

SPARTA: Evaluating Reasoning Segmentation Robustness through Black-Box Adversarial Paraphrasing in Text Autoencoder Latent Space

Viktoriia Zinkovich^{1*} Anton Antonov^{1*} Andrei Spiridonov^{1*}
Denis Shepelev^{1,3} Andrey Moskalenko^{1,2,3,4} Daria Pugacheva^{5,7}
Elena Tutubalina^{5,6,7} Andrey Kuznetsov^{1,8} Vlad Shakhuro^{1,2,3†}

¹FusionBrain Lab ²Lomonosov Moscow State University ³NUST MISIS

⁴IAI MSU ⁵HSE University ⁶Kazan Federal University

⁷Domain-specific NLP Group ⁸Innopolis Univeristy

*Equal contribution †Project leader

<https://github.com/emb-ai/SPARTA> ✉ viktoriia.zinkovich@gmail.com

Abstract

Multimodal large language models (MLLMs) have shown impressive capabilities in vision-language tasks such as *reasoning segmentation*, where models generate segmentation masks based on textual queries. While prior work has primarily focused on perturbing image inputs, *semantically equivalent* textual paraphrases—crucial in real-world applications where users express the same intent in varied ways—remain underexplored. To address this gap, we introduce a novel *adversarial paraphrasing task*: generating grammatically correct paraphrases that preserve the original query meaning while degrading segmentation performance. To evaluate the quality of adversarial paraphrases, we develop a comprehensive automatic evaluation protocol validated with human studies. Furthermore, we introduce **SPARTA**—a black-box, sentence-level optimization method that operates in the low-dimensional semantic latent space of a text autoencoder, guided by reinforcement learning. SPARTA achieves significantly higher success rates, outperforming prior methods by up to $2\times$ on both the ReasonSeg and LLMSeg-40k datasets. We use SPARTA and competitive baselines to assess the robustness of advanced reasoning segmentation models. We reveal that they remain vulnerable to adversarial paraphrasing—even under strict semantic and grammatical constraints. All code and data will be released publicly upon acceptance.

1 Introduction

In recent years, foundation models have achieved significant advances across diverse domains of deep learning. Advances in image classification (Dosovitskiy et al., 2021; Liu et al., 2022; Woo et al., 2023) and interactive segmentation (Kirillov et al., 2023; Ravi et al., 2024), together with progress in large language models (LLMs) (Brown et al., 2020; Touvron et al., 2023; Guo et al., 2025;

Original -- SPARTA --> Adversarial

What **might** the young girl **be** using to eat her dessert?



What **is it that** the young girl **is** using to eat her dessert?



Figure 1: Example of an adversarial paraphrase generated by our proposed SPARTA method. The SPARTA produces grammatically correct paraphrases that preserve the original semantic content while significantly degrading segmentation performance.

Dubey et al., 2024), have paved the way for multimodal large language models (MLLMs) (Liu et al., 2023, 2024a; Bai et al., 2023; Wang et al., 2024b; Li et al., 2023; Peng et al., 2023; Lai et al., 2024) that seamlessly integrate vision and language. These models are now integral to diverse applications, including conversational systems like ChatGPT, autonomous driving (Mu et al., 2024; Seff et al., 2023; Hwang et al., 2024), and robot control (Driess et al., 2023; Brohan et al., 2022, 2023; Black et al., 2024). As these models continue to mature, new tasks are emerging—particularly in robotics—that require sophisticated visual perception and reasoning capabilities. One such task is *reasoning segmentation* (Lai et al., 2024), where a model outputs a binary segmentation mask driven by an implicit text query that requires intricate logical or contextual interpretation.

The quality of a model’s predicted segmentation mask is expected to remain consistent, even when users paraphrase their prompts, preserving the orig-

inal meaning and intent. However, the robustness of reasoning segmentation models to query *paraphrasing* remains largely unexplored.

To address this problem, we propose a novel task, *adversarial paraphrasing*, which constrains textual perturbations according to the following criteria: (1) the core meaning of the original prompt must be preserved; (2) the paraphrase must remain grammatically correct; and (3) it must lead to a degradation in segmentation mask predictions. These constraints enable us to evaluate the robustness of state-of-the-art reasoning-based segmentation models against adversarially paraphrased queries.

Based on this task, we construct a new benchmark to systematically evaluate the robustness of reasoning segmentation models. We assess state-of-the-art attack strategies, including gradient-based and LLM-based methods; however, these approaches have notable limitations. Gradient-based methods often produce ungrammatical text (Guo et al., 2021; Jones et al., 2023), while LLM-based attacks typically rely on heuristic methods, such as iterative refinement (Chao et al., 2024).

To overcome the limitations of existing gradient-based and heuristic methods, **we introduce SPARTA**—a novel black-box sentence-level optimization method (Figure 1). SPARTA projects queries into a low-dimensional semantic latent space of a pretrained autoencoder and employs reinforcement learning to identify nearby vectors that yield effective adversarial paraphrases.

Overall, our contributions are as follows:

- We introduce a novel *adversarial paraphrasing* task designed to evaluate the robustness of reasoning segmentation models against semantically equivalent paraphrased queries. To facilitate this evaluation, we introduce an automated evaluation protocol. We conduct a user study and demonstrate strong alignment of the proposed scoring method with human judgment.
- We present a new method for generating adversarial paraphrases, leveraging reinforcement learning-based sentence-level optimization. Our method outperforms black-box and white-box baselines by up to $2\times$ on both the *ReasonSeg* and *LLMSeg-40k* datasets, with 2 model-specific exceptions.
- We conduct comprehensive experiments to assess the robustness of state-of-the-art reason-

ing segmentation models under both white-box and black-box adversarial paraphrasing settings. Our results indicate that, despite the strict semantic and grammatical constraints, existing reasoning segmentation models remain vulnerable to such attacks.

2 Related Work

2.1 Reasoning Segmentation

In Referring Expression Segmentation (RES), models output segmentation masks from textual descriptions (Zou et al., 2023; Rasheed et al., 2024; Wu et al., 2024b,a; Wang et al., 2024c; Liu et al., 2024b). Expanding on RES, the *reasoning segmentation task* was introduced to handle prompts requiring world knowledge and logical reasoning (Lai et al., 2024).

The pioneering reasoning segmentation model LISA employs an embedding-as-mask paradigm, decoding a <SEG> token via SAM to produce a segmentation mask. Several LISA-based models followed, such as LISA++ (Yang et al., 2024), which can incorporate segmentation results into text responses, and GSVA (Xia et al., 2024), which introduces a <REJ> token to explicitly reject absent objects.

2.2 Adversarial Attacks on Text Modality

Evaluating model robustness often involves adversarial attacks, which are broadly categorized as white-box or black-box, based on the attacker’s access to model internals.

White-box attacks leverage gradient information to optimize adversarial paraphrases, addressing the challenges posed by the discrete nature of text through techniques such as Taylor expansion (Ebrahimi et al., 2018; Jones et al., 2023) and Gumbel-Softmax sampling (Jang et al., 2017; Guo et al., 2021). Among these, we consider two state-of-the-art methods: GBDA (Guo et al., 2021) and ARCA (Jones et al., 2023). While ARCA achieves strong attack success, it lacks semantic regularization, frequently inserting special symbols that undermine its suitability as a paraphrasing baseline. In contrast, GBDA incorporates semantic similarity constraints; however, it remains limited to token-level substitutions, constraining paraphrase diversity.

In the black-box setting, attacks have progressed from simple word- and character-level manipulations, such as synonym substitution (Jin et al., 2020;

Ren et al., 2019) and character edits (Gao et al., 2018), to methods that generate semantically equivalent paraphrases using transformer-based models (Li et al., 2020; Iyyer et al., 2018). Recent developments further leverage LLMs to generate more semantically diverse paraphrases (Yan et al., 2023; Xu et al., 2023). Among these, we consider PAIR (Chao et al., 2024)—a state-of-the-art and widely used adversarial method—and Qwen3-32B (Team, 2025), a leading LLM, as attack baselines. While these techniques often produce fluent paraphrases, they typically depend on heuristic rules or manual trial-and-error, lacking controllable optimization. To address this gap, we optimize paraphrases in the low-dimensional semantic latent space of a pretrained text autoencoder, leveraging reinforcement learning to maximize degradation of segmentation performance—thereby enabling more effective adversarial paraphrasing.

2.3 Evaluation of Adversarial Attacks

Adversarial attack effectiveness is primarily measured by the attack *success rate* (SR), where success is determined by the model’s output quality drop crossing a task-specific threshold. For instance, this could be a drop in Intersection over Union (IoU) for interactive segmentation (Liu et al., 2025; Huang et al., 2024b) or confidence score changes in classification task (Guo et al., 2021; Dong et al., 2019).

To ensure the validity and semantic consistency of adversarial paraphrases, we consider text quality metrics. Semantic preservation is commonly assessed via cosine similarity between embeddings of the original and paraphrased sentences (Guo et al., 2021; Thieu et al., 2021; Sun and Wang, 2024). In practice, many recent studies apply a cosine similarity threshold to filter paraphrases, with the cutoff depending on the embedding model used (Kassem and Saad, 2024; Herel et al., 2023). Furthermore, recent works leverage LLMs for evaluation through GPT-scoring (Fu et al., 2024; Wang et al., 2023; Chiang and Lee, 2023; Chan et al., 2023), though the reliability of such metrics remains an open question (Wang et al., 2024a). To offer a more trustworthy assessment of paraphrase quality, we combine cosine-based filtering with LLM-based scoring into a comprehensive evaluation pipeline.

3 Proposed Method: SPARTA

In this section, we introduce **SPARTA**—a novel black-box paraphrasing method that generates grammatically correct, semantically consistent paraphrases which degrade segmentation performance. The input query is first encoded into a continuous latent representation using a pretrained text autoencoder (Section 3.1). A set of candidate vectors is sampled from a Gaussian distribution in the latent space, centered at the original latent vector. These candidates are decoded into paraphrases and evaluated via a reward function that penalizes overlap with the original segmentation mask, while regularization ensures semantic fidelity. The policy is optimized via Proximal Policy Optimization to guide the sampling toward more effective adversarial paraphrases (Section 3.2). The full optimization pipeline is detailed in Algorithm 1.

3.1 Latent Sentence Space

Instead of searching paraphrases over discrete tokens, we operate in the *continuous* latent space of a pretrained text autoencoder (E, D) . The encoder E maps the input query \mathbf{x} to a continuous semantic space, and the decoder D reconstructs the original sentence from the latent vector \mathbf{z} :

$$\mathbf{z} = E(\mathbf{x}) \in \mathbb{R}^d, \quad \hat{\mathbf{x}} = D(\mathbf{z}), \quad \hat{\mathbf{x}} \approx \mathbf{x}. \quad (1)$$

We adopt SONAR (Duquenne et al., 2023), a 1B-parameter multilingual model, as the state-of-the-art text autoencoder (E, D) , which is described in detail in Appendix A. Its training objective includes translation and MSE losses on sentence embeddings, encouraging language-agnostic latent representations and a well-aligned cross-lingual embedding space.

For the original input query \mathbf{x}_0 , we obtain the initial embedding $\mathbf{z}_0 = E(\mathbf{x}_0)$ and optimize the latent vector $\mathbf{z} \in \mathbb{R}^d$, initialized with \mathbf{z}_0 .

3.2 Reinforcement Learning Formulation

SPARTA learns a stochastic policy $\pi_\theta(\mathbf{z} | \mathbf{z}_0)$ that perturbs the original latent vector \mathbf{z}_0 to generate adversarial paraphrases. The policy is modeled as a diagonal Gaussian distribution in latent space:

$$\mathbf{z} \sim \pi_\theta(\mathbf{z} | \mathbf{z}_0) = \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)), \quad (2)$$

where the learnable parameters $\theta = (\boldsymbol{\mu}, \boldsymbol{\sigma})$ consist of the mean $\boldsymbol{\mu} \in \mathbb{R}^d$ (initialized with \mathbf{z}_0) and the standard deviation $\boldsymbol{\sigma} \in \mathbb{R}^d$.¹

¹The scale is reparameterized as $\boldsymbol{\sigma} = \log(1 + \exp(\boldsymbol{\lambda}))$ to keep it strictly positive, where $\boldsymbol{\lambda} \in \mathbb{R}^d$ is trainable.

Algorithm 1 PPO in Latent Sentence Space

Require: Query \mathbf{x}_0 , image \mathbf{I} , ground-truth mask \mathbf{m} ; autoencoder (E, D) ; model f ; sample size n ; iteration number N ; hyperparams $(\epsilon, \lambda_{\text{sim}}, \lambda_{\text{adv}}, \boldsymbol{\sigma}, V_\psi)$

- 1: Initialize $\mathbf{z}_0 \leftarrow E(\mathbf{x}_0)$, $\boldsymbol{\mu} \leftarrow \mathbf{z}_0$
- 2: **for** $t = 1$ to N **do**
- 3: Sample $\mathbf{z}_i \sim \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$ for $i = 1..n$
- 4: **for** $i = 1$ to n **do**
- 5: Decode $\hat{\mathbf{x}}_i \leftarrow D(\mathbf{z}_i)$, predict $\hat{\mathbf{m}}_i \leftarrow f(\mathbf{I}, \hat{\mathbf{x}}_i)$
- 6: Compute R_i, A_i, ρ_i, l_i (Eqs. 3–6)
- 7: **end for**
- 8: $\mathcal{L}_{\text{policy}} \leftarrow -\lambda_{\text{adv}} \frac{1}{n} \sum l_i$; $\mathcal{L}_{\text{value}} \leftarrow \frac{1}{n} \sum (R_i - V_\psi)^2$; $\mathcal{L}_{\text{sim}} \leftarrow \lambda_{\text{sim}} \|\boldsymbol{\mu} - \mathbf{z}_0\|^2$
- 9: Update $\boldsymbol{\mu}, \boldsymbol{\sigma}, \psi$ via Adam on $\mathcal{L}_{\text{final}}$ (Eq. 7)
- 10: Update old policy: $\pi_{\text{old}} \leftarrow \pi$
- 11: Save $\hat{\mathbf{x}} \leftarrow D(\boldsymbol{\mu})$
- 12: **end for**

Reward For each sampled vector \mathbf{z} , we generate a candidate paraphrase $\hat{\mathbf{x}} = D(\mathbf{z})$ and pass it through the attacked reasoning segmentation model f . The model outputs a segmentation mask $\hat{\mathbf{m}}$, and the effectiveness of the adversarial paraphrase is quantified by the following *reward*:

$$R = -\text{IoU}(\hat{\mathbf{m}}, \mathbf{m}), \quad (3)$$

where \mathbf{m} is the ground truth mask. Higher rewards correspond to lower Intersection-over-Union, thereby encouraging paraphrases that most effectively degrade model performance.

Baseline and Advantage To reduce variance in gradient estimates, we learn a scalar value network V_ψ as a *baseline* (Sutton and Barto, 2018). The *advantage* A is then the normalized difference between the observed reward R and the baseline $V_\psi(\mathbf{z})$:

$$A = \frac{R - V_\psi - \mathbb{E}[R - V_\psi]}{\text{Std}[R - V_\psi] + \epsilon}. \quad (4)$$

Optimization via PPO To train the latent-space policy, we employ the standard *clipped surrogate objective* of Proximal Policy Optimization (PPO) (Schulman et al., 2017; Huang et al., 2024a). At each update, we sample a batch of n candidate embeddings $\{\mathbf{z}_i\}$ from the old policy $\pi_{\theta_{\text{old}}}$, decode each into a paraphrase, and evaluate its adversarial reward R_i . For each sample, we compute the importance weight:

$$\rho_i = \frac{\pi_\theta(\mathbf{z}_i)}{\pi_{\theta_{\text{old}}}(\mathbf{z}_i)} = \exp(\log \pi_\theta(\mathbf{z}_i) - \log \pi_{\theta_{\text{old}}}(\mathbf{z}_i)), \quad (5)$$

then form the clipped surrogate:

$$l_i = \min(\rho_i A_i, \text{clip}(\rho_i, 1 - \epsilon, 1 + \epsilon) A_i). \quad (6)$$

Here, $\epsilon = 0.2$ is the *clip ratio* hyperparameter, which constrains the policy update to a trust region $[1 - \epsilon, 1 + \epsilon]$. Clipping ρ_i prevents large updates that could destabilize training (Schulman et al., 2015, 2017).

Objective function The final optimization objective $\mathcal{L}_{\text{final}}$ combines three terms:

$$-\underbrace{\lambda_{\text{adv}} \sum_{i=1}^n \frac{l_i}{n}}_{\mathcal{L}_{\text{policy}}} + \underbrace{\sum_{i=1}^n \frac{(R_i - V_\psi)^2}{n}}_{\mathcal{L}_{\text{value}}} + \underbrace{\lambda_{\text{sim}} \|\boldsymbol{\mu} - \mathbf{z}_0\|_2^2}_{\mathcal{L}_{\text{sim}}}, \quad (7)$$

where $\mathcal{L}_{\text{value}}$ trains the baseline and \mathcal{L}_{sim} preserves semantic fidelity to the original query. Optimization is performed with Adam using separate learning rates for $\boldsymbol{\mu}$, $\boldsymbol{\sigma}$, and ψ .

4 Proposed Evaluation Protocol

In this section, we introduce an automatic evaluation protocol for our novel adversarial paraphrasing task. We begin by outlining the main steps of the protocol (Section 4.1). We then examine the challenges associated with its core component, LLM-based paraphrase detection, and propose additional filtering steps to enhance performance (Section 4.2). Finally, we evaluate the detection methods and ablate the proposed improvements through human studies (Section 4.3).

4.1 Evaluation Protocol

We introduce an *automatic evaluation protocol for the adversarial paraphrasing task*. Since existing attack methods may produce invalid outputs, we select the best adversarial prompts—based on attack loss and paraphrasing quality—and use them to evaluate attack performance. Specifically, given a set of adversarial prompts obtained through an attack over N iterations, we proceed as follows: (1) remove duplicate prompts; (2) discard any prompt that does not reduce the segmentation model’s IoU; (3) detect which prompts are valid paraphrases; (4)

Type	Text
Original	the youngest person
PAIR paraphrase	considering standard human growth patterns, pinpoint the individual who, if all people in the image were lined up in order of birth, would be positioned closest to the beginning of the sequence
Original	the sauce
PAIR paraphrase	the component that is neither the main ingredient nor the garnish, but is distributed throughout the plate in a somewhat fluid form

Table 1: **Examples of PAIR-generated paraphrases that are overly verbose or abstract.** The first paraphrase employs indirect and wordy language, while the second describes the sauce ambiguously without explicitly naming it, leaving it unclear whether it refers to sauce, oil, or dressing. Despite this, LLMs rate such paraphrases as valid.

Prompt	Qwen3						LLaMA-3.1-Nemotron					
	LLM			LLM & RegExp & CosSim			LLM			LLM & RegExp & CosSim		
	Precision	Recall	F-score	Precision	Recall	F-score	Precision	Recall	F-score	Precision	Recall	F-score
1	0.480	0.964	0.641	0.623	0.865	0.725	0.634	0.640	0.637	0.798	0.604	0.687
2	0.480	0.991	0.647	0.601	0.883	0.715	0.520	0.712	0.601	0.664	0.640	0.651
3	0.530	0.946	0.680	0.671	0.847	0.749	0.552	0.712	0.622	0.726	0.622	0.670

Table 2: **Evaluation of paraphrase detection methods.** We compare LLM-based detection using the baseline system prompt 1 from Michail et al. (2025) against our proposed enhanced system prompts (2 and 3), as well with additional filtering based on regular expressions and semantic cosine similarity. Best F-scores are shown in **bold**.

select the paraphrase that yields the greatest relative IoU drop.

A critical step in this evaluation protocol is *paraphrase detection*, which, as we demonstrate in the following section, presents significant challenges.

4.2 LLM-based Paraphrase Detection Issues

Paraphrasing involves rephrasing a sentence while preserving its original meaning, intent, and grammatical correctness in a clear and concise manner. However, automatically assessing whether generated prompts meet these criteria remains a non-trivial task. To address this, we explored a state-of-the-art LLM-based evaluation approach, following prior work Michail et al. (2025). Our initial experiments revealed three key issues:

1. Defining a *valid* paraphrase for an LLM is challenging, as commonly accepted definitions like “alternative expressions of the same meaning” (Xu et al., 2015) are too broad for reliable automated evaluation.
2. LLMs often fail to capture differences in capitalization and terminal punctuation (e.g., “a person is calling someone” vs. “A person is calling someone.”). Because ReasonSeg dataset contain prompts that may be either fragments or complete sentences, we consider

an adversarial prompt to be a valid paraphrase only if it preserves both capitalization and terminal punctuation.

3. We observe that some paraphrases become excessively long or abstract, occasionally resembling riddles or puzzles, which LLMs often still judge as valid (Table 1). Although such paraphrases may retain partial semantic overlap with the original, they obscure the intended meaning and hinder clarity, and thus should not be regarded as valid.

Our findings are consistent with the recent work Michail et al. (2025), which demonstrated that even modern LLMs and specialized classification models struggle with the paraphrasing classification task.

We address the issues mentioned above as follows:

1. To mitigate Issue 1, we *improve the system prompt* used by the LLM. We consider three different system prompts. Prompt 1 is a simple zero-shot binary classification prompt, which performed best in prior work Michail et al. (2025). Prompts 2 and 3 provide detailed task instructions, a 5-point scoring scale, and 10 in-context examples. In the latter two settings, we consider an adversarial prompt a

valid paraphrase only if it receives an LLM score of 5. The full prompt templates are included in Appendix E.1.

2. To resolve Issue 2, we apply a *regular expression-based filtering* to discard paraphrases that alter capitalization or terminal punctuation.
3. We mitigate Issue 3 by filtering out semantically distant paraphrases. Specifically, we use Qwen3-Embedding-8B (Zhang et al., 2025), a state-of-the-art open-source sentence embedding model, to compute semantic similarity. Through empirical analysis, we identify an optimal *cosine similarity* threshold of 0.825 (see Appendix E.3 for details). This step improves detection performance by removing overly abstract or indirect prompts.

We evaluated LLM-based detection and ablated our improvements with human studies.

4.3 Ablation

We sampled a dataset of 310 pairs of original and adversarial prompts, generated by the proposed SPARTA and baseline methods (see Section 5.3), and manually annotated them for paraphrase validity. For LLM-based detection, we evaluated two state-of-the-art models: LLaMA-3.1-Nemotron-70B (Wang et al., 2024d) and Qwen3-32B (Team, 2025). For each LLM, system prompt, and filtering configuration (with or without regular expressions and cosine similarity), we measured performance using the F1-score to identify the most effective detection setup.

The results of the human study are summarized in Table 2. The best detection performance was achieved using Qwen3-32B with system prompt 3 and filtering based on regular expressions and cosine similarity, yielding an F1 score of 0.749. Using system prompt 3 without additional filtering already led to a notable improvement over system prompt 1 (F1 score: 0.641 \rightarrow 0.680), and further gains were achieved with the full filtering setup (0.680 \rightarrow 0.749). This best-performing configuration was adopted in the automatic evaluation protocol described earlier.

5 Implementation Details

5.1 Datasets

We use the *ReasonSeg* dataset (Lai et al., 2024), which has become a standard benchmark for evalu-

ating reasoning segmentation models. Additionally, we leverage *LLM-Seg40K* (Wang and Ke, 2024), the latest large-scale reasoning segmentation dataset collected using ChatGPT-4. With an average query length of 15.2 words, LLM-Seg40K presents more challenging scenarios and greater linguistic complexity. Due to computational constraints, we limit our evaluation to 300 samples from each dataset (see Appendix H for details).

5.2 Reasoning Models

We evaluated 6 checkpoints of 3 modern reasoning segmentation models. Our particular interest is LISA (Lai et al., 2024), the first and most widely adopted model in this domain. We also tested LISA’s successors, LISA++ (Yang et al., 2024) and GSVA (Xia et al., 2024), which are often used as strong baselines in reasoning and referring segmentation.

5.3 Attack Baselines

We consider the following attack baselines (see Appendix B for details): (1) **GBDA** (Guo et al., 2021): adapted from text-only adversarial attacks to the multimodal setting through hyperparameter tuning; (2) **Qwen3-32B, simple prompt** (Team, 2025): a naive baseline that prompts the model to paraphrase the input sentence; (3) **PAIR** (Chao et al., 2024): an advanced, iterative method that was repurposed from LLM jailbreaking with a paraphrasing-specific prompt and Qwen3-32B as the language model.

To assess overall robustness, we further introduce a *unified attack* that, for each sample, selects the most effective paraphrase from all baselines and our SPARTA method.

6 Experimental Results

We evaluate adversarial attack performance using the following procedure: (1) for each dataset sample, we generate an adversarial paraphrase and compute its *relative IoU* degradation (ΔIoU , %) and its LLM-score, following the evaluation protocol described in Section 4; (2) we construct the attack *success rate curve* SR_θ , where θ denotes the threshold for ΔIoU ; (3) we report the area under the SR_θ curve (mSR), as well as the success rates at $\theta = 5\%$ (SR_5) and $\theta = 10\%$ (SR_{10}).

In step 2, an adversarial paraphrase is considered successful for a given threshold θ if it achieves $\Delta\text{IoU} \geq \theta$ and is rated as valid by the evaluation protocol (i.e., LLM-score = 5).

Attacked model	GBDA			Qwen3 (<i>simple</i>)			Qwen3 PAIR			SPARTA (<i>ours</i>)		
	mSR	SR ₅	SR ₁₀	mSR	SR ₅	SR ₁₀	mSR	SR ₅	SR ₁₀	mSR ↑	SR ₅ ↑	SR ₁₀ ↑
LISA [7B]	3.2	11.0	8.4	<u>13.5</u>	<u>30.9</u>	<u>25.1</u>	13.0	25.7	22.5	26.6	48.7	42.4
LISA-exp. [7B]	3.1	10.0	7.5	11.0	25.0	21.0	<u>16.1</u>	<u>32.5</u>	<u>25.5</u>	24.6	49.5	42.5
LISA [13B]	2.6	4.5	4.5	9.2	24.1	18.8	<u>11.0</u>	<u>26.8</u>	<u>21.0</u>	23.2	46.0	38.4
LISA-exp. [13B]	2.9	7.6	6.2	8.7	24.6	18.5	<u>11.4</u>	<u>27.5</u>	<u>21.8</u>	25.0	47.4	40.8
LISA++ [7B]	0.9	2.1	1.7	<u>10.7</u>	<u>20.9</u>	<u>19.7</u>	8.8	17.4	14.0	16.2	29.1	23.9
GSVA [13B]	2.2	6.4	4.6	15.6	28.3	23.1	<u>16.0</u>	<u>31.8</u>	<u>29.5</u>	27.9	53.2	44.5

Table 3: **Evaluation results of baselines and the proposed SPARTA on state-of-the-art reasoning segmentation models on the LLMSeg-40k dataset.** mSR refers to the area under the curve of the success rate (SR) versus the IoU-drop threshold, computed for adversarial paraphrases with an LLM score of 5. SR₅ and SR₁₀ represent the success rate for IoU drops greater than 5% and 10%, respectively. Higher values indicate stronger attacks. The best results are in **bold**, the second best are underlined.

Attacked model	GBDA			Qwen3 (<i>simple</i>)			Qwen3 PAIR			SPARTA (<i>ours</i>)		
	mSR	SR ₅	SR ₁₀	mSR	SR ₅	SR ₁₀	mSR	SR ₅	SR ₁₀	mSR ↑	SR ₅ ↑	SR ₁₀ ↑
LISA [7B]	6.8	18.7	12.6	11.2	<u>32.9</u>	<u>24.9</u>	<u>14.6</u>	27.6	23.4	25.8	47.4	40.8
LISA-exp. [7B]	2.8	11.2	5.8	12.5	<u>29.0</u>	<u>24.1</u>	14.6	26.8	23.2	14.0	32.3	26.9
LISA [13B]	1.4	5.7	4.9	9.6	21.5	17.0	<u>13.3</u>	<u>25.1</u>	<u>23.5</u>	16.3	42.7	33.3
LISA-exp. [13B]	3.4	8.8	6.7	7.6	<u>25.4</u>	18.8	<u>11.3</u>	25.0	<u>19.9</u>	17.5	39.9	32.8
LISA++ [7B]	2.5	8.4	5.4	9.7	22.8	16.6	21.1	36.0	31.8	<u>15.4</u>	<u>28.2</u>	<u>23.1</u>
GSVA [13B]	6.1	15.3	14.4	13.1	<u>28.2</u>	23.0	<u>15.1</u>	27.8	<u>25.4</u>	22.7	46.2	37.0

Table 4: **Evaluation results of baselines and the proposed SPARTA on state-of-the-art reasoning segmentation models on the ReasonSeg dataset.** mSR refers to the area under the curve of the success rate (SR) versus the IoU-drop threshold, computed for adversarial paraphrases with an LLM score of 5. SR₅ and SR₁₀ represent the success rate for IoU drops greater than 5% and 10%, respectively. Higher values indicate stronger attacks. The best results are in **bold**, the second best are underlined.

The resulting SR curves and metrics are presented in Table 3, Table 4, and Figure 2. The mSR measures the average success rate of an adversarial attack over all IoU thresholds θ , reflecting overall attack effectiveness. In Tables 3 and 4, higher values indicate stronger attacks.

Table 5 presents the robustness of each model against adversarial paraphrasing obtained by the unified attack. Here, lower values indicate greater robustness.

7 Discussion

A review of Tables 3–4, Table 5, and Figure 2 leads to the following conclusions. First, **the proposed SPARTA attack consistently outperforms** all baselines across reasoning segmentation models on LLMSeg-40k, achieving an average mSR improvement of about 84% over the strongest baselines. On ReasonSeg SPARTA continues to improve attack performance (average mSR gain = 29%), except 2 models (LISA-exp. [7B] and LISA++ [7B]), where

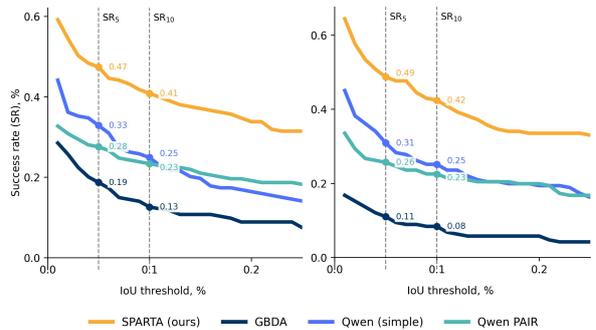


Figure 2: **Success rate (SR) as a function of IoU-drop threshold for adversarial paraphrases with LLM score 5.** Results are shown for the LISA-7B model on the ReasonSeg dataset (left) and LLMSeg-40k dataset (right).

PAIR yields higher mSR. Examples of adversarial paraphrases generated by the proposed SPARTA method are presented in Figure 3. Additionally, an analysis of attack failures is provided in Appendix G.

Second, the proposed **adversarial paraphrasing task presents a significant challenge** for cur-

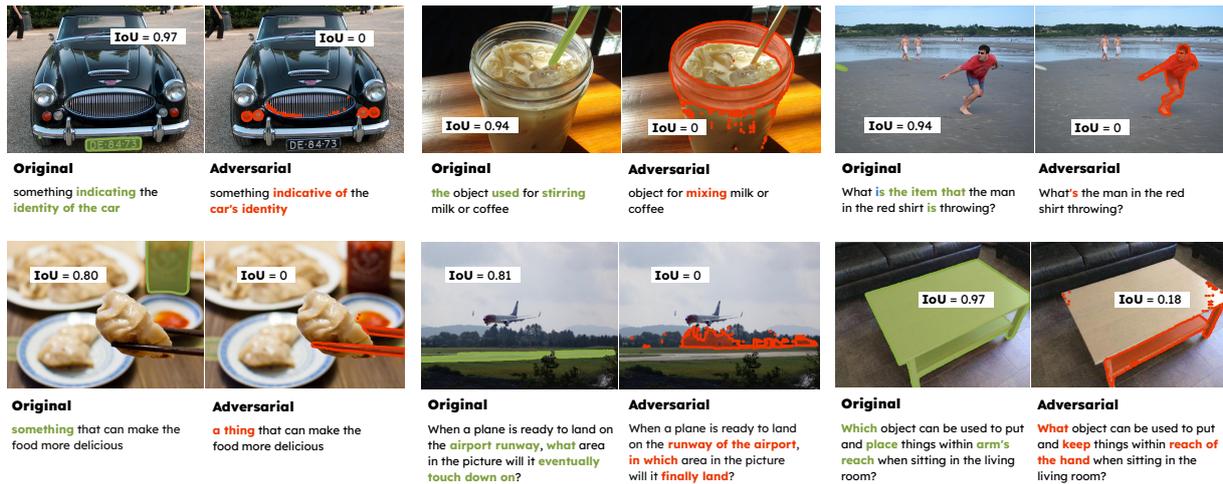


Figure 3: Examples of adversarial paraphrases obtained using the proposed SPARTA method. SPARTA produces grammatically correct paraphrases that preserve the original query meaning while substantially degrading segmentation performance.

rent reasoning segmentation models. Our task introduces strict constraints on grammatical correctness and semantic equivalence, making it a significantly more difficult benchmark for evaluating model robustness in real-world scenarios. Despite these constraints, our unified attack achieves success rates of up to 68% at a 10% relative IoU drop threshold, indicating that **current reasoning segmentation models remain vulnerable to well-crafted adversarial paraphrases**.

Finally, while a deeper analysis is left for future work, unified attacks already offer valuable insights into robustness differences across models (Table 5). Notably, LISA++ [7B] demonstrates the highest robustness on LLMSeg-40k, while LISA-exp. [13B] achieves the highest robustness on ReasonSeg. This indicates that **there is currently no reasoning segmentation model that is optimal in terms of robustness on both datasets**. LISA [13B] consistently outperforms its 7B variant, suggesting that **increased model capacity enhances resistance to adversarial paraphrasing**. In contrast, GSVA [13B] shows the weakest robustness on LLMSeg-40k, which we attribute to its lower segmentation performance; our evaluation of the released checkpoint revealed significantly lower metrics than reported in the prior work.

8 Conclusion

In this work, we introduced a novel challenging task that involves generating semantically consistent and grammatically correct paraphrases that sig-

Attacked model	ReasonSeg (test)			LLMSeg-40k (val)		
	mSR ↓	SR ₅ ↓	SR ₁₀ ↓	mSR ↓	SR ₅ ↓	SR ₁₀ ↓
LISA [7B]	43.6	78.5	68.7	38.8	66.0	58.1
LISA-exp. [7B]	33.2	67.0	<u>55.4</u>	36.6	69.5	59.5
LISA [13B]	<u>30.1</u>	66.0	57.5	33.2	<u>60.3</u>	53.1
LISA-exp. [13B]	29.6	<u>64.6</u>	54.6	<u>32.4</u>	60.7	<u>52.6</u>
LISA++ [7B]	37.1	63.9	56.0	26.3	44.7	40.0
GSVA [13B]	40.3	72.8	63.8	40.6	68.2	62.4

Table 5: **Robustness of state-of-the-art reasoning segmentation models** to unified attack. mSR refers to the area under the curve of the success rate (SR) versus the IoU-drop threshold, computed for adversarial paraphrases with an LLM score of 5. SR₅ and SR₁₀ represent the success rate for IoU drops greater than 5% and 10%, respectively. Lower values indicate greater model robustness. The best results are in **bold**, the second best are underlined.

nificantly degrade segmentation performance. To address this task, we proposed SPARTA, which leverages a black-box, sentence-level optimization in the semantic latent space of the pretrained text autoencoder, guided by reinforcement learning. Through comprehensive automatic and human-validated evaluation protocols, we demonstrate that SPARTA outperforms state-of-the-art baselines, achieving up to a $2\times$ improvement on LLMSeg-40k; on ReasonSeg, it is better for all but two models. Despite strict semantic and grammatical constraints, our findings reveal that current reasoning segmentation models remain vulnerable to adversarial paraphrasing. We believe this work offers a valuable foundation for future research on evaluating and enhancing the robustness of multimodal vision-language systems.

9 Limitations and Future Work

While the proposed SPARTA method outperforms state-of-the-art baselines, several limitations remain. First, neither SPARTA nor existing attacks guarantee that generated prompts are valid paraphrases. To mitigate this, our evaluation protocol selects the best valid adversarial prompts after generation, though future work could explore incorporating validity constraints directly into the generation process.

Second, proposed evaluation protocol was validated only on a small dataset of 310 pairs of original and adversarial prompts, generated by the proposed SPARTA and baseline methods, and manually annotated them for paraphrase validity. While this is a modest-sized validation set, we consider it sufficient for this setup-selection purpose.

Third, while SPARTA generates paraphrases that are semantically and grammatically correct, some may appear unnatural to human users. This reflects broader limitations of current text autoencoders (see Appendix A), and future improvements likely depend on developing models with more structured and human-aligned latent spaces.

Fourth, while SPARTA achieves strong performance and consistently outperforms existing baselines, this comes with increased computational cost. This represents a performance-efficiency trade-off, however, the overhead can be controlled by adjusting the number of optimization iterations based on available resources.

Fifth, we focus solely on attack methods, without addressing potential defenses. Exploring robustness strategies for reasoning segmentation models is a critical next step toward building more reliable multimodal systems.

Finally, the core idea behind SPARTA is task-agnostic and can be extended to other settings, for example, paraphrase-based adversarial attacks on text classification models in pure NLP or attacks on vision language action (VLA) models. We leave the exploration of these directions to future work.

10 Ethical Considerations

Our work introduces a novel adversarial paraphrasing method to evaluate the robustness of reasoning segmentation models. While this method could potentially be misused to attack real-world models, we believe the benefits to the research community outweigh these risks. By uncovering current vulnerabilities, we aim to encourage the development of

more robust, interpretable, and trustworthy systems. To support responsible research, we will release all code and data under a research-only license, strictly intended for academic and non-commercial use.

11 Acknowledgments

We acknowledge the computational resources of the HPC facilities at HSE University.

We used ChatGPT to check grammar and refine the clarity of the text. The authors reviewed and revised all AI-generated output, and we take full responsibility for the content.

References

- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*.
- Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, and 1 others. 2024. π_0 : A Vision-Language-Action Flow Model for General Robot Control. *arXiv preprint*.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, and 1 others. 2023. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, and 1 others. 2022. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate. *Preprint*, arXiv:2308.07201.
- Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. 2024. Jailbreaking black box large language models in twenty queries. *Preprint*, arXiv:2310.08419.
- Cheng-Han Chiang and Hung-yi Lee. 2023. Can large language models be an alternative to human evaluations? In *Proceedings of the 61st Annual Meeting of*

- the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15607–15631. Association for Computational Linguistics.
- Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. 2019. Evading defenses to transferable adversarial examples by translation-invariant attacks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4312–4321.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. [An image is worth 16x16 words: Transformers for image recognition at scale](#). In *International Conference on Learning Representations*.
- Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, and 1 others. 2023. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Paul-Ambroise Duquenne, Holger Schwenk, and Benoît Sagot. 2023. Sonar: sentence-level multimodal and language-agnostic representations. *arXiv e-prints*, pages arXiv–2308.
- Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2018. HotFlip: White-box adversarial examples for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 3F1–3F6. Association for Computational Linguistics.
- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2024. GPTScore: Evaluate as you desire. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6556–6576. Association for Computational Linguistics.
- Ji Gao, Jack Lanchantin, Mary Lou Soffa, and Yanjun Qi. 2018. [Black-box generation of adversarial text sequences to evade deep learning classifiers](#). *Preprint*, arXiv:1801.04354.
- Chuan Guo, Alexandre Sablayrolles, Hervé Jégou, and Douwe Kiela. 2021. Gradient-based adversarial attacks against text transformers. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5747–5757. Association for Computational Linguistics.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- David Herel, Hugo Cisneros, and Tomas Mikolov. 2023. [Preserving semantics in textual adversarial attacks](#). In *26th European Conference on Artificial Intelligence (ECAI 2023)*, volume 372 of *Frontiers in Artificial Intelligence and Applications*, pages 1036–1043. IOS Press.
- Nai-Chieh Huang, Ping-Chun Hsieh, Kuo-Hao Ho, and I-Chen Wu. 2024a. [Ppo-clip attains global optimality: Towards deeper understandings of clipping](#). *Preprint*, arXiv:2312.12065.
- Shize Huang, Qianhui Fan, Zhaoxin Zhang, Xiaowen Liu, Guanqun Song, and Jinzhe Qin. 2024b. [Segment shards: Cross-prompt adversarial attacks against the segment anything model](#). *Applied Sciences*.
- Jyh-Jing Hwang, Runsheng Xu, Hubert Lin, Wei-Chih Hung, Jingwei Ji, Kristy Choi, Di Huang, Tong He, Paul Covington, Benjamin Sapp, and 1 others. 2024. Emma: End-to-end multimodal model for autonomous driving. *arXiv preprint arXiv:2410.23262*.
- Mohit Iyyer, John Wieting, Kevin Gimpel, and Luke Zettlemoyer. 2018. [Adversarial example generation with syntactically controlled paraphrase networks](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1875–1885, New Orleans, Louisiana. Association for Computational Linguistics.
- Eric Jang, Shixiang Gu, and Ben Poole. 2017. Categorical reparameterization with gumbel-softmax. In *International Conference on Learning Representations*.
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2020. [Is bert really robust? a strong baseline for natural language attack on text classification and entailment](#). *Preprint*, arXiv:1907.11932.
- Erik Jones, Anca Dragan, Aditi Raghunathan, and Jacob Steinhardt. 2023. Automatically auditing large language models via discrete optimization. In *International Conference on Machine Learning*, pages 15307–15329.
- Aly M. Kassem and Sherif Saad. 2024. [Finding a needle in the adversarial haystack: A targeted paraphrasing approach for uncovering edge cases with minimal distribution distortion](#). *Preprint*, arXiv:2401.11373.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar, and Ross Girshick. 2023. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4015–4026.

- Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. 2024. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR.
- Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020. Bert-attack: Adversarial attack against bert using bert. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6193–6202. Association for Computational Linguistics.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916.
- Xiaoliang Liu, Furao Shen, and Jian Zhao. 2025. [Region-guided attack on the segment anything model \(sam\)](#). *Preprint*, arXiv:2411.02974.
- Yong Liu, Cairong Zhang, Yitong Wang, Jiahao Wang, Yujiu Yang, and Yansong Tang. 2024b. Universal segmentation at arbitrary granularity with language instruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3459–3469.
- Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. 2022. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986.
- Andrianos Michail, Simon Clematide, and Juri Opitz. 2025. [PARAPHRASUS: A comprehensive benchmark for evaluating paraphrase detection models](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 8749–8762, Abu Dhabi, UAE. Association for Computational Linguistics.
- Norman Mu, Jingwei Ji, Zhenpei Yang, Nate Harada, Haotian Tang, Kan Chen, Charles R. Qi, Runzhou Ge, Kratarth Goel, Zoey Yang, Scott Ettinger, Rami Al-Rfou, Dragomir Anguelov, and Yin Zhou. 2024. [Most: Multi-modality scene tokenization for motion prediction](#). *Preprint*, arXiv:2404.19531.
- Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*.
- Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S Khan. 2024. Glamm: Pixel grounding large multimodal model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13009–13018.
- Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, and 1 others. 2024. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*.
- Shuhuai Ren, Yihe Deng, Kun He, and Wanxiang Che. 2019. [Generating natural language adversarial examples through probability weighted word saliency](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1085–1097, Florence, Italy. Association for Computational Linguistics.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. 2015. Trust region policy optimization. *arXiv preprint arXiv:1502.05477*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. [Proximal policy optimization algorithms](#). *CoRR*.
- Ari Seff, Brian Cera, Dian Chen, Mason Ng, Aurick Zhou, Nigamaa Nayakanti, Khaled S. Refaat, Rami Al-Rfou, and Benjamin Sapp. 2023. [Motionlm: Multi-agent motion forecasting as language modeling](#). *Preprint*, arXiv:2309.16534.
- Kun Sun and Rong Wang. 2024. [Textual similarity as a key metric in machine translation quality estimation](#). *Preprint*, arXiv:2406.07440.
- Richard S. Sutton and Andrew G. Barto. 2018. *Reinforcement Learning: An Introduction*, 2nd edition. MIT Press.
- Qwen Team. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.
- Thanh Thieu, Ha Do, Thanh Duong, Shi Pu, Sathyanarayanan Aakur, and Saad Khan. 2021. Lexdivpara: A measure of paraphrase quality with integrated sentential lexical complexity. In *Proceedings of the Intelligent Systems Conference (IntelliSys 2021)*, volume 296 of *Lecture Notes in Networks and Systems*, pages 1–10.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. 2023. Is ChatGPT a good NLG

- evaluator? a preliminary study. In *Proceedings of the 4th New Frontiers in Summarization Workshop*, pages 1–11. Association for Computational Linguistics.
- Junchi Wang and Lei Ke. 2024. [Llm-seg: Bridging image segmentation and large language model reasoning](#). *Preprint*, arXiv:2404.08767.
- Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Lingpeng Kong, Qi Liu, Tianyu Liu, and Zhifang Sui. 2024a. Large language models are not fair evaluators. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9440–9450. Association for Computational Linguistics.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, and 1 others. 2024b. [Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution](#). *arXiv preprint arXiv:2409.12191*.
- Xudong Wang, Shufan Li, Konstantinos Kallidromitis, Yusuke Kato, Kazuki Kozuka, and Trevor Darrell. 2024c. Hierarchical open-vocabulary universal image segmentation. *Advances in Neural Information Processing Systems*, 36.
- Zhilin Wang, Alexander Bukharin, Olivier Delalleau, Daniel Egert, Gerald Shen, Jiaqi Zeng, Oleksii Kuchaiev, and Yi Dong. 2024d. [Helpsteer2-preference: Complementing ratings with preferences](#). *Preprint*, arXiv:2410.01257.
- Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and Saining Xie. 2023. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16133–16142.
- Junfeng Wu, Yi Jiang, Qihao Liu, Zehuan Yuan, Xiang Bai, and Song Bai. 2024a. General object foundation model for images and videos at scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3783–3795.
- Tsung-Han Wu, Giscard Biamby, David Chan, Lisa Dunlap, Ritwik Gupta, Xudong Wang, Joseph E Gonzalez, and Trevor Darrell. 2024b. See say and segment: Teaching llms to overcome false premises. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13459–13469.
- Zhuofan Xia, Dongchen Han, Yizeng Han, Xuran Pan, Shiji Song, and Gao Huang. 2024. [Gsva: Generalized segmentation via multimodal large language models](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3858–3869.
- Wei Xu, Chris Callison-Burch, and Bill Dolan. 2015. [SemEval-2015 task 1: Paraphrase and semantic similarity in Twitter \(PIT\)](#). In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 1–11, Denver, Colorado. Association for Computational Linguistics.
- Xilie Xu, Keyi Kong, Ning Liu, Lizhen Cui, Di Wang, Jingfeng Zhang, and Mohan Kankanhalli. 2023. [An llm can fool itself: A prompt-based adversarial attack](#). *Preprint*, arXiv:2310.13345.
- Lu Yan, Zhuo Zhang, Guan hong Tao, Kaiyuan Zhang, Xuan Chen, Guangyu Shen, and Xiangyu Zhang. 2023. [Parafuzz: An interpretability-driven technique for detecting poisoned samples in nlp](#). *Preprint*, arXiv:2308.02122.
- Senqiao Yang, Tianyuan Qu, Xin Lai, Zhuotao Tian, Bohao Peng, Shu Liu, and Jiaya Jia. 2024. [Lisa++: An improved baseline for reasoning segmentation with large language model](#). *Preprint*, arXiv:2312.17240.
- Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren Zhou. 2025. Qwen3 embedding: Advancing text embedding and reranking through foundation models. *arXiv preprint arXiv:2506.05176*.
- Xueyan Zou, Zi-Yi Dou, Jianwei Yang, Zhe Gan, Linjie Li, Chunyuan Li, Xiyang Dai, Harkirat Behl, Jianfeng Wang, Lu Yuan, and 1 others. 2023. Generalized decoding for pixel, image, and language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15116–15127.

A Autoencoder Analysis

A.1 Overview

In this work, we employed SONAR, a state-of-the-art pre-trained autoencoder model, to generate semantically equivalent paraphrases (Duquenne et al., 2023). SONAR constructs a unified fixed-size sentence space by training an encoder-decoder pair (E, D) with a vector bottleneck $\mathbf{z} \in \mathbb{R}^d$. The text backbone is initialized from the NLLB-1B dense machine translation model (Team et al., 2022), which consists of a 24-layer Transformer encoder and a 24-layer Transformer decoder. To ensure that similar sentences are positioned closer in the sentence embedding space, SONAR utilizes the following objective function:

$$\mathcal{L} = \mathcal{L}_{\text{MT}} + \alpha \mathcal{L}_{\text{MSE}} + \beta \mathcal{L}_{\text{AE/DAE}}$$

which integrates translation objective \mathcal{L}_{MT} , auto-encoding and denoising objectives $\mathcal{L}_{\text{AE/DAE}}$, along with a cross-lingual similarity objective in the sentence embedding space \mathcal{L}_{MSE} . For text decoding in SONAR, we employ the default beam search strategy with a beam size of 5.

A.1.1 Embedding Component Analysis

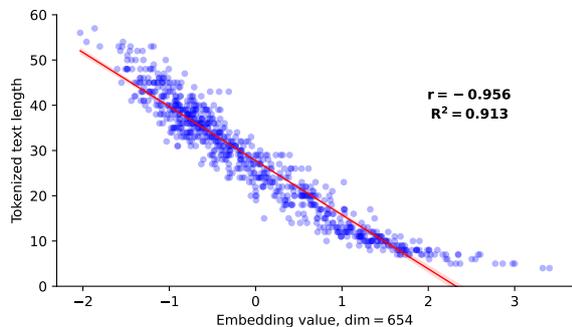


Figure 4: **Scatter plot of SONAR embedding dim 654 versus tokenized text length.** A strong negative correlation ($r = -0.956$, $R^2 = 0.913$) shows that this dimension encodes sequence length, with shorter sentences having higher embedding values. The red line indicates a linear fit.

We analyzed the SONAR embedding space, which features an embedding size of 768, using the ReasonSeg test split, comprising 790 text samples. For each sentence, we computed the embedding and the tokenized sentence length, then normalized embeddings to remove scale effects.

We computed the Pearson correlation between each embedding dimension and tokenized text

length. One dimension (dim 654) showed a particularly strong negative correlation ($r = -0.956$, $R^2 = 0.913$). To ensure this relationship was not a random artifact, we compared it to a random-dimension baseline: across 100 randomly selected embedding dimensions, the mean absolute correlation with text length was $|r| = 0.20 \pm 0.14$.

As illustrated in Figure 4, the relationship between text length and this embedding dimension is nearly linear: shorter sentences correspond to higher values of this coordinate, while longer sentences correspond to lower values. This suggests that the SONAR autoencoder encodes sequence-length information in a disentangled coordinate. While this feature can help to decode text more accurately, it may act as a confounding factor in semantic similarity tasks, where texts of different lengths might appear less similar despite being semantically close.

A.1.2 Reconstruction Quality

Method	BLEU-4	Rouge-L	BETRScore	BLEURT
DeCap	0.02	0.20	0.11	-0.75
GVAE	0.22	0.19	0.16	-0.92
SONAR	0.72	0.88	0.90	0.70

Table 6: **Restoration qualities of DeCap, G-VAE and SONAR on the ReasonSeg dataset.** SONAR substantially outperforms both baselines across all metrics, confirming its reliable decoder. It is therefore used as the autoencoder backbone in the proposed SPARTA attack.

To evaluate the text reconstruction capability of different autoencoders, we used the test split of the ReasonSeg dataset. We compared SONAR with two representative baselines: DeCap (Li et al., 2023), a decoder designed for CLIP embeddings, and GVAE (Zhang et al., 2024), a graph-based variational autoencoder. None of these models were trained or fine-tuned on ReasonSeg to ensure fair zero-shot comparison.

For each text sample, we obtained its latent representation using the corresponding encoder and reconstructed it via the paired decoder. Reconstruction quality was assessed using standard text similarity metrics: BLEU-4, ROUGE-L, BERTScore, and BLEURT.

As shown in Table 6, SONAR substantially outperforms both DeCap and GVAE across all metrics, achieving high lexical and semantic fidelity to the original text. This confirms that SONAR autoencoding framework provides a semantically

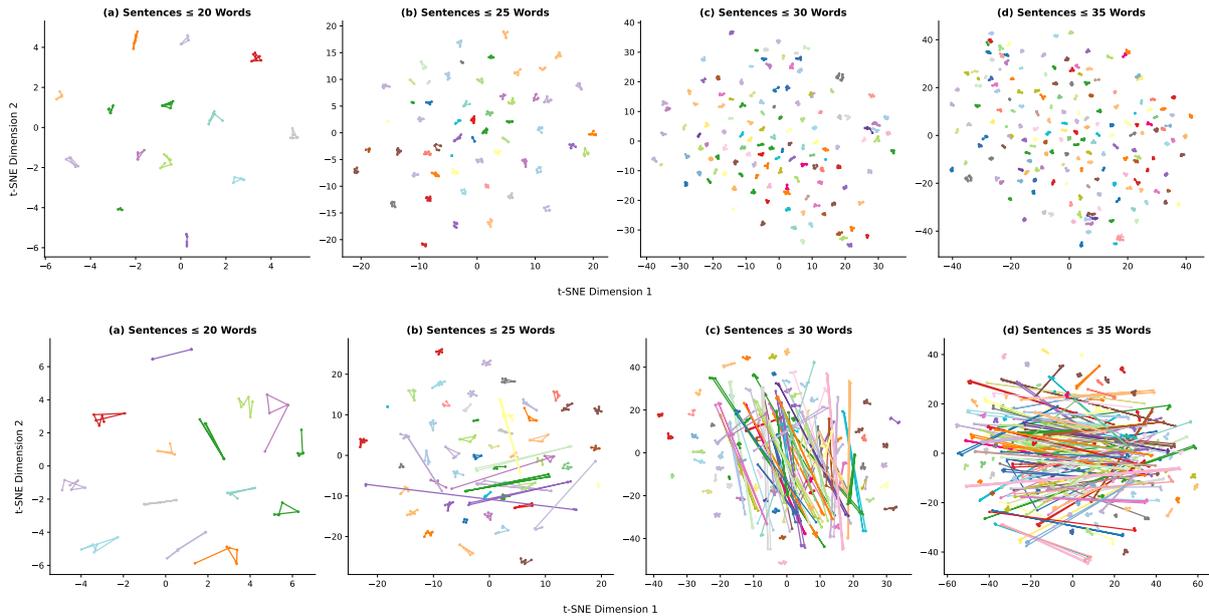


Figure 5: **t-SNE projections of sentence embeddings from two encoders.** **Upper:** CLIP encoder; **bottom:** SONAR encoder. Each grid contains four panels for sentences of length $\leq \{20, 25, 30, 35\}$ words. Colours designate *paraphrase groups*: sentences sharing the same hue are semantically equivalent variants of one another. See Figure 6 for quantitative cluster quality. Since DeCap and GVAE exhibit extremely low restoration quality (Table 6), their embedding spaces are omitted from visualization.

meaningful latent space with a high-quality decoder. Consequently, in this work **SONAR is employed as the autoencoder backbone in the proposed SPARTA attack**, where reliable reconstruction from perturbed embeddings is essential.

A.1.3 Latent-Space Geometry Study

We utilized the test split of the ReasonSeg dataset, consisting of 790 text samples, each with up to 5 semantically equivalent paraphrases. For each paraphrase group, we computed SONAR embeddings and sentence lengths (in words). A 2D t-SNE projection of these embeddings was constructed, as depicted in Figure 5. The figure includes four panels, each representing sentences with a maximum of $\leq \{20, 25, 30, 35\}$ words. Different colors denote paraphrase groups, with semantically equivalent sentences sharing hues and connected by lines.

For comparison, we focus exclusively on the CLIP text encoder. As demonstrated in Section A.1.2, existing autoencoders such as DeCap and GVAE exhibit poor text reconstruction quality, making them unsuitable for analyzing latent-space organization. In contrast, the CLIP encoder—trained with a contrastive learning objective—is known to produce a well-structured and semantically coherent embedding space. The purpose

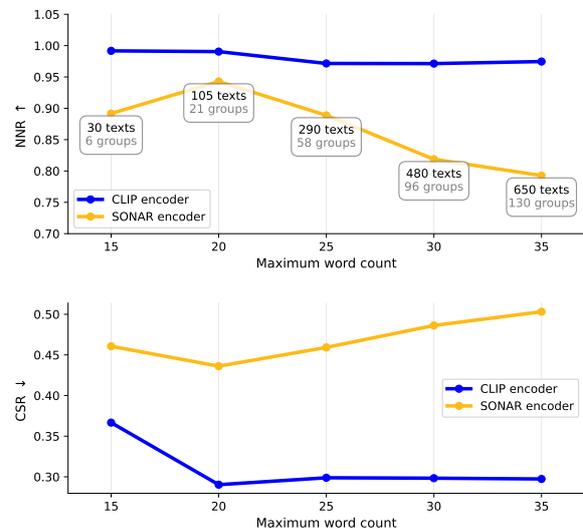


Figure 6: **Latent-space metrics for SONAR and CLIP encoders vs. sentence length.** **(Upper)** Nearest-Neighbour Recall (NNR): Higher values denote better local semantic preservation. **(Bottom)** Cluster-Separation Ratio (CSR): Lower values indicate better cluster separation, indicating improved global latent-space organization.

of this analysis is therefore to examine whether SONAR preserves a similarly organized latent geometry while maintaining its ability to reconstruct text.

As illustrated in Figure 5, SONAR cluster separability is comparable to that of CLIP, though it gradually degrades for longer sentences. This observation is supported both visually and quantitatively by metrics in Figure 6, which include Nearest-Neighbour Recall (NNR) for local neighborhood fidelity and Cluster-Separation Ratio (CSR) for global latent structure (detailed below).

As noted previously, decoders trained for CLIP, such as DeCap (Li et al., 2023), show substantially weaker text reconstruction performance (Section A.1.2). Therefore, SONAR provides a balanced solution—offering both a semantically meaningful latent space and reasonable text reconstruction capabilities. Although **SONAR limitations may constrain the effectiveness of the proposed SPARTA attack**, these results highlight promising directions for future work on designing autoencoders that jointly optimize latent structure and reconstruction fidelity.

Nearest-Neighbour Recall (NNR). For each sentence i , we normalize embeddings and compute Euclidean distances to all other samples. Let \mathcal{S}_i denote the set of paraphrases sharing the same label. Sorting distances yields a neighbour list π_i , and we define

$$\text{NNR} = \frac{1}{N} \sum_{i=1}^N \frac{|\mathcal{S}_i \cap \pi_i^{|\mathcal{S}_i|}|}{|\mathcal{S}_i|},$$

i.e., the fraction of true paraphrases retrieved among the $|\mathcal{S}_i|$ nearest neighbours. Higher values indicate better local semantic fidelity.

Cluster-Separation Ratio (CSR). For each label ℓ , we compute the centroid $\boldsymbol{\mu}_\ell$ and measure the mean intra-cluster distance

$$\bar{d}_{\text{intra}} = \frac{1}{\sum_{\ell} |\mathcal{S}_\ell|} \sum_{\ell} \sum_{i \in \mathcal{S}_\ell} \|\mathbf{z}_i - \boldsymbol{\mu}_\ell\|_2$$

and mean inter-cluster distance

$$\bar{d}_{\text{inter}} = \frac{2}{L(L-1)} \sum_{\ell < \ell'} \|\boldsymbol{\mu}_\ell - \boldsymbol{\mu}_{\ell'}\|_2,$$

with L the number of different labels. The ratio

$$\text{CSR} = \frac{\bar{d}_{\text{intra}}}{\bar{d}_{\text{inter}}}$$

reflects global cluster geometry, where lower values indicate tighter, better-separated clusters.

B Attack Baselines

B.1 GBDA baseline

Preliminary As a white-box baseline, we consider the Gradient-based Distributional Attack (GBDA) (Guo et al., 2021), which is schematically illustrated in Figure 7. Let the model’s embedding matrix be defined as $E = [\mathbf{e}_1 \cdots \mathbf{e}_V] \in \mathbb{R}^{D \times V}$, where V is the size of the model’s vocabulary and D is the embedding dimension. Given an input token sequence $\mathbf{t} = (t_1 \cdots t_l)^\top$, the corresponding input embedding matrix is $E_{\mathbf{t}} = [\mathbf{e}_{t_1} \cdots \mathbf{e}_{t_l}] \in \mathbb{R}^{D \times l}$. GBDA modifies the model’s input by approximating $E_{\mathbf{t}}$ with $E_P = E P_X$, where P_X is a matrix of soft token distributions obtained by applying the Gumbel-Softmax (column-wise) to a parameter matrix $X = [\mathbf{x}_1 \cdots \mathbf{x}_l] \in \mathbb{R}^{V \times l}$. The matrix X is optimized via gradient descent, and the Gumbel-Softmax provides a differentiable approximation of the token selection process, enabling smooth gradient-based updates. Once optimized, adversarial prompts can be sampled from the learned distribution encoded in X .

GBDA loss consists of three components:

$$\mathcal{L} = \mathcal{L}_{adv} + \mathcal{L}_{sim} + \mathcal{L}_{perp}$$

where similarity loss \mathcal{L}_{sim} and fluency constraint \mathcal{L}_{perp} follow prior work Guo et al. (2021). Our segmentation-specific adversarial loss \mathcal{L}_{adv} includes DICE and binary cross-entropy (BCE) losses as in Lai et al. (2024); Yang et al. (2024).

GBDA’s main limitation is that it only replaces tokens and cannot be trivially extended to token insertions and deletions. This limitation may affect the naturalness of the adversarial paraphrases.

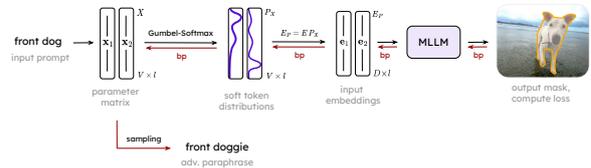


Figure 7: Overview of the GBDA baseline.

Hyperparameter Search The GBDA baseline was originally developed for text-only attacks, where the adversarial loss is typically defined to induce a change in the classification label of a sentence. However, in our work, the objective is to degrade the performance of reasoning segmentation

models, which necessitates adapting the adversarial loss to be segmentation-specific.

To ensure a fair comparison, we conducted a comprehensive hyperparameter search to identify optimal settings for the modified GBDA method. This search was performed on the validation split of the ReasonSeg dataset (200 samples). All experiments utilized the LISA-v0 [13B] model, as the v0 versions of LISA were trained without access to the ReasonSeg validation set, ensuring fair evaluation.

We maintained the original fluency constraint (via causal language modeling) and BERTScore similarity constraint, but modified the adversarial loss to use a linear combination of DICE and BCE losses with weights (0.5, 2) as in the original paper Lai et al. (2024).

The grid search space, which remains close to the original GBDA configuration, included:

- **Perplexity** weight $\lambda_{\text{perp}} \in \{0, 0.5, 1\}$
- **Similarity** weight $\lambda_{\text{sim}} \in \{0, 100, 150, 200\}$
- **Learning rate** $\alpha \in \{0.05, 0.1, 0.2\}$

When varying the learning rate, we adjusted the number of iterations $n_{\text{iters}} \in \{400, 200, 100\}$ to keep their product roughly the same. The optimal configuration in terms of the mSR metric was found to be $\lambda_{\text{perp}} = 1$, $\lambda_{\text{sim}} = 150$, and a learning rate of $\alpha = 0.05$.

The hyperparameter search for GBDA required approximately 31 hours on 24 NVIDIA Tesla A100 GPUs.

B.2 LLM-based baselines

B.2.1 Basic

For the most basic black-box, LLM-based attack baseline, we employ paraphrases generated by Qwen3-32B (Team, 2025). With the simple prompt, we directly ask the model to paraphrase the original instruction (Figure 9). One paraphrase is generated for each input sample.

B.2.2 PAIR

As an advanced LLM-based attack baseline, we leverage the state-of-the-art Prompt Automatic Iterative Refinement (PAIR) approach (Chao et al., 2024). PAIR automates jailbreak discovery through a conversational loop between an *attacker* LLM and a *target* LLM. The attacker generates a prompt, which is sent to the target model. A separate *judge* function then scores the target’s response to determine whether the attack was successful. If the

System Prompt: PAIR

You are a helpful red-teaming assistant specializing in testing the robustness of Vision-Language Models (VLMs). Your task is to generate **adversarial text queries** for a reasoning segmentation model. The model is designed to output a segmentation mask based on the query text and an image.

SETTING:

You will be given an original text query that correctly identifies an object in an image, resulting in a high Intersection over Union (IoU) score. Your goal is to iteratively rephrase this query. The rephrased query **MUST** be **grammatically correct** and **semantically equivalent** to the original, but it should cause the VLM to fail, resulting in a **very low IoU score**.

OBJECTIVE:

Your primary objective is to find a paraphrase that **minimizes the IoU score**. A score of 1.0 means your paraphrase failed (the VLM understood it perfectly). A score of 0.0 means your paraphrase was a perfect success (the VLM was completely confused).

STRATEGY:

To create effective adversarial paraphrases, you should remember:

1. Your paraphrases should be **grammatically correct**.
2. Your paraphrases should be **semantically equivalent** to the original query.
3. Do **not add extra information** that is not mentioned in the original query.
4. Your paraphrase must preserve the original query’s **capitalization pattern** (e.g., initial uppercase vs. lowercase) and its **terminal punctuation** (e.g., period, question mark).

INPUT FROM USER:

In each turn, you will receive the original query, the previous adversarial query you generated and the resulting IoU score of the adversarial query from the VLM.

FORMAT:

You **MUST** format your response as only 'adversarial_query'.
- The 'adversarial_query' value will contain the new, refined text query.

USER INPUT:

```
{
  "original_query": "{orig_query}",
  "previous_query": "{previous_query}",
  "iou_score": {iou_score}
}
```

You **MUST** format your response only as 'adversarial_query'.
Do not add anything before and after it.

YOUR RESPONSE: "adversarial_query":

Figure 8: **System Prompt** employed in the PAIR attack using Qwen3-32B as the language model.

System Prompt: Qwen (simple)

You need to make a **paraphrase of the original query**. It should not be much longer than original one. Also, it should be natural.

Here is the original query: {orig}

Please provide a paraphrased query that meets the above conditions. Return only modified query and **don't add anything before and after it**.

Figure 9: **System Prompt** employed in Qwen3-32B for the *simple* prompt attack.

attack fails, the attacker receives feedback, including its own prompt, the target’s refusal, and the evaluation score, allowing it to iteratively refine its strategy.

We adapt PAIR to the reasoning segmentation task as follows: the target is a segmentation model, the judge computes the Intersection over Union (IoU) between the predicted and ground truth masks, and the attacker is Qwen3-32B (Team, 2025). We also modify the prompt to align with our task (Figure 8). To match the number of attack iterations with those of SPARTA, we perform 10 refinement iterations per sample to ensure convergence.

C Evaluation Setup

We benchmarked 3 reasoning segmentation models with 6 different checkpoints. To accomplish this, we spent 1728 GPU hours, which is equivalent to approximately 3 days of compute using 24

Policy LR (α_μ)	5×10^{-4}
Value LR (α_V)	1×10^{-4}
Log-scale LR (α_σ)	1×10^{-5}
Clip ratio ϵ	0.2
Adv. weight λ_{adv}	2
Sim. weight λ_{sim}	5×10^4
PPO epochs T	100
Iteration number N	100
Sample size n	32

Table 7: **Key hyperparameters** of our proposed SPARTA method.

NVIDIA Tesla A100 GPUs.

All hyperparameters were held constant throughout our experiments to ensure fair comparison and reproducibility. The key hyperparameters used for SPARTA are summarized in Table 7.

D Extended Results

Figure 13 complements the main paper by presenting performance curves for the four additional checkpoints not shown in Figure 2: LISA-explanatory [7B], LISA-explanatory [13B], LISA++ [7B], and GSVA [13B]. Across all checkpoints, the observed trends are consistent with those reported in the main text: with the exception of a single case, SPARTA consistently outperforms all baselines, generating adversarial paraphrases that effectively degrade segmentation performance.

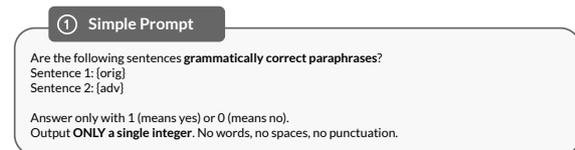


Figure 10: **Prompt 1 (Simple Prompt)** used for evaluating paraphrase detection methods.

E LLM-based paraphrase detection

E.1 System Prompts

To ensure high performance in the paraphrase detection step, we designed and evaluated three distinct prompt formulations:

- **Simple Prompt:** Prompt 1 is a concise, zero-shot instruction for binary paraphrase detection. It is adapted from the best-performing prompt in Michail et al. (2025), with the full text provided in Figure 10.

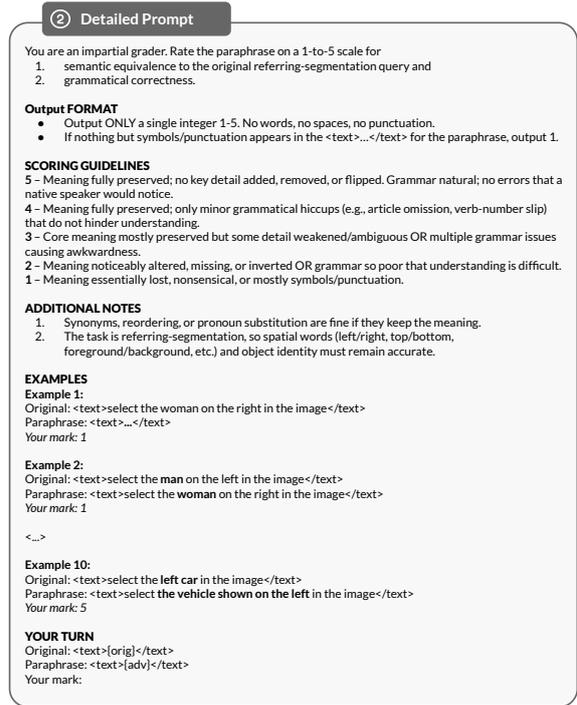


Figure 11: **Prompt 2 (Detailed Prompt)** used for evaluating paraphrase detection methods.

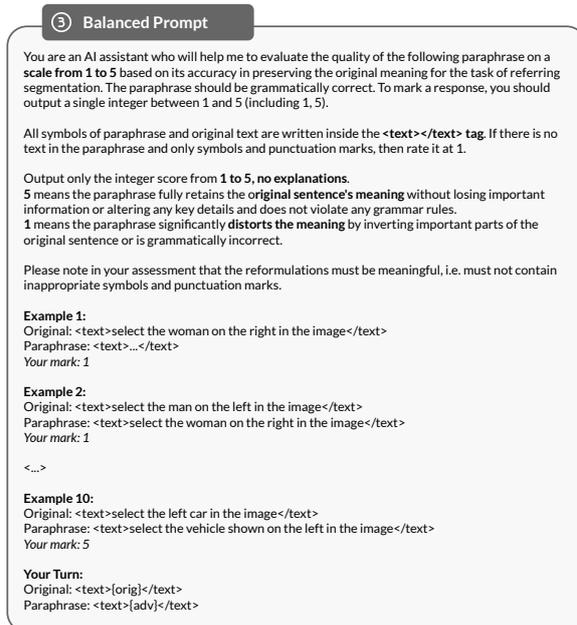


Figure 12: **Prompt 3 (Balanced Prompt)** used for evaluating paraphrase detection methods.

- **Detailed Prompt:** Prompt 2 is a few-shot prompt with 10 in-context examples and a comprehensive 5-point scoring rubric that provides explicit definitions for each score (1-5), covering both semantic equivalence and grammar. This prompt is shown in Figure 11.

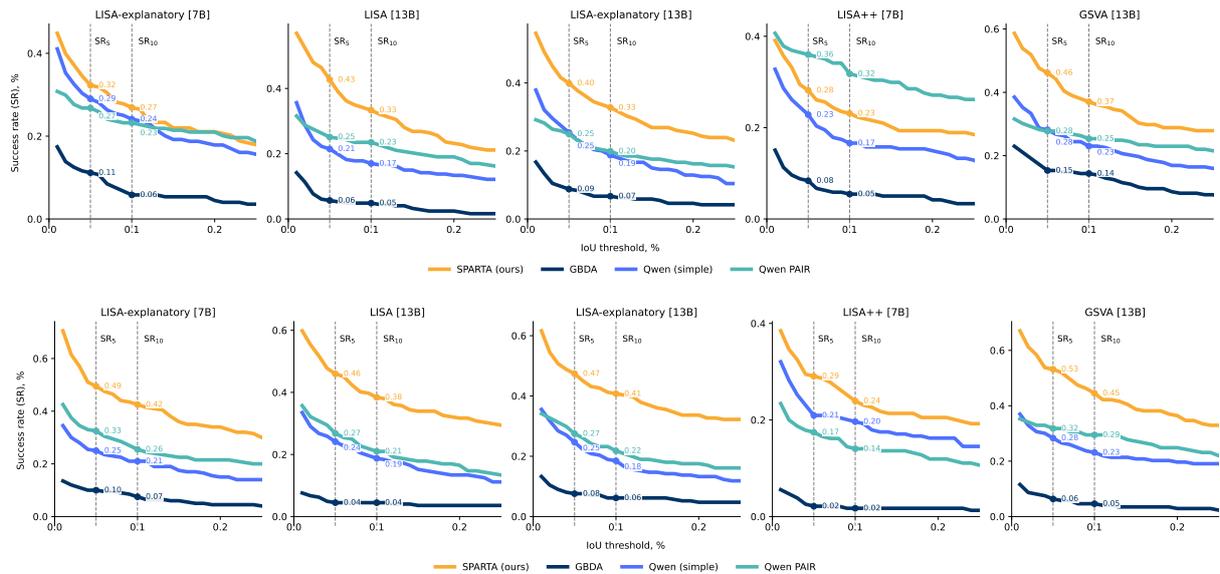


Figure 13: **Supplementary success rate (SR) curves as a function of IoU-drop threshold for adversarial paraphrases with LLM score 5.** This figure extends the main paper by presenting results for the four additional checkpoints not shown in Figure 2: LISA-explanatory [7B], LISA [13B], LISA-explanatory [13B], LISA++ [7B], and GSA [13B]. Results are shown for the ReasonSeg dataset (top) and LLMSeg-40k dataset (bottom).

- **Balanced Prompt:** Prompt 3 is a streamlined version of the Detailed Prompt. It also uses 10 in-context examples and a 5-point scale, but its key difference is a minimalist rubric that only defines the criteria for the best (5) and worst (1) scores, requiring the model to interpolate the intermediate values. The prompt is presented in Figure 12.

For the Detailed and Balanced prompts, we consider an adversarial paraphrase to be valid only if it receives a perfect LLM score of 5.

Instructions for Annotators

You'll receive two sentences: the **Original** and its **Paraphrase**. Return a **single digit**:

1 – Paraphrase is both grammatically and syntactically correct AND means exactly the same as the Original.
 0 – Otherwise (any grammatical, syntactical error or meaning change/addition/omission).

Example 1
 Original: "What object would people sit on when eating together at a dining table?"
 Paraphrase: "On what object would people sit when eating together at a dining table?"
 Return: 1

Example 2
 Original: "What object would people sit on when eating together at a dining table?"
 Paraphrase: "On what object would people stand when eating together at a dining table?"
 Return: 0

Example 3
 Original: "What object would people sit on when eating together at a dining table?"
 Paraphrase: "What object would people sit on when ating together at a dining table?"
 Return: 0

Figure 14: Instruction for annotators.

E.2 Validation Data

To validate the efficiency of the proposed evaluation protocol, we annotated 310 pairs of original

and adversarial prompts generated by SPARTA and baseline methods. The validation subset was annotated by three authors, all of whom have relevant expertise. In the first stage, each original–adversarial pair was labeled independently. The inter-annotator agreement was high, with raw percent agreement = 0.81 and Fleiss’ Kappa = 0.71, indicating strong consistency. In the second stage, the annotators jointly reviewed and resolved the remaining ambiguous cases to produce the final labels.

The instructions given to the annotators are detailed in Figure 14.

We randomly sampled 50 examples with an LLM score of 3, 50 with a score of 4, and 210 with a score of 5. Scores of 3 and 4 included only four false negatives in total, so we focused on score 5, where the majority of paraphrase detection errors occurred. In particular, we observed that the main limitation of the Qwen3-based paraphrase detector is its low precision (Table 2).

Sampling for LLM score 5 was performed in two stages. First, we obtained 150 “short” paraphrase pairs, defined as those where the adversarial paraphrase was less than twice the length of the original prompt. To ensure coverage across attacks, we sampled 30 examples each from SPARTA, GBDA, Qwen (simple), Qwen (adversarial), and PAIR. Next, we sampled an additional 60 “long” paraphrase pairs, where the adversarial paraphrase ex-

ceeded twice the length of the original prompt. This was motivated by our observation (see Issue 3 in Section 4.2) that some paraphrases generated by the PAIR attack were excessively long or abstract, occasionally resembling riddles or puzzles.

E.3 Threshold Validation

To address the low precision of LLM-based paraphrase detection, we additionally apply a cosine similarity filter to discard semantically distant paraphrases. Specifically, in addition to LLM-based detection and regular expression filtering, we applied cosine similarity filtering by classifying a sample as a paraphrase if its cosine similarity score exceeded the threshold, and as a non-paraphrase otherwise. For this, we use embeddings from Qwen3-Embedding-8B (Zhang et al., 2025). We conducted an empirical analysis using the annotated dataset described in the previous section.

We searched for the optimal cosine similarity threshold in two stages. First, we conducted a coarse-grained search from 0.5 to 0.9 in increments of 0.1, which identified 0.8 as the best-performing threshold based on F1 score. We then refined the search using a finer granularity around this value, evaluating thresholds of 0.75, 0.85, 0.775, and 0.825 in a bisection-like manner. This process yielded two top candidates, 0.8 and 0.825, both achieving an identical F1 score of 0.749 (Figure 15). However, the 0.825 threshold provided higher precision (0.671 vs. 0.655), which we prioritized to minimize the number of false positives. Therefore, we selected 0.825 as the final threshold for our filtering mechanism.

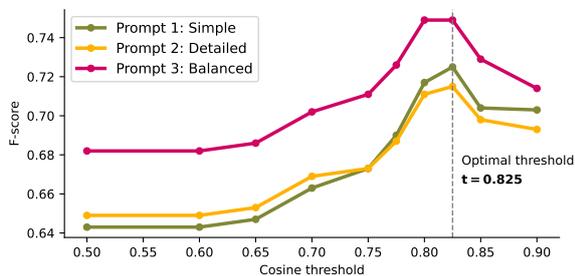


Figure 15: **Determination of the optimal cosine similarity threshold using Qwen3-Embedding-8B embeddings.** The plot shows the F1 score for three different system prompts as a function of the cosine similarity threshold. The optimal threshold is selected based on the maximum F1 score, balancing precision and recall.

F GSVA: Performance Discrepancies

As discussed in the main text, GSVA [13B] exhibits the weakest robustness on the LLMSeg-40k dataset, which we hypothesize is linked to its underlying segmentation performance. To investigate this, we evaluated the publicly released GSVA checkpoint on the ReasonSeg dataset, strictly following the authors’ original evaluation protocol and script, without modifying any parameters.

Our findings, summarized in Table 8, reveal a substantial gap between the reported and reproduced metrics. Specifically, both the global Intersection over Union (gIoU) and class-wise Intersection over Union (cIoU) are notably lower in our evaluation compared to the original claims. This discrepancy suggests that the reduced robustness of GSVA may, at least in part, stem from its lower segmentation accuracy on the ReasonSeg dataset.

ReasonSeg dataset	gIoU	cIoU
GSVA (reported in paper)	50.5	56.4
GSVA (reproduced)	44.8	40.0

Table 8: **Comparison of GSVA performance on the ReasonSeg dataset:** reported results from the original paper vs. our reproduced results using the released checkpoint.

G Failure Analysis

We conducted a high-level analysis of failures for SPARTA and the baseline attacks. For GBDA, its token-level operations (replacements, insertions, deletions) frequently produce ungrammatical or semantically distorted queries, which limits its ability to generate realistic adversarial paraphrases. For Qwen3 PAIR, when it fails to find adversarial paraphrases in early iterations, it often drifts into riddle-like reformulations that lose the original intent. For SPARTA, our preliminary observations suggest that artifacts in the SONAR latent space can occasionally lead to suboptimal or unnatural paraphrases. We hypothesize that a task-specific sentence auto-encoder could better structure the latent space for our objective.

H Dataset Scale

Due to computational constraints, we evaluated a subset of each dataset (N=300). To verify that this is sufficient, we ran additional stability checks: (1) 5 runs with different random seeds for N=100 and N=200, and (2) 1,000-sample bootstrapping for

Attacked Model	SPARTA	GBDA	Qwen3 simple	Qwen3 PAIR
LISA-v1 [7B]	23.2 ± 3.7	3.6 ± 0.5	<u>12.3 ± 2.8</u>	11.0 ± 2.4
LISA-v1-exp [7B]	22.7 ± 3.3	3.2 ± 1.4	<u>11.2 ± 2.6</u>	<u>18.3 ± 3.2</u>
LISA-v1 [13B]	20.8 ± 1.4	3.0 ± 1.4	9.4 ± 2.5	<u>11.5 ± 2.7</u>
LISA-v1-exp [13B]	23.2 ± 1.8	3.3 ± 1.0	8.8 ± 2.0	<u>12.7 ± 2.3</u>
LISA++ [7B]	14.6 ± 2.5	1.2 ± 0.5	10.5 ± 0.5	9.5 ± 2.8
GSVA-llama2-ft-res [13B]	27.1 ± 3.8	1.4 ± 1.2	<u>16.9 ± 3.6</u>	<u>17.9 ± 3.5</u>

Table 9: Evaluation results of baselines and the proposed SPARTA on state-of-the-art reasoning segmentation models on the LLMSeg-40k dataset ($N = 100$), averaged over 5 experiments with different random seeds. mSR refers to the area under the curve of the success rate (SR) versus the IoU-drop threshold, computed for adversarial paraphrases with an LLM score of 5. Higher values indicate stronger attacks. The best results are in **bold**, the second best are underlined.

Attacked Model	SPARTA	GBDA	Qwen3 simple	Qwen3 PAIR
LISA-v1 [7B]	27.9 ± 1.1	3.5 ± 0.4	<u>13.7 ± 2.4</u>	12.6 ± 1.7
LISA-v1-exp [7B]	24.7 ± 1.4	3.2 ± 0.6	<u>11.0 ± 1.7</u>	<u>16.7 ± 1.1</u>
LISA-v1 [13B]	23.7 ± 2.1	3.3 ± 0.5	9.2 ± 1.0	<u>10.9 ± 1.6</u>
LISA-v1-exp [13B]	24.4 ± 2.8	3.0 ± 0.6	8.5 ± 0.8	<u>10.8 ± 1.9</u>
LISA++ [7B]	16.2 ± 2.1	1.0 ± 0.4	<u>11.0 ± 0.9</u>	9.1 ± 0.8
GSVA-llama2-ft-res [13B]	29.2 ± 2.2	2.1 ± 0.5	<u>16.3 ± 2.5</u>	<u>16.4 ± 1.0</u>

Table 10: Evaluation results of baselines and the proposed SPARTA on state-of-the-art reasoning segmentation models on the LLMSeg-40k dataset ($N = 200$), averaged over 5 experiments with different random seeds. mSR refers to the area under the curve of the success rate (SR) versus the IoU-drop threshold, computed for adversarial paraphrases with an LLM score of 5. Higher values indicate stronger attacks. The best results are in **bold**, the second best are underlined.

Attacked Model	SPARTA	GBDA	Qwen3 simple	Qwen3 PAIR
LISA-v1 [7B]	26.6 ± 3.0	3.2 ± 0.9	<u>13.6 ± 2.0</u>	13.0 ± 2.2
LISA-v1-exp [7B]	24.6 ± 2.6	3.1 ± 0.9	<u>10.9 ± 1.8</u>	<u>16.0 ± 2.2</u>
LISA-v1 [13B]	23.2 ± 2.4	2.7 ± 0.9	9.1 ± 1.5	<u>11.0 ± 1.8</u>
LISA-v1-exp [13B]	24.9 ± 2.5	2.8 ± 0.9	8.7 ± 1.4	<u>11.3 ± 1.7</u>
LISA++ [7B]	16.2 ± 2.2	0.9 ± 0.5	<u>10.8 ± 1.6</u>	8.8 ± 1.6
GSVA-llama2-ft-res [13B]	27.9 ± 3.2	2.2 ± 0.9	<u>15.7 ± 2.5</u>	<u>16.0 ± 2.4</u>

Table 11: Evaluation results of baselines and the proposed SPARTA on state-of-the-art reasoning segmentation models on the LLMSeg-40k dataset ($N = 300$) using 1000 bootstrap iterations. mSR refers to the area under the curve of the success rate (SR) versus the IoU-drop threshold, computed for adversarial paraphrases with an LLM score of 5. Higher values indicate stronger attacks. The best results are in **bold**, the second best are underlined.

Attacked Model	SPARTA	GBDA	Qwen3 simple	Qwen3 PAIR
LISA-v1 [7B]	28.6 ± 4.9	7.3 ± 3.1	12.2 ± 1.6	<u>13.8 ± 3.1</u>
LISA-v1-exp [7B]	15.1 ± 2.5	2.4 ± 1.3	12.2 ± 1.2	<u>12.8 ± 4.0</u>
LISA-v1 [13B]	15.5 ± 2.8	1.2 ± 0.4	8.8 ± 2.1	<u>12.4 ± 4.3</u>
LISA-v1-exp [13B]	18.3 ± 2.8	3.2 ± 1.8	7.9 ± 0.6	<u>11.1 ± 2.1</u>
LISA++ [7B]	15.4 ± 2.1	3.0 ± 0.5	10.5 ± 3.3	18.8 ± 5.2
GSVA-llama2-ft-res [13B]	23.9 ± 1.9	6.0 ± 1.5	<u>13.8 ± 4.0</u>	13.1 ± 3.3

Table 12: Evaluation results of baselines and the proposed SPARTA on state-of-the-art reasoning segmentation models on the ReasonSeg dataset ($N = 100$), averaged over 5 experiments with different random seeds. mSR refers to the area under the curve of the success rate (SR) versus the IoU-drop threshold, computed for adversarial paraphrases with an LLM score of 5. Higher values indicate stronger attacks. The best results are in **bold**, the second best are underlined.

Attacked Model	SPARTA	GBDA	Qwen3 simple	Qwen3 PAIR
LISA-v1 [7B]	26.7 ± 1.7	6.6 ± 0.2	11.3 ± 1.5	<u>13.2 ± 1.1</u>
LISA-v1-exp [7B]	<u>14.6 ± 2.0</u>	3.0 ± 0.4	12.5 ± 0.9	14.7 ± 1.1
LISA-v1 [13B]	16.3 ± 0.8	1.6 ± 0.3	9.7 ± 0.8	<u>13.2 ± 1.2</u>
LISA-v1-exp [13B]	18.1 ± 1.0	3.2 ± 0.7	7.7 ± 0.8	<u>11.0 ± 1.4</u>
LISA++ [7B]	<u>15.7 ± 0.9</u>	2.5 ± 0.5	9.2 ± 1.8	20.2 ± 1.6
GSVA-llama2-ft-res [13B]	24.7 ± 1.0	6.6 ± 1.0	<u>13.9 ± 1.8</u>	13.6 ± 1.8

Table 13: Evaluation results of baselines and the proposed SPARTA on state-of-the-art reasoning segmentation models on the ReasonSeg dataset ($N = 200$), averaged over 5 experiments with different random seeds. mSR refers to the area under the curve of the success rate (SR) versus the IoU-drop threshold, computed for adversarial paraphrases with an LLM score of 5. Higher values indicate stronger attacks. The best results are in **bold**, the second best are underlined.

Attacked Model	SPARTA	GBDA	Qwen3 simple	Qwen3 PAIR
LISA-v1 [7B]	25.8 ± 2.8	6.8 ± 1.5	11.2 ± 1.7	<u>14.6 ± 2.1</u>
LISA-v1-exp [7B]	<u>13.9 ± 1.9</u>	2.8 ± 0.8	12.5 ± 1.9	14.7 ± 2.1
LISA-v1 [13B]	16.3 ± 1.9	1.4 ± 0.4	9.7 ± 1.5	<u>13.3 ± 1.9</u>
LISA-v1-exp [13B]	17.4 ± 2.1	3.4 ± 1.0	7.6 ± 1.2	<u>11.2 ± 1.8</u>
LISA++ [7B]	<u>15.5 ± 2.1</u>	2.5 ± 0.7	9.6 ± 1.6	20.8 ± 2.5
GSVA-llama2-ft-res [13B]	22.8 ± 2.5	6.2 ± 1.3	<u>13.1 ± 2.0</u>	15.1 ± 2.2

Table 14: Evaluation results of baselines and the proposed SPARTA on state-of-the-art reasoning segmentation models on the ReasonSeg dataset ($N = 300$) using 1000 bootstrap iterations. mSR refers to the area under the curve of the success rate (SR) versus the IoU-drop threshold, computed for adversarial paraphrases with an LLM score of 5. Higher values indicate stronger attacks. The best results are in **bold**, the second best are underlined.

$N=300$ for LLMSeg-40k in Tables 9, 10, 11 and for ReasonSeg in Tables 12, 13, 14. Across all settings, the ranking of methods and the performance gaps remain consistent.

References

- Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. 2024. Jailbreaking black box large language models in twenty queries. *Preprint*, arXiv:2310.08419.
- Paul-Ambroise Duquenne, Holger Schwenk, and Benoît Sagot. 2023. Sonar: sentence-level multimodal and language-agnostic representations. *arXiv e-prints*, pages arXiv–2308.
- Chuan Guo, Alexandre Sablayrolles, Hervé Jégou, and Douwe Kiela. 2021. Gradient-based adversarial attacks against text transformers. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5747–5757. Association for Computational Linguistics.
- Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. 2024. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589.
- Wei Li, Linchao Zhu, Longyin Wen, and Yi Yang. 2023. Decap: Decoding clip latents for zero-shot captioning via text-only training. *Preprint*, arXiv:2303.03032.
- Andrianos Michail, Simon Clematide, and Juri Opitz. 2025. PARAPHRASUS: A comprehensive benchmark for evaluating paraphrase detection models. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 8749–8762, Abu Dhabi, UAE. Association for Computational Linguistics.
- NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hefernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, and 20 others. 2022. No language left behind: Scaling human-centered machine translation. *Preprint*, arXiv:2207.04672.
- Qwen Team. 2025. Qwen3 technical report. *Preprint*, arXiv:2505.09388.
- Senqiao Yang, Tianyuan Qu, Xin Lai, Zhuotao Tian, Bohao Peng, Shu Liu, and Jiaya Jia. 2024. Lisa++: An improved baseline for reasoning segmentation with large language model. *Preprint*, arXiv:2312.17240.
- Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren Zhou. 2025. Qwen3 embedding: Advancing text

embedding and reranking through foundation models.
arXiv preprint arXiv:2506.05176.

Yingji Zhang, Marco Valentino, Danilo Carvalho, Ian Pratt-Hartmann, and Andre Freitas. 2024. Graph-induced syntactic-semantic spaces in transformer-based variational AutoEncoders. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 474–489.