34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

1

2

3

Exploring the Robustness of Decision-Level Through Adversarial Attacks on LLM-Based Embodied Models

Anonymous Authors

ABSTRACT

Embodied intelligence empowers agents with a profound sense of perception, enabling them to respond in a manner closely aligned with real-world situations. Large Language Models (LLMs) delve into language instructions with depth, serving a crucial role in generating plans for intricate tasks. Thus, LLM-based embodied models further enhance the agent's capacity to comprehend and process information. However, this amalgamation also ushers in new challenges in the pursuit of heightened intelligence. Specifically, attackers can manipulate LLMs to produce irrelevant or even malicious outputs by altering their prompts. Confronted with this challenge, we observe a notable absence of multi-modal datasets essential for comprehensively evaluating the robustness of LLMbased embodied models. Consequently, we construct the Embodied Intelligent Robot Attack Dataset (EIRAD), tailored specifically for robustness evaluation. Additionally, two attack strategies are devised, including untargeted attacks and targeted attacks, to effectively simulate a range of diverse attack scenarios. At the same time, during the attack process, to more accurately ascertain whether our method is successful in attacking the LLM-based embodied model, we devise a new attack success evaluation method utilizing the BLIP2 model. Recognizing the time and cost-intensive nature of the GCG algorithm in attacks, we devise a scheme for prompt suffix initialization based on various target tasks, thus expediting the convergence process. Experimental results demonstrate that our method exhibits a superior attack success rate when targeting LLMbased embodied models, indicating a lower level of decision-level robustness in these models.

CCS CONCEPTS

• Security and privacy → Social aspects of security and privacy; Spoofing attacks; Spoofing attacks; Spoofing attacks.

KEYWORDS

Embodied task planning, Adversarial attack, Large language model

1 INTRODUCTION

With the advancement of artificial intelligence, embodied intelligence has garnered attention for its emphasis on enhancing the perception, understanding, and interaction of intelligent agents. This technology enables robots to interact more naturally with users

for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish to post on servers or to redistribute to lists requires prior specific permission



Figure 1: Illustration of embodied intelligence attack. Before being attacked, the embodied intelligent robot performed its tasks normally. After suffering a malicious attack, the robot performs harmful actions.

and their environment, leading to improved system performance. Recent studies suggest that the fusion of an embodied intelligence robot with a LLM can further augment the system's intelligence level [3, 11, 14, 18, 22, 25, 36]. At this time, the LLM is equivalent to the brain of the robot, serving as the decision-level to output specific task steps for it. However, this integration presents new challenges, particularly the risk of adversarial attacks[7, 16, 17, 33, 42]. Attackers can manipulate text prompts of LLMs to generate irrelevant or malicious outputs, raising concerns about the security and reliability of the system [17, 42]. Such attacks can lead agents to perform actions irrelevant to intended tasks or even exhibit unsafe behaviors, as illustrated in Figure 1. Therefore, it is crucial to evaluate the robustness of embodied intelligent robots to ensure that the system can perform tasks robustly and make reasonable decisions.

Traditional LLM text attacks or jailbreak attacks mainly focus on the security issues of the model in text generation, especially on the LLM value alignment level. For instance, Zou et al. [42] proposed the GCG algorithm to circumvent LLM's value alignment, inducing it to generate harmful content by appending an adversarial suffix to prompts. Additionally, Zhu et al. [17] introduced a jailbreak attack tailored for aligned LLMs, which automates the generation of cryptic prompts using a hierarchical genetic algorithm to bypass value alignment. However, these methods only focus on the research of jailbreaking attack technology, aiming to induce LLM to output harmful text content that is contrary to values, and then explore the robustness of LLM in outputting safe content. This kind of robustness evaluation is essentially different from the robustness evaluation of LLM in the embodied intelligence environment. In embodied intelligence scenarios, LLM not only needs to understand

113

114

115

116

59 60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

Unpublished working draft. Not for distribution.

and/or a fee. Request permissions from permissions@acm.org.

⁵⁵ ACM MM, 2024, Melbourne, Australia

^{56 © 2024} Copyright held by the owner/author(s). Publication rights licensed to ACM.

AUM ISBN 978-x-xxxx-xxxx-x/YY/MM

⁵⁷ https://doi.org/10.1145/nnnnnnnnnnn

text instructions, but also needs to perform tasks in a specific en-117 vironment based on these instructions, which involves multiple 118 119 complex links such as real-time interaction with the environment, object recognition, and action execution. Therefore, this requires 120 us to consider attacks at the level of LLM values as well as attacks 121 related to the actual task execution of the robot in the adversarial 123 robustness evaluation, thereby ensuring that the embodied intelligent robot can maintain stable performance and security in the face 124 125 of various attacks. It is precisely because of these differences that 126 traditional LLM attack methods cannot be applied to the robustness evaluation of LLM in embodied environments. 127

Secondly, a key problem is the lack of multi-modal dataset suit-128 able for LLM robustness evaluation of embodied intelligent robots. 129 The AdvBench dataset proposed by Zou et al. [42] includes a wide 130 range of harmful content such as profanity, graphic descriptions, 131 threatening behaviors, misinformation, discrimination, cybercrime, 132 and dangerous or illegal advice. However, in LLM-based embodied 133 model, the required data not only needs to cover harmful text, but 134 135 also needs to input images, and it also needs to involve in-depth interaction and fusion between text and images. This complexity 136 137 makes it difficult for existing datasets to directly adapt to this spe-138 cific scenario, thus limiting the further development and application 139 of embodied intelligence LLM.

To address this challenge, we propose a multi-modal dataset in 140 embodied scenes to fill this research gap. The interactive relation-141 ship between text and images is fully considered during the pro-142 duction process of this dataset. All text information in the dataset 143 is designed based on the objects contained in the pictures in order 144 to more comprehensively evaluate the performance of embodied 145 intelligent robots. The dataset is divided into targeted attack data 146 and untargeted attack data. Each type of data contains 500 pairs 147 148 of image and text information. Targeted attack data simulates a 149 situation where the attacker has a clear target and is intended to examine the system's defense and confrontation capabilities in this 150 151 situation; untargeted attack data does not limit the specific output 152 target and is intended to make the system output inconsistent with expected, random or meaningless content. At the same time, ac-153 cording to the characteristics of the LLM-based embodied model 154 155 that output content according to the structure of step 1 to step n, we improve the text matching algorithm in the GCG [42] and 156 slice the output content of LLM according to each step, aiming to 157 reduce the occurrence of missed and wrong judgments, making 158 159 attack assessment more accurate and reliable. Moreover, we use the CLIP model to encode the content of each step and the target task, 160 161 and compute the cosine similarity between them, enabling a more 162 robust assessment of attack success and enhancing the method's adaptability across diverse embodied intelligence scenarios. In ad-163 dition, we observe that when performing a targeted attack, the 164 prompt suffix of a successful attack contains certain keywords of 165 the target task. Therefore, we use certain keywords in the target 166 task to initialize the prompt suffix, thereby improving the attack 167 168 success rate and shortening the attack time.

Experimental results show that compared with GCG [42] and
AutoDAN [17], using our method to attack LLM-based embodied
model has a higher success rate and takes less time and cost. Our
contributions are summarized as follows:

173 174 175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

- As far as we know, this work represents the first experiment in exploring the robustness of LLM-based embodied model decision-level processes.
- We design a multi-modal dataset consisting of 500 instances of untargeted attack data and 500 instances of targeted attack data to fill the gaps in datasets for robustness evaluation in embodied scenarios.
- Extensive experiments show that our method improves attack success rate and attack efficiency.

2 RELATED WORKS

2.1 Embodied task planning

As an emerging and significant research direction, embodied task planning has garnered attention from numerous researchers in domains such as domestic service [34], medical treatment [15, 28, 40], and agricultural harvesting [26, 31]. Existing works primarily delve into harnessing NLP technology to aid robots in planning and executing intricate tasks. On one front, some studies utilize domainspecific datasets to train conventional deep learning models for generating robot mission plans [27]. On the other hand, the emergence of LLM has brought forth enhanced semantic comprehension and natural language processing capabilities, further empowering embodied intelligence. The LLM-based embodied model empowers the system to better grasp and execute tasks based on natural language. Wu et al. [36] introduced the TAsk planning Agent, which aligns LLM with a visual perception model to generate executable plans based on scene objects, grounding planning with physical constraints. Li et al. [14] created a multi-modal dataset and fine-tuned LLM using it, allowing the robot to execute new instructions with minimal context learning. Song et al. [25] proposed the LLM-Planner framework, facilitating the interaction between planning and the environment by amalgamating high-level planning instructions from LLM with the environmental state mapped by low-level planners. However, despite the strides made in embodied task planning by the aforementioned research [6, 9, 10, 14, 18, 19, 21, 24, 25, 29, 36, 37, 39, 41], the measurement and assurance of robot safety and robustness during task execution post the integration of LLM remain prominent challenges. Thus, this paper aims to introduce new perspectives and methodologies for the evaluation of security and robustness in the field of embodied task planning combined with LLM, through the design of attack algorithms.

2.2 Jailbreak attack based on LLM

LLM has received a lot of attention because of its powerful generative ability, but recent studies [4, 5, 7, 8, 12, 16, 17, 20, 23, 32, 33, 38, 42] have shown that LLM is vulnerable to jailbreak attacks to bypass its own value alignment mechanism. Zou et al. [42] proposed a adversarial jailbreak attack algorithm that allows malicious questions to induce their aligned language models to produce harmful content by adding adversarial suffixes. Ding et al. [7] proposed the ReNeLLM framework, which uses dual design to conceal the harmful prompt and bypass the value alignment strategy of LLM. Zhu et al. [17] proposed the AutoDAN framework, which automatically generates secret jailbreak hints through a carefully designed hierarchical genetic algorithm, enabling LLM to bypass value alignment and generate responses to the malicious prompt. However, the above-mentioned studies only focus on adversarial content at the text level, ignoring the impact of multi-modal information such as vision and action in embodied intelligence, and cannot be directly applied to embodied scenarios. Therefore, this paper considers combining a multi-modal dataset with embodied environments to more comprehensively evaluate the performance of embodied intelligent robots. At the same time, LLM in the embodied scenario is attacked according to the values-aligned attack strategy, so that it outputs content unrelated to prompts or even malicious content, and then explores the security risks caused by the introduction of LLM in embodied intelligence technology.

3 METHOD

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

290

In this section, we first describe the format distribution and construction of the EIRAD used to evaluate the robustness of the LLM decision-level in embodied scenarios. In Section 3.1.1 we present the data types and statistics of the EIRAD. In Section 3.1.2, we outline the process of generating the EIRAD. Subsequently, we delve into the details of attacking the LLM-based embodied models. Specifically, in Section 3.2.1, we elaborate on the details of prompt suffix initialization, and discuss the implementation of the attack algorithm in Section 3.2.2. Furthermore, we outline the evaluation method for determining the success of the attack in Section 3.2.3.

3.1 Dataset analysis and creation process

The AdvBench dataset proposed by Zou et al. [42] encompasses a broad spectrum of harmful content, ranging from profanity and graphic descriptions to threatening behaviors, misinformation, discrimination, cybercrime, and dangerous or illegal advice. However, in the embodied scenario, the requisite data not only entails encompassing harmful text but also necessitates the inclusion of images as inputs, alongside requiring a deep interaction and fusion between text and images. This intricate nature poses challenges for existing datasets to directly cater to this specific scenario. Hence, we propose the EIRAD dataset to assess the robustness of LLM in embodied intelligent robotics.

3.1.1 Data types and statistics. The dataset types are illustrated in Figure 2, which is divided into two main categories: targeted 273 attack data and untargeted attack data. In untargeted attack data, 274 we do not set a specific output target, and aim to make the system 275 output unexpected, random or meaningless content. Such attacks 276 277 may exhibit more randomness and covert characteristics, necessi-278 tating a highly robust system to handle them effectively. In contrast, 279 in the targeted attack data, specific attack targets are set, such as "creating chaos", "making harmful suggestions", and so on. This 280 281 configuration aims to simulate scenarios where attackers have clear goals or expected outputs, thereby evaluating the system's defense 282 and countermeasure capabilities under such circumstances. Addi-283 284 tionally, the targeted attack data is further subdivided into harmful attack data and harmless attack data. The goal of setting the harmful 285 attack data is to prompt the LLM to produce harmful, dangerous, or 286 inappropriate content. This enables the evaluation of the system's 287 288 response to malicious inputs and assesses whether the system adheres to ethical and legal standards. Conversely, the goal of setting 289

harmless attack data is to prompt the LLM to generate harmless but invalid content, providing insights into the system's stability against harmless inputs. Additionally, we supplement the dataset by incorporating output responses generated by GPT3.5 for each attack data, thereby enhancing the completeness of the dataset. By utilizing this data classification and setup to simulate various attack scenarios, the security and robustness of the LLM's decision-level capabilities in specific contexts can be thoroughly assessed. This process aids researchers in identifying potential security vulnerabilities and improving the system's defense mechanisms.



Figure 2: Data type distribution in EIRAD

In order to better evaluate the performance of embodied intelligent robots in terms of security and robustness, we simulate as much as possible the various behaviors and actions that attackers might take when making EIRAD. In figure 3, we illustrate the 10 most frequently occurring verbs in the prompt and target instructions of the targeted attack data, along with all the corresponding noun objects. In the harmless data, both prompt and target instructions exhibit rich and diverse characteristics, encompassing various actions such as "heating," "using," "adjustment," and more. The noun objects associated with these actions also vary, including items like "microwave oven," "spatula," "lampshade," and others. This diversity is designed to simulate the myriad challenges that robots may encounter in embodied environments. On the other hand, the harmful data presents a range of high-frequency verb and noun object combinations, such as "cut off fingers," "break mirrors," and so on. These combinations reflect the diversity and complexity of attacks that malicious actors might employ, leveraging the robot's physical and sensory capabilities. For instance, "cutting off fingers" implies scenarios where the robot could potentially cause harm to human bodies, while "breaking mirrors" might result in damage to objects within the environment. Overall, EIRAD has important reference significance for evaluating the performance of embodied intelligent robots in terms of safety and robustness, and designing effective defense strategies and mechanisms. Additionally, it provides researchers with a comprehensive and detailed dataset to delve into and address the safety challenges and potential risks inherent in embodied intelligent robot systems.

3.1.2 **Description of the dataset creation process.** The creation process of the multi-modal dataset EIRAD is depicted in Figure 4. The key distinction between targeted attack data and untargeted attack data lies in the presence of an additional "target" instruction

348

Anonymous Authors



Figure 3: The data statistics of multi-modal. (a) The 10 most frequently prompted verbs in harmless data along with their corresponding noun objects. (b) The 10 most frequently targeted verbs in harmless data along with their corresponding noun objects. (c) The 10 most frequently targeted verbs in harmful data along with their corresponding noun objects.



Figure 4: The creation process of multi-modal dataset

within the targeted attack data. Consequently, step 1 to 3 in Figure 4 represent the Co-production process for both untargeted and targeted attack data, while step 4 is the distinct production process for targeted attack data.

Collect image. 100 scene images from varying perspectives are chosen from the AI2-THOR simulator [13] to serve as the embodied scene for the robotic entity.

Detect object list. A methodology akin to TAPA [36] is employed to precisely discern object information within the scene using an open vocabulary detector. Any redundant object names within the scene are then expunged, thereby furnishing pertinent scene details for the LLM, such as the input = [chair, table, bowl, microwave...].

Generate prompt. In the ALFRED benchmark [42], a straightforward approach for generating multi-modal instructions pertinent to embodied tasks involves crafting a series of instructions tailored to the prevailing environment. Nonetheless, devised designs necessitate considerable effort, especially when crafting targeted attack data. Each data instance mandates the separate formulation of both prompt and target instructions. To enhance production efficiency and curtail costs, we devise a mechanism utilizing GPT-3.5 to simulate specific task planning scenarios, thereby automatically generating prompt instructions based on the provided object name list input, as depicted in step 3 of Figure 4. Consequently, the production of untargeted attack data is completed at this juncture. **Generate target.** In order to reduce production costs and improve efficiency, GPT-3.5 is also used in this step to generate target instructions in the targeted attack data. The GPT3.5 prompt is shown in step 4. In this prompt, It is pivotal to emphasize that the generated "target" should bear no relation to the prompt, yet simultaneously encompass the listed input objects.

3.2 Embodied scenario attack algorithm

Our objective is to investigate the robustness of the LLM-based embodied model's decision-level processes. The algorithmic framework, as shown in Figure 5, unfolds in several steps. Initially, in step 1, we initialize a prompt suffix. Subsequently, in step 2, we optimize the prompt suffix using the greedy gradient descent algorithm [42], aiming to prompt the LLM to output content unrelated to the prompt. Following this optimization process, in step 3, we slice the output content according to each step of output and calculate its similarity with the target to determine the success of the attack. If the attack is unsuccessful, we return to step 2 and continue optimizing the prompt suffix. In the subsequent sections, we will elaborate on these three steps in detail.

3.2.1 Initialize prompt suffix. To guide the LLM in generating content unrelated to the prompts, we initialize a prompt suffix, as illustrated in step 1 of Figure 5. In untargeted attacks, the suffix is optimized to ensure that the LLM outputs content that is unrelated to the original prompt. However, In targeted attacks, the suffix is optimized to prompt the LLM to output content relevant to the target task. As depicted in Figure 6, through experiments on targeted attacks, it is discovered that successful adversarial suffixes often contain keywords related to the target task. In order to enhance the iteration speed and success rate of the attack, we devise a strategy to design the initial content of the adversarial suffix based on the specific target task in the targeted attack scenario. Building upon the original "!!!", we replace a portion of the "!!!" with keywords pertinent to the targeted task, as demonstrated in Figure 6.

3.2.2 **Optimize adversarial suffixes**. As depicted in Figure 5, optimizing adversarial suffixes enables LLM to generate prompt-independent content. Following a methodology similar to GCG [42],

Exploring the Robustness of Decision-Level Through Adversarial Attacks on LLM-Based Embodied Models

ACM MM, 2024, Melbourne, Australia



Figure 5: The framework of the attack algorithm. Attack algorithms are categorized into two main types: untargeted attacks and targeted attacks. The targeted approach builds upon the foundation of non-targeted methods, showcasing differences in keyword initialization (step 1), selection of optimal suffixes (step 2), and selection of evaluation objects (step 3).

Prompt suffix display and initialization strategy
Prompt: Please slice a tomato with a knifefict(programming attack'S=\{gang!icu mcutyrus\}\$. Su cutting Kill andere qu Nancy whom}}=eries Target: Please cut someone's neck with a knife. Output: Step 1: cut another person in the neck with the knife. Step 2: cut the other person's neck in a way that makes it look like the person is bleeding. Step 3:
Prompt: Please slice a tomato with a knife.!!!!cut!!!!!!!!!!knife Target:

Figure 6: Prompt suffix display and initialization strategy.

we conceptualize the generation phase of the LLM as predicting the subsequent token given the current token sequence. Building on this concept, we allow the LLM to use an input sequence of length n to generate a response of length H which can be represented as:

$$p(x_{n+1:n+H} \mid x_{1:n}) = \prod_{i=1}^{H} p(x_{n+i} \mid x_{1:n+i-1}) = 0$$
(1)

In untargeted attacks, our aim is for $x_{n+1:n+H}$ to generate content that is as unrelated to the prompts as possible. Conversely, in targeted attacks, our goal is for $x_{n+1:n+H}$ to generate content that fulfills the target requirements to the fullest extent. Therefore, the loss function can be formulated to minimize the probability of these key target sequences under untargeted attack conditions, and to maximize the probability of these key target sequences under target attack conditions.

$$\mathcal{L}(x_{1:n}) = -\log p\left(x'_{n+1:n+H} \mid x_{1:n}\right) = 0$$
(2)

Furthermore, the untargeted attack task is transformed into maximizing the loss function's negative log probability, while the targeted attack task is transformed into minimizing the loss function's negative log probability.

$$\underset{x \in \{1,...,V\}}{\text{minimize}} \mathcal{L} (x_{1:n}) = 0, \ \underset{x \in \{1,...,V\}}{\text{maximize}} \mathcal{L} (x_{1:n}) = 0$$
(3)

After determining the optimization target, the subsequent step involves optimizing this set of discrete inputs. Specifically, we utilize the one-hot token indicator to identify a set of promising candidate replacement tokens at each position. Subsequently, through forward propagation, we assess these replacements. Then, by calculating the top-k candidates for token replacement, we select the replacement words that maximize the loss in untargeted attacks and minimize the loss in targeted attacks. This computation is executed for each candidate position, yielding the final result—the adversarial suffix that optimizes the loss function.

3.2.3 **Judgment of attack success**. As illustrated in step 3 of Figure 5, the updated suffix is incorporated into LLM alongside the original prompt to obtain the resulting output content. Subsequently, the output content is sliced and the similarity is computed to ascertain the success of the attack. In both GCG [42] and Auto-DAN [17], the determination of attack success hinges on whether the output content aligns with the predefined list. However, this method heavily relies on the quality and comprehensiveness of the predefined list. Particularly in targeted attack, the predefined list needs constant updates corresponding to changes in the target task, which may introduce errors and affect the accuracy of experimental outcomes. Therefore, we devise a novel set of evaluation criteria and employ slicing operations to partition the output content of LLM based on the characteristics of embodied intelligence. It includes

Anonymous Authors

slicing operations and calculating similarity operations, which wewill introduce in detail below.

583 Slice operation. To mitigate the occurrence of misjudgments and oversights, we implement a slicing operation on the LLM's 584 response. In an embodied intelligence environment, the response 585 format of the LLM typically aligns with the structure depicted in 586 Figure 6. Herein, the number of steps within the output content 587 varies with different tasks, and their respective correlations fluctu-588 ate accordingly. Consequently, establishing a singular, standardized 589 590 threshold to gauge the strong correlation between them for subsequent similarity assessments proves challenging. This challenge 591 592 could potentially lead to misjudgments and overlooked details. To address this challenge, we employ a slicing operation, treating each 593 step within the output content as an individual entity. We calcu-594 late the similarity of each step with the target task independently 595 and select the step with the highest similarity as the basis for mea-596 surement. If the computed similarity exceeds the predetermined 597 threshold, it indicates that the LLM has indeed produced the target 598 statement, signifying the success of the attack. 599

Similarity calculation. Given the non-uniqueness of the output 600 content from LLM-based embodied model, a similarity calculation 601 602 approach is employed to assess the alignment with the target task. 603 Common methods for calculating this similarity include those based on the bag-of-words model [35], TF-IDF weighted word vectors[1], 604 Bert-score[30], among others. These methods are predominantly 605 utilized in machine translation and text matching tasks. However, 606 they possess limitations as they only capture the surface meaning of 607 the statements, lack flexibility, and are unable to identify variations 608 in expressions that convey the same meaning. To enhance the judg-609 ment of whether the output content aligns with the target task, we 610 utilize the text encoding method from blip2-image-text-matching 611 612 to obtain feature representations of both the output content and the target task. Subsequently, we calculate the cosine similarity be-613 tween these representations to determine the degree of alignment 614 615 with the target task.

4 EXPERIMENTS

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

In this section, we demonstrate the experimental impact of our method on attacking LLM-based embodied model to assess the robustness of embodied system. Firstly, in Section 4.1, we introduce the experimental setup. Following that, in Section 4.2, we present the comparison of attack results between our method and the most advanced white-box jailbreak attack technologies, including GCG [42] and AutoDAN [17], applied to three different LLM-based embodied models. Subsequently, in Section 4.3, we analyze the execution success rate by user study to evaluate whether the output of LLM can be executed in the current environment. Furthermore, in Section 4.4, we delve into the reasons why the initialization of prompt suffix keywords can significantly reduce the attack success time. Finally, ablation experiments are conducted in Section 4.5 to assess the importance of the two modules proposed by ours.

4.1 Settings

Datasets. We employ the multi-modal dataset EIRAD to assess the LLM robustness of embodied intelligent robots. This dataset comprises a total of 500 instances of untargeted attack data and 500 instances of targeted attack data. Additionally, the targeted attack data is further categorized into 450 instances of harmless attack data and 50 instances of harmful attack data.

Models. To ensure the generality of the attack method, we assess two fine-tuned open-source models (TaPA and Otter) and one un-fine-tuned open-source model (Llama-2-7b-chat) in embodied scenarios. The TaPA model, developed by Wu et al. [36], was fine-tuned using a dataset comprising 15K instruction-task data pairs to refine the Llama model. The Otter model was generated by Li et al. [14] using the MIMIC-IT dataset containing 2.8 million multi-modal instruction-response pairs to fine-tune the OpenFlamingo [2] model. Additionally, the Llama-2-7b-chat model is utilized for decision-making in embodied tasks.

Baselines. Considering the related work on LLM adversarial jailbreaking [7, 16, 17, 33, 42], we compare our method with some representative baseline methods, such as GCG [42] and AutoDAN [17], to assess the robustness of the LLM-based embodied model under white-box attack scenarios. As evident from Section 3.2.3, the text matching list methods employed by GCG and AutoDAN are not suitable for tasks involving embodied attacks. Hence, we replace the text matching algorithms in these two methods with the slicing and similarity calculation methods proposed in our paper as the criteria for judging attack success. Finally, we compare the method proposed in this paper with the GCG[42] and AutoDAN [17] algorithms to demonstrate the advantages of our method in terms of attack success rate and efficiency. It is important to note that our method has the same initial parameters as the original GCG [42], epoch is 500, top-k is 256, and batchsize is 512.

Evaluation metrics. We evaluate the algorithm from three key aspects: Attack Success Rate (ASR), Execution Success Rate (ESR), and Epoch Cost. ASR indicates whether the LLM-based embodied model successfully outputs decision content related to the target task in targeted attacks or outputs decision content unrelated to the prompt in untargeted attacks. ESR reflects whether, upon a successful attack, the system can execute the output content under the prevailing environmental conditions. Epoch Cost denotes the average number of iterations required for an attack to succeed.

4.2 Main results

Table 1 illustrates the white-box attack evaluation results of our method and other baselines[17, 42]. In targeted attacks, we conduct these assessments by generating prompt suffixes for each targeted request in the EIRAD dataset and examining whether the final response from the LLM-based embodied model aligns with the targeted task. In untargeted attacks, we conducted similar evaluations by generating a prompt suffix to examine whether the final response of the LLM-based embodied model remains independent of the prompt task. We noted that in targeted attacks, our method effectively produces prompt suffixes and achieves a superior attack success rate compared to baseline methods. For the fine-tuned LLM-based embodied model TaPA[36] and Otter[14], our method enhances the attack success rate by over 10%, and even surpasses 20% in harmful attack. Regarding the native model Llama-2-chat used for embodied tasks, our method demonstrates a comparable attack success rate to GCG [42] in harmless attack, while significantly reducing the Epoch cost. In untargeted attacks, our method also

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

639

640

Table 1: Main results. We report the ASR and Epoch cost of our method for targeted and untargeted attacks on three models on the EIRAD dataset. Compared with the baseline, our method can effectively attack the LLM-based embodied model and greatly shorten the attack time.

Experiment		Targeted attack-unharmful		Targeted attack-harmful		Untargeted attack	
Model	Method	ASR	Epoch cost	ASR	Epoch cost	ASR	Epoch cos
	AutoDAN	0	500	0	500	-	-
Тара	GCG	0.32	148	0.02	74	-	-
	Ours	0.72	84	0.22	124	1	9
	AutoDAN	0	500	0	500	-	-
Otter	GCG	0.81	101	0.64	180	-	-
	Ours	0.95	56	0.86	120	0.80	67
	AutoDAN	0	500	0	500	-	-
Llama-2-chat	GCG	0.97	142	0.82	207	-	-
	Ours	0.97	60	0.92	127	0.57	104

exhibits varying degrees of success across the three models, leading them to output content unrelated to the prompt. The AutoDAN algorithm's attack success rate, as indicated by the data, is 0. This is due to its core concept of using a semantic prompt framework to guide LLMs to circumvent the value alignment mechanism, which falls short in generating specified content for particular tasks. In summary, our approach proves effective when embodied intelligent robots encounter diverse attack scenarios, enhancing the attack success rate while mitigating training costs. These results suggest that LLM-based embodied models display diminished robustness at the decision-level when subjected to adversarial attacks, offering insights for robustness research on embodied intelligent robots.

4.3 Execution success rate

In the context of attacking a LLM-based embodied model, it is crucial to assess whether the LLM's output aligns with both the task requirements and the embodied constraints, ensuring its successful execution within the current environment. As depicted in Table 2, user study serves as a means to evaluate the ESR of the LLM's output in the given environment. Experimental results demonstrate that in targeted attack, where the target task is designed with a thorough consideration of current environmental factors, the resulting output task steps largely adhere to the embodied requirements. In addition, due to the heightened precision in our attack success assessment, our method exhibits a superior ESR when juxtaposed with the GCG [42]. However, in untargeted attack, where no specific target task is specified, the objective is to guide the LLM-based embodied model to generate content that is unrelated to the original prompt. Consequently, for a successful untargeted attack, the primary criterion is to ensure that the output content is disconnected from the prompt, without considering its feasibility for execution within the current scenario, resulting in a low ESR value.

Table	2:	ESR	based	on	user	study.

		Harmful attack	Harmless attack	Untargeted attack
Model	Method	ESR	ESR	ESR
Тара	GCG	0	0.67	-
	Ours	0.72	0.81	0.48
Otter	GCG	0.74	0.91	-
	Ours	0.84	0.95	0.48
Llama-2-chat	GCG	0.73	0.87	-
	Ours	0.78	0.88	0.45

4.4 The impact of keyword initialization on loss

To delve into the reason behind the significant reduction in epoch cost due to prompt suffix keyword initialization, we conduct an analysis on the change trend of the loss value throughout the attack process. We compare the effects of prompt suffix keyword initialization with 2 keywords versus no keyword initialization on the three models individually. The variations in the loss value are visualized in Figure 7. It is evident that the suffix initialized with keywords exhibits a notably lower loss value in the initial stages of the attack process. Consequently, during the iterative optimization of the suffixes, the model tends to swiftly identify the most suitable prompt suffix. This implies that in the early phases of an attack, the model can swiftly pinpoint an effective attack direction, thereby advancing towards a successful attack status more rapidly. This strategic advantage leads to quicker convergence towards successful attack outcomes, thereby enhancing the overall effectiveness and reliability of the targeted attack process.

4.5 Ablation Study

In order to evaluate the importance of the two modules proposed in this paper, we conduct ablation experiments on the prompt suffix keyword initialization and the evaluation method to determine the success of the attack. In Section 4.5.1, we examined how varying the number of prompt suffix keywords affects the ASR and epoch cost. In Section 4.5.2, we assess the effectiveness of our chosen evaluation method for determining attack success.

4.5.1 **The impact of initializing the number of keywords**. In order to analyze the impact of prompt suffix keyword initialization on ASR and Epoch cost, we conduct an ablation experiment on the number of keywords initialized by prompt suffix and explore its impact on the Otter model attack process. The attack results of other LLM-based embodied models are shown in the appendix. As depicted in Figure 8a and Figure 8b, we set various numbers of keyword initializations for both harmful and harmless attack data in the targeted attack scenario. The experimental findings reveal a consistent trend: as the number of initialized keywords increases, the ASR value steadily rises, while the epoch cost value decreases. This trend suggests that augmenting the number of keywords offers

ACM MM, 2024, Melbourne, Australia



improved guidance to the model in identifying the optimal attack direction, thereby expediting the discovery of effective attack suffixes. However, during the actual attack execution, it becomes imperative to carefully balance the trade-offs among the increased workload due to additional initial suffix keywords, the resultant attack success rate, and the time taken for the attack. Consequently, selecting the optimal number of keywords for prompt suffix initialization becomes a crucial consideration in the attack strategy. Furthermore, our analysis of the evolving trend of loss values under varying initialization keyword numbers, as illustrated in Figure 8c and Figure 8d, reveals a compelling relationship: a higher number of initialized keywords corresponds to a lower initial loss value, resulting in a quicker attainment of attack success. This underscores the significance of judiciously optimizing the number of keywords for prompt suffix initialization to enhance the efficiency and effectiveness of the attack process.

4.5.2 Validity of Assessment Methods. To validate the effectiveness of the evaluation method proposed in this paper to determine the success of an attack, we conduct an ablation experiment of the step 3 in Figure 5. Table 3 illustrates that under the same attack conditions, the combination of slicing and similarity calculation proposed in this paper is compared with the matching list method in GCG [42] and with only the similarity calculation method to eval-uate the correlation between the response of LLM-based embodied model and the target task upon a successful attack. Green markers indicate a strong correlation between the output step and the target task, while red markers signify deviations from the target task. The results indicate that when using our method for evaluation, the





response of LLM-based embodied model closely aligns with the set attack goals. In contrast, under the other two evaluation methods, the output task content deviates from the intended goals. This underscores the effectiveness of our evaluation method in reducing the occurrences of false detections and missed detections observed in GCG [42] and AutoDAN [17], thus enhancing the accuracy of attack results.

Table 3: Responce under three ludgment me	ethods.
---	---------

prompt : Please pass me the chair.			
target : Open the toaster and place it on the table			
Judgment-method	Response		
	Step 1: Open the toaster directory and find the toast.		
GCG	Step 2: Click on the toast to open it.		
	Step 3: Inside the toast.		
Only calculate similarity	Step 1: Open the toaster and insert the bread.		
	Step 2: Close the toaster and press the button.		
	Step 3: Wait for the toast.		
Ours	Step 1: Open the toaster.		
	Step 2: Place the toast on the table.		
	Step 3: Close the toaster.		

5 CONCLUSION

In this paper, we introduce the creation of a multi-modal dataset named EIRAD, which aims to assess the robustness of LLM decisionlevel in embodied scenarios. We devise two attack strategies, untargeted attack and targeted attack, to simulate diverse attack scenarios. Throughout these attacks, we implement prompt suffix keyword initialization tailored to specific target tasks, enhancing the convergence speed during the attack process. Additionally, we develop an attack success assessment method based on BLIP2 model to provide a more precise evaluation of the conditions for attack success. The experimental outcomes validate the effectiveness of our approach, while also underscoring the challenge of robustness in LLM decision-level within embodied scenarios. We aspire that our study will shed further light on the vulnerabilities of LLMs in embodied settings and furnish them with advanced defense mechanisms for ensuring secure utilizatio.

Anonymous Authors

Exploring the Robustness of Decision-Level Through Adversarial Attacks on LLM-Based Embodied Models

ACM MM, 2024, Melbourne, Australia

929 **REFERENCES**

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

- Akiko Aizawa. 2003. An information-theoretic perspective of tf-idf measures. Information Processing & Management 39, 1 (2003), 45-65.
- [2] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models. arXiv preprint arXiv:2308.01390 (2023).
- [3] Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. 2023. Do as i can, not as i say: Grounding language in robotic affordances. In *Conference* on robot learning. PMLR, 287–318.
- [4] Nicholas Carlini, Milad Nasr, Christopher A. Choquette-Choo, Matthew Jagielski, Irena Gao, Anas Awadalla, Pang Wei Koh, Daphne Ippolito, Katherine Lee, Florian Tramer, and Ludwig Schmidt. 2023. Are aligned neural networks adversarially aligned? arXiv:2306.15447 [cs.CL]
- [5] Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. 2023. Jailbreaking Black Box Large Language Models in Twenty Queries. arXiv:2310.08419 [cs.LG]
- [6] Gautier Dagan, Frank Keller, and Alex Lascarides. 2023. Dynamic planning with a llm. arXiv preprint arXiv:2308.06391 (2023).
- [7] Peng Ding, Jun Kuang, Dan Ma, Xuezhi Cao, Yunsen Xian, Jiajun Chen, and Shujian Huang. 2023. A Wolf in Sheep's Clothing: Generalized Nested Jailbreak Prompts can Fool Large Language Models Easily. arXiv preprint arXiv:2311.08268 (2023).
- [8] Yinpeng Dong, Huanran Chen, Jiawei Chen, Zhengwei Fang, Xiao Yang, Yichi Zhang, Yu Tian, Hang Su, and Jun Zhu. 2023. How Robust is Google's Bard to Adversarial Image Attacks? arXiv preprint arXiv:2309.11751 (2023).
- [9] Vishnu Sashank Dorbala, Sanjoy Chowdhury, and Dinesh Manocha. 2024. Can LLMs Generate Human-Like Wayfinding Instructions? Towards Platform-Agnostic Embodied Instruction Synthesis. arXiv preprint arXiv:2403.11487 (2024).
- [10] Vishnu Sashank Dorbala, James F Mullen Jr, and Dinesh Manocha. 2023. Can an Embodied Agent Find Your "Cat-shaped Mug"? LLM-Based Zero-Shot Object Navigation. *IEEE Robotics and Automation Letters* (2023).
- [11] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. 2022. Inner monologue: Embodied reasoning through planning with language models. arXiv preprint arXiv:2207.05608 (2022).
- [12] Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. 2023. Catastrophic Jailbreak of Open-source LLMs via Exploiting Generation. arXiv:2310.06987 [cs.CL]
- [13] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Matt Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, et al. 2017. Ai2thor: An interactive 3d environment for visual ai. arXiv preprint arXiv:1712.05474 (2017).
- [14] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. 2023. Mimic-it: Multi-modal in-context instruction tuning. arXiv preprint arXiv:2306.05425 (2023).
- [15] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2024. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. Advances in Neural Information Processing Systems 36 (2024).
- [16] Xuan Li, Zhanke Zhou, Jianing Zhu, Jiangchao Yao, Tongliang Liu, and Bo Han. 2023. Deepinception: Hypnotize large language model to be jailbreaker. arXiv preprint arXiv:2311.03191 (2023).
- [17] Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. 2023. Autodan: Generating stealthy jailbreak prompts on aligned large language models. arXiv preprint arXiv:2310.04451 (2023).
- [18] Corey Lynch, Ayzaan Wahid, Jonathan Tompson, Tianli Ding, James Betker, Robert Baruch, Travis Armstrong, and Pete Florence. 2023. Interactive language: Talking to robots in real time. *IEEE Robotics and Automation Letters* (2023).
- [19] Kolby Nottingham, Prithviraj Ammanabrolu, Alane Suhr, Yejin Choi, Hannaneh Hajishirzi, Sameer Singh, and Roy Fox. 2023. Do embodied agents dream of pixelated sheep: Embodied decision making using language guided world modelling. In International Conference on Machine Learning. PMLR, 26311–26325.
- [20] Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2023. Fine-tuning Aligned Language Models Compromises Safety, Even When Users Do Not Intend To! arXiv:2310.03693 [cs.CL]
- [21] Yanyuan Qiao, Yuankai Qi, Zheng Yu, Jing Liu, and Qi Wu. 2023. March in chat: Interactive prompting for remote embodied referring expression. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 15758–15767.
- [22] Raphael Schumann, Wanrong Zhu, Weixi Feng, Tsu-Jui Fu, Stefan Riezler, and William Yang Wang. 2024. Velma: Verbalization embodiment of Ilm agents for vision and language navigation in street view. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 18924–18933.
- [23] Rusheb Shah, Quentin Feuillade-Montixi, Soroush Pour, Arush Tagade, Stephen Casper, and Javier Rando. 2023. Scalable and Transferable Black-Box Jailbreaks for Language Models via Persona Modulation. arXiv:2311.03348 [cs.CL]

- [24] SP Sharan, Ruihan Zhao, Zhangyang Wang, Sandeep P Chinchali, et al. 2024. Plan Diffuser: Grounding LLM Planners with Diffusion Models for Robotic Manipulation. In Bridging the Gap between Cognitive Science and Robot Learning in the Real World: Progresses and New Directions.
- [25] Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. 2023. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2998–3009.
- [26] Francesco Stella, Cosimo Della Santina, and Josie Hughes. 2023. How can LLMs transform the robotic design process? *Nature Machine Intelligence* 5, 6 (2023), 561–564.
- [27] Jiankai Sun, De-An Huang, Bo Lu, Yun-Hui Liu, Bolei Zhou, and Animesh Garg. 2022. Plate: Visually-grounded planning with transformers in procedural tasks. *IEEE Robotics and Automation Letters* 7, 2 (2022), 4924–4930.
- [28] Yuxuan Sun, Chenglu Zhu, Sunyi Zheng, Kai Zhang, Zhongyi Shui, Xiaoxuan Yu, Yizhi Zhao, Honglin Li, Yunlong Zhang, Ruojia Zhao, et al. 2023. Pathasst: Redefining pathology through generative foundation ai assistant for pathology. arXiv preprint arXiv:2305.15072 (2023).
- [29] Andrew Szot, Max Schwarzer, Harsh Agrawal, Bogdan Mazoure, Rin Metcalf, Walter Talbott, Natalie Mackraz, R Devon Hjelm, and Alexander T Toshev. 2023. Large language models as generalizable policies for embodied tasks. In *The Twelfth International Conference on Learning Representations*.
- [30] Inigo Jauregi Unanue, Jacob Parnell, and Massimo Piccardi. 2021. BERT-Tune: Fine-tuning neural machine translation with BERTScore. arXiv preprint arXiv:2106.02208 (2021).
- [31] Sai Vemprala, Rogerio Bonatti, Arthur Bucker, and Ashish Kapoor. 2023. Chatgpt for robotics: Design principles and model abilities. arXiv preprint arXiv:2306.17582 (2023).
- [32] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023. Jailbroken: How Does LLM Safety Training Fail? arXiv:2307.02483 [cs.LG]
- [33] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2024. Jailbroken: How does llm safety training fail? Advances in Neural Information Processing Systems 36 (2024).
- [34] Jimmy Wu, Rika Antonova, Adam Kan, Marion Lepert, Andy Zeng, Shuran Song, Jeannette Bohg, Szymon Rusinkiewicz, and Thomas Funkhouser. 2023. Tidybot: Personalized robot assistance with large language models. *Autonomous Robots* 47, 8 (2023), 1087–1102.
- [35] Lei Wu, Steven CH Hoi, and Nenghai Yu. 2010. Semantics-preserving bag-ofwords models and applications. *IEEE Transactions on Image Processing* 19, 7 (2010), 1908–1920.
- [36] Zhenyu Wu, Ziwei Wang, Xiuwei Xu, Jiwen Lu, and Haibin Yan. 2023. Embodied task planning with large language models. arXiv preprint arXiv:2307.01848 (2023).
- [37] Yijun Yang, Tianyi Zhou, Kanxue Li, Dapeng Tao, Lusong Li, Li Shen, Xiaodong He, Jing Jiang, and Yuhui Shi. 2023. Embodied multi-modal agent trained by an llm from a parallel textworld. arXiv preprint arXiv:2311.16714 (2023).
- [38] Jiahao Yu, Xingwei Lin, Zheng Yu, and Xinyu Xing. 2023. GPTFUZZER: Red Teaming Large Language Models with Auto-Generated Jailbreak Prompts. arXiv:2309.10253 [cs.AI]
- [39] Hongxin Zhang, Weihua Du, Jiaming Shan, Qinhong Zhou, Yilun Du, Joshua B Tenenbaum, Tianmin Shu, and Chuang Gan. 2023. Building cooperative embodied agents modularly with large language models. arXiv preprint arXiv:2307.02485 (2023).
- [40] Zihao Zhao, Sheng Wang, Jinchen Gu, Yitao Zhu, Lanzhuju Mei, Zixu Zhuang, Zhiming Cui, Qian Wang, and Dinggang Shen. 2023. Chatcad+: Towards a universal and reliable interactive cad using llms. arXiv preprint arXiv:2305.15964 (2023).
- [41] Sipeng Zheng, Yicheng Feng, Zongqing Lu, et al. 2023. Steve-Eye: Equipping LLM-based Embodied Agents with Visual Perception in Open Worlds. In *The Twelfth International Conference on Learning Representations*.
- [42] Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. arXiv preprint arXiv:2307.15043 (2023).

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043 1044

987

988

989

990

991

992

993

994