

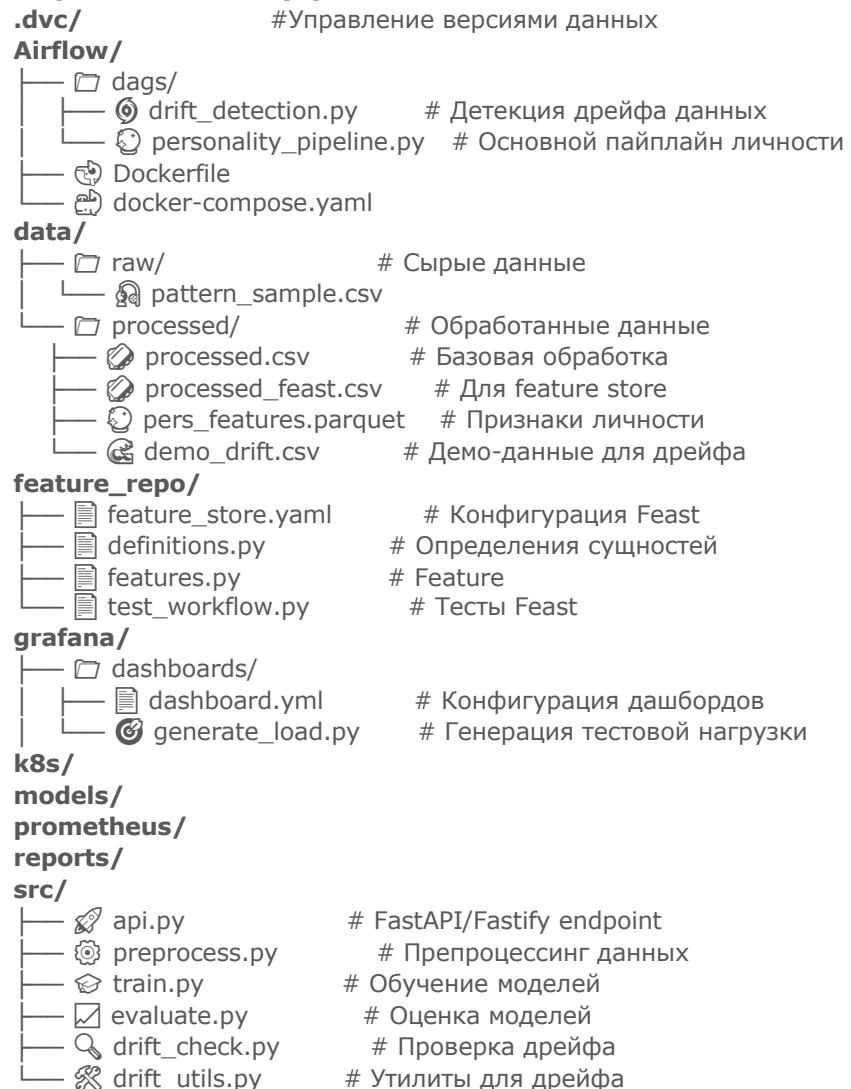
Предсказание экстраверта/ интроверта

MLOps Pipeline

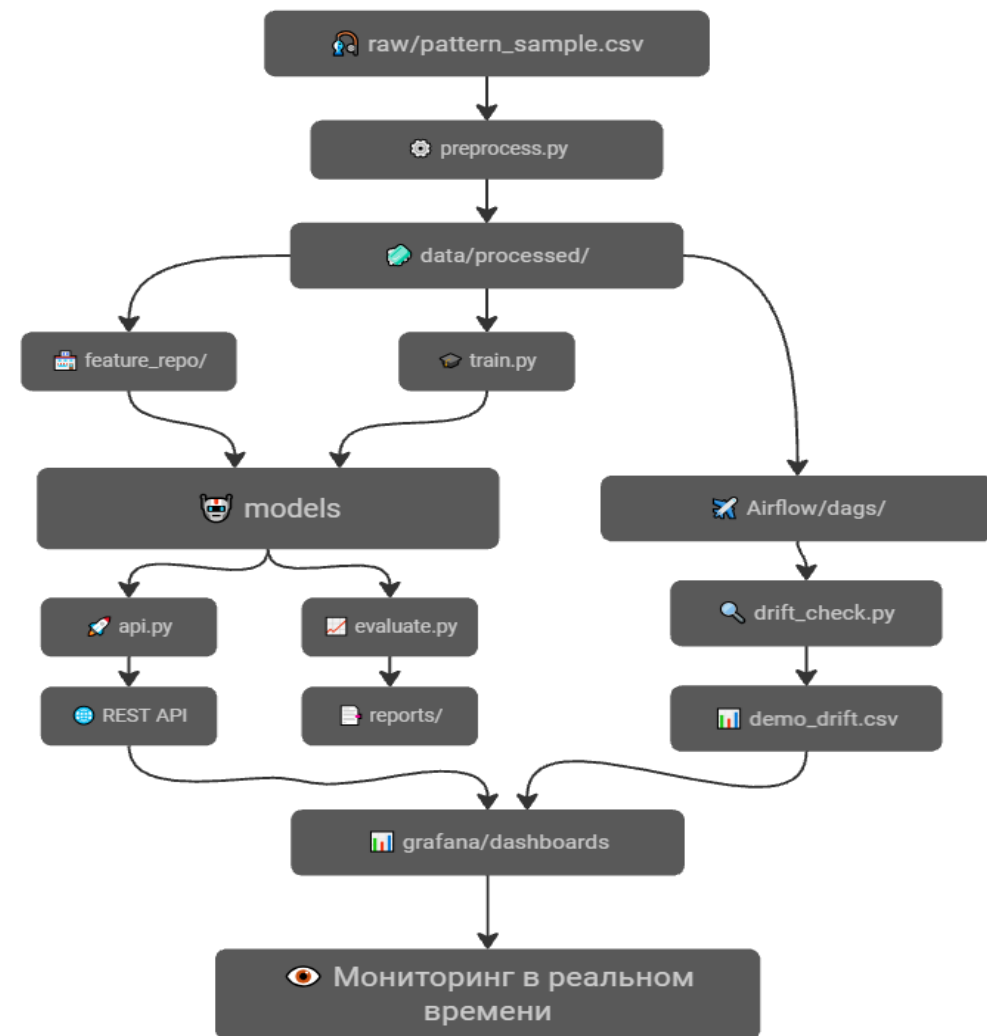
Выполнили:
Король А.Ю.
Раздоркина Д.С.
Банников Д.Ю.
Аврамчук И.А.

Структура и архитектура проекта

Архитектура



Путь данных



Данные и обработка признаков

Пайплайн обработки

- **Очистка:** удаление дубликатов и выбросов.
- **Нормализация:** приведение численных данных к единой шкале
- **Кодирование:** One-hot и Label encoding для категорий.

Целевая переменная и фичи

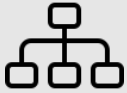
- **Target:** Extrovert/Introvert
- **Features:** Время проведения в одиночестве, страх сцены, посещение мероприятий, частота выхода из дома, чувство истощения после общения, круг друзей, частота постов в соц сетях.
- **Хранение:** Feast Feature Store обеспечивает доступность фичей для обучения и инференса.

Источник данных: [kaggle.com](https://www.kaggle.com)

Объем данных: 3000 строк описаний личности

Train/Test: 80% / 20%

Инструменты разработки и качества



Pre-commit Hooks

Автоматический линтинг и форматирование кода:

- Black
- Flake8
- Check YAML



Feature store (Feast)

Управление признаками и версионирование данных



Orchestration

Airflow DAGs для регулярных задач:

- Обучение модели
- Ежедневный дрифт чек

Feature store: Feast

Консистентность признаков

Feast используется для унификации признаков между этапами обучения и инференса

Offline Store:

- Хранение набора данных для обучения модели RandomForest.
- Регистрация: хранение метаданных о признаках (версии, источники).

Признаки (всего 10):

time_spent_alone

friends_circle_size

stage_fear

post_frequency

social_event_attend
ance

personality

going_outside

personality_encoded

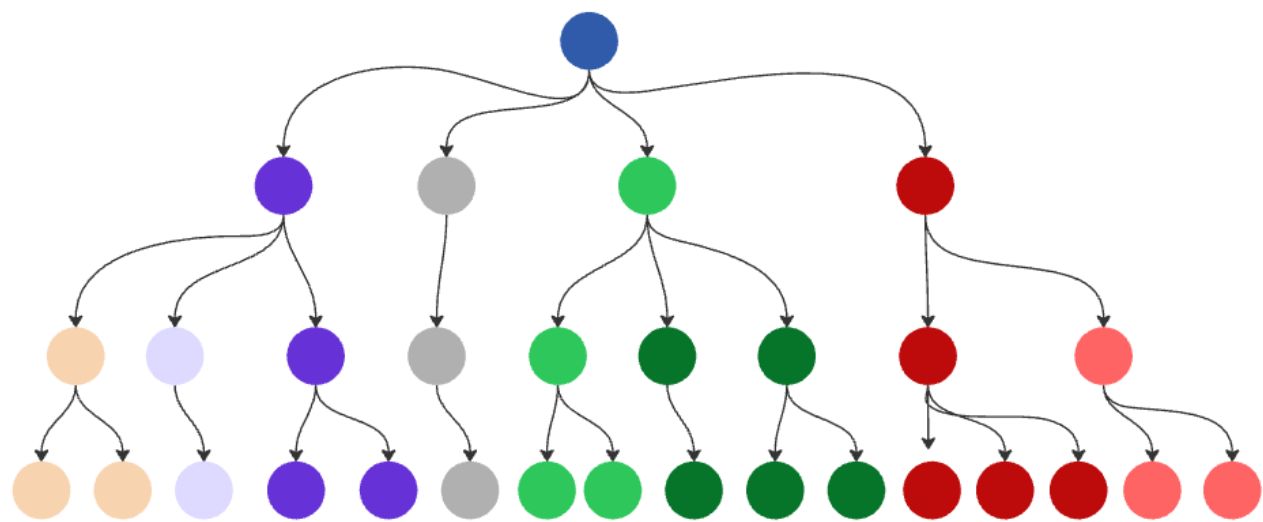
drained_after_sociali
zing

user_id

Модель

RandomForest

Выбрана за интерпритируемость и удобства при работе с категориальными признаками

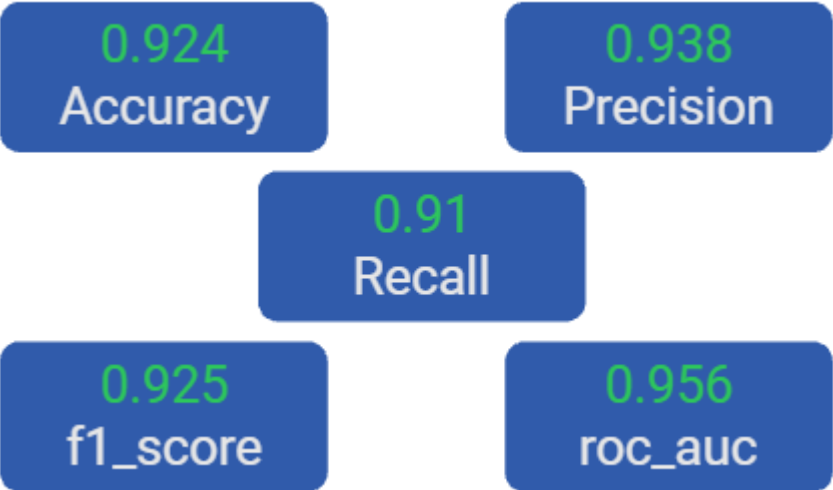


Параметры дерева

Кол-во деревьев	n_estimators=	100
Макс глубина дерева	max_depth=	10
Кол-во сэмплов для сплита	min_sam_split=	5
Кол-во сэмплов в листе	min_sam_leaf=	2
Стратегия выборки признаков	max_features=	sqrt
Критерий разделения	criterion=	gini
Бутстраппинг	bootstrap=	True
Веса классов	class_weight=	balanced
RANDOM_STATE	random_state=	42
Параллелизация	n_jobs=	-1

Результаты обучения

Метрики



Матрица ошибок

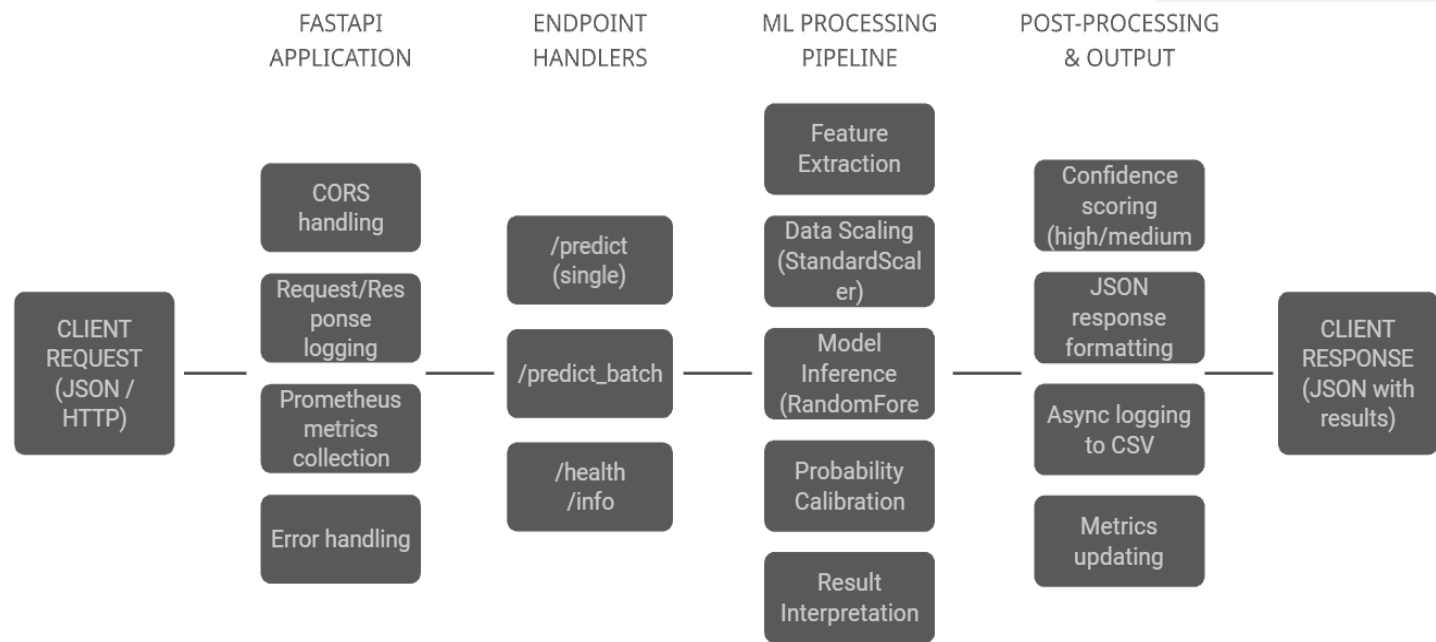
Матрица ошибок			
	Предсказано		
	Интроверт	Экстраверт	
Истина интров.	272	26	298
Истина экстрав.	18	264	282
	290	290	

Описание API (FastAPI)

Основные Эндпоинты

Документация доступна через Swagger UI.

Архитектура API



Personality Classifier API 1.0.0 OAS 3.1

/openapi.json

default

GET / Root

GET /health Health Check

GET /metrics Metrics

GET /model_info Model Info

POST /predict Predict

POST /predict_batch Predict Batch

API + UI

Сколько времени вы предпочитаете проводить в одиночестве?

5

Далее

Вы боитесь выступать перед людьми?

Да, боюсь

Да, боюсь

Нет, мне комфортно

Назад

Далее

Как часто вы публикуете посты в соцсетях?

5

Назад

Узнать результат



Introvert

Уверенность алгоритма: 65.4%

Пройти заново

Автоматизация и оркестрация Apache Airflow

Инструменты

- Настроен DAG для ежедневного обнаружения дрифта данных (Data Drift)
- Мониторинг: сравнение распределения признаков между тренировочным и текущим production - набором.
- Триггер: если дрифт превышает установленный порог, Airflow может автоматически запустить DAG переобучения.

manual__2025-12-15T17:25:26+00:00

Success

Add a note

Logical Date

Run Type

Start Date

End Date

Duration

2025-12-15 20:25:26

▶ manual

2025-12-15 20:25:27

2025-12-15 20:25:41

00:00:13.799

Task Instances

Asset Events

Audit Log

Code

Details

Q Search Tasks

All States

Task ID

Map Index

State

Start Date

Try Number

Operator

register

Success

2025-12-15 20:25:40

1

BashOperator

evaluate

Success

2025-12-15 20:25:38

1

BashOperator

train

Success

2025-12-15 20:25:29

1

BashOperator

preprocess

Success

2025-12-15 20:25:27

1

BashOperator

manual__2025-12-20T16:38:47+00:00

Success

Logical Date

Run Type

Start Date

End Date

Duration

2025-12-20 19:38:47

manual

2025-12-20 19:38:50

2025-12-20 19:38:56

00:00:05.494

Task Instances

Asset Events

Audit Log

Code

Details

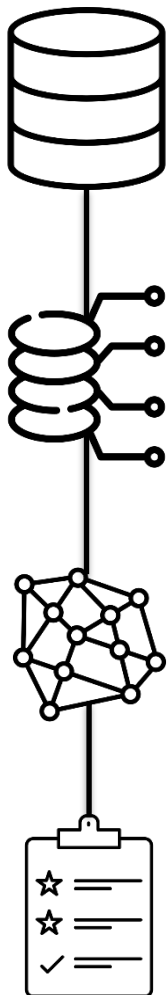
Q Search Tasks

All States

Task ID	Map Index	State	Start Date	Try Number	Operator
complete		Success	2025-12-20 19:38:56	1	EmptyOperator
trigger_retraining		Success	2025-12-20 19:38:54	1	TriggerDagRunOperator
no_drift_detected		Skipped	2025-12-20 19:38:54	0	BashOperator
check_drift		Success	2025-12-20 19:38:53	1	BranchPythonOperator
simulate_drift_for_demo		Success	2025-12-20 19:38:52	1	PythonOperator
create_reference_data		Success	2025-12-20 19:38:50	1	PythonOperator

manual__2025-12-17T14:44:16+00:00 Success					
Logical Date	Run Type	Start Date	End Date	Duration	Triggering User Name
2025-12-17 17:44:16	manual	2025-12-17 17:44:18	2025-12-17 17:44:35	00:00:17.010	airflow
Task Instances Asset Events Audit Log Code Details					
Q Search Tasks		All States			
Task ID	Map Index	State	Start Date	Try Number	Operator
save_report		Success	2025-12-17 17:44:34	1	BashOperator
evaluate		Success	2025-12-17 17:44:32	1	BashOperator
train_with_feast		Success	2025-12-17 17:44:26	1	BashOperator
feast_materialize		Success	2025-12-17 17:44:19	1	PythonOperator
preprocess		Success	2025-12-17 17:44:18	1	BashOperator

ML-пайплайн (DVC):



preprocess

Очистка данных, кодирование, нормализация. Выход:
`data/processed/processed.csv`

split

Разделение датасета на тренировочный и тестовый
(сохраняется в DVC).

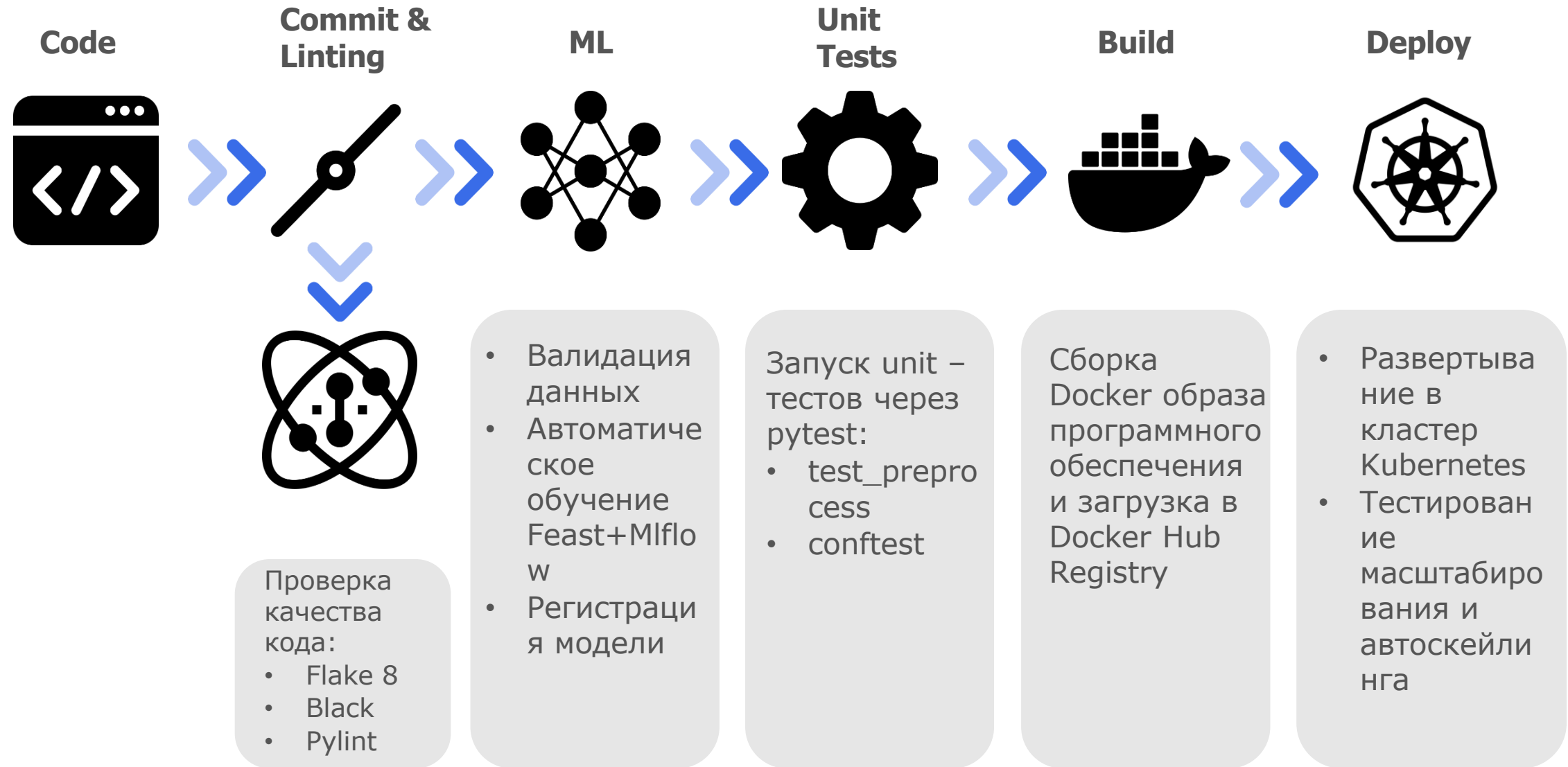
train

Обучение модели RandomForest. Выход:
`models/personality_model.joblib`

evaluate

Оценка модели, логирование метрик и артефактов в MLflow.
Выход: `models/classification_metrics.json`

CI/CD Pipeline



Мониторинг и observability

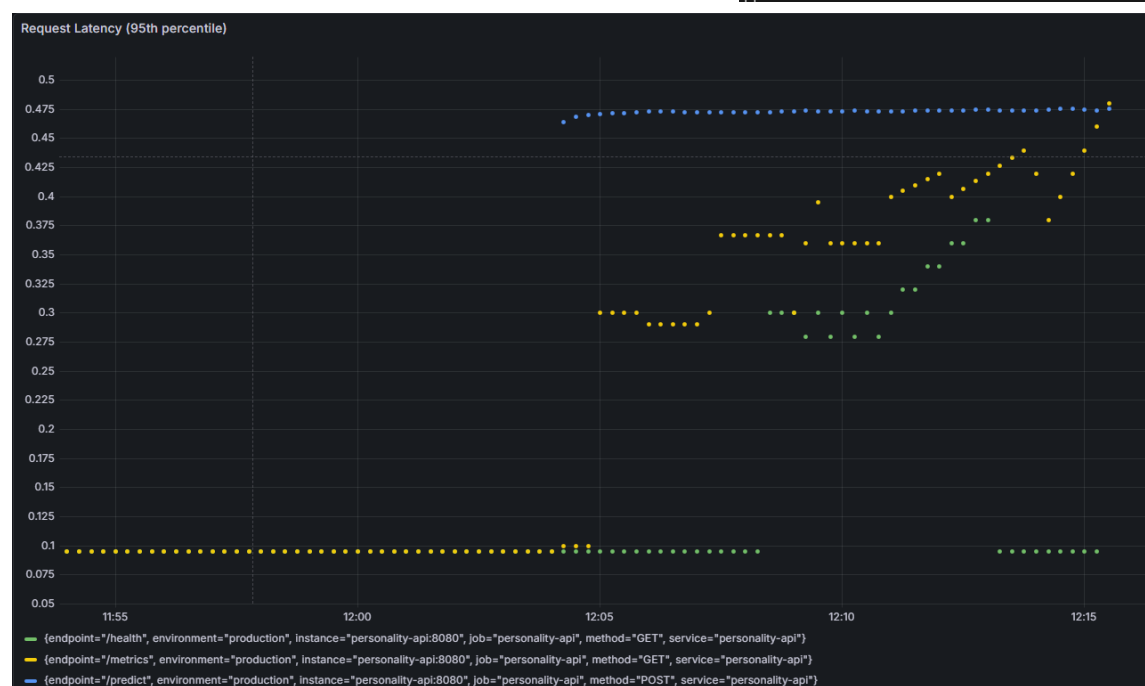
Инструменты

- Prometheus: сбор метрик с микросервисов
- Grafana: визуализация дашбордов

Отслеживаемые показатели

- Сколько запросов API обрабатывает в каждую секунду
- Среднюю нагрузку на систему

```
Worker 0: 300 requests, latency: 0.108s, status: 200
Worker 2: 300 requests, latency: 0.133s, status: 200
Worker 4: 300 requests, latency: 0.115s, status: 200
Worker 3: 300 requests, latency: 0.101s, status: 200
Worker 1: 320 requests, latency: 0.117s, status: 200
Worker 0: 320 requests, latency: 0.407s, status: 200
Worker 4: 320 requests, latency: 0.115s, status: 200
Worker 2: 320 requests, latency: 0.154s, status: 200
Worker 3: 320 requests, latency: 0.169s, status: 200
Worker 0: 340 requests, latency: 0.102s, status: 200
Worker 1: 340 requests, latency: 0.186s, status: 200
Worker 4: 340 requests, latency: 0.161s, status: 200
Worker 2: 340 requests, latency: 0.117s, status: 200
Worker 3: 340 requests, latency: 0.118s, status: 200
Worker 0: 360 requests, latency: 0.106s, status: 200
Worker 1: 360 requests, latency: 0.109s, status: 200
Worker 4: 360 requests, latency: 0.123s, status: 200
Worker 3: 360 requests, latency: 0.143s, status: 200
Worker 2: 360 requests, latency: 0.118s, status: 200
```



Команда

Король Алексей

Настроил Git-репозиторий с полной структурой проекта (src/, data/, models/, tests/)

Инициализировал DVC с локальным remote для версионирования данных

Реализовал DVC пайплайн с stage для предобработки данных

Добавил CI-реплику в виде базового GitHub Actions workflow

Разработал REST API на FastAPI с эндпоинтом /predict для онлайн-предсказаний

Раздоркина Дарья

Разработала unit-тесты для функций предобработки (tests/test_preprocess.py)

Создала тесты API (tests/test_api.py)
Настроила pytest конфигурацию и фикстуры в conftest.py

Реализовала Airflow DAG для оркестрации полного ML пайплайна

Настроила задачи Airflow - preprocess, train, evaluate, register

Добавила мониторинг выполнения - логи, ретраи для DAG задач

Банников Дмитрий

Разработал ML пайплайн обучения (src/train.py) с RandomForestClassifier

Интегрировал MLflow для логирования экспериментов, параметров и метрик

Реализовал оценку модели с вычислением ROC AUC, precision/recall, confusion matrix

Настроил MLflow Model Registry

Создал фронтенд-интерфейс (Swagger UI/ReDoc) для тестирования API

Аврамчук Игорь

Интегрировал Docker-сборку в GitHub Actions workflow

Реализовал деплой в Minikube - создал k8s манифесты

Добавил мониторинг Prometheus - инструментировал API метриками (/metrics эндпоинт)

Развернул Grafana dashboard для визуализации метрик

Реализовал детекцию дрейфа

Настроил автоматический retraining при обнаружении дрейфа