

Яндекс Go



# Метрики работы DWH: от проблематики до реализации

Евгений Ермаков, руководитель Data Office

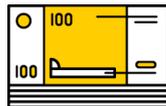
# Что такое Яндекс Go?



**700**

**ТЫС.**

активных водителей,  
которые сделали более  
одного заказа в месяц



**18**

стран присутствия, в  
том числе Гана и Кот-  
д'Ивуар



**1000**

городов, из них 300  
крупных

# Хранилище данных Яндекс Go



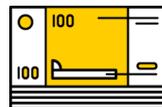
**>500**

уникальных  
пользователей  
данных в месяц



**>900**

отчетов по различным  
тематикам



**4**

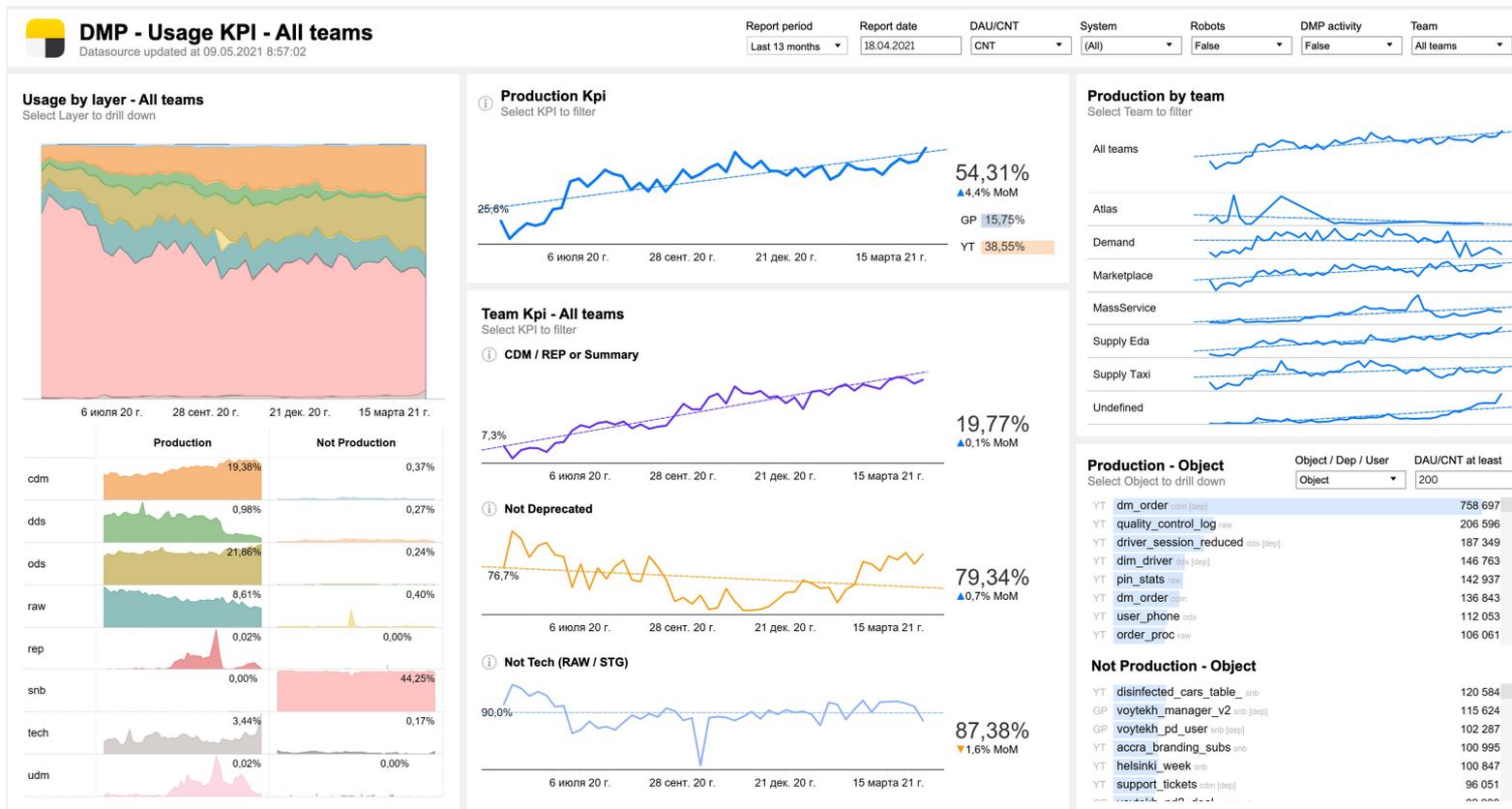
крупных бизнес-  
юнита: Такси, Еда,  
Лавка, Драйв



**3 ПБ**

накопленных данных  
по четырем бизнес-  
юнитам

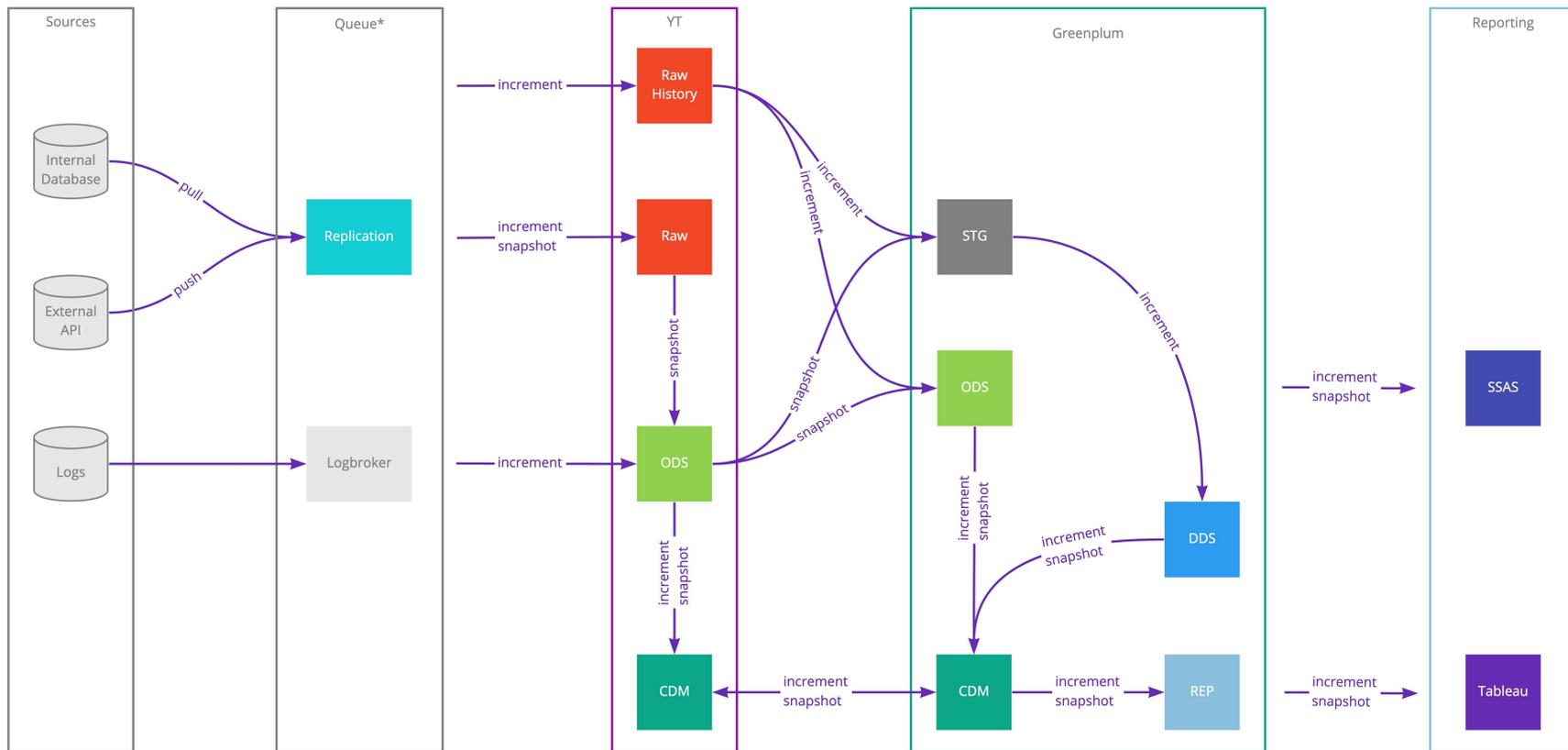
# Метрики хранилища



# Содержание

- I Зачем DWH метрики?
- II Как мы их реализовали?
- III Что получили?
- IV Стоило ли того?

# Платформа данных Яндекс GO



# Почему так сложно?

I. Зачем DWH метрики?

# Архитектура слоев данных



# Архитектура слоев данных



## Цель

- › Захватить сигналы источника

## Задачи

- › Собрать данные с источника as-is
- › Преобразовать их в объекты с понятным описанием и методом доступа

# Архитектура слоев данных



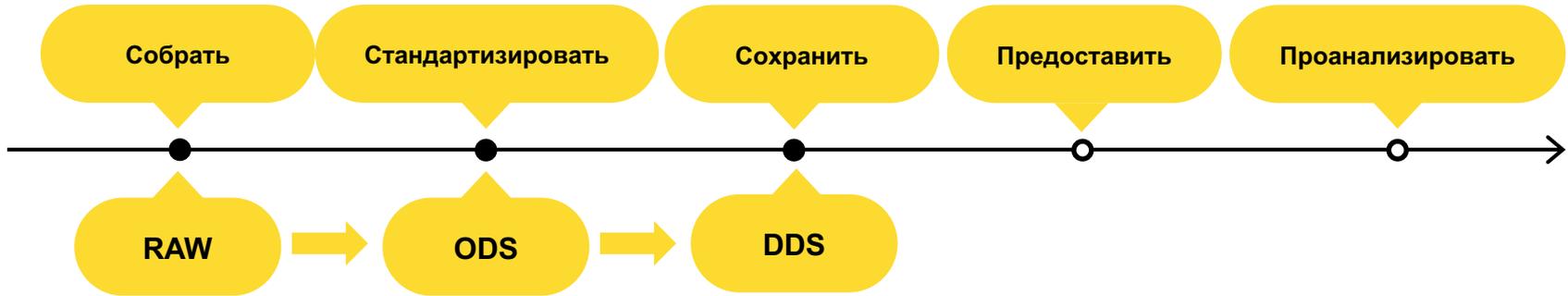
## Цель

- › Хранить операционные данные источника

## Задачи

- › Сформировать набор сущностей источника
- › Разложить данные по сущностям
- › Предоставить стандартный интерфейс доступа к данным

# Архитектура слоев данных



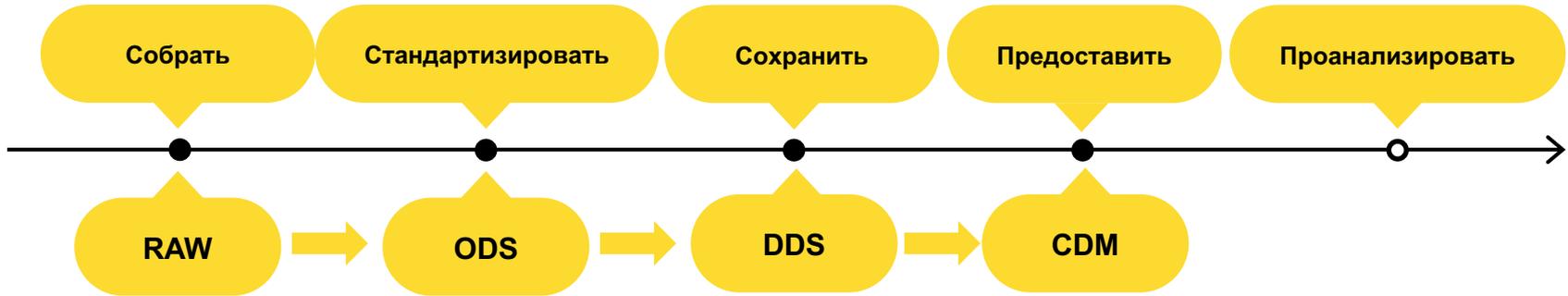
## Цель

- › Накапливать данные о сущностях доменной модели

## Задачи

- › Хранить детальную историю изменений
- › Консолидировать данные между источниками

# Архитектура слоев данных



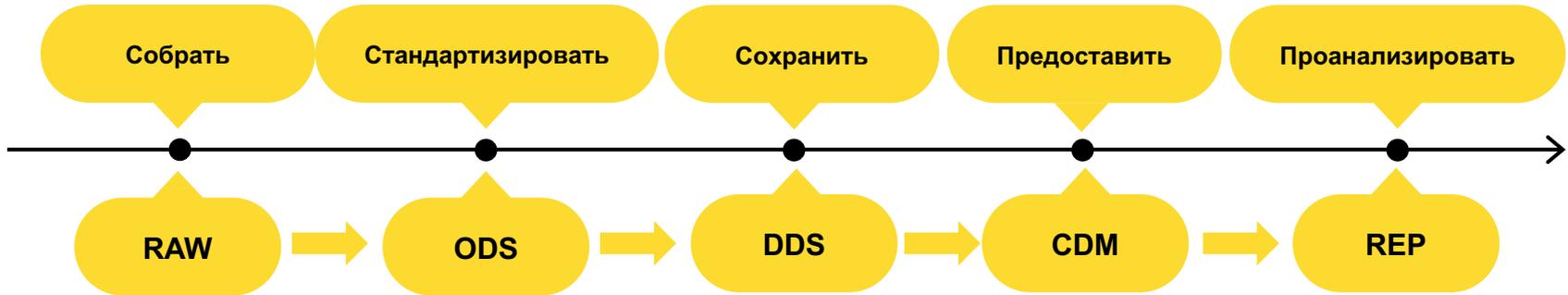
## Цель

- › Предоставлять витрины данных для анализа

## Задачи

- › Формировать данные в контексте бизнес-потребностей
- › Оптимизировать доступ на чтение

# Архитектура слоев данных



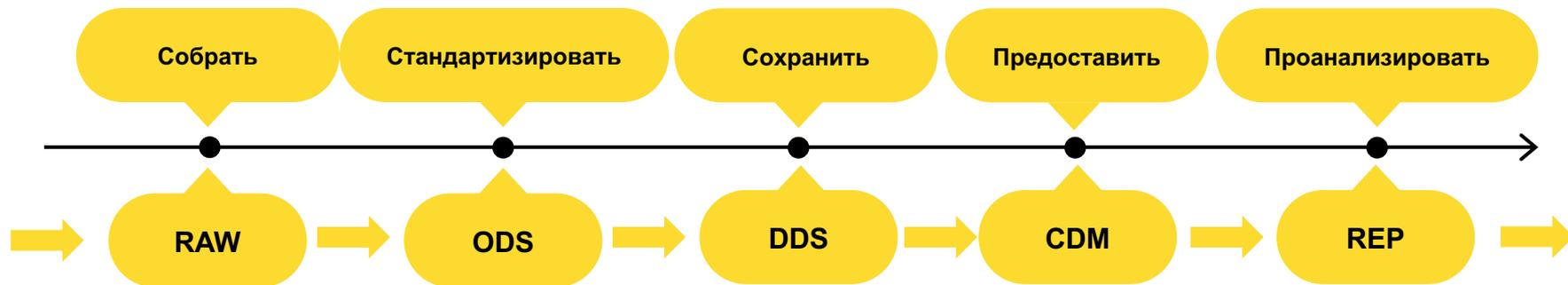
## Цель

- › Хранить отчетные срезы

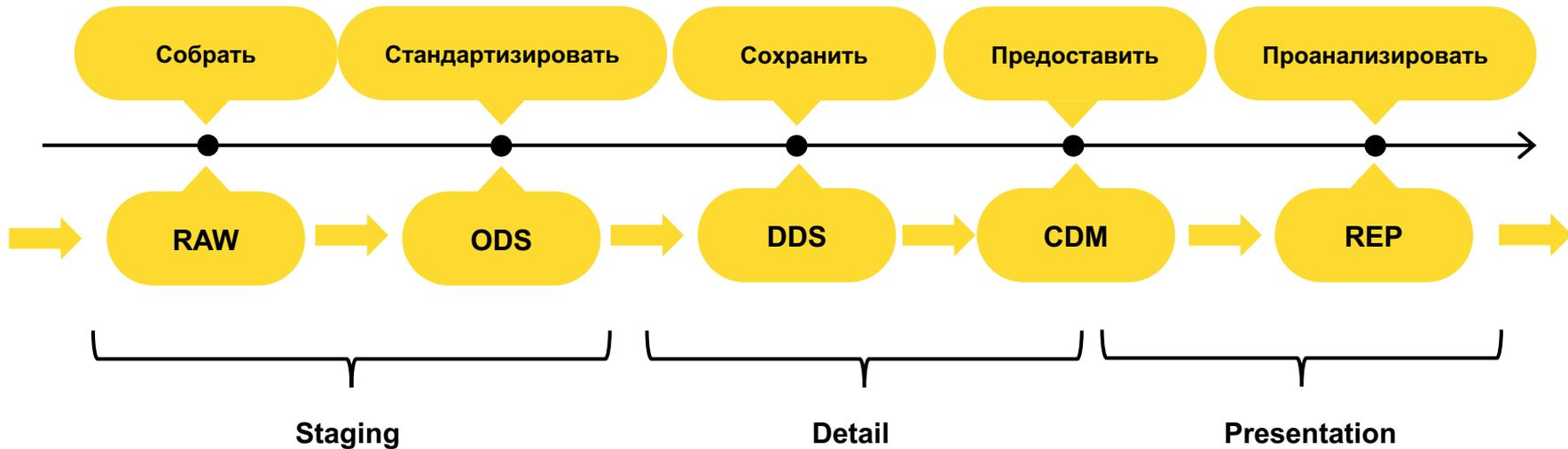
## Задачи

- › Формировать данные в контексте бизнес-потребностей
- › Готовить агрегированные отчеты

# Архитектура слоев данных



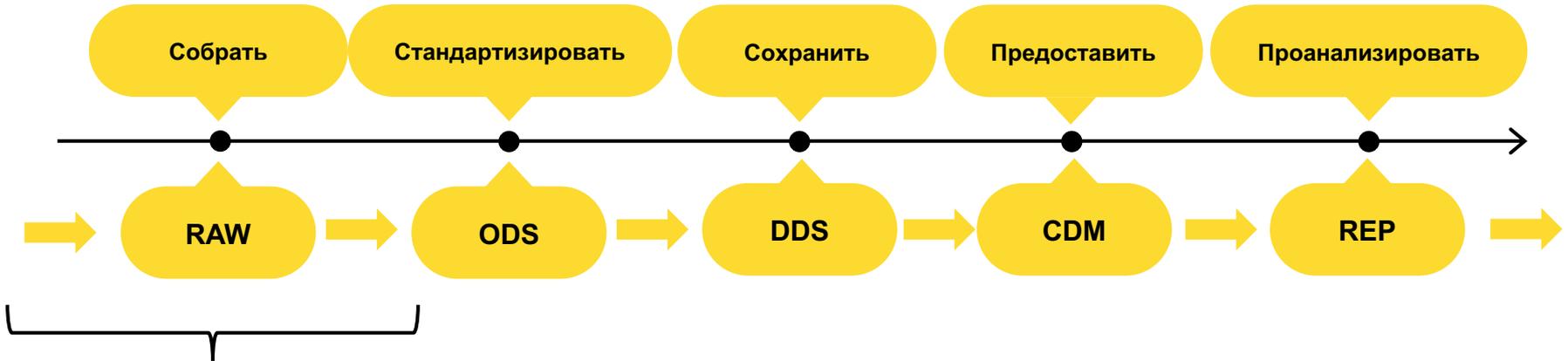
# Архитектура слоев данных



# Как слои связаны с системами?

I. Зачем DWH метрики?

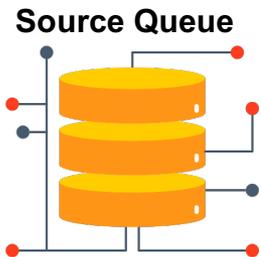
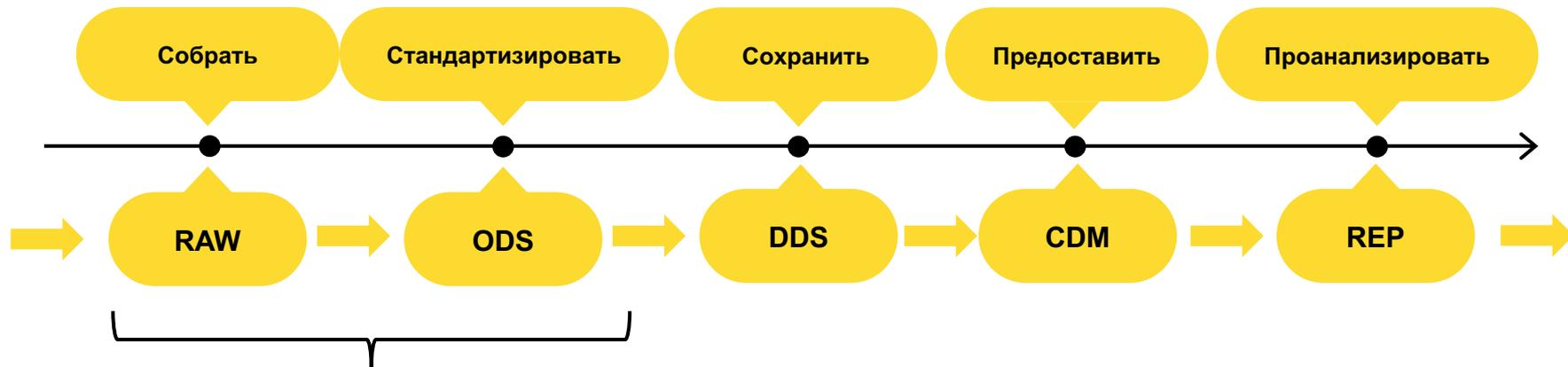
# Архитектура слоев данных



## Source Queue

- › Забирает инкременты и снэпшоты с источников различных типов
- › Преобразовывает данные в устойчивый к изменениям формат

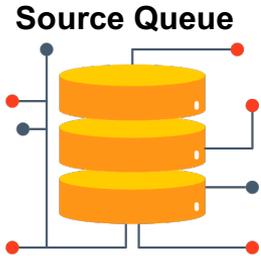
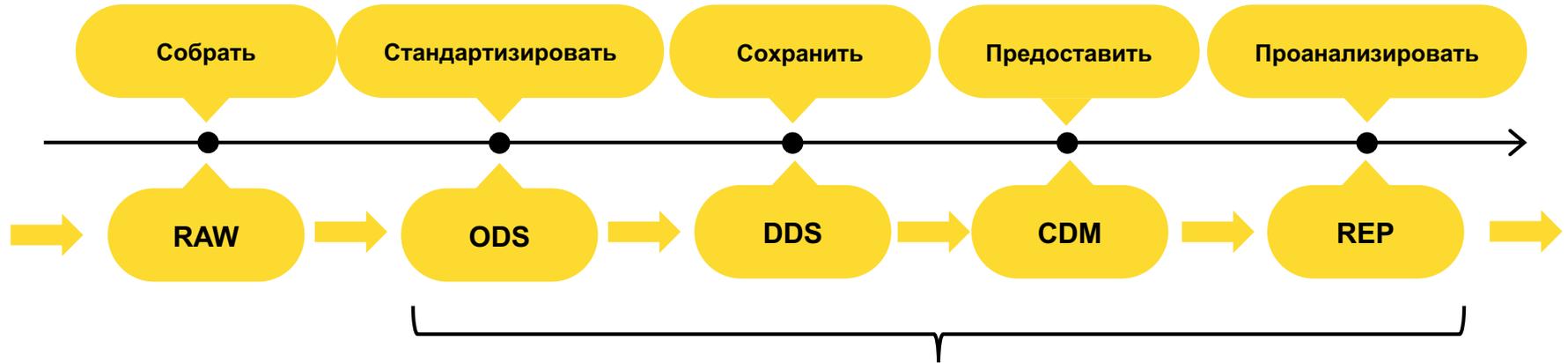
# Архитектура слоев данных



## YT (Data Lake)

- > Полуструктурированные данные
- > Каркас MapReduce
- > Аналогии экосистемы hadoop

# Архитектура слоев данных



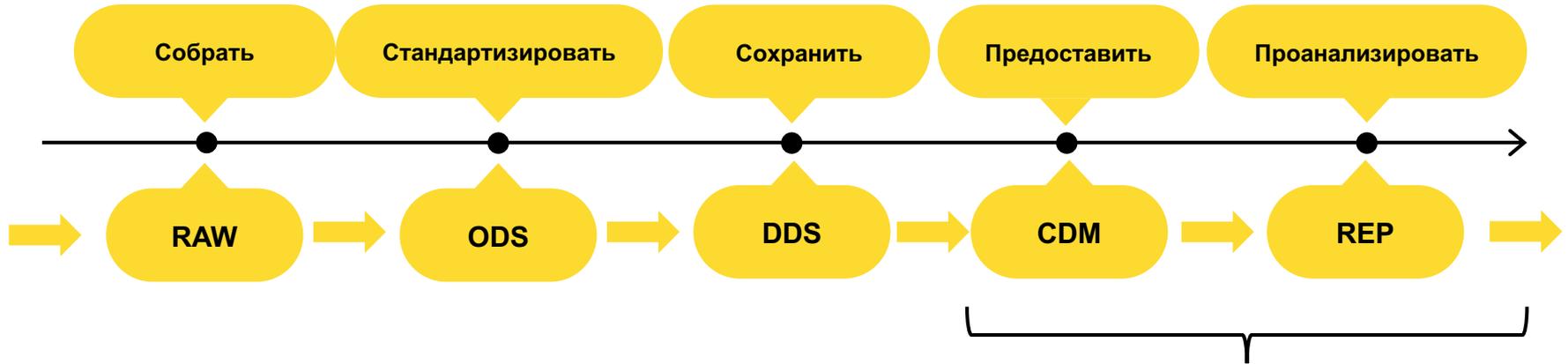
YT (Data Lake)



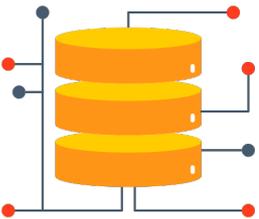
Greenplum (Data warehouse)

- > Различные ad-hoc-запросы
- > Большое количество join
- > Малое время отклика

# Архитектура слоев данных



Source Queue



YT (Data Lake)



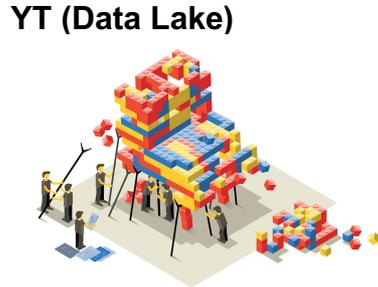
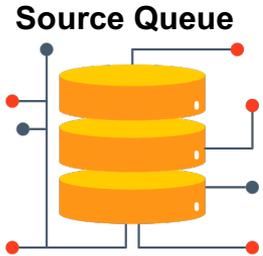
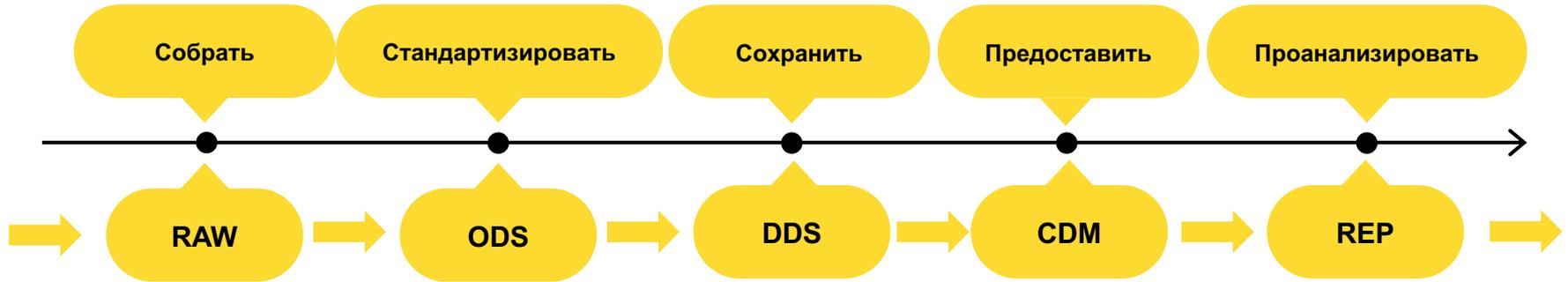
Greenplum (Data warehouse)



Reporting

- > Кубы данных
- > Отчеты и дашборды

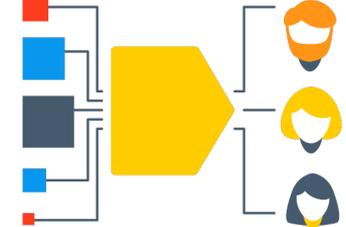
# Архитектура слоев данных



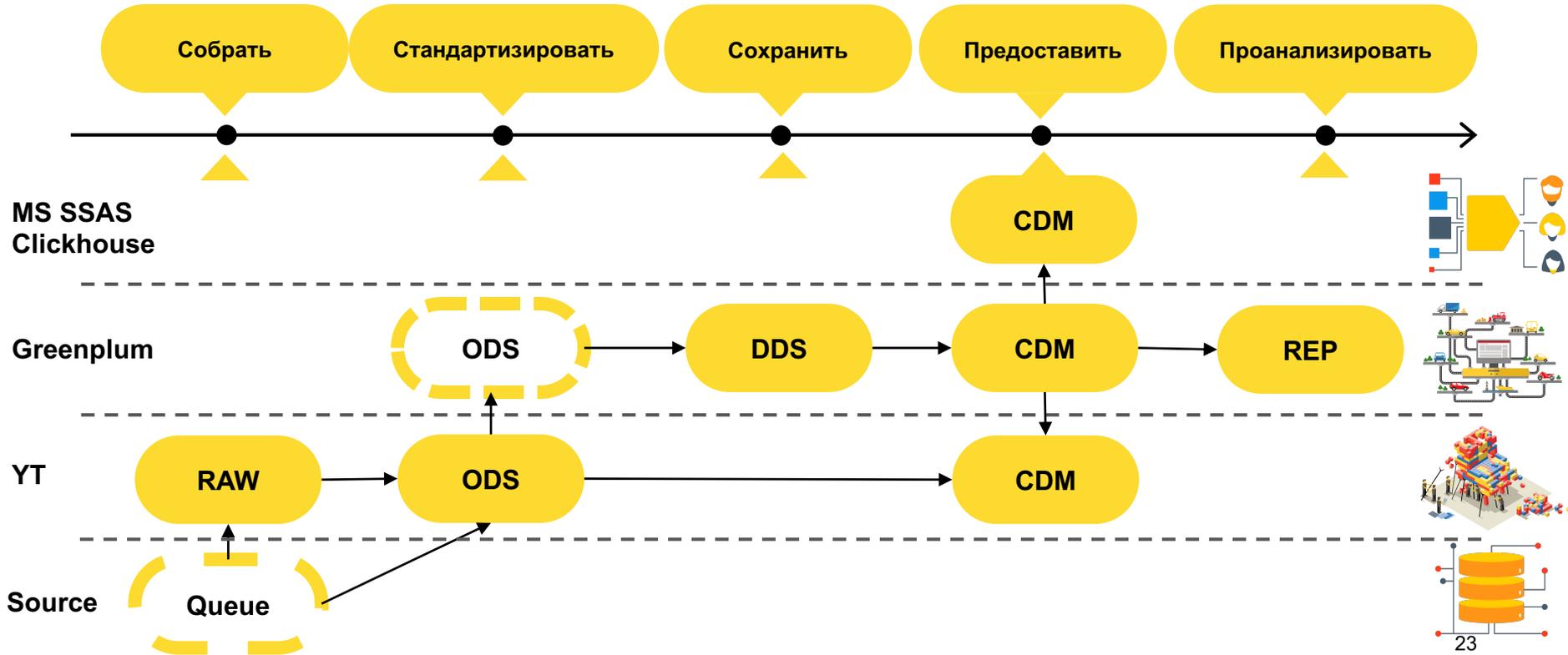
Greenplum (Data warehouse)



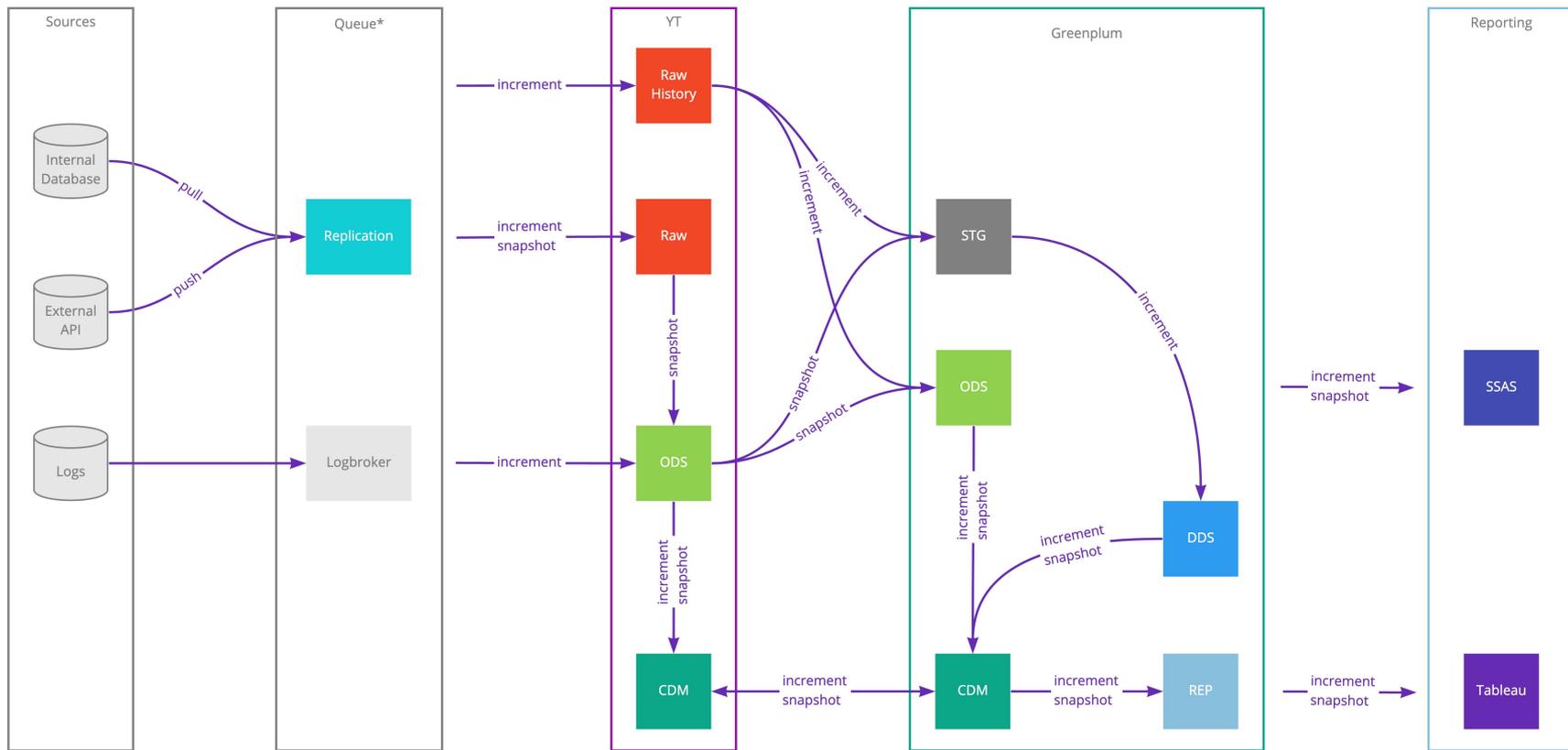
Reporting



# Архитектура слоев данных



# Платформа данных Яндекс GO



# Как развивать?

I. Зачем DWH метрики?

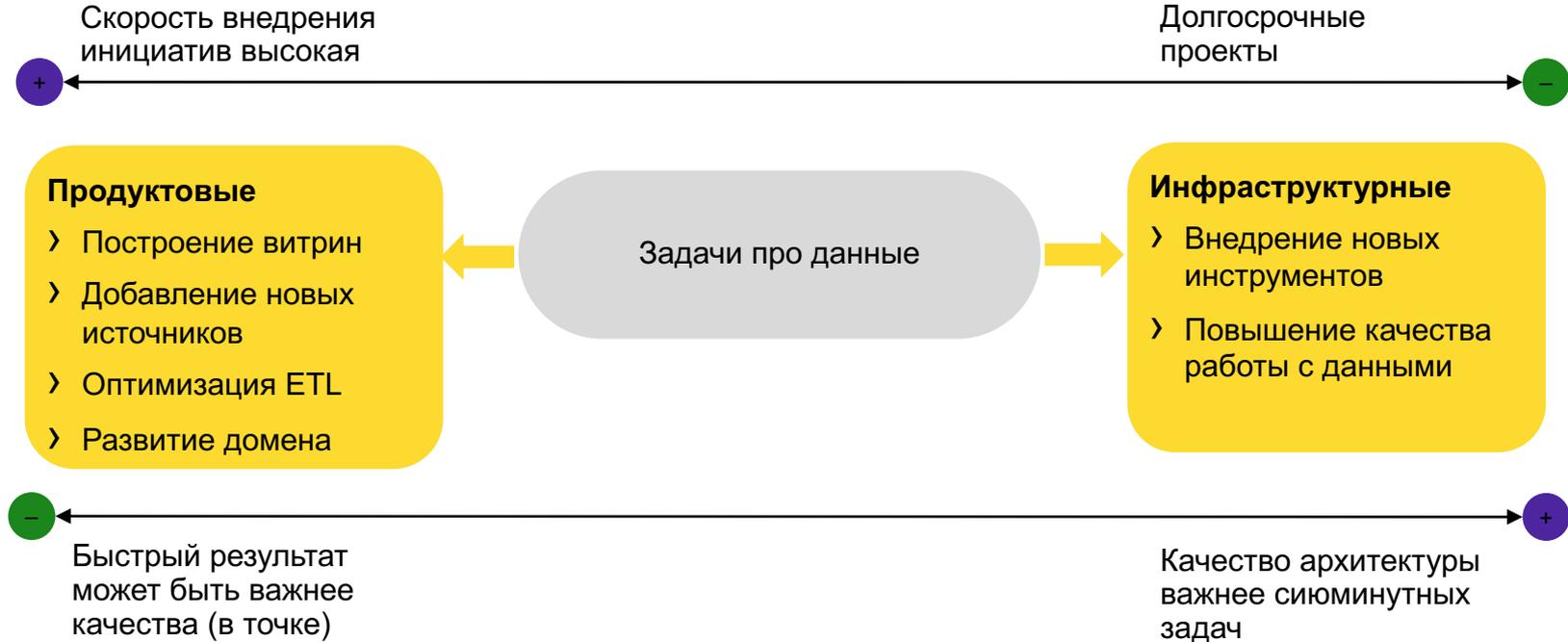
# Организация развития

Задачи про данные

# Организация развития



# Организация развития



# Организация развития

Троевластие: организационные, инфраструктурные и продуктовые команды



# Домены данных

**Данные сгруппированы по предметной области – домену (Domain)**

- › В одном домене может быть несколько объектов (таблиц)
- › За несколько доменов отвечает одна команда
- › Домены могут быть разных типов

# Домены данных

Данные сгруппированы по предметной области – домену (Domain)

- › В одном домене может быть несколько объектов (таблиц)
- › За несколько доменов отвечает одна команда
- › Домены могут быть разных типов

## Source Domain

---

- › Связаны с источником данных
- › Структура подогнана под источник
- › Включают в себя очистку, дедубликацию, приведение к стандартам и т.п.

## Core Domain

---

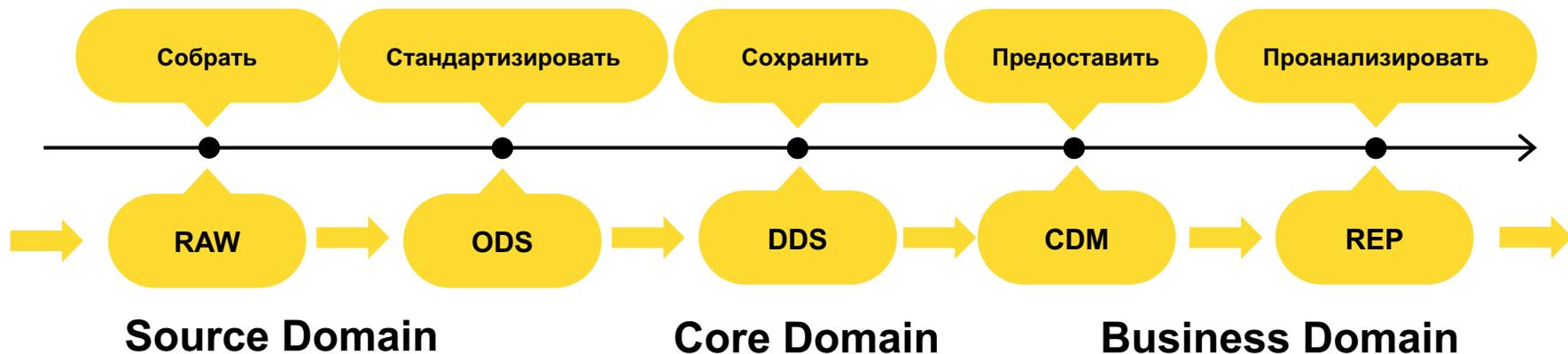
- › Связаны с крупной областью бизнеса
- › Структура подогнана под минимизацию изменений
- › Включает объединение данных из разных источников, генерацию суррогатных ключей и т.п.

## Business Domain

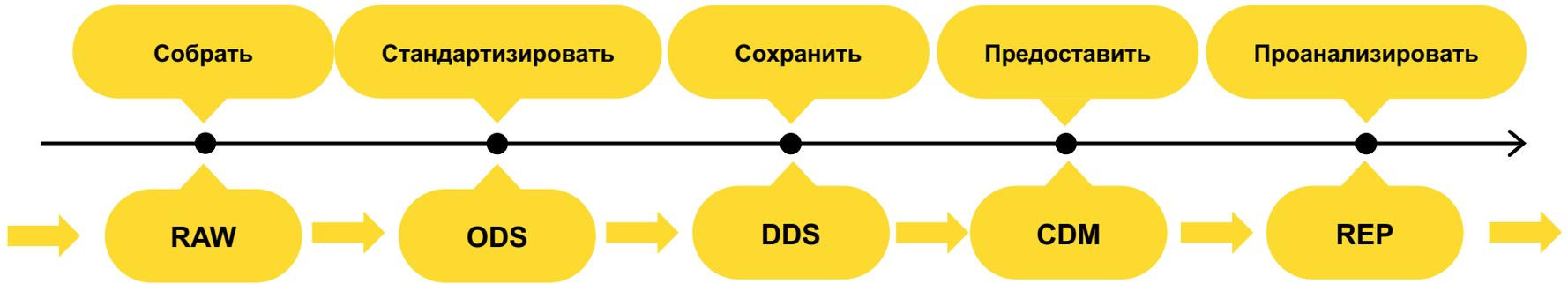
---

- › Связаны с потребителями данных
- › Структура подогнана под удобства использования
- › Фактически представляет собой специализированные витрины и/или отчеты

# Домены данных



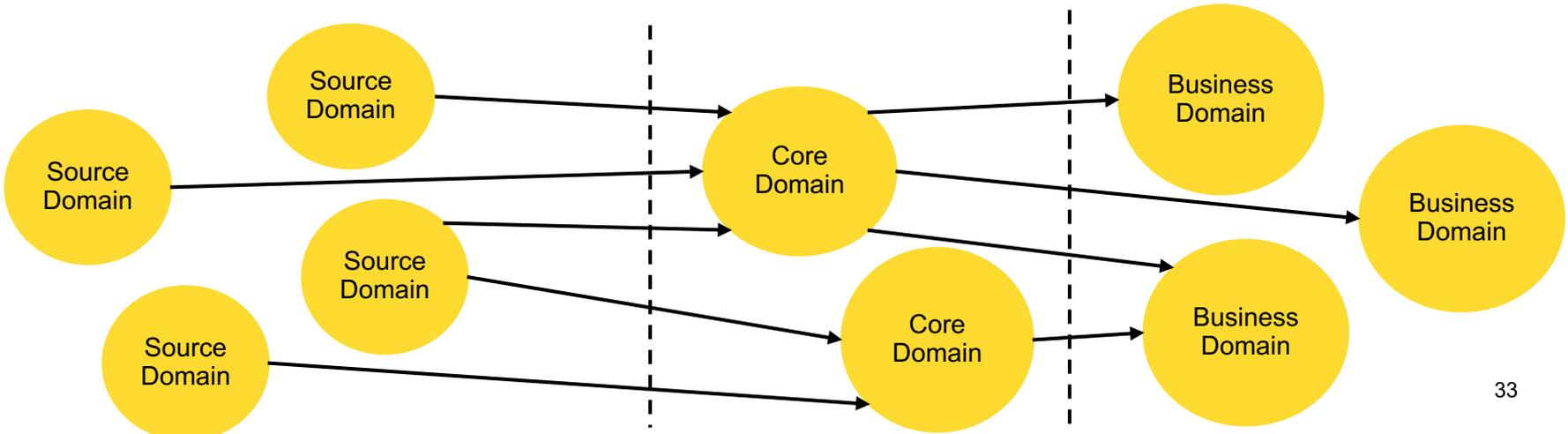
# Домены данных



**Source Domain**

**Core Domain**

**Business Domain**



# Как управлять (бес)порядком?

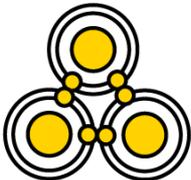
I. Зачем DWH метрики?

# Почему (бес)порядок?



**>500**

уникальных  
пользователей  
данных в месяц



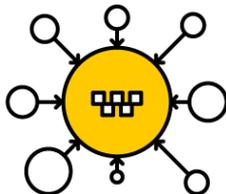
**>200**

доменов данных



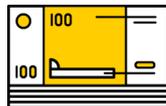
**>900**

отчетов по различным  
тематикам



**>3500**

объектов хранилища



**4**

крупных бизнес-  
юнита: Такси, Еда,  
Лавка, Драйв



**200**

коммитов в день



**3 Пб**

накопленных данных  
по четырем бизнес-  
юнитам



**300**

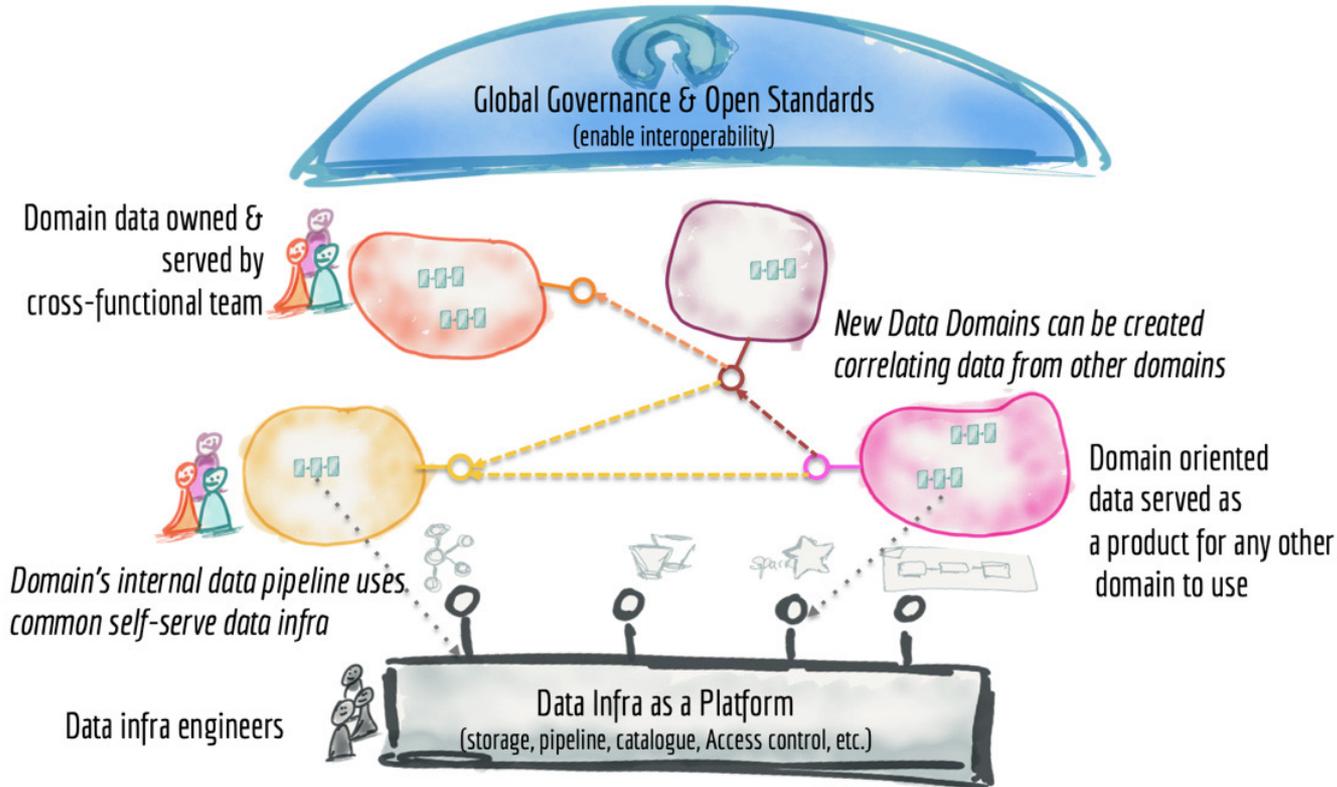
merged Pull Request в  
месяц

# **I. Проблема:**

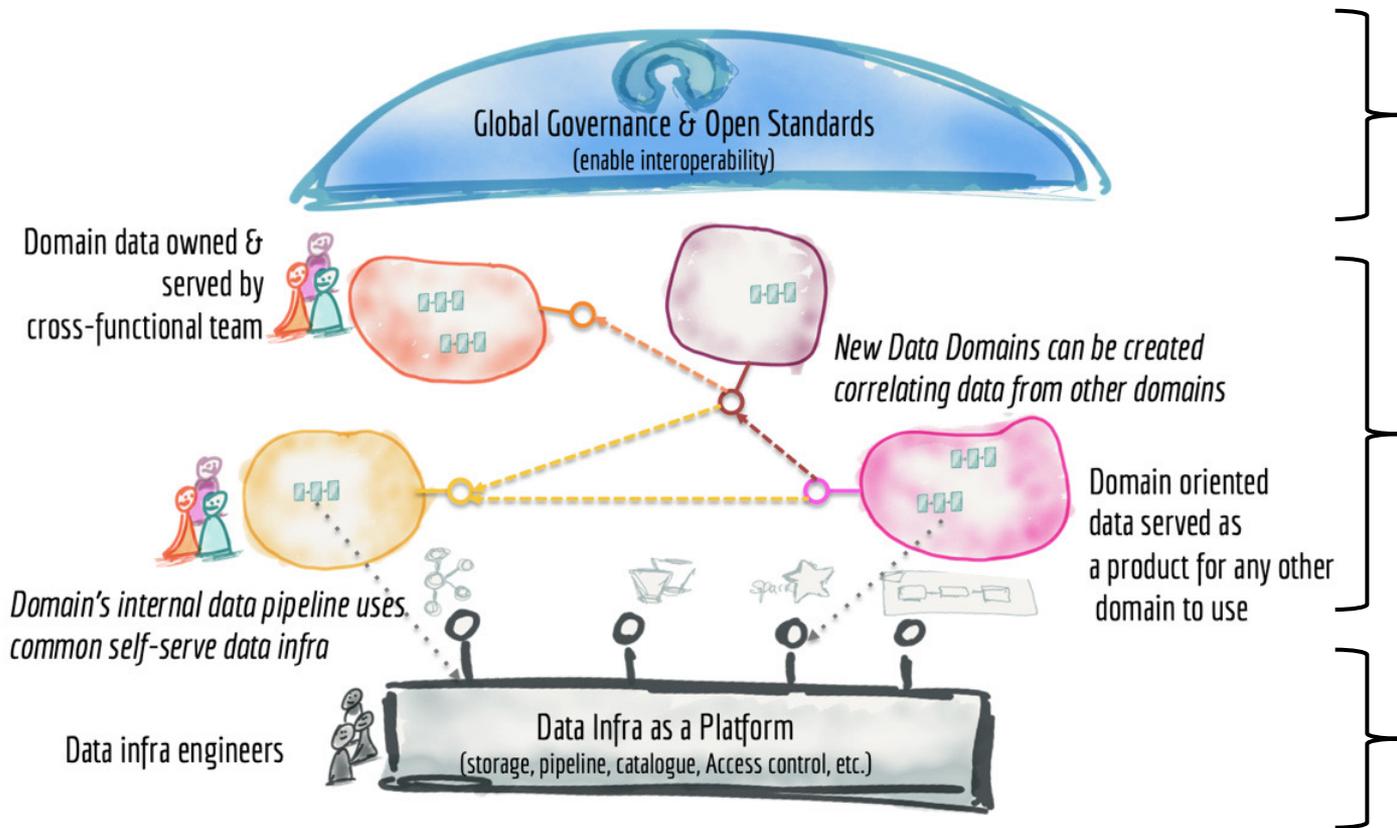
---

развитием крупного  
DWH сложно управлять

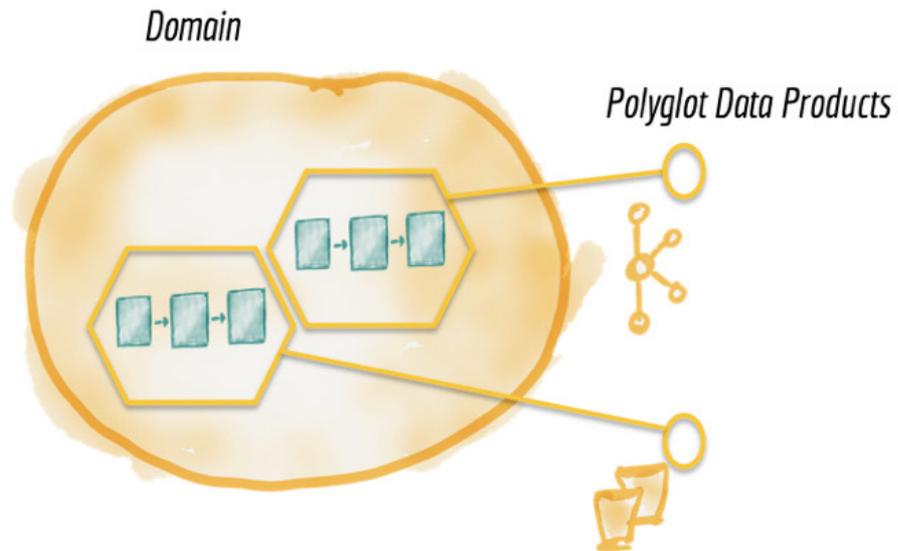
# Data Mesh



# Data Mesh



# Данные – это продукт



- DISCOVERABLE 
- ADDRESSABLE 
- TRUSTWORTHY  
(DEFINED & MONITORED SLOs) 
- SELF-DESCRIBING 
- INTER OPERABLE  
(GOVERNED BY OPEN STANDARDS) 
- SECURE 

## **I. Проблема:**

развитием крупного  
DWH сложно управлять

## **II. Решение:**

покрыть работу DWH  
метриками

# Данные – это продукт

Витрины, измерения, любые наборы данных – это продукт

Аналитики, DS, ML-специалисты, менеджеры – пользователи продукта

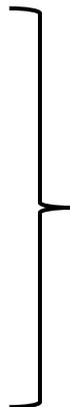
- › Удобство использования
- › Частота использования
- › Легкость обнаружения
- › Качество данных
- › Понятное описание семантики
- › Интегрируемость данных и стандарты

# Данные – это продукт

Витрины, измерения, любые наборы данных – это продукт

Аналитики, DS, ML-специалисты, менеджеры – пользователи продукта

- › Удобство использования
- › Частота использования
- › Легкость обнаружения
- › Качество данных
- › Понятное описание семантики
- › Интегрируемость данных и стандарты



Покроем метриками

# Данные – это продукт

Продуктовая команда – независимая единица поставки счастья

## Data Partner



- Владелец данных (= продукта)
  - › Коммуникации с пользователем
  - › Управление требованиями
  - › Развитие домена
  - › Постановка задач
  - › Создание метаданных витрин/отчетов

## Data Engineer



- Разработчик данных (= продукта)
  - › Выполнение задач на разработку
  - › Реализация ETL/ELT на базе платформы
  - › Создание сложных алгоритмов агрегации данных/подсчетов
  - › Физическая реализация метаданных на доступных инструментах

# Данные – это продукт

Продуктовая команда – независимая единица поставки счастья

**Data Partner**



- Владелец данных (= продукта)
  - › Коммуникации с пользователем
  - › Управление требованиями
  - › Развитие домена
  - › Постановка задач
  - › Создание метаданных витрин/отчетов

**Data Engineer**



- Разработчик данных (= продукта)
  - › Выполнение задач на разработку
  - › Реализация ETL/ELT на базе платформы
  - › Создание сложных алгоритмов агрегации данных/подсчетов
  - › Физическая реализация метаданных на доступных инструментах

Работу продуктовых команд будем оценивать через метрики

## **I. Проблема:**

развитием крупного  
DWH сложно управлять

## **II. Решение:**

покрыть работу DWH  
метриками

## **III. Идея:**

использовать данные  
систем DWH в самом  
DWH

(«DWH для DWH»)

# DHW для DWH

Почему бы не рассмотреть DWH как источник информации для самого DWH?

## Транзакционная информация

---

Что происходит?

- › Логи обращения к Greenplum
- › Логи обращения к YТ
- › Логи обращения к Tableau
- › Логи обращения к MS SSAS
- › Логи ошибок по объектам
- › Информация об отставании данных

## Статическая информация

---

С чем происходит?

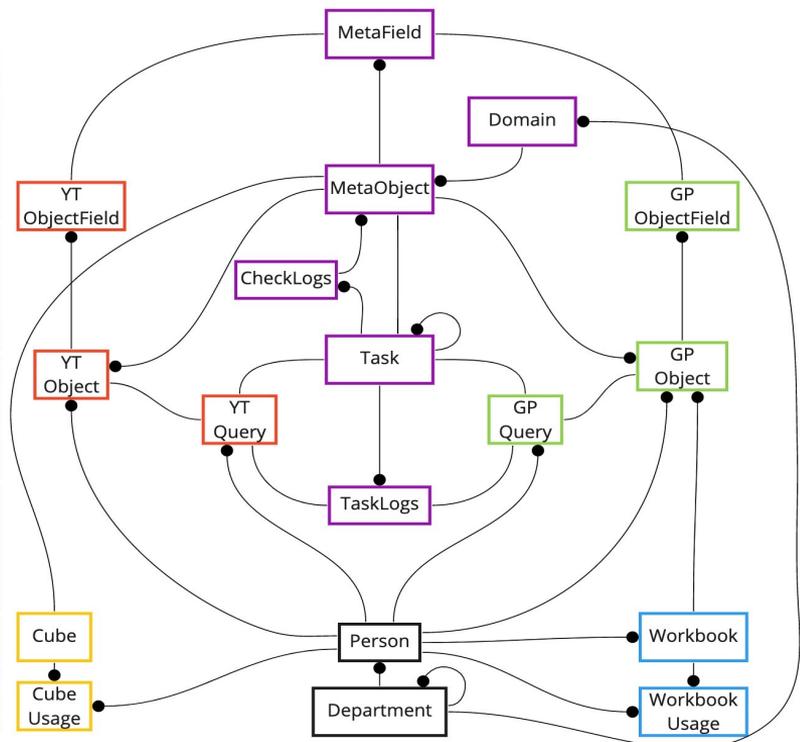
- › Метаданные систем
- › Данные из системы учета пользователей
- › Метаданные из нашего репозитория метаданных
- › Граф связей между задачами ETL-процессов

# MetaDWH

## Source Domain



## Core Domain



## Business Domain

### Техническая информация



### Использование объектов



### Витрины с метаданными



# Greenplum

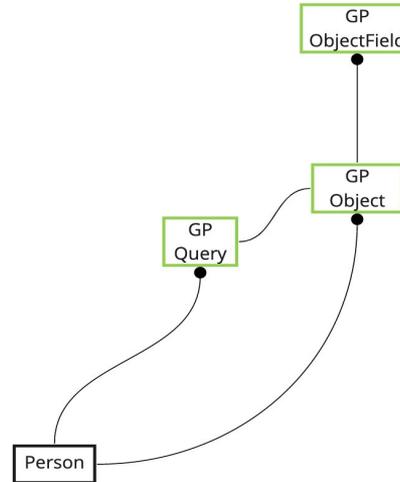
## Source Domain

### Greenplum

Логи  
использования

Метаданные  
объектов

## Core Domain



## Business Domain

# Staff

## Source Domain

### Greenplum

Логи  
использования

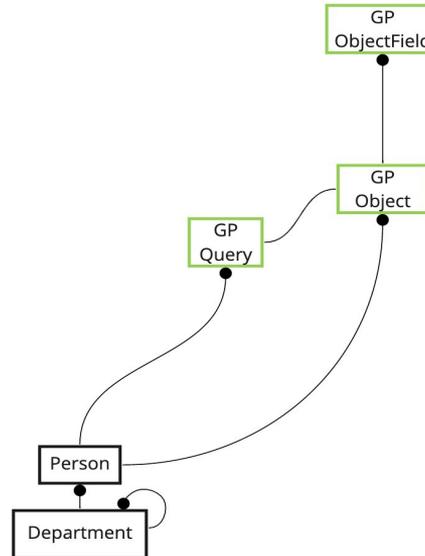
Метаданные  
объектов

### Staff

Профиль  
пользователя

Орг. Структура

## Core Domain



## Business Domain

# YT

## Source Domain

### Greenplum

Логи  
использования

Метаданные  
объектов

### YT

Логи  
использования

Метаданные  
объектов

### Staff

Профиль  
пользователя

Орг. Структура

## Core Domain

YT  
ObjectField

YT  
Object

YT  
Query

GP  
Query

GP  
ObjectField

GP  
Object

Person

Department

## Business Domain

# MS SSAS

## Source Domain

### Greenplum

Логи  
использования

Метаданные  
объектов

### YT

Логи  
использования

Метаданные  
объектов

### MS SSAS

Логи  
использования

Метаданные  
объектов

### Staff

Профиль  
пользователя

Орг.Структура

## Core Domain

YT  
ObjectField

YT  
Object

YT  
Query

GP  
Query

GP  
ObjectField

GP  
Object

Cube

Cube  
Usage

Person

Department

## Business Domain

# Потребление ресурсов

## Source Domain

### Greenplum

Логи использования

Метаданные объектов

### YT

Логи использования

Метаданные объектов

### MS SSAS

Логи использования

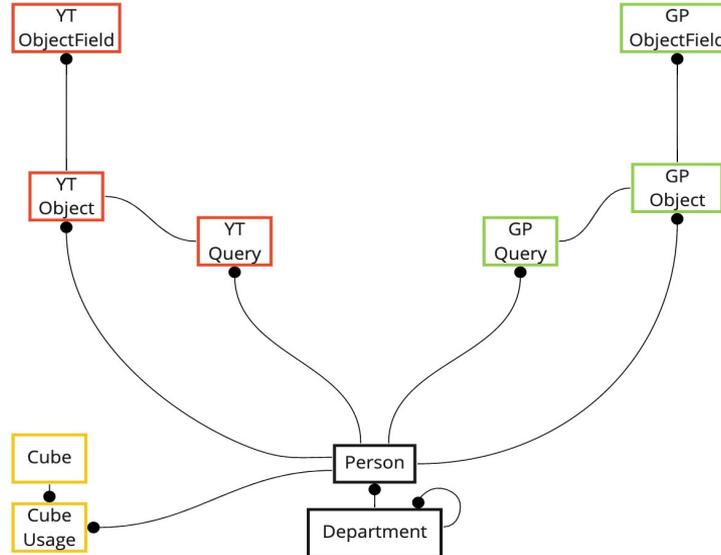
Метаданные объектов

### Staff

Профиль пользователя

Орг. Структура

## Core Domain



## Business Domain

### Техническая информация

Витрина по размеру данных

Витрина по потреблению ресурсов

# Tableau

## Source Domain

### Greenplum

Логи использования

Метаданные объектов

### Tableau

Логи использования

Метаданные объектов

### Staff

Профиль пользователя

Опр. Структура

### YT

Логи использования

Метаданные объектов

### MS SSAS

Логи использования

Метаданные объектов

## Core Domain

YT  
ObjectField

YT  
Object

YT  
Query

GP  
Query

GP  
ObjectField

GP  
Object

Cube

Cube  
Usage

Person

Department

Workbook

Workbook  
Usage

## Business Domain

### Техническая информация

Витрина по размеру данных

Витрина по потреблению ресурсов

# Использование объектов

## Source Domain

### Greenplum

Логи использования

Метаданные объектов

### Tableau

Логи использования

Метаданные объектов

### Staff

Профиль пользователя

Орг. Структура

### YT

Логи использования

Метаданные объектов

### MS SSAS

Логи использования

Метаданные объектов

## Core Domain

YT  
ObjectField

YT  
Object

YT  
Query

GP  
Query

GP  
ObjectField

GP  
Object

Cube

Cube  
Usage

Person

Department

Workbook

Workbook  
Usage

## Business Domain

### Техническая информация

Витрина по размеру данных

Витрина по потреблению ресурсов

### Использование объектов

Витрина по использованию объектов

Витрина по использованию отчетов

# Метаданные объектов

## Source Domain

### Greenplum

Логи использования

Метаданные объектов

### Tableau

Логи использования

Метаданные объектов

### Staff

Профиль пользователя

Орг. Структура

### YT

Логи использования

Метаданные объектов

### MS SSAS

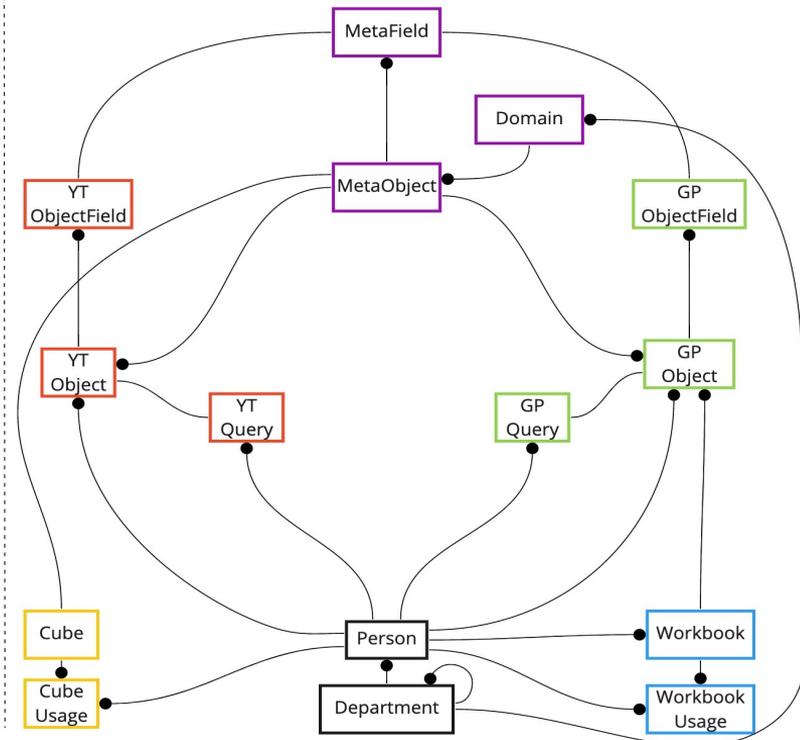
Логи использования

Метаданные объектов

### Platform

Метаданные объектов

## Core Domain



## Business Domain

### Техническая информация

Витрина по размеру данных

Витрина по потреблению ресурсов

### Использование объектов

Витрина по использованию объектов

Витрина по использованию отчетов



# Знания о запусках ETL

## Source Domain

### Greenplum

- Логи использования
- Метаданные объектов

### Tableau

- Логи использования
- Метаданные объектов

### Staff

- Профиль пользователя

### Platform

- Метаданные тасок
- Логи запусков

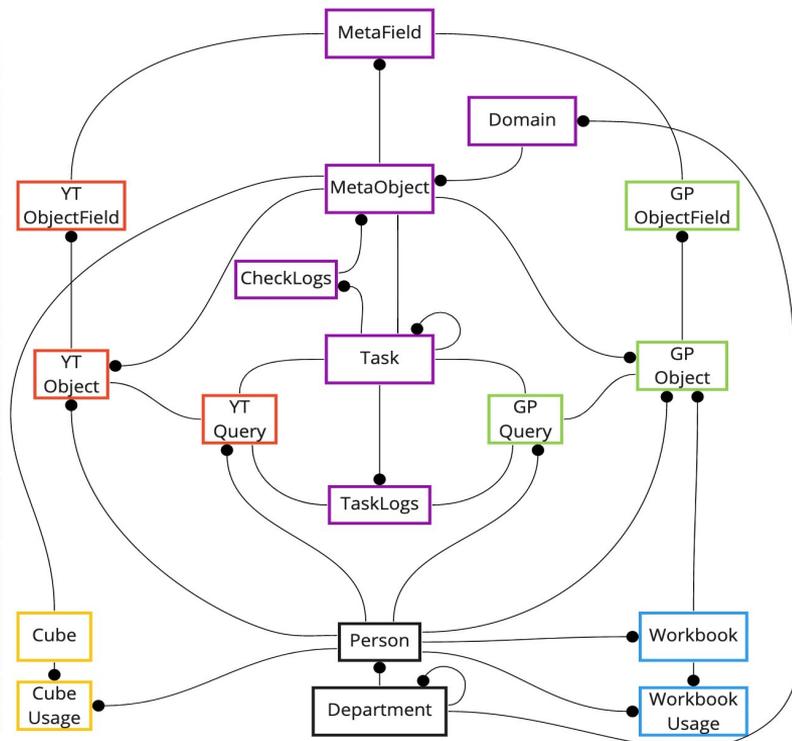
### YT

- Логи использования
- Метаданные объектов

### MS SSAS

- Логи использования
- Метаданные объектов

## Core Domain



## Business Domain

### Техническая информация

- Витрина по размеру данных
- Витрина по потреблению ресурсов

### Использование объектов

- Витрина по использованию объектов
- Витрина по использованию отчетов

### Витрины с метаданными

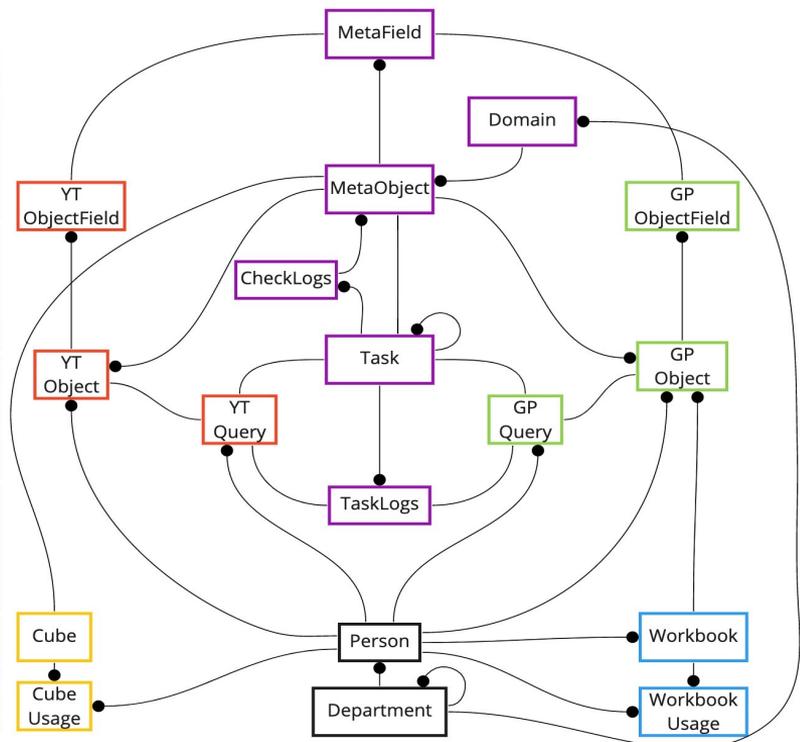
- Метаданные объектов

# MetaDWH

## Source Domain



## Core Domain



## Business Domain

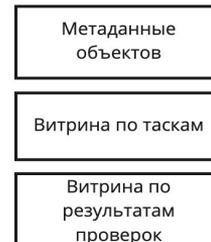
### Техническая информация



### Использование объектов



### Витрины с метаданными



## **I. Проблема:**

развитием крупного  
DWH сложно управлять

## **II. Решение:**

покрыть работу DWH  
метриками

## **III. Идея:**

использовать данные  
систем DWH в самом  
DWH

(«DWH для DWH»)

## **IV. Результат:**

аналитика по работе и  
развитию самого DWH

# Чем пользуются в хранилище?

III. Что получили?

# Профиль объекта

## Object profile

Datasource updated at 15.09.2020 12:26:12

Object usage  
Last 3 months

Include Robots  
True

DMP activity  
False

Sort by  
Requests

System  
YT

Object  
eda\_etl\_new / ...

Apply

---

### General info

**Object** dm\_order

**Path** //home/eda-dwh/cdm/order/dm\_order

**System** ■ yt

**Layer** ■ cdm

**Domain** [Business Domain] eda\_etl\_new.order

**DMP Team**

**Production flag**  True

**CDM or REP flag**  True

**Deprecated flag**  False

**Tech flag**  False

### Department

Name

	Usage	Users	Daily Usage
Центр организационного управления Проект 24	1 344	1	[Bar chart]
Центр Бизнес-анализов	878	3	[Bar chart]
Центр организационного управления Проект	850	1	[Bar chart]
Центр организационного управления Проект	728	1	[Bar chart]
Центр организационного управления Проект	677	5	[Bar chart]
Центр организационного управления Проект	555	1	[Bar chart]
Центр организационного управления Проект	514	2	[Bar chart]
Центр организационного управления Проект	270	3	[Bar chart]
Центр организационного управления Проект 24	86	1	[Bar chart]
Центр организационного управления Проект 24	55	1	[Bar chart]
Центр организационного управления Проект 24	50	1	[Bar chart]
Центр организационного управления Проект 24	29	1	[Bar chart]

### Fields

Name	Type	Description
cancel_reason_group_..	String	Название группы
cancel_reason_id	Int	Идентификатор прич..
cancel_reason_name	String	Описание причины от..
cancel_reasons_syste..	Boolean	Флаг автоматической..
cancelled_lat	Double	Широта координат ку..
cancelled_lon	Double	Долгота координат ку..
commission_value_w_v..	Double	Сумма комиссии, кот..
confirmed_flg	Boolean	Индикатор подтверж..
cooking_type	String	Тип готовки. *Опреде..
corp_order_flg	Boolean	Индикатор корпорати..
country_id	Int	Идентификатор стра..
country_name	String	Наименование страны
courier_assigned_lat	Double	Широта координат ку..
courier_assigned_lon	Double	Долгота координат ку..
courier_balance_id	Int	Идентификатор курь..
courier_delay_sec	Int	Опоздание курьера в..
courier_id	Int	Идентификатор курь..
courier_selfemployed_..	Boolean	Индикатор самозанят..
courier_service_id	Int	Идентификатор курь..
courier_service_income..	Int	Идентификатор дохо..
courier_service_name	String	Наименование курье..
courier_type	String	Тип курьера (пеший, ..
courier_type_code	String	Тип передвижения ку..
courier_username	String	Фамилия и имя курье..

6 077  
Usage

35  
Users

**Daily usage dynamic**

### User

Staff login

Staff login	Usage	Daily Usage
admin	1 344, 22,1%	[Bar chart]
admin	872, 14,3%	[Bar chart]
admin	850, 14,0%	[Bar chart]
admin	728, 12,0%	[Bar chart]
admin	555, 9,1%	[Bar chart]
admin	475, 7,8%	[Bar chart]
admin	376, 6,2%	[Bar chart]
admin	192, 3,2%	[Bar chart]
admin	100, 1,6%	[Bar chart]
admin	87, 1,4%	[Bar chart]
admin	86, 1,4%	[Bar chart]
admin	83, 1,4%	[Bar chart]
admin	55, 0,9%	[Bar chart]

# Профиль объекта

## Object profile

Datasource updated at 15.09.2020 12:26:12

Object usage  
Last 3 months

Include Robots  
True

DMP activity  
False

Sort by  
Requests

System  
YT

Object  
eda\_etl\_new / ...

Apply

### General info

**Object** dm\_order

**Path** //home/eda-dwh/cdm/order/dm\_order

**System** ■ yt

**Layer** ■ cdm

**Domain** [Business Domain] eda\_etl\_new.order

**DMP Team** .....

**Production flag**  True

**CDM or REP flag**  True

**Deprecated flag**  False

**Tech flag**  False

6 077  
Usage

35  
Users

Daily usage dynamic

### Department

Name

	Usage	Users	Daily Usage
Центр производственного анализа	1 344	1	
Центр Бизнес-анализа	878	3	
Центр автоматизированного анализа	850	1	
Центр интеллектуального анализа данных	728	1	
Центр анализа данных	677	5	
Аналитический центр	555	1	
Аналитический центр	514	2	
Аналитический центр	270	3	
Служба анализа данных	86	1	
Центр интеллектуального анализа	55	1	
Центр анализа данных	50	1	
Центр анализа и прогнозирования	29	1	

### User

Staff login

	Usage	Daily Usage
...	1 344 22,1%	
...	872 14,3%	
...	850 14,0%	
...	728 12,0%	
...	555 9,1%	
...	475 7,8%	
...	376 6,2%	
...	192 3,2%	
...	100 1,6%	
...	87 1,4%	
...	86 1,4%	
...	83 1,4%	

### Fields

Name	Type	Description
cancel_reason_group_..	String	Название группы
cancel_reason_id	Int	Идентификатор прич..
cancel_reason_name	String	Описание причины от..
cancel_reasons_syste..	Boolean	Флаг автоматической..
cancelled_lat	Double	Широта координат ку..
cancelled_lon	Double	Долгота координат ку..
commission_value_w_v..	Double	Сумма комиссии, кот..
confirmed_flg	Boolean	Индикатор подтверж..
cooking_type	String	Тип готовки. *Опреде..
corp_order_flg	Boolean	Индикатор корпорати..
country_id	Int	Идентификатор стра..
country_name	String	Наименование страны
courier_assigned_lat	Double	Широта координат ку..
courier_assigned_lon	Double	Долгота координат ку..
courier_balance_id	Int	Идентификатор курь..
courier_delay_sec	Int	Опоздание курьера в..
courier_id	Int	Идентификатор курь..
courier_selfemployed_..	Boolean	Индикатор самозанят..
courier_service_id	Int	Идентификатор курь..
courier_service_income..	Int	Идентификатор дохо..
courier_service_name	String	Наименование курье..
courier_type	String	Тип курьера (пеший, ..
courier_type_code	String	Тип передвижения ку..
courier_username	String	Фамилия и имя курье..

# Профиль объекта

## Object profile

Datasource updated at 15.09.2020 12:26:12

Object usage  
Last 3 months

Include Robots  
True

DMP activity  
False

Sort by  
Requests

System  
YT

Object  
eda\_etl\_new / ...

Apply

### General info

**Object** dm\_order

**Path** //home/eda-dwh/cdm/order/dm\_order

**System** ■ yt

**Layer** ■ cdm

**Domain** [Business Domain] eda\_etl\_new.order

**DMP Team**

**Production flag**  True

**CDM or REP flag**  True

**Deprecated flag**  False

**Tech flag**  False

### Department

Name

	Usage	Users	Daily Usage
Центр организационного управления Проект 26	1 344	1	
Центр Бизнес-аналитики	878	3	
Центр организационного управления Проект	850	1	
Центр организационного управления Проект	728	1	
Центр организационного управления Проект	677	5	
Центр организационного управления Проект	555	1	
Центр организационного управления Проект	514	2	
Центр организационного управления Проект	270	3	
Служба организационного управления Проект 26	86	1	
Центр организационного управления Проект 26	55	1	
Центр организационного управления Проект 26	50	1	
Центр организационного управления Проект 26	29	1	

### Fields

Name	Type	Description
cancel_reason_group_...	String	Название группы
cancel_reason_id	Int	Идентификатор прич...
cancel_reason_name	String	Описание причины от...
cancel_reasons_syste..	Boolean	Флаг автоматической..
cancelled_lat	Double	Широта координат ку...
cancelled_lon	Double	Долгота координат ку...
commission_value_w_v..	Double	Сумма комиссии, кот...
confirmed_flg	Boolean	Индикатор подтверж...
cooking_type	String	Тип готовки. *Опреде...
corp_order_flg	Boolean	Индикатор корпорати...
country_id	Int	Идентификатор стра...
country_name	String	Наименование страны
courier_assigned_lat	Double	Широта координат ку...
courier_assigned_lon	Double	Долгота координат ку...
courier_balance_id	Int	Идентификатор курь...
courier_delay_sec	Int	Опоздание курьера в..
courier_id	Int	Идентификатор курь...
courier_selfemployed_..	Boolean	Индикатор самозанят...
courier_service_id	Int	Идентификатор курь...
courier_service_income..	Int	Идентификатор дохо...
courier_service_name	String	Наименование курье...
courier_type	String	Тип курьера (пеший, ..
courier_type_code	String	Тип передвижения ку...
courier_username	String	Фамилия и имя курье...

**6 077 Usage**

**35 Users**

**Daily usage dynamic**

### User

Staff login

Staff login	Usage	Daily Usage
ivanov	1 344 22,1%	
ivanov	872 14,3%	
ivanov	850 14,0%	
ivanov	728 12,0%	
ivanov	555 9,1%	
ivanov	475 7,8%	
ivanov	376 6,2%	
ivanov	192 3,2%	
ivanov	100 1,6%	
ivanov	87 1,4%	
ivanov	86 1,4%	
ivanov	83 1,4%	

# Профиль объекта

## Object profile

Datasource updated at 15.09.2020 12:26:12

Object usage  
Last 3 months

Include Robots  
True

DMP activity  
False

Sort by  
Requests

System  
YT

Object  
eda\_etl\_new / ...

Apply

### General info

**Object** dm\_order

**Path** //home/eda-dwh/cdm/order/dm\_order

**System** ■ yt

**Layer** ■ cdm

**Domain** [Business Domain] eda\_etl\_new.order

**DMP Team**

**Production flag**  True

**CDM or REP flag**  True

**Deprecated flag**  False

**Tech flag**  False

### Department

Name

	Usage	Users	Daily Usage
Центр организационного управления Проект 26	1 344	1	
Центр Бизнес-аналитики	878	3	
Центр организационного управления Проект	850	1	
Центр организационного управления Проект	728	1	
Центр организационного управления Проект	677	5	
Центр организационного управления Проект	555	1	
Центр организационного управления Проект	514	2	
Центр организационного управления Проект	270	3	
Центр организационного управления Проект 26	86	1	
Центр организационного управления Проект 26	55	1	
Центр организационного управления Проект 26	50	1	
Центр организационного управления Проект 26	29	1	

### Fields

Name	Type	Description
cancel_reason_group_..	String	Название группы
cancel_reason_id	Int	Идентификатор прич..
cancel_reason_name	String	Описание причины от..
cancel_reasons_syste..	Boolean	Флаг автоматической..
cancelled_lat	Double	Широта координат ку..
cancelled_lon	Double	Долгота координат ку..
commission_value_w_v..	Double	Сумма комиссии, кот..
confirmed_flg	Boolean	Индикатор подтверж..
cooking_type	String	Тип готовки. *Опреде..
corp_order_flg	Boolean	Индикатор корпорати..
country_id	Int	Идентификатор стра..
country_name	String	Наименование страны
courier_assigned_lat	Double	Широта координат ку..
courier_assigned_lon	Double	Долгота координат ку..
courier_balance_id	Int	Идентификатор курь..
courier_delay_sec	Int	Опоздание курьера в..
courier_id	Int	Идентификатор курь..
courier_selfemployed_..	Boolean	Индикатор самозанят..
courier_service_id	Int	Идентификатор курь..
courier_service_income..	Int	Идентификатор дохо..
courier_service_name	String	Наименование курье..
courier_type	String	Тип курьера (пеший, ..
courier_type_code	String	Тип передвижения ку..
courier_username	String	Фамилия и имя курь..

6 077

Usage

35

Users

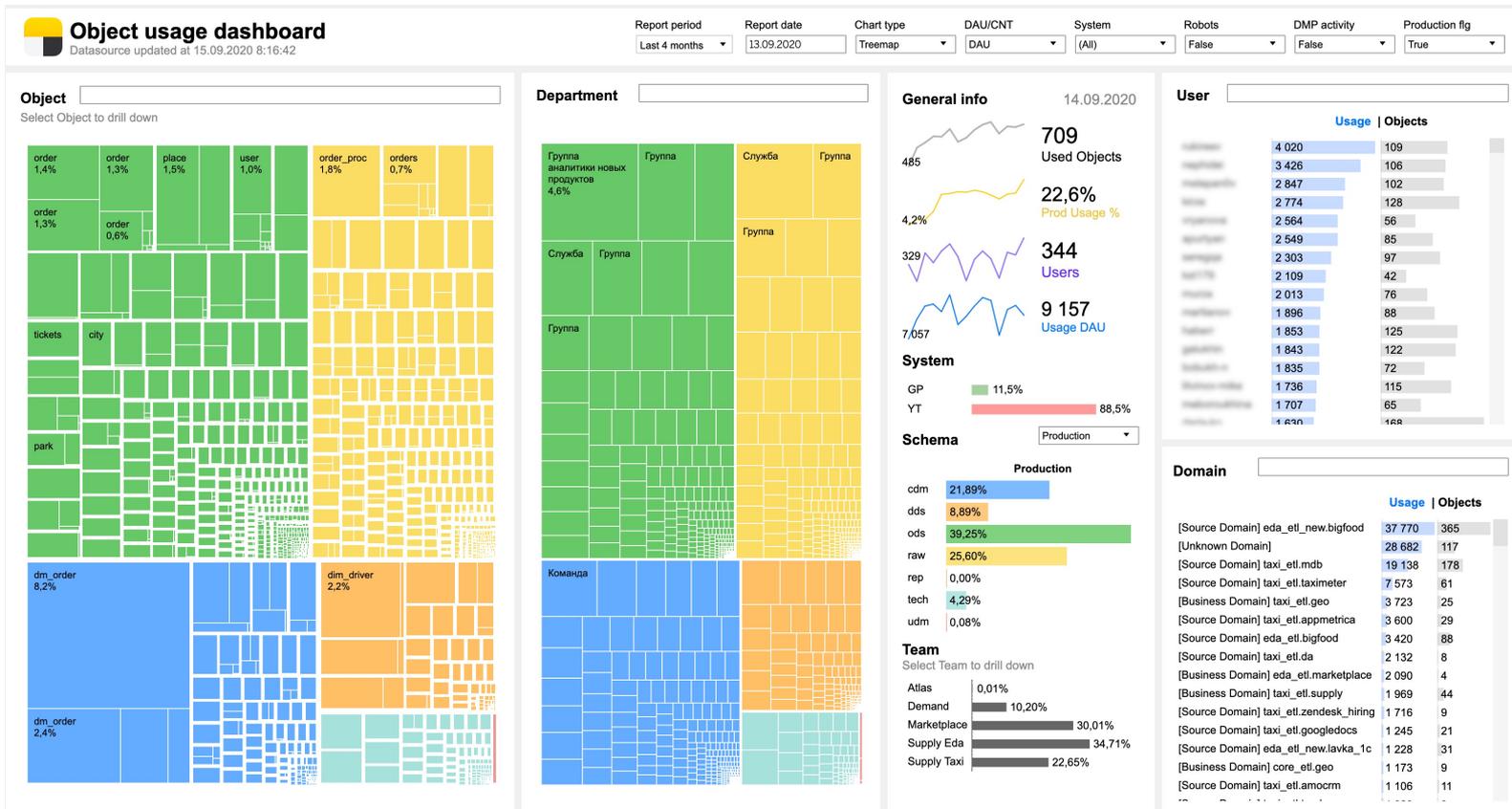
Daily usage dynamic

### User

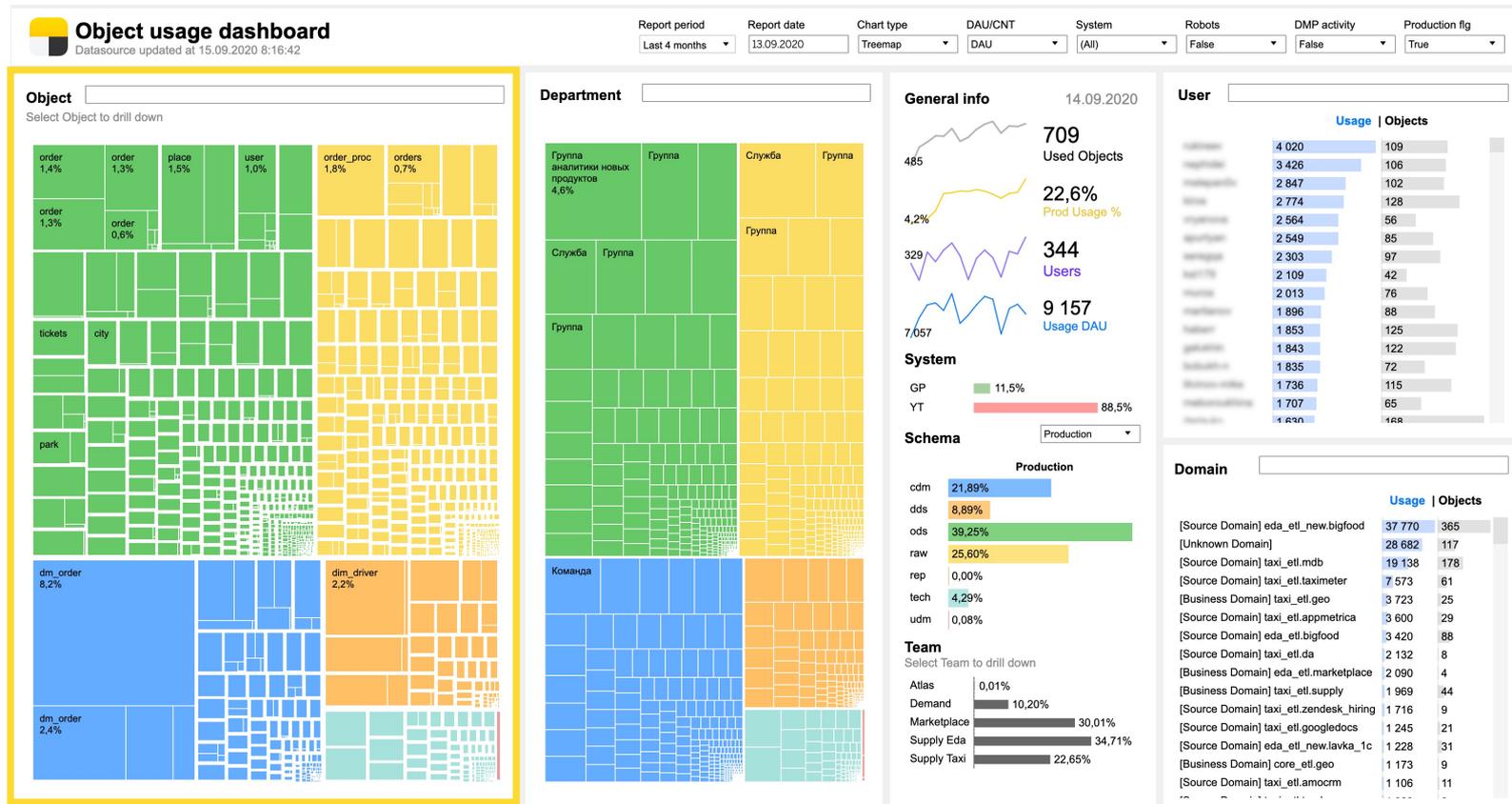
Staff login

	Usage	Daily Usage
1344	22,1%	
872	14,3%	
850	14,0%	
728	12,0%	
555	9,1%	
475	7,8%	
376	6,2%	
192	3,2%	
100	1,6%	
87	1,4%	
86	1,4%	
83	1,4%	
55	0,9%	

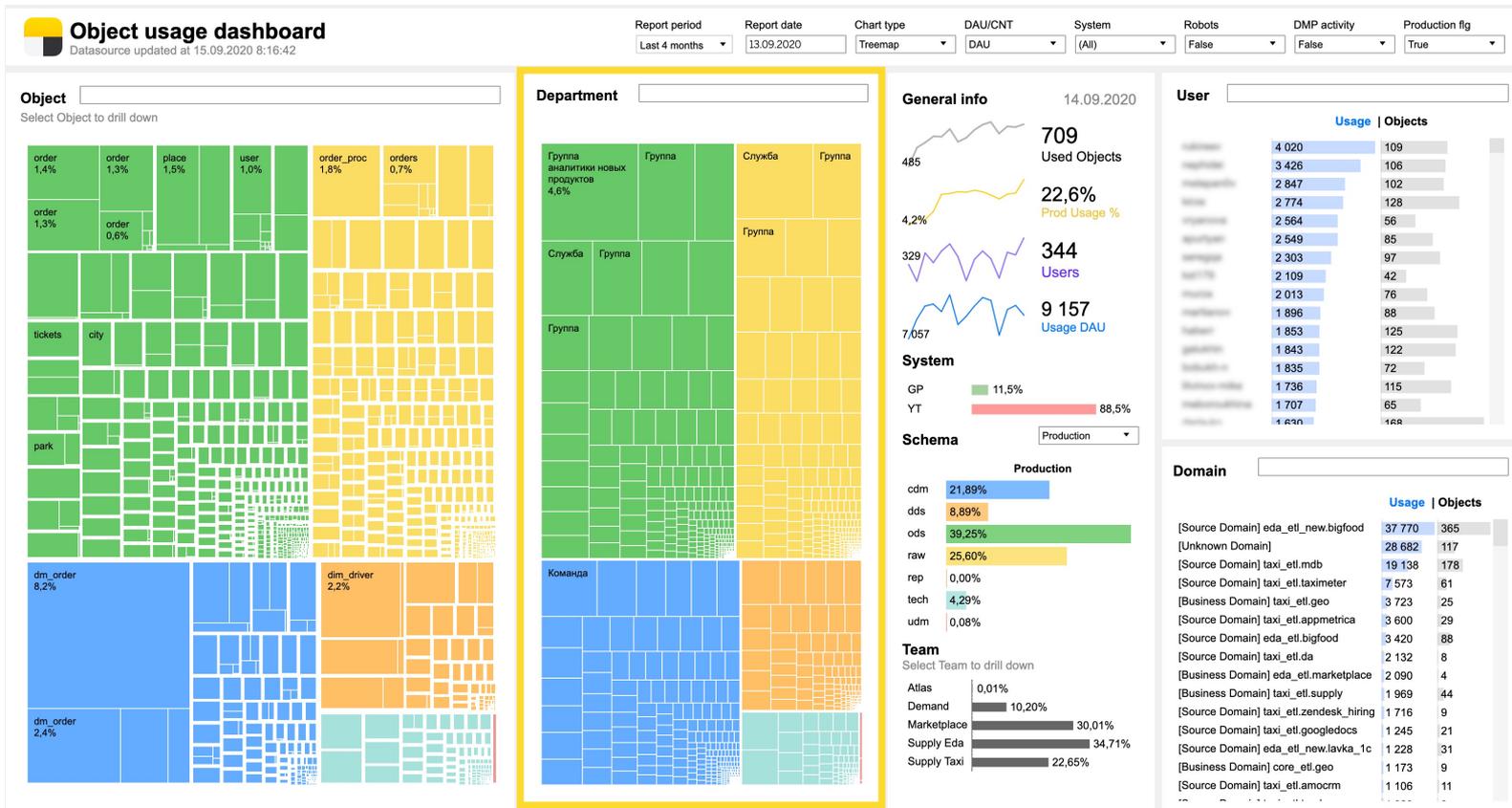
# Использование объектов



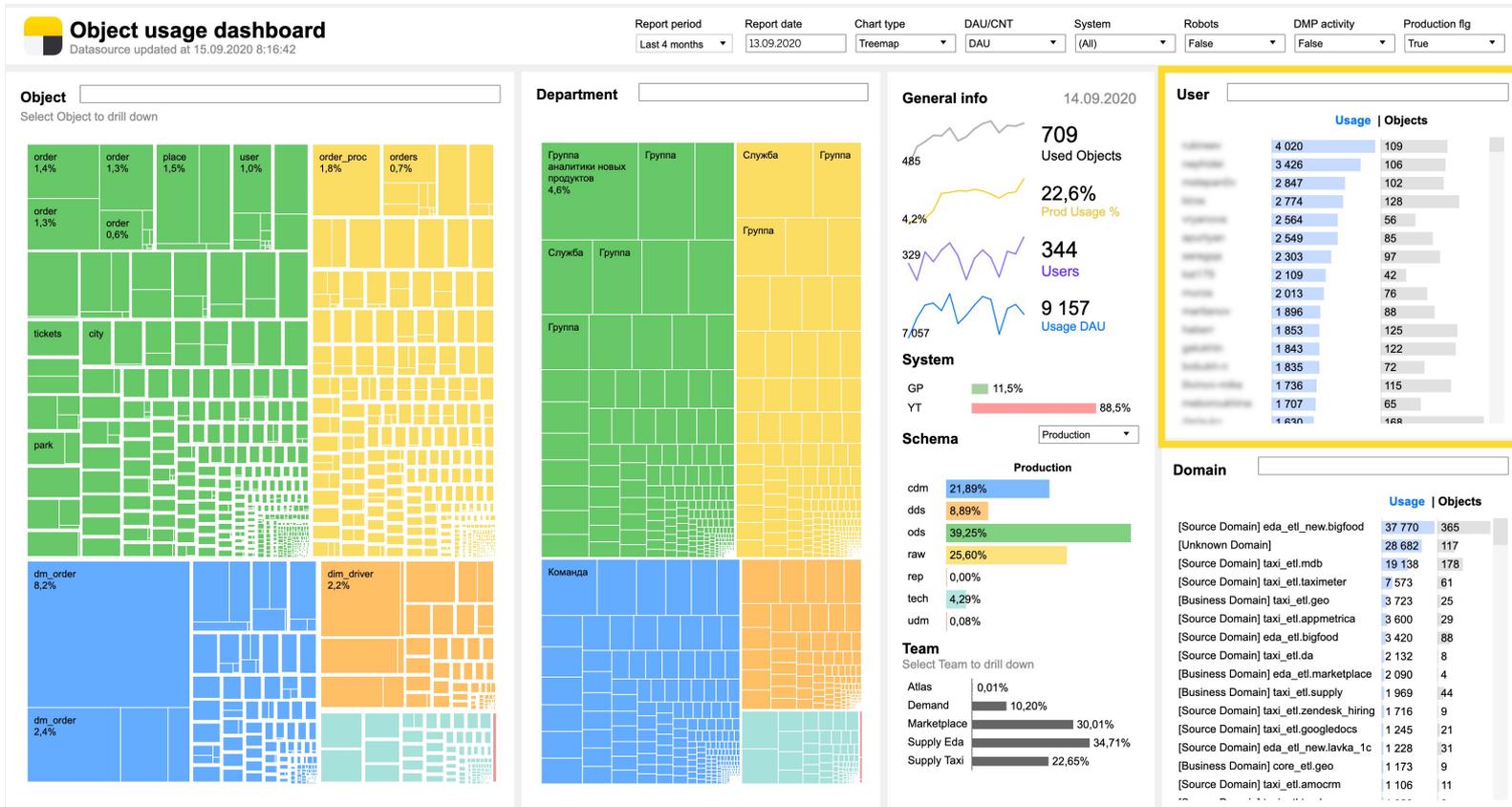
# Использование объектов



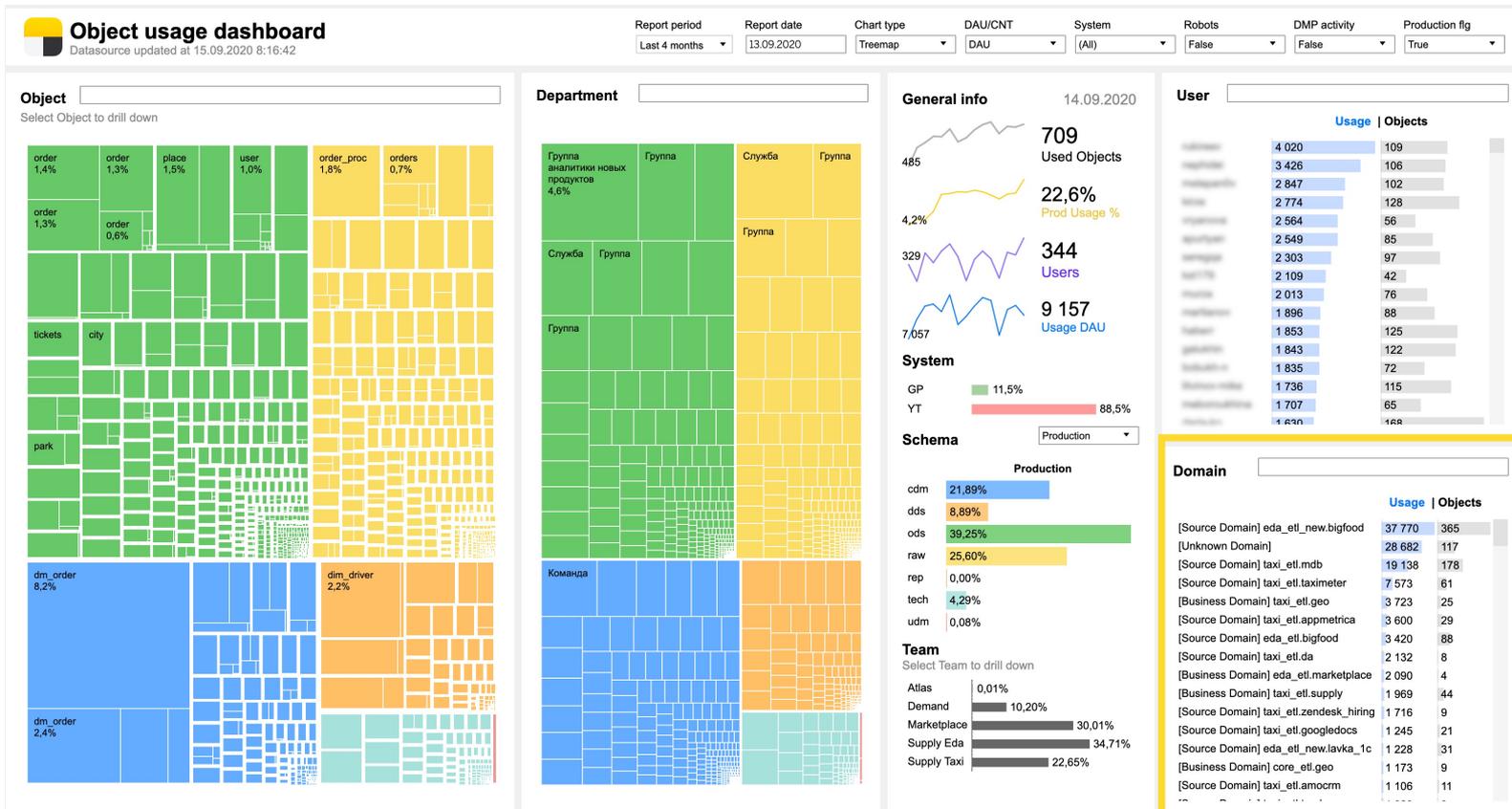
# Использование объектов



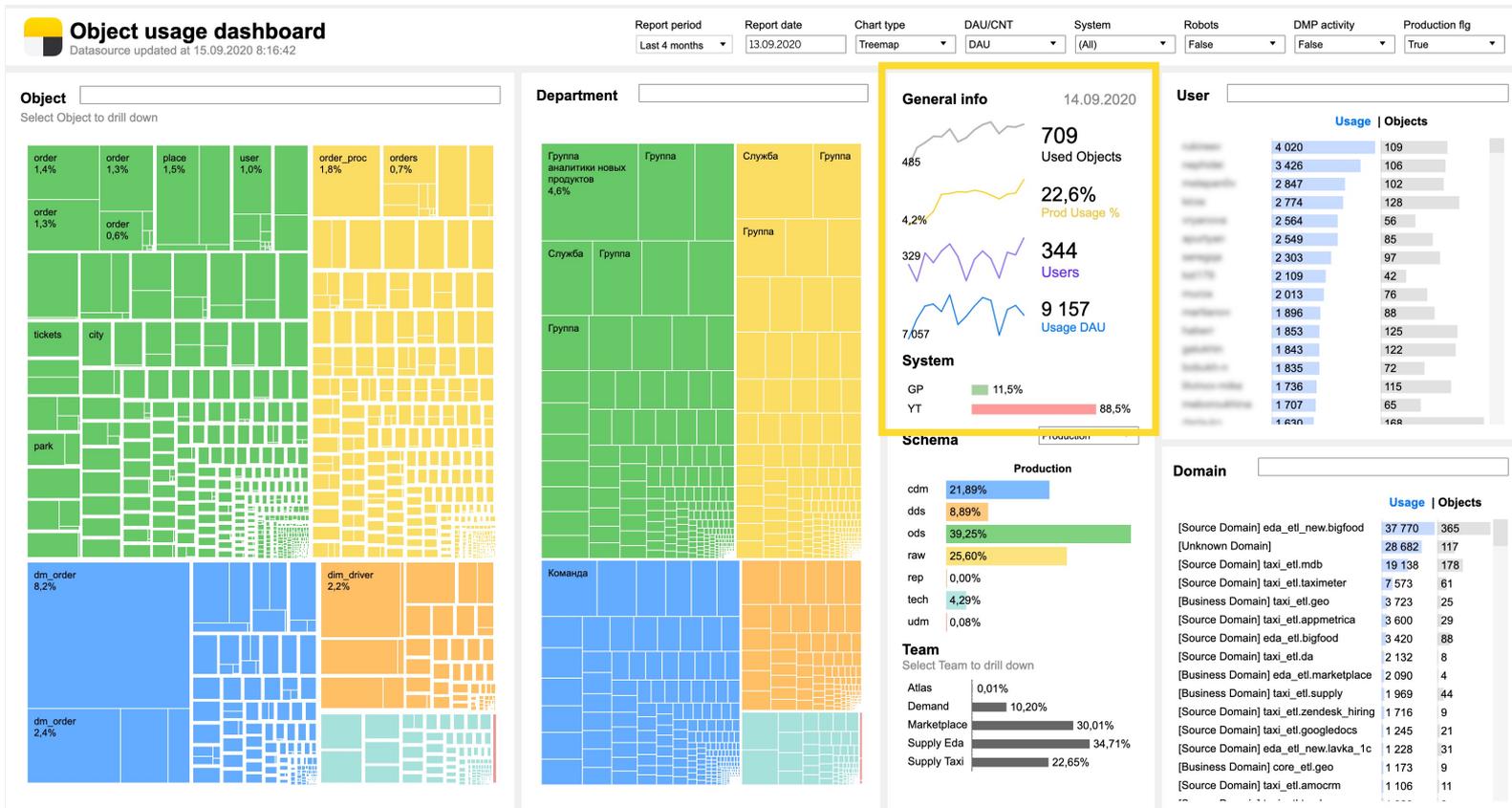
# Использование объектов



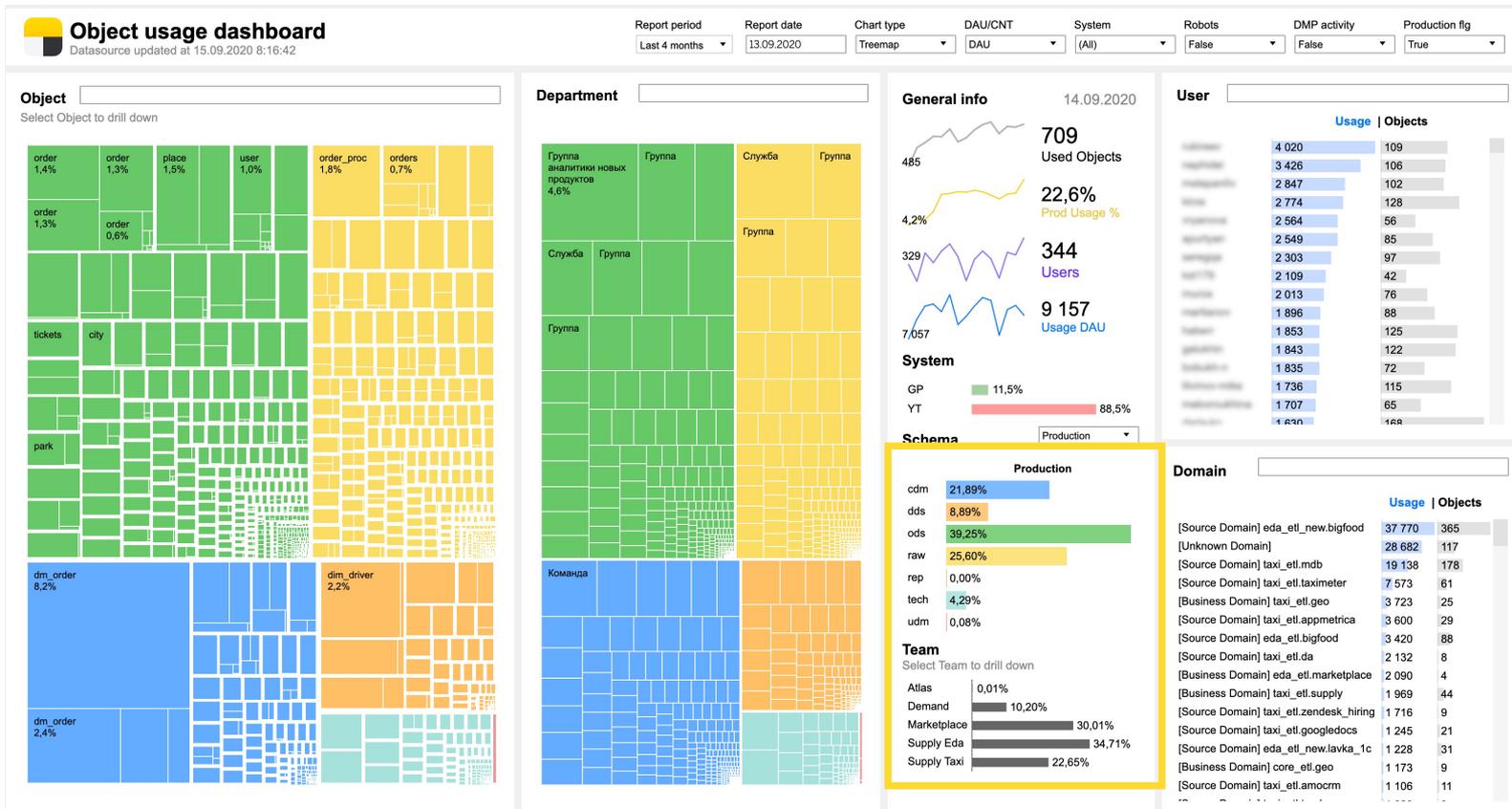
# Использование объектов



# Использование объектов



# Использование объектов



# Нотификация про изменения

При изменении объекта мы знаем, кто им пользовался, и можем точно уведомить про изменения в конкретном объекте в любом из доступных каналов коммуникации.

Пример письма:

**Действие:**

удалить

**Объект:**

//home/taxi-dwh/dds/driver\_session\_geoposition/

**Причина действия:**

В конце февраля 2021 мы остановим загрузку и перестанем поддерживать объект dds.driver\_session\_geoposition

Альтернативные объекты:

[fct\\_supply\\_state\\_hist](#) - водительские сессии, в которых собрано большинство атрибутов по активности водителей

[fct\\_taxi\\_tracker\\_position\\_enriched\\_log](#) - замена driver\_session\_geoposition, построенная на логге водительских геопозиций и fct\_supply\_state\_hist

Пожалуйста, переведите ваши процессы на новые объекты и сообщите нам о сроках, когда вы сможете запланировать переезд.

Вики как переезжать - <https://wiki.yandex-team.ru/taxi/dwh/data/business/driver-session/kak-perejiti-s-driversession-na-fctsupplystatehist/>

**Удаление запланировано на конец февраля 2021**

**Ссылка на тикет, в котором мы ведём работу над удалением/изменением:**

 **TAXIDWH-5913** [Открыт](#) [Удалить driver\\_session\\_geoposition](#) avbekker

**В этот день мы удалим/изменим объект:**

2021-02-28

# Как оценить работу продуктовой команды?

III. Что получили?

# Что нам важно?

**Результатом работы продуктовых команд (Объектами DWH) пользуются**

Считаем уникальные пары (пользователь, используемый DWH-объект) за каждый день и убираем те объекты, которыми пользуется ровно один пользователь (его личная песочница). На результирующих данных можем посчитать:

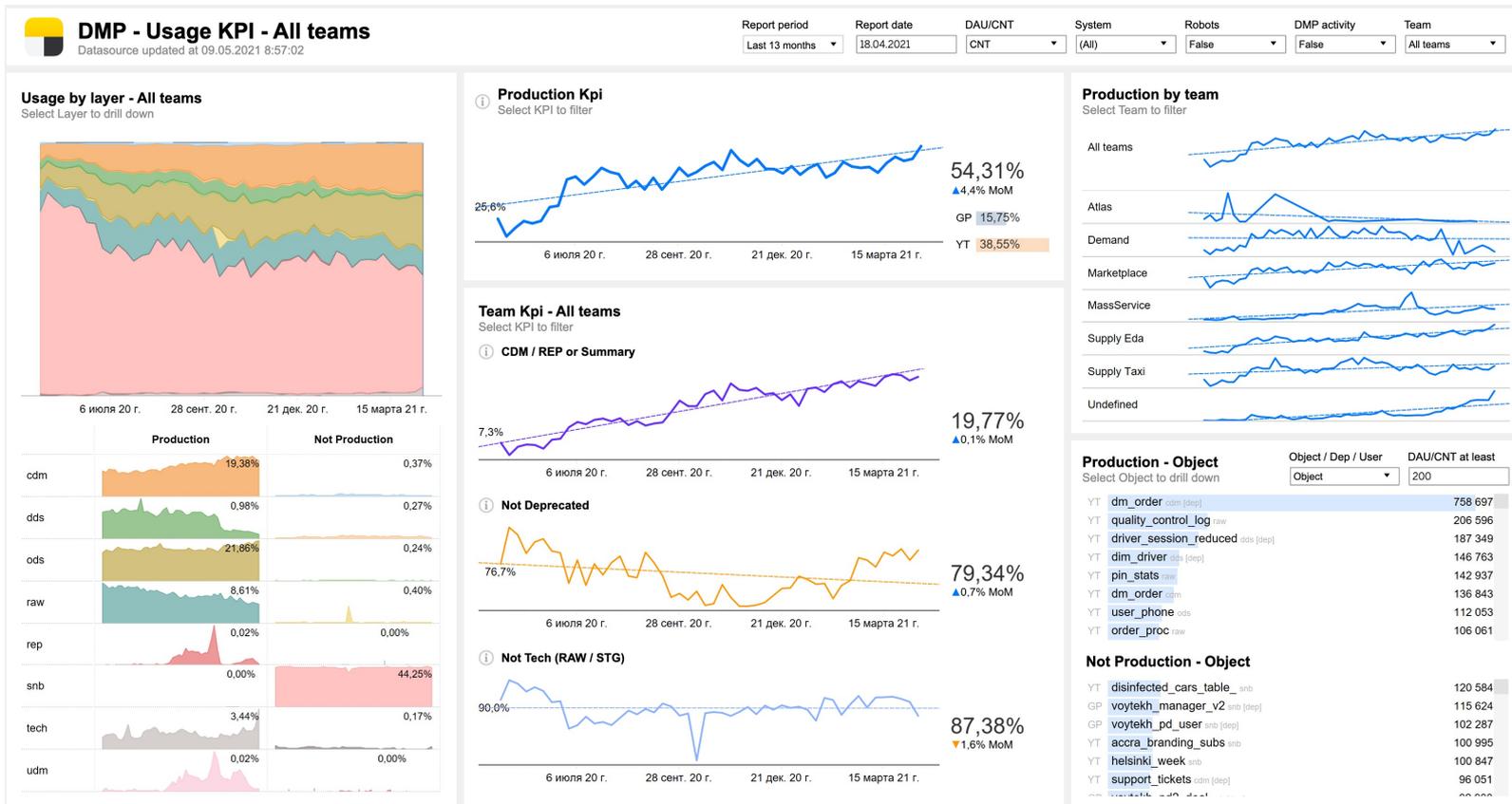
# Что нам важно?

## Результатом работы продуктовых команд (Объектами DWH) пользуются

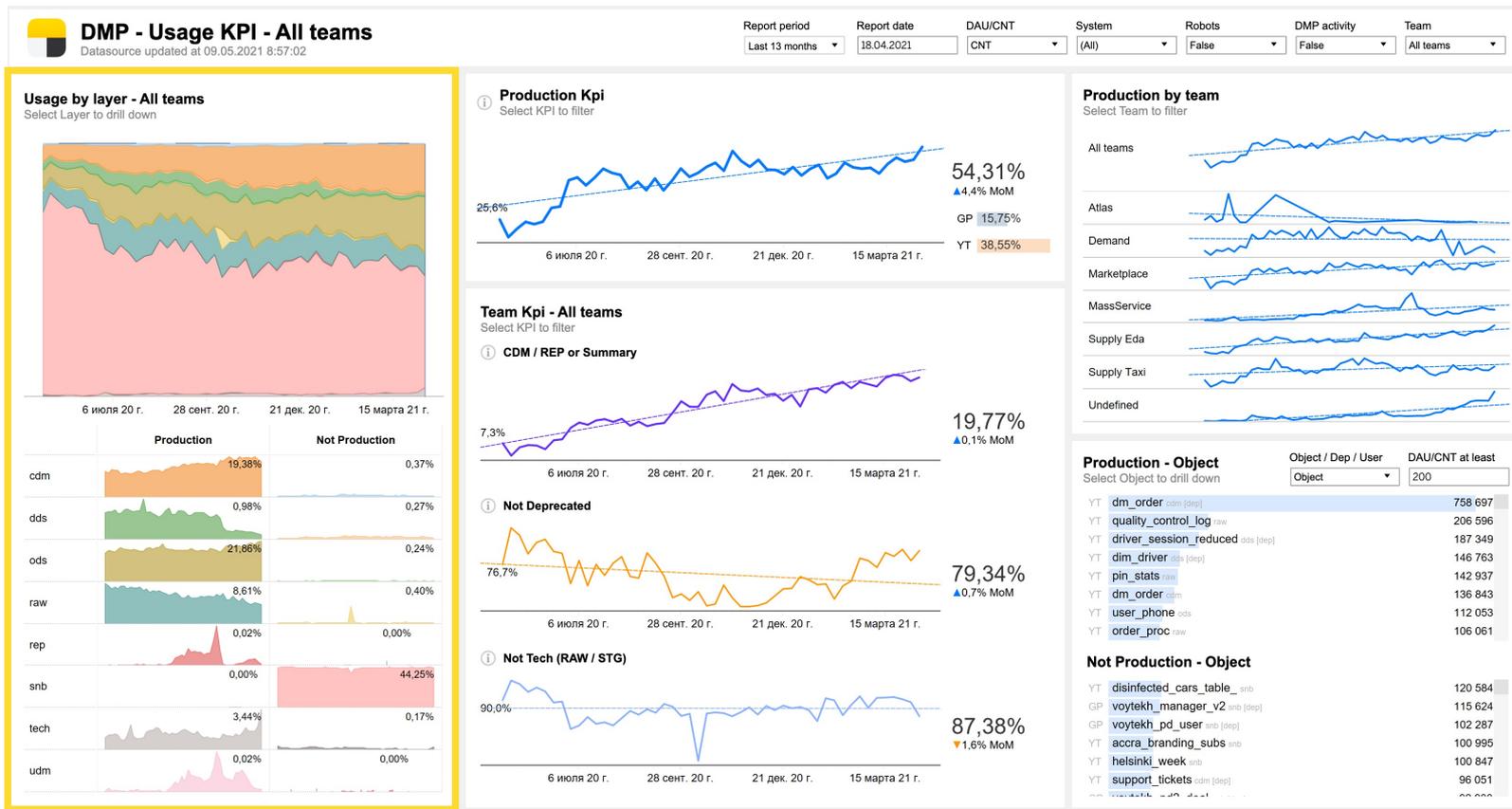
Считаем уникальные пары (пользователь, используемый DWH-объект) за каждый день и убираем те объекты, которыми пользуется ровно один пользователь (его личная песочница). На результирующих данных можем посчитать:

- соотношение обращений к prod- и не prod-объектам  
показывает, насколько пользователи смотрят в prod объекты
- соотношение обращений к deprecated- и не deprecated-объектам  
показывает, насколько мы избавляемся от легаси
- соотношение обращений tech vs all  
показывает, насколько мы быстро расшифровываем новые данные
- соотношение обращений CDM+REP vs all  
показывает, насколько наши целевые объекты удобны пользователям

# Количественные KPI команд



# Количественные KPI команд



### Production - Object

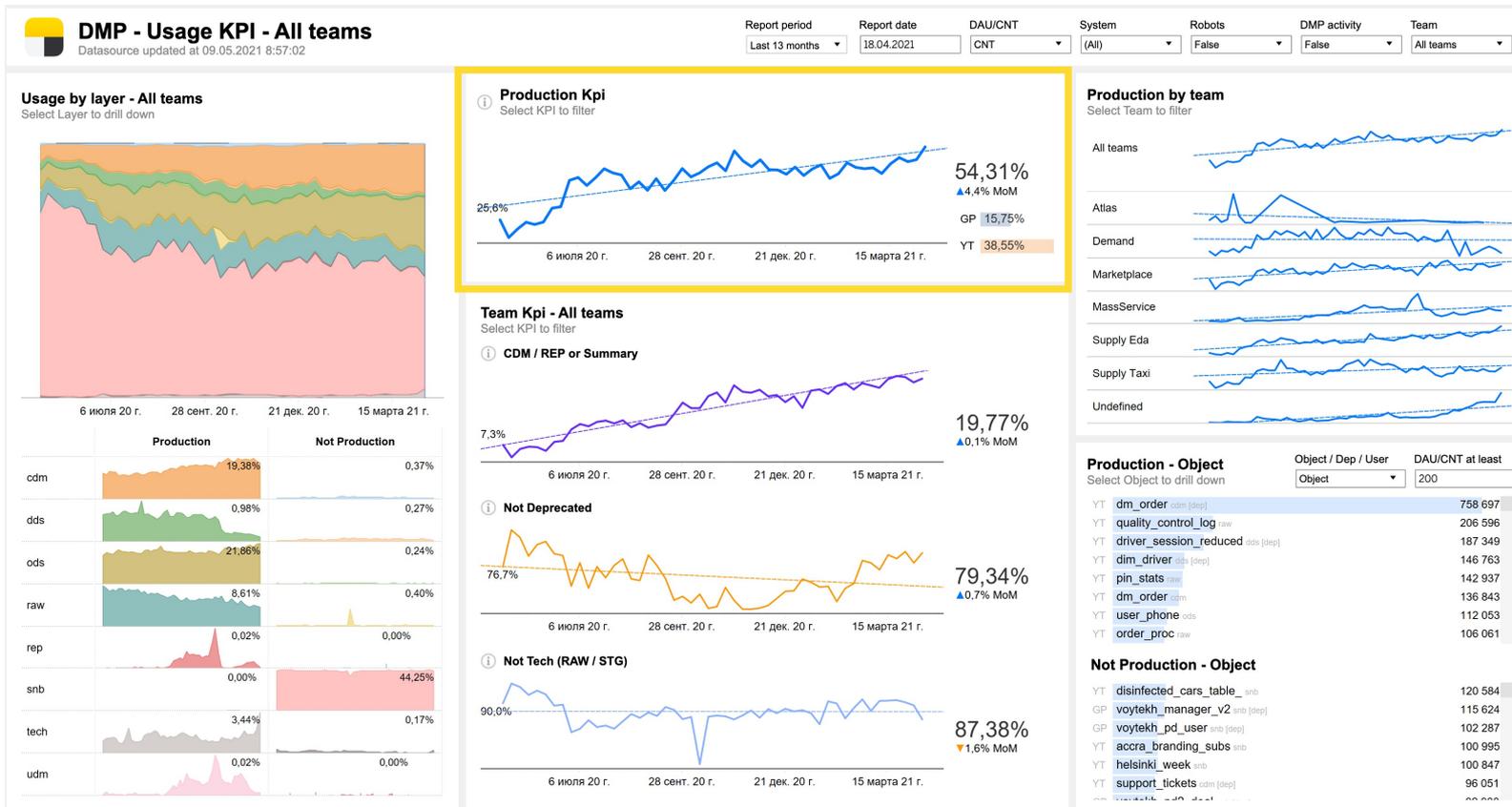
Select Object to drill down

Object / Dep / User	DAU/CNT at least
YT dm_order cdm [dep]	758 697
YT quality_control_log raw	206 596
YT driver_session_reduced dds [dep]	187 349
YT dim_driver ods [dep]	146 763
YT pin_stats raw	142 937
YT dm_order cdm	136 843
YT user_phone ods	112 053
YT order_proc raw	106 061

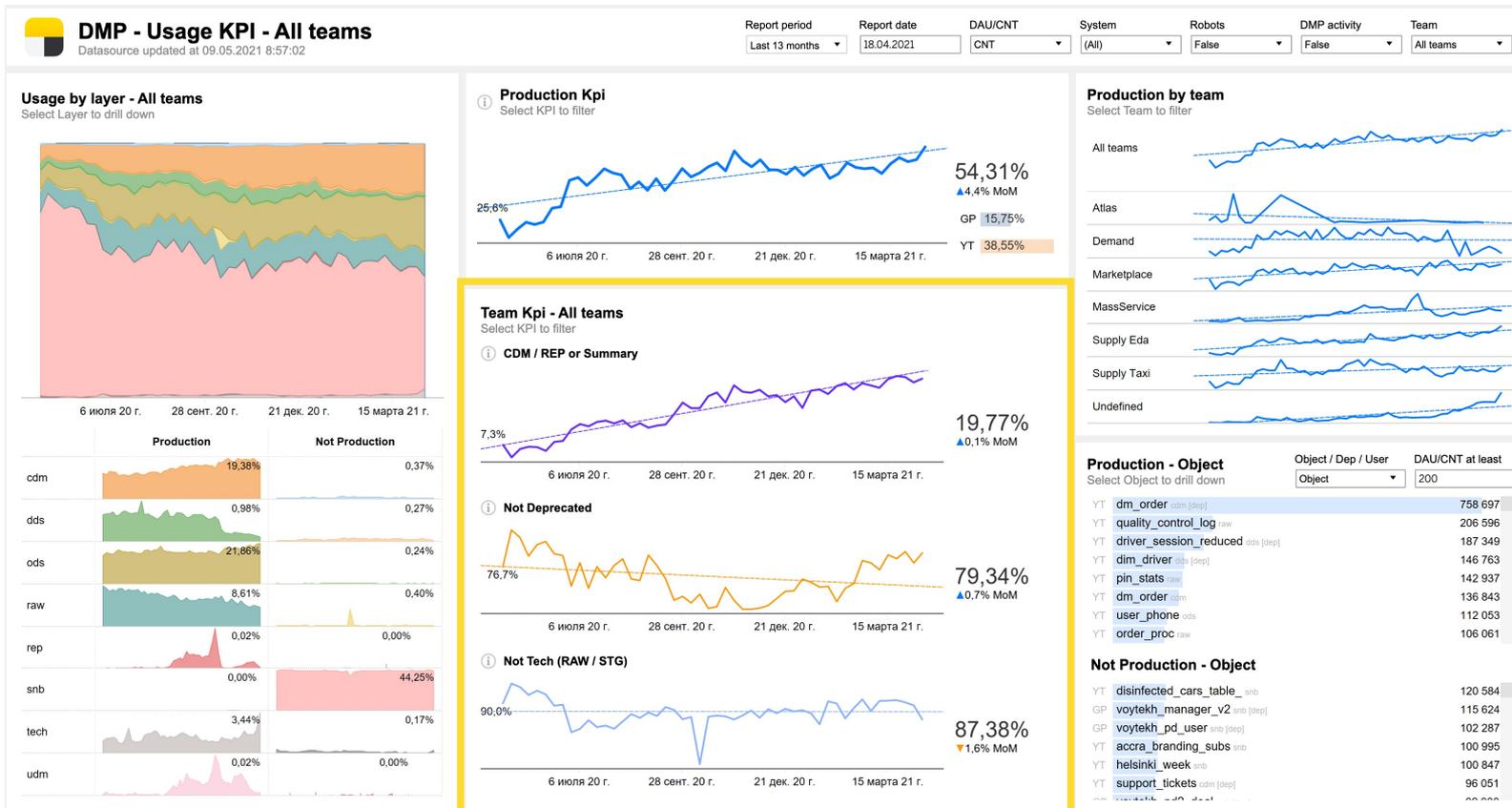
### Not Production - Object

YT disinfected_cars_table_snb	120 584
GP voytekh_manager_v2_snb [dep]	115 624
GP voytekh_pd_user_snb [dep]	102 287
YT accra_branding_subs_snb	100 995
YT helsinki_week_snb	100 847
YT support_tickets cdm [dep]	96 051

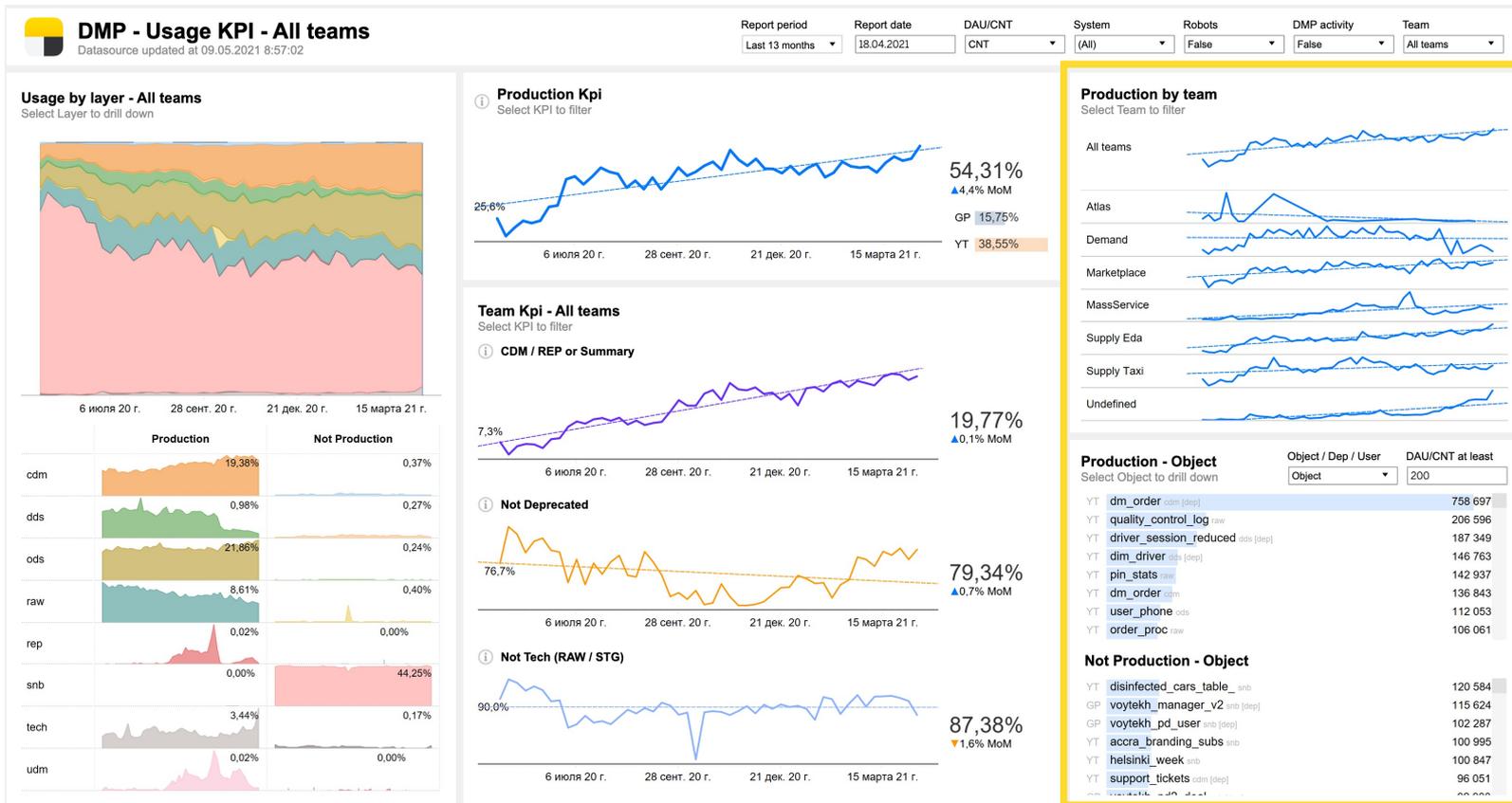
# Количественные KPI команд



# Количественные KPI команд



# Количественные KPI команд



# Оценка качества домена

Можем ввести метрики, косвенно оценивающие качество доменов

## Архитектура

- › Соблюдение naming convention
- › Использование legacy-объектов
- › Доля витрин, построенных на базе source domain (RAW|ODS), а не core domain (DDS|CDM)

## Качество данных

- › Отсутствие ПД
- › Скорость поставки данных
- › Качество документации
- › Покрытие данных проверками качества

## Качество расчетов

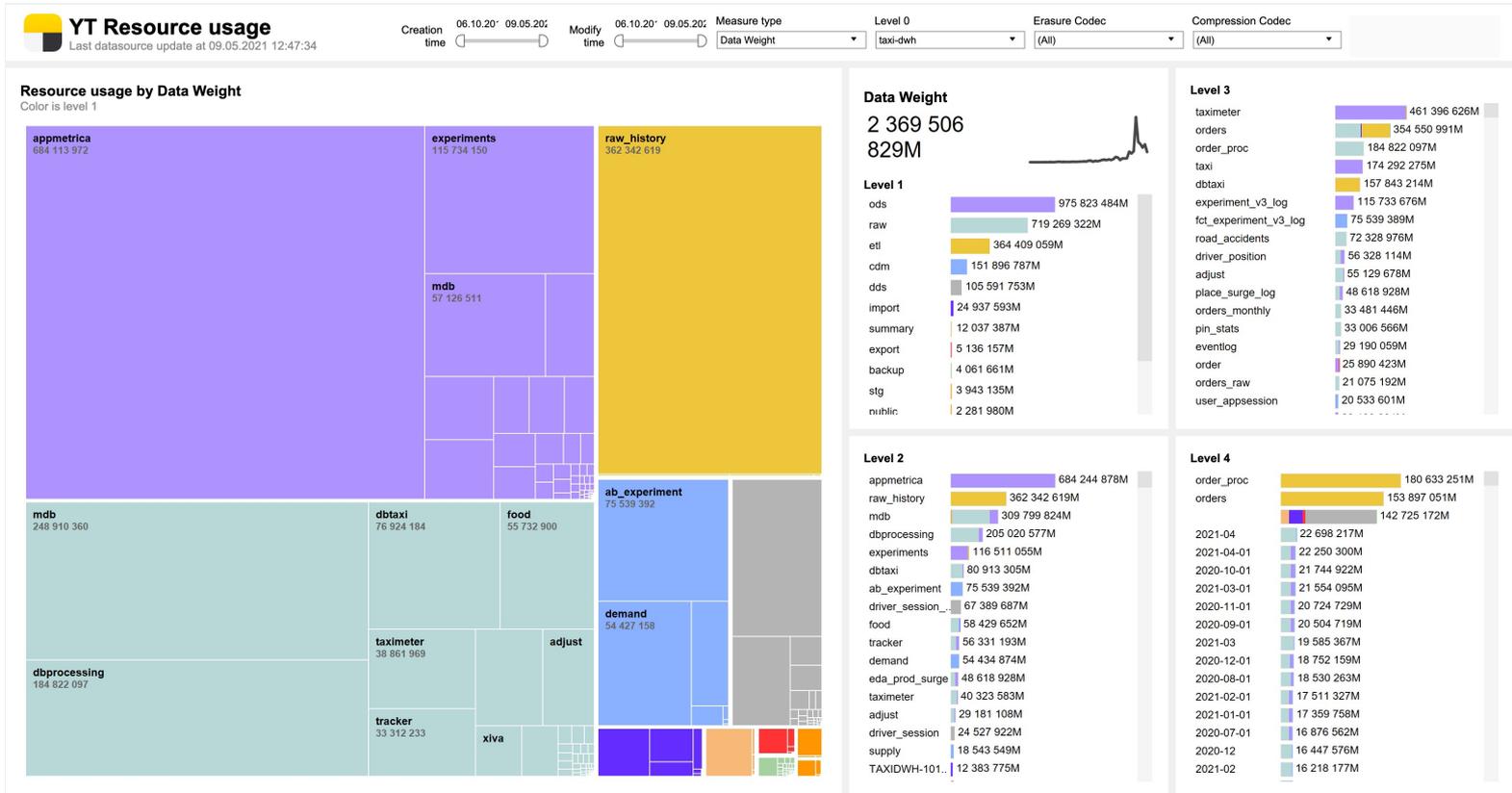
- › Использование последних инструментов платформы
- › Оптимальность ETL-процессов
- › Недоступность (downtime) объектов

Итоговая оценка качества домена как взвешенная сумма критериев

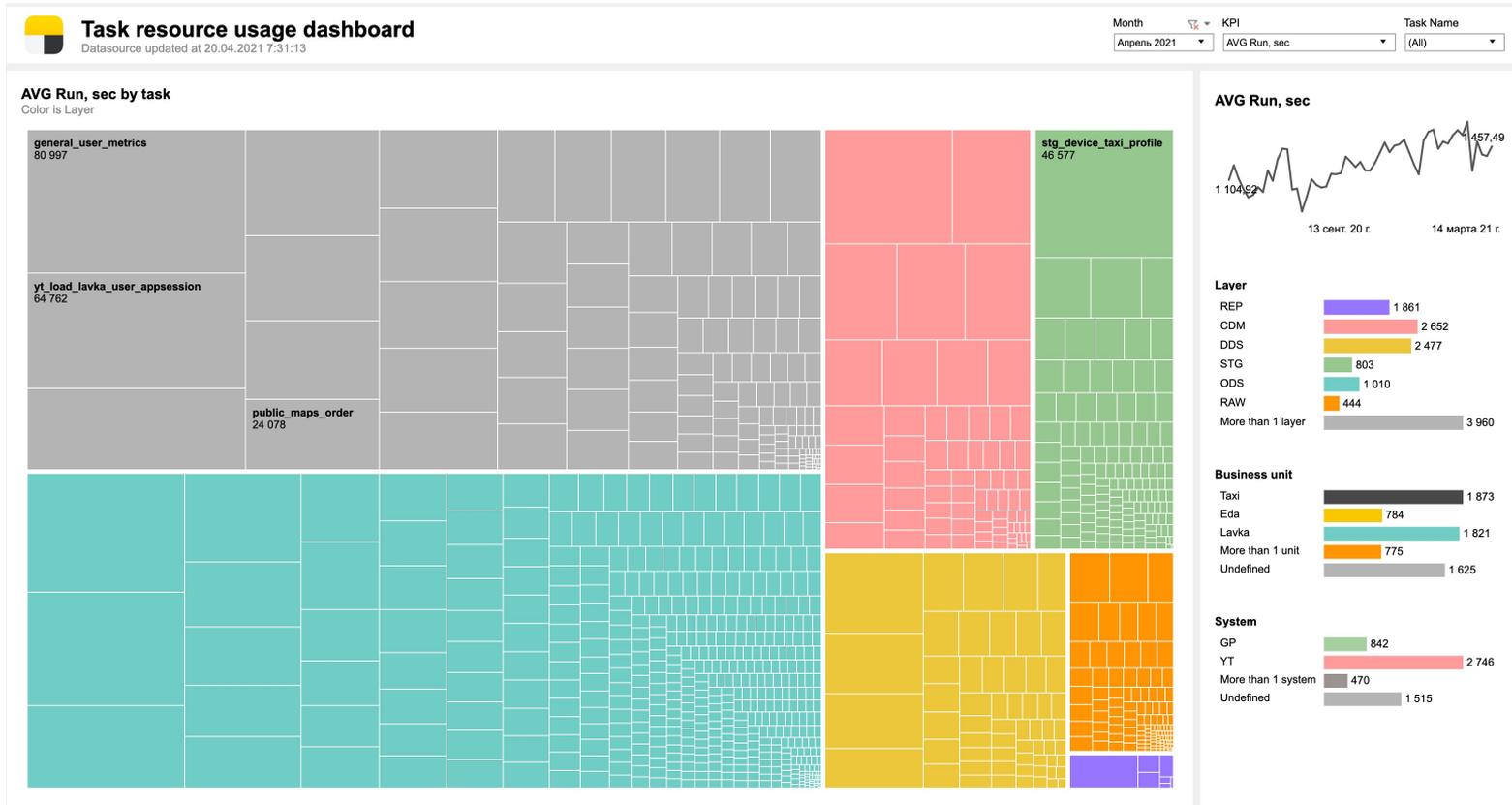
Есть ли более техническое  
применение?

III. Что получили?

# Размер объектов DWH



# Потребление ресурсов



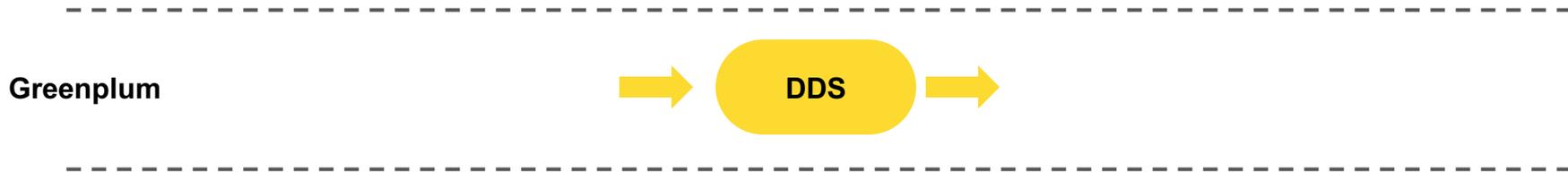
# Можно ли применить знания для оптимизации?

III. Что получили?

# Детальный слой

Детальный слой – ключевой для построения доменной модели

- › Хранить историю изменений сущностей
- › Отвечает за консолидацию данных между источниками
- › Устойчив к изменению в бизнесе
- › Модульный и масштабируемый



# Подходы к проектированию

сложность эксплуатации, простота внесения изменений

## Никакого

- › Денормализация
- › Можно использовать без подготовки
- › Неустойчиво к изменениям
- › Дублирование информации
- › Нет join

## Звезда и снежинка

- › Нормализация
- › Можно использовать с минимальной подготовкой
- › Неудобно перестраивать
- › Минимальное дублирование информации
- › Приемлемое количество join

## Data Vault

- › Строгая нормализация
- › Нельзя использовать без подготовки
- › Не надо перестраивать
- › Нет дублирования информации
- › Большое количество join

## Anchor modeling

- › Ультранормализация
- › Нельзя использовать без подготовки
- › Не надо перестраивать
- › Нет дублирования информации
- › Ультраколичество join

легкость эксплуатации, сложность внесения изменений

# Highly Normalized Hybrid Model

## Выбирать оптимальный формат хранения для каждого конкретного случая

- › Высокая нормализация
- › Параллельная загрузка из разных источников
- › Устойчив к изменению в бизнесе
- › Идемпотентный к повторной загрузке
- › Модульный и масштабируемый
- › Может эмулировать как Data Vault, так и Anchor Modeling



- › Атрибуты группируются в таблицы-**спутники** по принципам совместности: изменения и/или источника и/или использования
- › Есть специальные таблицы **Point-in-Time** и **Bridge**



- › На каждую сущность создается **anchor** – таблица с суррогатным ключом
- › **Связи** только через отдельные таблицы, никаких атрибутов – только хардкор

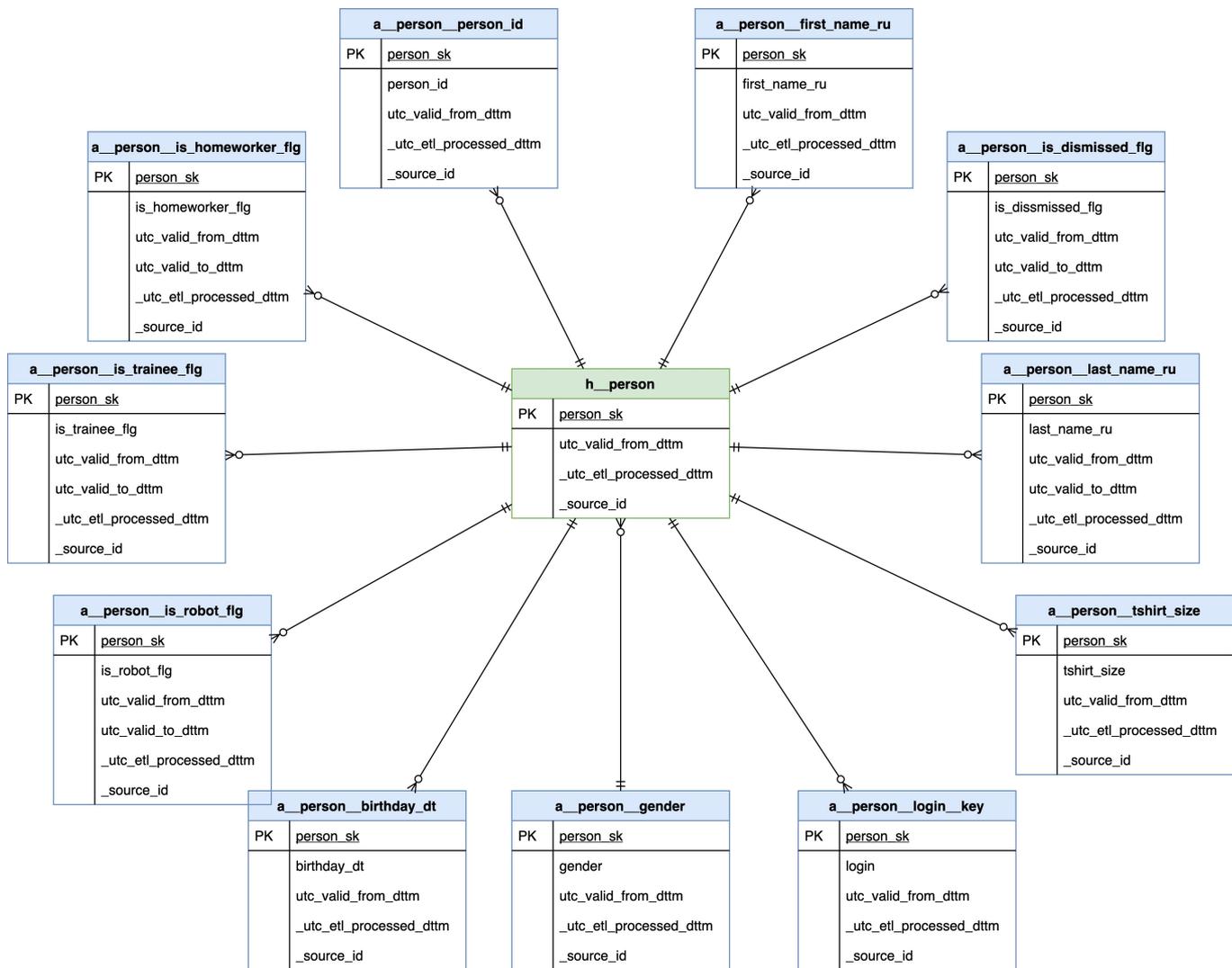
# Объявление сущности

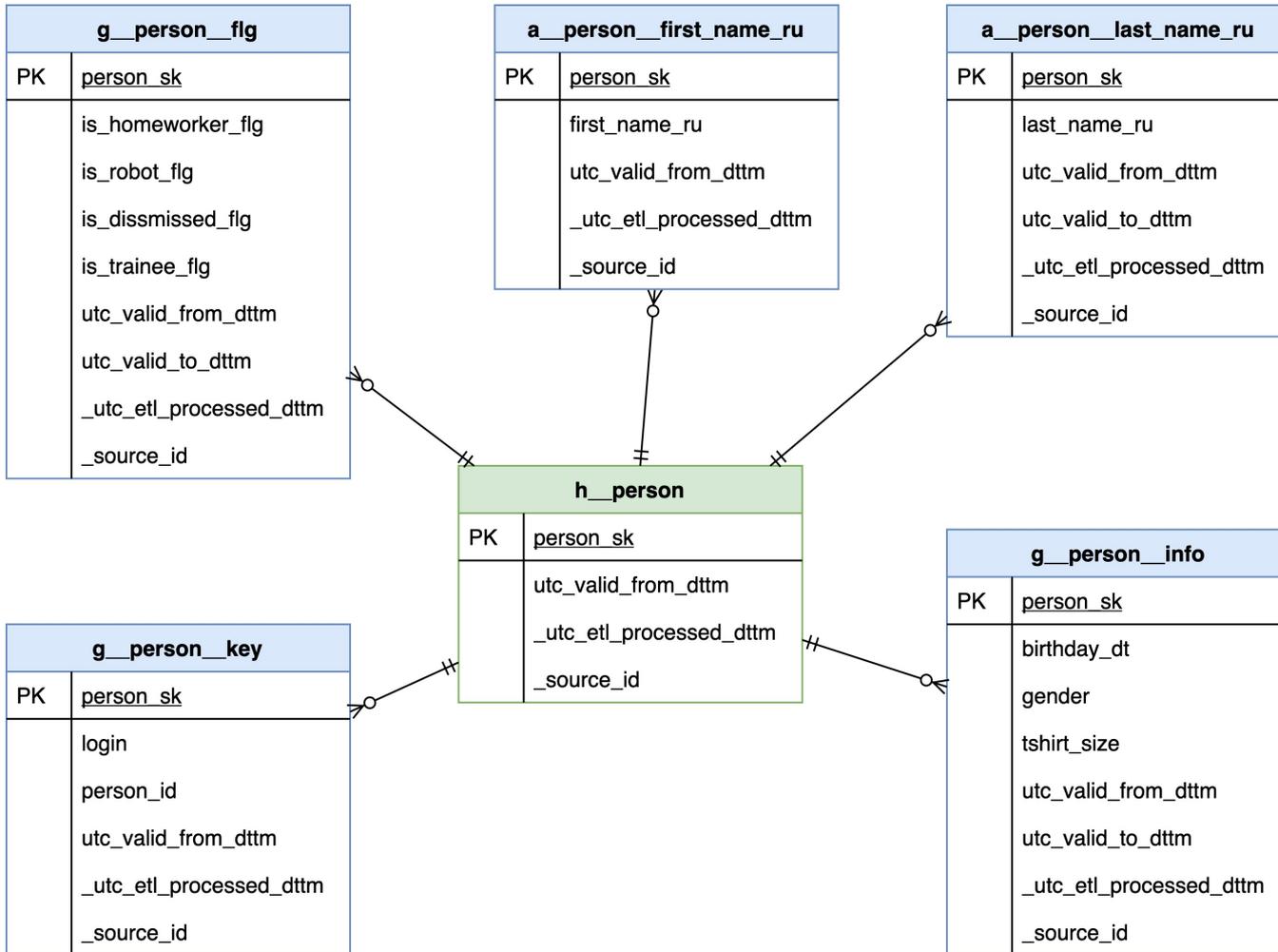
```
class Person(HnhmEntity):
    """Сотрудник со staff.yandex-team.ru"""

    __layout__ = DdsLayout(name='person', group='staff')

    person_id = Int(comment='ID в Стаффе', change_type=IGNORE)
    first_name_ru = String(comment='Имя сотрудника', change_type=UPDATE)
    last_name_ru = String(comment='Фамилия сотрудника', change_type=NEW)
    login = String(comment='Рабочий login', change_type=IGNORE)
    gender = String(comment='Пол', change_type=UPDATE)
    tshirt_size = String(comment='Размер футболки', change_type=UPDATE)
    birthday_dt = Date(comment='Дата рождения', change_type=UPDATE)
    is_dismissed_flg = Boolean(comment='Был уволен', change_type=NEW)
    is_homeworker_flg = Boolean(comment='Надомник', change_type=NEW)
    is_robot_flg = Boolean(comment='Робот', change_type=NEW)
    is_trainee_flg = Boolean(comment='Стажер', change_type=NEW)

    __keys__ = [login]
```





# Оптимизационная задача

**Вопрос:** как оптимально разбить данные по группам?

**Дано (и есть в metaDWH):**

- › Метаданные объектов
- › Маппинги полей и загрузки
- › Информация о количестве строк в объекте

**Ограничения**

- › Набор полей в метаданных объектов
- › Маппинги полей и загрузки (группа должна загружаться из одного источника)

**Оптимальность**

- › Будем минимизировать занимаемое место на диске

# Синтетический пример

business_dttm
key
Field 1
Field 2
Field 3

## Сущность состоит из 3х полей

- › Key – 16 байт
- › Field1 – 8байт
- › Field2 – 128байт
- › Field3 – 32байт

Итого 192 байт на строку

## Дополнительные поля

- › \_sk – 16 байт
- › \_dttm – 8 байт
- › \_source\_id – 2байт

Итого 42 байт на строку

# Синтетический пример

business_dttm
key
Field 1
Field 2
Field 3

Сущность состоит из 3х полей

- › Key – 16 байт
- › Field1 – 8байт
- › Field2 – 128байт
- › Field3 – 32байт

Итого 192 байт на строку

Дополнительные поля

- › \_sk – 16 байт
- › \_dttm – 8 байт
- › \_source\_id – 2байт

Итого 42 байт на строку

F1F2F3

sk
Field 1
Field 2
Field 3
valid_from_dttm
valid_to_dttm
_source_id
_etl_dttm

F1 | F2F3

sk	sk
Field 1	Field 2
valid_from_dttm	Field 3
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

F2 | F1F3

sk	sk
Field 2	Field 1
valid_from_dttm	Field 3
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

F3 | F1F2

sk	sk
Field 3	Field 2
valid_from_dttm	Field 1
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

F1 | F1 | F3

sk	sk	sk
Field 1	Field 1	Field 1
valid_from_dttm	valid_from_dttm	valid_from_dttm
valid_to_dttm	valid_to_dttm	valid_to_dttm
_source_id	_source_id	_source_id
_etl_dttm	_etl_dttm	_etl_dttm

# Синтетический пример

business_dttm
key
Field 1
Field 2
Field 3

Сущность состоит из 3х полей – 1 000 000 строчек

- › Key – 16 байт – 500 000 сущностей
- › Field1 – 8байт – 2 изменения на ключ
- › Field2 – 128байт – 2 изменения на ключ
- › Field3 – 32байт – 2 изменения на ключ

Итого 183 Мбайт

Дополнительные поля

- › \_sk – 16 байт
- › \_dttm – 8 байт
- › \_source\_id – 2байт

Итого 42 байт

F1F2F3

sk
Field 1
Field 2
Field 3
valid_from_dttm
valid_to_dttm
_source_id
_etl_dttm

200 Мбайт

F1 | F2F3

sk	sk
Field 1	Field 2
valid_from_dttm	Field 3
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

240 Мбайт

F2 | F1F3

sk	sk
Field 2	Field 1
valid_from_dttm	Field 3
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

240 Мбайт

F3 | F1F2

sk	sk
Field 3	Field 2
valid_from_dttm	Field 1
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

240 Мбайт

F1 | F1 | F3

sk	sk	sk
Field 1	Field 1	Field 1
valid_from_dttm	valid_from_dttm	valid_from_dttm
valid_to_dttm	valid_to_dttm	valid_to_dttm
_source_id	_source_id	_source_id
_etl_dttm	_etl_dttm	_etl_dttm

280 Мбайт

# Синтетический пример

business_dttm
key
Field 1
Field 2
Field 3

Сущность состоит из 3х полей – 1 000 000 строчек

- › Key – 16 байт – 31 250 сущностей
- › Field1 – 8байт – 2 изменения на ключ
- › Field2 – 128байт – **32 изменения** на ключ
- › Field3 – 32байт – 2 изменения на ключ

Итого 183 Мбайт

Дополнительные поля

- › \_sk – 16 байт
  - › \_dttm – 8 байт
  - › \_source\_id – 2байт
- Итого 42 байт

F1F2F3

sk
Field 1
Field 2
Field 3
valid_from_dttm
valid_to_dttm
_source_id
_etl_dttm

200 Мбайт

F1 | F2F3

sk	sk
Field 1	Field 2
valid_from_dttm	Field 3
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

174 Мбайт

F2 | F1F3

sk	sk
Field 2	Field 1
valid_from_dttm	Field 3
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

167 Мбайт

F3 | F1F2

sk	sk
Field 3	Field 2
valid_from_dttm	Field 1
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

195 Мбайт

F1 | F1 | F3

sk	sk	sk
Field 1	Field 1	Field 1
valid_from_dttm	valid_from_dttm	valid_from_dttm
valid_to_dttm	valid_to_dttm	valid_to_dttm
_source_id	_source_id	_source_id
_etl_dttm	_etl_dttm	_etl_dttm

169 Мбайт

# Синтетический пример

business_dttm
key
Field 1
Field 2
Field 3

Сущность состоит из 3х полей – 1 000 000 строчек

- › Key – 16 байт – 15 625 сущностей
- › Field1 – 8байт – 2 изменения на ключ
- › Field2 – 128байт – **32 изменения** на ключ
- › Field3 – 32байт – **64 изменения** на ключ

Итого 183 Мбайт

Дополнительные поля

- › \_sk – 16 байт
  - › \_dttm – 8 байт
  - › \_source\_id – 2байт
- Итого 42 байт

**F1F2F3**

sk
Field 1
Field 2
Field 3
valid_from_dttm
valid_to_dttm
_source_id
_etl_dttm

200 Мбайт

**F1 | F2F3**

sk	sk
Field 1	Field 2
valid_from_dttm	Field 3
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

155 Мбайт

**F2 | F1F3**

sk	sk
Field 2	Field 1
valid_from_dttm	Field 3
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

159 Мбайт

**F3 | F1F2**

sk	sk
Field 3	Field 2
valid_from_dttm	Field 1
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

194 Мбайт

**F1 | F1 | F3**

sk	sk	sk
Field 1	Field 1	Field 1
valid_from_dttm	valid_from_dttm	valid_from_dttm
valid_to_dttm	valid_to_dttm	valid_to_dttm
_source_id	_source_id	_source_id
_etl_dttm	_etl_dttm	_etl_dttm

**150 Мбайт**

# Наше решение

Вводим атомарные операции, меняющие схему,  
но не меняющие логику

- › Объединение групп/атрибутов
- › Соединение групп/атрибутов

В hNhM с точки зрения использования сущности логической модели все варианты физического хранения ниже одинаковы.

**F1F2F3**

sk
Field 1
Field 2
Field 3
valid_from_dttm
valid_to_dttm
_source_id
_etl_dttm

**F1 | F2F3**

sk	sk
Field 1	Field 2
valid_from_dttm	Field 3
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

**F2 | F1F3**

sk	sk
Field 2	Field 1
valid_from_dttm	Field 3
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

**F3 | F1F2**

sk	sk
Field 3	Field 2
valid_from_dttm	Field 1
valid_to_dttm	valid_from_dttm
_source_id	valid_to_dttm
_etl_dttm	_source_id
	_etl_dttm

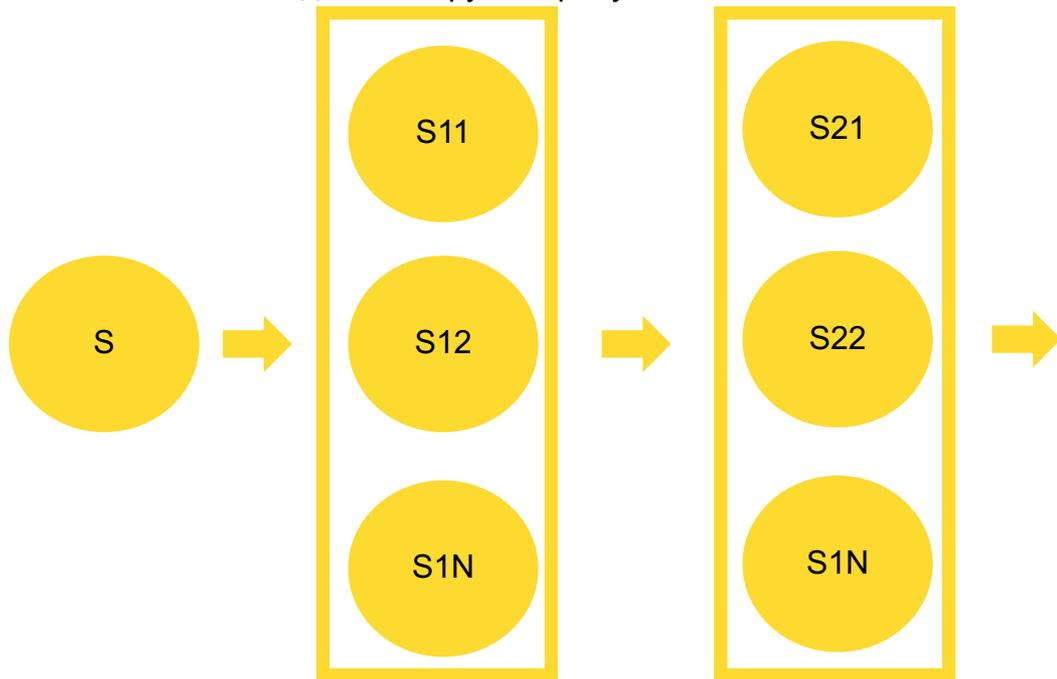
**F1 | F1 | F3**

sk	sk	sk
Field 1	Field 1	Field 1
valid_from_dttm	valid_from_dttm	valid_from_dttm
valid_to_dttm	valid_to_dttm	valid_to_dttm
_source_id	_source_id	_source_id
_etl_dttm	_etl_dttm	_etl_dttm

# Наше решение

Вводим атомарные операции, меняющие схему,  
но не меняющие логику

- › Объединение групп/атрибутов
- › Соединение групп/атрибутов



## Генетический алгоритм

- › Из текущего состояния мутациями (=атомарными операциями) создаем стартовую популяцию
- › Производим скрещивания и новые мутации
- › Каждое состояние оцениваем на оптимальность (в нашем случае по месту)
- › При подозрениях на сходимость останавливаемся

## Результат

- › Получаем итоговое состояние, которое лучше текущего
- › Сравниваем метаданные между состояниями и генерируем скрипт миграции
- › Миграция – отдельный вопрос

**Резюме**

## **I. Проблема:**

развитием крупного  
DWH сложно управлять

## **II. Решение:**

покрыть работу DWH  
метриками

## **III. Идея:**

использовать данные  
систем DWH в самом  
DWH

(«DWH для DWH»)

## **IV. Результат:**

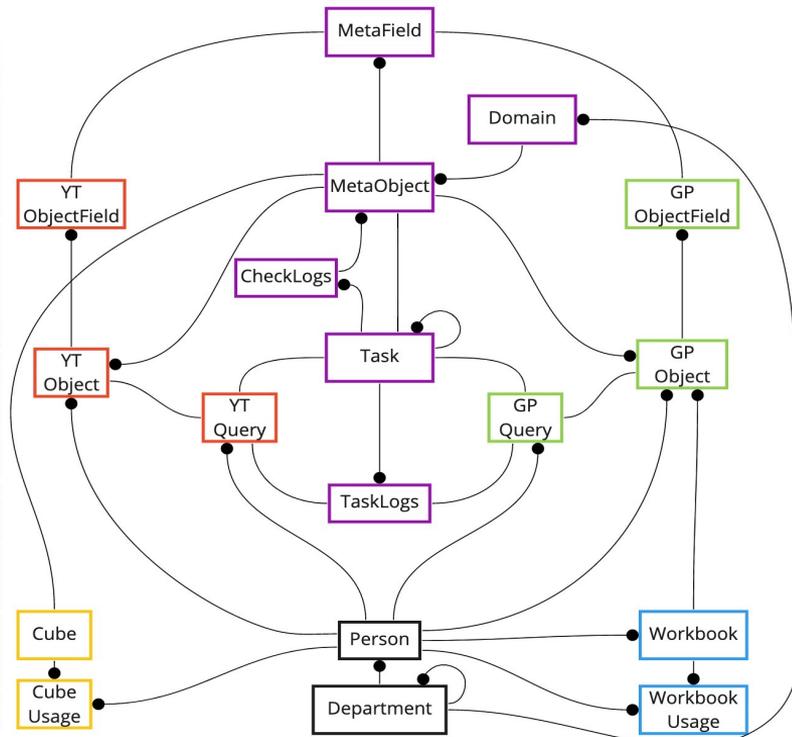
аналитика по работе и  
развитию самого DWH

# MetaDWH

## Source Domain



## Core Domain



## Business Domain

### Техническая информация



### Использование объектов



### Витрины с метаданными



# Затраты

## Стоимость реализации

- › Никаких дополнительных внедрений, исключительно существующие системы
- › Необходимы разноплановые специалисты (infra, de, dp, bi), которые уже есть в DWH
- › Абстрактно в вакууме 2FTE на Q
- › Фактически порядка 10 специалистов с  $\pm 20\%$  загрузкой

# Результат

## Стоимость реализации

- › Никаких дополнительных внедрений, исключительно существующие системы
- › Необходимы разноплановые специалисты (infra, de, dp, bi), которые уже есть в DWH
- › Абстрактно в вакууме 2FTE на Q
- › Фактически порядка 10 специалистов с  $\pm 20\%$  загрузкой

## Аналитика по ключевым аспектам

- › Целевые метрики и принятие стратегических решений
- › управление приоритизацией через KPI команд
- › Ad-hoc-запросы по использованию объектов хранилища
- › Поиск технически узких мест и оптимизация
- › Интеллектуальная нотификация пользователей

# Результат

## Стоимость реализации

- › Никаких дополнительных внедрений, исключительно существующие системы
- › Необходимы разноплановые специалисты (infra, de, dp, bi), которые уже есть в DWH
- › Абстрактно в вакууме 2FTE на Q
- › Фактически порядка 10 специалистов с  $\pm 20\%$  загрузкой



## Аналитика по ключевым аспектам

- › Целевые метрики и принятие стратегических решений
- › управление приоритизацией через KPI команд
- › Ad-hoc-запросы по использованию объектов хранилища
- › Поиск технически узких мест и оптимизация
- › Интеллектуальная нотификация пользователей

Возможно реализовать на любом отлаженном DWH

# Резюме

**DWH может быть источником данных для DWH**

**Создать MetaDWH – не слишком трудоемкая задача (при наличии рабочего DWH)**

**Обработка только логов запросов позволяет получить дашборды для анализа поведения пользователей**

**Более сложная систематизация (домены, слои, команды) позволяет ставить продуктовые метрики командам**

**Пример технической реализации: поиск узких мест среди объектов/тасок и модификация схемы в детальном слое**



# Спасибо

**Евгений Ермаков**

Руководитель Data Office



[jkermakov@yandex-team.ru](mailto:jkermakov@yandex-team.ru)



[@iJKos](https://www.t.me/iJKos)

[iJKos.com](http://iJKos.com)