# BRAINAV: INCORPORATING HUMAN BRAIN ACTIV ITY TO ENHANCE ROBUSTNESS IN EMBODIED VISUAL NAVIGATION

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Recent research shows that standard navigation agents significantly underperform and even fail in the presence of various visual corruptions. Unlike embodied agents, the human brain's visual system can robustly perceive the environment and extract the necessary information to complete the visual tasks. In this paper, we propose a two-phase **Brain**-Machine integration **Nav**igation method called **BraiNav**, which incorporates neural representations derived from human brain activity to enhance robustness against visual corruptions. In the first phase, a brain encoder, built upon a recently advanced self-supervised pretrained model, is trained on a large-scale human brain activity dataset and then frozen for downstream visual navigation. In the second phase, neural representations harboring high-level cognitive information from the human brain are constructed based on the pretrained frozen brain encoder. Additionally, we propose a multimodal fusion method based on cross-attention to obtain more consistent brain-visual joint representations, which are then used to learn the navigation policy. Sufficient experiments demonstrate that the proposed method exhibits higher robustness against various visual corruptions compared to standard navigation agent and multiple computer vision-enhanced agents. Our study pioneers the incorporation of human brain activity into embodied AI, aiming to catalyze further cross-disciplinary collaboration with computational neuroscience.

030 031 032

033

005 006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

#### 1 INTRODUCTION

Embodied visual navigation (Anderson et al., 2018; Batra et al., 2020), one of the most researched topics in embodied AI (Duan et al., 2022), requires agents to make action decisions based on ego-centric observations to achieve their goals. Despite the remarkable progress in embodied visual navigation, efforts (Wijmans et al., 2019; Zhao et al., 2021; Zhang et al., 2023) have primarily focused on training agents to generalize to unseen environments, assuming similarities between training and testing environments. However, a major challenge in this field is ensuring agent generalization across environments with different visual appearances (Chattopadhyay et al., 2021; Rajič, 2022).

041 Following standard protocol, agents are trained on a set of scenes and evaluated on unseen scenes, 042 which entails two types of evaluation. The first type (Figure 1 (a)) assesses generalization per-043 formance on clean observations, where existing navigation methods have demonstrated high perfor-044 mance. The second type (Figure 1 (b)) introduces visual corruption, such as defocus blur, simulating real-life challenges. Assessing robustness to such corruption alongside generalization performance poses a greater challenge. While previous studies have utilized standard deep learning techniques, 046 such as data augmentation and self-supervised adaptation (Chattopadhyay et al., 2021), to enhance 047 robustness, there remains significant room for improvement in fully recovering lost navigation per-048 formance. 049

Unlike deep models, the human brain's visual system possesses superior capabilities in processing
 high-level semantic information from images, going beyond basic features like color, shape, and
 texture. As shown in Figure 1, humans can achieve comparable performance in navigation tasks un der corrupted observations, indicating the robustness of the human visual system. Recent researches
 have demonstrated improved performance in deep learning models by incorporating neural repre-

sentations (Fong et al., 2018; Li et al., 2019; Nishida et al., 2020; Dapello et al., 2020; Fel et al., 2022; Liu et al., 2023; Shah et al., 2024). Notably, to the best of our knowledge, no prior study has utilized activity data collected from human brains as guidance for embodied visual navigation.

057 Motivated by these discussions, we investigated 058 the potential of leveraging human brain ac-059 tivity for embodied visual navigation, propos-060 ing a two-phase Brain-Machine integration 061 Navigation method, called BraiNav. In the 062 first phase, a brain encoder model, built upon 063 a recently advanced self-supervised pretrained 064 model, is trained on a large-scale functional magnetic resonance imaging (fMRI) dataset. 065 Following pretraining, the brain encoder is 066 frozen for downstream visual navigation tasks. 067 In the second phase, the navigation agent 068 learns through deep reinforcement learning 069 (DRL) (Sutton & Barto, 2018), comprising four major components: a brain module, a 071 visual-target module, a fusion module, and a 072 policy module. At each timestep, the agent 073 obtains RGB observation and target localiza-074 tion. The RGB observation is initially pro-075 cessed by the brain module to extract neural representations, constructed based on the pre-076 trained frozen brain encoder. Simultaneously, 077 the visual-target module generates visual-target representations from RGB observation and tar-079



Figure 1: Comparison of human brain's visual system and embodied navigation agents in the presence of visual corruption. The human brain's visual system demonstrates superior robustness compared to navigation agents when observations are visually corrupted.

get localization. Then, a cross-attention (Vaswani et al., 2017) based multimodal fusion module is
proposed to jointly learn the neural and visual-target representations to obtain the brain-visual joint
representations. Finally, the joint representations are fed into the GRU-based (Cho et al., 2014) policy module, and the agent is trained using the Proximal Policy Optimization (PPO) algorithm (Schulman et al., 2017). Our results demonstrate that the BraiNav shows higher robustness against various
visual corruptions compared to standard navigation agent and multiple computer vision-enhanced
agents. To summarize, our main contributions are as follows:

- We introduce BraiNav, the first embodied navigation method incorporating human brain activity to enhance agent robustness against visual corruptions.
- We propose a multimodal fusion method based on cross-attention, yielding more consistent brain-visual joint representations.
- Sufficient experiments demonstrate that BraiNav exhibits superior robustness against various visual corruptions compared to standard navigation agent and multiple computer vision-enhanced agents. Our source code will be made available online post-publication.

# 2 Methodology

098 099

100

087

090

091 092

093

094

095 096

2.1 TASK DEFINITION

In this paper, we focus solely on PointNav. The navigation agent is trained using deep reinforcement learning. At each timestep t, the agent obtains the egocentric RGB observation  $o_t$  and target location  $l_t$ . The agent then takes a predicted action to reach the target in as few timesteps as possible. There are four available actions for the agent: move forward (0.25m), turn left (30°), turn right (30°), and stop. An episode is considered successful if the agent stops within 0.2m of the target and executes a maximum of 300 steps. Unlike previous configurations where the RGB observation is clean during evaluation, in our task, the observation is subject to various visual corruptions to evaluate the robustness of the navigation agent.



139 140 141

Figure 2: Overall Framework of BraiNav. Phase 1: The brain encoder model receives stimulus and outputs predicted fMRI responses, which are then supervised pretrained with experimentally measured fMRI responses. Phase 2: At each timestep, the RGB observation and target localization are fed into the brain module and visual-target module to obtain neural representation and visual-target representation, respectively. The multimodal fusion module combines these two different modality representations and outputs a consistent brain-visual joint representation, which is used to train the navigation agent.

150

149

# 2.2 OVERALL FRAMEWORK OF BRAINAV

BraiNav is designed with two sequential phases, as outlined in the Figure 2. In phase 1, a brain
BraiNav is designed with two sequential phases, as outlined in the Figure 2. In phase 1, a brain
encoder model, built upon a self-supervised pretrained frozen model, is trained on a large-scale
fMRI dataset. The learned brain encoder model will be frozen to guide the visual navigation process
in the next phase. The architecture and training details of the brain encoder model will be provided
in Section 2.3.

In phase 2, our agent consists of four major components: a brain module, a visual-target module, a fusion module, and a policy module. The brain module  $E_b$  contains the pretrained frozen brain encoder, which is used to construct neural representations  $x_t^b$  harboring high-level cognitive information from the human brain(to be discussed in Section 2.3), expressed as follows:

$$x_t^b = \mathcal{E}_b\left(o_t\right). \tag{1}$$

162 The visual-target module is not elaborately designed, as the main purpose of this paper is to demon-163 strate the effectiveness of human brain activity for robust navigation. Specifically, the visual-target 164 module  $E_{v}$  comprises a frozen visual encoder, a target encoder, a compressor network, and a com-165 biner network, with reference to (Chattopadhyay et al., 2021). The visual encoder is a frozen 166 ResNet-18 (He et al., 2016) pretrained on ImageNet (Deng et al., 2009), which extracts visual features from images with a 512×7×7 output. The compressor network consists of two convo-167 lutional layers, each followed by ReLU activation, transforming the output of the visual encoder to 168  $512 \times 7 \times 7 \rightarrow 128 \times 7 \times 7 \rightarrow 32 \times 7 \times 7$ . The target localization is a 2-dimensional polar coordinate  $(r, \theta)$  mapped to a 32-dimensional target representation by the target encoder (a single-layer MLP), 170 which is then expanded to  $32 \times 7 \times 7$  dimensions. The outputs of the compressor network and 171 target representation are concatenated and further processed by the combiner network, consisting of 172 two convolutional layers that transform the output to  $64 \times 7 \times 7 \rightarrow 128 \times 7 \times 7 \rightarrow 32 \times 7 \times 7$ . 173 Finally, the output of the combiner network is flattened to obtain a 1568-dimensional visual-target 174 representation  $x_t^v$ , expressed as follows: 175

176 177

183

$$x_t^v = \mathcal{E}_v\left(o_t, l_t\right). \tag{2}$$

The neural representation and visual-target representation belong to different modalities, and the intuitive concatenate method cannot yield promising result because inter-modal interactions are not fully exploited (Du et al., 2023; Chen et al., 2023). To address this problem, we propose a multimodal fusion module  $E_f$  (to be discussed in Section 2.4) based on cross-attention (Vaswani et al., 2017) to obtain a more consistent brain-visual joint representation  $x_t^{bv}$ , expressed as follows:

$$x_t^{bv} = \mathcal{E}_f\left(x_t^b, x_t^v\right). \tag{3}$$

The joint representation is fed into a GRU with 512 hidden units, along with the previous hidden state. The GRU outputs a 512-dimensional vector and the next hidden state, followed by two independent MLPs(marked as Actor and Critic) that receive the 512-dimensional vector and output 4-dimensional action logits and a 1-dimensional scalar value, respectively. At timestep *t*, the reward received by the agent can be expressed as,

$$r_t = \mathcal{R}_{\text{success}} - \triangle_t^{\text{Geo}} + \lambda, \tag{4}$$

where  $R_{success} = 10$  denotes the reward obtained for a successful episode,  $\Delta_t^{Geo}$  denotes the change in geodesic distance to the target at a single timestep interval, and  $\lambda = -0.01$  denotes the movement penalty. The agent is trained using the Proximal Policy Optimization (PPO) (Schulman et al., 2017) algorithm.

195 196 197

190 191

#### 2.3 BRAIN ENCODER PRETRAINING

The brain encoder model follows the method proposed in (Adeli et al., 2023) and consists of three 199 components: an encoder, a decoder, and prediction heads, as shown in Figure 2. In the encoder, 200 the input image is first divided into fixed  $14 \times 14$  patches. These image patches are then processed 201 by DINOv2 (Oquab et al., 2023), a recent advanced foundation model capable of extracting highperformance visual features. Specifically, the DINOv2 used is a self-supervised pretrained frozen 202 ViT-B/14 (Dosovitskiy et al., 2020) model, which consists of a linear projection layer and a Trans-203 former encoder. Each image patch is projected into a 768-dimensional patch embedding by the 204 linear projection layer, and an additional learnable patch embedding is prepended to the sequence, 205 denoted as CLS. Position embeddings are added to the patch embeddings before they are fed into 206 the Transformer encoder, which consists of 12 stacked Transformer blocks (Vaswani et al., 2017). 207

The decoder is a single-layer Transformer with a feed-forward dimension of 2048 with 16 attention heads. The brain ROI queries consist of M learnable positional encodings corresponding to different brain ROIs in each hemisphere. Specifically, 16 queries are used, with 8 for each hemisphere (7 streams add 1 for all the vertices labeled 'unknown'). The Transformer decoder utilizes the output from the encoder to transform these queries into output tokens.

The prediction heads comprise 16 MLPs that map the output tokens from the decoder to the fMRI responses of the corresponding ROIs. Each output token is mapped to a vector by the MLP with the same number of voxels as the corresponding hemisphere (19,004 for the left hemisphere and 20,544 for the right hemisphere), with voxels not belonging to that ROI set to 0 using a mask. The prediction loss is computed by calculating the Mean Square Error (MSE) between the predicted and
 measured fMRI responses voxel by voxel.

After pretraining, we construct neural representations for downstream visual navigation using the pretrained frozen brain encoder, as shown in the brain module in Figure 2. Specifically, we use the low-dimensional and noise-reduced output tokens from the brain encoder to construct neural representations, effectively exploiting the high-level information embedded in human brain activities. Each 768-dimensional output token from each ROI of the anatomical streams is mapped to 48 dimensions using an MLP, and the resulting vectors are concatenated and fed into another MLP to obtain 768-dimensional neural representation.

#### 2.4 BRAIN-VISUAL MULTIMODAL LEARNING



Figure 3: Architecture of the multimodal fusion module. The module is primarily composed of N stacked computational blocks, each consisting of two branches. Each branch processes neural representations and visual-target representations of different dimension, which are effectively fused at the end using cross-attention mechanisms.

The architecture of the multimodal fusion module is shown in Figure 3, consisting of N stacked computational blocks. Each computational block includes multiple Cross-Attention (CA) layers and MLPs, and contains two branches.

For simplicity, the timestep t is omitted. In the brain branch (top), the Key embedding  $K_{i-1}^{v} \in \mathbb{R}^{V_{\text{dim}}}$  and Value embedding  $V_{i-1}^{v} \in \mathbb{R}^{V_{\text{dim}}}$  are first obtained by identity mapping the visual-target representation  $x_{i-1}^{v} \in \mathbb{R}^{V_{\text{dim}}}$ ,

$$\mathbf{K}_{i-1}^{v}, \mathbf{V}_{i-1}^{v} = \text{Identity}\left(x_{i-1}^{v}\right).$$
(5)

To align the dimensions, the Query embedding  $Q_{i-1}^b \in \mathbb{R}^{V_{dim}}$  is obtained by applying a singer-layer MLP to the neural representation  $x_{i-1}^b \in \mathbb{R}^{B_{dim}}$ ,

$$\mathbf{Q}_{i-1}^{b} = \mathrm{MLP}\left(x_{i-1}^{b}\right). \tag{6}$$

Then, the brain Query interacts with the visual-target Key and Value through CA, mathematically expressed as:

$$\widetilde{x}_{i}^{b} = CA \left( Q^{b}, K^{v}, V^{v} \right) 
= \operatorname{softmax} \left( \mathbf{q} \mathbf{k}^{T} / \sqrt{B_{\dim} / V_{\dim}} \right) \mathbf{v}, \qquad (7) 
\mathbf{q} = Q^{b} \mathbf{W}_{q}, \mathbf{k} = K^{v} \mathbf{W}_{k}, \mathbf{v} = V^{v} \mathbf{W}_{v},$$

where  $\mathbf{W}_q$ ,  $\mathbf{W}_k$ , and  $\mathbf{W}_q$  are learnable parameters. Finally, to align the dimensions,  $x_i^b \in \mathbb{R}^{B_{\text{dim}}}$  is obtained by mapping the output of CA with a residual shortcut through a single-layer MLP,

$$x_i^b = \mathrm{MLP}\left(\tilde{x}_i^b + \mathbf{Q}_{i-1}^b\right). \tag{8}$$

For the visual-target branch (bottom), the same procedure as the brain branch is performed, simply swapping Query with Key and Value. After N computational blocks, the final neural representation  $x_N^b$  and visual-target representation  $x_N^v$  are concatenated to produce the joint representation  $x_N^{bv}$ . Experimentally, the fusion module uses only one computational block (N = 1) for convergence considerations, where the cross-attention is a Transformer model with 3 heads.

275 276 277

278 279

280 281

282

283

284

285

286

299

300

301

302

303 304

# 3 EXPERIMENTS

# 3.1 DATASET AND ENVIRONMENT

The brain encoder is trained using the Natural Scenes Dataset (NSD) dataset (Allen et al., 2022), which is currently the largest and richest neuroimaging dataset. It provides high-quality whole-brain 7T fMRI responses from 8 subjects viewing  $\sim$ 73,000 different natural scenes while performing a continuous recognition task. The color image stimuli viewed by the subjects come from the Common Objects in Context (COCO) dataset (Lin et al., 2014), which is richly annotated and widely used in computer vision.



Figure 4: Navigation scene and visual corruptions. (a) An example of an indoor scene from RoboTHOR dataset, with each scene sized at  $8.8m \times 3.9m$ . (b) The top-left shows a clean RGB observation, while the remaining images display its corresponding corrupted observations. The illustrations of the corruptions are adapted from (Chattopadhyay et al., 2021).

The visual navigation agent will be trained and evaluated on the ROBUSTNAV benchmark (Chattopadhyay et al., 2021), built on top of the AI2-THOR simulator (Kolve et al., 2017) and the RoboTHOR dataset (Deitke et al., 2020). RoboTHOR contains 75 indoor scenes with different layouts(Figure 4(a) shows an example scene), with 60 used for training and 15 for evaluation. BraiNav is evaluated on 7 visual corruptions provided in ROBUSTNAV (shown in the Figure 4(b)): Defocus Blur, Motion Blur, Spatter, Camera Crack, Low Lighting, Lower FOV, and Speckle Noise. A detailed description of these visual corruptions can be found in (Chattopadhyay et al., 2021).

311 312 313

314

# 3.2 EXPERIMENTS CONFIGURATION

For brain encoder pretraining, the learning rate is set to 1e-4, the batch size to 32, the number of training epochs to 20, the gradient clip to 0.1, the weight decay to 1e-4, and dropout to 0.1. The subsequent brain encoder is trained using the fMRI data from subject 01, with results from other subjects presented shown in Appendix B.1. Additionally, the original  $425 \times 425$  resolution images are resized to  $224 \times 224$  to match the observation image resolution for downstream visual navigation tasks.

For BraiNav training, the learning rate is set to 3e-4, the discount factor to 0.99, the rollout length to 128, and the total training frames to 75M. The corruption level is consistently set to 5, as in ROBUSTNAV, where level 5 indicates the most severe corruptions. Training is conducted on a single NVIDIA GeForce RTX 4090 GPU.

# 324 3.3 EVALUATION METRICS

We use two common metrics in visual navigation: Success Rate (SR) and Success weighted by Path Length(SPL). SR indicates the percentage of successful episodes, defined as follows:

$$SR = \frac{1}{K} \sum_{i=1}^{K} \Pi_i, \tag{9}$$

(10)

Where K is the number of episodes in the evaluation and  $\Pi_i$  is the binary indicator whether the *i*-th episode is successful (1 if successful, otherwise 0). SPL is the percentage of the path length of successful episodes to the shortest path, defined as follows:

326

327

328

330 331 332

333

334

337 338

339

340

where l is the shortest path length and p is the agent's path length. Higher SR and SPL indicate that the navigation agent is more effective and efficient.

 $SPL = \frac{1}{K} \sum_{i=1}^{K} \prod_{i} \frac{l_i}{\max(p_i, l_i)},$ 

# 341342 3.4 PERFORMANCE OF BRAIN ENCODER

343 After pretraining, we evaluate the prediction accuracy of the brain encoder on the ROIs used to con-344 struct the neural representations. We compute the noise ceiling for each vertex using response data 345 from the subjects' three trials of the same stimulus image, and then average these values across each 346 ROI. Details of the calculation methodology and results are provided in Appendix B.2. Next, we 347 evaluate the brain encoder's performance by calculating the *Pearson correlation coefficient* between 348 its predicted fMRI responses and the actual values. To determine the noise-normalized prediction 349 accuracy, we divide the encoder's prediction accuracy by the corresponding noise ceiling, as summarized in Table 1. 350

The results indicate that the brain encoder effectively predicts the fMRI response to visual stimuli.
 High prediction accuracy ensures that the neural representations used for learning the navigation policy contain cognitive processing information from the human brain.

354 355

356

# 3.5 COMPARISON WITH STANDARD NAVIGATION AGENT

We compare the proposed method with the standard navigation agent (Chattopadhyay et al., 2021), as shown in Table 2. BraiNav outperforms the standard navigation agent across 6 visual corruptions, achieving varying degrees of performance improvement. Notably, on Speckle Noise and Defocus Blur, BraiNav shows significant improvements, with absolute improvements of (15.77%SPL, 17.75%SR) and (11.89%SPL, 11.37%SR), respectively. For other corruptions, the absolute improvements are (4.93%SPL, 1.37%SR) for Motion Blur, (2.76%SPL, 3.18%SR) for Spatter, (1.8%SPL, 0.73%SR) for Low Lighting, and (0.83%SPL, 0.46%SR) for Lower FOV.

It is worth noting that our method exhibits remarkable performance improvement on Defocus Blur.
 Neuroscience studies (Mon-Williams et al., 1998; Webster et al., 2002; Zhu et al., 2013) suggest the
 existence of a perceptual mechanism in the human brain that regulates image defocus. In conclusion,
 these experimental results demonstrate that our proposed navigation method, BraiNav, surpasses the
 standard navigation agent across various visual corruptions.

- 369
- 370 3.6 Comparison with Computer Vision-Enhanced Agents

We also compare our proposed method with several computer vision-enhanced agents outlined in (Chattopadhyay et al., 2021), as shown in Table 2.

**Standard Agent+AP**: Standard Agent introduces an auxiliary action prediction task.

Standard Agent+AP+SS-Adapt: Standard Agent+AP introduces self-supervised adaptation on specific corruptions.

**Standard Agent+RP**: Standard Agent introduces an auxiliary rotation prediction task.

Table 1: Noise-normalized prediction accuracy of the brain encoder.

Hemisphere			ROI				
	early	midventral	midlateral	midparietal	ventral	lateral	parietal
Left Hemisphere	0.7358	0.8254	0.8116	0.8432	0.7239	0.9179	0.7572
Right Hemishpere	0.7144	0.7910	0.7950	0.7239	0.8213	0.8761	0.6816

Table 2: Comparison with standard and computer vision-enhanced navigation agents. We compare BraiNav with six approaches proposed in (Chattopadhyay et al., 2021). The results are highlighted with best and second.

Approach							Visual C	orruption	L							
	Cl	ean	Spa	atter	Speckl	e Noise	Camer	a Crack	Lowe	r FOV	Defoc	us Blur	Motic	on Blur	Low L	ighting
	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL
Standard Agent	98.82	83.13	33.58	24.72	67.42	48.57	82.07	63.83	42.49	31.73	75.89	53.55	95.72	73.37	94.36	75.15
Standard Agent +AP	98.45	83.28	20.38	15.70	65.61	47.03	72.70	56.82	45.68	35.14	83.35	61.51	94.81	74.30	92.17	76.11
Standard Agent +AP+SS-Adapt	37.31	31.03	14.19	10.29	\	\	57.87	46.72	32.94	26.09	40.95	33.35	\	\	\	\
Standard Agent +RP	98.73	82.53	23.48	18.63	78.98	55.92	67.06	53.70	44.95	32.74	32.21	22.47	91.63	65.27	89.81	67.38
Standard Agent +RP+SS-Adapt	94.63	77.25	61.06	47.16	\	\	60.42	49.37	50.59	36.10	79.16	62.74	\	\	\	\
Standard Agent +Data Aug	98.45	81.08	23.93	18.41	77.25	57.95	88.44	71.57	71.70	54.54	81.26	61.32	96.91	75.97	97.27	78.74
BraiNav	97.73	80.88	36.76	27.48	85.17	64.34	77.80	60.25	42.95	32.56	87.26	65.44	97.09	78.30	95.09	76.95

**Standard Agent+RP+SS-Adapt**: Standard Agent+RP introduces self-supervised adaptation on specific corruptions.

406 Standard Agent+Data Aug: Standard Agent introduces various data augmentation during training.

BraiNav outperforms all the aforementioned computer vision-enhanced agents across 3 visual corruptions. Specifically, the absolute improvements are (6.39%SPL, 6.19%SR) for Speckle Noise, (2.7%SPL, 3.91%SR) for Defocus Blur, and (2.33%SPL, 0.18%SR) for Motion Blur. Additionally, BraiNav achieves the second best results on Spatter and Low Lighting, demonstrating competitive performance on the remaining visual corruptions. Additional experiments and detailed results of the comparison methods are provided in the Appendix B.3 and Appendix B.4.

3.7 ABLATION STUDY

BraiNav consists of two key components: the neural representation from the pretrained brain encoder and the multimodal fusion module. In this section, we first analyze the impact of the neural representation, followed by an evaluation of the multimodal fusion module's effectiveness.

For the first ablation experiment, we concatenate the original DINOv2 CLS representation with the visual-target representation to form the joint representation, labeled as BraiNav w/o NR in Table 3.

Except for Lower FOV, BraiNav consistently outperforms BraiNav w/o NR across clean and the
other 6 visual corruptions, achieving performance improvements to varying degrees. Specifically,
the absolute improvements are (14.6%SPL, 18.56%SR) for Speckle Noise, (9.08%SPL, 13.01%SR)
for Spatter, (7.42%SPL, 9.83%SR) for Low Lighting, (3.84%SPL, 5.37%SR) for Motion Blur,
(1.9%SPL, 8.64%SR) for Defocus Blur, and (2.55%SR) for Camera Crack. These results highlight that neural representation derived from human brain activity significantly enhances BraiNav's
robustness against visual corruptions, beyond the deep representation from DINOv2.

For the second ablation experiment, we replace the multimodal fusion module in BraiNav with
 the concatenate method, and the experimental results are labeled BraiNav w/o MF as shown in
 Table 3.

Table 3: Contribution of each component. **BraiNav w/o NR** presents experimental results without the neural representations (NR) derived from human brain activity. **BraiNav w/o MF** presents experimental results excluding the multimodal fusion (MF) module.

Approach						V	isual Co	orruptio	n							
	Cle	ean	Spa	itter	Speckl	e Noise	Camera	a Crack	Lowe	r FOV	Defoci	ıs Blur	Motio	n Blur	Low L	ighting
	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL
BraiNav w/o NR	95.45	79.74	23.75	18.40	66.61	49.74	75.25	61.40	48.95	38.63	78.62	63.54	91.72	74.46	85.26	69.53
BraiNav w/o MF	99.18	83.14	23.38	18.38	54.69	42.54	79.25	62.36	38.71	30.43	76.34	53.45	92.54	71.96	91.63	72.52
BraiNav	97.73	80.88	36.76	27.48	85.17	64.34	77.80	60.25	42.95	32.56	87.26	65.44	97.09	78.30	95.09	76.95

Except for Camera Crack, the multimodal fusion module enhances BraiNav's robustness across the other 6 visual corruptions. Specifically, the absolute improvements are (21.8%SPL, 30.48%SR) for Speckle Noise, (11.99%SPL, 10.92%SR) for Defocus Blur, (9.1%SPL, 13.38%SR) for Spatter, (6.34%SPL, 4.55%SR) for Motion Blur, (4.43%SPL, 3.46%SR) for Low Lighting, and (2.13%SPL, 4.24%SR) for Lower FOV. Overall, these findings demonstrate that the multimodal fusion module further enhances the robustness of BraiNav. Additional ablation experiments and detailed results are provided in the Appendix B.5.

#### 

# 4 DISCUSSION AND CONCLUSION

In this paper, we introduce BraiNay, a novel framework designed to address the robustness chal-lenges in embodied visual navigation. Our two-phase method leverages human brain activity to enhance the navigation agent's resilience to visual corruption. In the first phase, we pretrain a brain encoder model with DINOv2 as the backbone on a large-scale fMRI dataset. In the second phase, we utilize the pretrained frozen brain encoder to construct neural representations that encapsulate high-level cognitive information from the human brain. Additionally, we develop a multimodal fusion module based on cross-attention to facilitate the learning of consistent brain-visual joint representations for navigation policy acquisition. We evaluate BraiNav's navigation performance across multiple visual corruptions, demonstrating its superior robustness compared to standard visual navi-gation agent and multiple computer vision-enhanced agents. More importantly, our research bridges embodied AI and neuroscience, showcasing the potential for translating insights from neuroscience into advancements in embodied AI. 

To create a brain-like representation for the navigation agent, BraiNav first employs a brain encoder pretrained on fMRI data to obtain neural representations, followed by a multimodal fusion module to achieve more consistent joint representations. Thus, enhancements in both the brain encoder and the multimodal fusion module are expected to yield better performance. Future research could ex-plore different brain encoder architectures (Yang et al., 2024) and multimodal fusion methods (Mao et al., 2023). Furthermore, the brain module in BraiNav is decoupled from the specific task; it can receive images and output brain-like representations to improve the agent's robustness to visual cor-ruptions for other embodied tasks. Future studies could apply BraiNav's brain-like representations to a broader range of embodied tasks (Wan et al., 2023). 

# References

- Hossein Adeli, Sun Minni, and Nikolaus Kriegeskorte. Predicting brain activity using transformers. *bioRxiv*, pp. 2023–08, 2023.
- Emily J Allen, Ghislain St-Yves, Yihan Wu, Jesse L Breedlove, Jacob S Prince, Logan T Dowdle,
  Matthias Nau, Brad Caron, Franco Pestilli, Ian Charest, et al. A massive 7t fmri dataset to bridge
  cognitive neuroscience and artificial intelligence. *Nature neuroscience*, 25(1):116–126, 2022.
- Peter Anderson, Angel Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta,
   Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, et al. On
   evaluation of embodied navigation agents. *arXiv preprint arXiv:1807.06757*, 2018.

486 487 488	Dhruv Batra, Aaron Gokaslan, Aniruddha Kembhavi, Oleksandr Maksymets, Roozbeh Mottaghi, Manolis Savva, Alexander Toshev, and Erik Wijmans. Objectnav revisited: On evaluation of embodied agents navigating to objects. <i>arXiv preprint arXiv:2006.13171</i> , 2020.
489	
490	Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva,
491	Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor
492	environments. arXiv preprint arXiv:1709.06158, 2017.
493	Prithvijit Chattonadhyay, Judy Hoffman, Roozheh Mottaghi, and Aniruddha Kembhayi, Robustnay,
494	Towards benchmarking robustness in embodied navigation. In <i>Proceedings of the IEEE/CVF</i>
495	International Conference on Computer Vision, pp. 15691–15700, 2021.
496	
497	Xiaoyu Chen, Changde Du, Qiongyi Zhou, and Huiguang He. Auditory attention decoding with
498 499	task-related multi-view contrastive learning. In <i>Proceedings of the 31st ACM International Con-</i> <i>ference on Multimedia</i> , pp. 6025–6033, 2023.
500	Kyunghyun Cho, Bart Van Merriänhoer, Caglar Gulcehre, Dzmitry Rahdanau, Fethi Bougares, Hol
501	ger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder
502	
503	Joel Dapello, Tiago Marques, Martin Schrimpf, Franziska Geiger, David Cox, and James J DiCarlo.
504	Simulating a primary visual cortex at the front of cnns improves robustness to image perturbations.
505	Advances in Neural Information Processing Systems, 33:13073–13087, 2020.
507	Matt Deitke Winson Han Alvaro Herresti Aniruddha Kembhavi Eric Kolve Doorbeb Motteebi
507	Jordi Salvador, Dustin Schwenk, Eli VanderBilt, Matthew Wallingford, et al. Robothor: An
509	computer vision and nattern recognition pp. 3164, 3174, 2020
510	computer vision and pattern recognition, pp. 5104–5174, 2020.
511	Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-
512	erarchical image database. In 2009 IEEE conference on computer vision and pattern recognition,
513	pp. 248–255. leee, 2009.
514	Alexey Dosovitskiy Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
515	Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An im-
516 517 518	age is worth 16x16 words: Transformers for image recognition at scale. In <i>International Conference on Learning Representations</i> , 2020.
510	Changeda Dy Kajahang Ey Jinnang Li and Huigyang Ha Decoding visual neural representations
520 521	by multimodal learning of brain-visual-linguistic features. <i>IEEE Transactions on Pattern Analysis</i> and Machine Intelligence, 2023
522	una machine metagenee, 2023.
523	Jiafei Duan, Samson Yu, Hui Li Tan, Hongyuan Zhu, and Cheston Tan. A survey of embodied
524	ai: From simulators to research tasks. IEEE Transactions on Emerging Topics in Computational
525	Intelligence, 6(2):230–244, 2022.
526	Thomas Fel Ivan F Rodriguez Rodriguez Drew Linsley and Thomas Serre Harmonizing the object
527	recognition strategies of deep neural networks with humans. Advances in neural information
528	processing systems, 35:9432–9446, 2022.
529	
530	Ruth C Fong, Walter J Scheirer, and David D Cox. Using human brain activity to guide machine
531	learning. Scientific reports, 8(1):5397, 2018.
532	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
533	nition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp.
534	770–778, 2016.
535	V is in H. V. L'Olas Gills V's Value I' D's D 11/2 (D. C. L') Martin
536	Kaiming He, Ainlei Chen, Saining Aie, Yanghao Li, Piotr Dollar, and Ross Girshick. Masked au-
537	vision and pattern recognition pp 16000-16000 2022
538	vision and pattern recognition, pp. 10000-10009, 2022.
539	Kendrick N Kay, Thomas Naselaris, Ryan J Prenger, and Jack L Gallant. Identifying natural images from human brain activity. <i>Nature</i> , 452(7185):352–355, 2008.

578

579

581

588

589

590

540	Meenakshi Khosla, Gia Ngo, Keith Jamison, Amy Kuceyeski, and Mert Sabuncu. Neural encoding
541	with visual attention. Advances in Neural Information Processing Systems, 33:15942-15953,
542	2020.

- Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Matt 544 Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, et al. Ai2-thor: An interactive 3d environment for visual ai. arXiv preprint arXiv:1712.05474, 2017. 546
- 547 Jonas Kubilius, Martin Schrimpf, Kohitij Kar, Rishi Rajalingham, Ha Hong, Najib Majaj, Elias Issa, Pouya Bashivan, Jonathan Prescott-Roy, Kailyn Schmidt, et al. Brain-like object recognition with 548 high-performing shallow recurrent anns. Advances in neural information processing systems, 32, 549 2019. 550
- 551 Eun Sun Lee, Junho Kim, SangWon Park, and Young Min Kim. Moda: Map style transfer for self-552 supervised domain adaptation of embodied agents. In European Conference on Computer Vision, 553 pp. 338-354. Springer, 2022. 554
- Zhe Li, Wieland Brendel, Edgar Walker, Erick Cobos, Taliah Muhammad, Jacob Reimer, Matthias 555 Bethge, Fabian Sinz, Zachary Pitkow, and Andreas Tolias. Learning from brains how to regularize 556 machines. Advances in neural information processing systems, 32, 2019.
- 558 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr 559 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, 560 Proceedings, Part V 13, pp. 740–755. Springer, 2014. 561
- Dongjun Liu, Weichen Dai, Hangkui Zhang, Xuanyu Jin, Jianting Cao, and Wanzeng Kong. Brain-563 machine coupled learning method for facial emotion recognition. IEEE Transactions on Pattern 564 Analysis and Machine Intelligence, 2023. 565
- Yuxin Mao, Jing Zhang, Mochu Xiang, Yiran Zhong, and Yuchao Dai. Multimodal variational 566 auto-encoder based audio-visual segmentation. In Proceedings of the IEEE/CVF International 567 Conference on Computer Vision, pp. 954–965, 2023. 568
- 569 Mark Mon-Williams, James R Tresilian, Niall C Strang, Puja Kochhar, and John P Wann. Improving 570 vision: neural compensation for optical defocus. Proceedings of the Royal Society of London. Series B: Biological Sciences, 265(1390):71-77, 1998. 571
- 572 Thomas Naselaris, Kendrick N Kay, Shinji Nishimoto, and Jack L Gallant. Encoding and decoding 573 in fmri. *Neuroimage*, 56(2):400–410, 2011. 574
- 575 Satoshi Nishida, Yusuke Nakano, Antoine Blanc, Naoya Maeda, Masataka Kado, and Shinji Nishimoto. Brain-mediated transfer learning of convolutional neural networks. In Proceedings of the 576 AAAI Conference on Artificial Intelligence, volume 34, pp. 5281–5288, 2020. 577
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning 580 robust visual features without supervision. arXiv preprint arXiv:2304.07193, 2023.
- Ruslan Partsey, Erik Wijmans, Naoki Yokoyama, Oles Dobosevych, Dhruv Batra, and Oleksandr 582 Maksymets. Is mapping necessary for realistic pointgoal navigation? In Proceedings of the 583 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 17232–17241, 2022. 584
- 585 Maytus Piriyajitakonkij, Mingfei Sun, Mengmi Zhang, and Wei Pan. Tta-nav: Test-time adaptive re-586 construction for point-goal navigation under visual corruptions. arXiv preprint arXiv:2403.01977, 2024.
  - Frano Rajič. Robustness of embodied point navigation agents. In European Conference on Computer Vision, pp. 193–204. Springer, 2022.

Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied 592 ai research. In Proceedings of the IEEE/CVF international conference on computer vision, pp. 9339-9347, 2019.

614

- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Muhammad Shah, Aqsa Kashaf, and Bhiksha Raj. Training on foveated images improves robustness
   to adversarial attacks. *Advances in Neural Information Processing Systems*, 36, 2024.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
   Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.
- Weikang Wan, Haoran Geng, Yun Liu, Zikang Shan, Yaodong Yang, Li Yi, and He Wang. Unidexgrasp++: Improving dexterous grasping policy learning via geometry-aware curriculum and iterative generalist-specialist learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3891–3902, 2023.
- Michael A Webster, Mark A Georgeson, and Shernaaz M Webster. Neural adjustments to image
   blur. *Nature neuroscience*, 5(9):839–840, 2002.
- Haiguang Wen, Junxing Shi, Yizhen Zhang, Kun-Han Lu, Jiayue Cao, and Zhongming Liu. Neural encoding and decoding with deep learning for dynamic natural vision. *Cerebral cortex*, 28(12): 4136–4160, 2018.
- Erik Wijmans, Abhishek Kadian, Ari Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis Savva, and Dhruv Batra. Dd-ppo: Learning near-perfect pointgoal navigators from 2.5 billion frames. In *International Conference on Learning Representations*, 2019.
- Huzheng Yang, James Gee, and Jianbo Shi. Brain decodes deep nets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23030–23040, 2024.
- Jiazhao Zhang, Liu Dai, Fanpeng Meng, Qingnan Fan, Xuelin Chen, Kai Xu, and He Wang. 3d aware object goal navigation via simultaneous exploration and identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6672–6682, 2023.
- Xiaoming Zhao, Harsh Agrawal, Dhruv Batra, and Alexander G Schwing. The surprising effectiveness of visual odometry techniques for embodied pointgoal navigation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 16127–16136, 2021.
- Kiaoying Zhu, Neville A McBrien, Earl L Smith, David Troilo, and Josh Wallman. Eyes in various species can shorten to compensate for myopic defocus. *Investigative ophthalmology & visual science*, 54(4):2634–2644, 2013.

# 648 A RELATED WORK

# A.1 EMBODIED VISUAL NAVIGATION

652 Two widely studied tasks in embodied visual navigation are PointGoal Navigation (PointNav) (An-653 derson et al., 2018) and ObjectGoal Navigation (ObjectNav) (Batra et al., 2020). PointNav involves 654 navigating to a specified goal coordinate in a global reference frame, while ObjectNav requires the 655 agent to find an instance of a specified object. Thanks to high-quality simulators (Kolve et al., 2017; 656 Savva et al., 2019) and datasets (Chang et al., 2017; Deitke et al., 2020), significant advancements has been made in embodied visual navigation (Wijmans et al., 2019; Zhao et al., 2021; Partsey et al., 657 2022; Zhang et al., 2023), especially in the PointNav task, which is often considered "solved" (Wij-658 mans et al., 2019). 659

660 However, existing studies have largely overlooked the robustness of agents, a critical aspect for 661 real-world applications. To address this gap, (Chattopadhyay et al., 2021) proposed ROBUSTNAV, 662 a framework for analyzing the robustness of navigation agents. ROBUSTNAV quantifies the per-663 formance of embodied navigation agents under various common visual and dynamic corruptions. Extensive experiments demonstrated that navigation agents trained in simulation often exhibit sig-664 nificant performance drops when evaluated in corrupted environments. Furthermore, (Rajič, 2022) 665 conducted a robustness analysis of two successful agents from the 2021 Habitat Challenge, reveal-666 ing varying degrees of performance deterioration in corrupted settings. (Lee et al., 2022) proposed a 667 self-supervised domain adaptation method with map style transfer to boost agent robustness against 668 visual and dynamic perturbations. (Piriyajitakonkij et al., 2024) introduced TTA-Nav, which en-669 hances navigation performance across visual corruptions using a top-down decoder. Our approach 670 differs from theirs in that we exploits the capabilities of the human brain visual system by introduc-671 ing brain-like representations for agents to enhance robustness.

672 673

674

650

651

# A.2 DEEP LEARNING AND BRAIN ACTIVITY INTEGRATION

675 Recent studies have explored techniques to enhance computer vision models by integrating neural 676 network features with human brain activity. (Fong et al., 2018) improved image classification in 677 Convolutional Neural Networks (CNNs) by incorporating voxel responses from fMRI during train-678 ing, making the decision surface of the classifier more consistent with brain representations. (Li 679 et al., 2019) regularized CNNs using cortical representations from neuroscience data, enhancing 680 robustness against adversarial attacks. (Nishida et al., 2020) introduced a brain-mediated transfer learning (TL) method, transforming the feature representation of audiovisual input in CNNs into 681 brain representations, achieving superior performance in estimating human cognitive and behavioral 682 labels compared to standard TL. Additionally, (Dapello et al., 2020) developed VOneNet to simu-683 late the primary visual cortex, thereby enhancing CNN classification robustness. (Fel et al., 2022) 684 presented a neural harmonizer that aligns deep neural networks with human visual strategies, result-685 ing in improved classification accuracy. More recently, (Liu et al., 2023) proposed a brain-machine 686 coupled learning method that utilized visual images and electroencephalogram (EEG) signals for 687 training models in facial emotion recognition, demonstrating improved generalization. (Shah et al., 688 2024) developed an image transform that simulates peripheral vision, boosting DNN robustness 689 against adversarial attacks. These studies indicate that combining neural network features with 690 human brain activity can yield improvements. However, all of these works focus on image classifi-691 cation tasks, while our approach tackles the more complex domain of embodied visual navigation, a topic that has not been extensively explored. Furthermore, we leverage a much larger and higher-692 quality fMRI dataset. 693

694 695

696

#### A.3 BRAIN ENCODING MODEL

Brain encoding models are capable of predicting fMRI data from humans viewing visual stimuli,
which is crucial for understanding how information is represented in the brain (Naselaris et al.,
2011; Wen et al., 2018). (Kay et al., 2008) developed a linear encoding model based on a Gabor
wavelet pyramid to predict brain responses to stimulus images. (Khosla et al., 2020) proposed
an encoding model that incorporates visual attention, resulting in significant improvements in predicting neural responses. More recent studies (Adeli et al., 2023; Yang et al., 2024) have utilized

Table 4: Noise-normalized prediction accuracy of the brain encoder for subject 02.

Hemisphere			ROI				
	early	midventral	midlateral	midparietal	ventral	lateral	parietal
Left Hemisphere	0.7356	0.8213	0.8692	0.6995	0.9106	0.8831	0.7074
Right Hemishpere	0.7491	0.8166	0.8250	0.7682	0.8588	0.9891	0.7396

Table 5: Comparison with standard and computer vision-enhanced navigation agents for subject 02. We compare BraiNav with six approaches proposed in (Chattopadhyay et al., 2021). The results are highlighted with best and second.

Approach							Visual C	orruption								
	Cl	ean	Spa	itter	Speckl	e Noise	Camer	a Crack	Lowe	r FOV	Defoc	us Blur	Motic	on Blur	Low Lighting	
	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL
Standard Agent	98.82	83.13	33.58	24.72	67.42	48.57	82.07	63.83	42.49	31.73	75.89	53.55	95.72	73.37	94.36	75.15
Standard Agent +AP	98.45	83.28	20.38	15.70	65.61	47.03	72.70	56.82	45.68	35.14	83.35	61.51	94.81	74.30	92.17	76.11
Standard Agent +AP+SS-Adapt	37.31	31.03	14.19	10.29	\	\	57.87	46.72	32.94	26.09	40.95	33.35	\	\	\	\
Standard Agent +RP	98.73	82.53	23.48	18.63	78.98	55.92	67.06	53.70	44.95	32.74	32.21	22.47	91.63	65.27	89.81	67.38
Standard Agent +RP+SS-Adapt	94.63	77.25	61.06	47.16	\	\	60.42	49.37	50.59	36.10	79.16	62.74	\	\	\	\
Standard Agent +Data Aug	98.45	81.08	23.93	18.41	77.25	57.95	88.44	71.57	71.70	54.54	81.26	61.32	96.91	75.97	97.27	78.74
BraiNav	99.36	84.14	39.03	28.36	80.78	60.92	80.98	62.18	48.86	37.26	85.71	63.70	94.63	73.95	94.38	75.47

self-supervised pretrained Visual Transformer (ViT) (Dosovitskiy et al., 2020) models as backbones to extract features from stimulus images, achieving excellent performance.

# **B** MORE EXPERIMENTS

# B.1 DIFFERENT SUBJECT

We also train the brain encoder using the fMRI data from subj 02 and utilized the pretrained frozen brain encoder for downstream navigation tasks. We compute the noise ceiling for each vertex using response data from the subjects' three trials and then average these values across each ROI. Details of the calculation methodology and results are provided in Appendix B.2. Next, we evaluate the brain encoder's performance by calculating the *Pearson correlation coefficient* between its predicted fMRI responses and the actual values. To determine the noise-normalized prediction accuracy, we divide the encoder's prediction accuracy by the corresponding noise ceiling, as summarized in Table 4. The results indicate that the brain encoder effectively predicts the fMRI responses to visual stimuli. High prediction accuracy ensures that the neural representations used for navigation policy learning incorporate relevant cognitive processing information from the human brain.

Next, we compare the proposed method with standard and computer vision-enhanced agents, as shown in Table 5. BraiNav outperforms the standard navigation agent across clean and six vi-sual corruptions, achieving varying degrees of performance improvements. Notably, for Speckle Noise and Defocus Blur, BraiNav achieves significant improvements, with absolute improvements of (12.35%SPL, 13.36%SR) and (10.15%SPL, 9.82%SR), respectively. For other corruptions, the ab-solute improvements are (5.53%SPL, 6.37%SR) for Lower FOV, (3.64%SPL, 5.45%SR) for Spatter, (0.58%SPL) for Motion Blur, and (0.32%SPL, 0.02%SR) for Low Lighting. Additionally, BraiNav outperforms all computer vision-enhanced agents on clean and 2 visual corruptions. In detail, the absolute improvements are (2.97%SPL, 1.8%SR) for Speckle Noise and (0.96%SPL, 2.36%SR) for Defocus Blur. Additionally, BraiNav achieves the second best results on Spatter, Lower FOV, and Low Lighting, demonstrating competitive performance on the remaining visual corruptions.

		Table 6: Not	ise ceiling fo	r subject 01.			
Hemisphere			ROI				
	early	midventral	midlateral	midparietal	ventral	lateral	parietal
Left Hemisphere	0.5965	0.5967	0.5971	0.5942	0.5857	0.5881	0.5870
Right Hemishpere	0.5945	0.5967	0.5888	0.5910	0.5863	0.5870	0.5873

#### Table 7: Noise ceiling for subject 02.

Hemisphere			ROI				
	early	midventral	midlateral	midparietal	ventral	lateral	parietal
Left Hemisphere	0.6006	0.6056	0.5992	0.6047	0.5905	0.5921	0.5980
Right Hemishpere	0.5991	0.5983	0.5950	0.6013	0.5908	0.5941	0.5984

#### **B.2** Noise ceiling Across Different Subjects

For each stimulus image, subjects view it three times, resulting in 3-trial fMRI responses. For stimulus image i, let the three response values for vertex j be denoted as rep $0_i^i$ , rep $1_i^i$ , rep $2_i^i$ . Across N stimulus images, vertex j accumulates three sets of trial responses:  $\operatorname{rep}_{i}^{j} = [\operatorname{rep}_{i}^{1}, ..., \operatorname{rep}_{i}^{N}]$  $\operatorname{rep1}_{j} = [\operatorname{rep1}_{j}^{1}, ..., \operatorname{rep0}_{j}^{N}], \operatorname{rep2}_{j} = [\operatorname{rep2}_{j}^{1}, ..., \operatorname{rep2}_{j}^{N}].$  The average of these 3-trial responses is denoted as repm<sub>*i*</sub>. The noise ceiling for vertex j is then calculated as:

$$nc_j = \frac{\text{pearsonr}\left(\text{rep0}_j, \text{repm}_j\right) + \text{pearsonr}\left(\text{rep1}_j, \text{repm}_j\right) + \text{pearsonr}\left(\text{rep2}_j, \text{repm}_j\right)}{3}, \quad (11)$$

where pearsonr represents the Pearson correlation coefficient operator. After computing the noise ceilings for all vertices, we average them by ROI. The results are presented in Table 6 for subject 01 and Table 7 for subject 02.

#### **B.3** MORE COMPUTER VISION-ENHANCED AGENT

We develop a computer vision-enhanced agent leveraging the advanced self-supervised model, Masked Autoencoder (MAE) (He et al., 2022). The agent is first deployed in RoboTHOR indoor scenes to freely explore and collect 60,000 egocentric images as training data for MAE. These images are then masked and reconstructed for self-supervised pretraining. Specifically, the MAE encoder is implemented as a 12-layer Vision Transformer (ViT), while the decoder consists of an 8-layer ViT, with a masking ratio of 0.75. After pretraining, the encoder is retained as the agent's vi-sual backbone, and the decoder is discarded. We compare the performance of our proposed BraiNav agent against the MAE-based agent under seven visual corruption scenarios. As shown in Table 8, BraiNav consistently outperforms the MAE-based agent across all corruption types. 

#### **B.4** Comparison with Agent based on Brain-like Representations

We have replaced the brain encoder with CORnet-S (Kubilius et al., 2019), a compact, recurrent artificial neural network model that aligns closely with the anatomical structure and dynamic responses of the primate ventral visual stream. CORnet-S not only achieves high biological fidelity

				Tal	ble 8:	Con	nparis	on wi	th MA	AE ag	ent.					
Approach						V	Visual Co	orruptio	n							
	Cl	ean	Spa	atter	Speckl	e Noise	Camer	a Crack	Lowe	r FOV	Defoc	ıs Blur	Motio	n Blur	Low L	ighting
	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL
MAE Agent	98.54	83.31	8.55	5.92	10.10	7.66	48.95	37.67	31.57	23.09	67.42	51.95	79.34	60.14	36.12	28.27
BraiNav	97.73	80.88	36.76	27.48	85.17	64.34	77.80	60.25	42.95	32.56	87.26	65.44	96.90	78.30	95.09	75.95

				Table	e 9: C	Compa	arison	with	COR	net ag	gent.						
Approach						١	isual C	orruptio	n								
	Cle	Clean Spatter Speckle Noise Camera Crack Lower FOV Defocus Blur Motion Blur												n Blur	Low Lighting		
	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	
CORnet Agent	97.27	81.63	30.30	23.28	78.07	61.53	77.80	61.45	61.60	48.17	75.98	57.16	94.18	73.74	87.99	70.08	
BraiNav	97.73	80.88	36.76	27.48	85.17	64.34	77.80	60.25	42.95	32.56	87.26	65.44	96.90	78.30	95.09	75.95	

#### Table 10: Impact of different ROIs.

Approach						١	isual Co	orruptio	n							
	Cle	ean	Spa	tter	Speckle Noise Camera Crack					r FOV	Defocus Blur		Motion Blur		Low Lighting	
	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL	SR	SPL
Low-Level ROI	96.54	79.71	34.03	25.17	78.14	59.57	77.43	61.09	44.59	34.35	73.97	55.30	90.72	72.45	94.44	72.73
High-Level ROI	97.27	80.53	37.03	27.61	87.52	65.81	79.25	63.05	44.31	34.08	86.90	66.03	96.45	78.52	96.81	75.82

but also excels in both neuroscience and machine learning benchmarks. CORnet-S extracts visual features from images with a 1000-dimensional brain-like representations, followed by a linear layer converted to 768 dimensions. The results are presented in Table 9. BraiNav outperforms the CORnet Agent on Spatter, Speckle Noise, Defocus Blur, Motion Blur, and Low Lighting. Conversely, the CORnet Agent demonstrates strong performance on Camera Crack and Lower FOV. These results highlight that brain-like representations can enhance the robustness of embodied navigation agent.

#### **B.5** IMPACT OF DIFFERENT ROIS

Early brain regions are primarily responsible for processing low-level visual information, while
ventral, lateral, and parietal regions handle high-level visual information (Allen et al., 2022). To
investigate the contribution of features from different brain regions, we categorize the brain into
two groups: low-level regions (early) and high-level regions (ventral, lateral, and parietal). Neural
representations for visual navigation are then constructed using output tokens from each group. The
results are presented in Table 10.

Except for the Lower FOV corruption, representations derived from high-level brain regions significantly outperform those from low-level regions across all other corruptions, demonstrating strong
consistency. These findings further confirm that the brain encoder effectively generates neural representations containing cognitive information from the human brain. Moreover, representations from
higher brain regions exhibit greater robustness.