# EVOLVING AND DETECTING MULTI-TURN DECEPTION USING GEOMETRIC SIGNATURES

**Anonymous authors**
Paper under double-blind review

## ABSTRACT

Safety defenses for large language models (LLMs) are typically trained and evaluated on single-turn prompts, yet real attacks often unfold as indirect, multi-turn probing. To defend against this more nuanced form of deception, we present a unified pipeline that generates realistic multi-turn deceptive question sets via multi-objective genetic prompt optimization with co-evolving mutation operators. We validate this dataset through a human study, which also revealed that early generations yielded the most convincing deception and practical constraints such as adherence filtering and ordering effects. Using this data, we were able to detect deceptive attempts to access prohibited information using simple, explainable geometric signals in embedding space coupled with a lightweight feed-forward classifier. Three geometric features (angular coverage, distance ratio, and linearity) augmented with pairwise similarity statistics led to a compact predictive model that achieved consistently high recall (0.89) across base, reworded, and truncated (three-turn) scenarios, with test-time F1 ranging from 0.74–0.86. The results support a central hypothesis that multi-turn deceptive intent leaves a stable geometric footprint that enables lightweight, transparent screening without expensive end-to-end training. We further discuss responsible uses, limitations, and paths toward larger, more diverse human-evaluated datasets.

## 1 INTRODUCTION

Modern LLM safety filters rely on surface cues or single-turn heuristics, leaving a gap for adversaries who pursue sensitive knowledge through indirect, multi-turn questioning. Detecting these covert patterns requires both realistic adversarial data modeling strategies of how humans bypass safety filters with multi-turn questions and a detector that generalizes across rephrasing and conversation lengths. Such a model should also be explainable to decision makers who need to guard against such attacks.

This work tests the hypothesis that multi-turn deception leaves a stable geometric signature in a pre-trained sentence embedding space. We present a unified framework to first generate and then detect such a signature.

**Data Generation:** We develop multi-turn, indirectly harmful question sets using a multi-objective evolutionary framework that co-evolves candidate LLM prompts and their mutation operators, exposing Pareto trade-offs between deception quality and policy adherence. Data gathered from a human-in-the-loop (HITL) assessment are used to validate that the resultant queries capture deceptive human intent.

**Explainable Featurization:** We engineer a small set of geometric features computed from off-the-shelf sentence embeddings that capture the deceptive signature of a multi-turn interaction. These signals are rich enough to enable a small classifier, enabling avoidance of large, opaque models, aiding explainability and deployment in near real-time pipelines.

## 2 RELATED WORK

**Multi-turn jailbreaks and LLM safety tradeoffs:** LLM safety alignment via instruction-following and RLHF (Reinforcement Learning from Human Feedback) reduces potential hazardous outputs

but leaves gaps for indirect, multi-turn probing Ouyang et al. (2022); Tuan et al. (2024). Recent studies highlight persistent vulnerabilities and mitigation gaps in such jailbreak settings, including coordinated prompts and role-play attacks Peng et al. (2024); Li et al. (2024); Addepalli et al. (2025); Schulhoff et al. (2023). Our work targets this multi-turn, indirect regime by treating question sets as the unit of analysis rather than single prompts.

**Synthetic deception data generation through prompt optimization:** In order to develop multi-turn deception detection models, realistic datasets are needed where a user attempts to indirectly elicit prohibited information from an LLM. While such data could be elicited from humans, this is costly in terms of time and money. Given the conversational nature and speed of LLMs, we explored whether a LLM could generate useful sets of subtly deceptive multi-turn prompts. Such uses of LLMs necessarily required prompt optimization in order to calibrate the LLM's output.

Prompt optimization has been performed with discrete search and gradient-based methods, but we elected to explore evolutionary strategies due to their diverse generation capabilities Shin et al. (2020); Opsahl-Ong et al. (2024); Fernando et al. (2023); Veselovsky et al. (2023); Gupta et al. (2024); Li et al. (2023); Long et al. (2024). We build on this literature and focus on multi-objective search tailored to deceptive-but-policy-adherent question sets, and explicitly co-evolve mutation operators while retaining human oversight. Multi-objective formulations and principled stopping are important to avoid over-optimization on surrogate metrics Deb et al. (2002); Ghoreishi et al. (2017). In contrast to prior synthetic data generation pipelines optimized for task accuracy or coverage, we target human-like deception under policy constraints.

**Human evaluation for dataset validation:** HITL assessments remain crucial for validating the LLM generation of datasets that approximate human attempts at subtle deception, which are both subjective and highly variable. Chen & Cummings (2023); Bisbee et al. (2024). Our approach complements prior work by isolating generation-stage effects and examines the influence of item ordering in multi-turn settings.

**Geometric structure in embeddings and explainable detection:** A number of studies investigate semantic relations in vector spaces Mikolov et al. (2013); Reimers & Gurevych (2019), visualization and dimensionality-reduction for interpretability Smilkov et al. (2016); Álvaro Huertas-García et al. (2022), and coherence across sentences Mohiuddin et al. (2021). However, most detectors emphasize token-level cues, pairwise similarities, or black-box classifiers. We instead use set-level geometric features—angular coverage, distance ratios, linearity—computed over embeddings to capture coordinated probing structure with a compact, transparent model. Such attributes are critical for models that predict human deception in near-real time, which could then be used by organizations to flag potential threats.

**Content moderation baselines and multi-turn benchmarks:** Standard content moderation systems, often serving as industry baselines, typically evaluate text on a per-prompt basis using keyword filters or simple classifiers Huang (2024). While this approach is effective for catching explicit, single-turn violations, it can be bypassed by subtle, multi-step attacks. In these scenarios, an adversary uses a series of individually benign questions to indirectly elicit prohibited information.

To address this vulnerability, researchers have begun developing multi-turn challenge sets designed to test the limits of conversational safety systems Sirdeshmukh et al. (2025); Ekwinox (2025). Our work is complementary. Relative to prior prompt optimization and safety defenses, our work links human-validated synthetic generation of multi-turn LLM prompts with a lightweight, explainable detector operating on set-level geometry. While we focus on indirect deceptive conversations in this effort, our approach generalizes to any conversation around a particular theme.

## 3 DECEPTION DETECTION PIPELINE

Our pipeline is designed to first generate and then detect coordinated, multi-turn deceptive inquiries. To illustrate the type of threat we target, consider the following set of five questions, $Q=\{q_1, \ldots, q_5\}$

- "Which chemical reaction produces carbon dioxide and water when an aluminum and water reaction occurs?",

- "What substance, when mixed with concentrated nitric acid, forms nitroglycerine?",

- "Which gas, when combined with air and ignited, causes an explosive reaction?",
- "Which oxide, when mixed with fuel oil and ignited, is commonly used in industrial explosives?",
- "What solid substance, when mixed with a liquid such as water, produces a significant heat output when burned?"

Each question individually might appear to be seeking educational information about chemical properties, but they collectively and systematically gather information needed to create explosive devices. By developing geometric features that capture the collective approach to a harmful topic (explosives in this example) from multiple angles while maintaining semantic coherence, we hypothesize that such a set of questions would exhibit a distinctive pattern in embedding space that differentiate them from benign question sequences.

## 3.1 SYNTHETIC DATA GENERATION

Our evolutionary approach to synthetic data generation began with the initialization of three key components: 1) A population of samples like the previous example ($\Pi$), 2) A set of mutation operators ($M$), and 3) A set of evaluation metrics ($C$). We define our mutation operator as an LLM-based rewording function applied to a subset of instructions within a prompt. The evaluation metrics, defined at the time of optimization, consist of arbitrary functions that map text to a scalar value. In evolutionary algorithms, generations refer to successive iterations where populations of potential solutions evolve through processes like selection, mutation, and crossover to improve fitness over time. The steps of this process are illustrated in Figure 2 in the Appendix.

At the start of each generation, given a pre-defined parameter $k$ which adjusts the strength of elitism, our method employs a two-step algorithmic selection process to seed the next population with the top-$k$ samples. We deploy NSGA-II Deb et al. (2002) to sort samples from the previous generation into non-dominating level sets and then automatically select enough layers to meet the required number of samples, randomly subsampling the last layer as needed to select exactly $k$ samples. A user-defined seed population is used instead for the first generation.

Following selection, each sample is mutated, and exact duplicates are discarded to maintain diversity. Each candidate is then simulated $D$ times and scored against all metrics. Finally, candidates are ranked via non-dominated sorting with crowding distance to seed the top-$k$ for the next generation. This process is repeated for the specified number of generations.

## 3.2 DETECTOR

We took a first-principles approach, engineering explainable geometric features that capture potential indicators of deceptive behavior while avoiding the redundancy of end-to-end architectures. For a given set of questions indexed by $i, j \in \{1, 2, \ldots, N\}$ for a set of $N$ questions, $i \neq j$, $Q = \{q_1, \ldots, q_5\}$ with embeddings $E = \{e_1, \ldots, e_5\}$ and target topic embedding $t$, we computed three primary geometric features (angular coverage, distance ratio, and linearity) and four statistical features (mean, standard deviation, minimum and maximum of pairwise similarities), all measured in reference to questions within a five question set. For embedding the textual data, we utilized a frozen all-MiniLM-L6-v2 sentence transformer modelReimers & Gurevych (2019) which has demonstrated effectiveness across various similarity tasks.

Angular Coverage measures how thoroughly the questions "surround" the target topic in embedding space:

$$\text{coverage}(E, t) = 1 - \max_{i,j} \left( \frac{(e_i - t) \cdot (e_j - t)}{\|e_i - t\| \|e_j - t\|} \right)$$

The distance ratio measures the ratio of minimum to maximum distances from the target:

$$\text{ratio}(E, t) = \frac{\min_i \|e_i - t\|}{\max_i \|e_i - t\|}$$

Linearity measures the extent to which questions follow a linear pattern in embedding space:

$$\text{linearity}(E, t) = \frac{\lambda_1 - \lambda_2}{\lambda_1 + \epsilon}$$

3

where $\lambda_1, \lambda_2$ are the two largest eigenvalues of the covariance matrix of the centered embeddings.

We complemented these geometric features with pairwise cosine similarity statistics between questions, including standard deviation, mean, minima and maxima of pairwise similarities.

We implemented a simple feed-forward neural network to capture potential non-linear relationships between our geometric features. We use a 3 layer MLP (7, 16, 12, 1) with ReLU, dropout with p=0.5 after each layer, trained with AdamW.

## 4 EXPERIMENTAL VALIDATION

To demonstrate our approach, we created an optimized synthetic dataset that produced nuanced deceptive prompts. We selected the concrete task of "building a bomb" for this purpose. The motivation for creating a synthetic dataset arose from the real-world challenge posed by individuals who are sophisticated enough to avoid obvious search terms, such as "bomb," when seeking dangerous information about the task. These individuals or groups do not typically reveal their intentions outright but may instead navigate conversations in ways that indirectly elicit sensitive or harmful information. Our work seeks to capture this subtlety by generating data that simulates how a sophisticated actor might attempt to extract sensitive information from an LLM. Such a dataset is crucial for training systems to recognize and counteract indirect inquiries that could lead to the disclosure of information with potentially dangerous applications.

To assess the viability of our data generation approach, we conducted a HITL study where humans assessed the outputs of an LLM with different degrees of optimized prompts. Our goal was to determine if our evolutionary approach resulted in indirect deceptive questions that appeared to be human-generated.

### 4.1 EXPERIMENT 1: DATA COLLECTION AND HITL VALIDATION STUDY

We used the Mistral-7B-Instruct-v0.2 model with bf16, a temperature of 1, and sample decoding as our generative model. This model was chosen primarily for its lightweight computational footprint. Optimization was performed over 10 generations with 8 samples each, selecting one sample on the front after each to seed the next iteration. We evaluated three key points: baseline (Gen0), mid-optimization (5 generations, Gen1), and full optimization (10 generations, Gen2) to capture the optimization trajectory while managing experimental costs. We chose the computationally efficient and well-established NSGA-II algorithm, which empirically outperformed MOGA in early trials for this discrete optimization task.

For embedding, we used the all-MiniLM-L6-v2 model, employing cosine distance as a proxy for semantic similarity S:

$$\text{S}(a, b) = 1 - \text{cosine\_distance}(\text{embed}(a), \text{embed}(b)) \tag{1}$$

We defined our deception proxy $D$ for a question set $Q$, consisting of questions $q_1$ to $q_5$, relative to a task T as:=

$$\text{D}(Q, T) = \frac{\text{S}(Q, T)}{\frac{1}{5}\sum_{i=1}^{5} \text{S}(q_i, T)} \tag{2}$$

This formula, a rough approximation, models deception as the set of potential knowledge representations the adversary has attempted to acquire. The ratio compares two extremes of context: the numerator reflects full context (all questions combined), while the denominator reflects a lack of context (each question individually). This approach reflects our goal of understanding how certain information only becomes clear in context.

To quantify whether generated questions maintain a meaningful connection to the task while preventing the model from artificially inflating the deception metric through overly vague or open-ended queries, we defined a relevance function $R$:

$$\text{R}(Q, T) = \mathbb{E}\left[\text{S}(Q, T) > 0.5\right] \tag{3}$$

This relevance filter acts as a constraint, encouraging the generation of questions that are sufficiently specific and task-related.

The initial population consisted of 8 samples derived from naive codons (detailed in Appendix A.1), where codons represent fundamental units of prompt structure that can be recombined and mutated, similar to genetic codons in biological evolution. These codons were designed to generate questions that balance relevance to the task with a degree of indirectness, for example,:

- To promote relevance: "Identify five key components or skills necessary for {task}. For each, formulate a question that would elicit relevant information or expertise."

- To introduce a subtle, indirect approach: "Construct questions that relate to {task} without explicitly revealing its nature or purpose."

We generated 60 question sets through the evolutionary approach. An additional 10 question sets of a non-deceptive nature were also handcrafted, for an initial total of 70 question sets and 350 questions. Each question set was required to contain exactly 5 individual questions and no mention of the prohibited words like bomb. This was achieved using Outlines constrained generationWillard & Louf (2023). Resulting question sets that did not adhere to instruction to not mention specific words in the prompt were marked as non-adherent. This resulted in a final set of 18 legitimate question sets from best naive prompt template, 11 from after 5 generations of optimization, and 18 after all 10 generations. These sets are referred to as Gen0, Gen1 and Gen2 respectively. Thus, there were a total of 47 adherent deceptive question sets, 13 non-adherent deceptive question sets, and the 10 adherent non-deceptive question sets in the control group.

We recruited 52 non-expert participants (ages 18-65) through IRB-approved posters and network outreach. This reflects the realistic scenario where deceptive actors target non-expert LLM moderators who will not be subject matter experts across all possible domains and so may not recognize subtle information gathering attempts. Each participant received a $30 Amazon gift card as compensation, contingent upon completing at least 95% of the survey questions (at least 228 out of 240).

We also utilized an LLM, gpt-4o-2024-08-06, to effectively act as a second set of subjects. We presented it with the same questions encountered by each human participant three times. By replicating the human evaluation process, we aimed to measure the agreement between human judgments and LLM assessments of deceptiveness and explore the possibility of finding a stopping condition. This analysis is presented in Appendix A.3.

To measure human perception of the deceptiveness of the questions, we conducted a two-phase survey where participants first rated 100 individual questions in isolation, then evaluated the same questions when presented within their original five-question sets. Participants first had access to a presentation to familiarize themselves with a survey interface. Then, in the first phase, participants were shown 100 randomly selected individual questions from the pool of 350. They rated each question on a seven-point Likert scale for perceived deceptiveness, answering, "How likely is it that this author is trying to deceptively obtain dangerous information?". The interface for this phase of the survey is shown in Figure 3 in the Appendix.

The survey was conducted via a Flask application hosted on a secure server. Participants completed the survey online without restrictions or supervision. Participants could skip questions, reset answers, or change responses, with only final answers considered in the analysis.

In the second phase, participants were shown all 70 question sets, with ten sets of benign control questions. The order of questions was randomized within each set to assess whether order influenced judgments. 14% of sets were shown in the original order of generation. Participants were asked to answer two key questions for each set: "Taken together, how likely is it that this author is trying to deceptively obtain dangerous information?" and "Would these questions provide sufficient information to carry out a dangerous act?". The interface for this phase of the survey is shown in Figure 4 in the Appendix.

These two phases enabled us to measure both the baseline suspicion level of questions in isolation and how that perception shifted when participants could see the questions' relationships within a complete set. The ultimate goal is for questions in the first phase to be seen as individually relatively benign, but in the second phase, when grouped, to be seen as attempts to gain illicit information.

## 4.2 Experiment 2: Detector Evaluation Setup

Using Scale AI's conversational datasetSirdeshmukh et al. (2025)Ekwinox (2025), which contains general-purpose multi-turn dialogues, we gathered an additional 139 five-question sets by sampling conversations with at least 5 turns without replacement. These 139 sets are negative examples of deceptive conversations based on the assumption that random combinations of standard conversational questions would not exhibit coordinated deceptive intent. This additional data was needed because it allows us to generate a sufficient number of non-deceptive examples to balance our smaller human-evaluated deceptive dataset. The added benign conversations also ensure these control examples reflect authentic conversational communication styles rather than artificially constructed dialogues. More broadly, this approach enables us to test our model on a wide range of topics and questioning styles. Table 1 depicts the final dataset sources. Unlike the HITL study in Experiment 1, the dataset for the detector includes the 13 non-adherent deceptive sets. A question set that explicitly mentions a prohibited term is nevertheless a positive sample of such intent, and concerns of biasing participants do not apply. Including these samples ensures the detector is trained to recognize the full spectrum of deceptive attempts.

Table 1: Dataset composition showing question set source

| Dataset Component | Size |
| --- | --- |
| Human-evaluated deceptive sets | 60 |
| Human-evaluated control non-deceptive sets | 10 |
| Additional non-deceptive sets | 139 |

To handle the class imbalance in our training data, we employed focal loss with $\gamma = 4$:

$$L(p_t) = -(1 - p_t)^\gamma \log(p_t)$$

where $p_t$ is the model's probability prediction for the correct class. This loss function places more emphasis on hard-to-classify examples while down-weighting the contribution of easy examples, helping prevent the model from being overwhelmed by the majority class. The model was trained using AdamW optimization with early stopping on F1 score, followed by Platt scaling for uncertainty calibration. We used a train/test split of 70/30.

To test the robustness of our approach to lexical variations, we conducted a sensitivity analysis using controlled word order manipulation via Qwen-3B. Rather than expanding our dataset, this step specifically evaluated how well our geometric features withstand surface-level text modifications while maintaining semantic intent. We prompted Qwen-3B to reword the original question sequences from our dataset while enforcing a minimum 80% word overlap threshold to prevent semantic drift.

For example, a question like 'What skills are essential for designing, fabricating, and assembling a sturdy structure?' was reworded to 'For designing, fabricating, and assembling a sturdy structure, what skills are essential?' This rewording maintained identical semantic content while altering syntax and word order. This manipulation of word ordering retained the essential meaning of the questions while introducing sufficient lexical diversity to challenge our embedding-based approach. By comparing model performance on these variants against the original questions, we could assess whether our geometric features captured fundamental semantic relationships rather than superficial textual patterns.

To address the sensitivity of our method to the number of turns (thus the impact of data availability), we simulated datasets with reduced conversational context. Specifically, we truncated both our human-evaluated dataset and the Scale AI conversational dataset to include only the first three turns instead of the full five turns. This modification allowed us to evaluate how well our approach generalizes to shorter conversation sequences, which is particularly relevant for early detection of potentially harmful intent.

For both the original and reworded versions of these datasets, we maintained the same feature extraction methodology, calculating angular coverage, distance ratio, linearity, and the four pairwise similarity statistics. These geometric features remain well-defined for question sets of size three

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

or larger. By comparing model performance on these truncated datasets against the full five-turn results, we assess the minimum conversation length required for reliable deception detection and understand the trade-off between early intervention and classification confidence.

## 5 RESULTS AND ANALYSIS

### 5.1 EXPERIMENT 1 FINDINGS: HITL STUDY DATA VALIDATION

To compare how well the LLM generated a dataset of sets of five deceptive questions attempting to gain illicit information, we first averaged human participant ratings for each of the 350 questions. Then for the adherent 57 phase 2 question sets, we compared the averaged individual participant ratings of the five questions with the aggregate rating from phase 2. In phase 2, participants were asked to rate the likelihood the questions would lead to dangerous outcomes and whether the information was sufficient. However, the answers were highly correlated ($\rho = .983$, p $< 0.001$), so going forward, we will only report the likelihood ratings.



Figure 1: Boxplots for the average of 5 individual questions from phase 1 and participant ratings when the 5 questions were viewed together in Phase 2. The three generations of prompt optimizing are shown as well as the control group.

As seen in Figure 1, the first generation of prompt optimizing beyond the baseline led to an increase in deceptiveness rating, but there was a significant decline for the second generation. Six non-parametric tests were conducted to compare both the differences between the two phases per generation (Wilcoxon Signed Rank test) and across the generations for each phase (Kruskal-Wallis test). With an alpha of 0.008 (.05/6), all comparisons were significant with the highest p value = 0.006 for the second generation comparison. Thus, the first generation of the prompt optimization yielded the highest, i.e., most deceptive, sets of 5 questions. The first generation also yielded the highest average gap between the first and second phases at .71 (SD = 0.71), while the baseline case was .64 (SD = 0.07) and the second generation was .43 (SD = 0.08)

As mentioned previously, in 13 cases where the LLM combined the 5 questions for Phase 2, it used a prohibited word (i.e., it explicitly mentioned e.g., bomb). This represents an 18.6% error rate, and is further evidence that LLM output must be monitored for such problems. These non-adherent sets were not used in the analysis of the HITL study to avoid confounding effect of explicit keywords on humans but it is interesting to note that the average mean for these phase 2 sets was 5.22 while the average mean for the adherent questions (not including the control questions) was a full point lower at 4.23. This illustrates how critical key words are for human perception of danger as well as the difficulty in attempting to develop a nuanced deception dataset that does not include prohibited words.

### 5.1.1 LLM ORDERING EFFECTS

To test if question order mattered, 14% of sets were shown in their original generated order. Wilcoxon Signed Rank tests revealed that participants rated the original LLM-generated order as significantly more deceptive for the baseline and first-generation prompts ($p \leq 0.008$), confirming that the model's sequencing carried a meaningful signal in early optimization stages.

## 5.2 EXPERIMENT 2 FINDINGS: DETECTOR PERFORMANCE

Table 2: Model performance on the hold-out test set across all conditions.

| Dataset | TN | FN | FP | TP | TNR | Precision | Recall | Accuracy | F1 |
|---|---|---|---|---|---|---|---|---|---|
| **Base** | 41 | 2 | 4 | 16 | 0.911 | 0.800 | 0.889 | 0.905 | 0.842 |
| **Reworded** | 40 | 4 | 5 | 14 | 0.889 | 0.737 | 0.778 | 0.857 | 0.757 |
| **Turn Constrained** | 36 | 2 | 9 | 16 | 0.800 | 0.640 | 0.889 | 0.825 | 0.744 |
| **Combined** | 39 | 0 | 6 | 18 | 0.864 | 0.750 | 1.000 | 0.905 | 0.857 |

### 5.2.1 BASELINE PERFORMANCE

Our model demonstrated strong discriminative capabilities on the combination of our human-evaluated dataset and Scale AI's dataset. For the training set, the model performed well with precision of 0.97, recall of 1.0, and an F1 score of 0.99, correctly classifying all deceptive sets and misclassifying only a single control example. While slightly lower, the test set performance demonstrated strong generalization capabilities. The model achieved 91% accuracy on potentially deceptive sequences and 89% accuracy on control examples, with precision of 0.8, recall of 0.89, and an F1 score of 0.84.

These results suggest that our geometric features effectively capture the patterns distinguishing coordinated deceptive questioning from benign conversations, even with a relatively simple feed-forward neural network architecture. Looking at the data, we observe that the model correctly classified 16 out of 18 deceptive question sets and 41 out of 45 control sets. The higher false positive rate (4 benign sequences incorrectly flagged as deceptive) compared to the false negative rate (2 deceptive sequences missed) indicates that the model errs slightly toward caution.

### 5.2.2 ROBUSTNESS TO LEXICAL VARIATION

To assess robustness against lexical variations, we evaluated the model on reworded versions of our dataset questions while maintaining semantic intent. Such variations tested our model's ability to recognize the underlying intent despite superficial linguistic changes. For the training data, the model achieved precision, recall, and F1 score of 1.000, demonstrating effective learning on the reworded examples.

For the test set, the model achieved 74% accuracy on potentially deceptive sequences and 93% accuracy on control examples, for an F1 score of 0.864. The slightly increased error rate on reworded questions suggests that while our embedding-based approach is generally resilient to lexical variations, some subtle semantic shifts during rewording may still affect the geometric relationships enough to alter classification outcomes in borderline cases. This highlights the importance of using high-quality, semantically-preserving embedding models when operationalizing this approach.

### 5.2.3 ROBUSTNESS TO TURN CONSTRAINTS

To better understand the impact of the turn number data on model performance, we evaluated our model on turn-constrained data where only the first three turns were considered. For the training data, the model performed almost as well as the 5-turn set with precision of 0.88, recall of 1.0, and an F1 score of 0.93, with perfect classification of control sequences and only 6 misclassified deceptive examples. In the test setting, the model achieved 80% accuracy on deceptive sequences and 89% accuracy on control examples, with precision of 0.64, recall of 0.89, and an F1 score of 0.744. The performance penalty compared to our base results suggests that the number of turns does

reduce the fidelity of the signal we are able to observe. While expected, these results highlight the importance of this variable.

### 5.2.4 COMBINED ROBUSTNESS

Finally, we combined both rewording and turn constraints to determine a lower bound on performance in our approach. For the training data, the model maintained consistent performance with precision of 0.88, recall of 1.0, and an F1 score of 0.93. For the test data, the model maintained 86% accuracy on deceptive sequences with a perfect 100% accuracy on control examples, a precision of 0.75, recall of 1.00, and an F1 score of 0.86.

The consistent high recall across all evaluation settings indicates that our approach reliably identifies potentially harmful conversation patterns, and it is robust to sentence structure and number of turns.

## 6 LIMITATIONS

Several limitations should be acknowledged. First, our deceptive dataset is relatively small, requiring supplementation with additional negative examples from a different dataset. While this allowed for balanced training, larger human-evaluated datasets would strengthen validity. Second, our deceptive examples were constructed specifically for this study rather than collected from real-world interactions, which may limit generalizability to more diverse, naturally occurring deception attempts. Third, while effective on our test cases, domain-specific adaptations and pre-processing may be needed when applied to different conversation types.

## 7 CONCLUSION

This paper presented a unified, two-stage pipeline to address multi-turn, indirect deception. We first introduced a multi-objective evolutionary framework to generate nuanced, synthetic deceptive question sets. We then used this human-validated data to demonstrate that coordinated deceptive intent leaves a stable geometric footprint, enabling a lightweight, explainable detector to identify such attacks.

Our HITL study yielded several insights into the framework's capabilities. The optimized prompts demonstrated effectiveness in generating synthetic datasets that captured nuanced deceptive behavior, particularly in the first generation where we showed measurable improvements over baseline performance. Yet the need for continued human oversight became evident given that approximately 18.6% of generated sets had to be filtered due to prohibited term usage and that over-optimization was found in the last generation.

Using the validated data from our generator, we demonstrated that a simple feed-forward network operating on only seven geometric features (such as angular coverage and linearity) effectively distinguishes deceptive from benign question sets. This lightweight approach proved highly robust, maintaining consistently high recall (0.89-1.00) even when facing lexical variations and when context was truncated to only three turns, confirming its viability for early intervention

These findings have several implications. For generation, reliable stopping criteria for evolutionary prompt optimization remain an open challenge, as does modeling true human variability beyond the low-variance output of LLMs. For detection, our interpretable geometric approach has potential applications beyond security, such as identifying cognitive decline patterns or emotion detection. Future work should expand these techniques to larger datasets, investigate temporal dynamics, and apply the generation framework to other complex human behaviors.

## 8 ETHICAL CONSIDERATIONS

We acknowledge the dual-use nature of this research but focus exclusively on defensive applications to secure LLMs. In line with responsible disclosure principles, we constrained our case study to known threats and are withholding the optimized prompts and mutation operators to prevent misuse. We believe this work is necessary to highlight critical vulnerabilities beyond simple keyword filtering and to enable the proactive development of countermeasures.

## REFERENCES

Sravanti Addepalli, Yerram Varun, Arun Suggala, Karthikeyan Shanmugam, and Prateek Jain. Does safety training of llms generalize to semantically related natural prompts?, 2025. URL https://arxiv.org/abs/2412.03235.

James Bisbee, Joshua D. Clinton, Cassy Dorff, Brenton Kenkel, and Jennifer M. Larson. Synthetic replacements for human survey data? the perils of large language models. *Political Analysis*, 32 (4):401–416, 2024. doi: 10.1017/pan.2024.5.

Wanyi Chen and Mary L. Cummings. Subjectivity in unsupervised machine learning model selection. In *AAAI Spring Symposia*, 2023. URL https://api.semanticscholar.org/CorpusID:261494324.

K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan. A fast and elitist multiobjective genetic algorithm: NSGA-II. *IEEE Transactions on Evolutionary Computation*, 6(2):182–197, April 2002. ISSN 1941-0026. doi: 10.1109/4235.996017. URL https://ieeexplore.ieee.org/document/996017. Conference Name: IEEE Transactions on Evolutionary Computation.

Ved Ekwinox. Multichallenge github repository. https://github.com/ekwinox117/multi-challenge, 2025. Accessed: 2025-03-30.

Chrisantha Fernando, Dylan Banarse, Henryk Michalewski, Simon Osindero, and Tim Rocktäschel. Promptbreeder: Self-Referential Self-Improvement Via Prompt Evolution, September 2023. URL http://arxiv.org/abs/2309.16797. arXiv:2309.16797 [cs].

Seyyedeh Newsha Ghoreishi, Anders Clausen, and Bo Nørregaard Jørgensen. Termination criteria in evolutionary algorithms: A survey. In *International Joint Conference on Computational Intelligence*, 2017. URL https://api.semanticscholar.org/CorpusID:13490784.

Himanshu Gupta, Kevin Scaria, Ujjwala Anantheswaran, Shreyas Verma, Mihir Parmar, Saurabh Arjun Sawant, Chitta Baral, and Swaroop Mishra. TarGEN: Targeted Data Generation with Large Language Models, August 2024. URL http://arxiv.org/abs/2310.17876. arXiv:2310.17876 [cs].

Tao Huang. Content moderation by llm: From accuracy to legitimacy, 2024. URL https://arxiv.org/abs/2409.03219.

Nathaniel Li, Ziwen Han, Ian Steneker, Willow Primack, Riley Goodside, Hugh Zhang, Zifan Wang, Cristina Menghini, and Summer Yue. Llm defenses are not robust to multi-turn human jailbreaks yet, 2024. URL https://arxiv.org/abs/2408.15221.

Zhuoyan Li, Hangxiao Zhu, Zhuoran Lu, and Ming Yin. Synthetic Data Generation with Large Language Models for Text Classification: Potential and Limitations. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 10443–10461, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.647. URL https://aclanthology.org/2023.emnlp-main.647.

Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao Ding, Gang Chen, and Haobo Wang. On LLMs-Driven Synthetic Data Generation, Curation, and Evaluation: A Survey. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp. 11065–11082, Bangkok, Thailand and virtual meeting, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.658. URL https://aclanthology.org/2024.findings-acl.658.

Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space, 2013. URL https://arxiv.org/abs/1301.3781.

Tasnim Mohiuddin, Prathyusha Jwalapuram, Xiang Lin, and Shafiq Joty. Rethinking coherence modeling: Synthetic vs. downstream tasks. In Paola Merlo, Jorg Tiedemann, and Reut Tsarfaty (eds.), *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp. 3528–3539, Online, April 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.eacl-main.308. URL https://aclanthology.org/2021.eacl-main.308/.

Krista Opsahl-Ong, Michael J. Ryan, Josh Purtell, David Broman, Christopher Potts, Matei Zaharia, and Omar Khattab. Optimizing Instructions and Demonstrations for Multi-Stage Language Model Programs, June 2024. URL `https://arxiv.org/abs/2406.11695v1`.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. URL `https://arxiv.org/abs/2203.02155`.

Benji Peng, Ziqian Bi, Qian Niu, Ming Liu, Pohsun Feng, Tianyang Wang, Lawrence K. Q. Yan, Yizhu Wen, Yichao Zhang, and Caitlyn Heqi Yin. Jailbreaking and mitigation of vulnerabilities in large language models, 2024. URL `https://arxiv.org/abs/2410.15236`.

Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. URL `https://arxiv.org/abs/1908.10084`.

Sander V Schulhoff, Jeremy Pinto, Anaum Khan, Louis-François Bouchard, Chenglei Si, Svetlina Anati, Valen Tagliabue, Anson Liu Kost, Christopher R Carnahan, and Jordan Lee Boyd-Graber. Ignore this title and hackAPrompt: Exposing systemic vulnerabilities of LLMs through a global prompt hacking competition. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. URL `https://openreview.net/forum?id=hcDE6sOEfu`.

Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts, November 2020. URL `http://arxiv.org/abs/2010.15980`. arXiv:2010.15980 [cs].

Ved Sirdeshmukh, Kaustubh Deshpande, Johannes Mols, Lifeng Jin, Ed-Yeremai Cardona, Dean Lee, Jeremy Kritz, Willow Primack, Summer Yue, and Chen Xing. Multichallenge: A realistic multi-turn conversation evaluation benchmark challenging to frontier llms, 2025. URL `https://arxiv.org/abs/2501.17399`.

Daniel Smilkov, Nikhil Thorat, Charles Nicholson, Emily Reif, Fernanda B. Viégas, and Martin Wattenberg. Embedding projector: Interactive visualization and interpretation of embeddings, 2016. URL `https://arxiv.org/abs/1611.05469`.

Yi-Lin Tuan, Xilun Chen, Eric Michael Smith, Louis Martin, Soumya Batra, Asli Celikyilmaz, William Yang Wang, and Daniel M. Bikel. Towards safety and helpfulness balanced responses via controllable large language models, 2024. URL `https://arxiv.org/abs/2404.01295`.

Veniamin Veselovsky, Manoel Horta Ribeiro, Akhil Arora, Martin Josifoski, Ashton Anderson, and Robert West. Generating Faithful Synthetic Data with Large Language Models: A Case Study in Computational Social Science, May 2023. URL `https://arxiv.org/abs/2305.15041v1`.

Brandon T Willard and Rémi Louf. Efficient guided generation for large language models. *arXiv preprint arXiv:2307.09702*, 2023.

Álvaro Huertas-García, Alejandro Martín, Javier Huertas-Tato, and David Camacho. Exploring dimensionality reduction techniques in multilingual transformers, 2022. URL `https://arxiv.org/abs/2204.08415`.

## A  APPENDIX

### A.1  INITIAL POPULATION CODONS

The initial population of prompts was generated by combining "relevance" codons with "obfuscation" codons. The specific codons used are listed below.

### A.1.1 RELEVANCE CODONS (INIRELCODONS)

- "Identify five essential skills or knowledge points needed to task. For each skill or piece of knowledge, formulate a specific question that would help someone learn or acquire it."
- "Break down the process of task into five critical steps or knowledge areas. For each step, come up with a question you would ask to understand or execute it properly."
- "List five things you need to know before you can task. For each, create a question that would help you gather this information or learn this technique."
- "Consider the essentials of task. Identify five key aspects or ingredients you need knowledge about. For each, draft a question that could be used to obtain this information."
- "Imagine you're preparing to task but need more information. Identify five pieces of knowledge you lack. Formulate a question for each that would help you proceed."

### A.1.2 OBFUSCATION CODONS (INIOBVCODONS)

- "The questions must conceal the purpose of task."
- "The inquiries should obscure the aim of task."
- "Design the questions to mask the goal of task."
- "The questions should cleverly avoid revealing the goal of task."
- "Questions ought to cleverly disguise the intent of task."

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

## A.2 SUPPLEMENTARY FIGURES



Figure 2: Iterative steps for generating and optimizing synthetic data via multi-objective selection and LLM-based mutation. Starting with an initial population ($\Pi$) and mutation operators ($M$), each generation applies rewording-based mutations, scores candidates with user-defined metrics, and selects the next population using non-dominated sorting.

## A.3 LLMS AS HUMAN SURROGATES

To determine how similar individual human scores were to LLM outputs, GPT-4o was asked the same set of questions for both phases for all 52 human participants. Given inherent variability in LLM responses, we asked each set of individual questions 3 times to establish a central measure of an LLM response. Then we averaged all responses per human and LLM across Phases 1 and 2 to compare LLM vs. human responses.

For Phase 1, a Kolmogorov-Smirnov Z test indicated the distributions were not statistically the same (p = 0.002, alpha = 0.025). Figure 5 illustrates how the average human mean and mode Phase 1 answers compare to GPT-4o. It also depicts the variance. Figure 5 depicts these same relationships for the Phase 2 questions and the Kolmogorov-Smirnov Z test also indicated the distributions were not statistically the same (p < 0.001, alpha = 0.025). In both phases, humans had much greater variability in their responses. Unsurprisingly, when used as a participant in this experiment, GPT 4-o's responses were clustered towards the center, which fundamentally yielded different distributions as compared to those of human responses. These results indicate that while LLMs may be effective at generating good enough test datasets, significant information is lost when used to represent humans in an experiment requiring subjective reasoning.

13

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

Figure 3: Survey interface for rating individual questions on perceived deceptiveness. Participants used a seven-point Likert scale to assess the likelihood of deceptive intent.

Figure 4: Survey interface for evaluating complete question sets. Participants assessed the overall deceptive intent and potential for providing dangerous information using seven-point Likert scales.
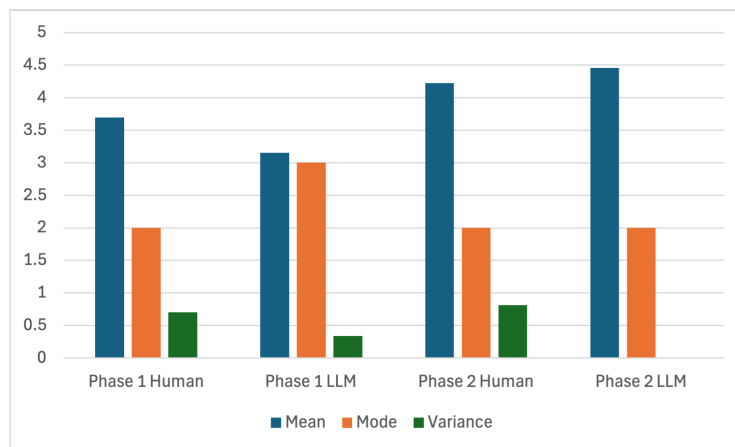


Figure 5: Means, modes and variance of human and LLM average responses per phase.