





для построения инфраструктуры высоконагруженных корпоративных и государственных информационных систем

скала р

ЛЕТ серийного выпуска

скала р

680

комплексов в промышленной эксплуатации 10

ТЫС. + вычислительных узлов

Продуктовые направления Скала^р



решения для высоконагруженных корпоративных и государственных систем



Динамическая инфраструктура

Машины динамической инфраструктуры Скала^р МДИ

на основе решений BASIS для создания динамической конвергентной и гиперконвергентной инфраструктуры ЦОД и виртуальных рабочих мест пользователей



Управление большими данными

Машины больших данных Скала^р МБД.8

на основе решений ARENADATA и PICODATA для создания инфраструктуры хранения, преобразования, аналитической, статистической обработки данных, а также распределенных вычислений



Высокопроизводительные базы данных

Машины баз данных Скала р МБД

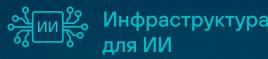
на основе решений Postgres Pro для замены Oracle Exadata в высоконагруженных системах с обеспечением высокой доступности и сохранности критически важных данных



Интеллектуальное хранение данных

Машины хранения данных Скала^р МХД

на основе технологии объектного хранения S3 для геораспределенных катастрофоустойчивых систем с сотнями миллионов объектов различного типа и обеспечения быстрого доступа к ним



Машина искусственного интеллекта Скала^р

на основе оптимизированного программноаппаратного стека для максимальной производительности при работе с моделями ИИ

- Использование опыта технологических лидеров (гиперскейлеров)
- Использование самых зрелых и перспективных технологий в кооперации с технологическими лидерами российского рынка в каждом из сегментов

Модульная платформа Скала р



Использование опыта технологических лидеров — гиперскейлеров

Единый принцип модульной компоновки и платформенный подход

Единая облачная система управления сервисами



laaS



PaaS



DBaaS

Единая система управления ресурсами и эксплуатацией



Разделение ресурсов



Мультитенантность



Автоматизация

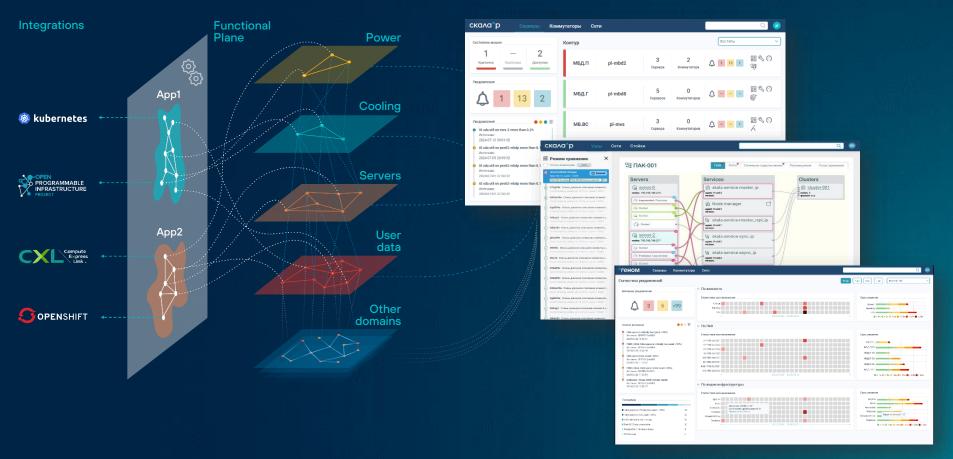
Серверная виртуализация осоверная виртуализация осоверная особерная осоверная особерная особерн



Программная платформа Скала^р



Объединения различных доменов управления в единую функциональную графовую CMDB Комплексное решение для эксплуатации инфраструктуры уровня ЦОД



- Единая точка обзора состояния контура
- Обозримость и удобство управления ЦОД
- Цифровой двойник инфраструктуры
- Контроль изменений быстроменяющихся топологий
- Моделирование изменений в инфраструктуре
- Высокая степень автоматизации
- Построение Al-Copilot для управления ЦОД

Скала^р – Secure by Design

Скала МБД.ИИ





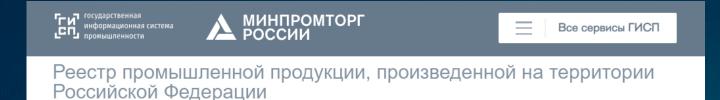
ПАК Скала р в Реестрах РФ



Машины

Модули

Компоненты



Машины

Модули

Программное обеспечение





Евразийский

Машины

Модули

Программное обеспечение

Соответствуют критериям доверенного ПАК

- Технологическая независимость
- Информационная безопасность
- Функциональная устойчивость

Импортозамещение: сложность выбора Отсутствие технологического лидерства





Вычислительная

инфраструктура





№ ГРАВИТОН

AQUARIUS

OpenYard

// kraftway*

GAGAR✓ N









Проблемы отсутствия ИТ-лидеров на российском рынке

СУБД

ORACLE

Microsoft

- Отсутствие информации и практического подтверждения совместимости продуктов
- Время и ресурсы для подтверждения соответствия заявленной функциональности

- Проблема совместимости с продуктами из разных классов.
- Размывание понятия «лидер»: в каждом сегменте существуют десятки на первый взгляд равноценных продуктов

Независимость: варианты реализации



Покомпонентное замещение

- Время на изучение вариантов, тестирование и выбор
- Лавина взаимосвязанных проектов по внедрению
- Сложность синхронизации дорожных карт развития
- Рост сроков внедрения и рисков на стыках



Создание целевой доверенной ИТ-инфраструктуры

- Последовательный перевод систем на целевую доверенную ИТ-инфраструктуру
- Снижение нагрузки с текущей инфраструктуры и необходимости ее масштабирования
- Сокращение сроков внедрения и снижение рисков



Почему ПАК Скала р?



- Гарантированно совместимые компоненты
- Отказоустойчивость на уровне архитектуры
- Оптимизация производительности
- Ответственность одного производителя за функционал и показатели назначения
- Решенные вопросы интеграции, эксплуатации, мониторинга, обеспечения ИБ, резервного копирования
- Поддержка и сервис из одного окна
- Серийность и преемственность
- Управляемая дорожная карта развития



Конкурентные преимущества оптимизированных решений Скала^р



Производительность



чем решения, использующие сопоставимые аппаратные средства за счет оптимизации ввода-вывода и интерконнекта и за счет разгрузки ЦПУ



чем решения в виртуальной среде, использующие сопоставимые аппаратные средства за счет снижения латентности



для систем с большим количеством сессий за счет использования специализированных пулеров и балансировщиков

RPO/RTO



время выполнения резервного копирования и восстановления за счет специализированного встроенного модуля резервного копирования



время полного восстановления узла в случае отказа за счет использования встроенной системы развертывания и цифрового двойника системы

Доступность

Кратное сокращение инцидентов

связанных с ошибками эксплуатации и существенное увеличение доступности за счет использования специализированной системы управления ресурсами

ПАК — Машины Скала[^]р — преимущества перед самостоятельными проектами



Высокая отказоустойчивость

За счет специализированной модульной и кластерной архитектуры решений

Высокая производительность

Встречная оптимизация и устранение узких мест по всему стеку применимых технологий

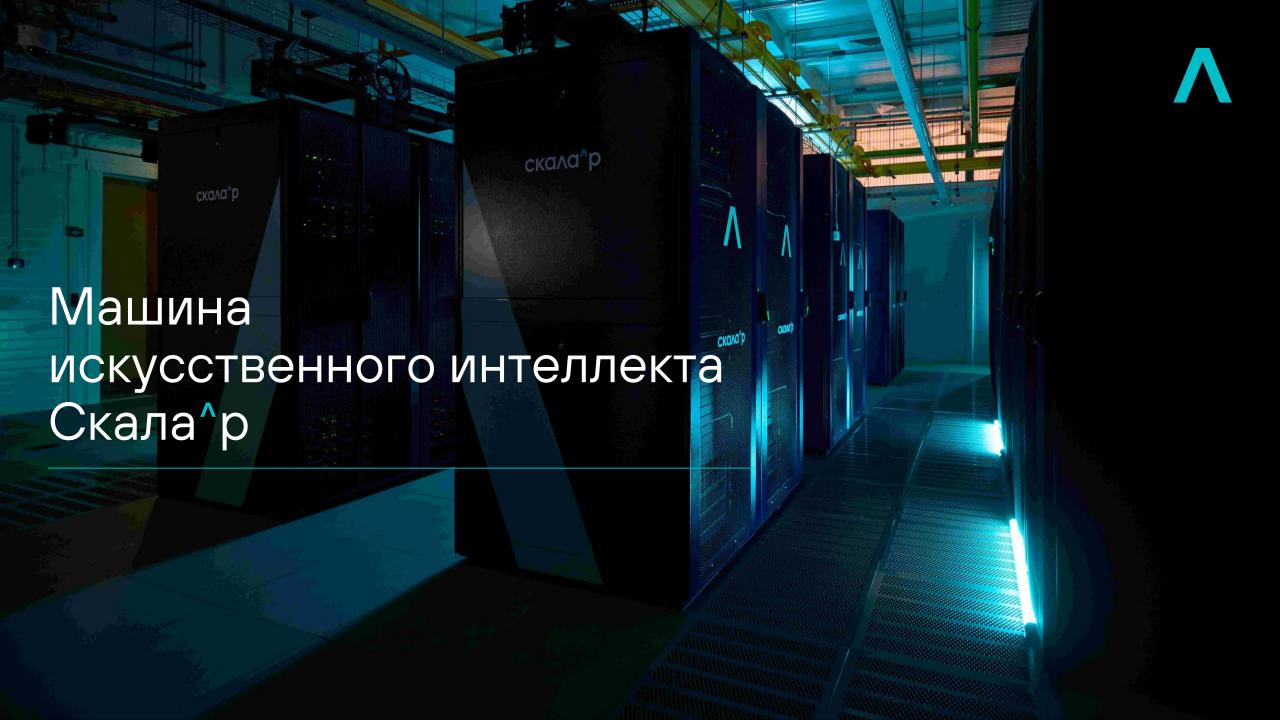
Премиальный сервис

Гарантированно работоспособное решение



Прямое взаимодействие с технологическими партнерами по развитию необходимого Заказчикам функционала

ПАК — Программно-аппаратный комплекс и модули платформы — включены в Единый реестр российской радиоэлектронной продукции и реестр Минцифры



Риски и сложности внедрения ИИ



Внедрение ИИ кардинально отличается от внедрения готовых программных продуктов (ERP, CRM, ITSM и т.п.) и от разработки ПО на заказ

- В процессе реализации ИИ-проекта доступен широкий спектр продуктов на рынке ИИ
- Список ИИ продуктов еженедельно меняется и появляются новые как решения так и технологии
- В процессе реализации ИИ-проекта подбираются и апробируются разные ИИ продукты

Для реализации ИИ инициатив нужны разные ИИ-специалисты (DataScience, ML-инженер, DevOps, Аналитик и т.п.)

- Специалистов очень мало на рынке и стоимость высокая
- Высокая конкуренция за ИИ-специалистов
- Долгий срок подготовки и взращивания специалистов даже при наличии наставников
- Зависимость ИИ-проектов от носителя знания или компетенции

ИИ инфраструктура дорогая

- Основная стоимость это графические ускорители и интерконнект
- Требуются специализированные аппаратные решения
- Потенциальные сложности и ограничения при масштабировании ИИ инфраструктуры

Комплексность ИИ решений усложняет организацию ИБ

- Во время интеграции в контур компании есть высокий риск получить дыру в безопасности
- Во время эксплуатации возникают дополнительные риски утечки данных

ПАК ИИ



- ПАК стандартизирует аппаратно-программные решения, гарантируя их совместимость, универсальность, сохранение уровня доверия к ИИ и возможность модернизации
- Готовый набор проверенных моделей и ИИ продуктов, входящих в состав Pandora
- Сокращение срока разработки и внедрения ИИ до 80%
- ПО автоматизации управления ЖЦ Машины ИИ (Pandora) позволяет развивать ИТ-специалистов, что снижает зависимость от рынка труда в сфере ИИ
- IT-специалисты -> MLDevOps -> ML инженеры -> DataScience
- Экономия на дальнейшей масштабируемости за счёт применения НРС технологий и устранении узких мест
- Отработанный план модернизации продукта от тестовой среды до ИИ-фабрики
- Обеспечена вариативность между NVIDIA, азиатскими вендорами и другими
- Оптимизация требуемых ресурсов от задачи к задачи
- Обеспечение безопасности ИИ сервисов за счёт подхода TRISM/secure by design и каталога доверенных контейнеров, готовых к запуску на on-premise инфраструктуре, в закрытом контуре
- ПАК может функционировать в изолированном контуре, предотвращая утечки данных
- Использование доверенных ИИ компонентов из его каталога и листа совместимости, минимизирует риски компрометации конечного продукта



Token

/TOPS

/FLOPS

Как итог улучшения показателя рубль за единицу вычисления

Импортозамещение: Соответствие требованиям регуляторов и уменьшение санкционных рисков

Машина искусственного интеллекта Скала р



ИИ решения



















Cotype

GigaChat

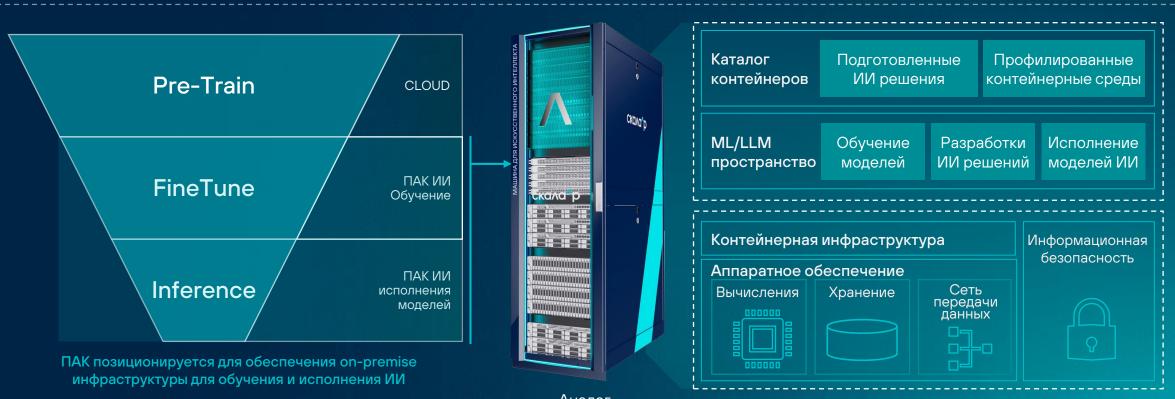
Llama

Сайбокс

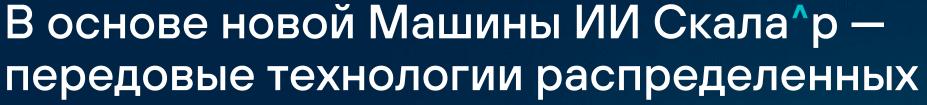
DeepSeek

YandexGPT 5

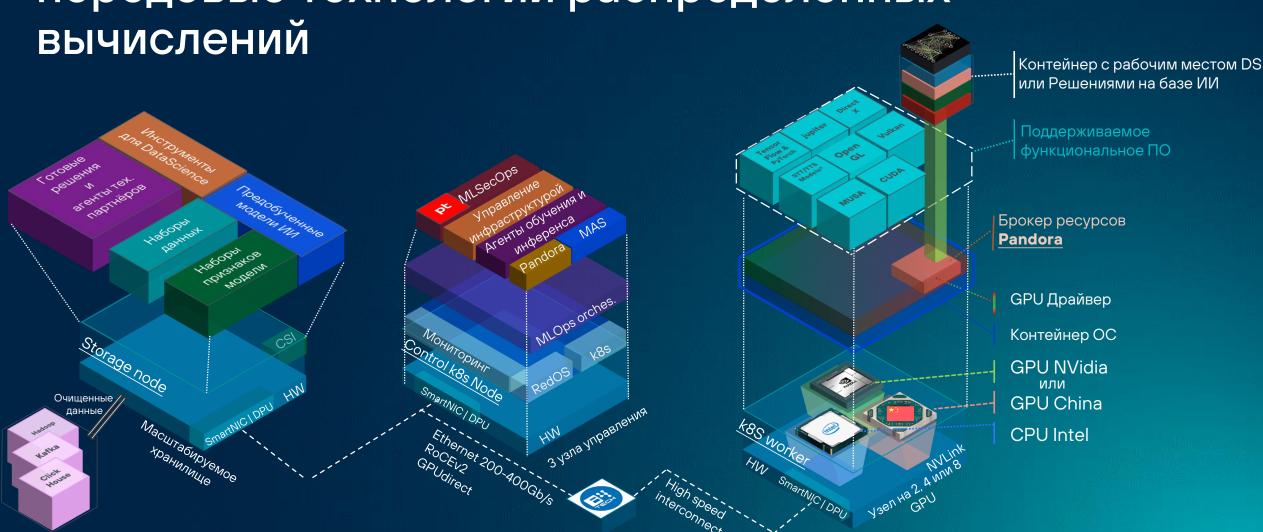
Смарт Платформа Другие ИИ решения



Аналог Huawei Atlas 900 Pod / NVIDIA DGX SuperPOD

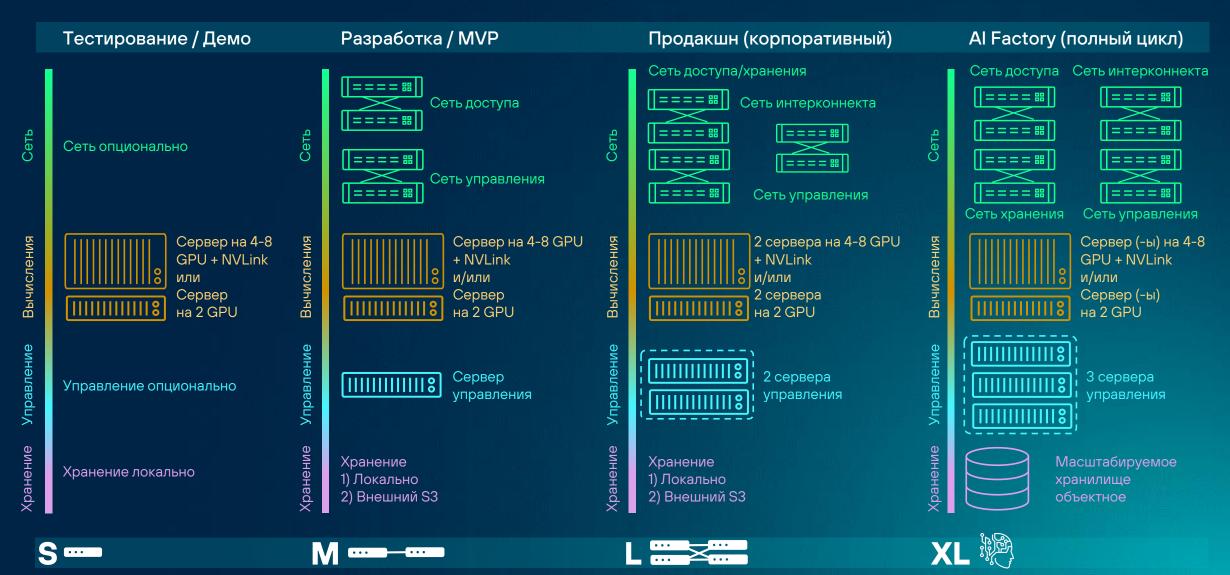






Размеры Машины Скала^р МИИ





Параметры продуктовой линейки



Функция / Характеристика	S (Small)	M (Medium)	L (Large)	XL (Extra Large)
Количество узлов	1 (Модуль инференса)	2+ (Минимальный кластер)	4+ (Отказоустойчивый кластер)	10+ (Масштабируемый кластер)
Поддержка инференса	Да (локальный)	Да (кластерный*)	Да (оптимизированный)	<mark>✓</mark> Да (масштабируемый)
Поддержка GPU/TPU	<u> </u>	✓ Да (несколько GPU)	✓ Да (кластер GPU)	Да (оптимизированные фермы)
Мониторинг и метрики	<u> </u>	✓ Prometheus + Grafana	Расширенная аналитика	✓ Al-аналитика + предсказания
Kubernetes (k8s) Management	Х Нет	Да (базовое управление)	Да (продвинутое управление)	Да (полный контроль + мониторинг)
Отказоустойчивость	Х Нет	<u>Л</u> Частично	Да (автовосстановление)	✓ Да (высокая доступность)
Создание ИИ-агентов	Х Нет	<u>∧</u> Базовые сценарии	✓ Да (сложные агенты)	<mark>✓</mark> Да (автономные агенты)
Масштабируемость	X Нет	<u>Л</u> Ручное масштабирование	✓ Да (автоматическое)	<mark>✓</mark> Да (гибкое + балансировка)
ИИ-ассистенты	X Нет	🗶 Нет	<u>Л</u> Простые интеграции	Да (многомодальные ассистенты)
Обучение моделей	X Нет	Х Нет	<u>М</u> Ограничено	У Да (распределённое обучение)
Целевой сценарий	Тестирование / Демо	Разработка / MVP	Продакшн (корпоративный)	Al Factory (полный цикл)

^{*} С добавлением модуля управления, можно кластеризировать модули инференса

Примеры исполняемых задач



с применением GPU платформы с NVLink



YandexGPT (LLM)



SciBox (Интсрумент)



RAG (Инструмент)



Cotype (LLM)



DeepSeek (LLM)



ValueAl (Инструмент)



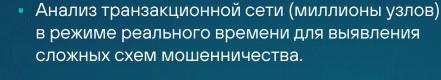
Llama (LLM)



GigaChat (LLM)



Другие ИИ решения



- Анализ кредитной истории + текстовых данных (договоры, переписка).
- Детекция аномалий в потоке транзакций (~100К TPS)
- Распознавание и верификация голоса в коллцентрах банка
- Парсинг договоров, регламентов, сканов документов и выявление рисков
- Обработка тысяч источников для прогноза волатильности рынка
- Автоматическое формирование отчетов по регуляторике на основе внутренних данных



Примеры исполняемых задач



с применением типовых серверов 2RU



YandexGPT (LLM)



SciBox (Интсрумент)



RAG (Инструмент)



Cotype (LLM)



DeepSeek (LLM)



ValueAl (Инструмент)



Llama (LLM)



GigaChat (LLM)



Другие ИИ решения

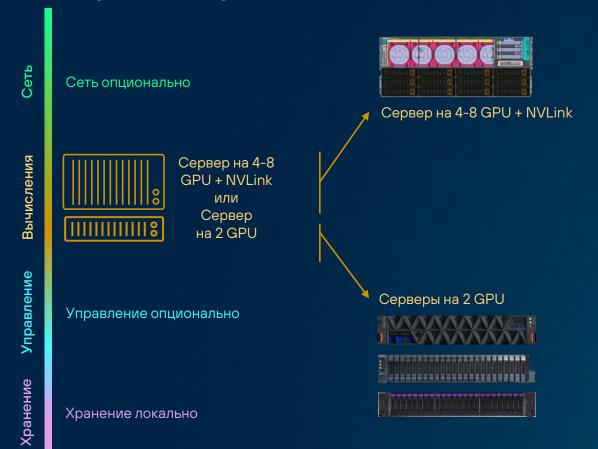


- Извлечение данных из документов
- Прогнозирование оттока клиентов
- Классический кредитный скоринг с фичами из транзакций.
- Выявление подозрительных транзакций (но не в реальном времени)
- Анализ клиентских профилей
- Ответы на типовые вопросы клиентов (без сложного RAG)
- Автоматическое категоризирование расходов.
 Разметка транзакций
- Проверка паспортов, договоров через компьютерное зрение

Инференс-узлы

s

тестирование/демо



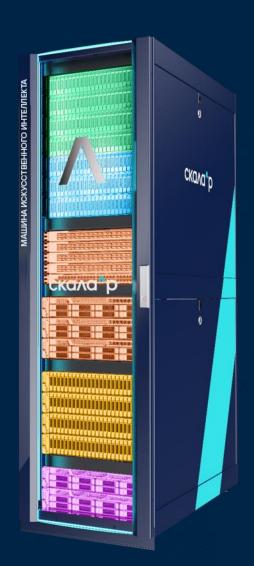
Юниты	4 RU		
Процессоры	2x Intel Xeon 4/5 Gen		
RAM	DDR5		
PCle слоты	Передняя панель: поддерживает максимум один слот PCle 5.0 Задняя панель: поддерживает максимум 10 слотов PCle 5.0 и 8 видеокарт двойной ширины		
Электропитание	~3,6КВт		

Юниты	2 RU
Процессоры	2x Intel Xeon 4/5 Gen
RAM	DDR5
PCIe слоты	До четырёх PCle 5.0 x16 и до трёх PCle 5.0 x8
Электропитание	~1,3КВт

Используемые GPU NVIDIA	NVIDIA H100	NVIDIA H200	NVIDIA A100	NVIDIA L40s	NVIDIA T4/L4
Используемые GPU Азия	16GB GDDR	32GB GDDR	48 GDDR	Аналог NVLink органичено, PCle 4.0 и PCle 5.0	

Машина Скала^р МИИ — Модули





Модуль полезной нагрузки Машины МИИ

- Bare metal узлы, выступающие в качестве Worker нод кластера Deckhouse Kubernetes Platform. Количество этих узлов можно варьировать от 3 до 16 (в некоторых случаях возможна конфигурация от 1 узла).
- Вычислительные мощности узла 64 физических ядра СРU, от 128ГБ до 4ТБ ОЗУ при оптимальной конфигурации памяти.
- От 1 до 8 GPU типа H100 в один узел
- Диски в этих узлах (от 4 штук в каждом узле в базовой конфигурации с возможностью масштабирования до 16 дисков на узел) можно
 использовать для организации хранения данных контейнеров, на сегодня это опции local path provisioner и SDS local volume в терминологии
 Deckhouse Kubernetes Platform.

Базовый модуль

Коммутационный модуль Машины МИИ

- Два коммутатора 100GbE или 400GbE на 32 порта(каждый) в отказоустойчивой конфигурации для сети интерконнекта Машины.
- Два коммутатора от 25GbE по 48 портов в отказоустойчивой конфигурации для организации доступа к сервисам Машины МИИ из сети заказчика.
- Два коммутатора от 25GbE на 48 портов(каждый) для организации сети хранения данных Машины.
- 🔻 Два коммутатора 1GbE на 48 портов (каждый) для организации управляющей сети (out-of-band управление и in-band управление).

Модуль управления Машиной МИИ

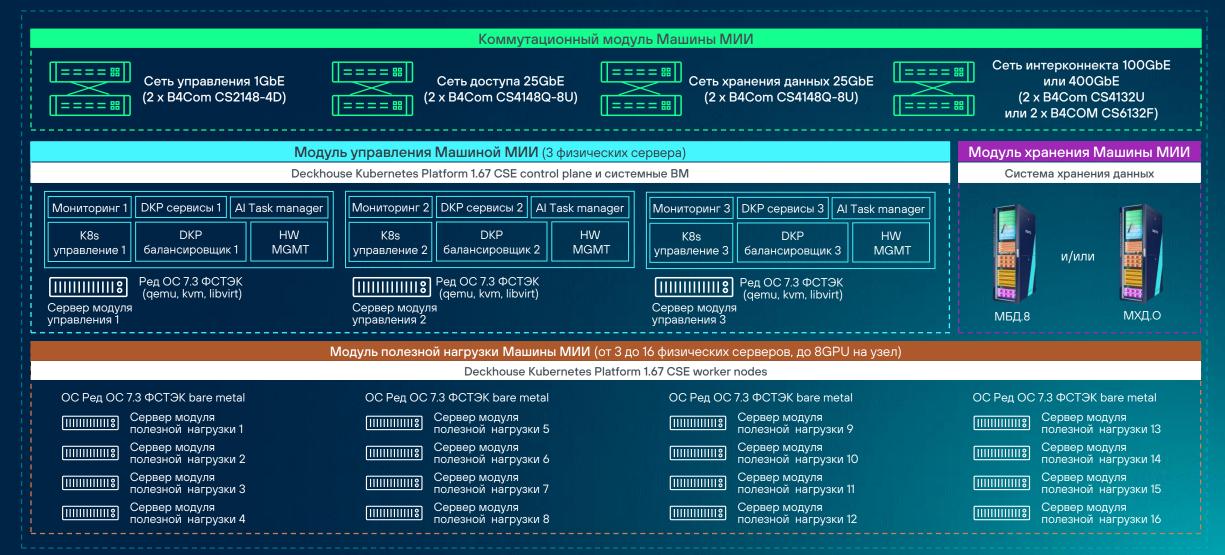
- Три сервера для размещения управляющих компонент Машины управляющих и служебных узлов Deckhouse Kubernetes Platform, сервисов Скала^Р.
- Диски в этих узлах (по 4 штуки в каждом узле в базовой конфигурации с возможностью масштабирования до 16 дисков на узел) можно
 использовать для организации различных вариантов хранилищ.

Модуль хранения Машины МИИ

- Подключаемый к кластеру DKP Машины контейнерной инфраструктуры посредством CSI драйвера.
- Поддержка распределенных вычислений
- Поддерживает многопоточную загрузку/выгрузку (например, через s5cmd, rclone)

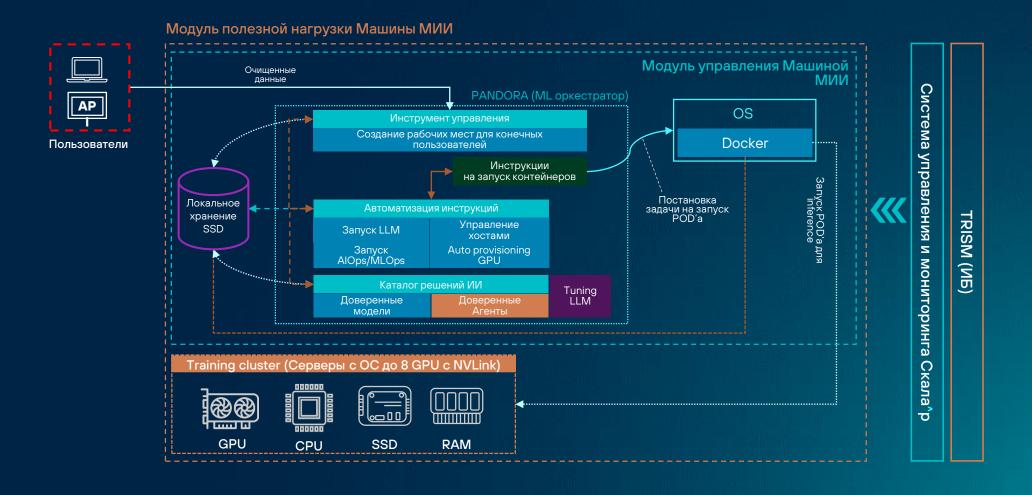
Машина Скала[^]р МИИ XL — компоненты





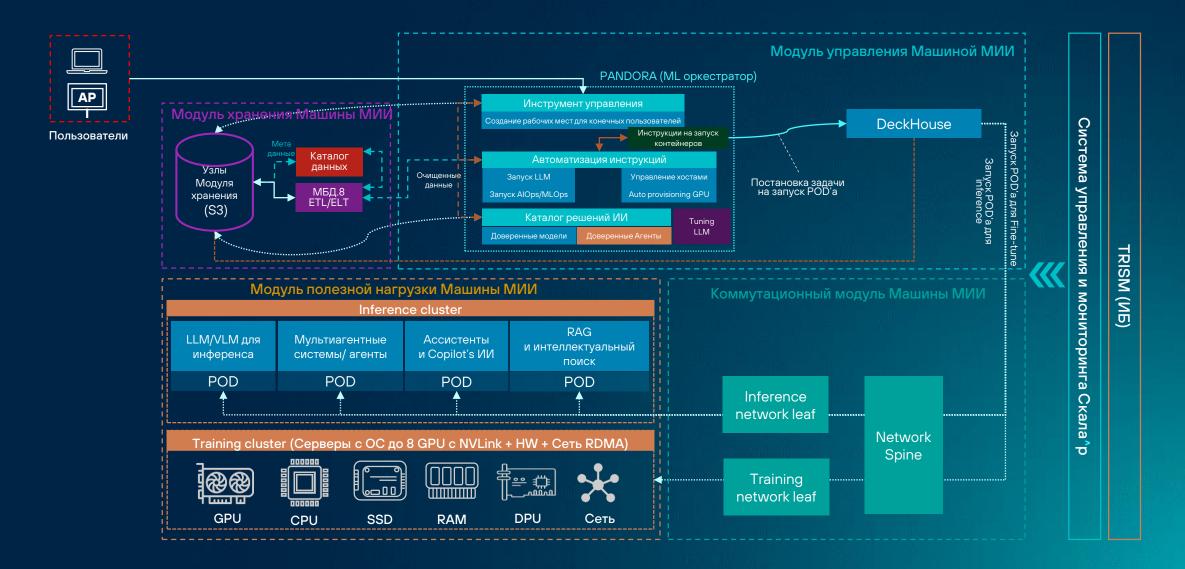
Машина Скала[^]р МИИ «S» — логическая схема





Машина Скала[^]р МИИ «XL» — логическая схема

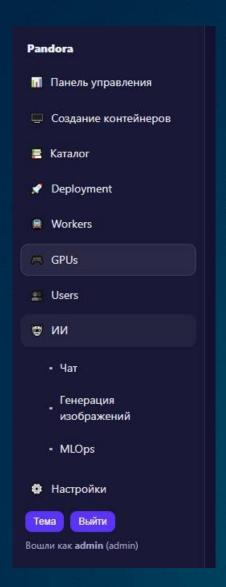


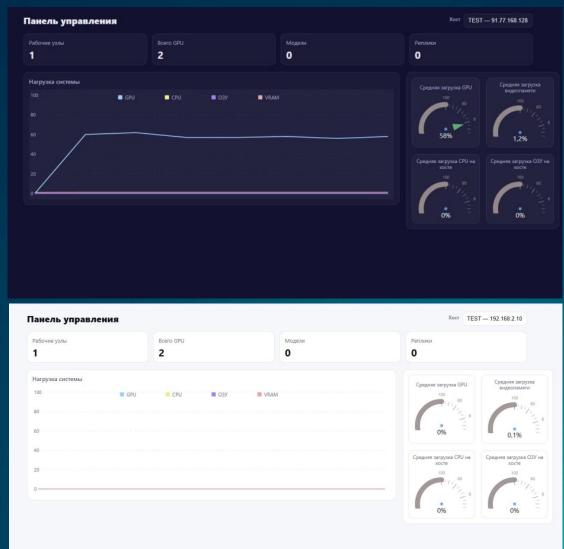


PANDORA (главная страница с мониторингом)



- Мониторинг ресурсов (GPU, CPU, RAM, vRAM)
- Можно увидеть количество хостов в кластере и кол-во GPU
- Кол-во запущенных ПОДов
- Кол-во используемых токенов (всего)
- Кол-во одновременных сессий в сторону модели/моделей (RPS)
- Кол-во ассистентов/агентов ИИ
- Тёмная и белая темы везде





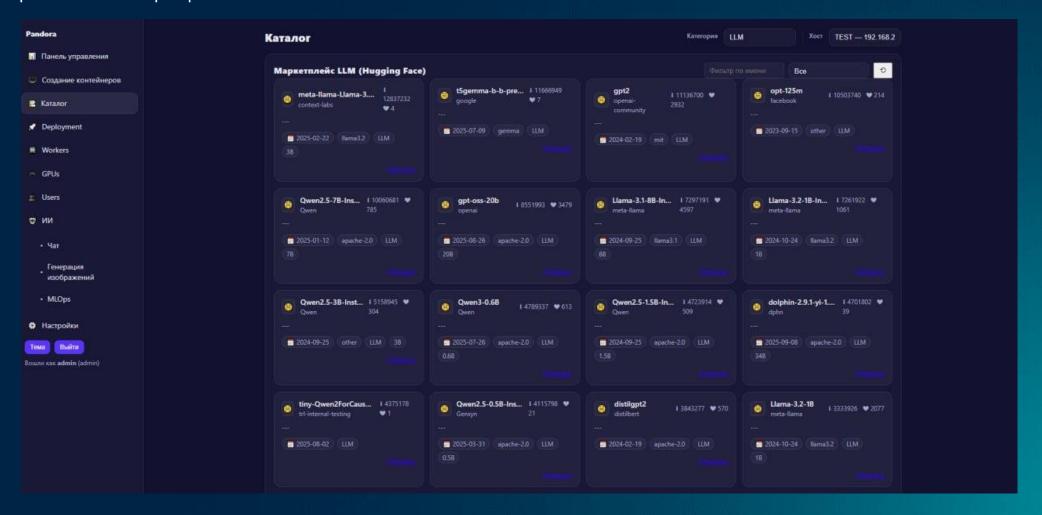


Создание рабочих мест с GPU ресурсами для конечных пользователей (разработчики, ML/DS-инженеры, бухгалтеры, юристы, сервисная поддержка и т.д.)

Pandora	Создать контейнер			
📊 Панель управления	Создать контейнер		Рабочий стол	
🥃 Создание контейнеров	Хост	GPU	Обновить	
■ Каталог	TEST — 91.77.168 128.2226	— без GPU —	redos8-mtt-latest-moore_threads-084639	
	СРИ (ядра)	RAM (ГБ)	redos8-mtt:latest @usertest, @port6082	
■ Workers	4	8	running	
	Образ redos8-mtt-kde:latest (9.14GB)	Пользователь — выберите —	Подключиться Перезапуск Остановить Удалить	
M GPUs	Протоколы	— выперите —	nvidia-redos7-latest-nvidia-083702	
4. Users	иротоколы ✓ noVNC ■ RDP ■ HTTP		nvidia-redos7:latest @usentest, @port6081	
Ф ии	Порты: noVNC 6080+, RDP 3390+, HTTP 8080+ (св	зободные на хосте)	running Подключиться Перезапуск Остановить Удалить	
• Настройки	Запустить контейнер			
Тема Выйти			nvidia-redos7-latest-nvidia-083428 nvidia-redos7:latest	
Вошли как admin (admin)			@usenalik @port6080 running	
			Подключиться Перезапуск Остановить Удалить	

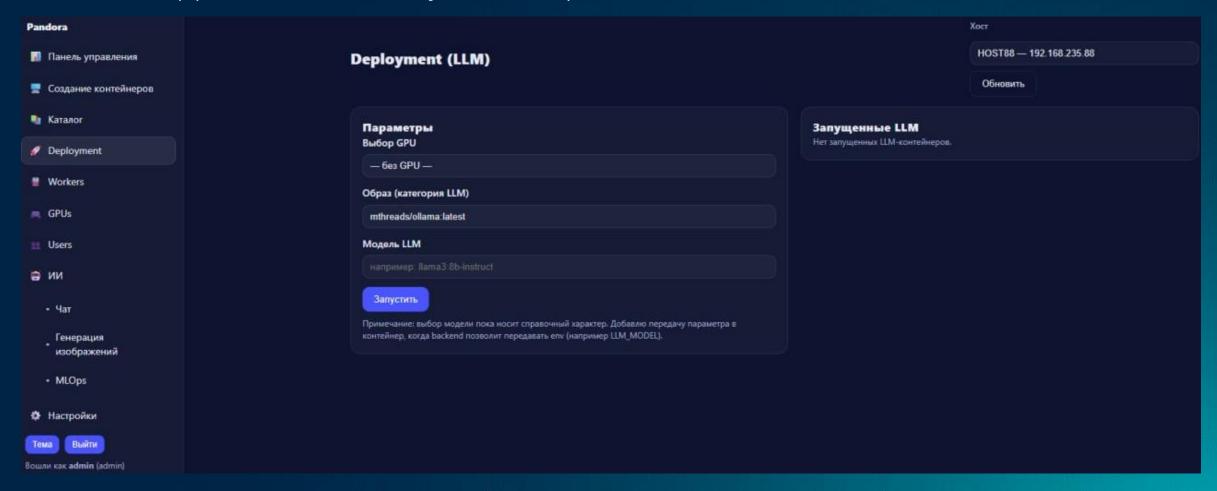


Каталог контейнеров и партнёрских решений, готовых к быстрому запуску на ПАК ИИ. Так же может пополняться решениями и разработками заказчика.



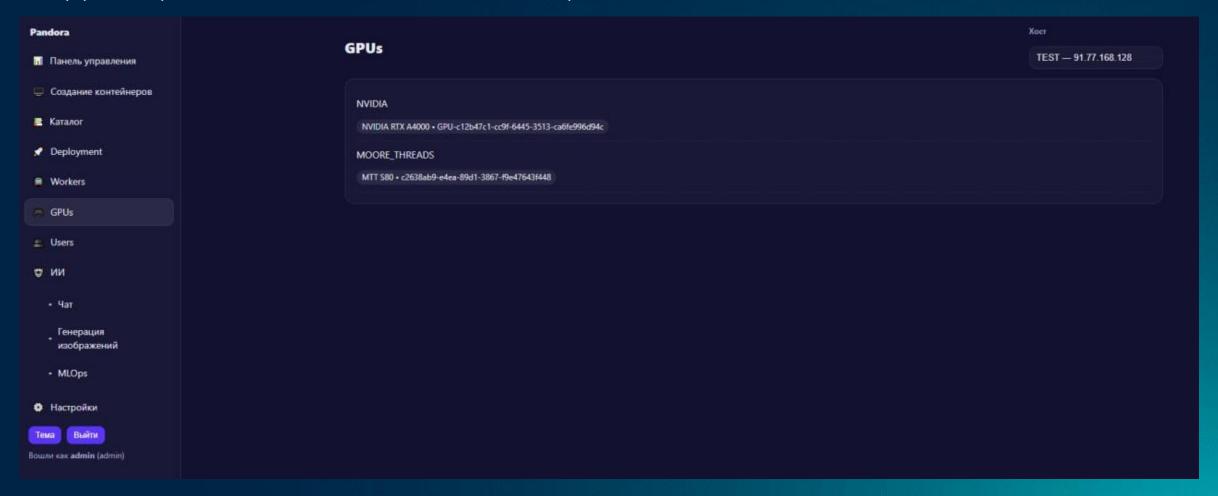


Отдельный интерфейс для создания и запуска контейнера с LLM моделью на ПАК ИИ



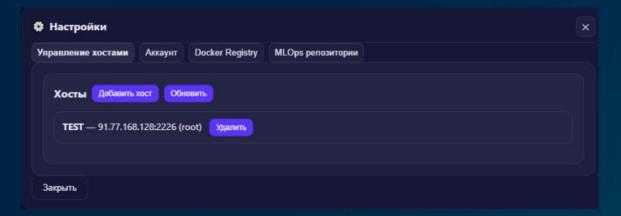


Интерфейс с представлением GPU на ПАК ИИ для хоста/блока или ПАК ИИ.





Модальное окно управления хостом (-ами)



Модальное окно добавления репозиториев docker

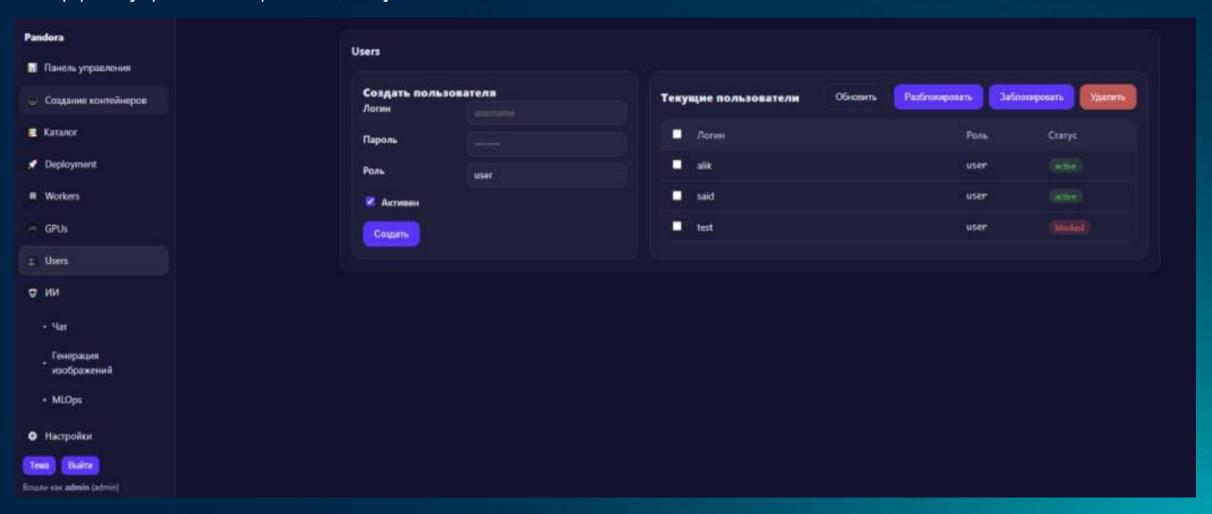


Модальное окно добавления репозиториев MLOps для LLM

Ф Настройки				×
Управление хостами	Аккаунт	Docker Registry	MLOps репозитории	
Название				
URL				
Токен (опц.)				
Добавить репозитори	й			
Репозиториев нет.				
Закрыть				

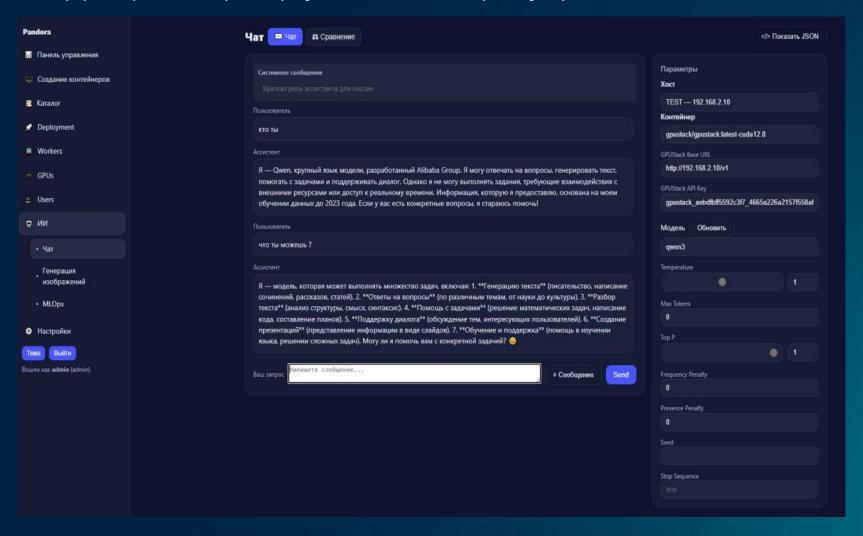


Интерфейс управления правами доступа





Интерфейс работы с развёрнутой LLM на хосте/модуле/ПАК ИИ , а так же Fine-tune модели



Метрики по ПАКу ИИ



Производительность

>= 6 Pflops

на один рабочий узел (TF32)

>= **400** Tflpos

на один рабочий узел (FP32)

Масштабирование сети с применением RDMA увеличивает производительность и уменьшает задержки более чем в 2 раза, что критично для моделей ИИ, работающих на распределённом кластере со сложной топологией

Максимальный размер LLM

до 188В параметров

в один рабочий узел без квантования

Использование NVLink для LLM увеличивает TPS примерно

в 2-5 раз

в зависимости от количества пользователей и типа запросов

Экосистема данных Lakehouse + ML/Al





Единое управление данными

объединяет структурированные и неструктурированные данные для обучения моделей ИИ, поддерживая разнообразные рабочие нагрузки (например, BI, ML, генеративный ИИ)

Расширенные возможности для ИИ

применение инновационных форматов с поддержкой транзакционности и версионирования, гарантии надёжности данных для ИИ

Масштабируемость и производительность

оптимизировано для крупномасштабного использования ИИ с инструментами, поддерживающими аналитику в реальном времени

Управление и безопасность

качество данных и соответствие требованиям для приложений ИИ

Интеграция генеративного ИИ

обеспечивает инновационные варианты использования, такие как агенты и системы рекомендаций

Преимущества Машины ИИ Скала^р





Платформенные решения позволяют сократить

- в 15 раз время подготовки среды разработки
- в 5 раз время работы дата-инженеров и дата-аналитиков*



Реализация каталога ИИ решений от валидированных партнеров на базе ПАК ИИ



Надежная мультивендорная Enterpriseинфраструктура с оптимальной конфигурацией и стабильным программно-аппаратным стеком на основе проведенных тестов и лучших практик.



Соответствие требованиям соблюдения принципов отказоустойчивости, масштабируемости на уровне архитектуры для использования в критичных и высоконагруженных корпоративных и государственных информационных системах



Исключение инцидентов на стыке технологий и высококвалифицированная поддержка Скала^р





Расширение возможностей как вертикального, так и горизонтального масштабирования



Предсказуемые характеристики, метрики функционирования платформенных решений



Управление жизненным циклом корпоративных ИИ решений



Поддержка регуляторных требований, отраслевых стандартов



Увеличение производительности

- в 3 раза при обучении ML моделей
- в 4 раза обученных ИИ моделей



Безопасное использование популярных языковых моделей LLM в закрытом контуре

^{*} Показатели могут варьироваться в зависимости от задачи

Примеры использования ИИ для корпоративных задач*



Совершенствование процессов технической поддержки продуктов компании IT_ONE

Автономная система для классификации, маршрутизации поступающих обращений клиентов по разным каналам связи на корректную линию технической поддержки.

Построена на основе обработки естественного языка с применением адаптированных языковых моделей LLM.

Повышение эффективности клиентского сервиса

Чат-бот технической поддержки клиентов для информирования, ответов на общие вопросы, уточнения дополнительной

информации.

Построен на основе технологии обработки естественного языка и дообученных языковых моделях LLM.

Совершенствование внутренних процессов по повседневной работе сотрудников

Расшифровка аудиозаписей встреч с суммаризацией итогов, определения решений и поручений по аудиозаписи: на основе обработки естественного языка, транскрибация, применение адаптированных

языковых моделей LLM.

Создание единого связанного пространства данных из разнородной информации документов ограниченного доступа, приходящих в ответ на запросы контролирующих органов государственной власти федерального уровня

Автономное (on-premise) ИИ-решение на основе LLM, в формате ПАК для автоматического извлечения данных из неструктурированных документов и автоматического формирования фабулы документа с гибкой настройкой правил извлечения данных.

Повышение эффективности разработки и тестирования программных продуктов

компании

Чат-ботов для разработчиков и тестировщиков, с поддержкой используемых языков программирования с учетом кодовой базы клиентских продуктов (ПО) во внутреннем контуре компании.

Создание изолированной ИТ-инфраструктуры для эксплуатации результатов инициатив ИИ.

Формирование у сотрудников компетенций, позволяющих использовать доверенные технологии ИИ

Средства обучения сотрудников промпт-инжинирингу и мотивации использования ИИ на основе.

Построены на больших фундаментальных языковых моделях (облачных) для выполнения текущих задач.

Повышение эффективности процессов управления проектами компании

Интеллектуальный помощник (чат-бот), повышающий эффективность повседневной работы руководителей проектов с внутренней документацией, базой знаний и регламентами компании, хранящимися в разнородных внутренних корпоративных сервисах компании.

Построен на основе адаптированных языковых моделей LLM, интеллектуального алгоритма для контекстного поиска, агрегации данных и предоставления структурированных ответов через интуитивный интерфейс чата.

Совершенствование процессов подбора сотрудников

на соответствие требованиям позиции (вакансии).

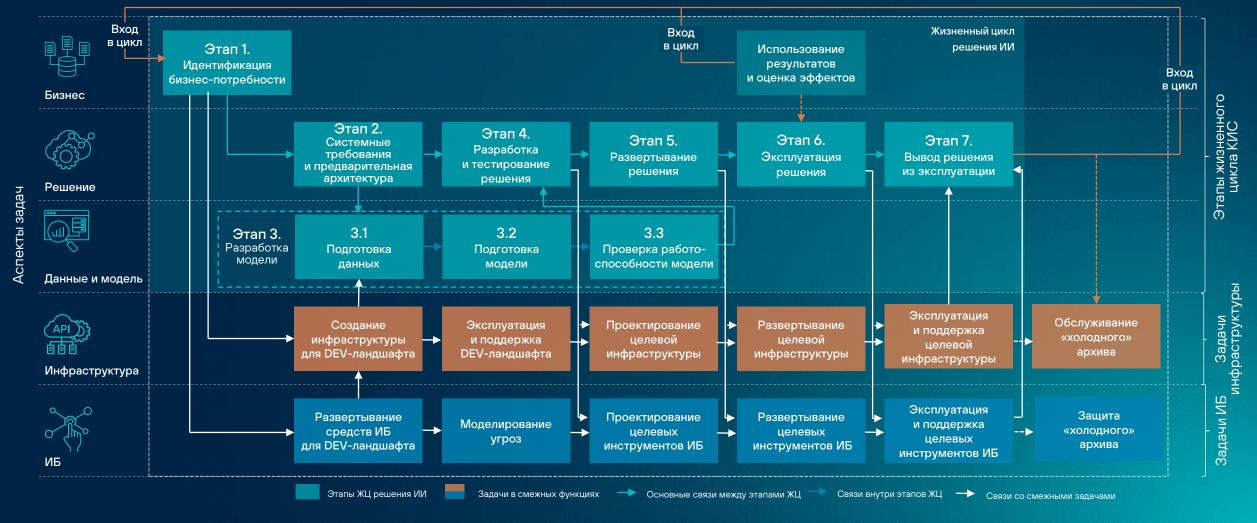
Построена на основе технологий NLP и применения адаптированных языковых моделей LLM.

^{*} Типовые задачи для инфраструктуры Машины ИИ Скала^р

Жизненный цикл КИС с ИИ

Общий взгляд на этапы и задачи

Современный жизненный цикл КИС с ИИ имеет специфические черты, связанные с работой с данными и моделями и тесную связь с задачами инфраструктуры и ИБ



История организационного и технологического развития



От импортозамещения Highload-стека к доверенной ИТ-инфраструктуре на ПАК

2025 2023 2014-2015 2016 2017 2019 2021 Замещение VMware и Citrix Динамическая инфраструктура Машины виртуализации Скала[^]р MB Virtuozzo + VDI Замещение Oracle Exadata Высокопроизводительные СУБД Машины баз данных Скала р МБД.П Postgres Professional Управление большими данными Замещение Teradata Машины больших данных Скала р МБД.8 Arenadata Масштабные S3-хранилища Замещение Scality, EMC/Netapp Машины хранения данных Скала р МХД.О Virtuozzo SDS + S3 **Модульная платформа Скала^р** 22 мая 2015 Приоритетное направления развития первое публичное представление платформы Скала^р Инфраструктура для ИИ Машины искусственного интелекта Скала р

