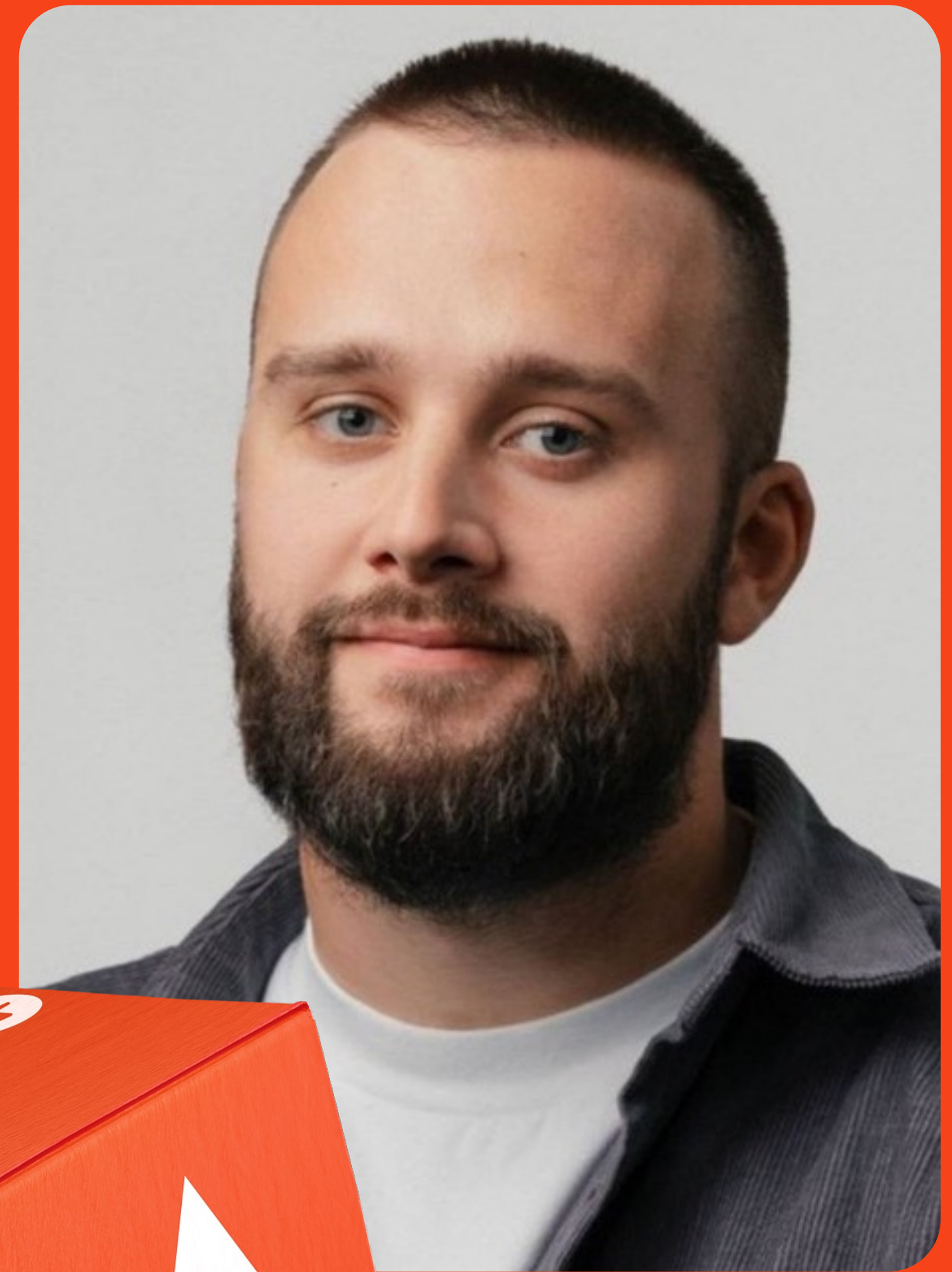


# критерии хорошей платформы данных

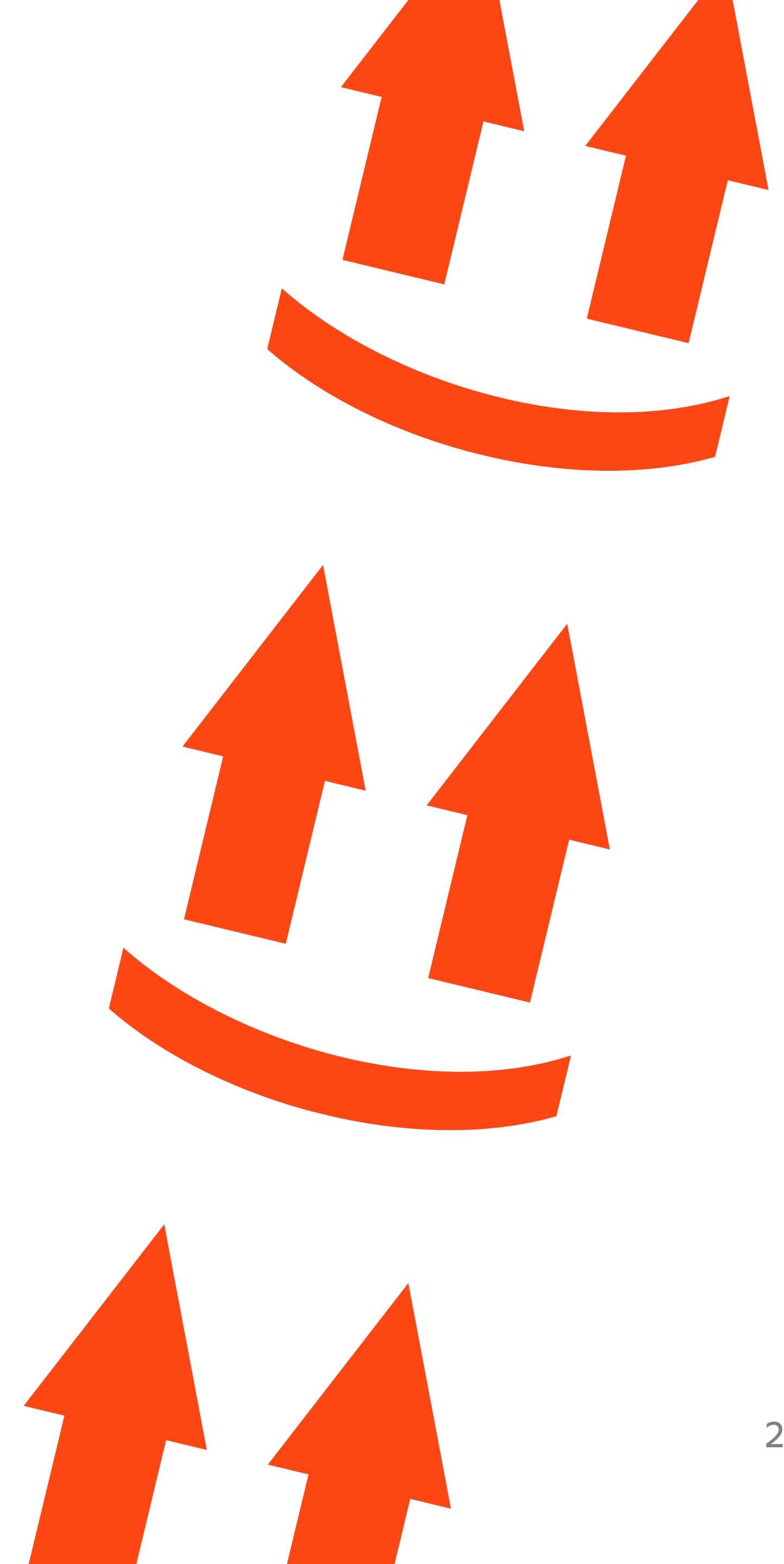
**Влад Гоцуляк**

Руководитель платформы данных в Яндекс Доставке



# ПОЗНАКОМИМСЯ!

- 01 Выпускник физтеха
- 02 Дата-инженерил в прошлом, теперь менеджер большой команды
- 03 Руководжу платформой данных в Яндекс Доставке
- 04 До этого работал CDO Eapteka, MLE в X5, DE в Sber
- 05 Люблю кайтсерфинг, сноуборд и вейкборд

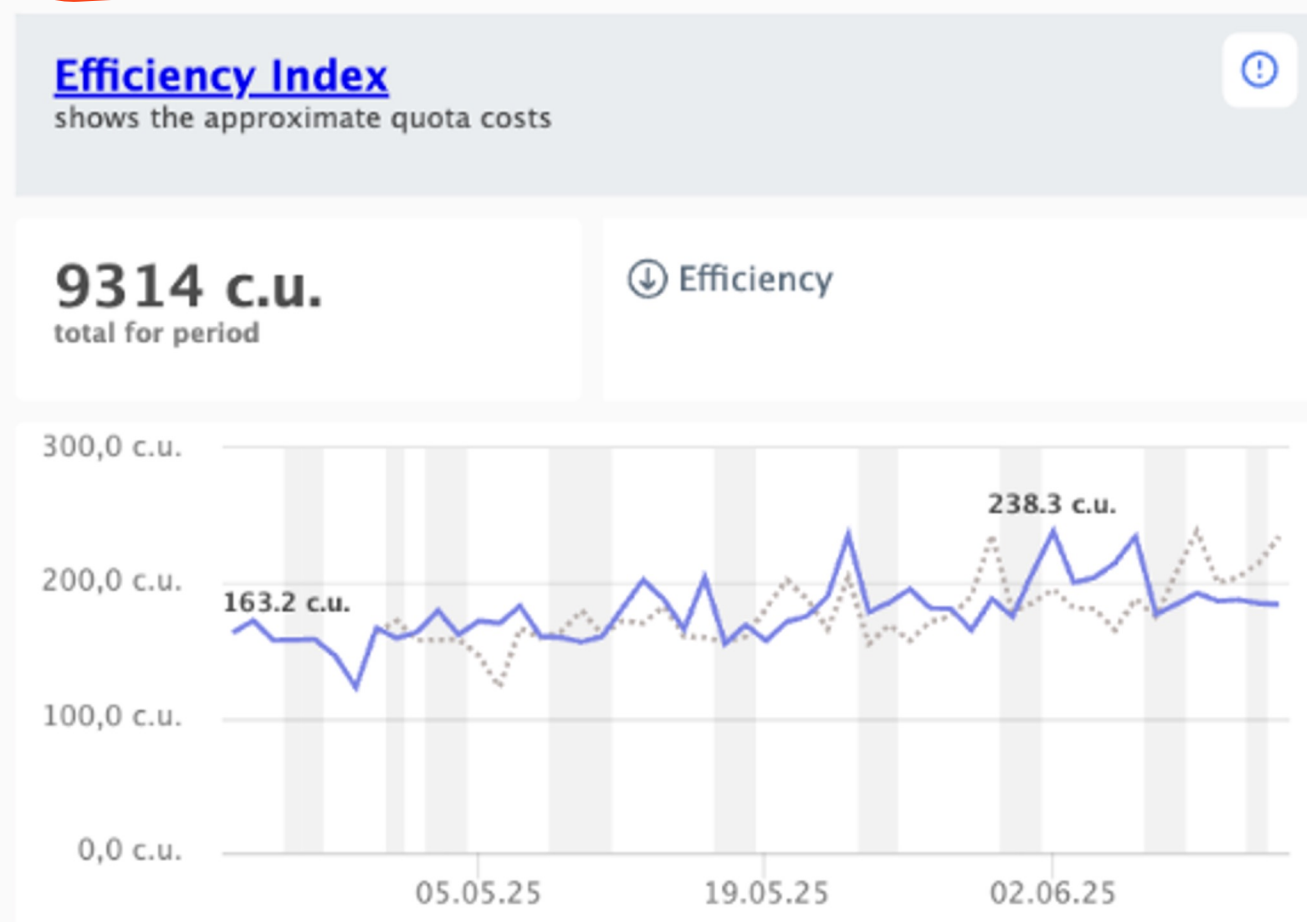
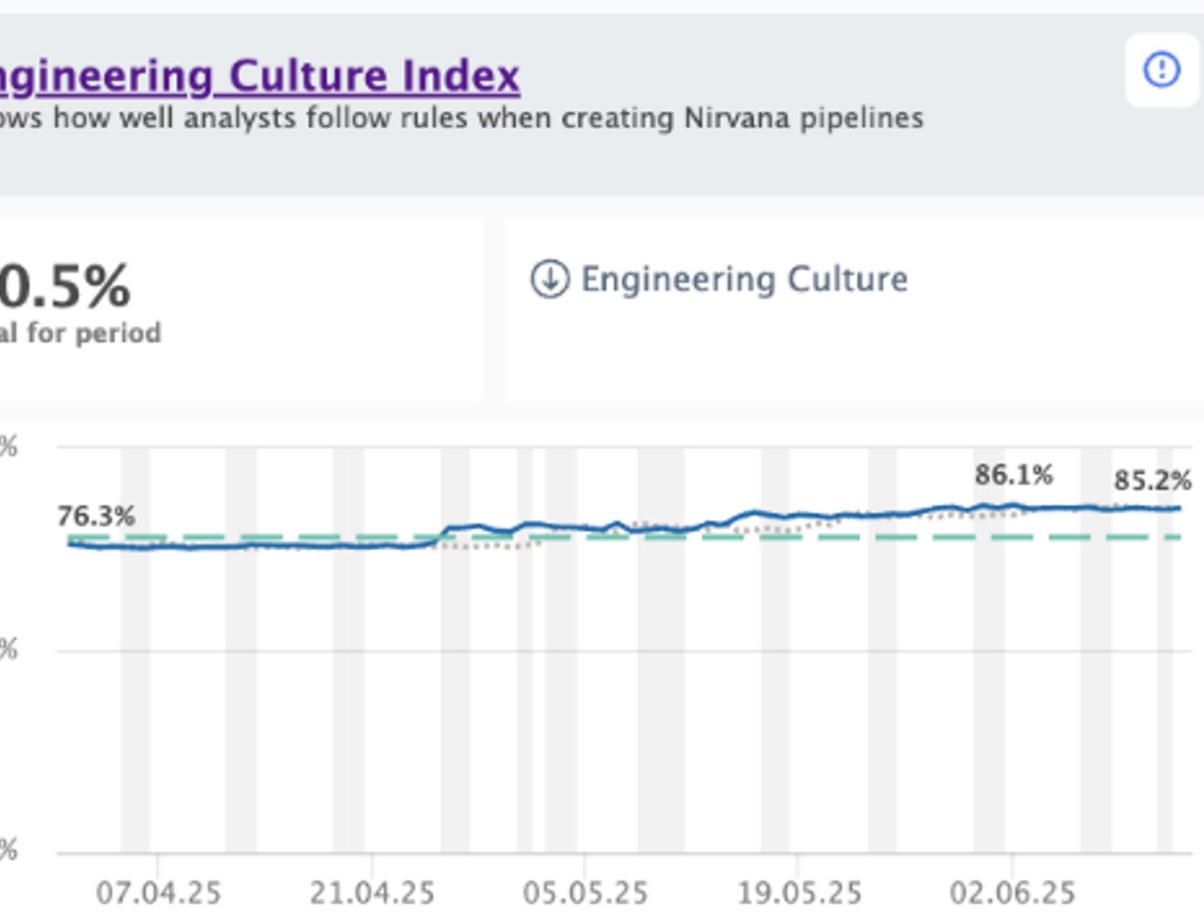
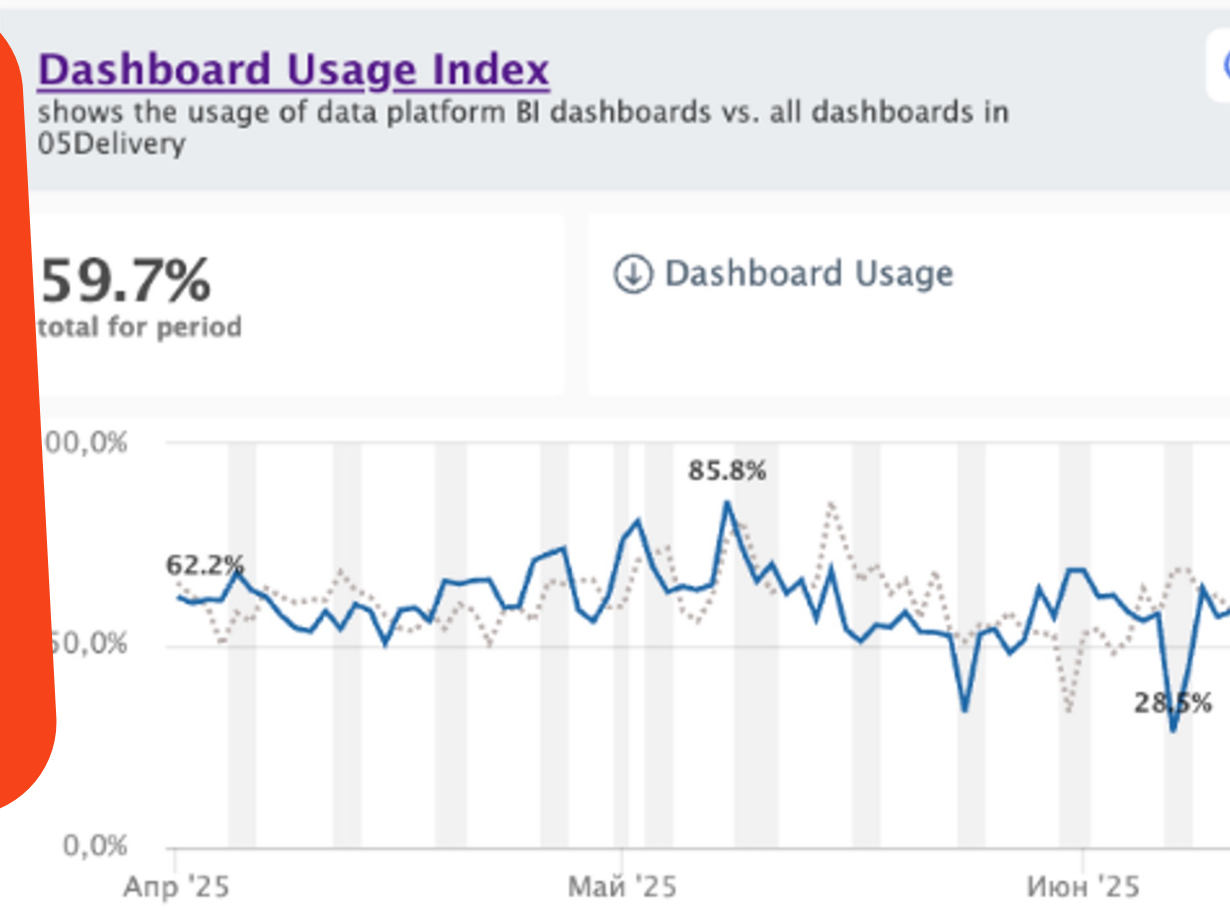
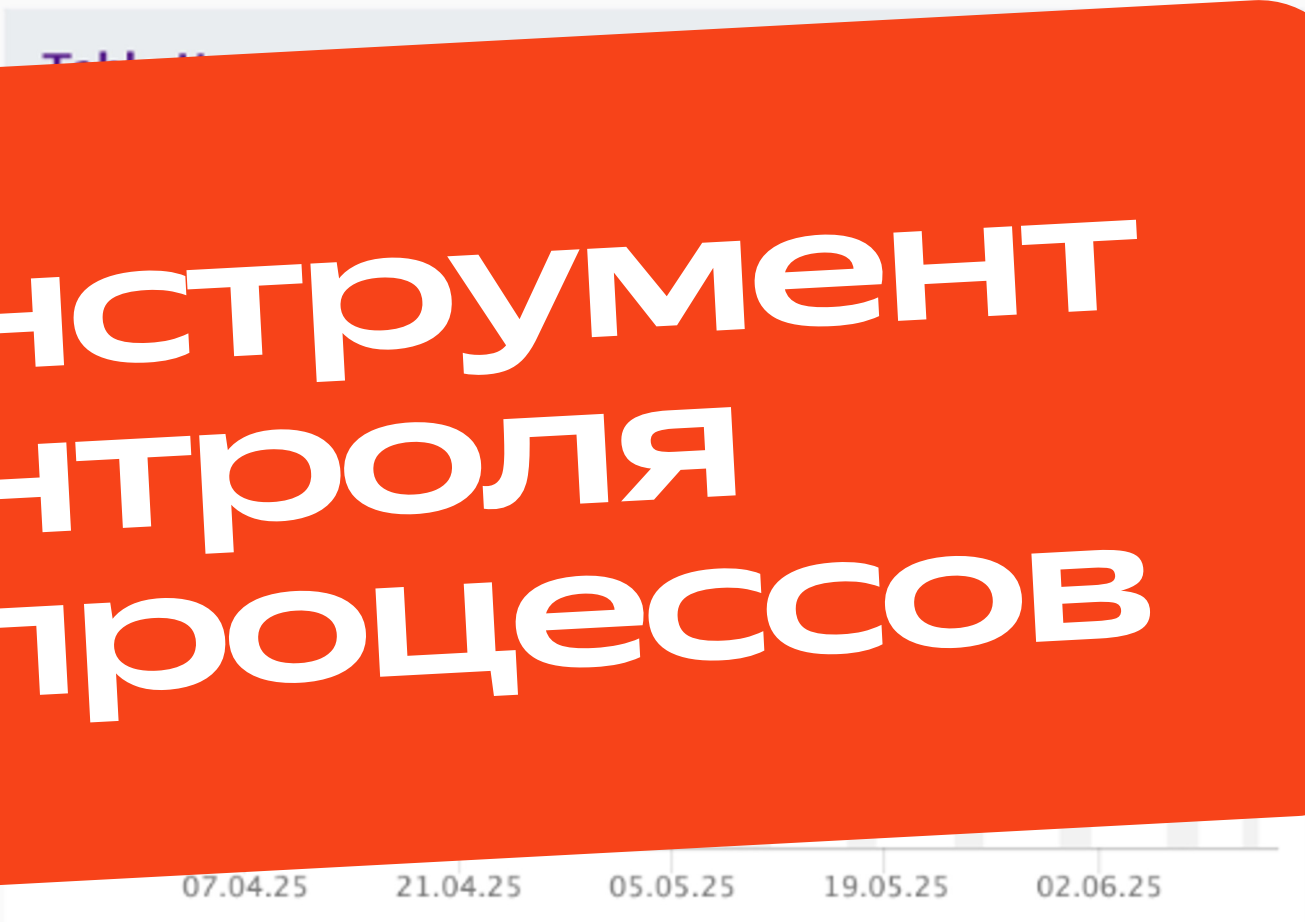
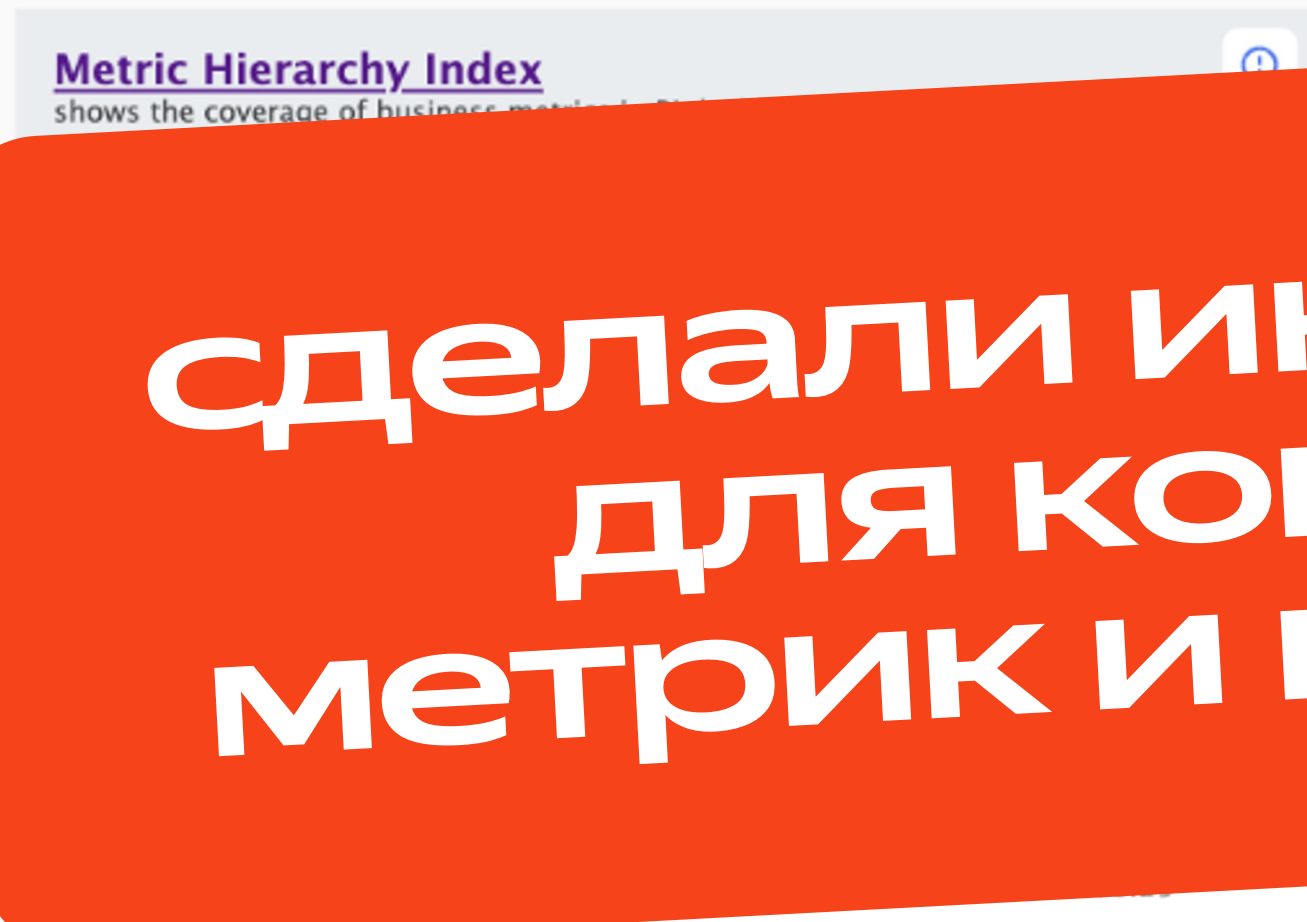
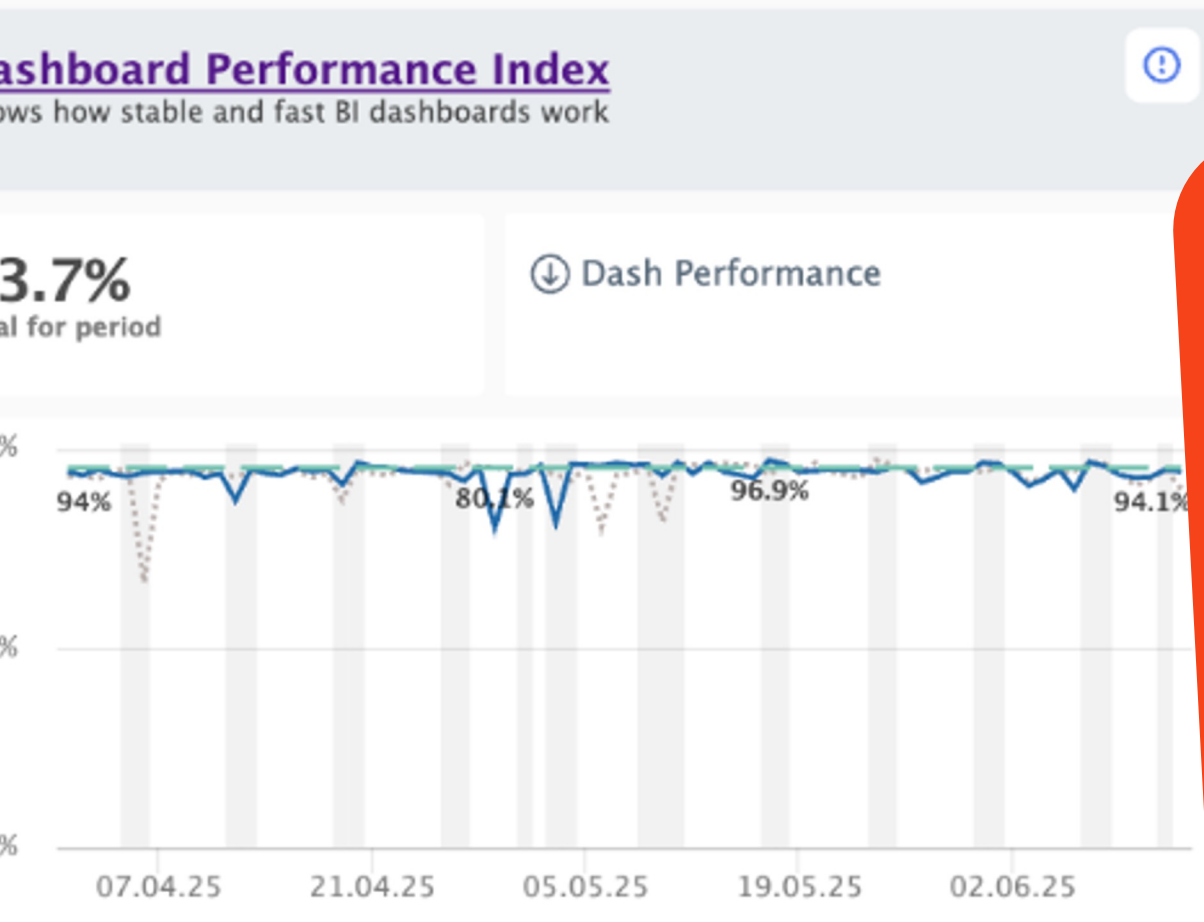
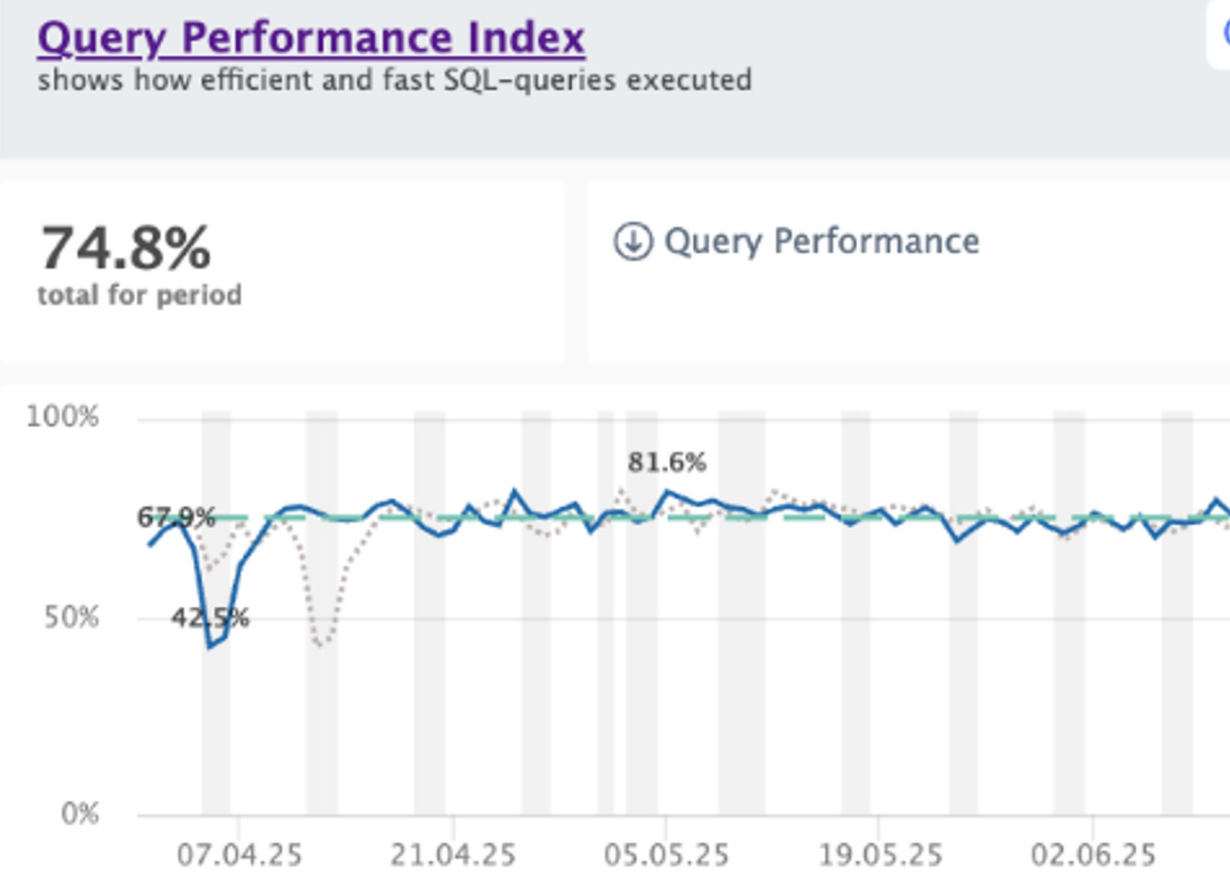
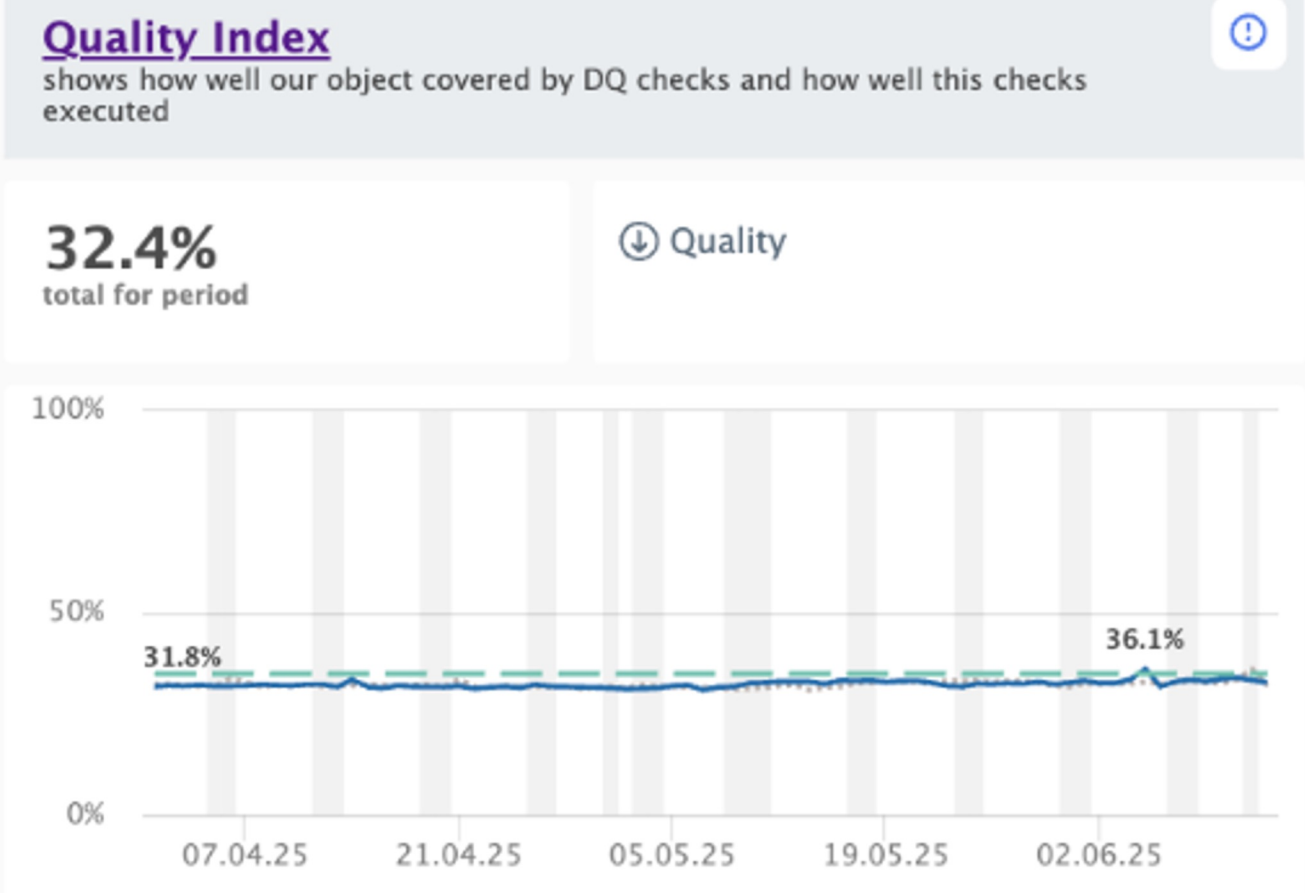
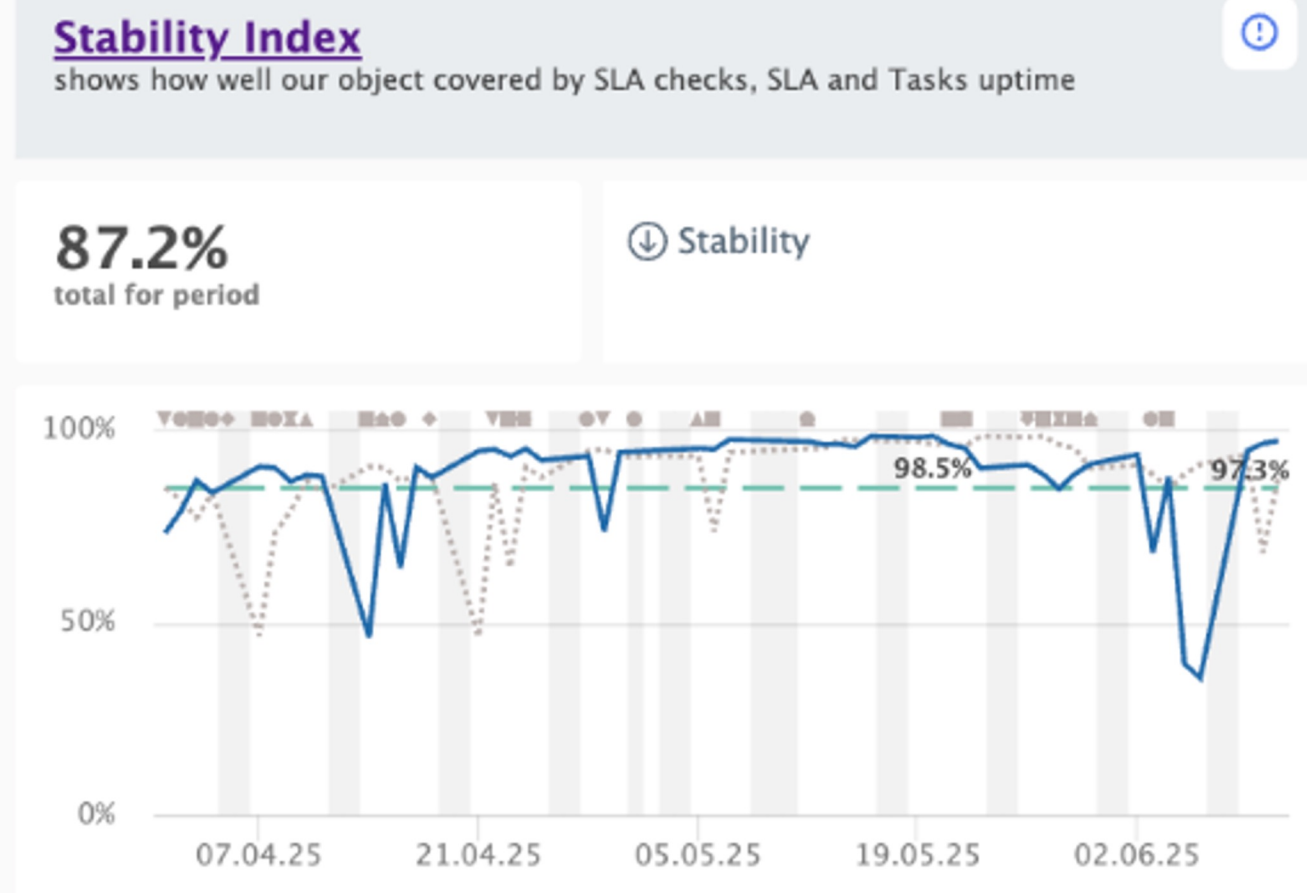
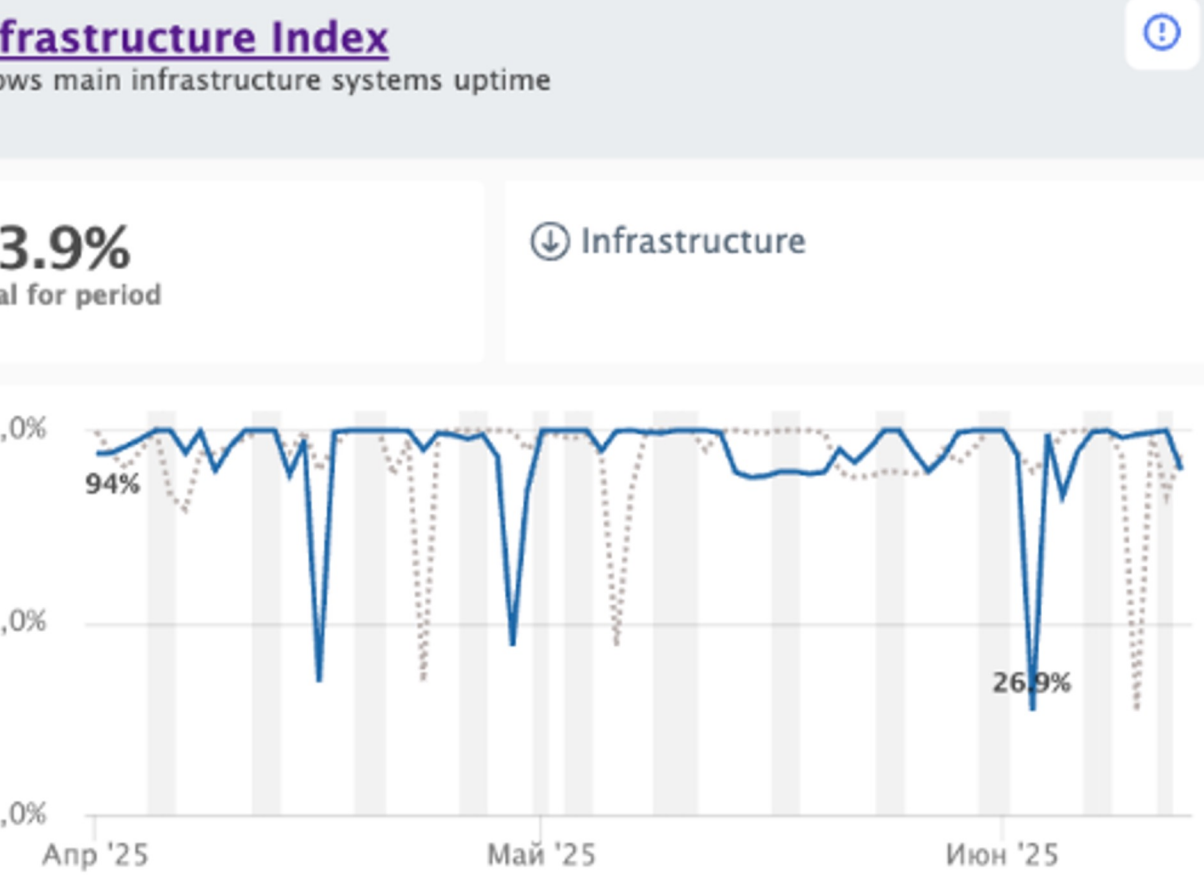


# О ЧЁМ СЕГОДНЯ ПОГОВОРИМ



# критерии хорошей платформы данных

- 01 Стабильность инфраструктуры и платформы
- 02 Качество данных
- 03 Скорость отчётов и аналитических запросов
- 04 Полнота покрытия бизнесовых метрик
- 05 Высокая частотность использования бизнесом и аналитикой
- 06 Эффективность использования железа
- 07 Инженерная культура аналитиков

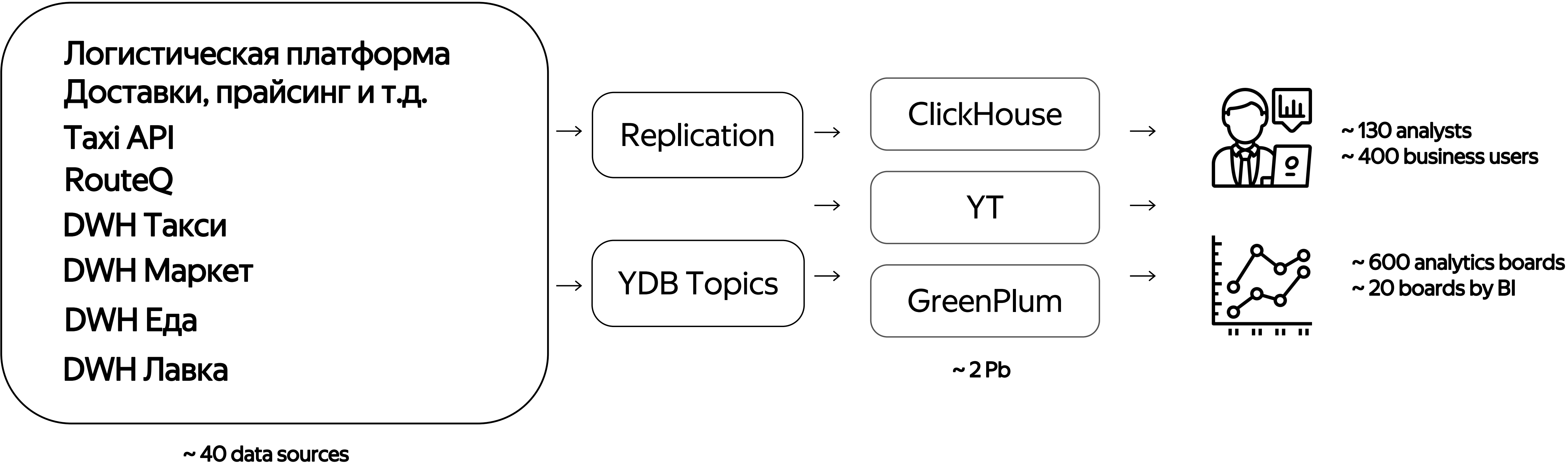


### Indexes by Domains

summary by domains, shows total values for selected period

DOMAIN	Stability Index	Quality Index	Query Performance Index	Dashboard Performance Index	Metrics Hierarchy Index	Tables Usage Index	Dashboard Usage Index	Engineering Culture Index	Efficiency Index c.u.
supply	89,2%	26,4%	61,7%	63,6%	66,1%	42,0%	66,6%	81,3%	739,4 c.u.
demand	88,1%	39,0%	67,7%	63,9%	32,0%	48,4%	18,9%	89,5%	1 817,6 c.u.
finance	87,9%	47,8%	81,3%	82,6%	76,6%	48,7%	56,0%	68,8%	638,4 c.u.
efficiency	87,8%	19,7%	76,2%	62,8%	53,8%	35,9%	18,8%	71,0%	3 490,4 c.u.
support	60,0%	31,0%	85,4%	76,5%	0,0%	44,7%	0,0%	null	1 180,5 c.u.
bets	80,9%	51,0%	71,4%	68,8%	3,6%	46,6%	0,0%	75,8%	1 429,6 c.u.
core	88,5%	26,9%	null	92,2%	0	null	98,7%	90,2%	18,7 c.u.

# внедрили для архитектуры всей платформы Яндекс Доставка

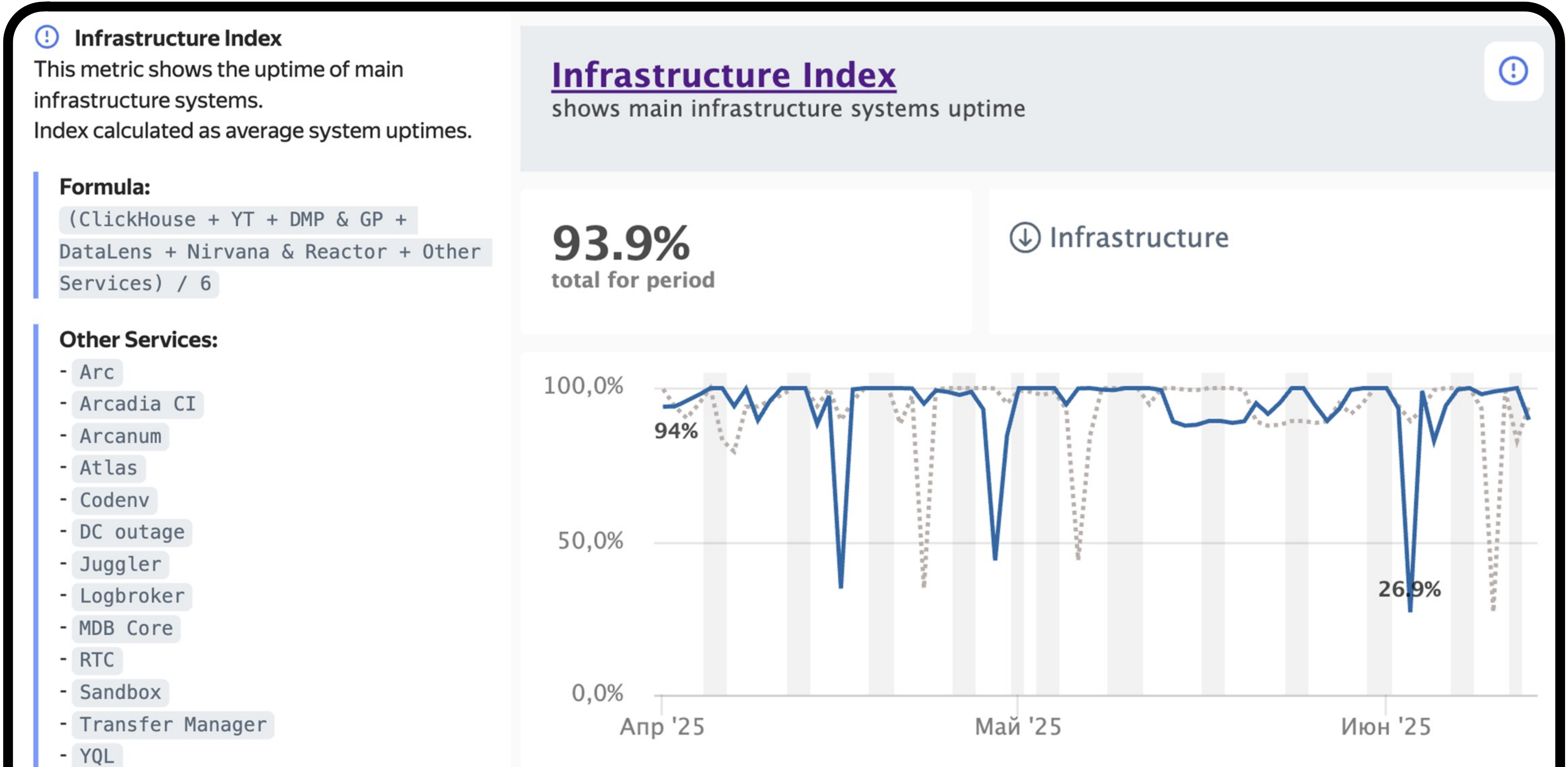




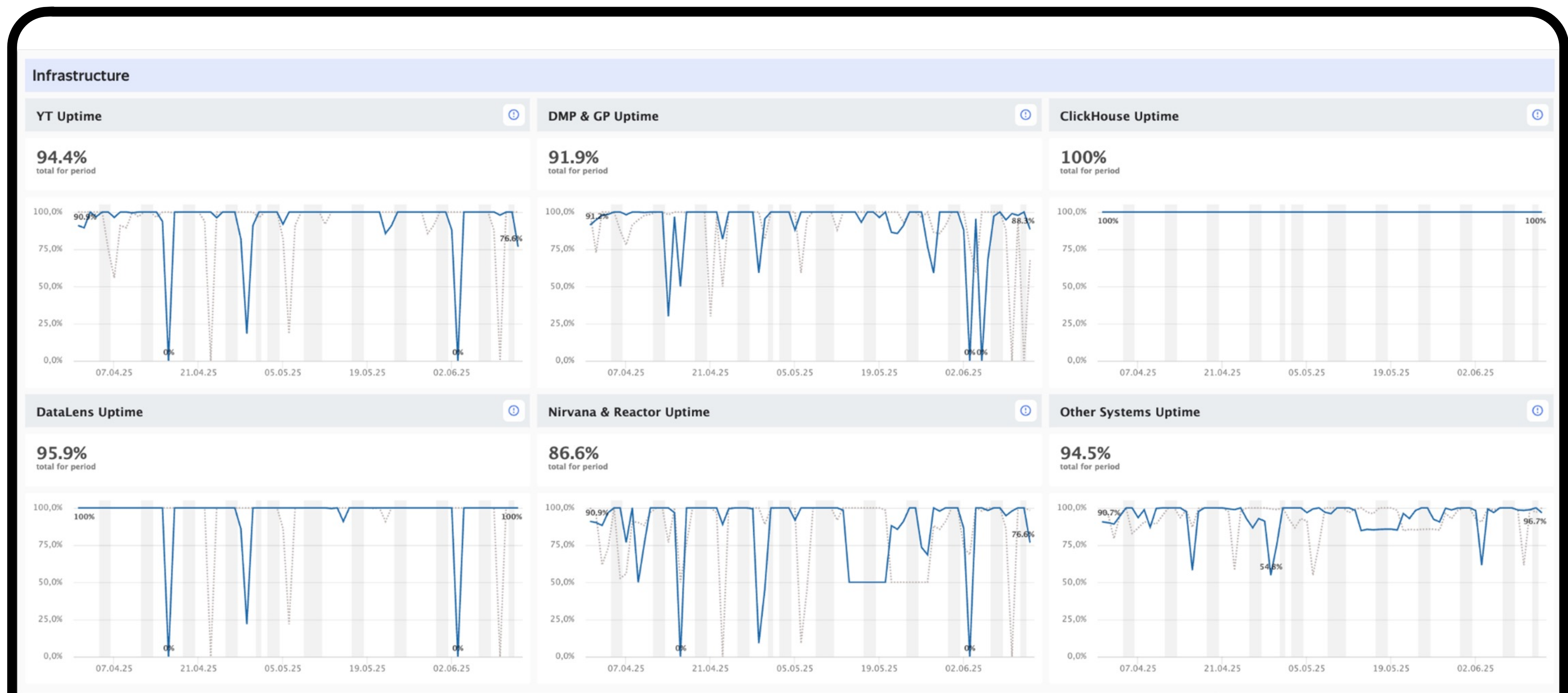
01

# метрики стабильности инфраструктуры

# следим за стабильностью нашей инфраструктуры через Infrastructure Index



# и делаем это в разрезе каждого сервиса, используемого в платформе данных Доставки



# как создать для своей платформы метрику `infrastructure Index`

- 01 Определи ключевые архитектурные компоненты своей системы
- 02 Подключи метрики аптайма/health-check для каждого компонента  
(Prometheus exporters, встроенные метрики из логов систем)
- 03 Сохраняй значения в витрину `infrastructure_uptime`  
(Service, day, uptime\_pct)
- 04 Считай индекс как средний uptime по всем сервисам за день  
(HDFS + GP + Kafka + Spark + Airflow + Metabase)
- 05 Усредняй за период -> Infrastructure Index
- 06 В дашборде (Metabase/Superset) строй график тренда и выделяй проблемные сервисы с командой



02

# метрики стабильности платформы данных (SLA)

# следим за выполнением гарантий обновления через Stability Index

## ! Stability Index

This metric shows the timeliness and success rate of object updates. It represents how well our objects are covered by SLA checks, whether our objects are updated within the established SLA, and the success rate of our data processing and update pipelines.

Index calculated as AVG stability by each table-day.

### Formula:

Share of time when SLA and Tasks was  
OK \* has SLA flag

*\*Index calculated for CDM,REP.  
Selector-like tables not included.*

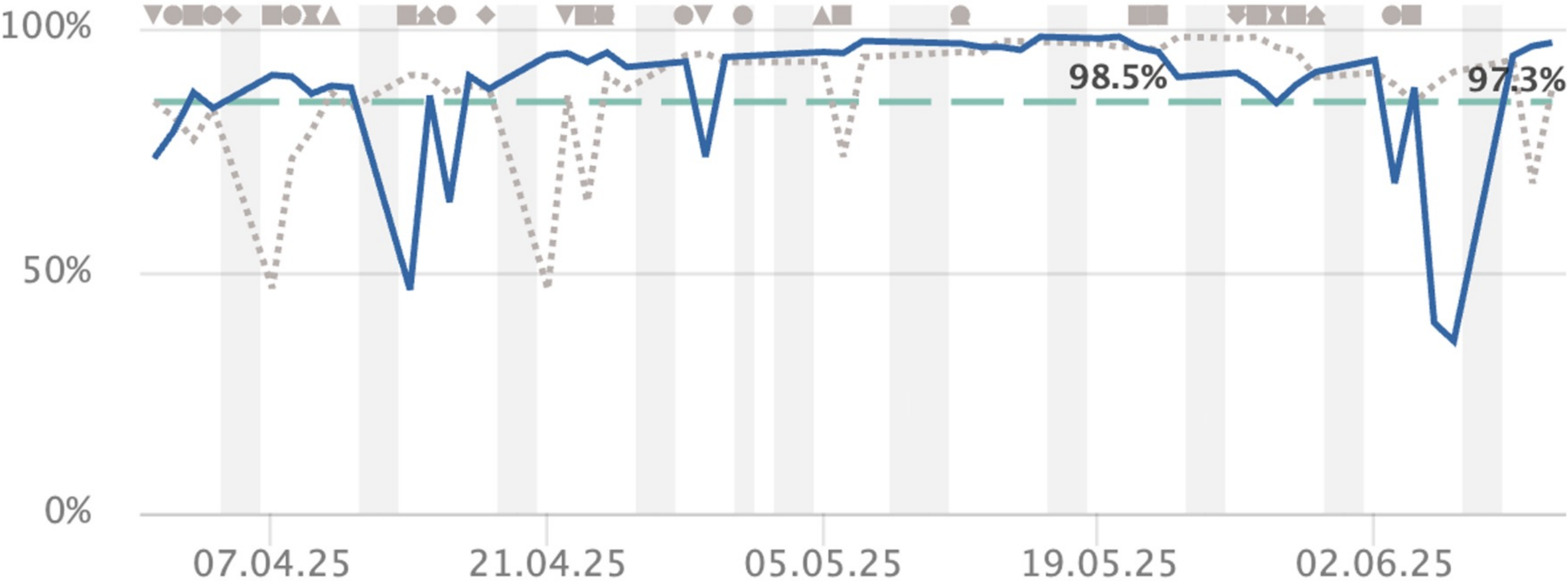
## Stability Index

shows how well our object covered by SLA checks, SLA and Tasks uptime

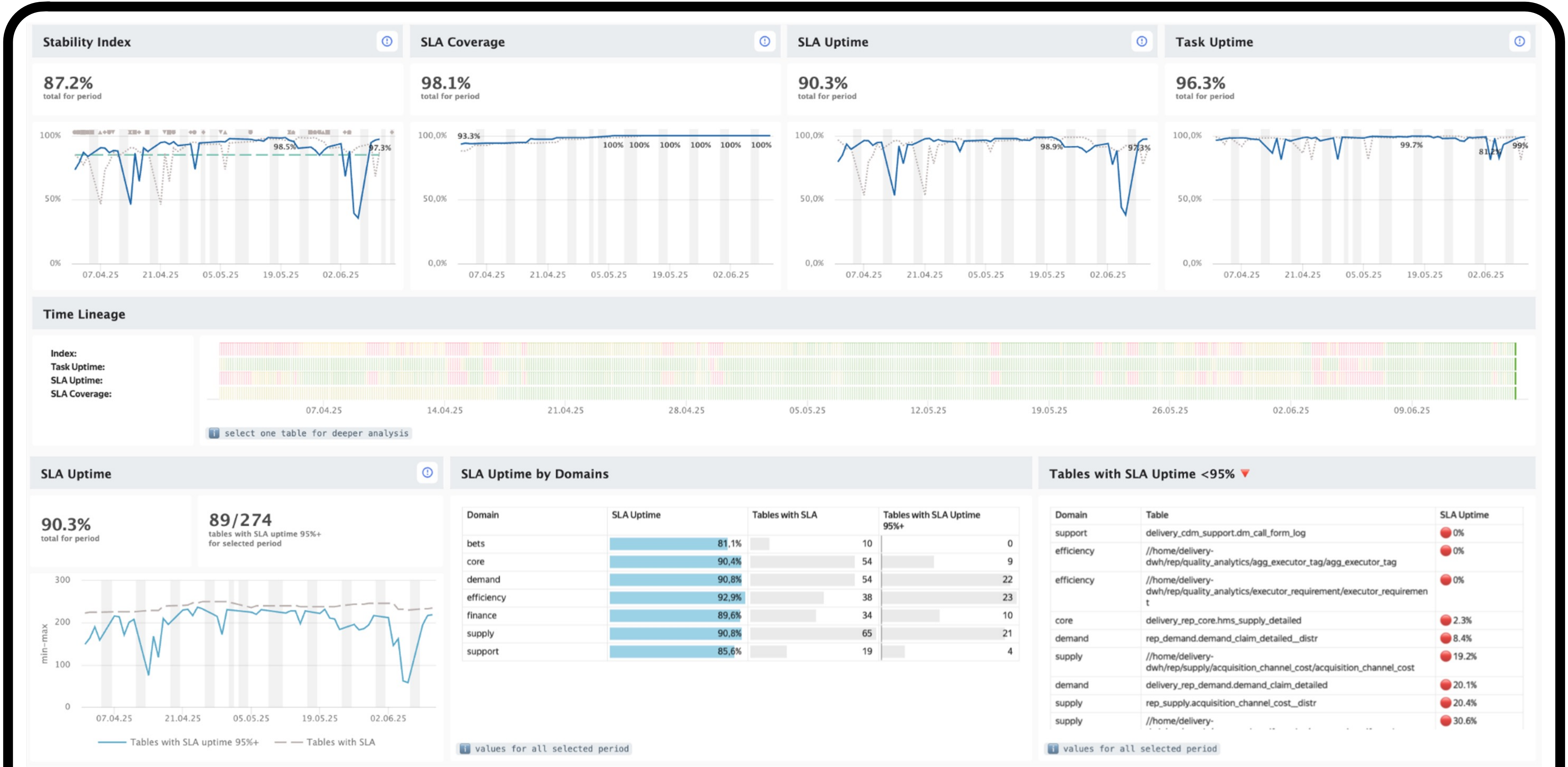
87.2%

total for period

⬇ Stability



# ключевой фокус в стабильности платформы — это прогнозируемое время обновления данных (SLA)



# как создать для своей платформы метрику **Stability Index**

- 01 Определи SLA-контракты с бизнесом и аналитикой
- 02 Собери данные SLA и аптайма витрин из оркестраторов (Airflow + DBT, Cron, DMP)
- 03 Для таблиц с `has_sla = true` фиксируй периоды, когда обновления выполнялись в срок и успешно
- 04 Считай долю времени, когда SLA и задачи были OK (`Stability = share_of_time(SLA OK & Tasks OK)`)
- 05  $SI = AVG(Stability)$  по каждой таблице-дню -> Stability Index
- 06 В BI отобрази тренд индекса и список таблиц с отклонением от SLA



03

# метрики качества данных

# следим за качеством данных в нашей платформе через

# Quality Index

## ! Quality Index

This metric shows the overall state of data quality. It represents how well our objects are covered by DQ checks and how successfully our objects pass these checks.

Index calculated as AVG value by each table-day.

### Formula per 1 table:

```
SUM (DQ Existing Measure Weight * DQ  
Check OK)
```

**For example:** table\_1 has 2 DQ checks

1 check weight = 0.2, check ok = 95%;

2 check weight = 0.05, check ok = 80%

**Index** =  $(0.2 * 0.95 + 0.05 * 0.8) = 23\%$

\* more about check weight read in DQ  
Coverage description.

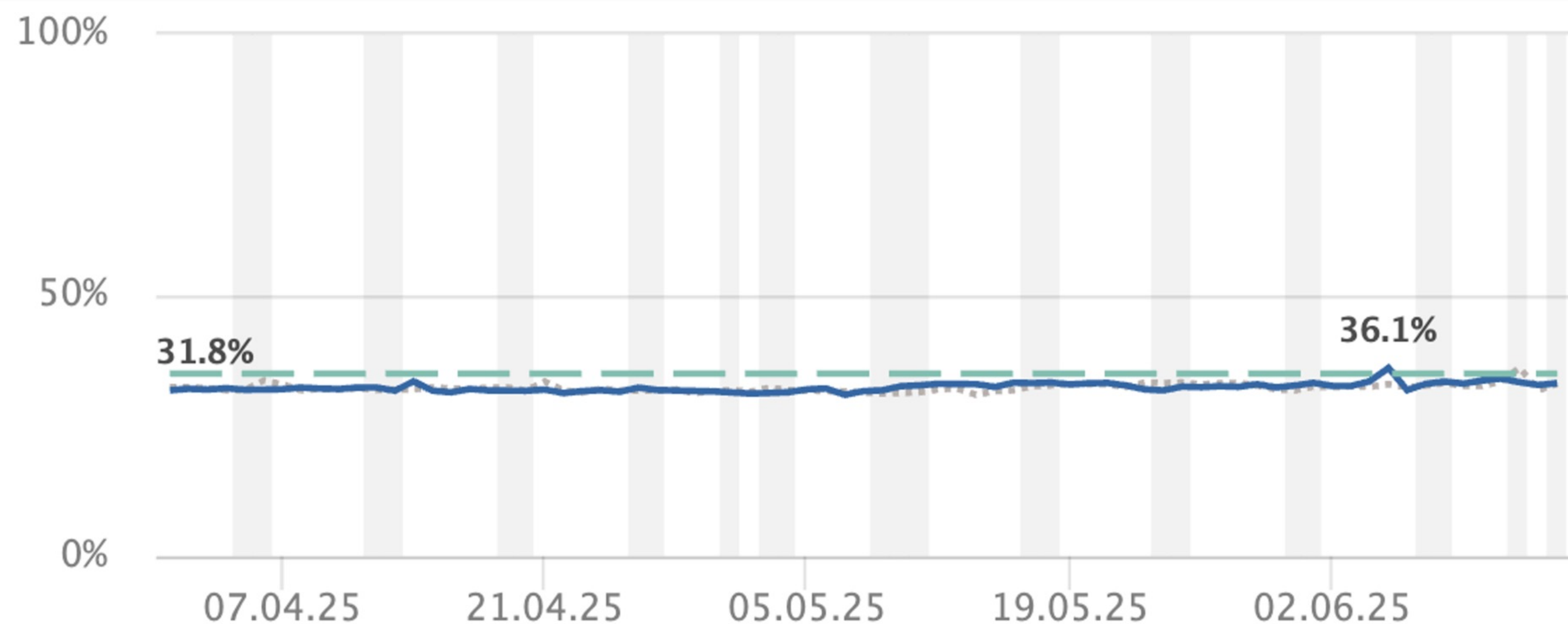
\* Index calculated for CDM, REP.  
Selector-like tables not included.  
Replicated tables not included

## Quality Index

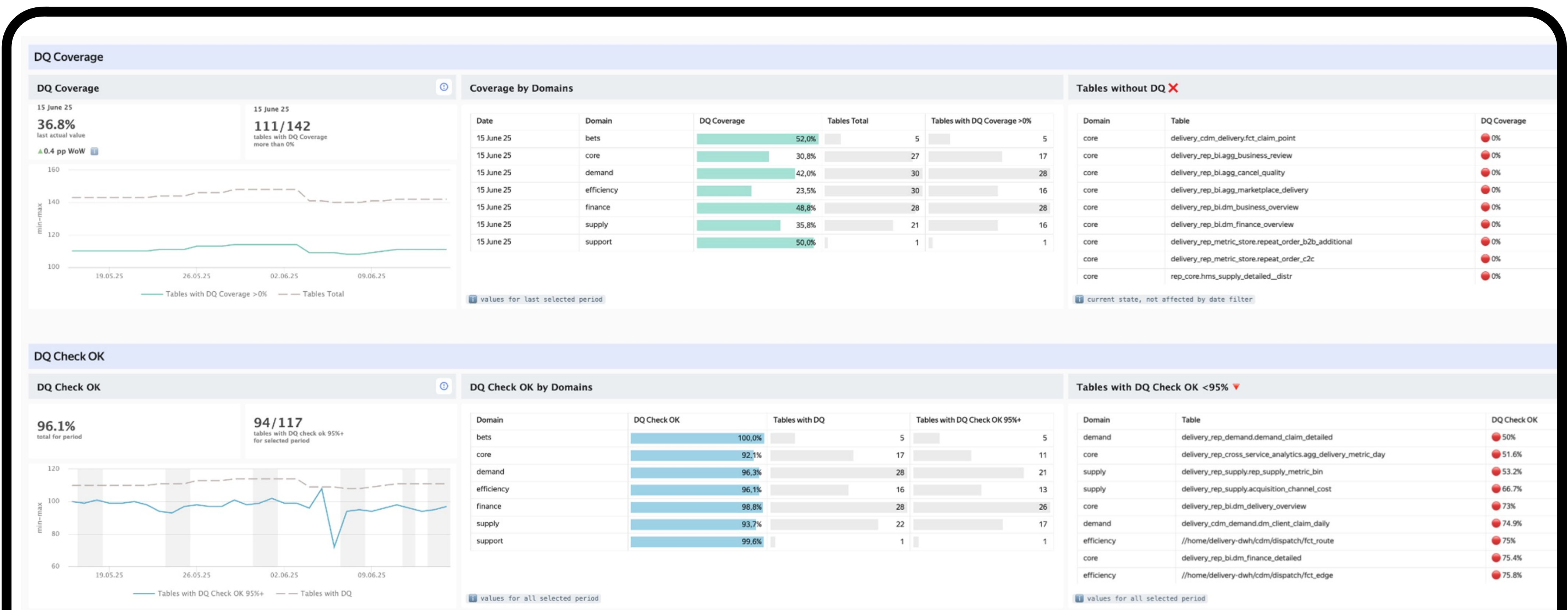
shows how well our object covered by DQ checks and how well this checks executed

**32.4%**  
total for period

↓ Quality



# С ПОМОЩЬЮ МЕТРИК КАЧЕСТВА ДАННЫХ МЫ ПРОЗРАЧНО УПРАВЛЯЕМ ЗАДАЧАМИ НА ИНЖЕНЕРОВ В ДОМЕНАХ



# как создать для своей платформы метрику Quality Index

- 01 Установи Great Expectations и подключи свой DWH/S3 как datasource
- 02 Создай expectation suites для таблиц, укажи вес проверки в `meta.dq_weight`
- 03 Настрой checkpoint для запуска проверок
- 04 Сохраняй результаты в SQL-таблицу `dq_results` (`check`, `вес`, `success_pct`)
- 05 Считай индекс:  $QI = \text{SUM}(\text{weight} \times \text{success\_pct})$  по таблице-дню, потом среднее
- 06 Собери дашборд: Quality Index, Coverage, Check OK, проблемные таблицы
- 07 Следи в динамике за этой метрикой и обсуждай с командой в разрезе доменов на регулярной основе



04

# метрики скорости отчётов

следим за скоростью отчетов платформы  
данных в DataLens через метрику

# Dashboard Performance Index

## ! Dashboard Performance Index

This metric shows the share of fast and successful events during BI production dashboard usage by users.

### Formula:

Fast and Successful Events/ Total  
Events

\*Event - one chart load, time for query  
execution in database

## Dashboard Performance Index

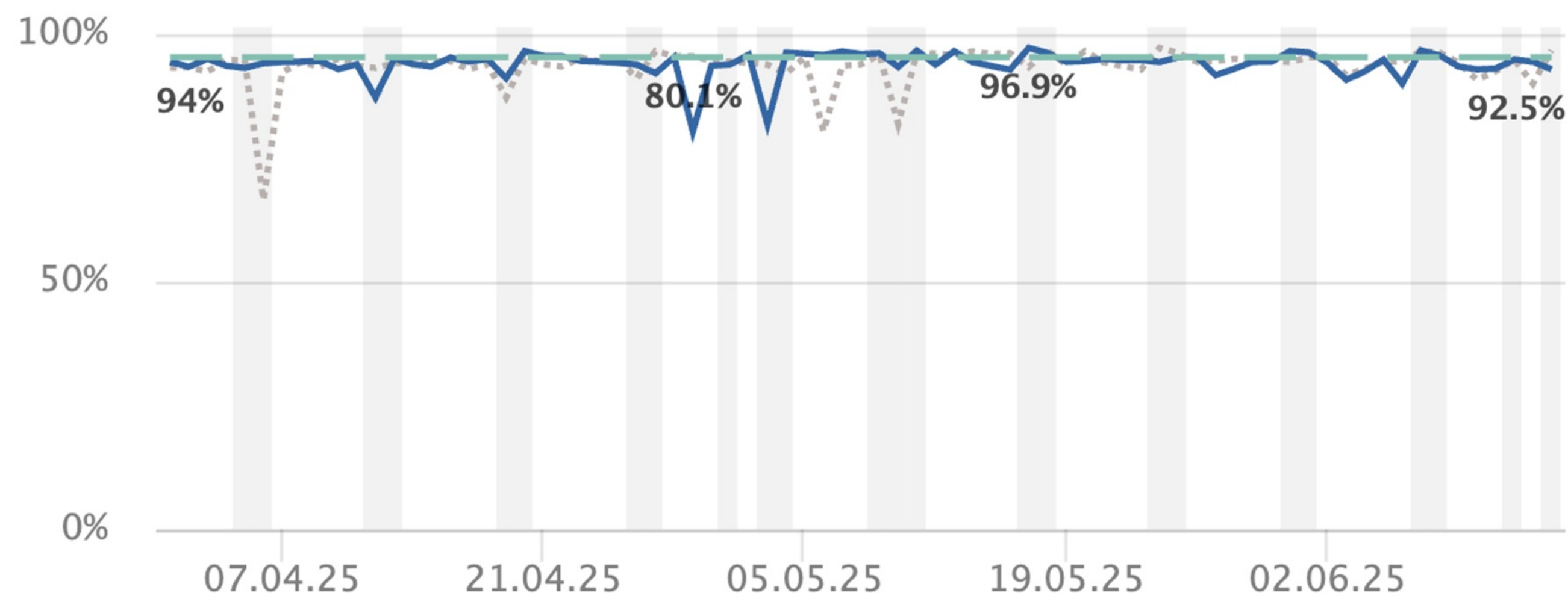
shows how stable and fast BI dashboards work



93.6%

total for period

⬇ Dash Performance



# как создать для своей платформы метрику Dashboard Performance Index

- 01 Настрой в своей BI системе логирование запросов
- 02 Собери витрину событий: `event_id`, `dashboard_id`, `chart_id`, `start_ts`, `end_ts`, `success_flag`, `type`
- 03 Задай пороги по типам: селекторы  $\leq 0.5-3s$ , чарты  $\leq 1-5s$  (определи с командой свои)
- 04 Рассчитай метрики по дню дашборду: `total_events`, `fast_events`, `success_events`  
( $Fast\% = fast\_events / total\_events$ ,  $Success\% = success\_events / total\_events$ )
- 05 Считай индекс:  $DPI = AVG (Fast\% , Success\%)$
- 06 Сохраняй в `dashboard_perf_index` (`dashboard_id`, `day`, `fast_pct`, `success_pct`, `dpi`)
- 07 В BI собери дашборд: общий DPI, тренд Fast/Success, список проблемных дашбордов (низкий Fast%)
- 08 Следи в динамике за этой метрикой и обсуждай с командой в разрезе доменов на регулярной основе



05

# метрики скорости аналитических запросов

следим за скоростью аналитических запросов  
в разрезе доменов через метрику

# Query Performance Index

## ! Query Performance Index

This index shows how fast and effective analytical scripts executed during ad-hocs or regular runs.

### Formula:

Fast and Effective Queries / Total  
Queries

## Query Performance Index

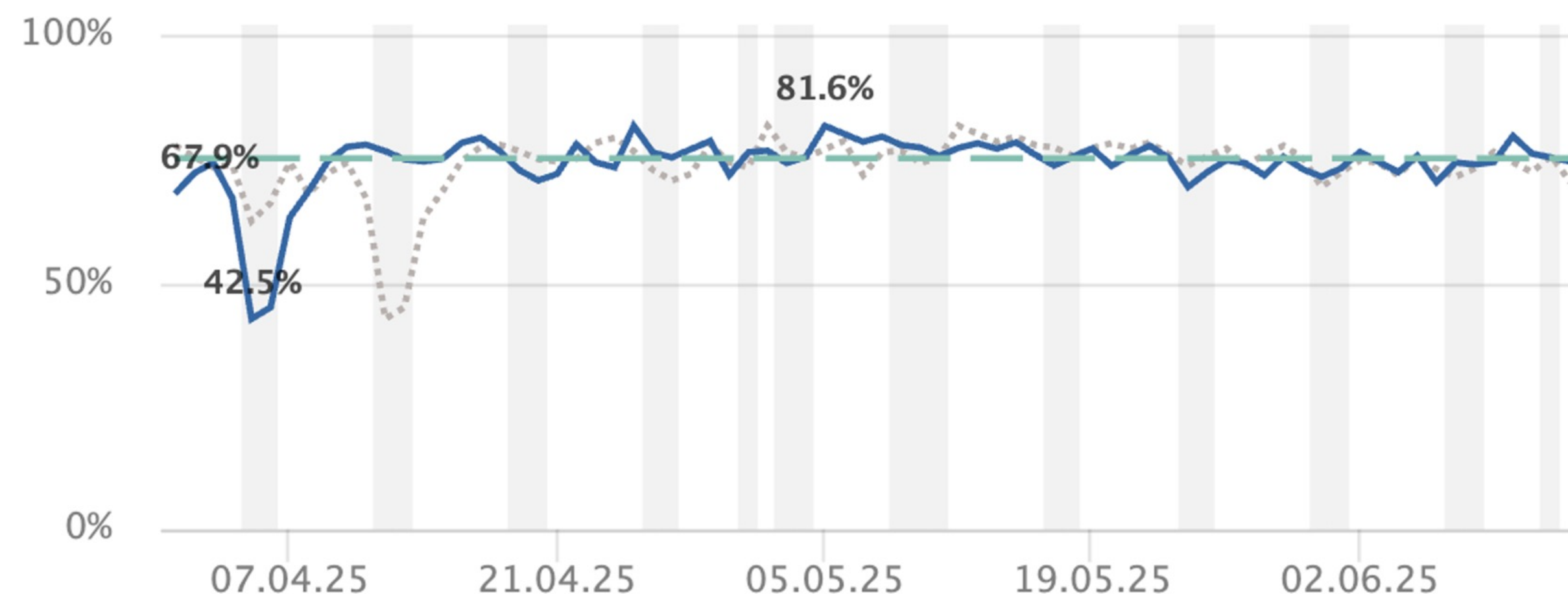
shows how efficient and fast SQL-queries executed



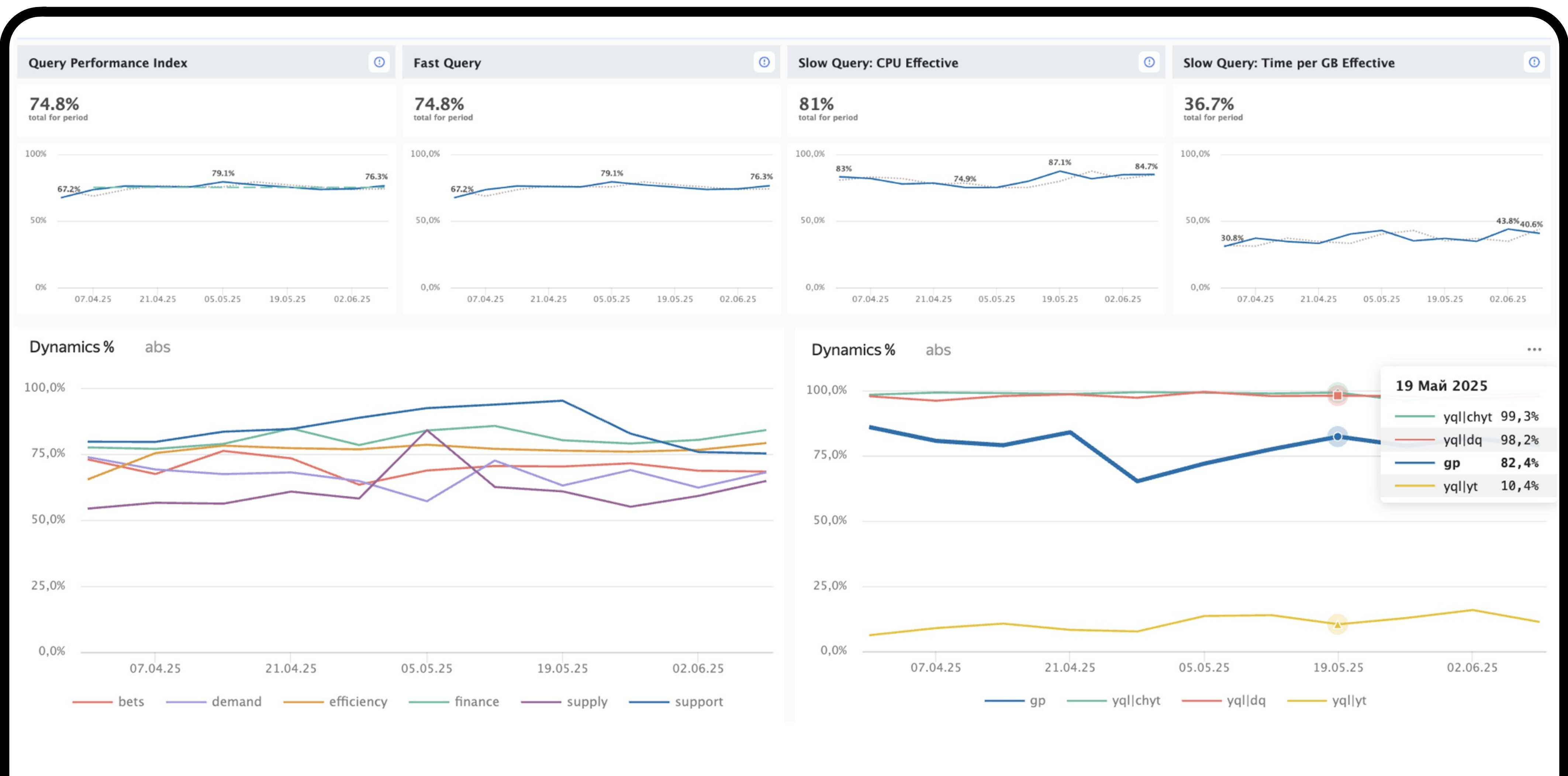
74.8%

total for period

↓ Query Performance



# С ПОМОЩЬЮ метрик скорости запросов выстраиваем линии поддержки аналитиков в каждом домене



# как создать для своей платформы метрику Query Performance Index

- 01 Собирай на основе логов все пользовательские запросы в разрезе доменов/команд
- 02 Добавь метрики:  $\text{duration} = \text{end\_ts} - \text{start\_ts}$ ,  $\text{success\_flag}$
- 03 Задай правила Fast and Effective  
(например:  $\text{runtime} \leq 30\text{s}$  для ad-hoc,  $\leq 5\text{m}$  для batch, статус=успех)
- 04 Считай показатели по домену/дню : TotalQueries, FastEffectiveQueries, SystemType
- 05  $\text{QPI} = \text{FastEffectiveQueries} / \text{TotalQueries}$
- 06 Собери дашборд: Query Performance Index, Check OK, проблемные таблицы
- 07 Следи в динамике за этой метрикой и обсуждай с командой в разрезе доменов на регулярной основе



06

# метрики использования продуктов платформы данных бизнесом и аналитикой

следим за использованием отчётов в разрезе  
пользователей и доменов через метрику

# Dashboard Usage Index

## ! Dashboard Usage Index

This index shows proportion of weighted views to production dashboards created by BI team vs. all other dashboards in 05 Delivery.

### Formula:

BI Dashboard Weighted Views / Total  
Weighted Views

### Weighted views - view \* coef:

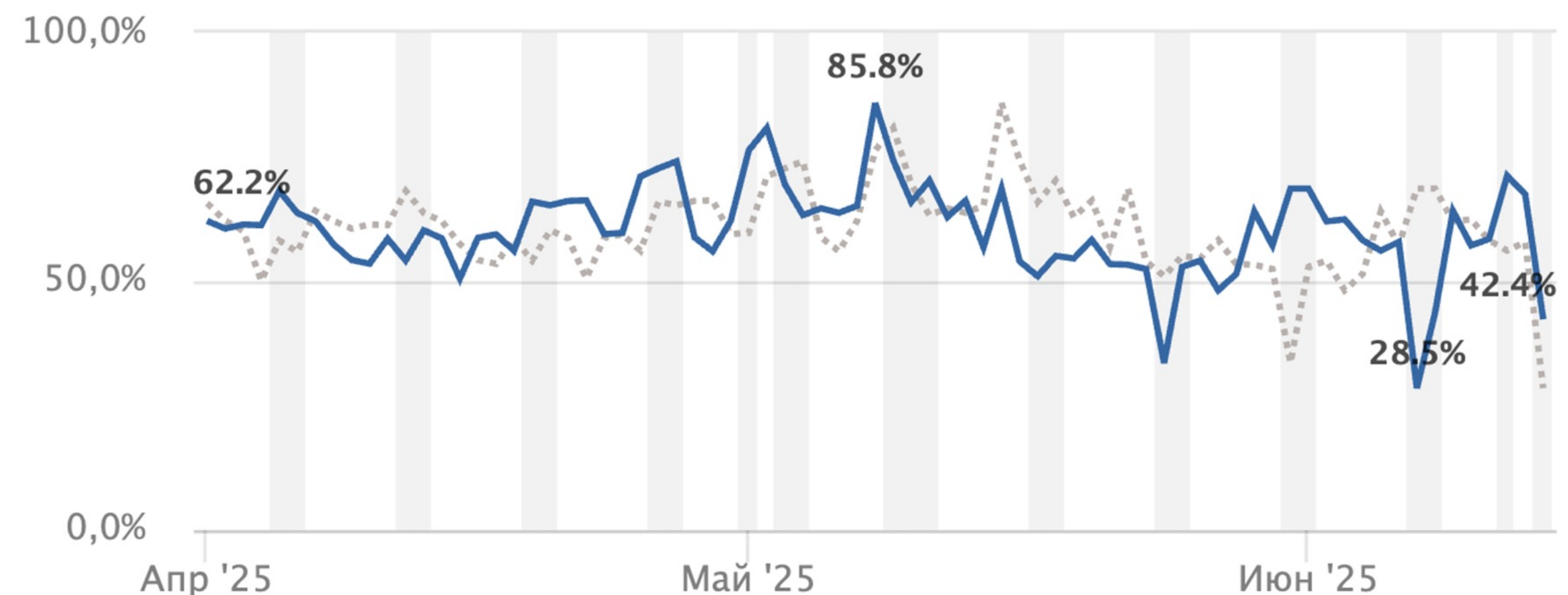
- CEO = 100
- CEO -1 = 50
- CEO -2 = 10
- Leads = 3
- Worker = 1
- Support = 0.3
- Robots = 0

## Dashboard Usage Index

shows the usage of data platform BI dashboards vs. all dashboards in 05Delivery

**59.7%**  
total for period

⬇ Dashboard Usage



# как создать для своей платформы метрику Dashboard Usage Index

- 01 Собери логи просмотров дашбордов (Metabase, Superset, PowerBI и так далее):  
user\_id, dashboard\_id, view\_ts
- 02 Добавь веса ролей  
(пример: CEO=100, Head=50, Lead=10, Analyst=3, Worker=1, Support=0.3, Robots=0)
- 03 Посчитай Weighted Views: views × weight для каждой роли
- 04 Выбери те дашборды, которые достойны внимания и пользуются бизнесом
- 05  $DUI = BI\_Dashboard\_Weighted\_Views / Total\_Weighted\_Views$
- 06 Сохраняй в витрину dashboard\_usage\_index (day, dui) и обновляй эти данные в t-1
- 07 Следи в динамике за этой метрикой и обсуждай с командой в разрезе доменов на регулярной основе

следим за использованием таблиц в разрезе  
пользователей и доменов через метрику

# Table Usage Index

## ! Table Usage Index

This metric shows the usage of DWH objects by analysts. It reflects the proportion of DWH tables used in analysts' work compared to all other objects, the share of active users, and a comparison of our DWH usage with DWHs from other business units.

Index calculated as weighted sum of 3 main measures.

### Formula:

Delivery DWH Usage x0.6 + Active  
Users% x0.3 + Delivery vs Other DWH  
Usage x0.1

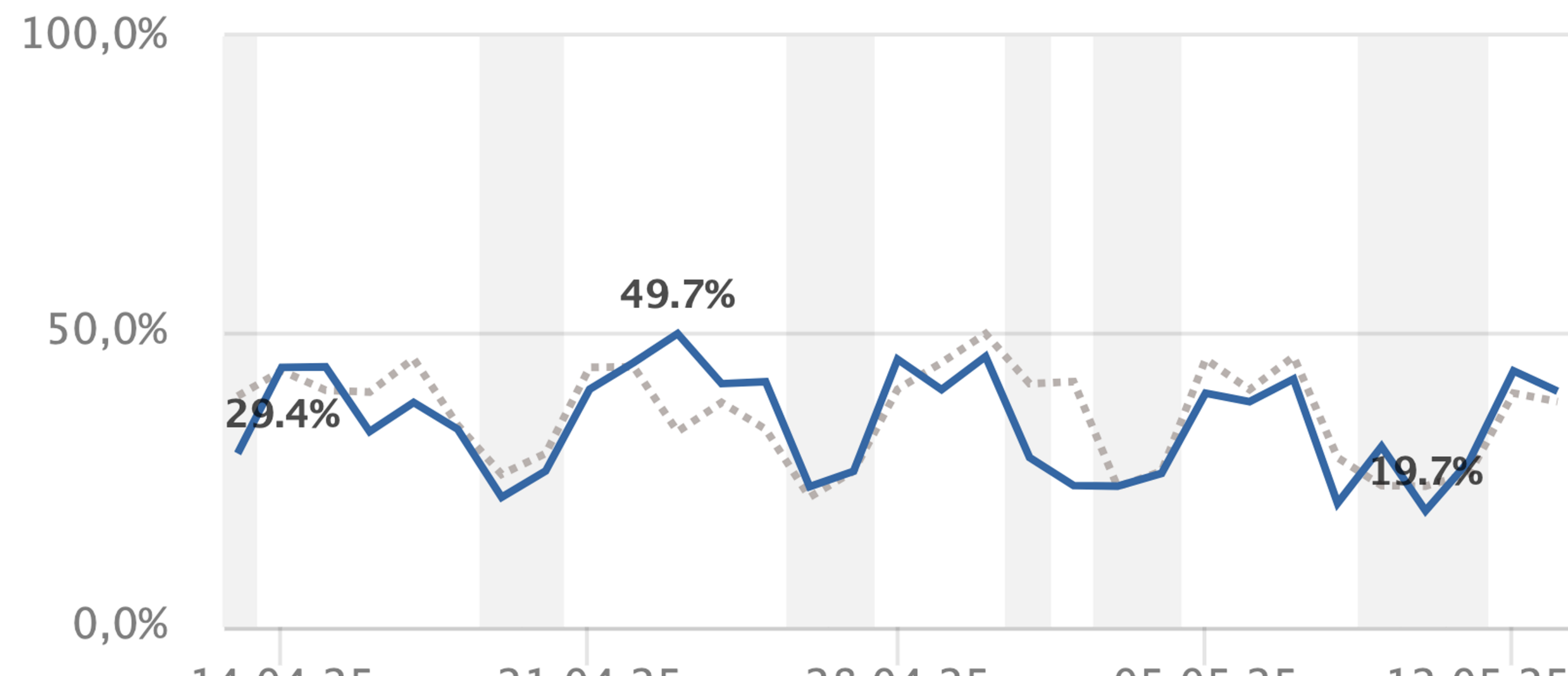
## Table Usage Index

shows the usage of DWH objects by analysts

46.7%

total for period

⬇ Table Usage



# как создать для своей платформы метрику Table Usage Index

- 01 Собери события использования DWH: кто, когда и какую таблицу использовал
- 02 Посчитай долю запросов к вашему DWH / запросы ко всем аналитическим данным в компании
- 03 Посчитай долю активных пользователей, определяй новичков/старичков и сравнивай их
- 04 Сравнивай сценарии использования, следи за тем, какие слои используют твои пользователи
- 05 Сохраняй в витрину dashboard\_usage\_index (day, dui) и обновляй эти данные в t-1
- 06 Следи в динамике за этой метрикой и обсуждай с командой в разрезе доменов на регулярной основе



07

# метрики эффективности использования ресурсов

следим за эффективностью использования ресурсов в разрезе доменов через метрику

# Efficiency Index

## ! Efficiency Index

This index shows approximate quota costs.  
([docs](#))

The cost of ownership is a guideline for estimating the approximate cost of a quota. It shows the monthly cost of the current quota volume at the moment.

The metric is displayed in conventional units (c.u.)

**c.u.** - conventional units are a virtual currency for estimating and comparing the cost of ownership of resources.

**i** To make the metric more flexible, we recalculated it by days, where 1 c.u. month = 1/30 c.u. per day

### Formula:

SUM (Total c.u.)

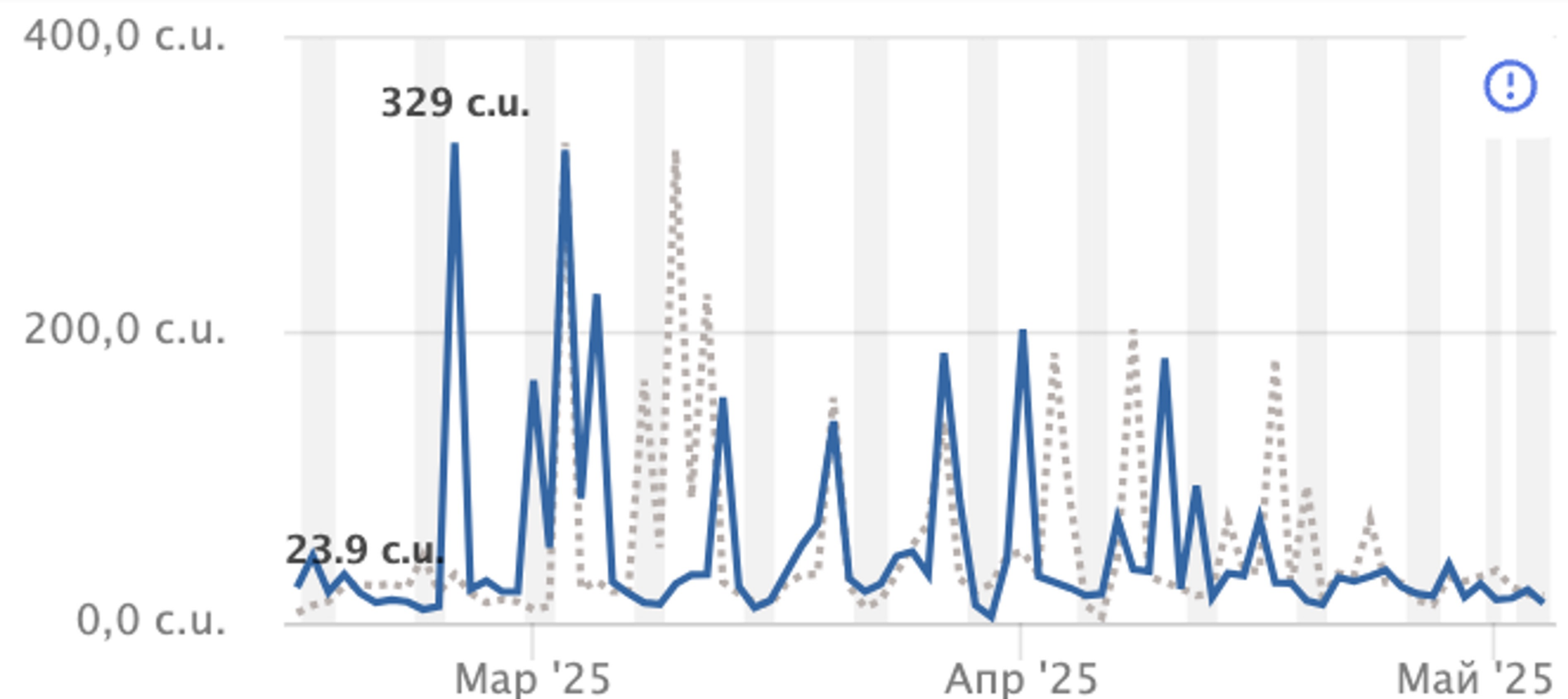
## Efficiency Index

shows the approximate quota costs

**23746 c.u.**

total for period

↓ Efficiency



# как создать для своей платформы метрику Efficiency Index

- 01 Определи квоты ресурсов (DWH, кластеры, хранилище, BI): CPU, RAM, Storage, Queries
- 02 Назначь «стоимость» квоты в условных единицах (с.у.) — внутренняя виртуальная валюта для сравнения
- 03 Собери фактическое потребление квот по дням
- 04 Пересчитай в с.у.: если 1 квота = 30 с.у./месяц → тогда  $\approx 1$  с.у./день
- 05  $\text{EfficiencyIndex} = \text{sum}(\text{Total} \setminus \text{с.у.})$
- 06 Сохрани в витрину `efficiency_index (day, total_cu)` и обновляй эти данные в `t-1`
- 07 Следи в динамике за этой метрикой и обсуждай с командой в разрезе доменов на регулярной основе



08

# метрики инженерной культуры аналитиков

# следим за инженерной культурой в разрезе доменов через метрику Engineering Culture Index

## ! Nirvana Engineering Culture Index

This index shows how analysts fits riles when creating pipelines in Nirvana.

Read more about rules here ([wiki](#))

### Formula:

$\text{AVG}(\text{Number of rules followed} / 7)$   
per pipeline

\* 7- targeted number of rules

## Engineering Culture Index

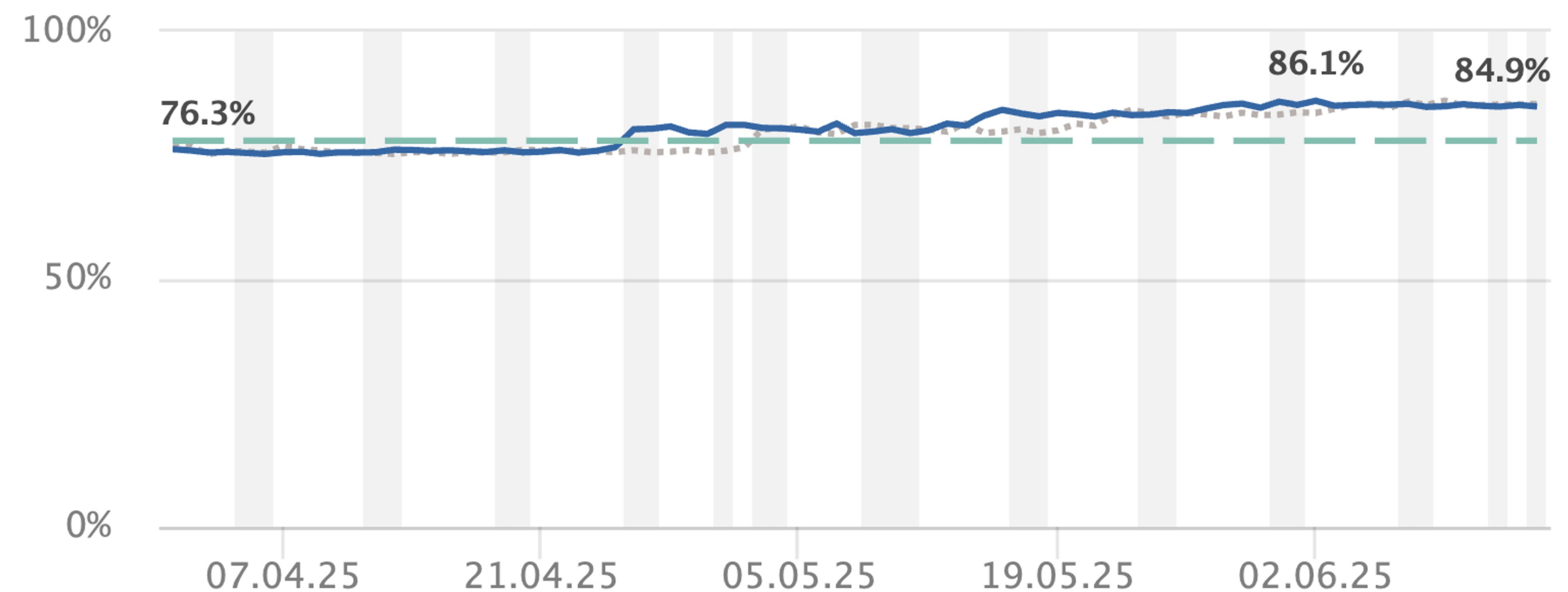
shows how well analysts follow rules when creating Nirvana pipelines



# 80.6%

total for period

⬇ Engineering Culture



# как создать для своей платформы метрику Engineering Culture Index

- 01 Определи правила качества кода/пайплайнов  
(например: тесты, логирование, мониторинг, документация, алерты, SLA, кодревью)
- 02 Логируй факт выполнения правил при создании ETL/ML/DA/BI пайплайнов
- 03 Для каждой задачи/пайплайна считай долю выполненных правил
- 04  $ECI = AVG(\text{pipeline} \backslash \text{score})$
- 05 Сохрани в витрину `engineering_culture_index` (`domain`, `day`, `eci`) и обновляй эти данные в `t-1`
- 06 Следи в динамике за этой метрикой и обсуждай с командой в разрезе доменов на регулярной основе
- 07 Расширяй индекс — добавь использование инженерных инструментов  
(Git PR, CI, тесты, мониторинг)

# покрытие иерархии метрик



# фреймворк необходимых действий по шагам

01

анализ бизнес-модели

02

декомпозиция  
бизнес-модели  
на доменные области

03

создание иерархии  
метрик в каждом домене

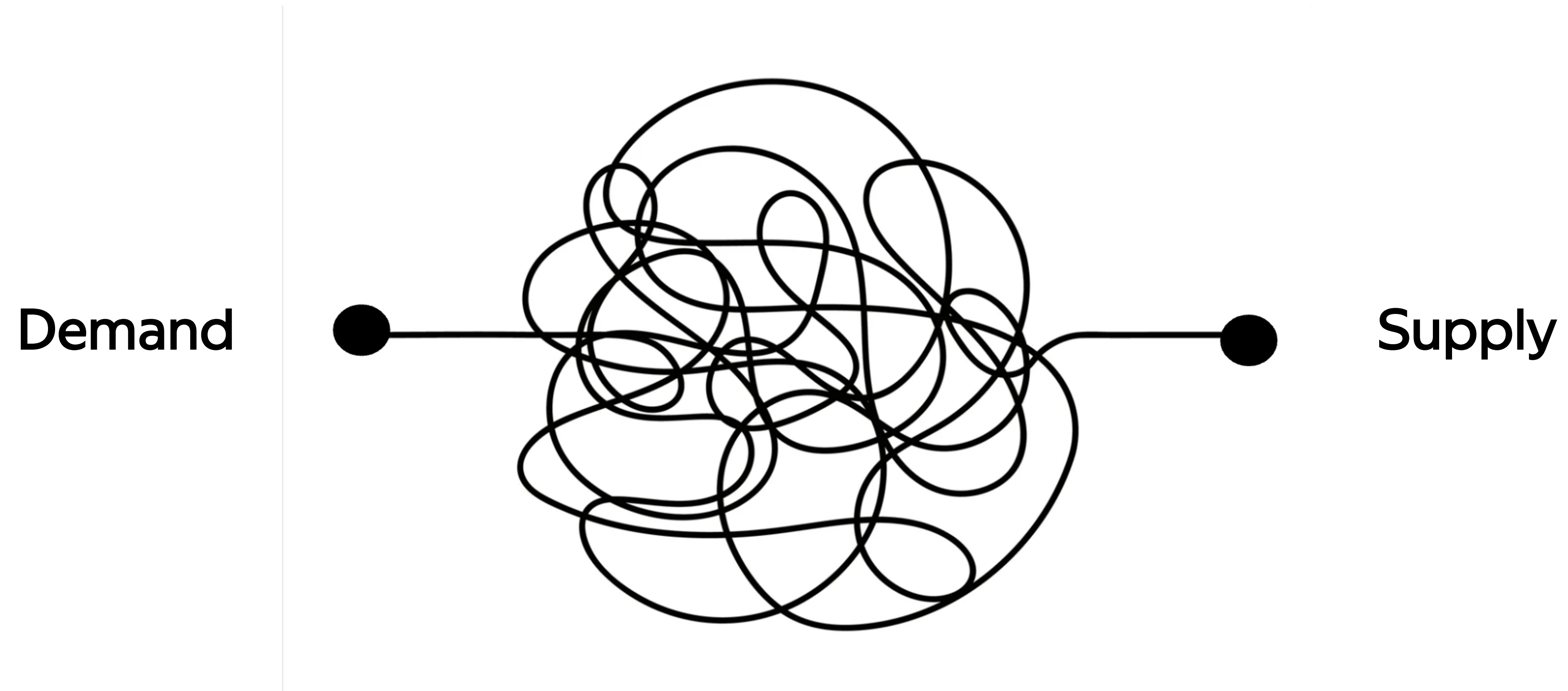
04

проектирование модели  
данных и системы отчётности

05

развитие через % покрытия  
иерархии метрик

# анализ бизнес-модели



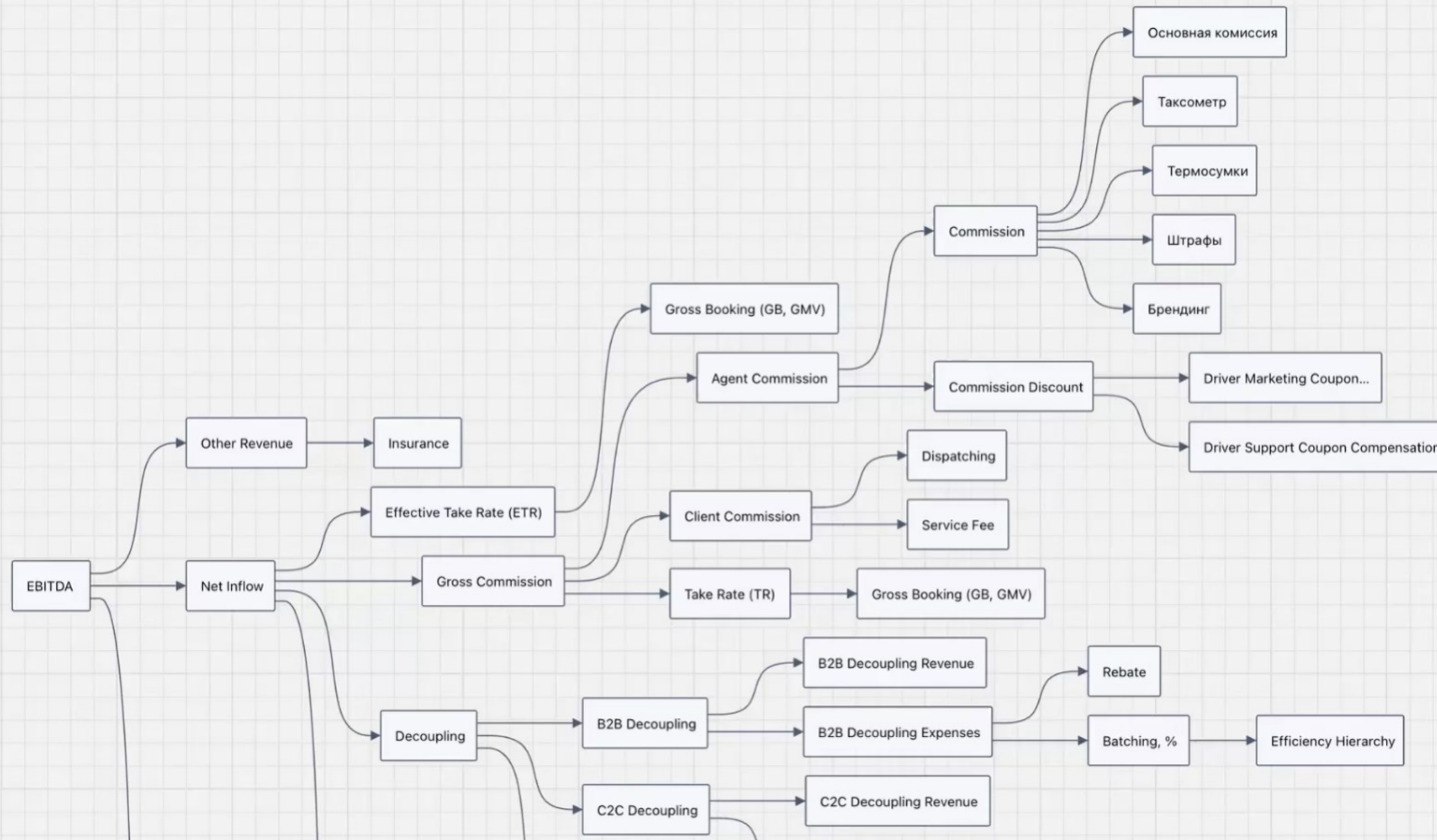
# декомпозиция бизнес-модели на доменные области

Demand

Efficiency

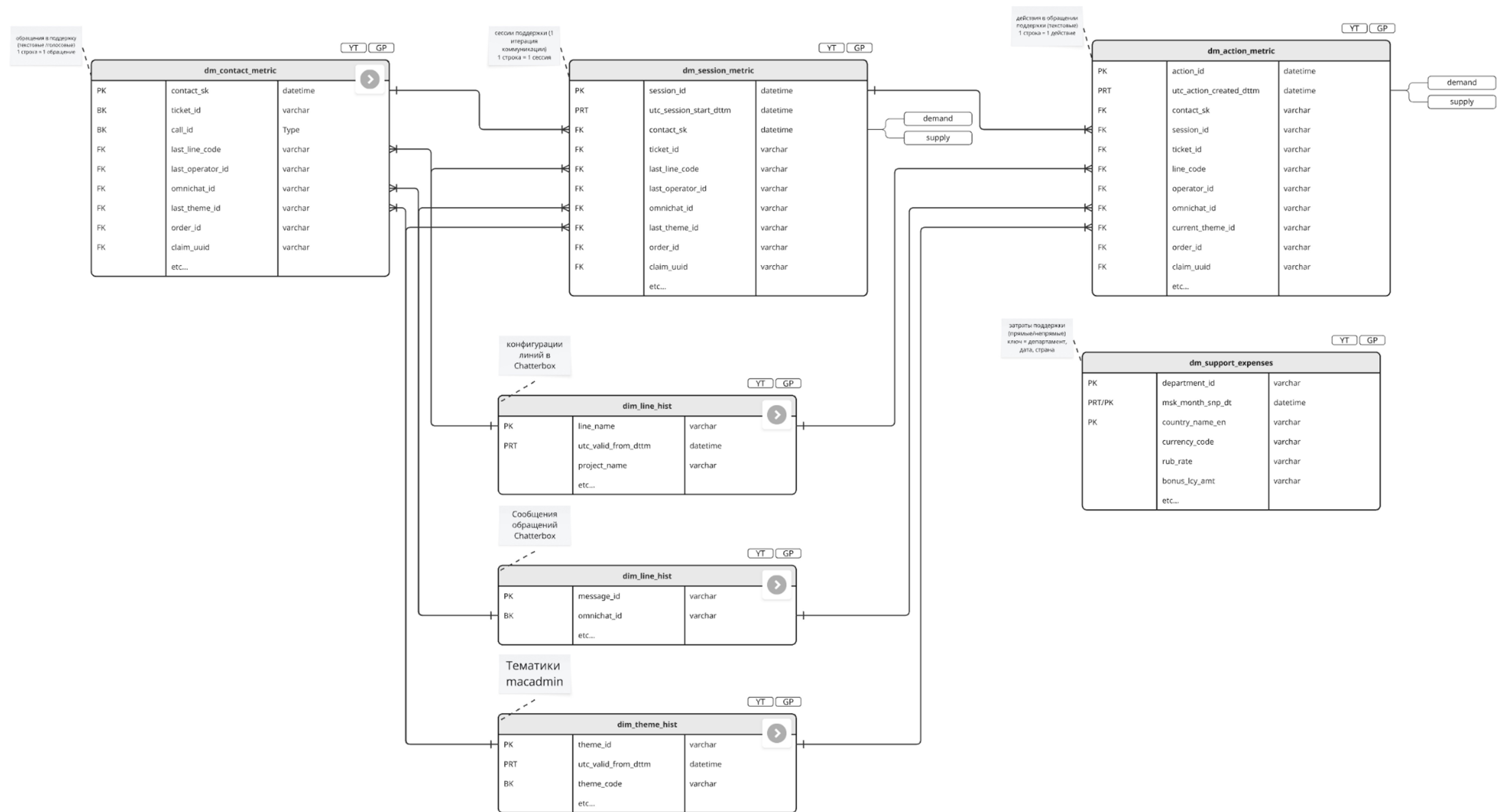
Supply

# Иерархии метрик в каждом домене бизнеса



# проектирование gold-слоя

## для аналитиков на основании метрик из иерархии бизнес-домена



# оценка качества модели аналитического слоя домена

через сложность и количество  
лишних соединений

CR Executor

```
SELECT
    date_trunc('month', lcl_snp_dt) AS month_dt,
    SUM(completed_delivery_cnt) / SUM(support_contacts_cnt) AS CR
FROM
    delivery_cdm_supply.agg_supply_data
WHERE
```

CR Client

```
SELECT
    date_trunc('month', utc_created_dt) AS month_dt,
    SUM(completed_delivery_cnt) / SUM(support_contacts_cnt) AS CR
FROM
    delivery_cdm_demand.dm_client_dayly_snp
WHERE
    1=1
    AND utc_created_dt BETWEEN '2020-01-01' AND '2024-12-31'
    AND
```

Theme

```
-- топ тематик по кол-ву обращений
SELECT
    main_theme_id,
    count(*)
FROM
    delivery_cdm_support.dm_contact_metric
GROUP BY
    1
ORDER BY
    2 DESC;

-- количество обращений в поддержку по финансовым вопросам от клиентов и курьеров
SELECT
    author_group_category_name,
    COUNT(*)
FROM
    delivery_cdm_support.dm_contact_metric
WHERE
    main_theme_name_ru = 'финансовые вопросы';
```

# проектирование системы отчётности домена на основании метрик из иерархии

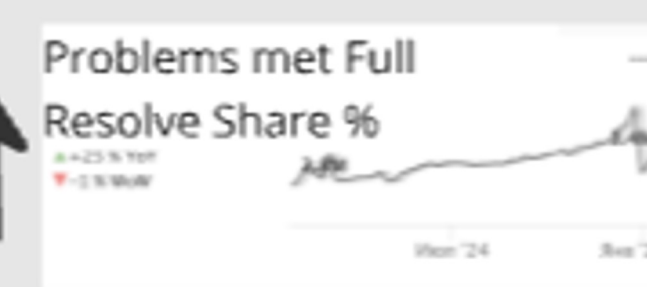
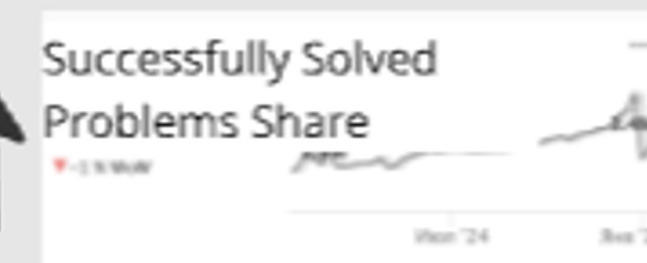
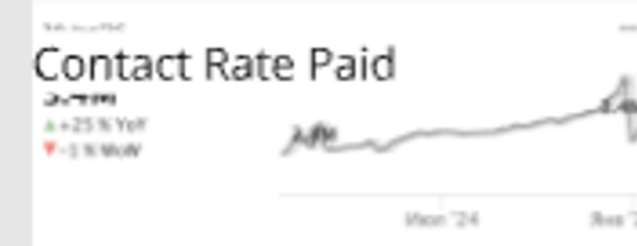
Фактоиды с динамикой

Support  
Requests

Support  
Quality

Support  
Expenses

Expenses не  
показываем  
КПБ

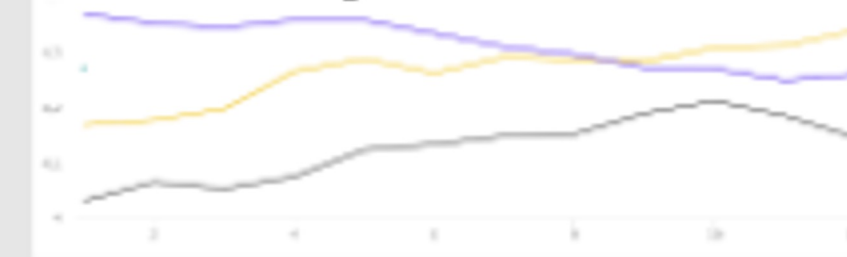


Пока  
непонятно  
как  
считать

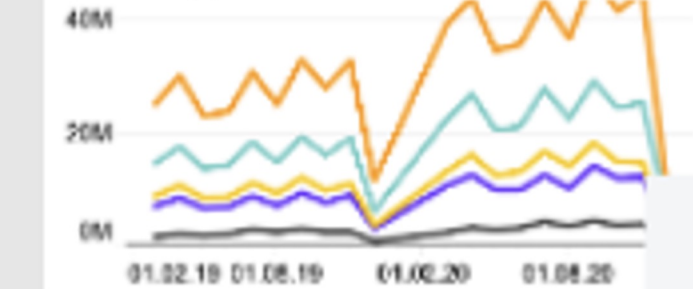
# проектирование системы отчётности домена на основании метрик из иерархии

## Contact Rate

CR & Incoming CR & CR Paid



Incoming CR по категориям



Format ④ Charts ④ Tables ④

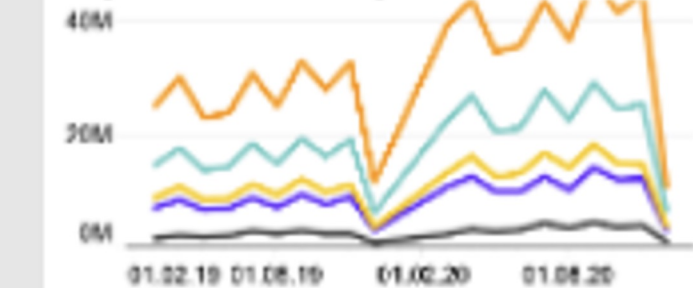
Incoming CR TOP стран

robot_dnp	140.71M
robot_catdnp	67.63M
robot_catdnp	23.7M
robot_cat	2.69M
robot_catdnp	5.70M
robot_catdnp	18.10M
robot_cat	10.00M
robot_cat	1.00M
robot_cat	1.00M
robot_cat	1.00M

Completed deliveries



CR paid по категориям



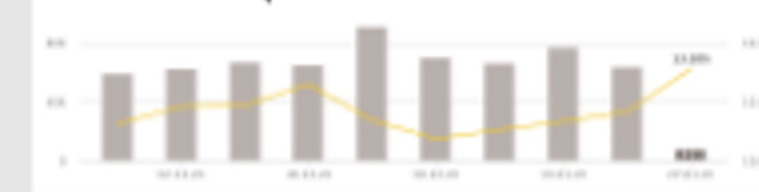
CR Paid TOP стран

robot_dnp	140.71M
robot_catdnp	67.63M
robot_catdnp	23.7M
robot_cat	2.69M
robot_catdnp	5.70M
robot_catdnp	18.10M
robot_cat	10.00M
robot_cat	1.00M
robot_cat	1.00M
robot_cat	1.00M

By language:  
ROM, RU,  
TURK, HSB

## CSAT

CSAT & Response rate



Tickets CSAT asked & Ticket (rated)



Format ④ Charts ④ Tables ④

Bad CSAT share (DR)

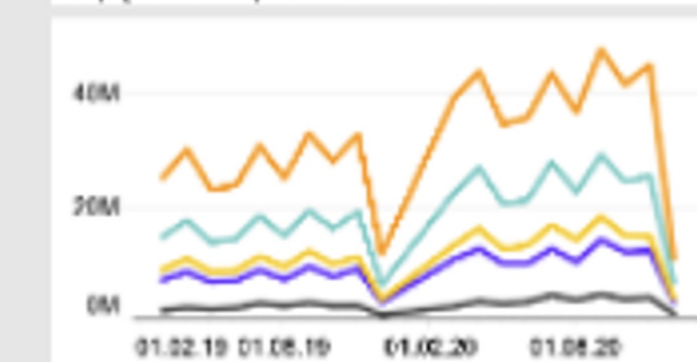


## Support Expenses

Support Expenses Share



Support Expenses Share



Format ④ Charts ④ Tables ④

Support Expenses



GMV



Support Expenses Share & Expenses no странам

robot_dnp	140.71M	140.71M
robot_catdnp	67.63M	67.63M
robot_catdnp	7.57M	7.57M
robot_cat	7.49M	7.49M
robot_catdnp	5.70M	5.70M
robot_catdnp	18.10M	18.10M
robot_cat	10.00M	10.00M
robot_cat	1.00M	1.00M
robot_cat	1.00M	1.00M
robot_cat	1.00M	1.00M

# следим за развитием покрытия иерархии метрик в системе отчётности через Metric Hierarchy Index

## ! Metric Hierarchy Index

This index shows the coverage of business metrics from hierarchy that available for users in BI domain dashboards in DataLens

### Formula:

$$\frac{\text{Metrics available in DataLens}}{\text{Total metrics in hierarchy}}$$

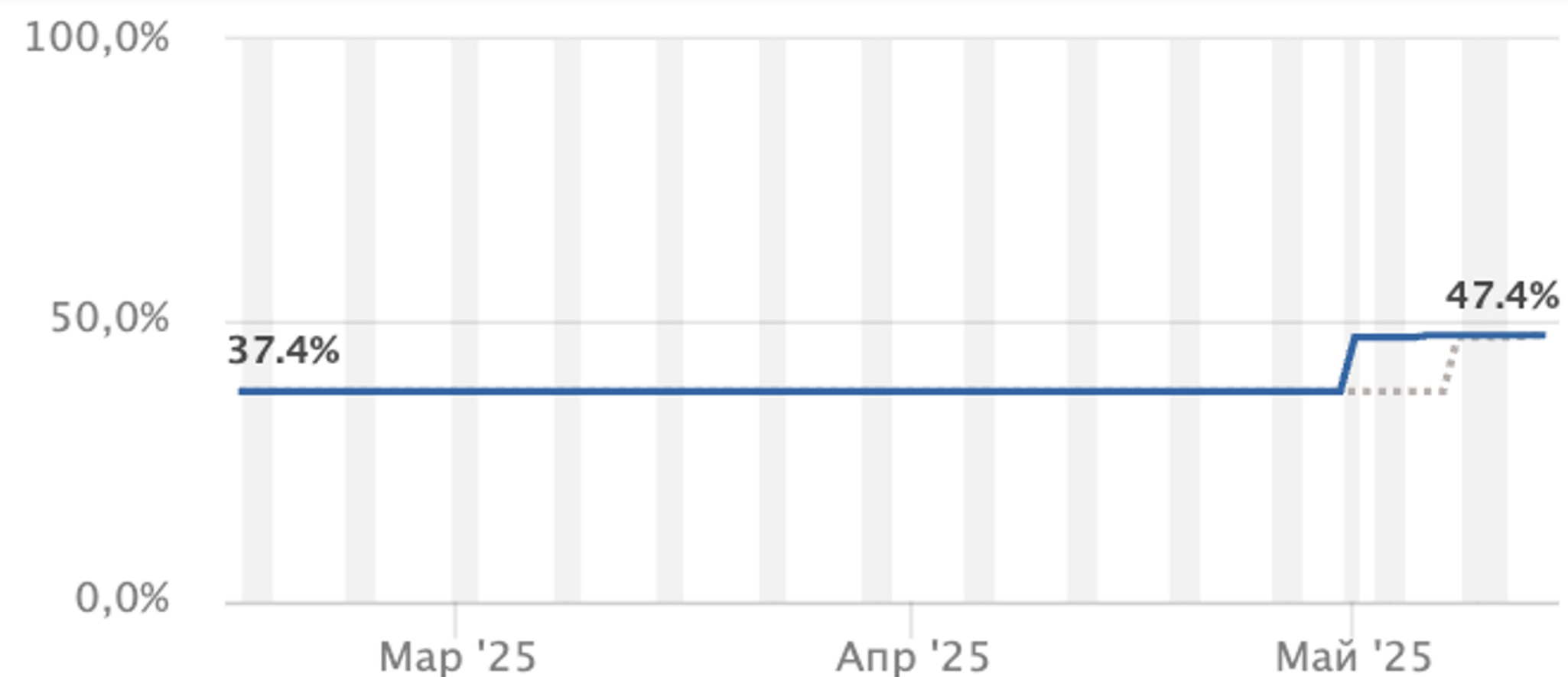
## Metric Hierarchy Index

shows the coverage of business metrics in BI dashboards

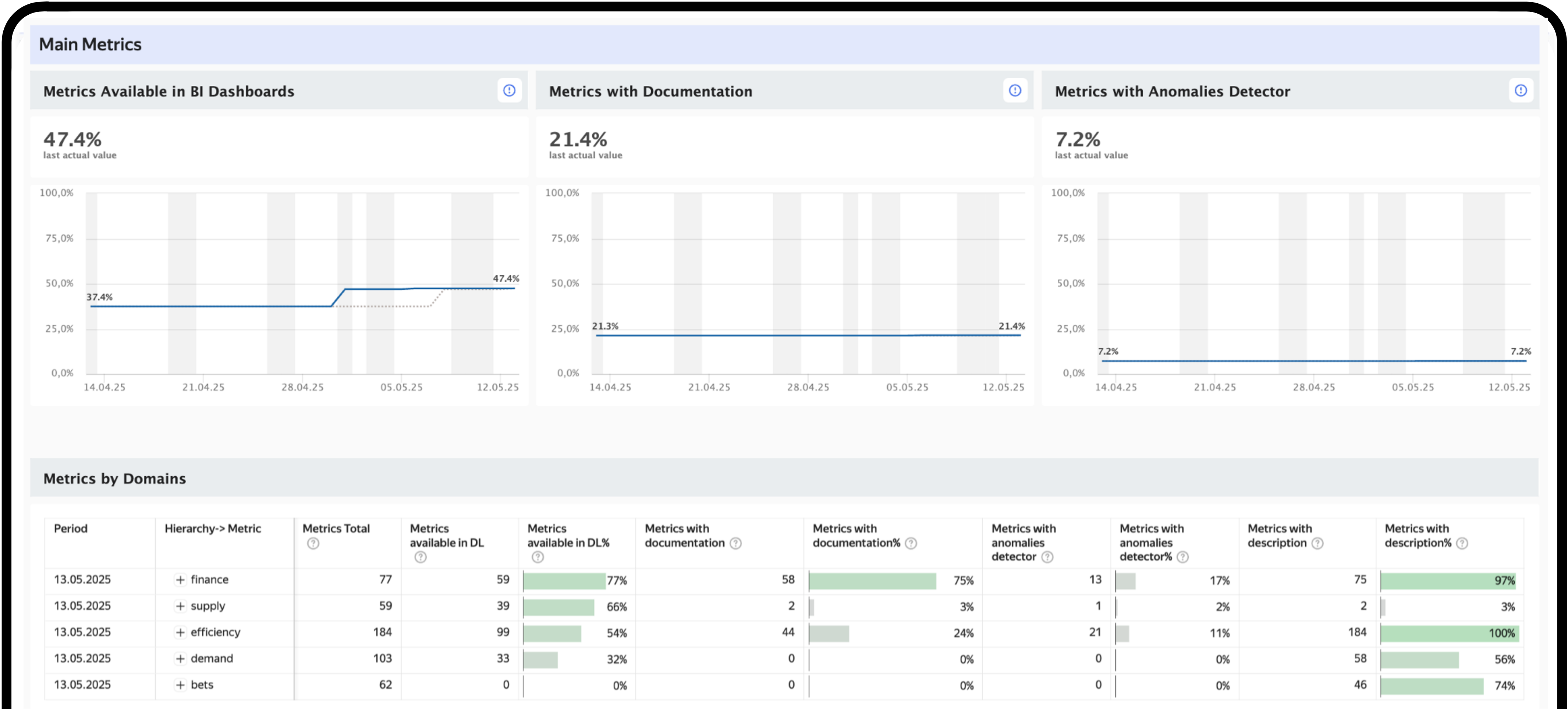
47.4%

last actual value

⬇ Metric Hierarchy



# И делаем это в разрезе каждого домена



# как создать для своей платформы метрику Metric Hierarchy Index

- 01 Определи эталонную иерархию бизнес-метрик для каждого домена
- 02 Собери список метрик в BI (DataLens/Superset/Metabase), доступных пользователям в дашбордах
- 03 Сопоставь метрики BI с эталонной иерархией
- 04 Формула индекса:  $MHI = \text{Metrics available in BI} \setminus \text{Total metrics in hierarchy}$
- 05 Агрегируй по своим доменам
- 06 Сохрани в витрину `metric_hierarchy_index (domain, day, mhi)`
- 07 Следи в динамике за этой метрикой и обсуждай с командой в разрезе доменов на регулярной основе

# ВЫВОДЫ



## ВЫВОДЫ

01

**оцифровка компонентов  
позволяет создавать понятную  
для всех в бизнесе стратегию  
развития платформы данных**

## ВЫВОДЫ

02

**для каждого из направлений  
платформы появляются  
оцифрованные цели и результаты**

## ВЫВОДЫ

03

**ИНЖЕНЕРЫ ПОНИМАЮТ, НА КАКИЕ  
МЕТРИКИ И КАК ИМЕННО ОНИ ВЛИЯЮТ  
ЧЕРЕЗ СВОИ ДОРАБОТКИ И ИЗМЕНЕНИЯ  
В СЦЕНАРИЯХ ОБРАБОТКИ ДАННЫХ**

## ВЫВОДЫ

04

**системные аналитики понимают,  
как их объектная модель используется  
и могут делать регулярные выводы,  
как её улучшить**

## ВЫВОДЫ

05

**bi-аналитики могут следить  
за развитием своей системы отчётности  
через покрытие иерархий метрик  
и метрик использования дашбордов**

## ВЫВОДЫ

06

**всё это работает вместе  
как единая система с общими целями  
и создает правильную культуру  
внутри команд платформы данных**

## ВЫВОДЫ

07

**В ЭТУ СИСТЕМУ МЕТРИК ЗАЕЗЖАЮТ  
ПЛАТФОРМЫ ДАННЫХ ТАКСИ, ЕДЫ,  
ЛАВКИ, ФИНТЕХ, СИБ**

**КАК МЫ ЕЁ ПРОКАЧАЛИ,  
РАССКАЖУ В СЛЕДУЮЩИХ СЕРИЯХ**

Smart Data

**ГОТОВ ОТВЕТИТЬ  
НА ВОПРОСЫ**

**Влад  
Гоцуляк**

Руководитель платформы данных в Яндекс Доставке

