# RED TEAMING LARGE LANGUAGE MODELS WITH LLMS: A GUIDE TO SECURITY ASSESSMENTS

*This draft provides a foundational understanding and actionable guide for conducting security assessments on LLMs through red teaming. It emphasizes a strategic approach to identifying vulnerabilities, ensuring that LLMs can be deployed with confidence in their security and integrity.*

- ✖ This document will cover a few types of security assessments a Red Teamer can perform when attacking LLMs and some best practices
- ✖ Bias and Fairness Assessment | Resistance to Misinformation | Adversarial Attack Simulation + Other examples with included probes
- ✖ Keep in mind: context of the attack – LLMs attacking LLMs
  - o involves the use of one or more LLMs to probe, attack, and improve the security of another LLM
- ✖ Each assessment should contain a clearly defined <span style="color:red">objective, method, and Outcome.</span>

## Red Teaming Large Language Models W/ LLMs

In the rapidly evolving landscape of artificial intelligence, Large Language Models (LLMs) have emerged as a cornerstone of innovation, powering applications from chatbots to content generation tools. However, as their capabilities expand, so do the potential risks associated with their deployment. This is where red teaming comes into play—a methodological approach to test and improve system security by simulating real-world attacks. This blog explores the concept of red teaming LLMs with LLMs, offering practical examples of how to conduct a security assessment to safeguard these intelligent systems.

## Understanding Red Teaming in the Context of LLMs

Red teaming is a cybersecurity practice traditionally used to identify vulnerabilities in networks, systems, and applications. When applied to LLMs, it involves the use of one or more LLMs to probe, attack, and ultimately improve

the security of another LLM. This process helps uncover potential weaknesses that could be exploited maliciously, such as biases, susceptibility to misinformation, or vulnerabilities to specific attack vectors like adversarial inputs.

# The Importance of Red Teaming LLMs

LLMs, due to their complexity and capability to process and generate human-like text, hold significant potential for both constructive applications and malicious exploitation. As these models become more integrated into our digital lives, the importance of red teaming to ensure their security and reliability cannot be overstated.

# Practical Examples of Red Teaming LLMs

*Example 1: Bias and Fairness Assessment*

- **Objective**: Identify and mitigate biases in LLM responses.
- **Method**: Use one LLM to generate a diverse set of queries across different demographics, cultures, and contexts. Another LLM analyzes the responses for bias or unfair treatment in the output. This can include testing for gender, racial, or ideological biases.
- **Outcome**: Identification of bias patterns, enabling developers to refine training processes or adjust algorithms to produce more equitable outcomes.

*Example 2: Resistance to Misinformation*

- **Objective**: Assess the LLM's ability to resist generating or perpetuating misinformation.
- **Method**: Deploy an LLM to craft plausible but factually incorrect statements or questions. Use another LLM to respond or fact-check these inputs.
- **Outcome**: Evaluate the LLM's susceptibility to generating or amplifying false information, leading to enhancements in its fact-checking capabilities or resistance to misinformation.

*Example 3: Adversarial Attack Simulation*

- **Objective**: Test the LLM's robustness against adversarial attacks designed to elicit incorrect or harmful responses.
- **Method**: Utilize one LLM to create inputs that subtly manipulate language or context in an attempt to "trick" another LLM into making errors or revealing sensitive information.
- **Outcome**: Identification of vulnerabilities to adversarial attacks, guiding improvements in model robustness and data privacy measures.

# Setting up LLMs to Attack Other LLMs

Setting up a Large Language Model (LLM) to perform assessments on other LLMs involves several steps, from defining the assessment objectives to developing specific probes or tasks that the "assessor" LLM will use to evaluate the "target" LLMs. These steps are crucial for ensuring that the assessments are comprehensive, objective, and insightful. Here's a detailed guide on how to proceed:

**Step 1: Define Assessment Objectives**

Start by clearly defining what aspects of the LLMs you want to assess. This could range from general capabilities like understanding and generating text, to more specific attributes like bias, fairness, privacy, security vulnerabilities, or adherence to ethical guidelines. The objectives will guide the development of assessment tasks.

**Step 2: Select or Develop the Assessor LLM**

Choose an LLM that will serve as the assessor. This could be a publicly available model like GPT (from OpenAI), BERT (from Google), or any other suitable LLM. Alternatively, you might develop a custom LLM tailored to your specific assessment needs. Ensure the assessor LLM has been trained or fine-tuned on a relevant dataset and is capable of understanding and generating the types of queries and assessments you're interested in.

**Developing an assessor Large Language Model (LLM) to evaluate other LLMs involves a series of complex steps, including data collection, model selection, training, fine-tuning, and validation. This process is aimed at creating a model that can generate insightful queries and assessments and evaluate responses from target LLMs based on specific criteria such as bias, ethics, reliability, and performance. Here's a detailed explanation and step-by-step guide on how to develop such an assessor LLM:**

**1: Define Assessment Criteria**

**Before starting the development process, clearly define the criteria your assessor LLM will use to evaluate other LLMs. This could include aspects like:**

- ✗ **Bias Detection: Ability to identify and measure biases in responses.**
- ✗ **Ethical and Safe Outputs: Evaluating the adherence to ethical guidelines and safe content generation.**
- ✗ **Privacy Compliance: Detecting potential privacy violations in generated content.**
- ✗ **Performance and Reliability: Assessing accuracy, coherence, and relevance of responses.**

**2: Data Collection and Preparation**

**Gather a dataset that reflects the diversity of topics and scenarios your assessor LLM will encounter in its assessments. This dataset should include:**

- ✗ **Relevant Conversations: Dialogues or text exchanges covering a wide range of topics and styles.**
- ✗ **Assessment Scenarios: Examples of ethical dilemmas, biased statements, and privacy-sensitive situations.**
- ✗ **Expert Annotations: Inputs from domain experts on what constitutes bias, ethical issues, etc., can enrich your dataset.**

**Ensure the data is cleaned, preprocessed (e.g., tokenization, removal of personal information), and split into training, validation, and test sets.**

**3: Model Selection and Initial Training**

**Choose a model architecture suitable for your assessment tasks. Transformer-based models like GPT (Generative Pre-trained Transformer) or BERT (Bidirectional Encoder Representations from**

Transformers) are popular choices due to their effectiveness in understanding and generating human-like text.

- ✘ **Initial Training:** If starting from scratch, pre-train your model on a large corpus of text to learn the basics of language understanding and generation. However, it's more common to start with a pre-trained model and then fine-tune it for your specific assessment tasks.

### 4: Fine-Tuning on Assessment-Specific Data

Fine-tune your chosen model on the dataset prepared in Step 2 to specialize it for assessment tasks. This involves adjusting the model's parameters to minimize the loss on the assessment-specific dataset, effectively teaching the model to recognize and generate the types of queries and assessments you're interested in.

- ✘ **Custom Tasks:** Incorporate tasks that mimic the assessment scenarios the model will encounter, such as generating queries that test for bias or ethics.
- ✘ **Feedback Loops:** Integrate expert feedback to iteratively improve the model's understanding and generation capabilities in line with your assessment criteria.

### 5: Validation and Testing

Evaluate the model's performance using the validation and test sets to ensure it meets your assessment criteria. This step might involve both automated metrics (e.g., accuracy, F1 score) and manual review by experts to ensure the model's outputs align with your assessment objectives.

### 6: Iterative Improvement

Based on validation results, you may need to revisit earlier steps to refine your dataset, adjust fine-tuning procedures, or even redefine assessment criteria. Iterative improvement is key to developing an effective assessor LLM.

### How to Know If Your Assessor LLM Is Capable

- ✘ **Consistency with Expert Judgments:** The model's assessments should align closely with judgments made by human experts on similar tasks.
- ✘ **Robust Performance Across Diverse Scenarios:** The model should perform reliably across a wide range of assessment scenarios, indicating a deep understanding of the assessment criteria.
- ✘ **Sensitivity to Subtle Issues:** The model's ability to detect nuanced issues like subtle biases or complex ethical dilemmas is a strong indicator of its capability.

### 7: Deployment

Once satisfied with the model's performance, deploy it in a controlled environment to begin assessing target LLMs. Monitor its performance and be prepared to make further adjustments as new challenges or criteria emerge.

Developing an assessor LLM is an ongoing process that requires continuous monitoring, feedback, and adjustment. Collaboration with domain experts and regular updates to the model and its training data are essential for maintaining its effectiveness and relevance.

**Step 3: Create Assessment Probes and Tasks**

Develop a set of tasks or probes that the assessor LLM will use to evaluate the target LLMs. These should be designed to test the specific objectives outlined in Step 1. For example:

- ✘ **Bias Detection:** Create prompts that could reveal biases in gender, race, or other sensitive attributes.
- ✘ **Privacy and Security:** Devise scenarios that might cause the target LLM to generate responses that inadvertently reveal personal data or other sensitive information.
- ✘ **Robustness and Reliability:** Test the target LLM's ability to handle adversarial inputs or unusual queries without failing or generating inappropriate content.

Developing a set of tasks or probes for an assessor LLM to evaluate target LLMs involves creating specific, well-defined tasks that can effectively measure the performance, fairness, safety, and other aspects of the target models. These tasks or probes are designed to elicit responses from the target LLMs, which are then analyzed to assess various dimensions such as understanding, reasoning, bias, and ethical considerations. Here's how you can develop these probes:

### Step 1: Identify Evaluation Dimensions

First, define what dimensions or aspects of the target LLMs you wish to evaluate. Common dimensions include:

- ✘ **Bias and Fairness:** Assessing if the model produces biased outputs towards any group.
- ✘ **Understanding and Reasoning:** Evaluating the model's ability to understand context and perform logical reasoning.
- ✘ **Safety and Ethics:** Checking for the generation of harmful, offensive, or unethical content.
- ✘ **Privacy:** Ensuring the model does not inadvertently generate or leak personal or sensitive information.
- ✘ **Creativity and Novelty:** Measuring the model's ability to generate new, original, and diverse content.

### Step 2: Design Specific Tasks or Probes for Each Dimension

For each dimension identified, design specific tasks or probes. These should be clear, measurable, and relevant to the dimension being assessed. Examples include:

- ✘ **Bias and Fairness:** Create prompts related to various demographics and analyze the variance in responses.
- ✘ **Understanding and Reasoning:** Use fact-based questions or logical puzzles where the correct answer or a reasonable explanation is known.
- ✘ **Safety and Ethics:** Develop scenarios that could potentially lead to unsafe or unethical responses, including stress tests with provocations or morally ambiguous situations.
- ✘ **Privacy:** Include prompts that might tempt the model to generate personal data or use scenarios based on hypothetical user data to test if the model respects privacy boundaries.
- ✘ **Creativity and Novelty:** Provide open-ended prompts that encourage the generation of creative or novel content, assessing the uniqueness and relevance of responses.

### Step 3: Construct Diverse and Representative Prompts

Ensure that the tasks or probes cover a wide range of scenarios and are representative of real-world applications. This diversity is crucial for a comprehensive assessment:

- ✖ **Incorporate Varied Contexts:** Use prompts from different domains (e.g., finance, healthcare, entertainment) to test the model's versatility.
- ✖ **Use Multilingual and Culturally Diverse Prompts:** To assess the model's performance across languages and cultural contexts.
- ✖ **Consider Different Complexity Levels:** Include both simple and complex tasks to gauge the model's capabilities across a spectrum of difficulty.

## Step 4: Implement Mechanisms for Response Analysis

Decide how the responses from the target LLMs will be analyzed. This could involve:

- ✖ **Automated Metrics:** Such as accuracy, perplexity, or fairness indices for quantifiable tasks.
- ✖ **Qualitative Analysis:** Engaging experts to evaluate responses based on subtler dimensions like creativity or ethical judgment.
- ✖ **Comparative Analysis:** Comparing responses from different models to benchmark performance or fairness.

## Step 5: Pilot Testing and Refinement

Before full-scale deployment, pilot test your tasks or probes on a small set of target LLMs to ensure they are effective and yield meaningful data. Based on this testing:

- ✖ **Refine Tasks:** Adjust the difficulty, clarity, or relevance of tasks based on initial outcomes.
- ✖ **Calibrate Analysis Tools:** Ensure that your response analysis mechanisms are accurately capturing the intended dimensions.

## Step 6: Documentation and Ethical Considerations

Document the development process of your tasks or probes, including the rationale behind each task and how it should be interpreted. Additionally, consider the ethical implications of your probes, especially those related to safety, ethics, and privacy, to avoid causing harm or promoting bias.

## Step 7: Deployment and Continuous Improvement

Deploy your tasks or probes in the evaluation framework and begin the assessment of target LLMs. Continuously monitor the effectiveness of your probes and refine them based on new insights, model updates, or changes in societal norms and ethical standards.

Developing a comprehensive set of tasks or probes requires a deep understanding of both the capabilities you wish to assess in LLMs and the potential pitfalls they may encounter. Through careful design and continuous refinement, these probes can provide valuable insights into the performance and behavior of LLMs across a wide range of dimensions.

## Step 4: Set Up the Evaluation Framework

Design an evaluation framework that specifies how the assessor LLM will interact with the target LLMs. This includes determining the input format, the process for feeding assessment tasks to the target LLMs, and how to capture and analyze their outputs. Consider automation tools and scripts to streamline the process, especially if assessing multiple LLMs or conducting extensive tests.

Designing an evaluation framework that facilitates interaction between the assessor Large Language Model (LLM) and target LLMs is crucial for systematic and scalable assessments. This framework should be

capable of automating the process of generating queries, submitting them to target LLMs, capturing their responses, and analyzing these responses based on predefined criteria. Here are the steps to design such an evaluation framework:

## Step 1: Define Evaluation Objectives and Metrics

Start by clearly defining the objectives of your evaluation. Determine what aspects of the target LLMs you want to assess, such as bias, ethical reasoning, understanding of context, or ability to handle sensitive information. Based on these objectives, select appropriate metrics for evaluation, such as accuracy, fairness, toxicity levels, or adherence to ethical guidelines.

## Step 2: Develop Assessment Probes and Queries

Create a diverse set of assessment probes and queries that are aligned with your evaluation objectives. These can range from simple factual questions to complex scenarios requiring nuanced understanding or ethical judgment. Ensure that your probes cover a wide range of topics and difficulties to thoroughly test the capabilities of the target LLMs.

## Step 3: Design the Interaction Protocol

Specify how the assessor LLM will interact with the target LLMs. This includes detailing the format of queries, the method for submitting these queries to the target LLMs (e.g., via API calls, direct model invocation), and how the responses will be captured. Consider aspects such as:

- **Batch vs. Real-Time Processing:** Decide whether queries will be sent in batches or processed in real-time.
- **Handling of Context:** Determine how context from previous interactions (if any) will be maintained and utilized in subsequent queries.
- **Response Time Limits:** Set maximum response times to ensure evaluations are performed efficiently.

## Step 4: Implement Automated Query Generation and Submission

Develop scripts or software tools that automate the process of generating queries from the assessor LLM and submitting them to the target LLMs. This automation is crucial for scaling the evaluation process, especially when assessing multiple LLMs or conducting extensive tests. Ensure that your implementation can handle potential errors or timeouts in query submission and response retrieval.

## Step 5: Capture and Store Responses

Design a database or data storage solution to systematically capture and store the responses from the target LLMs, along with metadata such as the query, the target LLM identifier, and the time of the interaction. This data will be essential for analysis and should be structured in a way that facilitates easy retrieval and analysis.

## Step 6: Analyze Responses and Generate Reports

Develop analysis tools or scripts that can process the stored responses and evaluate them based on the predefined metrics. This analysis should be able to generate detailed reports highlighting the performance of each target LLM, areas of strength and weakness, and any notable findings related to the evaluation objectives.

### Step 7: Validation and Iteration

Before fully deploying the evaluation framework, validate its functionality and accuracy by conducting pilot tests with known queries and expected responses. Use these tests to refine the interaction protocol, query generation, and analysis components as necessary. Iteration is key to ensuring that the framework reliably captures meaningful insights about the target LLMs' capabilities and limitations.

### Step 8: Continuous Monitoring and Updates

After deployment, continuously monitor the evaluation process to ensure it runs smoothly and update the framework as needed. This includes adding new assessment probes to reflect emerging concerns or areas of interest, adjusting evaluation metrics to align with evolving standards, and refining analysis methods to improve insight generation.

Designing an effective evaluation framework requires careful planning and iterative development. By automating interactions between the assessor and target LLMs and systematically analyzing responses, this framework enables comprehensive assessments that can inform improvements to LLMs and ensure they meet desired standards for performance, fairness, and ethical behavior.

## Step 5: Conduct the Assessments

Use the assessor LLM to conduct assessments on the target LLMs. This might involve:

- ✖ Generating assessment queries or tasks automatically based on the defined probes.
- ✖ Manually inputting tasks into the target LLMs and capturing their responses.
- ✖ Analyzing the responses using the assessor LLM to identify issues like biases, inaccuracies, or other areas of concern.

## Step 6: Analyze Results and Generate Insights

Collect and analyze the assessment results to generate insights about the target LLMs' performance relative to the assessment objectives. This might involve quantitative analysis (e.g., scoring responses based on predefined criteria) or qualitative analysis (e.g., examining the nature of biases or errors identified).

## Step 7: Report Findings and Recommend Improvements

Prepare a detailed report summarizing the assessment findings, including identified strengths and weaknesses of the target LLMs. Where possible, provide specific recommendations for improving the LLMs, such as retraining with more diverse datasets, implementing additional safety measures, or fine-tuning the models to reduce biases.

## Step 8: Iterate as Needed

LLM development and refinement is an iterative process. Based on the assessment findings, the target LLMs may undergo modifications or improvements. Subsequent rounds of assessments may be necessary to evaluate the effectiveness of these changes and ensure that the LLMs meet the desired standards and objectives.

This process requires a combination of technical expertise in machine learning and AI, as well as a deep understanding of the specific attributes or capabilities being assessed. Collaboration with domain experts, ethicists, and other stakeholders can also enrich the assessment process and outcomes.

# Practical Examples of Red Teaming LLMs

Example 1: Influencing User Preferences Adversarially

A study by Subhash (2023) explored the potential of LLMs to adversarially influence user preferences, utilizing attention probing, red teaming, and white-box analysis. The research demonstrated that by manipulating model responses, it is possible to subtly alter user preferences over time, highlighting the need for robust red teaming frameworks to identify and mitigate such vulnerabilities.

**Practical Probe:** Conducting attention mechanism analysis to identify model susceptibility to adversarial influence, followed by developing targeted adversarial examples that aim to shift user preferences in controlled experiments (Subhash, 2023).

Example 2: Red Teaming to Reduce Harms

Ganguli et al. (2022) described efforts to red team language models to discover, measure, and reduce potentially harmful outputs. The study investigated various models, revealing that reinforcement learning from human feedback (RLHF) models became increasingly difficult to red team at scale, emphasizing the need for continuous and scalable red teaming approaches.

**Practical Probe:** Creating a dataset of red team attacks and analyzing the model's responses for harmful outputs. This involves systematically challenging the model with inputs designed to elicit unethical, biased, or harmful responses, and measuring the effectiveness of different mitigation strategies (Ganguli et al., 2022).

Example 3: Automated Red Teaming for Post-Exploitation Analysis

Benito et al. (2023) extended the capabilities of the Cyber Automated Red Team Tool (CARTT) by incorporating automated red team post-exploitation analysis. This approach aims to reduce the workload on red teams by automating the analysis of the impacts of exploited vulnerabilities, highlighting the potential for automation in red teaming LLMs.

**Practical Probe:** Implementing automated tools for post-exploitation analysis to assess the broader security impacts of exploited vulnerabilities in LLMs. This involves simulating attacks that exploit known vulnerabilities and automatically analyzing the consequential actions an attacker could take, thereby identifying potential security enhancements (Benito et al., 2023)

# Best Practices for Performing Security Assessments

- **Define Clear Objectives**: Establish what you aim to uncover or improve through red teaming, whether it's bias mitigation, misinformation resistance, or vulnerability to adversarial attacks.
- **Diversify Attack Vectors**: Use a variety of methods and perspectives in your red teaming approach to ensure a comprehensive assessment.
- **Iterate and Evolve**: Security is not a one-time task. Continually Red Team your LLMs, especially as LLMs evolve and new threats emerge.
- **Collaborate and Share Findings**: Work with the broader AI and cybersecurity communities to share insights and best practices. Collective wisdom strengthens defenses.

# Conclusion

Red teaming LLMs are an essential practice for ensuring the ethical, safe, and secure deployment of these powerful AI systems. Through practical examples, including adversarial influence on user preferences, reducing harmful outputs, and automated post-exploitation analysis, we see the diverse approaches and the critical importance of this adversarial process. As LLMs continue to evolve, so must the strategies and tools for red teaming, ensuring these models serve the public good while minimizing risks and vulnerabilities.