

Selectel

# Инфраструктура для ЯЗЫКОВЫХ МОДЕЛЕЙ

Размер имеет значение



Владислав Кирпинский

Директор по облачной интеграции

# Selectel сегодня

40+

продуктовых  
решений

6

собственных дата-  
центров

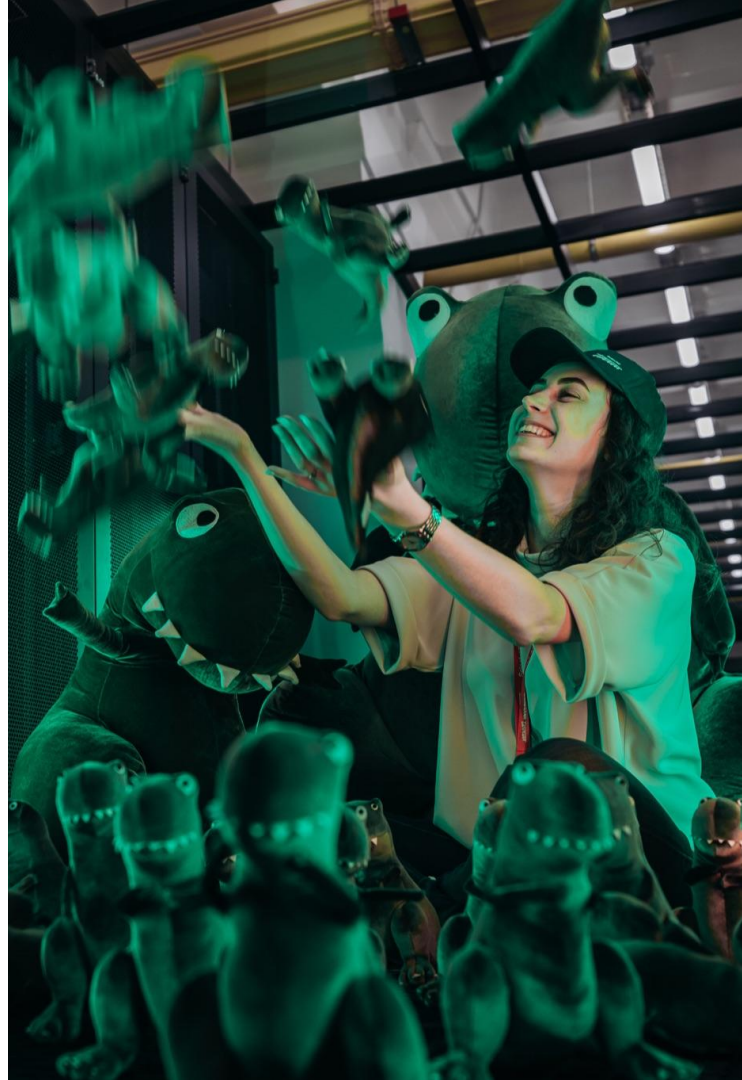
1 000+

сотрудников

24 000+





клиентов

С 2008 года помогаем компаниям решать бизнес-задачи, создавая надежную IT-инфраструктуру для проектов любой сложности.







# Более 40 готовых инфраструктурных решений

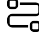
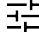
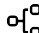

## Серверы и вычисления

-  Выделенные серверы готовых и произвольных конфигураций
-  Облачные серверы с моментальным запуском
-  Серверы с GPU
-  Экспериментальное железо




## Облачная платформа

-  Облачные базы данных
-  Объектное хранилище
-  Файловое хранилище
-  Managed Kubernetes и Container Registry




## Организация сети

-  Сеть доставки контента (CDN)
-  Балансировщики нагрузки
-  Selectel Connect и сети L3 VPN
-  Облачный DNS




## Облако на базе VMware

-  Публичное облако
-  Частное облако
-  Удаленные рабочие столы (VDI)

## ML и обработка данных

-  ML-платформа
-  Платформа обработки данных
-  Data Science & Analytical Virtual Machine

## Безопасность

-  Аттестованный сегмент ЦОД
-  Соответствие 152-ФЗ
-  Защита от DDoS и WAF

Единая панель управления  
и система биллинга

Документация к API  
и база знаний

Система управления  
ролями (IAM)

Базовая защита от DDoS  
по умолчанию

Техническая  
поддержка 24/7

# Содержание



---

Сколько нужно GPU, чтобы внедрить NLP?

---



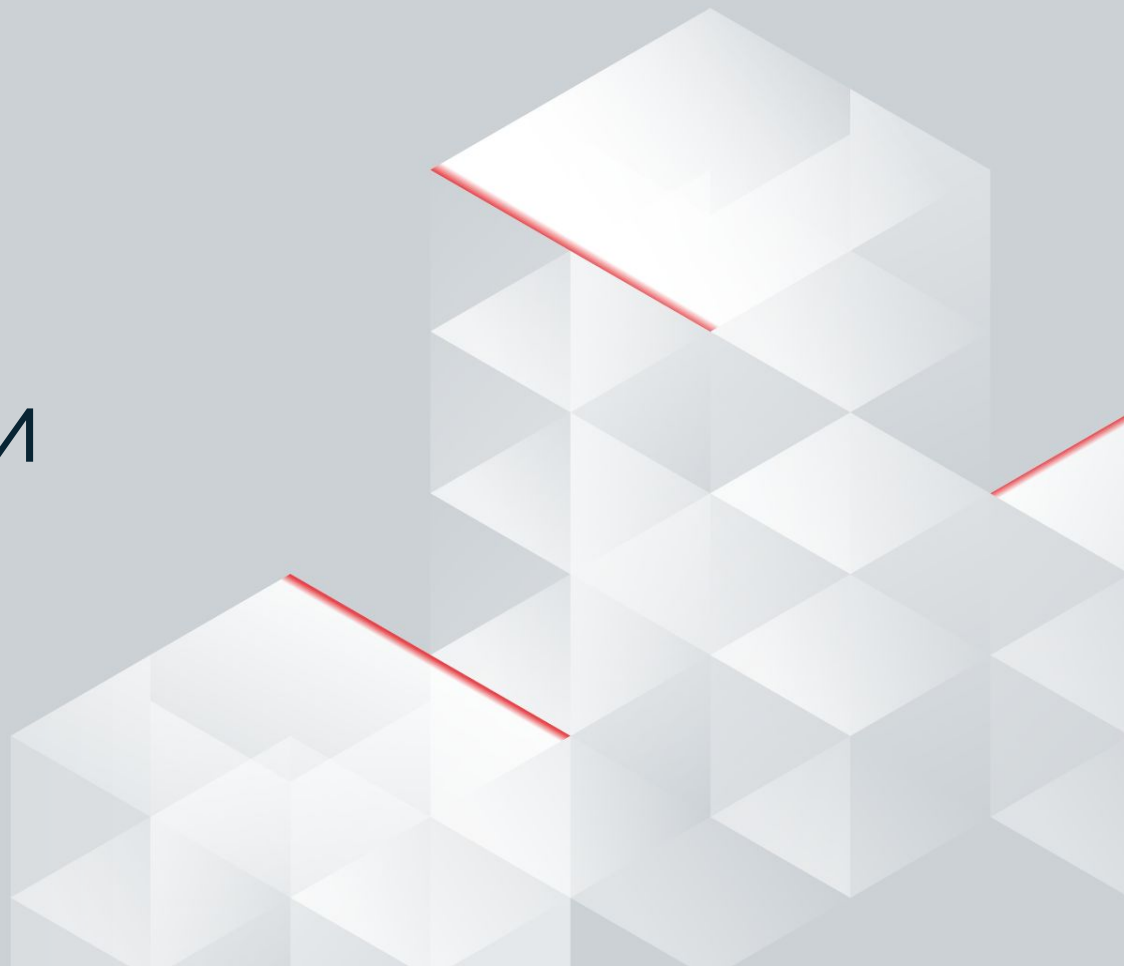
Инфраструктура для голосовых роботов: можно и на CPU

---



Inference-платформа: глобальный подход к экономии или лишние затраты?

Немного теории



# Эй, Кембриджский словарь, что такое...

## Inference

- A guess that you make or an opinion that you form based on the information that you have.

## Training

- The process of learning the skills you need to do a particular job or activity.

# Inference vs Training

IoT Data Input to ML Models



Raw IoT Data From IoT Endpoints

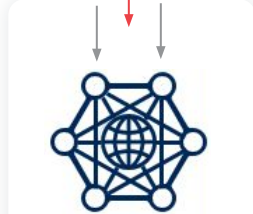
On-Premises or Cloud-Hosted

## Training

Learning a New Capability From Existing Data

Deep-Learning Framework

## Training Dataset



"cat"

"dog"



Trained Model



## Inference

Applying This Capability to New Data

Edge Device, On-Premises or Cloud-Hosted

## New Data



App or Service Featuring Capability

Logical Data Warehouse



The background features a series of overlapping circles in shades of gray and white, creating a sense of depth and movement. A prominent red arc curves across the right side of the image, adding a vibrant touch to the otherwise muted palette.

Что такое NLP в 2024?





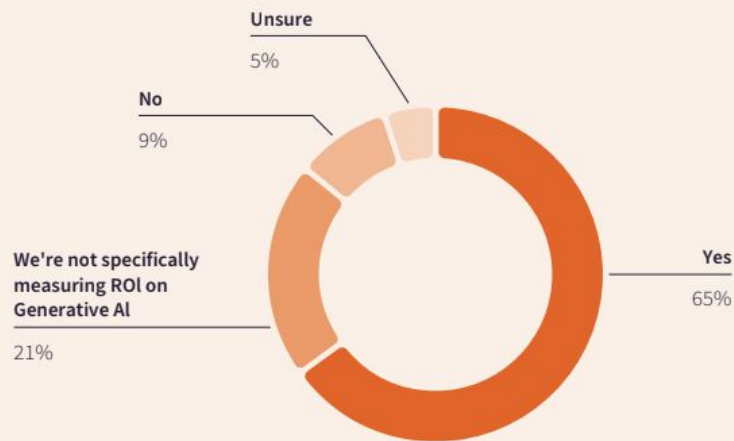
# Large Language Models

**ПОСЛЕ ДАННОГО СЛАЙДА  
ПРЕЗЕНТАЦИЯ В ПРОЦЕССЕ  
РЕДАКТИРОВАНИЯ**

# Автоматизация нагруженных процессов = экономия

## Experiencing Positive ROI From GenAI Use Cases

Q20new: Are you seeing positive ROI from GenAI use cases in production?



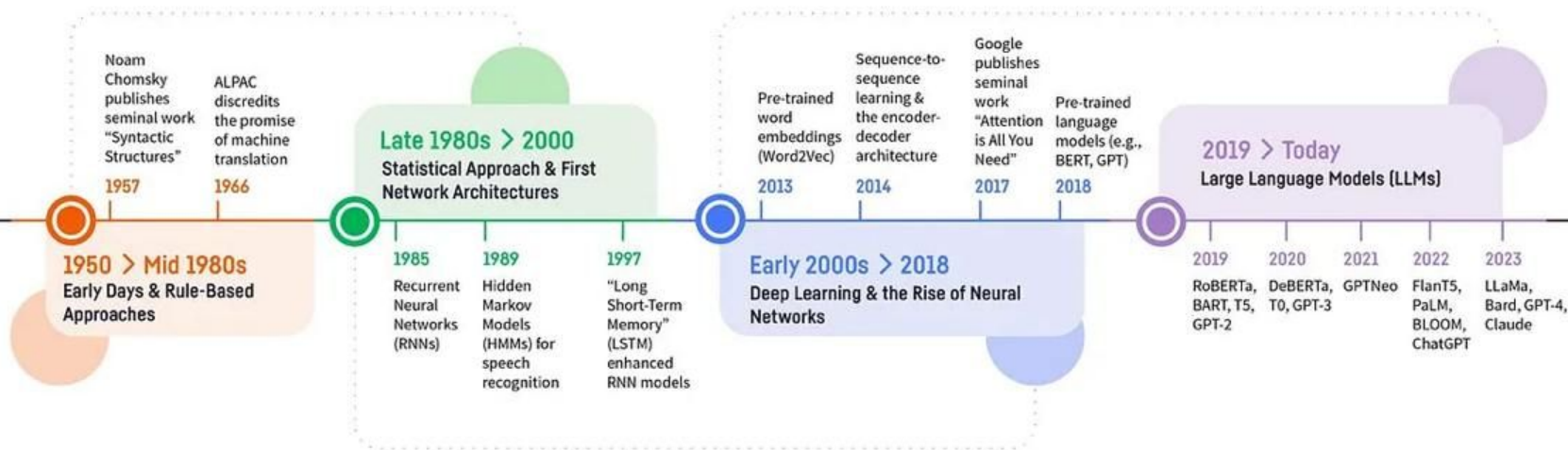
## Type of GenAI Models / LLMs Currently Using

Q19new: Are you currently using any of the following GenAI models / LLMs?



# Но до них тоже все было неплохо...

## The History of NLP





NLP-инфраструктура

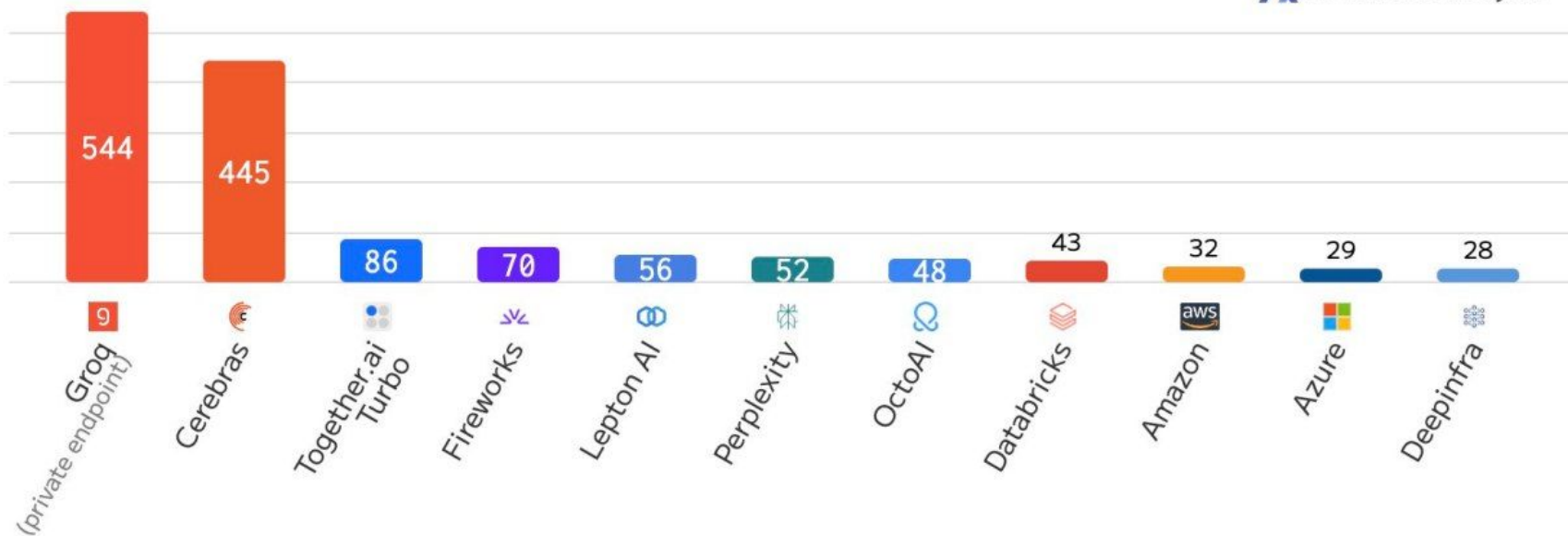
# Inference: битва за токены

## Output Speed: Llama 3.1 70B

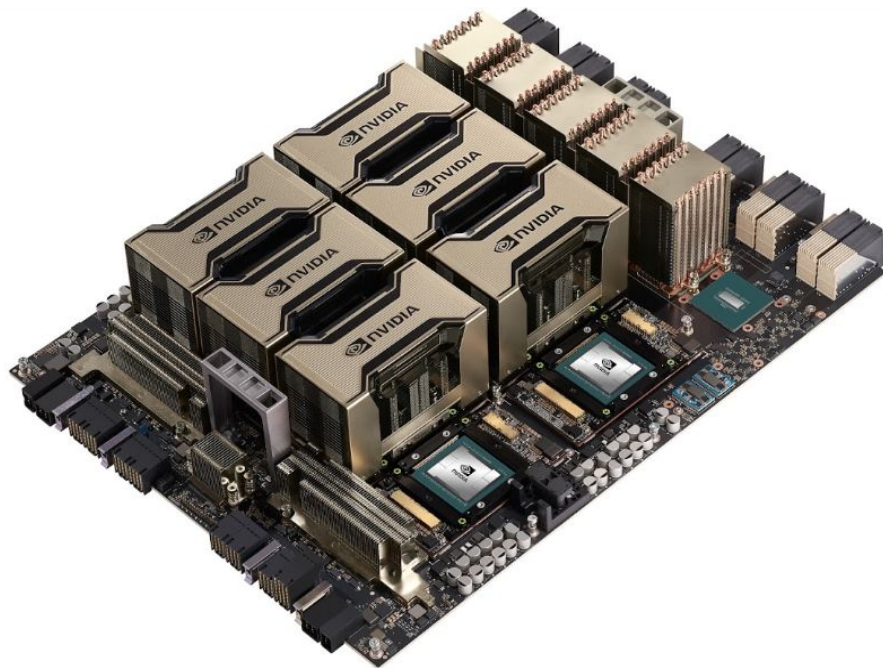
Output Tokens per Second; Higher is better; 1,000 Input Tokens

Groq (private endpoint): 8k context, N=100 requests; Tested: 9 September 2024

Artificial Analysis



ChatGPT: 3 500+ серверов NVIDIA HGX A100



## Занимательный факт из реальной жизни

2 голосовых робота могут одновременно работать на **одном** ядре **CPU** с задержкой 1-2 секунды.

Источник



Training: битва за обьем памяти

Llama 3.1 405B

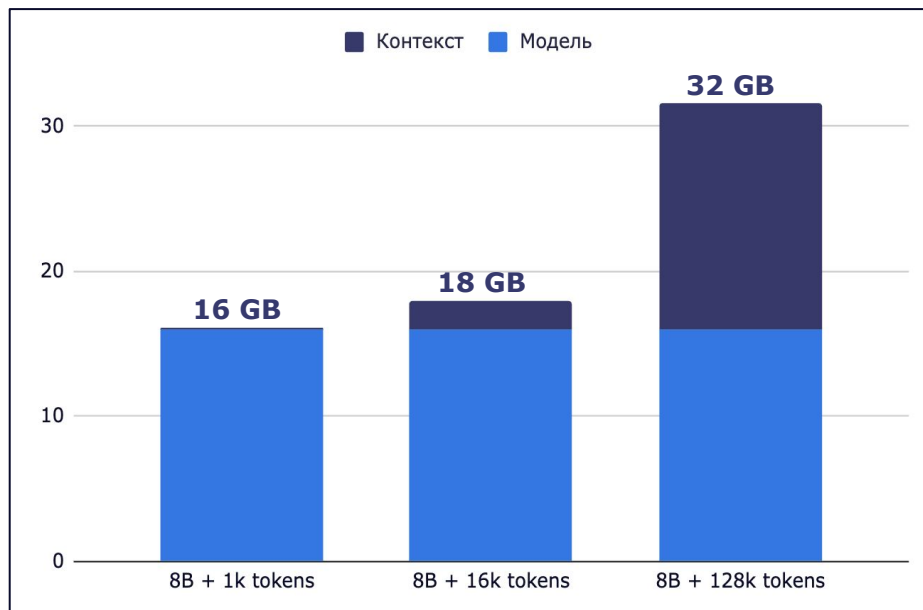


16,384  
NVIDIA H100  
80GB

Источник

# Еще один интересный факт из реальной жизни

## Требования к памяти LLaма 3.1 8B + контекст

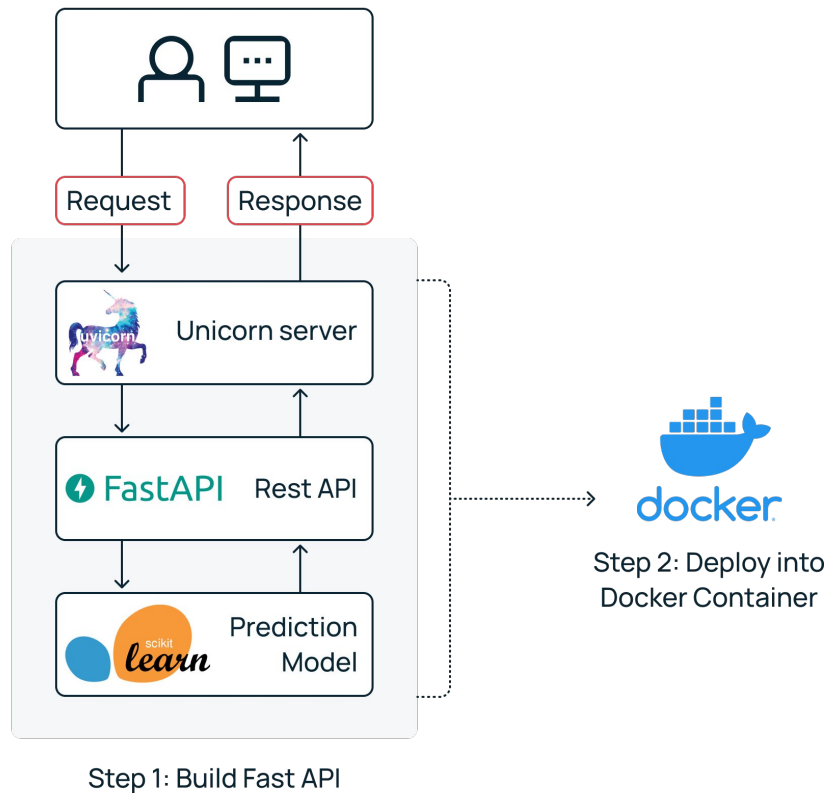


Для использования LLaма 8B достаточно одной NVIDIA Tesla A100 40Gb

Источник

Inference-платформа

# «Классический» путь ML-модели до Inference



# Спектры задач

Первый деплой

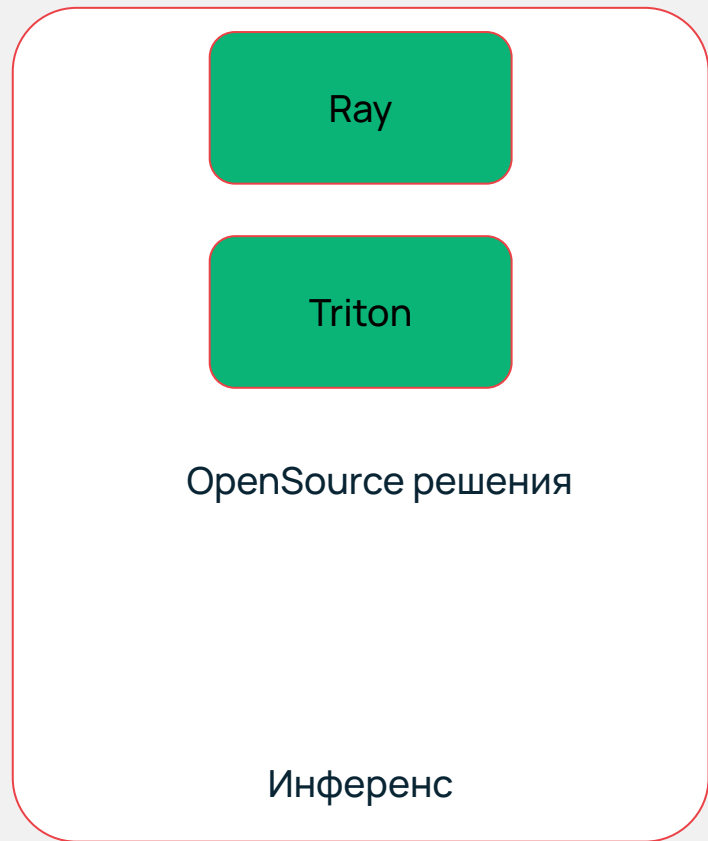
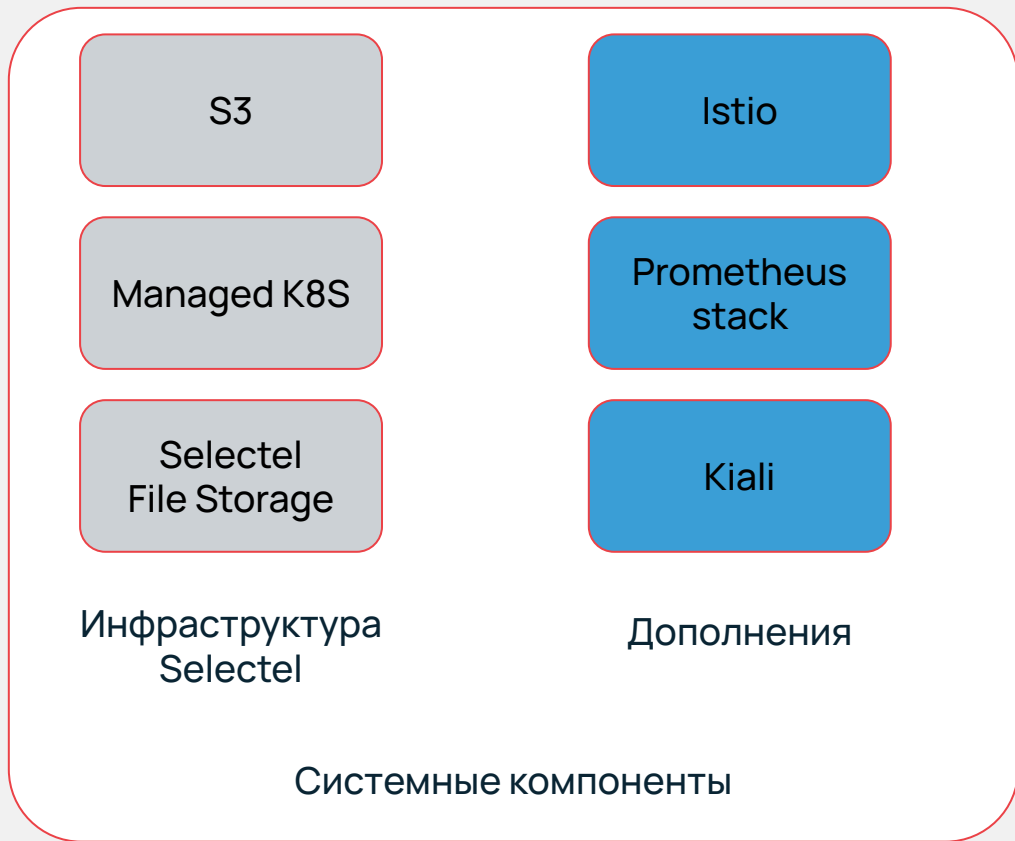
Деплой новой версии

Рост нагрузки

Мониторинг

Инференс-граф (LLM)

# Внутри inference-платформы



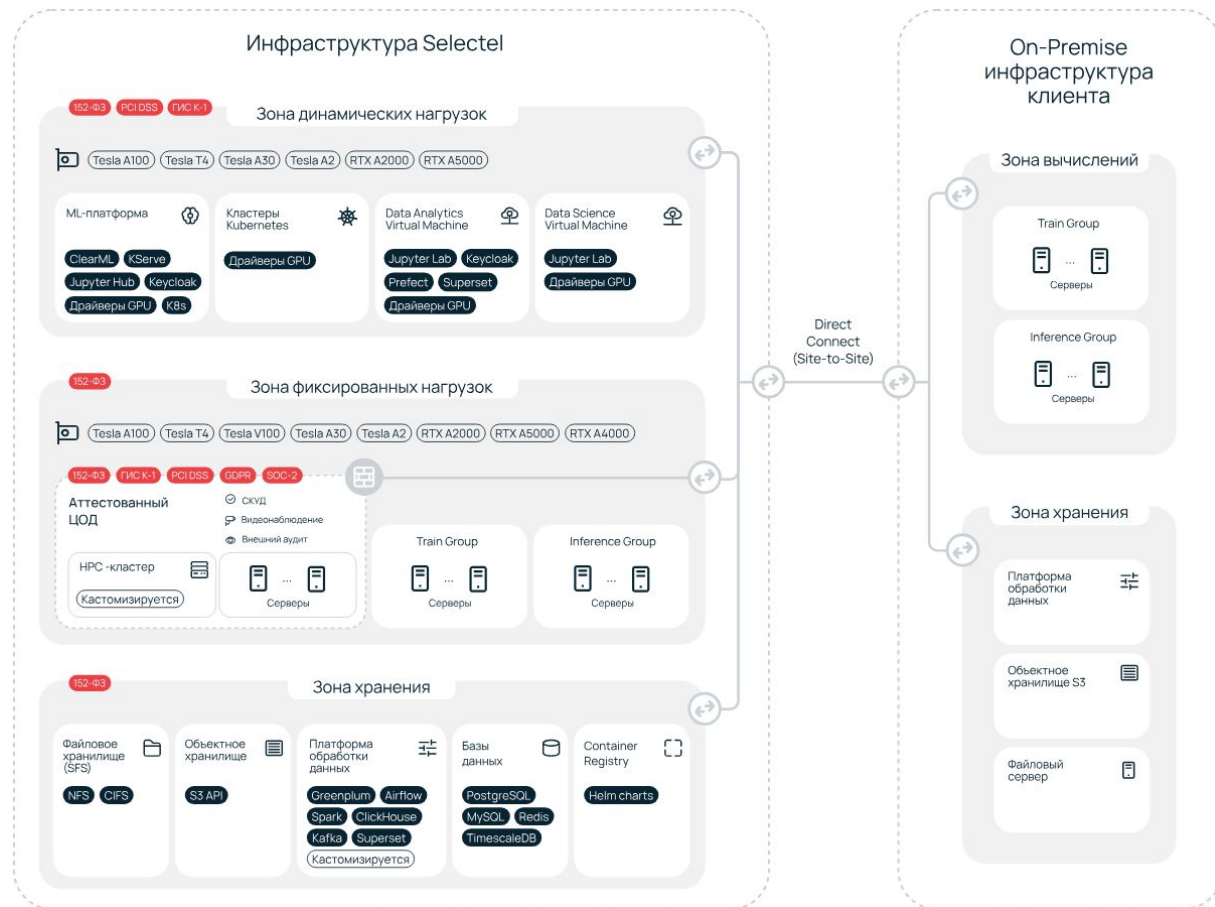
# Infrastructure чек-лист

Важно подумать над вопросами:

1. Какую NLP-технологию решено использовать?
2. Есть ли обоснование внедрения или нужен PoC?
3. Есть ли ответственные от IT и целевого бизнес-подразделения?
4. Какова предполагаемая нагрузка на итоговый сервис?
5. Нужно ли дообучать ML-модели?
6. Подходит ли формат собираемых бизнес-данных для дообучения ML-моделей?
7. Нужны ли изолированные контуры для training и inference?
8. Где провести тесты разных GPU под задачу?
9. Купить GPU-серверы или арендовать?
10. Требуется ли соблюдать 152-ФЗ?
11. Какой стек инструментов будет использоваться?
12. Есть ли необходимые компетенции?





# Selectel – провайдер для ML



# Selectel

 [selectel.ru](https://selectel.ru)

 [selectel](https://t.me/selectel)

 [selectel](https://vk.com/selectel)

Остались вопросы?  
Обращайтесь!