# ClickHouse для инженеров и архитекторов БД

# Меня хорошо видно && слышно?

# Clickhouse и dbt

## Андрей Поляков

**Senior Java/Groovy Developer/ Data Engineer**

**Об опыте:**

- В отрасли бекенд разработки на java я более 8 лет.
- Занимался fullstack разработкой приложений, разработкой высоко нагруженных compute-grid систем, а также микросервисов и etl-пайплайнов.
- Сейчас в роли старшего разработчика работаю над сервисами платежных систем в Unlimint.

# Правила вебинара

Активно участвуем

Задаем вопрос в чат или голосом

Вопросы вижу в чате, могу ответить не сразу

## Условные обозначения

Индивидуально

Время, необходимое на активность

Пишем в чат

Говорим голосом

Документ

Ответьте себе или задайте вопрос

# Маршрут вебинара

Мотивация: SQL в ETL/ELT

Конфигурация dbt, структура проекта

Специфичная для Clickhouse конфигурация

Подготовка инфраструктуры

Рефлексия

# Цели вебинара

1. Познакомиться с Data Build Tool – мультитул для работы с DWH;

2. Рассмотреть основные возможности и принципы dbt

3. Понять, как инструменты подобные dbt могут помочь инженерам и аналитикам при работе с clickhouse;
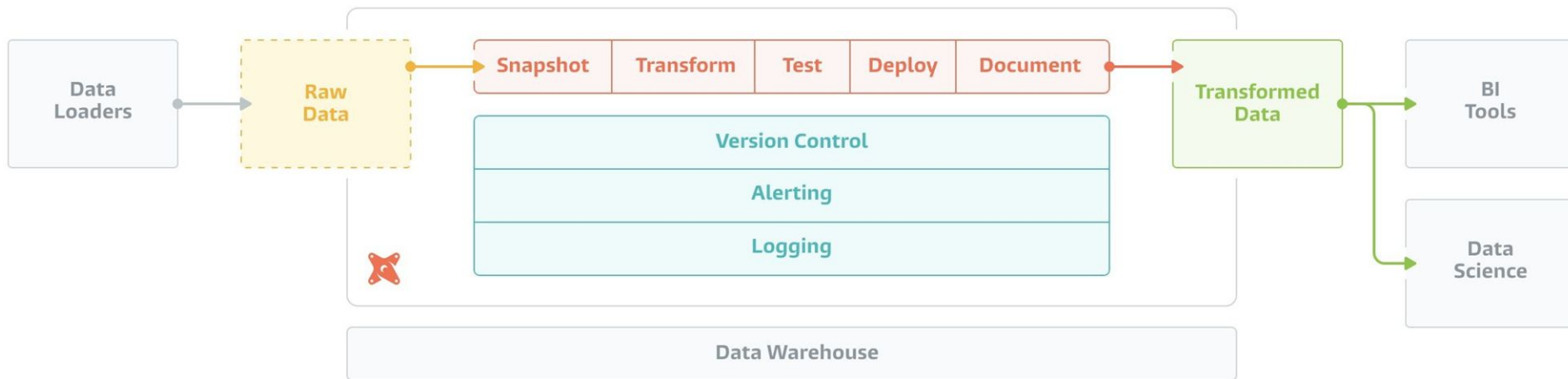
# Data Build Tool (dbt)

# Мотивация: SQL в ETL/ELT

```sql
CREATE PROCEDURE etl_example AS
BEGIN
-- Extract data from the source table
SELECT * INTO #temp_table FROM source_table;
-- Transform data
UPDATE #temp_table
SET column1 = UPPER(column1),
column2 = column2 * 2;
-- Load data into the target table
INSERT INTO target_table
SELECT * FROM #temp_table;
END
```

# Мотивация: ETL-инструменты

```python
from airflow import DAG
from airflow.operators.clickhouse_operator import ClickhouseOperator
from datetime import datetime, timedelta

default_args = {
    'owner': 'me',
    'start_date': datetime(2022, 1, 1),
    'depends_on_past': False,
    'retries': 1,
    'retry_delay': timedelta(minutes=5)}
dag = DAG('simple_dag',
    default_args=default_args,
    schedule_interval=timedelta(hours=1))
task1 = ClickhouseOperator(task_id='get_data',
    sql='select * from table',
    dag=dag)
task2 = ClickhouseOperator(task_id='insert_data',
    sql='insert into table values()',
    retries=3,
    dag=dag)
task1 >> task2
```

# Data Build Tool – T in ELT



Data Loaders → Raw Data → [ Snapshot | Transform | Test | Deploy | Document ] → Transformed Data → BI Tools, Data Science

Version Control
Alerting
Logging

Data Warehouse

# Модели – всё есть SELECT

# SQL + Jinja

Что делает этот код?

```
select
*

from event_tracking.events
{% if target.name == 'dev' %}
where created_at >= dateadd('day', -3, current_date)
{% endif %}
```

⏱️ **Сроки выполнения: 1 мин (пишем в чат)**

# Графы исполнения моделей DAGs

# Node selection syntax

```
# multiple arguments can be provided to --models
$ dbt run --models my_first_model my_second_model

# these arguments can be projects, models, directory paths, tags, or sources
$ dbt run --models tag:nightly my_model finance.base.*

# use methods and intersections for more complex selectors
$ dbt run --models path:marts/finance,tag:nightly,config.materialized:table
```

```
$ dbt run --models my_model+          # select my_model and all children
$ dbt run --models +my_model          # select my_model and all parents
$ dbt run --models +my_model+         # select my_model, and all of its parents and children
```

# Tagging models and running subgraphs

```
1. dbt run-operation stage_external_sources --vars 'ext_full_refresh: true'

2. dbt seed

3. dbt run-operation create_udf

4. dbt run --exclude cars_positions_zones tag:dq

5. dbt snapshot
```

```
1. dbt test --schema --exclude f_chauffeurs_sessions_corrected

2. dbt test --data
```

```
1. dbt run -m tag:dq --full-refresh
```

```yaml
wheely:
    +materialized: view
    staging:
        +schema: staging
        +tags: ["staging"]
        braze:
            +schema: braze
            +tags: ["braze"]
    flatten:
        +schema: flatten
        +materialized: incremental
        +unique_key: _id
        +dist: _id
        +sort: _id
        wheely_prod:
            +tags: ["flatten", "wheely_prod"]
        receipt_prod:
            +tags: ["flatten", "receipt_prod"]
    intermediate:
        +schema: intermediate
        +tags: ["intermediate"]
    marts:
        +tags: ["marts"]
    snapshots:
        +tags: ["snapshots"]
    braze:
        +schema: braze
        +materialized: table
        +tags: ["braze"]
```

# Конфигурация dbt

# dbt project - metadata

```
name: acme corp
profile: acme corp
version: '1.0'

require-dbt-version : ">=0.14.0"

source-paths : ["models"]
analysis-paths : ["analysis"]
test-paths : ["tests"]
data-paths : ["data"]
macro-paths : ["macros"]

target-path : "target"
clean-targets :
    - "target"
    - "dbt_modules"
```

# Настройки подключения

```yaml
acme corp:
  outputs:
    dev:
      type: postgres
      threads: 8
      host: [hostname]
      user: [username]
      pass: [password]
      port: 5439
      dbname: [database name]
      schema: dbt_[username] # e.g. dbt_alice
  target: dev
```

# Полная и инкрементальная загрузка

```
1   {{
2       config (
3       materialized='incremental',
4       sql_where='true',
5       unique_key='id',
6           dist="call_id",
7           sort="min_event_ts_msk",
8       )
9   }}
```

# CTE

```
WITH cte_name (column1, column2, ..., columnN) AS ( ❶
    -- Query definition goes here                    ❷
)
SELECT column1, column2, ..., columnN                ❸
FROM cte_name                                        ❸
-- Additional query operations go here               ❹
```

Что такое CTE? Для чего они нужны?

# CTE

## Без CTE

```sql
SELECT pb.book_id,
       pb.title,
       pb.author,
       s.total_sales
FROM (
    SELECT book_id,
           title,
           author
    FROM books
    WHERE rating >= 4.6
) AS pb
JOIN sales s ON pb.book_id = s.book_id
WHERE s.year = 2022
ORDER BY s.total_sales DESC
LIMIT 5;
```

## C CTE

```sql
WITH popular_books AS (
    SELECT book_id,
           title,
           author
    FROM books
    WHERE rating >= 4.6
),
best_sellers AS (
    SELECT pb.book_id,
           pb.title,
           pb.author,
           s.total_sales
    FROM popular_books pb
    JOIN sales s ON pb.book_id = s.book_id
    WHERE s.year = 2022
    ORDER BY s.total_sales DESC
    LIMIT 5
)
SELECT *
FROM best_sellers;
```

# Модели Stage

```sql
/* This should be file stg_books.sql, and it queries the raw table to create
the new model */

SELECT
    book_id,
    title,
    author,
    publication_year,
    genre
FROM
    raw_books
```

# Модели Intermediate

```sql
-- This should be file int_book_authors.sql

-- Reference the staging models
WITH
  books AS (
    SELECT *
    FROM {{ ref('stg_books') }}
  ),
  authors AS (
    SELECT *
    FROM {{ ref('stg_authors') }}
  )

-- Combine the relevant information
SELECT
  b.book_id,
  b.title,
  a.author_id,
  a.author_name
FROM
  books b
JOIN
  authors a ON b.author_id = a.author_id
```

# Модели Mart

```sql
-- This should be file mart_book_authors.sql

{{
  config(
    materialized='table',
    unique_key='author_id',
    sort='author_id'
  )
}}

WITH book_counts AS (
  SELECT
    author_id,
    COUNT(*) AS total_books
  FROM {{ ref('int_book_authors') }}
  GROUP BY author_id
)
SELECT
  author_id,
  total_books
FROM book_counts
```

# Параметры dbt-clickhouse

# Database ClickHouse configurations | dbt Developer Hub

## Models

| Type | Supported? | Details |
|---|---|---|
| view materialization | YES | Creates a view. |
| table materialization | YES | Creates a table. See below for the list of supported engines. |
| incremental materialization | YES | Creates a table if it doesn't exist, and then writes only updates to it. |
| ephemeral materialized | YES | Creates a ephemeral/CTE materialization. This does model is internal to dbt and does not create any database objects |

## Experimental models

The following are experimental features in Clickhouse:

| Type | Supported? | Details |
|---|---|---|
| Materialized View materialization | YES, Experimental | Creates a materialized view. |
| Distributed table materialization | YES, Experimental | Creates a distributed table. |
| Distributed incremental materialization | YES, Experimental | Incremental model based on the same idea as distributed table. Note that not all strategies are supported, visit this for more info. |
| Dictionary materialization | YES, Experimental | Creates a dictionary. |

# Table configurations

| Option | Description | Required? |
|---|---|---|
| `materialized` | How the model will be materialized into ClickHouse. Must be `table` to create a table model. | Required |
| `engine` | The table engine to use when creating tables. See list of supported engines below. | Optional (default: `MergeTree()`) |
| `order_by` | A tuple of column names or arbitrary expressions. This allows you to create a small sparse index that helps find data faster. | Optional (default: `tuple()`) |
| `partition_by` | A partition is a logical combination of records in a table by a specified criterion. The partition key can be any expression from the table columns. | Optional |

# Incremental table configurations

| Option | Description | Required? |
|---|---|---|
| `materialized` | How the model will be materialized into ClickHouse. Must be `table` to create a table model. | Required |
| `unique_key` | A tuple of column names that uniquely identify rows. For more details on uniqueness constraints, see here. | Required. If not provided altered rows will be added twice to the incremental table. |
| `engine` | The table engine to use when creating tables. See list of supported engines below. | Optional (default: `MergeTree()` ) |
| `order_by` | A tuple of column names or arbitrary expressions. This allows you to create a small sparse index that helps find data faster. | Optional (default: `tuple()` ) |
| `partition_by` | A partition is a logical combination of records in a table by a specified criterion. The partition key can be any expression from the table columns. | Optional |
| `inserts_only` | (Deprecated, see the `append` materialization strategy). If True, incremental updates will be inserted directly to the target incremental table without creating an intermediate table. | Optional (default: `False` ) |
| `incremental_strategy` | The strategy to use for incremental materialization. `delete+insert` , `append` and `insert_overwrite` (experimental) are supported. For additional details on strategies, see here | Optional (default: 'default') |
| `incremental_predicates` | Incremental predicate clause to be applied to `delete+insert` materializations | Optional |

# Table engines

## Supported table engines

| Type | Details |
|------|---------|
| MergeTree (default) | https://clickhouse.com/docs/en/engines/table-engines/mergetree-family/mergetree/. |
| HDFS | https://clickhouse.com/docs/en/engines/table-engines/integrations/hdfs |
| MaterializedPostgreSQL | https://clickhouse.com/docs/en/engines/table-engines/integrations/materialized-postgresql |
| S3 | https://clickhouse.com/docs/en/engines/table-engines/integrations/s3 |
| EmbeddedRocksDB | https://clickhouse.com/docs/en/engines/table-engines/integrations/embedded-rocksdb |
| Hive | https://clickhouse.com/docs/en/engines/table-engines/integrations/hive |

## Experimental supported table engines

| Type | Details |
|------|---------|
| Distributed Table | https://clickhouse.com/docs/en/engines/table-engines/special/distributed. |
| Dictionary | https://clickhouse.com/docs/en/engines/table-engines/special/dictionary |

# Вопросы?

Ставим "+",
если вопросы есть

Ставим "–",
если вопросов нет

# Список материалов для изучения

1. [dbt Getting Started Tutorial](#)
2. [dbt Documentation](#)
3. [dbt FAQ](#)
4. [How we structure our dbt projects](#)
5. [The Modern Data Stack: Past, Present, and Future](#)
6. [Five principles that will keep your data warehouse organized](#)
7. [The Analytics Engineering Guide](#)

Делитесь своими материалами в telegram

# Рефлексия

# Рефлексия

С какими впечатлениями уходите с вебинара?

Как будете применять на практике то,
что узнали на вебинаре?

# Приходите на следующие вебинары

**Алексей Железной**

***Senior Data Engineer*** *в Wildberries*
*Магистратура - ФКН ВШЭ*

*Руководитель курсов **DWH Analyst, ClickHouse для инженеров и архитекторов БД** в OTUS*

*Преподаватель курсов **Data Engineer, DWH Analyst, PostgreSQL** и пр. в OTUS*

LinkedIn