



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 *Робототехники и комплексной автоматизации*

КАФЕДРА *Системы автоматизированного проектирования (РК-6)*

## **ОТЧЕТ О ВЫПОЛНЕНИИ ЛАБОРАТОРНОЙ РАБОТЫ**

по дисциплине:

### **Архитектура параллельных вычислительных систем**

|                          |                                     |
|--------------------------|-------------------------------------|
| Студент                  | Абидоков Рашид Ширамбиевич          |
| Группа                   | РК6-11М                             |
| Тип задания              | лабораторная работа                 |
| Тема лабораторной работы | Применение технологий OpenMP и CUDA |

|               |       |  |
|---------------|-------|--|
| Студент       | _____ | <b><u>Абидоков Р. Ш.</u></b><br><i>подпись, дата</i><br><i>фамилия, и.о.</i> |
| Преподаватель | _____ | <b><u>Спасенов А. Ю.</u></b><br><i>подпись, дата</i><br><i>фамилия, и.о.</i> |

*Москва, 2020 г.*

## Оглавление

|  |    |
|--|----|
| Цель выполнения лабораторной работы..... | 3  |
| Задание на лабораторную работу .....     | 3  |
| 1. Реализация с помощью OpenMP .....     | 4  |
| 2. Реализация с помощью CUDA .....       | 5  |
| 3. Сравнение времени работы .....        | 7  |
| Заключение .....                         | 10 |

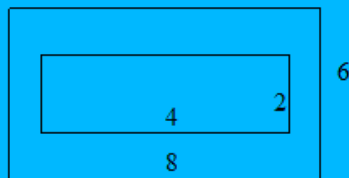
## Цель выполнения лабораторной работы

Цель выполнения лабораторной работы – получение практических навыков применения технологий OpenMP и CUDA, оценка времени работы программы.

## Задание на лабораторную работу

### Задание 1

Методом конечных разностей решить задачу распространения тепла в стенках трубы, изображенной на рис. Выполнить два варианта расчета с кратными шагами сетки. Уточнить решение с помощью экстраполяции Ричардсона.



Размеры сторон показаны на рисунке. Граничные условия : внутри трубы протекает жидкость с температурой 100 градусов, на внешней границе задано условие:

$$dT/dn = T.$$

Отчет должен содержать: текст программы, таблицы результатов в совпадающих узлах сетки, рисунок объекта с распределением фазовой переменной, полученной после экстраполяции, сравнение результатов расчета с результатами, полученными с помощью пакета ANSYS.

Рис. 1 Исходное задание

Однако для большей наглядности применения методов распараллеливания условие задачи было изменено – требуется решить нестационарную задачу, т.е. получить распределение температуры в зависимости от времени. Начальной температурой пластины считается 20 градусов Цельсия.

## 1. Реализация с помощью OpenMP

Поскольку практически всё время работы программы занимает решение СЛАУ методом Гаусса, распараллеливанию как в данном пункте, так и при использовании CUDA подвергался именно метод Гаусса. Исходный код реализующей его функции приведен в Листинге 1.

Листинг 1.

```
118 template <typename T>
119 vector<T>& gauss(const Matrix<T>& l_matr, const vector<T>& r_vect) {
120     unsigned size_ = r_vect.size();
121     if (l_matr.get_cols() != l_matr.get_rows()) {
122         cout << "not square matrix in gauss" << endl;
123         throw ("not square matrix in gauss");
124     }
125     if (size_ != l_matr.get_rows()) {
126         cout << "l_matr size != r_vect size in gauss" << endl;
127         throw ("l_matr size != r_vect size in gauss");
128     }
129     Matrix<T> l_matr_ = l_matr;
130     vector<T> r_vect_ = r_vect;
131     vector<T>& ret_vect = *(new vector<T>(size_));
132
133     // Прямой ход
134     // Внешний нельзя, важен порядок
135     for (unsigned i = 0; i < size_; i++) {
136         #pragma omp parallel for
137         for (unsigned j = 0; j < i; j++) {
138             #pragma omp atomic
139             r_vect_[i] -= r_vect_[j] * l_matr_[i][j];
140             l_matr_.plus_row(i, j, -l_matr_[i][j]);
141         }
142         r_vect_[i] /= l_matr_[i][i];
143         l_matr_.set_diag_to_one_r(i);
144     }
145
146     // Обратный ход
147     for (unsigned i = 0; i < size_; i++) {
148         T sum = 0;
149         #pragma omp parallel for
150         // Здесь нельзя трогать наружный цикл
151         for (unsigned j = 0; j < i; j++) {
152             #pragma omp atomic
153             sum += l_matr_[size_-i-1][size_-j-1] * ret_vect[size_-j-1];
154         }
155         ret_vect[size_-i-1] = r_vect_[size_-i-1] - sum;
156     }
```

В явном виде количество потоков задается в функции main, фрагмент которой приведен в Листинге 2.

Листинг 2.

```
24  int main(int argc, char *argv[]) {
25      // Выставляем количество потоков
26      omp_set_dynamic(0);
27      omp_set_num_threads(4);
```

## 2. Реализация с помощью CUDA

Поскольку прямой ход метода Гаусса имеет сложность  $O(n^3)$ , а обратный, в свою очередь,  $O(n^2)$ , при большом количестве узлов (когда распараллеливание является эффективным) наибольшее влияние на время выполнения оказывает именно прямой ход. Исходный код функции, реализующей метод Гаусса, приведен в Листинге 3.

Листинг 3.

```
197  vector<double>& gauss_cuda(const Matrix& l_matr, const vector<double>& r_vect) {
198      size_t size = r_vect.size();
199      // Делаем копию вектора правой части в виде массива
200      double* r_vec = new double[size];
201      for (size_t i = 0; i < size; i++) {
202          r_vec[i] = r_vect[i];
203      }
204      // Количество блоков и тредов
205      dim3 N_threads1(8);
206      dim3 N_blocks1(size / 8 + 1);
207      dim3 N_threads_once(1);
208      dim3 N_blocks_once(1);
209      // Выделяем память на девайсе
210      double* l_matr_dev;
211      double* r_vec_dev;
212      double* ret_vec_dev;
213      size_t l_matr_size = size * size * sizeof(double);
214      size_t vec_size = size * sizeof(double);
215      cudaMalloc((void**) &l_matr_dev, l_matr_size);
216      cudaMalloc((void**) &r_vec_dev, vec_size);
217      cudaMalloc((void**) &ret_vec_dev, vec_size);
218      // Зануляем правый вектор
219      fillVecKernel<<<N_blocks1, N_threads1>>>(ret_vec_dev, size, 0);
220      // Копируем туда данные
221      cudaMemcpy(l_matr_dev, l_matr.get_arr(), l_matr_size, cudaMemcpyHostToDevice);
222      cudaMemcpy(r_vec_dev, r_vec, vec_size, cudaMemcpyHostToDevice);
```

```

223 // Прямой ход
224 // Приводим матрицу к верхнетреугольной
225 double* mcoeff;
226 cudaMalloc((void**)&mcoeff, sizeof(double));
227 for (size_t i = 0; i < size; i++) {
228     for (size_t j = 0; j < i; j++) {
229         calculateMCoeffKernel<<<N_threads_once, N_blocks_once >>>(l_matr_dev, size, i, j, mcoeff);
230         plusRowKernel<<<N_blocks1, N_threads1>>>(l_matr_dev, r_vec_dev, size, i, j, mcoeff);
231     }
232 }
233 cudaFree(mcoeff);
234 cudaDeviceSynchronize();
235 // Выводим единицы на главной диагонали
236 // Вытаскиваем диагональные элементы, чтобы модифицировать матрицу
237 double* diag_vec_dev;
238 cudaMalloc((void**)&diag_vec_dev, vec_size);
239 diagElmsToVecKernel <<<N_blocks1, N_threads1 >>> (l_matr_dev, diag_vec_dev, size);
240 // Диагональные элементы строк в единицы
241 for (size_t i = 0; i < size; i++) {
242     setDiagToOneKernel <<<N_blocks1, N_threads1 >>> (l_matr_dev, r_vec_dev, diag_vec_dev, size, i);
243 }
244 // И правый вектор
245 rvecDiagDivKernel<<<N_blocks1, N_threads1 >>>(r_vec_dev, diag_vec_dev, size);
246 cudaFree(diag_vec_dev);
247 // Обратный ход
248 double* l_matr_host = new double[size*size];
249 double* r_vec_host = new double[size];
250 cudaMemcpy(l_matr_host, l_matr_dev, l_matr_size, cudaMemcpyDeviceToHost);
251 cudaMemcpy(r_vec_dev, r_vec_dev, vec_size, cudaMemcpyDeviceToHost);
252 vector<double>& ret_vect = *(new vector<double>(size));
253 for (size_t i = 0; i < size; i++) {
254     double sum = 0;
255     for (size_t j = 0; j < i; j++) {
256         sum += l_matr_host[(size - i - 1)*size + (size - j - 1)] * ret_vect[size - j - 1];
257     }
258     ret_vect[size - i - 1] = r_vec_host[size - i - 1] - sum;
259 }

```

Используемые Kernel-функции в Листингах 4 и 5.

Листинг 4.

```

153 __global__ void fillVecKernel(double* vec, size_t size, double val) {
154     size_t idx = threadIdx.x + blockIdx.x * blockDim.x;
155     if (idx < size) {
156         vec[idx] = val;
157     }
158 }
159 __global__ void calculateMCoeffKernel(double* l_matr, size_t size, size_t i, size_t j, double* mcoeff) {
160     *mcoeff = l_matr[i * size + j] / l_matr[j * size + j];
161 }
162
163 __global__ void plusRowKernel(double* l_matr, double* r_vec, size_t size,
164                             size_t mr_idx, size_t pr_idx, double* mcoeff) {
165     size_t idx = threadIdx.x + blockIdx.x * blockDim.x;
166     if (idx < size) {
167         l_matr[mr_idx * size + idx] -= l_matr[pr_idx * size + idx] * (*mcoeff);
168     }
169 }

```

Листинг 5.

```

171 __global__ void diagElemsToVecKernel(double* l_matr, double* diag_vec, size_t size) {
172     size_t idx = threadIdx.x + blockIdx.x * blockDim.x;
173     if (idx < size) {
174         diag_vec[idx] = l_matr[idx*size + idx];
175     }
176 }
177
178 __global__ void setDiagToOneKernel(double* l_matr, double* r_vec, double* diag_vec, size_t size,
179                                     size_t mr_idx) {
180     size_t idx = threadIdx.x + blockIdx.x * blockDim.x;
181     if (idx < size) {
182         l_matr[mr_idx * size + idx] /= diag_vec[mr_idx];
183     }
184 }
185
186 __global__ void rvecDiagDivKernel(double* r_vec, double* diag_vec, double size) {
187     size_t idx = threadIdx.x + blockIdx.x * blockDim.x;
188     if (idx < size) {
189         r_vec[idx] /= diag_vec[idx];
190     }
191 }

```

### 3. Сравнение времени работы

Сравнение времени работы однопоточной и многопоточных реализаций проводилось на компьютере с процессором Intel Core i7-6500u, имеющем 2 ядра с 4 потоками, и видеокартой gtx-960m. Для автоматизации использовался скрипт, написанный на языке python, код которого приведен в Листинге 6.

Листинг 6.

```

1  import numpy as np
2  import pandas as pd
3  import subprocess
4
5
6  single_path = "single\\out\\main.exe"
7  openmp_path = "openmp\\out\\main.exe"
8  cuda_path = "cuda\\x64\\Release\\cuda.exe"
9
10 # Сравниваем при различных разбиениях при постоянном количестве шагов интегрирования
11 t_end = 10
12 h_t = 0.2
13 N_x_list = np.append(np.array([4]), [np.arange(8, 48, 8)])
14 N_y_list = np.append(np.array([3]), [np.arange(6, 36, 6)])
15 N_list = np.array([N_x_list, N_y_list]).T
16 print(N_list)
17

```

```

18 N_nodes_data = pd.DataFrame(columns=['N_nodes', 'cuda_total', 'cuda_avg'])
19 for i, [N_x, N_y] in enumerate(N_list):
20     print(i, N_x, N_y)
21     # Однопоточное решение
22     single_curr = subprocess.run([single_path, str(t_end), str(h_t), str(N_x), str(N_y), '-d'], stdout=subprocess.PIPE).stdout
23     n_nodes, single_total, single_avg = [float(s) for s in single_curr.split()]
24     print('single_total', single_total)
25     print('single_avg', single_avg)
26     # OpenMP решение
27     openmp_curr = subprocess.run([openmp_path, str(t_end), str(h_t), str(N_x), str(N_y), '-d'], stdout=subprocess.PIPE).stdout
28     n_nodes, openmp_total, openmp_avg = [float(s) for s in openmp_curr.split()]
29     print('openmp_total', openmp_total)
30     print('openmp_avg', openmp_avg)
31     # CUDA решение
32     cuda_curr = subprocess.run([cuda_path, str(t_end), str(h_t), str(N_x), str(N_y), '-d'], stdout=subprocess.PIPE).stdout
33     n_nodes, cuda_total, cuda_avg = [float(s) for s in cuda_curr.split()]
34     print('cuda_total', cuda_total)
35     print('cuda_avg', cuda_avg)
36     print('\n')
37     N_nodes_data.loc[i] = [n_nodes, single_total, single_avg, openmp_total, openmp_avg, cuda_total, cuda_avg]
38
39 print(N_nodes_data)
40 N_nodes_data.to_csv("N_nodes_compare.csv")

```

Всё время указано в мс, столбцы "шаг" – общее время, деленное на количество шагов интегрирования

Таблица 1.

| N узлов | Однопоточн.<br>всего | Однопоточн.<br>шаг | OpenMP<br>всего | OpenMP<br>шаг | CUDA<br>всего | CUDA<br>шаг |
|---------|----------------------|--------------------|-----------------|---------------|---------------|-------------|
| 40      | 31.0                 | 0.6                | 299.0           | 6.0           | 130.0         | 2.6         |
| 160     | 1354.0               | 27.08              | 1794.0          | 35.88         | 594.0         | 9.2         |
| 360     | 14783.0              | 295.66             | 10096.0         | 201.92        | 1975.0        | 35.22       |
| 640     | 82487.0              | 1649.74            | 48314.0         | 966.28        | 4955.0        | 92.08       |
| 1000    | 287670.0             | 5753.4             | 179644.0        | 3592.88       | 11382.0       | 214.68      |

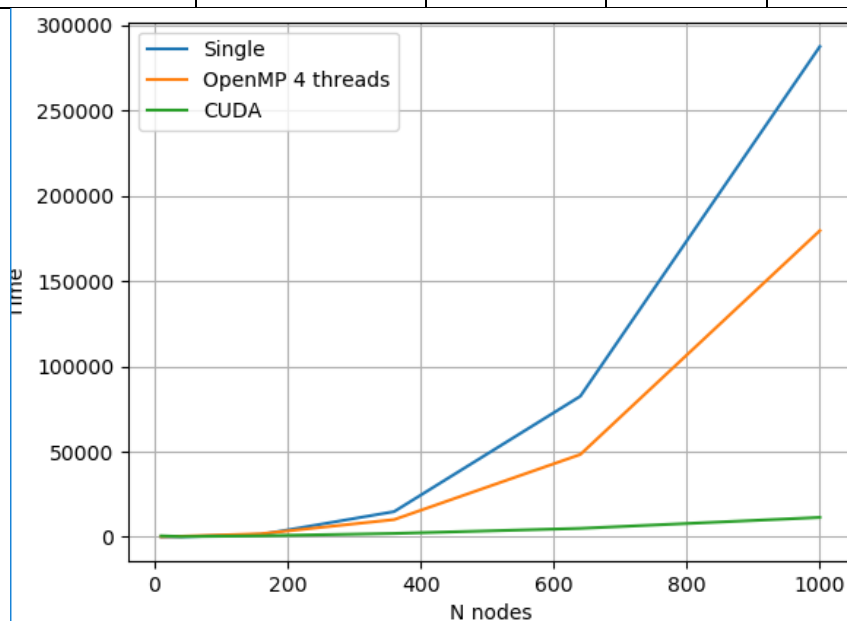


Рис 2.



Более наглядным является график кубического корня времени от количества узлов. Т.к. сложность алгоритма  $O(n^3)$ , графики похожи на прямые линии. Видно, что при малом количестве узлов (и, как следствие, малой сложности) параллельные реализации проигрывают в скорости однопоточной, что связано с накладными расходами при распараллеливании.

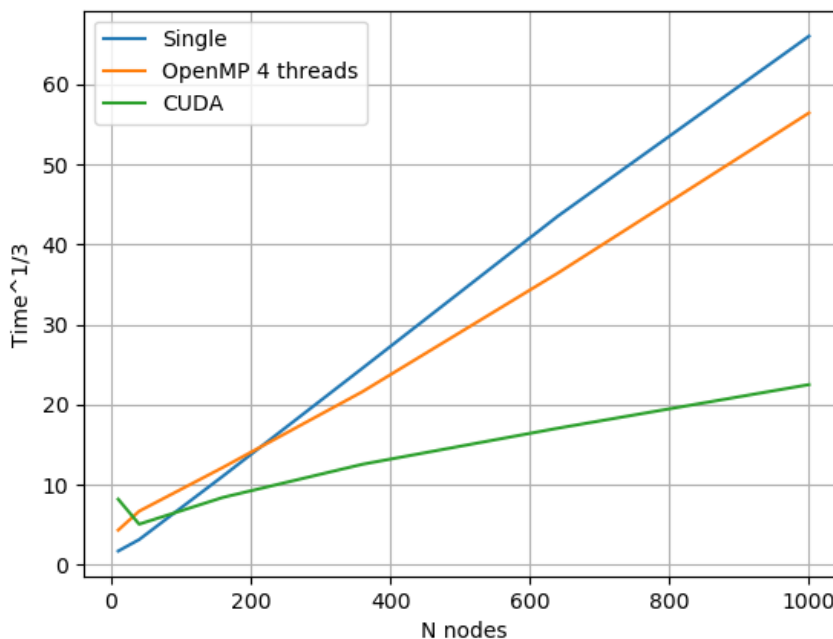


Рис 3.

Графики отношения времени однопоточной реализации к времени многопоточных

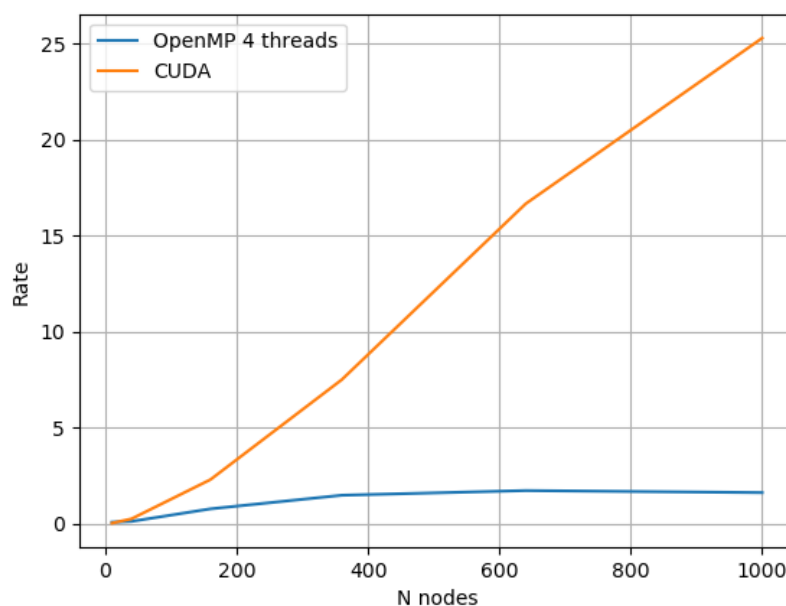


Рис 4.

## **Заключение**

Было реализовано распараллеливание метода Гаусса с использованием технологий OpenMP и CUDA. Проведено сравнение времени работы однопоточной и многопоточных реализаций, показано, что даже при локальном применении данных технологий прирост производительности является существенным.

При этом в случае, если вычислительная сложность задачи не очень велика, распараллеливание может не только не привести к ускорению работы программы, но даже замедлить её, что связано с возникающими накладными расходами, связанными с управлением памятью и синхронизацией потоков.