

# Data Lineage

Как настроить в зоопарке технологий и зачем это нужно

Булат Усманов, Архитектор

**купер**

# Об авторе

✦ Булат Усманов

✦ Архитектор в Купере

✦ 10+ лет в Fintech, E-Com, etc.;  
AWS Resilience Hub

✦ [usmanovbf@gmail.com](mailto:usmanovbf@gmail.com),  
✈ [@usmanovbf](https://www.instagram.com/usmanovbf)

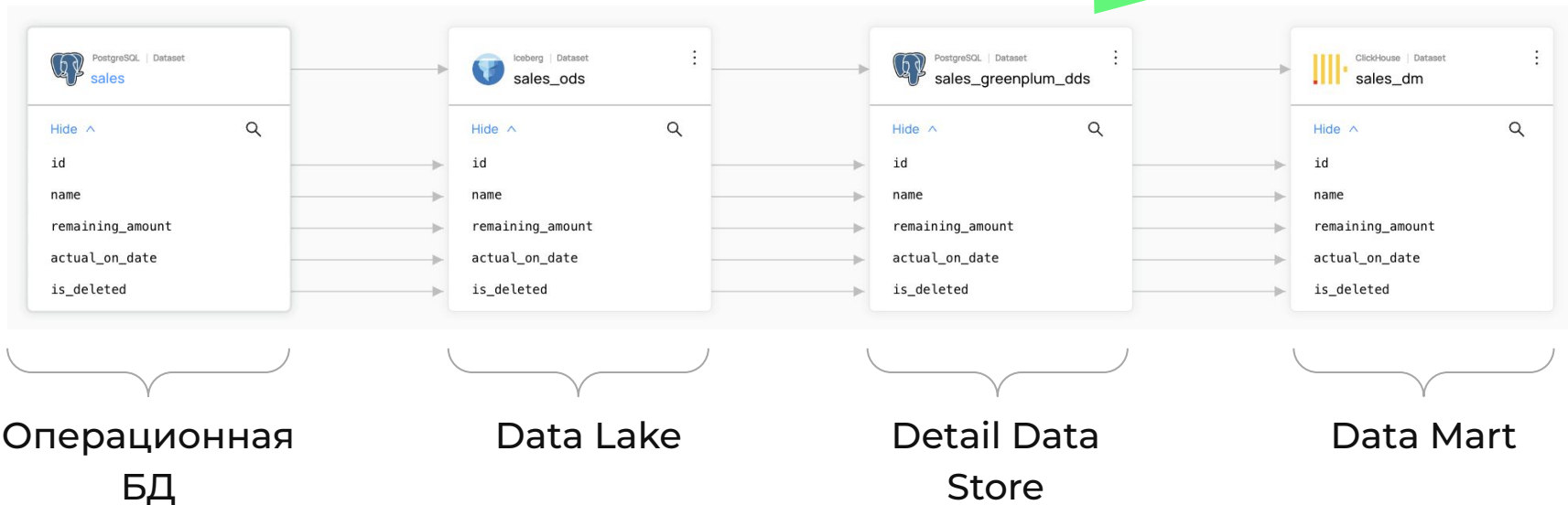


# План

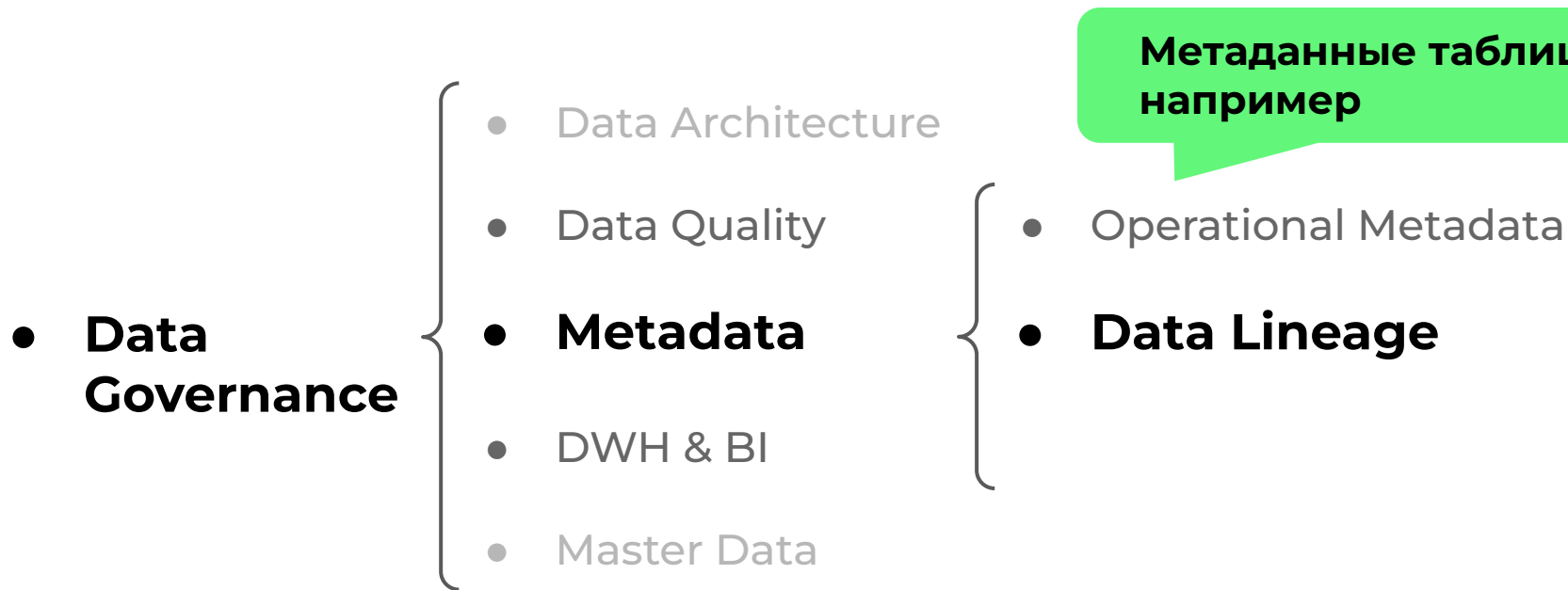
- 01 Что такое **Data Lineage** и его роль в **Data Governance**
- 02 Зачем нужен Data Lineage
- 03 **DataHub** — платформа для Metadata и Data Lineage
- 04 Простое и сложное внедрение Data Lineage
- 05 Решения на основе Data Lineage

# Что такое Data Lineage?

Это же tracing, ха!



# Data Lineage в Data Governance



Если упростить «[Data Management Body of Knowledge](#)»

# Зачем Data Lineage нужен?

**01**  Принятие обоснованных решений и снижение рисков

# Зачем Data Lineage нужен?

01 Принятие обоснованных решений и снижение рисков



02 Соответствие регуляторным требованиям

# Зачем Data Lineage нужен?

**01** Принятие обоснованных решений и снижение рисков

**02** Соответствие регуляторным требованиям

**03** Улучшение координации между коллегами, командами



# Зачем Data Lineage нужен?

01 Принятие обоснованных решений и снижение рисков

02 Соответствие регуляторным требованиям

03 Улучшение координации между коллегами, командами

04 Поддержание Data Mesh, Data Products

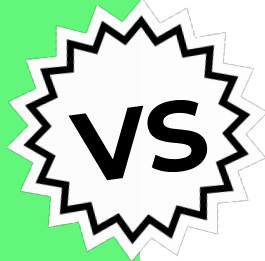


# Зачем Data Lineage нужен?

- 01 Принятие обоснованных решений и снижение рисков
- 02 Соответствие регуляторным требованиям
- 03 Улучшение координации между коллегами, командами
- 04 Поддержание Data Mesh, Data Products
- 05 Ускорение разработки, Time to Market



# DataHub



# OpenMetadata

- › Больше подключений к источникам
- › Богаче API
- › Больше community



Fork

2.8k



Star

9.5k



DataHub

- › Приятнее UX/UI
- › Нагляднее и богаче документация
- › Отслеживание заполняемости

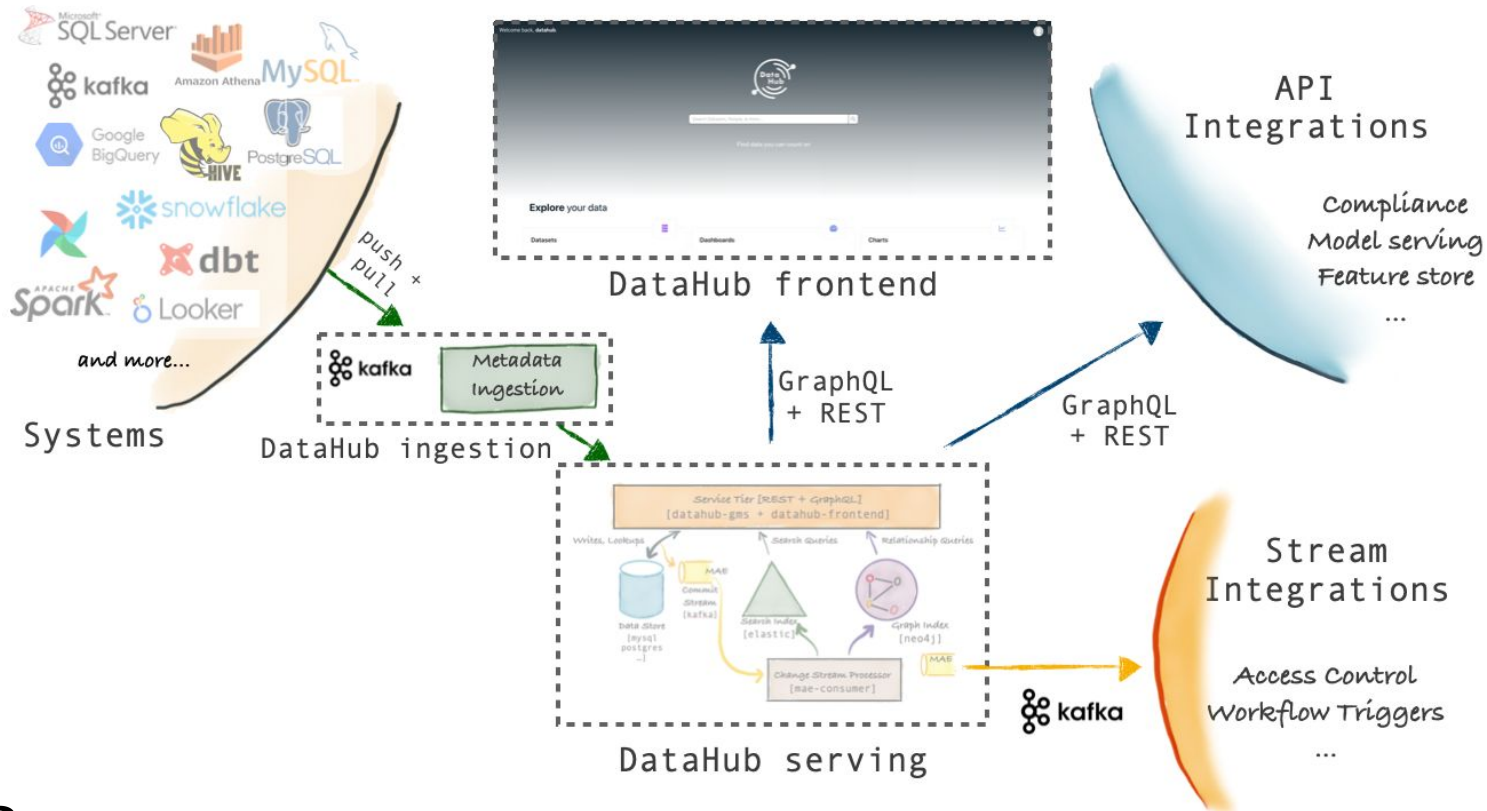


Open  
Metadata

# Приблизились к плюсам OpenMetadata

- › Приятнее UX/UI
  - Привыкли к DataHub, но открыты к улучшениям
  - Дополнительный UI с помощью отчетов Tableau
- › Нагляднее и богаче документация
  - Мы разве исходники не сможем прочесть?
- › Отслеживание заполняемости
  - Отчеты в Tableau, куда данные выгружаются с помощью Airflow

# DataHub. Взгляд сверху



# DataHub. Deploy

datahub-gms — Metadata Service. 3 instances при ~250 MAU, 150k items

datahub-ingestion — выгрузка метаданных с помощью “datahub ingest run -c recipe.yaml” по Cron

datahub-mce-consumer — приём новых метаданных из источников

datahub-mae-consumer — отправка готовых метаданных, для интеграций  
вовне

datahub-actions — фреймворк для кастомизации, для интеграций изнутри

Также ElasticSearch (Neo4j), MySQL, PostgreSQL, Kafka...



kubernetes

# DataHub. Пример метаданных

Datasets > Snowflake > long\_tail\_companions > ecommerce [Details](#) [Lineage](#) 1 upstream, 4 downstream

Table Snowflake > long\_tail\_companions > ecommerce [Share](#)

## users

[Schema](#) [Documentation](#) [Lineage](#) [Properties](#) [Queries](#) [Stats](#) [Quality](#) [Governance](#) [Incidents](#)

Q Search in schema...

Field	Description	Tags	Glossary Terms
pk <span>Struct</span> <span>Primary Key</span>	primary key of users table		
email <span>Struct</span> <span>Nullable</span>	email address currently associated with the user	<span>pii</span>	<span>Email</span>
created_at <span>Struct</span>	UTC timestamp when record was created		

● Last synchronized 4 hours ago

---

### About

this table contains records of all users that have ever made a purchase on the ecommerce site

[+ Add Link](#)

---

### Owners

Other

E Ecommerce X

[+ Add Owners](#)

---

### Part Of

users

# DataHub. Выгрузка метаданных

recipe.yml

```
source:  
  type: "mysql"  
  config:  
    username: "datahub"  
    password: "datahub"  
    host_port: "localhost:3306"  
sink:  
  type: "datahub-rest"  
  config:  
    server: 'http://localhost:8080'
```

```
datahub ingest -c recipe.yml
```

- Pull подход, но есть ещё Push



# DataHub. Кастомизация выгрузки

recipe.yml

```
source:
  type: "mysql"
  config:
    username: "datahub"
    password: "datahub"
    host_port: "localhost:3306"
sink:
  type: "datahub-rest"
  config:
    server: 'http://localhost:8080'
```

```
datahub ingest -c recipe.yml
```

А еще:

- SQL профилирование
- Удаление неактуальных метаданных
- Статистика использования таблиц
- Классификация персональных данных по регулярным выражениям
- Пост-трансформация

# Внедрение Data Lineage

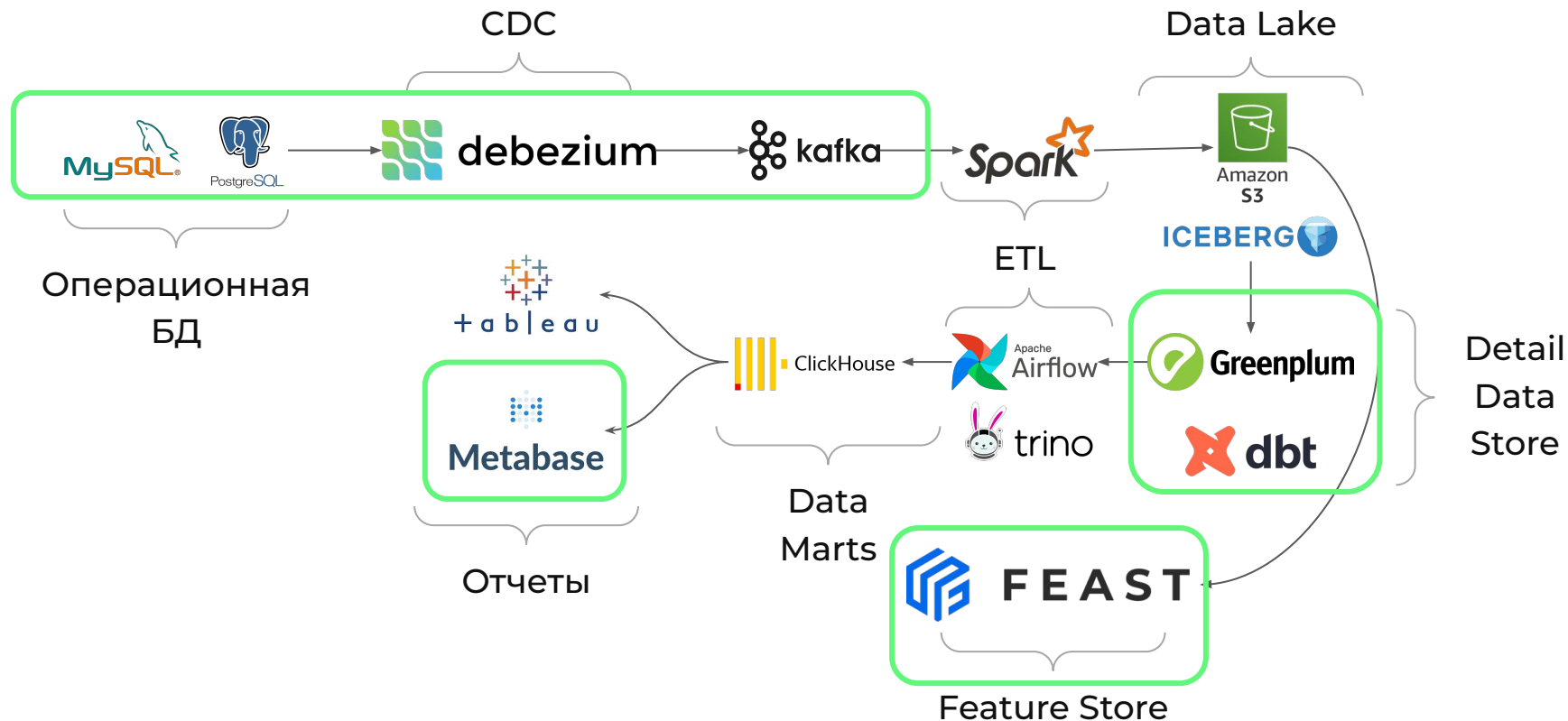
## Простое

1. Написать Ingestion recipe.
2. Запустить.
3. Отслеживать — готово!

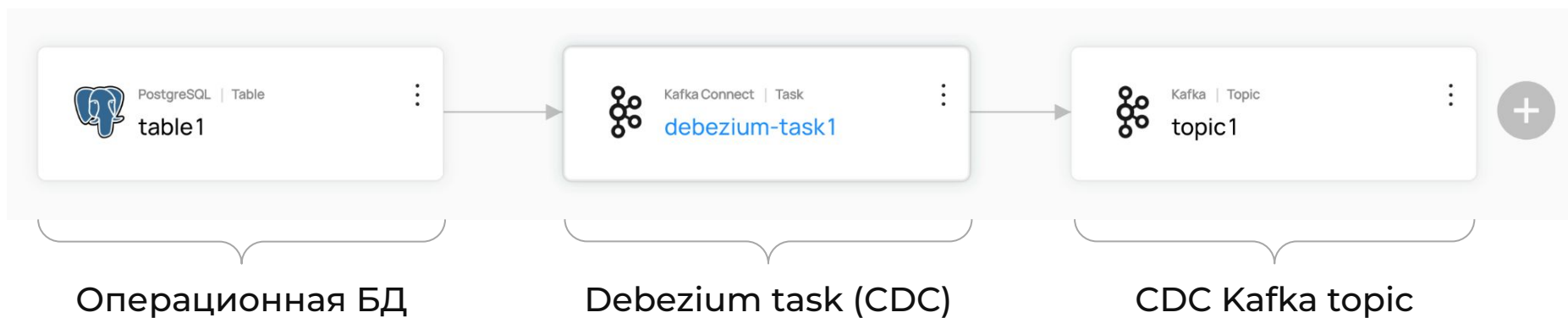
## Посложнее

1. Написать Ingestion recipe.
2. **Кастомизировать выгрузку Metadata и Data Lineage (Java, Python).**
3. Запустить.
4. Отслеживать — готово!

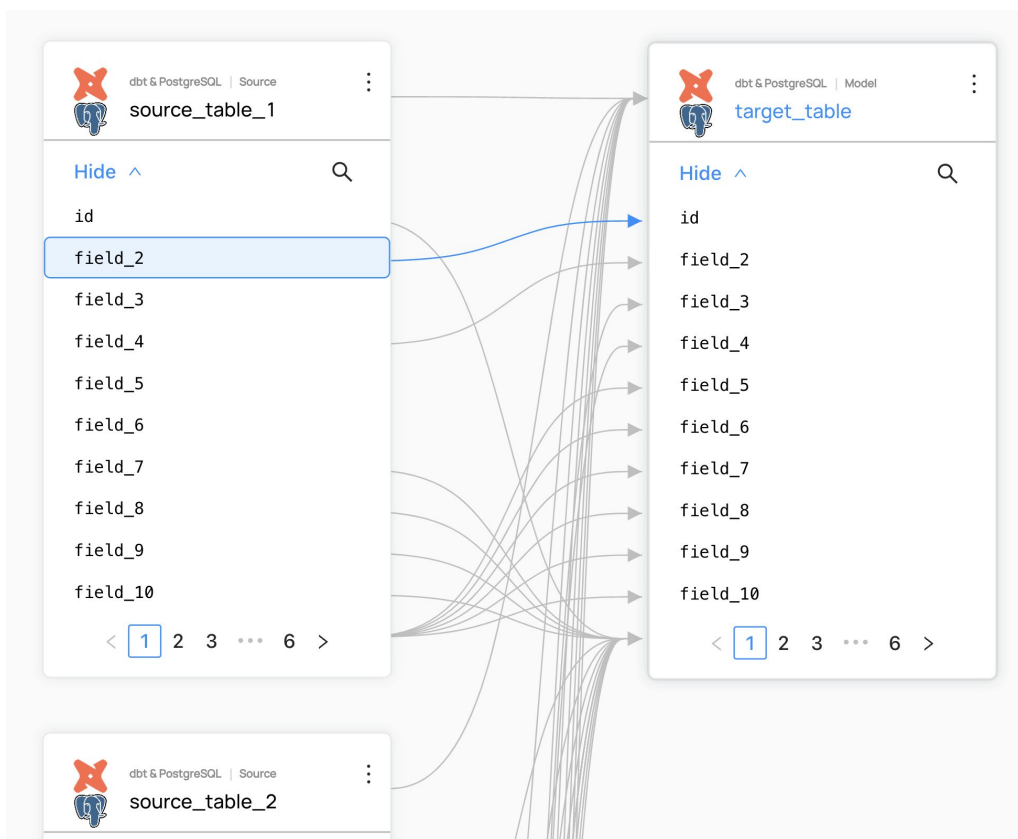
# Простое внедрение Data Lineage



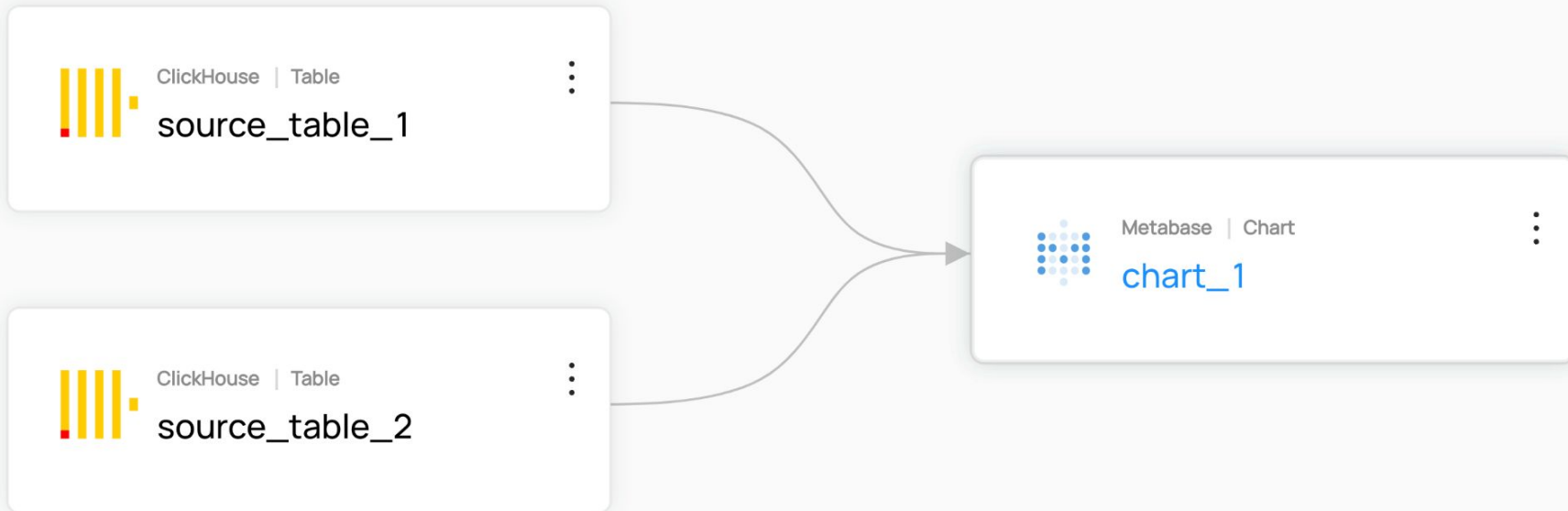
# Lineage для CDC от БД до Kafka



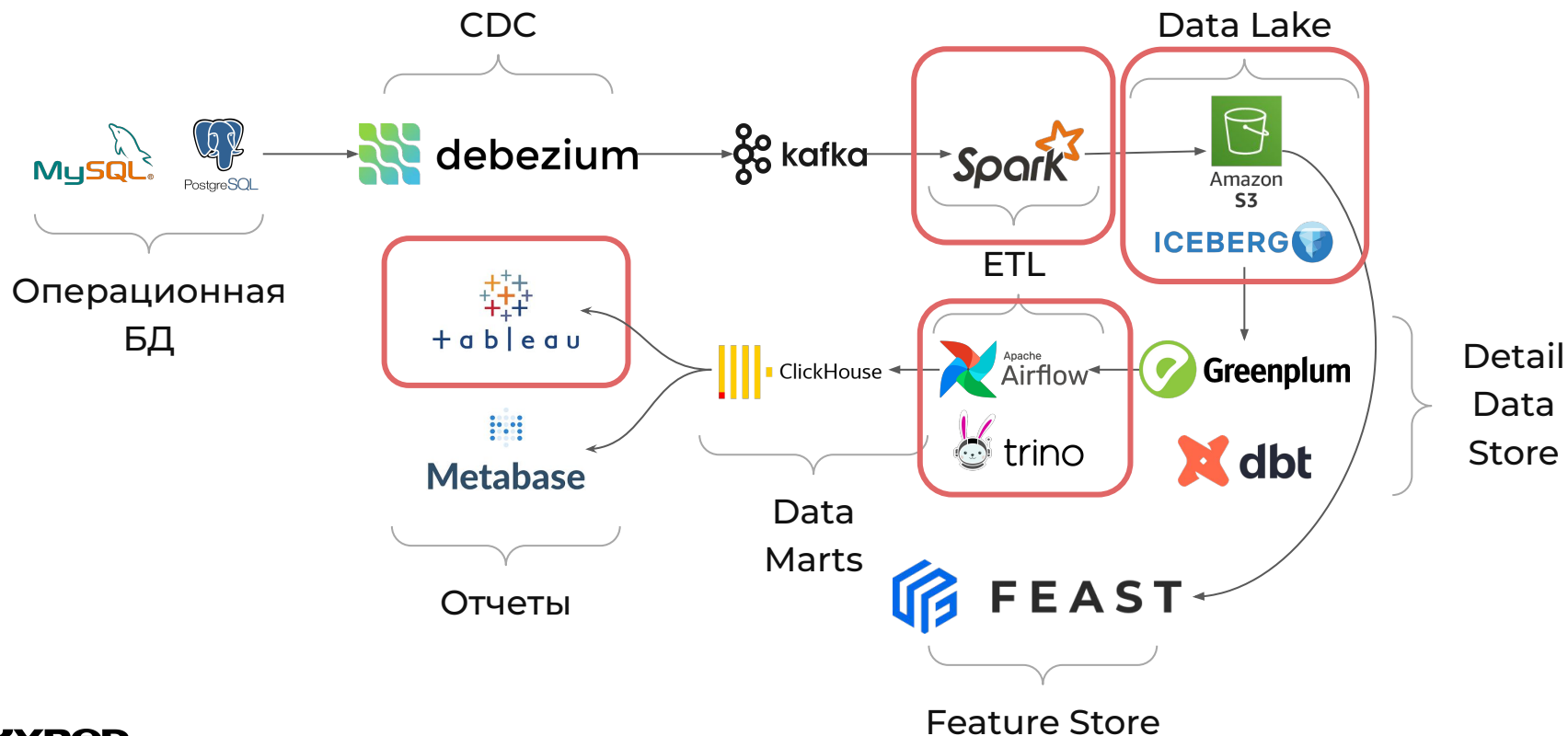
# Lineage для dbt (Greenplum)



# Lineage для Metabase



# Внедрение Data Lineage посложнее



# DataHub. Пост-трансформация. Проблема с S3

```
recipe.yaml

1 source:
2   type: s3
3   config:
4     path_specs:
5     -
6       include: 's3://bucket-name/{table}/{partition_key[0]}={partition[0]}/platform=ios/*.parquet'
7     -
8       include: 's3://bucket-name/{table}/{partition_key[0]}={partition[0]}/platform=android/*.parquet'
9   aws_config:
10    aws_access_key_id: '${ACCESS_KEY}'
11    aws_secret_access_key: '${SECRET_KEY}'
```



bucket-name.some-object для ios, android





# DataHub. Пост-трансформация. Решение

```
recipe.yaml

1 source:
2   type: s3
3   config:
4     path_specs:
5       -
6         include: 's3://bucket-name/{table}/{partition_key[0]}={partition[0]}/platform=ios/*.parquet'
7       -
8         include: 's3://bucket-name/{table}/{partition_key[0]}={partition[0]}/platform=android/*.parquet'
9     aws_config:
10      aws_access_key_id: '${ACCESS_KEY}'
11      aws_secret_access_key: '${SECRET_KEY}'
12
13 transform:
14   - type: "urn_modifier_transformer.urn_modifier_transformer:UrnModifierTransformer"
15     config:
16       regex_to_search: "/platform=(.*)?/"
17       urn_to_modify: "{parsed_urn.entity_ids[1]}-platform-{match.group(1)}"
```



# DataHub. Пост-трансформация. Результат

Dataset | AWS S3 > bucket-name

## some-object-platform-ios

Updated 2 months ago

[Schema](#) Relationships Documentation Lineage

Search in schema...

Field
field_1 Number
field_2 Number

Dataset | AWS S3 > bucket-name

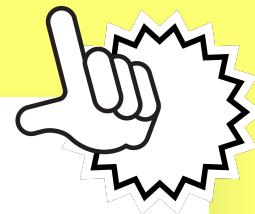
## some-object-platform-android

Updated 2 months ago

[Schema](#) Relationships Documentation Lineage Properties

Search in schema...

Field	Description
field_3 Number	
field_4 Number	



# Data Lineage из Spark. Из коробки

datahub\_spark\_listener\_out\_of\_the\_box.sh

Но не всегда как надо 😞

```
1 spark-submit \  
2   ... \  
3   --packages io.acryl:acryl-spark-lineage:0.2.16 \  
4   --conf "spark.extraListeners=datahub.spark.DatahubSparkListener" \  
5   your_spark_job_to_run.py # Или через JAR
```

Много лишних связей в Lineage, много метаданных

# Data Lineage из Spark. Кастомизация на Java

```
1 public class YourOwnDatahubSparkListener extends org.apache.spark.scheduler.SparkListener {
2
3     @Override
4     public void onApplicationStart(SparkListenerApplicationStart applicationStart) {
5         saveToContext(applicationStart); // Сохранение данных в контексте для дальнейшей обработки
6     }
7
8     @Override
9     public void onOtherEvent(SparkListenerEvent input) {
10        SparkListenerSQLExecutionStart event = (SparkListenerSQLExecutionStart) input;
11        QueryExecution execution = SQLExecution.getQueryExecution(event.executionId());
12        LogicalPlan logicalPlan = execution.logical();
13        extractLineage(logicalPlan); // Извлечение Lineage из текущего запроса в виде LogicalPlan
14    }
15
16    @Override
17    public void onApplicationEnd(SparkListenerApplicationEnd applicationEnd) {
18        emitToDatahub(applicationEnd); // Отправка извлеченного Lineage в DataHub
19    }
20 }
```

# Извлечение Lineage из LogicalPlan

```
LogicalPlan.scala
1 package org.apache.spark.sql.catalyst.plans.logical
2
3 abstract class LogicalPlan
4   extends QueryPlan[LogicalPlan]
5   with AnalysisHelper
6   with LogicalPlanStats
7   with LogicalPlanDistinctKeys
8   with QueryPlanConstraints
9   with Logging {}
10
11 abstract class QueryPlan[PlanType <: QueryPlan[PlanType]]
12   extends TreeNode[PlanType]
```

LogicalPlan похож  
на результат  
от EXPLAIN благодаря  
TreeNode

**Обойдём дерево!**

# Извлечение Lineage из LogicalPlan. Пример

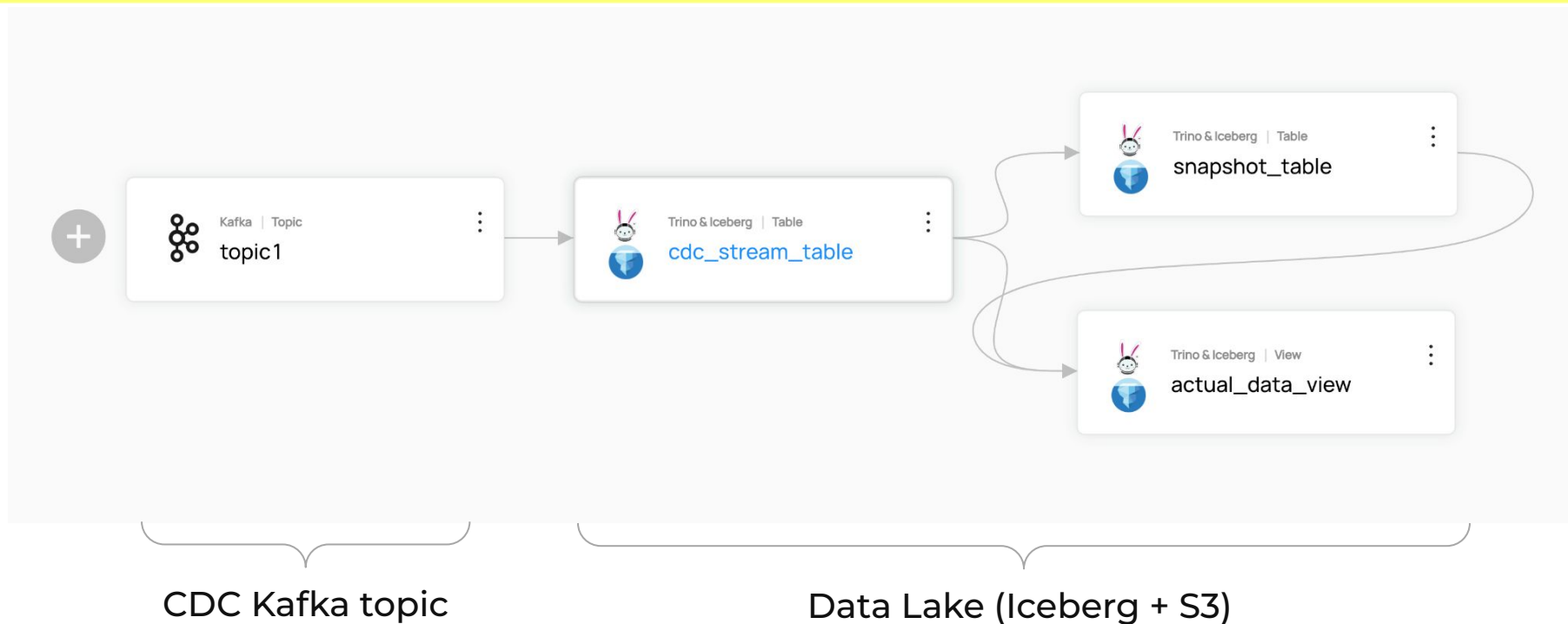


```
1 MERGE INTO iceberg_catalog.table AS target
2 USING iceberg_catalog.table_stream_from_cdc AS source ...
```

```
1 class ReplaceIcebergData extends LogicalPlan:
2
3 +- Sort [id ASC], false
4   +- RepartitionByExpression [id]
5     +- MergeRows[id, version_id, ...]
6       +- Join FullOuter (id = id)
7         :- BatchScan[id, version_id, ...] iceberg_catalog.db_version
8         +- BatchScan[id, version_id, ...] iceberg_catalog.db_version_cdc_stream
```



# Lineage для Spark. От Kafka до Data Lake



# Data Lineage из Airflow. Из коробки

```
simple airflow integration.sh

1 pip install 'acryl-datahub-airflow-plugin[plugin-v2]'
2 airflow connections add --conn-type 'datahub-rest' \
3   'connection_name' \
4   --conn-host 'https://datahub-gms:443' \
5   --conn-password 'password'
```

Но не всегда 😞

Работает для некоторых популярных Airflow операторов в Airflow 2.3+



# Data Lineage из Airflow. Inlets & outlets

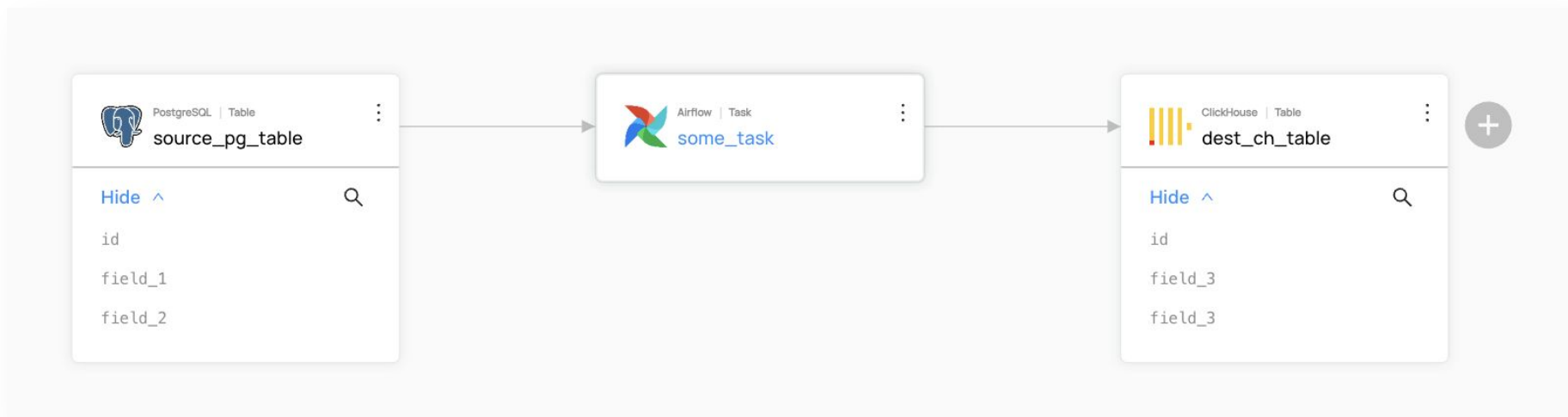
```
1 ParamsInletsOutletsOperator(  
2     task_id="run_data_task", dag=dag,  
3     inlets=[Dataset(platform="postgres", name="db.schema.table1")],  
4     outlets=[Dataset("clickhouse", "db.schema.table2")],  
5 )
```

```
1 class InnerInletsOutletsOperator(BaseOperator):  
2  
3     def execute(self, context):  
4         ...  
5         inlets, outlets = self._get_lineage()  
6         context['ti'].task.inlets = self.inlets  
7         context['ti'].task.outlets = self.outlets  
8  
9     def _get_lineage(self) -> Tuple[List[str], List[str]]:  
10         ...  
11         return inlets, outlets
```

# DataHub SDK внутри оператора со своим API

```
1 # urn:li:dataset:(urn:li:dataPlatform:postgres,instance-1.db.schema.table1,PROD)
2 src_entity = DataHubEntity(
3     # platform="postgres", platform_instance="instance-1"
4     platform_instance=datahub_platform_instances.POSTGRES_PLATFORM_INSTANCE,
5     # db.schema.table1
6     entity_name=full_src_table
7 )
8 dst_entity = DataHubEntity(
9     platform_instance=datahub_platform_instances.CLICKHOUSE_PLATFORM_INSTANCE,
10    entity_name=full_dst_table
11 )
12
13 lineage_items = [DataHubLineageItem(src=[src_entity], dst=dst_entity)]
14
15 datahub_lineage_operator = DataHubLineageOperator(
16     task_id=f"lineage_{dst_table}",
17     lineage=lineage_items,
18     dag=dag
19 )
```

# Data Lineage из Airflow. Результат



# А не лучше ли не прописывать руками?

**01** Не везде используется SQL – Python, Java ETL ...

**02** Парсер SQL запросов SQLGlot

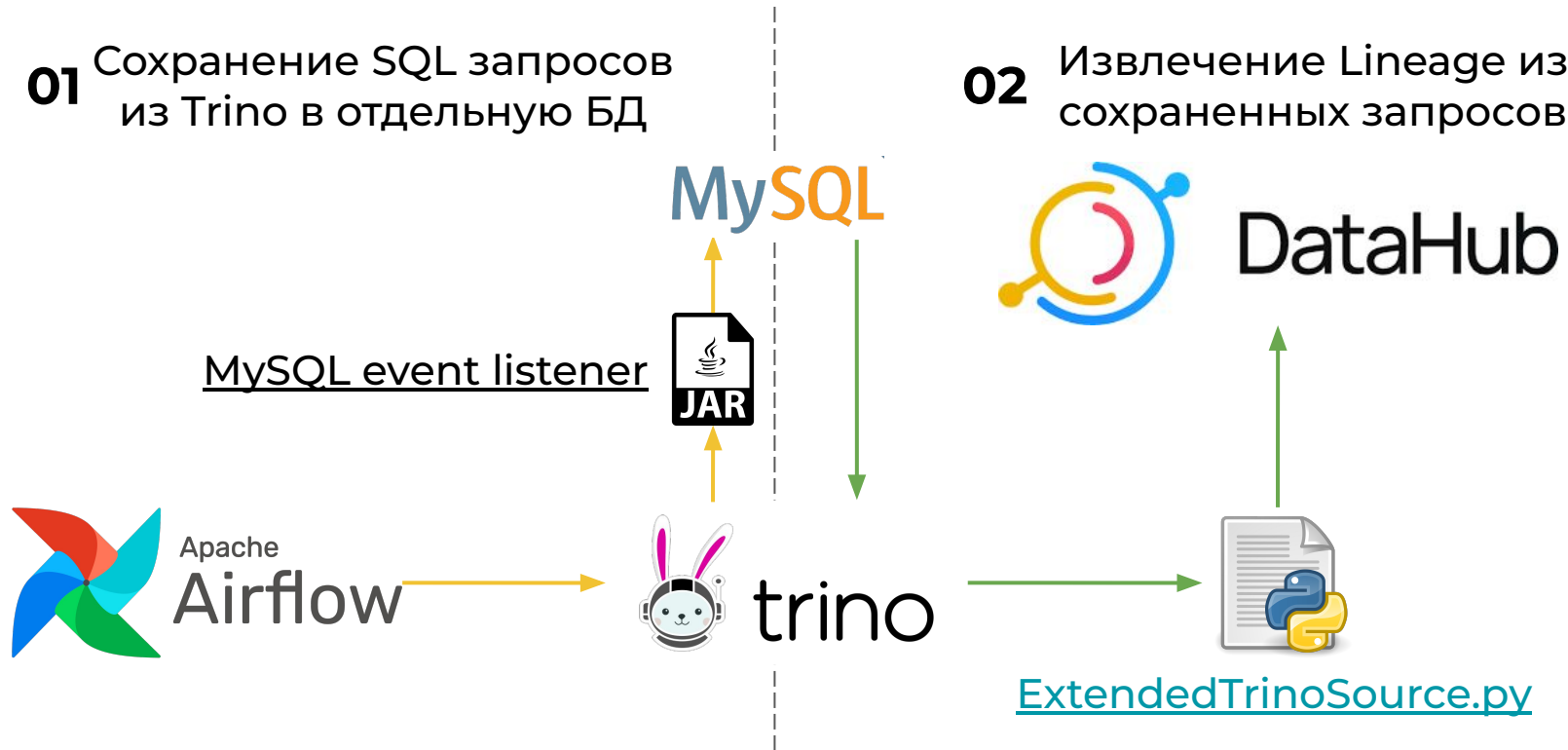


**03** Нужно расширять его диалекты для Postgres, Clickhouse, etc. — покажу на примере Trino

# Lineage из Trino

**01** Сохранение SQL запросов из Trino в отдельную БД

**02** Извлечение Lineage из сохраненных запросов



# Lineage из Trino. Кастомизация Trino ingestion

```
ExtendedTrinoSource.py

1 from datahub.ingestion.source.sql.trino import TrinoConfig, TrinoSource
2 from datahub.sql_parsing.sqlglot_lineage import sqlglot_lineage, SqlParsingResult
3
4 class ExtendedTrinoSourceConfig(TrinoConfig):
5     some_extra_config_param: str = Field(description="Write you own config param")
6
7 class ExtendedTrinoSource(TrinoSource):
8
9     def get_workunits_internal(self) -> typing.Iterable[MetadataWorkUnit]:
10         for inspector in self.get_inspectors():
11             sql_queries: List[str] = find_queries_from_mysql_db(inspector)
12
13             for sql_query in sql_queries:
14                 lineage: SqlParsingResult = sqlglot_lineage(sql=sql_query, ...)
15                 yield extract_metadata_work_units(lineage)
```

# Lineage из Trino. Использование SQLGlott

```
ExtendedTrinoSource.py

1 from datahub.ingestion.source.sql.trino import TrinoConfig, TrinoSource
2 from datahub.sql_parsing.sqlglott_lineage import sqlglott_lineage, SqlParsingResult
3
4 class ExtendedTrinoSourceConfig(TrinoConfig):
5     some_extra_config_param: str = Field(description="Write you own config param")
6
7 class ExtendedTrinoSource(TrinoSource):
8
9     def get_workunits_internal(self) -> typing.Iterable[MetadataWorkUnit]:
10         for inspector in self.get_inspectors():
11             sql_queries: List[str] = find_queries_from_mysql_db(inspector)
12
13             for sql_query in sql_queries:
14                 lineage: SqlParsingResult = sqlglott_lineage(sql=sql_query, ...)
15                 yield extract_metadata_work_units(lineage)
```

# Расширение диалектов SQLGlot. Пример



presto\_ddl\_query.sql

```
1 CREATE VIEW some_view SECURITY DEFINER AS
2 SELECT id, name FROM some_table
```



# Расширение диалекта SQLGlot. Решение

parser.py

```
1 class Parser(metaclass=_Parser):
2     PROPERTY_PARSERS: t.Dict[str, t.Callable] = {
3         ...
4         "SECURITY": lambda self: self._parse_security(),
5     }
6
7     def _parse_security(self) -> t.Optional[exp.SecurityProperty]:
8         if self._match_texts(("DEFINER", "INVOKER")):
9             security_specifier = self._prev.text.upper()
10            return self.expression(exp.SecurityProperty, this=security_specifier)
11            return None
```

expressions.py

```
1 class SecurityProperty(Property):
2     arg_types = {"this": True}
```

# Расширение диалекта SQLGlot. Решение



[sqlglot PR #4008](#) – Merged

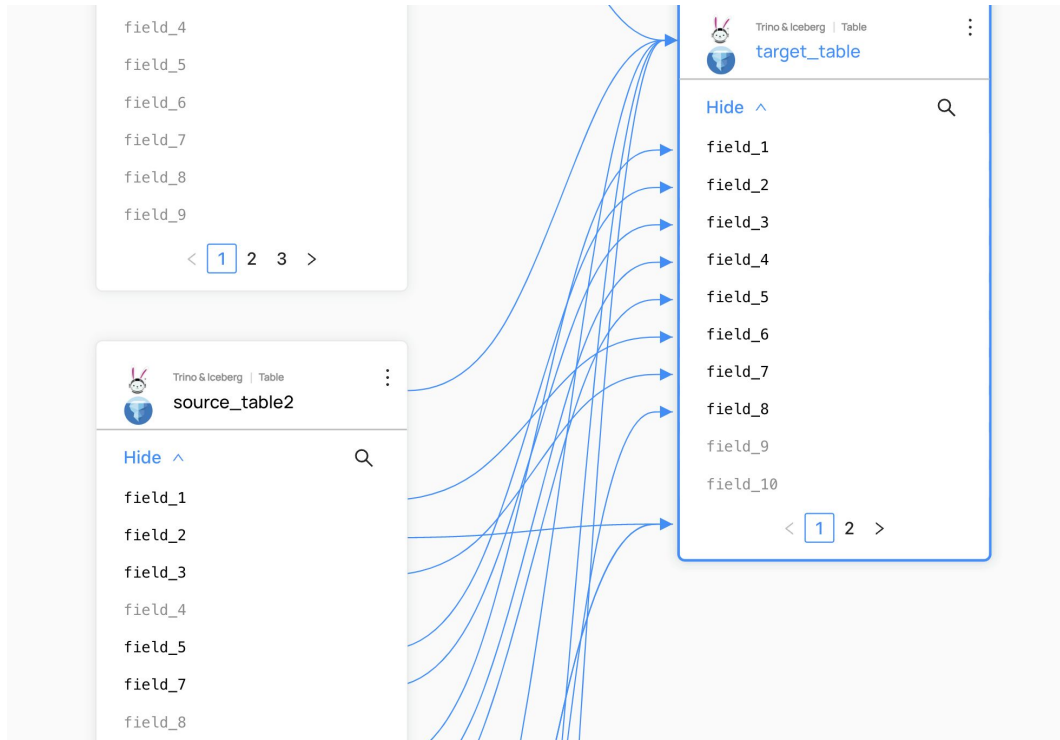
generator.py

```
1 class Generator(metaclass=_Generator):
2     TRANSFORMS: t.Dict[t.Type[exp.Expression], t.Callable[..., str]] = {
3         ...
4         exp.SecurityProperty: lambda self, e: f"SECURITY {self.sql(e, 'this')}",
5     }
6
7     PROPERTIES_LOCATION = {
8         ...
9         exp.SecurityProperty: exp.Properties.Location.POST_SCHEMA,
10    }
```

expressions.py

```
1 class SecurityProperty(Property):
2     arg_types = {"this": True}
```

# Lineage из Trino. Результат



# Решения на основе Data Lineage

01 Предупреждение ломающих DDL изменений



[Kuper-Tech/ddl-breaking-change](#)

Скоро будет код 😊

02 Автоматическое проставление тегов



[Kuper-Tech/datahub-tags-distribution](#)

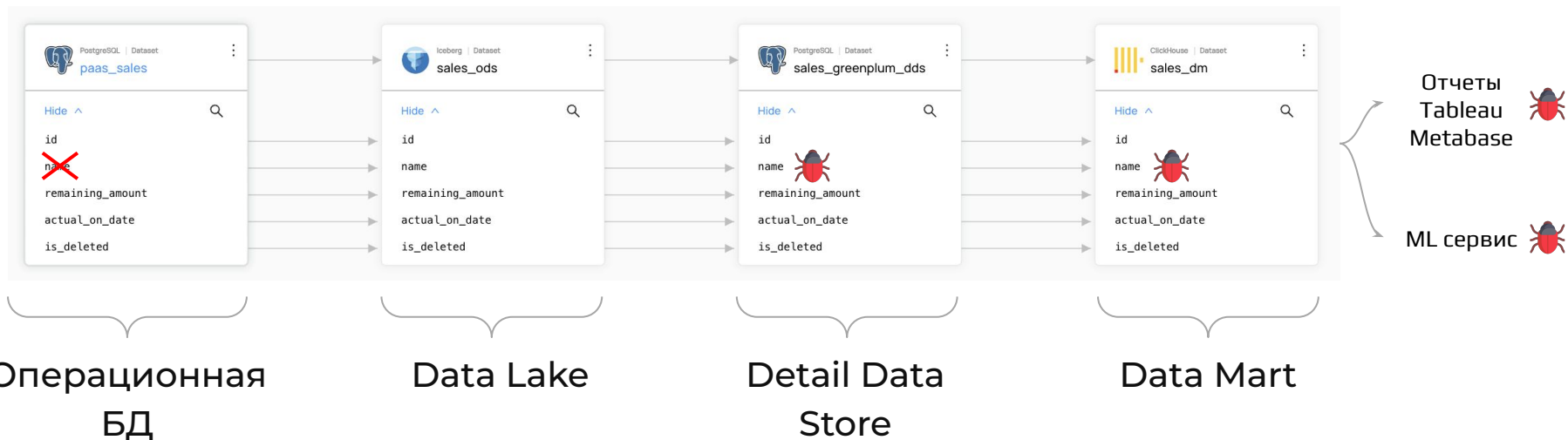
# DDL breaking change. Проблема

01 DROP TABLE ...

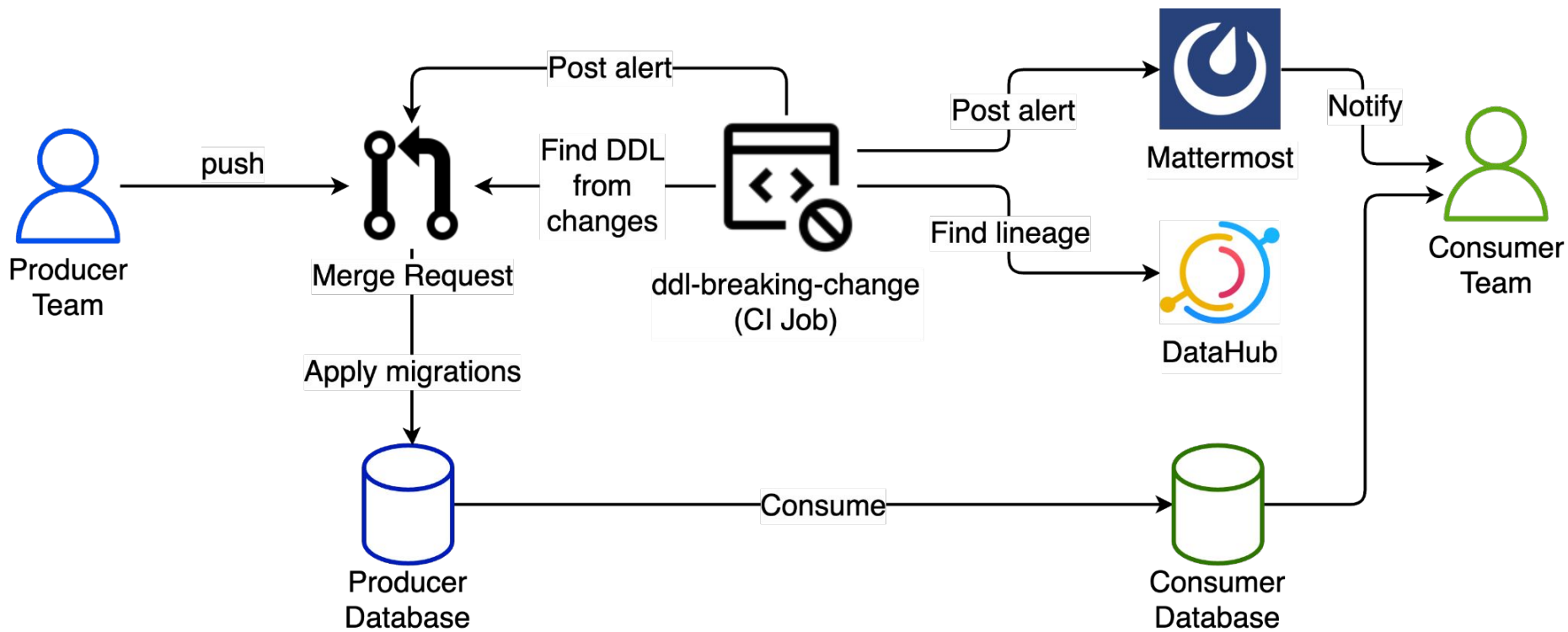
02 ... DROP COLUMN ...

03 ... RENAME COLUMN ...

04 ... ALTER COLUMN ... TYPE



# DDL breaking change. Решение



# Автоматическое проставление тегов



Tier 1, Перс.



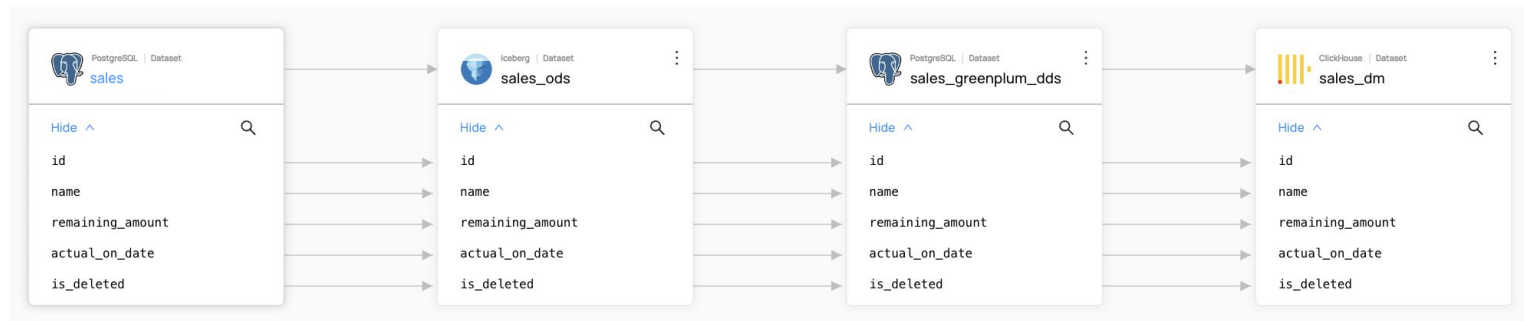
Tier 1, Перс.



Tier 1, Перс.



Tier 1, Перс.



Операционная  
БД

Data Lake

Detail Data  
Store

Data Mart

# Спасибо! Вопросы?

✧ Булат Усманов

✧ Архитектор в Купере

✧ [usmanovbf@gmail.com](mailto:usmanovbf@gmail.com),  
✈ [@usmanovbf](#)

