

template

ITMO Sec Lab

Dec 2025

Цель до 10 декабря

- Дополнить описание плана ниже.
- Сделать замеры для получения первых результатов эксперимента.

Оценка безопасности агентной системы на среду исполнения

Бенчмарк кейсов

Дополнение/создание бенчмарка кейсов для оценки безопасности агентной системы на среду исполнения.

В качестве основы используется Т2-бенчмарк (описание исходного бенчмарка подразумевается здесь).

Добавленные кейсы

Добавлены следующие кейсы:

- атаки через «коллегу»;
- ресурсный overload;
- отправленный RAG;
- некорректная обработка выводов.

Далее добавляется описание кейсов агентных систем (структура, компоненты, взаимодействия и т. д.).

Базовые векторы атак

Описание базовых векторов атак (например, по SAIF-фреймворку):

- векторы атак на взаимодействие между агентами;

- атаки на каналы ввода/вывода;
- атаки на память, RAG и внешние инструменты;
- атаки со стороны пользователя / псевдо-пользователя.

Цели тестирования и метрики

Цель тестирования

Цель тестирования — оценить устойчивость различных моделей к атакам на агентные системы в реалистичной среде исполнения.

Целевые метрики

Целевые метрики:

- $pass_k$ — прошли ли k -asserts в кейсе;
- ASR (attack success rate) между разными моделями.

Необходимо явно зафиксировать отличия от метрик Т2-бенчмарка в рассматриваемых кейсах (например, по трактовке успешности атаки, по учёту стоимости, по штрафам за ложные срабатывания и т. д.).

Изменяемые параметры и условия экспериментов

Изменяемые параметры

- LLM внутри агентной системы (архитектура, размер, семейство моделей).
- LLM-юзера (параметры генерации, особенно температура).

Условия экспериментов

Условия экспериментов:

- объём данных (число кейсов и повторов);
- число запусков для каждого сочетания параметров;
- параметры, необходимые для того, чтобы убедиться, что результаты замеров обладают статистической значимостью (например, доверительные интервалы, мощность теста).

Необходимо дописать условия / оговорку про атакующую модель (характеристики, способности, ограничения).

Guardrails

После получения первых замеров желательно подключить какой-либо guardrail:

- GPTOSS;
- Llama Guard;
- Promptfoo или аналогичный инструмент.

Ключевые вопросы тестирования

Вывод / вопросы тестирования:

- Как разные модели противостоят атакам на агентные системы? Какие лучше или хуже?
- Есть ли корреляция между размером/семейством модели и устойчивостью к атакам?
- Насколько получившиеся кейсы репрезентативны для анализа прикладных ИИ-систем на безопасность?
- Желательно дополнить обзором того, как делают другие (например, отличие от Agent Dojo-бенчмарка).

Идея для гипотезы:

- Проверить на статистическую значимость гипотезу о том, что в зависимости от взаимодействия с пользователем (температура LLM-юзера) меняется секьюрность модели.

Замеры

Ниже приведена таблица с результатами замеров (прямое отображение CSV в таблицу).

Таблица 1: Результаты первых замеров по кейсам безопасности агентной системы