



Сложность vs Доступность:

Как облако демократизирует искусственный интеллект



Михаил Соколов

Мировой рынок расходы на ИИ

млрд\$

1479

2023

988

2025

2026



Развитие инфраструктуры

4

Интеграция в потребительскую электронику



AlaaS — ИИ как услуга

интересно

×10

0,3 Вт энергоемкость стандартного

поискового запроса Google

VS

3,0 BT

энергоемкость поискового запроса ChatGPT

7

Стремительный рост ИИ → драйвер рынка видеокарт



2025

18,2% CAGR до 2032 Операторы дата-центров наращивают мощности

- У Обучение больших языковых моделей (LLM)
- **4** Инференс

Сценарии использования



Персонализированое клиентское обслуживание

Запуск сложных Аl-продуктов и исследований

Обработка и извлечение данных из документов

Ускорение разработки и аналитики разработки и разработки и аналитики и аналит

Сквозная автоматизация бизнес-процессов

Прогнозирование и предиктивная аналитика

7

Проблемы внедрения инфраструктуры с ИИ



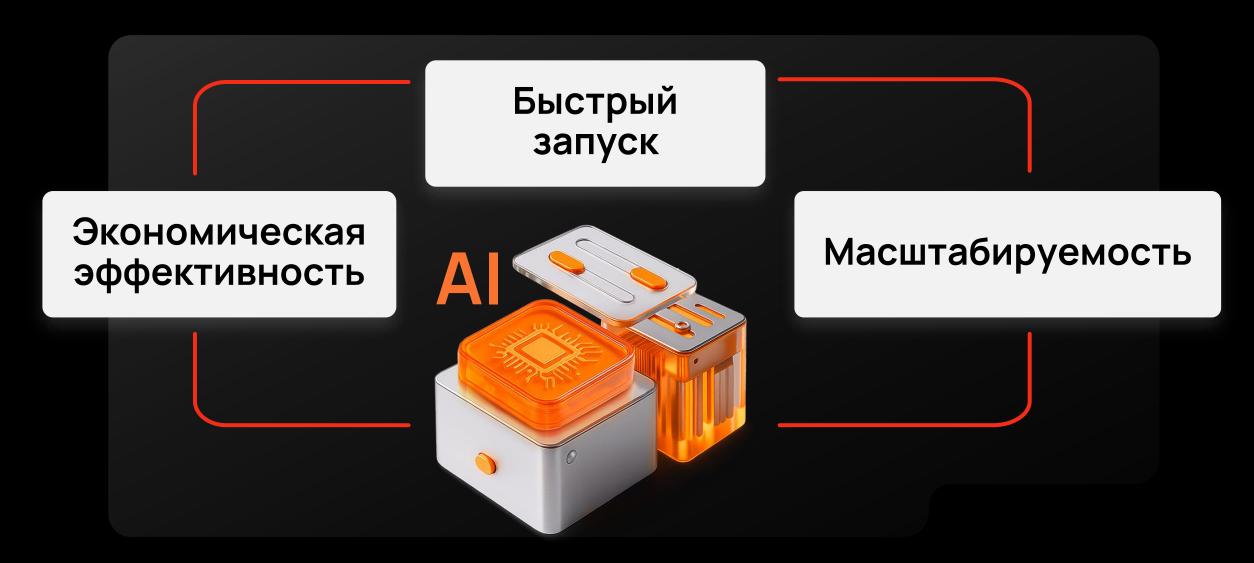
Очень дорогое «железо» — необходимы значительные вложения в инфраструктуру (On-premise)



Для оптимизации языковых моделей необходимы дополнительные капитальные вложения

Преимущества облачных сервисов





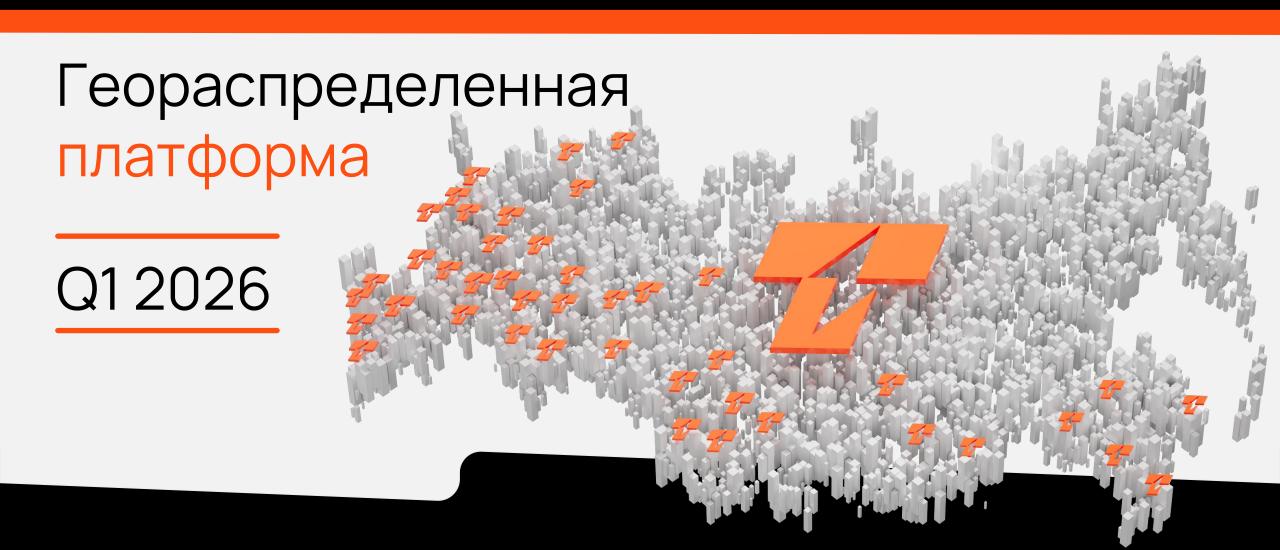


ИИ рядом

Вызовы географической распределенности



качество доступа



Лнициатива

Создание инфраструктурных сервисов для продуктов AI/ML



→ Потребность рынка в вычислениях на GPU

82

клиентских запроса в '22-'23 на 1000+ GPU разных классов

У большинства облачных → игроков продукты с GPU уже запущены

Предпосылки

2024

- → Предоставление вычислительных мощностей на базе GPU из облака с возможностью выбора оптимальной конфигурации: ML, 3D, VDI
- → Создание MVP сервисов AI/ML

Анализ опыта запуска GPU



Клиентам нужны не «сырые» GPU, а виртуализированные и преднастроенные для работы среды

Выводы

2025 ->

- Созданную инфраструктуру с GPU масштабируем
- Запускаем новый сервис «Виртуальные машины с GPU»
- Разворачиваем высокоуровневые сервисы для ML/AI

Сервисы

Сервисы для Al



Продукт	Серверы BareMetal	Виртуальная инфраструктура	Inference	Foundation models	Нейрошлюз
Формат	Выделенный физический сервер с GPU	ВМ или контейнер с GPU	Виртуальный инстанс	ML модель	Чат с ИИ-моделями для бизнес- пользователей
Тариф	За конфигурацию сервера	За виртуальные ресурсы	За инстансы	За токены	За пользователя
Доступно	от 150 340 ₽ / мес	от 11 900 ₽ / мес	Скоро	Скоро	Скоро

Производительность GPU

Гибкое управление конфигурацией

Минимальные накладные расходы

Изоляция от «шума»

Минимальная стоимость **BareMetal**











Виртуальная машина











Linux контейнеры











Выделенные серверы BareMetal



Без виртуализации для ресурсоемких задач, включая серверы с графическими картами для нейросетей, работы с графикой и машинного обучения

- Современные конфигурации серверов
- Инструменты для удаленного управления сервером
- Аппаратное резервирование компонентов сервера
- Поддержка NVLink
- _____

- Широкий выбор операционных систем
- Единый портал управления облаком и серверами BareMetal
- Возможна доукомплектация базовой конфигурации
- Нет переподписки по PCIe
- Доступна поддержка функционала виртуальных дисков

- Сетевая связанность с облачной инфраструктурой
- В составе сервера могут быть предоставлены GPU-карты

Выделенные серверы BareMetal



	Сервер Medium	Сервер Medium Plus	Сервер Medium Pro	Сервер Medium Ultra	Сервер Large GPU
CPU	AMD EPYC 7313 32 ядра, 3.0 GHz	Intel Xeon Gold 6240 36 ядер, 2.6 GHz	AMD EPYC 7443 48 ядер, 2.85 GHz	Intel Xeon Gold 6248R 48 ядер, 3.0 GHz	AMD EPYC 7543 64 ядра, 2.8 GHz
RAM	от 256 ГБ	от 512 ГБ	от 512 ГБ	от 1024 ГБ	от 512 ГБ
Хранение	SAS / SATA / NVMe SSD	SATA / SAS / SSD	SAS / SATA / NVMe SSD	SATA / SAS / SSD	SAS / SATA / NVMe SSD
Кол-во GPU	До 1-й на сервер	До 2-х на сервер	До 2-х на сервер	До 2-х на сервер	До 8-ми на сервер
Поддерживаемые GPU	L4, A6000	L4, A6000, A100	L4, L40S, A6000, A100, H100	L4, A6000, A100	L4, L40S, A6000, A100, H100
+Комплектация RAM, HDD, SSD	Да	Да	Да	Да	Да

Виртуальная инфраструктура с GPU



Позволяет делиться ресурсами физического GPU между несколькими BM или контейнерами, обеспечивая гибкость и эффективность без владения дорогостоящим оборудованием

- Выделенный профиль GPU на каждую BM/контейнер
- Отдельный образ для Data Science с популярными инструментами и фреймворками
- Масштабирование в пару кликов
- Инструменты для удаленного управления ВМ

- Подготовленные образы с оптимизацией под GPU
- Возможность заказаВМ и контейнерас несколькими GPU
- Современное железо с DDR5 и PCIe Gen5
- Единый портал управления

- TensorFlow, PyTorch и другие инструменты для вашей работы
- Гарантированная производительность без влияния «шумных соседей»

Конфигурации виртуальных машин c GPU



	Профиль	Профиль 2	Профиль 3	Профиль 4	Профиль 5	Профиль 6	Профиль 7
GPU	L4	L4	L4	L4	L40S	L40S	L40S
Профиль GPU	3 ГБ	6 ГБ	12 ГБ	24 ГБ	12 ГБ	24 ГБ	48 ГБ
vCPU, ядра	4	8	16	32	8	16	32
RAM	16 ГБ	32 ГБ	64 ГБ	128 ГБ	48 ГБ	96 ГБ	192 ГБ
SSD	60 ГБ	60 ГБ	60 ГБ	60 ГБ	60 ГБ	60 ГБ	60 ГБ

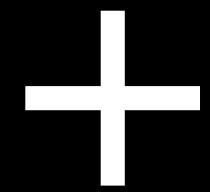
Конфигурации Linux-контейнеров с GPU



	Профиль 1	Профиль 2	Профиль 3	Профиль 4	Профиль 5
GPU	H200	H200	H200	H200	H200
Профиль MIG	1g.35gb	2g.35gb	3g.71gb	4g.71gb	7g.141gb
vCPU	14	14	28	28	56
RAM	64 ГБ	64 ГБ	128 ГБ	128 ГБ	256 ГБ
SSD	60 ГБ				

Новые Сервисы

2026







облако высоких скоростей

облачных сервисов

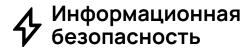




Платформенные сервисы



Виртуальная инфраструктура и сеть



Профессиональные сервисы















































Public preview — Q1 2026

скоро

Нейрошлюз

Private Preview — уже сейчас



Онлайн-платформа, где можно работать с множеством средств и сервисов AI в едином рабочем пространстве



Public preview — Q1 2026

скоро

Inference platform



Позволяет создавать инстансы с необходимыми библиотеками и инструментами для гибкого управления ML-моделями



Public preview — Q2 2026

скоро

Foundation models



Сервис обеспечивает доступ к актуальным ML-моделям, расположенным на доверенной инфраструктуре провайдера, с тарификацией по токенам

БЫСТРЫЙ CTAPT C TYPБО



Соколов Михаил



Turbo Cloud.ru