mathrm{r}(t_i)), \quad (4)$$

$$\alpha_i = \exp(-\sum_{j=1}^{i-1} \sigma_j \delta_j) \left(1 - \exp\left(\sigma_i \delta_i\right)\right), \quad (5)$$

where $\delta_j = t_{i+1} - t_i$ represents the distance between neighboring sampling points, and $N_s$ is the number of sampling points on each ray.

Uncertainty map $I_u$ can reflect the familiarity of NeRF with different regions of the scene in the current perspective. As shown in Fig. 3, we first use two cascaded residual structures to extract features from the uncertainty map $I_u$. Through the resnet network, we compress the input uncertainty map $I_u$ into a feature map $F_u \in \mathbb{R}^{16 \times W/16 \times H/16}$.

After compressing $I_u$ to $F_u$, we employ channel attention and spatial attention to adaptively optimize the feature map.

**Channel Attention.** Channel attention compresses the input feature maps into $F_{avg}^c$ and $F_{max}^c$ by aggregating the spatial information using average pooling and maximum pooling in the spatial dimension. The process is as follows:

$$M_C(F) = \mathrm{Sigmod}\left(\mathrm{MLP}\left(F_{avg}^c\right) + \mathrm{MLP}\left(F_{max}^c\right)\right), \quad (6)$$

**Spatial Attention.** Spatial attention complements channel attention by exploiting average pooling and maximum pooling operations on the channel dimension to compress the feature graph into $F_{avg}^c$ and $F_{max}^c$. The specific process is as follows:

$$M_S(F) = \mathrm{Sigmod}\left(\mathrm{Conv}\left(\left[F_{avg}^s, F_{avg}^s\right]\right)\right), \quad (7)$$

We multiply the two attention maps $M_C$, $M_S$ with the original feature map $F_u$. Finally, the feature map is further compressed into a perceptual feature vector $f_u$ of length 64 through an MLP. This process can be represented as follows:

$$f_u = \mathrm{MLP}(\mathrm{Flatten}(M_C \otimes M_S \otimes F_u)), \quad (8)$$

where $\otimes$ denotes element-wise multiplication

*2) Exploitative Thinking Generation:* To generate exploitative thinking, we extract the spatial features from NeRF. We first generate the volume density $F_\sigma \in \mathbb{R}^{1 \times W \times H \times N_s}$ and the feature map $F_c \in \mathbb{R}^{64 \times W \times H \times N_s}$ from the intermediate layer through NeRF. After obtaining the voxel density and

TABLE I: **Image-goal navigation results.** The data in bold in each column represents the optimal data for that column. The bottommost row is the test results of our full model, NUE.

| | | | | Easy | | | Medium | | | Hard | | | Total | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | SR↑ | SPL↑ | DTS↓ | SR↑ | SPL↑ | DTS↓ | SR↑ | SPL↑ | DTS↓ | SR↑ | SPL↑ | DTS↓ |
| **Baselines** | **FR** [17] | | | 22.37% | 0.213 | 1.635 | 8.96% | 0.085 | 4.242 | 8.18% | 0.075 | 2.976 | 12.54% | 0.012 | 3.049 |
| | **Nav-A3C** [18] | | | 44.83% | 0.318 | 1.327 | 37.72% | 0.320 | 2.138 | 23.22% | 0.200 | 3.245 | 33.90% | 0.272 | 2.355 |
| | **MSM** [19] | | | 47.32% | 0.260 | **1.097** | 25.60% | 0.114 | 2.220 | 18.08% | 0.117 | 3.440 | 30.50% | 0.167 | 2.259 |
| | **VGM** [20] | | | 47.13% | 0.295 | 1.313 | 42.04% | 0.278 | 2.478 | 36.22% | 0.272 | 3.702 | 41.37% | 0.295 | 2.590 |
| | **NRNS** [21] | | | 43.00% | 0.317 | 1.493 | 30.60% | 0.209 | 2.106 | 18.61% | 0.133 | 3.575 | 30.73% | 0.219 | 2.391 |
| | **SLING+DDPPO** [22, 23] | | | 40.00% | 0.273 | 2.007 | 36.50% | 0.272 | 2.136 | 24.00% | 0.205 | 3.191 | 33.50% | 0.250 | 2.445 |
| | **SLING+OVRL** [22, 24] | | | 63.18% | **0.477** | 1.385 | 55.00% | **0.395** | 1.749 | 51.75% | 0.267 | 3.269 | 56.64% | 0.379 | 2.134 |
| **Ablation** | $f_u$ | AT | CBAM | SR↑ | SPL↑ | DTS↓ | SR↑ | SPL↑ | DTS↓ | SR↑ | SPL↑ | DTS↓ | SR↑ | SPL↑ | DTS↓ |
| | ✘ | ✘ | ✘ | 45.74% | 0.340 | 1.577 | 34.02% | 0.278 | 2.371 | 30.87% | 0.251 | 3.052 | 36.48% | 0.287 | 2.391 |
| | ✔ | ✘ | ✘ | 55.09% | 0.404 | 1.464 | 39.41% | 0.326 | 2.396 | 37.86% | 0.316 | 2.957 | 43.53% | 0.316 | 2.346 |
| | ✘ | ✔ | ✔ | 52.13% | 0.367 | 1.605 | 42.31% | 0.319 | 2.203 | 44.40% | 0.324 | 3.048 | 46.25% | 0.336 | 2.232 |
| | ✔ | ✘ | ✔ | 54.41% | 0.368 | 1.546 | 42.21% | 0.272 | 2.096 | 42.10% | 0.292 | 3.032 | 45.49% | 0.307 | 2.204 |
| | ✔ | ✔ | ✘ | 62.23% | 0.422 | 1.406 | 51.48% | 0.352 | 1.847 | 47.15% | 0.328 | 2.794 | 53.15% | 0.365 | 2.190 |
| **Ours** | ✔ | ✔ | ✔ | **66.00%** | 0.458 | 1.289 | **59.50%** | 0.388 | **1.511** | **54.75%** | **0.388** | **2.684** | **59.65%** | **0.385** | **2.053** |

color features for all sampling points, we perform voxel rendering on the color features to generate a compressed spatial feature map $I_{cog} \in \mathbb{R}^{64 \times W \times H}$. The specific procedure is described as follows:

$$f_c = \text{MLP}_{\theta_2}(f, \gamma_d(\text{r}(t))), \quad (9)$$

$$I_{cog}(\text{r}) = \sum_{i=1}^{N_s} \alpha_i f_c(\text{r}(t_i)), \quad (10)$$

where $\text{MLP}_{\theta_2}$ is the same as in equation 2, $N_s$, $\alpha_i$ and $\text{r}(t_i)$ is the same as in equation 4.

The spatial feature map $I_{cog} \in \mathbb{R}^{64 \times W \times H}$ contains the structural information of the scene from the current viewpoint. We concatenate the target RGB image with $I_{cog}$ along the channel dimension, and the concatenated feature map undergoes two residual blocks for additional feature extraction. Subsequently, the compressed features are passed through CBAM structure to eliminate noise and amplify navigation-related information. Finally, the cognitive information is compressed into a 64-length feature vector $f_{cog}$ using an MLP. This step associates the spatial cognitive ability of the robot with target navigation, enhancing the robot's navigation capabilities during the exploitation phase.

**Auxiliary Task.** We aggregate cognitive information through an implicit process, ensuring an end-to-end characteristic of the network. To improve network interpretability, we introduce an auxiliary task predicting the angle between the robot's current orientation and the target direction. The operation procedure of the auxiliary task is to input the feature vector $f_{cog}$ into the two-layer MLP to obtain the predicted value $\hat{\alpha}$ for the angle of the pinch, and compare $\hat{\alpha}$ and the true value $\alpha$ using the $L1$ loss function, and optimize it by minimizing the loss.

### D. Multi-thinking Integration

We use adaptive feature fusion for multi-thinking integration. Before feature fusion, we generate real-time perception features through visual inputs, which are crucial for the robot's obstacle avoidance capability. Here, we utilize a visual encoder to extract structural scene information.

To enable the balancing of exploratory and exploitative behaviors, we concatenate feature vectors $f_{cog}$, $f_u$, and $f_p$ into a fused feature $f_{cat_1}$. Subsequently, an attention layer with perceptrons calculates weights $w$ for each feature, which is then element-wise multiplied with $f_{cat_1}$ to produce $f_{cat_2}$, allowing the network to allocate attention adaptively. Afterward, $f_{cat_2}$ is fed into the navigation policy network to generate navigation action. The specific process is as follows:

$$f_{cat_2} = \text{MLP}(w \otimes f_{cat_1}), \quad (11)$$

$$a = \text{Sample}(\text{Softmax}(\text{MLP}(f_{cat_2}))), \quad (12)$$

where $f_{cat_1} = \text{Concat}(f_{cog}, f_u, f_p)$, $w = \text{MLP}(f_{cat_1})$, and $\text{Sample}(\cdot)$ refers to the probability-based sampling.

## IV. EXPERIMENT

### A. Implementation

**Task Setup.** We conduct image-goal navigation tasks in iGibson [25]. NUE is trained using imitation learning on the Gibson dataset [26] with 21 scenes as training splits and 14 scenes as validation splits. These scenes are categorized into 3 difficulty levels based on the distance from the starting point to the target: 1) Easy: 1.5m - 3.0m; 2) Medium: 3.0m - 5.0m; 3) Hard: 5.0m - 10.0m. The robot has only access to the current RGBD observation and its current pose. The RGBD observation is obtained from a single RGBD camera with $180 \times 240$ resolution and $90°$ horizontal field of view. The maximum time step for each episode is set to 800. An event is considered successful when the robot reaches within a range of 0.8m of the target. Three evaluation metrics are used: success rate (SR), success weighted by path length (SPL), and distance to success (DTS).

**Baselines and Ablation Study.** We first compare the navigation efficiency of our model with the most basic models: Fully Reactive (**FR**) [17] and **Nav-A3C** [18]. In addition, to validate the effectiveness of our cognitive structure, we include **VGM** [20] and Multi-Storage Memory (**MSM**) [19] as baselines. We also compare our method with recent state-of-the-art image-goal navigation approaches. **NRNS** [21] constructs a topological map by predicting the distance between the target image and node images. **SLING** [22] enhances the robot's navigation capability in the "last mile" by predicting the relative pose of the target based on keypoint matching. We combine SLING with the reinforcement
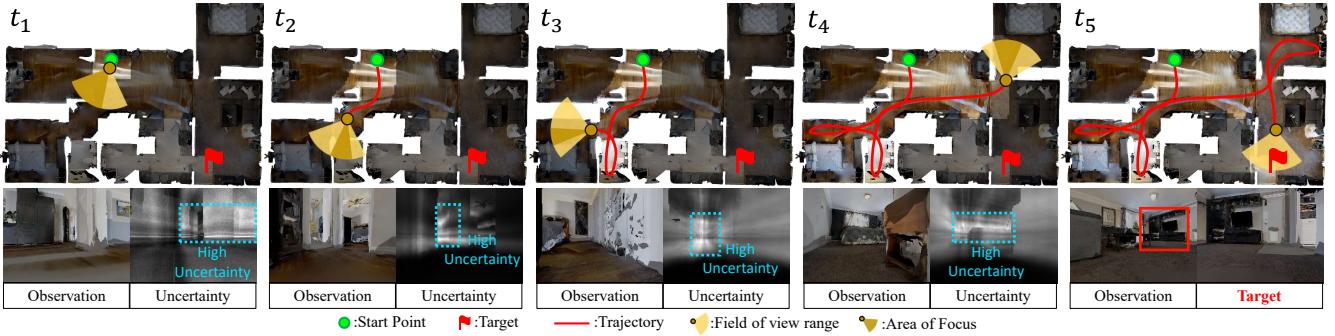
Fig. 4: **Visualization examples of image-goal navigation.** The visualized results showcase the behavioral logic of our model. Our model successfully utilizes uncertainty to explore the scene in the early stages of navigation, while also effectively leveraging cognitive information for navigation upon sighting the target.

learning-based **DDPPO** [23] and **OVRL** [24] as baselines. In ablation study, we conduct ablation on 3 parts: 1) $f_u$: we remove $f_u$, and only fuse $f_{cog}$ and $f_p$ to generate actions; 2) *auxiliary task* (AT): we delete the auxiliary task for target direction prediction from the full model; 3) CBAM: we remove the CBAM layer and rely solely on the residual network and MLPs for cognitive feature extraction.

### B. Image-goal navigation Results

Table I presents the average SR, SPL, and DTS for each method. Compared to FR and Nav-A3C, our model achieves significantly enhanced navigation performance, highlighting the significance of cognitive structure. Moreover, our model outperforms baselines with internal scene representation (MSM, VGM, NRNS), showcasing the effectiveness of NeRF as a scene representation method and the value of our information extraction approach. Furthermore, compared to reinforcement learning methods enhanced by last-mile navigation (SLING+DDPPO, SLING+OVRL), our model still exhibits superior performance in hard scenarios, which are more exploration-dependent. This underscores the benefits of using uncertainty to enhance exploration behaviors for robot navigation in complex environments.

The ablation study is also shown in Table I. The model only with the addition of $f_u$ demonstrated significantly better performance compared to the fully ablated model. Meanwhile, The removal of $f_u$ led to a significant decrease in the navigation performance of our model, demonstrating that our approach of enhancing exploration using uncertainty improves the robot's ability to establish cognitive understanding. Similarly, the ablation of AT also resulted in a decrease in navigation performance, indicating that our auxiliary task effectively improves the robot's reasoning capability for determining the target's position. Furthermore, after ablating the CBAM layer, we can observe a decrease in the navigation performance of our model, indicating that CBAM attention contributes to improving feature extraction effectiveness.

### C. Interpretability Experiments

*1) Visualization of Typical Navigation Behaviors:* In Fig. 4, we visualize the robot's trajectory in a typical testing episode. At the initial stage of this task ($t_1$), the cognitive
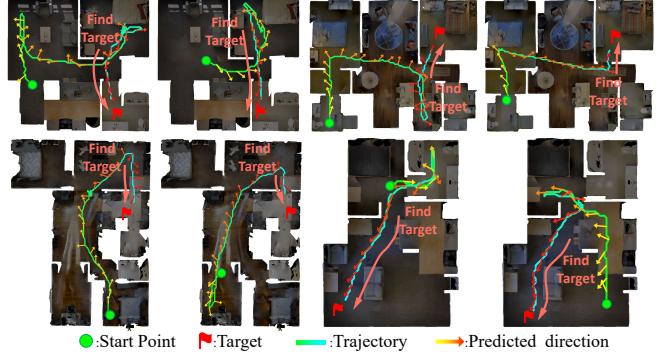


Fig. 5: **Predicted results of the auxiliary task.** The change in arrow color from yellow to red indicates the progress of navigation. During navigation tasks, auxiliary task accuracy steadily enhances, especially after observing the target.

structure has limited environmental memory, resulting in low signal-to-noise ratio outputs. At this phase, uncertainty information provides additional assistance to the robot's decision-making. As the robot observes a new room ($t_2$, $t_3$, $t_4$), higher uncertainty is observed in unexplored areas, guiding the robot to explore new regions and avoiding redundant paths. Finally, when the robot encounters an area similar to the target image ($t_5$), it decisively chooses to move toward it. This experiment confirms that our pipeline effectively achieves the transition from exploration to exploitation.

*2) Results of Auxiliary Task:* To assess the auxiliary task's efficacy, we select 8 navigation tasks across 4 scenes and visualize the predicted target direction in Fig. 5. Initially, with limited robot exploration range and lack of target-related cues, the cognitive structure lacks target memory, resulting in subpar prediction performance. A noticeable mismatch between predicted and actual robot actions suggests exploratory thinking guiding decision-making. However, as the robot's exploration expands, auxiliary task accuracy notably improves. At this stage, exploitative thinking prevails, effectively steering the robot towards the target. The visualization confirms the auxiliary task's role in enhancing the robot's reasoning for the target position.

*3) Interpretation of Uncertainty Extraction:* Fig. 6 showcases the gradient visualization of our uncertainty extraction structure using Grad-CAM. The color gradient from blue to
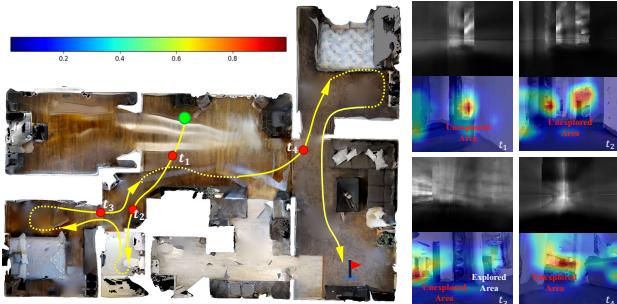
Fig. 6: **Interpretation results of uncertainty extraction.** The color in the image signifies the attention level of our uncertainty extraction structure on various regions. The top image in each pair shows the uncertainty map generated by NeRF, while the bottom image displays the gradient heatmap used for policy generation.

red represents the ascending order of the model's attention intensity. It can be observed that the model can effectively allocate high attention to regions with high uncertainty, which indicates that our uncertainty extraction structure can focus on unexplored regions through the uncertainty map, thus guiding the robot to explore these areas.

## V. CONCLUSION

We introduce an end-to-end visuomotor navigation framework, NUE, which applies the powerful scene representation capability of NeRF to the field of image-goal navigation, and specifically enhances the robot's exploratory behavior through uncertainty estimation and extraction. This innovative approach overcomes the problem of insufficient focus on exploratory behavior in traditional methods, enabling the robot to rapidly establish environmental cognition to provide more information for the subsequent exploitation phase. Our experiments provide compelling evidence that our framework has achieved active exploration behavior for the robot, thereby improving the efficiency of navigation.

## REFERENCES

[1] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[2] S. Zhi, T. Laidlow, S. Leutenegger, and A. J. Davison, "In-place scene labelling and understanding with implicit scene representation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 15 838–15 847.

[3] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Transactions on Graphics (ToG)*, vol. 41, no. 4, pp. 1–15, 2022.

[4] Y. Wei, S. Liu, J. Zhou, and J. Lu, "Depth-guided optimization of neural radiance fields for indoor multi-view stereo," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

[5] X. Pan, Z. Lai, S. Song, and G. Huang, "Activenerf: Learning where to see with uncertainty estimation," in *European Conference on Computer Vision*. Springer, 2022, pp. 230–246.

[6] L. Yen-Chen, P. Florence, J. T. Barron, A. Rodriguez, P. Isola, and T.-Y. Lin, "inerf: Inverting neural radiance fields for pose estimation. in 2021 ieee," in *RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021, pp. 1323–1330.

[7] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys, "Nice-slam: Neural implicit scalable encoding for slam," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 786–12 796.

[8] J. Ichnowski, Y. Avigal, J. Kerr, and K. Goldberg, "Dex-nerf: Using a neural radiance field to grasp transparent objects," *arXiv preprint arXiv:2110.14217*, 2021.

[9] Y. Li, S. Li, V. Sitzmann, P. Agrawal, and A. Torralba, "3d neural scene representations for visuomotor control," in *Conference on Robot Learning*. PMLR, 2022, pp. 112–123.

[10] M. Adamkiewicz, T. Chen, A. Caccavale, R. Gardner, P. Culbertson, J. Bohg, and M. Schwager, "Vision-only robot navigation in a neural radiance world," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 4606–4613, 2022.

[11] M. Kurenkov, A. Potapov, A. Savinykh, E. Yudin, E. Kruzhkov, P. Karpyshev, and D. Tsetserukou, "Nfomp: Neural field for optimal motion planner of differential drive robots with nonholonomic constraints," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 10 991–10 998, 2022.

[12] O. Kwon, J. Park, and S. Oh, "Renderable neural radiance map for visual navigation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 9099–9108.

[13] P. Marza, L. Matignon, O. Simonin, and C. Wolf, "Multi-object navigation with dynamically learned neural implicit representations," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 11 004–11 015.

[14] J. Zeng, Y. Li, Y. Ran, S. Li, F. Gao, L. Li, S. He, J. Chen, and Q. Ye, "Efficient view path planning for autonomous implicit reconstruction," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 4063–4069.

[15] S. Lee, L. Chen, J. Wang, A. Liniger, S. Kumar, and F. Yu, "Uncertainty guided policy for active robotic 3d reconstruction using neural radiance fields," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 12 070–12 077, 2022.

[16] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "Cbam: Convolutional block attention module," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 3–19.

[17] Y. Zhu, R. Mottaghi, E. Kolve, J. J. Lim, A. Gupta, L. Fei-Fei, and A. Farhadi, "Target-driven visual navigation in indoor scenes using deep reinforcement learning," in *2017 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2017, pp. 3357–3364.

[18] M. Zhang, Z. McCarthy, C. Finn, S. Levine, and P. Abbeel, "Learning deep neural network policies with continuous memory states," in *2016 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2016, pp. 520–527.

[19] H. Sang, R. Jiang, Z. Wang, Y. Zhou, and B. He, "A novel neural multi-store memory network for autonomous visual navigation in unknown environment," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 2039–2046, 2022.

[20] O. Kwon, N. Kim, Y. Choi, H. Yoo, J. Park, and S. Oh, "Visual graph memory with unsupervised representation for visual navigation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 15 890–15 899.

[21] Z. Al-Halah, S. K. Ramakrishnan, and K. Grauman, "Zero experience required: Plug & play modular transfer learning for semantic visual navigation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17 031–17 041.

[22] J. Wasserman, K. Yadav, G. Chowdhary, A. Gupta, and U. Jain, "Last-mile embodied visual navigation," in *Conference on Robot Learning*. PMLR, 2023, pp. 666–678.

[23] E. Wijmans, A. Kadian, A. Morcos, S. Lee, I. Essa, D. Parikh, M. Savva, and D. Batra, "Dd-ppo: Learning near-perfect pointgoal navigators from 2.5 billion frames," *arXiv preprint arXiv:1911.00357*, 2019.

[24] K. Yadav, R. Ramrakhya, A. Majumdar, V.-P. Berges, S. Kuhar, D. Batra, A. Baevski, and O. Maksymets, "Offline visual representation learning for embodied navigation," in *Workshop on Reincarnating Reinforcement Learning at ICLR 2023*, 2023.

[25] C. Li, F. Xia, R. Martín-Martín, M. Lingelbach, S. Srivastava, B. Shen, K. Vainio, C. Gokmen, G. Dharan, T. Jain *et al.*, "igibson 2.0: Object-centric simulation for robot learning of everyday household tasks," *arXiv preprint arXiv:2108.03272*, 2021.

[26] F. Xia, A. R. Zamir, Z. He, A. Sax, J. Malik, and S. Savarese, "Gibson env: Real-world perception for embodied agents," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 9068–9079.