



Платформа управления данными вокруг YTsaurus

Владимир Верстов
Наталья Савенкова

Яндекс Go

Yandex Go

Services for Ridetech, Foodtech,
E-commerce and Delivery

40 mln

MAU of Yandex Go



47 mln

SKU on
MarketPlace



570 k

deliveries
per day



Содержание

01 Платформа

02 Архитектура

03 YTsaurus

04 Выводы



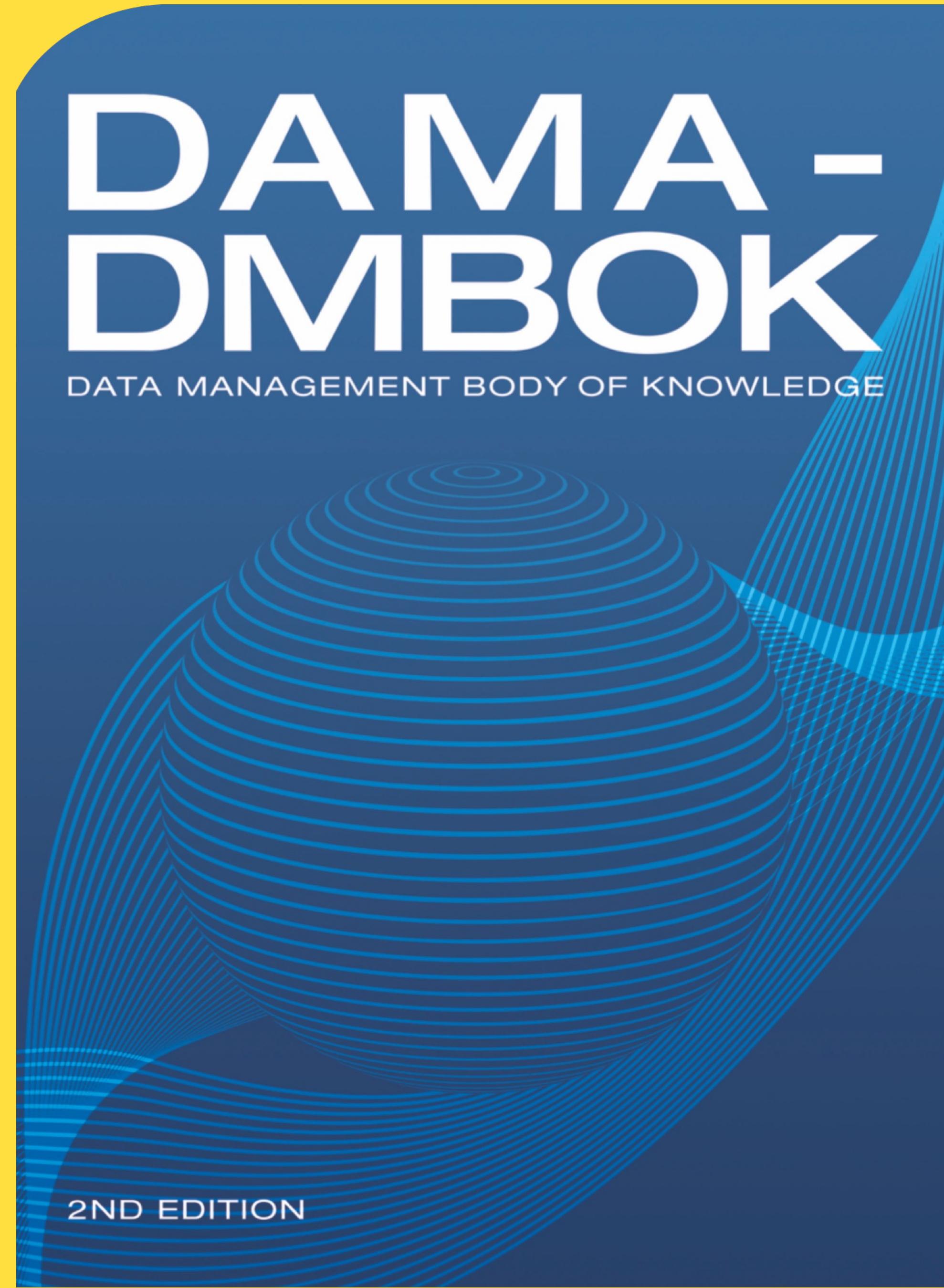
Платформа



Data Management

Data Management is the development, execution, and supervision of plans, policies, programs, and practices that **deliver, control, protect, and enhance the value of data assets throughout their lifecycles.**

The primary driver for data management is to enable organizations **to get value from their data assets.**



Data Management Goals

- Understanding and supporting the information needs of the enterprise and its stakeholders, including customers, employees, and business partners
- Capturing, storing, protecting, and ensuring the integrity of data assets
- Ensuring the quality of data and information
- Ensuring the privacy and confidentiality of stakeholder data
- Preventing unauthorized or inappropriate access, manipulation, or use of data and information
- Ensuring data can be used effectively to add value to the enterprise

Data Lifecycle

Creation and usage are the most critical points:

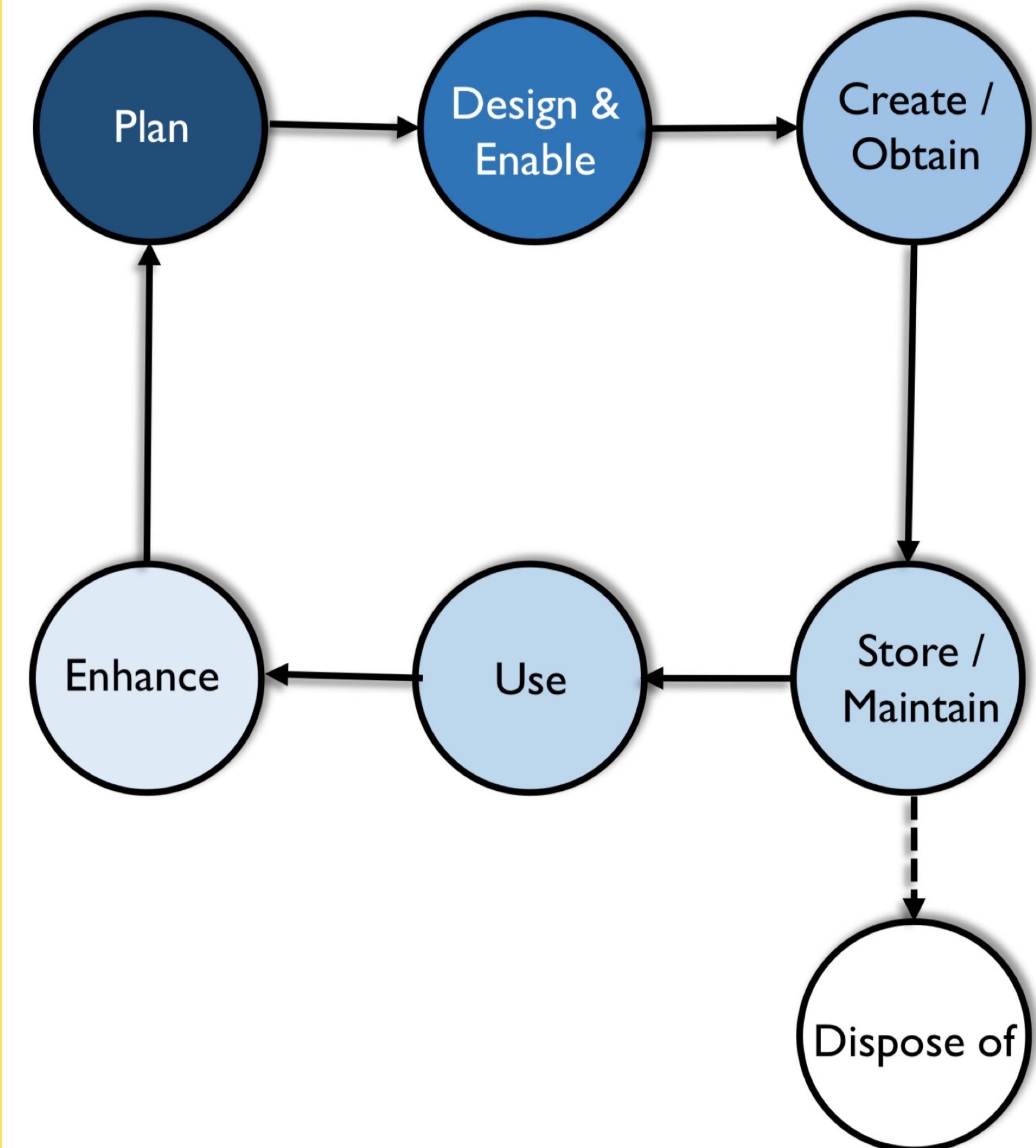
- It costs money to produce data
- Data is valuable only when it is consumed or applied

Management thought lifecycle:

- Data Quality
- Metadata
- Data Security

Efforts should focus on the most critical data:

- Organisations produce a lot of data
- Large portions of data is never actually used
- Manage every piece of data is not possible



Data Lifecycle

Creation and usage are the most critical points:

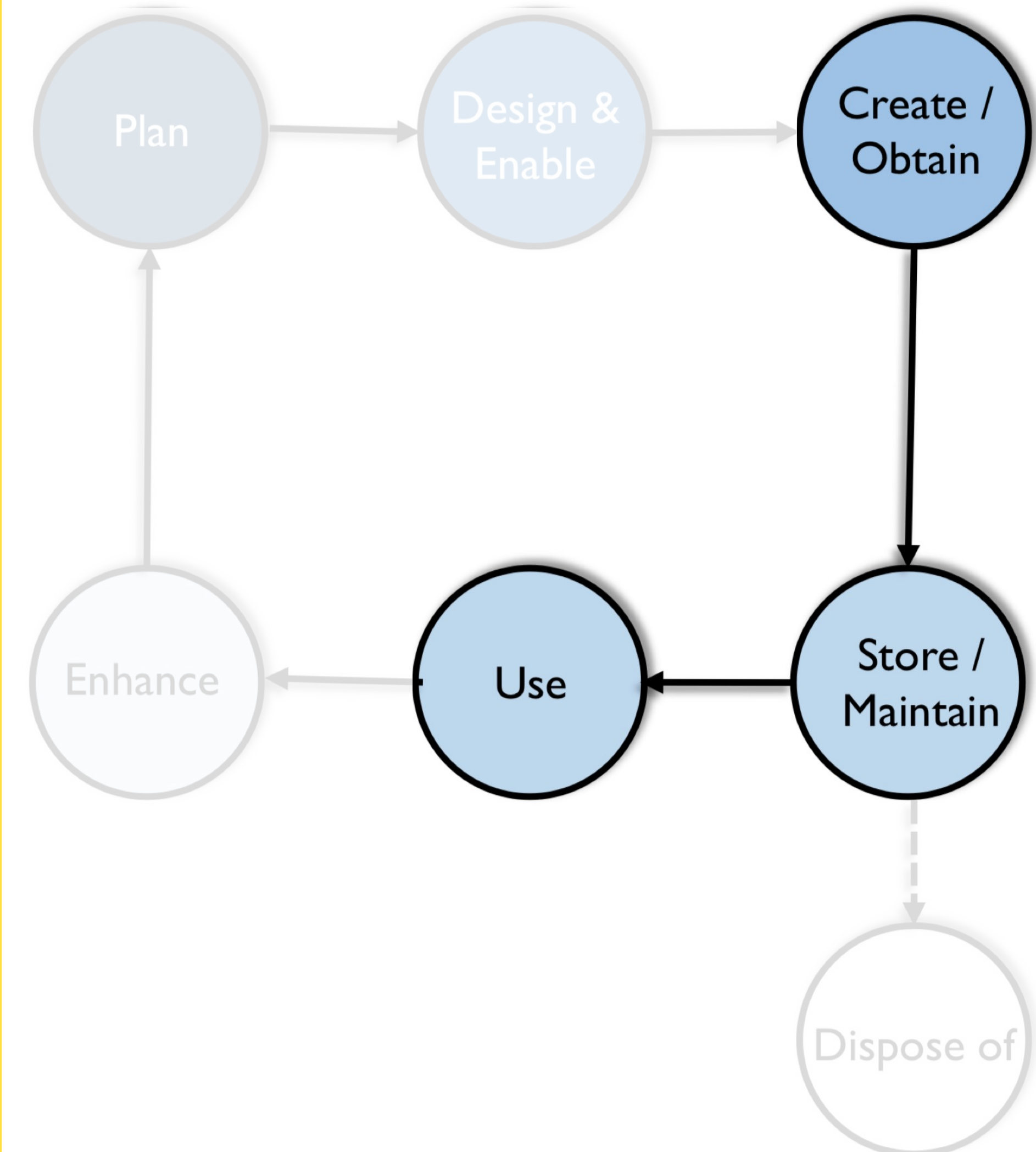
- It costs money to produce data
- Data is valuable only when it is consumed or applied

Management thought lifecycle:

- Data Quality
- Metadata
- Data Security

Efforts should focus on the most critical data:

- Organisations produce a lot of data
- Large portions of data is never actually used
- Manage every piece of data is not possible



Архитектура



2016



Нужны данные и отчеты





2016

2017

YTsaurus

Data Lake
Staging & Data Marts
Sandboxes



YTsaurus

Data Lake
Staging & Data Marts
Sandboxes

BI

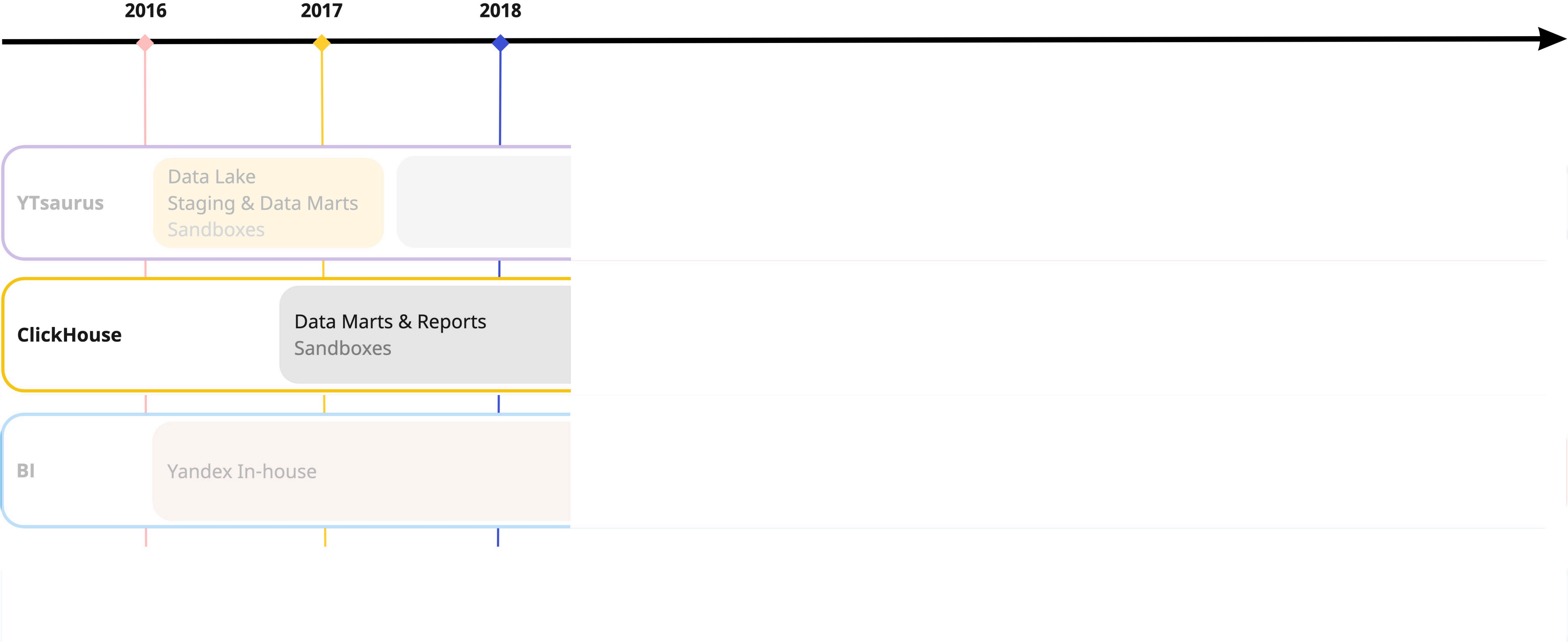
Yandex In-house

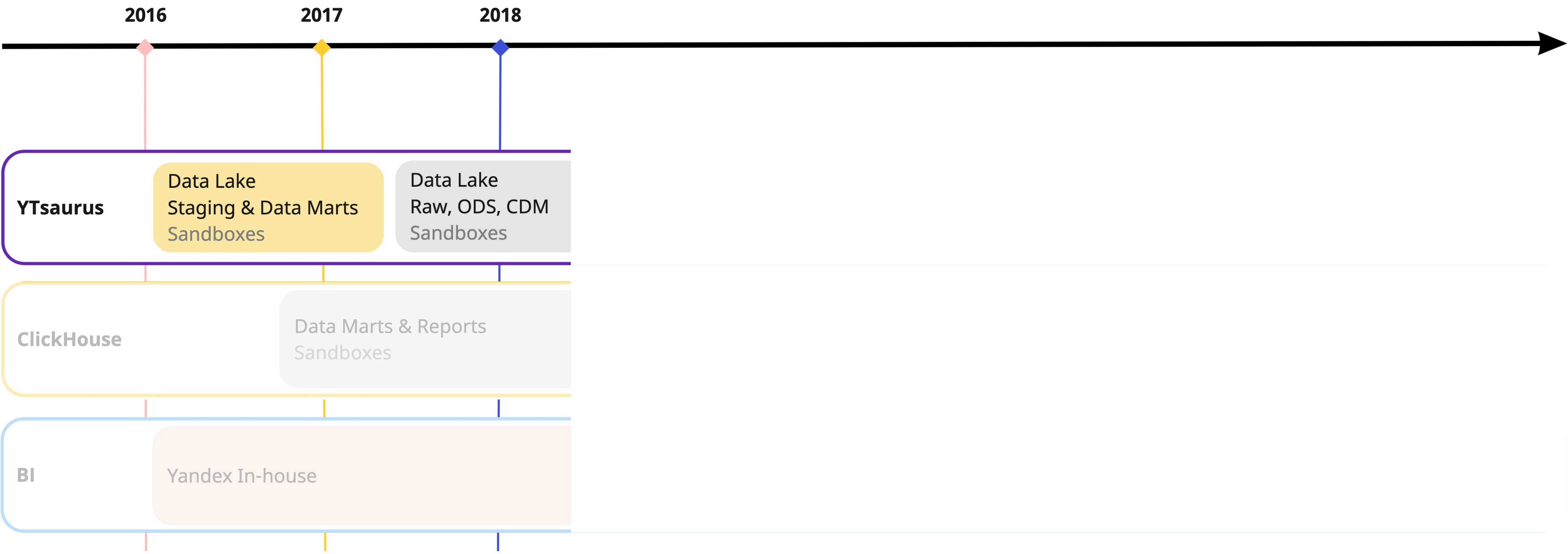
**Бизнес быстро
растет**

**Бэкенд активно
развивается**

**Нужно бежать
быстрее**







2016

2017

2018

YTsaurus

Data Lake
Staging & Data Marts
Sandboxes

Data Lake
Raw, ODS, CDM
Sandboxes

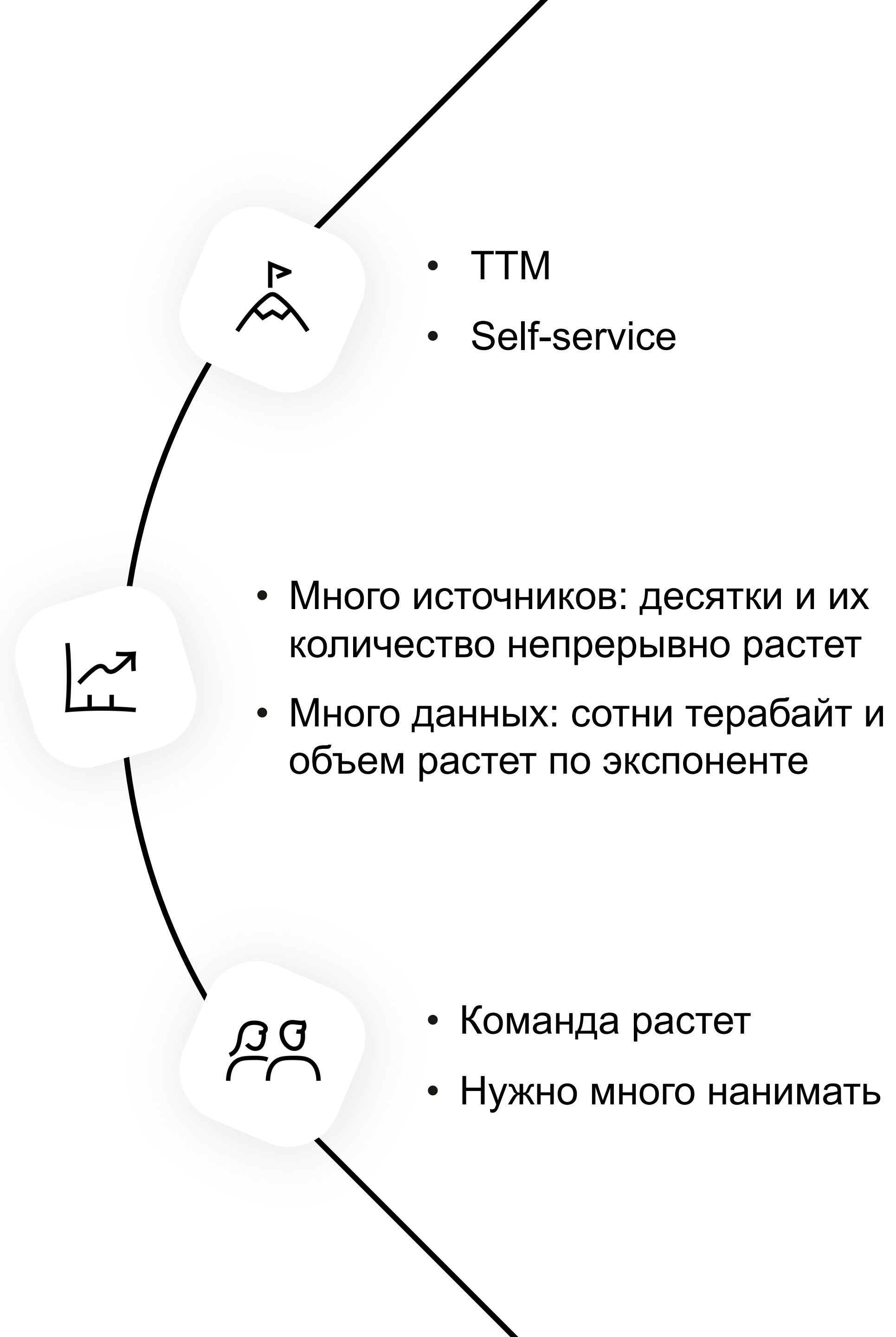
ClickHouse

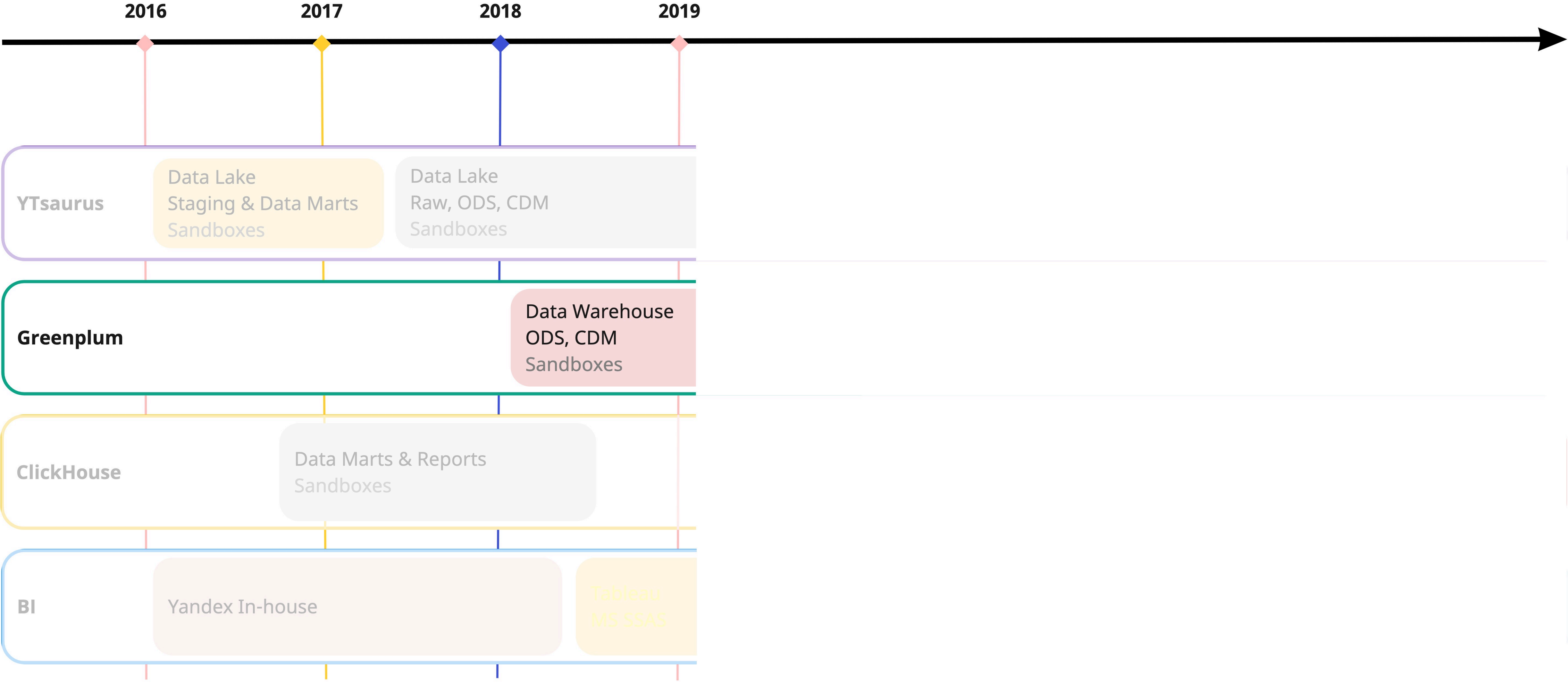
Data Marts & Reports
Sandboxes

BI

Yandex In-house

Компания стала большой





2016

2017

2018

2019

YTsaurus

Data Lake
Staging & Data Marts
Sandboxes

Data Lake
Raw, ODS, CDM
Sandboxes

Greenplum

Data Warehouse
ODS, CDM
Sandboxes

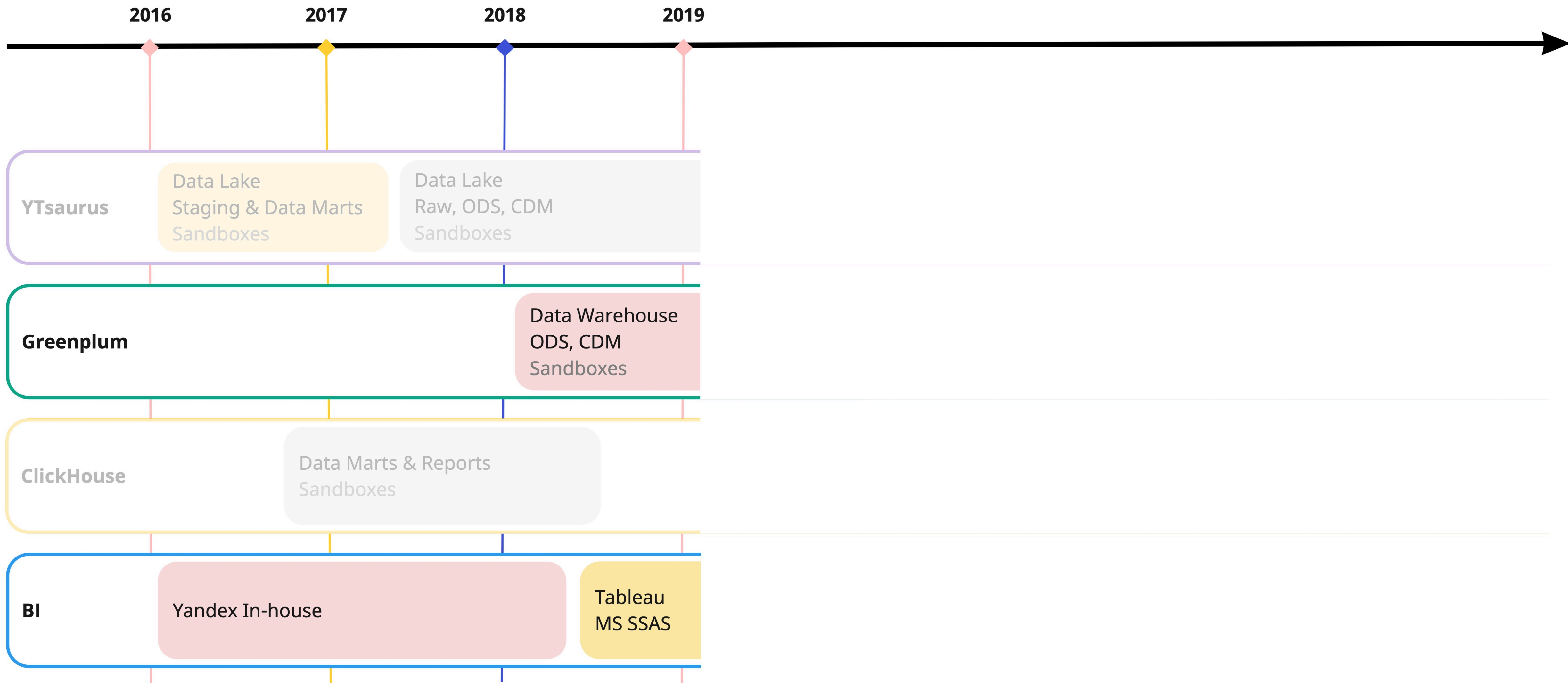
ClickHouse

Data Marts & Reports
Sandboxes

BI

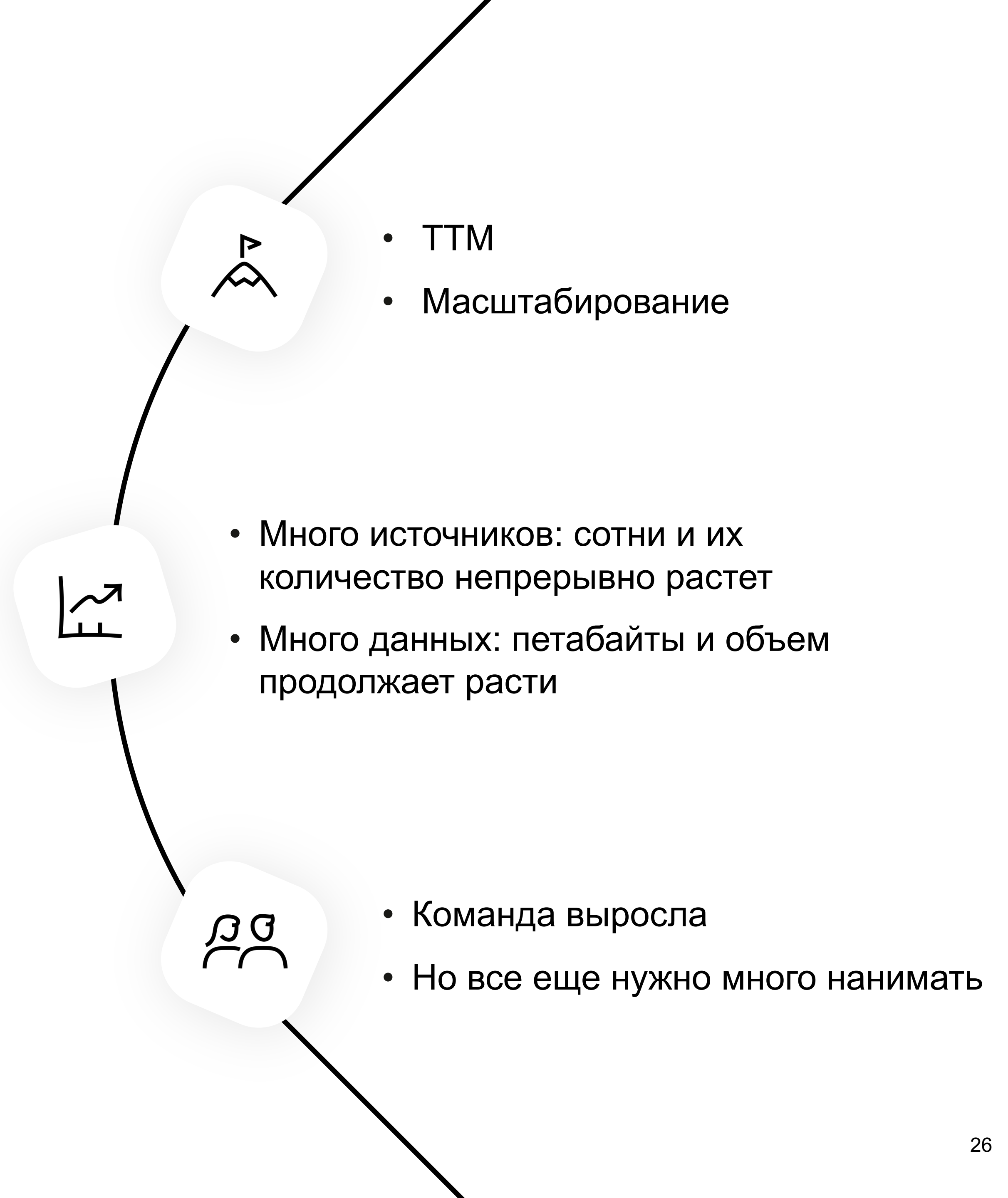
Yandex In-house

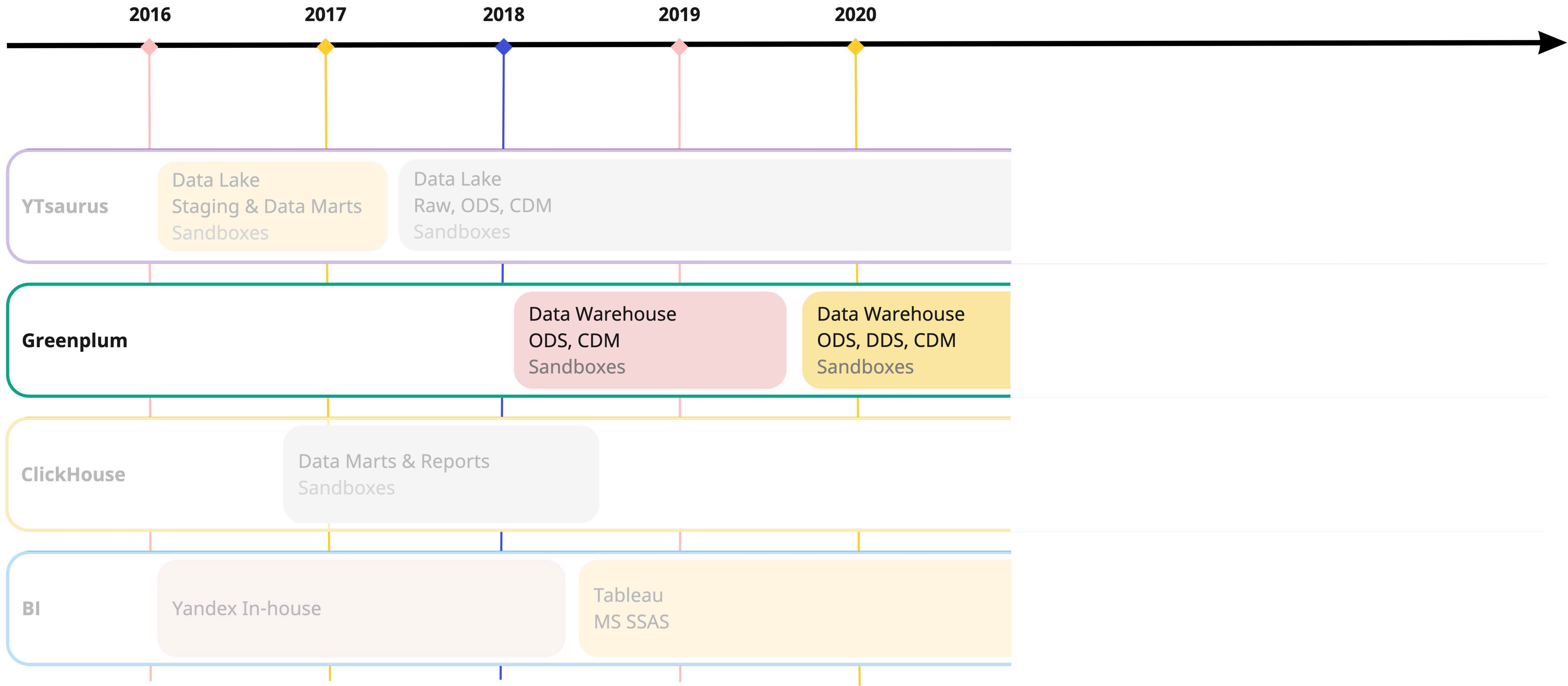
Tableau
MS SSAS



**Развиваются
новые бизнесы**

**Растет сложность
предметной
области**





2016

2017

2018

2019

2020

YTsaurus

Data Lake
Staging & Data Marts
Sandboxes

Data Lake
Raw, ODS, CDM
Sandboxes

Greenplum

Data Warehouse
ODS, CDM
Sandboxes

Data Warehouse
ODS, DDS, CDM
Sandboxes

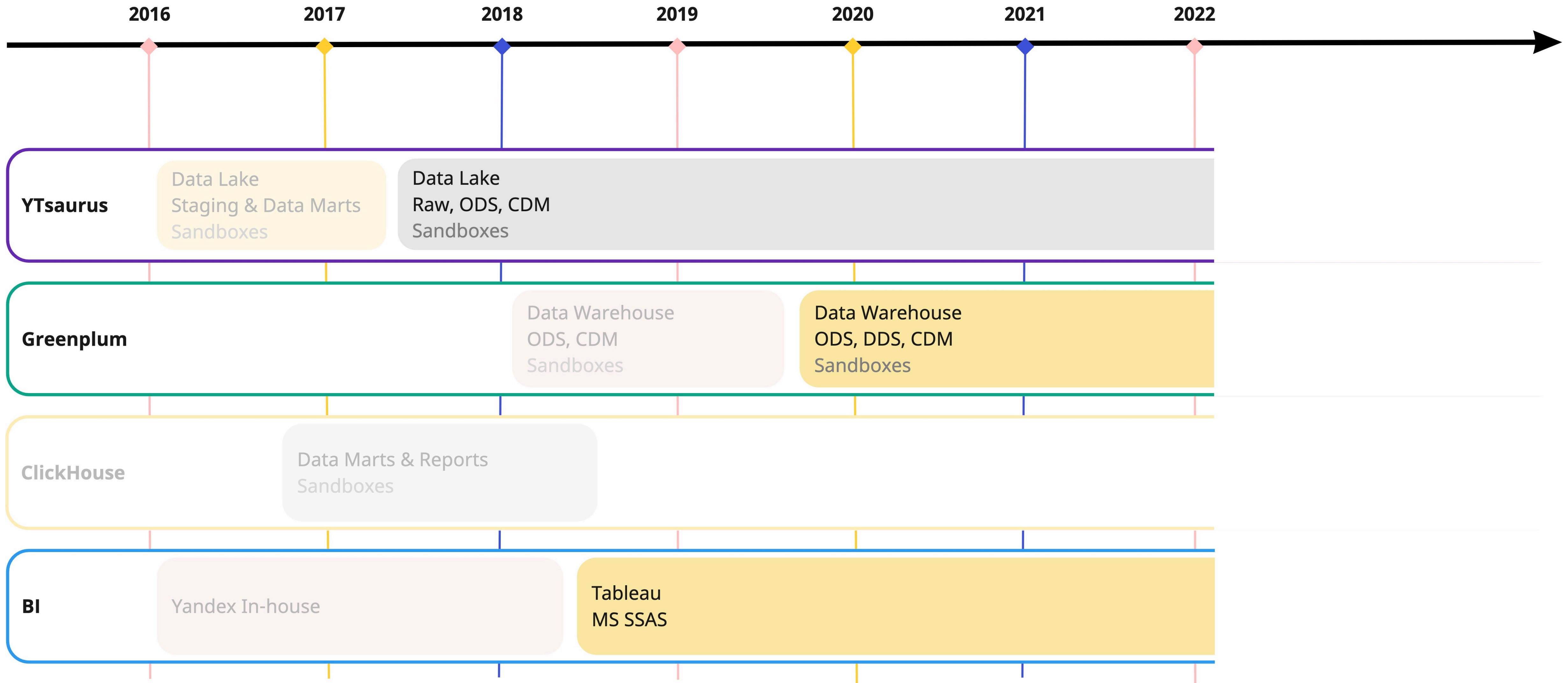
ClickHouse

Data Marts & Reports
Sandboxes

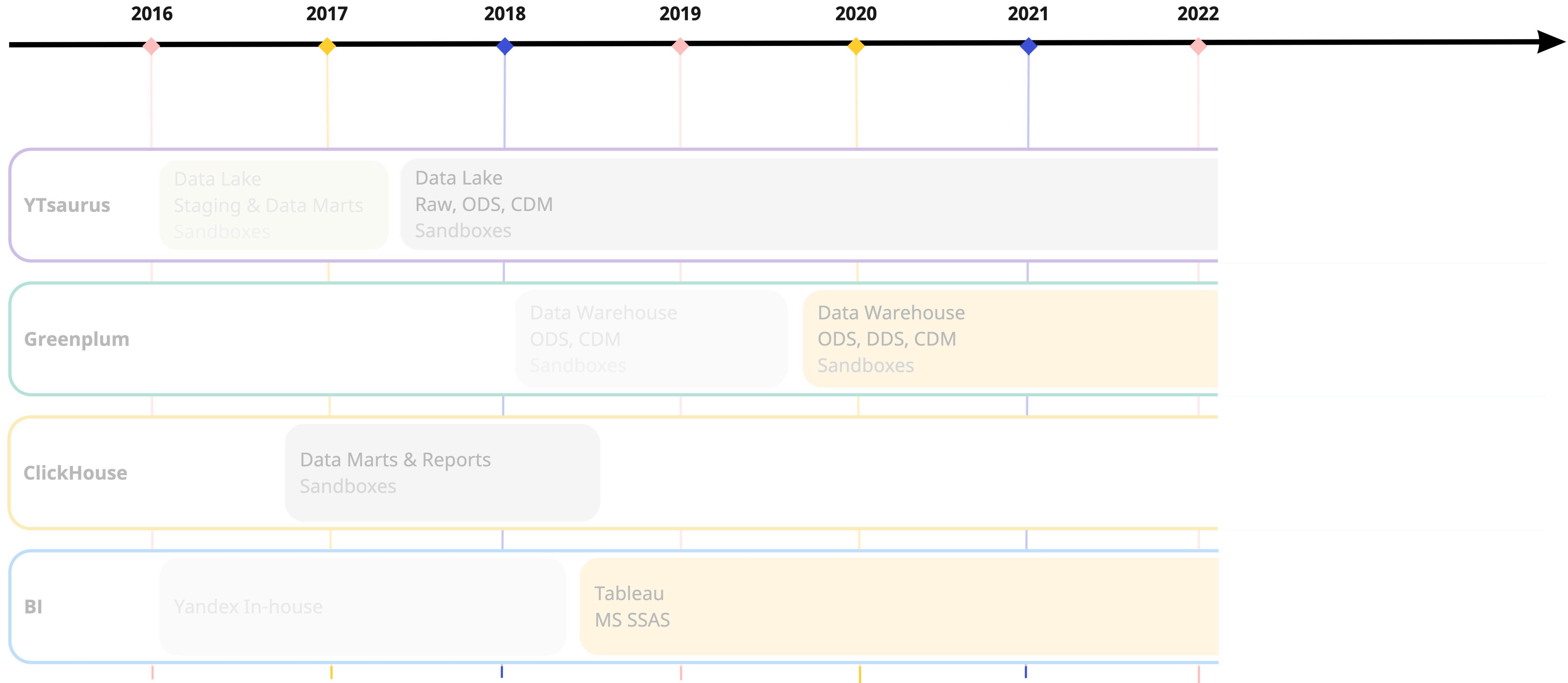
BI

Yandex In-house

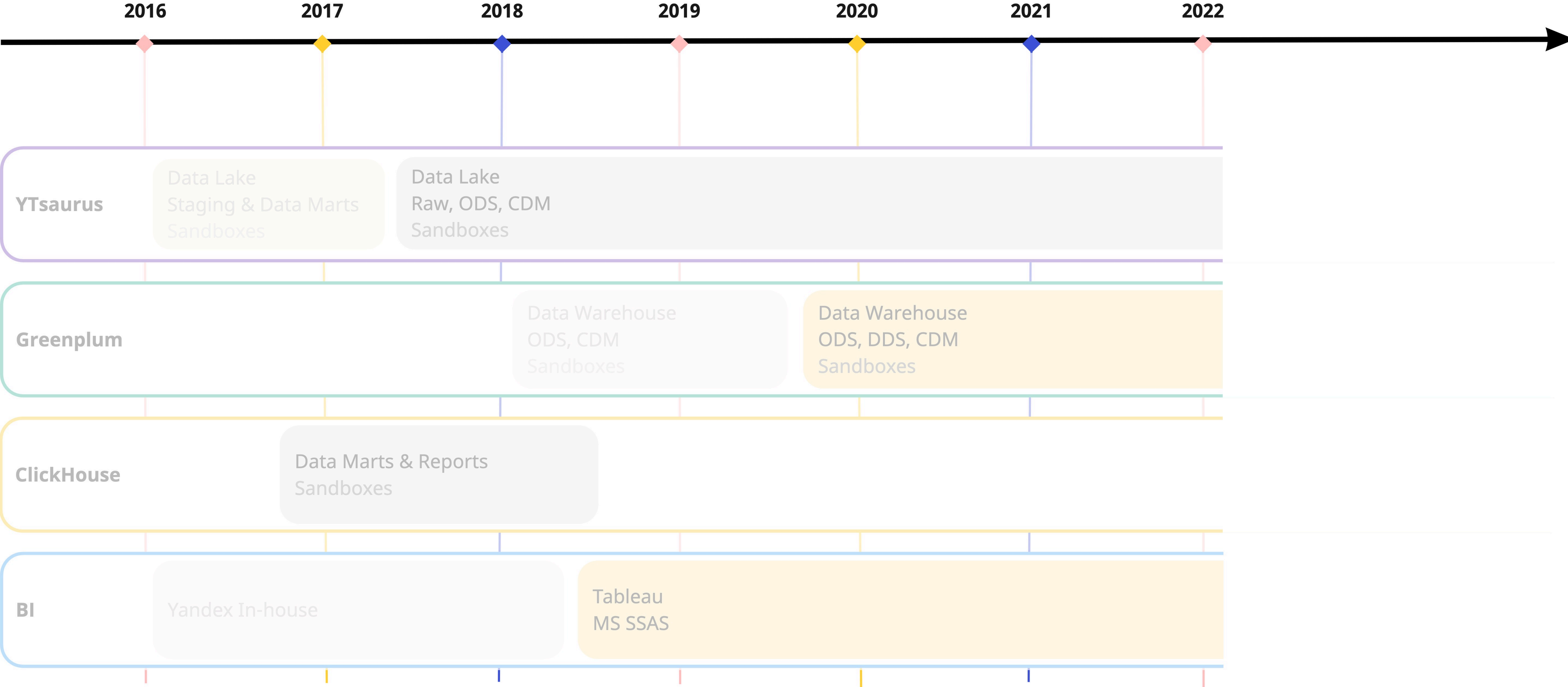
Tableau
MS SSAS



Apache Spark



Свой фреймворк



X PB на 2020-01-01

10X PB к 2022-12-31

2016

2017

2018

2019

2020

2021

2022



YTsaurus

Data Lake
Staging & Data Marts
Sandboxes

Data Lake
Raw, ODS, CDM
Sandboxes

Greenplum

Data Warehouse
ODS, CDM
Sandboxes

Data Warehouse
ODS, DDS, CDM
Sandboxes

ClickHouse

Data Marts & Reports
Sandboxes

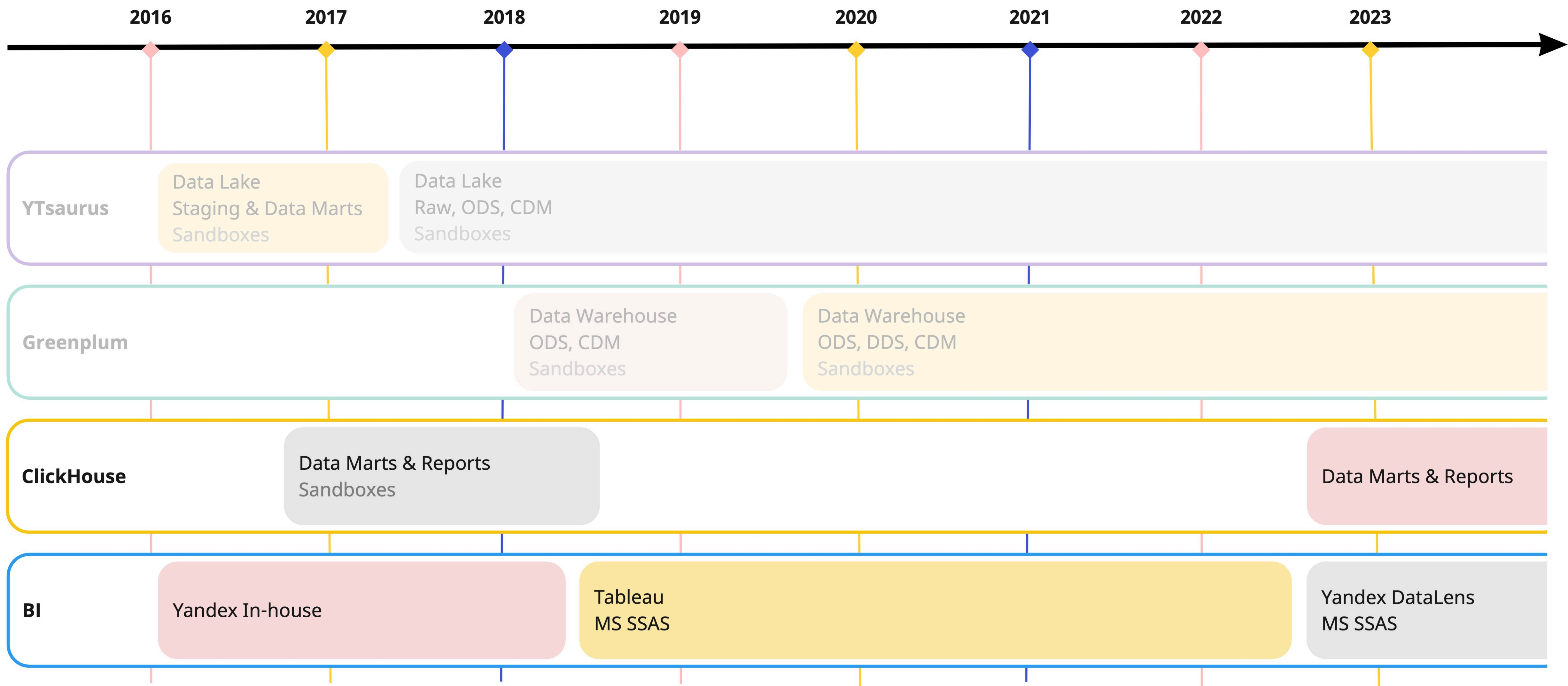
BI

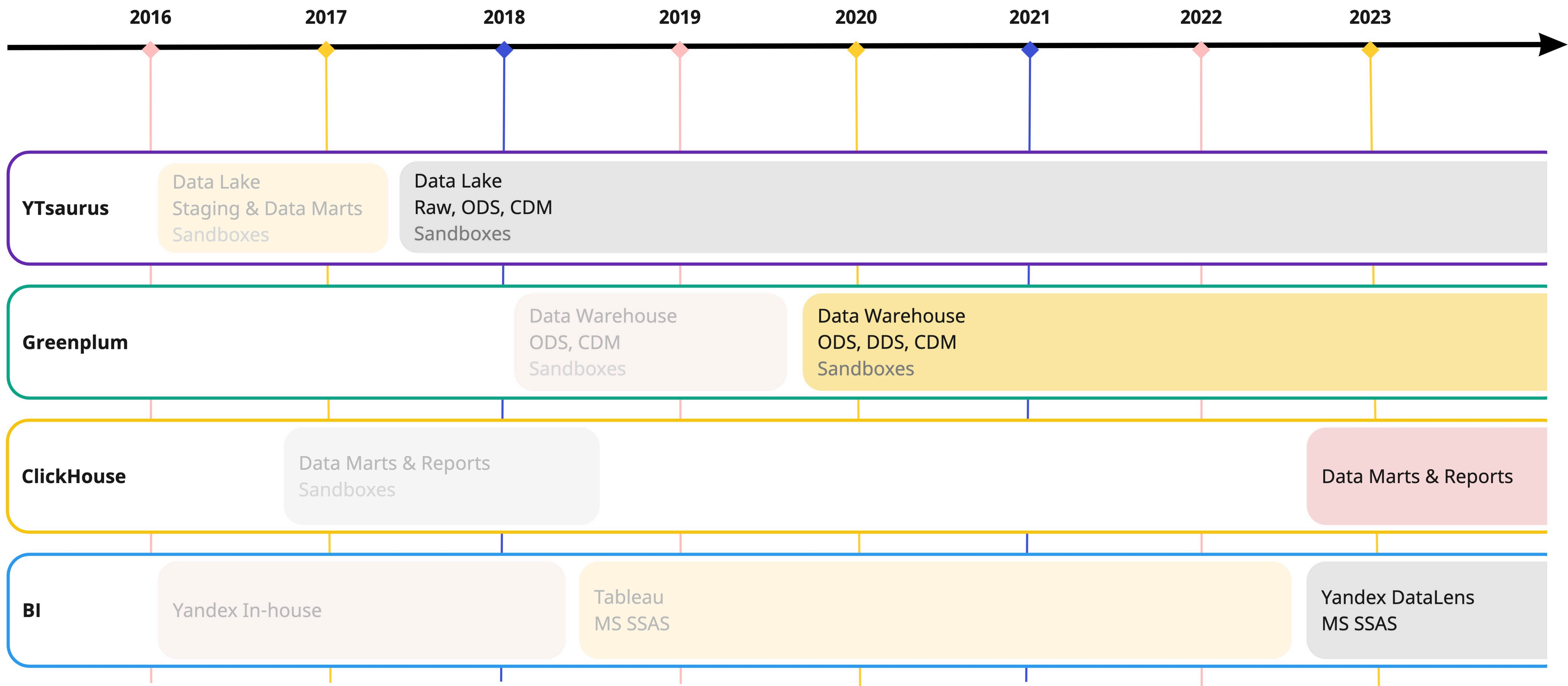
Yandex In-house

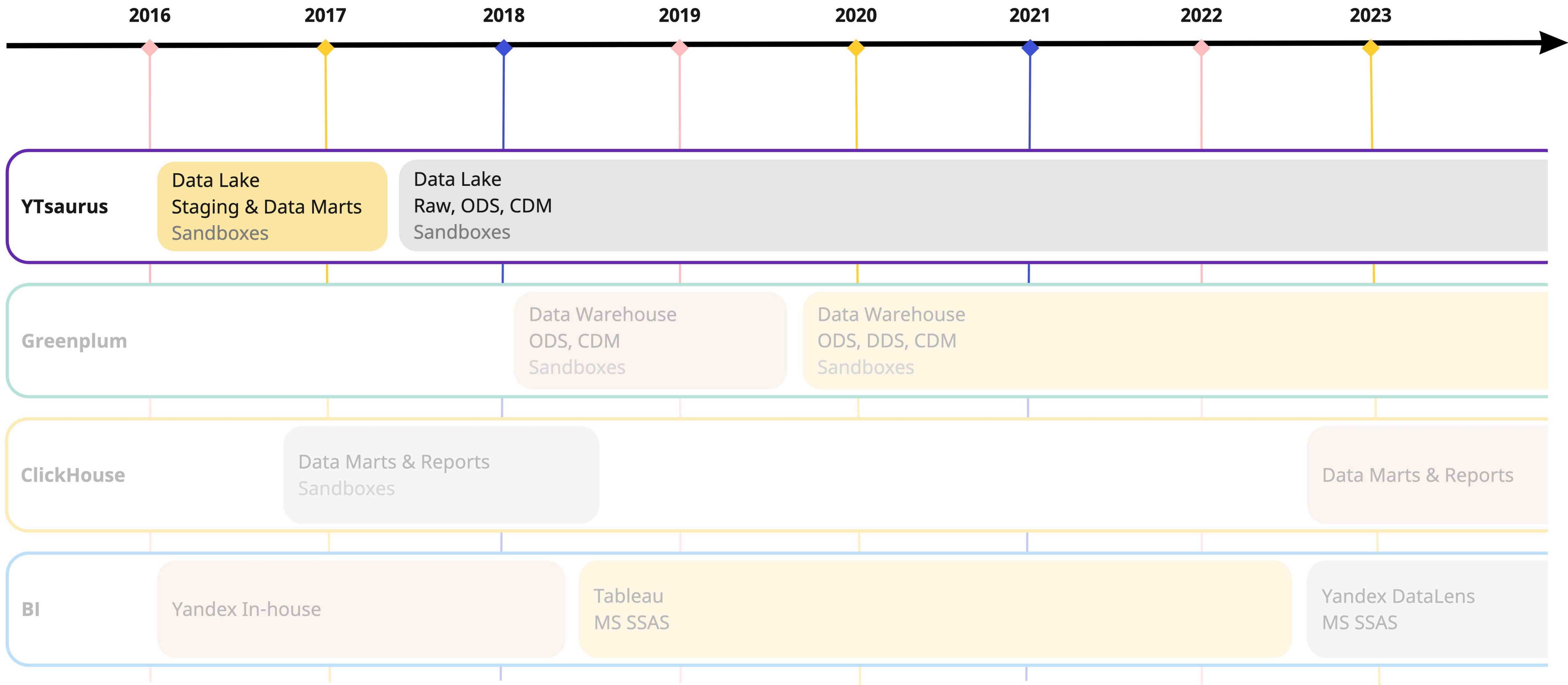
Tableau
MS SSAS

Внешние факторы









Можно проще

YTsaurus является необходимым и достаточным компонентом нашей платформы для множества сценариев.

Достаточно дополнить YTsaurus практически любым BI инструментом и это сразу закрывает огромное количество пользовательских сценариев в DWH, BI, аналитике и ML.

YTsaurus

Data Lake
Raw, ODS, CDM
Sandboxes

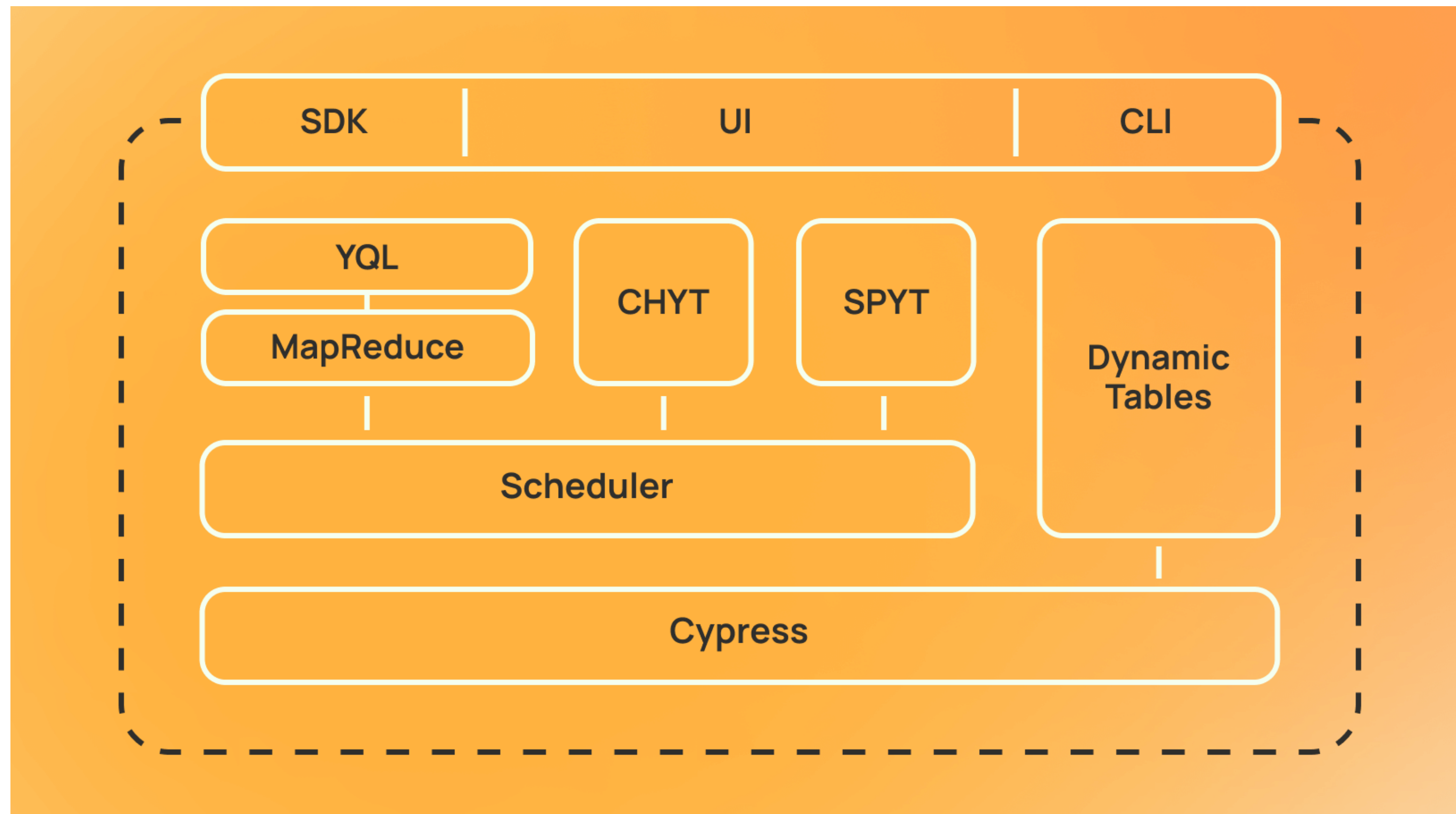
BI

BI with odbc, jdbc or
ClickHouse support

YTsaurus



Архитектура YTsaurus



1. Cypress

Кипарис - распределенная файловая система и хранилище метаданных. Хранит row- и column-oriented таблицы для batch и OLTP нагрузки

2. Scheduler

Планировщик с поддержкой парадигмы MapReduce и на основе принципа Dominant Resource Fairness

3. Compute

MapReduce и высокоуровневые вычислительные движки:

- YQL
- CHYT
- SPYT

Интерфейс

Удобный, понятный и с богатой функциональностью

The screenshot shows the Csupress interface with a table of resources. The table has columns for Name, Locks, Account, Disk space, Rows, Modification time, and Creation time. The data is as follows:

Name	Locks	Account	Disk space	Rows	Modification time	Creation time
account_balance		resource_planner	132.21 KiB	≈ 709	17 Feb 2023, 09:22	03 Feb 2023, 16:52
account_statistics		resource_planner	-	-	17 Feb 2023, 10:30	10 Feb 2023, 22:19
bundle_balance		resource_planner	24.21 KiB	≈ 134	17 Feb 2023, 11:57	10 Feb 2023, 18:21
bundle_statistics		resource_planner	-	-	13 Feb 2023, 22:13	12 Feb 2023, 11:57
equipment_delivery_dates		resource_planner	0 B	-	15 Feb 2023, 11:21	14 Feb 2023, 11:21
orders		resource_planner	1.16 MiB	2 746	17 Feb 2023, 16:52	10 Feb 2023, 16:33
orders_backup		resource_planner	228.22 KiB	≈ 351	10 Feb 2023, 16:33	01 Feb 2023, 11:57
orders_statistics		resource_planner	142.78 KiB	350	13 Feb 2023, 18:21	10 Feb 2023, 10:57
pool		resource_planner	-	-	12 Feb 2023, 11:57	05 Feb 2023, 19:57
pool_statistics		resource_planner	-	-	15 Feb 2023, 10:34	12 Feb 2023, 14:78

Csupress – распределённая файловая система и дерево метаданных

The screenshot shows the Scheduler interface with resource usage bars and a configuration table. The resource usage is as follows:

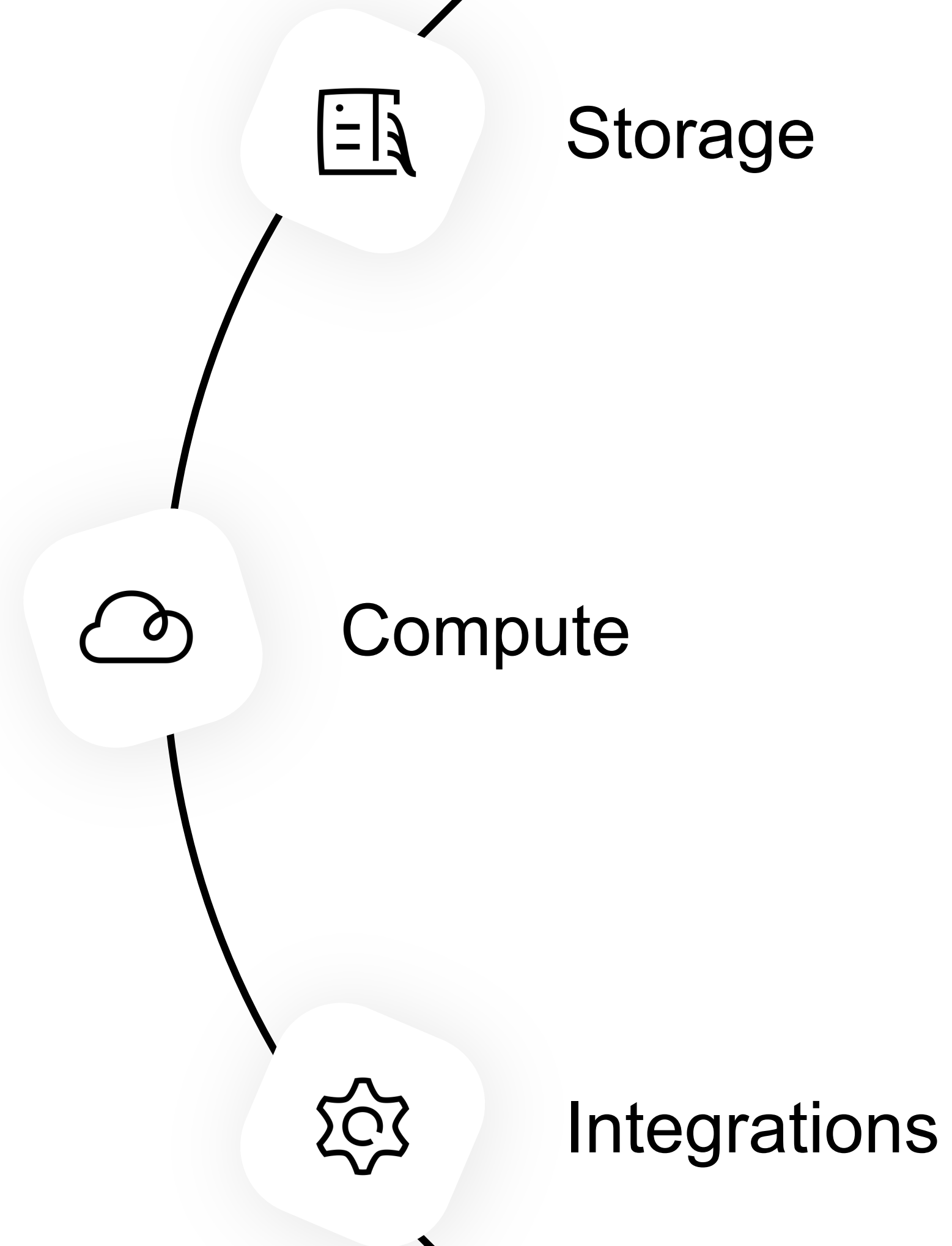
Resource	Usage
CPU	8 176 / 10 606
User memory	35.77 TiB / 66.52 TiB
Network	18.59 %
GPU	0 / 0
User slots	7 404 / 163 600

The configuration table is as follows:

Pool / Operation	FI	Weight	Usage/Fair share	Fair share	Usage	Demand	Guarantee	Operations	Dom. res.
babenko		1 000	-	0	0	0	0	0 / 0	CPU
cron		1		0.000233	0.000233	0.000233	0	1 / 1	Memory
cron_compression		1		0	0	0	0	0 / 0	CPU
cron_master_snapshot		1		0.000233	0.000233	0.000233	0	1 / 1	Memory
cron_merge		1		0	0	0	0	0 / 0	CPU
maps	Integral	1		0	0	0	0	0 / 0	CPU
b2b	Burst	1		0	0	0	0	0 / 0	CPU
production		1		0.065842	0.064074	0.065842	0	9 / 9	Memory

Scheduler – планировщик, который управляет несколькими видами ресурсов (CPU, RAM, GPU) на основе принципа Dominant Resource Fairness

О чем будем говорить



Storage

The screenshot shows a storage management interface with a table of resources. The table has columns for Name, Locks, Account, Disk space, Rows, Modification time, and Creation time. The resources listed include account_balance, account_statistics, bundle_balance, bundle_statistics, equipment_delivery_dates, orders, orders_backup, orders_statistics, pool, and pool_statistics.

Name	Locks	Account	Disk space	Rows	Modification time	Creation time			
account_balance		resource_planner	132.21 KiB	≈ 709	17 Feb 2023, 09:22	03 Feb 2023, 16:52	@	✎	🗑️
account_statistics		resource_planner	-	-	17 Feb 2023, 10:30	10 Feb 2023, 22:19	@	✎	🗑️
bundle_balance		resource_planner	24.21 KiB	≈ 134	17 Feb 2023, 11:57	10 Feb 2023, 18:21	@	✎	🗑️
bundle_statistics		resource_planner	-	-	13 Feb 2023, 22:13	12 Feb 2023, 11:57	@	✎	🗑️
equipment_delivery_dates		resource_planner	0 B	-	15 Feb 2023, 11:21	14 Feb 2023, 11:21	@	✎	🗑️
orders		resource_planner	1.16 MiB	2 746	17 Feb 2023, 16:52	10 Feb 2023, 16:33	@	✎	🗑️
orders_backup		resource_planner	228.22 KiB	≈ 351	10 Feb 2023, 16:33	01 Feb 2023, 11:57	@	✎	🗑️
orders_statistics		resource_planner	142.78 KiB	350	13 Feb 2023, 18:21	10 Feb 2023, 10:57	@	✎	🗑️
pool		resource_planner	-	-	12 Feb 2023, 11:57	05 Feb 2023, 19:57	@	✎	🗑️
pool_statistics		resource_planner	-	-	15 Feb 2023, 10:34	12 Feb 2023, 14:78	@	✎	🗑️

Таблицы

- Статические – основной тип таблиц
- Динамические – key-value хранилище
- Row- & column-oriented
- Богатая система типов

Транзакции

- ACID
- Могут длиться часами и даже днями
- Затрагивают Cypress, а не только таблицы

Квоты

- Диск: HDD, SSD
- Кол-во узлов в Cypress
- Кол-во чанков
- Кол-во таблеток (шардов) в динамических таблицах

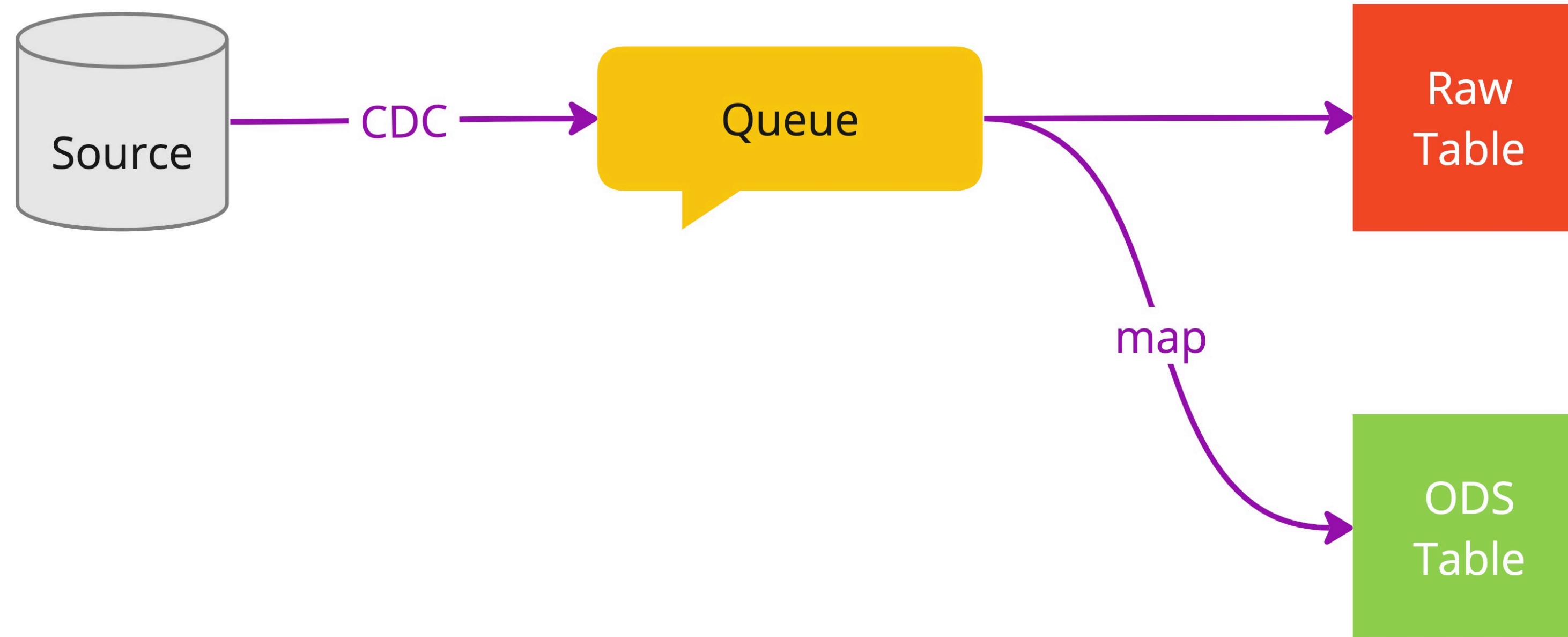
ETL в динамические таблицы

100+ MB/s

- Есть пиковые нагрузки
- Есть особо «жирные» потоки

1000+ сущностей

2+ ПБ



Таблицы в «Raw» слое

Ключи

- Primary key
- Partition key

Док

- Документ как есть в формате Yson

ETL дата

"id"↓	"doc"	"etl_updated"
"7d9d6b4f..."	<pre>{ "_id": ..., "": 7, "": [{ "": "...", "": "...", "": 7, "": "...", }], "": [{ "": { "": 0, "": 317, "": [...], }, "": ... }, { "": "...", "": 1, "": null, "": "...", "": "...", "": "...", "": "...", "": "...", "": "...", }], "": "...", }</pre>	"..."

Таблицы в ODS

Column-oriented

- Усложняет вставку
- Упрощает чтение
- Читаем данные мы чаще, чем пишем

"order_id"↓	"a"	"b"	"c"	"d"
"1"	"a"	666666666	66666667	null
"2"	"a"	5039999999995	6666667	null
"3"	null	null	null	null
"4"	"a"	4071975816	5262087	null
"5"	"a"	809826286	354255	null
"6"	"a"	9155188	1155529	null
"7"	"a"	3038841215	55557338553	"a"
"8"	"y"	65	999999995	null
"9"	"a"	974861	54665287532	null
"10"	"a"		5333334	null
"11"	"a"	3309841454	555529270184	null
"12"	null	null	null	null
"13"	null	null	null	null
"14"	"a"	666666667	55518	null
"15"	"a"	3333333	6666667	null

Схема

3NF, ну почти

- Разрешаем массивы и словари

Статические таблицы для логов и витрин

Sorting

- Может гарантировать уникальность ключа
- Ускоряет фильтрацию, например, по дате
- Ускоряет join-ы по ключу сортировки
- Life hack: может заметно экономить диск для column-oriented таблиц*

Schema

- Примитивы
- Сложные типы: массивы, словари, структуры и т.п.
- Есть constraint: required (not null)
- Может быть strong, weak или отсутствовать

Partitioning

- Нет нативной поддержки, но есть стандартный подход
- Эмулируется через папки: папку считаем таблицей, а таблице в ней – партициями
- Ускоряет обработку и экономит вычислительные ресурсы

Таблицы: управляем размером

Compression

- Большой выбор кодеков сжатия
- Можно выбирать разные кодеки для «горячих» и «холодных» данных

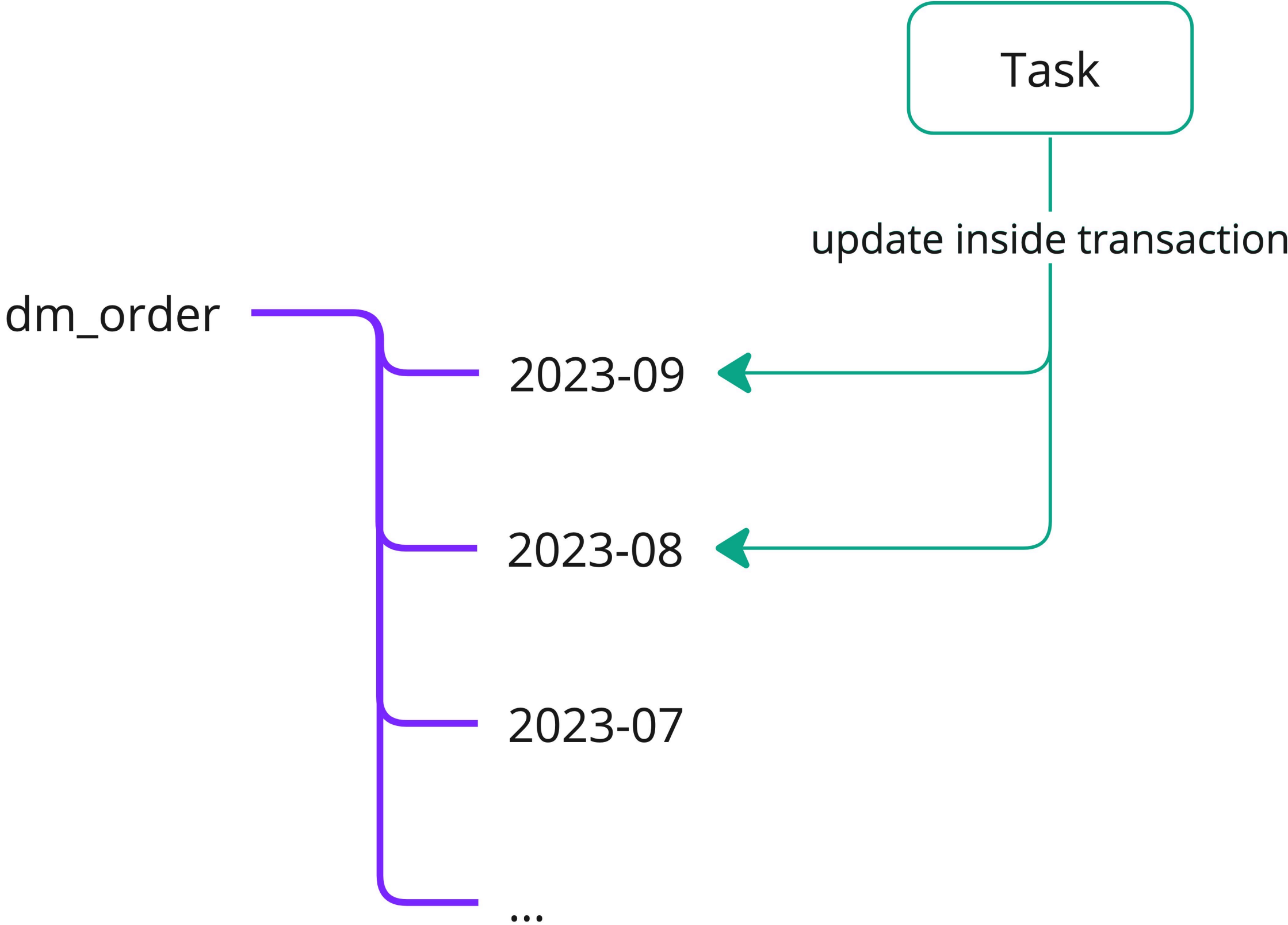
Time To Live

- Expiration Time – момент времени, когда надо удалить узел
- Expiration Timeout – время без обращений к узлу

Redundancy

- Replication: меняем диск на CPU, управляем количеством реплик
- Erasure Coding: меняем CPU на диск, подходит для «холодных» редко используемых данных, объем диска от 1.33x до 1.5x от сжатых данных

Транзакции: расчет за N дней



ACID

Isolation:

- Serializable для таблиц
- Read Committed для Cypress

N узлов Cypress

- Можно пересчитать всю историю и обновить в рамках одной транзакции
- Можно изменять структуру файловой системы

Незаметно для пользователя

- Snapshot Lock на чтение – read-only копия узла
- Exclusive Lock на модификацию узла
- Shared Lock на модификацию части состояния узла



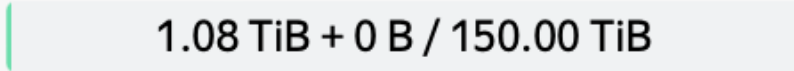
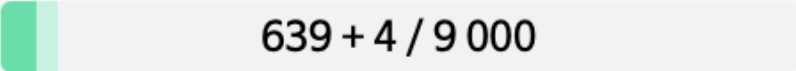
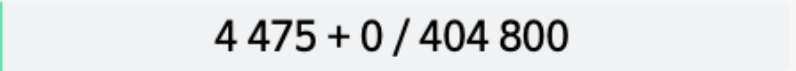


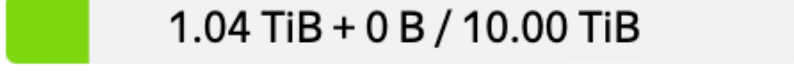
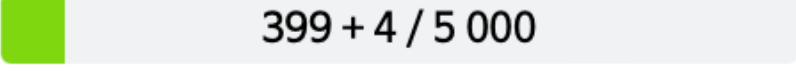
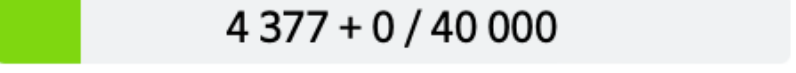


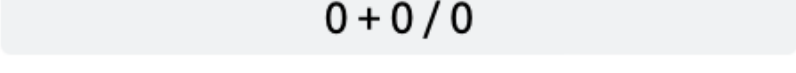
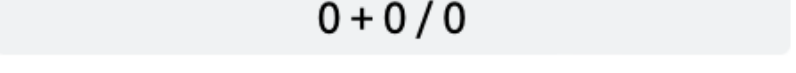


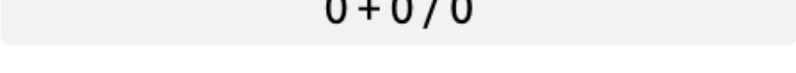
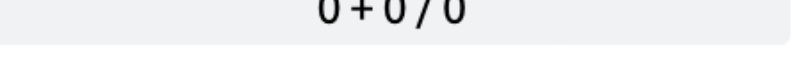


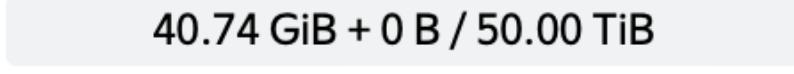
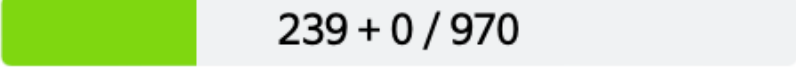
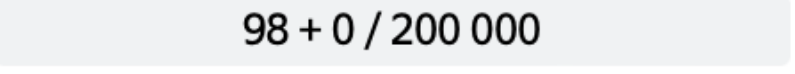

Изоляция: квоты и аккаунты

КВОТЫ:

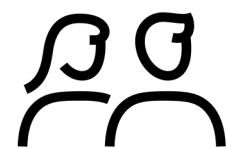
- Размер данных: HDD, SSD
- Кол-во объектов: nodes, chunks

Аккаунты:

- Аккаунты могут быть вложенными, это упрощает управление квотами
- Аккаунт может быть указан для любого узла
- Может быть over commit

Name ^	A ⌵	Disk space ⌵	Nodes ⌵	Chunks ⌵	
^ go-home	 	 1.08 TiB + 0 B / 150.00 TiB	 639 + 4 / 9 000	 4 475 + 0 / 404 800	
∨ dmp		 1.04 TiB + 0 B / 10.00 TiB	 399 + 4 / 5 000	 4 377 + 0 / 40 000	
dwh		-	 0 + 0 / 0	 0 + 0 / 0	
ml		-	 0 + 0 / 0	 0 + 0 / 0	
taxi-analytics		 40.74 GiB + 0 B / 50.00 TiB	 239 + 0 / 970	 98 + 0 / 200 000	

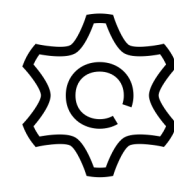
Управление доступами



Доступ можно выдать на

- Списки пользователей
- Группы

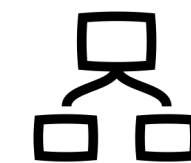
Группы могут быть созданы внутри самого YTsaurus.



Каждый объект в YTsaurus имеет **Access Control List (ACL)**.

Пример:

```
"acl": [  
  {  
    "action": "allow",  
    "subjects": ["user1"],  
    "permissions": ["read"],  
  },  
  ...  
]
```



ACL может наследоваться от родительского узла в Cypress.

Это позволяет гибко и удобно управлять доступами и минимизировать их количество.



Можно настраивать **Column Level Security**.

CLS позволяет делать меньше таблиц, но скрывать доступ к чувствительным данным, например, к финансовой информации.

Сравним Storage

YTsaurus

Плюсы:

- Нативная поддержка таблиц
- Нативный Key-Value Storage
- Транзакционность
- Enterprise Level управление квотами и правами
- Удобный и функциональный UI

Минусы:

- В самом начале пути в Open Source
- Нет готовых специалистов на рынке

s3 & hdfs

Плюсы:

- Давно в Open Source: всем известны плюсы, минусы и грабли
- «Много» готовых специалистов
- Огромная экосистема Hadoop & Open Source Big Data Stack
- Доступны в «любом» облаке

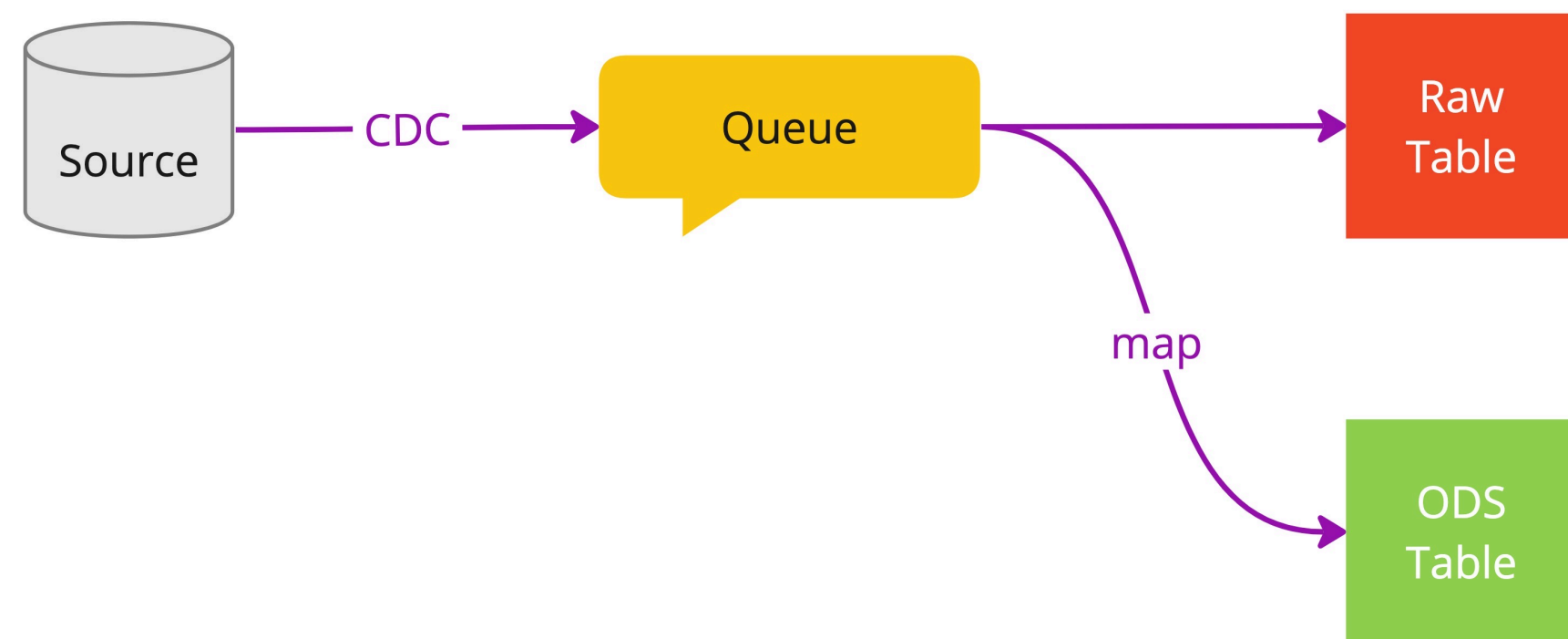
Минусы:

- Нет транзакций
- Для таблиц или key-value нужны отдельные «кубики»

MapReduce в 2023?!..

MapReduce в 2023?!..

Потоковая поставка данных в Raw и ODS

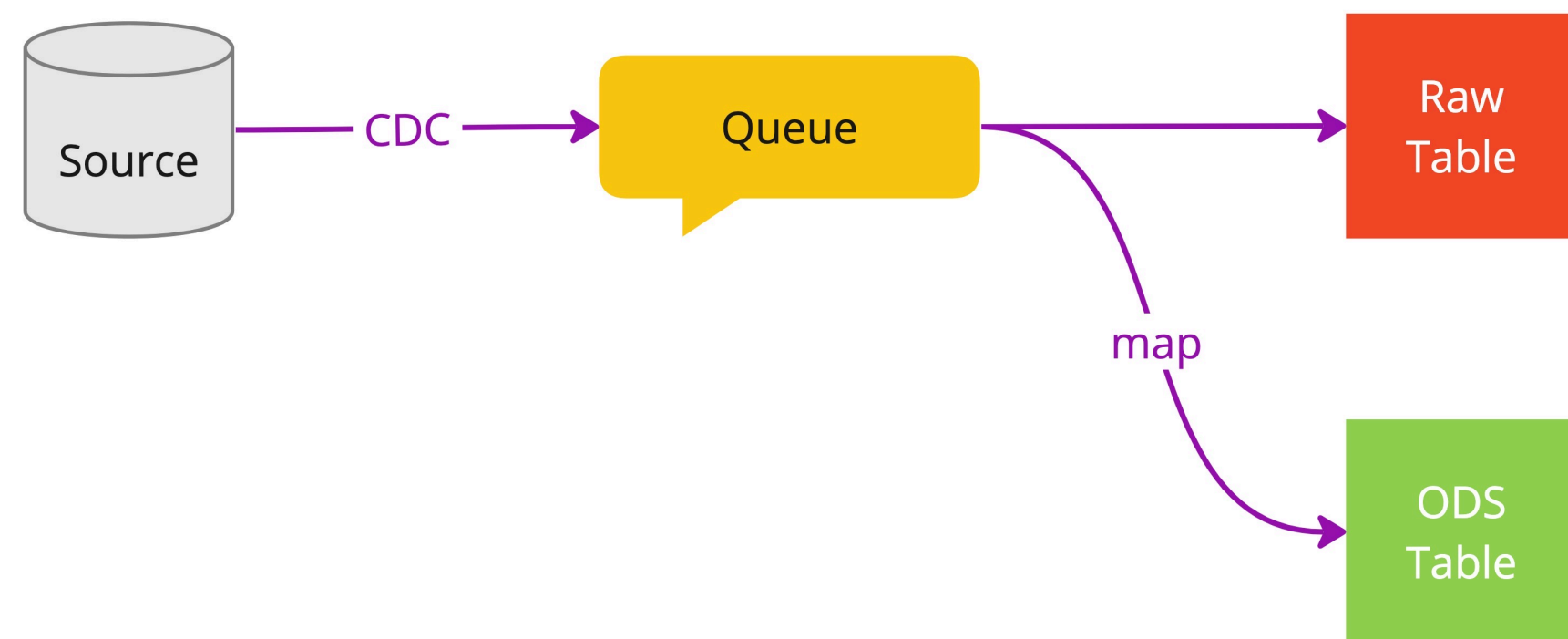


Проблемы пересчета истории:

- Прогонять все через очередь – долго и может тормозить поставку свежих данных
- Источник может не хранить историю
- Могут быть проблемы с огромной пиковой нагрузкой

MapReduce в 2023?!.. Для пересчета истории

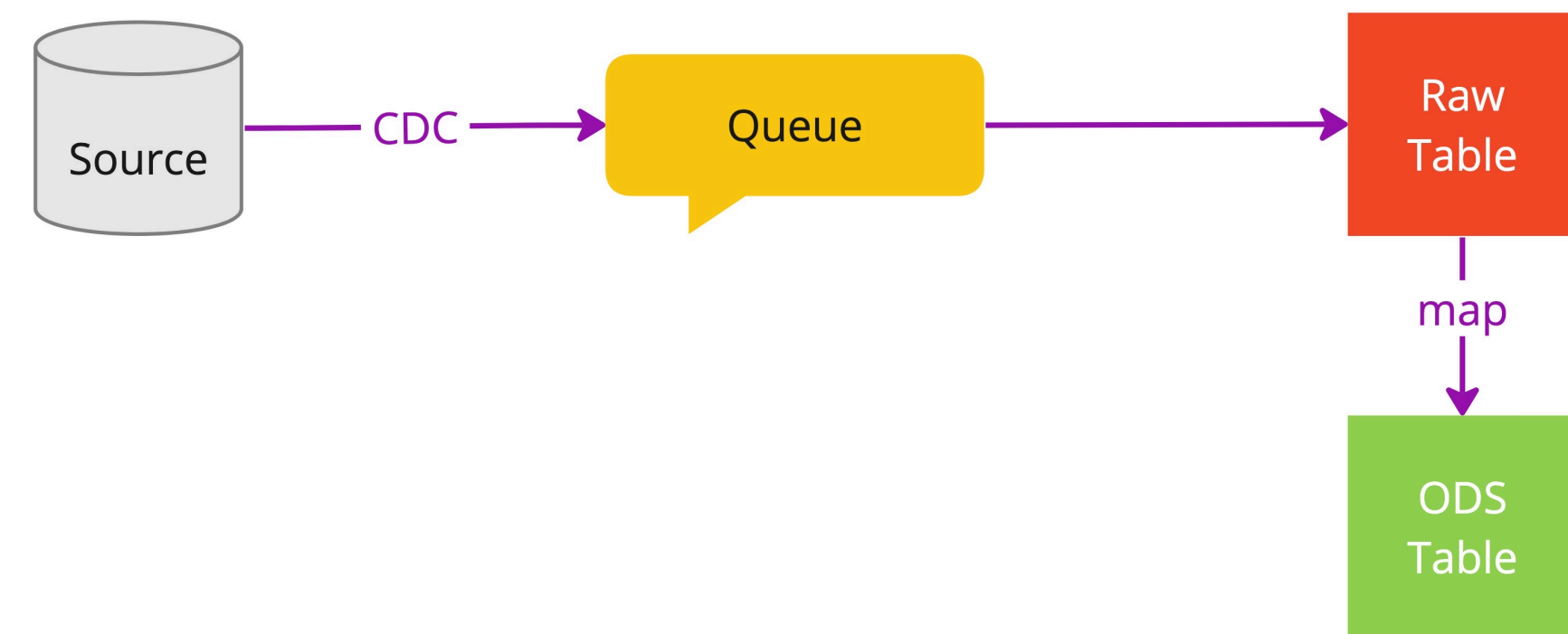
Потоковая поставка данных в Raw и ODS



Проблемы пересчета истории:

- Прогонять все через очередь – долго и может тормозить поставку свежих данных
- Источник может не хранить историю
- Могут быть проблемы с огромной пиковой нагрузкой

Backfill ODS



Почему Map из Raw данных удобен?

- Нет необходимости останавливать текущую поставку*
- Нет зависимости на источник данных
- Хорошо масштабируется, нативная поддержка Python, Java, Go и C++

Построение графа
вычислений на основе
MapReduce операций

The screenshot shows the YQL interface with a query editor on the right and a list of queries on the left. The query editor contains the following code:

```
1 use cluster-name;
2
3 PRAGMA yt.InferSchema = '1';
4
5 -----
6 -- JOIN clause allows to combine columns from multiple input tables in
7 -- the single resulting table.
8 --
9 -- Unless stated otherwise, rows are matched by equality of
10 -- the specified column values and the cartesian product is created if
11 -- multiple rows have the same value.
12 --
13 -- Task: append lastnames to users, which are listed in a separate table,
14 -- and omit those who are 30 years old or older.
15 -----
16 SELECT
17   a.name AS name,    -- "a" and "b" here are so called
18   b.lastname,        -- "correlation names", table aliases
19   a.age,             -- declared below with AS; they are used
20   a.last_time_on_site -- to resolve ambiguity of mentioned column names
21 FROM
```

The result table shows the following data:

#	a.age	a.last_time_on_site	b.lastname	name
1	15	0.5	"Ivanova"	"Anya"
2	25	5	null	"Petr"
3	17	10.5	"Petrova"	"Masha"
4	5	22.5	null	"Alena"
5	23	17.5	null	"Irina"
6	13	19.5	null	"Inna"
7	27	12	null	"German"

In-memory

При определенных условиях
весь граф вычислений или его
часть может быть выполнена
в памяти без MapReduce
операций

UDF

- C++
- Python
- JavaScript

YQL. MapReduce ≠ плохо

The screenshot shows the YQL interface with a query editor and a results table. The query is as follows:

```
1 use cluster-name;
2
3 PRAGMA yt.InferSchema = '1';
4
5 -- JOIN clause allows to combine columns from multiple input tables in
6 -- the single resulting table.
7
8 -- Unless stated otherwise, rows are matched by equality of
9 -- the specified column values and the cartesian product is created if
10 -- multiple rows have the same value.
11
12 -- Task: append lastnames to users, which are listed in a separate table,
13 -- and omit those who are 30 years old or older.
14
15
16 SELECT
17   a.name AS name, -- "a" and "b" here are so called
18   b.lastname, -- "correlation names", table aliases
19   a.age, -- declared below with AS; they are used
20   a.last_time_on_site -- to resolve ambiguity of mentioned column names
21 FROM
```

The results table shows the following data:

#	a.age	a.last_time_on_site	b.lastname	name
1	15	0.5	"Ivanova"	"Anya"
2	25	5	null	"Petr"
3	17	10.5	"Petrova"	"Masha"
4	5	22.5	null	"Alena"
5	23	17.5	null	"Irina"
6	13	19.5	null	"Inna"
7	27	12	null	"German"

Clickstream

- Очень много данных
- На первом этапе происходит простая обработка: фильтрация и очистка

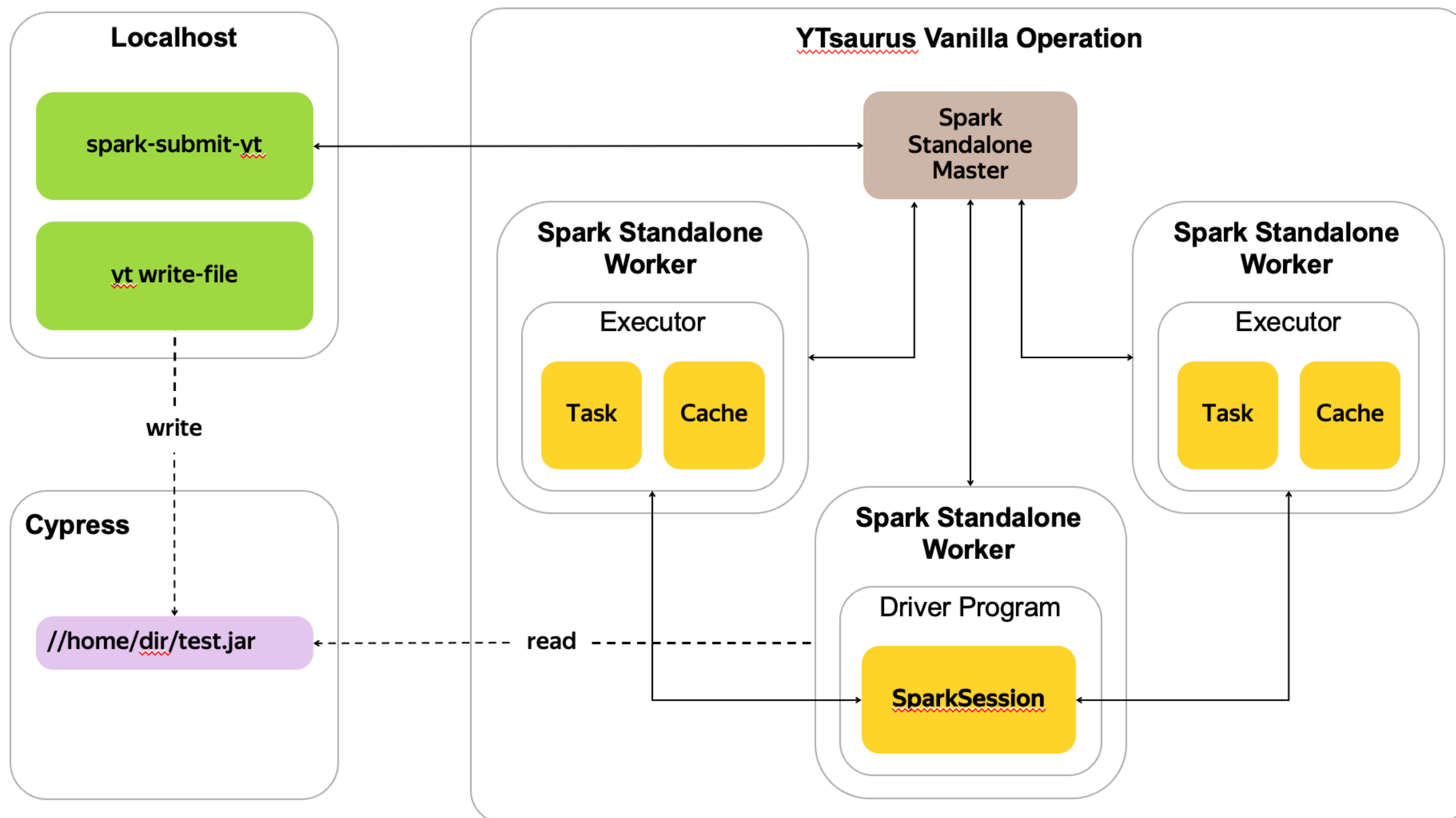
«Большие» логи

- YQL написан на C++ и весьма эффективен
- YQL пользуется фидами YT. Например, MapReduce вместо Map→Sort→Reduce

Витрины

- Простота разработки
- Хорошо масштабируется
- Часть вычислений может быть выполнена In-Memory

Spark over YTsaurus



Fork, но...

- Все изменения касаются только работы с YT и внутри YT
- Доступна вся функциональность «ванильного» Spark-a

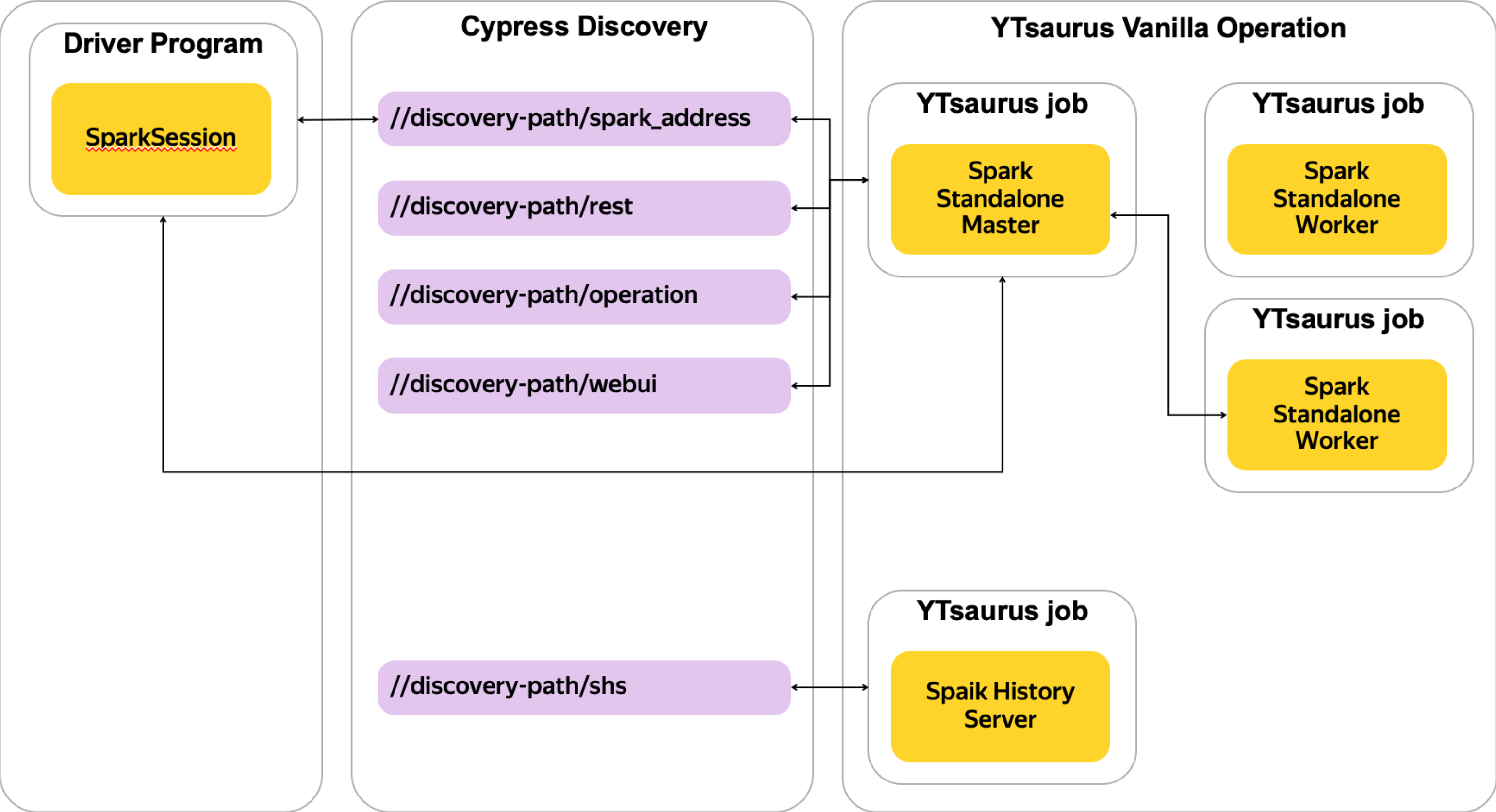
In-Memory

- Чем сложнее расчет, тем больше выгоды ускорения по сравнению с MapReduce
- Spill на SSD или в tmpfs

Внутри YTsaurus

- Легко переиспользуются вычислительные ресурсы между разными движками
- Нет необходимости поддерживать несколько систем

Spark over YTsaurus. Client Mode



Витрины

- Сложная логика: join-ы, агрегации, UDF и т.п.
- Идеально, когда данные помещаются в память

Ad-hoc

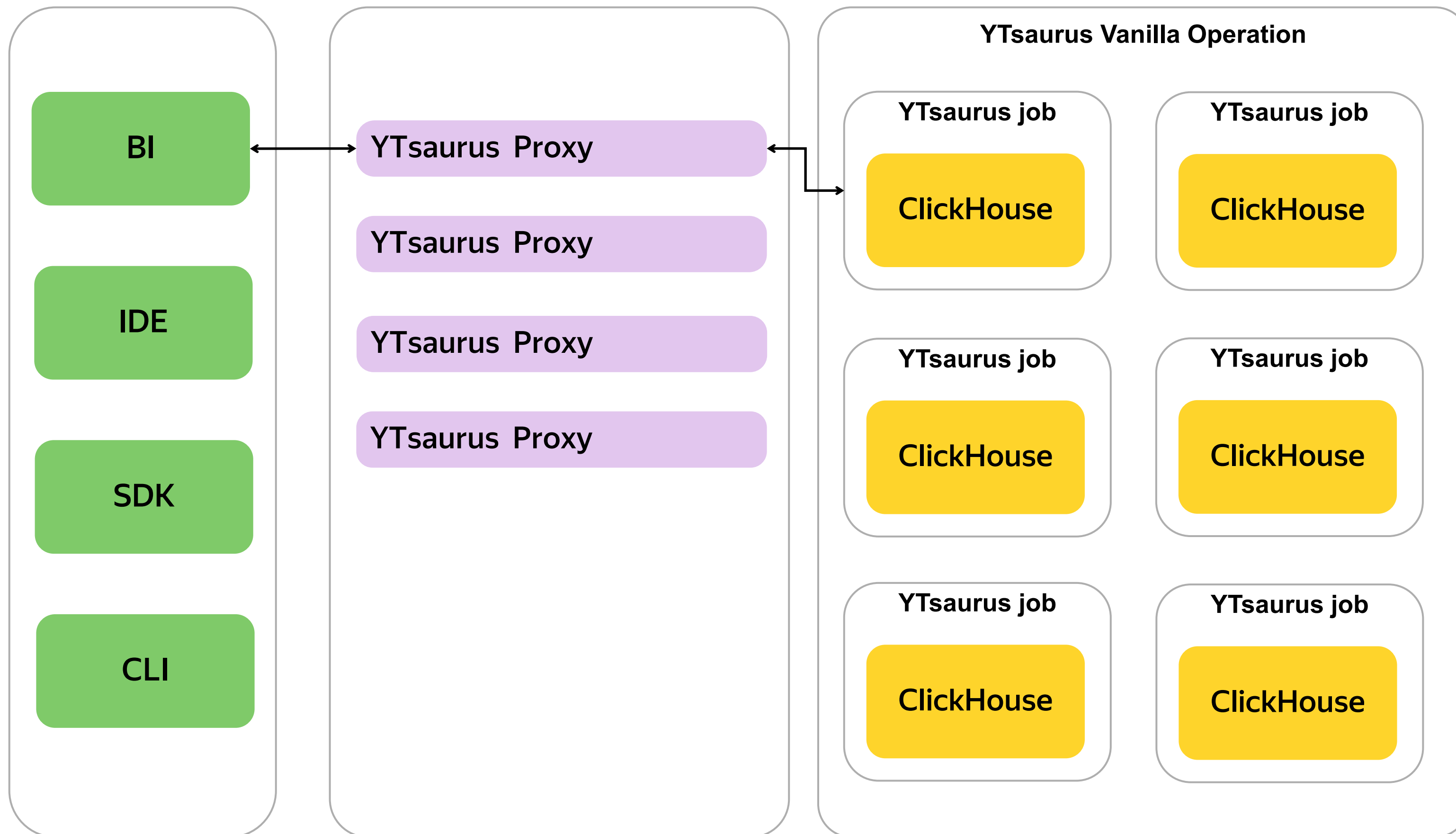
- PySpark прекрасен в Jupyter

SQL

Саша Белоусова про подключение Spark к YTsaurus: <https://highload.ru/spring/2021/abstracts/7266>
Леша Шишкин про оптимизации в SPYT: <https://highload.ru/spb/2023/abstracts/10222>

BI over YTsaurus?..

ClickHouse over YTsaurus



BI

Любой BI с поддержкой коннекторов:

- ClickHouse
- odbc
- jdbc

Ad-hoc

- Скорость работы уступает ClickHouse, но не сильно
- Значительно быстрее YQL и Spark*

Простота

- Не надо копировать данные
- Не надо поддерживать отдельные системы
- Нужно только выделить вычислительные ресурсы

... over YT???!111

Планировщик и пулы

КВОТЫ:

- CPU (Strong, Relaxed Integral, Burst Integral)
- RAM

Особенности:

- Fair Share & Preemption внутри операций пула, внутри поддерева пулов и в кластере в целом
- Поддержка MapReduce: Map, Reduce, MapReduce, Sort, Merge, Erase
- Копирование данных между кластерами: RemoteCopy
- Запуск произвольной нагрузки: Vanilla

Параметры:

- Weight
- Max & Running Operation Count
- ...

<https://ytsaurus.tech/docs/ru/user-guide/data-processing/scheduler/scheduler-and-pools>

Pool / Operation [△]

☰ root

☰ analytics

☰ bi-ch...

△ ☰ dwh

☰ backfill

☰ batch

☰ priority

☰ spyt

☰ ml

YTsaurus vs ...

- Cypress with static row- & column-oriented tables & ACID transactions
- Sorted Dynamic Tables
- Ordered Dynamic Tables

- YQL
- Spark over YT
- ClickHouse over YT

- Fair Share & Preemption планировщик

- hdfs, s3 with Parquet, ORC, Avro

- HBase
- Kafka

- Hive
- Spark
- Impala, Presto, Trino, ...

- YARN или даже Kubernetes

YTsaurus vs ...

- Готовое коробочное решение
- Enterprise level: масштабируемость, надежность, ролевая модель, изоляция хранения и вычислений, красивый UI, ...

- Конструктор с кучей кубиков
- Нужен напильник 😊 или покупка готовой сборки с батарейками

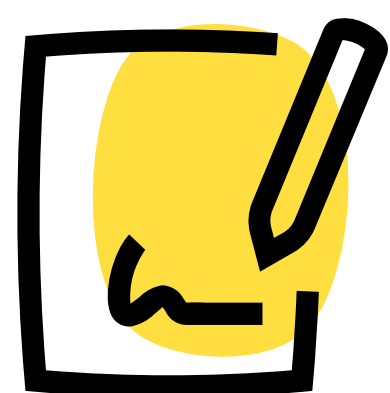
YTsaurus vs ...

- Готовое коробочное решение
- Enterprise level: масштабируемость, надежность, ролевая модель, изоляция хранения и вычислений, красивый UI, ...

- Конструктор с кучей кубиков
- Нужен напильник 😊 или покупка готовой сборки с батарейками

Выводы





Все доступно
Все Open Source

1

YTsaurus

<https://ytsaurus.tech/>

2

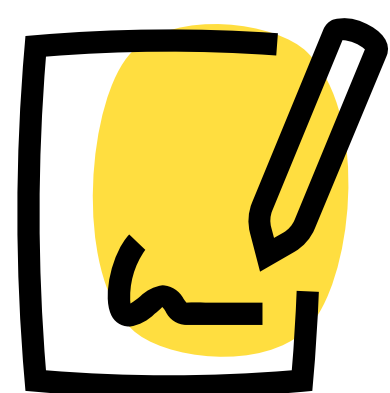
Greenplum

<https://greenplum.org>

3

ClickHouse

<https://clickhouse.com>



YTsaurus

**МОЖНО
попробовать**

1

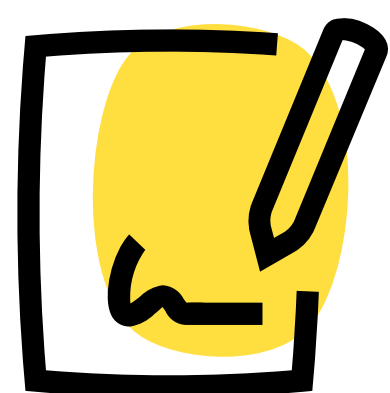
Docker

2

Демо-стенд

3

Kubernetes



Как строить платформу

1

Опираться на накопленный опыт и best practices

2

Идти от требований

3

Экспериментировать в поисках лучшего решения

Яндекс **Go**

Спасибо за внимание!

Вопросы?