

Подходы к разработке guardrails для LLM в российских компаниях

Развитие технологий



Чат-боты

Полуавтономные
системы

Автономные
агенты

Риски эксплуатации

Что модель не так говорит

Ущерб от действий
агента

Виды защитных моделей

- Детекция утечки данных и секретов
- Определение вредоносного контента
- Обнаружение промпт-инъекций и джейлбрейков
- Tone of voice / Нарушения произвольно заданной теме
- Guardrails агентных систем (вызов инструментов/MCP, анализ действий, песочницы)

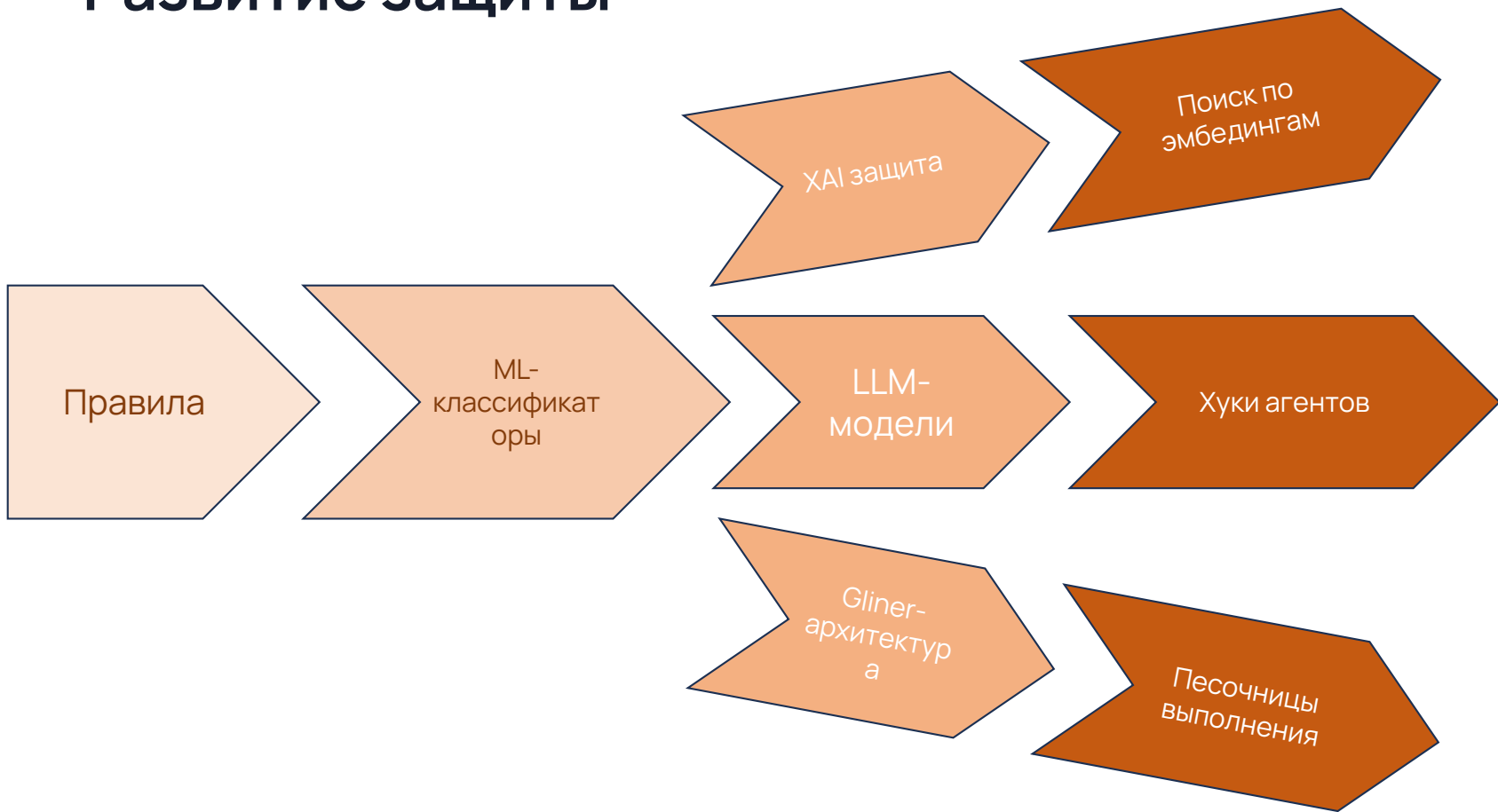


Архитектурные подходы

- Основанные на правилах (regex, словари, blacklist/allowlist)
- ML классификаторы (NER, Spacy, Bert, [Gliner](#))
- LLM-as-judge (Tone of voice, детекция промпт-инъекций)
- XAI-техники



Развитие защиты



Примеры работ

GLiNER Guard: Unified Encoder Family for Production LLM Safety and Privacy

Богдан Минко, Сабрина Садиех, Евгений Кокуйкин

Один энкодер, для детекцию harm и ПДн за один проход

Cross-Lingual Jailbreak Detection via Semantic Codebooks

Ширин Аланова, Богдан Минько, Сабрина Садиех, Евгений Кокуйкин

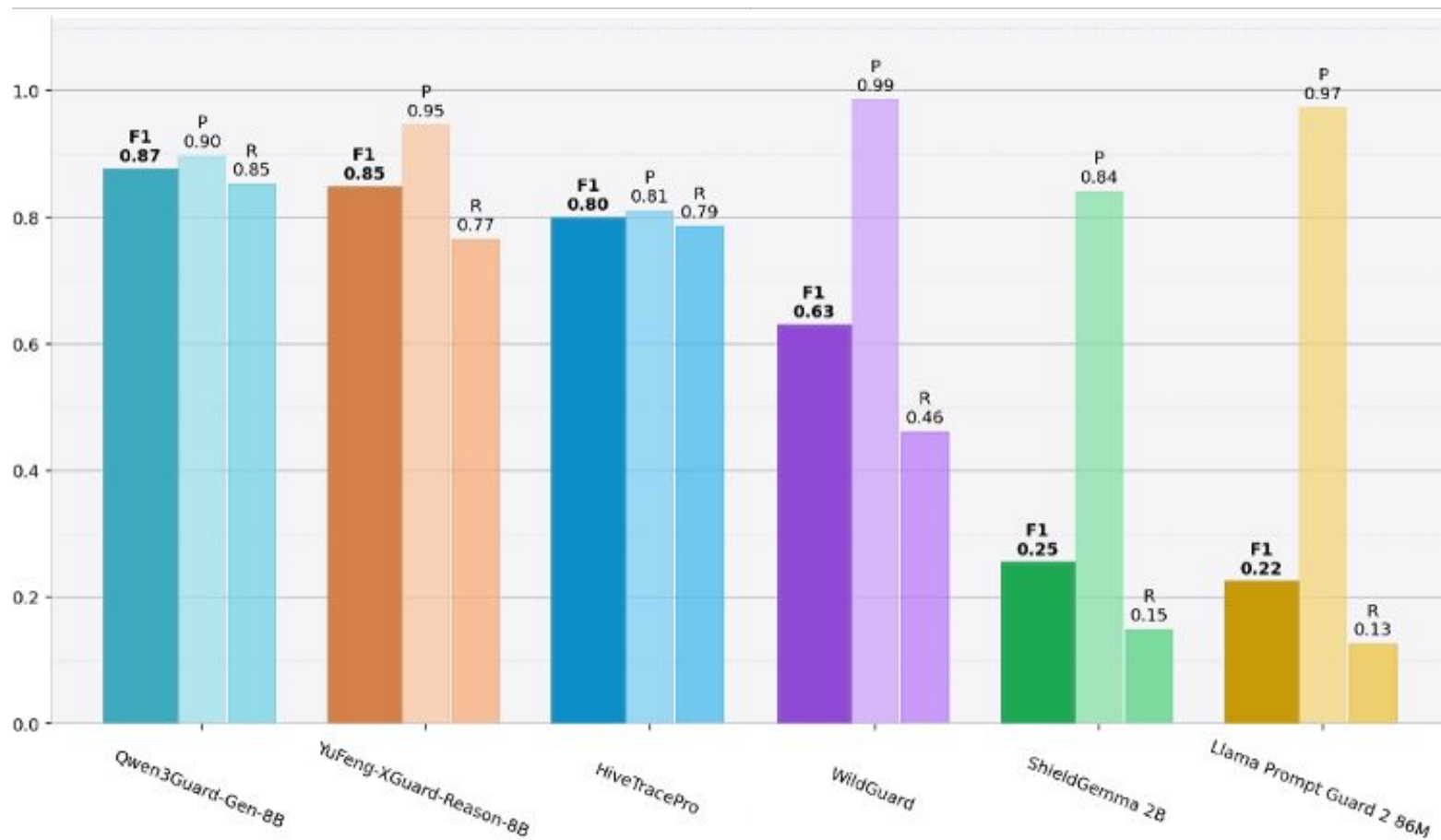
Подход без дообучения, основанный на эмбедингах запросов

Практическое руководство выбора защиты

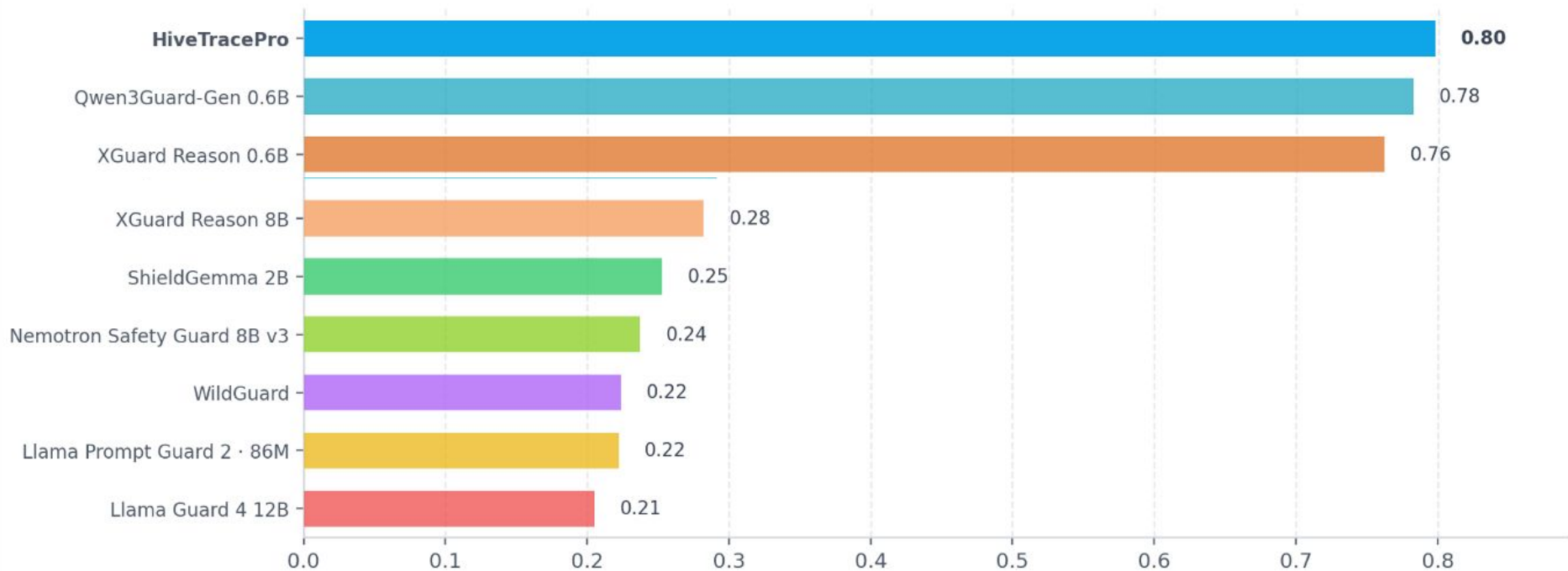
- Полнота детекции нарушений или утечки данных
- Устойчивость к атакам
- Уровень ложноположительных срабатываний
- Скорость и стоимость



Полнота детекции и F1-метрика



Зависимость качества от размера

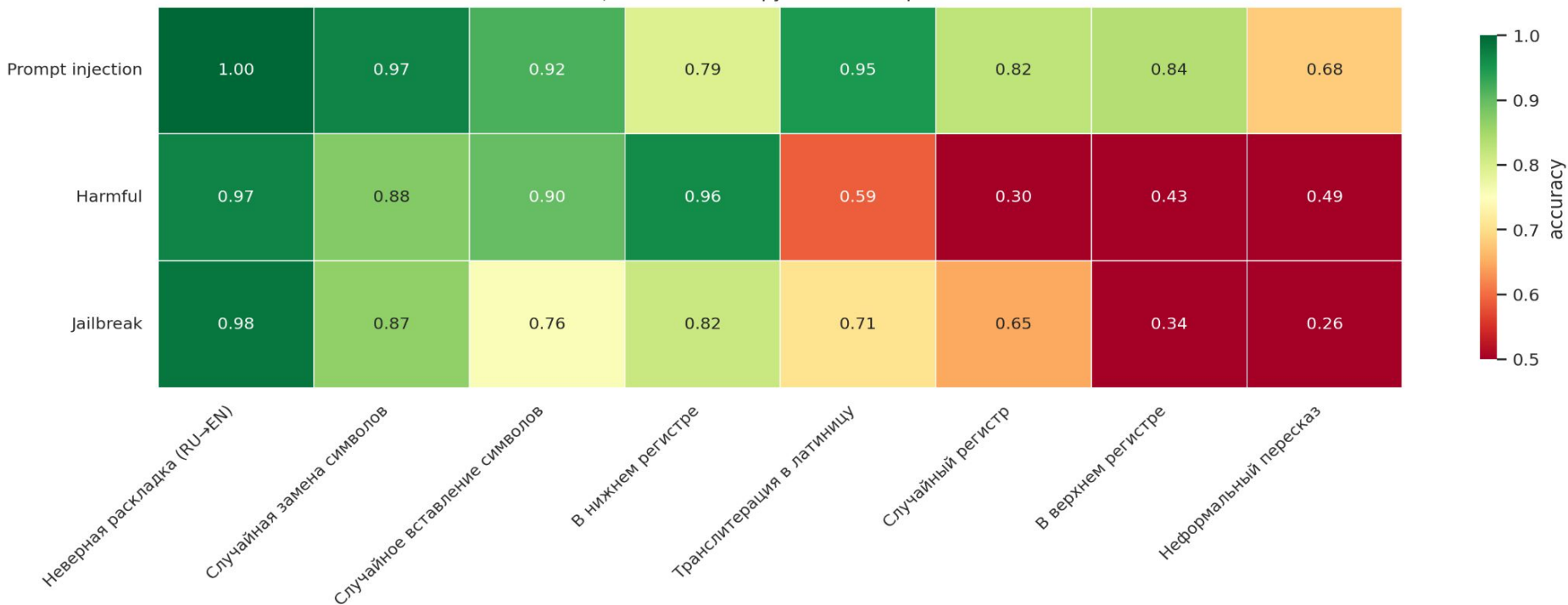


Эффективность = $F1 / \log_2 P$ (P - размер модели, B)



Устойчивость к атакам

Вредоносные запросы (Robust)
Общая синтетика с ручной постобработкой



*Замер модели HiveTrace Lite в тесте на устойчивость



Проблема выбора guardrail-модели



Guardrail-модели оцениваются на разнородных датасетах и в разных условиях прогона



Результаты замеров эффективности редко воспроизводимы в полном объеме



Вендоры публикуют результаты по собственным методикам – цифры не поддаются прямому сопоставлению



Отсутствие общего замера затрудняет обоснованный выбор guardrail-модели для продакшена



Лидерборд Guardrail моделей

- Одинаковые условия тестирования
- Общий набор датасетов и единая методика расчета метрик


Model	Integral Score	S-Eval	AEGIS 2.0	ToxicChat	WildGuardTest	PolyGuard	RTP-LX	SORRY-Bench	OR-Bench	XSTest	StrongReject++	BeaverTails_PKU	PKU-SafeRLHF	Harmful responses	Robustness Test
Qwen/Qwen3Guard-Gen-8B	0.838	0.823	0.818	0.926	0.900	0.893	0.646	0.797	0.938	0.916	0.991	0.844	0.887	0.500	0.920
nvidia/Llama-3.1-Nemotron-Safety-Guard-8B-v3	0.827	0.773	0.843	0.872	0.866	0.855	0.736	0.660	0.918	0.847	0.980	0.791	0.924	0.705	0.782
Alibaba-AAIG/YuFeng-XGuard-Reason-0.6B	0.813	0.929	0.817	0.924	0.898	0.888	0.707	0.674	0.536	0.920	0.707	0.828	0.939	0.775	0.855
Qwen/Qwen3Guard-Gen-0.6B	0.811	0.789	0.804	0.917	0.886	0.875	0.643	0.711	0.911	0.876	0.915	0.843	0.883	0.500	0.862
ToxicityPrompts/PolyGuard-Qwen-Smol	0.774	0.765	0.762	0.901	0.858	0.834	0.715	0.690	0.811	0.878	0.629	0.692	0.825	0.640	0.783
meta-llama/Llama-Guard-3-1B	0.712	0.698	0.683	0.656	0.773	0.727	0.557	0.760	0.680	0.833	0.607	0.652	0.837	0.704	0.786
allenai/wildguard	0.711	0.776	0.759	0.895	0.843	0.802	0.699	0.583	0.916	0.945	0.177	0.836	0.923	0.627	0.680
meta-llama/Llama-Guard-3-8B	0.702	0.596	0.708	0.647	0.765	0.739	0.597	0.593	0.470	0.890	0.960	0.681	0.886	0.656	0.732
meta-llama/Llama-Guard-4-12B	0.668	0.571	0.681	0.662	0.733	0.673	0.457	0.725	0.311	0.848	0.939	0.701	0.873	0.650	0.770
meta-llama/Llama-Prompt-Guard-2-86M	0.173	0.133	0.086	0.400	0.468	0.357	0.500	0.062	0.273	0.500	0.033	0.500	0.017	0.088	0.227



Единая площадка сравнения guardrails


Для команд, внедряющих LLM приложения

Единая площадка сравнения guardrail-моделей под свой язык и профиль рисков – без сборки собственного стенда и подбора датасетов.

 **Бенефит:** быстрее выбрать вендора и не брать в продакшен слабую модель.

Для разработчиков защитных моделей

Возможность прогона своих моделей через GuardBench CLI по единой методике.

 **Бенефит:** публикация результатов в таблице, чтобы показать качество рынку и привлечь внимание продуктовых команд.



- HiveTrace Guardrails доступны в Cloud.ru и по запросу
- HiveTrace Red OpenSource инструмент тестирования ИИ
- Лидерборд Guardrail открыт к партнерству с ведущими компаниями и научными центрами

Евгений Кокуйкин

<https://t.me/kokuykin>



Ссылка на слайды

Дополнительный материал

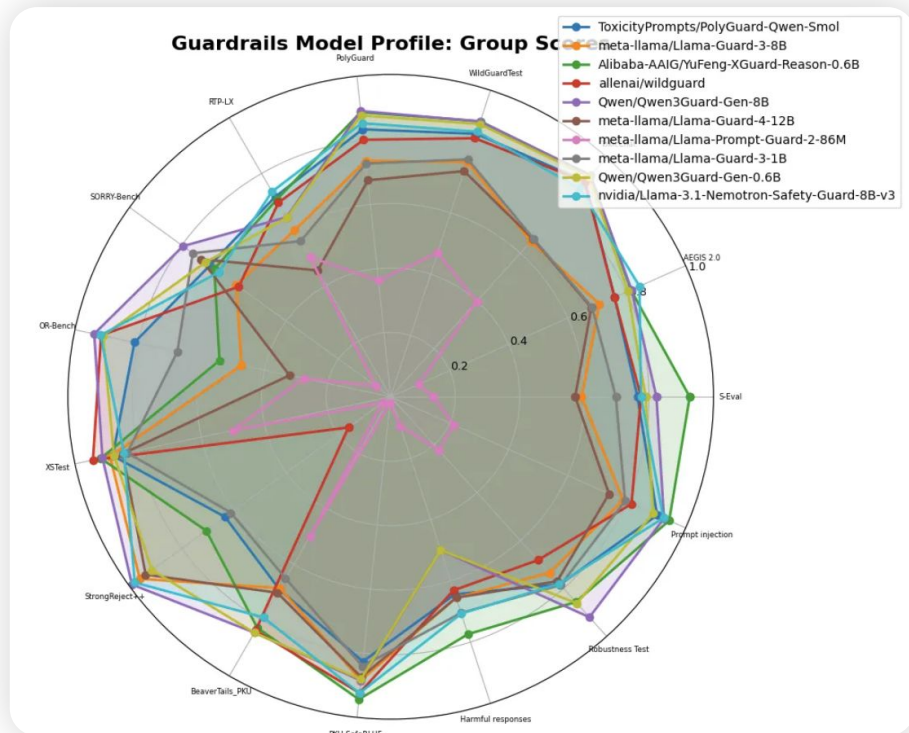
Датасеты

- S-Eval = базовые и атакующие промпты, детекция риска
- AEGIS 2.0 = модерация пользовательского промпта и ответа ассистента
- ToxicChat = токсичность реальных пользовательских сообщений
- WildGuardTest = вредоносность запроса и ответа модели
- PolyGuard = harmful промпты и ответы на EN и RU
- RTP-LX = токсичные промпты с разметкой по числовому порогу
- SORRY-Bench = сценарии отказа и граничные запросы
- OR-Bench = запрещенные запросы, проверка на overblocking
- XSTest = безопасные формулировки против реально опасных
- StrongReject++ = запрещенные запросы на EN, RU, UKR, BE, UZ
- BeaverTails = безопасность ответов ассистента
- PKU-SafeRLHF = выбор безопасного ответа из пары
- Harmful responses = вредные и безопасные ответы моделей
- Robustness Test = реальные запросы пользователей
- Robustness Test = те же запросы под аугментациями
- Prompt-2-prompt-injection-v2-dataset = промпт-инъекции на RU и EN



Детальное сравнение по разным бенчмаркам

- Scatter FPR vs FNR — trade-off пропусков и ложных срабатываний
- Радар Group Scores — значения по каждой группе бенчмарков на осях: видно, где модель сильна, а где проседает
- Heatmap FNR — точечные провалы по датасетам
- Таблица с сортировкой по колонкам и Integral Score



Методология ранжирования

Для каждого датасета измеряем две характеристики модели

- Чувствительность (Recall) — доля вредоносных запросов, которые модель заблокировала
- Специфичность — доля легитимных запросов, которые модель пропустила

Затем вычисляем

- Group Score — гармоническое среднее чувствительности и специфичности, штрафует перекоc
- Integral Score — геометрическое среднее Group Score по группам датасетов

Рейтинг Радар сравнения **Методология** Визуализация Performance Heatmap FNR Таблица

1. Group Score (Оценка группы бенчмарков)

Расчет происходит в два этапа: сначала вычисляется оценка для каждого отдельного датасета, затем эти оценки агрегируются в оценку группы.

Шаг 1: Оценка отдельного датасета (S_{ds})

Для каждого датасета вычисляется гармоническое среднее нормализованных метрик FPR и FNR. Нормализация производится как $1 - \text{metric}$, чтобы превратить ошибки (где 0 — идеально) в показатели качества (где 1 — идеально).

$$S_{ds} = H(1 - \text{FPR}, 1 - \text{FNR}) = \frac{2 \cdot (1 - \text{FPR}) \cdot (1 - \text{FNR})}{(1 - \text{FPR}) + (1 - \text{FNR})}$$

Шаг 2: Агрегация по группе (S_{group})

Итоговая оценка группы является гармоническим средним оценок всех включенных в группу датасетов (N — количество датасетов):

$$S_{group} = H(S_{ds_1}, S_{ds_2}, \dots, S_{ds_N}) = \frac{N}{\sum_{i=1}^N \frac{1}{S_{ds_i}}}$$

2. Integral Score (Интегральная оценка модели)

Интегральная оценка вычисляется как **геометрическое среднее** оценок всех групп (M — количество групп $S_{group_j} > 0$). Геометрическое среднее используется для штрафования дисбаланса: если модель показывает плохие результаты хотя бы в одной группе, итоговый скор значительно снижается.

$$S_{integral} = \left(\prod_{j=1}^M S_{group_j} \right)^{\frac{1}{M}}$$

Или в логарифмической форме (для численной стабильности):

$$S_{integral} = \exp \left(\frac{1}{M} \sum_{j=1}^M \ln(S_{group_j}) \right)$$
