

# Визуализация для ELT-процессов в DWH

Как мы готовили у себя DBT и какие  
ошибки допустили.

# Содержание

1. О нас
2. Инструмент в помощь доставке данных
3. Зачем тут DBT
4. Что такое DBT
5. Наши особенности



# СберМаркет – сервис доставки продуктов и товаров



> 32k

Магазинов

160

Ретейлеров

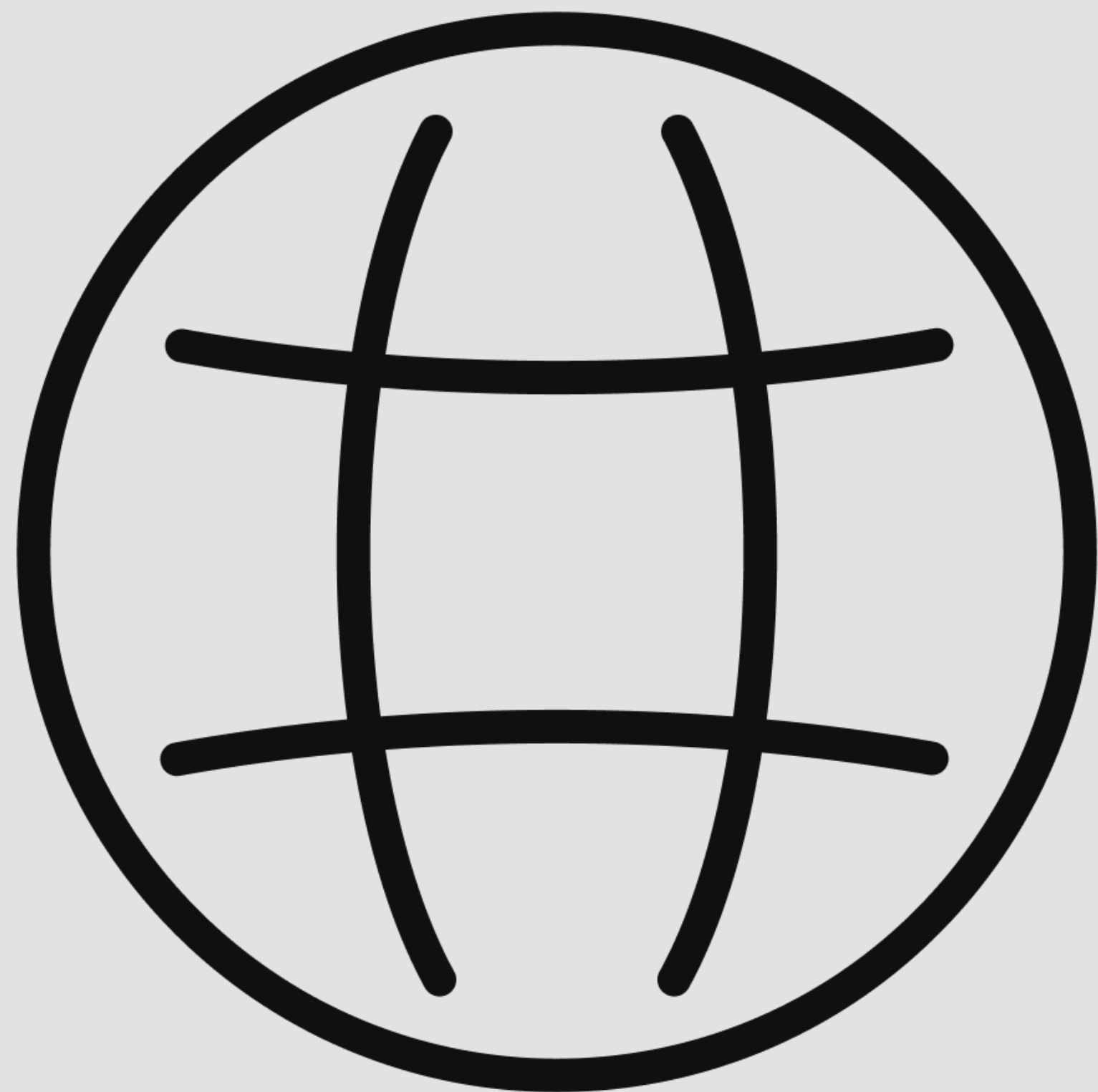
> 160

Городов

> 200kk

Количество SKU  
(товарных позиций)

# О данных



> 200T6

DataLake

> 150

Пользователей GP

> 10T6

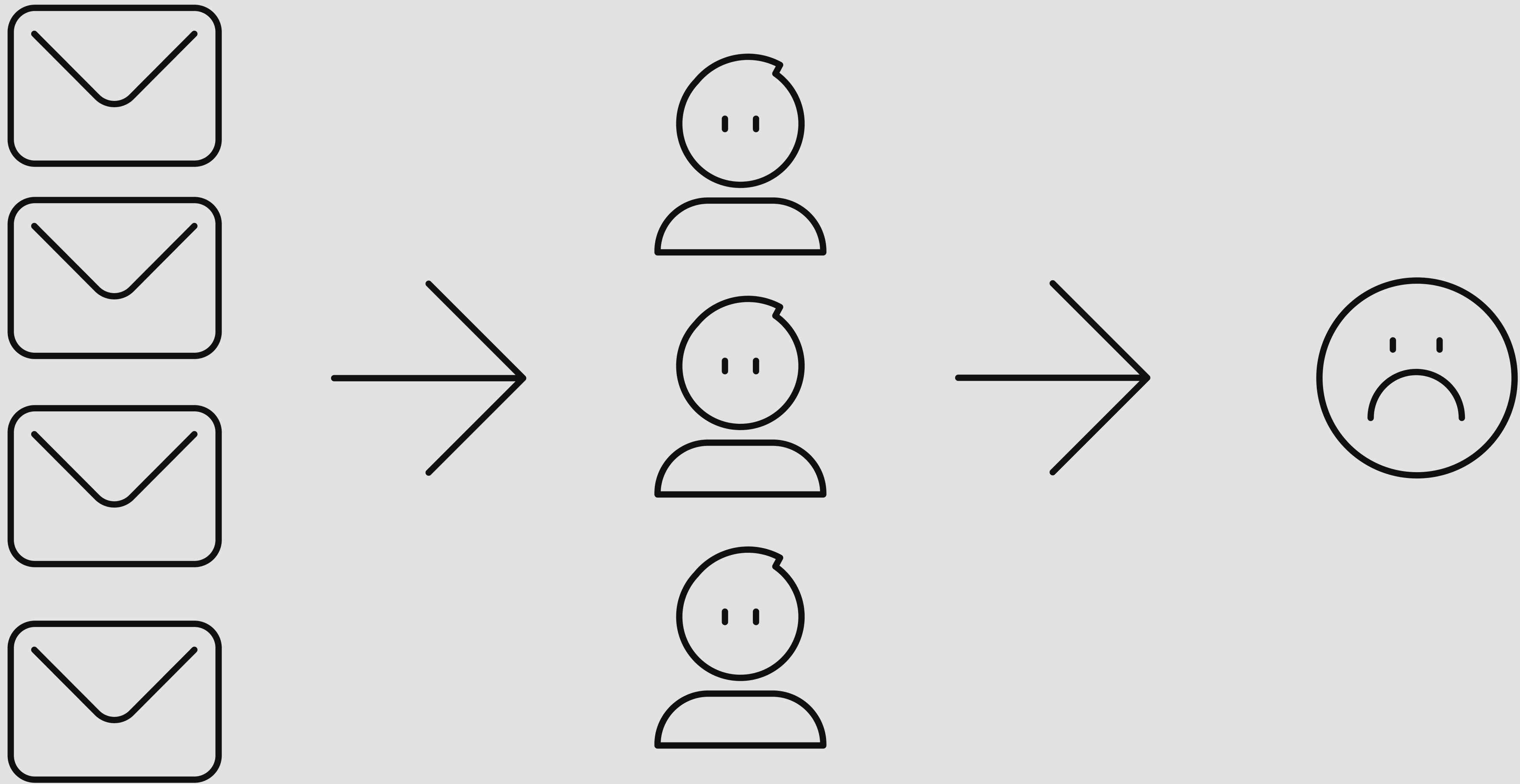
Хранилище

> 100

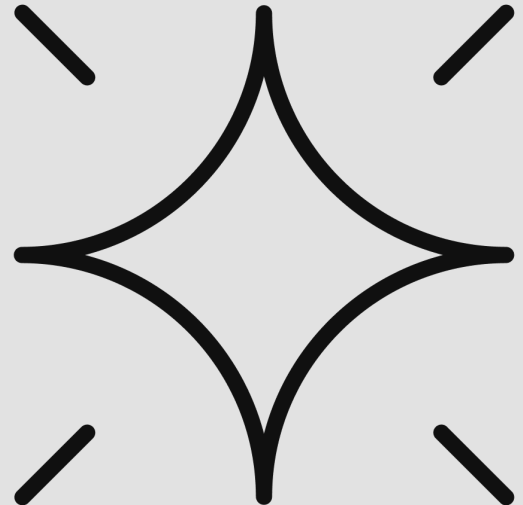
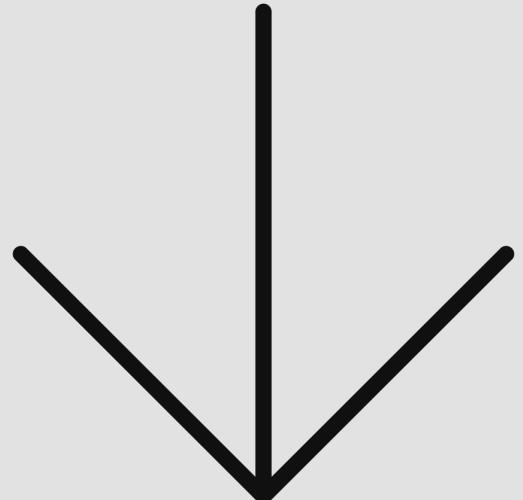
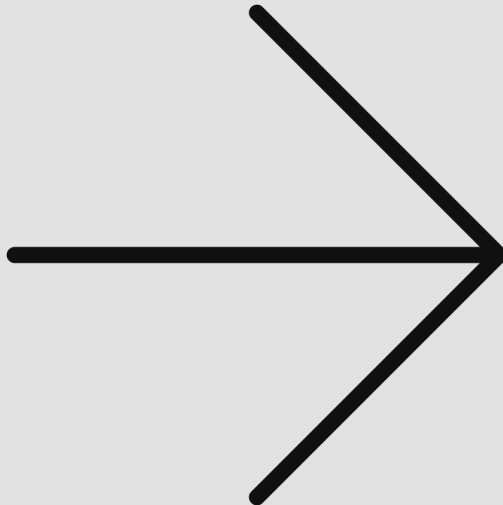
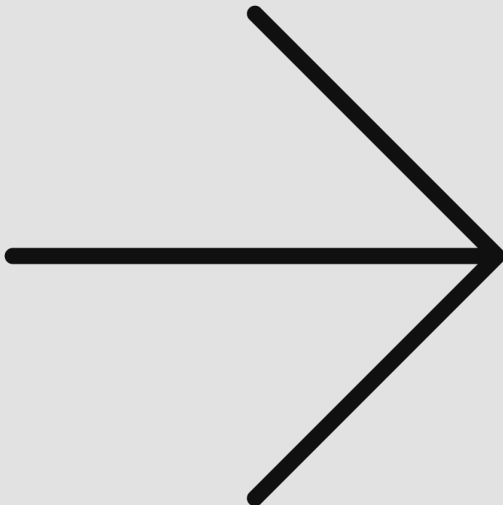
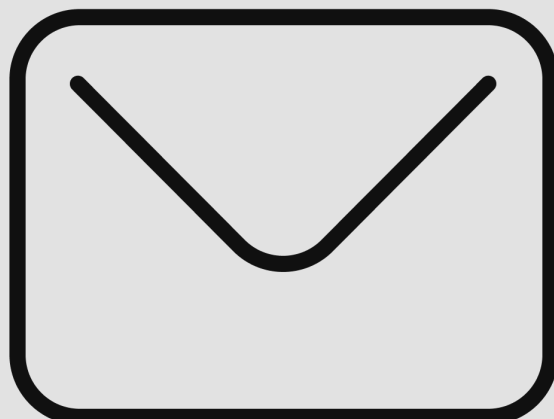
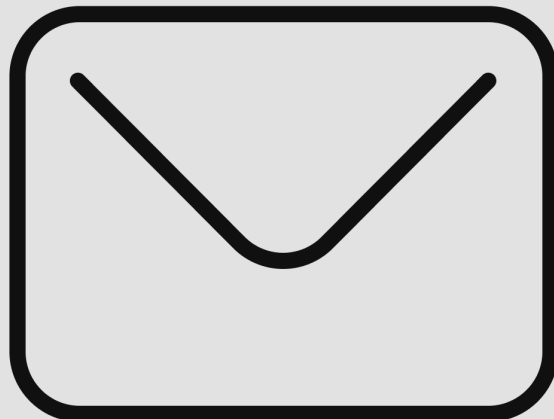
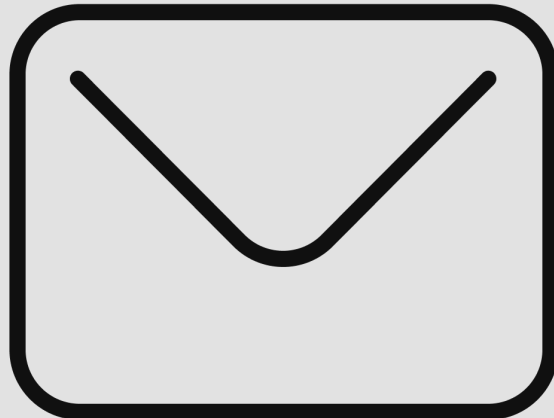
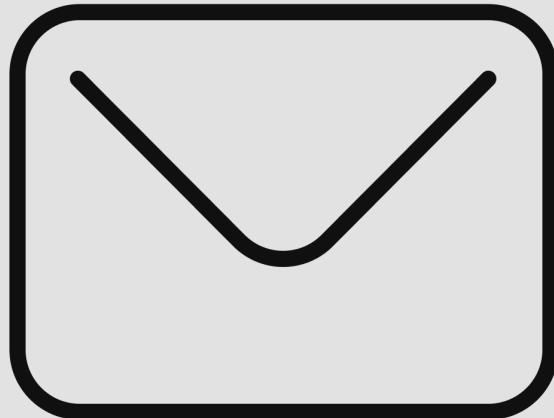
Витрин




# Если бы инженер сам доставлял витрину



# Предоставить инструмент



 Airflow

DAGs

Security ▾

Browse ▾

Admin ▾

Docs ▾

08:24 UTC ▾

AA ▾

Deleting DAG with id dwh\_gp\_preaggregate\_receipt. May take a couple minutes to fully disappear.

×

DAGs

All 0

Active 0

Paused 0

Filter DAGs by tag

Search DAGs

<div>ⓘ</div> DAG	Owner	Runs ⓘ	Schedule	Last Run ⓘ	Next Run ⓘ	Recent Tasks ⓘ	Actions	Links
No results								

« < > »

Showing 0-0 of 0 DAGs

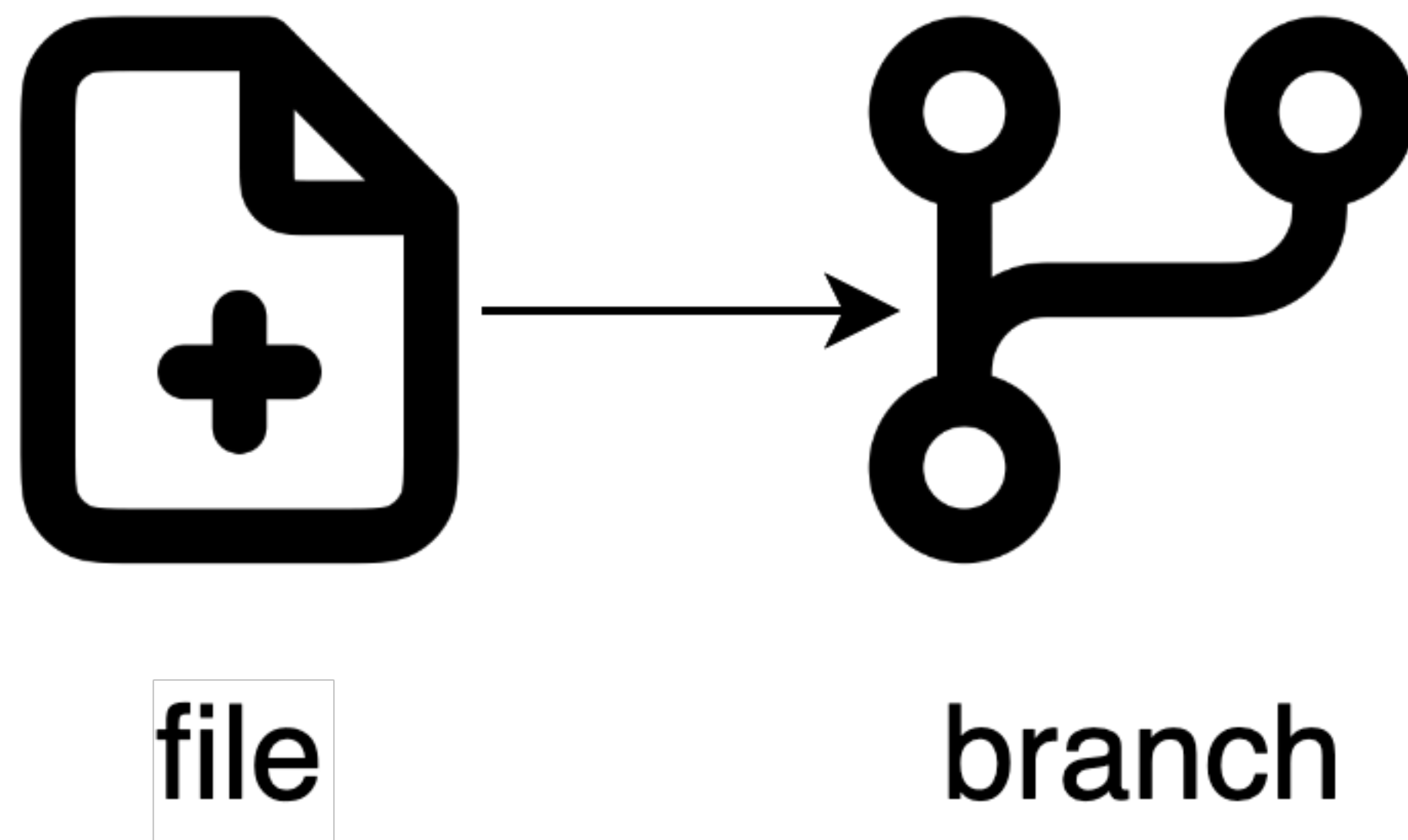


# Границы применения

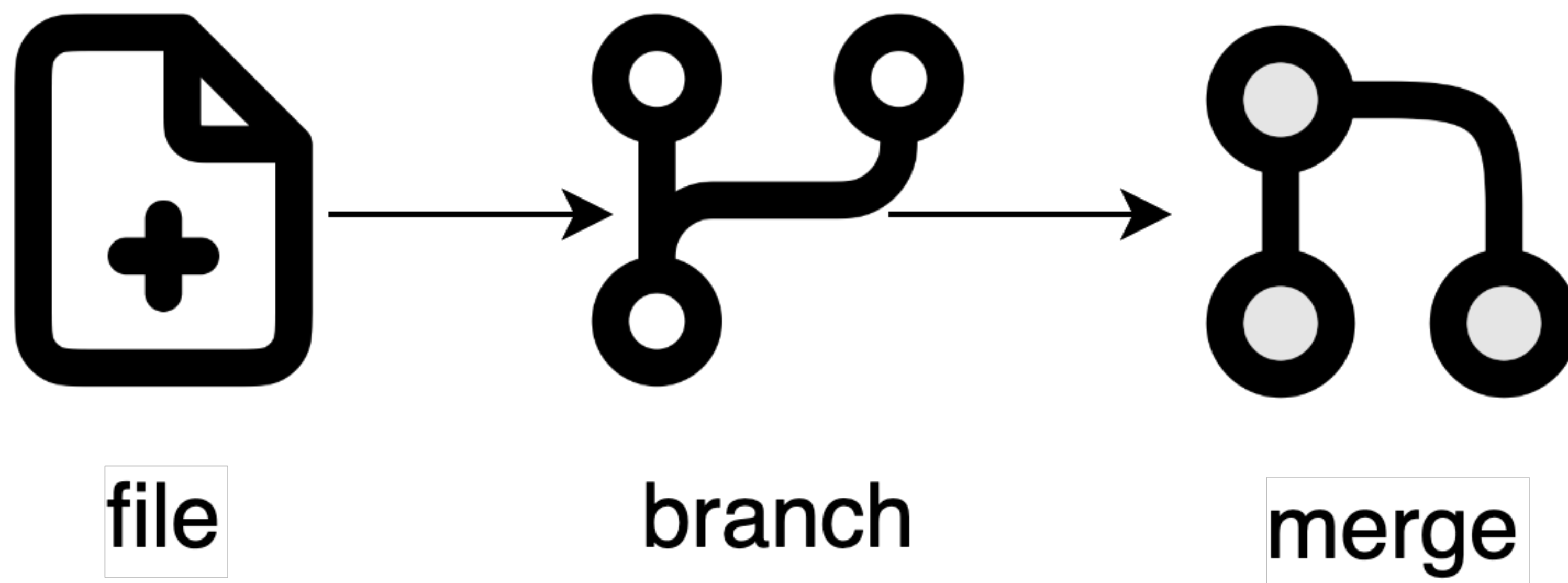
Не подойдет для

- Наполнение сырых данных
- Не хватает SQL

# Наше решение до DBT

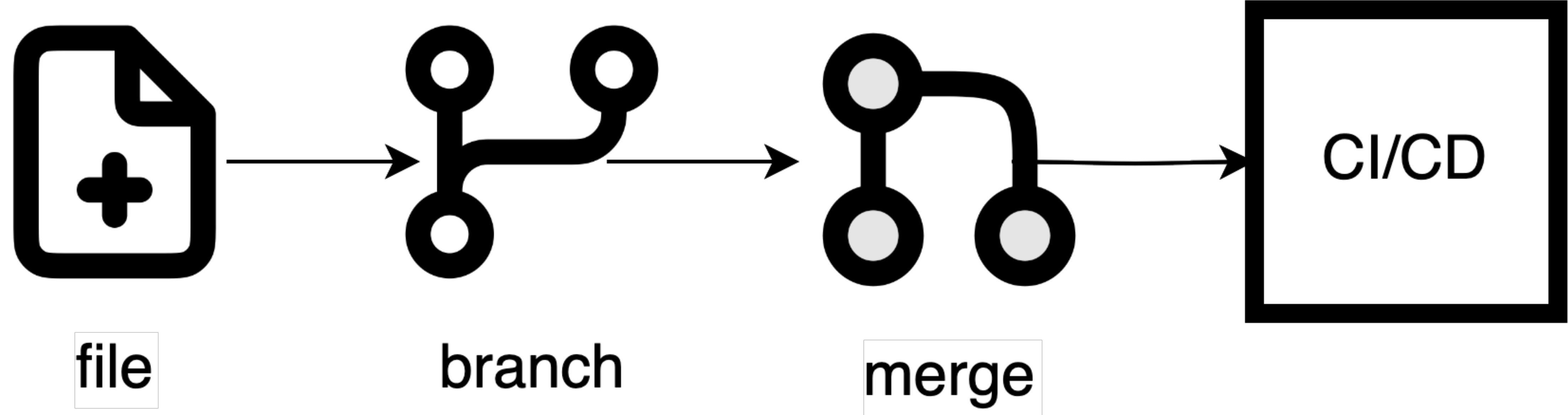


# Наше решение до DBT

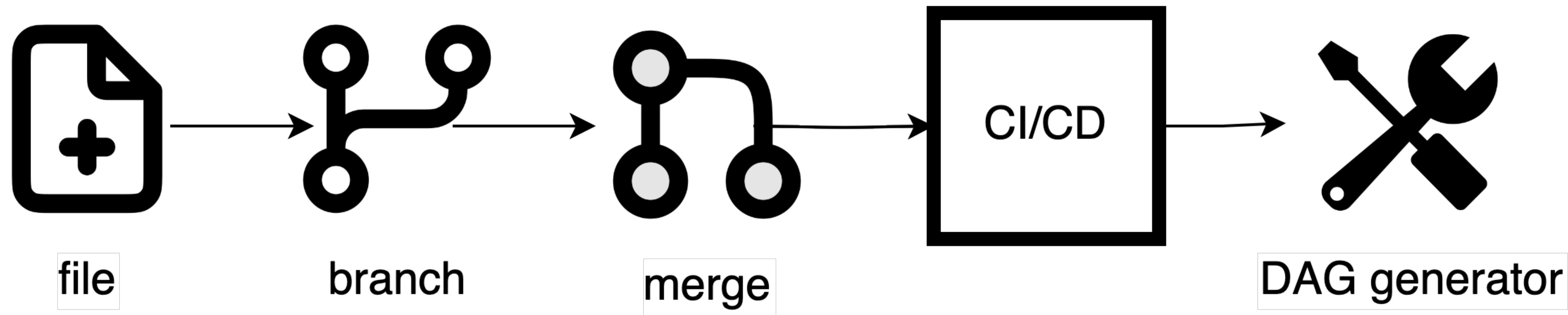




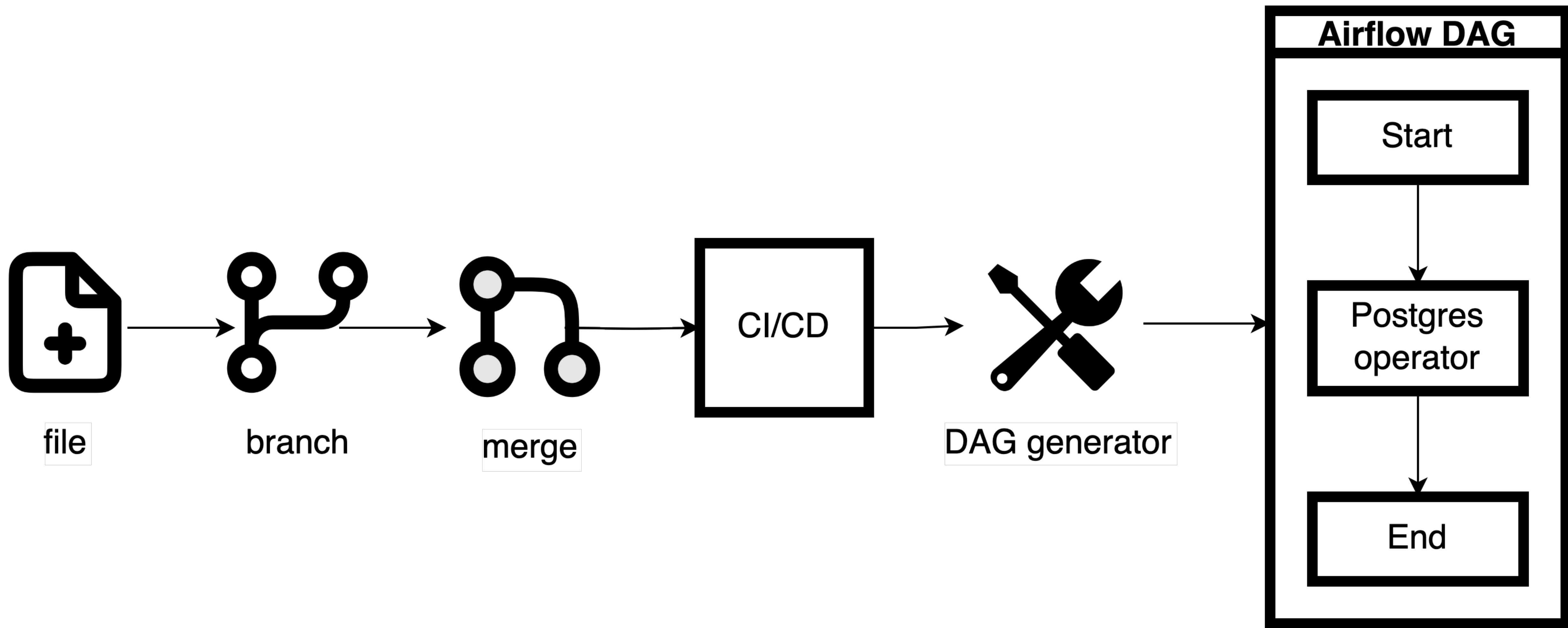
# Наше решение до DBT



# Наше решение до DBT

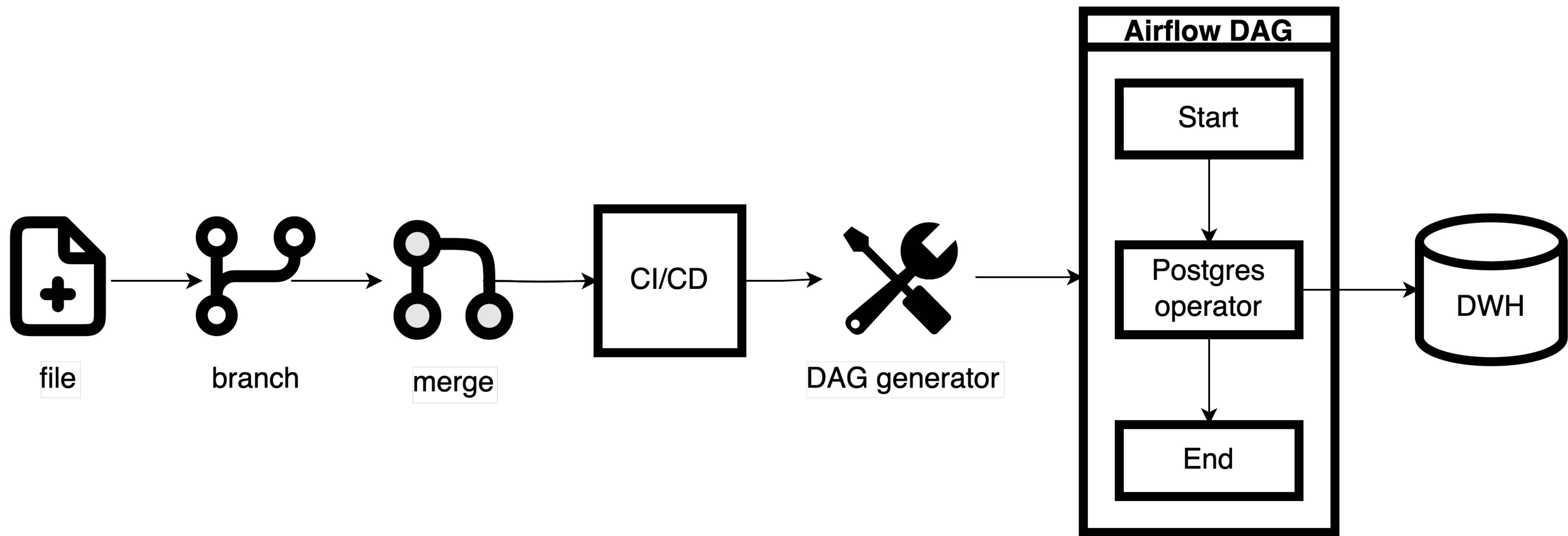


# Наше решение до DBT

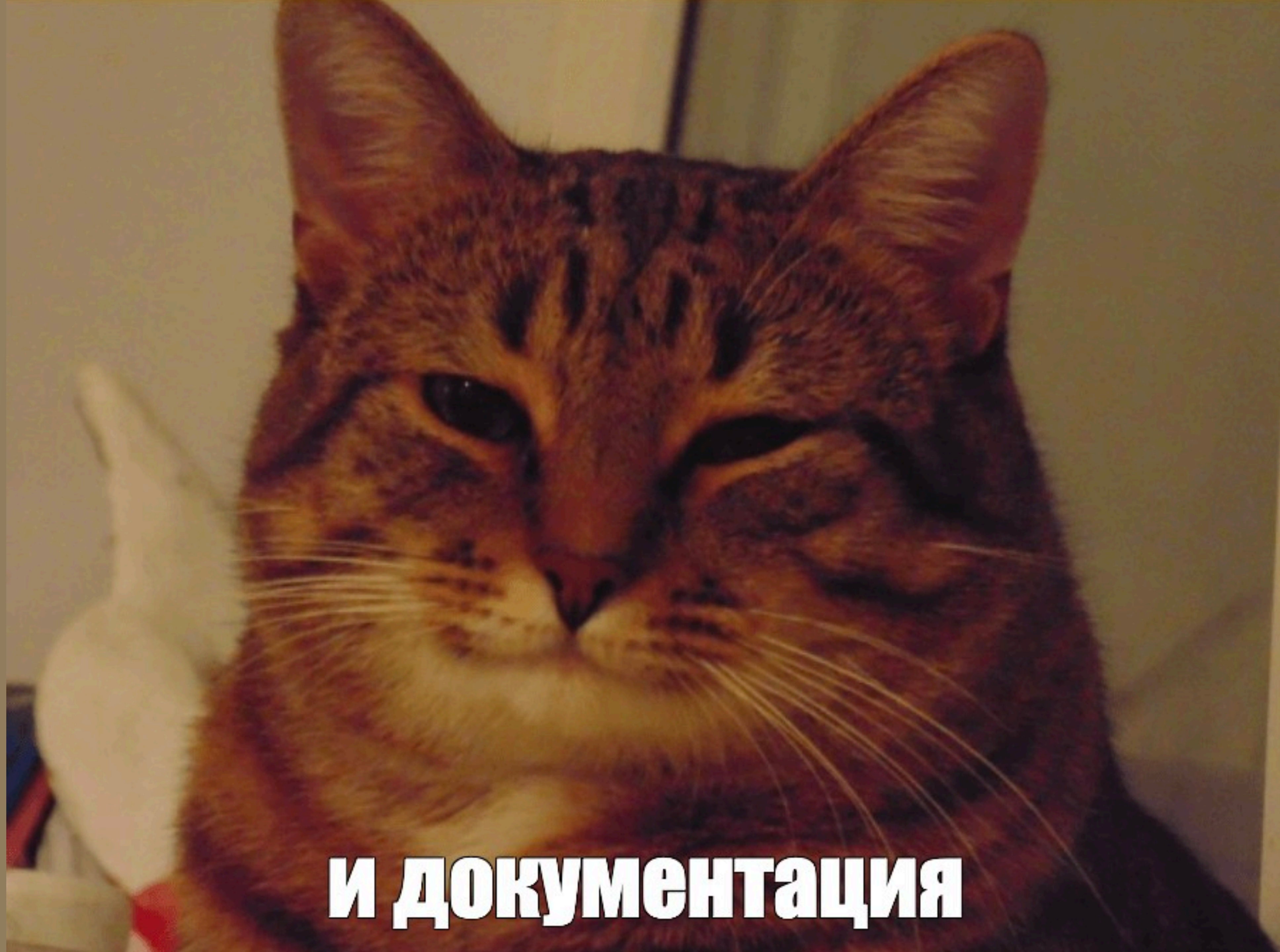




# Наше решение до DBT







**И ДОКУМЕНТАЦИЯ**





**ВЗЯТЬ ГОТОВОЕ**



**ПИСАТЬ  
САМИМ**



# Data build tool

- DBT — это инструмент командной строки с открытым исходным кодом. Призван помочь более эффективно преобразовывать данные в хранилищах данных. В акрониме ELT (Extract, Load, Transform) DBT занимает этап T.
- Целью dbt является предоставление аналитикам возможности работать как инженеры-программисты



# Чем привлек DBT

- Распространенность, развитие, обширное сообщество



# Чем привлек DBT

- Распространенность, развитие, обширное сообщества
- Конфигурация на .yaml



# Чем привлек DBT

- Распространенность, развитие, обширное сообщество
- Конфигурация на .yaml
- Возможность добавления тестов



# Чем привлек DBT

- Распространенность, развитие, обширное сообщества
- Конфигурация на .yaml
- Возможность добавления тестов
- Логика не в ХП



# Чем привлек DBT

- Распространенность, развитие, обширное сообщества
- Конфигурация на .yaml
- Возможность добавления тестов
- Логика не в ХП
- Метаданные



# Чем привлек DBT

- Распространенность, развитие, обширное сообщества
- Конфигурация на .yaml
- Возможность добавления тестов
- Логика не в ХП
- Метаданные
- Графы зависимостей



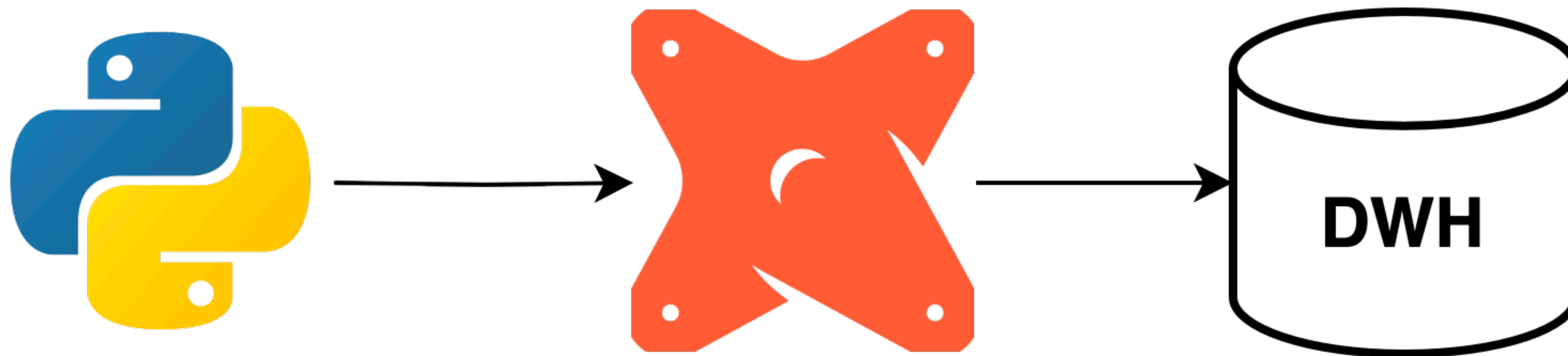


# Чем привлек DBT

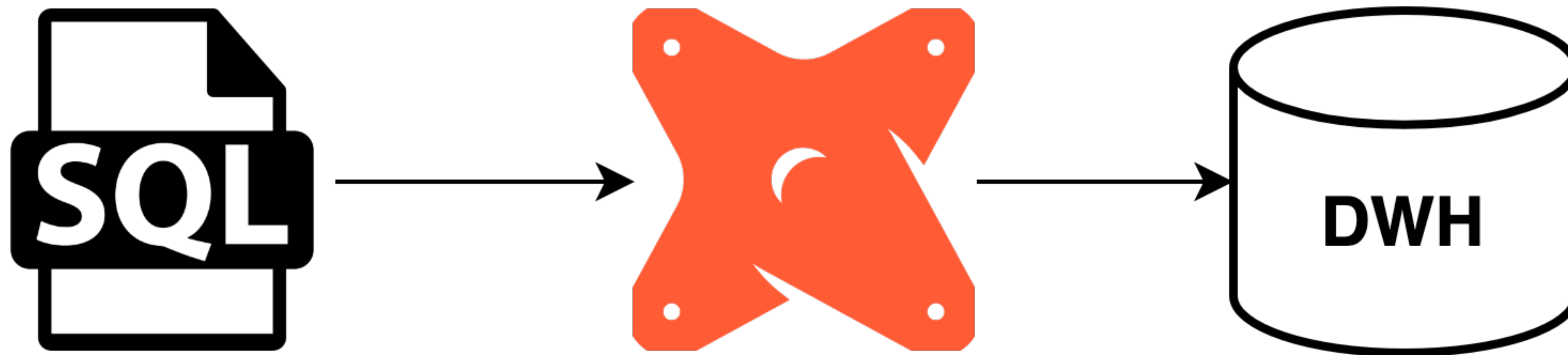
- Распространенность, развитие, обширное сообщества
- Конфигурация на .yaml
- Возможность добавления тестов
- Логика не в ХП
- Метаданные
- Графы зависимостей
- Визуализация и документация



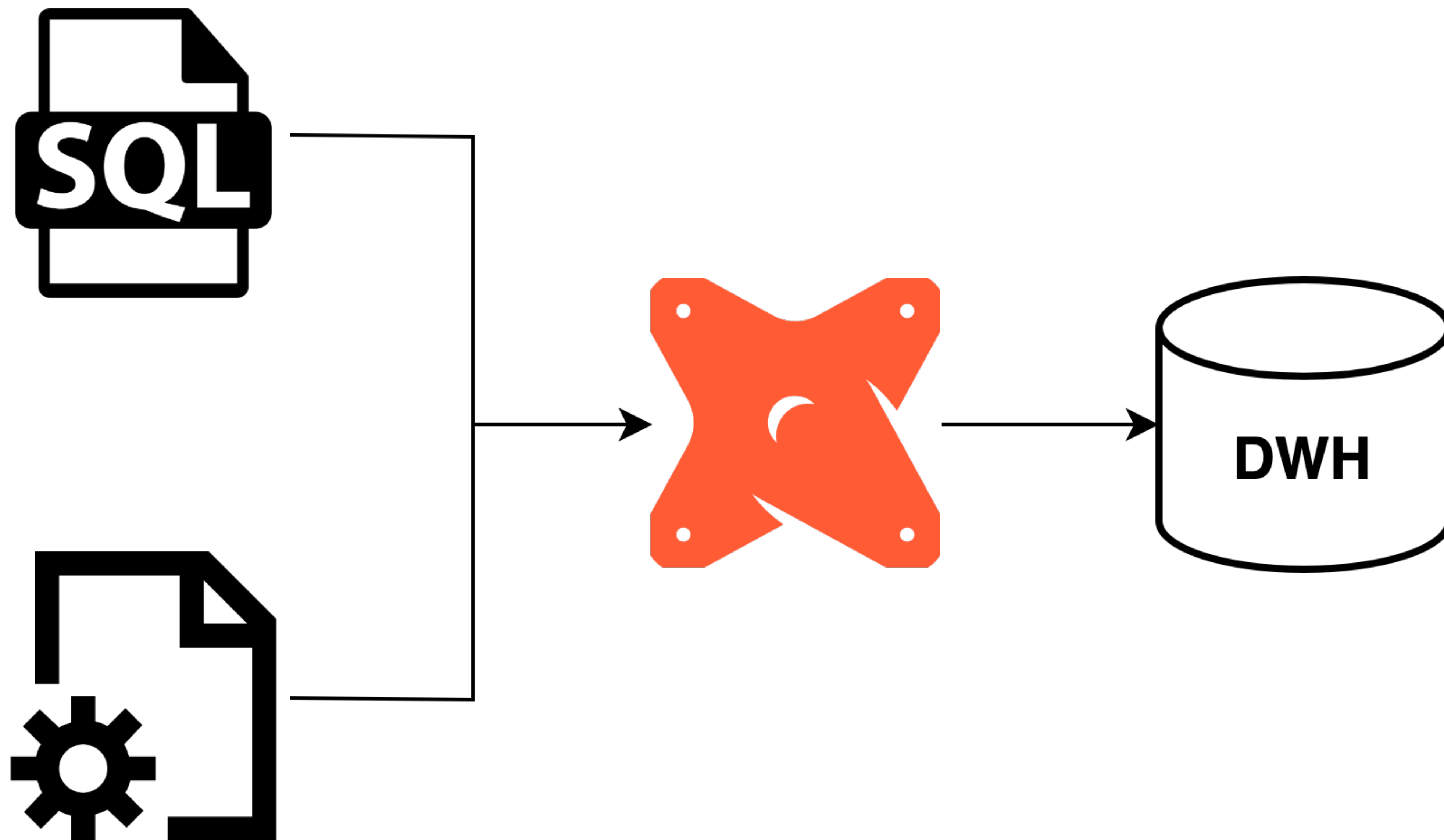
# Как DBT работает



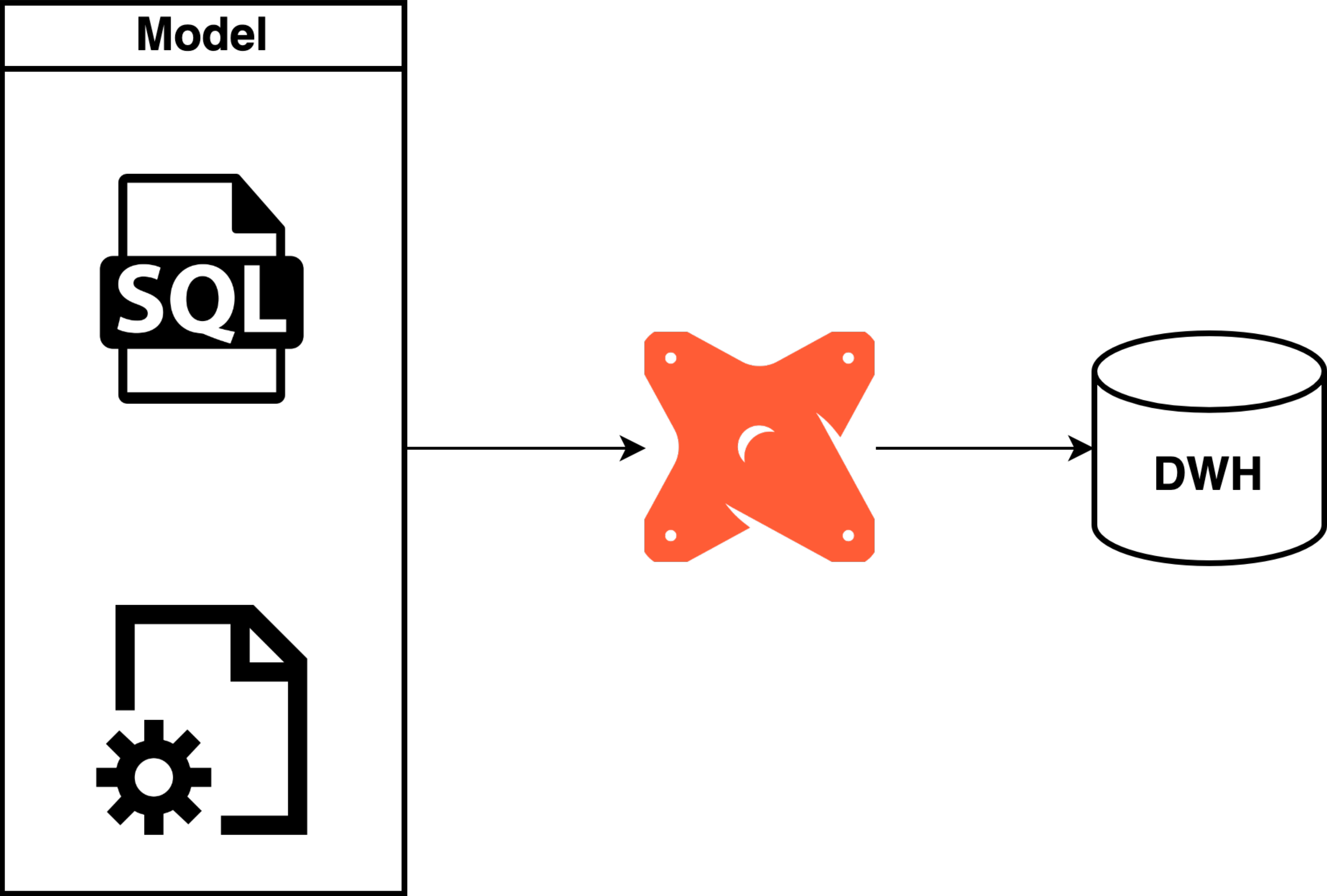
# Как DBT работает



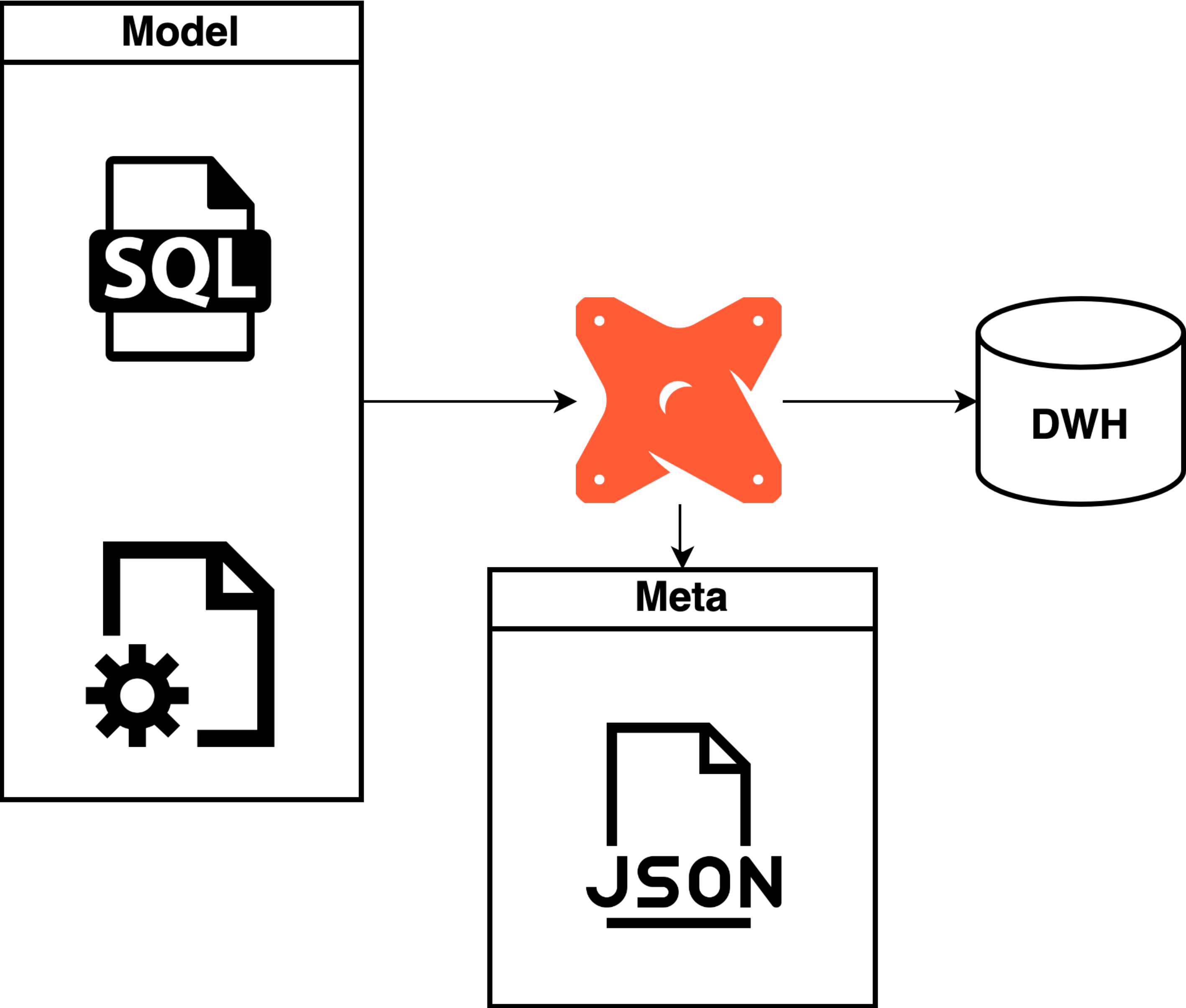
# Как DBT работает



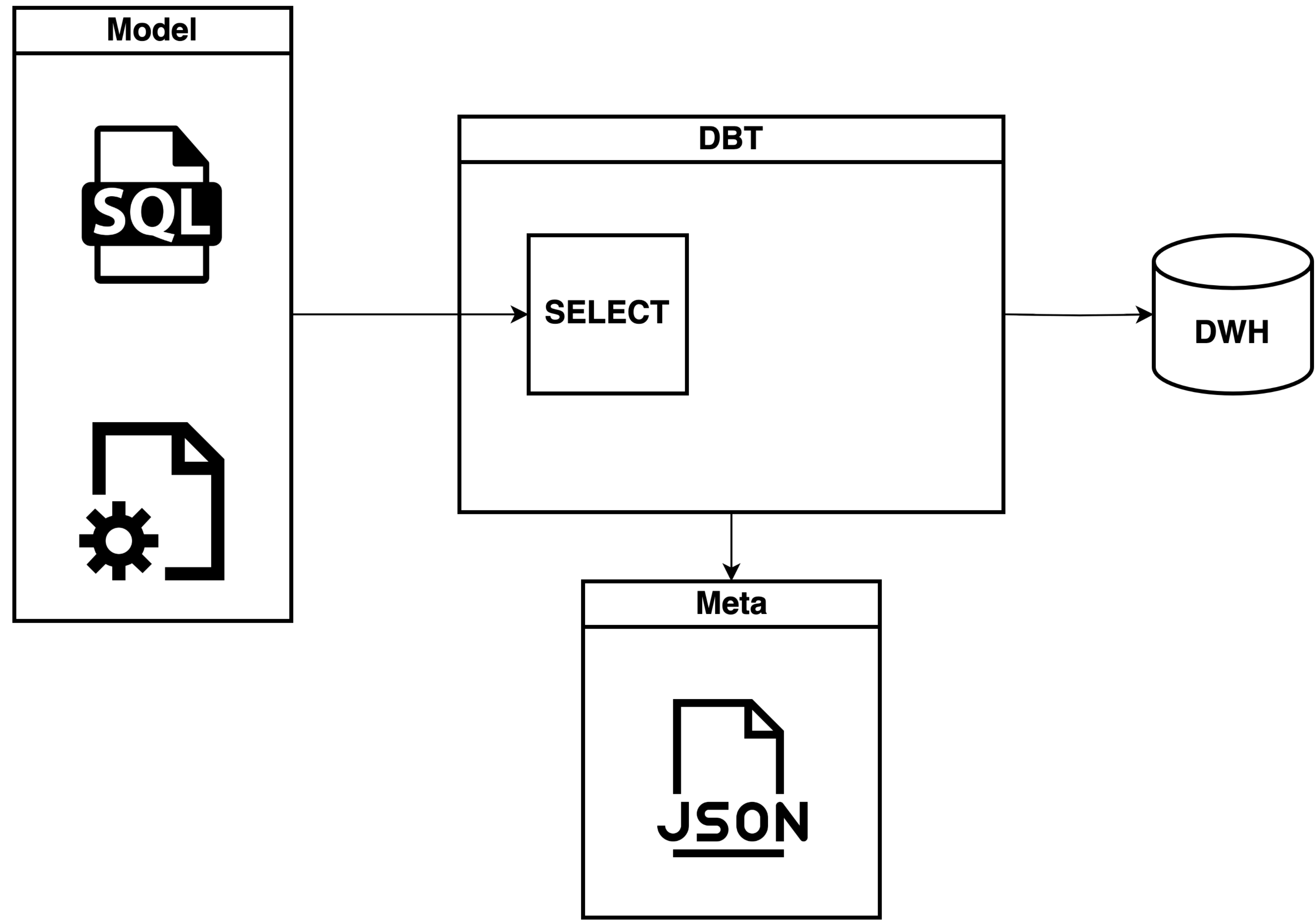
# Как DBT работает



# Как DBT работает

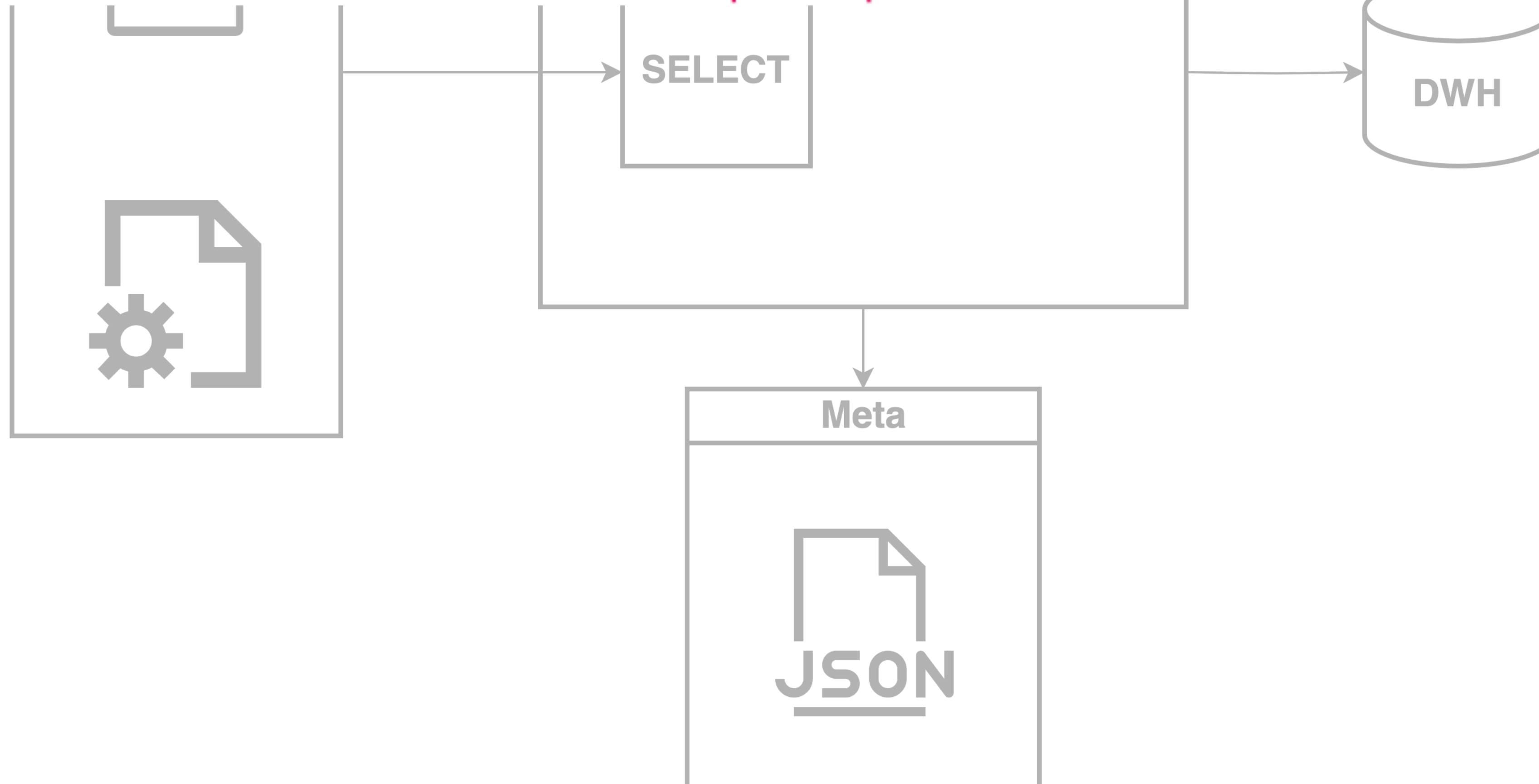


# Как DBT работает





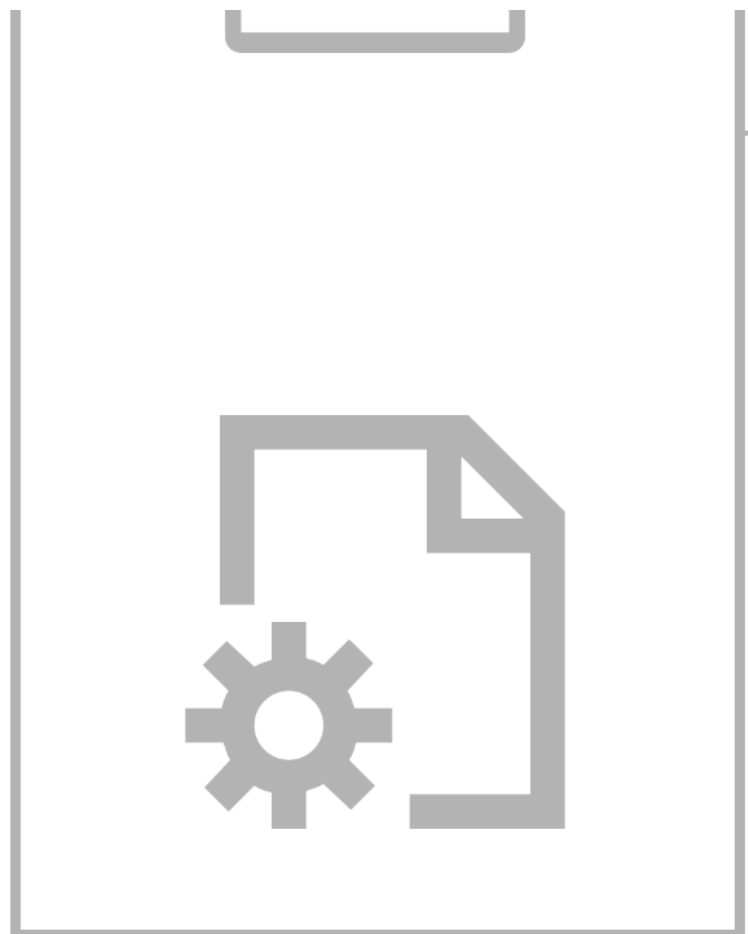
```
join {{ ref( 's_delivery_window' ) }} pwin  
on sdwch.previous_delivery_window_id = pwin.  
{{ actual_rows_by_source('pwin', 'INST') }}  
left join {{ source('meta', 'all_layers_params')  
on not_sm_reason.layer_cd = 'dds'  
and not_sm_reason.table_name = 'p_shipment'
```



```

join {{ ref( 's_delivery_window' ) }} pwin
on sdwch.previous_delivery_window_id = pwin.
{{ actual_rows_by_source('pwin', 'INST') }}
left join {{ source('meta', 'all_layers_params') }}
on not_sm_reason.layer_cd = 'dds'
and not_sm_reason.table_name = 'p_shipment'

```



SELECT

```

join "dwh"."dds"."s_delivery_window" pwin

```

```

on sdwch.previous_delivery_window_id = pwin.d

```

```

and pwin.mt_src_cd = 'INST'

```

```

and not pwin.mt_deleted_flg

```

```

and pwin.mt_active_flg

```

```

left join "dwh"."meta"."all_layers_params" not_sm_

```

```

on not_sm_reason.layer_cd = 'dds'

```

```

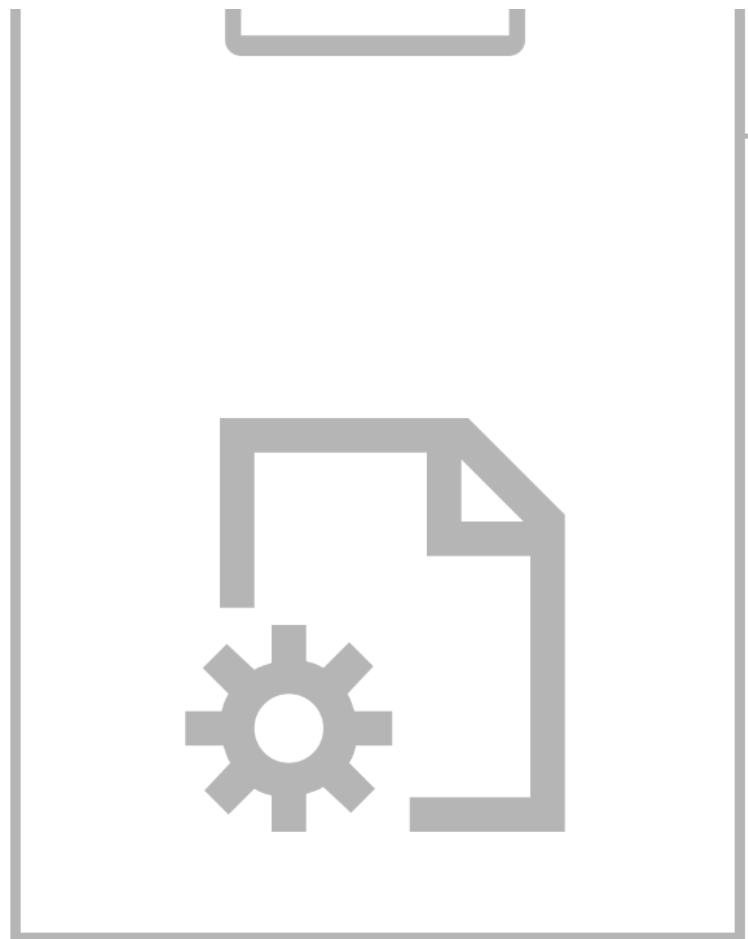
and not_sm_reason.table_name = 'p_shipment'

```

```

join {{ ref( 's_delivery_window' ) }} pwin
on sdwch.previous_delivery_window_id = pwin.
{{ actual_rows_by_source('pwin', 'INST') }}
left join {{ source('meta', 'all_layers_params') }}
on not_sm_reason.layer_cd = 'dds'
and not_sm_reason.table_name = 'p_shipment'

```



SELECT

```

join "dwh"."dds"."s_delivery_window" pwin
on sdwch.previous_delivery_window_id = pwin.d

```



```

and pwin.mt_src_cd = 'INST'
and not pwin.mt_deleted_flg
and pwin.mt_active_flg

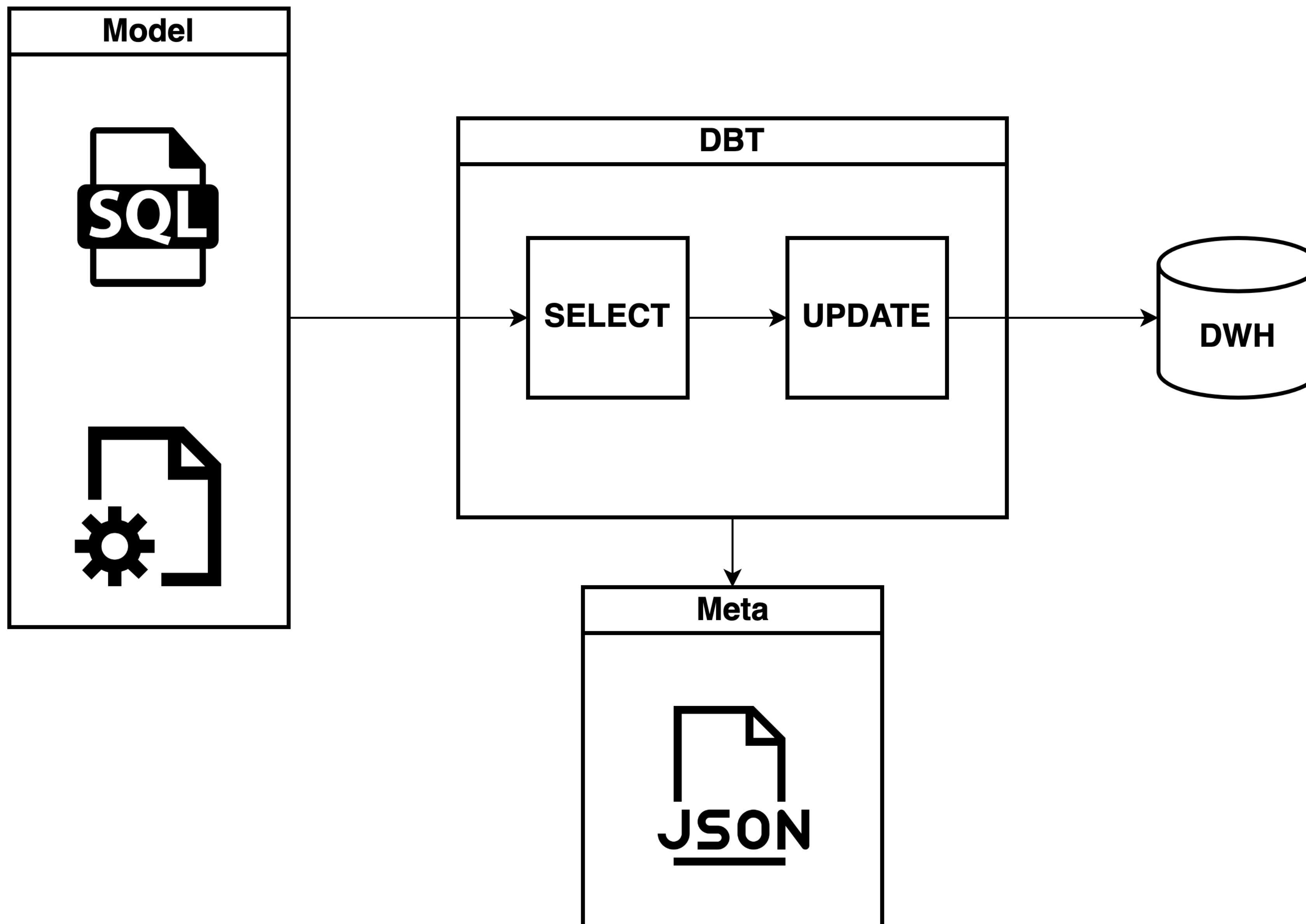
```

```

left join "dwh"."meta"."all_layers_params" not_sm_
on not_sm_reason.layer_cd = 'dds'
and not_sm_reason.table_name = 'p_shipment'

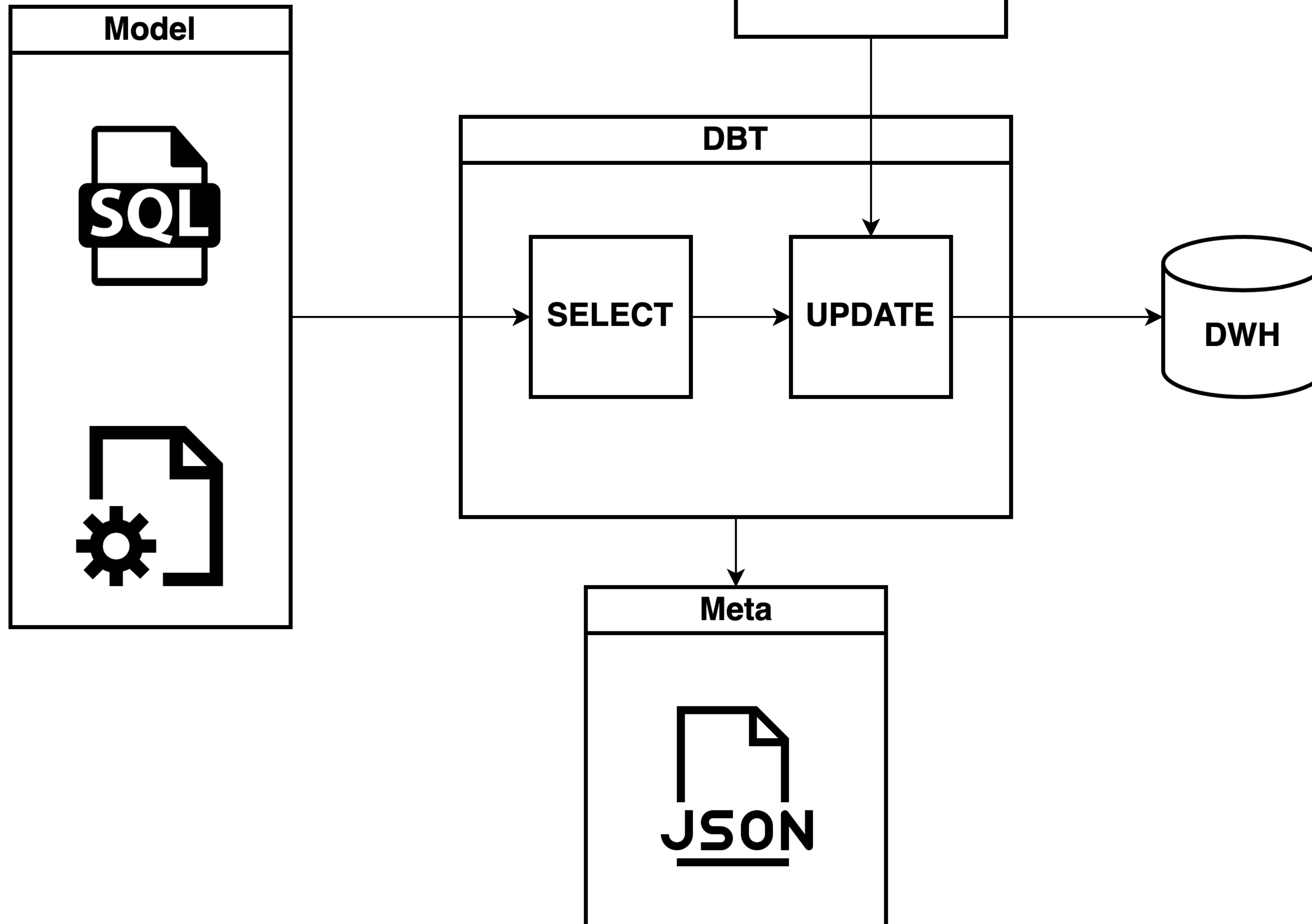
```

# Как DBT работает

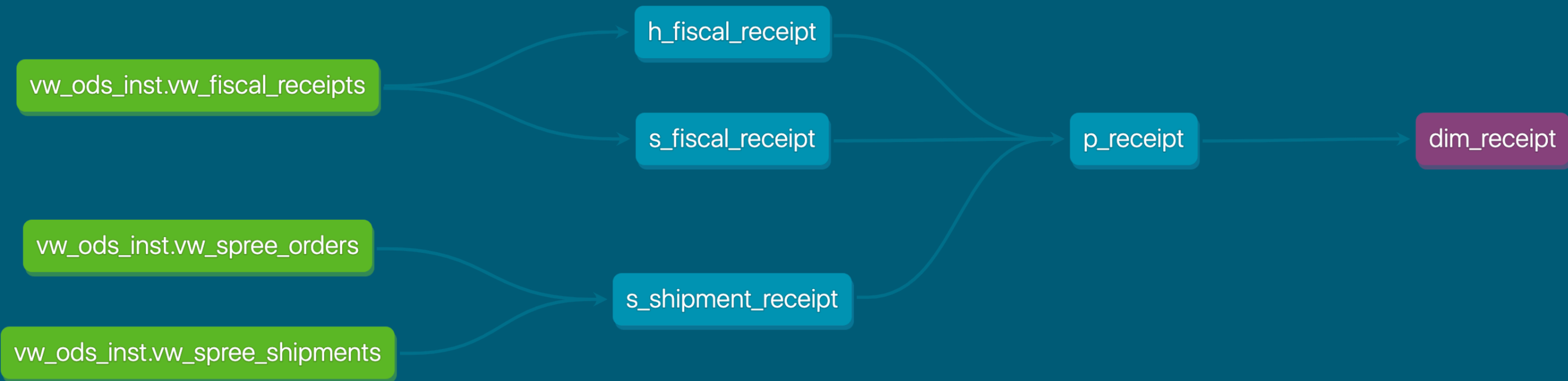




# Как DBT работает



# Граф зависимостей



## p\_receipt scd2\_materialization

[Details](#) [Description](#) [Columns](#) [Depends On](#) [SQL](#)

### Code

[Source](#) [Compiled](#)

```
1  with fiscal_receipt_ship_id as
2  (
3    select distinct s.shipment_id
4    from {{ ref('s_fiscal_receipt') }} s
5    where not s.mt_deleted_flg
6    and s.mt_active_flg
7    and s.mt_src_cd = 'INST'
8  ),
9  invoice_receipt_tmp
10 as
11 (
12   select
13     md5(ss.receipt_num::varchar) as receipt_id
```



aph



packages

tags

--select

--exclude

dbt\_dags

All selected

++

...

Update



# Раскатка и основные команды DBT

- Установка:

```
pip install dbt
```

# Раскатка и основные команды DBT

- Установка:

```
pip install dbt
```

- Запуск проекта:

```
dbt init
```

# Раскатка и основные команды DBT

- Установка:

`pip install dbt`

- Запуск проекта:

`dbt init`

- Выполнение моделей:

`dbt run`

- Выполнить конкретную модель:

`dbt run -m`

# Раскатка и основные команды DBT

- Установка:

`pip install dbt`

- Запуск проекта:

`dbt init`

- Выполнение моделей:

`dbt run`

- Выполнить конкретную модель:

`dbt run -m`

- Проверка:

`dbt parse`

- Компиляция:

`dbt compile`



# Раскатка и основные команды DBT

- Установка:

`pip install dbt`

- Запуск проекта:

`dbt init`

- Генерация документации:

`dbt docs generate`

- Запуск веб-клиента с докой:

`dbt docs serve`

- Выполнение моделей:

`dbt run`

- Выполнить конкретную модель:

`dbt run -m`

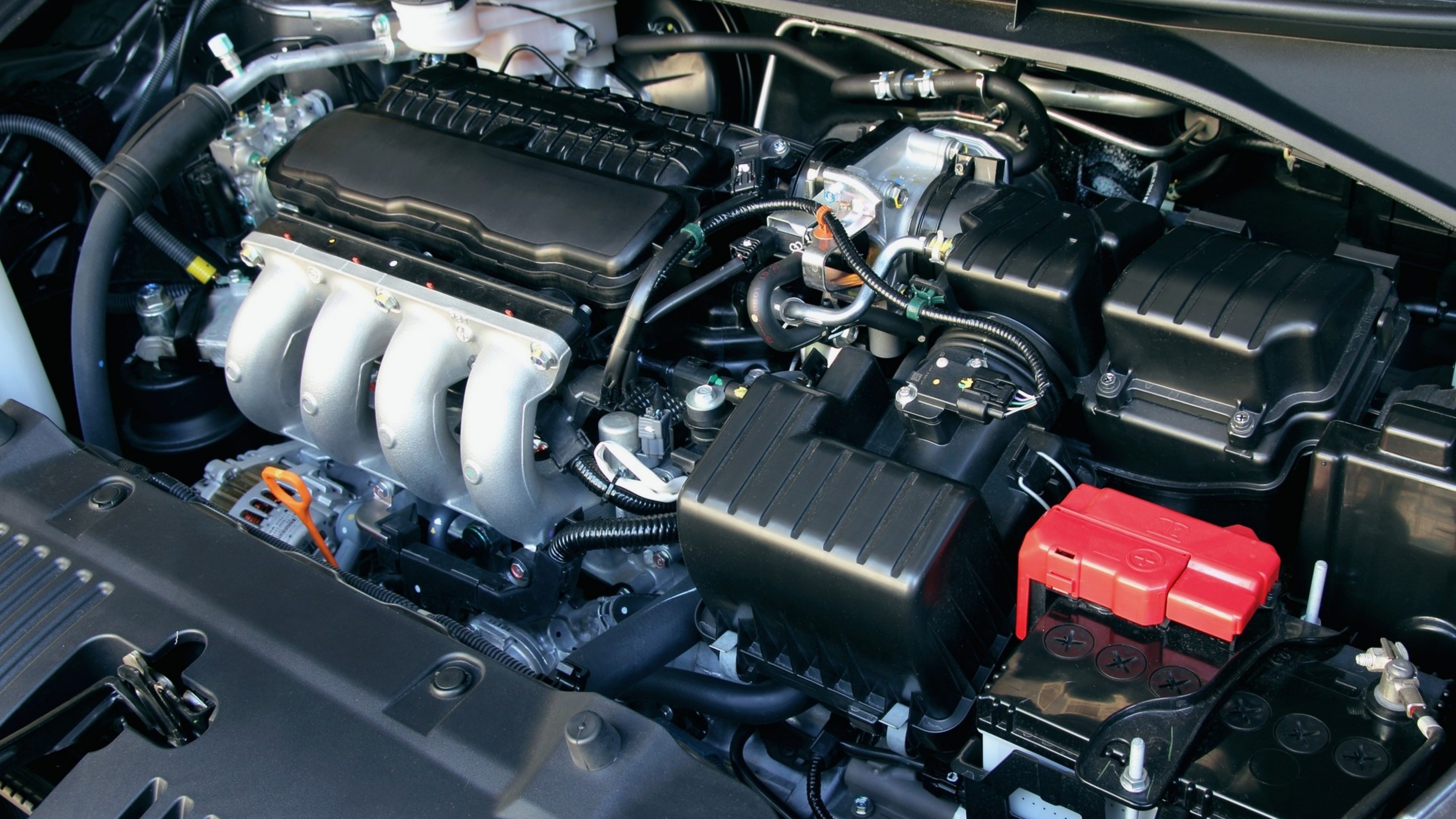
- Проверка:

`dbt parse`

- Компиляция:

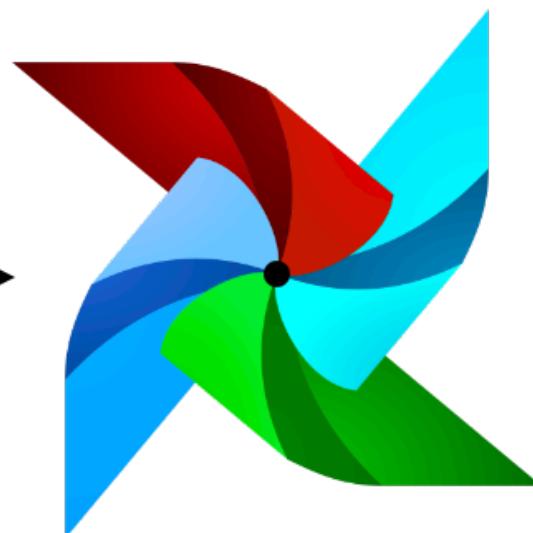
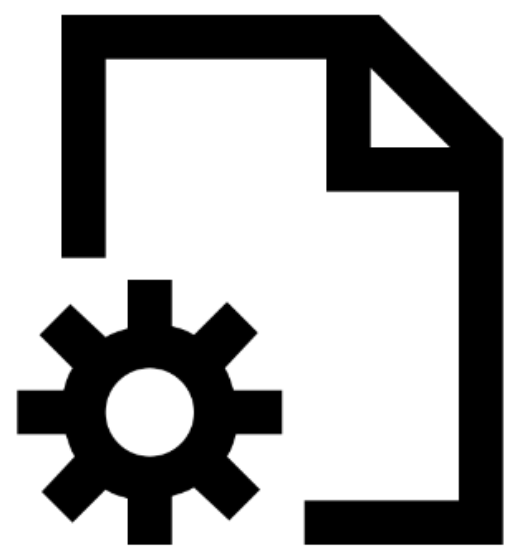
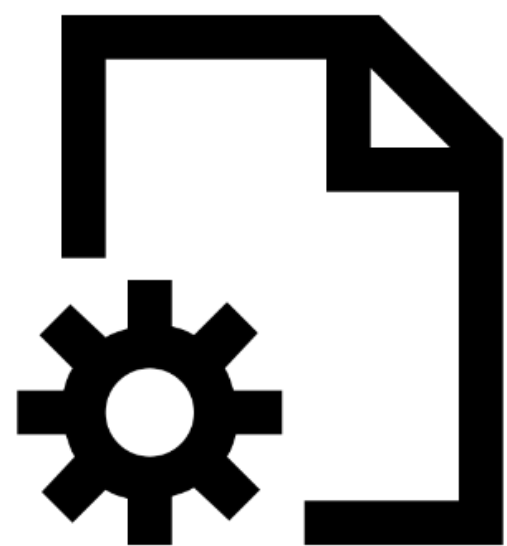
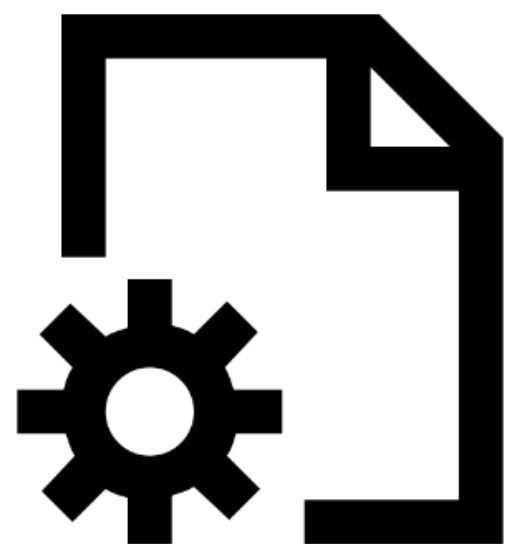
`dbt compile`



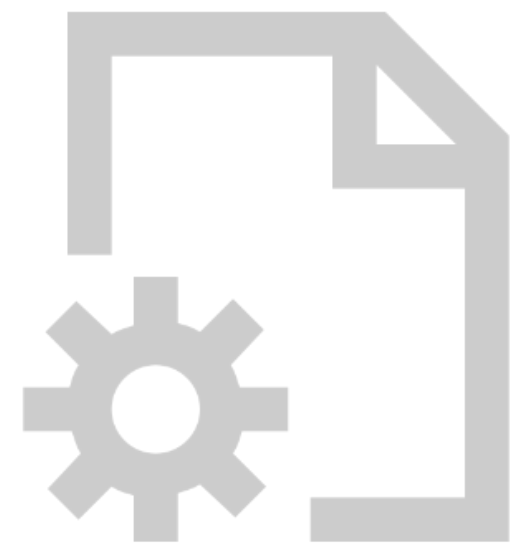
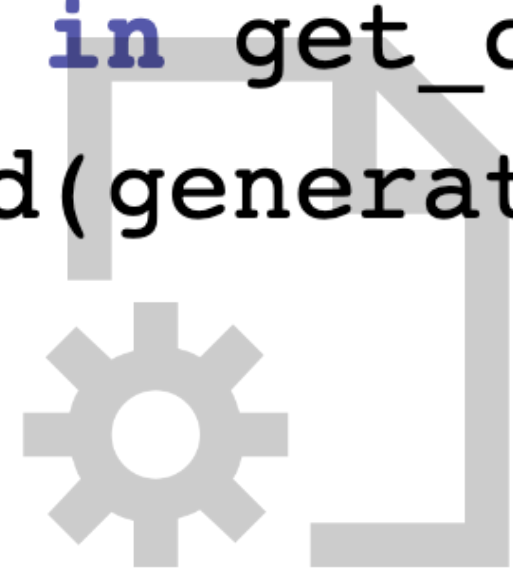




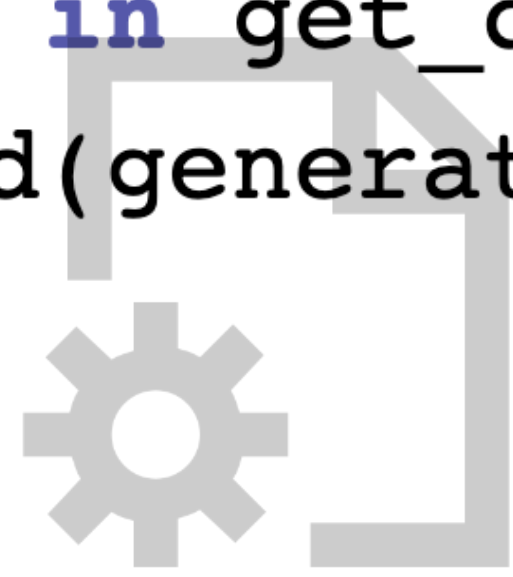
# Даг-генератор



```
for dag_config in get_dag_configs(SCHEMA_PATH, dag_configs):  
    dags.append(generate_dag_from_config(dag_config, connection_details))  
return dags
```



```
for dag_config in get_dag_configs(SCHEMA_PATH, dag_configs):  
    dags.append(generate_dag_from_config(dag_config, connection_details))  
return dags
```



```
def generate_dag_from_config(dag_config, connection_details) -> DAG:  
    with DAG(  
        dag_id=dag_id,  
        schedule_interval=schedule,  
        ...  
    ) as dag:  
        dag.doc_md = generate_doc_md_from_config(dag_id, dag_config.get('models'))  
  
    return generate_dag_branches(dag, main_configs, models, connection_details, process_type)
```





# Дэг-генератор



Tree



Graph



Calendar



Task Duration



Task Tries



Landing Times



Gantt



Details

< > Code

```
1 from airflow import DAG
2 from  import generate_dags_from_configs
3
4 dags = generate_dags_from_configs()
5
6 for dag in dags:
7     globals()[dag.dag_id] = dag
```

# Выбор оператора для DBT

```
def run_dbt( ):
    import subprocess
    subprocess.run( [ 'dbt', 'run' ] )

run_dbt_operator = PythonOperator(
    task_id='run_dbt',
    python_callable=run_dbt,
    dag=dag
)
```

# Выбор оператора для DBT

```
return DbtRunOperator(  
    task_id=f'dbt_run_{model_name}',  
    profiles_dir=DBT_DIR,  
    dir=DBT_DIR,  
    select=f'{" ".join(temp_models)} {model_name}',  
    vars={  
        "dag_id": dag_id,  
        "process_id": process_id,  
        "process_type": process_type,  
        "load_dttm": load_dttm,  
        "ext_vars": ext_vars,  
        **conn_params_dict  
    },  
    on_failure_callback=callback_error_table(process  
)
```



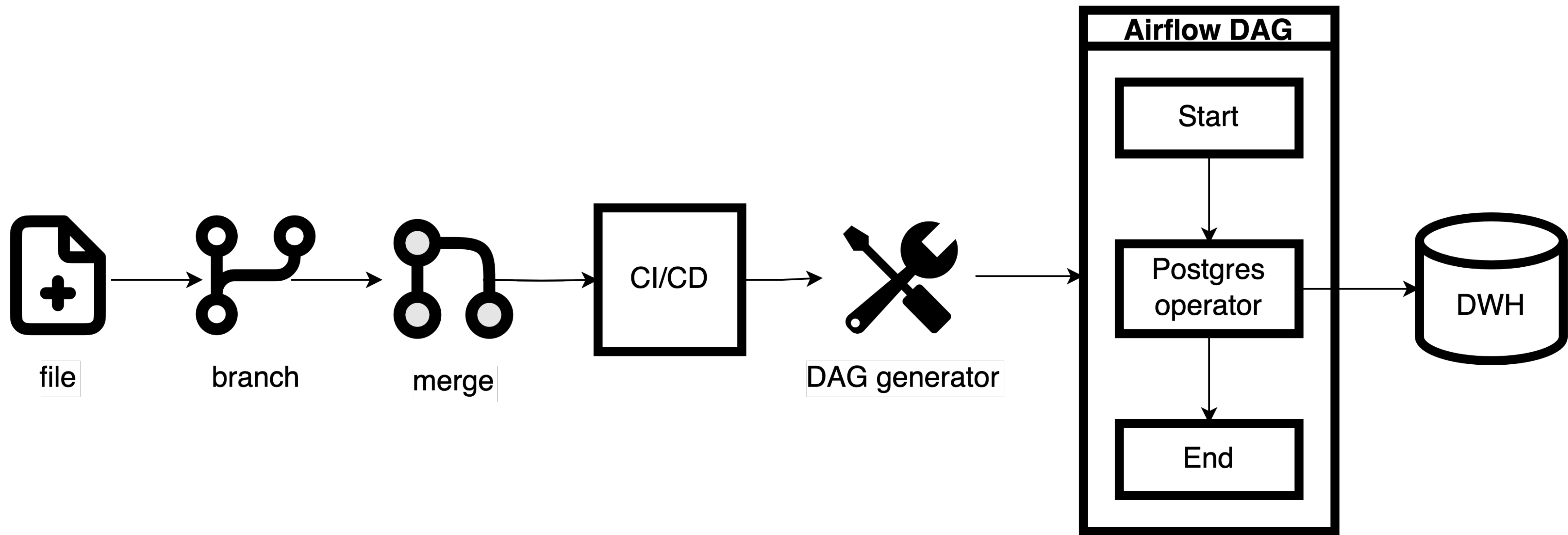
# Выбор оператора для DBT

```
return DbtRunOperator(  
    task_id=f'dbt_run_{model_name}',  
    profiles_dir=DBT_DIR,  
    dir=DBT_DIR,  
    select=dbt_run,                                ._name}',  
    vars={  
        "da    --profiles-dir   
        "pr    --vars {"dag_id": "  
        "pr    --select   
        "lc  
        "ext_vars : ext_vars,  
        **conn_params_dict  
    },  
    on_failure_callback=callback_error_table(process  
)
```

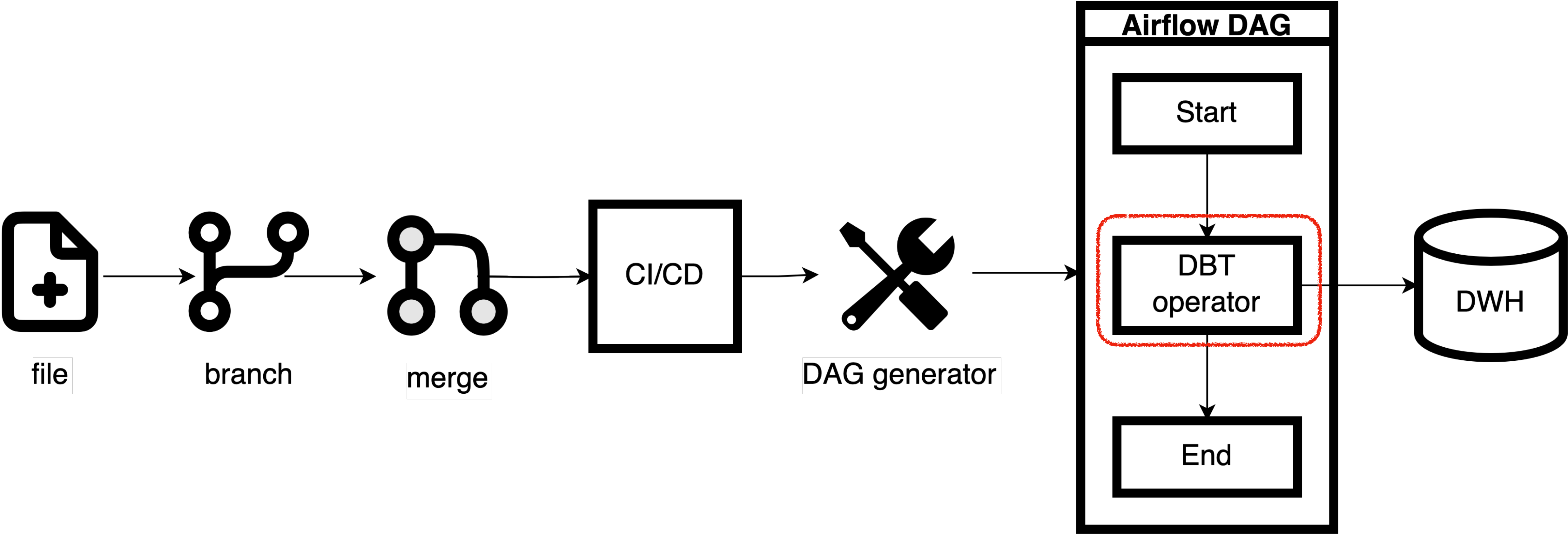




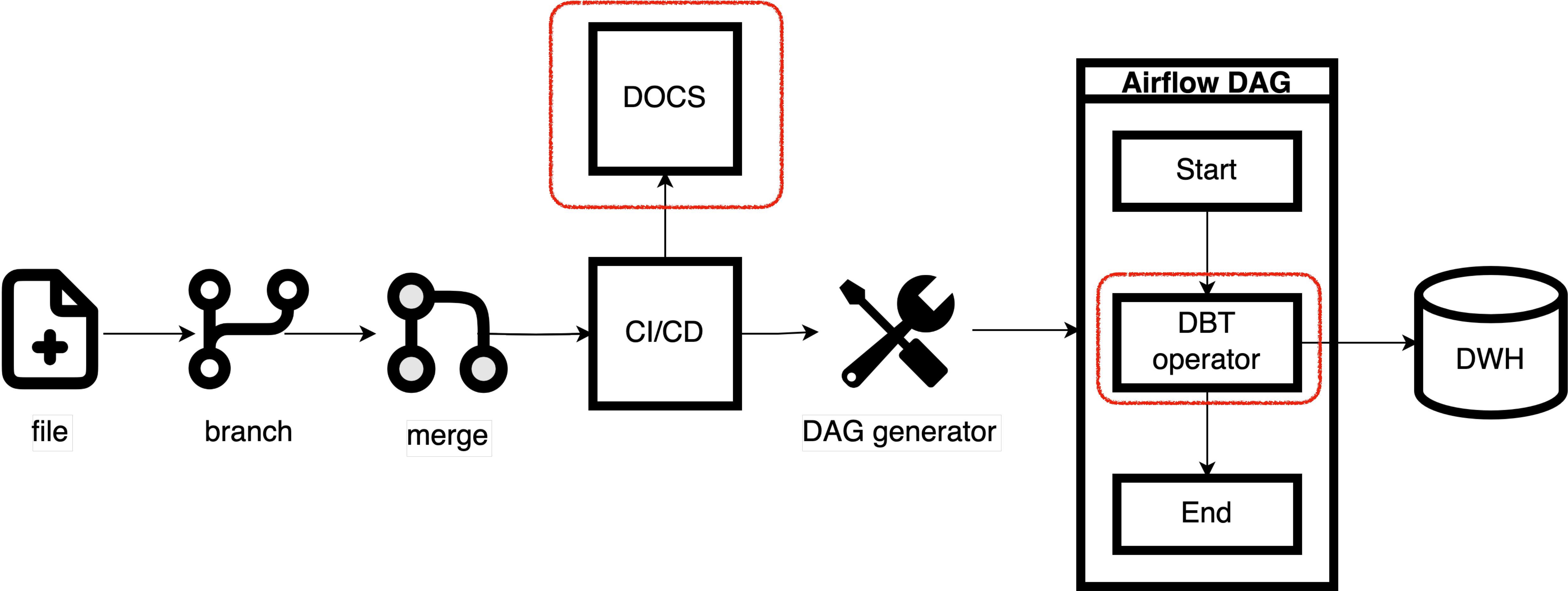
# Наше решение до DBT



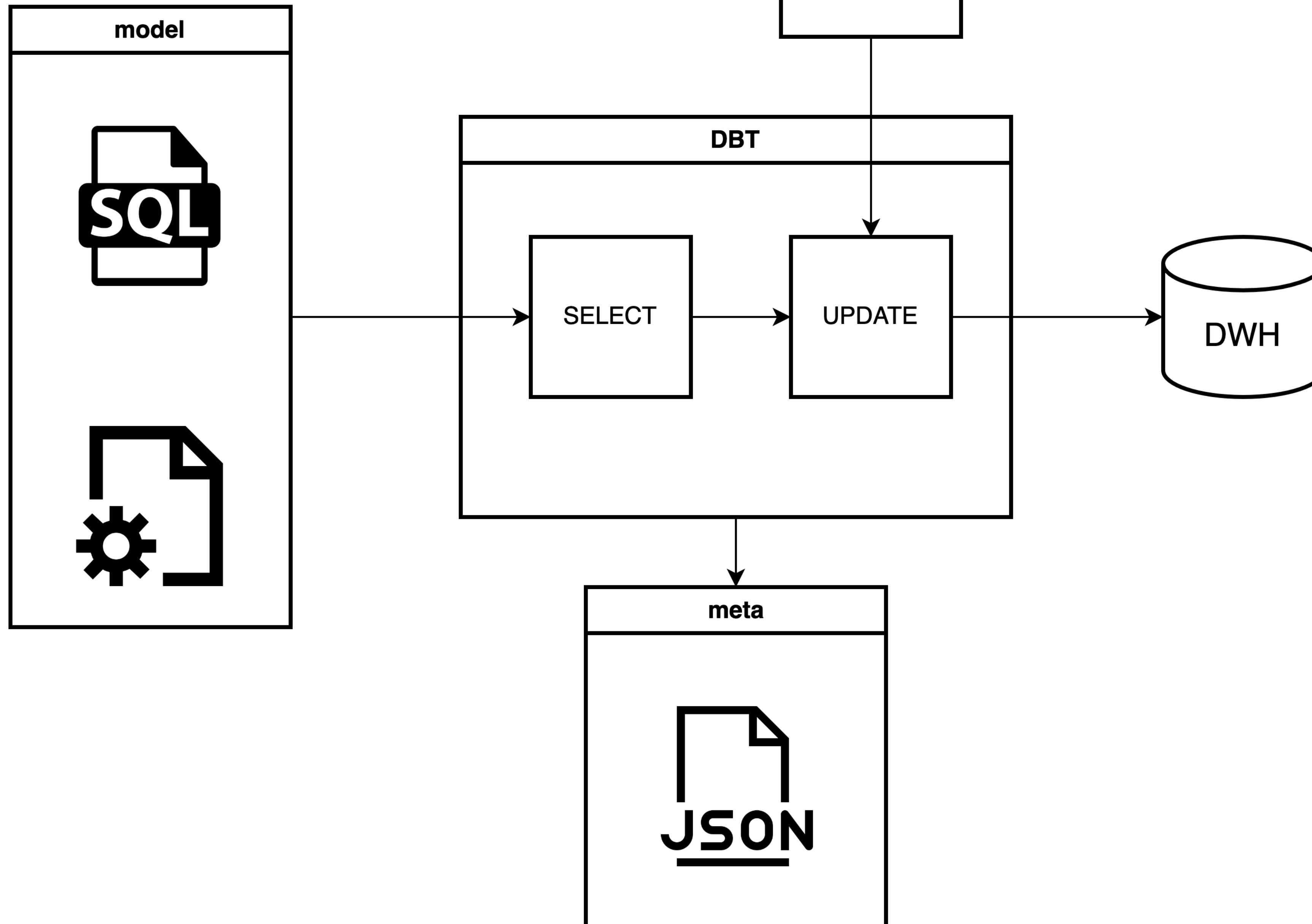
# Наше решение с DBT



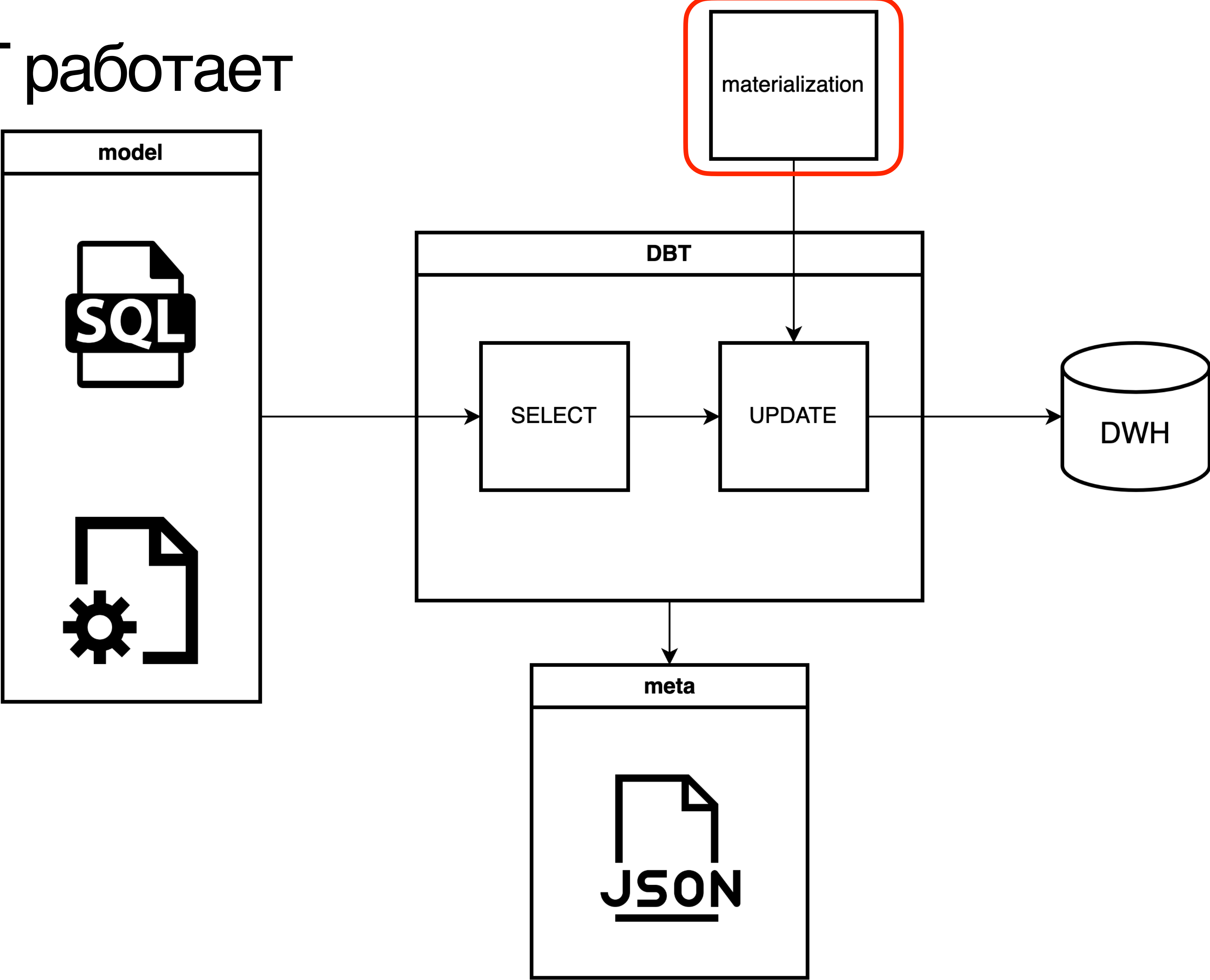
# Наше решение с DBT



# Как DBT работает



# Как DBT работает



```

{% materialization demo_materialization, default -%}
  {{ run_hooks(pre_hooks, inside_transaction=False) }}

  {% call statement("main") -%}
    insert into {{ target_relation }} (
      {{ columns_name }},
      {{ meta_columns_name }}
    )
    select
      {{ cast_columns }},
      {{ meta_columns }}
    from (sql) as source;
  {% - endcall %}

  {{ run_hooks(post_hooks, inside_transaction=True) }}

  {{ adapter.commit() }}
{% - endmaterialization %}

```



```
{% materialization demo_materialization, default -%}
```

```
from (sql) as source;
```

```
{%- endmaterialization %}
```

MODEL.sql

```
{% materialization demo_materialization, default -%}
```

```
{% call statement("main") -%}  
  insert into {{ target_relation }} (  
    {{ columns_name }},  
    {{ meta_columns_name }}  
  )  
  select  
    {{ cast_columns }},  
    {{ meta_columns }}  
  from (sql) as source;  
{%- endcall %}
```

```
{%- endmaterialization %}
```

## call statement

MODEL.sql

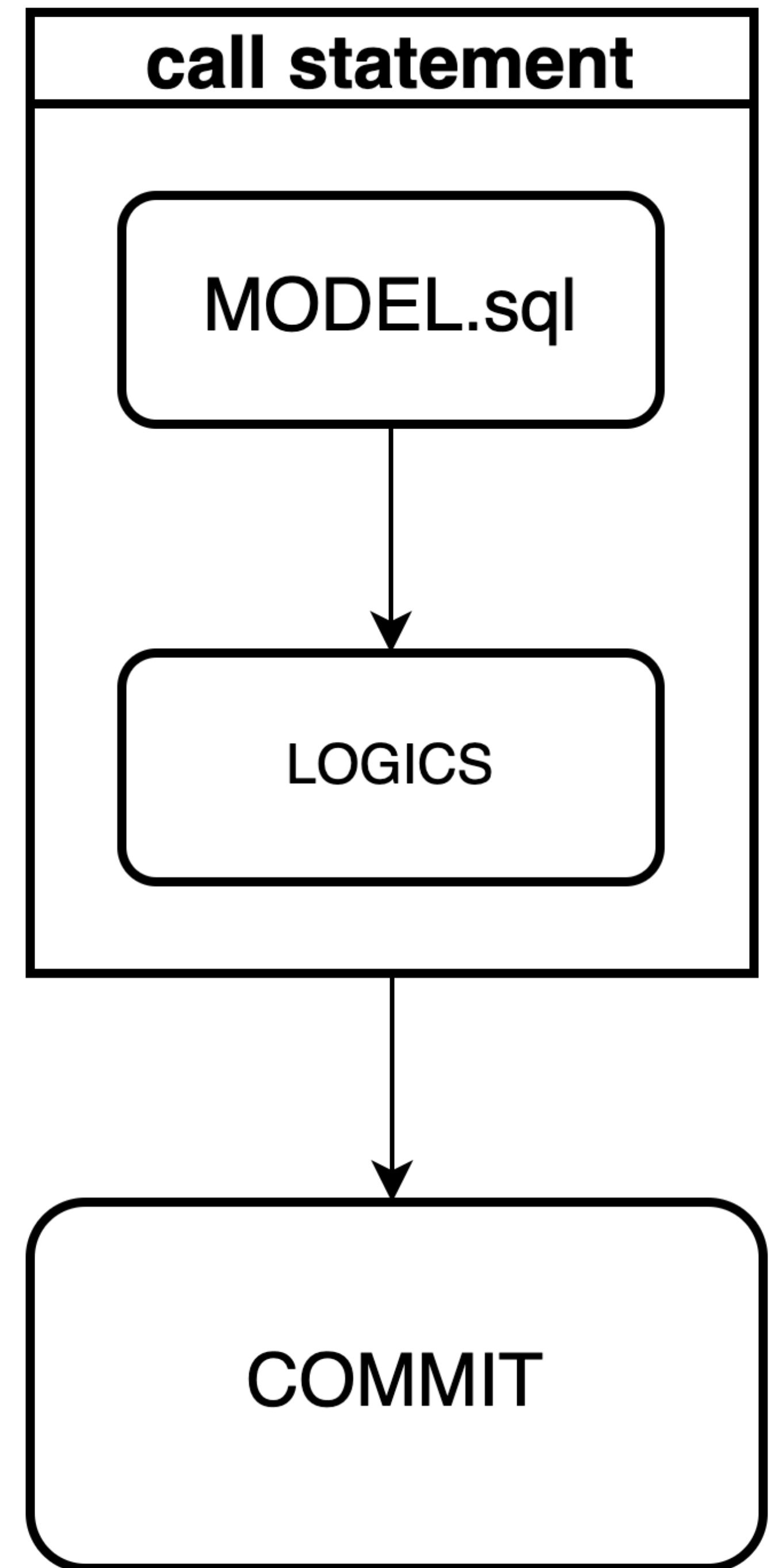


LOGICS

```
{% materialization demo_materialization, default -%}
```

```
{% call statement("main") -%}  
  insert into {{ target_relation }} (  
    {{ columns_name }},  
    {{ meta_columns_name }}  
  )  
  select  
    {{ cast_columns }},  
    {{ meta_columns }}  
  from (sql) as source;  
{%- endcall %}
```

```
  {{ adapter.commit() }}  
{%- endmaterialization %}
```



```

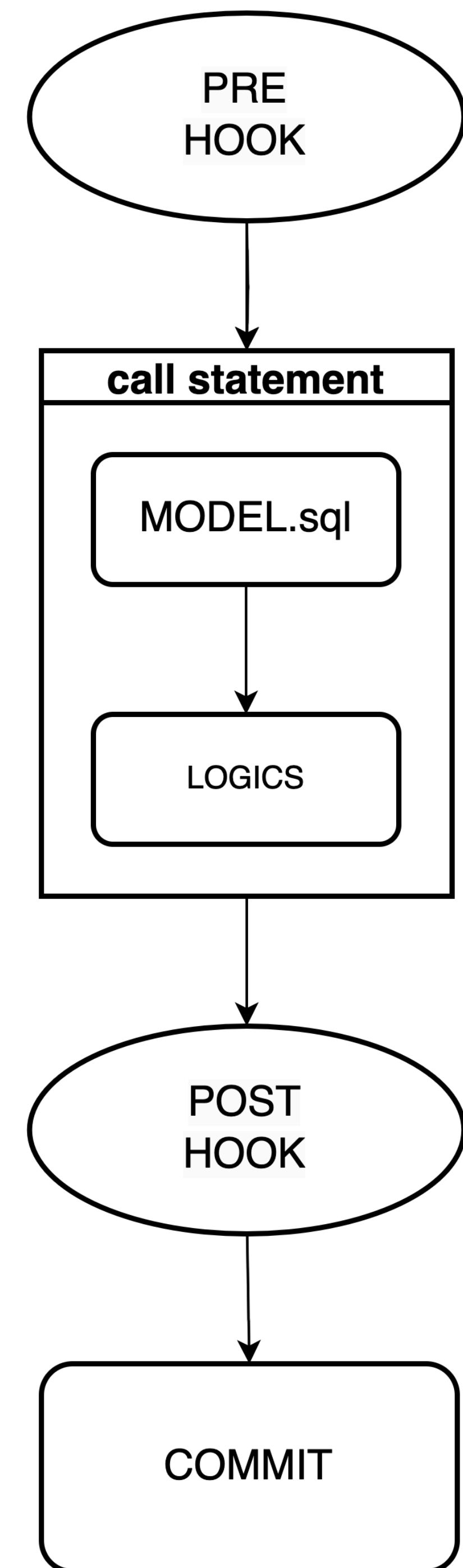
{% materialization demo_materialization, default -%}
  {{ run_hooks(pre_hooks, inside_transaction=False) }}

  {% call statement("main") -%}
    insert into {{ target_relation }} (
      {{ columns_name }} ,
      {{ meta_columns_name }}
    )
    select
      {{ cast_columns }} ,
      {{ meta_columns }}
    from (sql) as source;
  {% - endcall %}

  {{ run_hooks(post_hooks, inside_transaction=True) }}

  {{ adapter.commit() }}
{% - endmaterialization %}

```



```

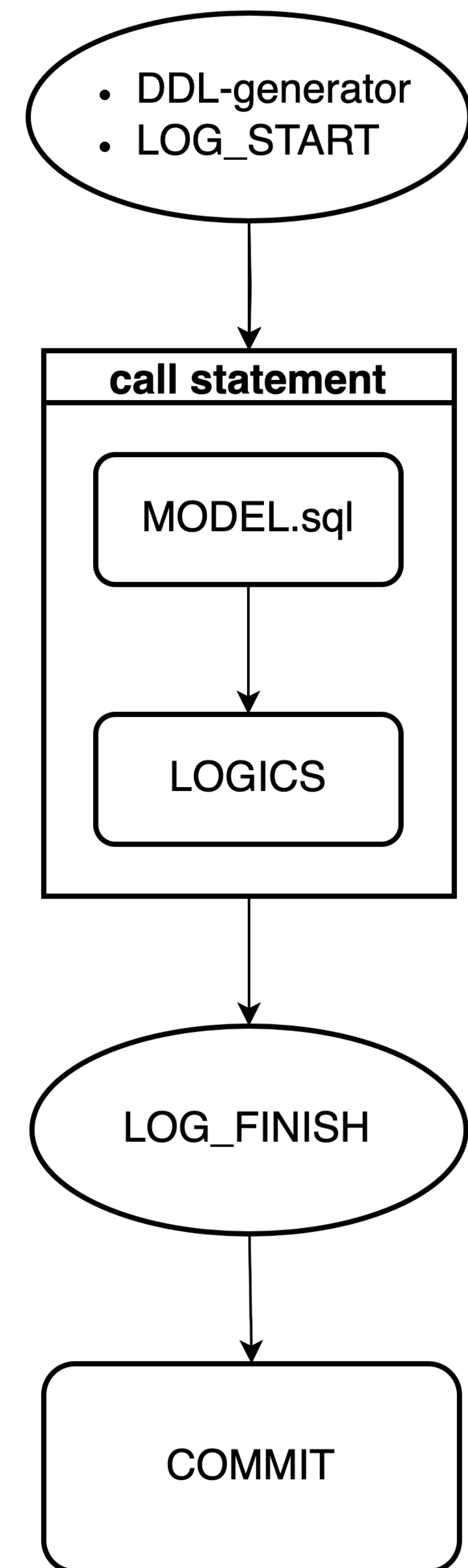
{% materialization demo_materialization, default -%}
  {{ run_hooks(pre_hooks, inside_transaction=False) }}

  {% call statement("main") -%}
    insert into {{ target_relation }} (
      {{ columns_name }} ,
      {{ meta_columns_name }}
    )
    select
      {{ cast_columns }} ,
      {{ meta_columns }}
    from (sql) as source;
  {% - endcall %}

  {{ run_hooks(post_hooks, inside_transaction=True) }}

  {{ adapter.commit() }}
{% - endmaterialization %}

```





```

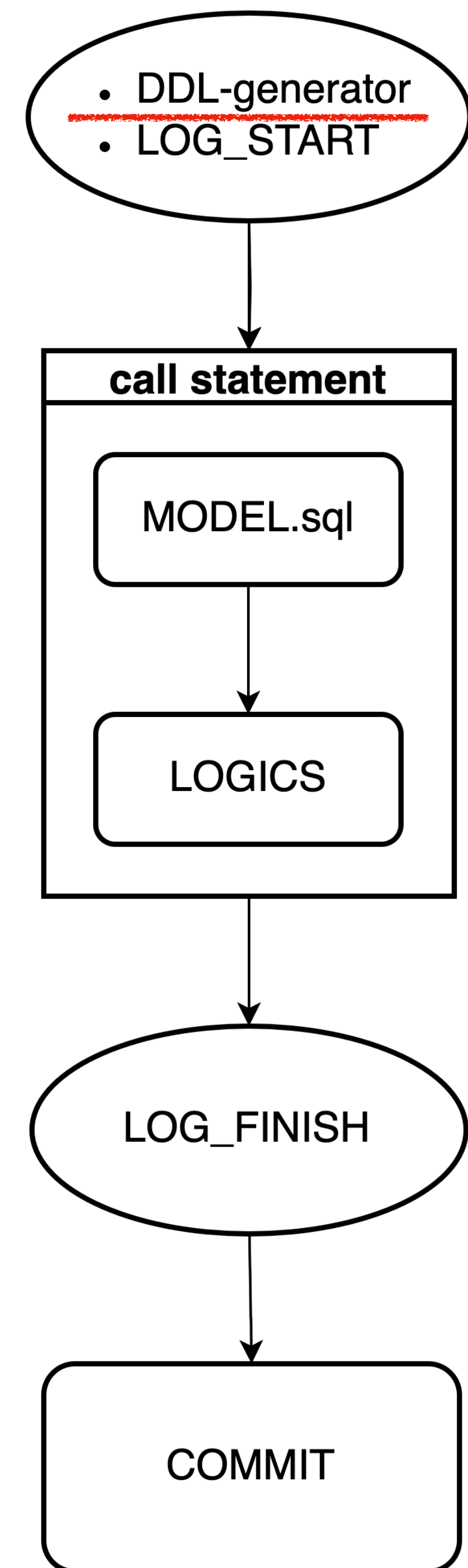
{% materialization demo_materialization, default -%}
  {{ run_hooks(pre_hooks, inside_transaction=False) }}

  {% call statement("main") -%}
    insert into {{ target_relation }} (
      {{ columns_name }} ,
      {{ meta_columns_name }}
    )
    select
      {{ cast_columns }} ,
      {{ meta_columns }}
    from (sql) as source;
  {% - endcall %}

  {{ run_hooks(post_hooks, inside_transaction=True) }}

  {{ adapter.commit() }}
{% - endmaterialization %}

```



# DDL генератор

```
{% macro create_table_ddl(relation) %}  
    CREATE TABLE {{ relation }} (  
        {%- for col in columns() %}{% if not loop.first %}, {% endif %}  
        {{ col }}{% endfor %}  
    )  
    {{ storage_parameters() }}  
    {{ distribution() }}  
    {{ partition() }}  
    ;  
  
    {%- for comment in comments(relation) %}  
    {{ comment }}{% endfor %}  
  
    {{ grant_table(relation) }}  
{% endmacro %}
```

# DDL генератор

```
{% macro create_table_ddl(relation) %}  
    CREATE TABLE {{ relation }} (  
        {%- for col in columns() %}{% if not loop.first %}, {% endif %}  
        {{ col }}{% endfor %}  
    )  
    {{ storage_parameters() }}  
    {{ distribution() }}  
    {{ partition() }}  
    ;  
  
    {%- for comment in comments(relation) %}  
    {{ comment }}{% endfor %}  
  
    {{ grant_table(relation) }}  
{% endmacro %}
```

# Промежуточные итоги

ДБТ установили.

Генератор airflow написали. Даги создаются по конфигам.

Материализации написали.

Генератор таблиц.

Логирование.

Бери и пользуйся?



— Мне надоело, у нее были слишком большие запросы.

— Например какие?

— Ну например update instance inner join  
(select **group.id** as group\_id, (select **message.id**  
from message inner join thread on **thread.id**  
= message.thread\_id where location\_id =  
@location\_id and language\_id = @language\_id  
and concat(group\_key, '.') like concat(group.`key`,  
'.%')) order by message.created desc limit 1) as  
last\_message\_id, (select count(\*) from thread  
where location\_id = @location\_id and language\_id  
= @language\_id and concat(group\_key, '.') like  
concat(group.`key`, '.%')) as thread\_count,  
(select if(sum(thread.message\_count) is null, 0,  
sum(thread.message\_count)) from thread where  
location\_id = @location\_id and language\_id =  
@language\_id and concat(group\_key, '.') like  
concat(group.`key`, '.%')) as message\_count from  
group where @channel\_key like concat(`key`, '.  
%')) as statistics on statistics.group\_id =







# Временные таблицы

```
on-run-end:
```

```
- "{{ drop_temp_tables(results) }}"
```

```
models:
```

```
  dbt_dags:
```

```
    scripts:
```

```
      tmp:
```

```
        +schema: tmp
```

```
        +materialized: temp_table_materialization
```



# Временные таблицы

```
models:
  - name: p_shipment
    config:
      temp_models:
        - p_shipment_num_items_cancel_replace_tmp

  - name: p_shipment_num_items_cancel_replace_tmp
    docs:
      show: false
    config:
      meta:
        temporary: true
        distribution: shipment_id
```







# Зависимости дагов airflow

sensor:

required\_tables:

detail:

- 'dds.d\_api\_client'
- 'dds.d\_city'
- 'dds.d\_delivery\_category'
- 'dds.d\_delivery\_method\_kind'
- 'dds.d\_delivery\_window\_change\_reason'
- 'dds.d\_loyalty\_transaction\_kind'
- 'dds.d\_loyalty\_transaction\_state'

# Как долго

- MVP (генератор + одна модель) - 1 мес
- Все модели на DBT - 6 мес
- К текущему виду - 3 мес
- Дальше - ?



# Выводы

1. Легко начать
2. Гибко можно настроить под свои алгоритмы
3. Управление возможно по-модельное
4. Экономится время доставки витрин
5. Наглядные графы зависимостей



Спасибо за внимание