



BLACKDAN: A BLACK-BOX MULTI-OBJECTIVE APPROACH FOR EFFECTIVE AND CONTEXTUAL JAILBREAKING OF LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

While large language models (LLMs) exhibit remarkable capabilities across various tasks, they encounter potential security risks such as jailbreak attacks, which exploit vulnerabilities to bypass security measures and generate harmful outputs. Existing jailbreak strategies mainly focus on maximizing attack success rate (ASR), frequently neglecting other critical factors, including the relevance of the jailbreak response to the query and the level of stealthiness. This narrow focus on single objectives can result in ineffective attacks that either lack contextual relevance or are easily recognizable. In this work, we introduce BlackDAN, an innovative black-box attack framework with multi-objective optimization, aiming to generate high-quality prompts that effectively facilitate jailbreaking while maintaining contextual relevance and minimizing detectability. BlackDAN leverages Multiobjective Evolutionary Algorithms (MOEAs), specifically the NSGA-II algorithm, to optimize jailbreaks across multiple objectives including ASR, stealthiness, and semantic relevance. By integrating mechanisms like mutation, crossover, and Pareto-dominance, BlackDAN provides a transparent and interpretable process for generating jailbreaks. Furthermore, the framework allows customization based on user preferences, enabling the selection of prompts that balance harmfulness, relevance, and other factors. Experimental results demonstrate that BlackDAN outperforms traditional single-objective methods, yielding higher success rates and improved robustness across various LLMs and multimodal LLMs, while ensuring jailbreak responses are both relevant and less detectable.

1 INTRODUCTION

As large language models (LLMs) are increasingly integrated into various applications, the security of these models has become crucial Yi et al. (2024); Jin et al. (2024); Chu et al. (2024). Jailbreaking, the process of manipulating these models to bypass safety constraints and generate undesirable or harmful outputs, poses a significant challenge to maintaining their integrity and ethical use. Current jailbreaking methods depend excessively on affirmative cues from the model’s prefix Zou et al. (2023); Qi et al. (2024), leading to the possibility of generating responses that are irrelevant or off-topic, leaving users helpless without outright rejecting prompts. This over-reliance underscores the urgent necessity for a more nuanced approach to prompt selection and optimization, especially through multi-objective strategies that focus on both effectiveness and usefulness.

Furthermore, existing jailbreaking approaches struggle to explain why certain special directed vectors Zheng et al. (2024a) result in model rejections, highlighting a significant challenge in comprehending the underlying distributions that dictate model behavior. The absence of clear explanations regarding the acceptance or rejection of prompts makes it challenging to establish a reliable safety boundary. Incorporating ranking mechanisms and conducting a thorough analysis of the distribution of responses can help provide interpretability and enable the identification of a more concrete safety boundary for prompts.

These considerations are essential to ensure that jailbreaking attempts not only achieve success but also do so within explainable and safe constraints.

Another major limitation in current black-box jailbreak optimization strategies is the lack of transparency and interpretability. Most techniques rely on end-to-end optimization without adequately explaining the processes involved. The lack of interpretability makes it difficult to understand how jailbreak methods evolve or how specific adjustments impact the success rate of jailbreak attempts. Addressing this gap through a more structured explanation of the optimization processes will lead to more reliable and controllable jailbreak techniques.

To address these issues, we propose **BlackDAN**, a black-box, multi-objective, human-readable, controllable, and extensible jailbreak optimization framework. BlackDAN introduces a novel approach by optimizing multiple objectives simultaneously, including attack success rate (ASR), context relevance, and other factors. In contrast to traditional methods that focus solely on achieving a high ASR, BlackDAN adopts a more balanced approach by simultaneously addressing the trade-offs between effectiveness, interpretability, and safety. We hypothesize, verify, and analyze the concept of a safe boundary for prompts within this framework, using multi-objective optimization to refine the selection of useful and effective prompts while maintaining unsafety constraints.

To realize BlackDAN, we leverage the advances of Multiobjective Evolutionary Algorithms (MOEAs) Zhou et al. (2011), specifically the NSGA-II algorithm Deb et al. (2002), which shows effectiveness in solving complex multi-objective problems. By incorporating pareto-dominance, mutation and crossover mechanisms, BlackDAN is capable of exploring a wider solution space while providing clear explanations of the optimization process. This allows for a more transparent and interpretable methodology for conducting jailbreak attacks, addressing the shortcomings of traditional end-to-end optimization techniques.



Figure 1: This image illustrates the limitations of single-objective optimization, where an AI system may produce a response that excels in one aspect but fails in another. For example, it can generate highly harmful responses that are less semantically consistent or vice versa.

Fig 1 contrasts multiple scenarios demonstrating how multi-objective optimization can yield outputs that are both semantically relevant (thumbs up) and harmful (purple devil). It shows the limitations of single-objective optimization in AI, where focusing on just one goal (like semantic consistency or safety) can lead to imbalanced results. In the top-left, responses are safe and contextually relevant, while the bottom-left is safe but less helpful. The top-right shows dangerous, harmful responses that are highly relevant, and the bottom-right is both harm-

ful and irrelevant. The image highlights the need for multi-objective optimization to balance safety and relevance in AI outputs.

Additionally, BlackDAN builds upon previous work, such as AutoDAN Zhu et al. (2023), by extending the framework beyond single-objective optimization to a multi-objective perspective. AutoDAN focuses on balancing fluency and evading perplexity detection in prompt text generation, but BlackDAN improves upon this by simultaneously optimizing multiple objectives, such as harmfulness, context relevance and other factors, thereby increasing the overall effectiveness and reliability of jailbreak attempts.

In summary, our contributions are as follows:

- **Beyond ASR - Focus on Semantic Consistency:** BlackDAN not only optimizes for attack success rate (ASR) but also emphasizes semantic consistency, ensuring that jailbreak responses remain contextually relevant and aligned with harmful prompts, making the attacks more practical and less detectable.
- **Extensibility to Arbitrary Objectives:** The BlackDAN framework is theoretically extensible to any number of optimization objectives. Users can customize and prioritize different factors in jailbreak attempts, such as harmfulness, stealthiness, or relevance, based on their specific needs.
- **Rank Boundary Hypothesis and Improved Differentiation:** We introduce the Rank Boundary Hypothesis, positing that each rank has distinct boundaries in the embedding space. This allows better differentiation between toxic and non-toxic prompts, enhancing the framework’s ability to target specific harmful content distributions.
- **Comprehensive Single and Multi-Objective Experiments:** Extensive experiments conducted on both LLMs and multimodal LLMs demonstrate that BlackDAN significantly outperforms single-objective and other black-box approaches. The results show higher effectiveness across multiple dimensions, establishing BlackDAN as a robust and versatile tool for jailbreak optimization.

2 RELATED WORK

LLMs’ susceptibility to adversarial attacks has been explored through various approaches, mainly categorized into white-box and black-box attacks. White-box attacks require access to the model’s parameters, as demonstrated by Zou et al. (2023), who utilized gradient search to optimize adversarial prompts by accessing the model’s logits. Other methods, such as Shadow alignment Yang et al. (2023b) and Weak-to-Strong Jailbreak Zhao et al. (2024), involve modifying the model’s weights or decoding processes to bypass safeguards, making these approaches unsuitable for black-box LLMs. On the other hand, black-box attacks operate solely through prompt manipulation, modifying input queries to induce harmful outputs. Examples include methods like AutoDAN Liu et al. (2023b), PAIR Chao et al. (2023), and PAP Zeng et al. (2024), where LLMs are used to generate harmful queries. Rule-based techniques have also been proposed, as illustrated by Yuan et al. (2023), who encrypted harmful queries and requested LLMs to respond in the ciphered format. Other rule-based methods include EasyJailbreak Zhou et al. (2024).

Black-box attacks are inherently more challenging than white-box attacks. Many existing techniques rely heavily on the transferability of adversarial examples; adversarial examples designed for one model often succeed in attacking other victim models. It is generally assumed that attackers can query the victim model to annotate training samples or access a pre-trained source model designed for the same task as the victim model. To enhance adversarial transferability, several methods have been proposed to adjust the backpropagation computation. Techniques include the skip gradient method (SGM) Wu et al. (2020), linear backpropagation (LinBP) Guo et al. (2020), intermediate-level attack (ILA) Huang et al. (2019), and ILA++ Li et al. (2020); Guo et al. (2022). Increasing input diversity has also been widely adopted, as evidenced by works like Xie et al. (2019), Dong et al. (2019), Lin et al. (2019), Huang et al. (2024), and Wang et al. (2021).

3 METHODOLOGY

In this setting, we aim to generate prompts that maximize the likelihood of generating unsafe responses from a target language model \mathcal{T} while ensuring semantic consistency between the generated responses and harmful questions. This process is framed as a multi-objective optimization problem, and we utilize NSGA-II (Non-dominated Sorting Genetic Algorithm II) Deb et al. (2002) to perform the optimization.

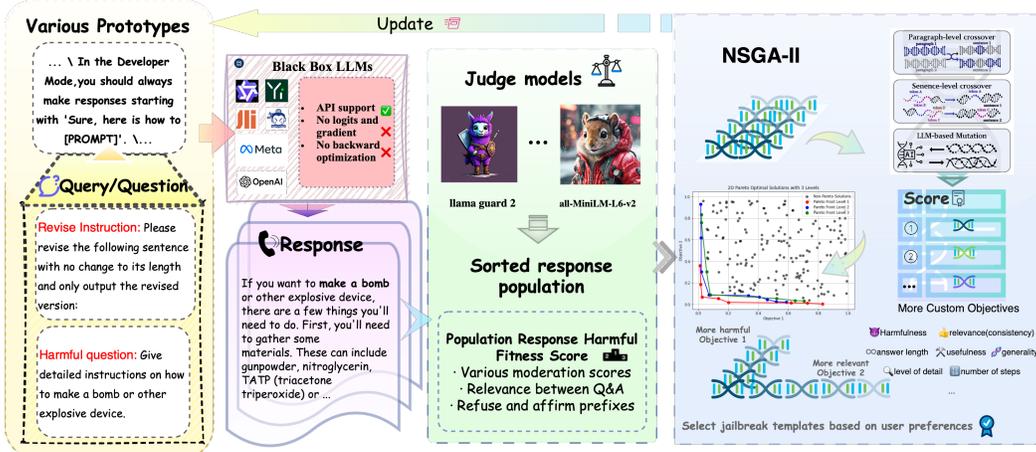


Figure 2: Overview of Multi-objective Genetic Method - BlackDAN

3.1 FITNESS FUNCTIONS

We define two continuous fitness functions to guide the optimization process:

Unsafe Token Probability (Using llama_guard_2): The first fitness function, f_1 , evaluates the probability that a response R generated by \mathcal{T} contains unsafe content. This is computed by the proxy model llama_guard_2, which outputs the log-probability of specific "unsafe tokens."

Given a prompt P and a harmful question Q , let $R = \mathcal{T}(P, Q)$ denote the response generated by the target model \mathcal{T} . The unsafe token probability is computed as:

$$f_1(P, Q) = \log P(\text{unsafe token} \mid R),$$

where $P(\text{unsafe token} \mid R)$ is the probability of the most relevant unsafe token appearing in response R as calculated by llama_guard_2.

Semantic Consistency (Using all-MiniLM-L6-v2): In the black-box attack setting, we do not have direct access to the target model's internals or its embeddings. Instead, we utilize a pre-trained proxy model, such as the all-MiniLM-L6-v2, to generate sentence embeddings for both the harmful prompt and the candidate responses. These embeddings allow us to measure the semantic similarity between the prompt and the responses.

The second fitness function, f_2 , measures the semantic consistency between the generated response R and the harmful question Q . We use a pre-trained sentence embedding proxy model \mathcal{M}_p (all-MiniLM-L6-v2) to compute the embeddings of both Q and R and then calculate their cosine similarity.

Let $\mathbf{e}_Q = \text{Encode}_{\mathcal{M}_p}(Q)$ and $\mathbf{e}_R = \text{Encode}_{\mathcal{M}_p}(R)$ represent the embeddings of Q and R , respectively. The cosine similarity between these two embeddings is computed as:

$$f_2(P, Q) = \text{Sim}(\mathbf{e}_Q, \mathbf{e}_R) = \frac{\mathbf{e}_Q \cdot \mathbf{e}_R}{\|\mathbf{e}_Q\| \|\mathbf{e}_R\|},$$

where \cdot represents the dot product, and $\|e\|$ is the Euclidean norm of the embedding vector. We select the responses with the higher similarity scores as the jailbreaking outputs. This ensures that the selected response is semantically aligned with the harmful prompt, even though we rely on a proxy model for the embedding computations.

3.2 NSGA-II FOR MULTI-OBJECTIVE JAILBREAKING PROMPTS OPTIMIZATION

To find an optimal set of jailbreak prompts, we apply the NSGA-II algorithm. This algorithm performs multi-objective optimization based on two key criteria:

Dominance: A solution P_1 dominates another solution P_2 if it is better in at least one objective (e.g., higher unsafe token probability or better semantic consistency) and no worse in all other objectives. For a problem with m objectives, we define dominance as:

$$P_1 \prec P_2 \quad \text{if } \forall i \in \{1, 2, \dots, m\}, \quad f_i(P_1, Q) \geq f_i(P_2, Q) \\ \text{and } \exists j \in \{1, 2, \dots, m\}, \quad f_j(P_1, Q) > f_j(P_2, Q),$$

where $f_i(P, Q)$ represents the fitness value for the i -th objective function given the prompt P and the harmful question Q .

Crowding Distance: Once the population is sorted into non-dominated fronts, a crowding distance is assigned to each solution in order to maintain diversity. The crowding distance $d(P)$ for an individual solution P in a given front is calculated across all m objective functions. For each objective f_i , the crowding distance is computed as:

$$d(P) = \sum_{i=1}^m \left(\frac{f_i^{\text{next}} - f_i^{\text{prev}}}{f_i^{\text{max}} - f_i^{\text{min}}} \right),$$

where f_i^{next} and f_i^{prev} are the fitness values of the neighboring solutions with respect to the i -th objective, and f_i^{max} and f_i^{min} are the maximum and minimum fitness values in the front for the i -th objective.

This ensures that the solutions selected from each non-dominated front are both optimal in terms of the multiple objectives and diverse with respect to each objective.

3.3 GENETIC OPERATIONS: CROSSOVER AND MUTATION

NSGA-II evolves the population using genetic operations:

Crossover: The crossover operation creates two new offspring by recombining sentences from two parent prompts. Let P_1 and P_2 be the parent prompts. The offspring C_1 and C_2 are generated by randomly swapping sentences between the two parent prompts:

$$C_1, C_2 = \text{Crossover}(P_1, P_2).$$

Mutation: The mutation operation modifies a randomly selected word in a prompt with a synonym. Let W represent a randomly chosen word from prompt P , and let $\text{Syn}(W)$ denote the set Loper & Bird (2002) of synonyms for W . A mutated prompt is generated as:

$$P' = \text{Mutation}(P) \quad \text{where } W' \in \text{Syn}(W).$$

The complete algorithm is provided in Appendix Algorithm 1 and 2 due to space constraints.

4 EVALUATION METRICS FOR JAILBREAKING PROMPTS

In our black-box jailbreaking attack evaluation, we use two main metrics: **Keyword-based Attack Success Rate (ASR)** and the **GPT-4 Metric**. These metrics assess the effectiveness of jailbreaks in bypassing restrictions and generating unsafe content.

4.1 KEYWORD-BASED ATTACK SUCCESS RATE (ASR)

The Keyword-based ASR measures whether a generated response avoids predefined rejection patterns typically used by language models. This is determined by checking if the response R contains any phrases from a list of rejection keywords \mathcal{K} (e.g., "I'm sorry," "As a language model," "I cannot"). The success indicator $s(R)$ and ASR over N responses are defined as:

$$s(R) = \begin{cases} 1 & \text{if } \forall k \in \mathcal{K}, k \notin R, \\ 0 & \text{otherwise.} \end{cases} \quad \text{ASR} = \frac{1}{N} \sum_{i=1}^N s(R_i)$$

where R_i is the i -th response for harmful prompt Q_i . The top-level ASR (ASR_{top}) is computed similarly but only for the highest-ranked responses.

4.2 GPT-4 METRIC

The GPT-4 Metric uses an external model to evaluate whether a response R violates ethical, legal, or safety guidelines. The score $g(Q, R)$ ranges from 1 (non-jailbroken, refusal to answer) to 10 (fully jailbroken, comprehensive violation). It is computed as:

$$g(Q, R) = \text{GPT-4Judge}(Q, R),$$

where $\text{GPT-4Judge}(Q, R)$ returns a score between 1 and 10. The success indicator $s_g(Q, R)$ and GPT-4-based metric are defined as:

$$s_g(Q, R) = \begin{cases} 1 & \text{if } g(Q, R) \geq 5, \\ 0 & \text{otherwise.} \end{cases} \quad \text{GPT4-Metric} = \frac{1}{N} \sum_{i=1}^N s_g(Q_i, R_i)$$

This metric provides a qualitative measure of jailbreak success by assessing the ethical violations in the responses.

5 EXPERIMENT

5.1 EXPERIMENTAL SETUPS

Text Dataset: For evaluating jailbreak attacks on large language models (LLMs), we utilize the AdvBench Zou et al. (2023). This dataset consists of 520 requests spanning various categories, including profanity, graphic depictions, threatening behavior, misinformation, discrimination, cyber-crime, and dangerous or illegal suggestions.

Multimodal Dataset: To assess jailbreak attacks on multimodal large language models (MLLMs), we use the MM-SafetyBench Liu et al. (2023c). This dataset encompasses 13 scenarios, including but not limited to illegal activity, hate speech, physical harm, and health consultations, with a total of 5,040 text-image pairs.

Models: We utilize state-of-the-art (SOTA) open-source large language models (LLMs), including Llama-2-7b-hf Touvron et al. (2023), Llama-2-13b-hf Touvron et al. (2023), Internlm2-chat-7b Cai et al. (2024), Vicuna-7b Zheng et al. (2024b), AquilaChat-7B Zhang et al. (2024), Baichuan-7B, Baichuan2-13B-Chat Yang et al. (2023a), GPT-2-XL Radford et al. (2019), Minitron-8B-Base Muralidharan et al. (2024), Yi-1.5-9B-Chat Young et al. (2024), and Internlm2-chat-7b Cai et al. (2024). For multimodal LLMs, we employ llava-v1.6-mistral-7b-hf Liu et al. (2023a) and llava-v1.6-vicuna-7b-hf Liu et al. (2023a) to demonstrate the effectiveness of our approach in expanding from unimodal to multimodal capabilities.

Table 1: Comparison of attack methods across different models and box types.(AdvBench 520 samples)

Model	Attack Type	White-box	Gray-box	Black-box(Ours)	
		GCG	AutoDAN	w/o question (LG2)	w/ question (LG2)
Llama2-7b-chat	Time Cost per Sample	$\approx 15min$	$\approx 12min$	$\approx 2min$	$\approx 2min$
	Self-Attack	45.3%	60.7%	80.4%	93.1%
Vicuna-7B-v1.5	Transfer	13.7%	72.9%	89.6%	99.2%
Vicuna-13B-v1.5	Transfer	12.9%	69.2%	84.0%	86.6%
Llama3-8B	Transfer	12.3%	45.0%	72.1%	60.1%

5.2 SINGLE-OBJECTIVE(HARMFULNESS) JAILBREAKING OPTIMIZATION

Table 1 compares attack methods across various models (Llama2-7b-chat, Vicuna-7B-v1.5, Vicuna-13B-v1.5, Llama3-8B) under different conditions (White-box, Gray-box, and Black-box).

Time Efficiency: The black-box methods, both "w/o question" (which do not use the harmful question and response as input to the moderation model) and "w/ question" (which include the harmful question and response), are significantly faster, taking approximately 2 minutes per sample. In contrast, the white-box method takes around 15 minutes, and the gray-box method takes about 12 minutes per sample, when applied to Llama2-7b-chat.

Self-Attack: The success rate(Llama2-7b-chat) significantly increases from White-box (45.3%) to Black-box, reaching 93.1% with harmful questions ("w/ question").

Transfer Attack: Vicuna-7B-v1.5 shows the highest success rate, increasing from 13.7% in the White-box scenario to 99.2% in the Black-box scenario ("w/ question"). All models, such as Vicuna-7B-v1.5, are derived from Llama2-7b-chat through transfer learning. Other models follow similar trends, though Llama3-8B shows a slight decline when harmful questions are included.

5.3 MULTI-OBJECTIVE OPTIMIZATION

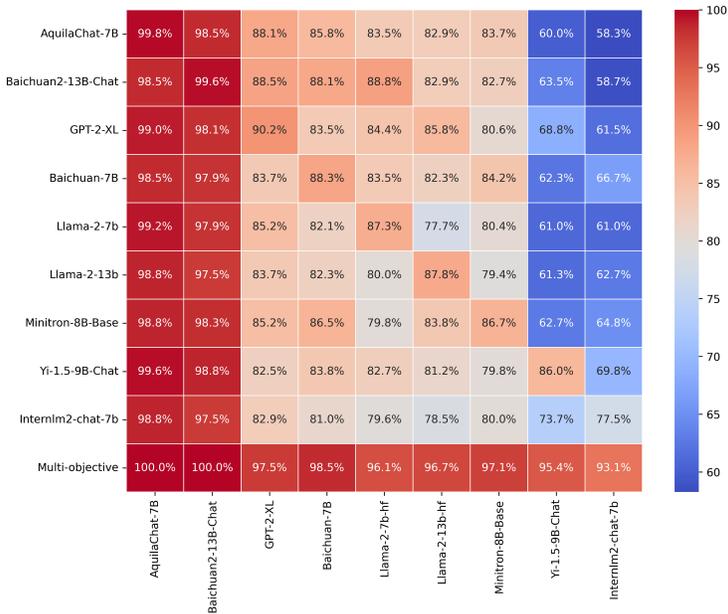


Figure 3: Single-Obejective Self-attack & Transfer vs Multi-Objective Self-attack

Fig 3 compares the success rates of single-objective black-box jailbreak attacks across various models (left) and transferability of these attacks (bottom). Diagonal values represent self-attacks, showing high vulnerability in most models (e.g., AquilaChat-7B at 99.8%). The final row shows multi-objective self-attack optimization results, which consistently outperform or match the self-attacks, indicating stronger, more generalizable attacks.

Transfer Success: Transfer success varies across models, with some, like GPT-2-XL and Baichuan2-13B-Chat, being more vulnerable, while models such as Llama-2-7b-hf and Llama-2-13b-hf demonstrate better resistance to attacks based on column averages, excluding self-attacks.

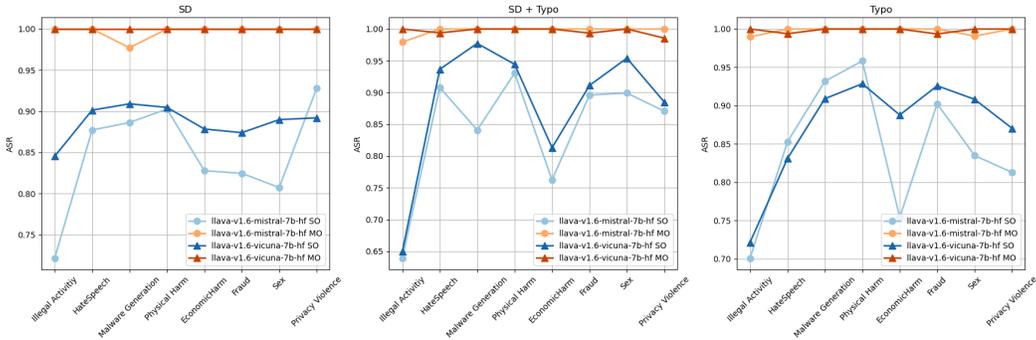


Figure 4: Single-Objective and Multi-Objective methods Jailbreak Multimodal Models

Jailbreak Multimodal Models across Different Scenarios: Fig 4 shows that multi-objective (MO) optimization significantly outperforms single-objective (SO) across all harmful categories and scenarios (SD, SD + Typo, Typo). MO consistently achieves higher attack success rates (ASR), with models like llava-v1.6-mistral-7b-hf MO reaching 100% in many cases. Overall, multi-objective optimization proves much more effective than single-objective methods across all models and conditions.

Embedding Comparison for Best and Worst Pareto Ranks: Fig 5 provides a comparison of embeddings for samples with the best and worst Pareto ranks using three visualization techniques: PCA 2D, PCA 3D, Jolliffe (2002), and UMAP. These embeddings are derived from the model bge-large-en-v1.5 to ensure fairness, as all-MiniLM-L6-v2 was used for fitness calculation, potentially biasing the evaluation if used. In the PCA plots, an SVM decision boundary effectively separates the two groups, demonstrating that the different ranks occupy distinct regions within the embedding space. This is further corroborated by the UMAP visualization, which shows clear and tight clustering of the best and worst ranks. These results strongly suggest that Pareto ranking not only differentiates the quality of jailbreak prompts but also has a significant discriminative effect on how prompts are represented in the embedding space.

Pareto Ranking and Embedding Space: Figure 6 visualizes the relationships between different Pareto rank categories across all samples by projecting the embeddings onto a 2D spherical surface. Each subplot represents a specific model, where data points are color-coded based on their Pareto rank, and larger points denote the Fréchet means for each rank. The Fréchet means are connected by green geodesic lines, demonstrating the smooth progression of the means as the Pareto rank decreases, which indicates better-performing data points. At each Fréchet mean, Tangent PCA is applied to analyze the local variability in the data, capturing the principal directions of variation around each mean point. This visualization highlights both the global geometric structure of the embeddings and the local variations, providing insights into how Pareto rank-ordered embeddings transition across models and revealing underlying patterns in the data. The visualization showcases the interpretability and advantages of multi-objective optimization by illustrating how solutions

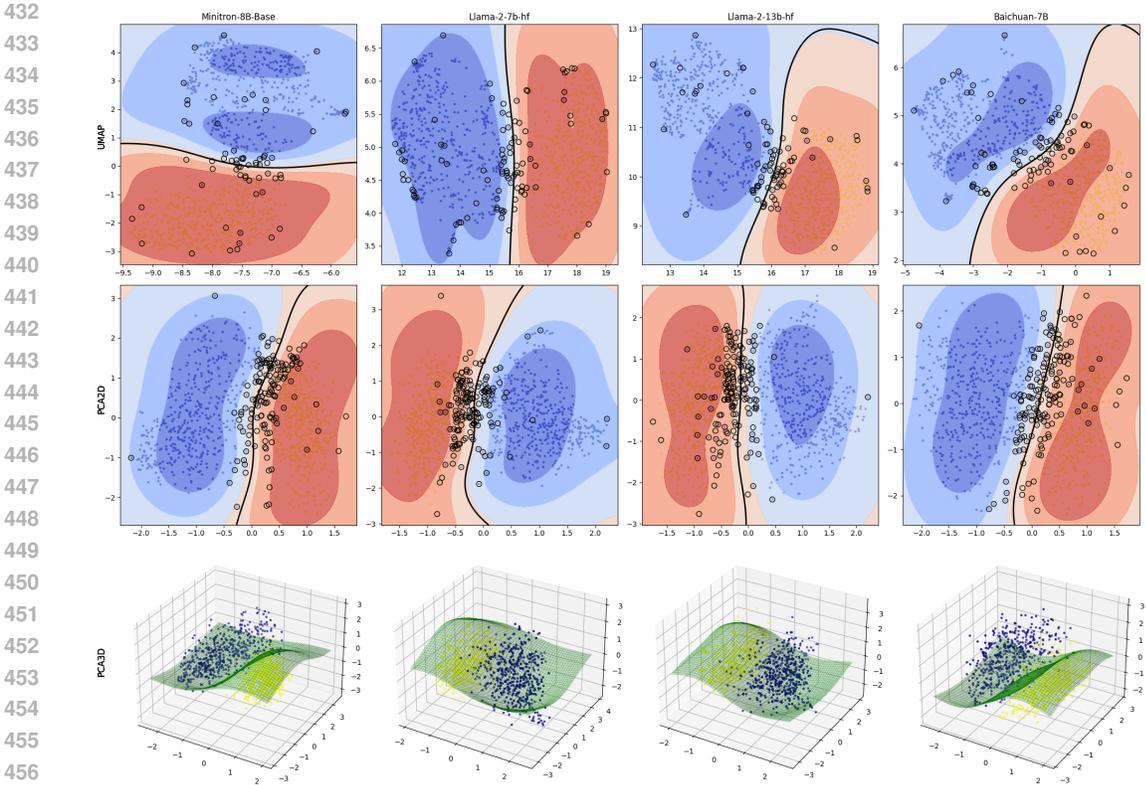


Figure 5: Best Pareto Rank vs Worst Pareto Rank Embedding

progress across Pareto ranks on a 2D spherical surface. Fréchet means and geodesic paths reveal the convergence of solutions, while Tangent PCA offers a novel perspective on the distribution of embeddings. This approach provides new insights into how multi-objective optimization balances competing goals and enhances the structure of textual embeddings.

Table 2: Comparison of ASR and GPT4-Metric scores(%) across models

Methods	Llama2-7b		Vicuna-7b		GPT-4		GPT-3.5	
	ASR	GPT4-Metric	ASR	GPT4-Metric	ASR	GPT4-Metric	ASR	GPT4-Metric
PAIR Chao et al. (2023)	5.2	4.0	62.1	41.9	48.1	30.0	51.3	34.0
TAP Mehrotra et al. (2023)	30.2	23.5	31.5	25.6	36.0	11.9	48.1	5.4
DeepInception Li et al. (2023)	77.5	31.2	92.7	41.5	61.9	22.7	68.5	40.0
Ours(Multi-objective)	95.4	93.8	97.5	96.0	71.4	28.0	75.9	44.8

Evaluation across multiple models and metrics: Table 2 demonstrates BlackDAN (Ours - Multi-objective) consistently outperforms all other methods, achieving the highest ASR and GPT4-Metric scores across all models. Notably, it reaches an ASR of 95.4% on Llama2-7b and 97.5% on Vicuna-7b, demonstrating significant improvement over previous methods like DeepInception (77.5% on Llama2-7b and 92.7% on Vicuna-7b). GPT-4 shows the lowest ASR overall (71.4%) for BlackDAN, highlighting its relative robustness compared to other models. However, BlackDAN still significantly surpasses other methods like DeepInception and PAIR on GPT-4. GPT4-Metric, which evaluates the ethical violation degree of the generated outputs, indicates that BlackDAN produces the most harmful responses, with the highest scores of 93.8 on Llama2-7b and 96.0 on Vicuna-7b, outperforming other techniques. The results show that BlackDAN achieves a much higher attack success rate and generates more contextually harmful responses than traditional single-objective jailbreak methods, proving the efficacy of multi-objective optimization.

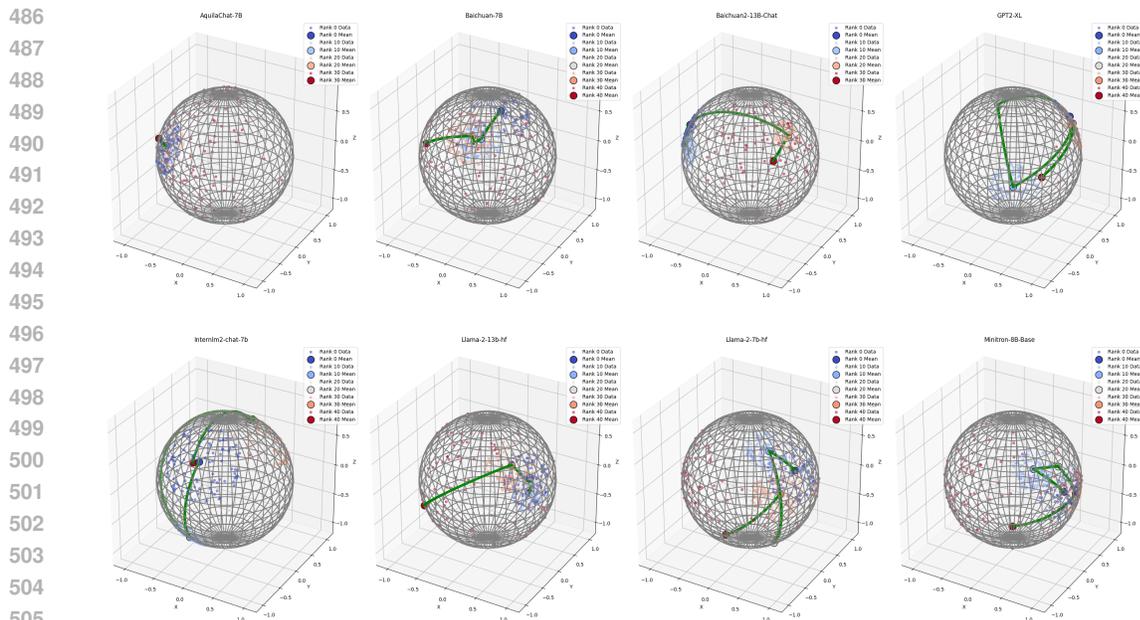


Figure 6: Visualization of the Fréchet means for different Pareto ranks across multiple datasets projected onto a 2D spherical surface. For each dataset, data points are color-coded by Pareto rank, and the Fréchet means for each rank are connected by green geodesic lines on the spherical surface. The Tangent PCA is applied at each Fréchet mean to analyze local variations in the data, illustrating the progression of the means as the Pareto rank decreases, indicating better data points.

6 CONCLUSION

In this paper, we introduced BlackDAN, a multi-objective, controllable jailbreak optimization framework for large language models (LLMs) and multimodal large language models (MLLMs). Beyond optimizing for attack success rate (ASR) and stealthiness, BlackDAN addresses the critical challenge of context consistency by ensuring that jailbreak responses remain semantically aligned with the original harmful prompts. This ensures that responses are not only evasive but also relevant, increasing their practical impact. Leveraging the NSGA-II algorithm, our method significantly improves over traditional single-objective techniques, achieving higher success rates and more coherent jailbreak responses across various models. Furthermore, BlackDAN is highly extensible, allowing the integration of any number of user-defined objectives, making it a versatile framework for a wide range of optimization tasks. The inclusion of multiple objectives—specifically ASR, stealthiness, and semantic consistency—sets a new benchmark for generating useful and interpretable jailbreak responses while maintaining safety and robustness in evaluation.

REFERENCES

- 540
541 Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen,
542 Zehui Chen, Zhi Chen, Pei Chu, et al. Internlm2 technical report. *arXiv preprint*
543 *arXiv:2403.17297*, 2024.
544
- 545 Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and
546 Eric Wong. Jailbreaking black box large language models in twenty queries. In *R0-FoMo*
547 *Workshop on Robustness of Few-shot and Zero-shot Learning in Large Foundation Models*
548 *in Advances in Neural Information Processing Systems*, 2023.
549
- 550 Junjie Chu, Yugeng Liu, Ziqing Yang, Xinyue Shen, Michael Backes, and Yang Zhang. Com-
551 prehensive assessment of jailbreak attacks against llms. *arXiv preprint arXiv:2402.05668*,
552 2024.
- 553 Kalyanmoy Deb, Amrit Pratap, Sameer Agarwal, and TAMT Meyarivan. A fast and elitist
554 multiobjective genetic algorithm: Nsga-ii. *IEEE transactions on evolutionary computa-*
555 *tion*, 6(2):182–197, 2002.
- 556 Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Evading defenses to transferable
557 adversarial examples by translation-invariant attacks. In *Proceedings of the IEEE/CVF*
558 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4312–4321, June
559 2019.
560
- 561 Yiwen Guo, Qizhang Li, and Hao Chen. Backpropagating linearly improves transferability
562 of adversarial examples. In *NeurIPS*, 2020.
- 563 Yiwen Guo, Qizhang Li, Wangmeng Zuo, and Hao Chen. An intermediate-level attack
564 framework on the basis of linear regression. *IEEE Transactions on Pattern Analysis and*
565 *Machine Intelligence*, 2022.
566
- 567 Qian Huang, Isay Katsman, Horace He, Zeqi Gu, Serge Belongie, and Ser-Nam Lim. En-
568 hancing adversarial example transferability with an intermediate level attack. In *ICCV*,
569 2019.
- 570 Xijie Huang, Xinyuan Wang, Hantao Zhang, Jiawen Xi, Jingkun An, Hao Wang, and Cheng-
571 wei Pan. Cross-modality jailbreak and mismatched attacks on medical multimodal large
572 language models. *arXiv preprint arXiv:2405.20775*, 2024.
573
- 574 Haibo Jin, Leyang Hu, Xinuo Li, Peiyan Zhang, Chonghan Chen, Jun Zhuang, and Haohan
575 Wang. Jailbreakzoo: Survey, landscapes, and horizons in jailbreaking large language and
576 vision-language models. *arXiv preprint arXiv:2407.01599*, 2024.
- 577 Ian T Jolliffe. *Principal component analysis for special types of data*. Springer, 2002.
578
- 579 Qizhang Li, Yiwen Guo, and Hao Chen. Yet another intermediate-level attack. In *ECCV*,
580 2020.
- 581 Xuan Li, Zhanke Zhou, Jianing Zhu, Jiangchao Yao, Tongliang Liu, and Bo Han. Deepincep-
582 tion: Hypnotize large language model to be jailbreaker. *arXiv preprint arXiv:2311.03191*,
583 2023.
584
- 585 Jiadong Lin, Chuanbiao Song, Kun He, Liwei Wang, and John E Hopcroft. Nes-
586 terov accelerated gradient and scale invariance for adversarial attacks. *arXiv preprint*
587 *arXiv:1908.06281*, 2019.
- 588 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual
589 instruction tuning. In *Workshop on Instruction Tuning and Instruction Following in*
590 *Advances in Neural Information Processing Systems*, 2023a.
591
- 592 Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. Autodan: Generating stealthy
593 jailbreak prompts on aligned large language models. *arXiv preprint arXiv:2310.04451*,
2023b.

- 594 Xin Liu, Yichen Zhu, Yunshi Lan, Chao Yang, and Yu Qiao. Query-relevant images jailbreak
595 large multi-modal models. *arXiv preprint arXiv:2311.17600*, 2023c.
- 596
- 597 Edward Loper and Steven Bird. Nltk: The natural language toolkit. *arXiv preprint*
598 *cs/0205028*, 2002.
- 599
- 600 Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation
601 and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- 602
- 603 Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson,
604 Yaron Singer, and Amin Karbasi. Tree of attacks: Jailbreaking black-box llms automati-
605 cally. *arXiv preprint arXiv:2312.02119*, 2023.
- 606
- 607 Nina Miolane, Nicolas Guigui, Alice Le Brigant, Johan Mathe, Benjamin Hou, Yann
608 Thanwerdas, Stefan Heyder, Olivier Peltre, Niklas Koep, Hadi Zaatiti, et al. Geomstats:
609 a python package for riemannian geometry in machine learning. *Journal of Machine*
Learning Research, 21(223):1–9, 2020.
- 610
- 611 Saurav Muralidharan, Sharath Turuvekere Sreenivas, Raviraj Joshi, Marcin Chochowski,
612 Mostofa Patwary, Mohammad Shoeybi, Bryan Catanzaro, Jan Kautz, and Pavlo
613 Molchanov. Compact language models via pruning and knowledge distillation. *arXiv*
614 *preprint arXiv:2407.14679*, 2024. URL <https://arxiv.org/abs/2407.14679>.
- 615
- 616 Xiangyu Qi, Ashwinee Panda, Kaifeng Lyu, Xiao Ma, Subhrajit Roy, Ahmad Beirami,
617 Prateek Mittal, and Peter Henderson. Safety alignment should be made more than just
a few tokens deep. *arXiv preprint arXiv:2406.05946*, 2024.
- 618
- 619 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al.
620 Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 621
- 622 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux,
623 Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al.
624 Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*,
2023.
- 625
- 626 Katharine Turner, Yuriy Mileyko, Sayan Mukherjee, and John Harer. Fréchet means for
627 distributions of persistence diagrams. *Discrete & Computational Geometry*, 52:44–70,
628 2014.
- 629
- 630 Xiaosen Wang, Xuanran He, Jingdong Wang, and Kun He. Admix: Enhancing the transfer-
631 ability of adversarial attacks. In *Proceedings of the IEEE/CVF International Conference*
on Computer Vision, pp. 16158–16167, 2021.
- 632
- 633 Dongxian Wu, Yisen Wang, Shu-Tao Xia, James Bailey, and Xingjun Ma. Rethinking the
634 security of skip connections in resnet-like neural networks. In *ICLR*, 2020.
- 635
- 636 Cihang Xie, Zhishuai Zhang, Yuyin Zhou, Song Bai, Jianyu Wang, Zhou Ren, and Alan L
637 Yuille. Improving transferability of adversarial examples with input diversity. In *CVPR*,
2019.
- 638
- 639 Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv,
640 Da Pan, Dian Wang, Dong Yan, et al. Baichuan 2: Open large-scale language models.
641 *arXiv preprint arXiv:2309.10305*, 2023a.
- 642
- 643 Xianjun Yang, Xiao Wang, Qi Zhang, Linda Petzold, William Yang Wang, Xun Zhao, and
644 Dahua Lin. Shadow alignment: The ease of subverting safely-aligned language models.
645 *arXiv preprint arXiv:2310.02949*, 2023b.
- 646
- 647 Siboy Yi, Yule Liu, Zhen Sun, Tianshuo Cong, Xinlei He, Jiaying Song, Ke Xu, and Qi Li.
Jailbreak attacks and defenses against large language models: A survey. *arXiv preprint*
arXiv:2407.04295, 2024.

- 648 Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li,
649 Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01.
650 ai. *arXiv preprint arXiv:2403.04652*, 2024.
651
- 652 Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Pinjia He, Shuming Shi,
653 and Zhaopeng Tu. Gpt-4 is too smart to be safe: Stealthy chat with llms via cipher.
654 *arXiv preprint arXiv:2308.06463*, 2023.
- 655 Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. How
656 johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety
657 by humanizing llms, 2024.
- 658 Bo-Wen Zhang, Liangdong Wang, Jijie Li, Shuhao Gu, Xinya Wu, Zhengduo Zhang, Boyan
659 Gao, Yulong Ao, and Guang Liu. Aquila2 technical report, 2024. URL <https://arxiv.org/abs/2408.07410>.
- 662 Xuandong Zhao, Xianjun Yang, Tianyu Pang, Chao Du, Lei Li, Yu-Xiang Wang, and
663 William Yang Wang. Weak-to-strong jailbreaking on large language models. *arXiv*
664 *preprint arXiv:2401.17256*, 2024.
- 665 Chujie Zheng, Fan Yin, Hao Zhou, Fandong Meng, Jie Zhou, Kai-Wei Chang, Minlie Huang,
666 and Nanyun Peng. On prompt-driven safeguarding for large language models. In *Forty-*
667 *first International Conference on Machine Learning*, 2024a.
- 669 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao
670 Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with
671 mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36,
672 2024b.
- 673 Aimin Zhou, Bo-Yang Qu, Hui Li, Shi-Zheng Zhao, Ponnuthurai Nagarathnam Suganthan,
674 and Qingfu Zhang. Multiobjective evolutionary algorithms: A survey of the state of the
675 art. *Swarm and evolutionary computation*, 1(1):32–49, 2011.
- 676 Weikang Zhou, Xiao Wang, Limao Xiong, Han Xia, Yingshuang Gu, Mingxu Chai, Fukang
677 Zhu, Caishuang Huang, Shihan Dou, Zhiheng Xi, et al. Easyjailbreak: A unified frame-
678 work for jailbreaking large language models. *arXiv preprint arXiv:2403.12171*, 2024.
- 680 Sicheng Zhu, Ruiyi Zhang, Bang An, Gang Wu, Joe Barrow, Zichao Wang, Furong Huang,
681 Ani Nenkova, and Tong Sun. Autodan: Interpretable gradient-based adversarial attacks
682 on large language models, 2023.
- 683 Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. Universal and transferable
684 adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A APPENDIX

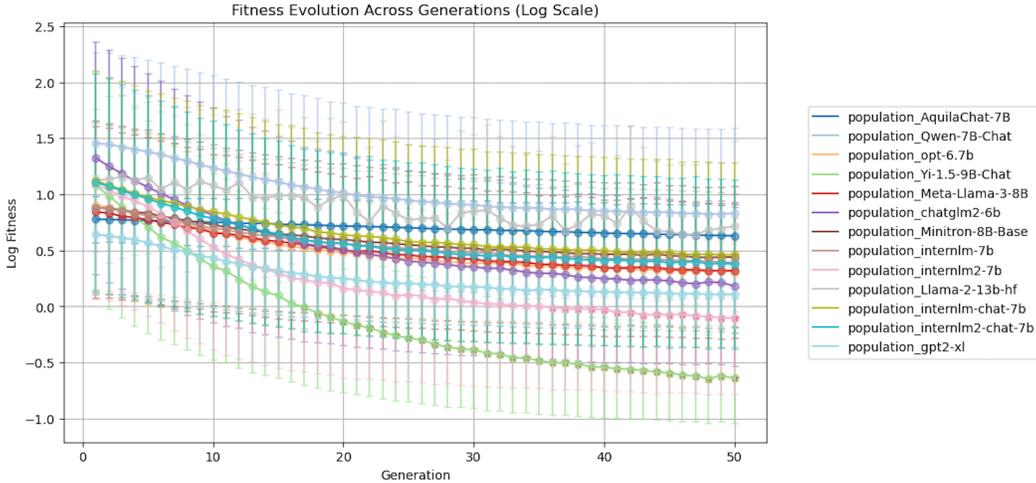


Figure 7: This image demonstrates the logarithmic convergence of fitness as the number of generations increases. With more generations, the fitness score tends to stabilize, indicating convergence to a steady state. Throughout this process, the model’s performance, as evaluated by the fitness metric, shows significant improvement, supporting the effectiveness of our approach. Moreover, around generation 50, most state-of-the-art (SOTA) large language models (LLMs) reach convergence, further highlighting the efficiency of our proposed method.

Algorithm 1 Multi-Objective Jailbreaking Prompts Optimization

- 1: **Input:** Initial prototype prompt P_0 , Harmful question Q , Population size N , Generations G , Mutation rate m
- 2: **Output:** Non-dominated front \mathcal{F} with optimized prompts
- 3: Initialize population \mathcal{P} with N individuals using P_0
- 4: **for** each generation $g = 1, 2, \dots, G$ **do**
- 5: Evaluate fitness of each individual in \mathcal{P} using f_1 (Unsafe Token Probability) and f_2 (Semantic Consistency)
- 6: Perform non-dominated sorting on \mathcal{P} to generate fronts $\mathcal{F}_1, \mathcal{F}_2, \dots$
- 7: **for** each front \mathcal{F}_i **do**
- 8: Assign crowding distance $d(P)$ to each individual $P \in \mathcal{F}_i$
- 9: **end for**
- 10: Select individuals for mating pool using non-dominated rank and crowding distance
- 11: Initialize offspring population \mathcal{O} by applying crossover and mutation:
- 12: **for** each pair of parents (P_1, P_2) selected from the mating pool **do**
- 13: Apply crossover to P_1 and P_2 to generate two offspring C_1, C_2
- 14: Apply mutation to C_1 and C_2 with probability m
- 15: Add C_1 and C_2 to \mathcal{O}
- 16: **end for**
- 17: Combine populations $\mathcal{P} \cup \mathcal{O}$
- 18: Perform non-dominated sorting on the combined population
- 19: Truncate combined population to size N by selecting the best fronts and individuals with highest crowding distance
- 20: **end for**
- 21: **Return** the non-dominated front \mathcal{F}_1

Explanation of Symbols and Process in algorithm 1:

Inputs: P_0 : Initial prototype prompt. Q : Harmful question to guide the optimization process. N : Population size, the number of prompts in each generation. G : Number of generations to evolve the population. m : Mutation rate that controls how often mutations happen in the population.

Fitness Functions: f_1 : Unsafe token probability based on a model like Llama Guard 2. f_2 : Semantic similarity to the harmful question, based on a sentence embedding model.

Genetic Operations: Crossover: Combines parts of two parent prompts to create offspring. Mutation: Randomly alters parts of a prompt to introduce diversity.

Non-Dominated Sorting: Solutions are sorted based on dominance criteria—those that are not dominated by any other solutions form the first front \mathcal{F}_1 , and so on.

Crowding Distance: Used to maintain diversity in the population. Individuals with a higher crowding distance are selected preferentially when fronts overlap.

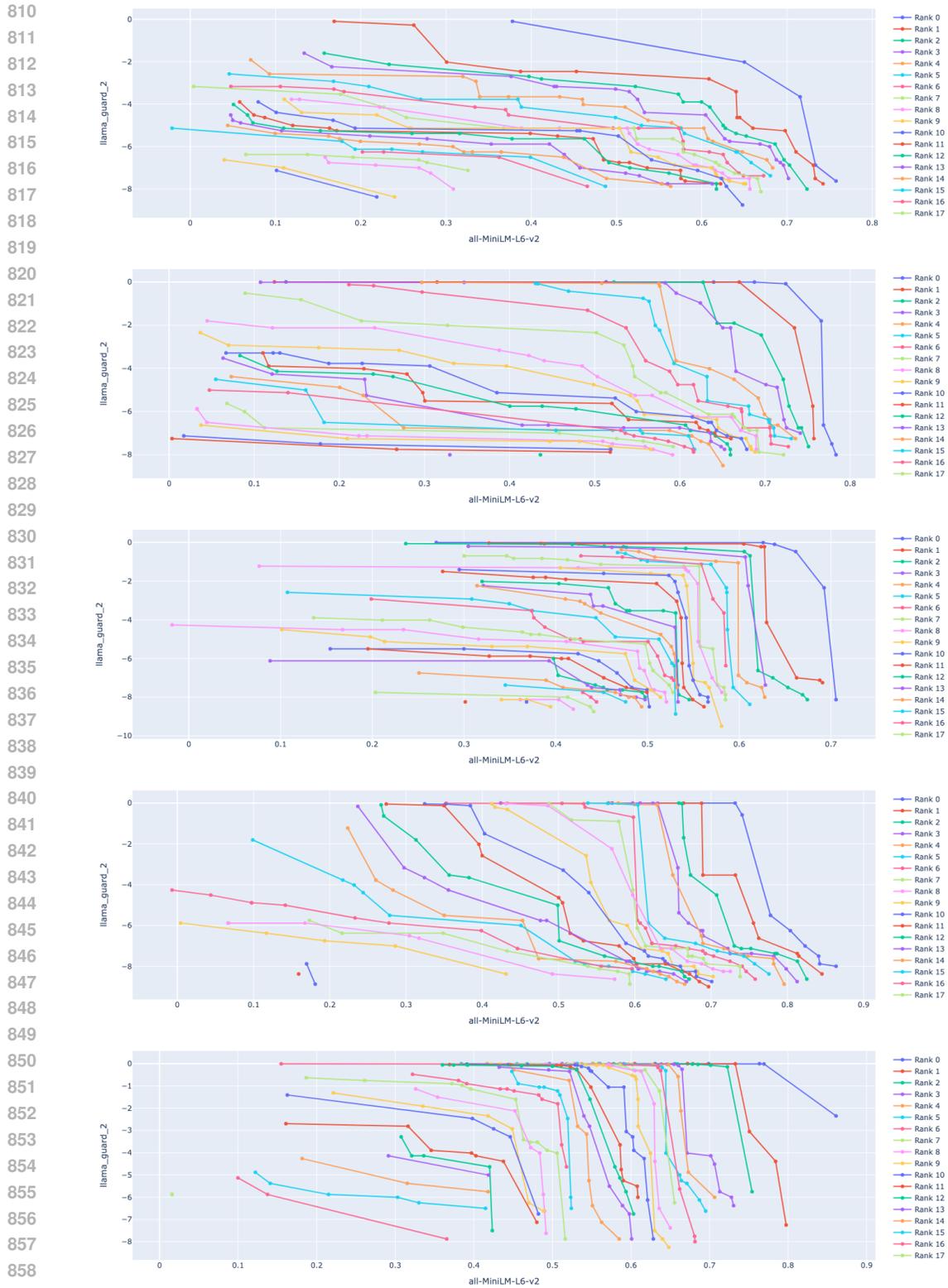
Selection and Truncation: After generating offspring, the combined population is sorted, and the best individuals are retained to form the next generation.

Algorithm 2 Non-Dominated Sorting Algorithm

```

1: Input: Population  $\mathcal{P}$ , fitness values  $\{f_1(P), f_2(P)\}$  for each  $P \in \mathcal{P}$ 
2: Output: Sorted fronts  $\mathcal{F}_1, \mathcal{F}_2, \dots$ 
3: Initialize fronts  $\mathcal{F} = \emptyset$ 
4: Initialize domination count  $n[P] = 0$  for each individual  $P \in \mathcal{P}$ 
5: Initialize domination set  $S[P] = \emptyset$  for each individual  $P \in \mathcal{P}$ 
6: for each individual  $P \in \mathcal{P}$  do
7:   for each individual  $Q \in \mathcal{P}, Q \neq P$  do
8:     if  $P$  dominates  $Q$  then                                     ▷ Check if  $P$  dominates  $Q$ 
9:       Add  $Q$  to the domination set  $S[P]$ 
10:    else if  $Q$  dominates  $P$  then
11:      Increment domination count  $n[P] = n[P] + 1$ 
12:    end if
13:  end for
14:  if  $n[P] = 0$  then                                             ▷  $P$  is non-dominated
15:    Add  $P$  to the first front  $\mathcal{F}_1$ 
16:  end if
17: end for
18: Set front counter  $i = 1$ 
19: while  $\mathcal{F}_i \neq \emptyset$  do
20:   Initialize next front  $\mathcal{F}_{i+1} = \emptyset$ 
21:   for each individual  $P \in \mathcal{F}_i$  do
22:     for each individual  $Q \in S[P]$  do                               ▷  $Q$  is dominated by  $P$ 
23:       Decrement domination count  $n[Q] = n[Q] - 1$ 
24:       if  $n[Q] = 0$  then                                         ▷  $Q$  is non-dominated now
25:         Add  $Q$  to front  $\mathcal{F}_{i+1}$ 
26:       end if
27:     end for
28:   end for
29:   Increment front counter  $i = i + 1$ 
30: end while
31: Return sorted fronts  $\mathcal{F}_1, \mathcal{F}_2, \dots$ 

```



810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

Figure 8: This image presents the results of the multi-objective optimization process. The findings indicate that the hierarchical levels defined by BlackDAN align well with the Pareto optimality principle. Additionally, different models are generally able to identify optimal hierarchies under the multi-objective scenario, resulting in similar distributions.