



Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Московский государственный технический университет
имени Н. Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н. Э. Баумана)

ФАКУЛЬТЕТ «Информатика и системы управления»

КАФЕДРА «Программное обеспечение ЭВМ и информационные технологии»

ОТЧЕТ

по лабораторной работе № 5

по курсу «Анализ алгоритмов»

на тему: «Организация асинхронного взаимодействия потоков вычисления на
примере конвейерных вычислений»

Студент ИУ7-51Б
(Группа)

(Подпись, дата)

Д. В. Шубенина
(И. О. Фамилия)

Преподаватель

(Подпись, дата)

Л. Л. Волкова
(И. О. Фамилия)

Преподаватель

(Подпись, дата)

Ю. В. Строганов
(И. О. Фамилия)

2023 г.

СОДЕРЖАНИЕ

ВВЕДЕНИЕ	3
1 Аналитическая часть	4
1.1 Конвейерная обработка данных	4
1.2 Нечеткий алгоритм кластеризации с-средних	4
1.3 Описание алгоритмов	5
2 Конструкторская часть	7
2.1 Требования к программному обеспечению	7
2.2 Разработка алгоритмов	7
3 Технологическая часть	14
3.1 Средства реализации	14
3.2 Сведения о модулях программы	14
3.3 Реализация алгоритмов	15
4 Исследовательская часть	21
4.1 Технические характеристики	21
4.2 Демонстрация работы программы	21
4.3 Временные характеристики	22
4.4 Вывод	22
ЗАКЛЮЧЕНИЕ	23
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ	24

ВВЕДЕНИЕ

Целью данной лабораторной работы является описание параллельных конвейерных вычислений.

Для достижения поставленной цели необходимо решить следующие задачи:

- 1) описать организацию конвейерной обработки данных;
- 2) описать алгоритмы обработки данных, которые будут использоваться в текущей лабораторной работе;
- 3) определить средства программной реализации;
- 4) реализовать программу, выполняющую конвейерную обработку с количеством лент не менее трех в однопоточной и многопоточной реализаций;
- 5) сравнить и проанализировать реализации алгоритмов по затраченному времени.

1 Аналитическая часть

В данном разделе приведена информация, касающаяся основ конвейерной обработки данных.

1.1 Конвейерная обработка данных

Конвейер — организация вычислений, при которой увеличивается количество выполняемых инструкций за единицу времени за счет использования принципов параллельности.

Конвейеризация в компьютерной обработке данных основана на разбиении выполнения функций на более мелкие этапы, называемые ступенями, и выделении отдельной аппаратуры для каждой из них. Это позволяет организовать передачу данных от одного этапа к следующему, что увеличивает производительность за счет одновременного выполнения нескольких команд на различных ступенях конвейера.

Хотя конвейеризация не уменьшает время выполнения отдельной команды, она повышает пропускную способность процессора, что приводит к более быстрому выполнению программы по сравнению с простой, не конвейерной схемой.

1.2 Нечеткий алгоритм кластеризации с-средних

Кластерный анализ — это ряд математических методов интеллектуального анализа данных, предназначенных для разбиения множества исследуемых объектов на компактные группы, называемые кластерами. Под объектами кластерного анализа подразумеваются предметы исследования, нуждающиеся в кластеризации по некоторым признакам. Признаки объектов могут иметь как непрерывные, так и дискретные значения [1].

Метод с-средних — итеративный нечеткий алгоритм кластеризации. В данном методе кластеры являются нечеткими множествами, и каждый объект из выборки исходных данных относится одновременно ко всем кластерам с различной степенью принадлежности. Таким образом, матрица принадлежности объектов к кластерам (или матрица разбиения) содержит не бинарные, а вещественные значения, принадлежащие отрезку $[0; 1]$ [1].

Пусть X — исходный набор данных размера N . Обновление матрицы

принадлежности и списка центров кластеров производится в 5 этапов:

- 1) инициализация матрицы центров кластеров W случайными значениями;
- 2) инициализация матрицы разбиения $U = (\mu_{ij})$ следующим образом:

$$\mu_{ij}^{(t)} = \frac{1}{C \sum_{l=1}^C \left(\frac{d_{ij}}{d_{il}} \right)^{\frac{2}{m-1}}}, i = 1, \dots, N; j, l = 1, \dots, C$$

$$\mu_{ij} = \begin{cases} 1, & \text{если } d_{ij} = 0 \\ 0, & \text{для } l \neq j, \end{cases}$$
(1.1)

где t — номер итерации,

C — количество кластеров,

m — показатель нечеткости, регулирующий точность разбиения,

d_{ij} — расстояние от x_i до w_j , $d_{ij} = \|x_i - w_j^{(t)}\|$;

- 3) увеличить t на 1 и рассчитать матрицу $W^{(t)}$ по формуле (1.2)

$$W_j^{(t)} = \frac{\sum_{i=1}^N \left(\mu_{ij}^{(t-1)} \right)^m x_i}{\sum_{i=1}^N \left(\mu_{ij}^{(t-1)} \right)^m}, j = 1, \dots, C;$$
(1.2)

- 4) вычислить матрицу разбиения $U^{(t)}$ согласно соотношению (1.1);

- 5) если $\|U^{(t)} - U^{(t-1)}\| \geq \varepsilon$ перейти на шаг 3 [2].

1.3 Описание алгоритмов

В качестве операций, выполняющихся на конвейере, взяты следующие:

- 1) обработка файла, описывающего датасет файлов MIDI [3];
- 2) извлечение биграмм нот;
- 3) кластеризация нечетким алгоритмом с-средних.

Ленты конвейера (обработчики) будут передавать друг другу заявки. Первый этап, или обработчик, будет формировать заявку, которая будет передаваться от этапа к этапу. Заявка содержит:

- 1) массив нот, извлеченных из файла MIDI;
- 2) массив биграмм нот;
- 3) результат кластеризации: матрица разбиения и список центроидов кластеров.

Вывод

В данном разделе было рассмотрено понятие конвейерной обработки, а также выбраны этапы для обработки датасета файлов MIDI. Также был рассмотрен и описан алгоритм кластеризации с-средних.

2 Конструкторская часть

В данном разделе будут представлены схемы последовательной и параллельной работы стадий конвейера.

2.1 Требования к программному обеспечению

К программному обеспечению предъявлен ряд требований:

- 1) наличие интерфейса для выбора действий;
- 2) возможность обработки файлов MIDI;
- 3) возможность выбора линейной или конвейерной реализаций алгоритма.

2.2 Разработка алгоритмов

На рисунке 2.1 представлена схема последовательного нечеткого алгоритма кластеризации с-средних.

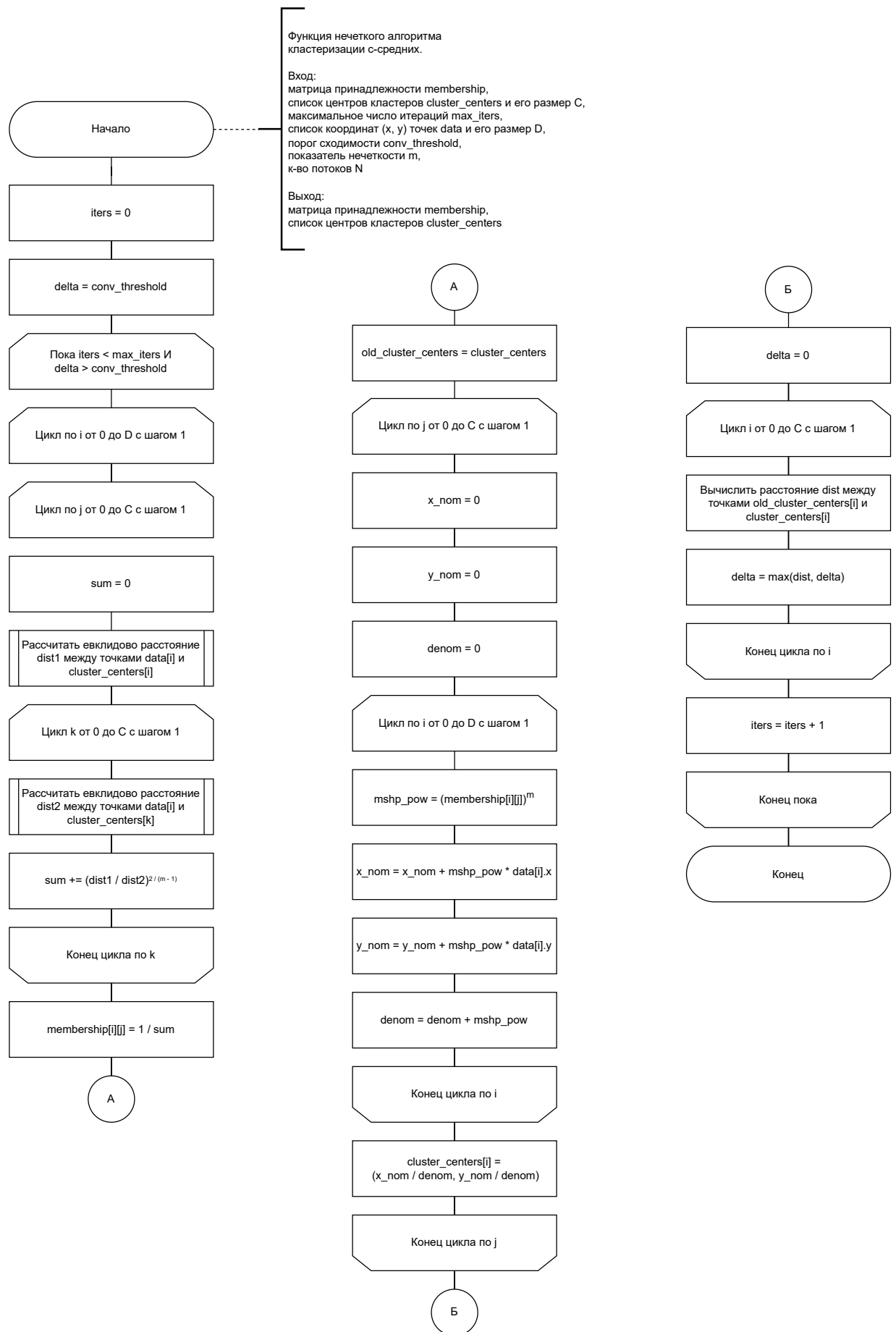


Рисунок 2.1 – Схема последовательного нечеткого алгоритма с-средних

На рисунке 2.2 представлена схема алгоритма линейной обработки датасетов файлов MIDI.

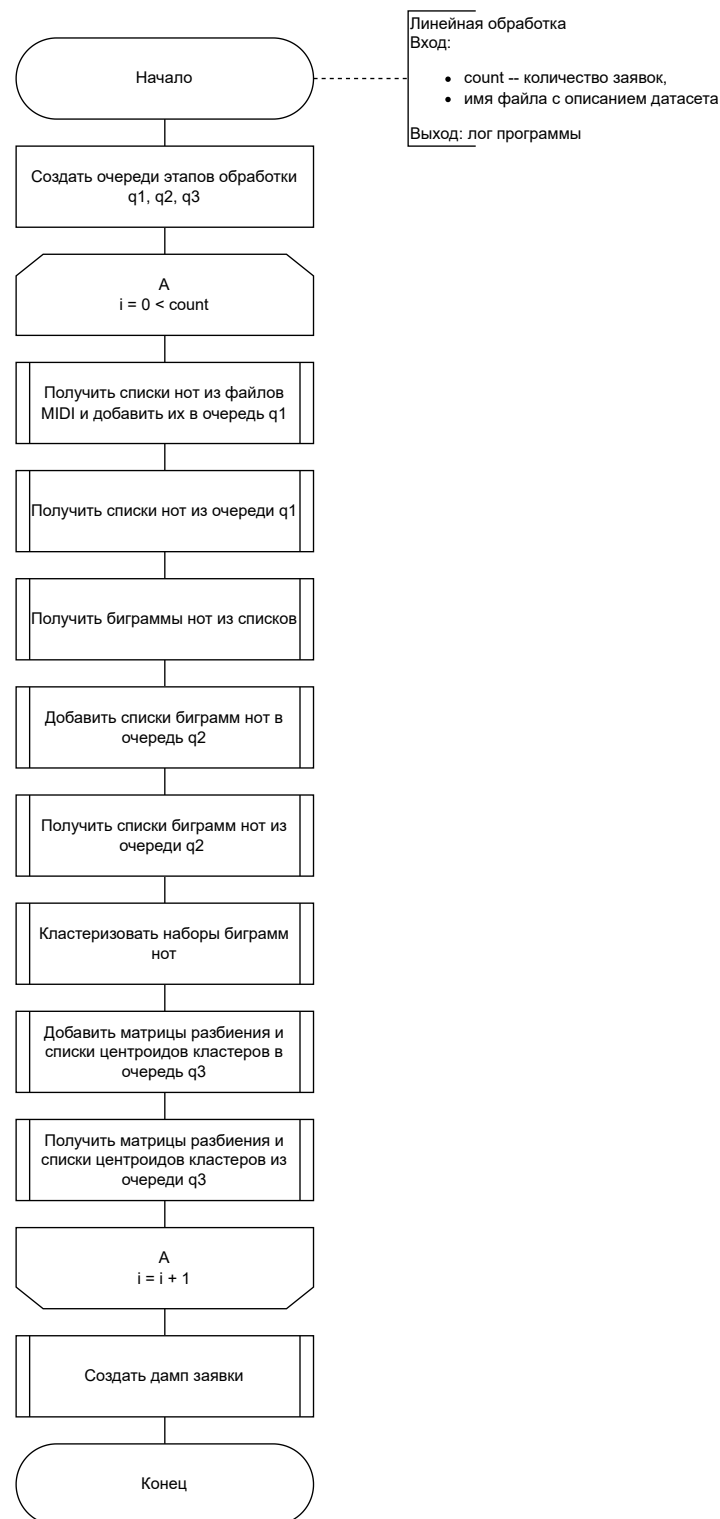


Рисунок 2.2 – Схема алгоритма линейной обработки датасетов файлов MIDI

На рисунке 2.3 представлена схема алгоритма главного потока конвейерной обработки датасетов файлов MIDI.

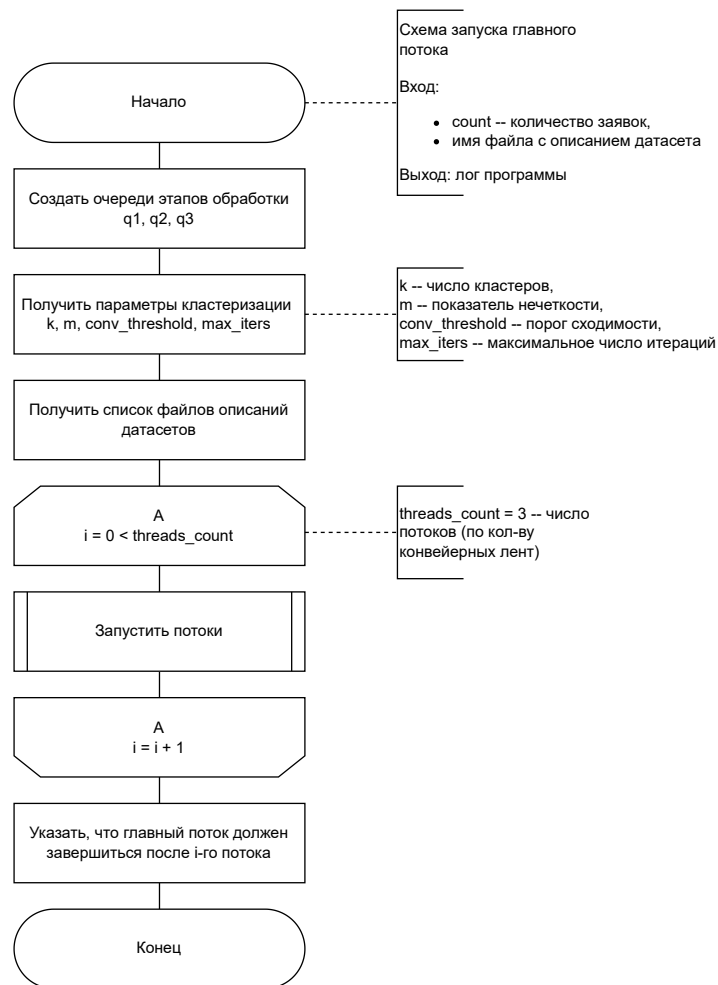


Рисунок 2.3 – Схема алгоритма главного потока конвейерной обработки датасетов файлов MIDI

На рисунке 2.4 представлена схема алгоритма потока, выполняющего извлечение нот из MIDI файлов.

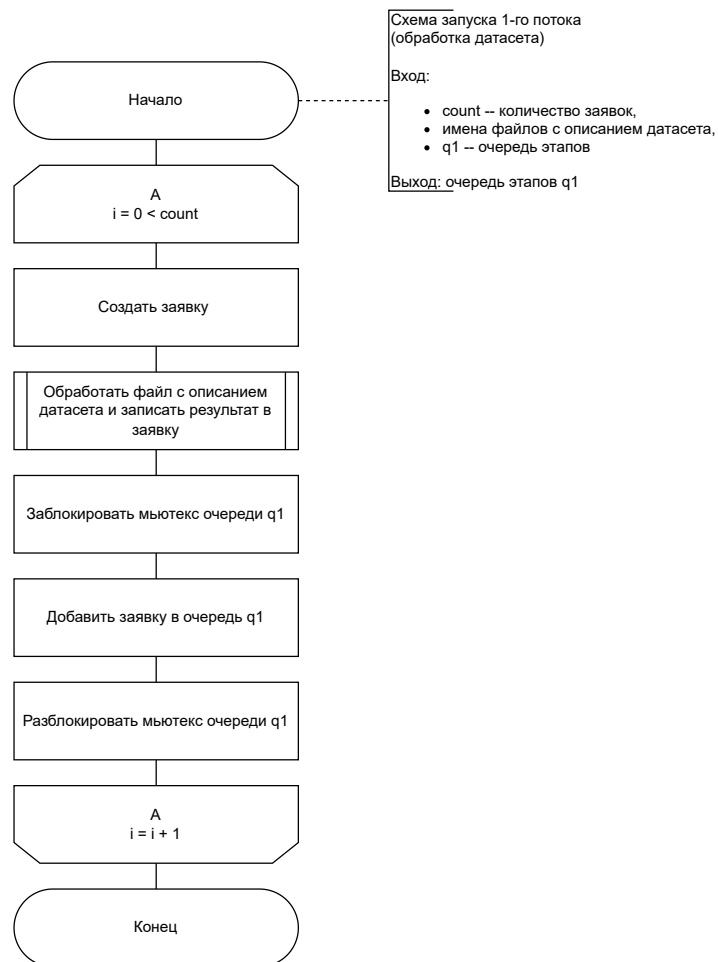


Рисунок 2.4 – Схема алгоритма потока, выполняющего извлечение нот из MIDI файлов

На рисунке 2.5 представлена схема алгоритма потока, выполняющего извлечение биграмм нот из обработанных на предыдущем этапе файлов.

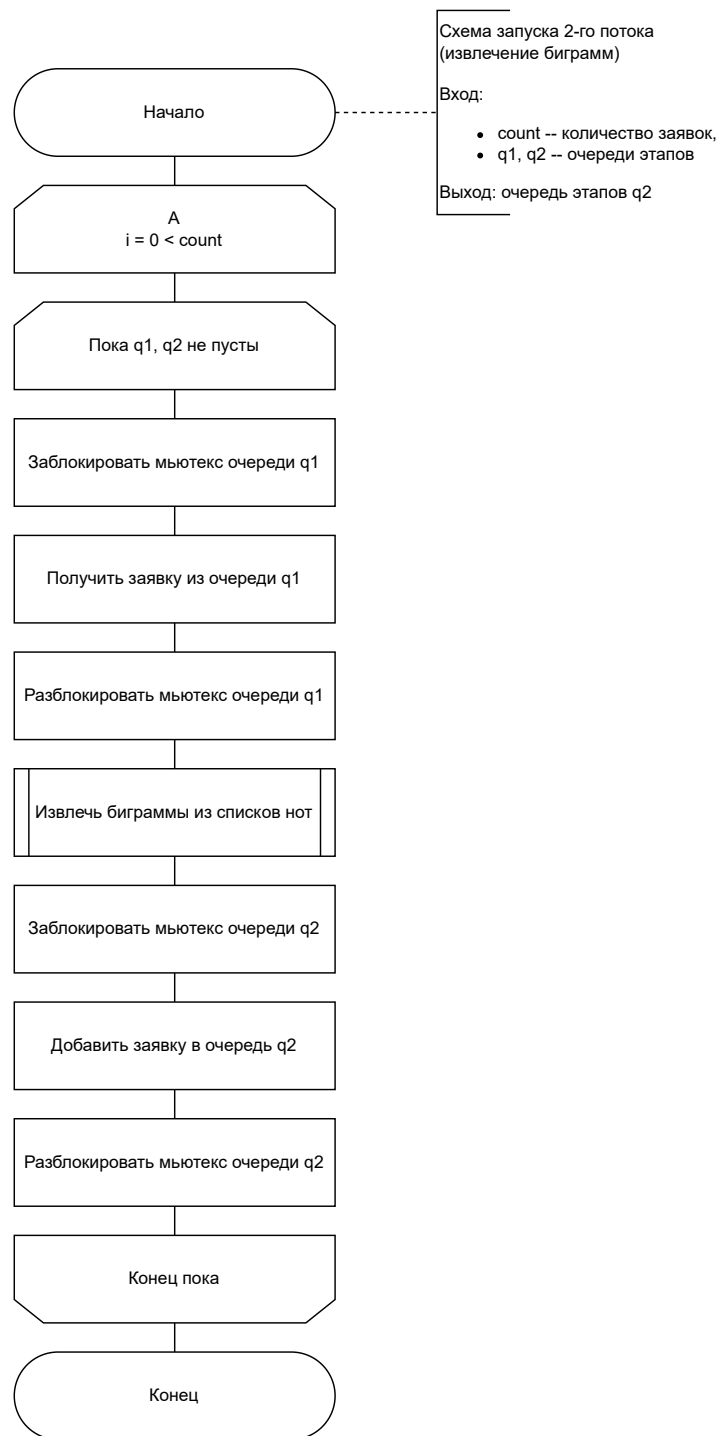


Рисунок 2.5 – Схема алгоритма потока, выполняющего извлечение биграмм нот из обработанных на предыдущем этапе файлов

На рисунке 2.6 представлена схема алгоритма потока, выполняющего кластеризацию биграмм нот.

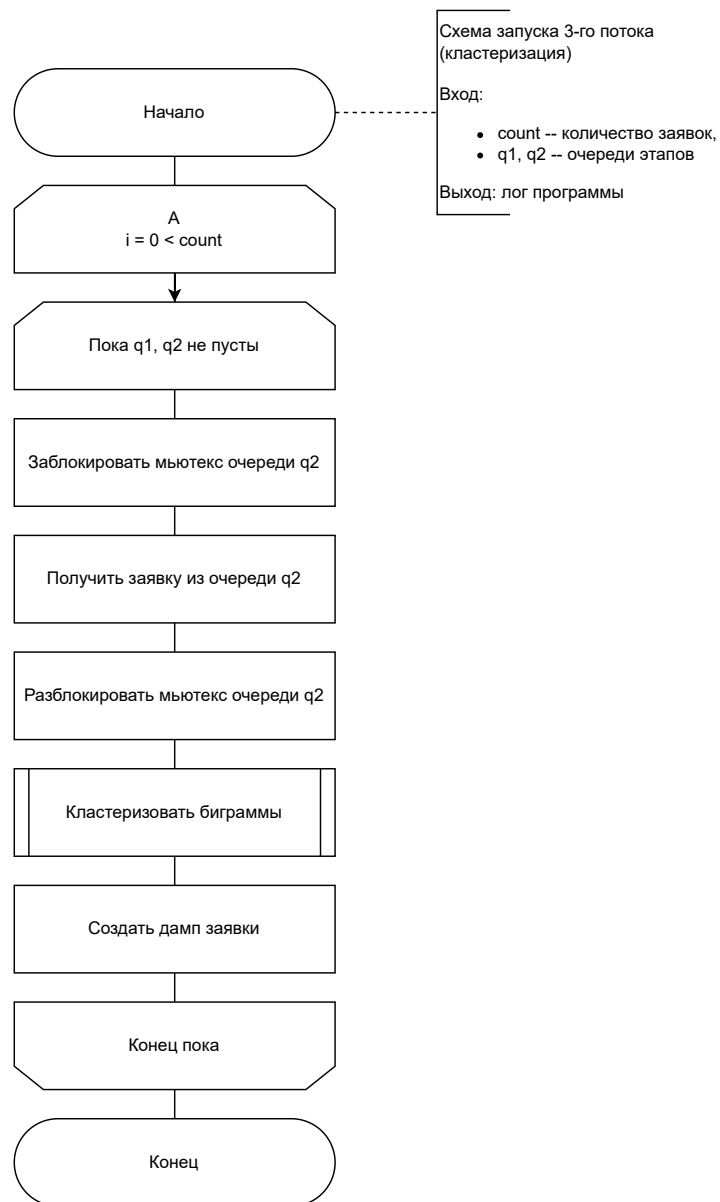


Рисунок 2.6 – Схема алгоритма потока, выполняющего кластеризацию биграмм нот

Вывод

В данном разделе были перечислены требования к программному обеспечению и построены схемы рассматриваемых алгоритмов.

3 Технологическая часть

В данном разделе описаны средства реализации программного обеспечения, а также листинги и функциональные тесты.

3.1 Средства реализации

В качестве языка программирования, используемого при написании данной лабораторной работы, был выбран C++ [4], так как в нем имеется контейнер `std::vector`, представляющий собой динамический массив данных произвольного типа, и библиотека `<ctime>` [5], позволяющая производить замеры процессорного времени. Также выбранный язык программирования предоставляет возможность работы с:

- потоками (класс `thread` [6]);
- мьютексами (класс `mutex` [7]);

Обработка файлов MIDI производилась с помощью библиотеки `MidiFile` [8].

3.2 Сведения о модулях программы

Данная программа разбита на следующие модули:

- `main.cpp` — файл, содержащий точку входа в программу;
- `algorithms.cpp` — файл, содержащий реализации алгоритмов, используемых на различных стадиях конвейера;
- `ts_queue.cpp` — файл, содержащий реализацию потокобезопасной очереди;
- `pipeline.cpp` — файл, содержащий функции конвейерной обработки;
- `utils.cpp` — файл, содержащий вспомогательные функции;
- `measure.cpp` — файл, содержащий функции, измеряющие процессорное время выполнения реализуемых алгоритмов.

3.3 Реализация алгоритмов

На листинге 3.1 представлены реализации алгоритмов извлечения нот из файлов MIDI, извлечения биграмм и кластеризации с-средних.

Листинг 3.1 – Реализации алгоритмов извлечения нот из файлов MIDI, извлечения биграмм и кластеризации с-средних

```
1  std::vector<note_vector_t> read_dataset(const std::string
    &descr_file)
2  {
3      auto [id, n, filenames] = read_descr_file(descr_file);
4      std::vector<note_vector_t> notes;
5      for (const auto &filename : filenames)
6      {
7          smf::MidiFile midifile;
8          midifile.read(filename);
9
10         note_vector_t temp_notes;
11         for (int track = 0; track < midifile.getTrackCount();
            track++)
12         {
13             for (int event = 0; event < midifile[track].size();
                event++)
14             {
15                 if (midifile[track][event].isNoteOn())
16                 {
17                     temp_notes.push_back(
18                         midifile[track][event][1]);
19                 }
20             }
21         }
22         notes.emplace_back(std::move(temp_notes));
23     }
24     return notes;
25 }
26
27 std::vector<point_vec_t> extract_bigrams(const
    std::vector<note_vector_t> &notes)
28 {
29     std::vector<point_vec_t> bigrams;
30     for (const auto &note_set : notes)
31     {
```

```

32     point_vec_t temp_bigrams;
33     for (size_t i = 1; i < note_set.size(); ++i)
34     {
35         std::vector<double> bigram{
36             (double)note_set[i - 1], (double)note_set[i]};
37         temp_bigrams.emplace_back(bigram);
38     }
39     bigrams.emplace_back(temp_bigrams);
40 }
41 return bigrams;
42 }
43
44 void c_means(
45     membership_t &membership, point_vec_t &cluster_centers,
46     const point_vec_t &data,
47     double m, double conv_threshold, int max_iters)
48 {
49     int iters = 0;
50     double delta = conv_threshold + 1.0;
51     while (iters < max_iters && delta > conv_threshold)
52     {
53         for (size_t i = 0; i < data.size(); ++i)
54         {
55             for (size_t j = 0; j < cluster_centers.size(); ++j)
56             {
57                 double sum = 0.0;
58                 double dist1 = sqrt(pow(data[i][0] -
59                                         cluster_centers[j][0], 2) +
60                                     pow(data[i][1] -
61                                         cluster_centers[j][1],
62                                         2));
63                 for (size_t k = 0; k < cluster_centers.size();
64                     ++k)
65                 {
66                     double dist2 = sqrt(pow(data[i][0] -
67                                             cluster_centers[k][0], 2) +
68                                           pow(data[i][1] -
69                                             cluster_centers[k][1],
70                                             2));
69                     sum += pow(dist1 / dist2, 2.0 / (m - 1.0));
70                 }

```



```

66         membership[i][j] = 1.0 / sum;
67     }
68 }
69 auto old_cluster_centers = cluster_centers;
70 for (size_t j = 0; j < cluster_centers.size(); ++j)
71 {
72     double x_nom = 0.0, y_nom = 0.0, denom = 0.0;
73     for (size_t i = 0; i < data.size(); ++i)
74     {
75         double membership_pow_m = pow(membership[i][j],
76             m);
77         x_nom += membership_pow_m * data[i][0];
78         y_nom += membership_pow_m * data[i][1];
79         denom += membership_pow_m;
80     }
81     cluster_centers[j] = {x_nom / denom, y_nom / denom};
82 }
83 delta = 0.0;
84 for (size_t i = 0; i < cluster_centers.size(); ++i)
85 {
86     double distance = sqrt(pow(old_cluster_centers[i][0]
87         - cluster_centers[i][0], 2) +
88         pow(old_cluster_centers[i][1]
89         - cluster_centers[i][1],
90         2));
91     if (distance > delta)
92         delta = distance;
93 }
94 ++iters;
95 }
96 }

```

На листинге 3.2 представлена реализация линейного алгоритма обработки набора файлов MIDI.

Листинг 3.2 – Реализация линейного алгоритма обработки набора файлов MIDI

```

1 void consequent()
2 {
3     int req_cnt = utils::get_request_count();
4     auto datasets = utils::pick_datasets(req_cnt);
5     auto [k, m, conv_threshold, max_iters] =

```

```

        utils::get_clust_params();
6
7    std::vector<std::unique_ptr<stages_t>> pool;
8    for (int i = 0; i < req_cnt; ++i)
9    {
10        stages_t *s = new stages_t;
11
12        clock_gettime(CLOCK_REALTIME, &s->parsed.op_start);
13        s->parsed.notes = read_dataset(datasets[i]);
14        clock_gettime(CLOCK_REALTIME, &s->parsed.op_end);
15
16        clock_gettime(CLOCK_REALTIME, &s->embedded.op_start);
17        s->embedded.embeddings =
            extract_bigrams(s->parsed.notes);
18        clock_gettime(CLOCK_REALTIME, &s->embedded.op_end);
19
20        clock_gettime(CLOCK_REALTIME, &s->clusterized.op_start);
21        for (const auto &embed : s->embedded.embeddings)
22        {
23            auto [mshp, cc] = init_structures(embed.size(), k);
24            c_means(mshp, cc, embed, m, conv_threshold,
                max_iters);
25            s->clusterized.results.emplace_back(std::move(mshp),
                std::move(cc));
26        }
27        clock_gettime(CLOCK_REALTIME, &s->clusterized.op_end);
28        pool.emplace_back(s);
29    }
30    dump_pool(pool, "cons.txt");
31 }

```

На листинге 3.3 представлена реализация конвейерного алгоритма обработки набора файлов MIDI.

Листинг 3.3 – Реализация конвейерного алгоритма обработки набора файлов MIDI

```

1 static void service_01(
2     int req_cnt,
3     const std::vector<std::string> &datasets,
4     ts_queue<stages_t *> &q1)
5 {
6     for (int i = 0; i < req_cnt; ++i)

```

```

7      {
8          stages_t *s = new stages_t();
9          clock_gettime(CLOCK_REALTIME, &s->parsed.op_start);
10         s->parsed.notes = read_dataset(datasets[i]);
11         clock_gettime(CLOCK_REALTIME, &s->parsed.op_end);
12         q1.push(s);
13     }
14 }
15
16 static void service_02(
17     int req_cnt,
18     ts_queue<stages_t *> &q1,
19     ts_queue<stages_t *> &q2)
20 {
21     for (int i = 0; i < req_cnt; ++i)
22     {
23         stages_t *s = q1.pop();
24
25         clock_gettime(CLOCK_REALTIME, &s->embedded.op_start);
26         s->embedded.embeddings =
27             extract_bigrams(s->parsed.notes);
28         clock_gettime(CLOCK_REALTIME, &s->embedded.op_end);
29         q2.push(s);
30     }
31 }
32
33 static void service_03(
34     int req_cnt,
35     std::tuple<int, double, double, int> cls_params,
36     ts_queue<stages_t *> &q2,
37     std::vector<std::unique_ptr<stages_t>> &pool)
38 {
39     for (int i = 0; i < req_cnt; ++i)
40     {
41         stages_t *s = q2.pop();
42         auto [k, m, conv_threshold, max_iters] = cls_params;
43
44         clock_gettime(CLOCK_REALTIME, &s->clusterized.op_start);
45         for (const auto &embed : s->embedded.embeddings)
46         {
47             auto [mshp, cc] = init_structures(embed.size(), k);

```

```

47         c_means(mshp, cc, embed, m, conv_threshold,
48                 max_iters);
49         s->clusterized.results.emplace_back(std::move(mshp),
50                                             std::move(cc));
51     }
52     clock_gettime(CLOCK_REALTIME, &s->clusterized.op_end);
53     pool.emplace_back(s);
54 }
55 void concurrent()
56 {
57     int req_cnt = utils::get_request_count();
58     auto datasets = utils::pick_datasets(req_cnt);
59     std::tuple cls_params = utils::get_clust_params();
60
61     std::vector<std::unique_ptr<stages_t>> pool;
62     ts_queue<stages_t *> q1;
63     ts_queue<stages_t *> q2;
64
65     std::thread t_01(service_01, req_cnt, std::cref(datasets),
66                     std::ref(q1));
67     std::thread t_02(service_02, req_cnt, std::ref(q1),
68                     std::ref(q2));
69     std::thread t_03(service_03, req_cnt, cls_params,
70                     std::ref(q2), std::ref(pool));
71
72     t_01.join();
73     t_02.join();
74     t_03.join();
75
76     dump_pool(pool, "conc.txt");
77 }

```

Вывод

В данном разделе были рассмотрены средства реализации, а также представлен листинг реализаций линейного и конвейерного алгоритмов обработки набора файлов MIDI.

4 Исследовательская часть

В данном разделе приведены технические характеристики устройства, на котором проводилось измерение времени работы программного обеспечения, а также результаты замеров времени.

4.1 Технические характеристики

Технические характеристики устройства, на котором выполнялись замеры по времени:

- процессор: AMD Ryzen 7 5800X @ 3.800 ГГц, 8 физ. ядер, 16 лог. ядер;
- оперативная память: 32 ГБайт.
- операционная система: Manjaro Linux x86_64 (версия ядра Linux 6.5.12-1-MANJARO).

Измерения проводились на стационарном компьютере. Во время проведения измерений устройство было нагружено только системными приложениями.

4.2 Демонстрация работы программы

На рисунке 4.1 продемонстрирована работа программы для случая, когда пользователь выбрал пункт 2 «Запустить последовательную обработку датасета файлов MIDI», выбрал файл с описанием датасета под номером 1 и ввел следующие значения:

- количество заявок — 10,
- число кластеров — 3,
- значение показателя нечеткости — 2,
- значение порога сходимости — 1,
- максимальное к-во итераций — 10.

```

        Меню
1. Запустить последовательную обработку датасета файлов MIDI.
2. Запустить конвейерную обработку датасета файлов MIDI.
3. Произвести замеры по времени реализуемых алгоритмов.
0. Выход.

Выберите опцию (0-3): 2

Введите количество заявок: 10
Выберите файлы описания датасета:
    1. "/home/daria/Документы/bmstu-aa/lab5/prog/descriptions/ds_02.txt"
    2. "/home/daria/Документы/bmstu-aa/lab5/prog/descriptions/ds_01.txt"
Номера файлов (0 для выхода): 2 1 2 2 1 2 2 2 1 2
Введите количество кластеров: 3
Введите показатель нечеткости: 2
Введите порог сходимости: 1
Введите максимальное количество итераций: 10

        Меню
1. Запустить последовательную обработку датасета файлов MIDI.
2. Запустить конвейерную обработку датасета файлов MIDI.
3. Произвести замеры по времени реализуемых алгоритмов.
0. Выход.

Выберите опцию (0-3): 0

```

Рисунок 4.1 – Демонстрация работы программы

4.3 Временные характеристики

Исследование временных характеристик реализуемых алгоритмов производилось 2 раза:

1)

Наборы данных генерировались из равномерного распределения.

4.4 Вывод

В результате исследования реализуемых алгоритмов по времени выполнения можно сделать следующие выводы:

1)

ЗАКЛЮЧЕНИЕ

В результате выполнения лабораторной работы по исследованию алгоритмов сортировок решены следующие задачи:

- 1) описан нечеткий алгоритм кластеризации с-средних;
- 2) разработана параллельная версия алгоритма;
- 3) определены средства программной реализации;
- 4) реализованы последовательная и параллельная версии алгоритма;
- 5) проведен сравнительный анализ процессорного времени выполнения реализованных алгоритмов:
 - при варьировании размера датасета последовательная версия нечеткого алгоритма кластеризации с-средних выполнялась в среднем в 92.6 раз дольше, чем параллельная;
 - при варьировании числа потоков последовательная версия нечеткого алгоритма кластеризации с-средних выполнялась в среднем 72.7 раз дольше, чем параллельная;
 - наименьшее время работы многопоточной реализации алгоритма достигается при 8 вспомогательных потоках; наибольшее же время выполнения алгоритма достигается при 4 потоках;

таким образом, рекомендуется использование 8 вспомогательных потоков, т. к. при таком количестве временные затраты на создание потоков, переключение аппаратного контекста и синхронизацию ниже, чем получаемая скорость обработки набора данных.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Лосев Д. Г. Разработка и сравнение параллельных реализаций итеративных алгоритмов кластеризации // Наука молодых - будущее России : сборник научных статей 6-й Международной научной конференции перспективных разработок молодых ученых (9-10 декабря 2021 года), в 5 томах. Т. 4. — Курск : Юго-Зап. гос. ун-т, 2021. — С. 71—74.
2. Hung M.-C., Yang D.-L. An efficient Fuzzy C-Means clustering algorithm // Proceedings 2001 IEEE International Conference on Data Mining. — 2001. — С. 225—232. — DOI: 10.1109/ICDM.2001.989523.
3. MIDI — Wikipedia, The Free Encyclopedia. — [Электронный ресурс]. — Режим доступа: <https://en.wikipedia.org/w/index.php?title=MIDI&oldid=1190344258> (дата обращения: 21.12.2023).
4. C++ language. — [Электронный ресурс]. — Режим доступа: <https://en.cppreference.com/w/cpp/language> (дата обращения: 21.12.2023).
5. Standard library header <ctime>. — [Электронный ресурс]. — Режим доступа: <https://en.cppreference.com/w/cpp/thread/thread> (дата обращения: 21.12.2023).
6. std::thread. — [Электронный ресурс]. — Режим доступа: <https://en.cppreference.com/w/cpp/thread/thread> (дата обращения: 22.12.2023).
7. std::mutex. — [Электронный ресурс]. — Режим доступа: <https://en.cppreference.com/w/cpp/thread/mutex> (дата обращения: 22.12.2023).
8. Midifile: C++ MIDI file parsing library. — [Электронный ресурс]. — Режим доступа: <https://github.com/craigsapp/midifile> (дата обращения: 22.12.2023).