

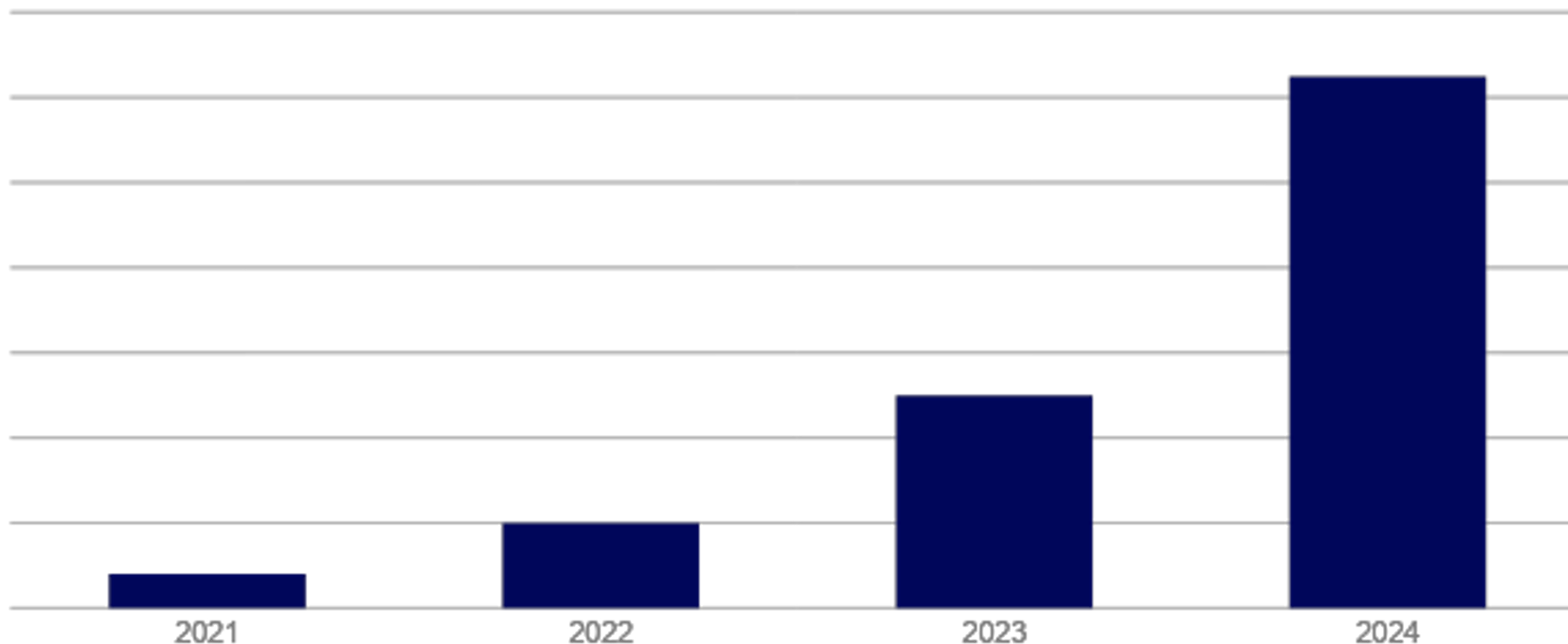


Aqueduct – open source для экономии железа в ML-инференсе



Растет потребность железа для МЛ

График роста потребности в оборудовании для МЛ



Потребность растет

Новые фреймворки и алгоритмы для МЛ-моделей всё более требовательны к потреблению железа



МЛ-сервера

Сервера включает от одного гри-устройства. Они дороже и их тяжелее закупать



С акведуком удастся уменьшить рост потребности железа

Акведук

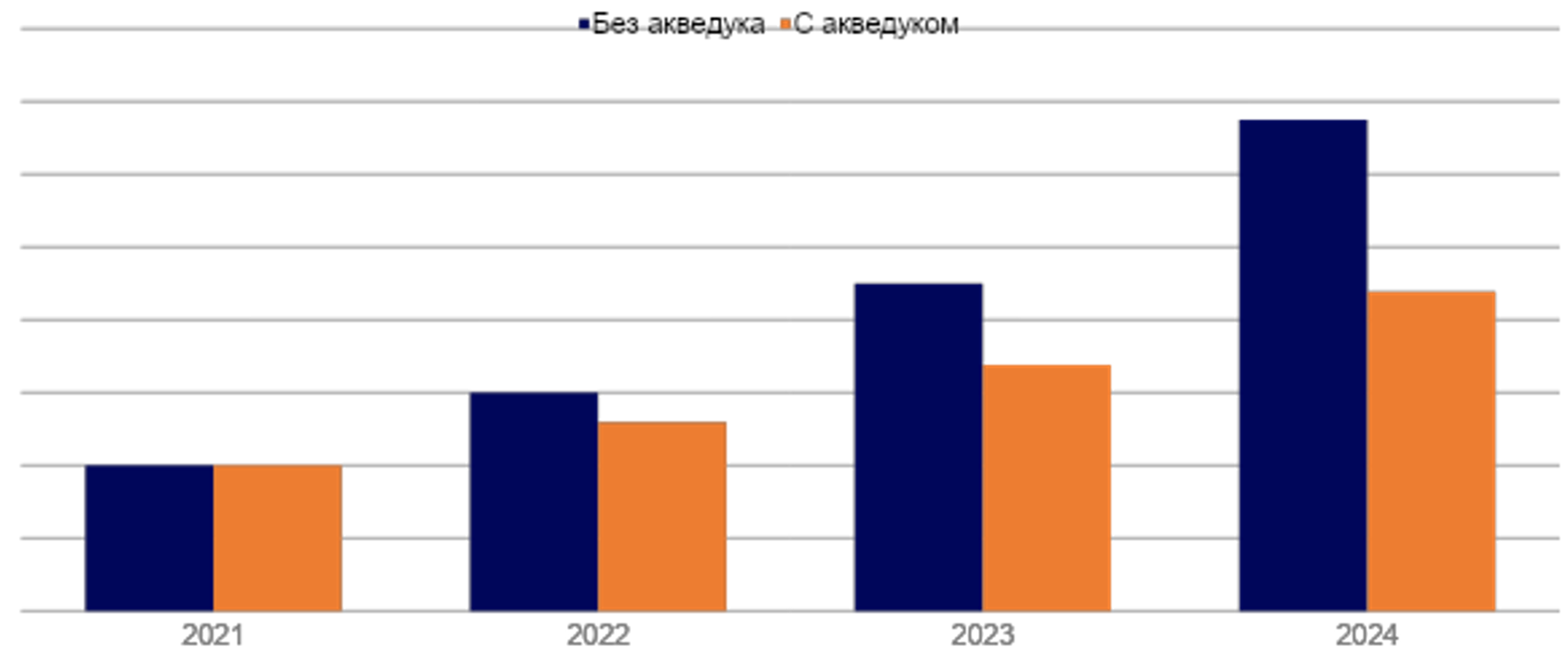
Позволяет уменьшить рост потребности в железе

Потребность растет

Новые фреймворки и алгоритмы для ML-моделей все более требовательны к потреблению железа

ML-сервера

Сервер включает от одного гри-устройства. Они дороже и их тяжелее закупать



Как работает акведук

РАБОТА МОДЕЛИ СОСТОИТ ИЗ ЭТАПОВ



Каждая МЛ-модель — это последовательность шагов вычислений.

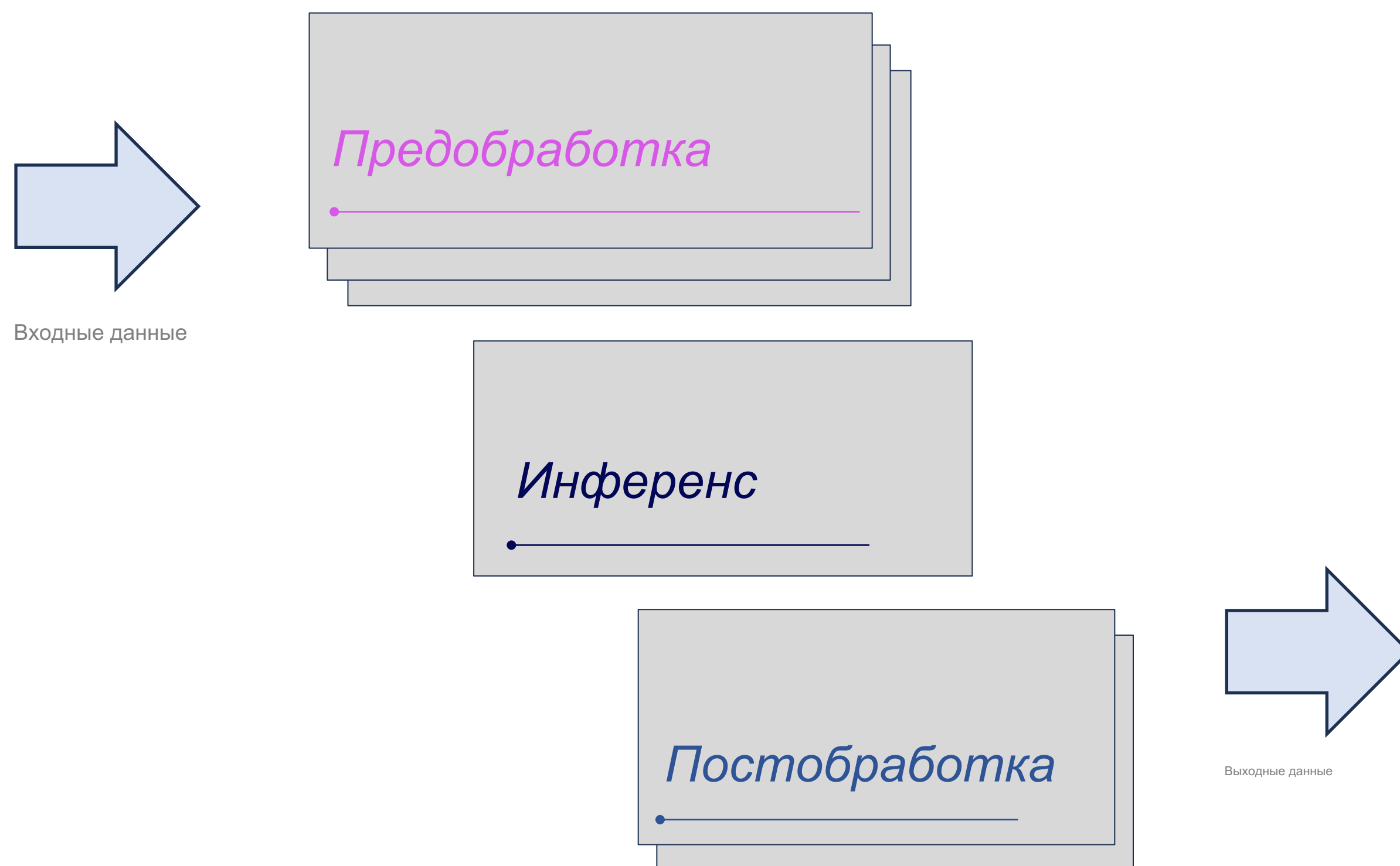
Обычно это три шага:

- 1) подготовка данных для отправки их в модель,
- 1) сам инференс, то есть работа уже самой модели,
- 1) постобработка полученных данных.



Как работает акведук

РАЗДЕЛЯЕМ С ПОМОЩЬЮ АКВЕДУКА

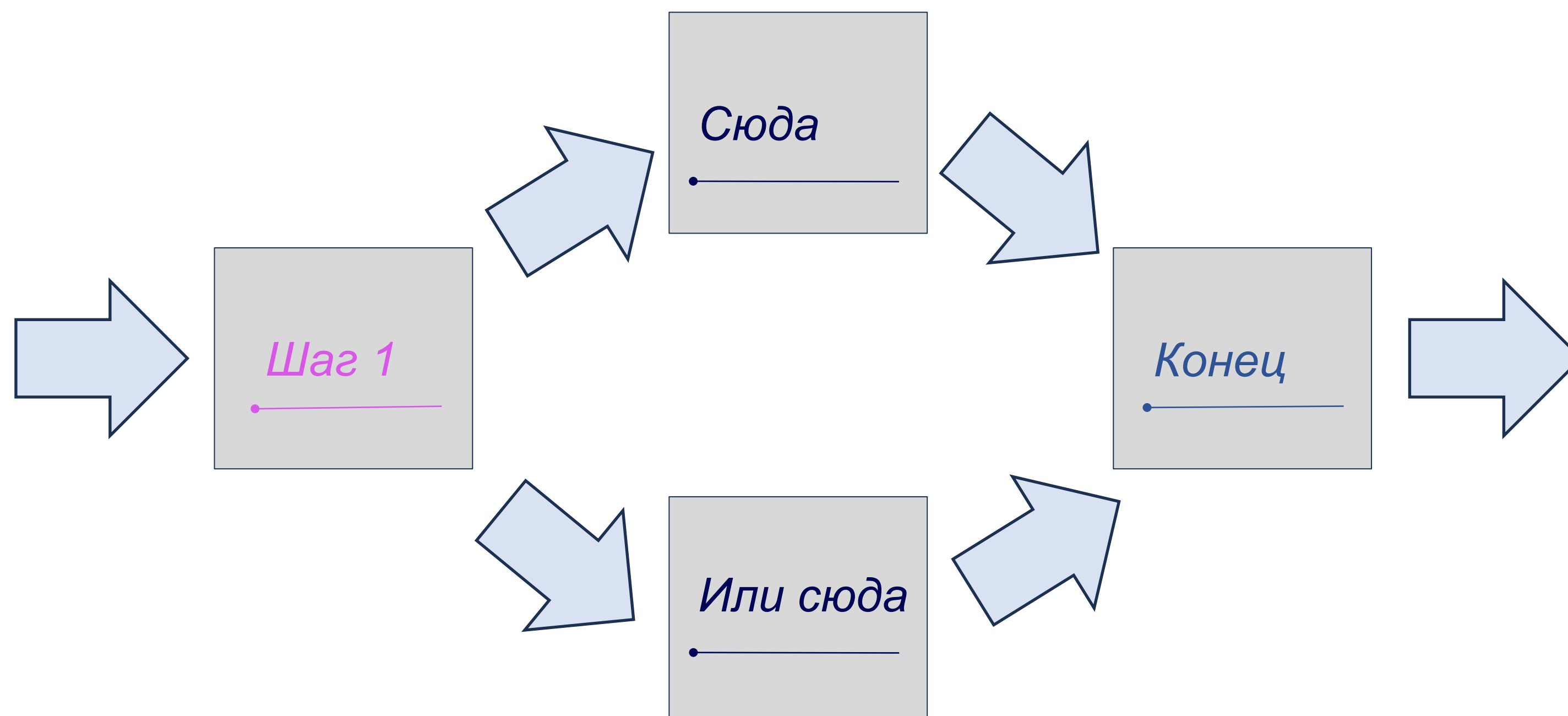


МЛ-модель

С помощью акведука можно разделить такую логику на разные этапы и сделать так, чтобы они выполнялись в разных процессах.



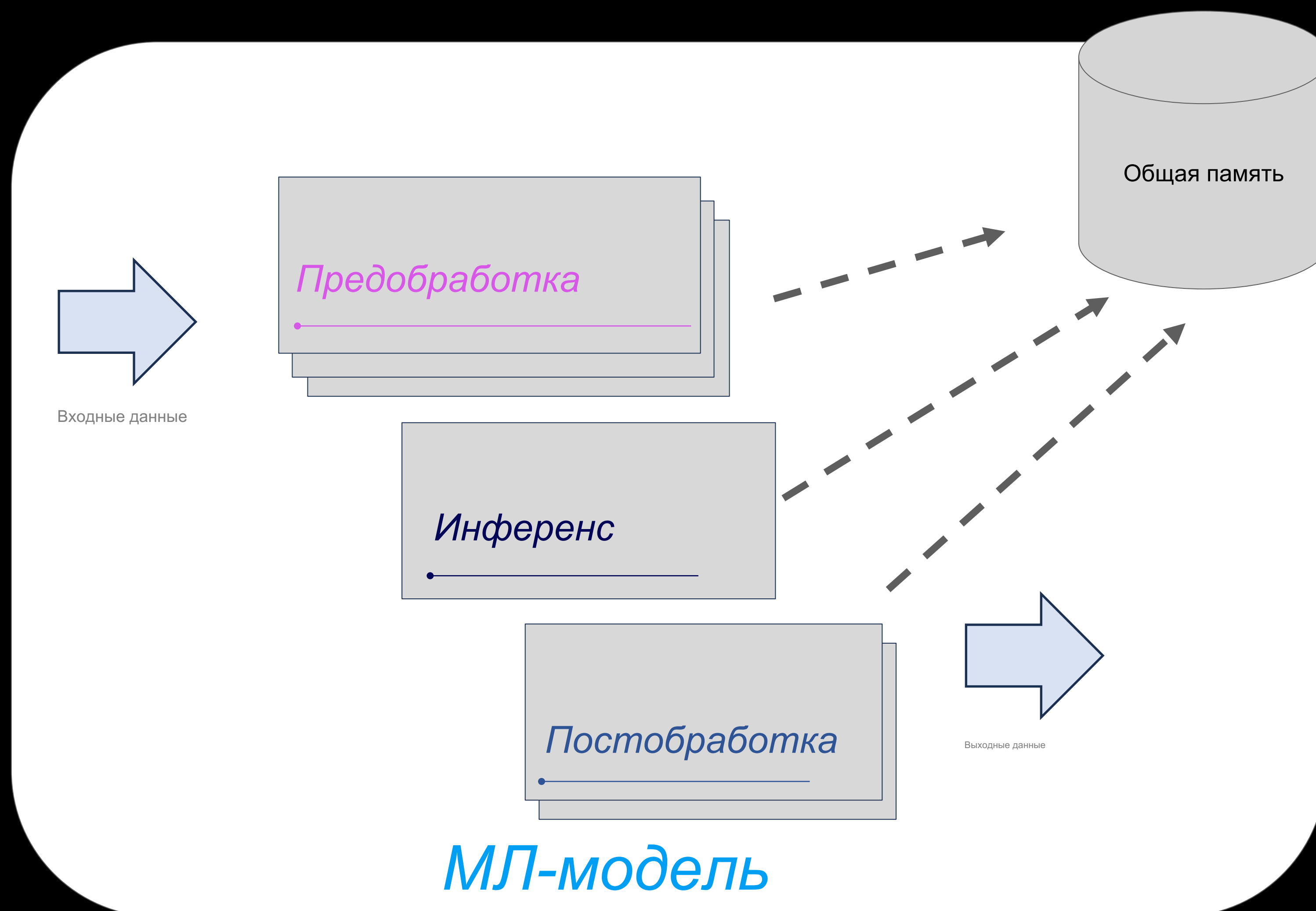
Условный переход между шагами



С помощью акведука можно реализовать условный переход между шагами



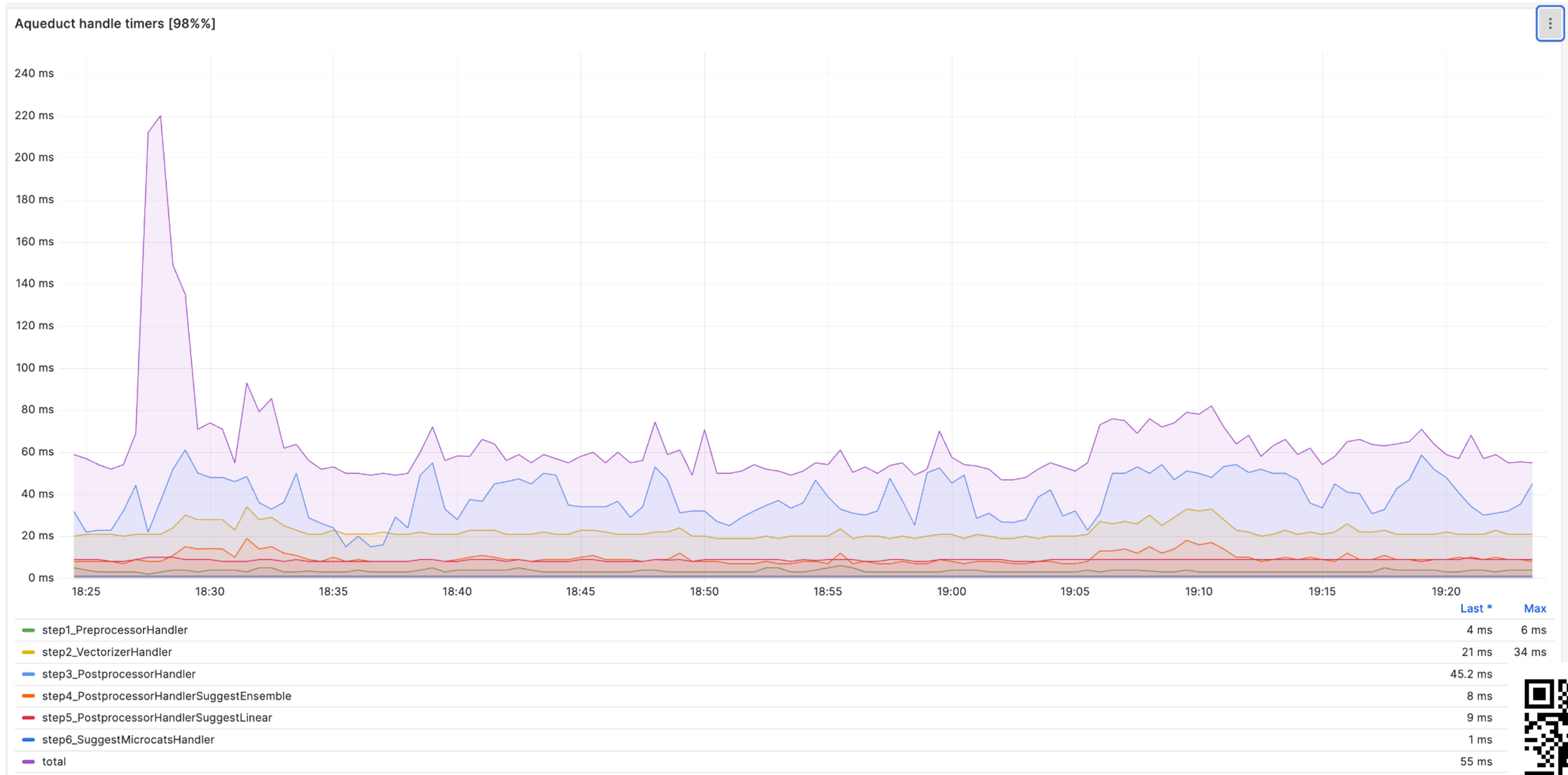
Шаренная память



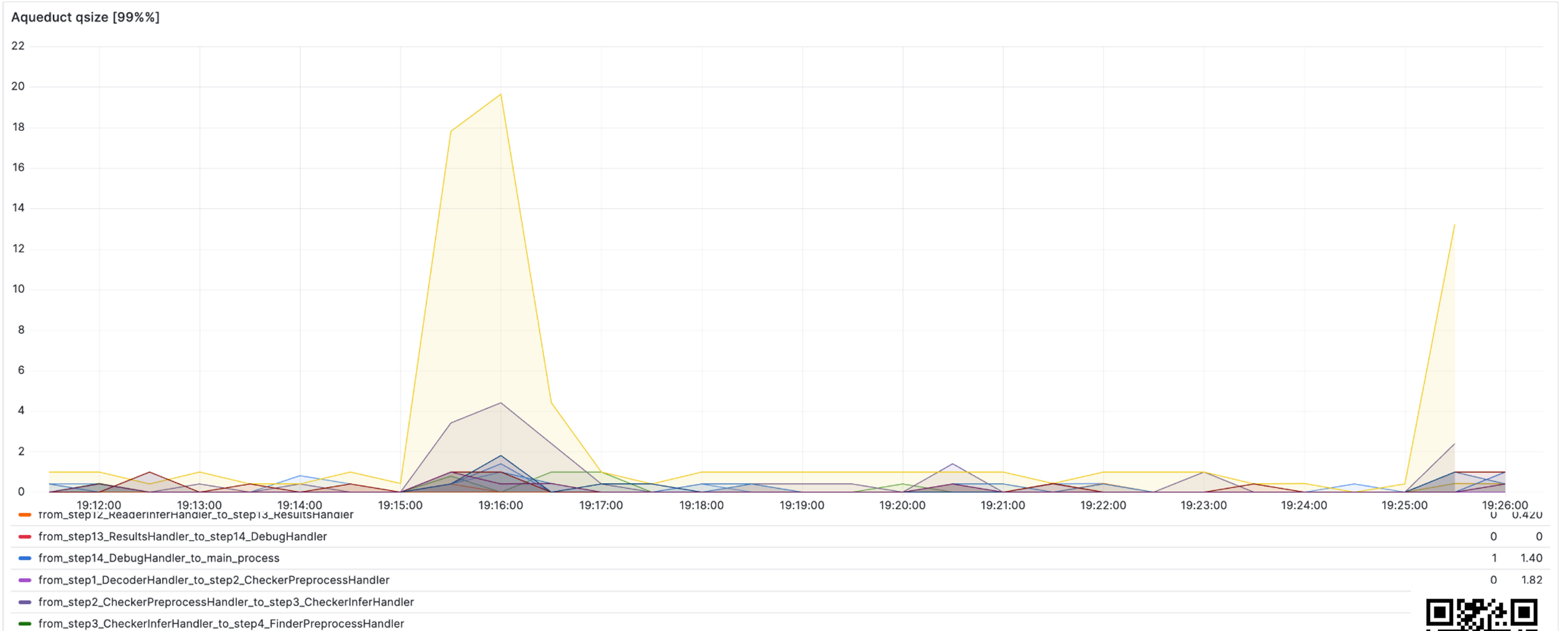
С помощью акведука можно реализовать использование общей памяти



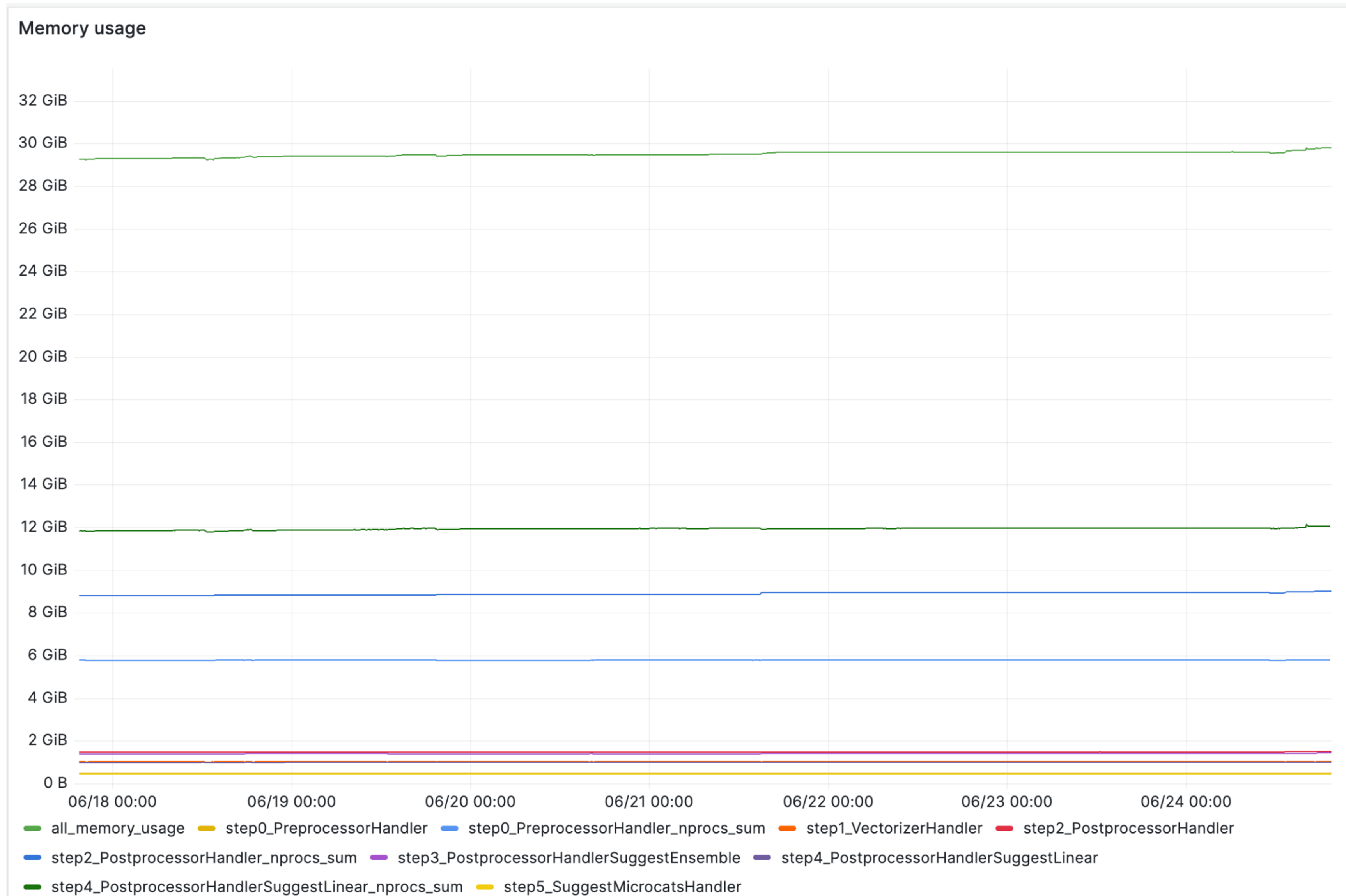
Метрики



Метрики



Метрики



Акведук — универсальное решение

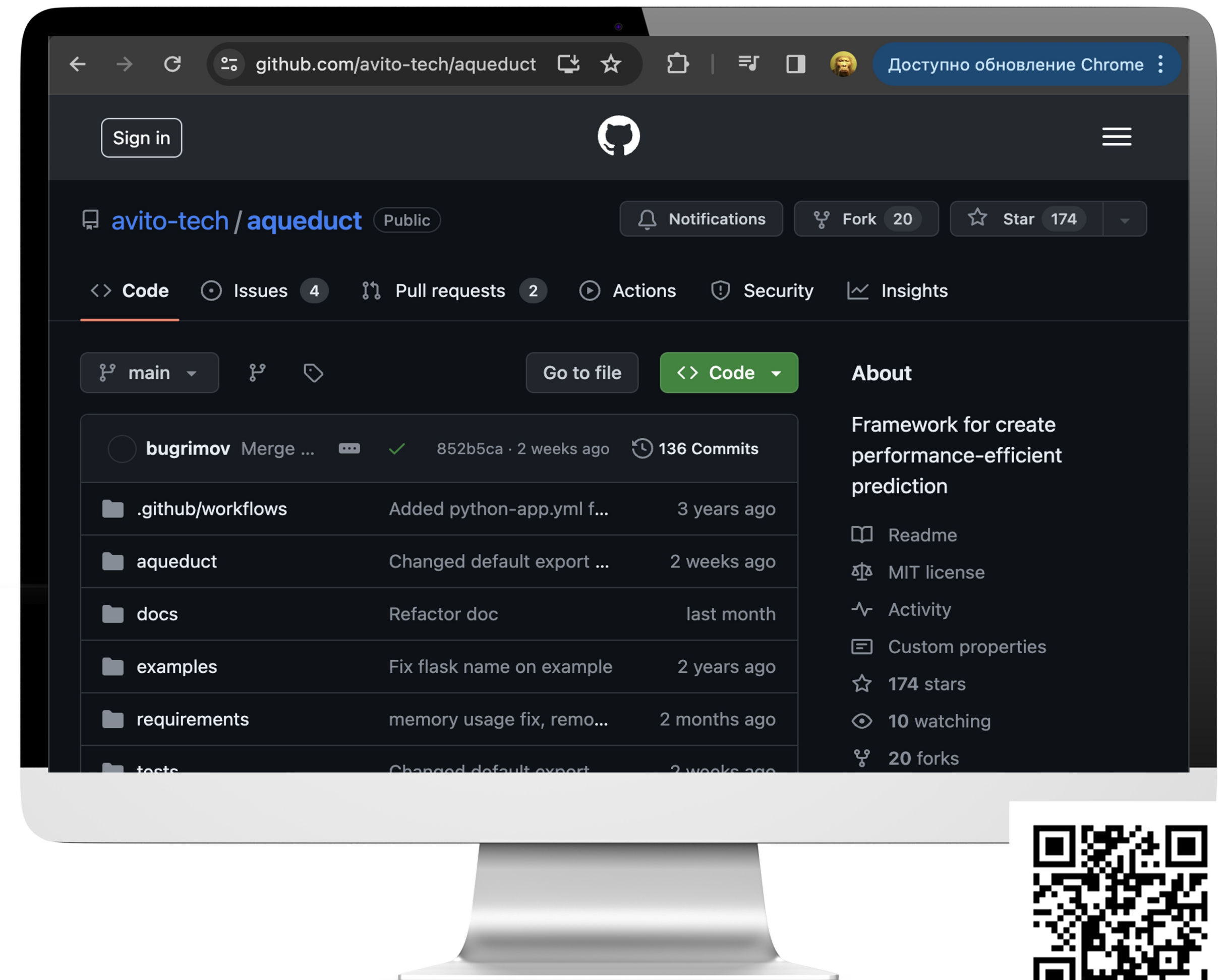
Преимущества от выкладывания в Open Source

- + изменения/фиксы, которые мы не планировали\ не замечали,
- + изменения от других команд внутри компании,
- + такая разработка заметно мотивирует команду
- + используют решение в своих продуктах

Недостатки от выкладывания в Open Source

- Нужно поддерживать решение, отделенное от внутренних костылей
- Больше ответственности на внесение изменений
- Больше задач на документацию

Работает для всех типов моделей.
Доступно на github, бесплатно под лицензией MIT.
Есть внедрения в других компаниях





Aqueduct – open source для экономии железа в ML-инференсе

