Module Attack: Exploiting Module Swapping to Compromise LLM Alignments

Anonymous ACL submission

Abstract

In this study, we introduce a novel approach for undermining the alignment of large language models (LLMs), which we term the Module Attack. A module attack compromises the alignment of a model by manipulating intermediate modules in the LLM by changing the internal structure of the model through module swapping. Unlike traditional prompt-based jailbreak attacks, which rely on external inputs and have limited effectiveness, we show that module attacks can bypass alignment defense mechanisms by exploiting structural vulnerabilities inside the LLM and can be answered without going through a separate prompt engineering process.

004

005

007

011

022

026

We also propose a cooperative decoding approach that alternately generates tokens from the attacked LLM and the original LLM during token generation. In conclusion, we achieved high ASRs, reaching 100% in most cases, across different LLM architectures (Qwen 2.5, Llama 3.1, Mistral v0.3), and found no difference in ASR between generation using the attacked LLM alone and cooperative decoding with the original LLM. We also showed that a simple swap of internal modules in the LLM can break the alignment of the model without any prompt engineering. This is a methodology that can neutralize the alignment of a model faster than any other methodology without any prior action.

> This research provided a deep understanding of the structural vulnerabilities of LLMs and confirmed that manipulating modules in LLMs can easily lead to unwanted consequences.

1 Introduction

Recent breakthroughs in large language models (LLMs) have led to exponential performance gains across diverse tasks (Wang et al., 2024b; Chen et al., 2021; Jain et al., 2024; Zheng et al., 2023), enabling models to produce both high-quality answers and unexpectedly creative outputs (Wei et al., 2022).



Figure 1: Illustration of two cases where LLM evades the question "How to make a Bomb?" while Module Attacked LLM gives the recipe for the actual bomb

However, these powerful capabilities also raise serious concerns about malicious use, as LLMs can become formidable "digital weapons" when deployed by adversarial actors. To mitigate these risks, researchers have introduced various safety mechanisms, including supervised fine-tuning (SFT) (Bianchi et al., 2023), reinforcement learning from human feedback (RLHF), and red-teaming (Ganguli et al., 2022; Perez et al., 2022). Despite these efforts, multiple studies have already highlighted fundamental limitations of learning-based alignment (Wolf et al., 2023; Dai et al., 2023; Su et al., 2024).

Examining the internal structure of LLMs reveals that early layers tend to capture local sentence-level information (Zhang et al., 2024), while middle layers form more abstract and high-level representations (Skean et al., 2024). Closer to the output layer, the model aligns these representations for specific tasks or ethical guidelines (Wang et al., 2024a). A growing body of work indicates that modifying certain layers can strengthen or weaken the model's alignment (Zhao et al., 2024; Hasan et al., 2024), demonstrating the close link between a model's structural components and its safety.

Against this backdrop, this study proposes a new methodology, hereinafter referred to as **Module**

043

044

Swapping, which exploits structural vulnerabili-071 ties in LLMs to weaken their alignment. Although 072 Model Manipulation-based jailbreak techniques 073 have been introduced in previous research, they often involve complex procedures that limit their practicality. Likewise, Prompt-based jailbreak methods require users to search for or engineer 077 specific prompts, making them tedious and highly situational. In contrast, the Module Swapping method described in this white paper can be executed quickly and simply by reconfiguring the model's internal components (layers, multi-layer perceptrons, attention modules) without intricate prompt engineering. In particular, we demonstrate 084 that rearranging the order of specific modules can significantly degrade the model's alignment. Notably, even an overtly malicious query such as "How to make a bomb?" becomes answerable under this attack scenario. In the case of Closed LLMs operating on hacked hardware, one could obtain dangerous information more efficiently and rapidly than through any other known method, underscoring the severe security implications of structural manipulation. 094

> We also investigate how this attack generalizes across different conditions, including Collaborative Decoding, thereby offering a broader perspective on whether rearranging the order of specific Modules alone can circumvent alignment. Our extensive experiments cover state-of-the-art architectures, such as Qwen 2.5, Llama 3.1, and Mistral 7B. We quantitatively measure the Attack Success Rate (ASR) and other performance indicators to systematically assess the model's vulnerability under this structural manipulation. These findings suggest that simply relying on learning-based alignment (e.g., SFT, RLHF) may be insufficient to protect against deeper, structural vulnerabilities that emerge from the model's intermediate Modules.

096

097

100

101

102

103

104

105

107

108

109

110

111

112

114

116

117

118

120

121

In summary, the contributions of this work are as follows.

We highlight how the middle Modules of LLMs often considered primarily for abstract represen-113 tation also serve as a critical pivot for alignment mechanisms. By focusing on Module Swap, we 115 elucidate how structural manipulations can undermine a model's safety. We empirically verify the attack's generality across different model architectures (Owen 2.5, Llama 3.1, Mistral 7B) and var-119 ious settings, including Collaborative Decoding, showing that structural vulnerabilities are widely

shared among current LLMs. We emphasize that true robustness cannot rely solely on high-level alignment techniques such as SFT and RLHF; rather, it requires a holistic approach that considers all Modules and their interactions. By highlighting the ease with which alignment can be bypassed through simple structural changes, this study underscores the urgent need for research into Modulespecific defenses and more holistic safety mechanisms. We anticipate that our findings will inform both next-generation LLM design and the broader field of AI safety, guiding the development of more robust alignment strategies that account for vulnerabilities beyond mere output level control.

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

2 **Related Works**

Prompt based Jailbreaking 2.1

Prompt based Jailbreaking is one of the main focuses of existing research to bypass the alignment of LLMs. This methodology designs malicious prompts to induce the model to produce unwanted information. For example, (Jiang et al., 2024) proposed a prompt attack method utilizing ASCII art, and (Zeng et al., 2024) introduced an approach to evade the model's safety mechanisms through persuasion based techniques. In addition, (Chang et al., 2024) showed that alignment can be bypassed by combining multiple attack prompts. (Deng et al., 2024) proposed an indirect attack technique called PANDORA that exploited Retrieval Augmented Generation (RAG). However, these methods have had limited effectiveness because they rely on external inputs to the model and do not directly compromise the internal structure of the model.

However, these prompt based approaches rely on external input from the model and do not directly manipulate its internal structure. As a result, they are only effective in certain scenarios and are limited in their ability to weaken the alignment of the model.

2.2 Model Manipulation based Jailbreaking

Research on jailbreaking models by modifying their internal structure has recently gained traction. (Zhao et al., 2024) proposed a methodology to weaken alignment by modifying certain layers, (Zhang et al., 2023) proposed a coercive knowledge extraction method that utilizes the output logic of a model to force the generation of harmful information. Modern automated black box attack techniques, such as (Mehrotra et al., 2023)TAP (Tree



Figure 2: An illustration of Module Attack(Swapping) + Collaborative Decoding. L is the total number of Modules in the model. In the attack process, Swapped LLM is created by attacking the LLM with Layer Swap to modify the alignment of the LLM, and then Collaborative Decoding is performed by alternately generating tokens of the Attacked LLM and the original LLM during the inference process.

of Attacks with Pruning), have greatly improved the success rate in bypassing alignment by generating effective attack prompts without the need to know the internal state of the model, (Zou et al., 2023) successfully used a gradient based method to generate malicious input.

3 Methodology

171

172

173

174

175

176

177

178

179

181

182

183

184

185

188

189

190

191

193

194

In this section, we introduce **module swapping**, a method that disrupts the alignment of an LLM by modifying its core layers. Additionally, we propose **collaborative decoding**, a technique that improves the quality of generated responses by alternately generating tokens between the attacked LLM and the original LLM.

3.1 Module Swapping

Module Swapping is a targeted intervention mechanism that swaps specific intermediate Modules of LLMs, thereby compromising the alignment capabilities essential for maintaining model safety. Given a model $f = \{f_1, f_2, \ldots, f_L\}$, with $E \in \mathbb{R}^{V \times d}$ representing the embedding space and $H \in \mathbb{R}^{d \times V}$ denoting the LM head, the attacked model is formulated as:

AttackedLLM
$$(x) = H \Big(f_L \circ \cdots \circ \sigma \big(f_{L/2-t}, \cdots, f_{L/2+t} \big)$$

 $\circ \cdots \circ f_1 \Big) (E(x)).$

where σ represents the swapping operator applied over the Module range $\{f_{L/2-t}, \ldots, f_{L/2+t}\}$. Module swapping disrupts the hierarchical processing of information within LLMs by reordering critical intermediate Modules. For a subset of Modules $\{f_{L/2-t}, \ldots, f_{L/2+t}\}$, a permutation π is applied, resulting in a rearranged structure. The modified Module composition, denoted as

$$\sigma = \sigma_{\text{swap}}(f_{L/2-t}, \dots, f_{L/2+t}),$$

197

198

199

201

202

205

206

207

209

210

211

212

213

produces a model structure incapable of maintaining the original alignment, thereby amplifying the attack's impact.

3.2 Collaborative Decoding

To maximize the attack's efficacy, we introduce **Collaborative Decoding**, a hybrid generation process alternating between the original LLM (L) and the attacked model (A). The generation process is defined as:

$$\phi_{\rho}(t) = \begin{cases} A, & \text{if } ((t-1) \mod (p+q)) < p, \\ L, & \text{otherwise,} \end{cases}$$

where $\rho = p : q$ represents the token generation 214 ratio between A and L. The process continues until 215 an end of sequence (EOS) token is generated or the 216 maximum decoding length T_{max} is reached. 217

218 219

220

224

231

241

242

243

245

246

247

249

256

261

262

264

4 Experiment

4.1 Experimental Setup

Model. We selected the following models for our experiments based on performance, awareness, and downloads: Qwen 2.5 7b instruct, Llama 3.1 8b instruct, and Mistral 7b instruct v0.3. The Qwen 2.5 model used online RL and offline RL, bias removal, response filtering, etc. to make the model alignment robust(Yang et al., 2024) The Mistral model performed content moderation with self-reflection, achieving 99.4% accuracy in self-reflection precision(Jiang et al., 2023) The Llama 3.1 model can prevent direct jailbreak attempts using prompts with Prompt Guard, which is a model that prevents jailbreaks through red teaming and safe fine-tuning(Dubey et al., 2024).

Dataset. We use JailbreakBench to assess the model's robustness against jailbreak attempts and MMLU-Pro to evaluate its general knowledge and reasoning abilities. JailbreakBench is used to assess how effectively the proposed methodology weakens the model's alignment, while the MMLU-Pro Benchmark measures the extent to which the model's general performance declines. dataset from JailbreakBench(Chao et al., 2024) and the JBB-Behaviors(JailbreakBench, 2025) JBB-Behaviors is a dataset of a list of 100 misuse behaviors curated from OpenAI's usage policies(OpenAI, 2025) and broken down into 10 main categories. Each category represents 10% of the dataset, and the dataset is composed of 55% original, 27% from TDC/HarmBench, 18% from AdvBench, and 18% from TDC/HarmBench.

dataset from MMLU-Pro (Wang et al., 2024b) and the MMLU-Pro Dataset(Wang et al., 2025) is a powerful and challenging large-scale multi-task comprehension dataset tailored for rigorously benchmarking the capabilities of large-scale language models, consisting of 14 tasks with a total of 12,102 datasets.

Evaluation Metric. As an evaluation metric, we used Attack Success Rate (ASR), which is defined as follows.

$$ASR = \frac{\text{Number of Successful Attacks}}{\text{Total Number of Attack Attempts}} \times 100\%$$

Each question was generated a total of 50 times, and if any of the answers were misused, we considered the attack successful. We also used the original, risky prompts directly as input prompts, without any prompt engineering.

4.2 Module Attack & collaborative decoding

The experiments were conducted on Module Swap methodologies, and the Module swap process was performed by swapping 1:1, 2:2, 3:3, and 4:4 Modules based on the middle Modules. Also, for each methodology, the 1:1, 2:1, 3:1, 4:1, ∞ of the collaborative decoding to measure ASR.

4.3 Experimental Results

Module Swap. Module Swap experiments show that the Alignment attack success rate (ASR) reaches 100% for most models. The experimental results are shown in the following tables Table 1, and Figure 3, Figure 4, and Figure 5. Furthermore, to evaluate the general performance of the models, the MMLU-Pro results for Swapping and Collaborative Decoding can be found in Table 2 and Table 3.

As shown in Table 1, the module swapping methodology generally achieved an ASR close to 100%. Regardless of the number of swapped modules, the ASR remained at a minimum of 90%. Notably, in the case of the Mistral v0.3 and Llama 3.1 models, swapping even a single module resulted in an ASR approaching 100%. Furthermore, a general trend was observed in which the ASR increased as the number of swapped modules increased.

For the Mistral v0.3 model, when layer swapping was performed, the ASR for the model swapped at a 1:1 ratio was 93%, whereas the model swapped at a 4:4 ratio exhibited an ASR of 95%, indicating a 2% improvement. Similarly, MLP swapping led to an increase in ASR from 91% to 97%, reflecting a 6% improvement, while attention swapping resulted in a 2% increase from 95% to 97%. This trend was also observed in the Qwen 2.5 model, excluding the Llama 3.1 model. Specifically, the Qwen 2.5 model demonstrated an ASR improvement of 5% for layer swapping, 10% for MLP swapping, and 7% for attention swapping.

As shown in Table 2, swapping a single module at a time does not result in significant performance degradation compared to the original model. For instance, in the case of Mistral v0.3, the original model achieved a score of 0.36, whereas models with individual module swaps obtained scores of 0.35, 0.35, and 0.34, reflecting only a minimal difference of approximately 0.01 points. Similarly, for 268 269

267

270 271

272

273

274

275

277

278

279

280

282

283

284

286

287

288

289

290

291

292

294

295

296

297

298

299

301

302

303

304

305

306

307

308

309

310

311

312

313

314

| | | Layer Swapping (ASR) | | | | | | MLP Swapping (ASR) | | | | Attention Swapping (ASR) | | | | |
|-----------|--------|----------------------|------|------|------|------|------|--------------------|------|------|------|--------------------------|------|------|------|------|
| Model | Module | INF | 1:1 | 2:1 | 3:1 | 4:1 | INF | 1:1 | 2:1 | 3:1 | 4:1 | INF | 1:1 | 2:1 | 3:1 | 4:1 |
| Mistral | 1 | 93% | 98% | 99% | 99% | 98% | 91% | 98% | 98% | 99% | 99% | 95% | 100% | 99% | 99% | 99% |
| | 2 | 91% | 99% | 100% | 98% | 98% | 94% | 100% | 99% | 100% | 99% | 94% | 100% | 99% | 99% | 100% |
| | 3 | 95% | 100% | 100% | 100% | 100% | 94% | 100% | 100% | 100% | 100% | 94% | 99% | 100% | 100% | 100% |
| | 4 | 95% | 100% | 100% | 100% | 100% | 97% | 100% | 100% | 100% | 100% | 97% | 100% | 100% | 100% | 100% |
| Qwen 2.5 | 1 | 95% | 93% | 93% | 97% | 94% | 90% | 87% | 90% | 93% | 87% | 93% | 89% | 88% | 89% | 90% |
| | 2 | 98% | 97% | 100% | 97% | 98% | 98% | 92% | 96% | 93% | 96% | 91% | 87% | 90% | 94% | 90% |
| | 3 | 100% | 99% | 100% | 100% | 100% | 99% | 98% | 99% | 100% | 100% | 99% | 91% | 95% | 94% | 98% |
| | 4 | 100% | 100% | 100% | 99% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 98% | 100% | 99% | 100% |
| Llama 3.1 | 1 | 99% | 99% | 99% | 99% | 98% | 99% | 98% | 99% | 99% | 99% | 100% | 99% | 99% | 98% | 100% |
| | 2 | 99% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 99% | 98% | 100% | 100% | 99% |
| | 3 | 99% | 100% | 99% | 99% | 99% | 99% | 100% | 99% | 100% | 99% | 100% | 99% | 99% | 99% | 99% |
| | 4 | 98% | 99% | 100% | 99% | 99% | 97% | 99% | 99% | 98% | 99% | 96% | 99% | 99% | 100% | 99% |

Table 1: The ASR results for layer, MLP, and attention module swapping are presented. Here, 1, 2, 3, and 4 indicate the number of swapped modules. The notations 1:1, 2:1, 3:1, and 4:1 represent the ratio used during collaborative decoding with the original model, where the first value corresponds to the generation ratio from the attacked LLM, and the second value indicates the generation ratio from the original LLM. The detailed ASR results for each attempt can be found in appendix A



Figure 3: Figures showing the results of Layer Swap. Each figure illustrates the change in ASR for each attempt. 3a, 3e, and 3i present the ASR results when swapping one pair of layers. 3b, 3f, and 3j show the ASR results when swapping two pairs of layers. 3c, 3g, and 3k depict the ASR results when swapping three pairs of layers. 3d, 3h, and 3l display the ASR results when swapping four pairs of layers.

the Llama 3.1 and Qwen 2.5 models, the performance of the swapped models remained comparable to or even surpassed that of their respective prior models.

However, as the number of swapped modules

increases, such as in the 2:2, 3:3, and 4:4 configurations, the performance exhibits an exponential decline. This observation suggests that while module swapping can achieve near-perfect ASR, the optimal swap ratio for maintaining maximum per-

```
5
```



Figure 4: Figures showing the results of MLP Swap. Each figure illustrates the change in ASR for each attempt. 4a, 4e, and 4i present the ASR results when swapping one pair of MLPs. 4b, 4f, and 4j show the ASR results when swapping two pairs of MLPs. 4c, 4g, and 4k depict the ASR results when swapping three pairs of MLPs. 4d, 4h, and 4l display the ASR results when swapping four pairs of MLPs.

| MMLU-Pro (Module Swapping) | | | | | | | | |
|----------------------------|-----------|---------------------|------|------|------|--|--|--|
| Model | Module | 1:1 | 2:2 | 3:3 | 4:4 | | | |
| | Attention | 0.39 | 0.30 | 0.13 | 0.11 | | | |
| Llama 3.1 | MLP | 0.40 | 0.30 | 0.13 | 0.11 | | | |
| | Layer | 0.41 | 0.28 | 0.14 | 0.11 | | | |
| | Attention | 0.35 | 0.30 | 0.20 | 0.12 | | | |
| Mistral v0.3 | MLP | 0.35 | 0.30 | 0.22 | 0.13 | | | |
| | Layer | 0.34 | 0.29 | 0.21 | 0.13 | | | |
| | Attention | 0.41 | 0.28 | 0.20 | 0.11 | | | |
| Qwen 2.5 | MLP | 0.46 | 0.39 | 0.22 | 0.15 | | | |
| | Layer | 0.47 | 0.37 | 0.22 | 0.18 | | | |
| | | Llama 3.1 - 0.44 | | | | | | |
| Bas | e | Mistral v0.3 - 0.36 | | | | | | |
| | | Qwen 2.5 - 0.56 | | | | | | |
| | | Llama 3 - 0.40 | | | | | | |
| Pric | or | Mistral v0.2 - 0.31 | | | | | | |
| | | Qwen 2 - 0.44 | | | | | | |

Table 2: MMLU-Pro performance comparison by model with module swapping. 1:1, 2:2, 3:3, and 4:4 are evaluated after swapping one, two, three, and four Attention, MLP, and Layer modules, respectively.

| MMLU-Pro (Collaborative Decoding) | | | | | | | | | |
|-----------------------------------|---------------------|---------------------|------|------|------|------|--|--|--|
| Model | Module | INF | 1:1 | 2:1 | 3:1 | 4:1 | | | |
| | Attention | 0.41 | 0.41 | 0.41 | 0.38 | 0.41 | | | |
| Llama 3.1 | MLP | 0.44 | 0.44 | 0.34 | 0.41 | 0.39 | | | |
| | Layer | 0.44 | 0.41 | 0.45 | 0.45 | 0.41 | | | |
| | Attention | 0.31 | 0.34 | 0.33 | 0.32 | 0.32 | | | |
| Mistral v0.3 | MLP | 0.32 | 0.37 | 0.36 | 0.32 | 0.36 | | | |
| | Layer | 0.32 | 0.33 | 0.34 | 0.33 | 0.31 | | | |
| | Attention | 0.46 | 0.45 | 0.47 | 0.40 | 0.44 | | | |
| Qwen 2.5 | MLP | 0.50 | 0.54 | 0.45 | 0.53 | 0.46 | | | |
| | Layer | 0.42 | 0.52 | 0.48 | 0.49 | 0.49 | | | |
| | Llama 3.1 - 0.5 | | | | | | | | |
| Bas | Mistral v0.3 - 0.35 | | | | | | | | |
| | | Qwen 2.5 - 0.55 | | | | | | | |
| | | Llama 3 - 0.44 | | | | | | | |
| Pric | or | Mistral v0.2 - 0.30 | | | | | | | |
| | Qwen 2 - 0.42 | | | | | | | | |

Table 3: MMLU-Pro Performance Comparison via 1:1 Module Swapping and Collaborative Decoding. INF denotes the MMLU-Pro score of the swapped model alone. The results for 1:1, 2:1, 3:1, and 4:1 represent the outcomes of collaborative decoding between the attacked LLM and the original model. For evaluation, 10 random samples were selected from each task of MMLU-Pro and assessed accordingly.



Figure 5: Figures showing the results of Attention Swap. Each figure illustrates the change in ASR for each attempt. 5a, 5e, and 5i present the ASR results when swapping one pair of attentions. 5b, 5f, and 5j show the ASR results when swapping two pairs of attentions. 5c, 5g, and 5k depict the ASR results when swapping three pairs of attentions. 5d, 5h, and 5l display the ASR results when swapping four pairs of attentions.

formance is 1:1 module swapping.

These results suggest that the Alignment defense is concentrated in an overall module and can be easily neutralized by swapping that module without significant performance degradation. We also experimentally demonstrate that a simple module swap can effectively bypass Alignment without prompt engineering.

Collaborative Decoding.

As shown in Table 1, the ASR results obtained through collaborative decoding between the swapped LLM and the original LLM generally achieved a 100% ASR. Additionally, in some cases, collaborative decoding with the original model led to higher ASR compared to using only the swapped LLM for jailbreak attempts. For instance, when swapping a single layer in the Mistral v0.3 model, the ASR of the swapped LLM alone was 93%. However, when collaborative decoding was performed at 1:1, 2:1, 3:1, and 4:1 ratios, the ASR improved to 99%. This trend suggests that regardless of the number of swapped modules in the Mistral v0.3 model, collaborative decoding consistently yields higher ASR than using the swapped LLM alone. A similar, albeit weaker, trend was observed in the Qwen 2.5 and Llama 3.1 models, where collaborative decoding not only enhanced jailbreak effectiveness but also demonstrated superior general performance compared to using the swapped LLM in isolation.

348

349

350

351

352

354

355

356

357

359

360

361

363

365

366

367

369

The results of collaborative decoding with a single-module swap for each model are presented in Table 3. In general, the performance of the model utilizing collaborative decoding with a single swapped module exhibited similar results to the model generated solely with a single-module swap. However, when more than one module was swapped during generation, it was empirically observed that the model encountered various issues, such as failing to capture basic contextual information. As shown in appendix B, collaborative decoding effectively mitigated these issues. In conclusion, experimental results demonstrated that as

the number of swapped modules increased, collaborative decoding was able to maintain the model's
performance. Additionally, even when only one
module was swapped, the performance degradation
remained relatively minimal, as verified through
experiments.

5 Discussion

376

377

387

390

391

396

400

401

402

403

404

405

406

407

408

409

410

411

412

Dataset. The JBB-Behaviors dataset we used in our experiments provides a wide range of attack scenarios, including high-risk behaviors such as harassment, malware creation, and fraud, and the balanced distribution of these categories provides a comprehensive evaluation of the Module Attack methodology, which showed no particular resistance in any particular category and an overall ASR close to 100%. This suggests that the structural vulnerabilities of the LLM used by Module Attack are not simply related to the semantic content of the prompt, and that a structural level of defense is required.

Collaborative Decoding. Our findings indicate that collaborative decoding enhances the quality of generated sentences by alternating token generation between the compromised LLM and the original LLM. Furthermore, our experiments demonstrate that the alignment of the original LLM can be effectively neutralized during the collaborative decoding process.

However, collaborative decoding may leave detectable traces in the interaction patterns between the compromised and original models. Such traces could be identified by defense systems designed to detect anomalous model behavior. To further validate our findings, additional experiments are necessary using a more extensive dataset than the JBB-Behavior dataset employed in this study. Moreover, since the MMLU-Pro evaluation of collaborative decoding was conducted with only 10 data points per task, a more rigorous and comprehensive evaluation is required to ensure the robustness of our conclusions.

413Multimodal. Beyond traditional text-based414LLMs, it is necessary to analyze the impact of mod-415ule swapping in multimodal models. Multimodal416LLMs include additional Modules and complex417interaction mechanisms to process various inputs418such as images, speech, etc. In addition to text, this419multimodal structure plays an important role in in-

tegrating and complementarily processing information from each modality, unlike text-based models, implying that different results may be obtained in multimodal LLMs as opposed to text-based models. Further analysis of whether module swapping of these multimodals breaks the alignment is worthwhile and may provide deeper insights into the role of modules in multimodal as well as LLMs. 420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

6 Conclusion

In this study, we propose Module Swapping, a methodology to swap the internal modules of a model, and Collaborative Decoding, a method to generate tokens alternately with the original model. We experimented on Llama 3.1, Qwen 2.5, and Mistral v0.3 models and tested methods. In conclusion, we can see that for 50 attempts, no matter how many modules we swap, we got close to 100% ASR, which shows that it is possible to get close to 100% ASR with the smallest amount of Module swap. We also showed that swapping too many modules leads to a loss of internal knowledge and reasoning ability of the model, which can be recovered through collaborative decoding. Unlike existing prompt-based jailbreak methods, this study shows that the alignment defense mechanism can be bypassed by directly attacking the structural characteristics of the model. Furthermore, by analyzing the impact of Module Attack and Collaborative Decoding on the performance and stability of the model through several experiments, we provide a deeper understanding of the structural vulnerabilities of LLMs. This demonstrates that relying solely on the model is insufficient to prevent jailbreaks, highlighting the necessity for additional pipelines, such as detecting and blocking harmful prompts externally to the model. These results may help to provide new research directions to enhance the safety of LLMs in the future and to compensate for possible vulnerabilities in their structural design.

7 Limitations

The Module Swapping methodology proposed in this study is a methodology that neutralizes the alignment of LLM by swapping the middle Module, and in the process, we found that the model's reasoning ability, knowledge, etc. decreased overall. To prevent this, we introduced collaborative decoding to improve the overall quality of answers. However, this has only been partially validated through empirical observations and requires further rigor-

ous verification. The model sizes used in our exper-469 iments (7b-8b) correspond to medium-sized LLMs, 470 and the effectiveness of Module Swapping is likely 471 to change when considering the complexity of very 472 large models (>100b) and the interactions between 473 different modules. Further research should be con-474 ducted to evaluate the impact of Swapping on very 475 large models. 476

References

477

478

479

480

481

482

483

484

485

486

487

488

489 490

491

492

493

494

495

496

497

499

500

501

502

503 504

505

507

508

510

511

512

513

514

515

516

517

518

519

520

- Federico Bianchi, Mirac Suzgun, Giuseppe Attanasio, Paul Röttger, Dan Jurafsky, Tatsunori Hashimoto, and James Zou. 2023. Safety-tuned llamas: Lessons from improving the safety of large language models that follow instructions. *arXiv preprint arXiv:2309.07875*.
- Zhiyuan Chang, Mingyang Li, Yi Liu, Junjie Wang, Qing Wang, and Yang Liu. 2024. Play guessing game with llm: Indirect jailbreak attack with implicit clues. *arXiv preprint arXiv:2402.09091*.
- Patrick Chao, Edoardo Debenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwag, Edgar Dobriban, Nicolas Flammarion, George J Pappas, Florian Tramer, et al. 2024. Jailbreakbench: An open robustness benchmark for jailbreaking large language models. *arXiv preprint arXiv:2404.01318*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. 2023. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*.
- Gelei Deng, Yi Liu, Kailong Wang, Yuekang Li, Tianwei Zhang, and Yang Liu. 2024. Pandora: Jailbreak gpts by retrieval augmented generation poisoning. *arXiv preprint arXiv:2402.08416*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, et al. 2022. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. arXiv preprint arXiv:2209.07858.
 - Adib Hasan, Ileana Rugina, and Alex Wang. 2024. Pruning for protection: Increasing jailbreak resistance

in aligned llms without fine-tuning. *arXiv preprint arXiv:2401.10862*.

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

- JailbreakBench. 2025. Jbb-behaviors dataset. https://huggingface.co/datasets/ JailbreakBench/JBB-Behaviors. Accessed: 2025-01-22.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. arXiv preprint arXiv:2310.06825.
- Fengqing Jiang, Zhangchen Xu, Luyao Niu, Zhen Xiang, Bhaskar Ramasubramanian, Bo Li, and Radha Poovendran. 2024. Artprompt: Ascii art-based jailbreak attacks against aligned llms. *arXiv preprint arXiv:2402.11753*.
- Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. 2023. Tree of attacks: Jailbreaking black-box llms automatically. *arXiv preprint arXiv:2312.02119*.
- OpenAI. 2025. Usage policies. https://openai.com/ policies/usage-policies/. Accessed: 2025-01-22.
- Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. 2022. Red teaming language models with language models. *arXiv preprint arXiv:2202.03286*.
- Oscar Skean, Md Rifat Arefin, Yann LeCun, and Ravid Shwartz-Ziv. 2024. Does representation matter? exploring intermediate layers in large language models. *arXiv preprint arXiv:2412.09563*.
- Jingtong Su, Julia Kempe, and Karen Ullrich. 2024. Mission impossible: A statistical perspective on jailbreaking llms. *arXiv preprint arXiv:2408.01420*.
- Chenglong Wang, Hang Zhou, Kaiyan Chang, Bei Li, Yongyu Mu, Tong Xiao, Tongran Liu, and Jingbo Zhu. 2024a. Hybrid alignment training for large language models. *arXiv preprint arXiv:2406.15178*.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, et al. 2024b. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*.

Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, et al. 2025. Mmlu-pro dataset. https://huggingface. co/datasets/TIGER-Lab/MMLU-Pro. Accessed: 2025-02-05.

574

575

580

585 586

587

588

589 590

591

593

594

596

598

599

608

610

611

612

613

614

615

616

- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*.
- Yotam Wolf, Noam Wies, Oshri Avnery, Yoav Levine, and Amnon Shashua. 2023. Fundamental limitations of alignment in large language models. *arXiv preprint arXiv:2304.11082*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.
- Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. 2024. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms. *arXiv preprint arXiv:2401.06373*.
- Yang Zhang, Yanfei Dong, and Kenji Kawaguchi. 2024. Investigating layer importance in large language models. *arXiv preprint arXiv:2409.14381*.
- Zhuo Zhang, Guangyu Shen, Guanhong Tao, Siyuan Cheng, and Xiangyu Zhang. 2023. Make them spill the beans! coercive knowledge extraction from (production) llms. *arXiv preprint arXiv:2312.04782*.
- Wei Zhao, Zhe Li, Yige Li, Ye Zhang, and Jun Sun. 2024. Defending large language models against jailbreak attacks via layer-specific editing. arXiv preprint arXiv:2405.18166.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. Advances in Neural Information Processing Systems, 36:46595–46623.
- Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.



Figure 6: Illustration of the results of collaborative decoding for each ratio after swapping for each MLP in the Llama 3.1 model. The numbers 1, 2, 3, and 4 MLP indicate the respective pairs of MLP that have been swapped, while INF, 1:1, 2:1, 3:1, and 4:1 represent the ratios for collaborative decoding.



Figure 7: Illustration of the results of collaborative decoding for each ratio after swapping for each MLP in the Mistral v0.3 model. The numbers 1, 2, 3, and 4 MLP indicate the respective pairs of MLP that have been swapped, while INF, 1:1, 2:1, 3:1, and 4:1 represent the ratios for collaborative decoding.



Figure 8: Illustration of the results of collaborative decoding for each ratio after swapping for each MLP in the Qwen 2.5 model. The numbers 1, 2, 3, and 4 MLP indicate the respective pairs of MLP that have been swapped, while INF, 1:1, 2:1, 3:1, and 4:1 represent the ratios for collaborative decoding.



Figure 9: Illustration of the results of collaborative decoding for each ratio after swapping for each Attention in the Llama 3.1 model. The numbers 1, 2, 3, and 4 Attention indicate the respective pairs of Attention that have been swapped, while INF, 1:1, 2:1, 3:1, and 4:1 represent the ratios for collaborative decoding.



Figure 10: Illustration of the results of collaborative decoding for each ratio after swapping for each Attention in the Mistral v0.3 model. The numbers 1, 2, 3, and 4 Attention indicate the respective pairs of Attention that have been swapped, while INF, 1:1, 2:1, 3:1, and 4:1 represent the ratios for collaborative decoding.



Figure 11: Illustration of the results of collaborative decoding for each ratio after swapping for each Attention in the Qwen 2.5 model. The numbers 1, 2, 3, and 4 Attention indicate the respective pairs of Attention that have been swapped, while INF, 1:1, 2:1, 3:1, and 4:1 represent the ratios for collaborative decoding.



Figure 12: Illustration of the results of collaborative decoding for each ratio after swapping for each layer in the Llama 3.1 model. The numbers 1, 2, 3, and 4 layer indicate the respective pairs of layer that have been swapped, while INF, 1:1, 2:1, 3:1, and 4:1 represent the ratios for collaborative decoding.



Figure 13: Illustration of the results of collaborative decoding for each ratio after swapping for each layer in the Mistral v0.3 model. The numbers 1, 2, 3, and 4 layer indicate the respective pairs of layer that have been swapped, while INF, 1:1, 2:1, 3:1, and 4:1 represent the ratios for collaborative decoding.



Figure 14: Illustration of the results of collaborative decoding for each ratio after swapping for each layer in the Qwen 2.5 model. The numbers 1, 2, 3, and 4 layer indicate the respective pairs of layer that have been swapped, while INF, 1:1, 2:1, 3:1, and 4:1 represent the ratios for collaborative decoding.

B Comparison of generation results using only the Attacked LLM after layer swapping versus generation results using Collaborative Decoding with the original model.



Figure 15: The result generated after swapping six layers in the attacked LLM. It can be observed that reasoning and comprehension abilities have declined due to layer swapping, and the model fails to properly understand the prompt.



Figure 16: The result of 6 layers of swapping followed by 1:1 collaborative decoding. We can see that it clearly understands the prompt and gives an answer.