

СКОЛЬКО «ВЕШАТЬ В ГРАММАХ»?

«ВЗЯТЬ С ЗАПАСОМ» — МОЖЕМ СЕБЕ ПОЗВОЛИТЬ, НО НУЖНО ЛИ?

КАК ОТП БАНК ПОДОШЕЛ К ОЦЕНКЕ LLM-ИНФЕРЕНСА

ДЛЯ КОГО ЭТОТ ДОКЛАД? ВЫ ТОЧНО ЗДЕСЬ:

1



Если вы планируете внедрение **GenAI**

2



Если вы платите за **GPU**

3



Если вас спрашивают «почему столько?»

ЧТО ВЫ СМОЖЕТЕ ПОСЛЕ ДОКЛАДА

1



Самостоятельно рассчитать потребность

2



Защитить цифру перед финансистами

3



Уточнять и калибровать расчеты



**AI Platform Center Lead
Богдан Гарбар**

1



Предпосылки. Как мы пришли к мысли разработать свою методику

2



Методика. Принципы расчета

3



Проверка. Уточнение и калибровка

4



Практическая польза

01

ПРЕДПОСЫЛКИ



Что мы хотим, и сколько GPU для этого нужно на 1-2 года вперед?



Ну что, ребят? Рассказывайте...
Чего, сколько, какое ROI?

РЕАЛЬНОСТЬ LLM-ИНФЕРЕНСА СЕГОДНЯ



Много и дорого

Для создания своего LLM инференс кластера нужно большое кол-во GPU. А GPU – это дорого.

1



Сложности с поставками в РФ

В текущей обстановке, наблюдаются сложности с поставками, что увеличивает сроки и повышает риски.

2



Ошибка в подсчетах - не дешево

GPU – это дорого, а поставки длительные. Подсчет «пальцем в небо» может дать существенную ошибку.

3



Быстрое устаревание железа

Новые поколения AI моделей и чипов появляются каждый квартал. Нужно выжать максимум из оборудования за конечное время его актуальности.

4



“Взять с запасом” – можем конечно себе позволить, но нужно ли?

Любое инженерное капиталоемкое решение должно быть объяснимо, воспроизводимо и защищено.



КАКИЕ ИНСТРУМЕНТЫ ОЦЕНКИ БЫЛИ ДОСТУПНЫ И ПОЧЕМУ ИХ ОКАЗАЛОСЬ НЕДОСТАТОЧНО

В процессе планирования мы использовали:

Публичные и вендорские калькуляторы
для первичной прикидки

Эвристические оценки и бенчмарки
как ориентиры

Экспертные статьи и кейсы
как источник гипотез

Почему этого оказалось недостаточно для production:

Эти подходы не дают воспроизводимой логики расчёта

Не всегда есть прямая связь с SLA и пользовательской нагрузкой

Сложно адаптировать расчёт под реальные сценарии эксплуатации

Сложно или невозможно защищать итоговую цифру перед CFO

Всё это недостаточно прозрачно для принятия обоснованных решений

А может сделаем расчет сами?
Можно будет объяснить, повторить и
уточнить по данным...

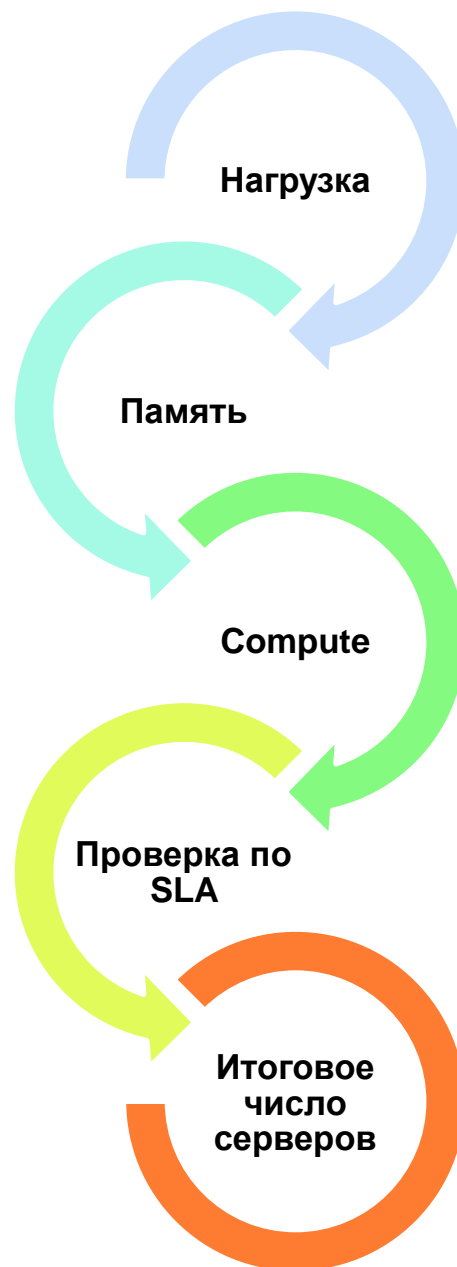
02

ПРИНЦИПЫ РАСЧЕТА

ЧТО И КАК МЫ СЧИТАЕМ?

Расчет потребления памяти
на KV кэш и веса моделей

Проверка
конфигураций по SLA



Расчет кол-ва
пользователей и
потребления токенов

Расчет пропускной
способности

Максимальное из расчетов по
Памяти и Пропускной
способности

ОБЪЯСНЕНИЕ МЕТОДИКИ ЧЕРЕЗ АНАЛОГИЮ С СОСТАВОМ МЕТРО

АНАЛОГИЯ

У города есть N жителей, которых надо перевозить каждый день.

Сколько составов нужно закупить, чтобы все доехали вовремя?

ДВА МОМЕНТА

- Хватает ли мест, чтобы все сели?
- Хватает ли скорости, чтобы всех перевезти за час?



В переводе на GPU:

сколько серверов нужно, чтобы обслужить целевое число пользователей с заданным SLA?

ЧТО ПОДАЕМ В КАЧЕСТВЕ ИСХОДНЫХ ДАННЫХ

Пользовательская нагрузка

- Количество пользователей
- Конкурентность (одновременные сессии)

Характер запросов

- Средний размер контекста (tokens)
- Средняя длина ответа (tokens)

Модель

- Размер и архитектура (вес модели)

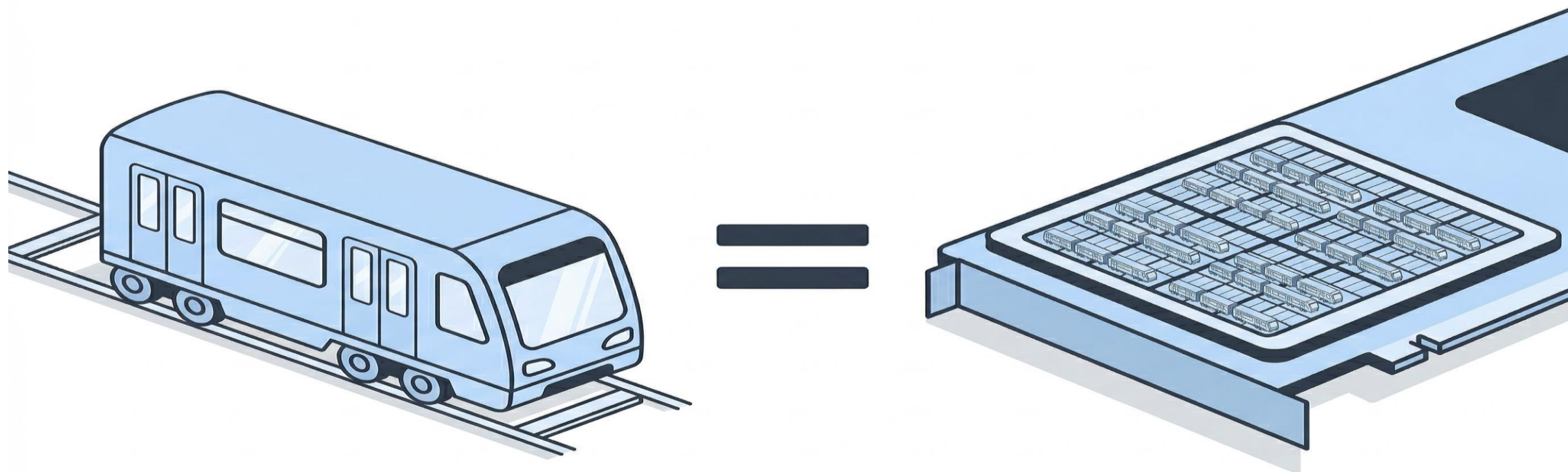
Требования

- Целевая латентность (SLA)
- Требуемый throughput (tokens/sec)

Инфраструктурные допущения

- Тип и конфигурация GPU (как входное ограничение)

ПРЕДСТАВИМ, ЧТО GPU = ВАГОН МЕТРО



GPU-чип

= один вагон метро
Память (HBM) = общий объём внутри вагона

GPU-сервер

= состав (несколько вагонов в сцепке)

Кластер GPU

= сеть метро города
(много линий,
много составов)

ТОГДА, ЧТО ЗАНИМАЕТ МЕСТО В ВАГОНЕ?



Кабина + Двигатель + Оборудование
= Weights модели (параметры нейросети)

Пассажирский салон
= KV-cache (контекст каждой сессии)

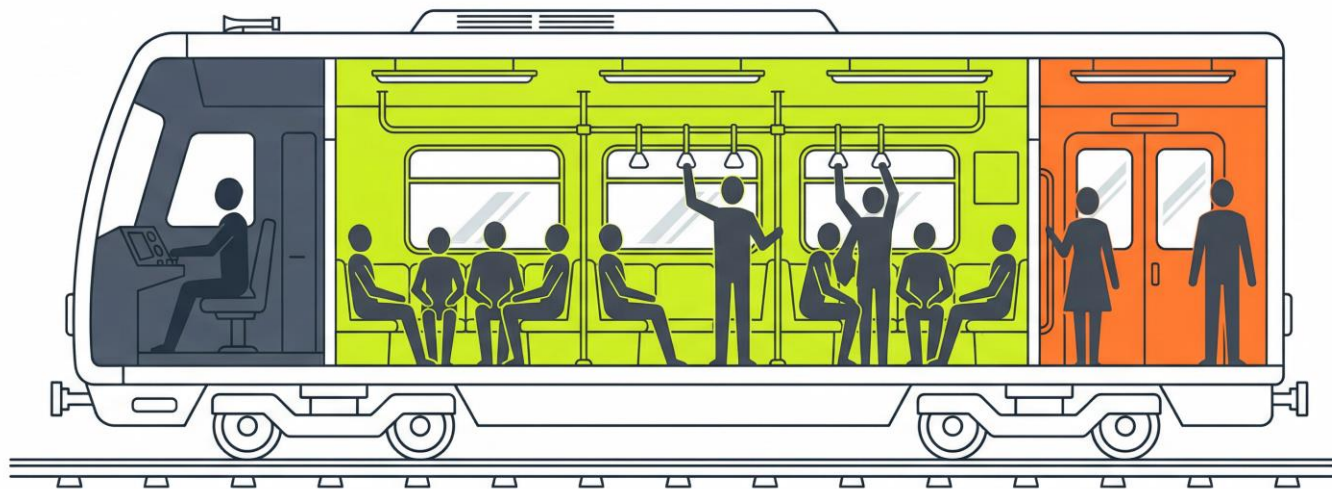
Проходы и тамбуры
= Reserve (буферы, активации, KV overhead)

Формула:

Память GPU = Кабина (Weights) + Салон (KV-cache) + Проходы (Reserve)

Чем больше кабина (модель) — тем меньше места для пассажиров (сессий)

ПАССАЖИРЫ И БАГАЖ (KV-CACHE)



Пассажир = сессия пользователя

Каждый пассажир занимает место (KV-cache) в салоне. Чем длиннее его маршрут (контекст), тем больше багажа.

ВМЕСТИМОСТЬ САЛОНА

- Размер вагона — память GPU (HBM)
- Минус кабина — веса модели (weights)
- Минус проходы — резерв на пики (reserve)

Формула:

$\text{Max сессий} = (\text{HBM} - \text{Weights} - \text{Reserve}) / \text{KV на сессию}$

Чем больше размер вагона — тем больше помещается пассажиров

СОСТАВ ИЗ НЕСКОЛЬКИХ ВАГОНОВ (TENSOR PARALLELISM)



Один машинист на весь состав

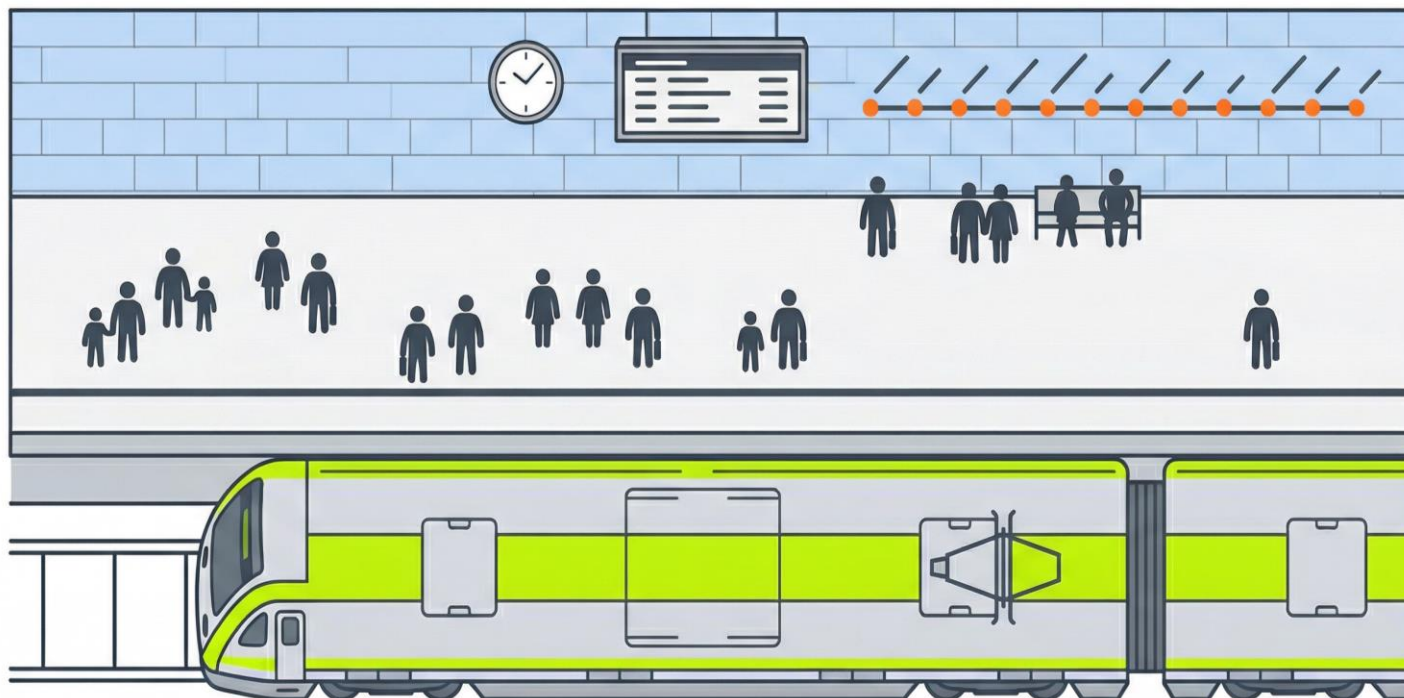
Модель (weights) распределена по всем вагонам (GPU), но работает как единое целое.

Зачем нужен состав?

Если модель не помещается в один вагон — сцепляем несколько.

Например, TP=2: два вагона на одну модель.

СКОЛЬКО ПОЕЗДОВ НАДО ЗАПУСТИТЬ?



Ограничение по местам:

Каждый поезд вмещает ограниченное число пассажиров одновременно.

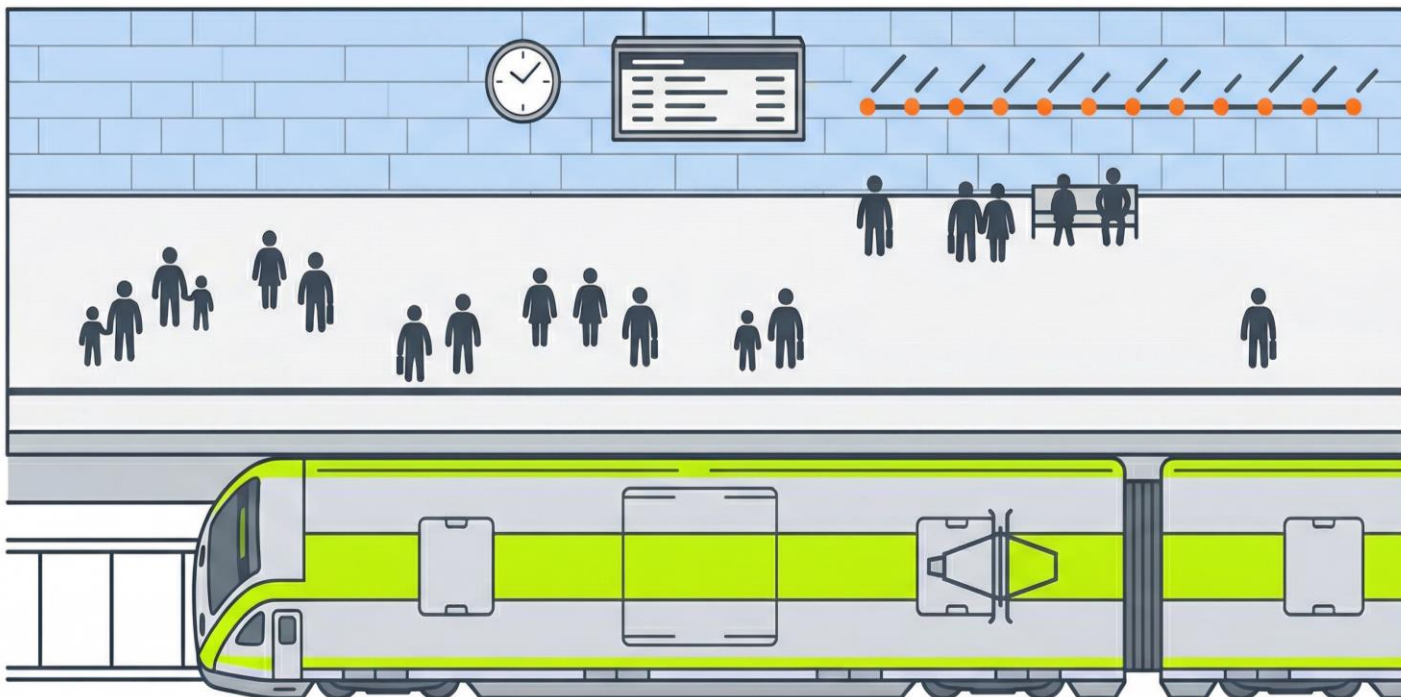
- Concurrent users — сколько пассажиров одновременно в метро
- Max sessions per train — вместимость одного состава (зависит от GPU-памяти, модели, контекста)

Формула:

$$N_{\text{поездов (по местам)}} = \text{Concurrent_users} / \text{Max_sessions_per_train}$$

Чем больше пассажиров в метро — тем больше требуется составов

РАСПИСАНИЕ: СКОРОСТЬ ПЕРЕВОЗКИ



TTFT

ОЖИДАНИЕ НА ПЛАТФОРМЕ

Сколько пассажир ждёт первый поезд

Time to First Token

ITL

ВРЕМЯ МЕЖДУ СТАНЦИЯМИ

Скорость, с которой генерируются следующие токены

Inter-Token Latency

THR

ПАССАЖИРОВ ЗА ЧАС

Общая пропускная способность линии

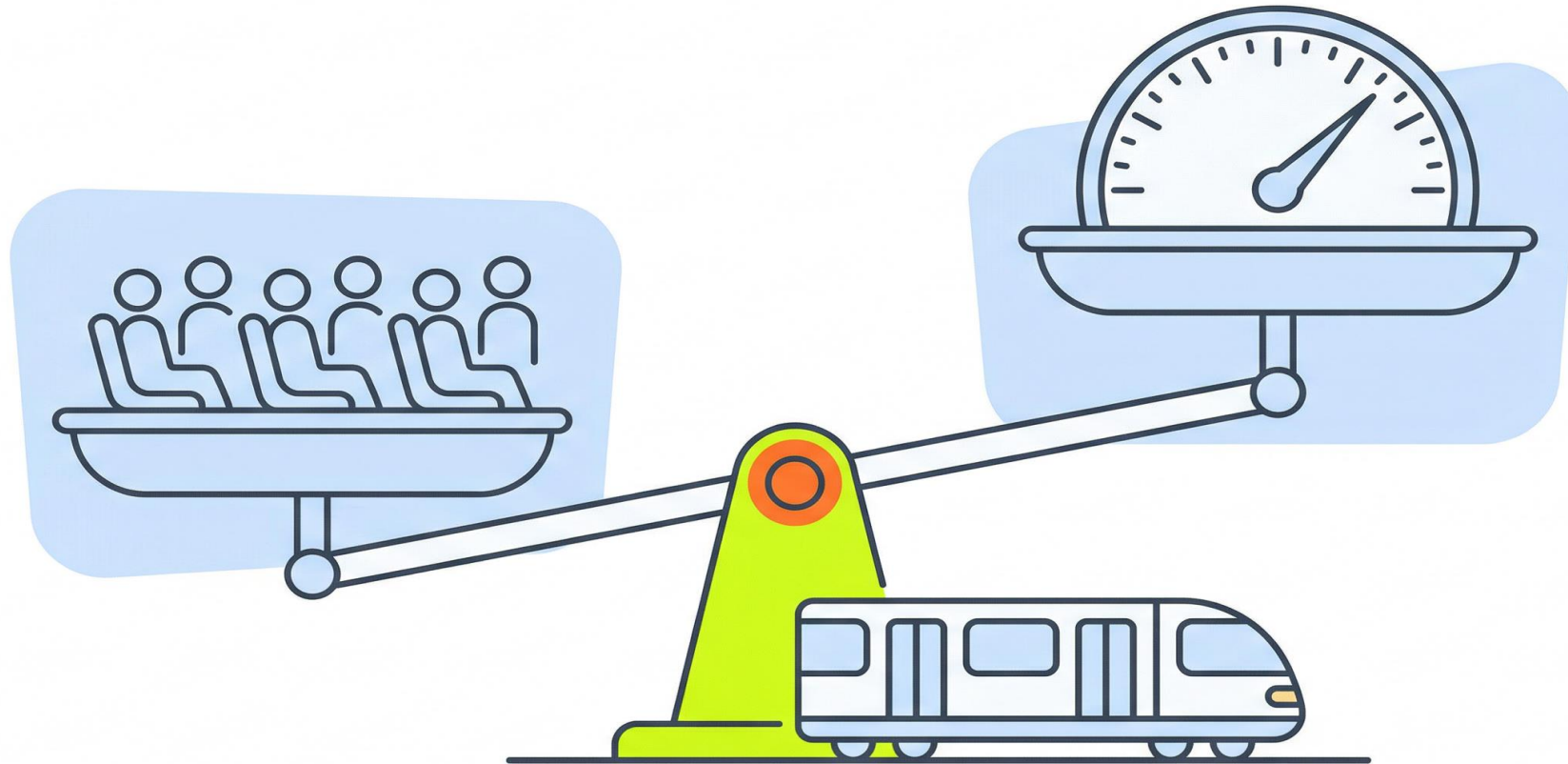
Throughput

Формула:

$$N_{\text{поездов (по скорости)}} = \text{Target_throughput} / \text{Throughput_per_train}$$

Чем больше пассажиров надо перевезти за единицу времени — тем больше требуется составов

ОТВЕТ НА ВОПРОС «СКОЛЬКО ВЕШАТЬ В ГРАММАХ»



Формула:

$$N_{\text{servers}} = \text{MAX}(N_{\text{по_местам}}, N_{\text{по_скорости}})$$

Берём максимум из двух ограничений — вместимость и пропускная способность

А что если мы не уложились в целевой SLA?
Что можно сделать, чтобы улучшить
показатели?

ОПТИМИЗАЦИЯ: СЕРИИ ВАГОНОВ (Н100 / Н200 / В200)



Н100

80 GB НВМ3

Классический вагон.
Достаточно для
большинства
моделей до 70В.

Н200

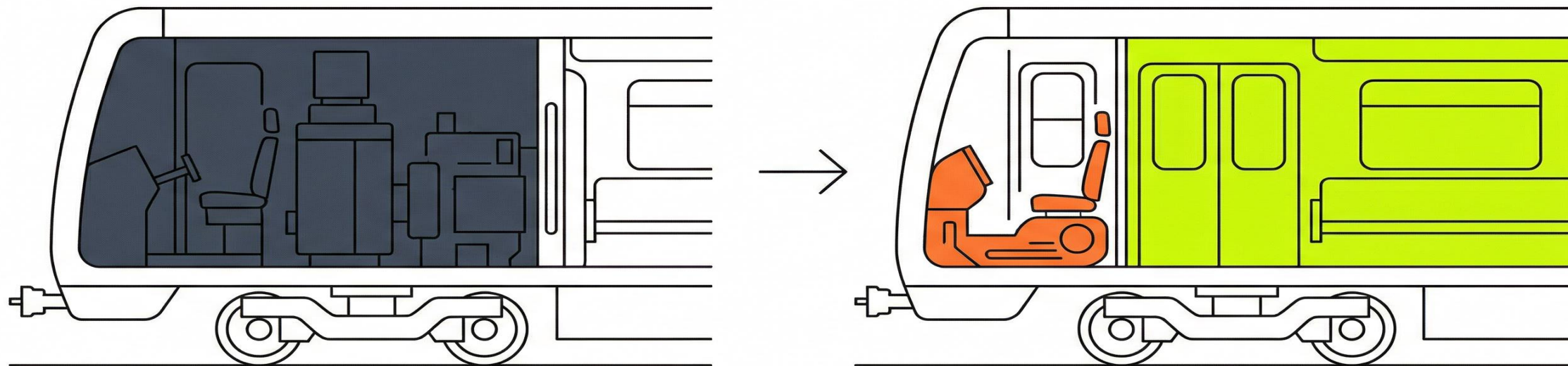
141 GB НВМ3е

Просторный вагон.
Вмещает больше
пассажиров
(длинные контексты).

В200

192 GB НВМ3е

Вагон нового поколения.
Больше места + лучше
технические условия для
размещения двигателя.



КВАНТОВАНИЕ

Компактное оборудование кабины

Уменьшаем «оборудование кабины» — освобождаем место для пассажиров.

2-4x

Менше памяти на веса модели

Точность падает на единицы процентов — пропускная способность кратно растёт



CONTINUOUS BATCHING

Подбираем пассажиров на ходу

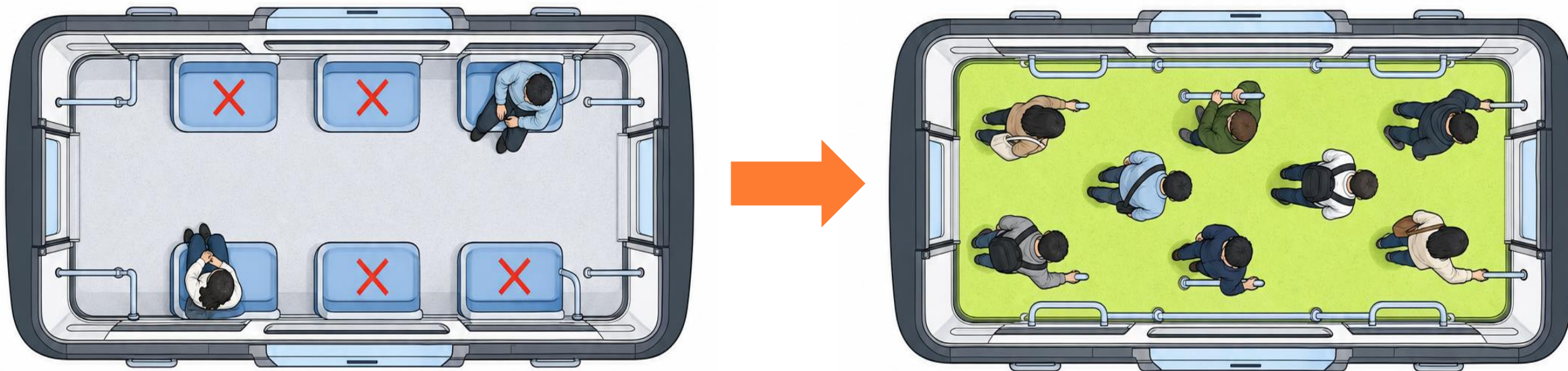
Один пассажир сошёл — на его место тут же садится следующий, не дожидаясь конечной станции.

2–3×

выше utilization GPU

Вагон всегда +/- полный.

ОПТИМИЗАЦИИ: КОМПАКТНОЕ ОБОРУДОВАНИЕ



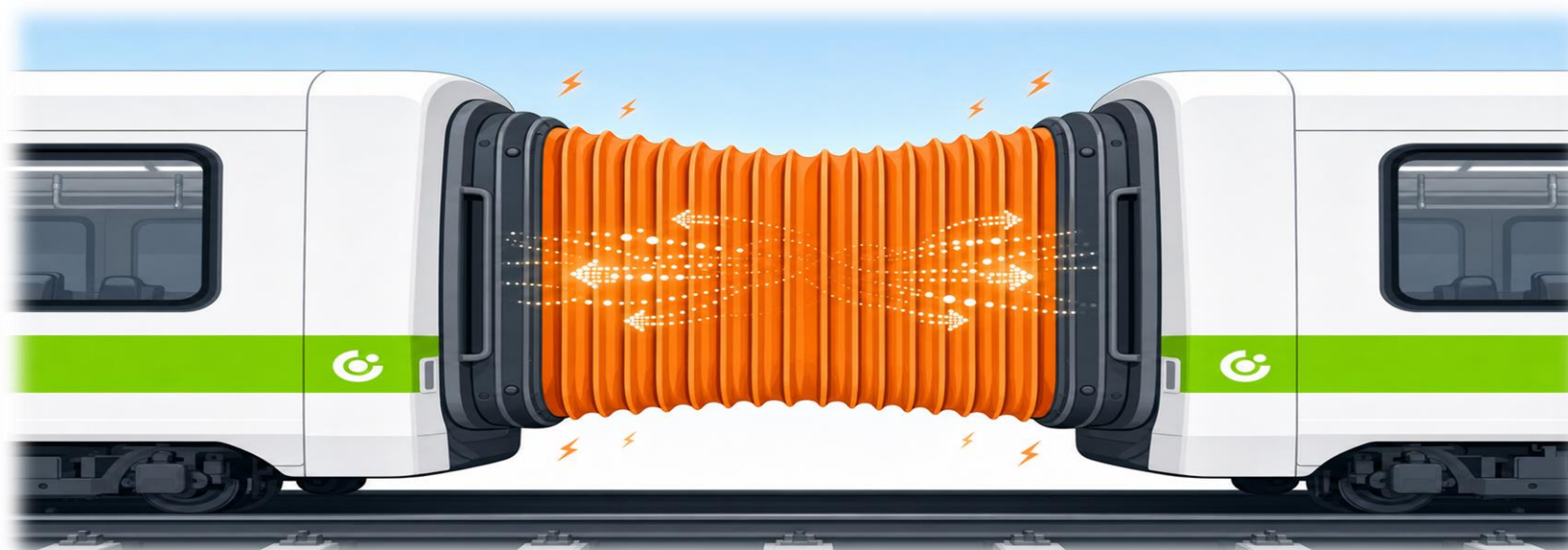
PAGED ATTENTION

Динамическая рассадка

Рассаживаем пассажиров гибко — без фиксированных, зарезервированных мест.

2–3×**Больше параллельных сессий**

на тот же GPU без потери качества



NVLINK

Мягкая сцепка между вагонами

Поток пассажиров внутри вагона осуществляется бесшовно и намного быстрее

5–10x

ПРОПУСКНАЯ СПОСОБНОСТЬ

выше, чем просто у PCIe

Новая модель — другие размеры

Пересели с Llama-70B на Llama-405B?
Пересчёт всех параметров.

Длинный маршрут

Пользователь с контекстом 128K
занимает места как 10 обычных.

Вывод:

Расчёт даёт хорошую стартовую точку, но нужна калибровка.

Формула + замер + корректировка = надёжная оценка.

03 ПРОВЕРКА И УТОЧНЕНИЕ

КАЛИБРОВКА

1. Расчёт

Формула даёт теоретическую оценку N серверов.

Входные данные: модель, пользователи, SLA.

2. Замер

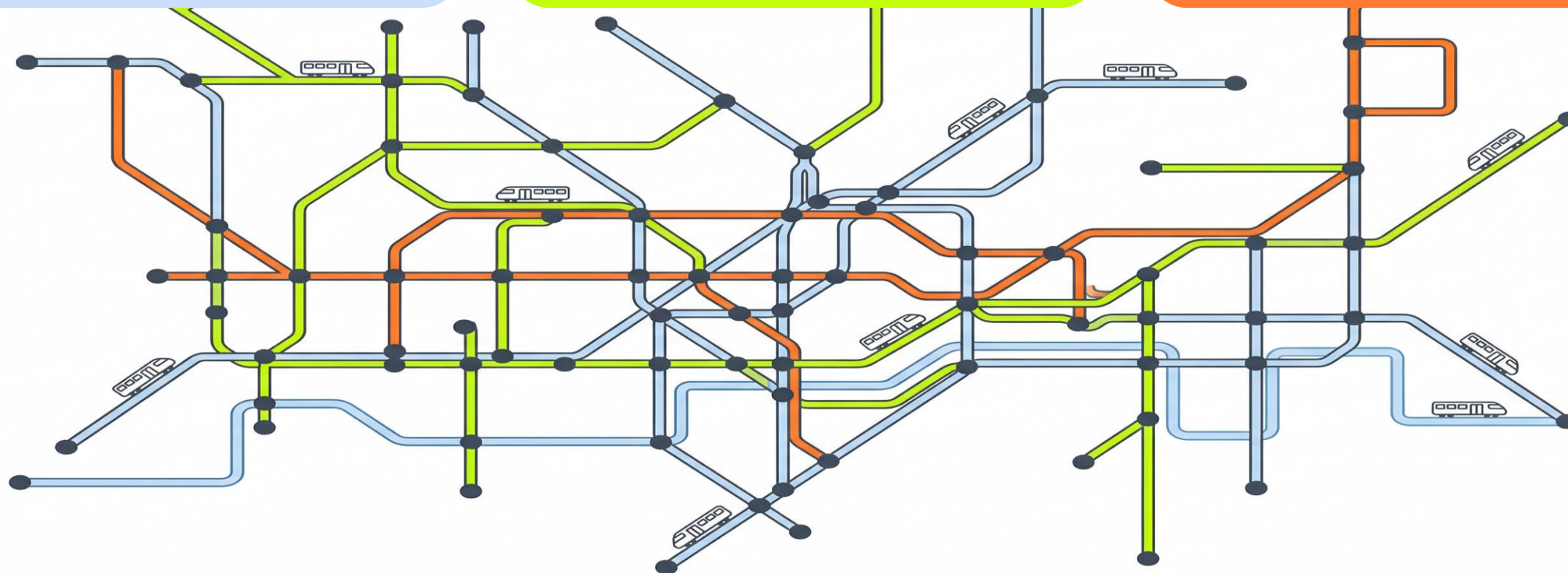
Бенчмарк на реальной инфраструктуре.

TTFT, ITL, throughput — сравниваем с расчётом.

3. Уточнение

Корректируем формулу по результатам.

Итеративный цикл:
расчёт → замер → уточнение.



НАГРУЗОЧНОЕ ТЕСТИРОВАНИЕ: ЧТО МЕРЯЕМ

Три параметра, которые отличают работающий стенд от красивого слайда

01

Latency

Укладываемся ли в SLA

TTFT и ITL — основные метрики качества для пользователя

02

Throughput

Запросов и токенов в секунду

Максимальная пропускная способность, после которой начинается деградация

03

Concurrency

Параллельные сессии

Как система ведёт себя при росте числа одновременных пользователей

ВАЖНЫЙ ПРИНЦИП:

Тестируем не одну точку, а поведение системы до предела

НАКОПЛЕНИЕ СТАТИСТИКИ ИСПОЛЬЗОВАНИЯ

Фактические токены

Реальное распределение длин запросов и ответов

Реальные пики активности

Пиковая нагрузка по времени суток

Реальная одновременность

Фактическое число параллельных сессий

Реальные профили пользователей

Поведенческие паттерны и частота обращений

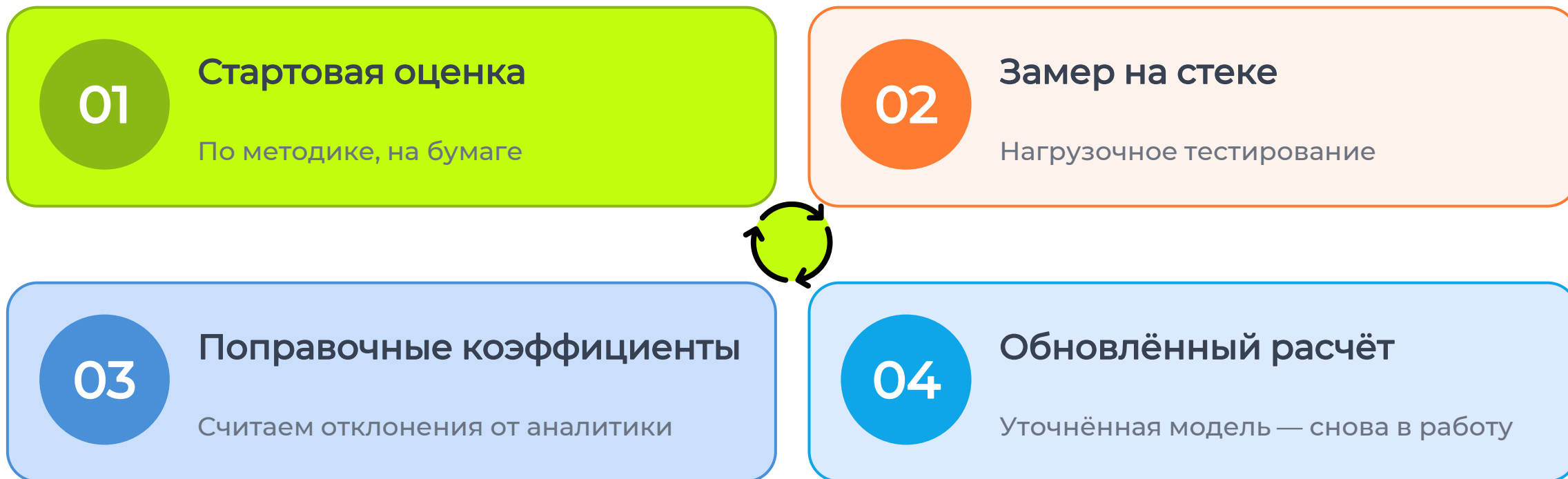


ВАЖНЫЙ ПРИНЦИП:

Теория → Практика: каждый параметр уточняется на реальных данных

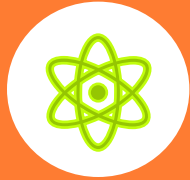
ЦИКЛ УТОЧНЕНИЯ

Расчёт — это не «посчитал и забыл», а непрерывный цикл уточнения по фактическим данным



НАБЛЮДЕНИЕ:

Через месяц-два production оценка из теоретической становится эмпирической



Идеальная ли это модель?

Это не идеальная модель

Но достаточная для оценки сценариев
и принятия взвешенных решений

И мы уверены, что она будет становиться точнее

04

ПРАКТИЧЕСКАЯ ПОЛЬЗА

БЕРИ И ПОЛЬЗУЙСЯ: OPEN SOURCE

2 Definition of Simultaneously Active Users

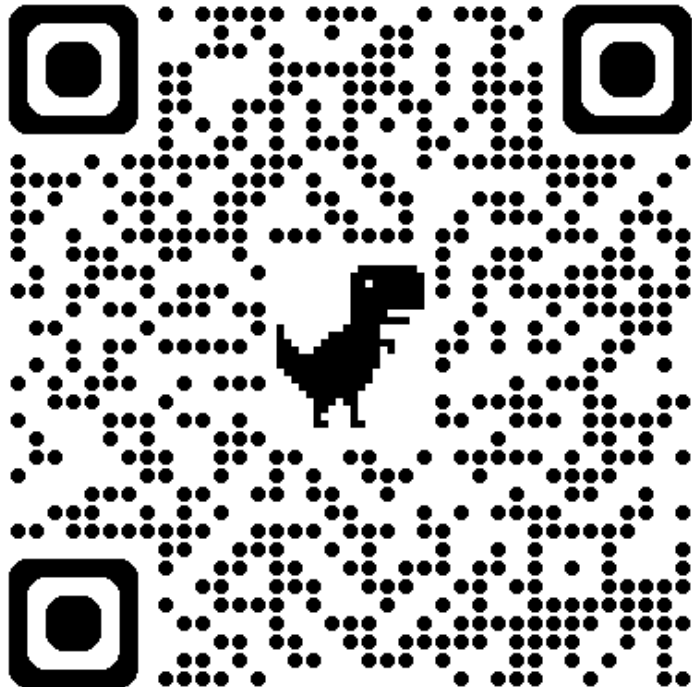
The calculation of the number of users is performed separately for each segment. In particular, if there is a division into external and internal users (for example, staff and contractors, staff specialists and vendor/customer representatives), the calculation is performed for the internal segment, then for the external segment.

$$N_{sim} = N_{users} \times K_{pen} \times K_{sim} \quad (1)$$

Used in: eq. (13), eq. (18)
where:

Variable	Definition	Unit	Note
N_{sim}	Number of employees simultaneously using the service per unit time	User/unit time	Calculations are performed for each user segment separately if needed (e.g., internal employees

Полный текст методики опубликован на Archive.org и в GitHub



Чтобы выйти из полноэкранного режима, нажмите и удерживайте Esc

AI Infrastructure Calculator

Find out how many servers and GPUs you need for your AI models

[Take a Tour](#)
[Star on GitHub](#)
[Documentation](#)

Configuration Parameters

Quick Presets

- 32B / 500K users (A100 80GB | 2 = 4)
- 7B / 2K users (A100 80GB | 2 = 1)
- 70B / 10K users (H100 80GB | 2 = 2)

Basic Configuration | Advanced

Users

Total Users:

Model

Search Model:

SELECTED MODEL: Qwen/Qwen3-VL-30B-A3B-Instruct

Hardware

GPU: NVIDIA H100 GPU accelerator (PCIe card)

SELECTED GPU: NVIDIA H100 GPU accelerator (PCIe card) (80 GB) | 756.4 TFLOPS

GPUs per Server:

Usable Memory Fraction:

Tensor Parallelism

TP Degree:

Calculation Results

COST ESTIMATE

CONCURRENT SESSIONS: 2.5K

SESSION CONTEXT: 4.0K tokens

SERVERS REQUIRED: 22
max(mem: 22, comp: 12)

SESSIONS PER SERVER: 116
2 inst × 58 sess each

SERVER THROUGHPUT: 5.3
prefill: 12.1K · decode: 9.0K

GPUS PER SERVER: 8

GPUS PER INSTANCE: 4

INSTANCES PER SERVER: 2

SLA Validation ✓ SLA Passed

Time To First Token (TTFT)		End-to-End Latency	
Calculated	0.33 sec	Calculated	0.38 sec
SLA Target	1.00 sec	SLA Target	2.00 sec
PASS		PASS	

GPU Memory per Instance

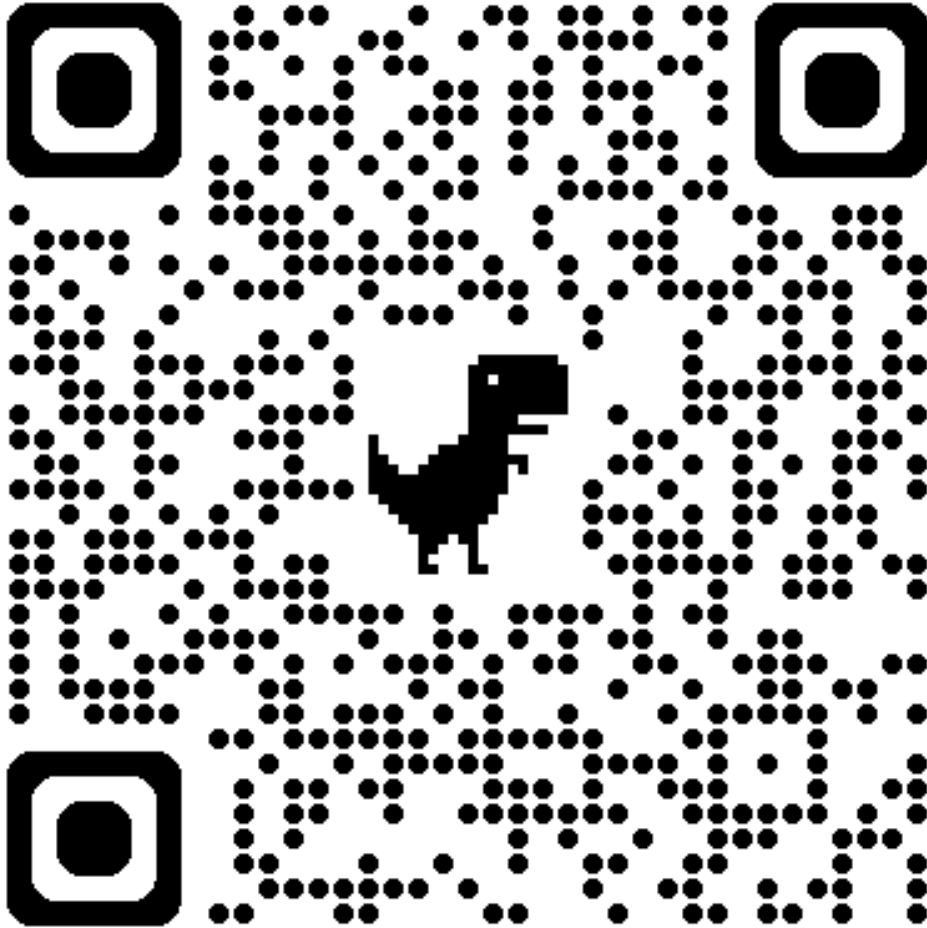
4 GPUs × 80 GiB = 320.0 GiB total

320.0 GiB

- Model Weights: 60.88 GiB
- Available for KV-cache: 227.12 GiB
- Reserved (1 - Kavail): 32.00 GiB

Detailed Results Memory Path | Compute Path

БЕРИ И ПОЛЬЗУЙСЯ: OPEN SOURCE REPOSITORY



The screenshot shows a GitHub repository page for 'AI.ServerCalculationApp' (Private). The repository has 2 branches and 0 tags. A notification at the top states 'Your main branch isn't protected'. The commit history shows several updates, including UI corrections and deployment scripts. The README section is visible, titled 'AI Server Calculator', and includes a description of the application and a preview of the 'AI Infrastructure Calculator' web interface.

File	Description	Time
app	Корректирование UI (с учетом ext users)	3 months ago
frontend	Корректирование UI (с учетом ext users)	3 months ago
nginx	Fix status check in nginx	4 months ago
LICENSE	Initial commit	4 months ago
METHODOLOGY_README.md	Add common readme	4 months ago
README.md	Update README.md	now
docker-compose.prod.yml	Redeploy script, separate docker-compose file	4 months ago
docker-compose.yml	Redeploy script, separate docker-compose file	4 months ago
redeploy.sh	Redeploy script, separate docker-compose file	4 months ago

AI Server Calculator

Полнофункциональное приложение для расчета требований к серверной инфраструктуре для AI/LLM моделей. Включает в себя FastAPI backend и React frontend с интерактивным интерфейсом.

AI Infrastructure Calculator
Find out how many servers and GPUs you need for your AI models

Configuration Parameters: Basic Configuration, Advanced

Calculation Results: Final Server Count, GPUs per Server

Contributors: PavelSozonov, PovarisSFEDU, Interloperok

Languages: Python 62.6%, JavaScript 34.3%, Dockerfile 1.2%, Other 1.9%

<https://github.com/Interloperok/AI.ServerCalculationApp>

Здесь вы найдете все что нужно!

3 ШАГА НА ЗАВТРА

1



Соберите вход и прогоните калькулятор

2



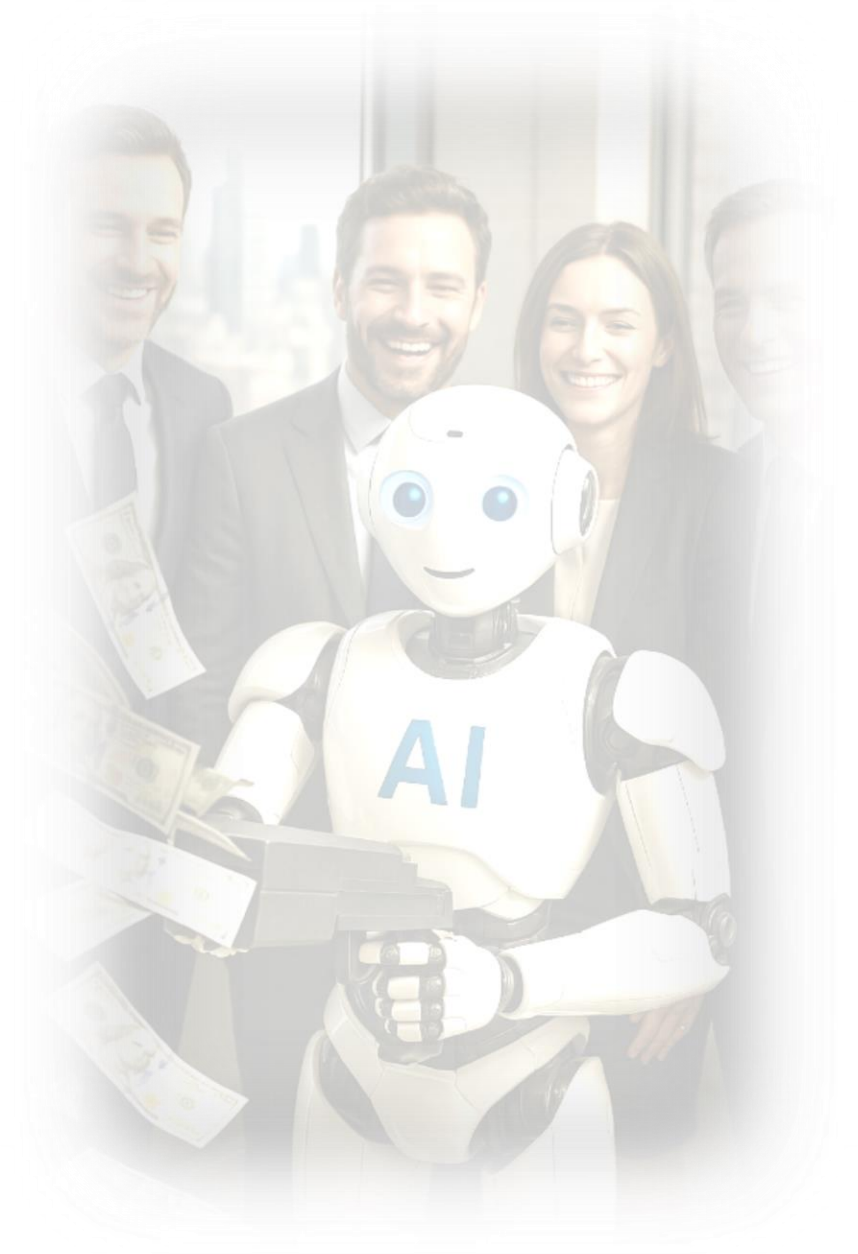
Защитите цифру перед CFO на 2-3 года

3



**Закупите минимум, подключите метрики
и собирайте факт с последующей
калибровкой**

05 ЗАКЛЮЧЕНИЕ



ТАК СКОЛЬКО «ВЕШАТЬ В ГРАММАХ»?

«ВЗЯТЬ С ЗАПАСОМ» — МОЖЕМ СЕБЕ ПОЗВОЛИТЬ, НО НУЖНО ЛИ?

*МЫ ДЛЯ СЕБЯ НА ЭТИ ВОПРОСЫ ОТВЕТИЛИ ЧЕРЕЗ
ФОРМАЛИЗОВАННУЮ И ПРОВЕРЯЕМУЮ ЛОГИКУ.*

И ПОКАЗАЛИ, КАК МОЖНО ПРИЙТИ К ТАКИМ ВЫВОДАМ.

СПАСИБО ЗА ВНИМАНИЕ



AI Platform Center Lead
Богдан Гарбар

