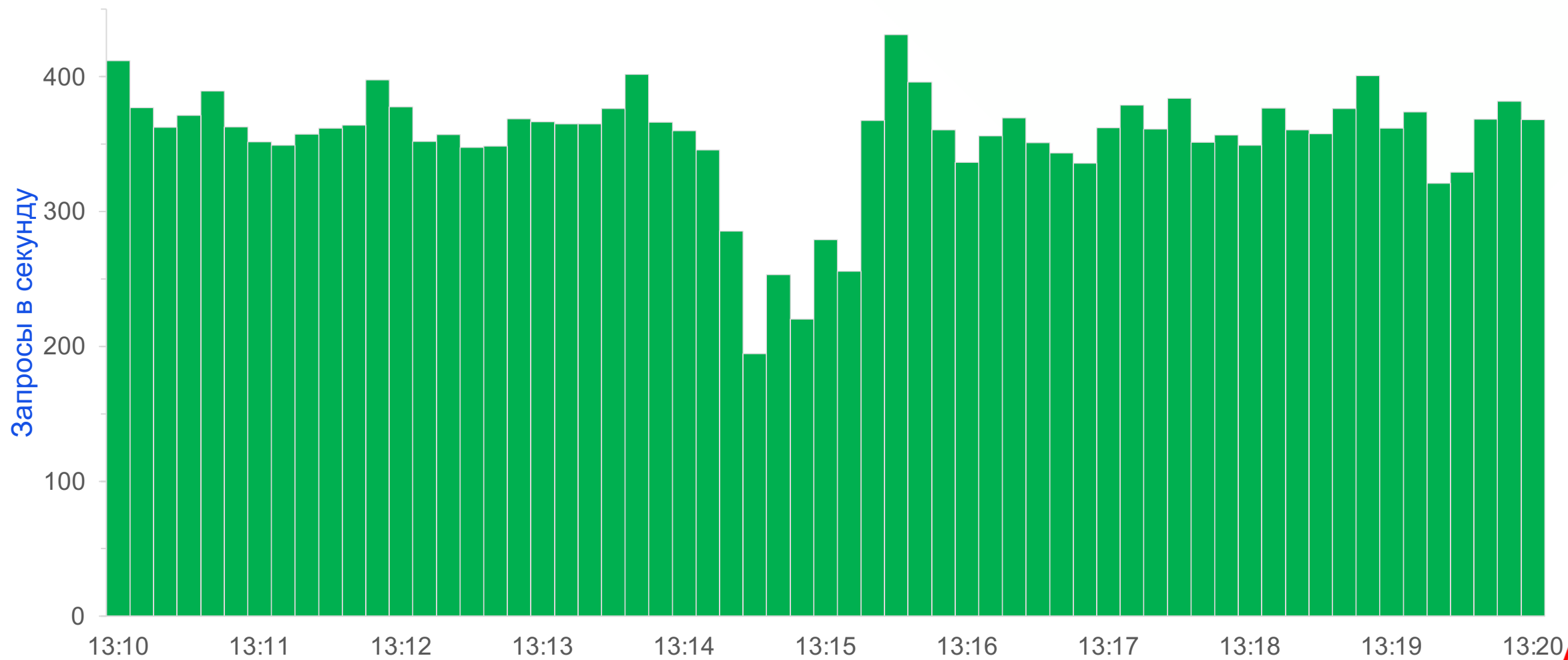


Как сделать деплой незаметным для пользователя

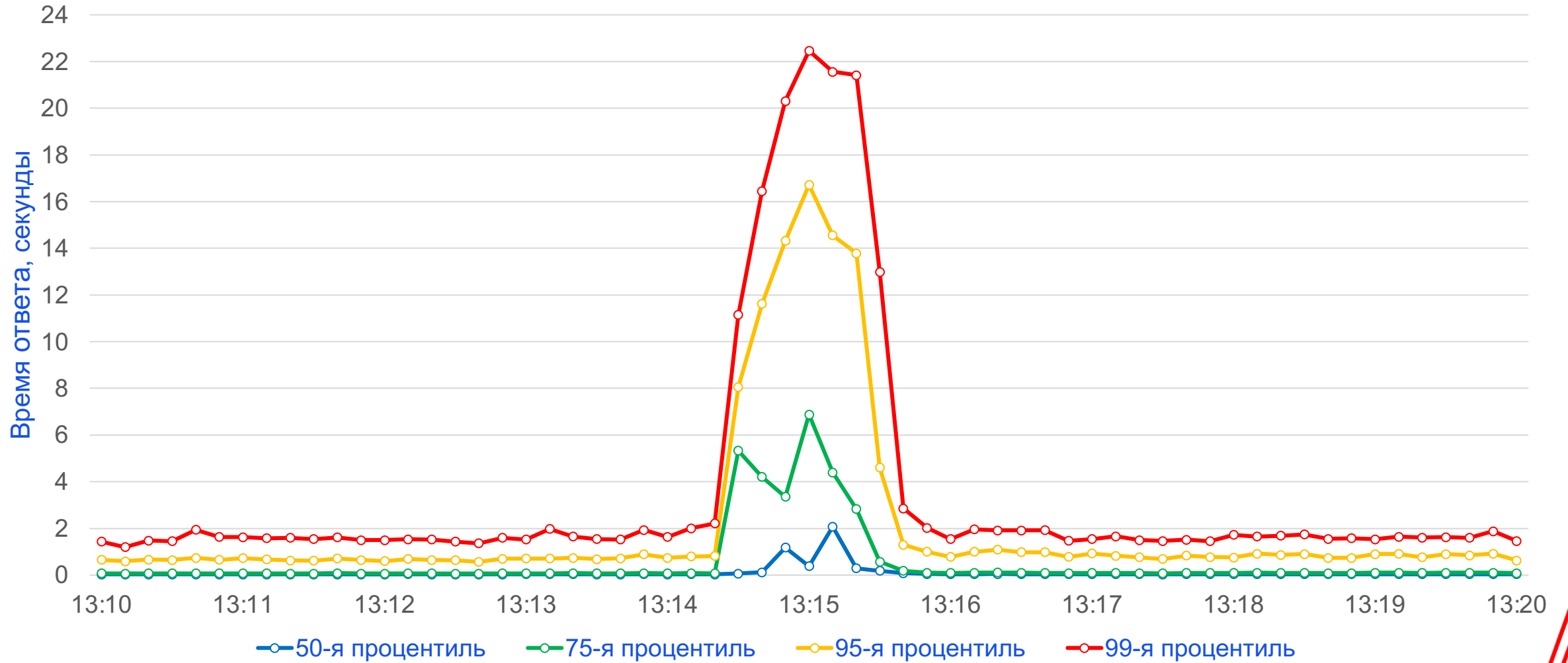


Деплой

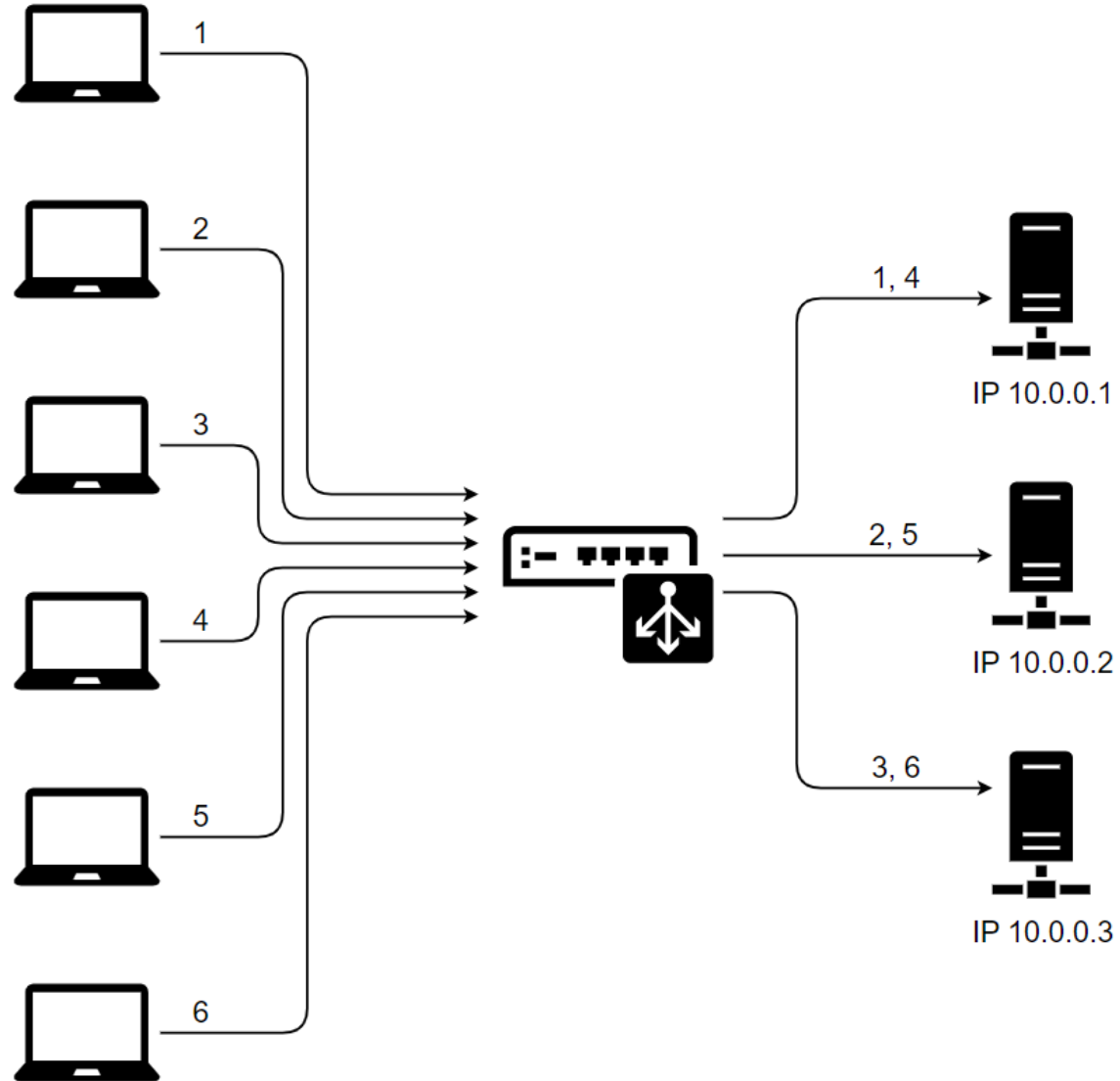




Деплой



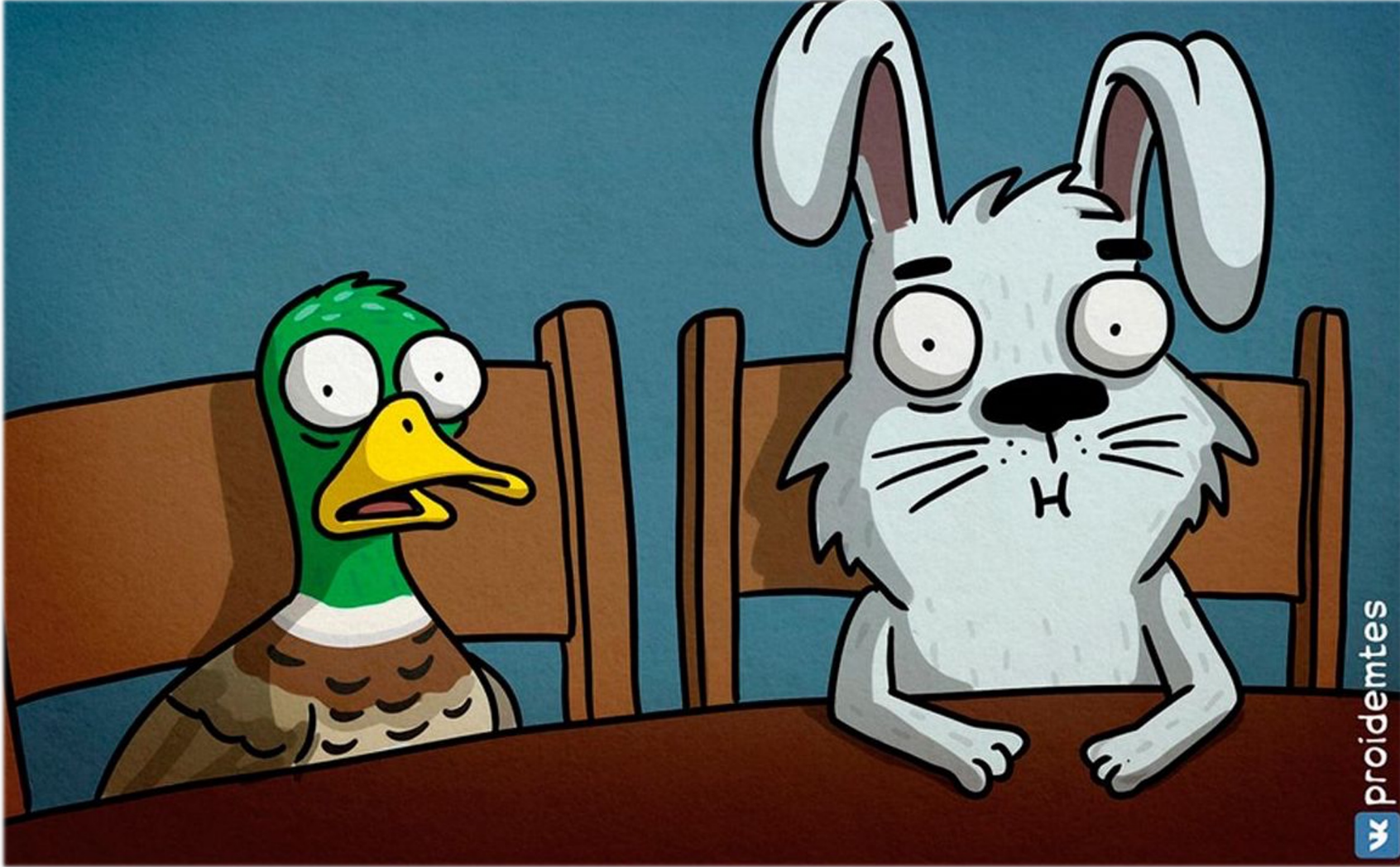
Round robin балансировка



Пространство решений

- уменьшить время прогрева
 - вертикальное масштабирование
 - горизонтальное масштабирование
 - оптимизация производительности
- предварительный прогрев
 - нагрузочными тестами
 - реальным трафиком
- использовать балансировку с липкими сессиями?!?!





EWMA балансировка

EWMA - *Exponentially Weighted Moving Average*.

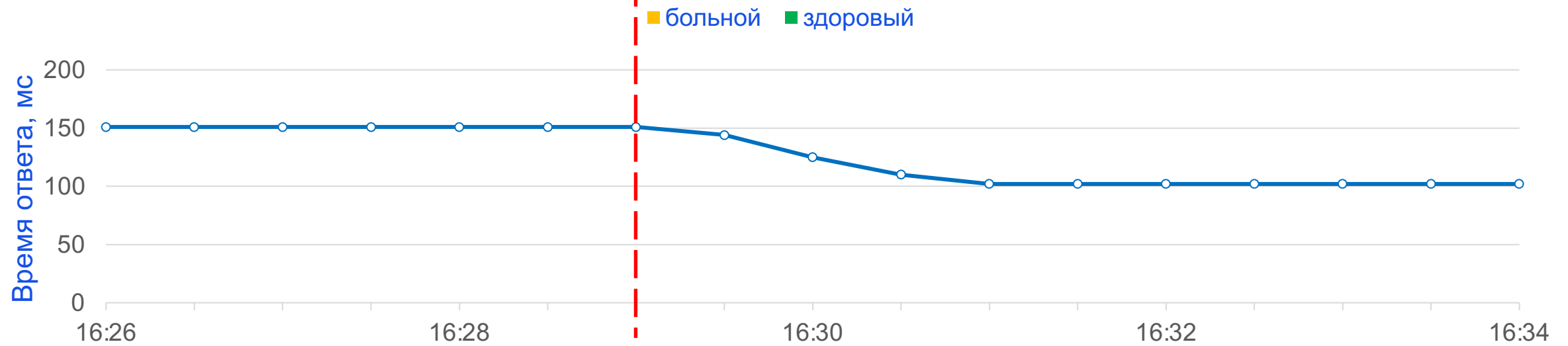
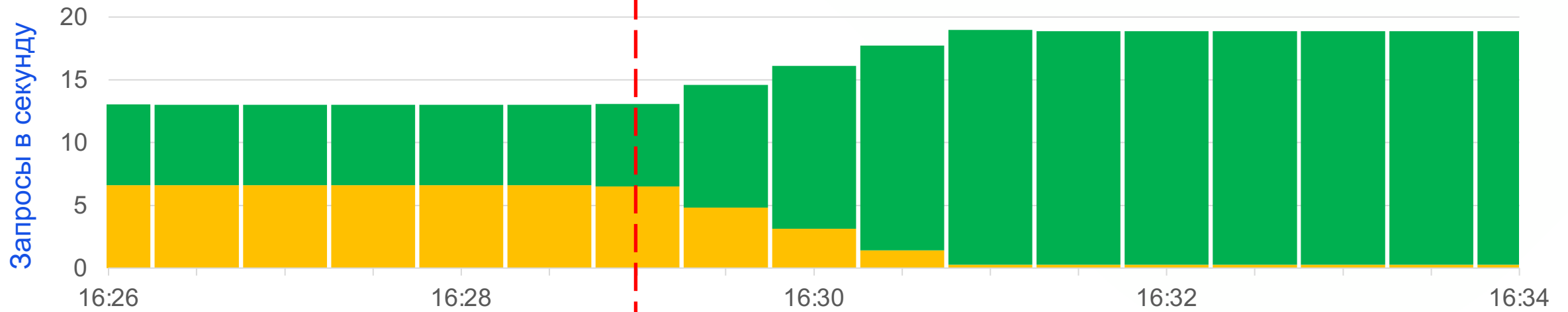
Запросы распределяются таким образом, чтобы выравнять **время ответа** на всех сервисах. Предпочтение отдаётся тому сервису, который быстрее отвечает.

Время ответа вычисляется как экспоненциально взвешенное скользящее среднее.

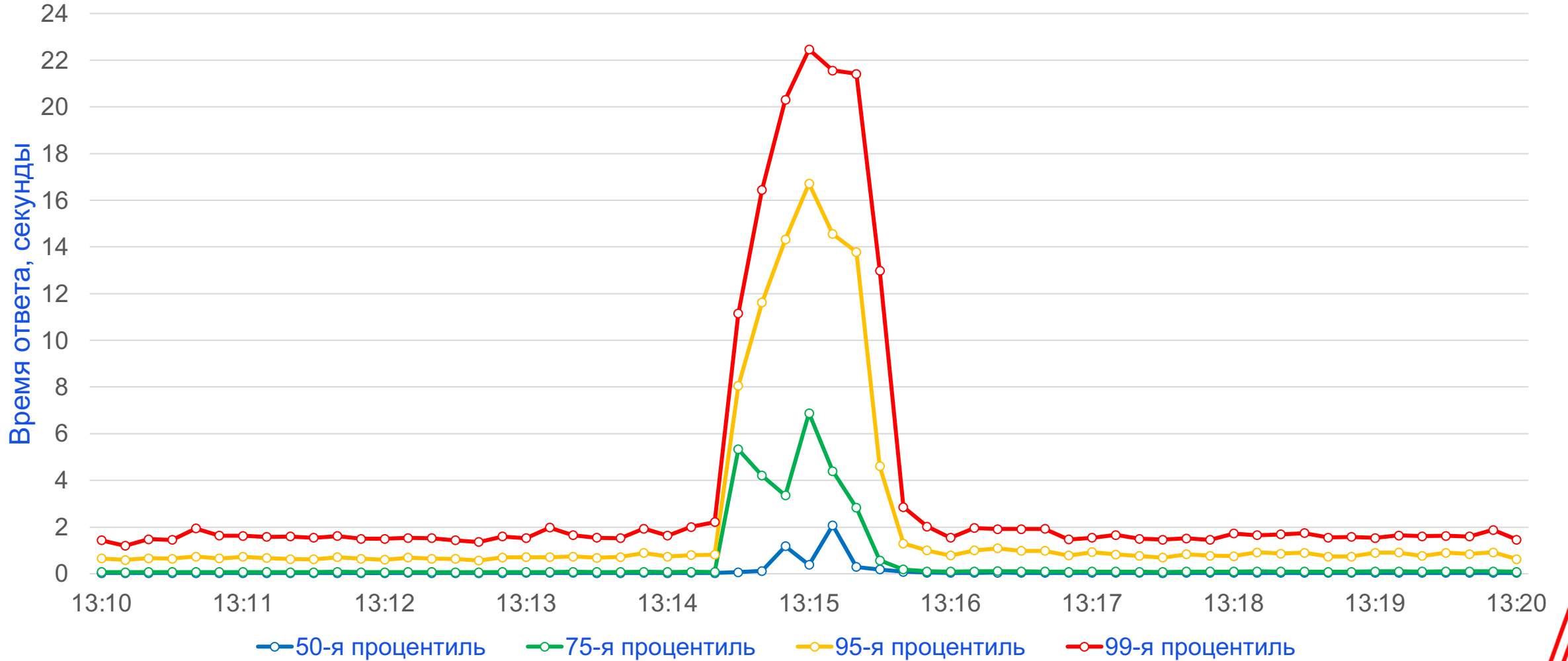
В production впервые начали использоваться в 2014г.

- Finagle - Twitter
- Linkerd - Service Mesh для Kubernetes

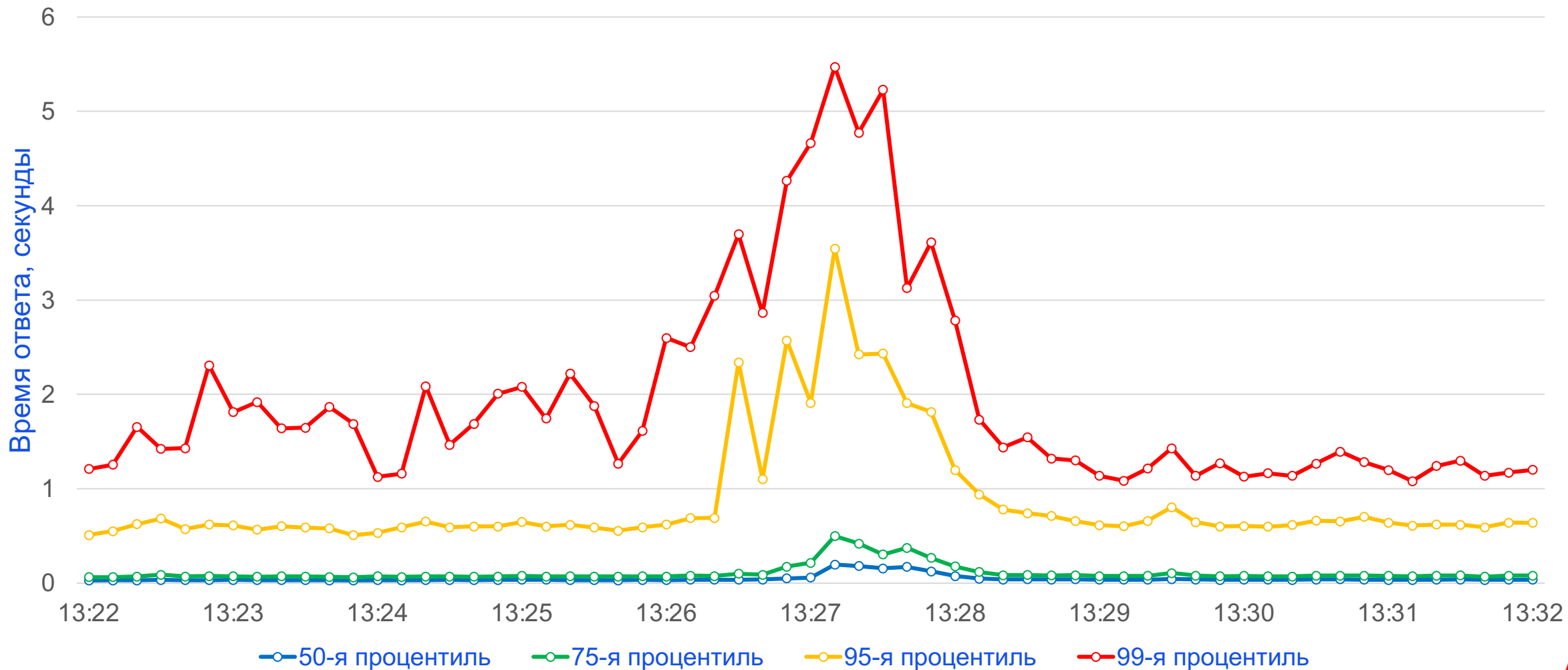
Синтетический тест



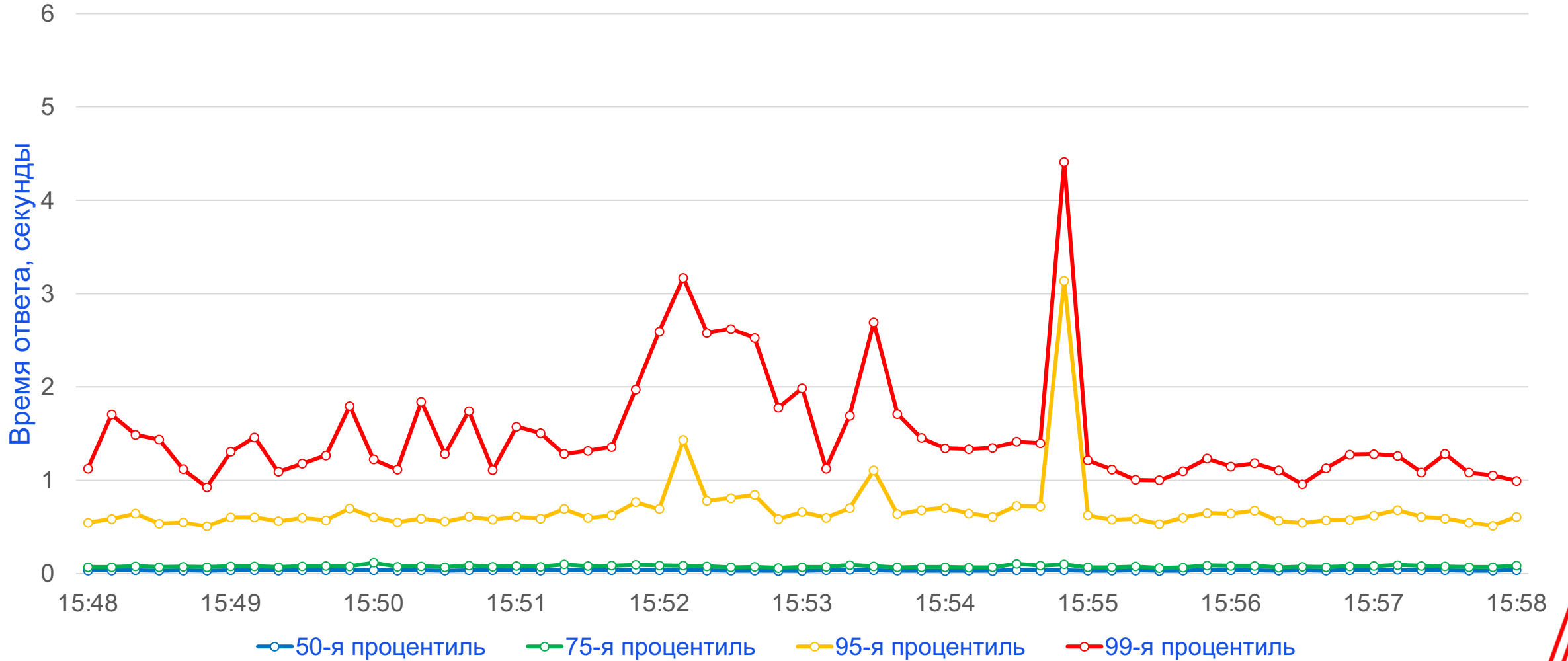
Деплой + Round Robin



Деплой + EWMA



Деплой + EWMA + прогрев







Технические детали

```
apiVersion: apps/v1
  kind: Deployment
spec:
  strategy:
    type: RollingUpdate
    rollingUpdate:
      maxUnavailable: 0
      maxSurge: 1
  minReadySeconds: 60
  readinessProbe:
    ...
```


Технические детали

```
apiVersion: networking.k8s.io/v1
kind: Ingress
metadata:
  annotations:
    nginx.ingress.kubernetes.io/load-balance: ewma
```

Вопросы?



Спасибо!

