

Оркестратор пайплайнов для небольшой команды инженеров и аналитиков: как мы используем dagster

Дмитрий Крылов
 mityakrylov

Алексей Завальский
 al_zav

Bestplace



Геоаналитическая платформа

- Геоданные (координаты, геометрии)
- Аналитика (визуализация, обработка)
- Машинное обучение



Компания

- Основана в 2016 году
- Является одним из лидеров рынка геоаналитики в России



Клиенты и партнеры

- Ритейл
- FMCG
- Сервисы поиска недвижимости



Контекст: данные

Различные источники:

- парсинг веба
- государственные открытые данные
- партнерские и платные апи
- OpenStreetMap

Ну и что ты тут накачал?



Контекст: данные

Различные источники:

- парсинг веба
- государственные открытые данные
- партнерские и платные апи
- OpenStreetMap

Храним в mongodb и S3

Ну и что ты тут накачал?



Контекст: данные

Различные источники:

- парсинг веба
- государственные открытые данные
- партнерские и платные api
- OpenStreetMap

Храним в mongodb и S3

Данных много разных по смыслу (домохозяйства, точки интереса, трафики и т.д.), но немного по объему (десятки Тб)

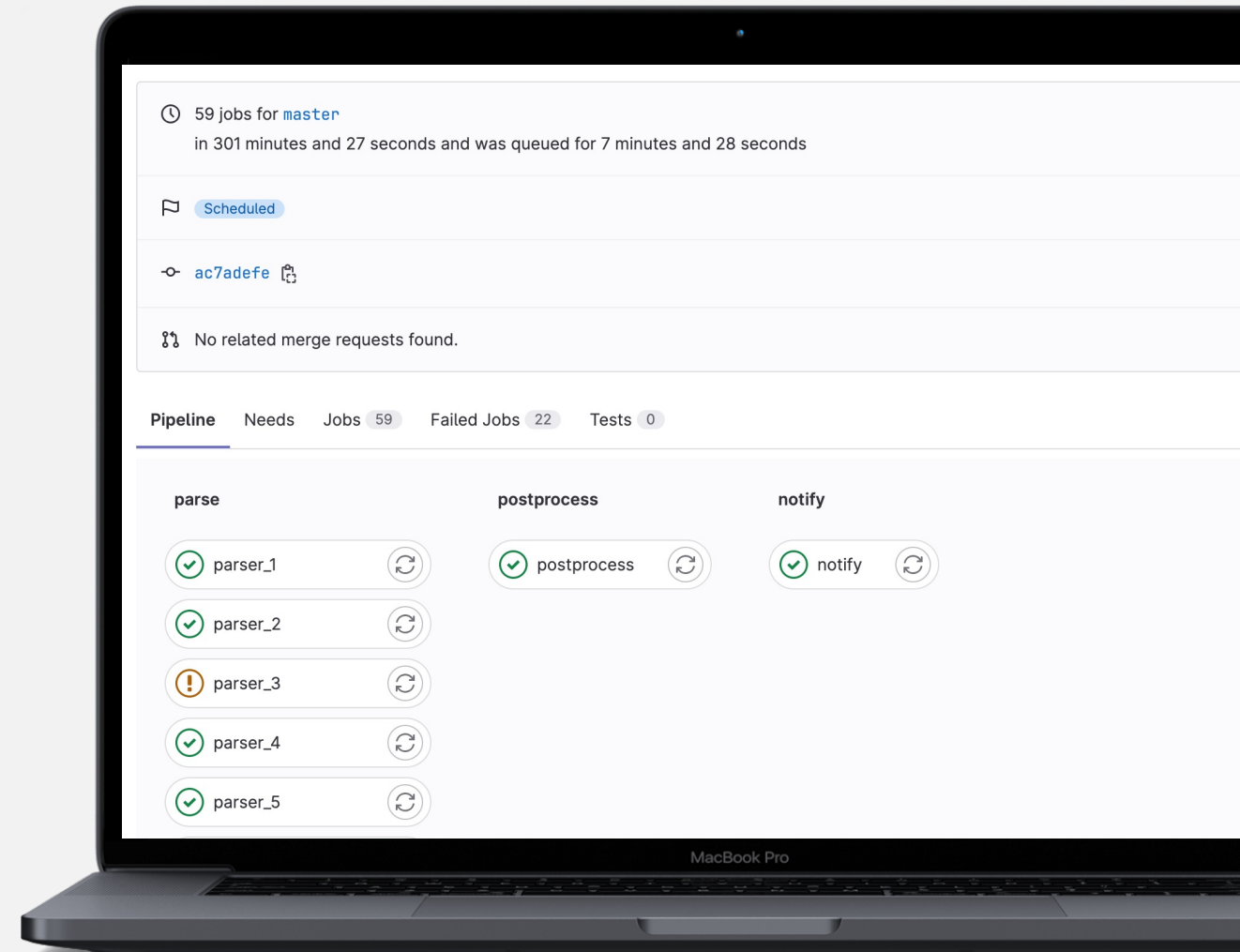
Ну и что ты тут накачал?



Контекст: пайплайны

Часть пайплайнов писали разработчики в gitlab ci:

- запуск по расписанию
- простые графы зависимостей

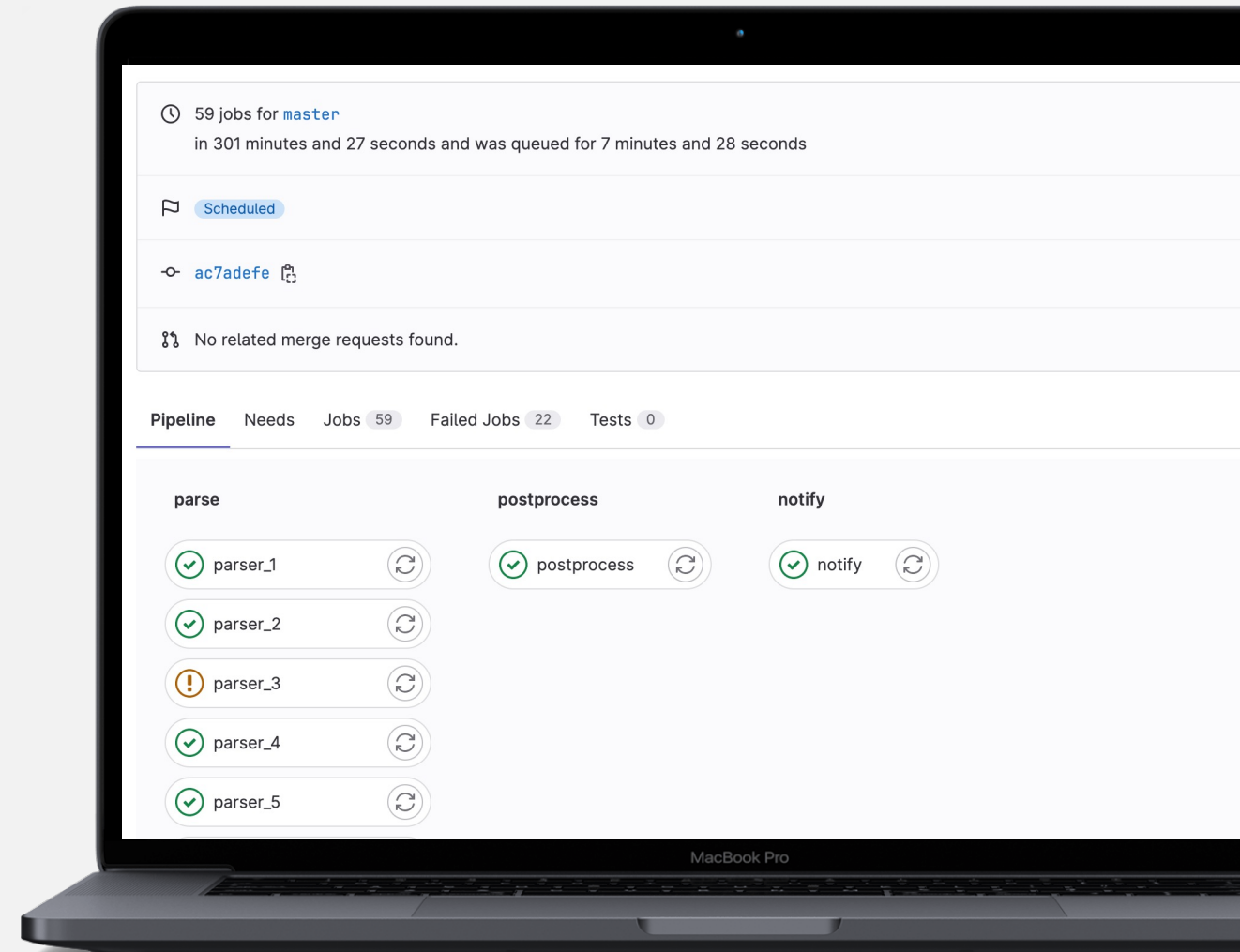


Контекст: пайплайны

Часть пайплайнов писали разработчики в gitlab ci:

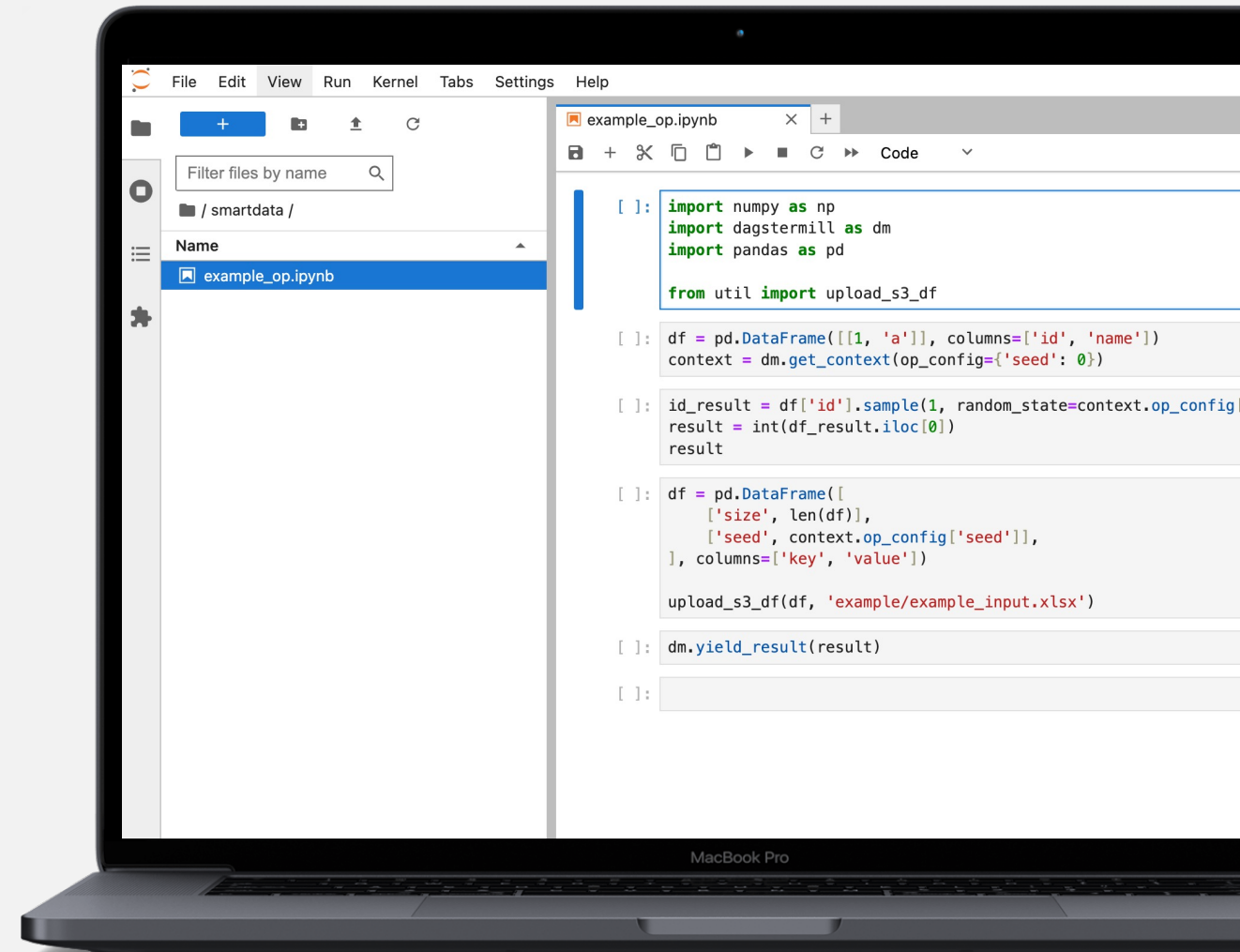
- запуск по расписанию
- простые графы зависимостей

Проблемы: неудобно следить и параметризировать



Контекст: пайплайны

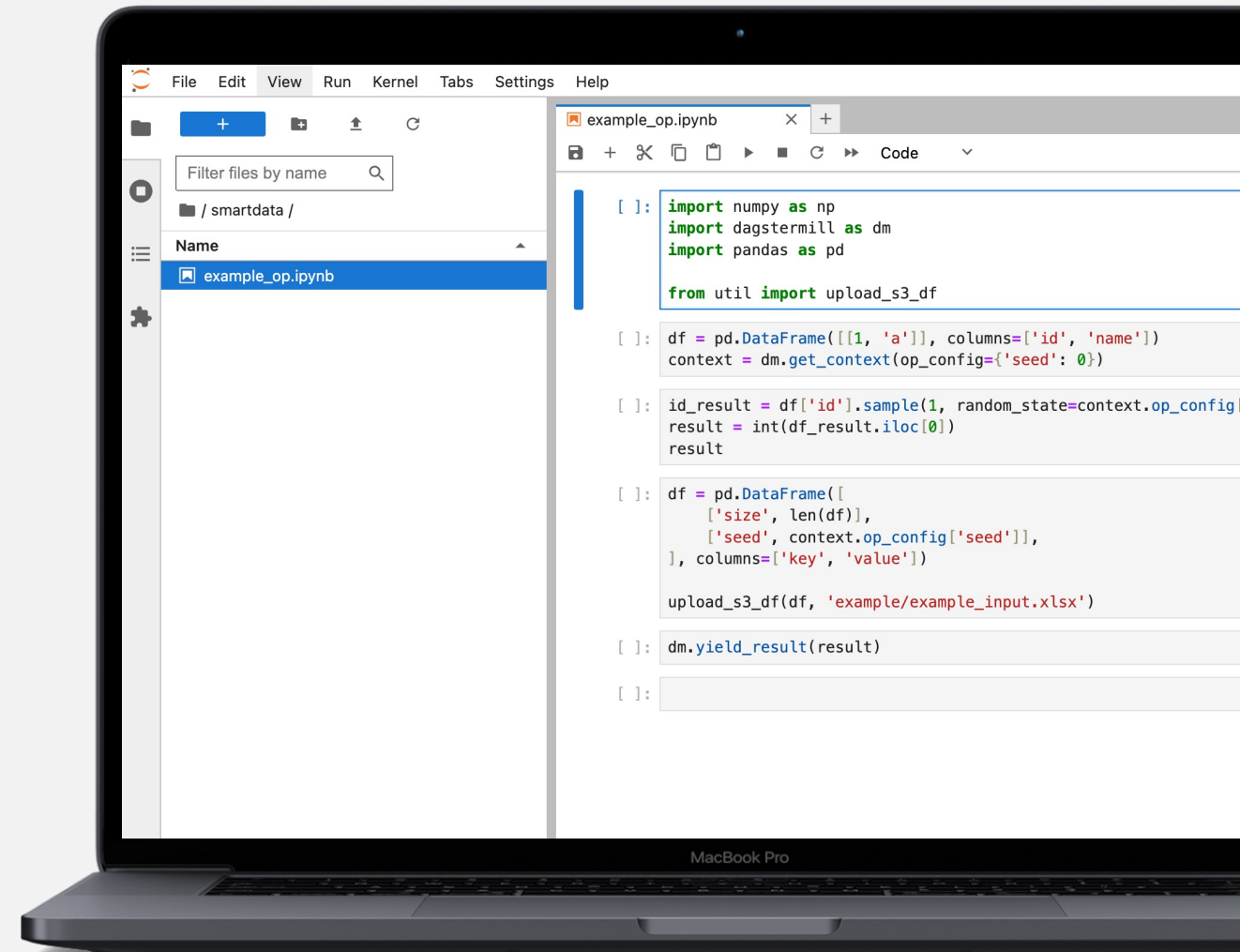
Часть пайплайнов писали аналитики и дата саентисты в jupyter ноутбуках



Контекст: пайплайны

Часть пайплайнов писали аналитики и дата саентисты в jupyter ноутбуках

Проблемы: неудобно воспроизводить



Выбор инструмента, 2020 год

Пожелания:

- Несложно разобраться и развернуть
- Параметризируемый запуск из интерфейса
- Интеграция с jupyter
- Симпатичный UI



Dagster

Что понравилось:

- Соответствовал всем пожеланиям
- Хорошая документация
- Скорость обновлений

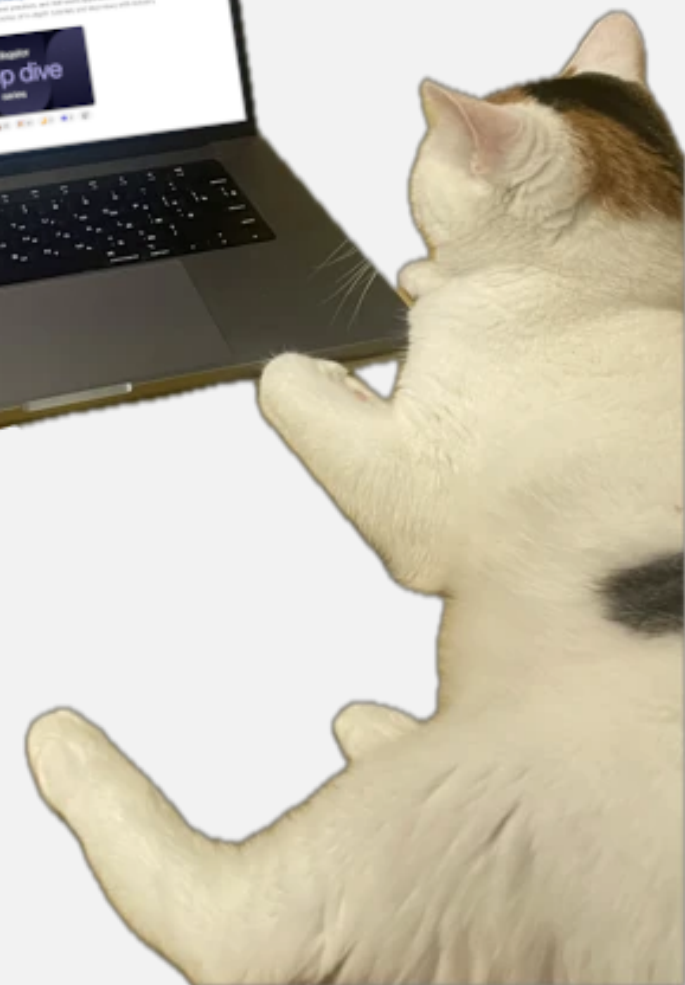


Коммуникация с командой dagster

github, slack, созвоны



*Ого, даже
отдельный канал?*

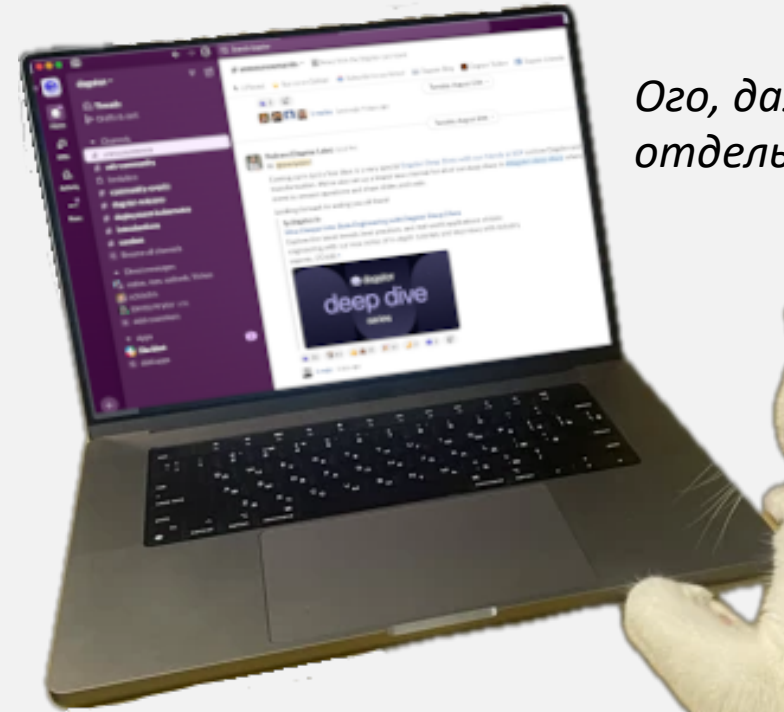


Коммуникация с командой dagster

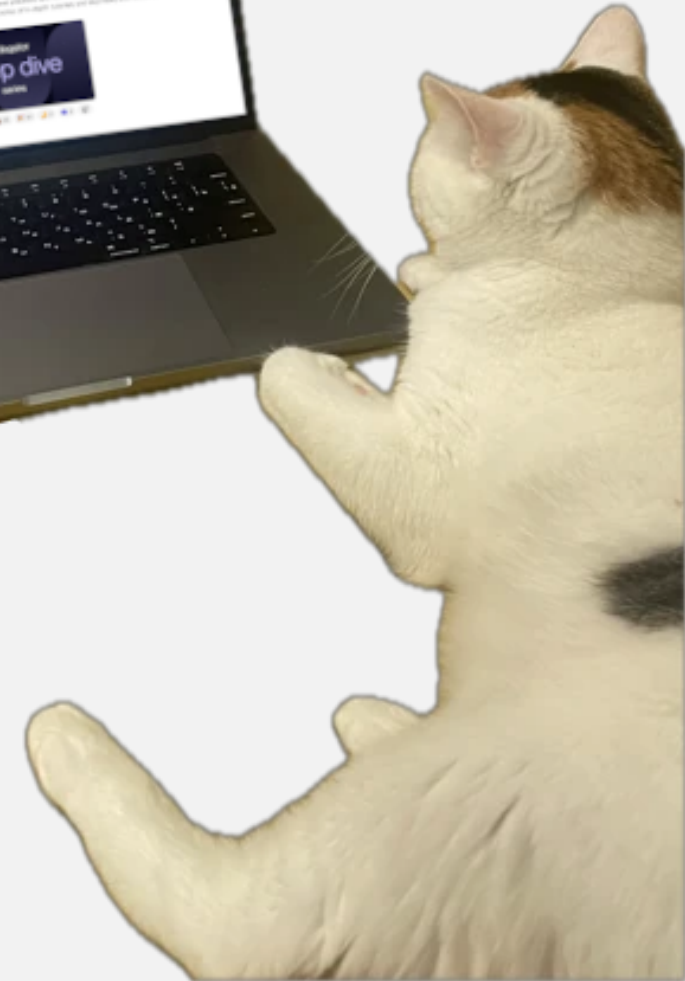
github, slack, созвоны

Совместные улучшения:

- docker in docker executor
- фикс взаимодействия с minio
- фиксы взаимодействия с jupyter
- и другие



*Ого, даже
отдельный канал?*

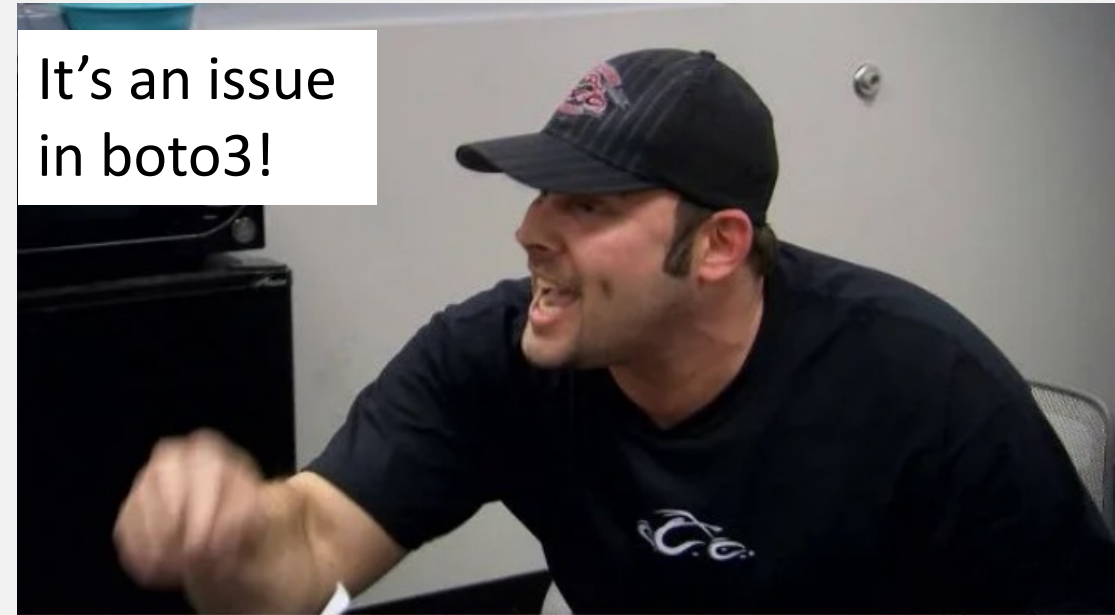


Коммуникация с командой dagster

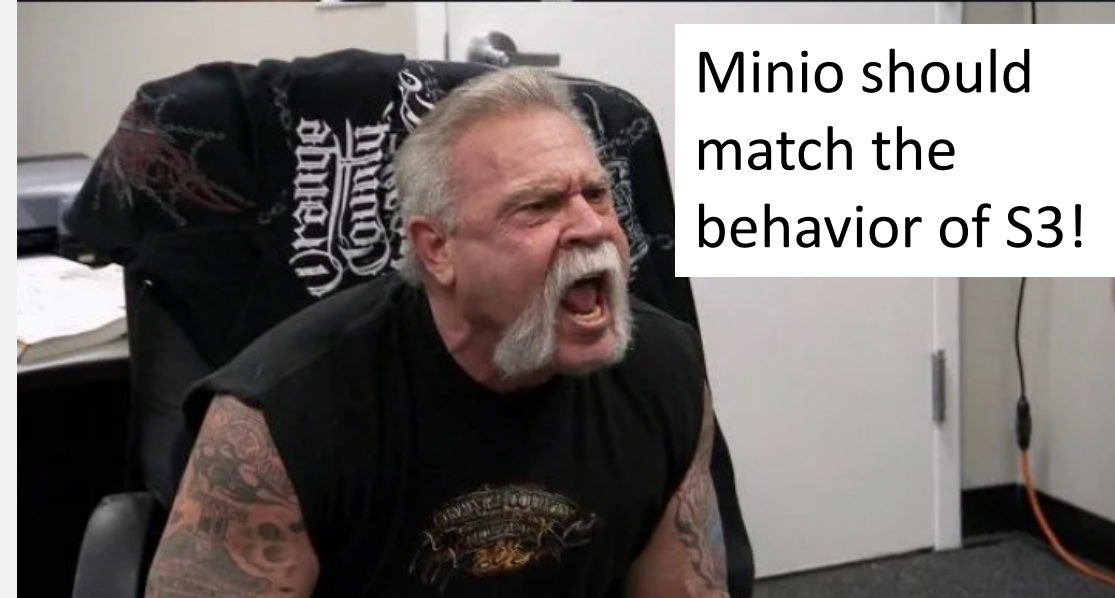
github, slack, созвоны

Совместные улучшения:

- docker in docker executor
- фикс взаимодействия с minio
- фиксы взаимодействия с jupyter
- и другие



It's an issue
in boto3!



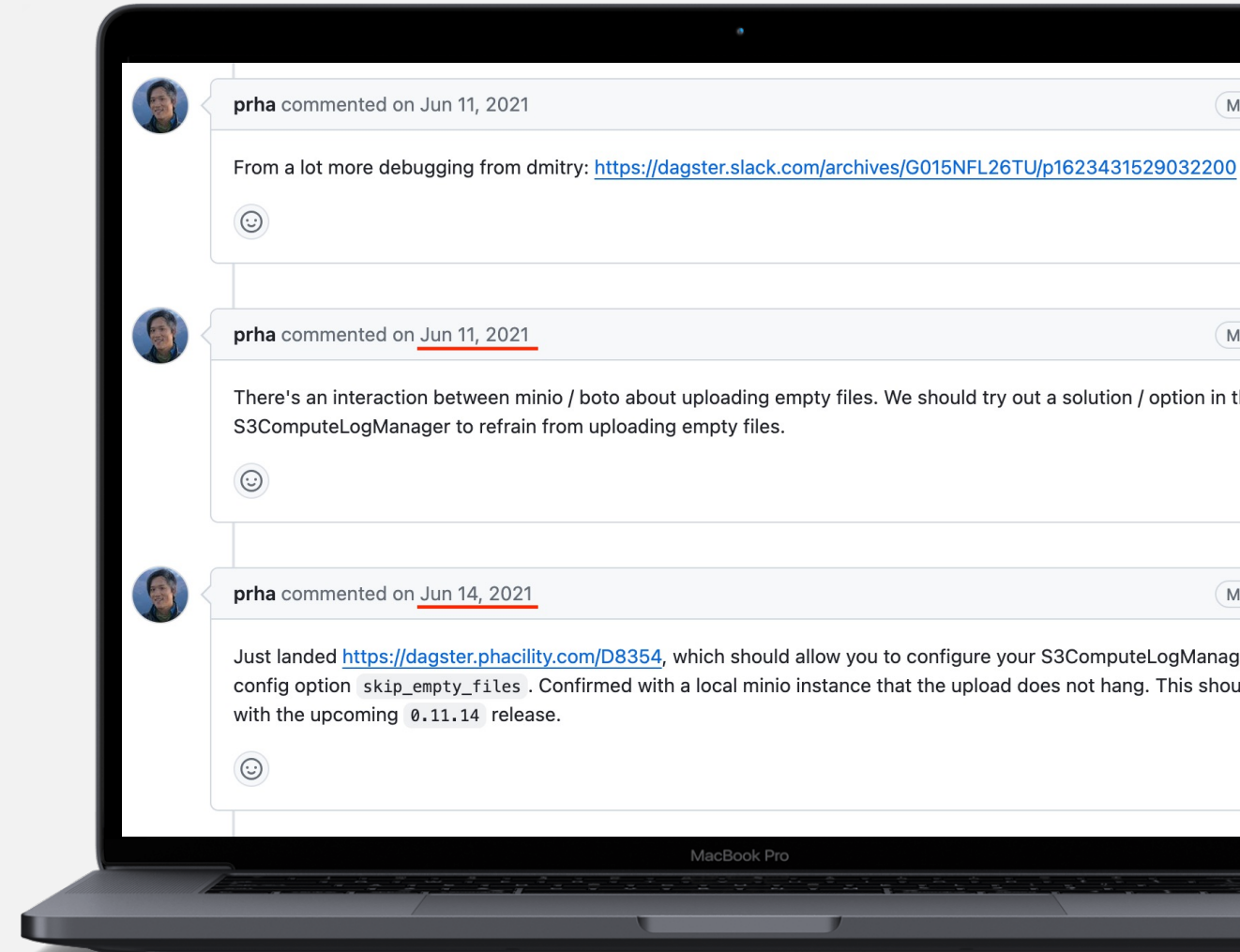
Minio should
match the
behavior of S3!

Коммуникация с командой dagster

github, slack, созвоны

Совместные улучшения:

- docker in docker executor
- фикс взаимодействия с minio
- фиксы взаимодействия с jupyter
- и другие



Коммуникация с командой dagster

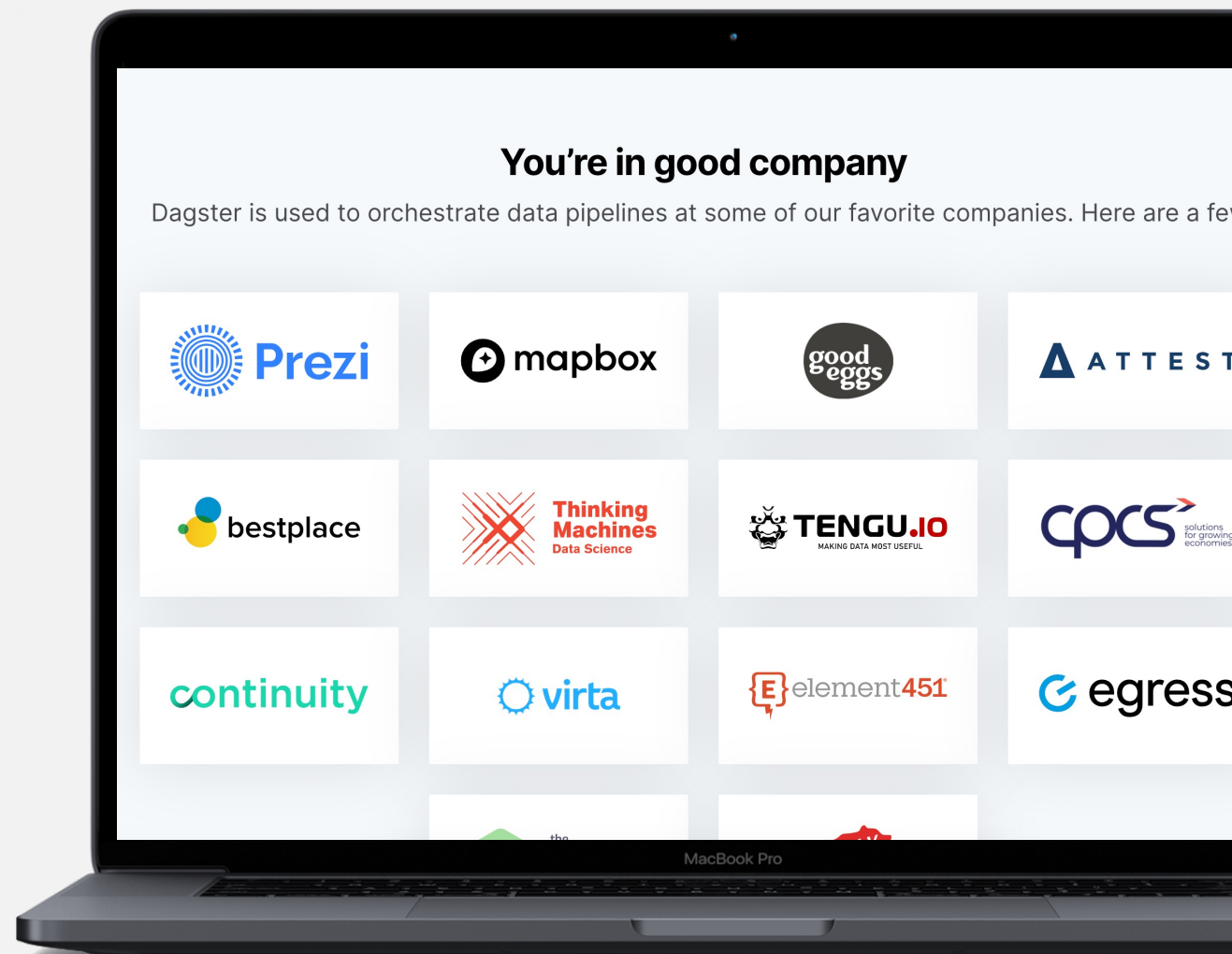
github, slack, созвоны

Совместные улучшения:

- docker in docker executor
- фикс взаимодействия с minio
- фиксы взаимодействия с jupyter
- и другие

Релизы в течение недели после исправлений.

Взаимодействуйте с open source командами!



Коммуникация с командой dagster

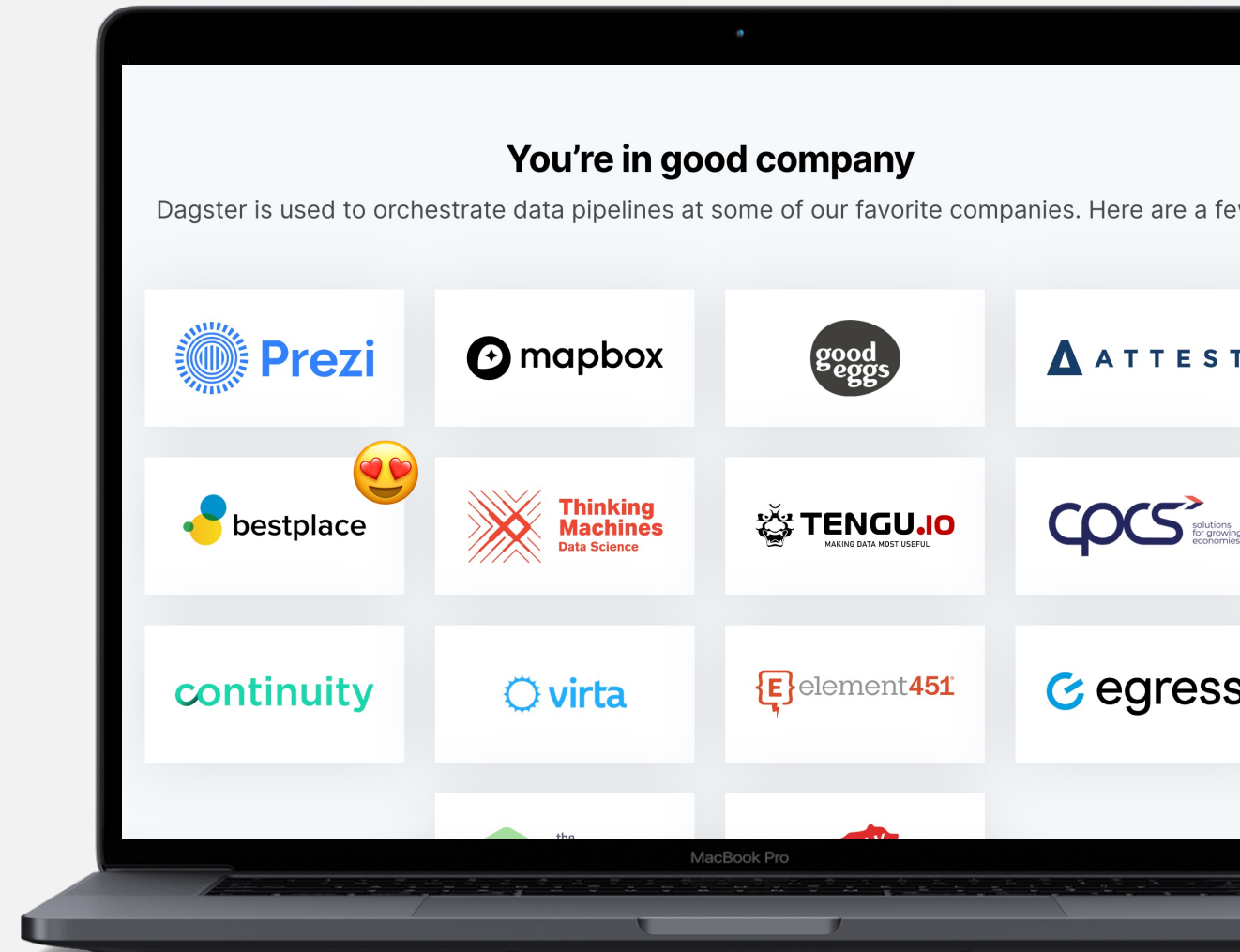
github, slack, созвоны

Совместные улучшения:

- docker in docker executor
- фикс взаимодействия с minio
- фиксы взаимодействия с jupyter
- и другие

Релизы в течение недели после исправлений.

Взаимодействуйте с open source командами!



Результат внедрения

Результат для дата-инженеров

Фреймворк для создания пайплайнов:

- op
- graph
- job
- schedule

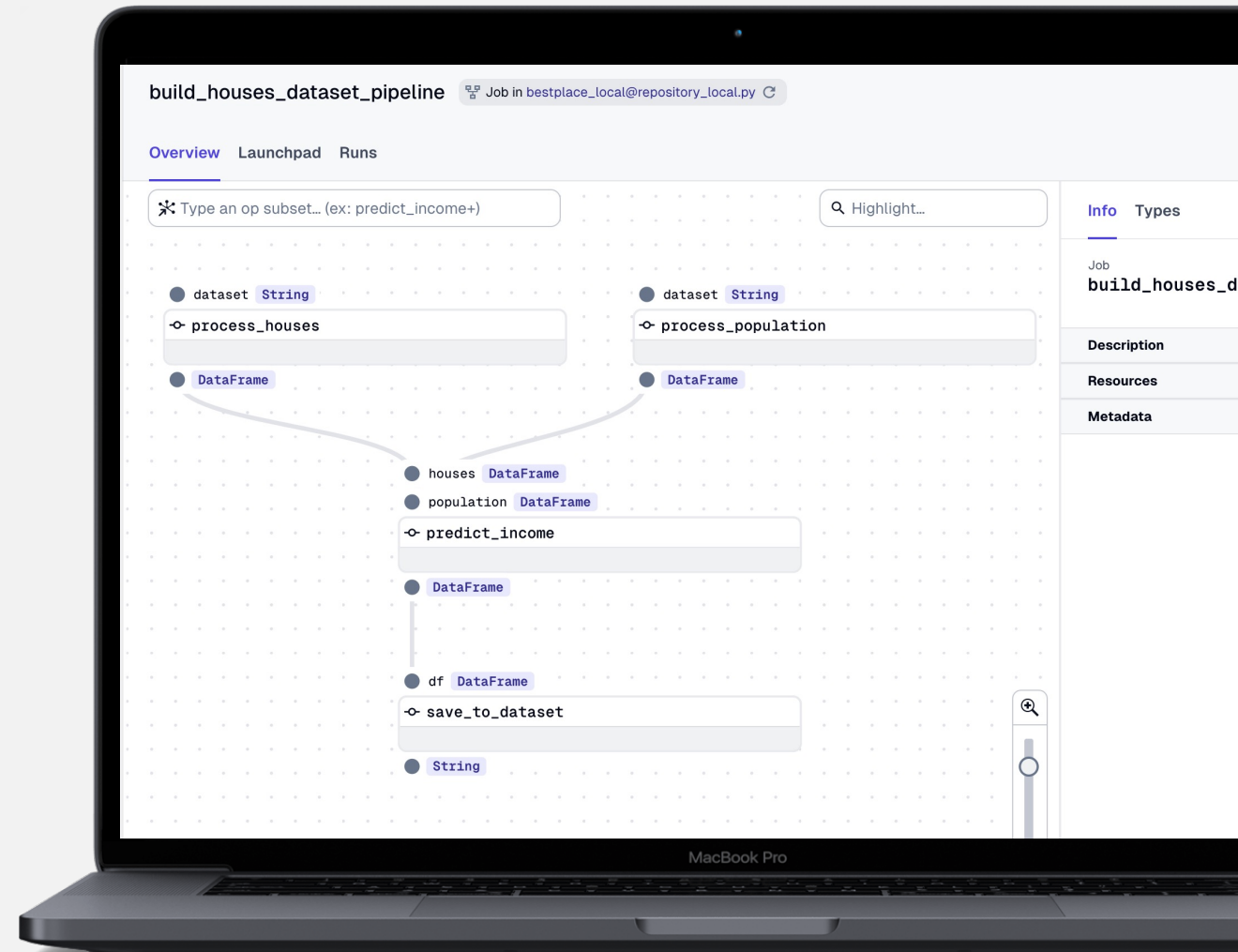
```
1  import pandas as pd
2  from dagster import op, graph
3
4  @op
5  def process_houses(dataset: str) -> pd.DataFrame: ...
10
11  @op
12  def process_population(dataset: str) -> pd.DataFrame: ...
16
17  @op
18  def predict_income(houses: pd.DataFrame, population: pd.DataFrame) -> pd.DataFrame: ...
22
23  @op
24  def save_to_dataset(df: pd.DataFrame) -> str: ...
27
28
29  @graph
30  def build_houses_dataset_pipeline():
31      dataset = predict_income(process_houses(), process_population())
32      save_to_dataset(dataset)
33
34
```

Результат для дата-инженеров

Фреймворк для создания пайплайнов:

- op
- graph
- job
- schedule

Удобное отслеживание статуса и логов запусков



Результат для аналитиков

Пайплайны из jupyter ноутбуков с помощью dagstermill (papermill)

```
1 import pandas as pd
2 import dagstermill as dm
3 from dagster import In, Out, Any
4
5
6 process_dataset_op = dm.define_dagstermill_op(
7     name='process_dataset',
8     notebook_path='notebooks/process_dataset_op.ipynb',
9     ins={
10         'df': In(pd.DataFrame, description='Input dataframe')
11     },
12     outs={
13         'result': Out(Any)
14     },
15 )
16
17
18
19
20
```

Результат для аналитиков

Пайплайны из jupyter ноутбуков с помощью dagstermill (papermill)

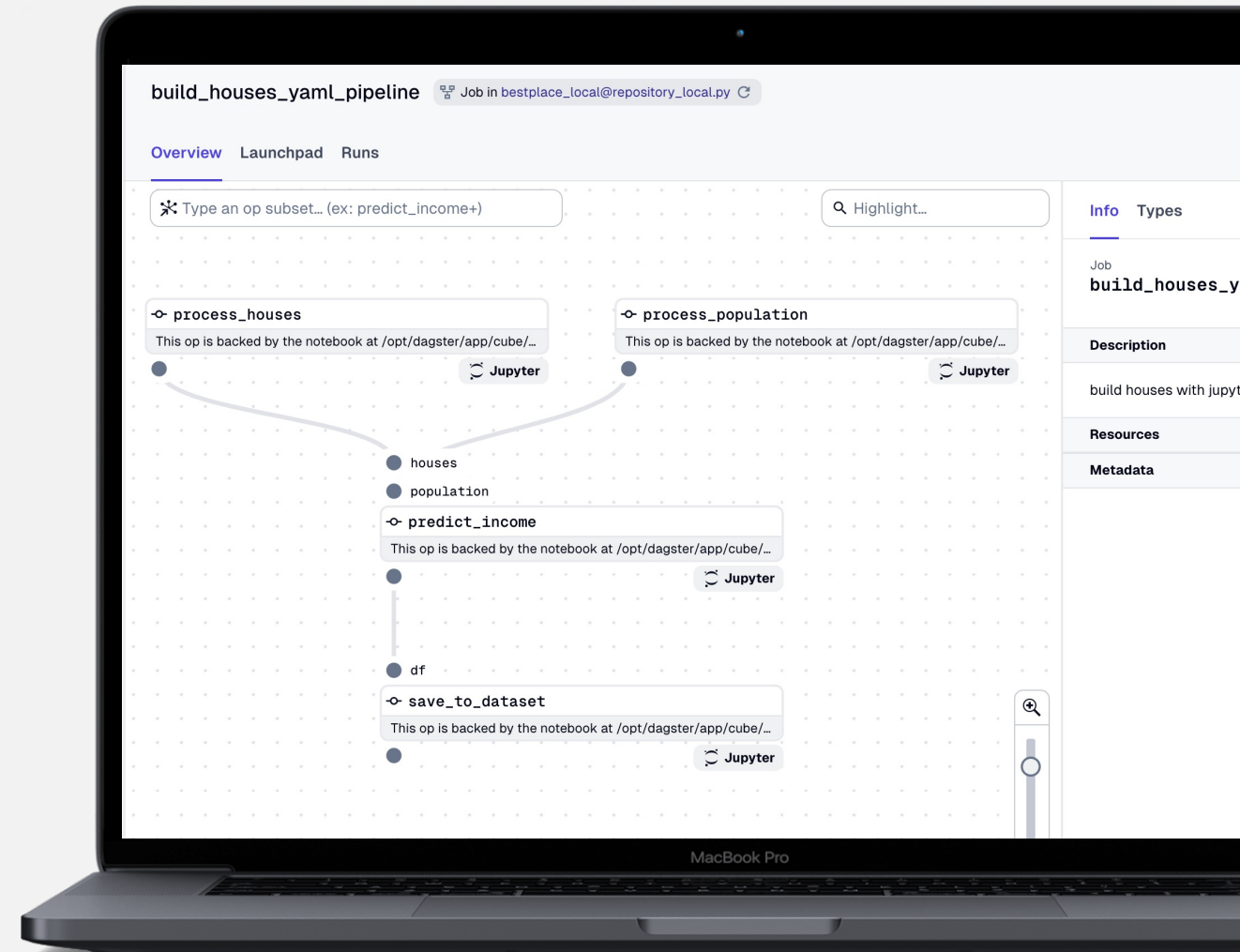
Надстройка над синтаксисом dagster: yaml пайплайны

```
1   name: build_houses_yaml_pipeline
2   description: build houses with jupyter notebooks as ops
3   ops:
4   - def: build_houses/process_houses.ipynb
5     alias: process_houses
6
7   - def: build_houses/process_population.ipynb
8     alias: process_population
9
10  - def: build_houses/predict_income.ipynb
11    alias: predict_income
12    deps:
13      houses: process_houses
14      population: process_population
15
16  - def: build_houses/save_to_dataset.ipynb
17    alias: save_to_dataset
18    deps:
19      df: predict_income
20
```

Результат для аналитиков

Пайплайны из jupyter ноутбуков с помощью dagstermill (papermill)

Настройка над синтаксисом dagster: yml пайплайны



Результат для аналитиков

Пайплайны из jupyter ноутбуков с помощью dagstermill (papermill)

Надстройка над синтаксисом dagster: уaml пайплайны

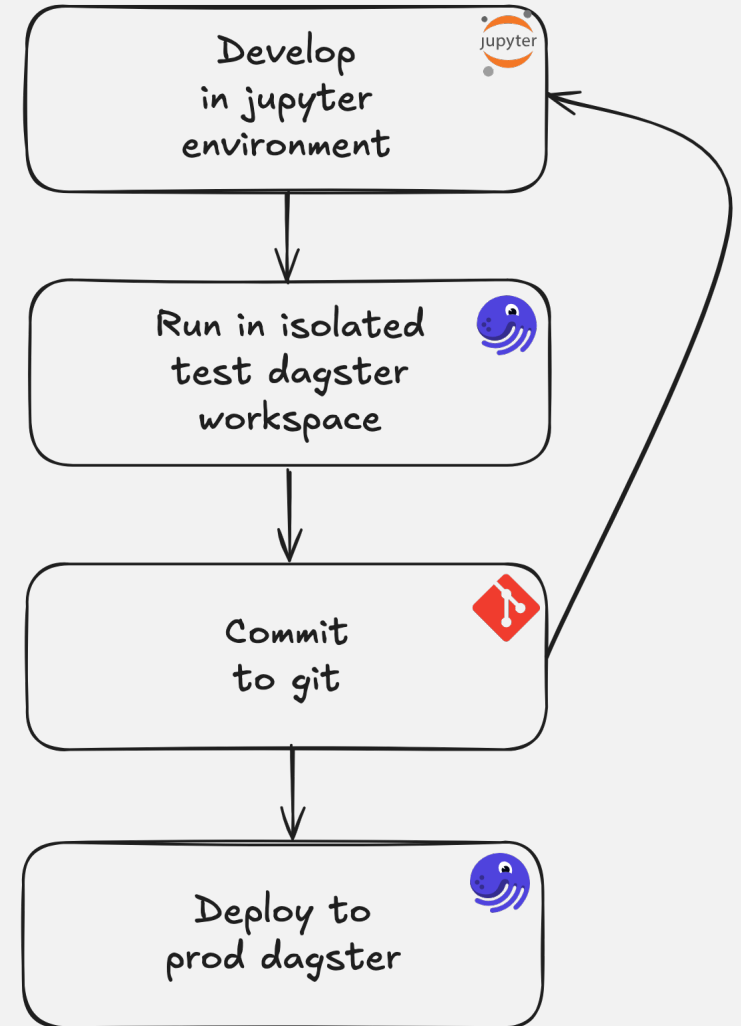


Результат для аналитиков

Пайплайны из jupyter ноутбуков с помощью dagstermill (papermill)

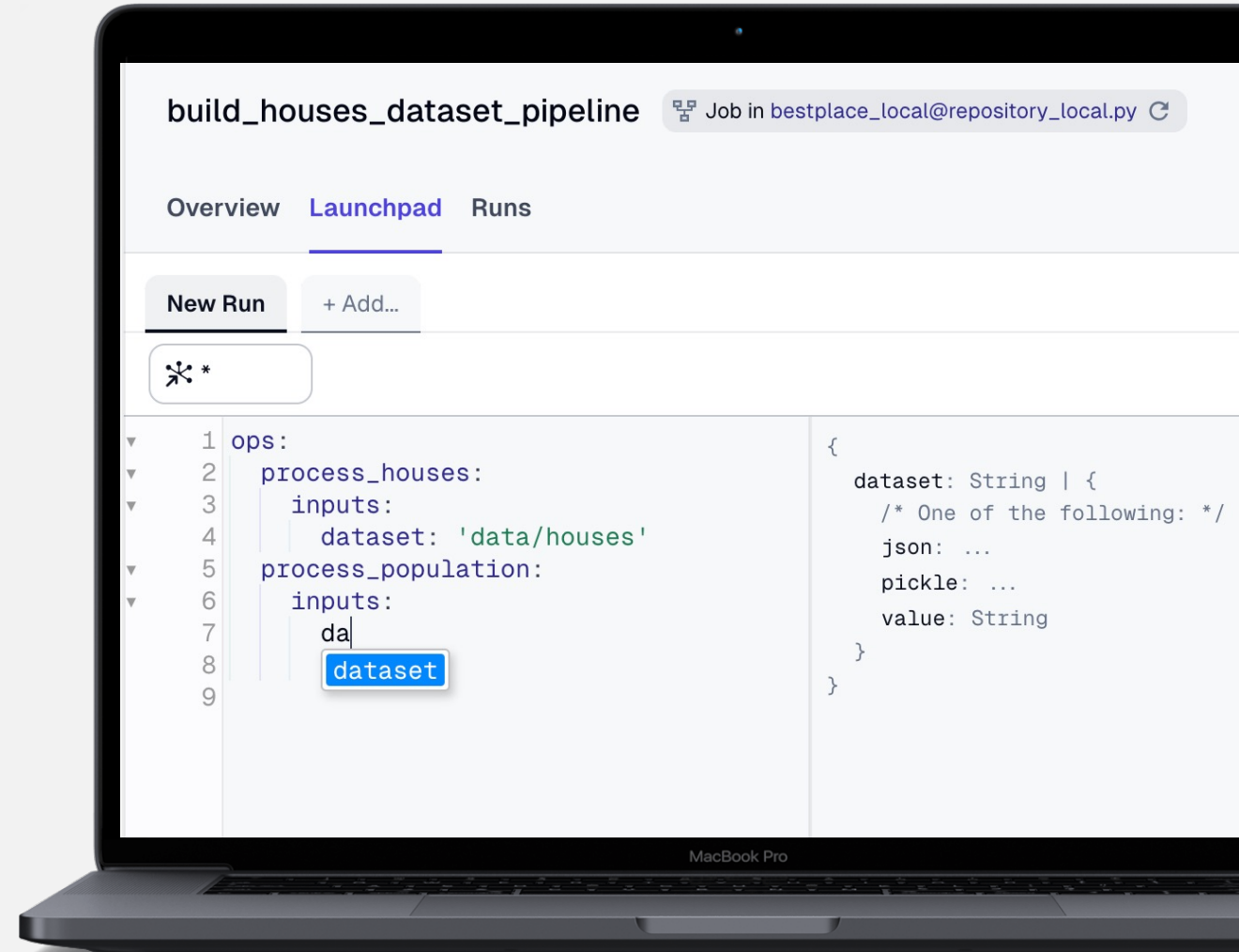
Настройка над синтаксисом dagster: yaml пайплайны

Тестовый контур



Результат для всех

Параметризируемый запуск по кнопке



Время историй

История про мониторинг запусков

История про мониторинг запусков

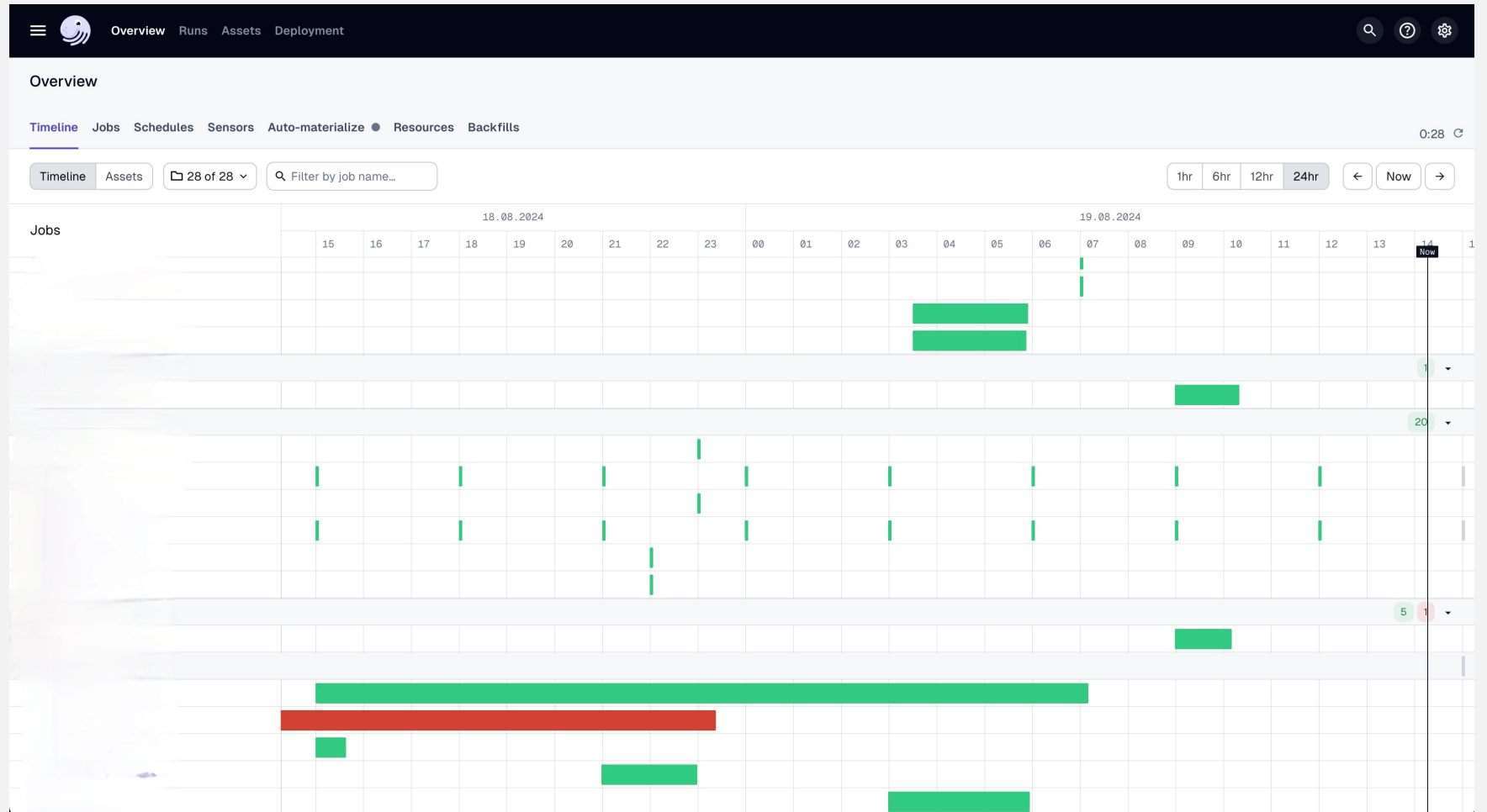
Пользователи дагстера:

- дата-инженеры
- аналитики
- **менеджеры**

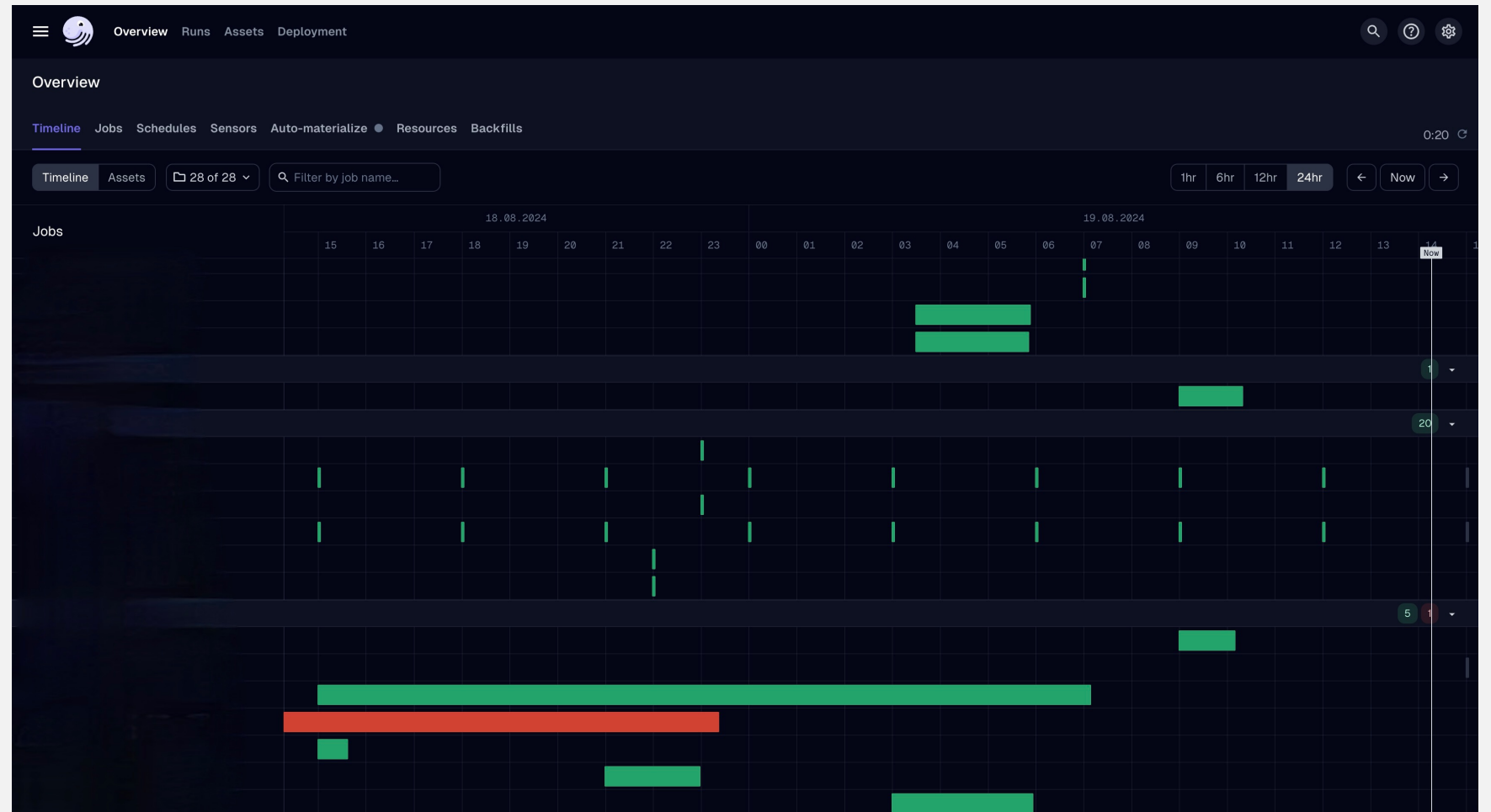
*Не вижу,
куда смотреть*



Мониторинг в UI



Мониторинг в UI



Мониторинг в telegram

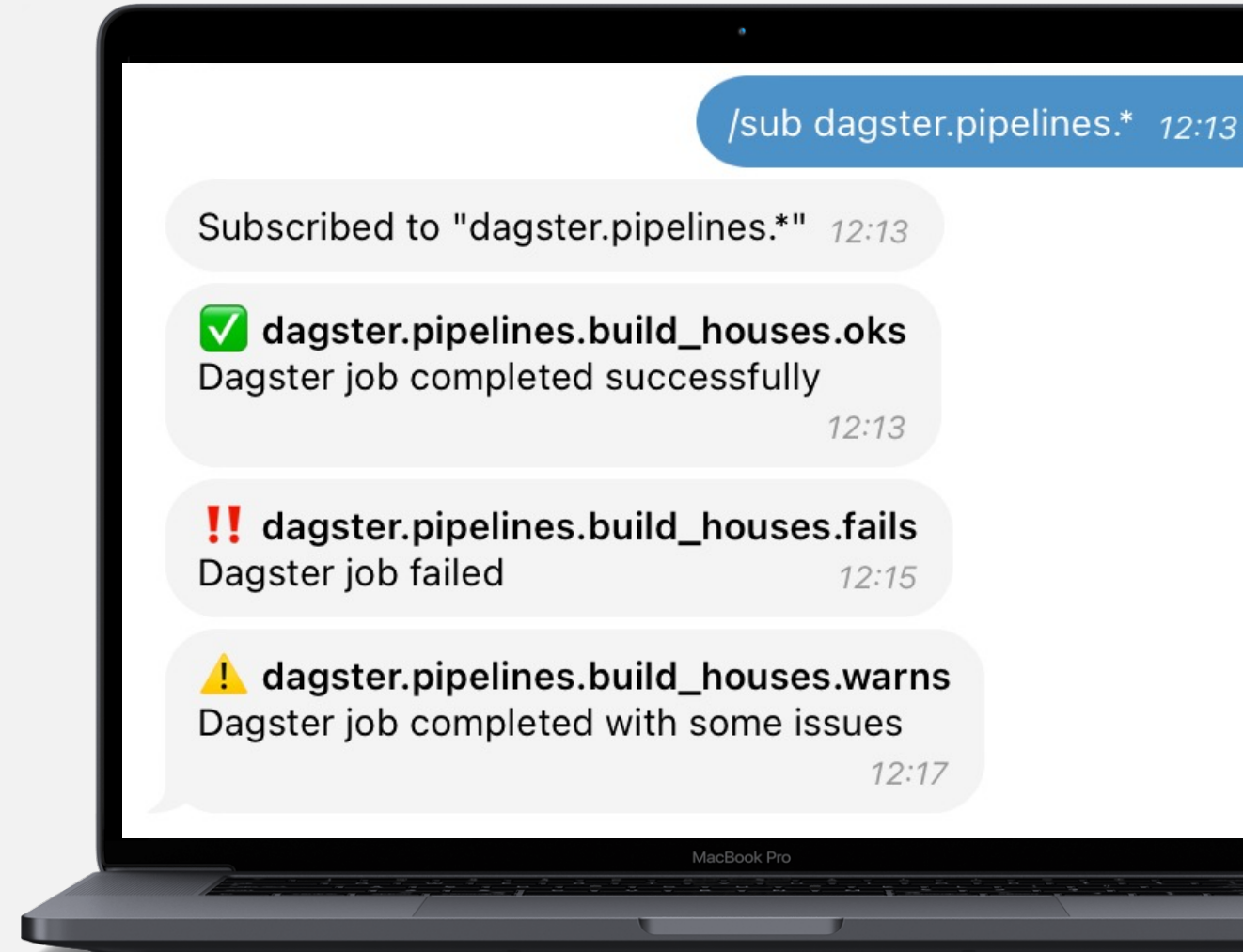
UI — это хорошо

Но хочется в телеграм



Мониторинг в telegram

TAAS (telegram as a service) — publish-subscribe сервис для отправки уведомлений через телеграм



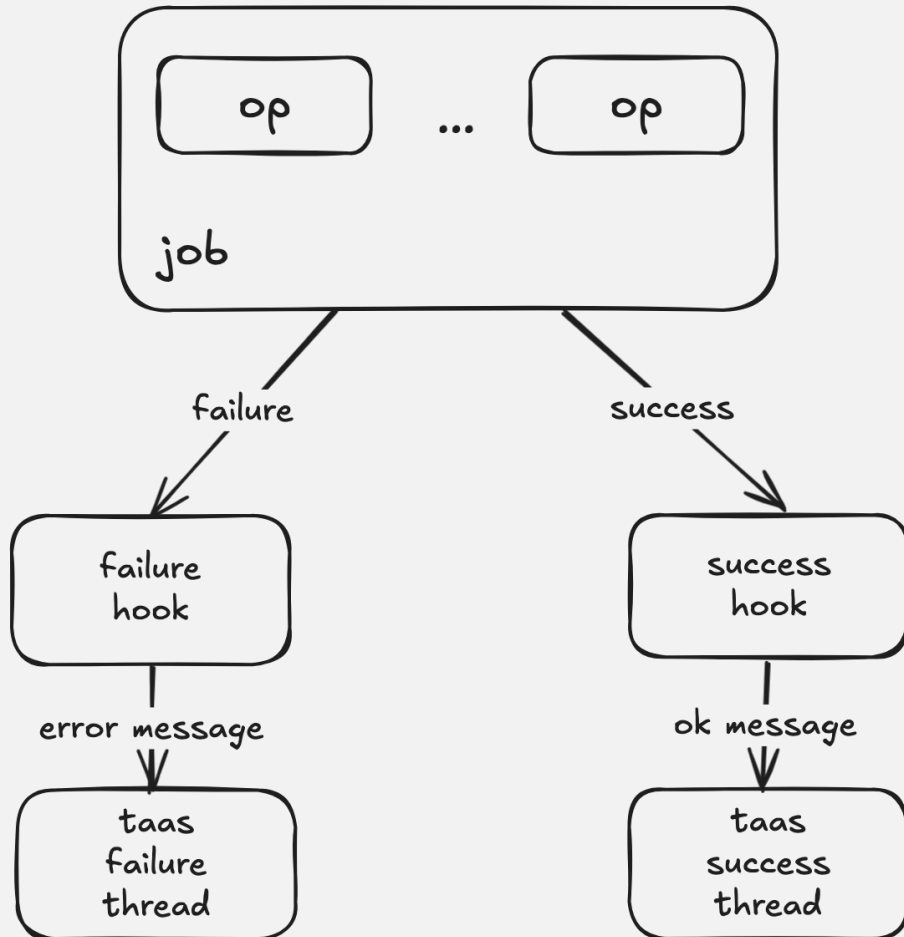
Мониторинг в telegram

TAAS (telegram as a service) — publish-subscribe сервис для отправки уведомлений через телеграм

Dagster hooks — механизм обработки завершившихся запусков



Мониторинг в telegram



```
1 from dagster import HookContext, job, op, success_hook
2 from taas import submit_event
3
4
5 @success_hook
6 def taas_success_hook(context: HookContext):
7     message = f'Dagster job {context.job_name} completed.'
8     submit_event(f'bp.dagster.oks', message)
9
10
11 @op
12 def process_spatial_data():...
13
14
15
16 @job
17 def spatial_job():
18     process_spatial_data.with_hooks({taas_success_hook})()
19
20
21
```

Итоги внедрения

Для менеджера и аналитика:

- уведомления в удобном формате

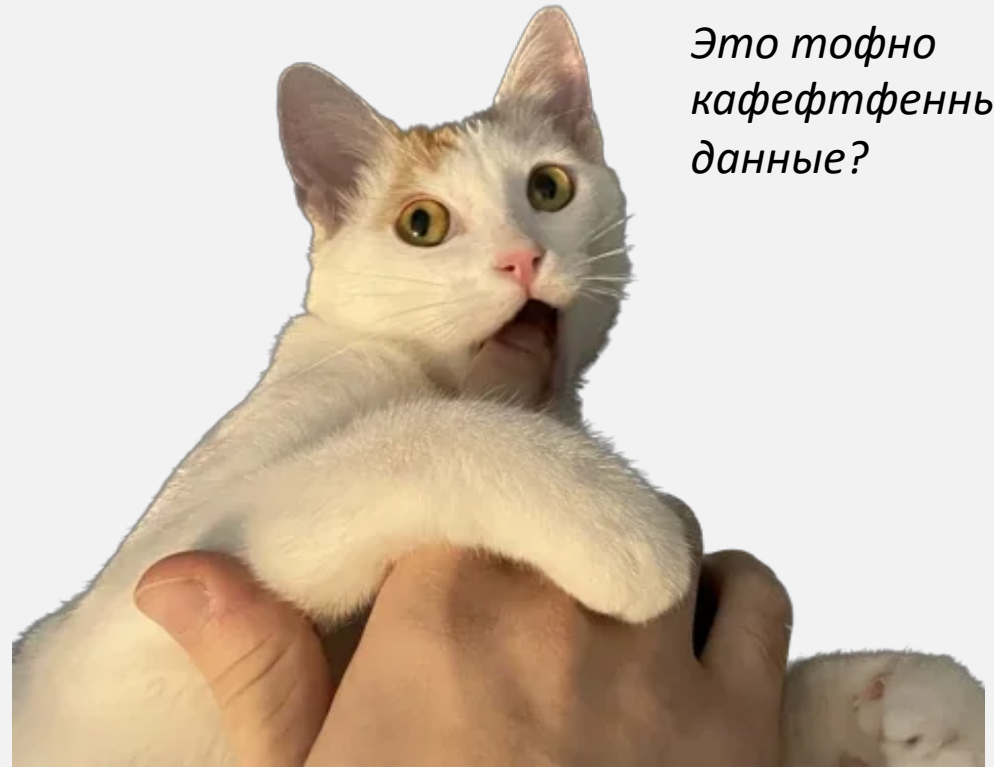
Для дата-инженера:

- простой способ добавлять уведомления

История про качество данных

История про качество данных

- Мониторинга недостаточно
- Пайплайны работают как надо, но есть нюанс
- Требуются проверки качества данных



*Это точно
кафефтенные
данные?*

Немного data quality

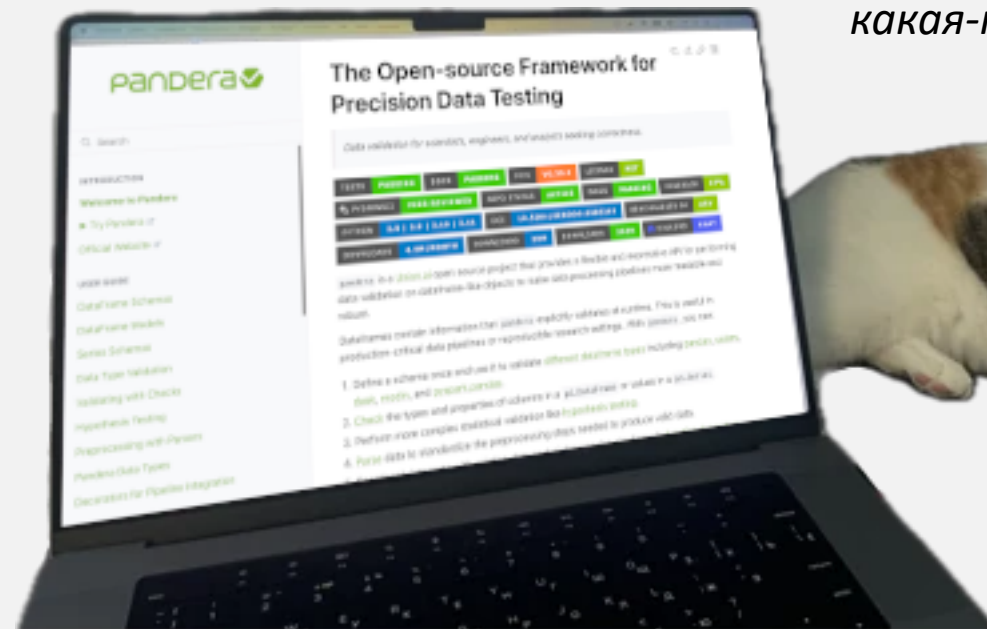
- great expectations
- **pandera**

*Ваши expectations —
ваши проблемы!*



Data quality: pandera

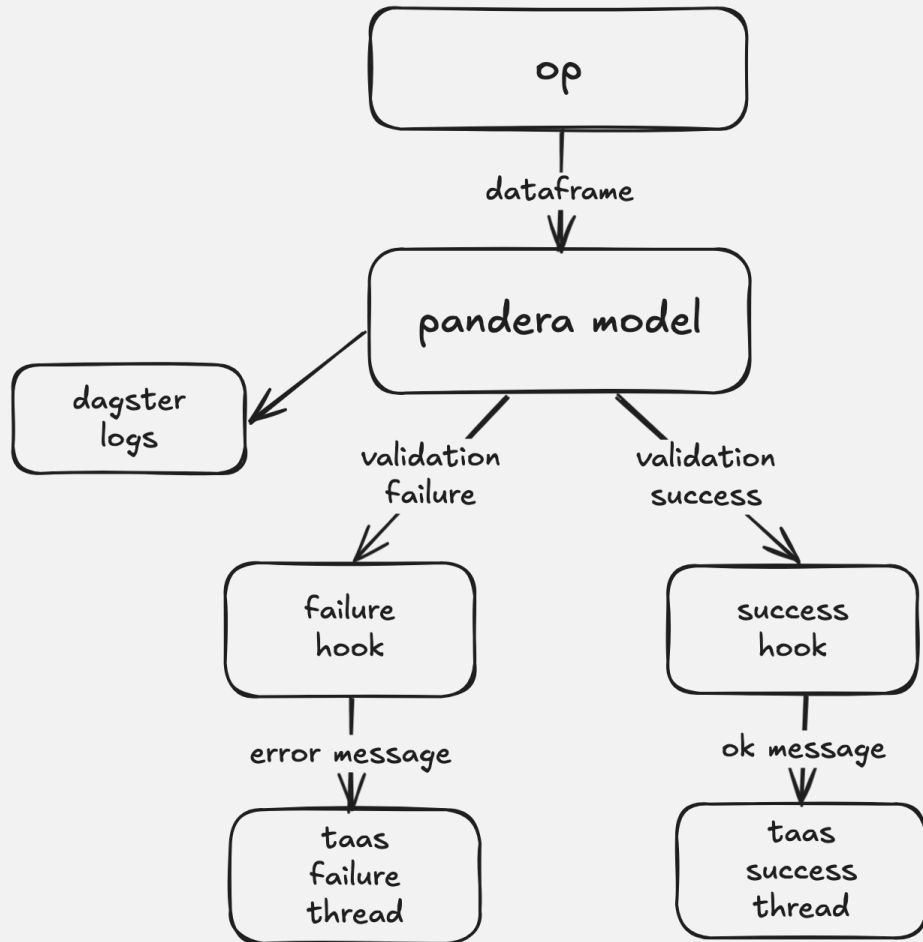
- Проще в использовании
- Изначально ориентирована на работу с pandas
- Схожесть с уже известным ruydantic



*Почему тебе пишет
какая-то пантера?*



Схема с валидацией и уведомлением



```
1 import pandera as pa
2 from pandera.typing import Series
3 from dagster_pandera import pandera_schema_to_dagster_type
4
5 class SomeDFWithGeo(pa.DataFrameModel):
6     name: Series[str] = pa.Field(description='POI name')
7     traffic: Series[int] = pa.Field(ge=0, nullable=False,
8                                     description='Traffic nearby')
9     geopos: Series[object] = pa.Field(nullable=False,
10                                       description='Point geometry')
11
12     @pa.check('geopos', name='Correct geojson')
13     def check_geopos(col: Series[object]) -> bool: ...
14
15     @op(out=Out(dagster_type=pandera_schema_to_dagster_type(SomeDFWithGeo)))
16     def process_spatial_data(): ...
17
18
19
20
21
22
23
24 @job
25 def spatial_job():
26     process_spatial_data.with_hooks({taas_success_hook})()
```

Итоги внедрения

Для всех:

- своевременные уведомления о проблемах в данных

История про непривычно много данных

История про непривычно много данных

- Встретили задачу, которая считалась привычными инструментами в руках аналитика 5 месяцев
- Данные лежат в паркетах в S3



*Откуда
столько?!*

Dagster + Spark

Интеграция с Apache Spark:

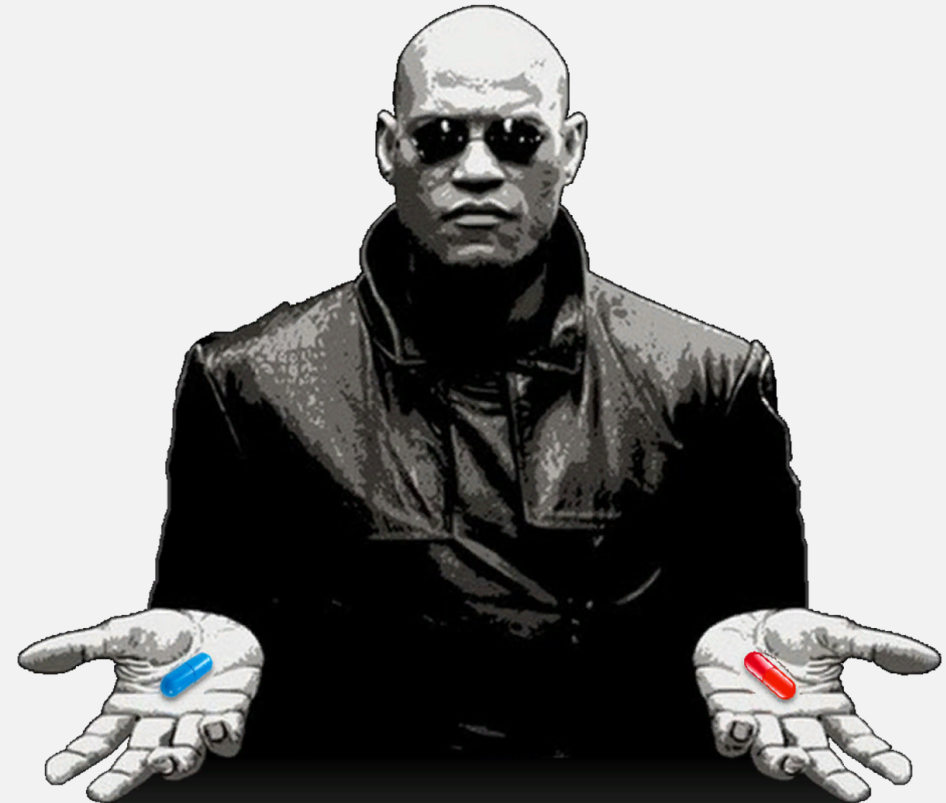
- Databricks
- AWS EMR
- YARN / spark standalone



Dagster + spark-operator

Проблемы, с которыми столкнулись при интеграции:

- сложности с запуском приложений через yam1 конфиг
- необходимость отслеживать прогресс задачи
- хочется читать логи



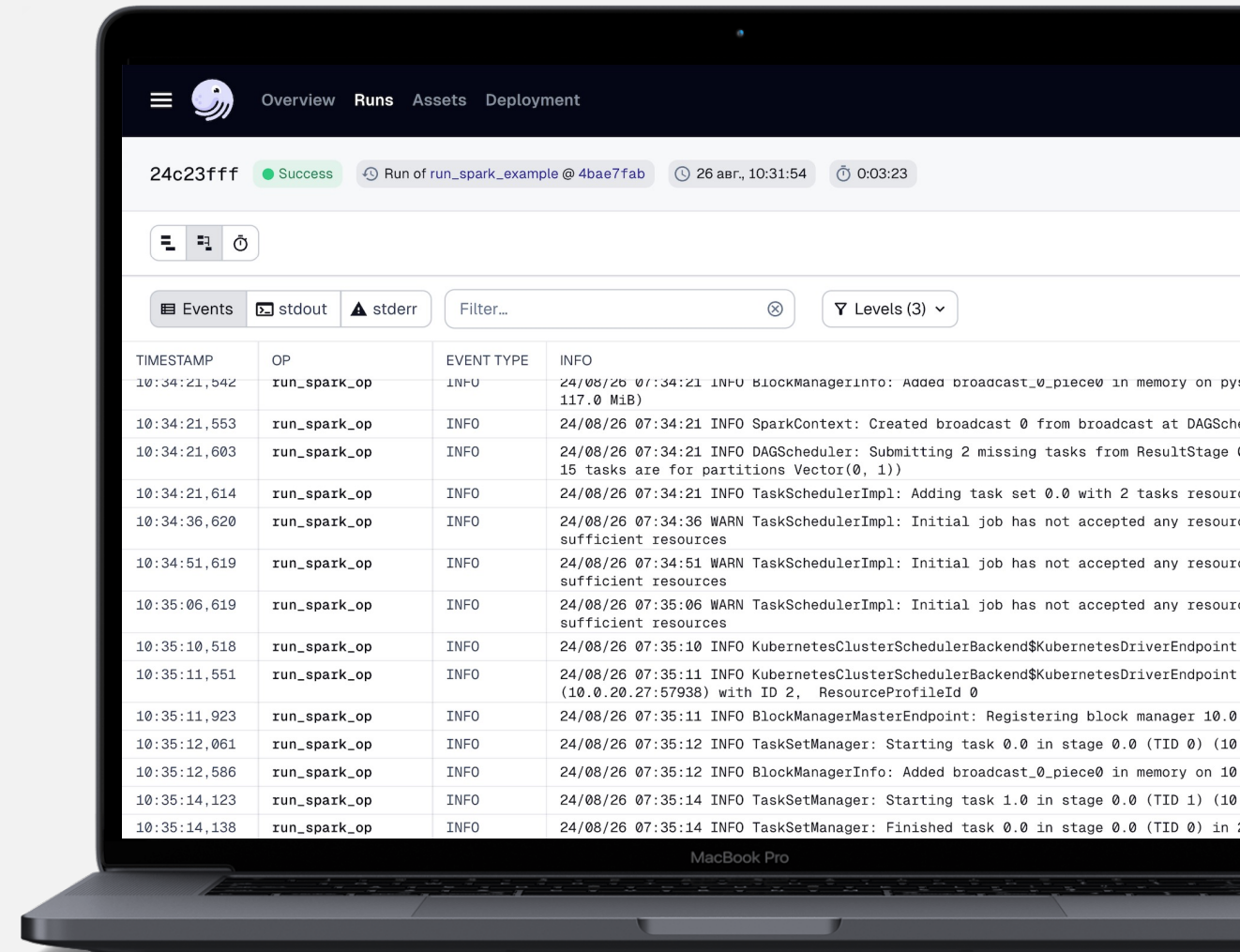
пишешь все сам

*пишешь все сам,
только клубничная*

Dagster + spark-operator

Собственная библиотека:

- задание конфигурации с помощью кода
- запуск и очистка ресурсов в k8s
- стриминг логов драйвера



spark8s

<https://github.com/bestplace/spark8s>



```
1  app_name = 'test-app'
2  # create manifest using static method
3  app = AppManifest.load_default_manifest(app_name)
4
5  app.set_driver_spec(1, '1000m', 2)
6  app.set_executor_spec(1, '1000m', 2, 4)
7  app.set_application_file('local:///somefile.py')
8
9  manifest = app.get_manifest()
10
11 # create kubernetes api client
12 api = client.ApiClient()
13 core_v1_api = client.CoreV1Api()
14 pod_name = f'{app_name}-driver'
15
16 # run app, stream logs, delete app and allocated resources
17 run_spark_app(api, manifest)
18 for line in stream_pod_logs(core_v1_api, pod_name, namespace):
19     log.info(line)
20 delete_spark_app(api, app_name)
21
```


Итоги внедрения

Для аналитика:

- новый инструмент обработки данных
- в конкретной задаче кратное ускорение

Для дата-инженера:

- инфраструктура для запуска spark приложений
- библиотека для шаблонизации конфигов и управления запусками

История про аналитику и визуализации

История про аналитику и визуализации

Проблема:

- аналитика и графики в jupyter ноутбуке
- отсутствие версионирования
- неочевидная логика преобразований



*Мне кажется,
или это не
best practice?*

Как же автоматизировать пайплайны?

- Настройка выгрузки данных по расписанию
- Добавление автоматизации преобразований



Как же автоматизировать пайплайны?

- Настройка выгрузки данных по расписанию
- Добавление автоматизации преобразований
- Используем dbt и его интеграцию с dagster



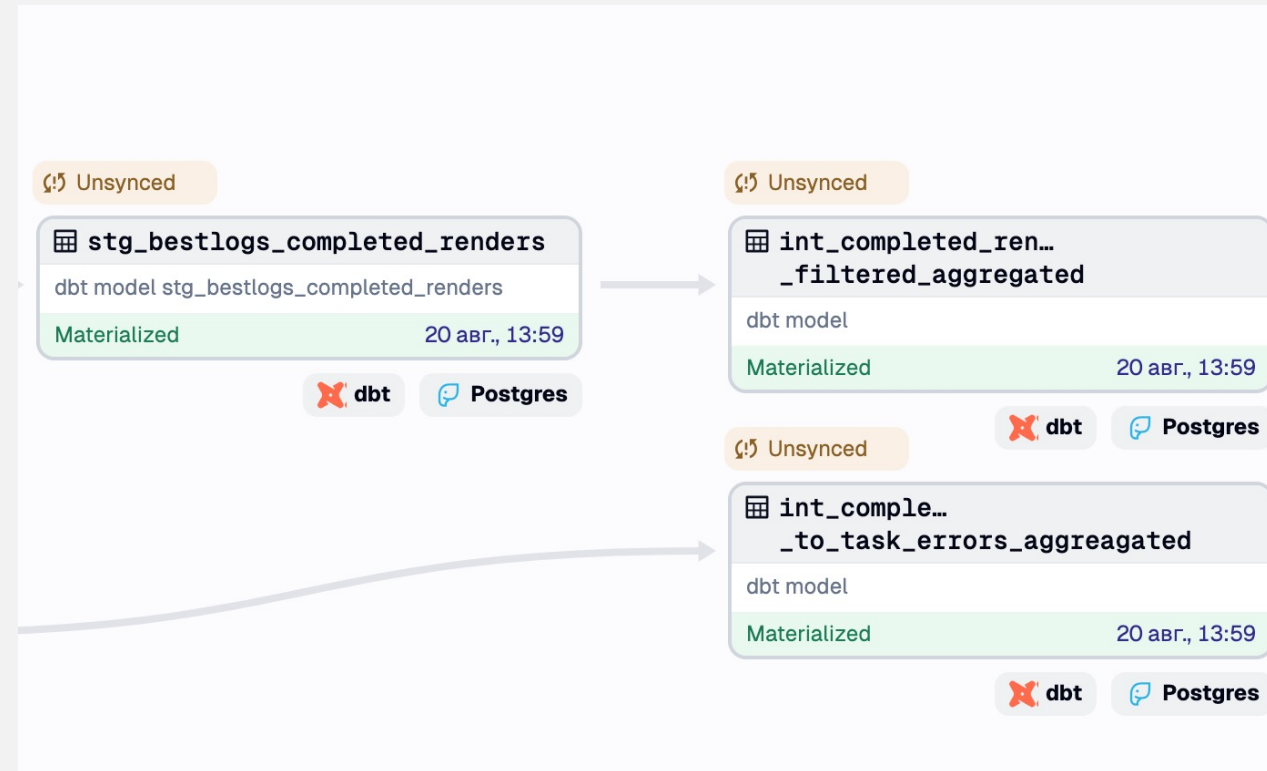
dagster + dbt

У дагстера есть **asset** — объект в постоянном хранилище (таблица, файл и т.п.)

У dbt есть **модели**

dagster-dbt дает нам:

- dbt_assets
- интеграцию в глобальный lineage
- возможность запуска и параметризация через UI



dagster + dbt

Проблемы, с которыми столкнулись при интеграции:

- ???



Итоги внедрения

Для менеджера:

- дашборды на свежих данных каждое утро
- не нужно заниматься поддержкой кода

Для дата-инженера:

- автоматизация создания витрин
- процессы работы с данными стали прозрачнее
- граф связей между данными из разных источников

ИТОГИ

Итоги

- Легко вкатиться как дата-инженеру, так и аналитику
- Большое количество интеграций
- Активная поддержка со стороны разработчиков и растущее англоязычное комьюнити
- Несложное развертывание и надежная работа



```
git clone https://github.com/dagster-io/dagster.git
cd dagster/examples/deploy_docker
docker compose up --build
```



Bestplace

Спасибо за внимание!

Дмитрий Крылов
✈️ mityakrylov



Алексей Завальский
✈️ al_zav



и Скрепка

