

# Теория Параллелизма

## Отчет

### Уравнение теплопроводности (cublas + openacc)

Выполнил 21932, Бабенко Егор Степанович

18.03.2022

## Цель работы

Реализовать решение уравнение теплопроводности методом Якоба (пятиточечный шаблон) для двумерной сетки. Произвести профилирование программы для GPU и оптимизацию кода. Произвести сравнение времени работы на CPU и GPU. Для CPU использовалась программа из прошлого задания.

Используемый компилятор: *pgc++*

Для компиляции версии с cublas использовалось:

- `pgc++ t3.cpp -o t3_GPU.pg -fast -acc=gpu -O2 -D OPENACC__ -Mcudalib=cublas`

Для дополнительной профилировки: `-D NVPROF_`

Используемый профилировщики: *nvprof*, *nsys*

Nsys использовался с OpenACC trace, cuBLAS trace, а также NVTX trace.

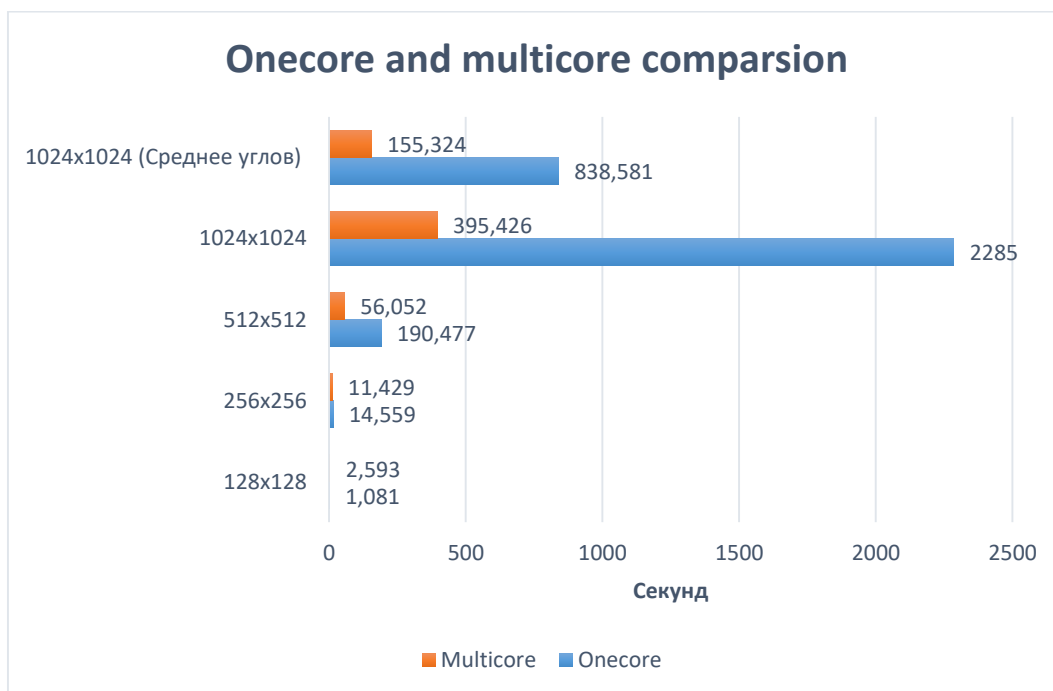
Замер времени работы всей программы производился с помощью команды `time`. Время высчитывалось как среднее между пятью запусками.

## CPU-onecore

Размер сетки	Время выполнения (сек)	Ошибка	Количество итераций
128x128	1.081	9.9998e-07	30074
256x256	14.559	9.9993e-07	102885
512x512	190.477	9.99984e-07	339599
1024x1024	2285.354	1.36929e-06	1000000
1024x1024 <i>(среднее значение углов сетки как стартовое значение)</i>	838.581	9.99993e-07	364620

## CPU-multicore

Размер сетки	Время выполнения (сек)	Ошибка	Количество итераций
128x128	2.593	9.9998e-07	30074
256x256	11.429	9.9993e-07	102885
512x512	56.052	9.99984e-07	339599
1024x1024	395.426	1.36929e-06	1000000
1024x1024 <i>(среднее значение углов сетки как стартовое значение)</i>	155.324	9.99993e-07	364620



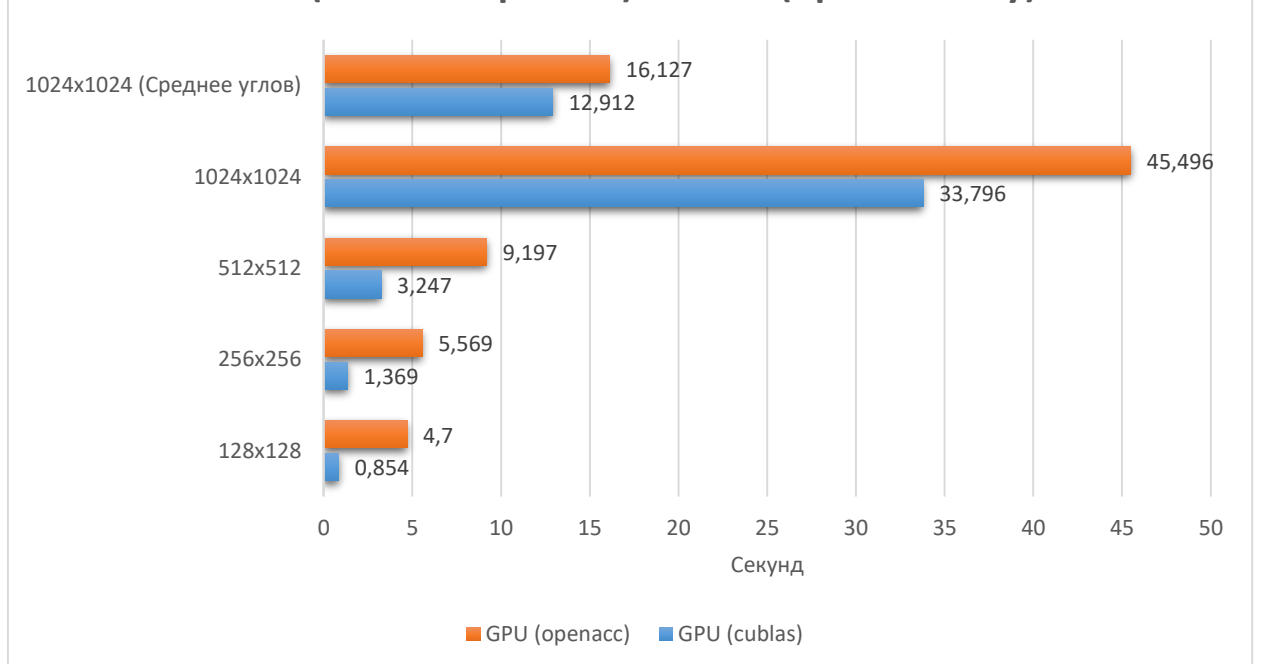
## GPU-optimized (openacc only)

Размер сетки	Время выполнения (сек)	Ошибка	Количество итераций
128x128	4.700	9.98757e-07	30078
256x256	5.569	9.99703e-07	102888
512x512	9.197	9.99965e-07	339600
1024x1024	45.496	1.36929e-06	1000000
1024x1024 <i>(среднее значение углов сетки как стартовое значение)</i>	16.127	9.99993e-07	364620

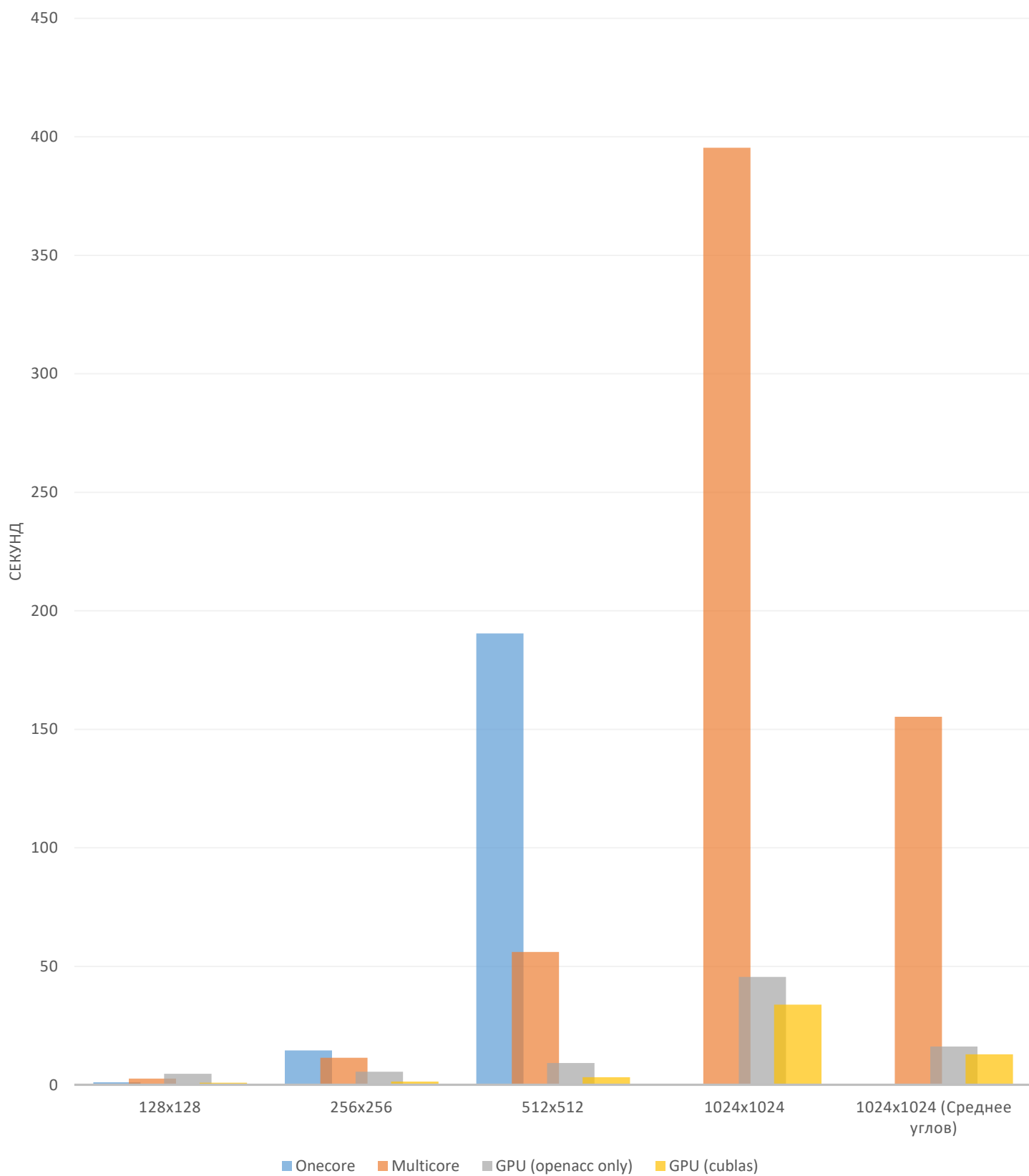
## GPU-optimized (cublas)

Размер сетки	Время выполнения (сек)	Ошибка	Количество итераций
128x128	0.854	9.98757e-07	30120
256x256	1.369	9.99703e-07	102910
512x512	3.247	9.99965e-07	339603
1024x1024	33.796	1.36929e-06	1000000
1024x1024 <i>(среднее значение углов сетки как стартовое значение)</i>	12.912	9.99993e-07	364703

### GPU (cublas + openacc) vs GPU (openacc only)



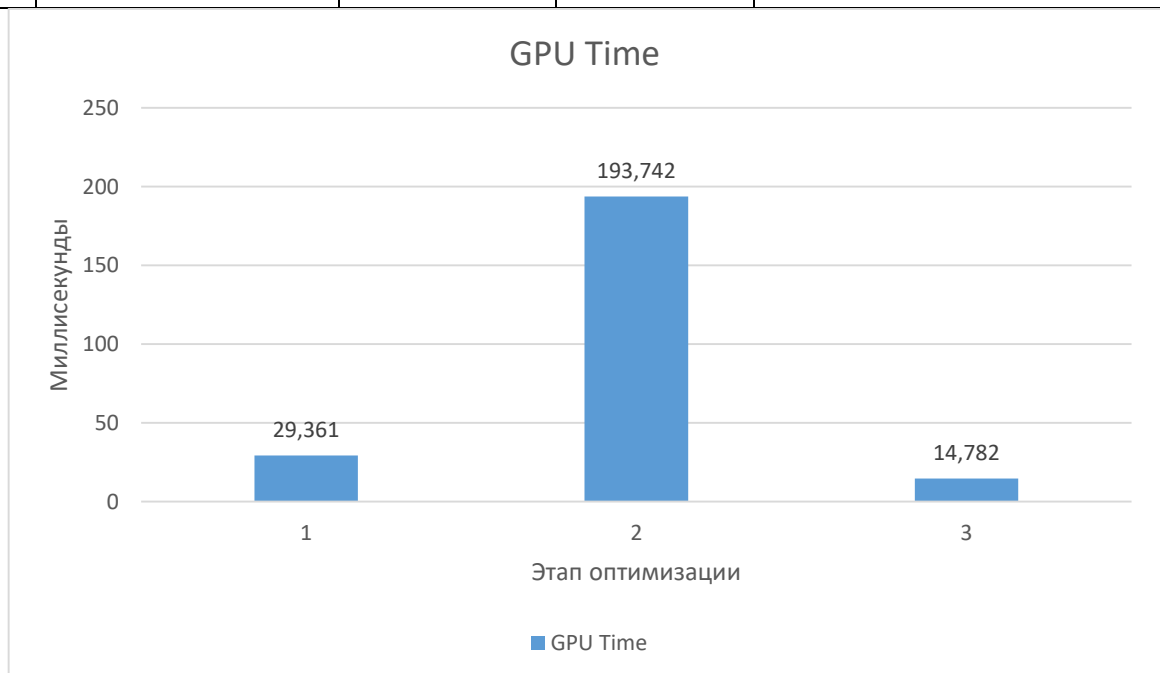
Onecore, multicore and GPU comparsion



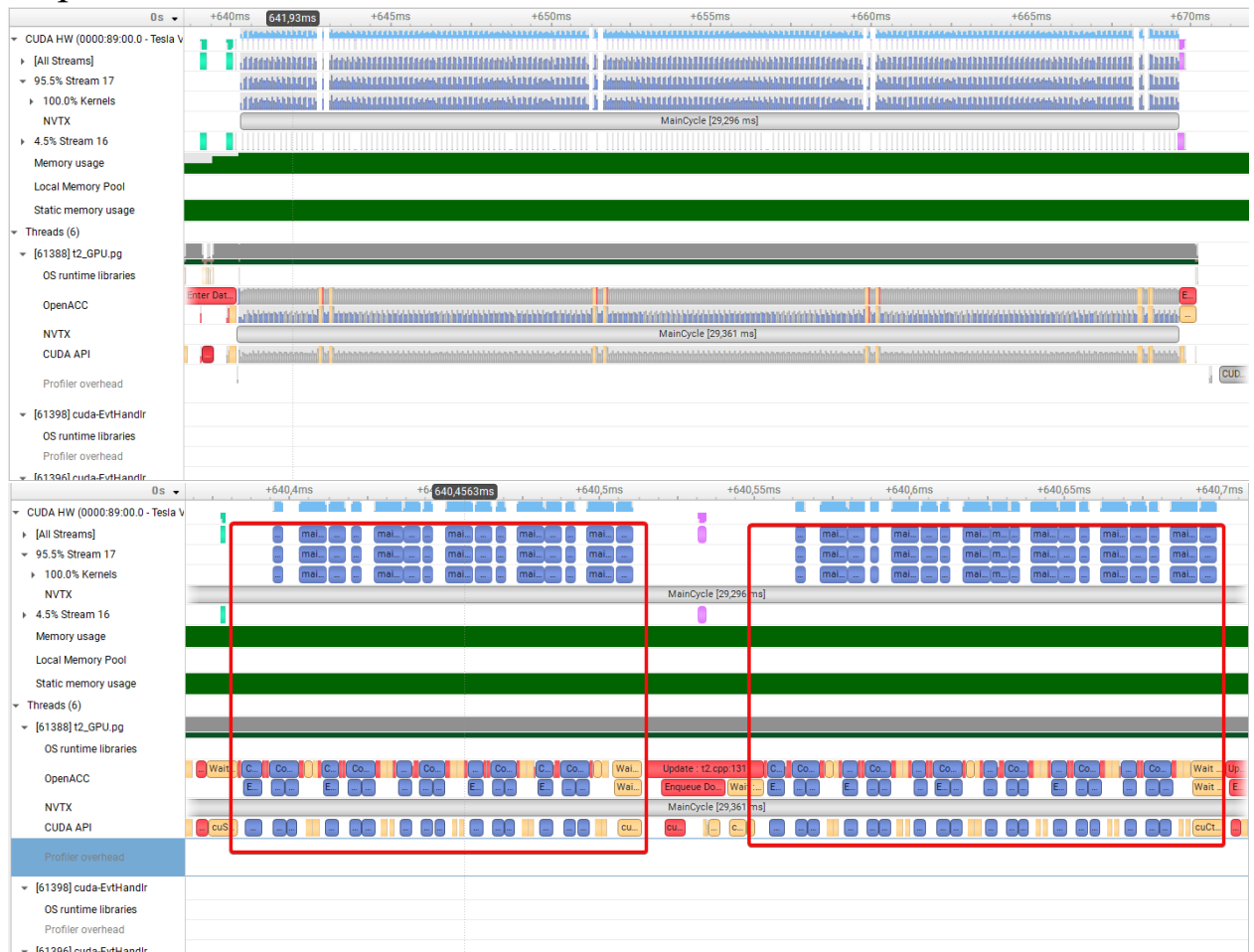
## Этапы оптимизации кода на GPU (cublas)

### Этапы оптимизации на сетке 512x512

№	Время выполнения основного цикла	Ошибка	Кол-во итераций	Комментарий
1	29,361 мс	0.0105525	1000	Код написанный только с использованием орепасс (reduction выполнено через прагмы орепасс, результат прошлого задания)
2	193,742 мс	0.0105524	1000	Выполнение расчета ошибки с помощью cublas (перерасчет каждую итерацию)
3	14,782 мс	0.0140183	1000	Перерасчет ошибки каждые 250 итераций. Значение подобрано эмпирически на основе данных профилировщика, как оптимальное для разных размеров сеток и с учетом времени на синхронизацию потоков.



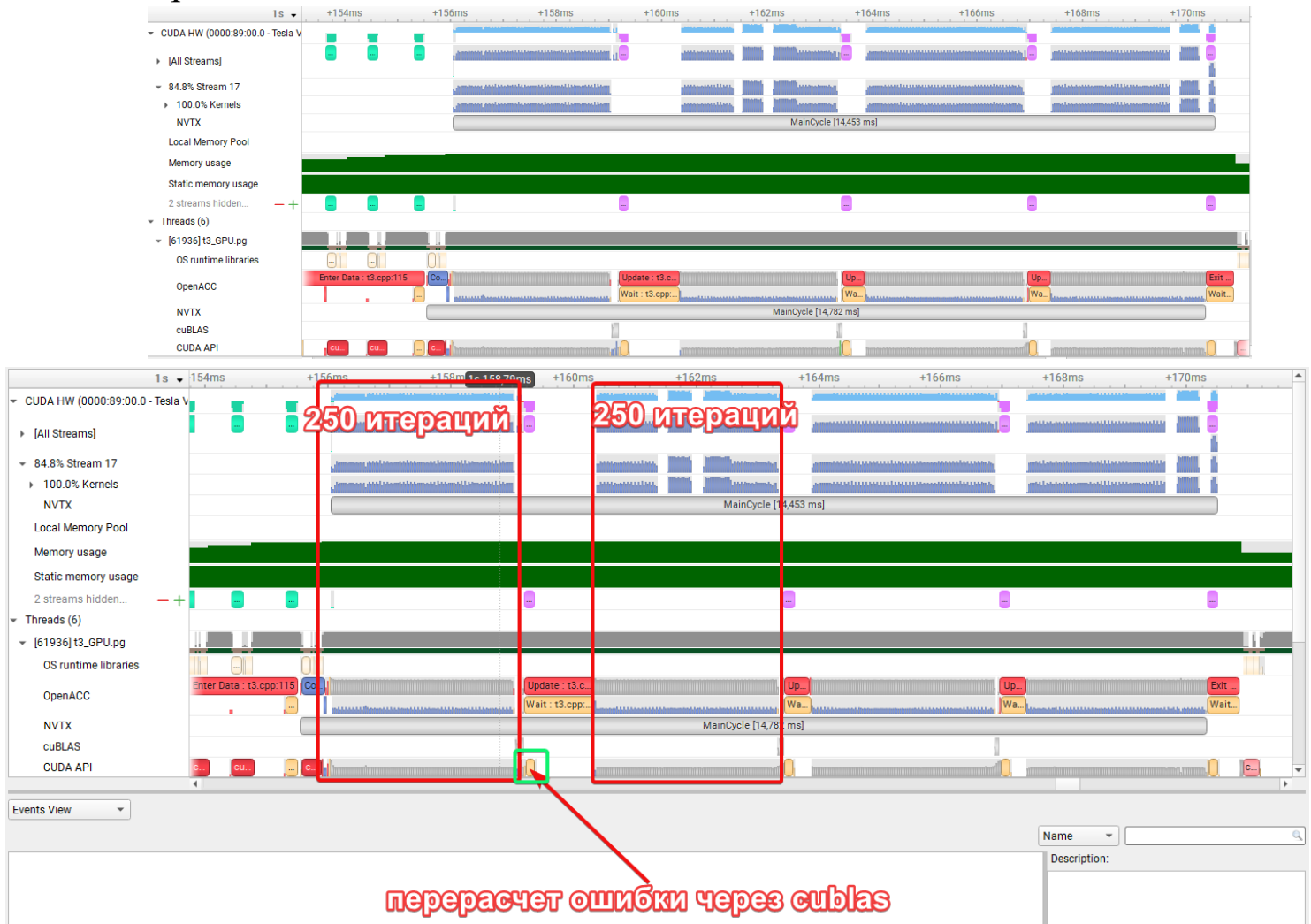
## Первый этап:



## Второй этап



## Третий этап





## Вывод:

Используя `cublas` мы можем добиться меньшего времени за счет того что производим синхронизацию вычислительных потоков как можно реже, а само вычисление ошибки производим достаточно быстро. В итоге получаем время вычисления уравнения на маленьких сетках меньше, чем на CPU, за счет грамотного использования ресурсов видеокарты и минимизации синхронизаций данных между CPU/GPU и вычислительными потоками.

Github: <https://github.com/JooudDoo/Parallelism-Tasks>

```
181 lines (144 sloc) | 5.1 KB
1 #include <iostream>
2 #include <string>
3 #include <sstream>
4 #include <math.h>
5 #include <cmath>
6
7 #ifdef OPENACC_
8 #include <openacc.h>
9 #endif
10 #ifdef NVPROF_
11 #include </opt/nvidia/hpc_sdk/Linux_x86_64/22.11/cuda/11.8/targets/x86_64-linux/include/nvtx3/nvToolsExt.h>
12 #endif
13
14 #include <cublas_v2.h>
15
16 #define at(arr, x, y) (arr[(x)*n+(y)])
17
18 constexpr int LEFT_UP = 10;
19 constexpr int LEFT_DOWN = 20;
20 constexpr int RIGHT_UP = 20;
21 constexpr int RIGHT_DOWN = 30;
22
23 void initArrays(double* mainArr, double* subArr, int& n, int& m, bool initUsingMean){
24     std::memset(mainArr, 0, sizeof(double)*(n)*(m));
25
26     for(int i = 0; i < n*n && initUsingMean; i++){
27         mainArr[i] = (LEFT_UP+LEFT_DOWN+RIGHT_UP+RIGHT_DOWN)/4;
28     }
29
30     at(mainArr, 0, 0) = LEFT_UP;
31     at(mainArr, 0, m-1) = RIGHT_UP;
32     at(mainArr, n-1, 0) = LEFT_DOWN;
33     at(mainArr, n-1, m-1) = RIGHT_DOWN;
34     for(int i = 1; i < n-1; i++){
35         at(mainArr,0,i) = (at(mainArr,0,m-1)-at(mainArr,0,0))/(m-1)*i+at(mainArr,0,0);
36         at(mainArr,i,0) = (at(mainArr,n-1,0)-at(mainArr,0,0))/(n-1)*i+at(mainArr,0,0);
37         at(mainArr,n-1,i) = (at(mainArr,n-1,m-1)-at(mainArr,n-1,0))/(m-1)*i+at(mainArr,n-1,0);
38         at(mainArr,i,m-1) = (at(mainArr,n-1,m-1)-at(mainArr,0,m-1))/(m-1)*i+at(mainArr,0,m-1);
39     }
40     std::memcpy(subArr, mainArr, sizeof(double)*(n)*(m));
41 }
42
43 template<typename T>
44 T extractNumber(char* arr){
45     std::stringstream stream;
46     stream << arr;
47     T result;
48     if (!(stream >> result)){
49         throw std::invalid_argument("Wrong argument type");
50     }
51     return result;
52 }
53
54 constexpr int ITERS_BETWEEN_UPDATE = 250;
55 constexpr double negone = -1;
56
57 int main(int argc, char *argv[]){
58
59     bool showResultArr = false;
60     bool initUsingMean = false;
61     double eps = 1E-6;
62     int iterations = 100;
63     int n = 10;
64     int m = n;
65
66     for(int arg = 0; arg < argc; arg++){
67         if(std::stricmp(argv[arg], "-eps") == 0){
68             eps = extractNumber<double>(argv[arg+1]);
69             arg++;
70         }
71         else if(std::stricmp(argv[arg], "-i") == 0){
72             iterations = extractNumber<int>(argv[arg+1]);
73             arg++;
74         }
75         else if(std::stricmp(argv[arg], "-s") == 0){
76             n = extractNumber<int>(argv[arg+1]);
77             m = n;
78             arg++;
79         }
80         else if(std::stricmp(argv[arg], "-show") == 0){
81             showResultArr = true;
82         }
83         else if(std::stricmp(argv[arg], "-o") == 0){
84             initUsingMean = true;
85         }
86     }
87
88     std::cout << "Current settings:" << std::endl;
89     std::cout << "\tEPS: " << eps << std::endl;
90     std::cout << "\tMax Iteration: " << iterations << std::endl;
91     std::cout << "\tSize: " << n << 'x' << m << std::endl << std::endl;
92
93     double* F = new double[n*m];
94     double* Fnew = new double[n*m];
95
96     double* inter = new double[n*m];
97
98     initArrays(F, Fnew, n, m, initUsingMean);
99
100     double error = 1;
101     int iteration = 0;
102
103     int itersBetweenUpdate = 0;
104
105     cublasHandle_t handle;
106
107     cublasCreate(&handle);
108
109     cublasStatus_t status;
110     int max_idx = 0;
111
112     #pragma acc enter data copyin(Fnew[:n*m], F[:n*m], inter[:n*m])
113
114     #ifdef NVPROF_
115         nvtxRangePush("MainCycle");
116     #endif
117     do {
118         #pragma acc parallel loop collapse(2) present(Fnew[:n*m], F[:n*m]) vector_length(128) async
119         for(int x = 1; x < n-1; x++){
120             for(int y = 1; y < m-1; y++){
121                 at(Fnew,x,y) = 0.25 * (at(F, x+1,y) + at(F,x-1,y) + at(F,x,y-1) + at(F,x,y+1));
122             }
123         }
124
125         double* swap = F;
126         F = Fnew;
127         Fnew = swap;
128
129     #ifdef OPENACC_
130         acc_attach((void**)F);
131         acc_attach((void**)Fnew);
132     #endif
133
134     if(itersBetweenUpdate >= ITERS_BETWEEN_UPDATE && iteration < iterations){
135
136         #pragma acc data present(inter[:n*m], Fnew[:n*m], F[:n*m]) wait
137         {
138             #pragma acc host_data use_device(Fnew, F, inter)
139             {
140
141                 status = cublasDcopy(handle, n*m, F, 1, inter, 1);
142                 if(status != CUBLAS_STATUS_SUCCESS) std::cout << "copy error" << std::endl, exit(30);
143
144                 status = cublasDaxpy(handle, n*m, &negone, Fnew, 1, inter, 1);
145                 if(status != CUBLAS_STATUS_SUCCESS) std::cout << "sum error" << std::endl, exit(40);
146
147                 status = cublasIdamax(handle, n*m, inter, 1, &max_idx);
148                 if(status != CUBLAS_STATUS_SUCCESS) std::cout << "abs max error" << std::endl, exit(41);
149             }
150         }
151
152         #pragma acc update self(inter[:n*m])
153         error = fabs(inter[max_idx]);
154
155         itersBetweenUpdate = -1;
156     }
157     iteration++;
158     itersBetweenUpdate++;
159 } while(iteration < iterations && error > eps);
160
161 #ifdef NVPROF_
162     nvtxRangePop();
163 #endif
164
165     #pragma acc exit data delete(Fnew[:n*m]) copyout(F[:n*m])
166     cublasDestroy(handle);
167
168     std::cout << "Iterations: " << iteration << std::endl;
169     std::cout << "Error: " << error << std::endl;
170     for(int x = 0; x < n && showResultArr; x++){
171         for(int y = 0; y < m; y++){
172             std::cout << at(F,x, y) << ' ';
173         }
174         std::cout << std::endl;
175     }
176
177     delete[] F;
178     delete[] Fnew;
179
180     return 0;
181 }
```