

Теория параллелизма

Отчёт

Уравнение теплопроводности

Выполнил Грищенко Александр Михайлович, 21932

03.2023

1 Цели работы

Реализовать решение уравнение теплопроводности (пятиточечный шаблон) в двумерной области на равномерных сетках.

Перенести программу на GPU используя директивы OpenACC.

Произвести профилирование программы и оптимизацию кода.

Сравнить время работы на CPU и GPU.

2 Используемый компилятор

pgc++

3 Используемый профилировщик

nsys (NVIDIA Nsight Systems) с OpenACC trace.

4 Как проводился замер времени работы

Для замера времени работы использовалась библиотека chrono.

Замер времени производился несколько раз, затем бралось среднее время.

5 Выполнение на CPU

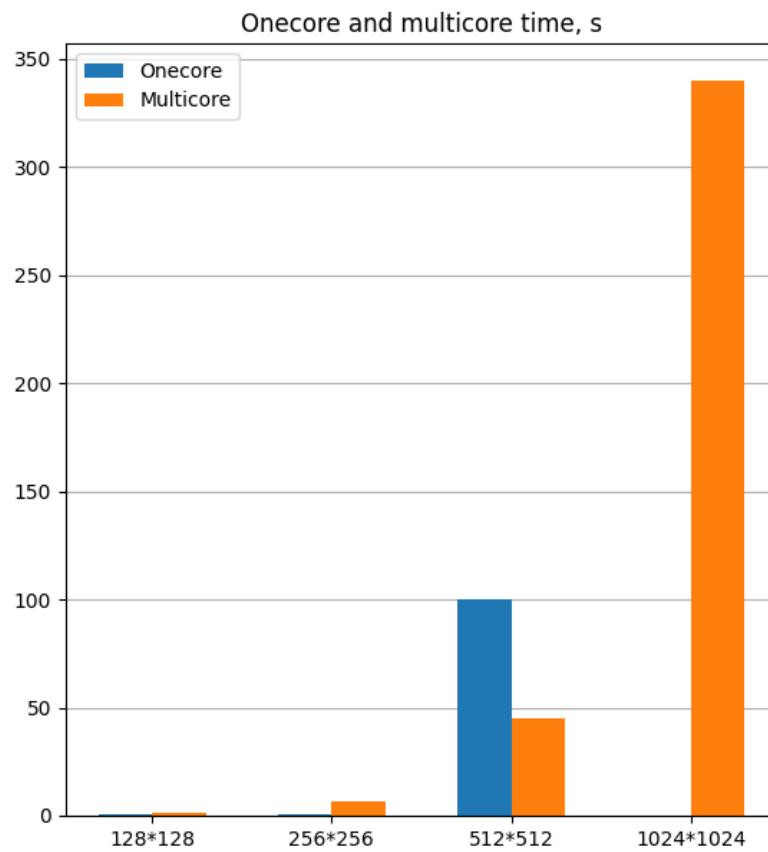
5.1 CPU-onecore

Размер сетки	Время выполнения, с	Точность	Количество операций
128*128	0.5	9.52e-07	11136
256*256	0.67	9.85e-07	37376
512*512	100	9.77e-07	120832

5.2 CPU-multicore

Размер сетки	Время выполнения, с	Точность	Количество операций
128*128	1.1	9.52e-07	11136
256*256	6.5	9.85e-07	37376
512*512	45	9.78e-07	120832
1024*1024	340	9.89e-07	365568

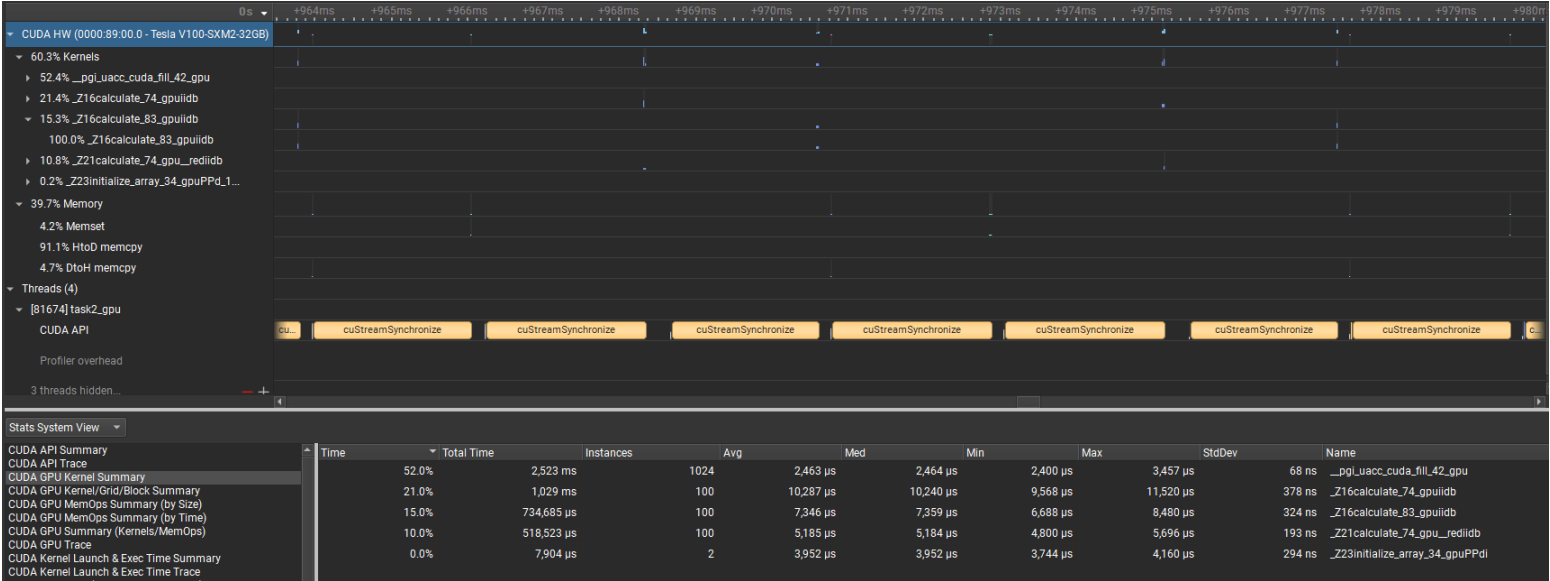
5.3 Диаграмма сравнения время работы CPU-onecore и CPU-multicore



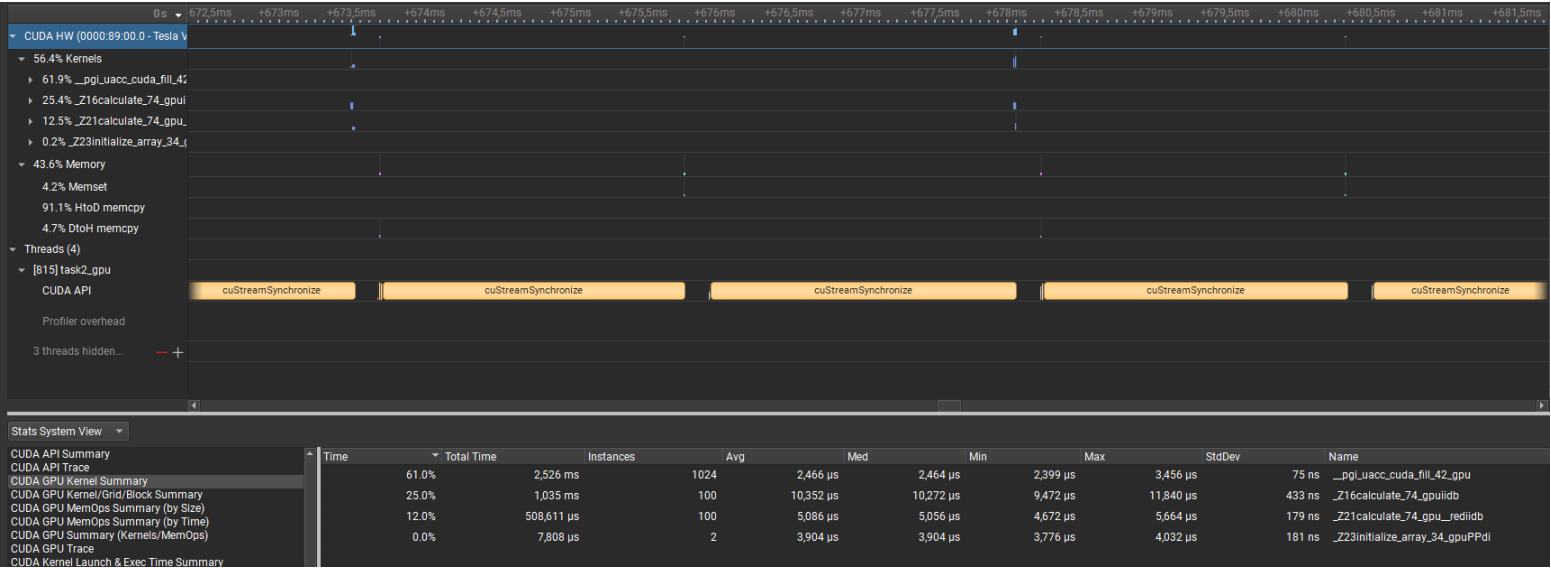
6 Выполнение на GPU

6.1 Этапы оптимизации на сетке 512*512
 (количество итераций при профилировании 100)

Этап №	Время вы- полнения, с	Точность	Количество операций	Комментарии (что было сделано)
1	0.95	0.107	100	Распараллелены циклы, reduction(max:error)
2	0.73	0.107	100	Замена копирования массива swar'ом через указатели
3	0.7	0.035	100	Изначальная инициализация массива значениями 20
4	0.5	0.036	100	Асинхронный запуск ядер
5	0.25	N/A	100	Подсчет ошибки не каждую итера- цию



Этап 1



Этап 2

Почему инициализация значениями 20? Для начала посмотрим на финальный вывод (для простоты взята сетка 8x8):

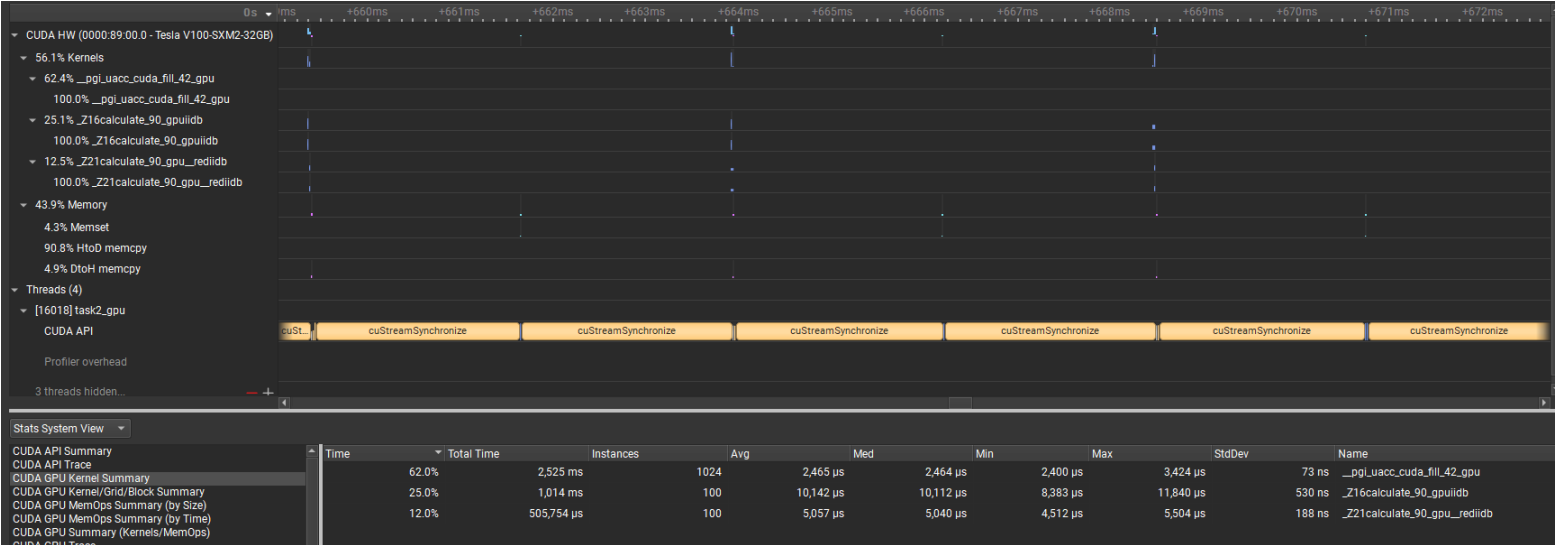
10	11.43	12.86	14.29	15.71	17.14	18.57	20
11.43	12.86	14.29	15.71	17.14	18.57	20	21.43
12.86	14.29	15.71	17.14	18.57	20	21.43	22.86
14.29	15.71	17.14	18.57	20	21.43	22.86	24.29
15.71	17.14	18.57	20	21.43	22.86	24.29	25.71
17.14	18.57	20	21.43	22.86	24.29	25.71	27.14
18.57	20	21.43	22.86	24.29	25.71	27.14	28.57
20	21.43	22.86	24.29	25.71	27.14	28.57	30

Как можно заметить значения лежат в отрезке [10, 30], значит, имеет смысл изначально инициализировать массив значениями из этого диапазона. Также стоит отметить, что модой в матрице является число 20 (а ещё и средним арифметическим границ отрезка). Значит изначальная инициализация массива 20-ю существенно приблизит ответ к правильному, ещё до запуска основного алгоритма.

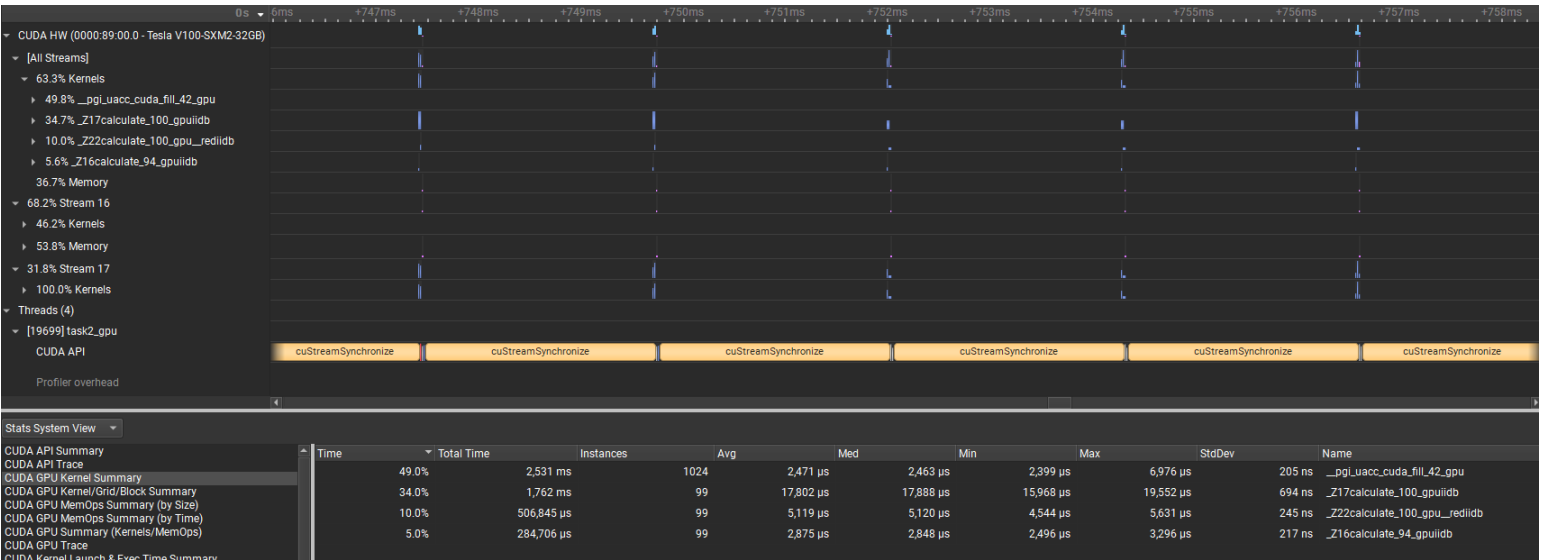
Сравнение поведения программы на сетке 128x128.

Инициализация	Время выполнения, с	Точность	Количество операций
Нет	1.1	1e-06	11073
Да	0.4	9.995e-07	30074

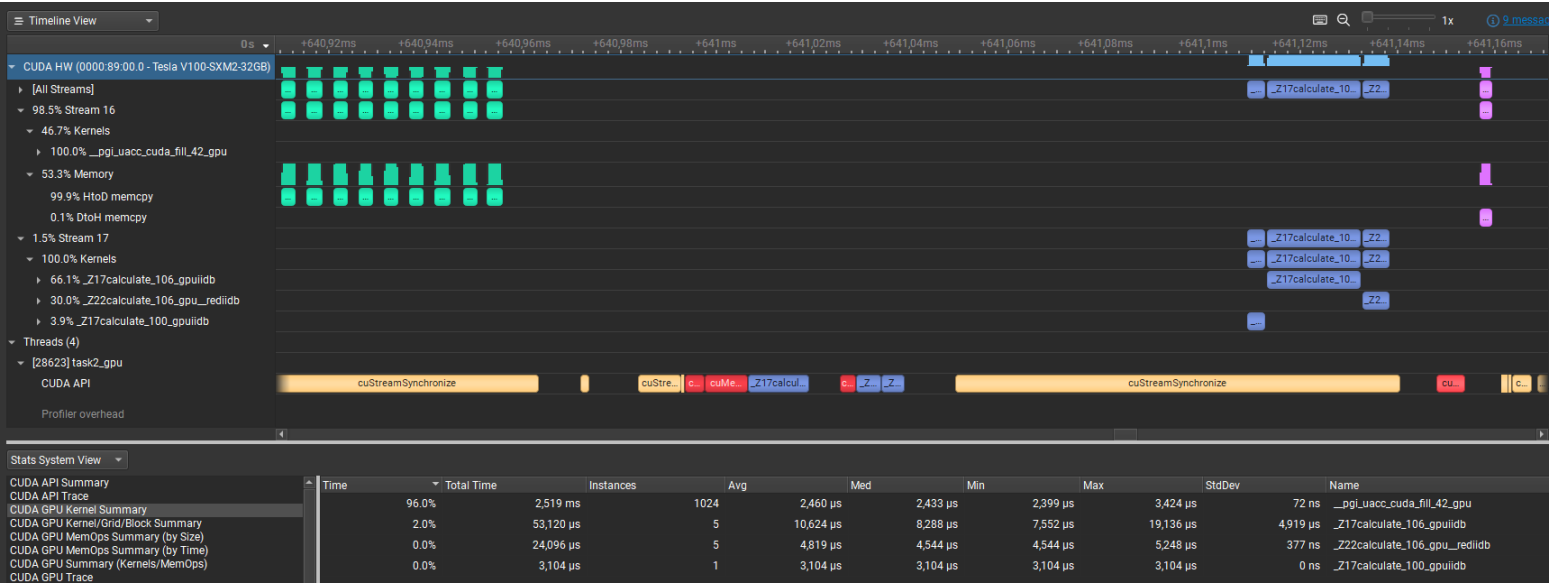
Как видно в таблице, изначальная инициализация дает хороший прирост производительности.



Этап 3

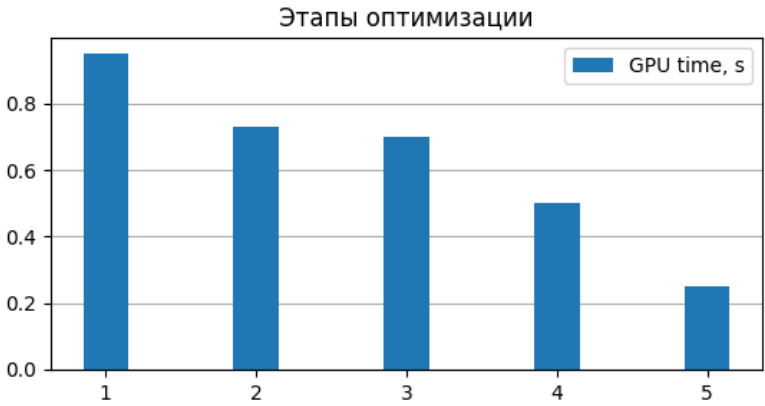


Этап 4



Этап 5

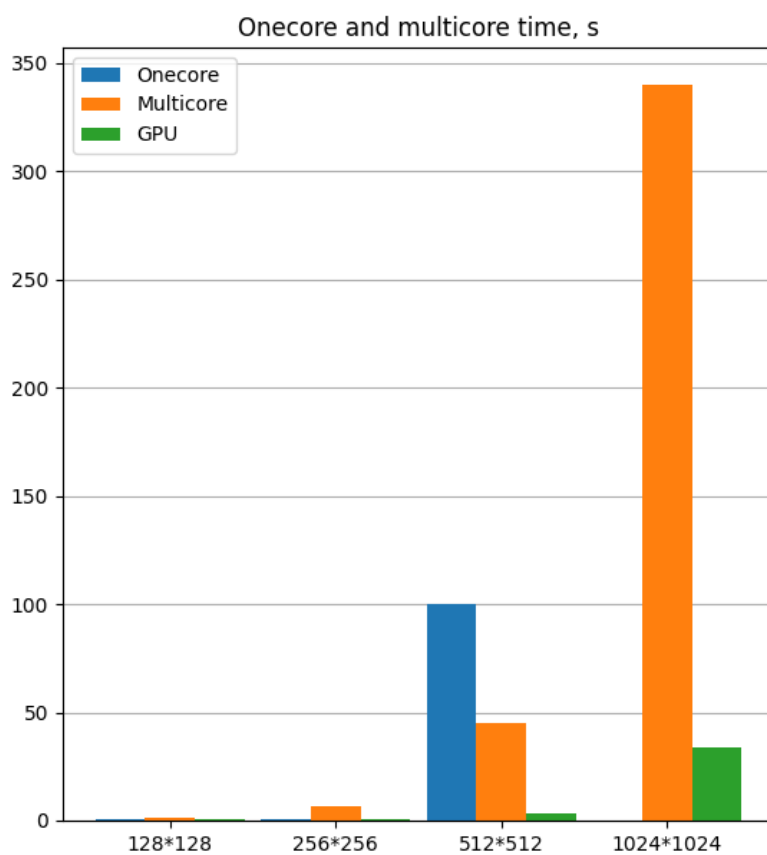
6.2 Диаграмма оптимизации (по горизонтали номер эта- па; по вертикали время работы)



6.3 GPU – оптимизированный вариант

Размер сетки	Время выполнения, с	Точность	Количество опреаций
128*128	0.5	9.52e-07	11136
256*256	0.8	9.85e-07	37376
512*512	3.1	9.78e-07	120832
1024*1024	34	9.89e-07	365568

7 Диаграмма сравнения времени работы CPU-one, CPU-multi, GPU (оптимизированный вариант) для разных размеров сеток



8 Вывод

При грамотном использовании ресурсов GPU можно достичь ускорения работы программы в десятки раз.

9 Приложение

9.1 Ссылка на GitHub

https://github.com/busyhedg03/ParallelismTheory/tree/master/task_2

```

1  #include <iostream>
2  #include <cmath>
3  #include <chrono>
4
5  #ifdef _FLOAT
6  #define T float
7  #define MAX std::fmaxf
8  #else
9  #define T double
10 #define MAX std::fmax
11 #endif
12
13 // Вывести значения двумерного массива
14 void print_array(T **A, int size)
15 {
16 #pragma acc update host(A[:size][:size])
17     std::cout.precision(4);
18     for (int i = 0; i < size; i += 1)
19     {
20         for (int j = 0; j < size; j += 1)
21             std::cout << A[i][j] << "\t";
22         std::cout << std::endl;
23     }
24     std::cout << std::endl;
25 }
26
27 // Инициализация матрицы, чтобы подготовить ее к основному алгоритму
28 void initialize_array(T **A, int size)
29 {
30     // Заполнение углов матрицы значениями
31     A[0][0] = 10.0;
32     A[0][size - 1] = 20.0;
33     A[size - 1][size - 1] = 30.0;
34     A[size - 1][0] = 20.0;
35
36     // Заполнение периметра матрицы
37     T step = 10.0 / (size - 1);
38     for (int i = 1; i < size - 1; i++)
39     {
40         T addend = step * i;
41         A[0][i] = A[0][0] + addend; // horizontal
42         A[size - 1][i] = A[size - 1][0] + addend; // horizontal
43         A[i][0] = A[0][0] + addend; // vertical
44         A[i][size - 1] = A[0][size - 1] + addend; // vertical
45     }
46
47     // Заполнение 2D-ю, чтобы сократить количество операций в основном алгоритме
48     for (int i = 1; i < size - 1; i++)
49         for (int j = 1; j < size - 1; j++)
50             A[i][j] = 20.0;
51 }
52
53 // Очистка динамической памяти, выделенной под двумерную матрицу
54 void delete_2d_array(T **A, int size)
55 {
56     for (int i = 0; i < size; i++)
57         delete[] A[i];
58     delete[] A;
59 }

```

```

61 // Основной алгоритм
62 void calculate(int net_size = 128, int iter_max = 1e6, T accuracy = 1e-6, bool res = false)
63 {
64     // Создание 2-х двумерных матриц, одна будет считаться на основе другой
65     T **Anew = new T *[net_size],
66       **A = new T *[net_size];
67     for (int i = 0; i < net_size; i++)
68     {
69         A[i] = new T[net_size];
70         Anew[i] = new T[net_size];
71     }
72
73     // Инициализация матриц
74     initialize_array(A, net_size);
75     initialize_array(Anew, net_size);
76
77     // Больше 30-и ошибки быть не должно
78     T error = 30;
79     // Счетчик итераций
80     int iter = 0;
81     // Указатель для swap
82     T **temp;
83     // Флаг обновления ошибки на хосте для обработки условием цикла
84     bool update_flag = true;
85 #pragma acc data copyin(A[:net_size][:net_size], Anew[:net_size][:net_size]) copy(error)
86     {
87         while (error > accuracy && iter++ < iter_max)
88         {
89             // Сокращение количества обращений к CPU. Чем больше сетка, тем реже стоит сверять значения.
90             update_flag = !(iter % net_size);
91             // напомнить о массивах
92             #pragma acc data present(A, Anew)
93             // асинхронно
94             #pragma acc kernels async(1)
95             {
96                 // вернуть ошибку на девайс и занулить
97                 if (update_flag)
98                 {
99                     #pragma acc update device(error)
100                     error = 0;
101                 }
102                 #pragma acc loop independent collapse(2) reduction(max \
103                     : error)
104                 for (int j = 1; j < net_size - 1; j++)
105                 {
106                     for (int i = 1; i < net_size - 1; i++)
107                     {
108                         Anew[i][j] = (A[i + 1][j] + A[i - 1][j] + A[i][j - 1] + A[i][j + 1]) * 0.25;
109                         if (update_flag)
110                             error = std::max(error, std::abs(Anew[i][j] - A[i][j]));
111                     }
112                 }
113             }
114
115             // swap(A, Anew)
116             temp = A, A = Anew, Anew = temp;
117
118             // синхронизировать и отправить хосту для проверки условия
119             if (update_flag)
120             {
121                 #pragma acc wait(1)
122                 #pragma acc update self(error)
123             }
124             if (res)
125                 print_array(A, net_size);
126             std::cout << "iter=" << iter << ",\terror=" << error << std::endl;
127         }
128     }
129
130     delete_2d_array(A, net_size);
131     delete_2d_array(Anew, net_size);
132 }

```

```

133
134 int main(int argc, char *argv[])
135 {
136     auto begin_main = std::chrono::steady_clock::now();
137     int net_size = 128, iter_max = (int)1e6;
138     T accuracy = 1e-6;
139     bool res = false;
140     for (int arg = 1; arg < argc; arg++)
141     {
142         std::string str = argv[arg];
143         if (!str.compare("-a"))
144         {
145             #ifdef _FLOAT
146                 accuracy = std::stof(argv[arg + 1]);
147             #else
148                 accuracy = std::stod(argv[arg + 1]);
149             #endif
150             arg++;
151         }
152         else if (!str.compare("-i"))
153         {
154             iter_max = (int)std::stod(argv[arg + 1]); // 1e6
155             arg++;
156         }
157         else if (!str.compare("-s"))
158         {
159             net_size = std::stoi(argv[arg + 1]);
160             arg++;
161         }
162         else if (!str.compare("-res"))
163         {
164             res = true;
165         }
166     }
167     calculate(net_size, iter_max, accuracy, res);
168     auto end_main = std::chrono::steady_clock::now();
169     int time_spent_main = std::chrono::duration_cast<std::chrono::milliseconds>(end_main - begin_main).count();
170     std::cout << "The elapsed time is:\nmain\t\t\t" << time_spent_main << " ms\n";
171     return 0;
172 }
173

```