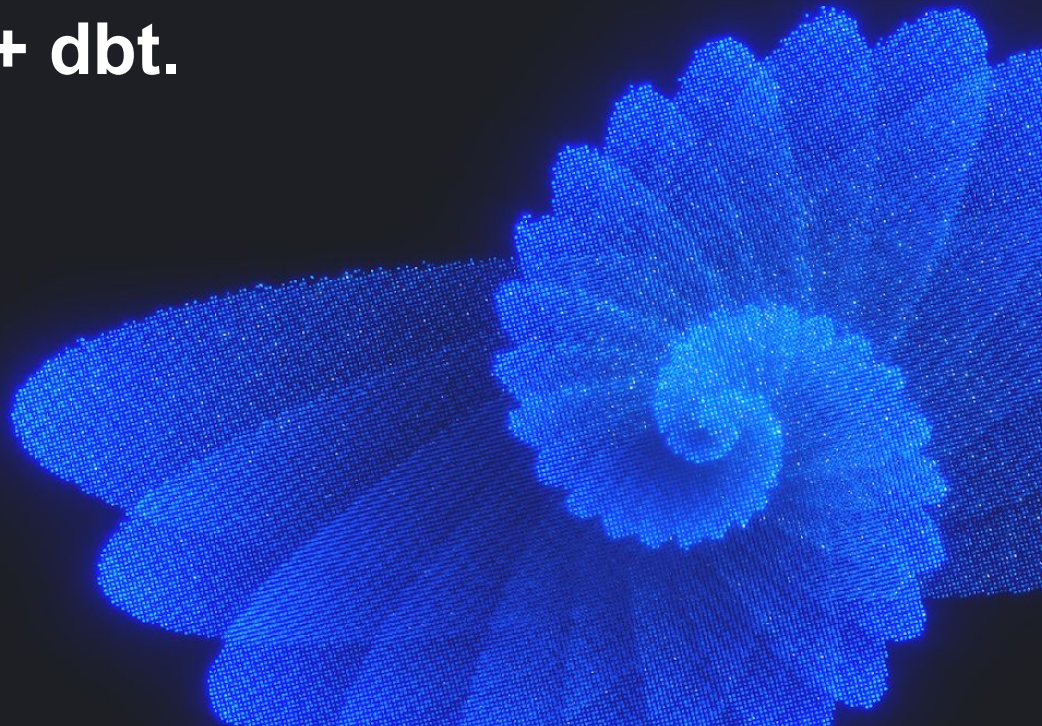




От хайпа до продакшена: DataMesh на Airflow + dbt.

@nikitayurasov
@lkozhinov



Никита Юрасов



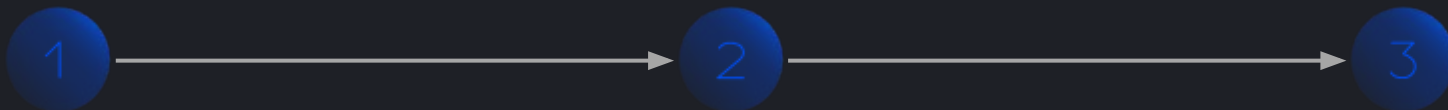
Разработчик по призванию, data engineer по профессии и математик в душе.

Кожинев Леонид



Старший разработчик мигрировавший в DevOps'а.
Фанат модных языков программирования, если это C++.

[TODO] Какой план?



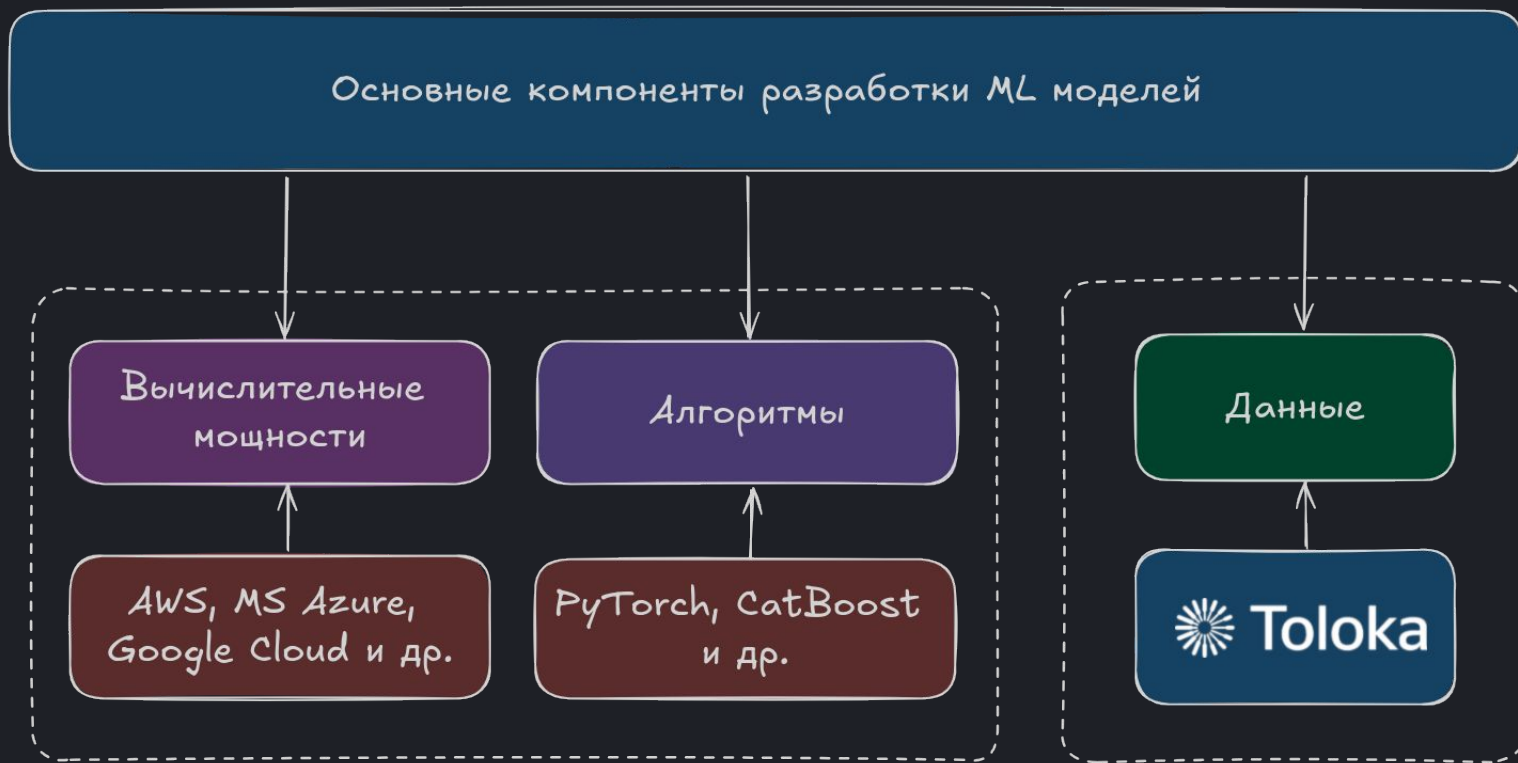
- О нас
- Data LakeHouse
- Data Mesh
- Наша архитектура

- Что мы хотим от dbt + Airflow
- dbt + Airflow = dbt-af
- dbt-af = ❤️

- Что получилось
- Что хорошо
- А что плохо

Глава I. Будем знакомы

Данные - это новая нефть ©

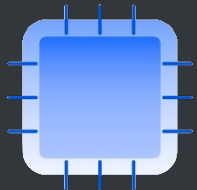


<https://toloka.ai/about-us/>

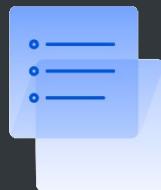
Наша прекрасная Data и ML платформа!



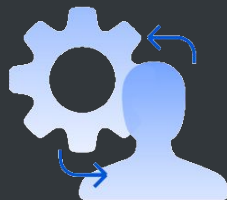
530 Тб



400 CPU



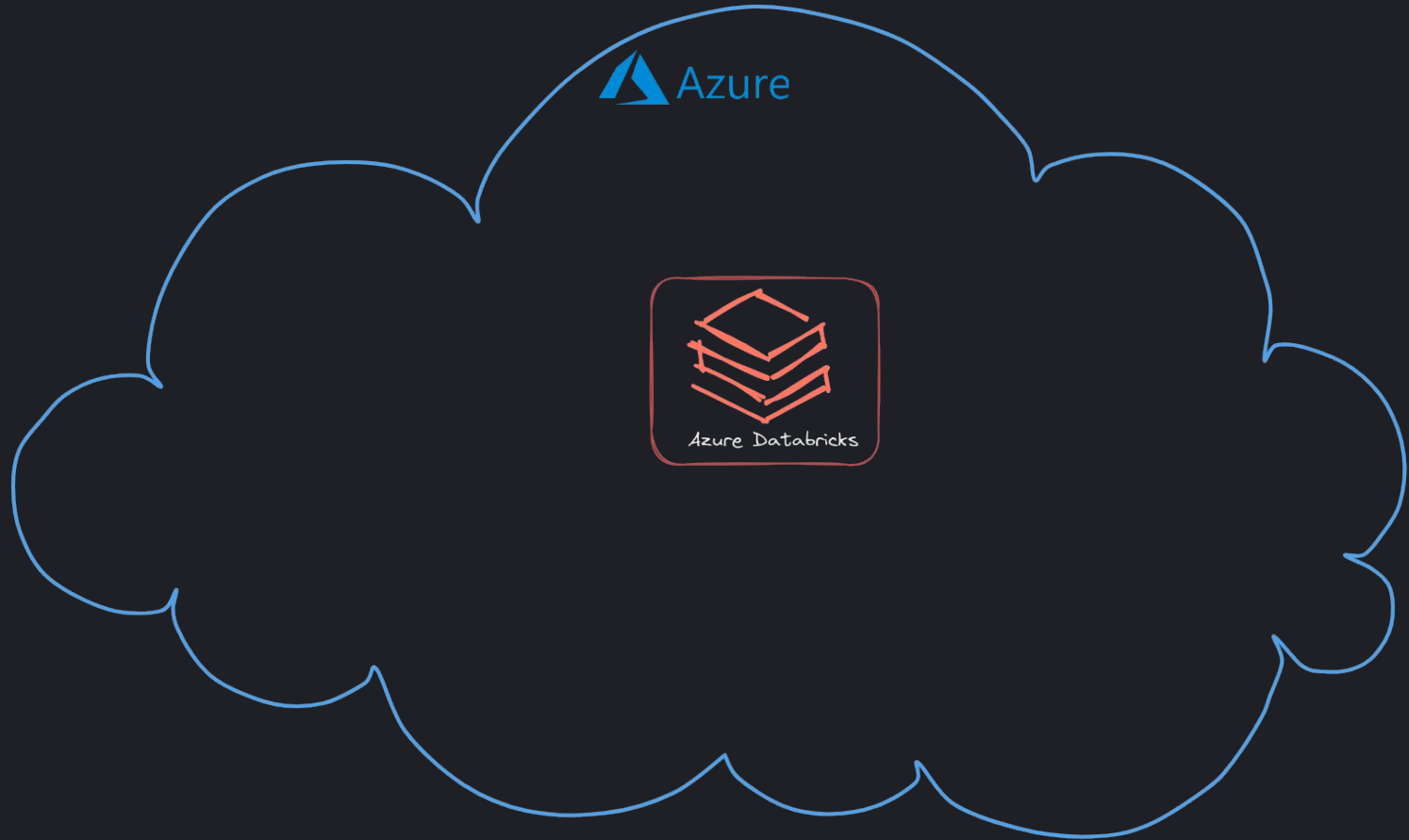
4700 таблиц



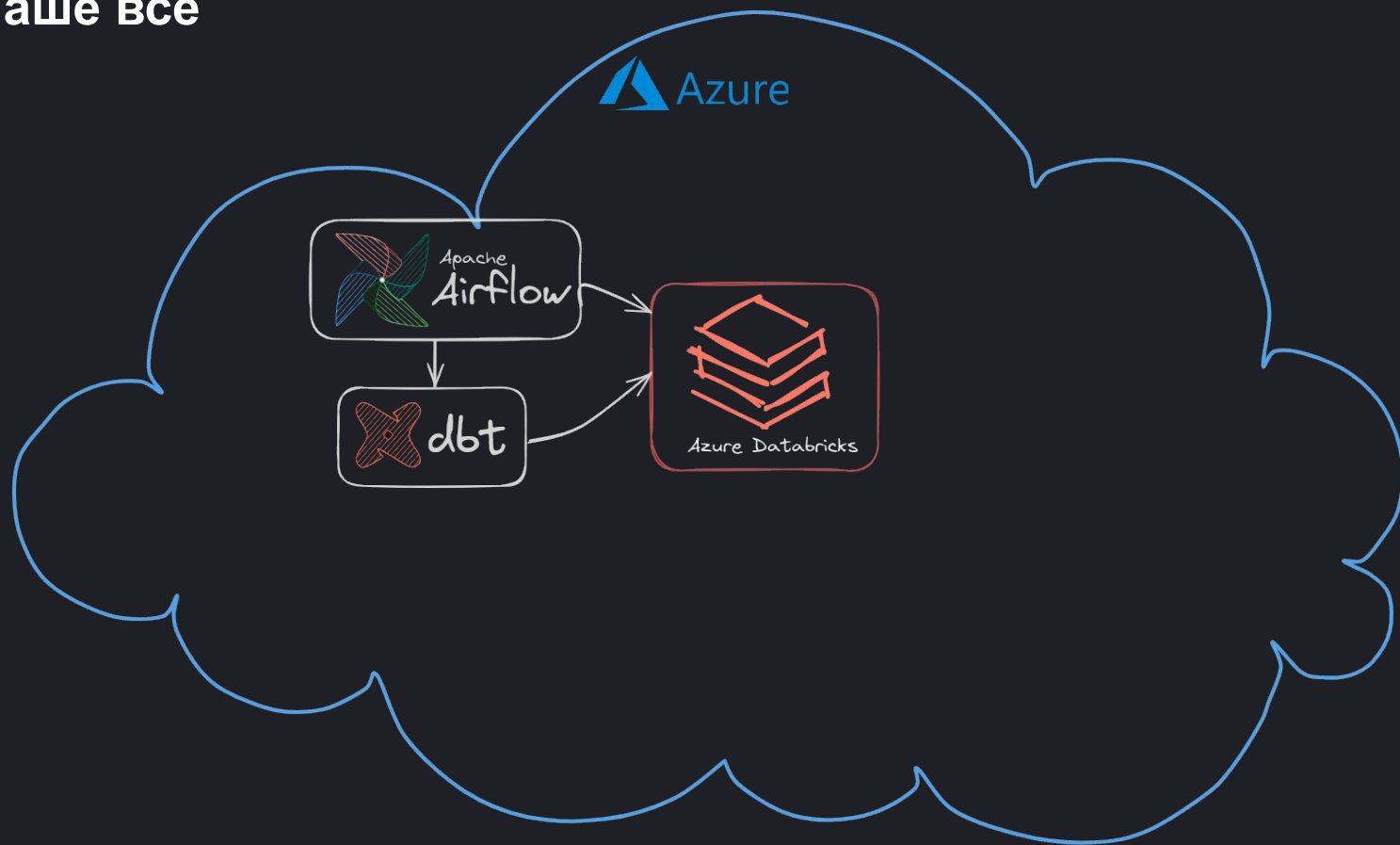
100 внутренних
пользователей



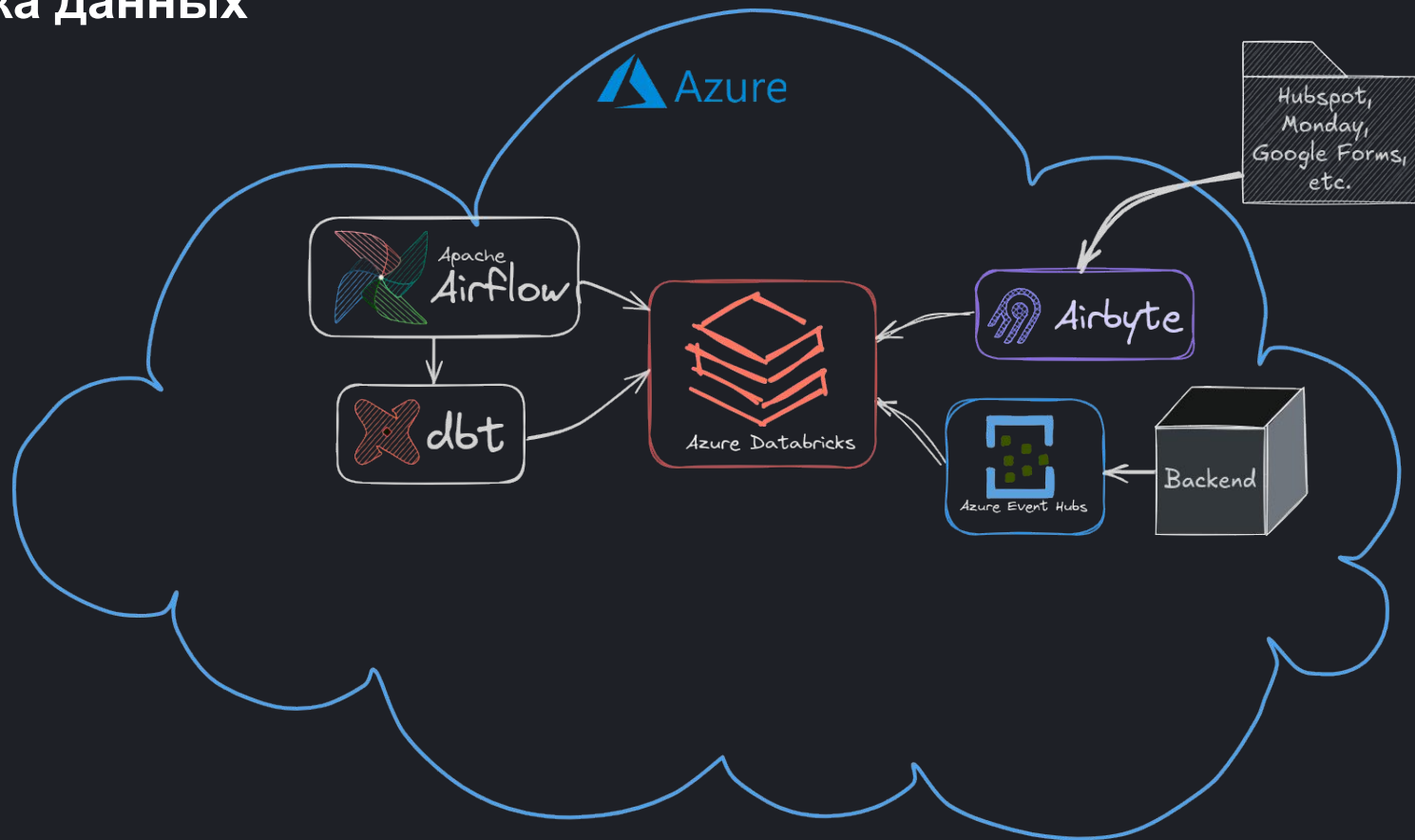
Databricks - ядро нашего DataLake house'a



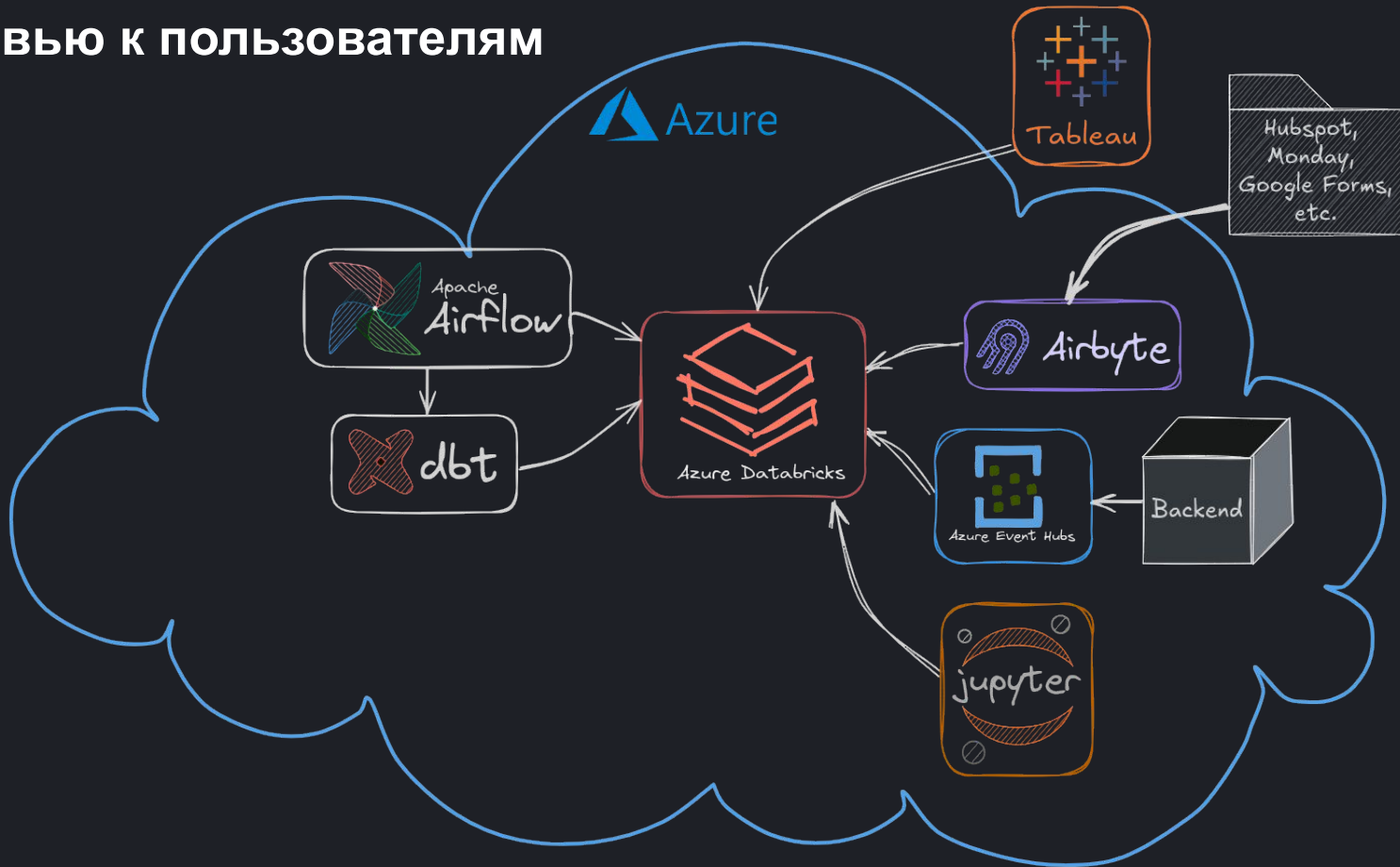
DBT - наше всё



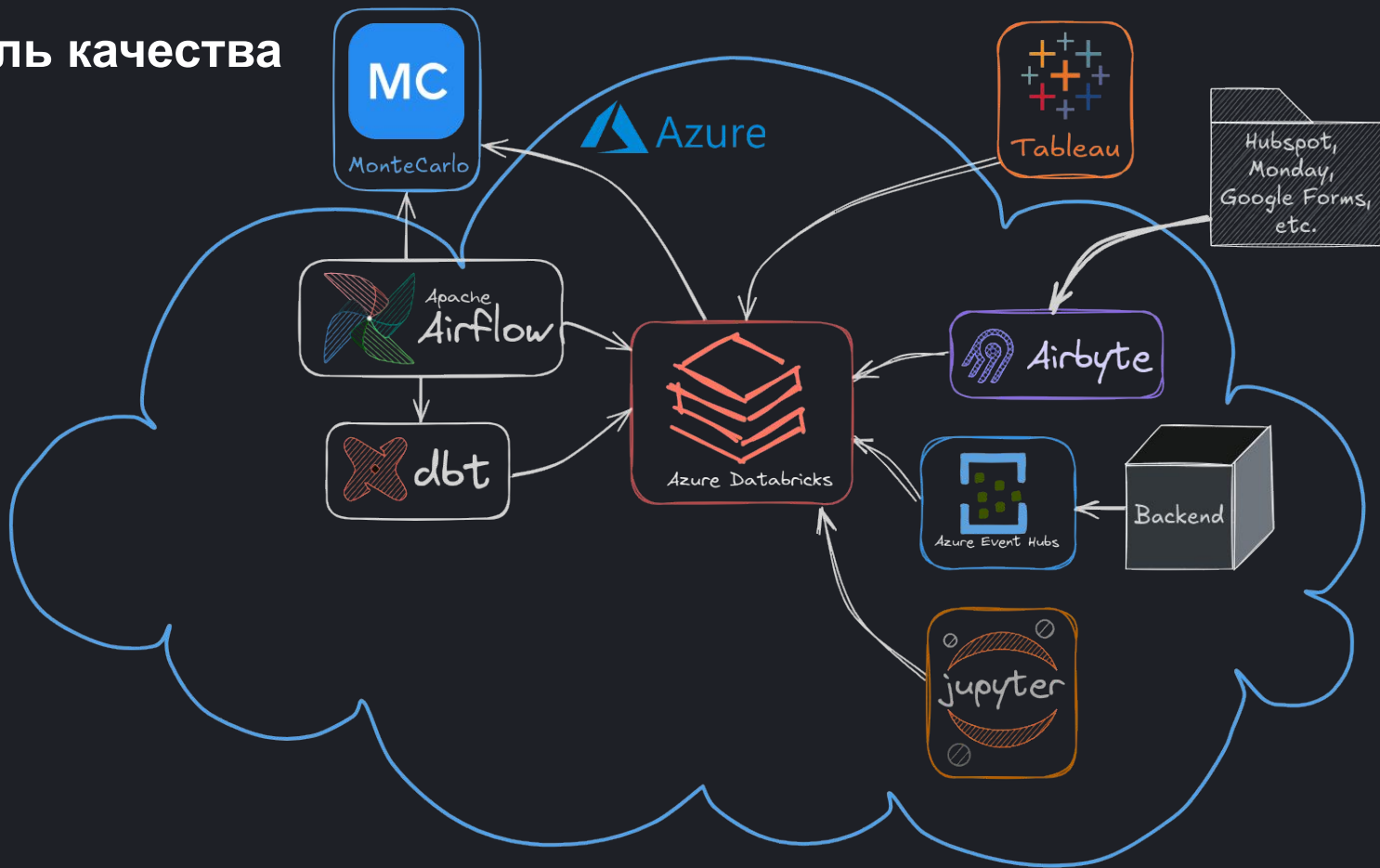
Загрузка данных



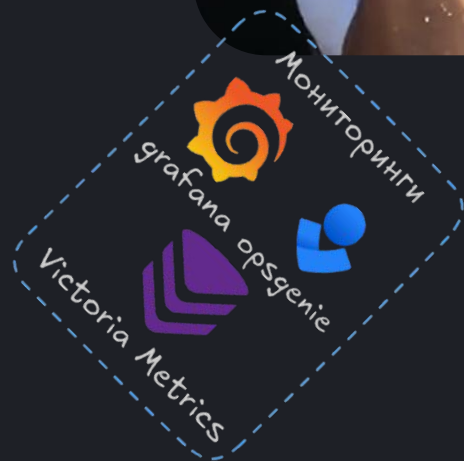
С любовью к пользователям



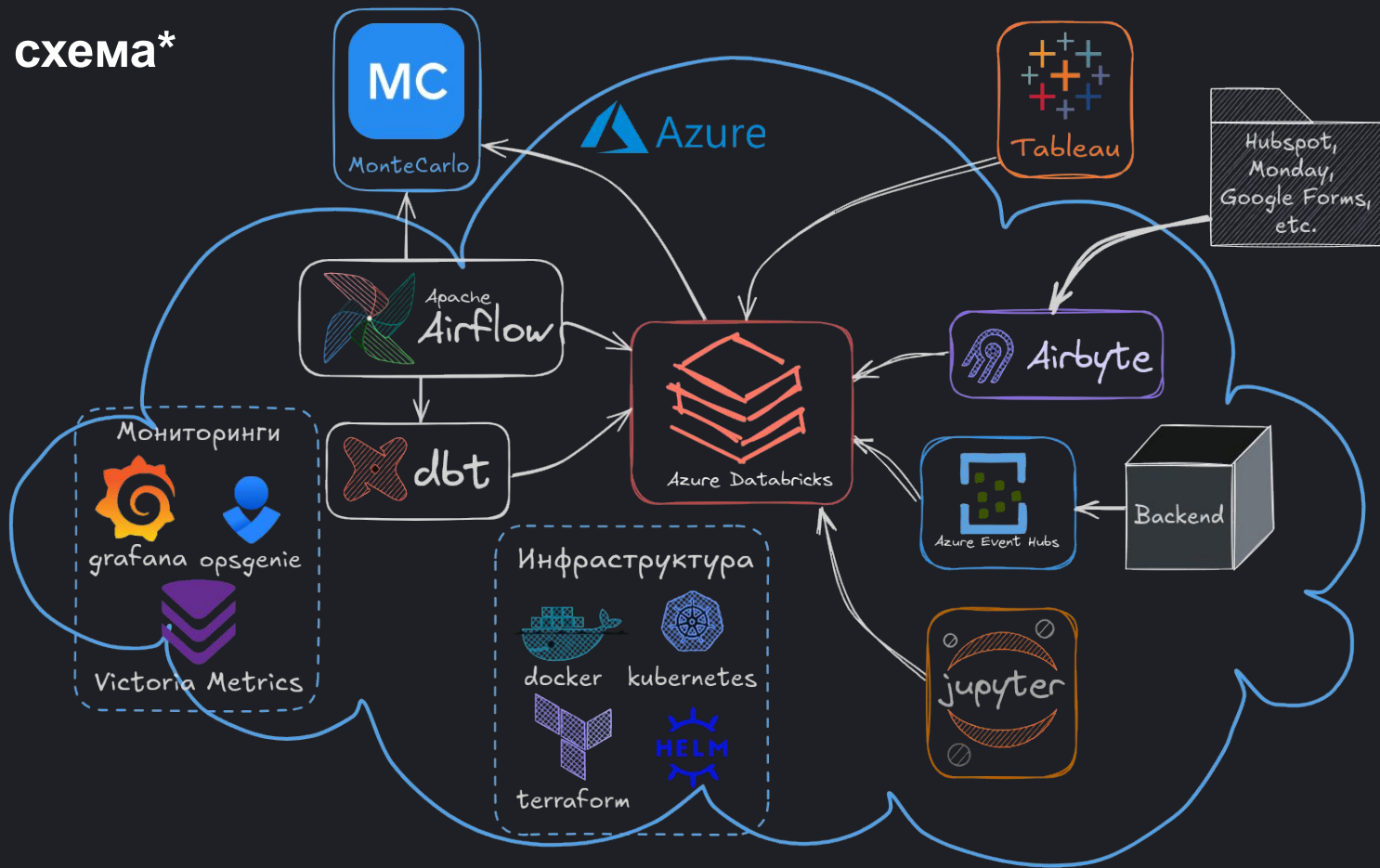
Контроль качества



Приправим мониторингами и инфрой



Полная схема*



Так что, dbt run и “поехали”, да?



Размер имеет значение

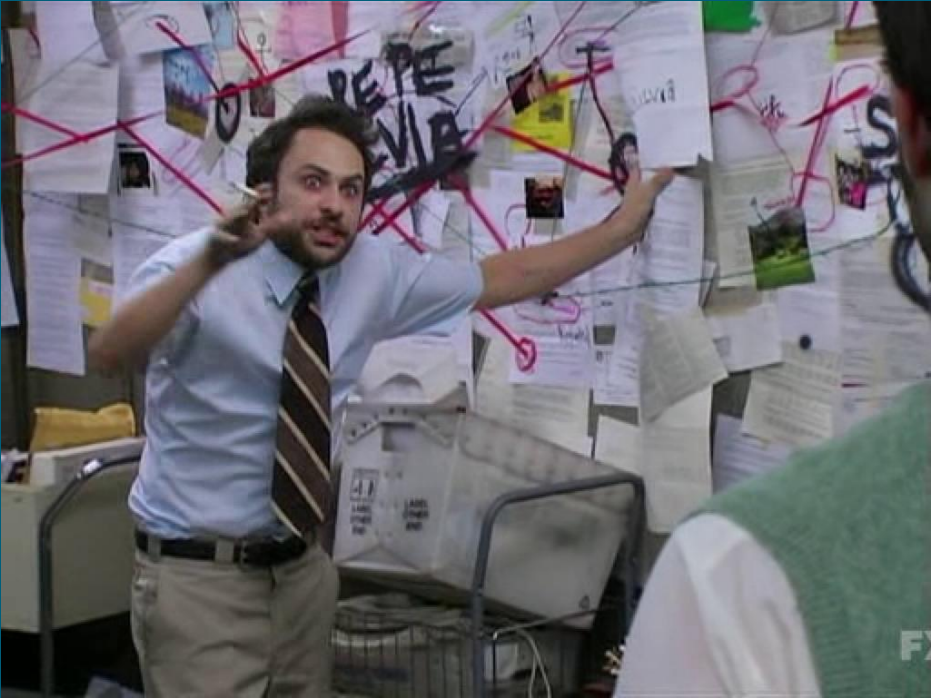


Toasty!



Размер имеет значение

Lineage Graph



resources All selected

packages tik_dwh

tags All selected


--select

--exclude

Update Graph X

И сказал лид “Декомпозируй” и стало хорошо

Lineage Graph



resources packages tags --select --exclude
All selected tik_dwh All selected Update Graph X

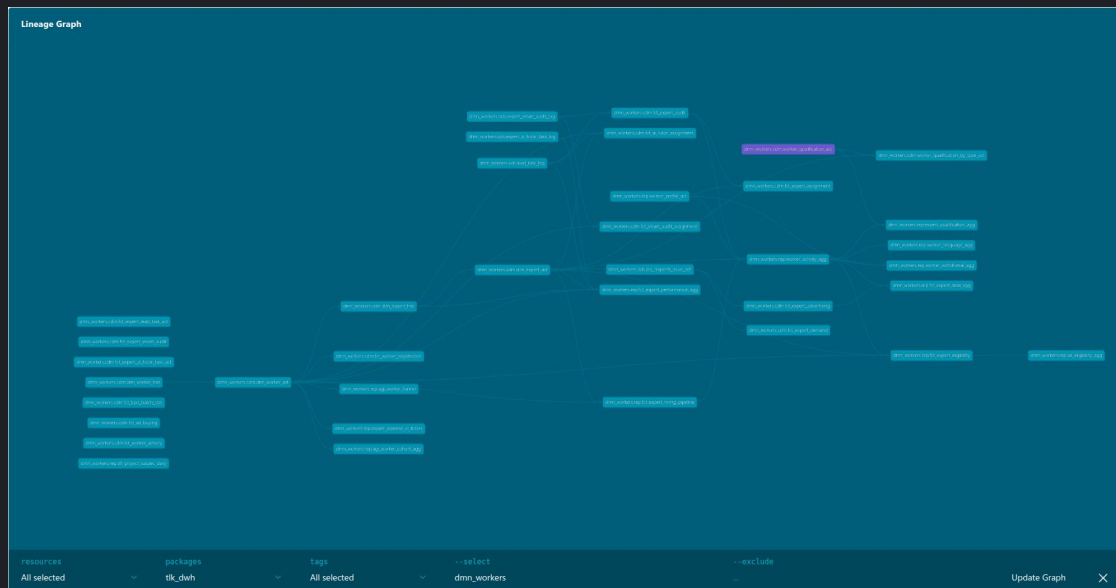
Домен 1

- Название: svc_crm
- Ответственные: ...
- Доменные RBAC роли



Домен 2

- Название: dmn_workers
- Ответственные: ...
- Доменные RBAC роли



[TODO] Слайды про нашу доменную структуру

[TODO] Слайд-подводка для dbt-af

Глава II. Решающая



Что мы хотим?


- Авто-генерация DAG
- Идемпотентность запусков dbt
- Поддержка разных уровней шедулинга

Дополнительно:

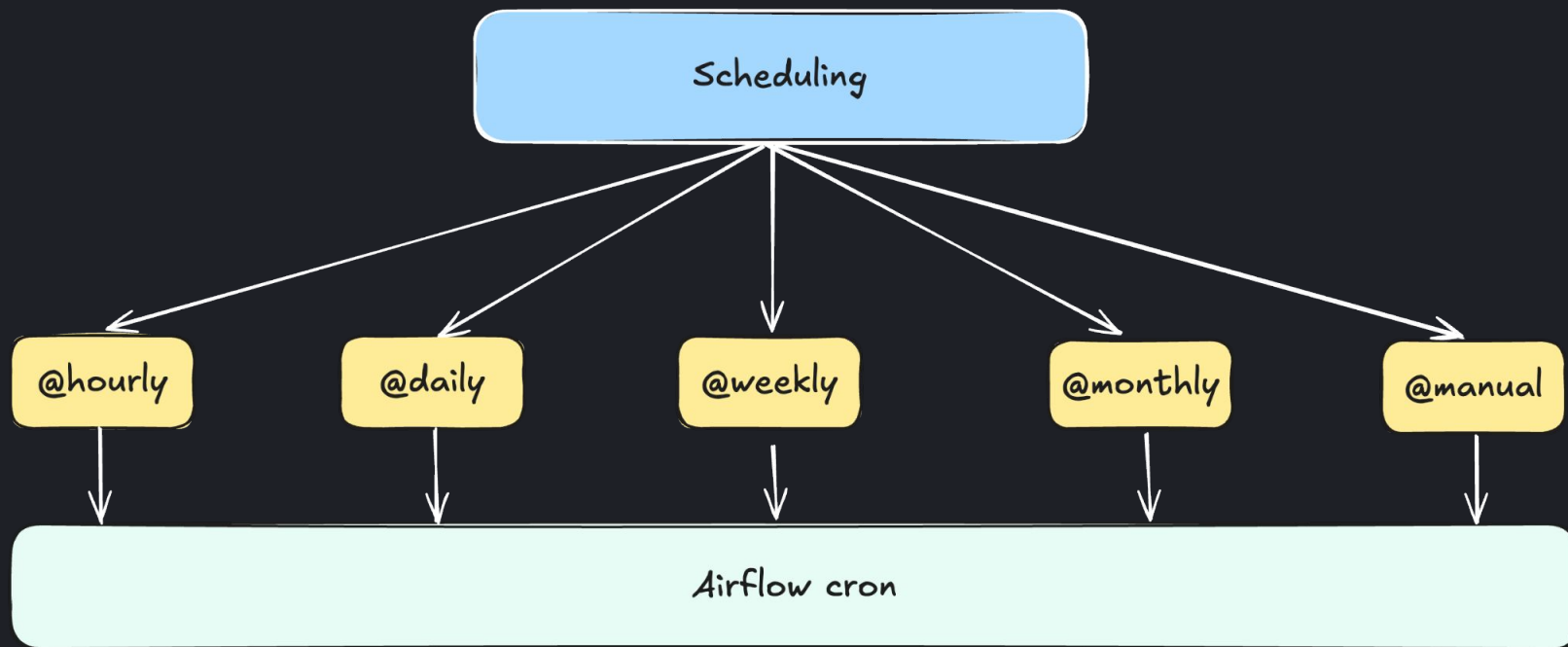
- Доменно-ориентированность (Data Mesh)
- Интеграции с другими системами (Airbyte, Monte Carlo, ...)



Велосипеды

	astronomer-cosmos	dbt-airflow	dbt Cloud	Dagster
Авто-генерация DAG'ов	+/-	+/-	-	++
Идемпотентность запусков	+	+	+	+
Разные уровни шедулинга	-	-	-	-/+
Управление нагрузкой, сценарии тестов	-	-	-	-/+

Шедулинг (от @monthly до @hourly)



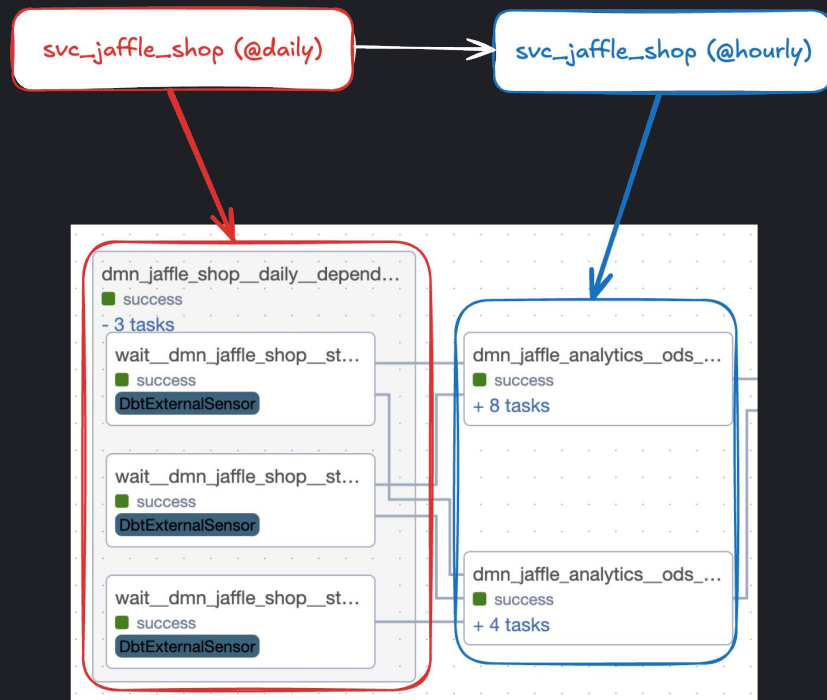
Кросс-доменные зависимости

- Зависимости указаны через функцию `ref`
- Не создается циклический граф из DAG'ов



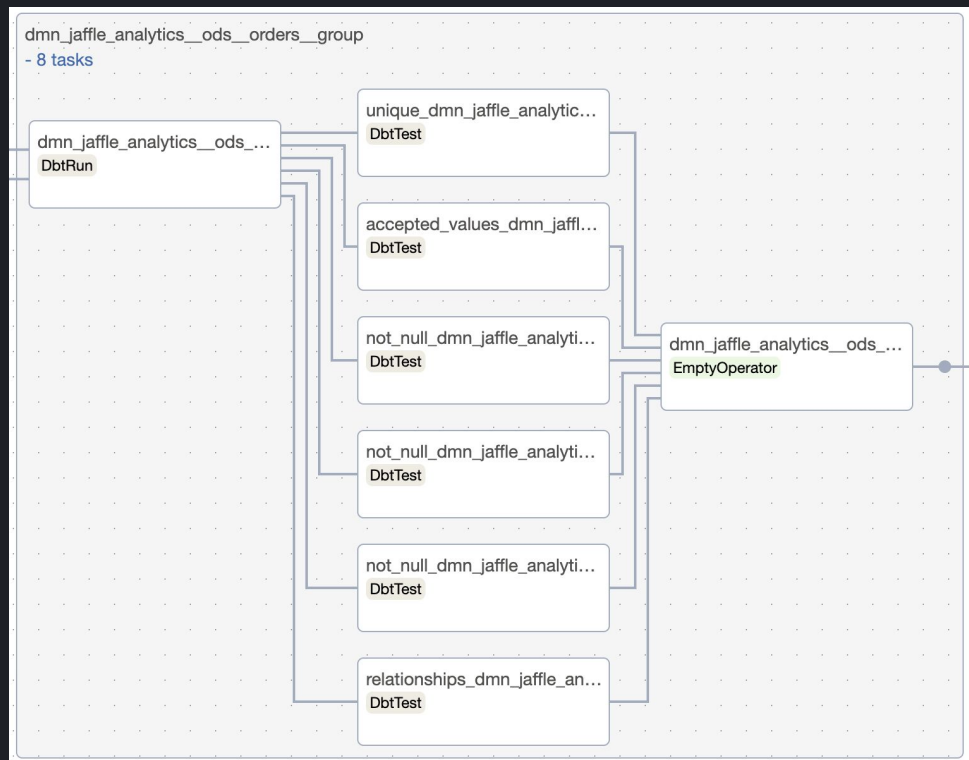
Кросс-доменные зависимости

- Зависимости указаны через функцию `ref`
- Не создается циклический граф из DAG'ов



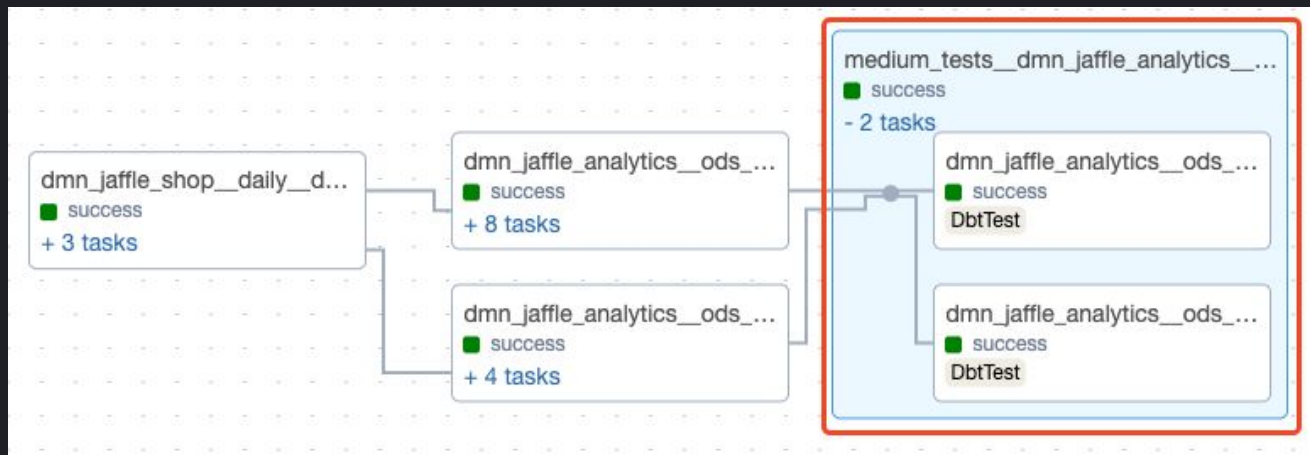
3 сценария запуска тестов

- @small
 - легкие
 - блокируют downstream
- @medium
- @large



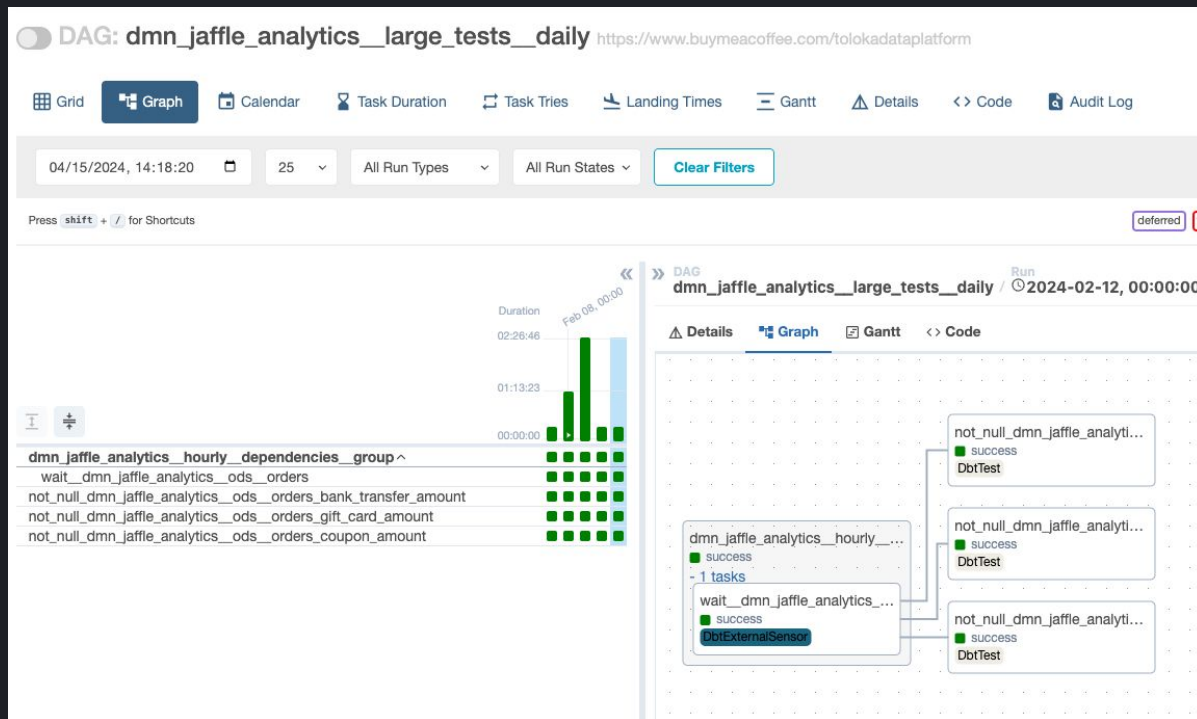
3 сценария запуска тестов

- @small
- @medium
 - Запускаются в конце DAG Run
- @large



3 сценария запуска тестов

- @small
- @medium
- @large
 - **особо тяжелые тесты**
 - **отдельный DAG**

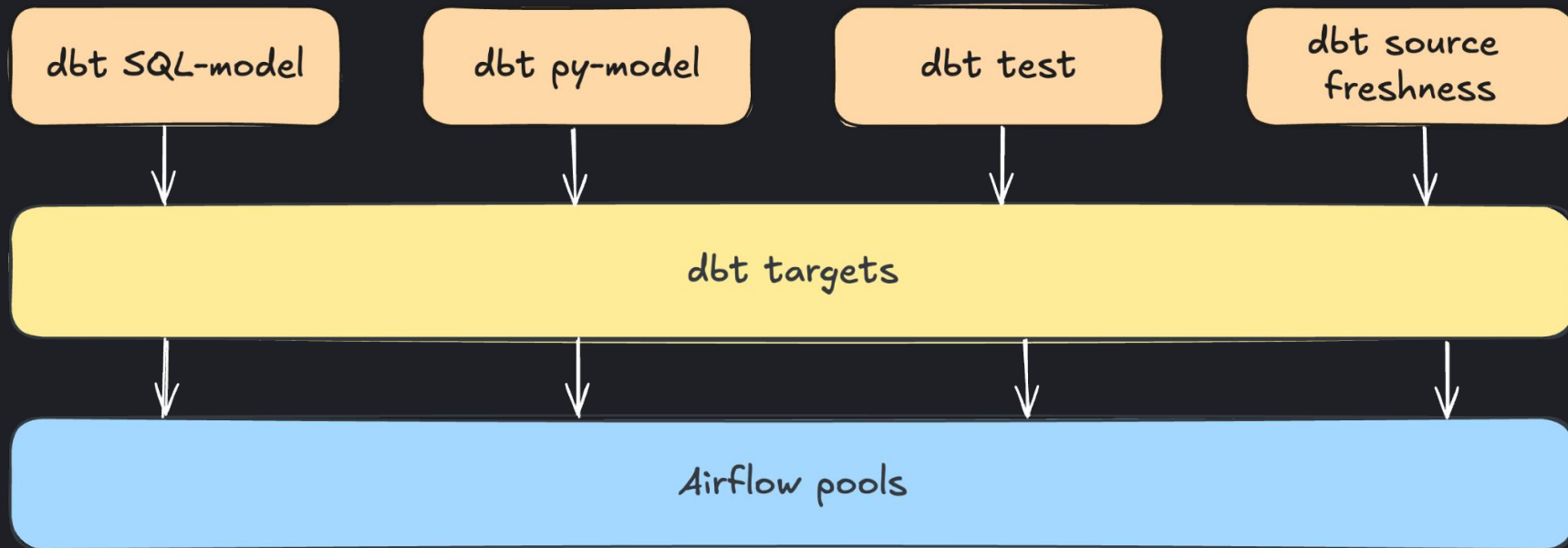


Дополнительные фичи

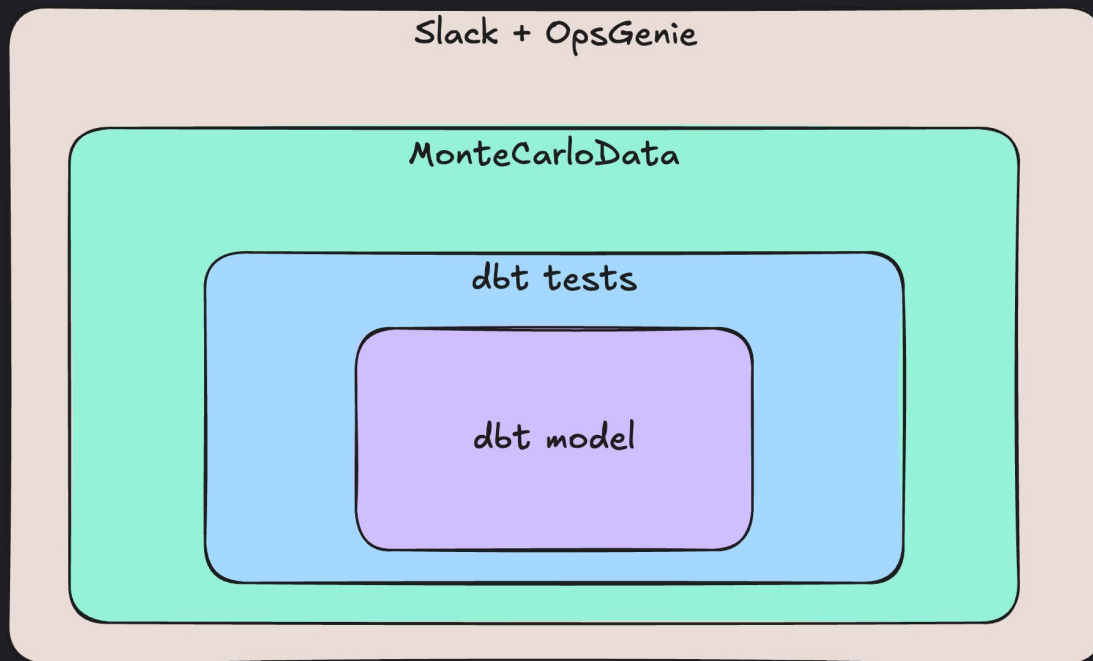
- explicit dbt targets
- source freshness
- Интеграция с другими системами (Monte Carlo, Tableau)
- произвольные k8s dbt-python-like таски
- minidbt (оптимизация компиляции dbt проекта)

[TODO] dbt-af в дикой природе

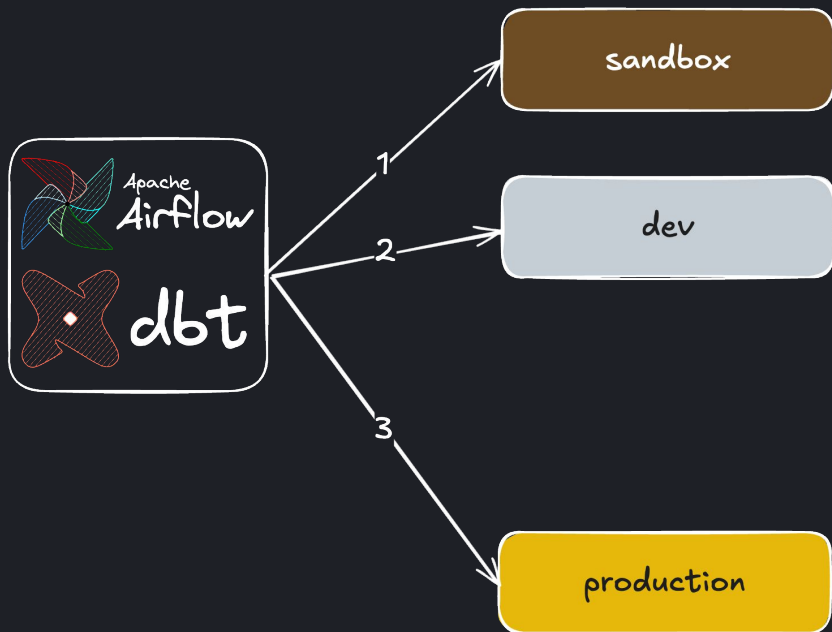
С большим продом приходит большая ответственность



[TODO] DQ: как мы узнаем, что оно развалилось



DX дата инженера

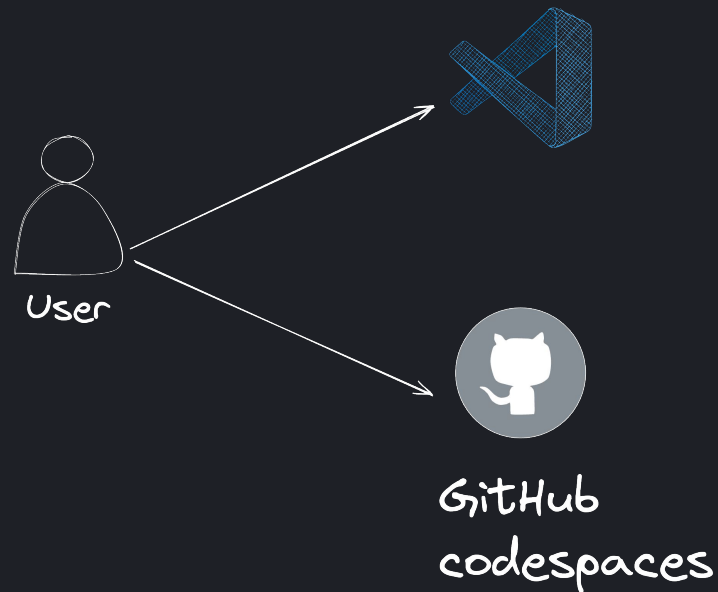


```
# LABELS: dag, airflow (it's required for airflow dag-processor)
from dbt_af.dags import compile_dbt_af_dags
from dbt_af.conf import Config, DbtDefaultTargetsConfig,
DbtProjectConfig
# specify here all settings for your dbt project
config = Config(
    dbt_project=DbtProjectConfig(
        dbt_project_name='my_dbt_project',
        dbt_project_path='/path/to/my_dbt_project',
        dbt_models_path='/path/to/my_dbt_project/models',
        dbt_profiles_path='/path/to/my_dbt_project',
        dbt_target_path='/path/to/my_dbt_project/target',
        dbt_log_path='/path/to/my_dbt_project/logs',
        dbt_schema='my_dbt_schema',
    ),
    dbt_default_targets=DbtDefaultTargetsConfig(default_target='dev'),
    is_dev=False, # set to True if you want to turn on dry-run mode
)

dags = compile_dbt_af_dags(
    manifest_path='/path/to/my_dbt_project/target/manifest.json',
    config=config,
)
for dag_name, dag in dags.items():
    globals()[dag_name] = dag
```

DX аналитика / ML researcher

```
models:  
  - name: "domain.schema.table"  
    description: "Some meaningful description"  
    config:  
      schedule: "@hourly"  
      dbt_target: "nightly_target"  
      dependencies:  
        other_domain.other_schema.other_table:  
          skip: true
```



[TODO] Про тесты, CI/CD для DX?

Глава III. Заключительная



Немного статистики

- 32 доменов
- ~8900 запусков dbt run в сутки
- ~310 уникальных связей между доменами
- ~320 тестов
- 59к строк sql-кода



Наши боли



Analysts



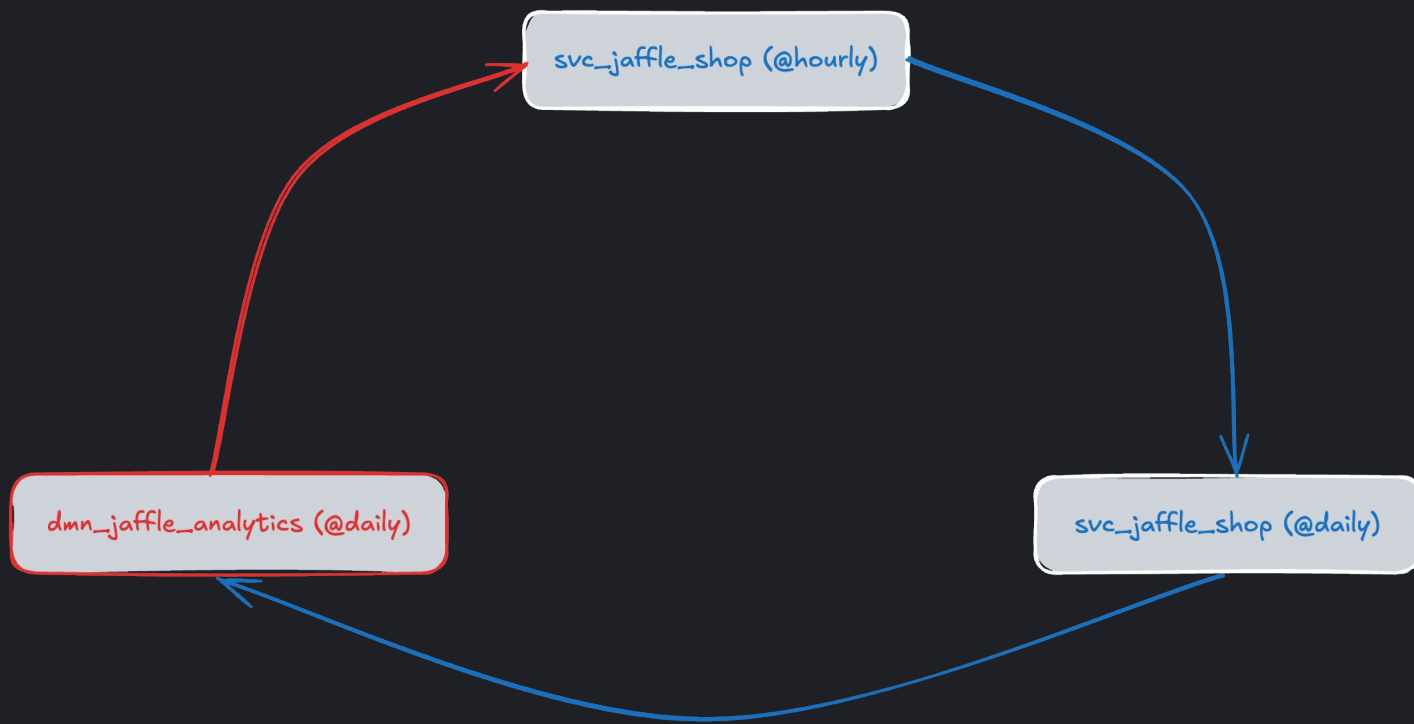
Researchers



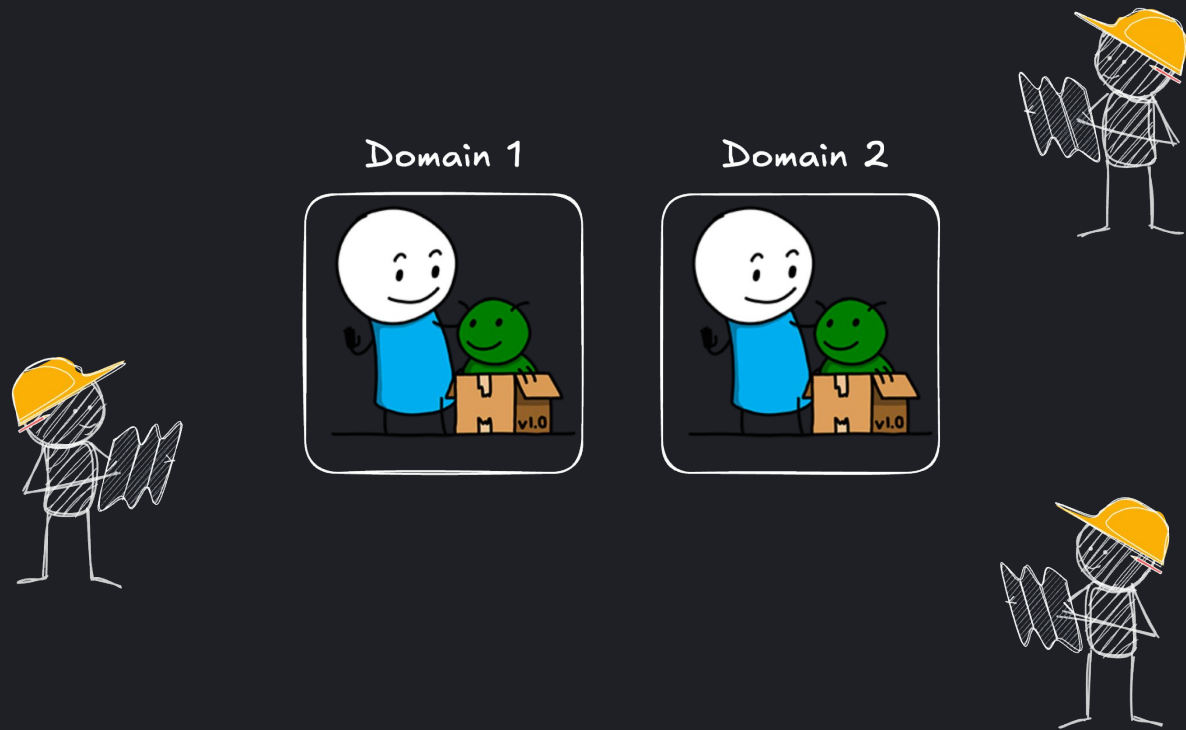
Data Engineers

Боль №0: на сколько сильно дробить предметную область на домены?

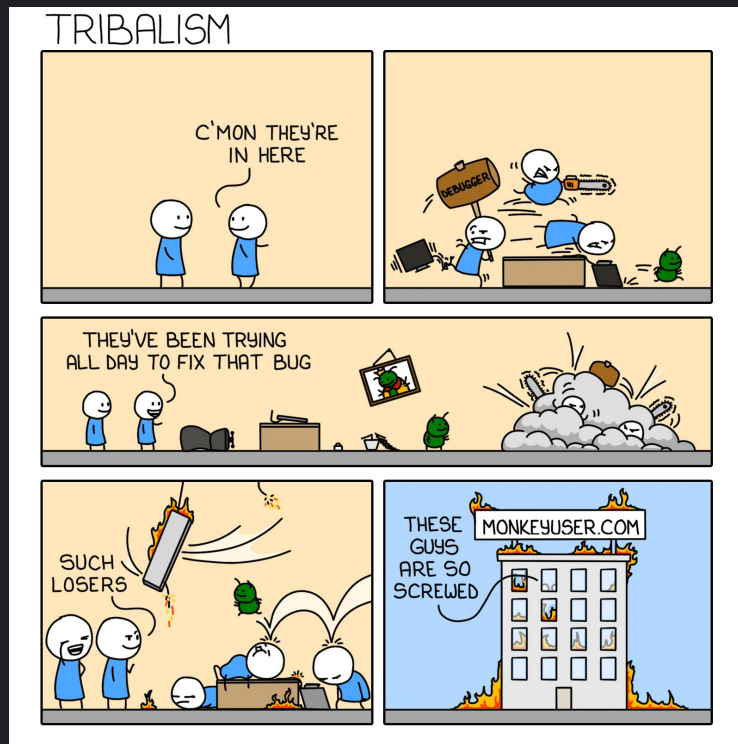
Боль №1: циклические зависимости



Боль №2: проблемные домены



Боль №3: восстановление после инцидентов



Плюсы Data Mesh

- Возможность управлять большим хранилищем силами небольшой команды
- Observability
- Низкий порог входа для пользователей дата платформы



Минусы Data Mesh

- Data Mesh требует от пользователей понимания принципов архитектуры
- Data Mesh приносит дополнительную сложность
- dbt-af требует определенной структуры dbt проекта



[TODO] Выводы

- Data Mesh топ, обратно на монолит не хотим
- Инфра местами топ, но не совсем



@lkozhinov



@nikitayurasov



dbt-af

