

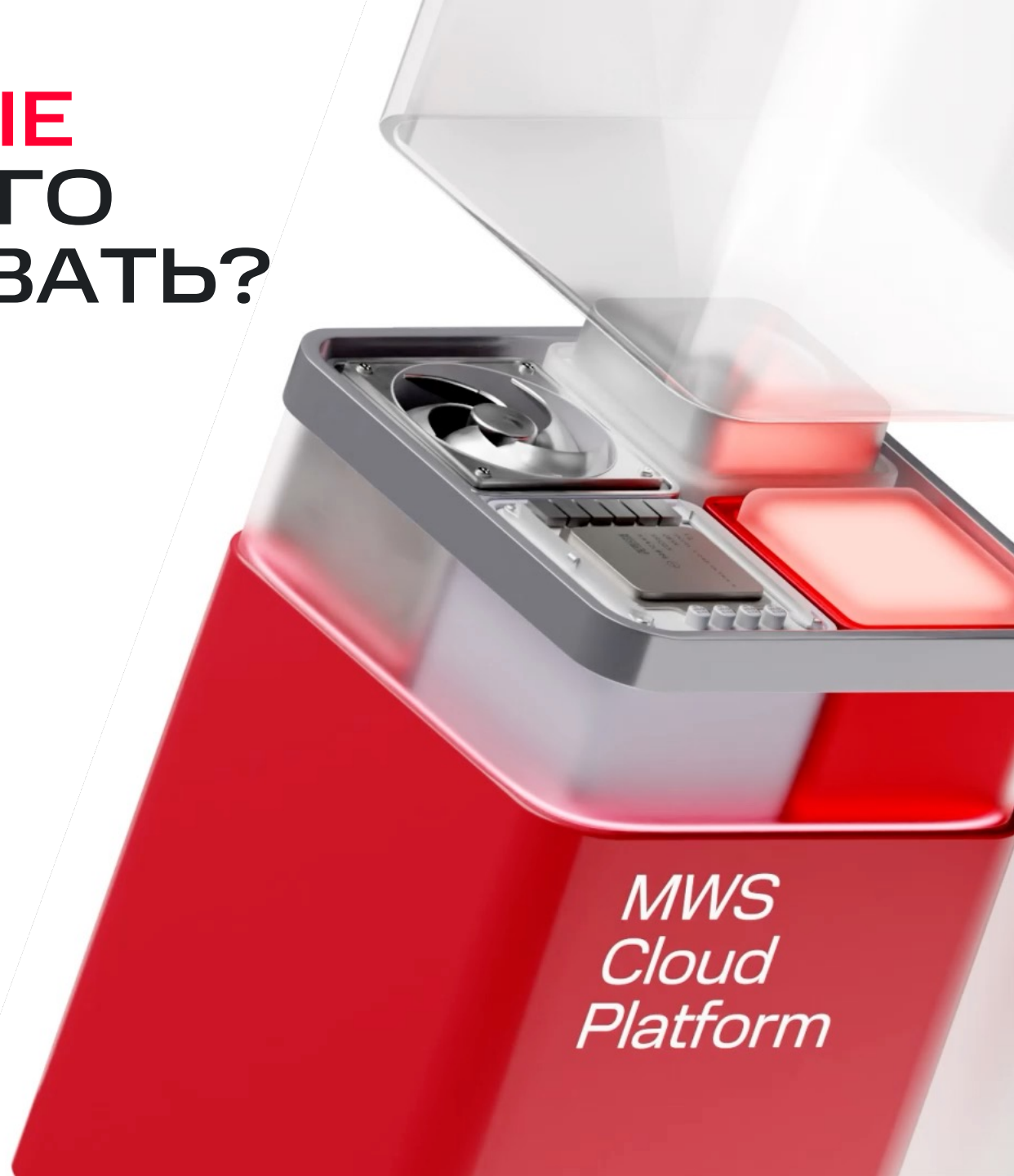
БЫСТРЫЕ ОБЪЕКТНЫЕ ХРАНИЛИЩА: ДЛЯ ЧЕГО И КАК ИХ ИСПОЛЬЗОВАТЬ?



**Кирилл
Беспалов**

Технический руководитель
MWS Cloud Platform

*MWS
Cloud
Platform*



MWS Cloud Platform Object Storage — это:

Erasure Encoding & Replication

Надёжность

KMS (шифрование)

Безопасность

?

Производительность

SLA 99,95%

Доступность

Общая архитектура



User



REST Storage API
(Golang)



postgresql
shards

Objects DB

Storage

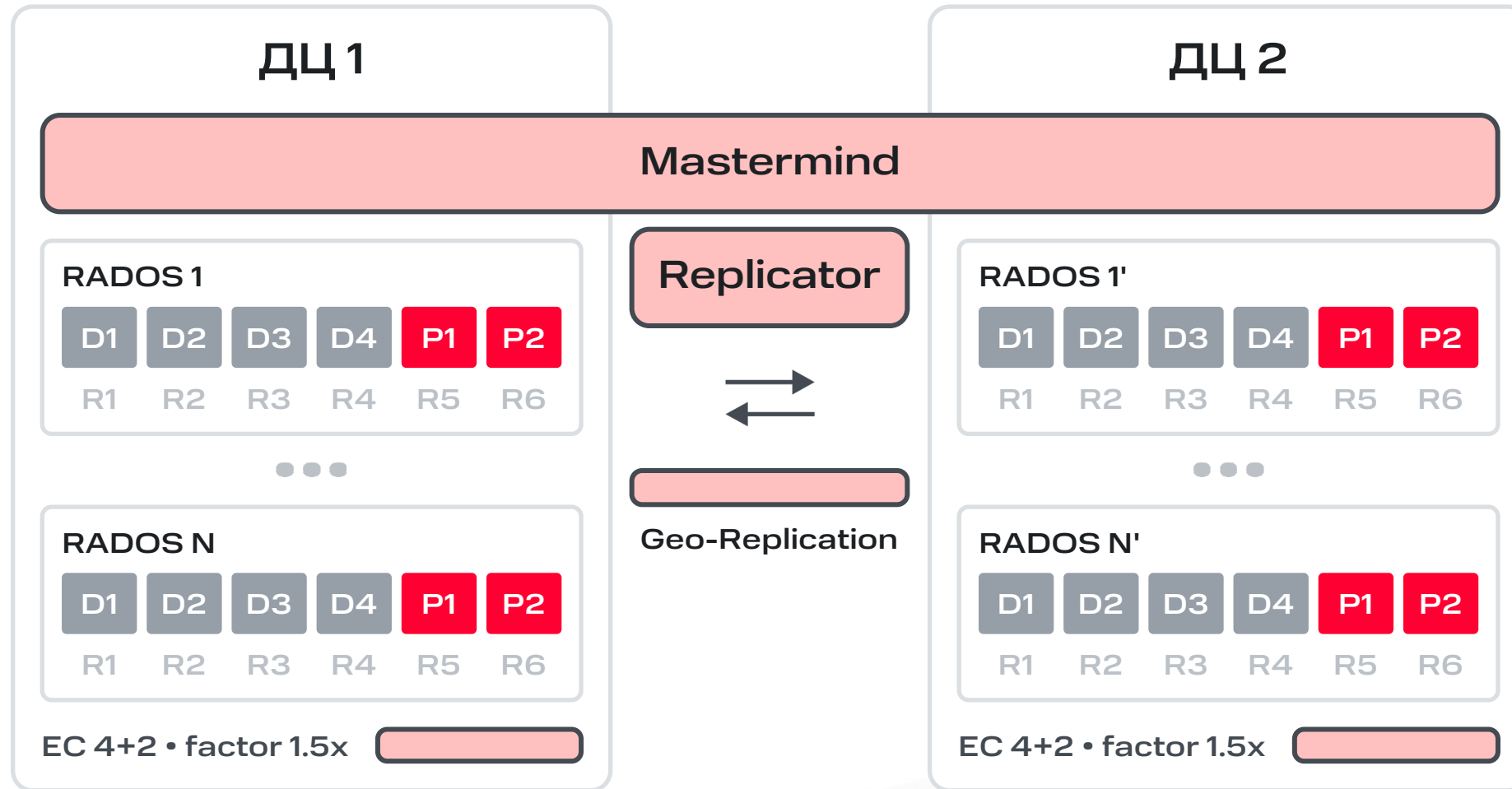


Ceph RADOS
Multiple Clusters

Building the Cloud.
Episode 3: Object Storage

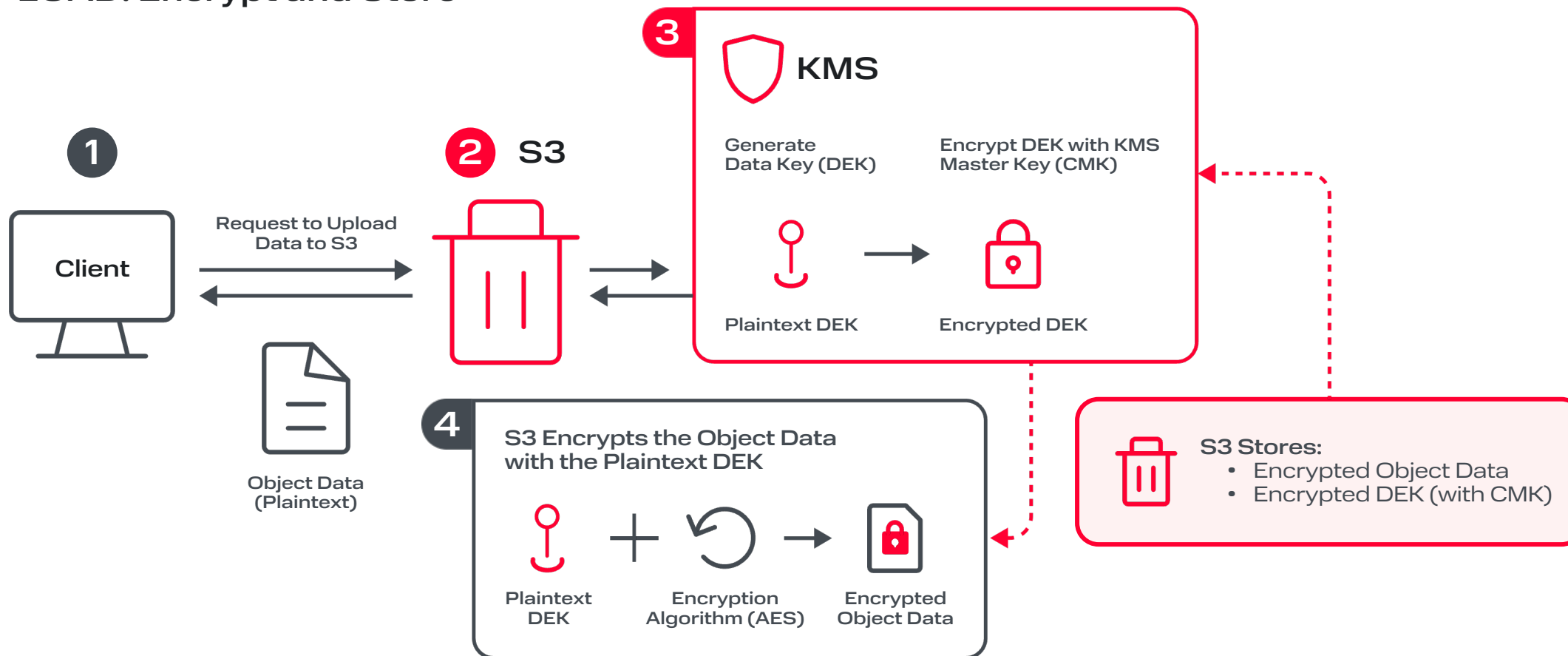


Отказоустойчивость



Безопасность

UPLOAD: Encrypt and Store



Производительность



Standard — основной класс хранения

Отлично подходит для:

- Бэкапов
- Резервных копий
- Артефактов CI/CD
- Статического контента и так далее

Standard Class

Cold

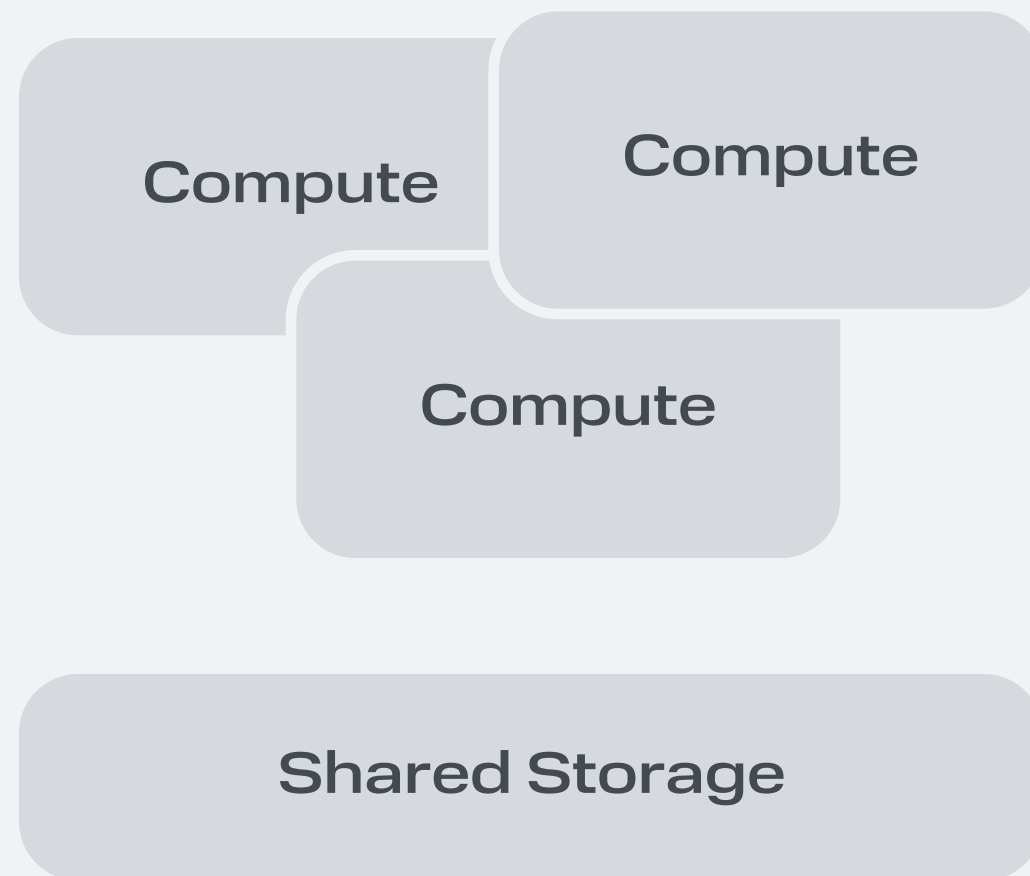
Ice



Требования к latency и throughput растут

Новые типы нагрузок

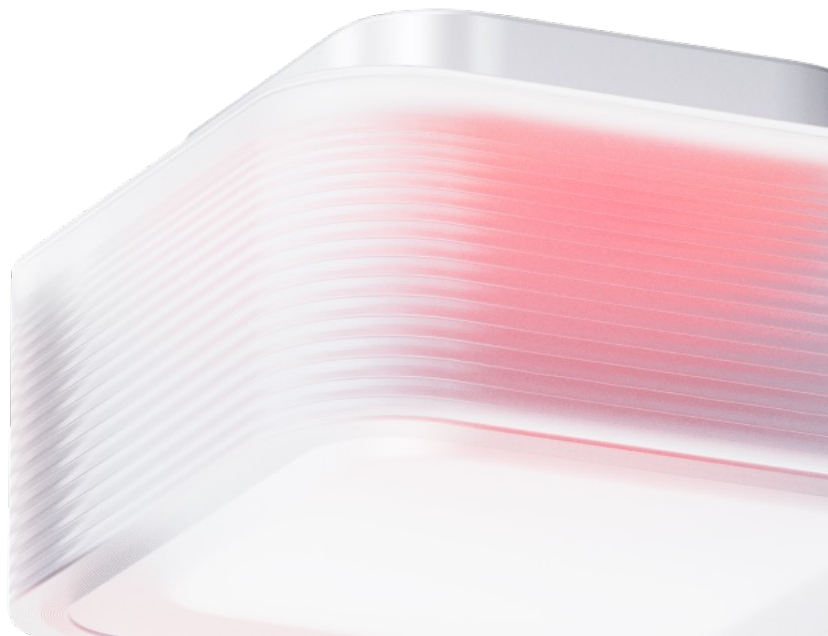
- Lakehouse & Serverless BI
- ML & Model Training
- Observability



MWS Cloud Platform — классы хранения

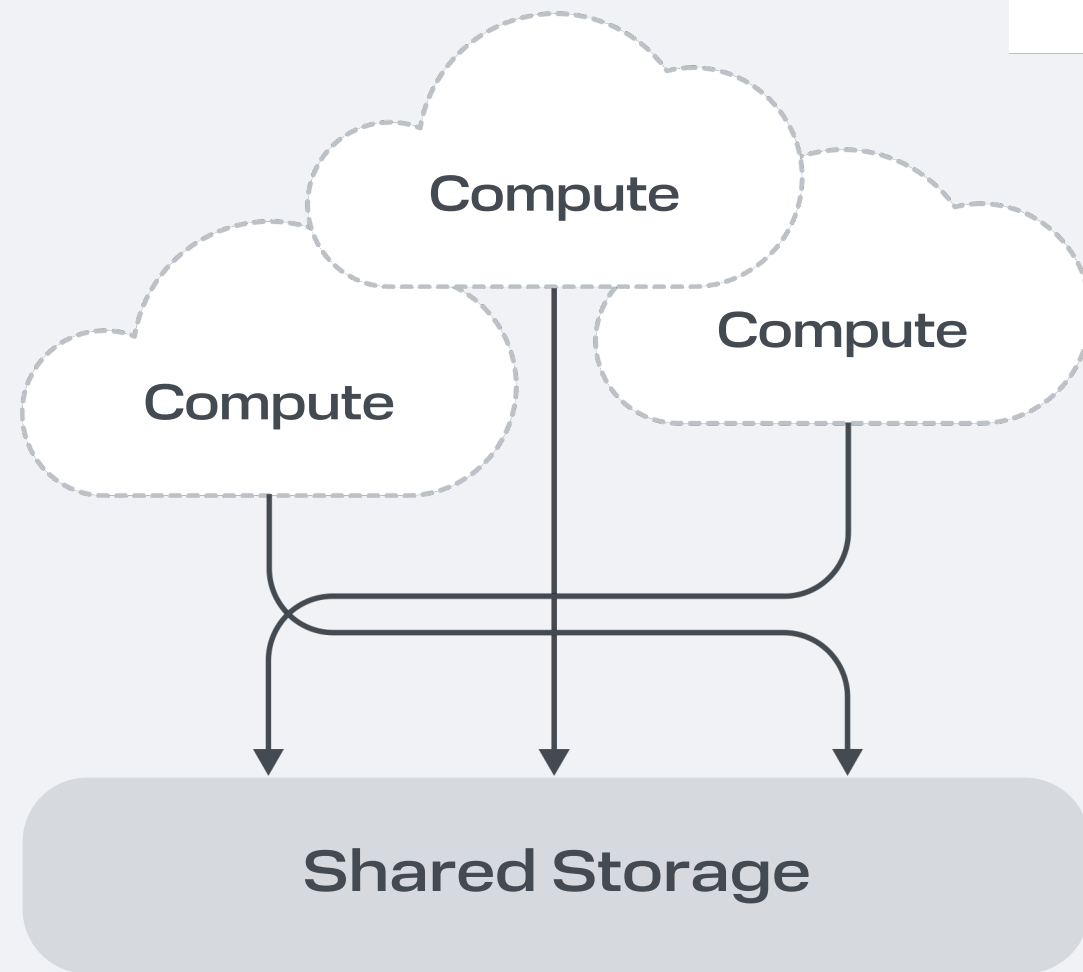
Новые типы нагрузок

- Lakehouse & Serverless BI
- ML & Model Training
- Observability



Что их объединяет?

Концепция
Compute & Storage
Separation



Базовый сценарий. Что измеряем?



VM



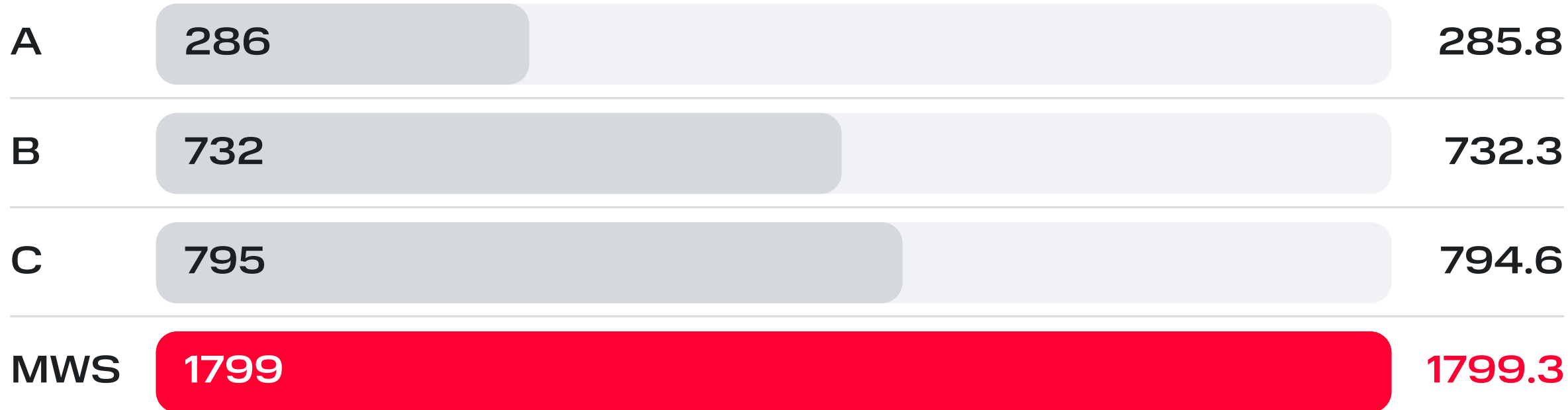
cmd: warp run scenario.yaml

PUT OBJECT 64 MB — MEDIAN THROUGHPUT

Warm storage class



MiB/s



GET OBJECT 64 MB — MEDIAN THROUGHPUT

Warm storage class

M W
S

MiB/s

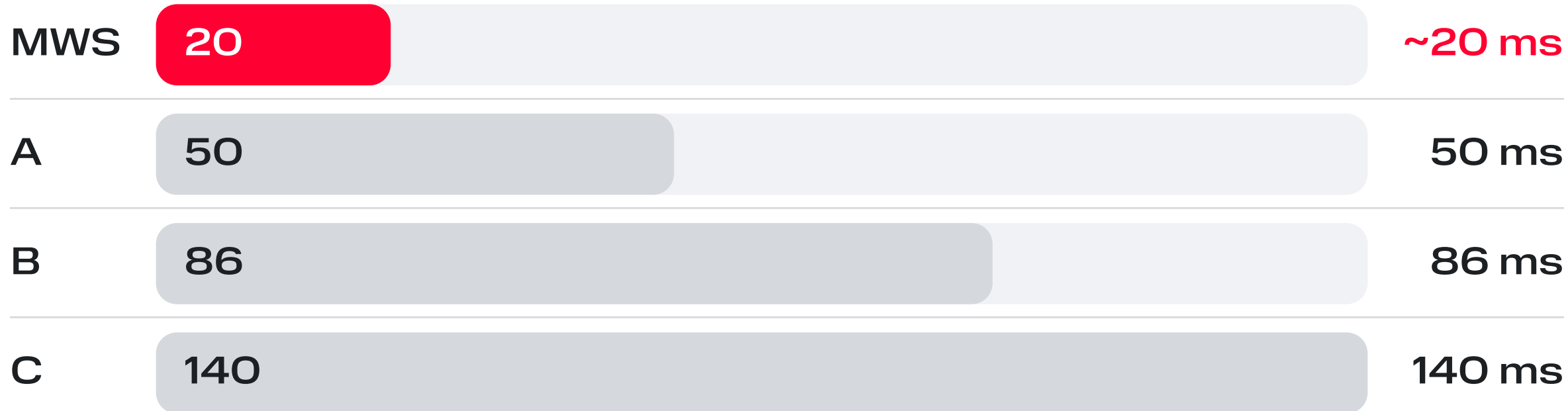
A	388	388.2
B	496	496.1
C	704	704.3
MWS	1648	1647.5

TIME TO FIRST BYTE — P90

Warm storage class



ms (меньше — лучше)



Почему так?

Высокопроизводительный
SDN

Up to 25 Gbps per link

DPDK

VPP

Производительность
meta и storage
слой Object Storage

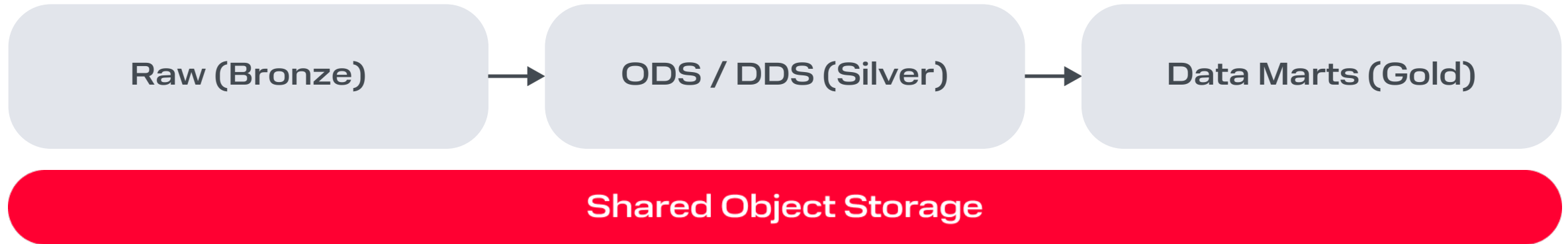
storage

storage

nvme

pg

Классический Lakehouse

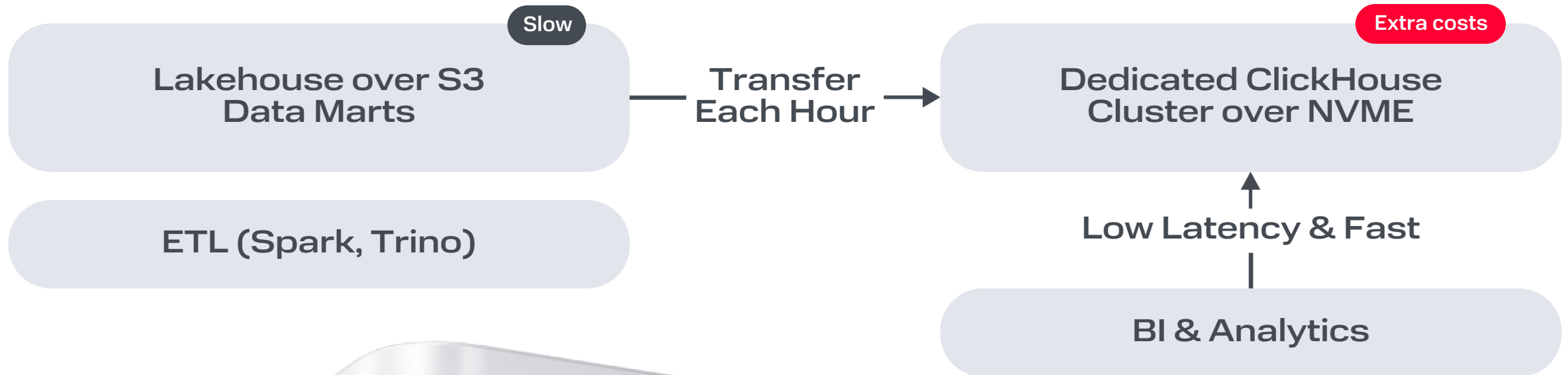


ETL over Airflow + Spark / Trino



Что это значит на практике?

Классический Lakehouse для интерактивной аналитики требует выделенного кластера, например с ClickHouse для работы аналитиков и BI



Что это значит на практике?



Пример на практике: выполнение набора аналитических запросов

Пример на практике:

**ClickHouse over
Network NVME Disk
(Native Merge Tree)**

nvme disk,
нативный формат хранения

**Diskless ClickHouse
over Object Storage
(S3 table + raw parquete)**

табличная функция s3 table
и «сырой» parquete

Пример на практике: выполнение набора аналитических запросов

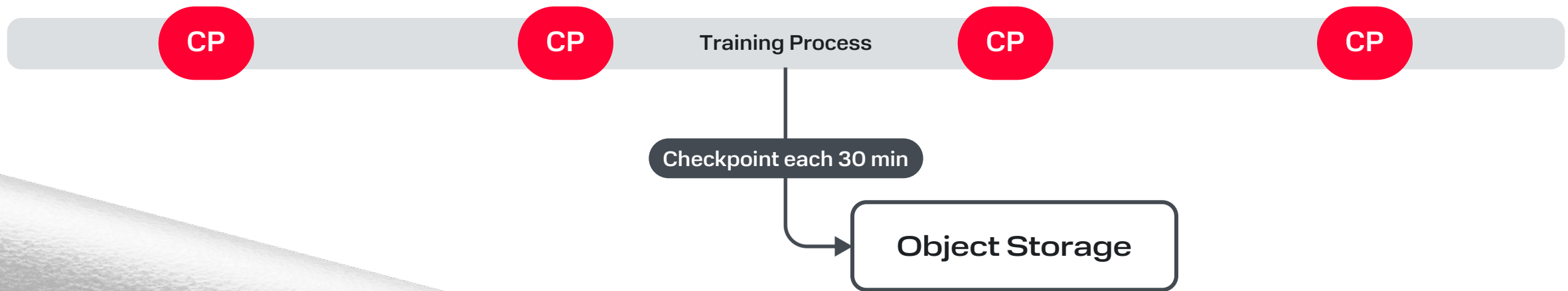
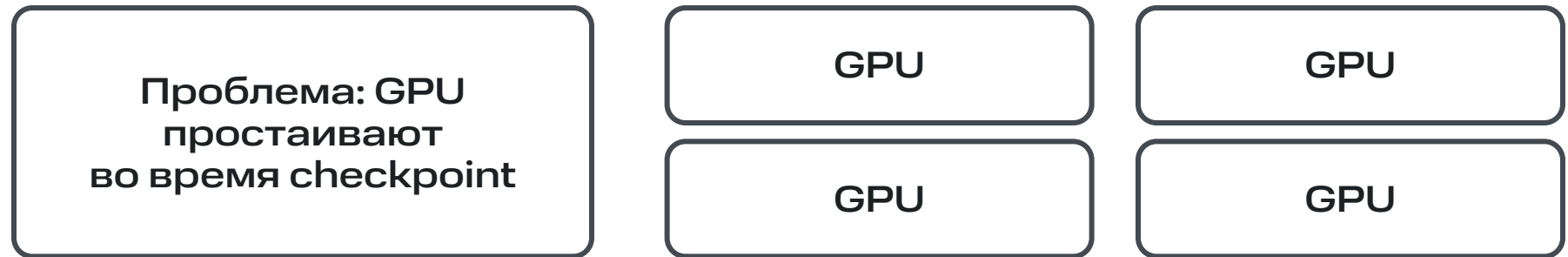
Сопоставимо по времени выполнения для аналитических запросов

В худших случаях overhead порядка ~1,5x

#	ЗАПРОС	NVME, C	OBJSTOR, C	РАЗНИЦА
1	GROUP BY UserID x SearchPhrase	1.818	1.656	x0.9
2	High-cardinality GROUP BY UserID	2.001	1.918	x1.0
3	Search + mobile / desktop split	1.009	1.246	x1.2
4	Top search phrases	0.577	0.733	x1.3
5	count(DISTINCT SearchPhrase)	0.567	0.779	x1.4
6	GROUP BY + uniqExact	0.583	0.842	x1.4
7	count(DISTINCT UserID)	0.404	0.636	x1.6
8	Subquery + session buckets	0.459	0.764	x1.7
9	LIKE '%google%' + domain()	0.938	2.117	x2.3
10	Hourly analytics за период	0.412	0.936	x2.3
11	Funnel + HAVING	0.318	0.868	x2.7
12	Daily stats по top-10 счётчикам	0.288	0.945	x3.3
13	Mobile phone models	0.185	0.765	x4.1



Другой пример: ML





Под капотом: инфраструктура

*MWS
Cloud
Platform*