

Санкт-Петербургский государственный университет
Факультет прикладной математики – процессов управления

ЭФФЕКТИВНАЯ ПАРАЛЛЕЛЬНАЯ РЕАЛИЗАЦИЯ МЕТОДОВ МАШИННОГО ОБУЧЕНИЯ С ИСПОЛЬЗОВАНИЕМ ГРАФИЧЕСКОГО ПРОЦЕССОРА

Щекалёв Александр Сергеевич

Научный руководитель, к.т.н., доцент Гришкин В.М.

Постановка задачи

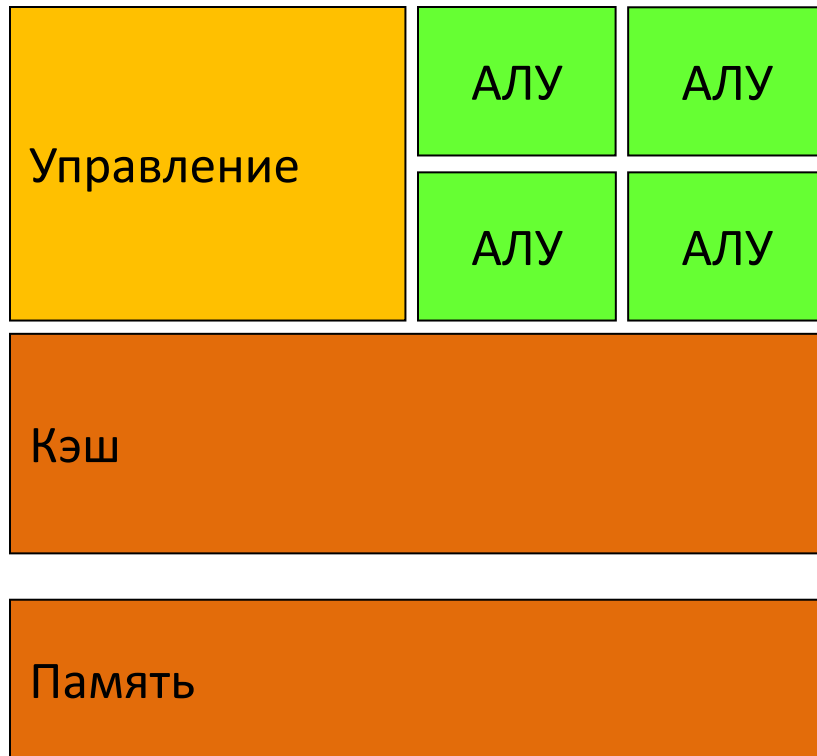
2

- Параллельная реализация метода опорных векторов и метода бустинга деревьев решений
 - Выделение вычислительно-сложных этапов каждого метода
 - Распараллеливание и перенос на графический процессор
 - Оптимизация с учётом архитектуры видеокарты

Графические процессоры

3

Центральный процессор



Графический процессор



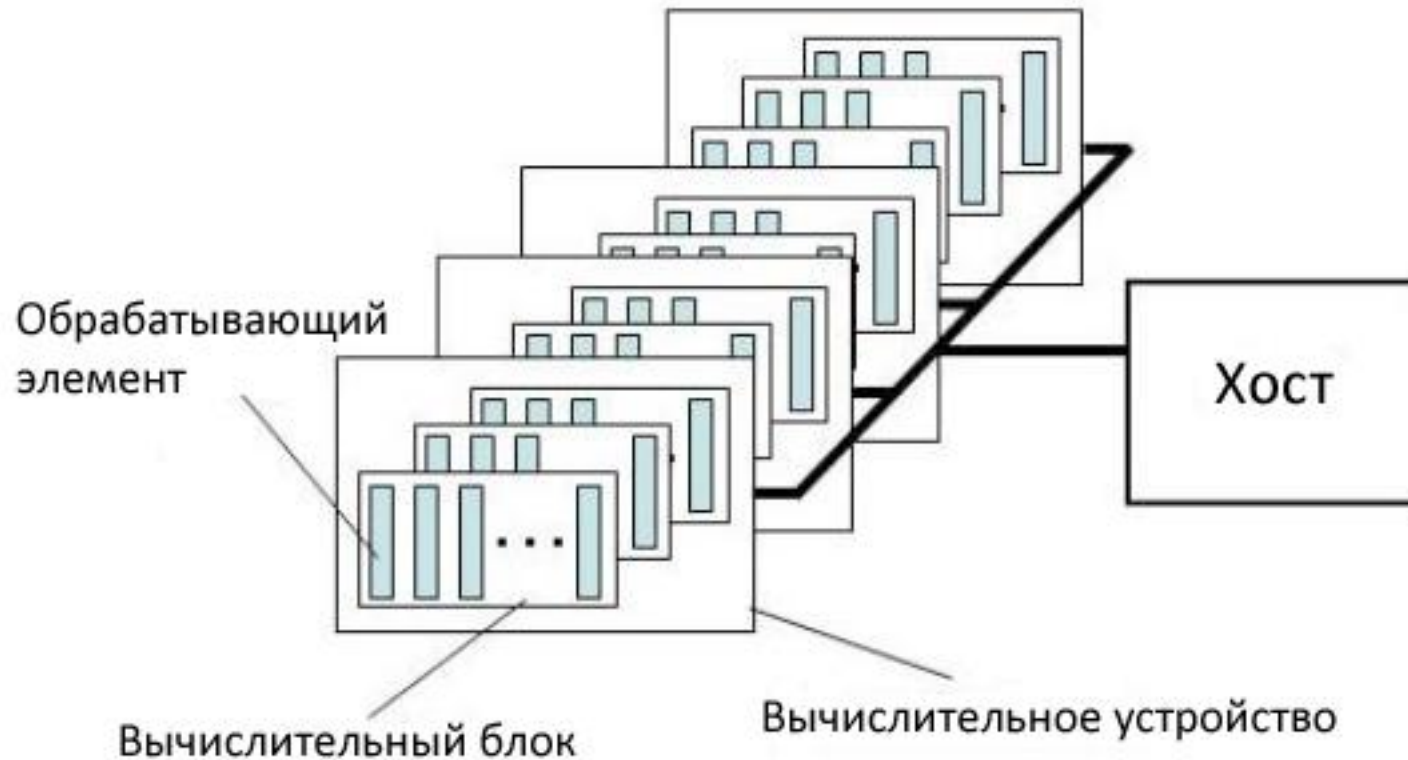
Доступные технологии

4

- Шейдеры
 - ▣ Решение задач, связанных с визуализацией графики
- Технологии от производителей видеокарт
 - ▣ NVIDIA CUDA
 - ▣ AMD Stream
- OpenCL
 - ▣ Универсальный стандарт

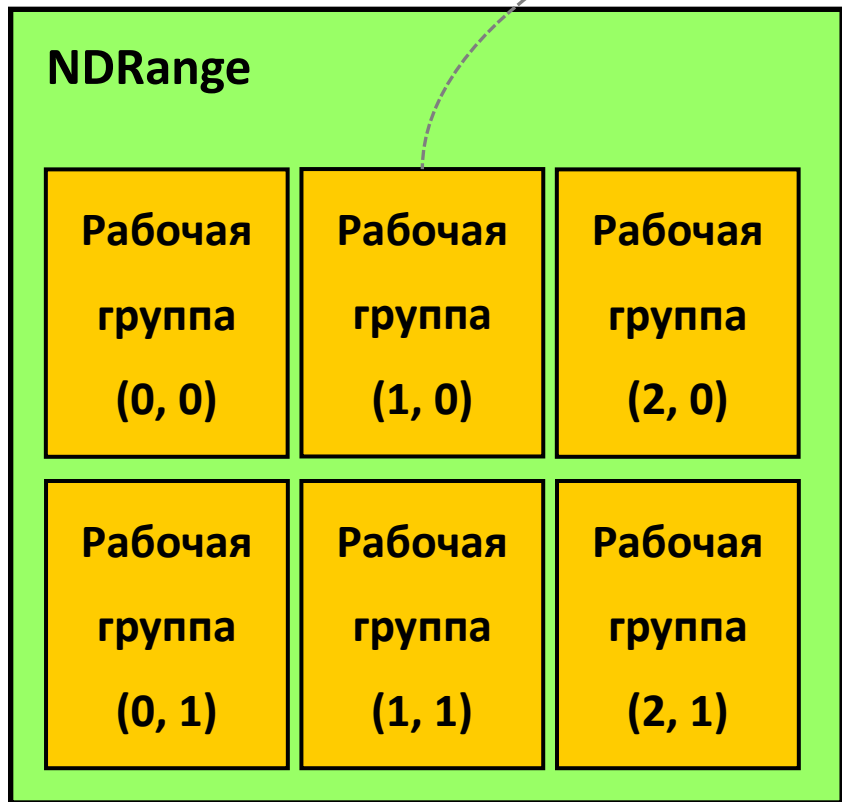
OpenCL. Модель платформы

5



OpenCL. Модель выполнения

6



OpenCL. Модель памяти

7

- Доступно четыре вида памяти
 - ▣ Глобальная память
 - ▣ Константная память
 - ▣ Локальная память
 - ▣ Собственная память

Метод опорных векторов

8

□ Задача:

$$\min_{\alpha} \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha$$

$$0 \leq \alpha_i \leq C \quad i = 1 \dots L$$

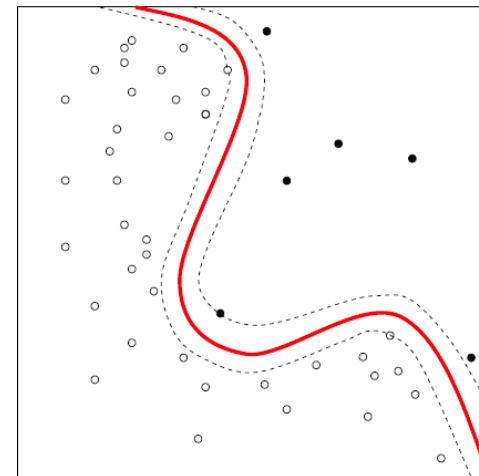
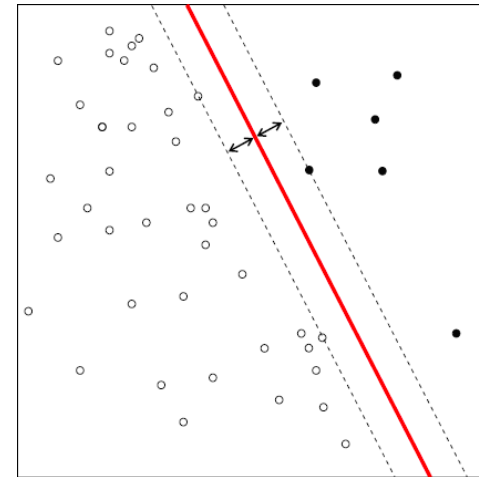
$$y^T \alpha = 0$$

$$Q_{ij} = y_i y_j \phi(x_i, x_j)$$

□ Примеры ядер:

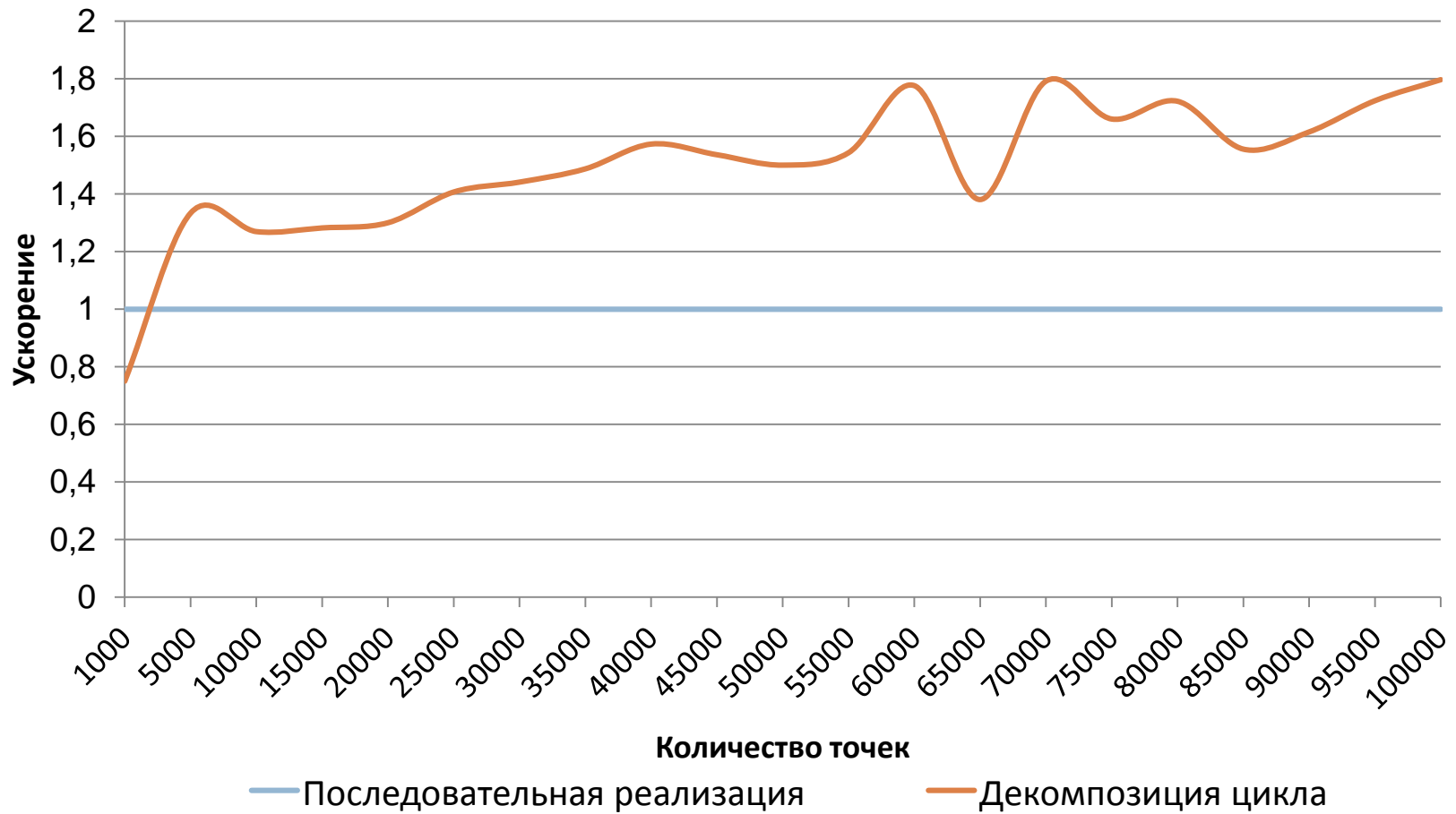
$$\phi(x_i, x_j) = (x_i \cdot x_j + a)^b$$

$$\phi(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$



Наивная реализация

9



Обращения к памяти

10

- Развёртывание по строкам

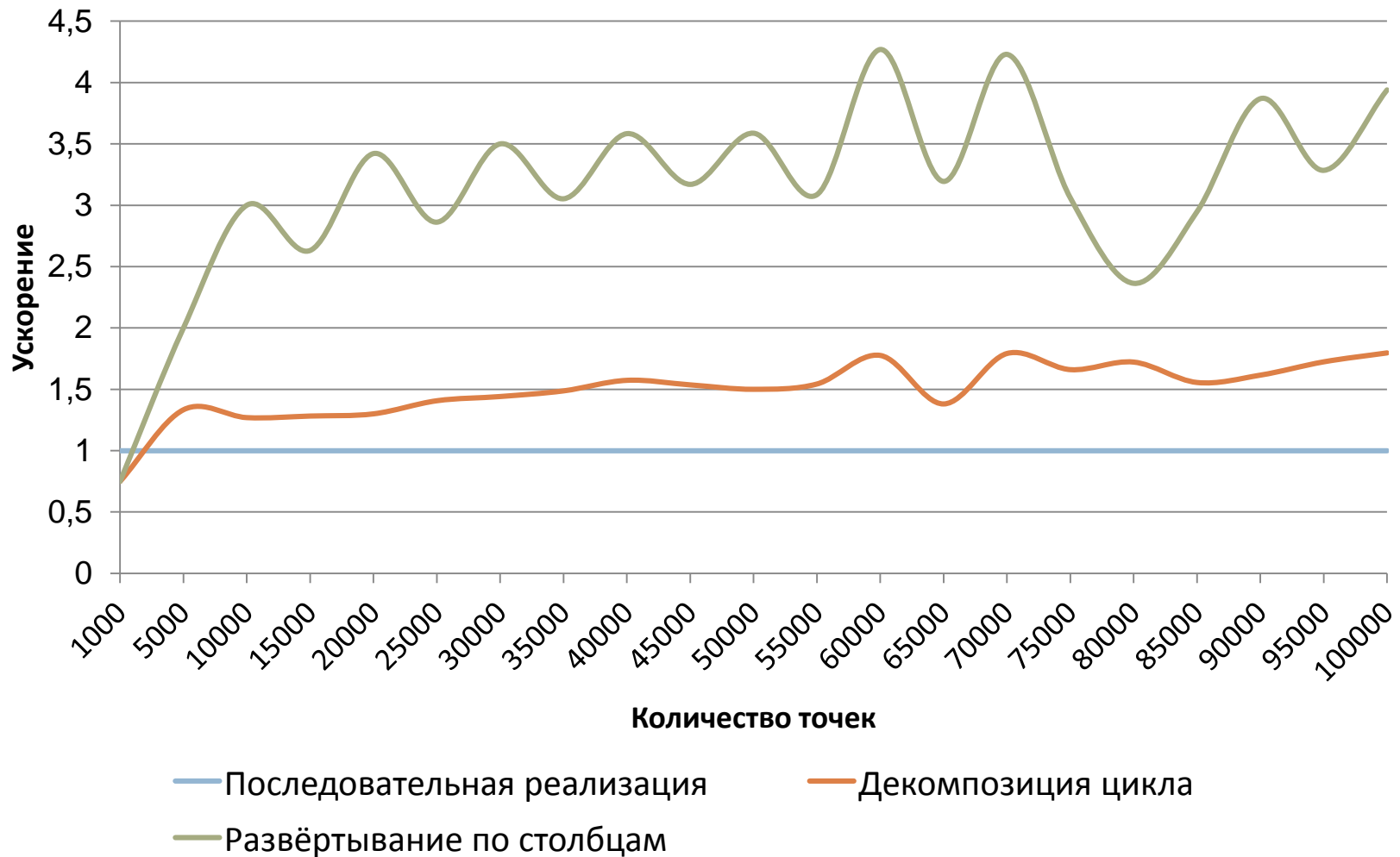
$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} = [\boxed{1} \quad \boxed{2} \quad \boxed{3} \quad \boxed{4} \quad \boxed{5} \quad \boxed{6}]$$

- Развёртывание по столбцам

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} = [\boxed{1} \quad \boxed{4} \quad \boxed{2} \quad \boxed{5} \quad \boxed{3} \quad \boxed{6}]$$

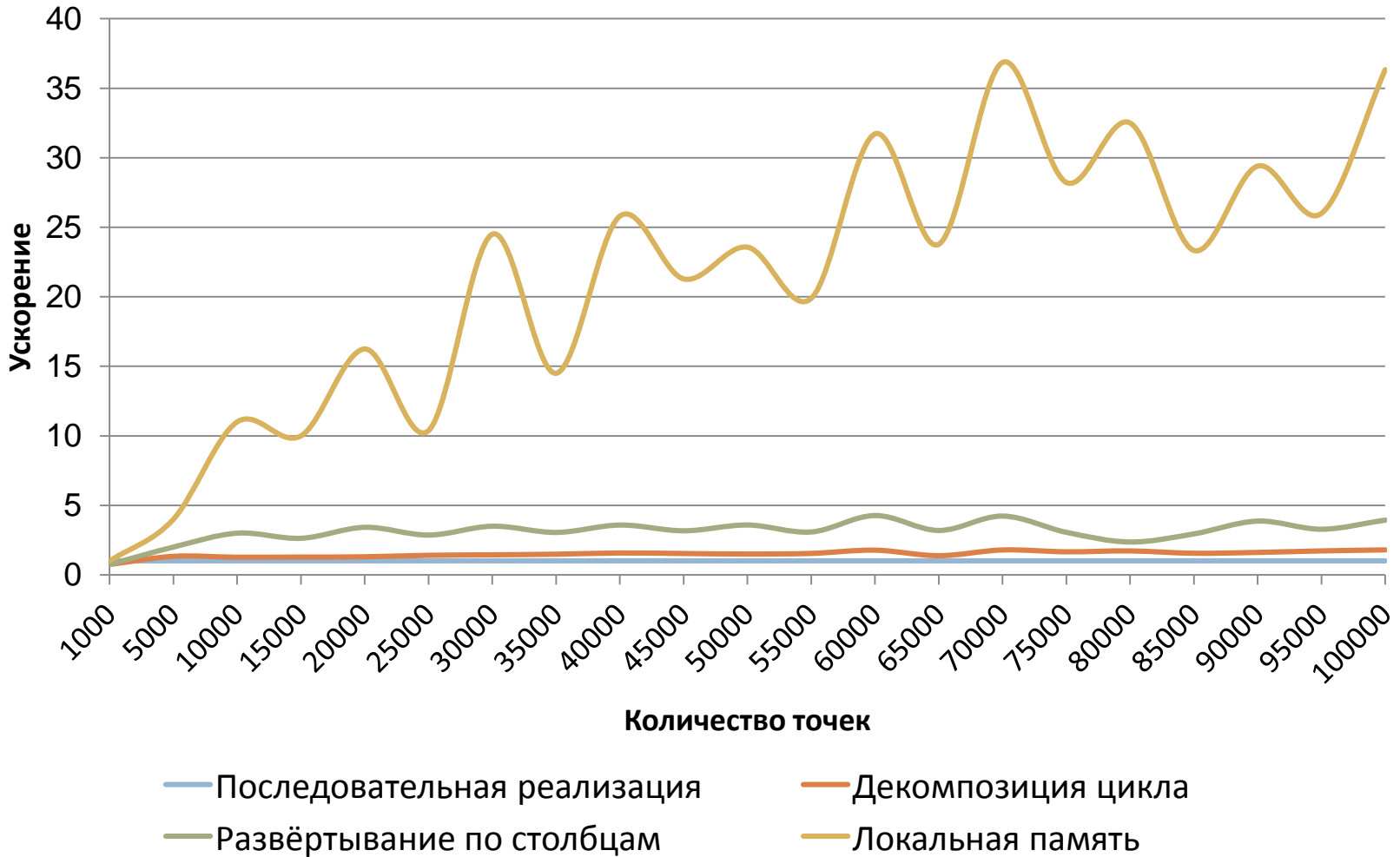
Обращения к памяти

11



Локальная память

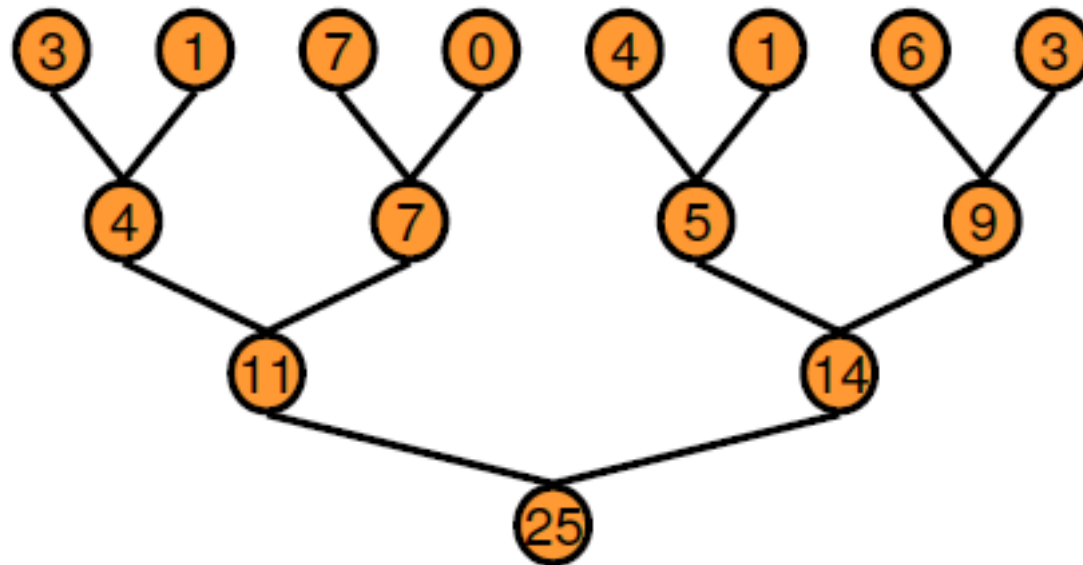
12



Сокращение

13

- Операции $\arg \min(\cdot)$ и $\arg \max(\cdot)$ реализуются с помощью алгоритма параллельного сокращения



Бустинг

14

- Построение сильного классификатора как набор слабых классификаторов



$$G(x) = \sum_{m=1}^M \alpha_m G_m(x)$$

Аддитивная модель

15

- Аддитивное разложение на базисные функции

$$f(x) = \sum_{m=1}^M \beta_m b(x; \gamma_m)$$

- Подбор модели:

$$\min_{\{\beta_m, \gamma_m\}_1^M} \sum_{i=1}^N L \left(y_i, \sum_{m=1}^M \beta_m b(x_i; \gamma_m) \right)$$

- Прямой пошаговый подбор

Градиентный бустинг

16

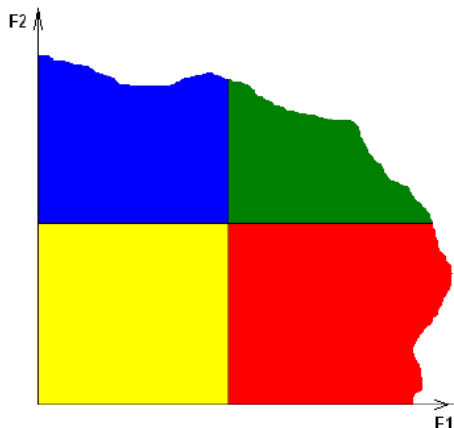
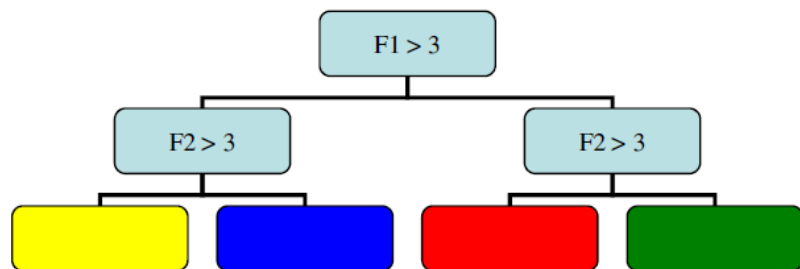
- Вывод слабого классификатора с учётом отрицательного градиента

$$g_{im} = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x_i)=f_{m-1}(x_i)}$$

- Подбор методом наименьших квадратов

Деревья решений

17



$$x \in R_j \Rightarrow f(x) = \gamma_j$$
$$T(x; \Theta) = \sum_{j=1}^J \gamma_j I(x \in R_j)$$
$$\Theta = \{R_j, \gamma_j\}_1^J$$

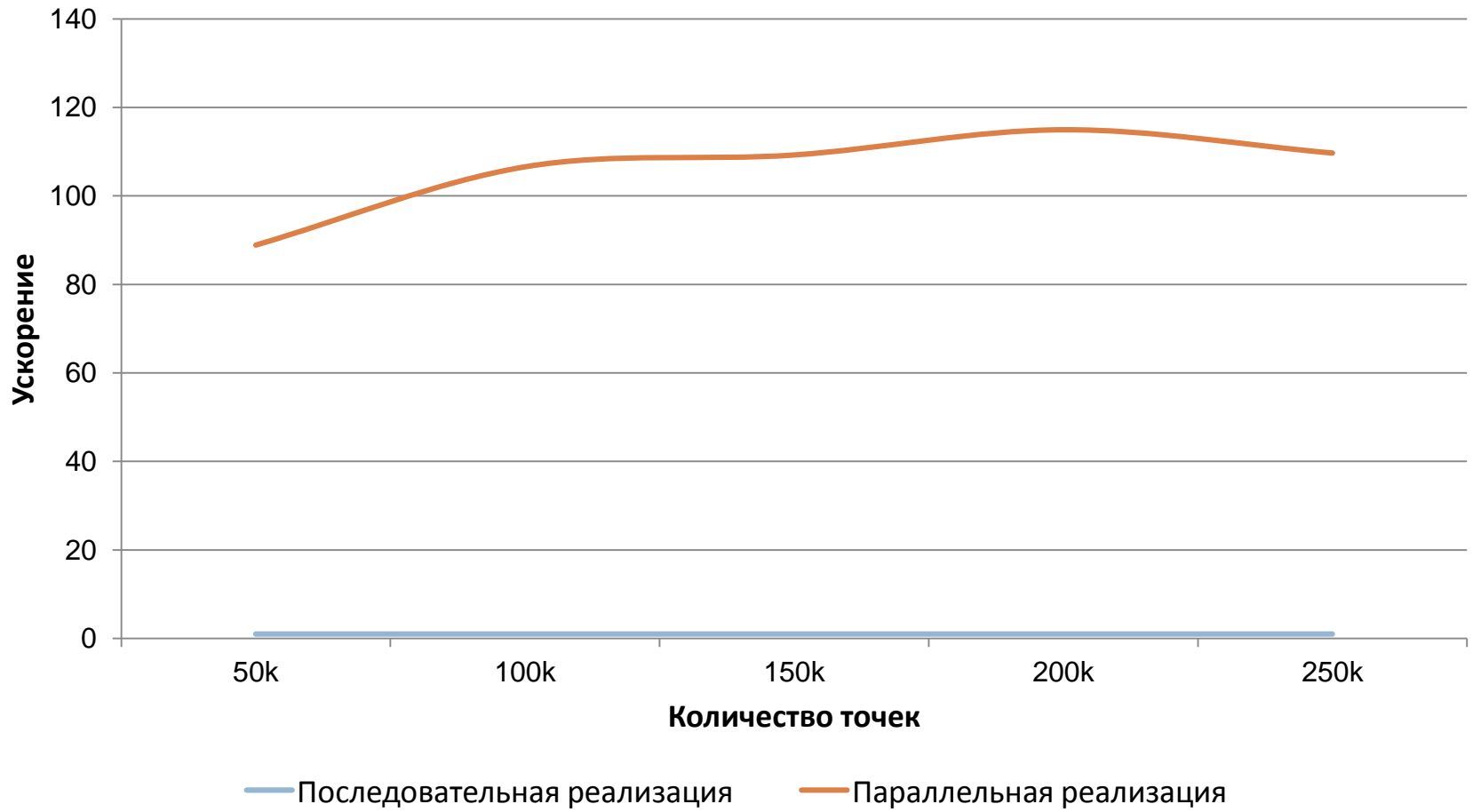
Бустинг на графическом процессоре

18

- Будем строить решающие деревья с 64 листьями
- Будем рассматривать 32 разбиения для каждого фактора
- Независимость разбиений – источник параллелизма

Бустинг на графическом процессоре

19



Результаты

20

- Метод бустинга деревьев решений полностью реализован на графическом процессоре
 - Ускорение в 70 раз
- Для метода опорных векторов реализованы вычислительно-сложные этапы
 - Ускорение в 36 раз на старом оборудовании
- На основе анализа производительности вычислений на графическом процессоре выявлены основные цели оптимизации

Результаты

21

- Полученная реализация бустинга применена на практике для решения задачи перекрёстной проверки
 - Кластер из нескольких десятков машин можно заменить одной машиной, использующей видеокарты

Спасибо за внимание!



Перекрёстная проверка

23

- Оценка обобщающей способности модели
- K-блочная перекрёстная проверка

