

Как небольшой компании  
запустить MLOps?

01  
Введение

О компании  
Что было  
и почему пошли в  
MLOps

02  
Итерация 0

Ресерч  
Стало  
Вывод

03  
Итерация I

Ресерч  
Стало  
Вывод

04  
Итерация II

Ресерч  
Стало  
Вывод

05  
Заключение

Что делать  
или roadmap  
по кварталам

CASE

# Продуктовый ритейлер

## Проблема

Трудоемкость сбора и обработки информации о ценах конкурентов в магазинах

## Решение

Приложение для оффлайн-мониторинга цен конкурентов с функцией планировщика заданий

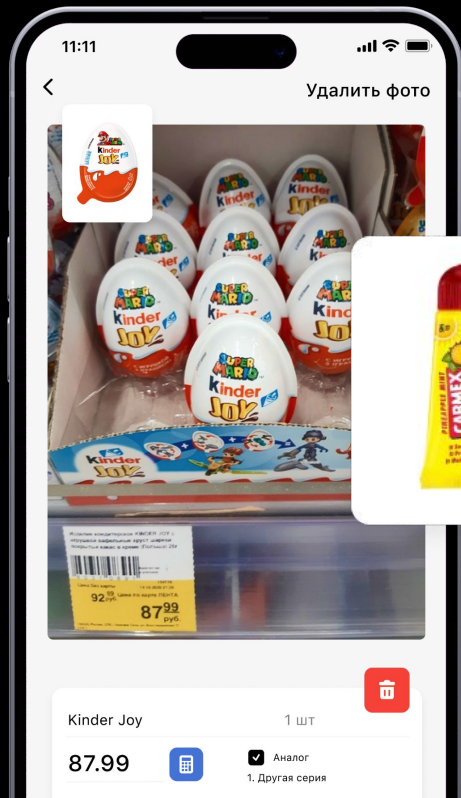
Супервайзер ставит работнику задачу по мониторингу цен на отдельные товарные позиции в магазинах конкурентов

Пользователь делает фотографии товара с ценником или всего стеллажа. Они автоматически распознаются, заносятся в реестр и передаются супервайзеру в режиме реального времени



CASE

# Показатели бизнес-кейса



87 → 33 МИН

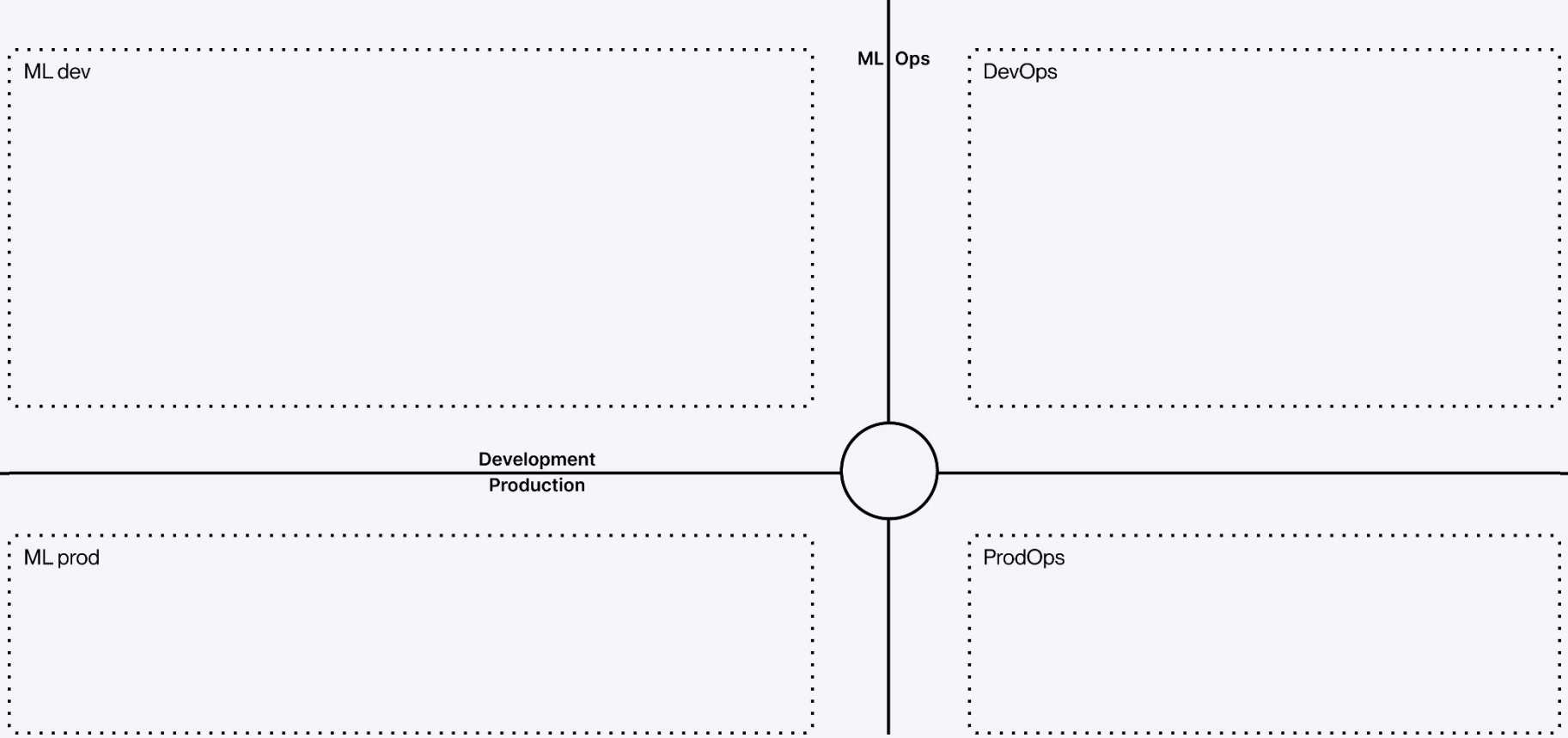
сокращение времени  
мониторинга в расчете  
на 100 SKU

20%

рост точности собираемых  
данных по ценам конкурентов

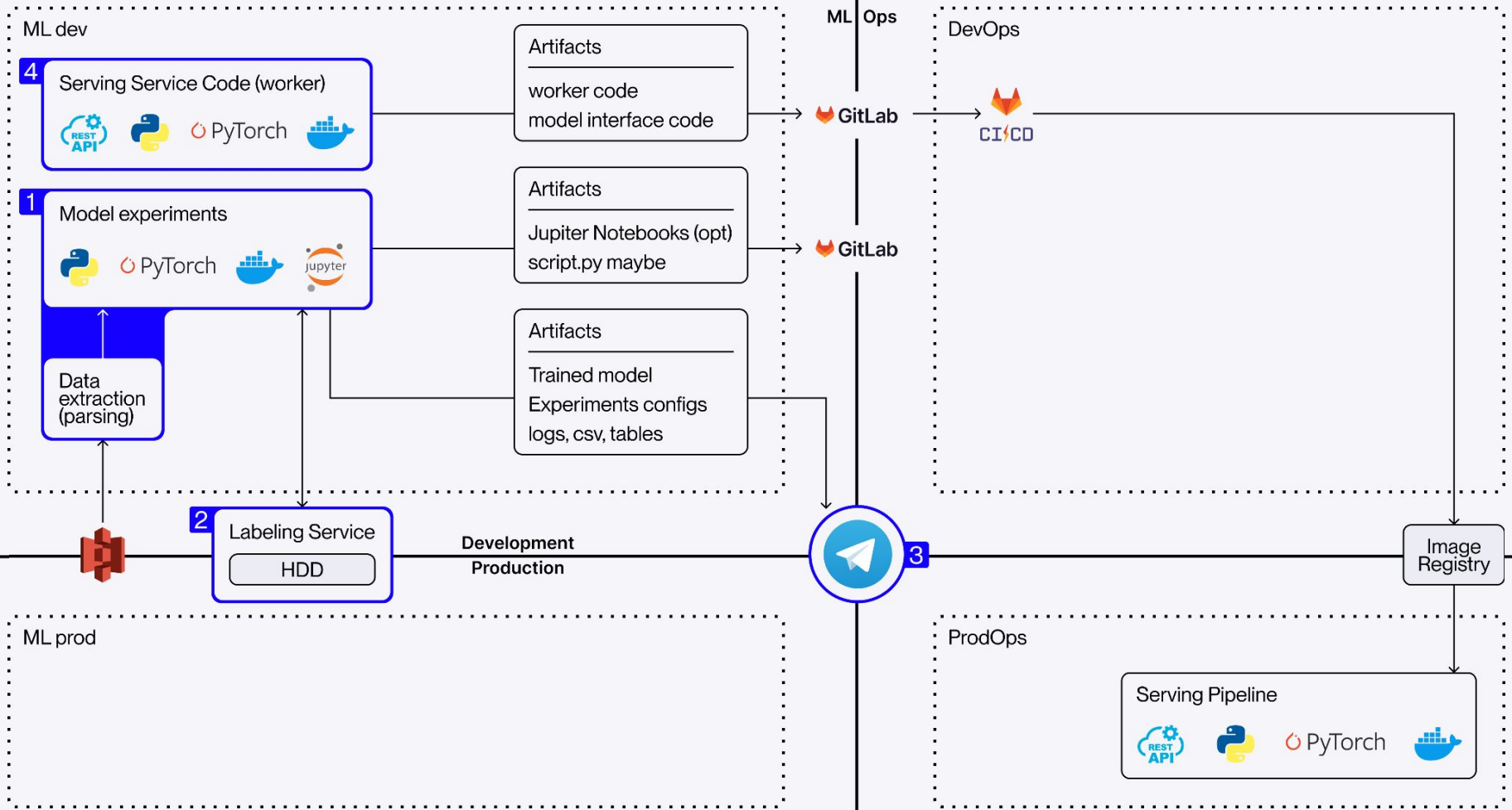


увеличение охвата рынка  
конкурентных магазинов  
в ходе мониторинга цен



MLOps в билайн: как катить машинное обучение в production без ML-инженеров

Построение MLOps платформы. Как мы обуздали хаос в головах и технике



- 01 Введение
- 02 Интеграция 0
- 03 Интеграция I
- 04 Интеграция II
- 05 Заключение

# Почему пошли в MLOps

- Увольнение разработчика
- А где та модель, которую мы получили 2 месяца назад?
- Ой, Я случайно сделал `sudo rm-rf` на датасет
- Проблема с утилизацией железа и деплоем моделей

90%

кода лежит локально

90%

артефактов лежат  
локально

100%

ручной однотипной работы  
для релиза

0%?

проекта удовлетворяет  
концепциям mlops

- 01 Введение
- 02 Интеграция 0
- 03 Интеграция I
- 04 Интеграция II
- 05 Заключение

Интеграция 0

# Ресерч

MLOps — методология оптимизации, ускорения и автоматизации процесса разработки и доставки решений Машинного обучения

DevOps  
MLOps  
CloudOps  
AIOps  
DataOps  
ITOps  
RevOps  
FinOps  
HROps  
LegalOps  
BizOps  
LLMOps  
ChatOps  
NoOps



# Ресерч. Требования

01

Нужно что-то, чтобы хранить и **версионировать артефакты**

03

Нужно что-то, чтобы хранить и **версионировать данные**

02

Нужно что-то, что сможет выступать как **оркестратор** для пайплайнов

04

Нужно что-то, что предоставляет функционал **разметки данных**

# Ресерч. Результат

The logo for mlflow, featuring the word "mlflow" in a blue, lowercase, sans-serif font with a small "TM" symbol.

Логирование  
экспериментов  
и моделей

The logo for DVC, consisting of the letters "DVC" in a bold, black, sans-serif font.

Версионирование  
данных

The logo for Apache Airflow, featuring a colorful pinwheel icon with four blades in red, blue, green, and yellow, followed by the text "Apache Airflow" in a black, sans-serif font.

Оркестратор

The logo for CVAT, featuring a black circle with a white dot inside, followed by the letters "CVAT" in a bold, black, sans-serif font.

Сервис разметки



[MLOps: Continuous delivery and automation pipelines in machine learning](#)

[The 2024 Mad \(Machine Learning, Artificial Intelligence & Data\) Landscape](#)

Интеграция 0

# Стало. Инструменты





# Вывод

- Единое окно: трудности в эксплуатации и поддержке такой инфраструктуры
- Сложности в масштабировании — нужно писать кучу обвязок вокруг этих инструментов (возвращаемся к проблеме выше)
- Инструментарий бигтеха не всегда лучший выбор для компаний/команд маленького/среднего размера

90% → 0%  
кода лежит локально

90% → 0%  
артефактов лежат  
локально

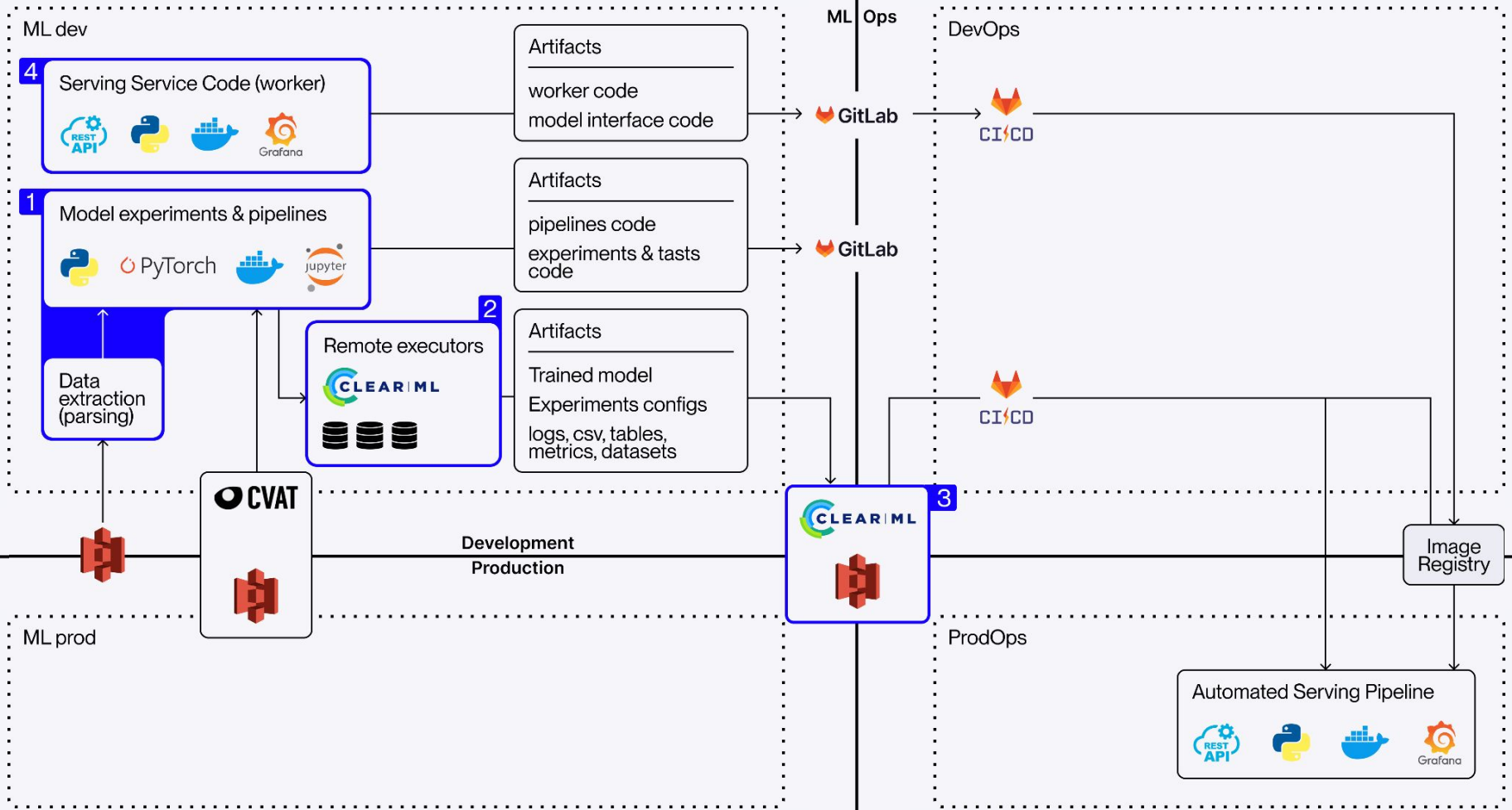
100% → 60%  
ручной однотипной работы  
для релиза

0% → 100%?  
проекта удовлетворяет  
концепциям MLOps

- 01 Введение
- 02 Интеграция 0
- 03 Интеграция I
- 04 Интеграция II
- 05 Заключение

- 01 Изменение метрик, по которым отслеживаем результат внедрения MLOps
- 02 Более точечный ресерч с учетом опыта
  - Более ML-ый оркестратор
  - Единое окно для рабочего места ML специалиста

|                    | ML Flow | Clear ML            | Wandb | DVC studio | V7                  |
|--------------------|---------|---------------------|-------|------------|---------------------|
| Интерфейс          | -       | +                   | +     | + -        | +                   |
| On Premise         | +       | +                   | +     | +          | -                   |
| Enterprise version | -       | +                   | +     | +          | +<br>ТОЛЬКО ПЛАТНАЯ |
| RBAC               | -       | -<br>ТОЛЬКО ПЛАТНАЯ | -     | - +        | +                   |
| Пайплайны          | + -     | +                   | -     | -          | +                   |
| Remote execution   | -       | +                   | -     | -          | -                   |
| Мониторинг         | -       | +                   | -     | -          | +                   |
| Версионирование    | -       | +                   | +     | +          | +                   |
| Разметка           | -       | +<br>enterprise     | -     | -          | +                   |



# Вывод

- ClearML прост в использовании и имеет большой функционал из коробки
- Обернуть в скрипты + remote executors + release\_pipeline уменьшаем время релиза и увеличиваем утилизацию железа на 35%
- CI/CD уменьшил кол-во факапов деплоя моделей на 30%

0% → 25%

уменьшено время релиза

—

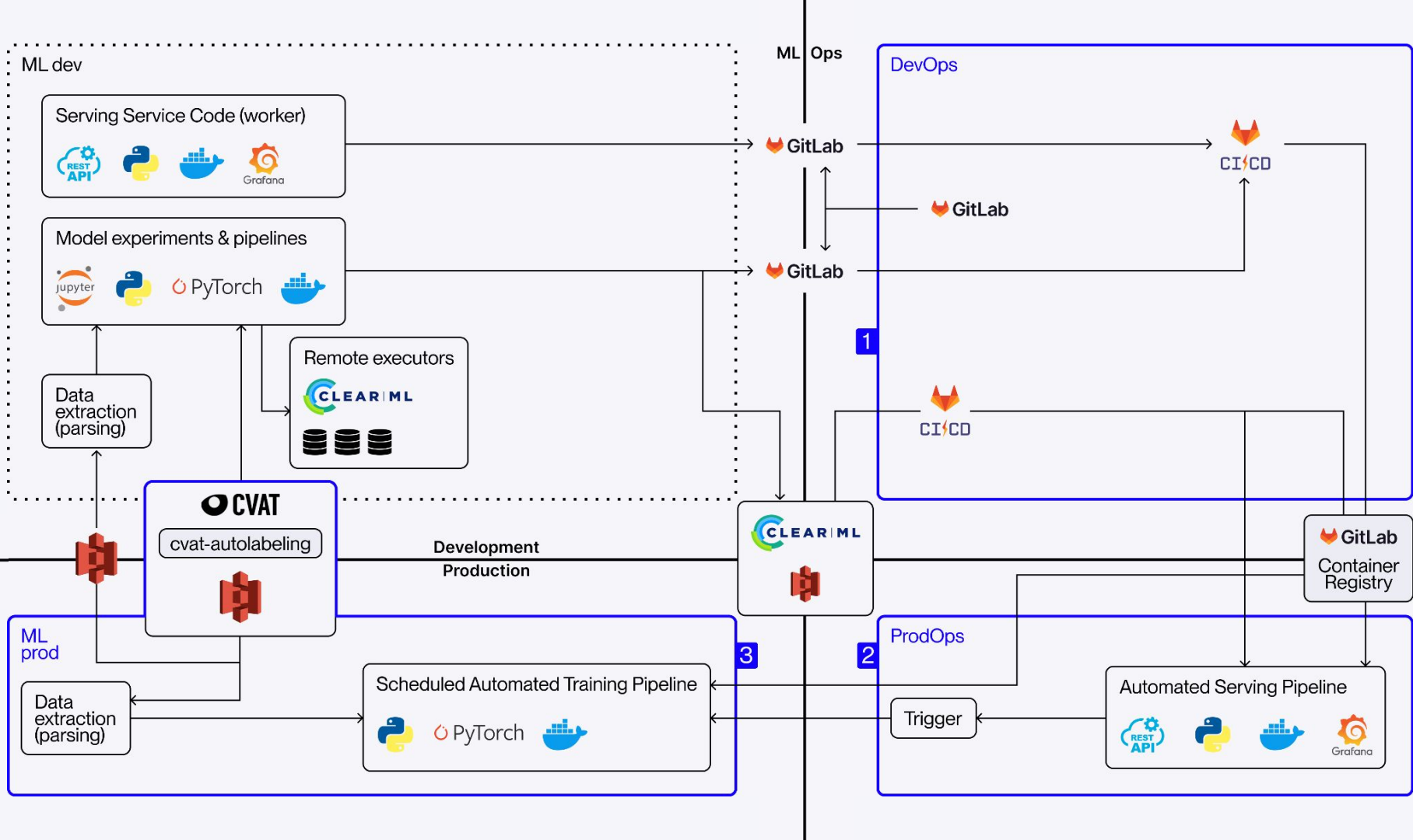
уменьшено время демо

0% → 35%

соответствие уровню  
MLOps lvl 2

- 01 Введение
- 02 Интеграция 0
- 03 Интеграция I
- 04 Интеграция II
- 05 Заключение

- А как организовать рабочее место не только ML Engineer но и Data Analyst?
- С dev разобрались, а что делать с prod?
- Можем ли мы еще ускорить процесс обучения?



# Выводы

- Авторазметка CVAT-а экономит ресурсы и уменьшает время разметки => релиза
- Добавление шедулеров на пайплайны уменьшает время релизов и загруженность разработчиков
- Даже костыльные триггеры на метрики в связке с автоматизацией дают прирост по метрикам

25% → 65%

уменьшено время релиза

0% → 20%

уменьшено время демо

20% → 65%

соответствие уровню  
MLOps lvl 2

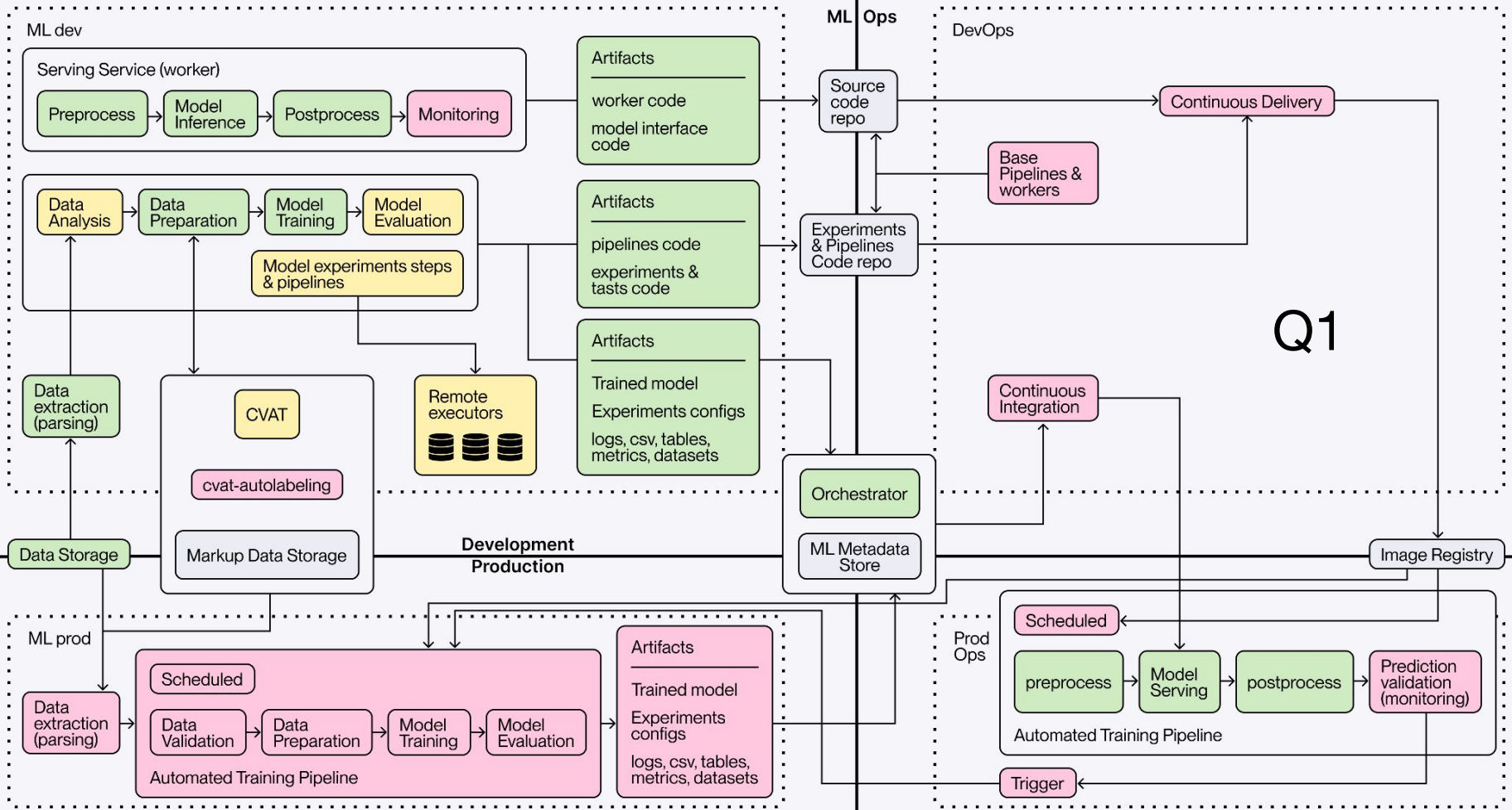
# Планы

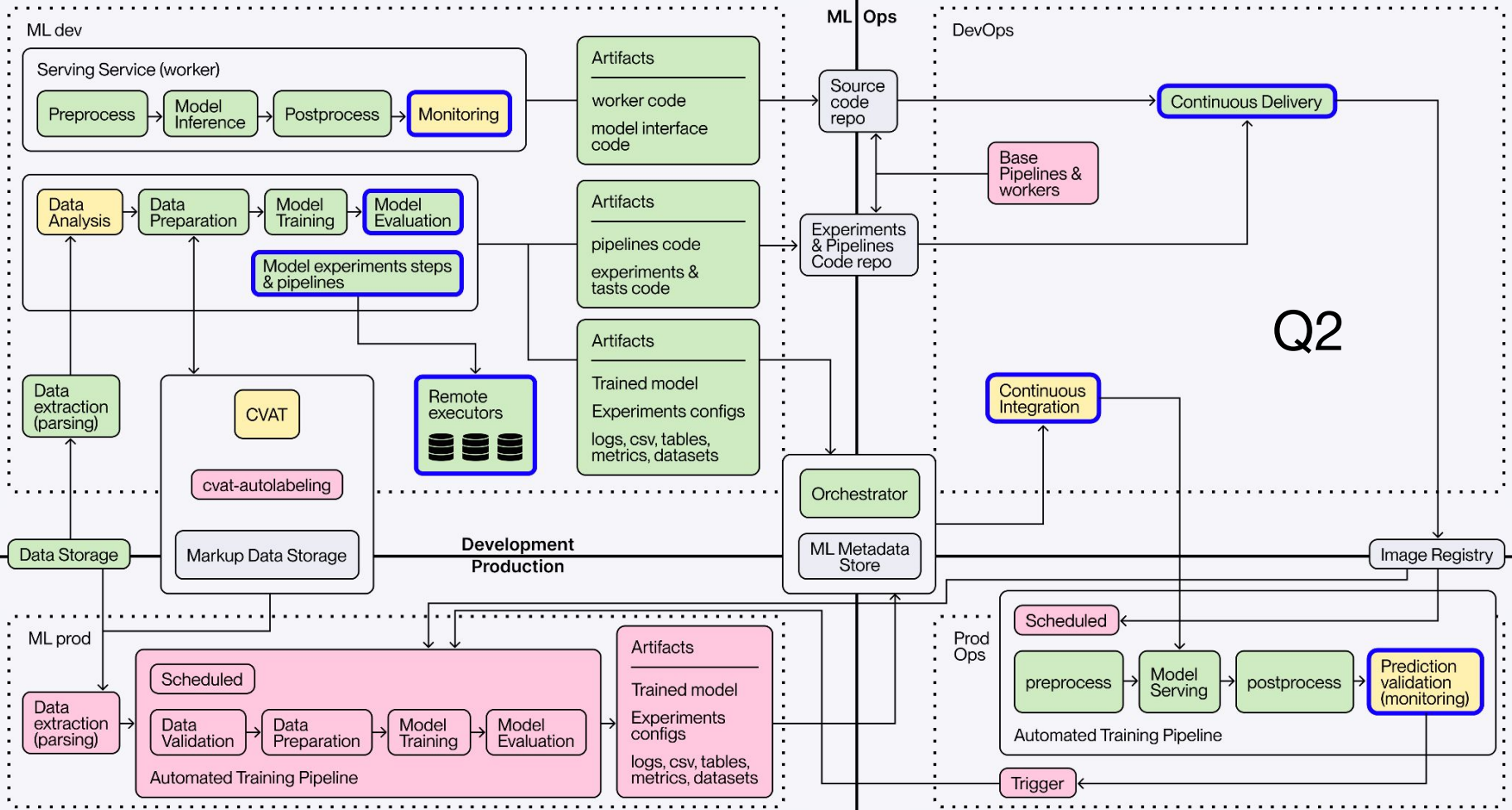
- Поддержка и развитие текущих инструментов
- Ресерч, ресерч и еще раз ресерч
- Оптимизация процесса деплоя и сервисов-оберток моделей (Seldon Core?)
- Рабочее место дата аналитика и дашборды (SuperSet? JupiterHub?)
- Доработка мониторингов и триггеров (Evidently?)
- Доработка валидации данных (Great Expectations?)

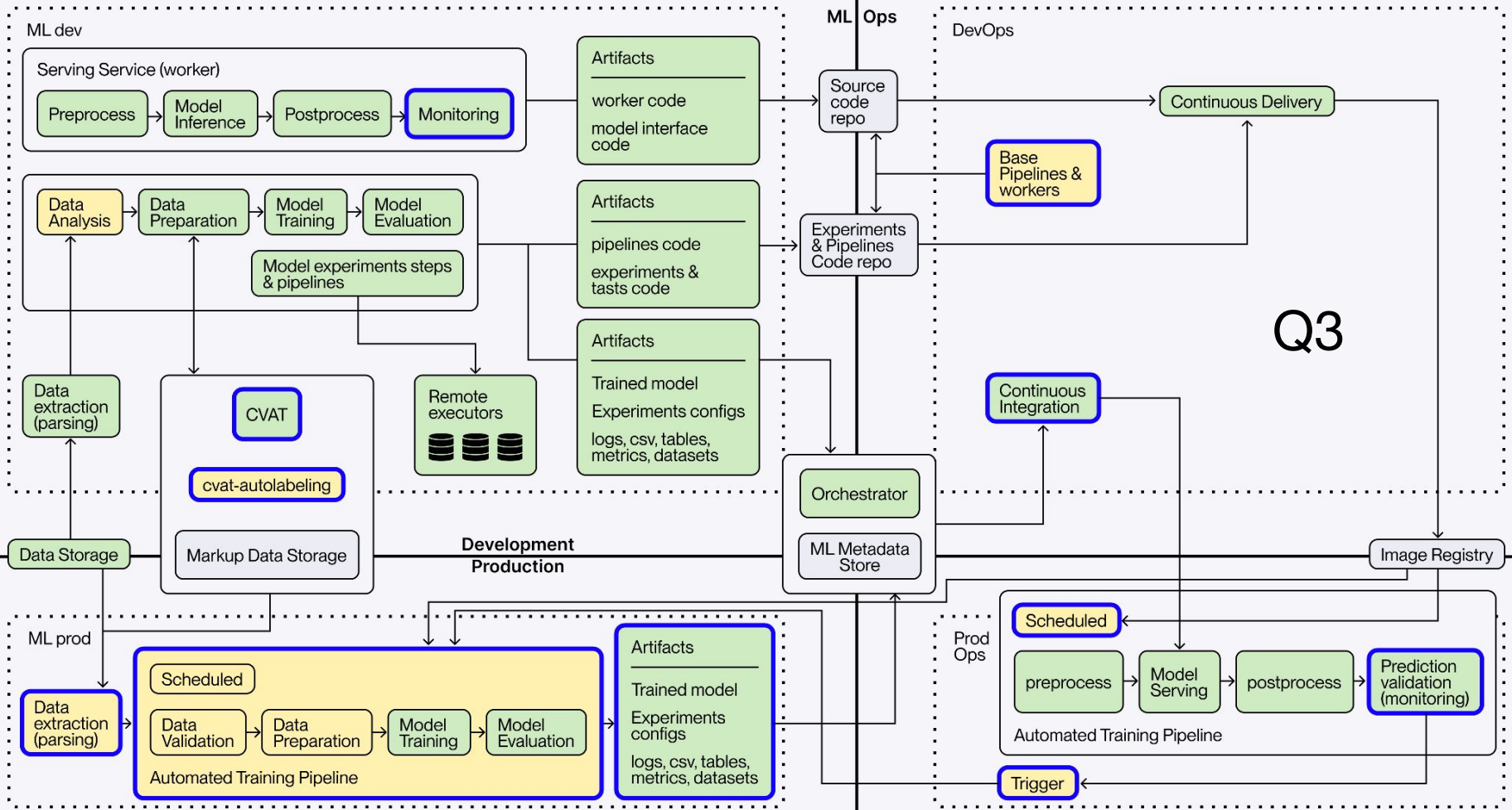
- 01 Введение
- 02 Интеграция 0
- 03 Интеграция I
- 04 Интеграция II
- 05 Заключение

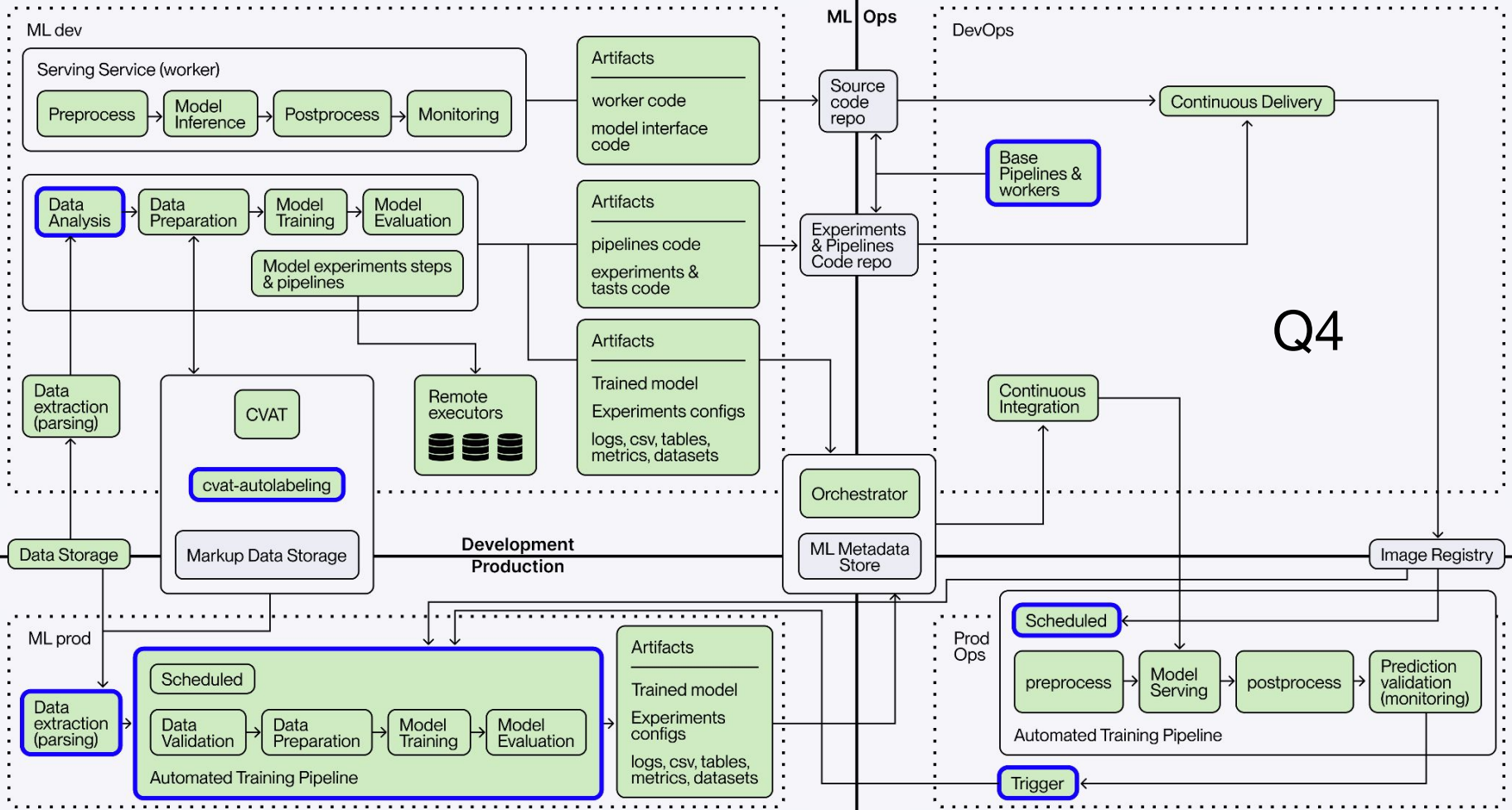
# Что делать или RoadMap по кварталам

- Проекты с неструктурированными данными
- 2 ML Engineer-a
- 1 DevOps
- Team Lead
- Бешенное желание страдать меньше











Игорь  
Терёхин

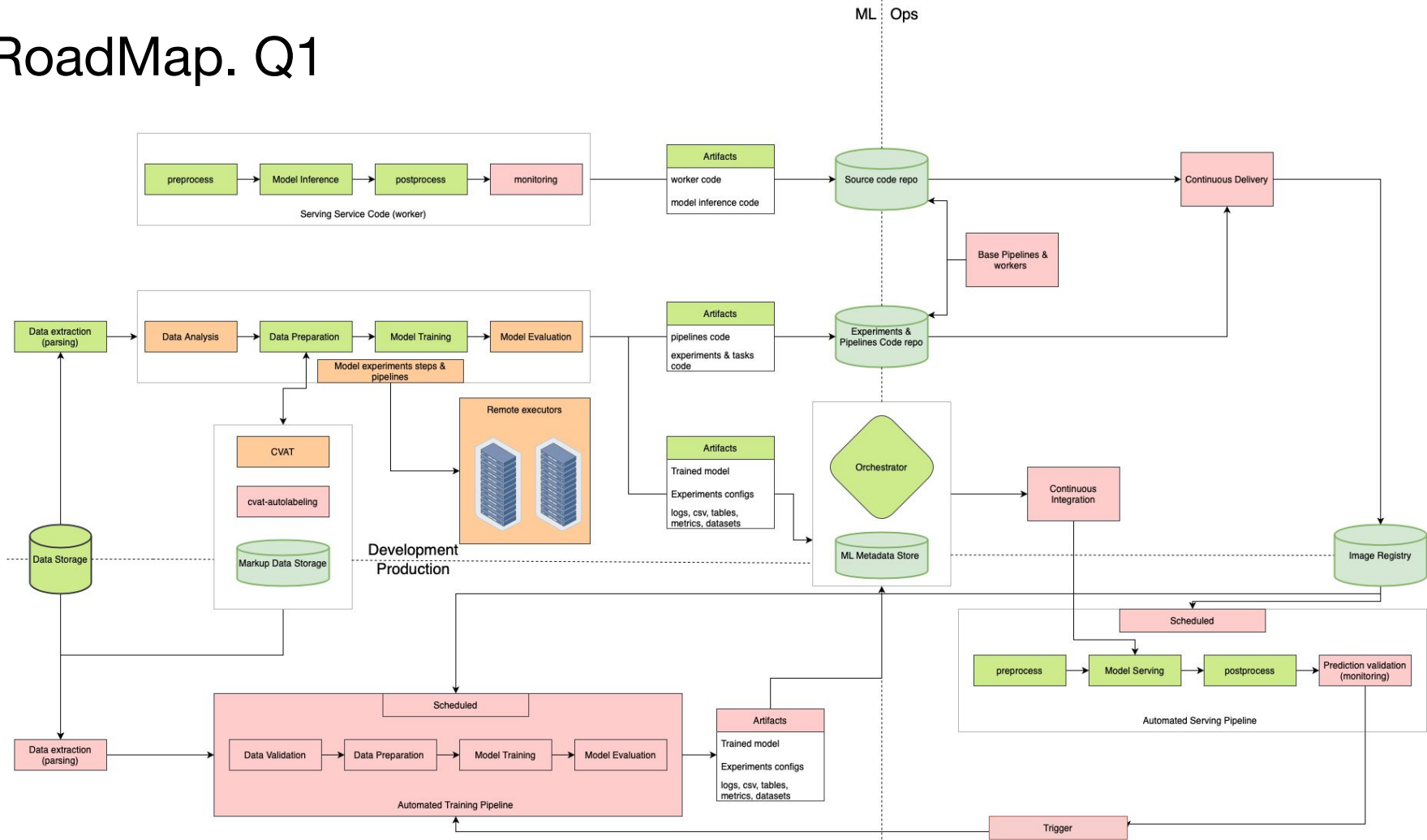
Руководитель отдела  
ML Napoleon IT

@spannenberger

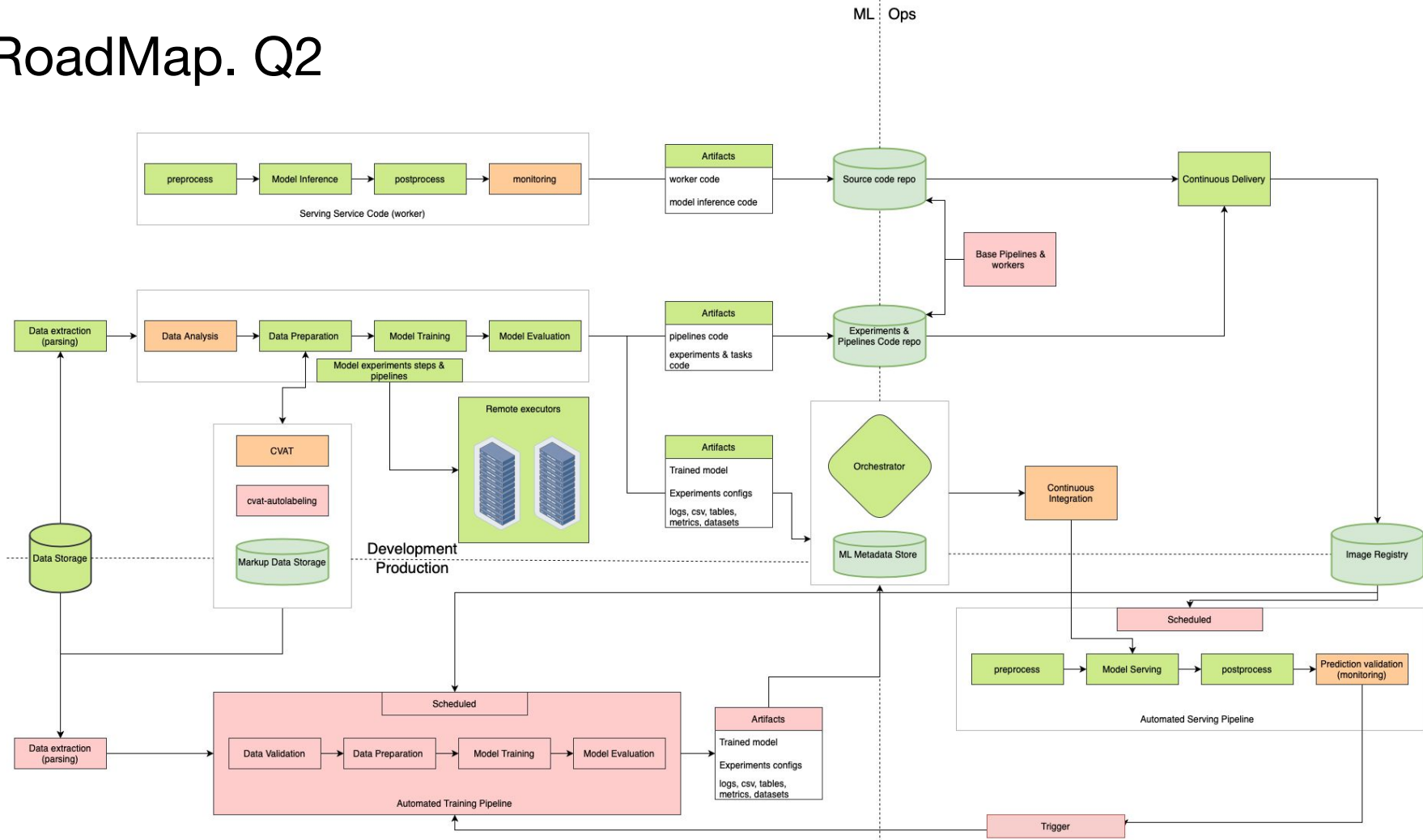
i.terekhin@napoleonit.ru



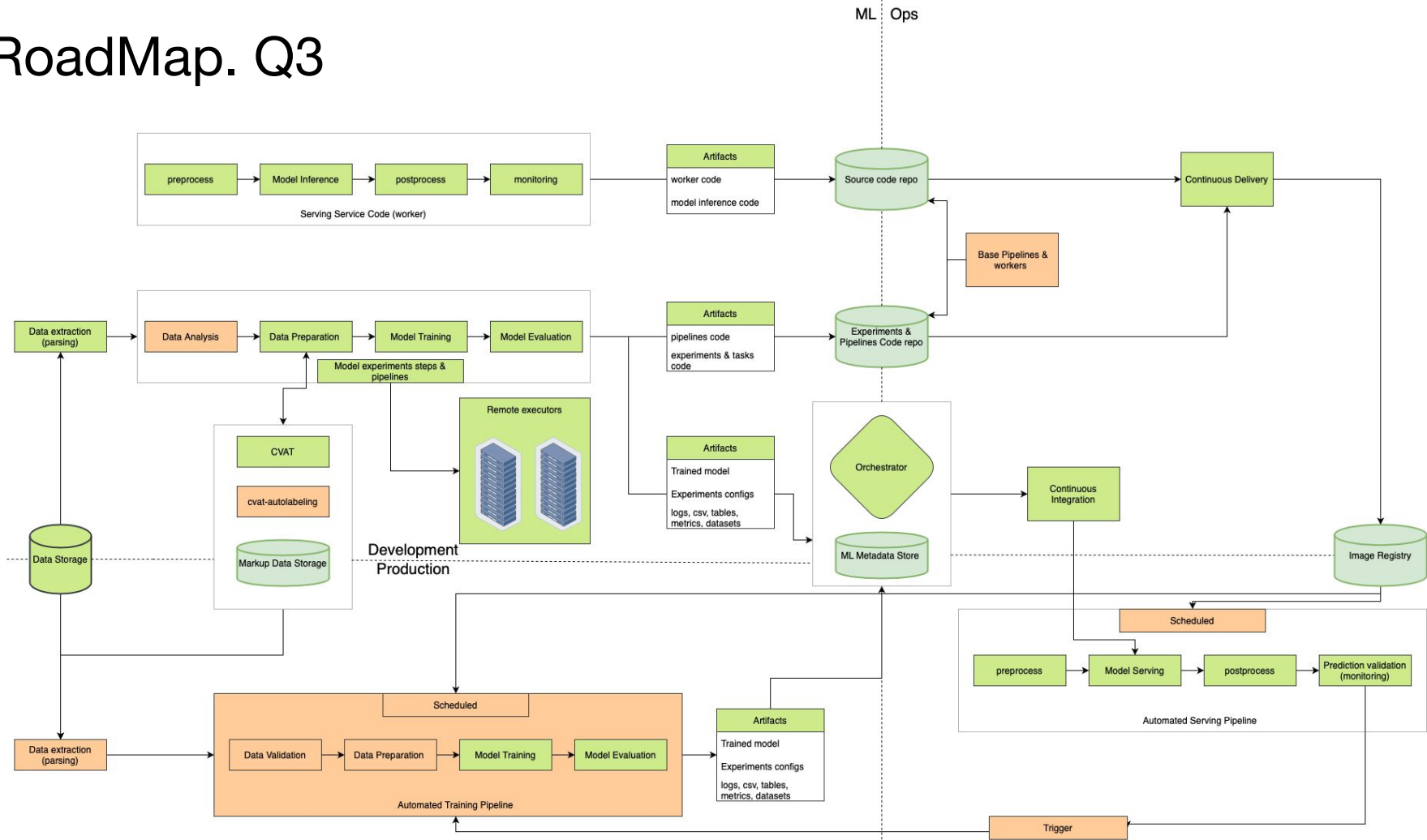
# RoadMap. Q1



# RoadMap. Q2



# RoadMap. Q3



# RoadMap. Q4

