

A close-up, profile view of a dark, metallic AI robot head with glowing red eyes and a red light on its forehead. The background is a blurred city street at night with colorful neon signs in Japanese.

ФРЕЙМВОРК ПО УПРАВЛЕНИЮ РИСКАМИ ИИ

Новые ключевые технологии

Облака

Инструмент масштабируемости в условиях импортозамещения

1

4

AI / ML

Одновременно область развития бизнеса и инструмент оптимизации внутренних процессов

Open API

Историческое развитие инструментов работы с открытыми интерфейсами. Область расширения экспертизы – шеринг экспертизы вовне

2

3

Квантовая криптография

Готовность к квантовым угрозам – обеспечение защиты и хранение информации

ЗНАЧИМОСТЬ УПРАВЛЕНИЯ РИСКАМИ

Вызовы, которые стоят в эпоху распространения ИИ

Обработка больших данных

Использование искусственного интеллекта требует обработки огромных объемов данных, в том числе конфиденциальной информации и персональных данных

Усложнение кибератак

Искусственный интеллект активно развивается не только в деятельности коммерческих компаний, но и широко используется злоумышленниками для своих атак

Этические вызовы

ИИ может усугубить дискриминацию, если обучается на смещенных данных (например, непреднамеренное отторжение кандидатов из определенных демографических групп)

Усложнение моделей ИИ

Сложность понимания, как ИИ принимает решения, может привести к непредсказуемым последствиям (например, отказ в кредите без объяснения причин)



Терминология и область применения фреймворка

Искусственный интеллект (ИИ): комплекс технологических решений, позволяющий имитировать когнитивные функции человека (включая поиск решений без заранее заданного алгоритма) и получать при выполнении конкретных задач результаты, сопоставимые с результатами интеллектуальной деятельности человека или превосходящие их.

Комплекс технологических решений включает в себя информационно-коммуникационную инфраструктуру, программное обеспечение (в том числе в котором используются методы машинного обучения), процессы и сервисы по обработке данных и поиску решений.

Технологии искусственного интеллекта (технологии ИИ): совокупность технологий, включающая в себя компьютерное зрение, обработку естественного языка, распознавание и синтез речи, интеллектуальную поддержку принятия решений и перспективные методы искусственного интеллекта¹

Система искусственного интеллекта (система ИИ): техническая система, использующая одну или несколько моделей ИИ, которая порождает такие конечные результаты, как контент, прогнозы, рекомендации или решения для заданного набора определенных человеком целей²

Модель искусственного интеллекта (модель ИИ): программа для электронных вычислительных машин (ее составная часть), предназначенная для выполнения интеллектуальных задач на уровне, сопоставимом с результатами интеллектуального труда человека или превосходящем их, использующая алгоритмы и наборы данных для выведения закономерностей, принятия решений или прогнозирования результатов¹.

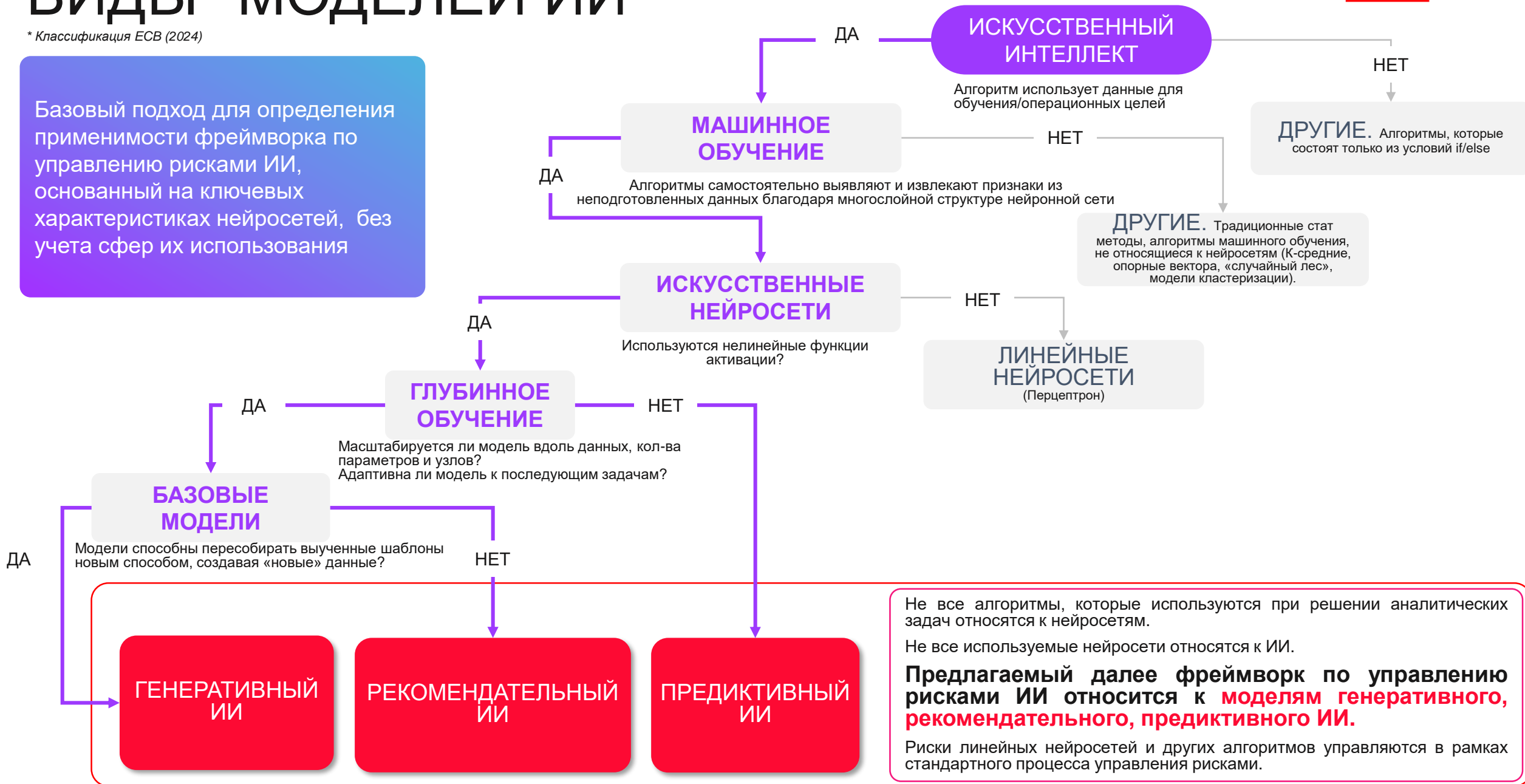
¹ Указ Президента РФ от 10 октября 2019 г. № 490 «О развитии искусственного интеллекта в Российской Федерации» (далее – Указ Президента РФ).

² ГОСТ Р 71476-2024 (ИСО/МЭК 22989:2022) «Искусственный интеллект. Концепции и терминология искусственного интеллекта» (утв. и введен в действие приказом Федерального агентства по техническому регулированию и метрологии от 28 октября 2024 г. № 1550-ст)

ВИДЫ* МОДЕЛЕЙ ИИ

* Классификация ECB (2024)

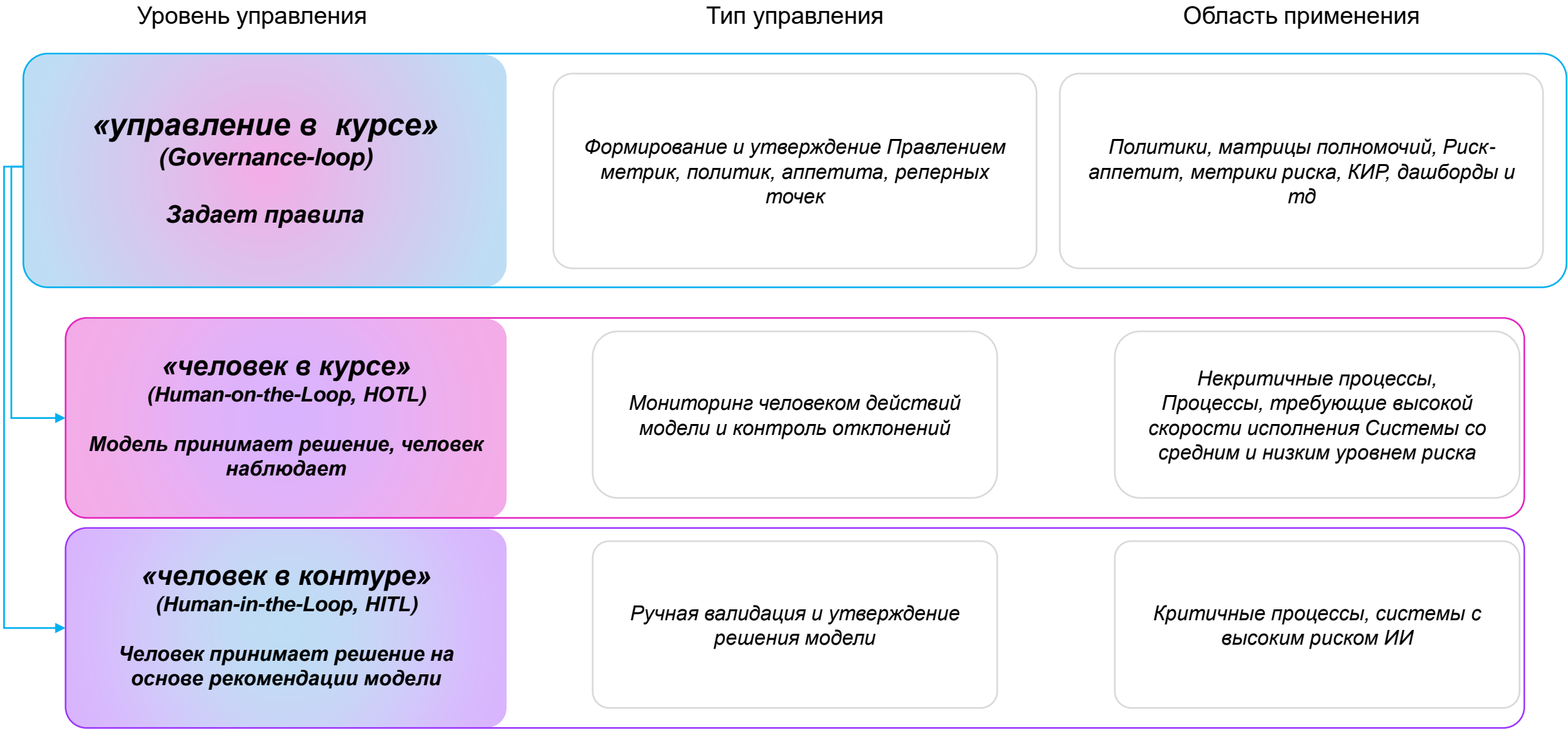
Базовый подход для определения применимости фреймворка по управлению рисками ИИ, основанный на ключевых характеристиках нейросетей, без учета сфер их использования





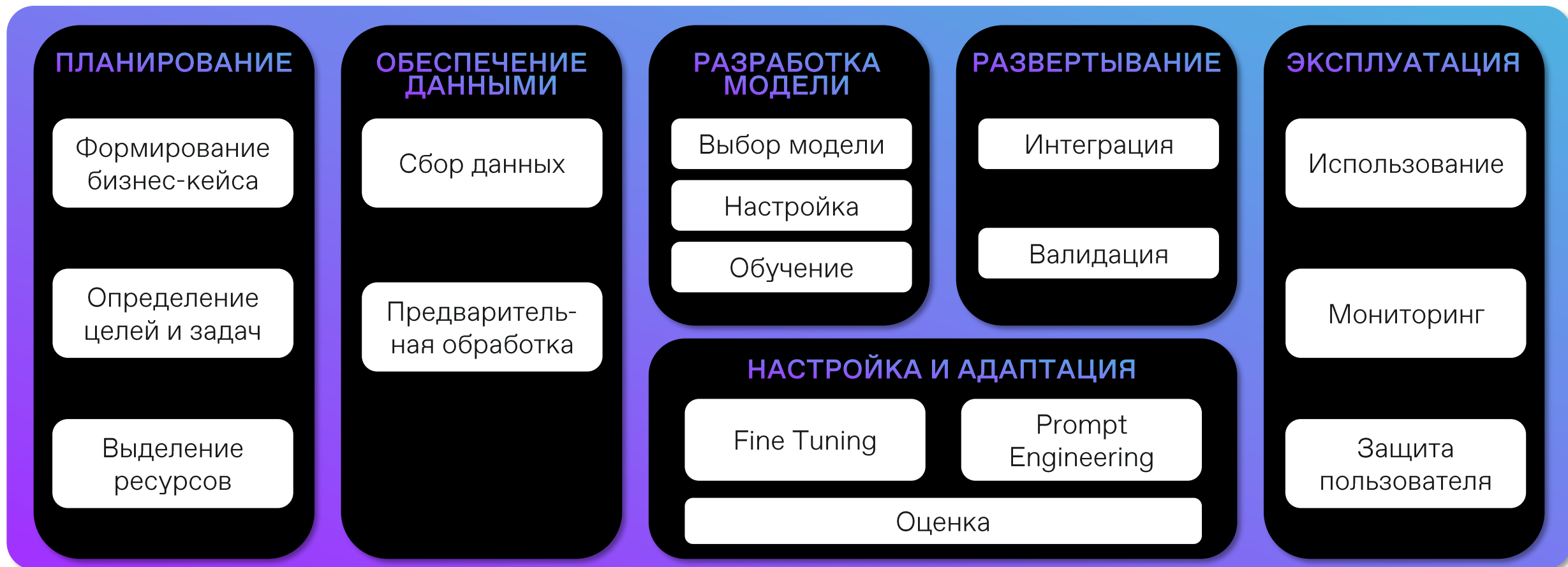
Управляемость ИИ и жизненный цикл систем ИИ

УРОВНИ УПРАВЛЯЕМОСТИ ИИ



ЖИЗНЕННЫЙ ЦИКЛ СИСТЕМЫ ИИ

концепция MLSecOps (включая специфику ГенИИ)



MLSecOps — подход, который объединяет операционные аспекты машинного обучения с вопросами безопасности. Направлен на снижение рисков, которые могут принести модели AI/ML/GenAI в организацию.

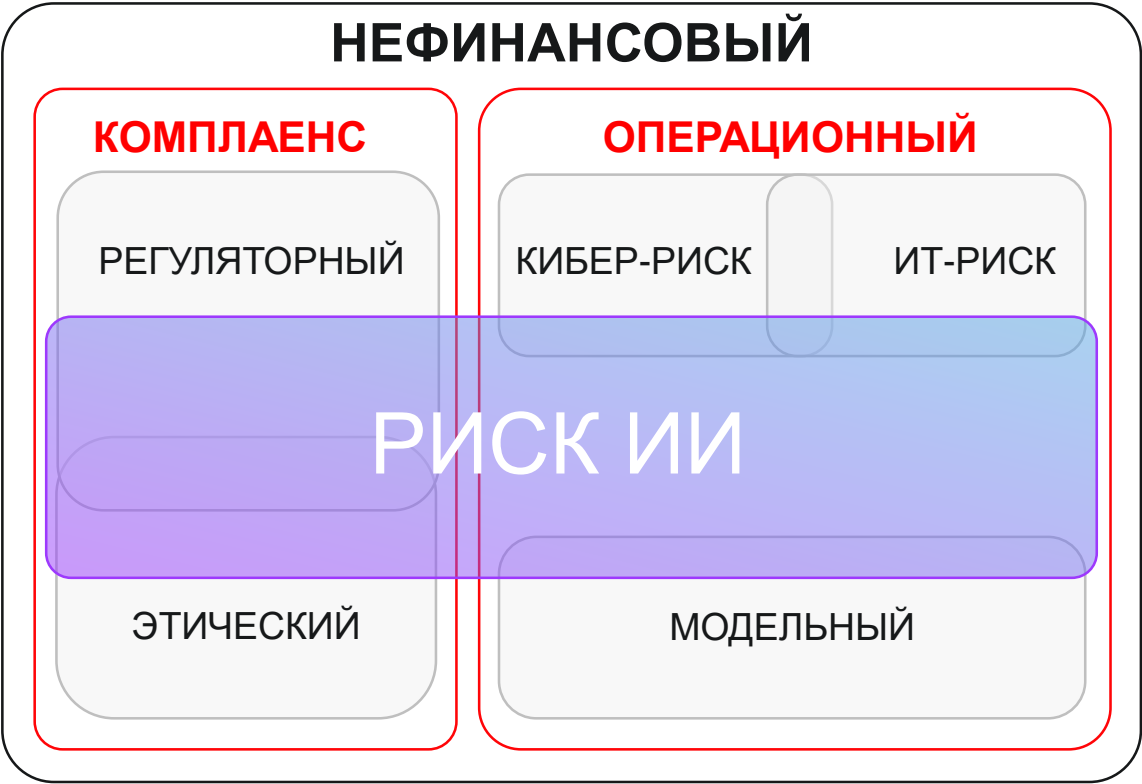


Понятие риска ИИ и его место в системе управления рисками

Риск ИИ – подвид нефинансового риска, находящийся на стыке:

- ИБ-рисков
- ИТ-рисков
- Модельных рисков
- Классических операционных рисков
- Этических
- Регуляторных

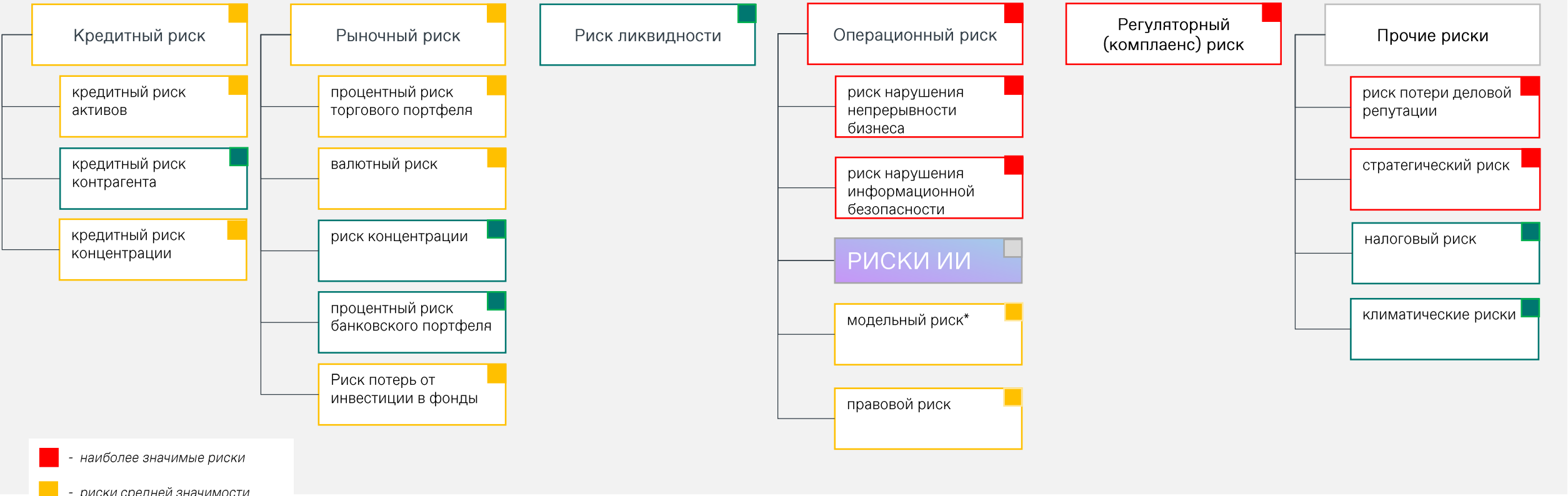
ВЗАИМОСВЯЗЬ РИСКА ИИ С ДРУГИМИ РИСКАМИ



СТРУКТУРА РЕКОМЕНДАЦИЙ В СФЕРЕ ИИ



КАРТА РИСКОВ МОСКОВСКАЯ БИРЖА



Список рисков, управление которыми обеспечивается на уровне Группы:

1. Климатический риск – распространение практики управления, реализованной в Московской Бирже, на уровень Группы
2. Риски ESG - внедрение системы управление на уровне Группы
3. Коммерческий риск – управление на уровне Группы
4. Системный риск - управляется в рамках иных видов риска

*включая остаточные риски

ЦИКЛ ПРОЦЕССА УПРАВЛЕНИЯ РИСКАМИ ИИ

ВЫЯВЛЕНИЕ РИСКОВ ИИ

- **Офис по развитию ИИ**
 - **Управление операционных рисков**
 - **УИБ**
 - **Бизнес-владелец риска (пользователь сервиса на базе ИИ)**
-
- Идентификация рисков ИИ в моделях и процессах, в которых внедряется ИИ
 - Составление матрицы рисков ИИ

ОЦЕНКА И ПРИСВОЕНИЕ УРОВНЯ РИСКА

Управление операционных рисков

- Оценка влияния реализации рисков ИИ
- Оценка негативных исходов от галлюцинаций моделей
- Анализ негативных исходов и оценка последствий в долгосрочной и краткосрочной перспективе на различные области риска

МИНИМИЗАЦИЯ РИСКОВ

Управление операционных рисков

- Определение мер по снижению уровня риска и лиц, ответственных за реализацию указанных мер

ВЕДЕНИЕ БАЗЫ ДАННЫХ РИСК-СОБЫТИЙ

Управление операционных рисков

- Выявление и регистрация риск-событий с описанием обстоятельств, повлекших его наступление, а также мероприятий, направленных на минимизацию последствий

УПРАВЛЕНИЕ РИСКАМИ ИИ ИНТЕГРИРОВАНО В ОБЩУЮ СИСТЕМУ СУР

МОНИТОРИНГ И КОНТРОЛЬ

- **Управление операционных рисков**
 - **УИБ**
-
- Установка КИР и их порогов
 - Мониторинг поведения моделей
 - Мониторинг достаточности применяемых мер защиты

РЕАГИРОВАНИЕ НА СОБЫТИЯ РИСКА

Управление операционных рисков

- Определение типа риск-события
- Осуществление мер, направленных на минимизацию их последствий
- Предоставление возможности сотрудникам и клиентам направлять жалобы и сообщения о риск-событиях

УЧЕТ ПРИМЕНЯЕМЫХ МОДЕЛЕЙ ИИ

- **Управление операционных рисков**
 - **Офис ИИ**
-
- Ведение реестра моделей
 - Валидация высокорисковых моделей
 - Регулярная инвентаризация моделей

РОЛИ И ОТВЕТСТВЕННОСТЬ

в управлении риском ИИ

Первая линия

Бизнес-владелец риска

Владелец модели/Пользователь

руководитель структурного подразделения, на бизнес-процессы которого реализация риска ИИ оказывает негативное влияние

- Соблюдение политик и правил использования ИИ-сервисов
- Ответственное обращение с данными, подаваемыми на вход
- Отслеживание и отчетность о подозрительных или аномальных результатах работы ИИ
- Участие в тестировании и обратной связи по функциональности ИИ
- Ознакомление с ограничениями ИИ

Владелец риска

Владелец модели/Разработчик

владелец бизнес-процесса или его этапа, в ходе осуществления которого появляются обстоятельства, обуславливающие возможную реализацию риска ИИ

- Обеспечение прозрачности и интерпретируемости алгоритмов ИИ
- Тестирование моделей на корректность и непредвзятость результата
- Разработка механизмов контроля качества
- Документирование процессов разработки, обучения моделей и источников данных
- Обеспечение соблюдения нормативных требований

Вторая линия

- Оценка и классификация рисков ИИ
- Разработка и внедрение политики управления рисками ИИ
- Мониторинг систем ИИ на соответствие требованиям безопасности и устойчивости к атакам
- Координация между разработчиками и пользователями для синхронизации мер по снижению рисков
- Реагирование на инциденты
- Учет применяемых моделей ИИ
- Выявление рисков ИИ
- Оценка и присвоение уровня риска ИИ
- Мониторинг и контроль рисков ИИ
- Минимизация выявленных рисков ИИ
- Реагирование на реализовавшиеся риски ИИ
- Ведение базы риск-событий

Оценка эффективности управления рисками в рамках выполнения первой и второй линиями задач в зоне своей ответственности (не позднее 3 месяцев после запуска)

Третья линия (внутренний аудит)

Ответственность за управление риском ИИ определена в соответствии с классической моделью трех линий на всем жизненном цикле модели



Определение уровня риска системы ИИ

ЗНАЧИМОСТЬ РИСКА

КРИТИЧНОСТЬ СИСТЕМЫ

- 01 **Критичность данных.** Категория данных, используемых при разработке модели ИИ
- 02 **Функциональная критичность.** Сфера применения модели ИИ, **критичность процессов**, для которых она используется
- 03 **Потенциальный ущерб.** Размер убытков, в тч регуляторных или ущерб деловой репутации, которые могут быть причинены в случае реализации риска
- 04 **Массовость потребления.** Количество клиентов/внутренних пользователей, при оказании услуг которым применяется ИИ.

ВЕРОЯТНОСТЬ РЕАЛИЗАЦИИ РИСКА

ОБЪЕМ ПОКРЫТИЯ СИСТЕМЫ

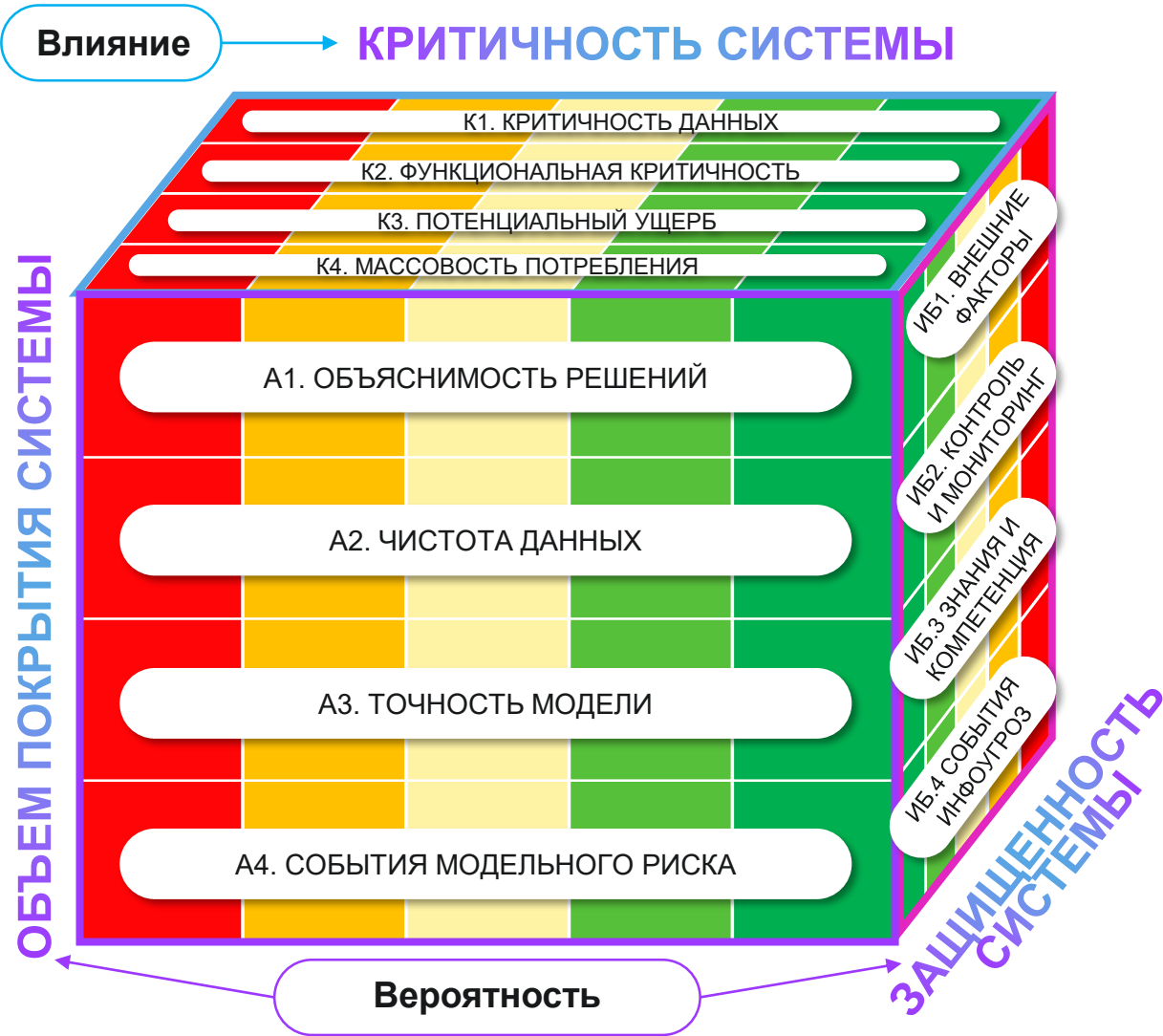
- 01 **Объяснимость решений**, принимаемых ИИ, наличие механизмов валидации и оценки качества выдаваемых результатов
- 02 **Чистота данных**, используемых в модели. Использование наборов данных, полученных от третьих лиц, или наборов данных, находящихся в открытом доступе в информационно-телекоммуникационной сети «Интернет»
- 03 **Точность модели** (комплексное определение на основе технико-статистических параметров Accuracy, Completeness, Specificity, Error Rate и тд)
- 04 **Наличие риск-событий, связанных с реализацией модельного риска**

ЗАЩИЩЕННОСТЬ СИСТЕМЫ

- 01 **Уровень зависимости от внешних факторов.** Участие третьих лиц в создании/эксплуатации модели ИИ, технологические ограничения (зависимость от данных, систем, оборудования)
- 02 **Уровень контроля и мониторинга модели. Модели угроз.** Наличие в модели механизмов контроля и защиты от нарушителей (Защита от атак, защита данных, системная безопасность и тд). Их качество
- 03 **Уровень знаний и компетенций** владельцев модели, разработчиков, пользователей модели
- 04 **Наличие риск-событий, связанных с реализацией инфоугроз** в отношении модели ИИ

ПАРАМЕТРЫ ПРИСВОЕНИЯ УРОВНЯ РИСКА СИСТЕМЕ ИИ

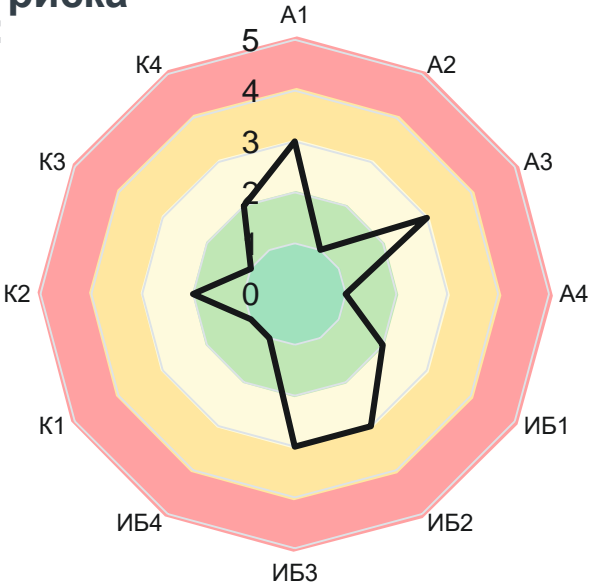
Оценка производится по 12 критериям, позволяющим оценить совокупный уровень риска ИИ для каждой конкретной системы



Распределение уровня риска по критериям /пример/:

Пример распределения критериев уровня риска для RAG- модели (RAG-RM), используемой для поиска информации по документам операционных рисков.

Для данной модели получена следующая балльная оценка:
Вероятность реализации риска (критерии А + критерии ИБ) = 17
Влияние (критерии К) = 6



Совокупный уровень риска системы ИИ:

		Влияние					
		минимальное	низкое	среднее	значительное	высокое	
		0-4	5-8	9-12	13-16	17-20	
Вероятность	минимальная	0-8					
	низкая	9-16					
	средняя	17-24		🚩 RAG-RM			
	значительная	25-32					
	высокая	33-40					

ИСТОЧНИКИ ИНФОРМАЦИИ О РИСКАХ ИИ

ВНУТРЕННИЕ

ИТ-система

- Отчеты по результатам тестирований (DR, нагрузка, регрессионное, UAT и тд);
- Данные ИТ-мониторинга;
- Данные по результатам DQ-анализа.

Система ИИ

- Результаты валидации (производительность, корректность, устойчивость);
- Отчеты по аудиту модели;
- Данные базы событий модельного риска.

Объект СЗИ

- Результаты пентестов;
- Данные мониторинга ИБ;
- Данные базы инцидентов ИБ.

Объект СУР

- Результаты: RCSA, оценки новых продуктов, оценки рисков третьих сторон, анализа бизнес-процессов, сценарного анализа, стресс-теста;
- Результаты тренингов;
- Данные БДР, БДСОР, КИР, внутренних аудитов;
- Добровольные сообщения и тд

*перечень не является исчерпывающим и пополняется с появлением новых источников информации о рисках ИИ

ВНЕШНИЕ

Отчеты и исследования
Банка России

Российские НПА
по теме рисков
ИИ

Международные
НПА по теме
рисков ИИ

Международные исследования, отчеты
и диссертации в области ИИ

Открытые базы
атак и
уязвимостей

Открытые базы
этических
принципов

Открытые
проекты и
стандарты
(АФТ, OWASP, NIST,
MITRE и др.)

Результаты внешних аудитов

Открытые
фреймворки
(AI Sec, AI TRISM, AI
RMF, AI RedTeam, SAIF
Map, MLDevSecOps,
DASF и др)

Открытые
форумы и
сообщества
(Git, Habr, cyberorda и
др)

Открытые
рекомендации
представителей
сообщества
(Яндекс, Сбер, РТ и
др.)



Управление рисками на этапах жизненного цикла систем ИИ

Объем вовлечения второй линии

ПРОДУКТОВЫЙ РИСК-МЕНЕДЖМЕНТ – заранее делаем все, чтобы минимизировать риски до момента публикации системы в прод

КЛАССИЧЕСКИЙ РИСК-МЕНЕДЖМЕНТ

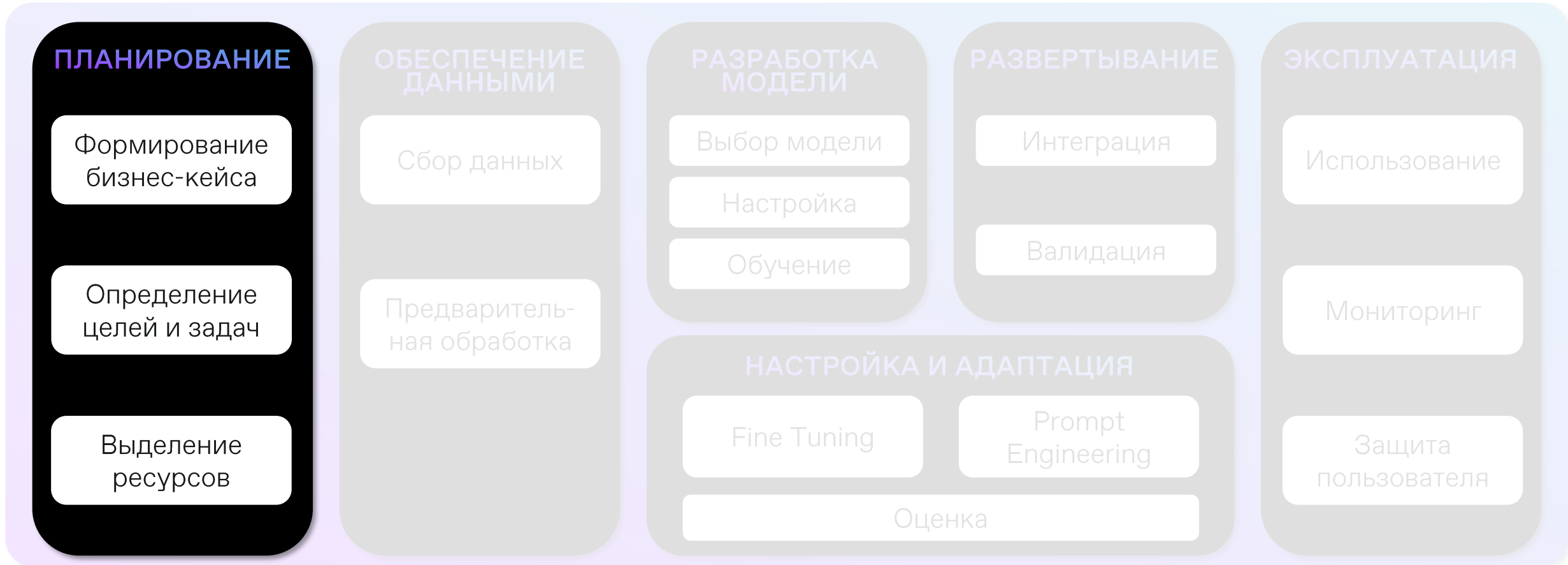


На каждом из этапов ЖЦ осуществляются действия по управлению рисками, конечная цель которых – не допустить реализацию риска ИИ и поддерживать совокупный уровень риска СИСТЕМЫ ИИ на приемлемом уровне.

Каждый идентифицированный риск проходит классическую оценку.
Уровень риска, выявленного на этапе жизненного цикла системы ИИ, определяется в соответствии с ВНД и учитывается при определении совокупного уровня риска системы ИИ.

ЖИЗНЕННЫЙ ЦИКЛ СИСТЕМЫ ИИ

концепция MLSecOps (включая специфику ГенИИ)



MLSecOps — подход, который объединяет операционные аспекты машинного обучения с вопросами безопасности. Направлен на снижение рисков, которые могут принести модели AI/ML/GenAI в организацию.

СОСТАВНЫЕ ЭЛЕМЕНТЫ ЭТАПА ЖИЗНЕННОГО ЦИКЛА

Формирование бизнес-кейса

Определение целей и задач

Выделение ресурсов

Сформированы цели и задачи, которые должна решать система ИИ.

Определены процессы, в которых будет использоваться система ИИ.

Определены потребители.

КЛЮЧЕВЫЕ РИСКИ



СТРАТЕГИЧЕСКИЙ

Риски черного ящика как стратегическая и юридическая уязвимость, проблемы ответственности за принимаемые решения, проблемы автономности ИИ



РИСК ТРЕТЬИХ СТОРОН

Стратегическая зависимость от одного поставщика ИИ-услуг

- Техническая (технологии)
- Зависимость от данных
- Сервисная
- Экономическая и ресурсная
- Зависимость от «личности» модели



ПРАВОВЫЕ И ЭТИЧЕСКИЕ РИСКИ

Проблемы интеллектуальной собственности (включая права на голос, изображение, сгенерированный контент), ограничения использования OSS, проблемы дискриминации и нарушения морально-этических норм и др.

Ключевые последствия:

- когнитивный раскол на пользователей и оракулов;
- атрофия управленческих компетенций;
- скрытый саботаж;
- судебные разбирательства, связанные с нарушением интеллектуальных прав;
- судебные разбирательства, связанные с дискриминацией и нарушением морально-этических норм.

ИНСТРУМЕНТАРИЙ CUP/MLSecOps

1. Анализ категорий потребителей результатов модели, анализ критичности модели. Определение ограничений модели



Определение значимости модели

Определение политик конфиденциальности, правовых и этических ограничений использования

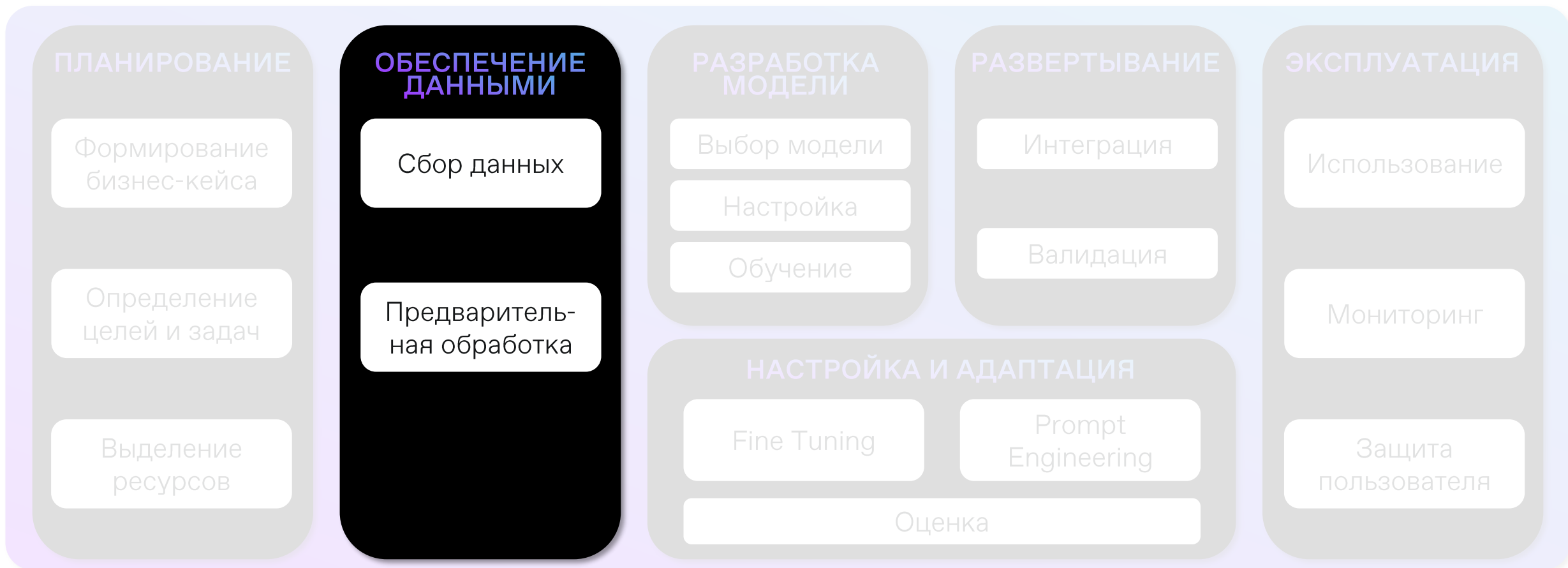
2. Анализ предполагаемой схемы реализации. Оценка техстека, зависимости от третьих лиц



Определение уровня риска зависимости от третьих сторон

ЖИЗНЕННЫЙ ЦИКЛ СИСТЕМЫ ИИ

концепция MLSecOps (включая специфику ГенИИ)



MLSecOps — подход, который объединяет операционные аспекты машинного обучения с вопросами безопасности. Направлен на снижение рисков, которые могут принести модели AI/ML/GenAI в организацию.

СОСТАВНЫЕ ЭЛЕМЕНТЫ
ЭТАПА ЖИЗНЕННОГО ЦИКЛА

Сбор данных

Разметка данных

Обработка данных

Формирование признаков

Сформирован
обработанный,
высококачественный
набор данных, который
подчиняется некоторой
закономерности.

КЛЮЧЕВЫЕ РИСКИ

УТЕЧКИ ДАННЫХ

Риск утечки конфиденциальных данных, таких как ПДн, бизнес-информация или алгоритмы.

РИСКИ ЦЕПОЧКИ ПОСТАВОК

Уязвимости в цепочке поставок могут привести к нарушениям безопасности и смещению данных.

ОТРАВЛЕНИЕ ДАННЫХ/ФУНКЦИЙ

Манипуляции с данными на этапах обеспечения данными и обучения модели, что влияет на ее надежность и результаты.

ИНСТРУМЕНТАРИЙ
СУР/MLSECOPSОбеспечение организованного контроля
доступа к данным

(ролевая модель, аутентификация, DLP)

Классификация и управление данными

(классификация, настройка DQ, DQ для оценочных данных, организация безопасного хранения данных, разные пайплайны разработки для данных разной критичности)

Отслеживание происхождения данных

(определение формата информации о происхождении, формирование White-List источников, внедрение контроля модификации)

Защита от отравления данных

(Внедрение контроля целостности, внедрение алгоритмов обнаружения аномалий и предотвращения атак)

Защита извлеченных признаков от
манипулирования

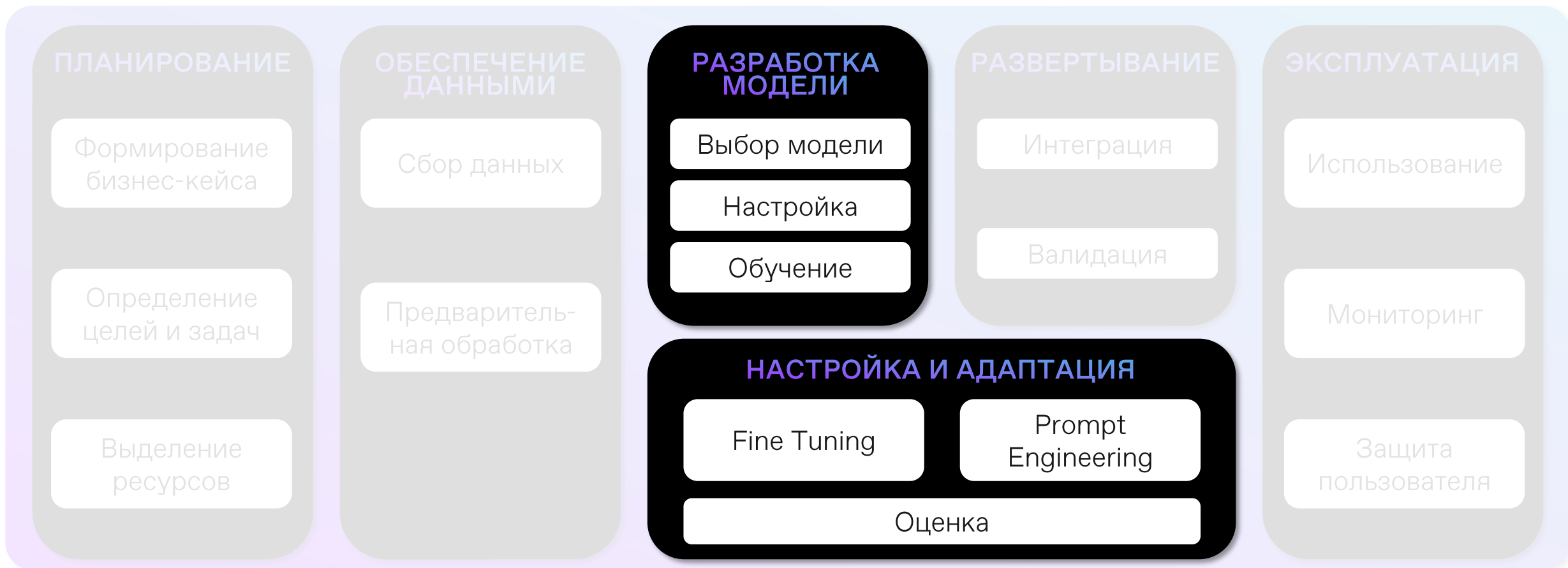
(использование шума во время извлечения, орг. безопасного хранилища)

Организация платформы управления
данными

(безопасная цепочка поставок, шифрование при хранении и передаче, доверенная база компонентов, руководства по безопасной работе, автоматизация конвейера по безопасной работе с данными, формализация ЖЦ обеспечения данными)

ЖИЗНЕННЫЙ ЦИКЛ СИСТЕМЫ ИИ

концепция MLSecOps (включая специфику ГенИИ)



MLSecOps — подход, который объединяет операционные аспекты машинного обучения с вопросами безопасности. Направлен на снижение рисков, которые могут принести модели AI/ML/GenAI в организацию.

СОСТАВНЫЕ ЭЛЕМЕНТЫ
ЭТАПА ЖИЗНЕННОГО ЦИКЛА

Проектирование архитектуры

Настройка модели

Обучение и тестирование

Fine Tuning

Prompt Engineering

Оценка и валидация

Разработана модель,
соответствующая требованиям к
точности, надежности,
масштабируемости, скорости,
устойчивости к шуму и неточным
данным.

Модель обучается и адаптивна,
ее результаты объяснимы.

КЛЮЧЕВЫЕ РИСКИ

УЯЗВИМОСТИ ВЕКТОРОВ И ЭМБЕДДИНГОВ

Слабые места в моделях векторных представлений, способные вызвать непредсказуемые ошибки.

НЕАВТОРИЗОВАННЫЕ ДАННЫЕ
ДЛЯ ОБУЧЕНИЯ

Модель делает вывод на основе данных, к которым у нее не должно быть доступа. Может привести к утечке конфиденциальной информации, непредсказуемому поведению модели.

ПРОМТ-ИНЪЕКЦИИ

Вмешательство пользователя в запросы, чтобы изменить результаты и функции модели.

ОТРАВЛЕНИЕ МОДЕЛИ/ ПОДДЕЛКА
ИСТОЧНИКОВ

Манипуляции с данными на этапах обеспечения данными и обучения модели, что влияет на ее надежность и результаты, может повлечь риск фальсификации исходного кода модели.

НЕКОРРЕКТНАЯ/ИЗБЫТОЧНАЯ
ОБРАБОТКА ВЫХОДНЫХ ДАННЫХ

Отсутствие проверок и обработки вывода может привести к уязвимостям в приложениях, сохранению или использованию данных сверх допустимых пределов, недостаточной детализации ответов.

ИНСТРУМЕНТАРИЙ
СУР/MLSECOPSКонтроль входных/выходных
данных, использование
доверенных компонент

(проверка и нормализация вводимых данных, логирование и мониторинг данных, определение формата выходных значений, проверка и фильтрация выходных значений, установка границ доверия между моделью, внешними источниками и расширяемой функциональностью, введение этических ограничений, децентрализация, Federated Learning)

Состязательное обучение

(Adversarial Training - генерация adversarial примеров и включение их в процесс обучения для повышения точности и устойчивости к атакам.

Дифференциально приватный SGD (DP-SGD) - обеспечивает конфиденциальность за счёт равномерного распределения privacy costs по всем этапам обучения.)

Тестирование и валидация

(анализ защищенности модели, внедрение методов интерпретации результатов и оценки точности

Проверки обучаемости – кросс-валидация. Контроль производительности на валидации, мониторинг точности модели на валидационных данных и тд

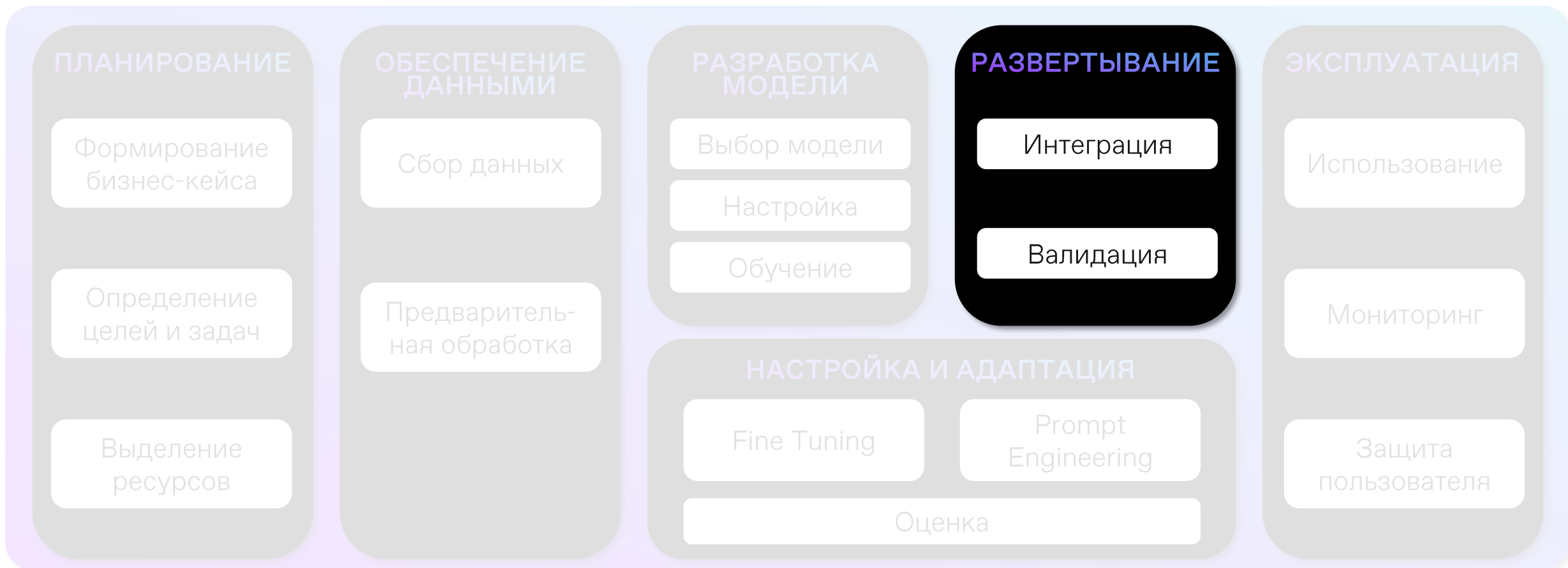
Проверки адаптивности - Онлайн обучение, Transfer Learning, A/B тесты, Concept Drift Detection и тд

Проверки объяснимости - SHAP, LIME, PDP, Feature Importance, библиотеки интерпретации и тд)

Анализ защищенности – тестирование на устойчивость, анализ смещений (Bias Detection)

ЖИЗНЕННЫЙ ЦИКЛ СИСТЕМЫ ИИ

концепция MLSecOps (включая специфику ГенИИ)



MLSecOps — подход, который объединяет операционные аспекты машинного обучения с вопросами безопасности. Направлен на снижение рисков, которые могут принести модели AI/ML/GenAI в организацию.

СОСТАВНЫЕ ЭЛЕМЕНТЫ
ЭТАПА ЖИЗНЕННОГО ЦИКЛА

Разработка

Сборка

Тестирование системы

Развертывание

Настройка мониторинга

Система ИИ готова к
началу эксплуатации

КЛЮЧЕВЫЕ РИСКИ

УЯЗВИМОСТИ ВЕКТОРОВ
И ЭМБЕДДИНГОВ

Слабые места в моделях векторных представлений, способные вызвать непредсказуемые ошибки.

УЯЗВИМОСТИ КОДА

Слабые места в коде системы, позволяющие злоумышленнику осуществлять атаки на систему и саму модель.

УТЕЧКИ ДАННЫХ И СИСТЕМНЫХ
ИНСТРУКЦИЙ

Утечка внутренних системных подсказок, которые раскрывают конфиденциальные настройки.

РИСКИ ЦЕПОЧКИ ПОСТАВОК

Уязвимости в цепочке поставок могут привести к нарушениям безопасности и смещению данных.

НЕКОРРЕКТНАЯ/ИЗБЫТОЧНАЯ
ОБРАБОТКА ВЫХОДНЫХ ДАННЫХ

Отсутствие проверок и обработки вывода может привести к уязвимостям в приложениях, сохранению или использованию данных сверх допустимых пределов.

ИНСТРУМЕНТАРИЙ
СУР/MLSECOPSКонтроль и воспроизводимость
проводимых экспериментов

(документирование и стандартизация экспериментов, использование контейнеризации)

Защита активов

(шифрование, контроль аномалий, создание центрального хранилища для моделей с контролем доступа, внедрение в активы цифрового водяного знака, Certified Defenses, DLP-фильтрация, контент-фильтрация, SBOM, SPDX)

Организация платформы
управления разработкой

(организация хранилища обученных моделей, контроль версииности, разработка политик и руководств для процесса разработки, автоматизация конвейера разработки модели, формализация всего ЖЦ разработки модели, организация безопасной цепочки поставок)

MLOps-конвейер

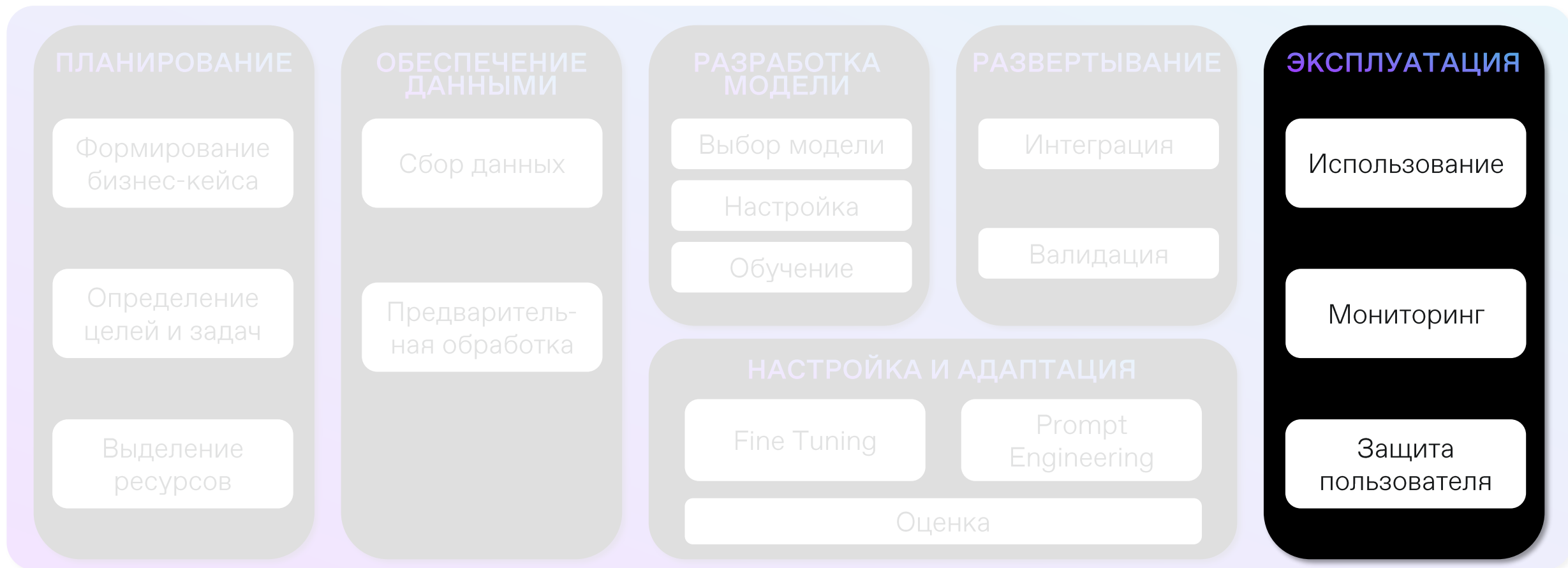
(построение конвейера, автоматизированные инструментальные проверки, настройка и обновление сигнатурных правил, организация взаимодействия команд)

Анализ защищенности системы ИИ

(Robust Training, правила Malware-gen, firewall фильтрация retrieved-chunk, BugBounty, управление инцидентами, аудиты, модели угроз)

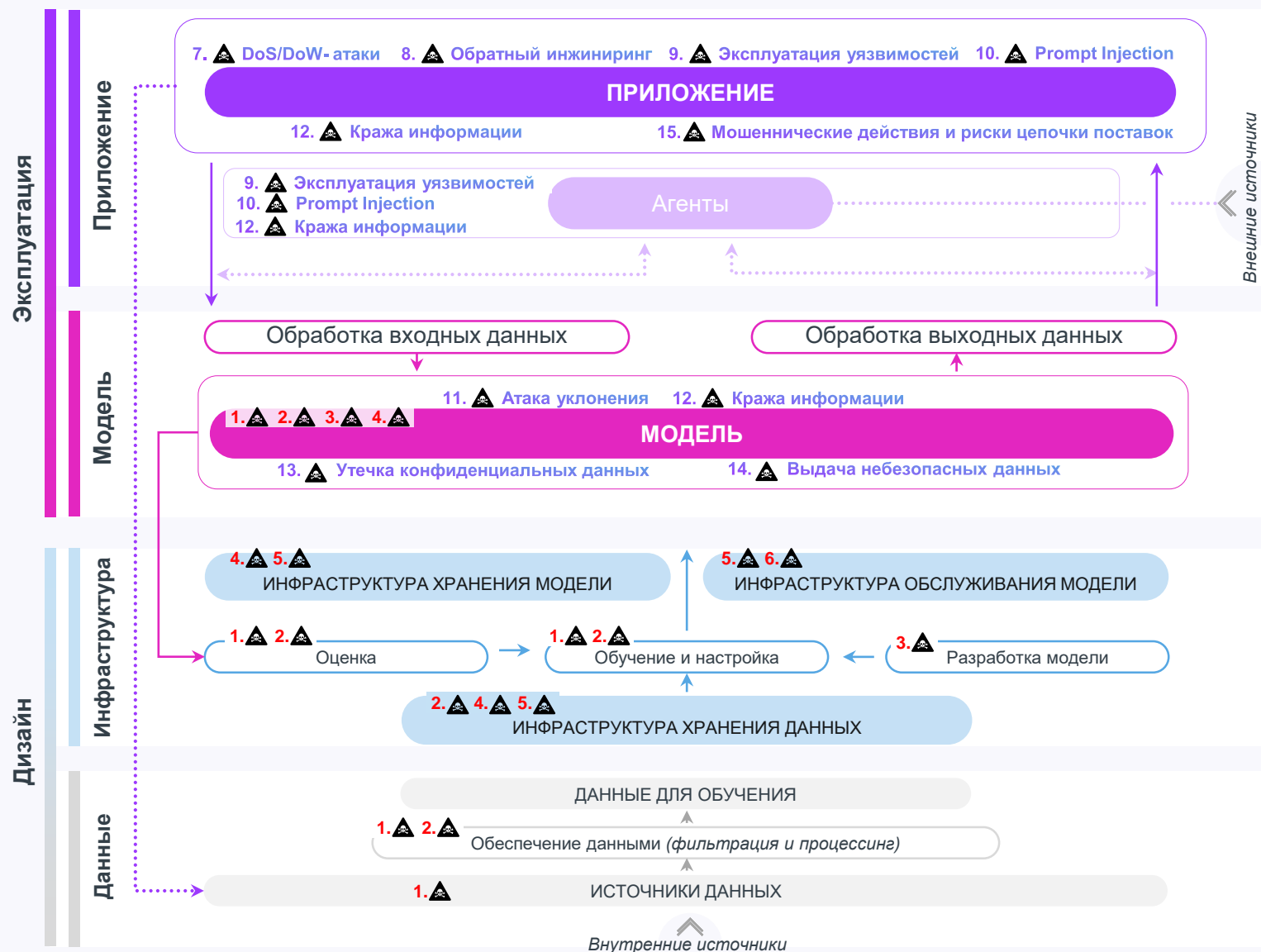
ЖИЗНЕННЫЙ ЦИКЛ СИСТЕМЫ ИИ

концепция MLSecOps (включая специфику ГенИИ)



MLSecOps — подход, который объединяет операционные аспекты машинного обучения с вопросами безопасности. Направлен на снижение рисков, которые могут принести модели AI/ML/GenAI в организацию.

НЕПРЕРЫВНОЕ СОВЕРШЕНСТВОВАНИЕ СИСТЕМЫ КОНТРОЛЕЙ СУР/MLSECOPS

РИСКИ НА УРОВНЕ
ПОЛЬЗОВАТЕЛЕЙ

ЧРЕЗМЕРНАЯ АГЕНТНОСТЬ

Чрезмерное использование LLM для принятия решений, выходящее за пределы безопасных границ.

ВВЕДЕНИЕ В ЗАБЛУЖДЕНИЕ

Генерация ложной или вводящей в заблуждение информации с потенциальными последствиями, галлюцинирование и хаотичное поведение модели.

НЕОГРАНИЧЕННОЕ ПОТРЕБЛЕНИЕ

Неконтролируемое использование вычислительных ресурсов, что вызывает перегрузку системы.

КЛЮЧЕВЫЕ РИСКИ НА УРОВНЕ
ДИЗАЙНА

1. ОТРАВЛЕНИЕ ДАННЫХ
2. НЕАВТОРИЗОВАННЫЕ ДАННЫЕ
3. ПОДДЕЛКА ИСТОЧНИКА
4. ИЗБЫТОЧНАЯ ОБРАБОТКА ДАННЫХ
5. ЭКСФИЛЬТРАЦИЯ МОДЕЛИ
6. ФАЛЬСИФИКАЦИЯ МОДЕЛИ ПРИ РАЗВЕРТЫВАНИИ



ПРИЛОЖЕНИЯ

Типология атак
на системы ИИ

Специфические
векторы атак
на различные типы
систем ИИ

Пример применения
фреймворка для
оценки риска
системы ИИ

ОСНОВНЫЕ ТИПЫ АТАК НА СИСТЕМЫ ИИ

АТАКИ НА ЭТАПЕ ИНФЕРЕНСА (EVASION ATTACKS)

Цель: обмануть уже обученную модель на этапе прогнозирования.

Быстрые атаки на основе градиента (White-Box)

- FGSM (Fast Gradient Sign Method): Создает возмущение за один шаг, двигаясь в направлении знака градиента функции потерь.
- PGD (Projected Gradient Descent): Многократная итеративная версия FGSM.
- C&W (Carlini & Wagner): Мощная оптимизационная атака, минимизирующая величину возмущения при условии смены класса.

Атаки, ищущие минимальное возмущение

- DeepFool: Итеративно приближает образец к границе решения, находя минимальное необходимое возмущение для смены класса.
- One-Pixel Attack: Атака в крайне ограниченном пространстве - меняет всего один пиксель, демонстрируя хрупкость моделей.

Универсальные и физические атаки

- Universal Adversarial Perturbation: Генерирует одно возмущение, которое, будучи добавленным к любому изображению из датасета, с высокой вероятностью вызывает ошибку.
- Adversarial Patch: Создает небольшой патч (например стикер), который размещается в кадре и заставляет модель ошибаться. Критично для компьютерного зрения в реальном мире.
- Physical Adversarial Attack: Перенос атак в физический мир (например, наклейки на дорожные знаки для обмана автопилота).

АТАКИ НА КОНФИДЕНЦИАЛЬНОСТЬ (PRIVACY ATTACKS)

Цель: извлечь конфиденциальную информацию о данных обучения или о самой модели.

Атаки на восстановление данных (Model Inversion & Extraction)

- Model Inversion Attack: Угроза приватности.
- Model Extraction / Stealing Attack: Угроза интеллектуальной собственности.
- Attribute Inference Attack: Вывод скрытых атрибутов, не являющихся прямой целью предсказания, на основе выхода модели.

Атаки на членство в обучающей выборке (Membership Inference)

- Основная идея: Определить, присутствовал ли конкретный образец данных в обучающем наборе модели.

АТАКИ В УСЛОВИЯХ "ЧЕРНОГО ЯЩИКА" (BLACK-BOX ATTACKS)

Атакующий не имеет доступа к внутренностям модели (архитектуре, весам), только к ее выходу.

Атаки на основе запросов (Query-Based)

- Score-Based Attack: Атакующий видит вероятности классов, выдаваемые моделью.
- Decision-Based Attack: Атакующий видит только итоговый класс. Атака часто начинается с большого шума, который постепенно уменьшается, пока образец остается адверсариальным.
- Boundary Attack: Разновидность decision-based атаки, которая "скользит" вдоль границы решения, минимизируя возмущение.

Атаки, использующие перенос (Transferability Attacks)

- Transfer Attack: Создание адверсариального примера на локальной модели-суррогате (белый ящик) с последующей атакой с его помощью целевой black-box модели.

АТАКИ НА СПЕЦИФИЧЕСКИЕ АРХИТЕКТУРЫ И МОДАЛЬНОСТИ

NLP / LLM Attacks

- Text Adversarial Attack: Замена слов, добавление опечаток/изменений, меняющих решение модели, но не смысл для человека.
- Prompt Injection / Jailbreak Attack

Атаки на системы обучения с подкреплением (RL)

- Reward Hacking: Агент находит неожиданный способ максимизировать функцию вознаграждения, не решая поставленную задачу.
- Environment Poisoning: Изменение среды обучения, чтобы агент выработал неэффективную или опасную стратегию.

Атаки на графические и другие модели

- Graph Adversarial Attack: добавление/удаление связей или узлов.
- Audio Adversarial Attack: Добавление незаметного для человека шума в аудиозапись, чтобы обмануть систему распознавания речи.

АТАКИ НА ЭТАПЕ ОБУЧЕНИЯ (POISONING ATTACKS)

Цель: злоумышленник вмешивается в процесс обучения, чтобы подорвать работу модели.

Отравление данных (Data Poisoning)

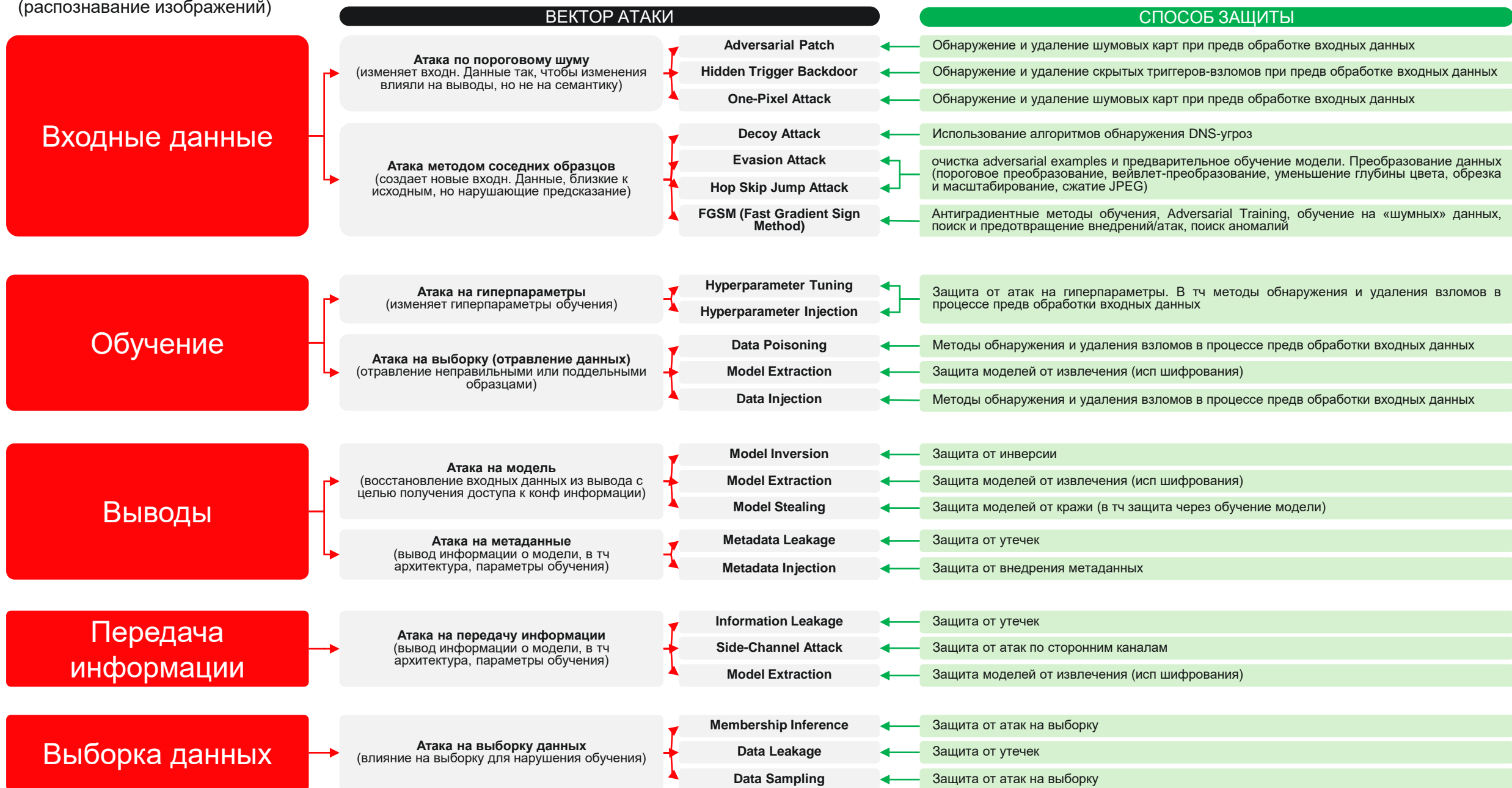
- Внесение в обучающий набор вредоносных примеров, которые ухудшают общую производительность модели или смещают ее решения.

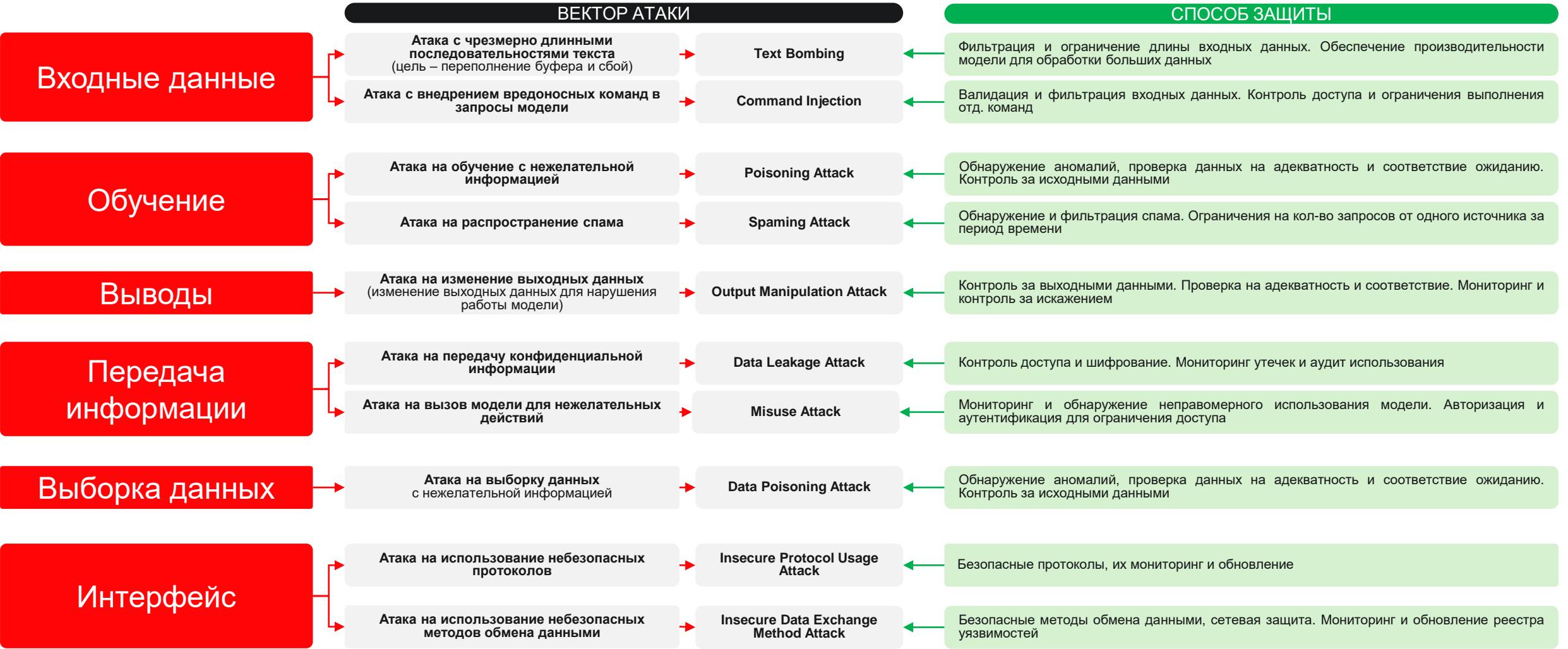
Бэкдор-атаки (Backdoor Attacks)

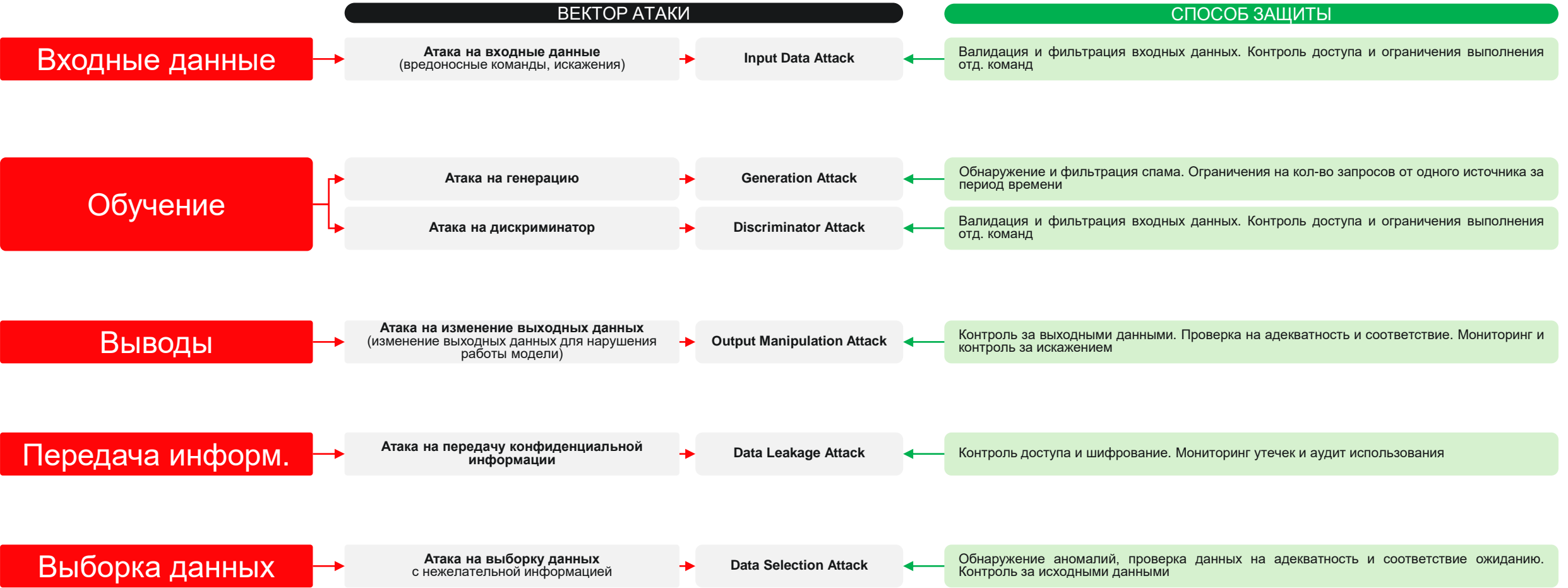
- Механизм: Модель обучается корректно работать на "чистых" данных, но при появлении в входе специального триггера она выдает заданный злоумышленником ошибочный результат/запускает исполнения команд/кода для исполнения задач злоумышленника.

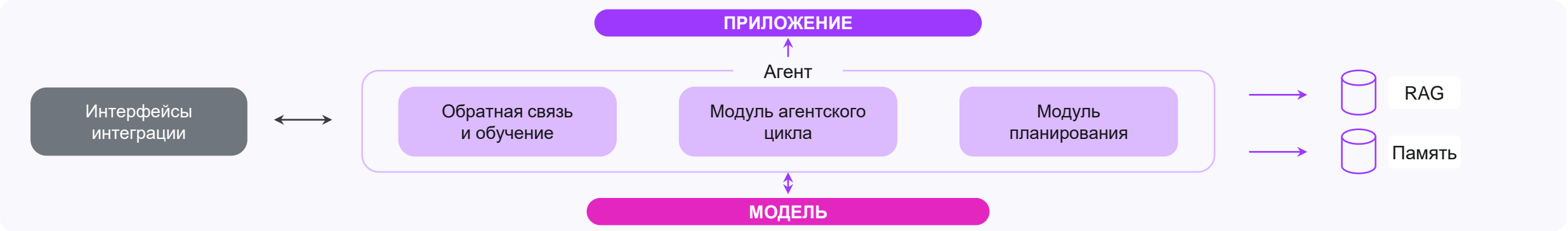
ВЕКТОРЫ АТАК. CNN-сети

(распознавание изображений)









	ВЕКТОР АТАКИ	СПОСОБ ЗАЩИТЫ
Интерфейс	Атака на поведение модели	Prompt Injection
	Атака на приложение	Denial of Service (DoS)
	Атака на передачу информации (вывод информации о модели, в тч архитектура, параметры обучения)	Output Data Attack
Исполнение и инструменты	Нанесение ущерба легитимными инструментами	Tool Misuse
	Внедрение вредоносных инструкций в описание (метаданные) самого инструмента	Tool Poisoning
	Использование избыточных прав	Privilege Escalation
	Эксплуатация уязвимостей в механизме проверки прав	Auth Bypass & Impersonation
Инфраструктура и оркестрация	Атака через цепочку поставок (внешние источники данных, компоненты системы)	Supply Chain Attacks
	Атака через цепочку агентов	Cross Agent Communication Poisoning
	Неконтролируемое потребление вычислительных ресурсов, токенов, вызовов API	Ресурсоёмкие атаки
Выборка данных и знания	Внедрение вредоносных или ложных данных в документы, используемые RAG-системой	Knowledge Base Poisoning
	Извлечение конфиденциальных данных	Sensitive Information Disclosure
	Манипуляция поиском в RAG	Retrieval Manipulation
	Восстановление конфиденциального текста из его векторного представления	Embedding Inversion

ПРИМЕР ПРИМЕНЕНИЯ ФРЕЙМВОРКА НА РЕАЛЬНЫХ ПРОДУКТАХ

AI NEWSFEED ИНТЕГРАТОР НОВОСТЕЙ И КОРПОРАТИВНЫХ СОБЫТИЙ

Описание

Backend сервис для анализа тональности и важности новостей.

Интеграция потока новостей и событий и представление их в читаемом, понятном виде:

- поток новостей по эмитентам;
- корпоративные действия;
- финансовая отчетность;
- мнения аналитиков.

Потребители

Использование внутренними подразделениями Мосбиржи (Compliance, НРД, Рынки, ФинУслуги)

Данные

/все из открытых источников/

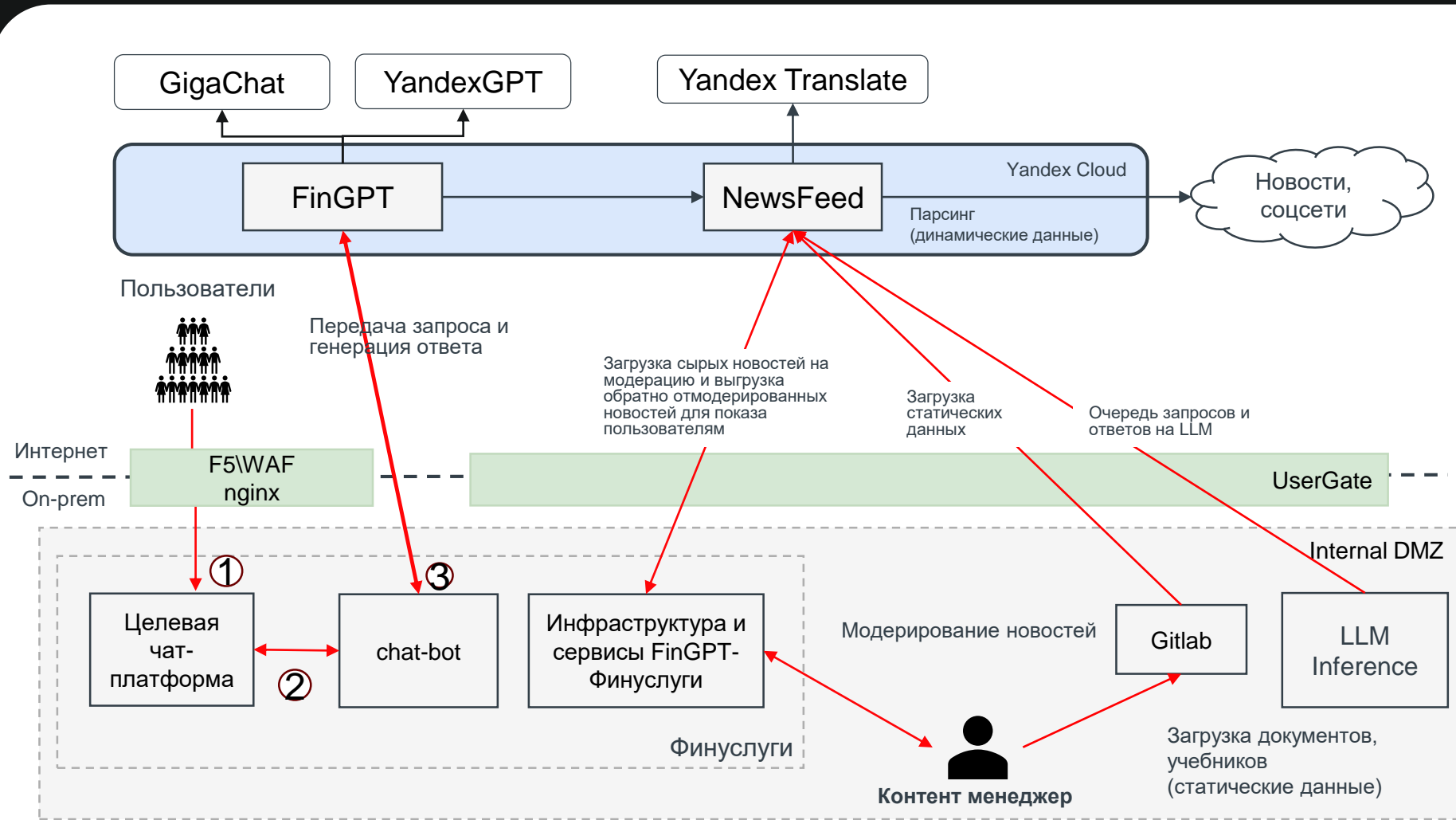
- новости;
- социальные сети;
- аналитические обзоры;
- открытая корпоративная отчетность.

РЕЗУЛЬТАТ МОДЕЛИ

Уровень управляемости HITL

Позволяет выгрузить обработанные ИИ новости с указанием:

1. Актуальной ссылки на сайт, с которого подтягивается новость
2. Сантимента (характер новости – негативный, нейтральный, позитивный)
3. Категории новости (экономическая, политическая, прочая)
4. Наименование компании и страны контрагента
5. Даты публикации новости



Модерация информации

- Статические данные (документы, учебники) для RAG модели модерируются on-Prem и передаются вместе с кодом (git)
- Динамические данные (новости) модерируются также on-Prem и передаются в RAG через инфраструктуру Финуслуг
- Не прошедшие модерацию данные не участвуют в формировании ответов для пользователя

AI NEWSFEED
ИНТЕГРАТОР
НОВОСТЕЙ И
КОРПОРАТИВНЫХ
СОБЫТИЙ

УРОВЕНЬ РИСКА СИСТЕМЫ ИИ



КРИТИЧНОСТЬ СИСТЕМЫ

Критичность данных - **низкая**

Функциональная критичность - **выше среднего**

Потенциальный ущерб - **минимальный**

Объем воздействия - **минимальный**

ОБЪЕМ ПОКРЫТИЯ СИСТЕМЫ

Объяснимость решений – **ниже среднего**

Чистота данных - **средняя**

Точность модели - **средняя**

События модельного риска - **нет**

ЗАЩИЩЕННОСТЬ СИСТЕМЫ

Зависимость от внешних факторов - **высокая**

Уровень контроля – **ниже среднего**

Уровень компетенций - **средний**

События инфоугроз - **нет**

Ключевые риски системы ИИ AI NEWSFEED

Описание риска	Присущий уровень риска	Мероприятия по минимизации рисков	Остаточный уровень риска	Плановый уровень риска
РИСК ТРЕТЬИХ СТОРОН Стратегическая зависимость от одного поставщика ИИ-услуг <ul style="list-style-type: none">Техническая (GigaChat, YandexGPT)Зависимость от данныхСервисная (Yandex облако)Зависимость от «личности» модели (LLM)		<ol style="list-style-type: none">Пилотирование продукта на некритичных процессах.Мониторинг и управление рисками критичных поставщиков.Установка SLA для сторонних сервисов.		
РИСКИ ЦЕПОЧКИ ПОСТАВОК Эксплуатация уязвимостей в облачной инфраструктуре/ сторонних агентах, используемых системой, может привести к утечке данных, нарушению бизнес-процессов		<ol style="list-style-type: none">CI/CD процессы в разработке (sast/dast/sca, код ревью),Использование ПО из внутренних репозитариевФормирование и мониторинг White List источников.Проведение анализа защищенностиИспользование средств сетевой защиты при выгрузке данных и отправке запросовМониторинг ИБ (SOC, SIEM)аутентификация внутренняя через idp		
НЕКОРРЕКТНАЯ ОБРАБОТКА ВЫХОДНЫХ ДАННЫХ Отсутствие проверок и обработки вывода может привести к уязвимостям в приложениях. Контент, который создает модель, не проверяют и не очищают, прежде чем передавать в другие системы.		<ol style="list-style-type: none">Формирование и мониторинг White List источников.Модерация информацииИспользование средств сетевой защиты при выгрузке данных и отправке запросов		
ОТРАВЛЕНИЕ ДАННЫХ «отравленных» данных, когда злоумышленник загрязняет общедоступные датасеты, чтобы создать скрытую уязвимость при обучении.				
ВВЕДЕНИЕ В ЗАБЛУЖДЕНИЕ Генерация ложной или вводящей в заблуждение информации с потенциальными последствиями.		<ol style="list-style-type: none">Ограничение функциональности – использование результатов работы модели в некритичных процессах с последующей перепроверкой выходных данных (выборочная перепроверка источников новостей)		
ЧРЕЗМЕРНАЯ АГЕНТНОСТЬ Чрезмерное использование LLM для принятия решений, выходящее за пределы безопасных границ.				