

Теория параллелизма

Отчёт

Оптимизированные библиотеки

Выполнил: Лейсле Александр Геннадьевич, группа 21932

28.03.2023

Цели работы:

1. Реализовать решение уравнение теплопроводности (пятиточечный шаблон) в двумерной области на равномерных сетках:
 - 128x128
 - 256x256
 - 512x512
 - 1024x1024
2. Перенести программу на GPU используя директивы OpenACC
3. Операцию редукции (вычисление максимального значения ошибки) на графическом процессоре реализовать через вызовы функций из библиотеки cuBLAS
4. Сравнить скорость работы для разных размеров сеток на:
 - Нескольких ядрах CPU
 - GPU (с реализацией с библиотекой cuBLAS и без неё)
5. Произвести профилирование программы с использованием NsightSystems.

Используемый компилятор: `pgc++`

Используемый профилировщик: `nsys`

Как производили замер времени работы:

В начале и в конце программы производилась фиксация текущего времени с использованием команд из `chrono`, разница этого времени выводилась в стандартный поток вывода.

Реализация без библиотеки cuBLAS

CPU-multicore

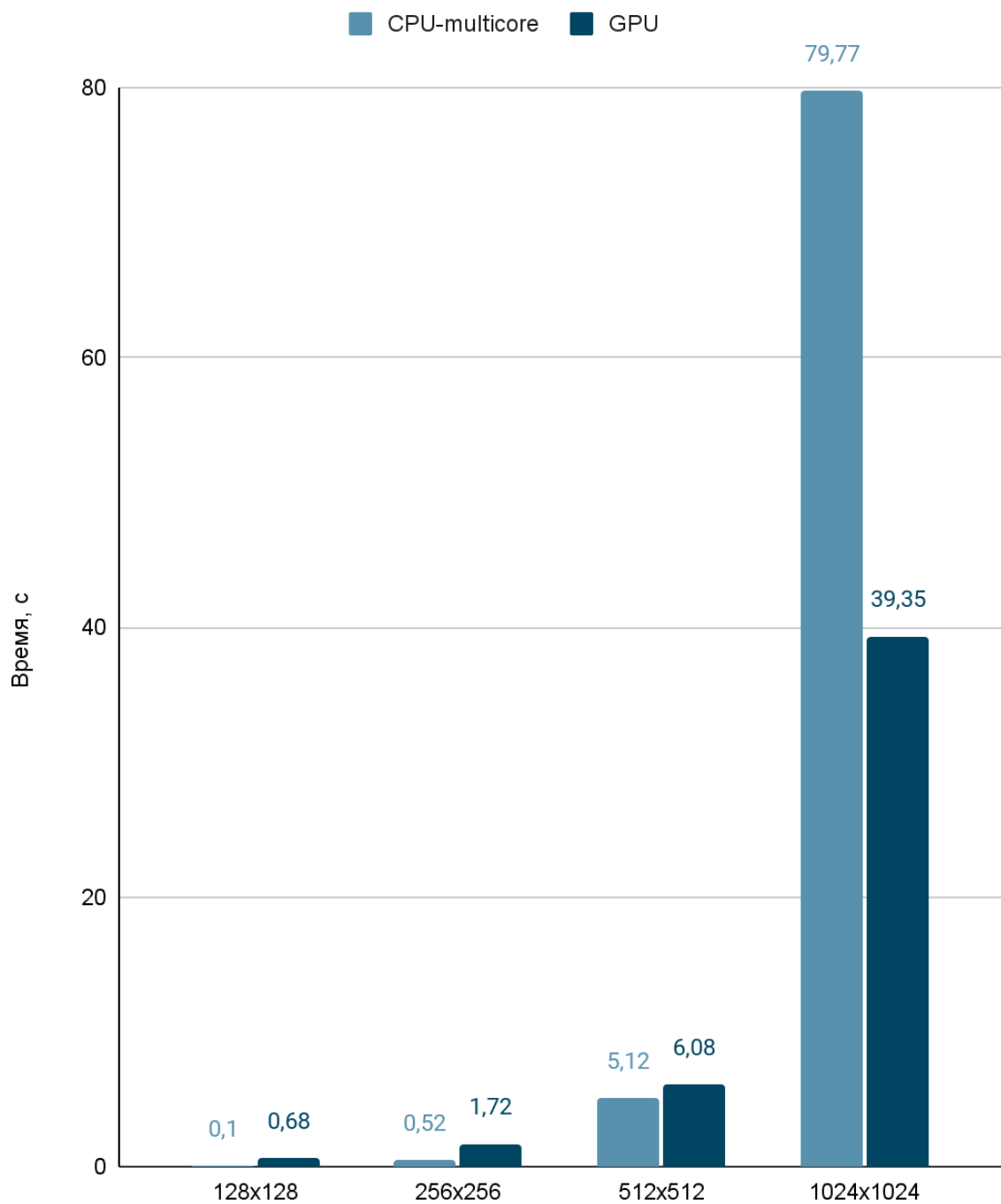
Размер сетки	Время выполнения	Точность	Количество итераций
128x128	0.10с	$9.9 * 10^{-7}$	30080
256x256	0.52с	$9.9 * 10^{-7}$	102912
512x512	5.12с	$9.9 * 10^{-7}$	339712
1024x1024	79.77с	$1.37 * 10^{-6}$	1000000

GPU

Размер сетки	Время выполнения	Точность	Количество итераций
128x128	0.68с	$9.9 * 10^{-7}$	30080
256x256	1.72с	$9.9 * 10^{-7}$	102912
512x512	6.08с	$9.9 * 10^{-7}$	339712
1024x1024	39.35с	$1.37 * 10^{-6}$	1000000

Диаграмма сравнения времени работы CPU-multicore и GPU (реализация без cuBLAS)

Общее время работы



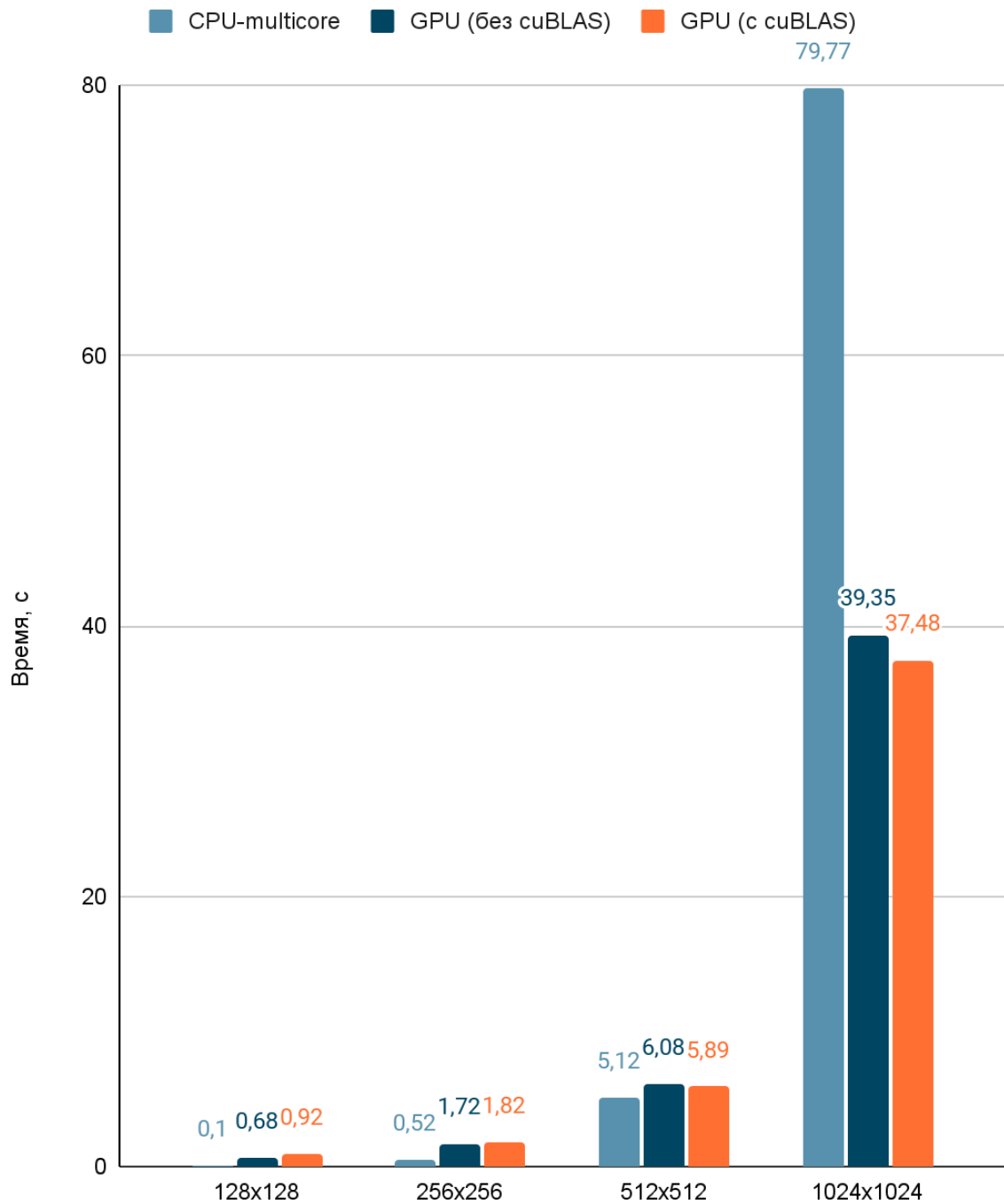
Реализация с библиотекой cuBLAS

GPU

Размер сетки	Время выполнения	Точность	Количество итераций
128x128	0.92с	$9.9 * 10^{-7}$	30080
256x256	1.82с	$9.9 * 10^{-7}$	102912
512x512	5.89с	$9.9 * 10^{-7}$	339712
1024x1024	37.48с	$1.37 * 10^{-6}$	1000000

Диаграмма сравнения времени работы CPU-multicore и GPU (реализация с cuBLAS и без cuBLAS)

Общее время работы



Выполнение на GPU

Этапы оптимизации на сетке 512x512

Этап №	Время выполнения	Точность	Количество итераций	Комментарии (что было сделано)
1	0.20с	0.11	100	<p>Реализация только с использованием директив OpenACC;</p> <p>Сетка представляется в виде одномерного массива;</p> <p>Выделение памяти под сетки на текущем и новом шагах происходит на CPU, затем обе сетки отправляются на GPU;</p> <p>Под обе сетки выделяется память, и происходит копирование массивов (которые заполнены нулями) с помощью асс data copyout. По выходе из секции данные массивов копируются обратно на CPU;</p> <p>В CPU происходит выделение памяти под массив из 4 элементов,</p>

			<p>представляющих собой размеры шагов для линейной интерполяции между углами сетки по границам. Под данный массив происходит выделение памяти на GPU с помощью асс <code>data create</code>;</p> <p>В память GPU копируются значения: общего количества элементов сетки, размера стороны сетки и максимальная ошибка, которая высчитывается между итерациями (в начальный момент превышает заданную ошибку на 1);</p> <p>Обращение к индексам массива происходит с помощью макроса <code>ID</code>, который позволяет перевести индексы двумерного массива в индекс одномерного;</p> <p>Заполнение границ вспомогательной и основной сеток происходит на GPU асинхронно при помощи директивы асс <code>asunc</code>. При этом циклы, в которых происходит</p>
--	--	--	---

			<p>заполнение, распараллеливаются с помощью асс parallel loop independent. В конце заполнения происходит синхронизация для обеих сеток;</p> <p>Новое состояние каждой точки сетки рассчитывается через среднее 4х соседних точек. Это происходит в 2х вложенных циклах, которые распараллеливаются с помощью директивы асс parallel loop collapse(2) independent present(grid, grid_new).</p> <p>Для вычисления максимальной ошибки между двумя итерациями используется редукция по функции максимума для переменной max_error с помощью директивы асс reduction(max:max_error)</p> <p>Обновление основной сетки происходит с помощью смены указателей сеток на текущем и новом шагах</p>
--	--	--	--

				<p>на CPU, а затем с помощью функции <code>асс_attach</code> происходит пересвязывание адресов этих сеток в памяти GPU;</p> <p>Происходит вычисление значения максимальной ошибки, а затем обновление в памяти CPU, если прошедшее количество итераций кратно половине размера стороны сетки или пройдено заданное количество итераций, с помощью <code>асс update host</code></p>
2	0.53с	0.11	100	<p>Реализация с использованием директив <code>OpenACC</code> и функций <code>cuBLAS</code>;</p> <p>На CPU создается контекст <code>cuBLAS</code> и инициализируется с помощью функции <code>cublasCreate</code>;</p> <p>Режим обращения к указателям для функций <code>cuBLAS</code> устанавливается в режим обращения на GPU;</p>

				<p>Если прошедшее количество итераций кратно половине размера стороны сетки или пройдено заданное количество итераций, то вычисляется максимальная ошибка с помощью функций cuBLAS:</p> <p>Текущее состояние сетки сохраняется во временный массив путем копирования сетки с помощью функции cublasDcopy;</p> <p>Из нового состояния сетки вычитается скопированный массив с помощью cublasDaxpy, результат записывается на место скопированного массива;</p> <p>В конце с помощью функции cublasIdamax в переменную id записывается адрес самого большого по модулю значения, которое представляет собой значение наибольшей ошибки между двумя итерациями.</p>
--	--	--	--	--

				<p>В переменную max_err г записывается значение наибольшей ошибки</p> <p>После обновления текущего состояния сетки происходит обновление максимальной ошибки в памяти CPU.</p>
--	--	--	--	--

Вывод

В ходе работы было реализовано решение уравнения теплопроводности на сетках размера 128x128, 256x256, 512x512, 1024x1024. Программа была перенесена на GPU с использованием директив OpenACC. Затем была реализована версия этой программы, где вычисление максимальной ошибки производилось с использованием функций cuBLAS. Затем было выполнено сравнение скорости выполнения на разных размерах сеток, и в конце было выполнено профилирование на Nsight System.

Как видно, использование функций cuBLAS позволяет производить вычисления, связанные с матрицами и векторами, за короткое время (нахождение индекса максимального элемента в массиве размерностью 262144 происходит примерно за 32 мкс), что позволяет ускорить выполнение программы на больших сетках, однако инициализация и удаление контекста cuBLAS занимает достаточно большое время, что нужно учитывать при запуске на небольшом количестве итераций или на небольших сетках.

Приложение

https://github.com/HerrPhoton/Optimized_libraries

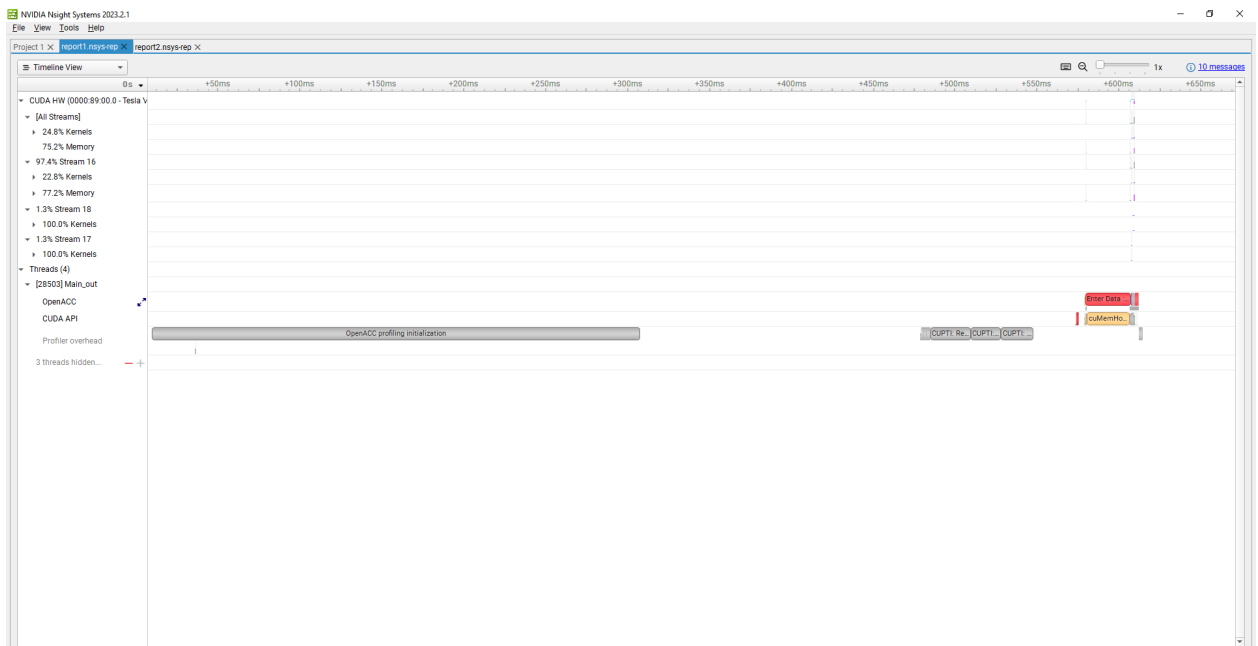
```
● sanya@sanya-desktop:~/Lab3$ ./out -side 15 -iters 1000000 -error 0.000001 -show
Passed iterations: 539/1000000
Maximum error: 9.52953e-07/1e-06
Total execution time: 80.759 ms
10 10.7143 11.4286 12.1429 12.8571 13.5714 14.2857 15 15.7143 16.4286 17.1429 17.8571 18.5714 19.2857 20
10.7143 11.4286 12.1429 12.8571 13.5714 14.2857 15 15.7143 16.4286 17.1428 17.8571 18.5714 19.2857 20 20.7143
11.4286 12.1429 12.8571 13.5714 14.2857 15 15.7143 16.4286 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4286
12.1429 12.8571 13.5714 14.2857 15 15.7143 16.4285 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4286 22.1429
12.8571 13.5714 14.2857 15 15.7143 16.4285 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4286 22.1429 22.8571
13.5714 14.2857 15 15.7143 16.4285 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4286 22.1428 22.8571 23.5714
14.2857 15 15.7143 16.4285 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4285 22.1428 22.8571 23.5714 24.2857
15 15.7143 16.4286 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4285 22.1428 22.8571 23.5714 24.2857 25
15.7143 16.4286 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4285 22.1428 22.8571 23.5714 24.2857 25 25.7143
16.4286 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4285 22.1428 22.8571 23.5714 24.2857 25 25.7143 26.4286
17.1429 17.8571 18.5714 19.2857 20 20.7143 21.4285 22.1428 22.8571 23.5714 24.2857 25 25.7143 26.4286 27.1429
17.8571 18.5714 19.2857 20 20.7143 21.4286 22.1428 22.8571 23.5714 24.2857 25 25.7143 26.4286 27.1429 27.8571
18.5714 19.2857 20 20.7143 21.4286 22.1428 22.8571 23.5714 24.2857 25 25.7143 26.4286 27.1429 27.8571 28.5714
19.2857 20 20.7143 21.4286 22.1429 22.8571 23.5714 24.2857 25 25.7143 26.4286 27.1429 27.8571 28.5714 29.2857
20 20.7143 21.4286 22.1429 22.8571 23.5714 24.2857 25 25.7143 26.4286 27.1429 27.8571 28.5714 29.2857 30
○ sanya@sanya-desktop:~/Lab3$ █
```

Вывод сетки 15x15 на реализации без cuBLAS

```
● sanya@sanya-desktop:~/Lab3$ ./out -side 15 -iters 1000000 -error 0.000001 -show
Passed iterations: 539/1000000
Maximum error: 9.52953e-07/1e-06
Total execution time: 754.439 ms
10 10.7143 11.4286 12.1429 12.8571 13.5714 14.2857 15 15.7143 16.4286 17.1429 17.8571 18.5714 19.2857 20
10.7143 11.4286 12.1429 12.8571 13.5714 14.2857 15 15.7143 16.4286 17.1428 17.8571 18.5714 19.2857 20 20.7143
11.4286 12.1429 12.8571 13.5714 14.2857 15 15.7143 16.4286 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4286
12.1429 12.8571 13.5714 14.2857 15 15.7143 16.4285 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4286 22.1429
12.8571 13.5714 14.2857 15 15.7143 16.4285 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4286 22.1429 22.8571
13.5714 14.2857 15 15.7143 16.4285 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4286 22.1428 22.8571 23.5714
14.2857 15 15.7143 16.4285 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4285 22.1428 22.8571 23.5714 24.2857
15 15.7143 16.4286 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4285 22.1428 22.8571 23.5714 24.2857 25
15.7143 16.4286 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4285 22.1428 22.8571 23.5714 24.2857 25 25.7143
16.4286 17.1428 17.8571 18.5714 19.2857 20 20.7143 21.4285 22.1428 22.8571 23.5714 24.2857 25 25.7143 26.4286
17.1429 17.8571 18.5714 19.2857 20 20.7143 21.4285 22.1428 22.8571 23.5714 24.2857 25 25.7143 26.4286 27.1429
17.8571 18.5714 19.2857 20 20.7143 21.4286 22.1428 22.8571 23.5714 24.2857 25 25.7143 26.4286 27.1429 27.8571
18.5714 19.2857 20 20.7143 21.4286 22.1428 22.8571 23.5714 24.2857 25 25.7143 26.4286 27.1429 27.8571 28.5714
19.2857 20 20.7143 21.4286 22.1429 22.8571 23.5714 24.2857 25 25.7143 26.4286 27.1429 27.8571 28.5714 29.2857
20 20.7143 21.4286 22.1429 22.8571 23.5714 24.2857 25 25.7143 26.4286 27.1429 27.8571 28.5714 29.2857 30
○ sanya@sanya-desktop:~/Lab3$ █
```

Вывод сетки 15x15 на реализации с cuBLAS

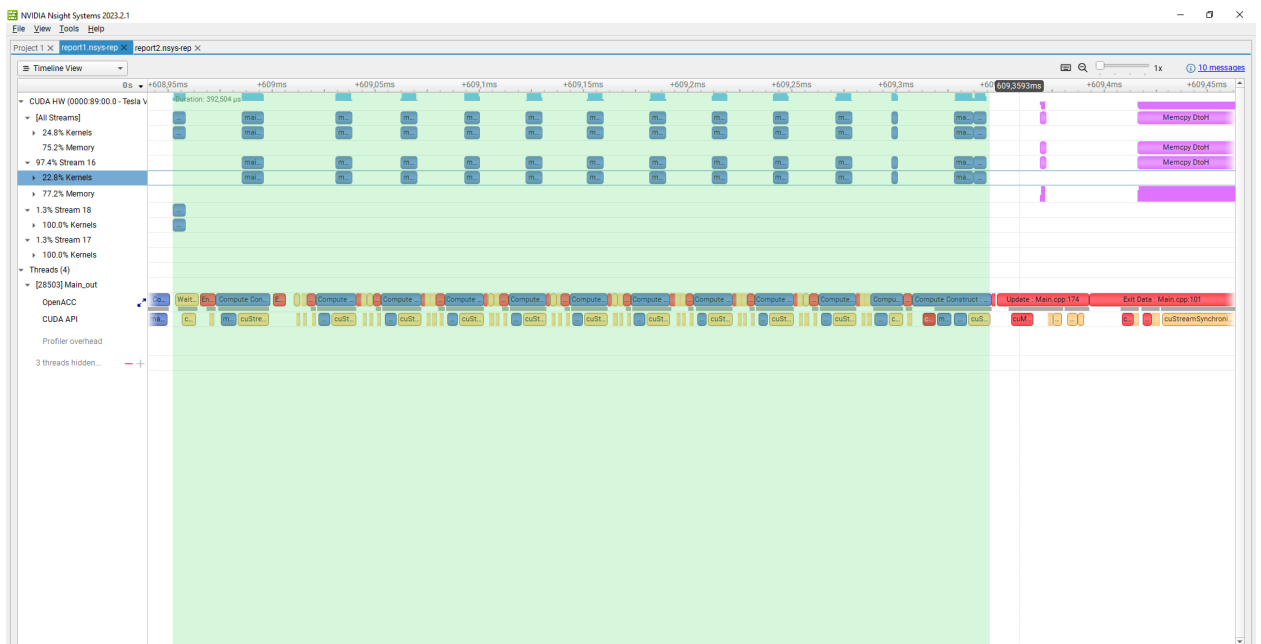
Профилирование реализации без cuBLAS



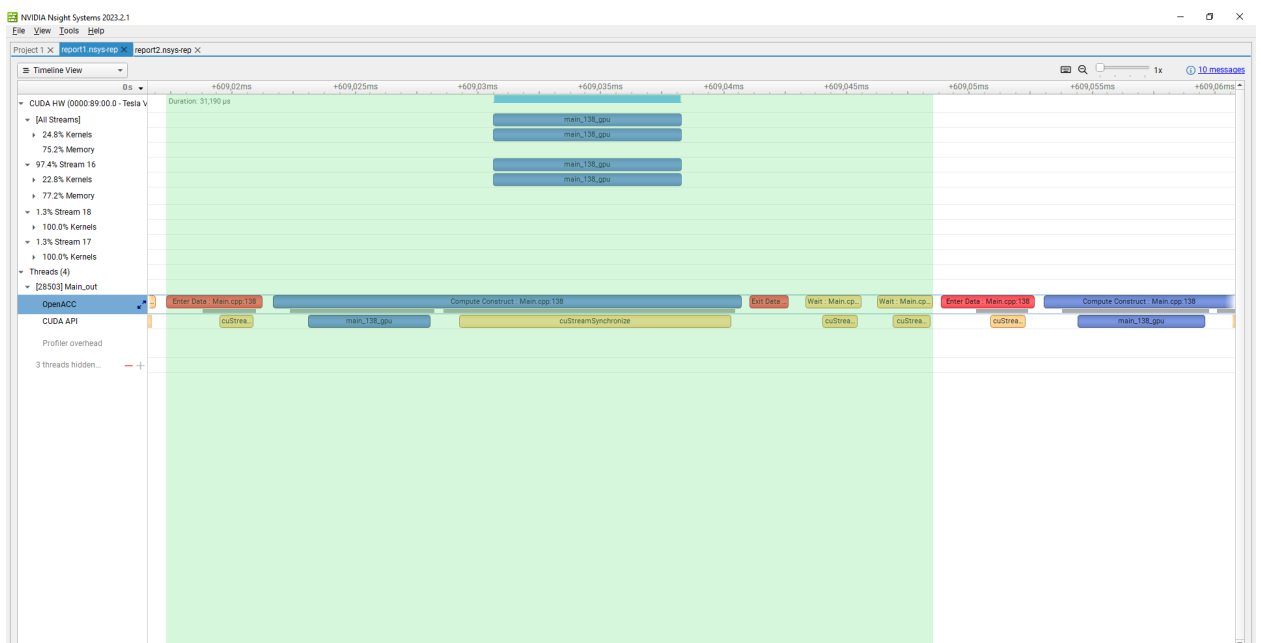
Общий вид профилирования



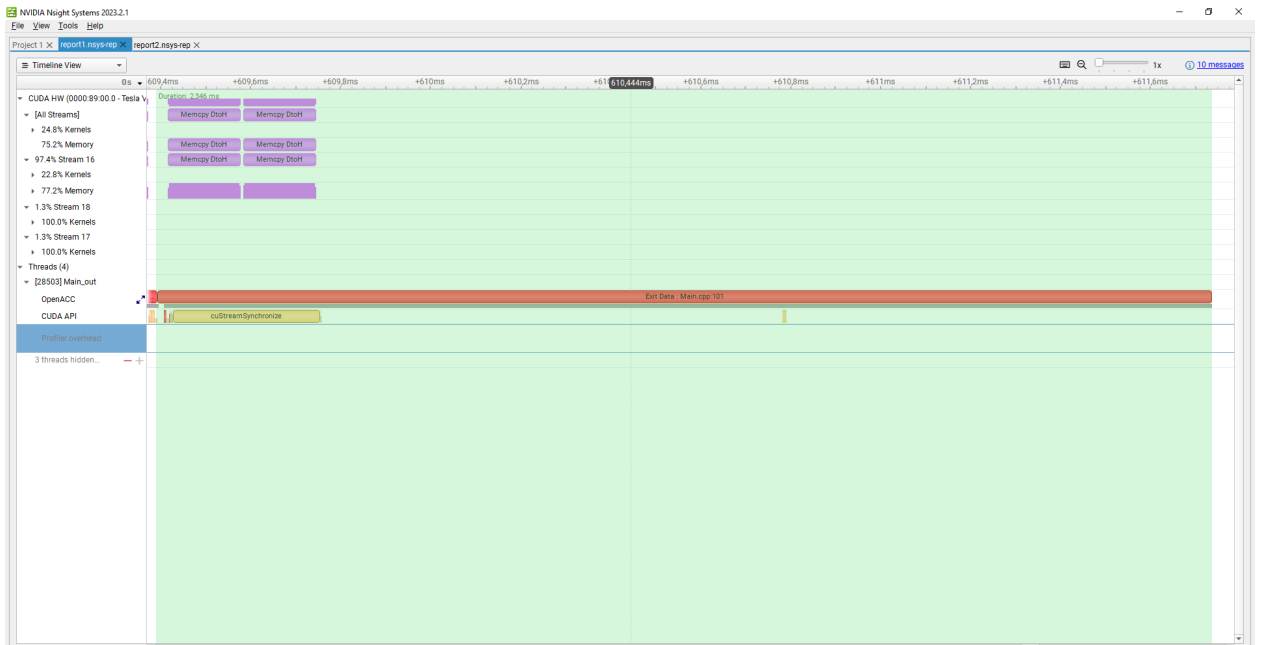
Перенос данных на GPU + выделение памяти



Общий вид процесса расчета сетки 512x512 с 10ю итерациями



Одна итерация расчета сетки

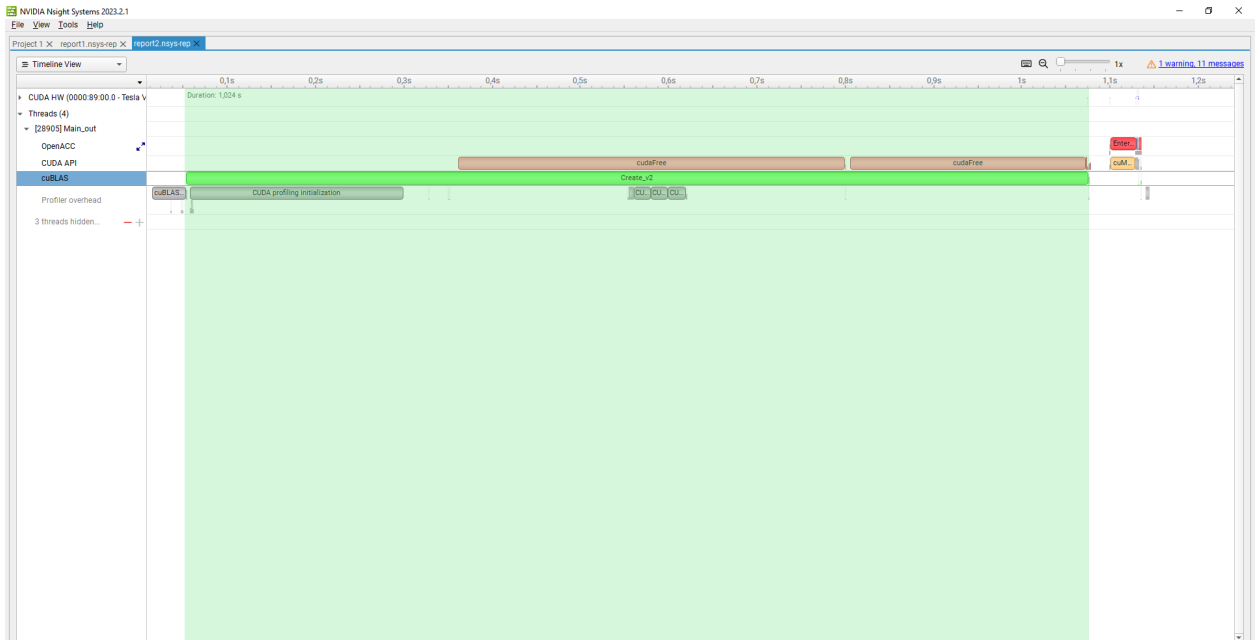


Копирование данных с памяти GPU

Профилирование реализации с cuBLAS



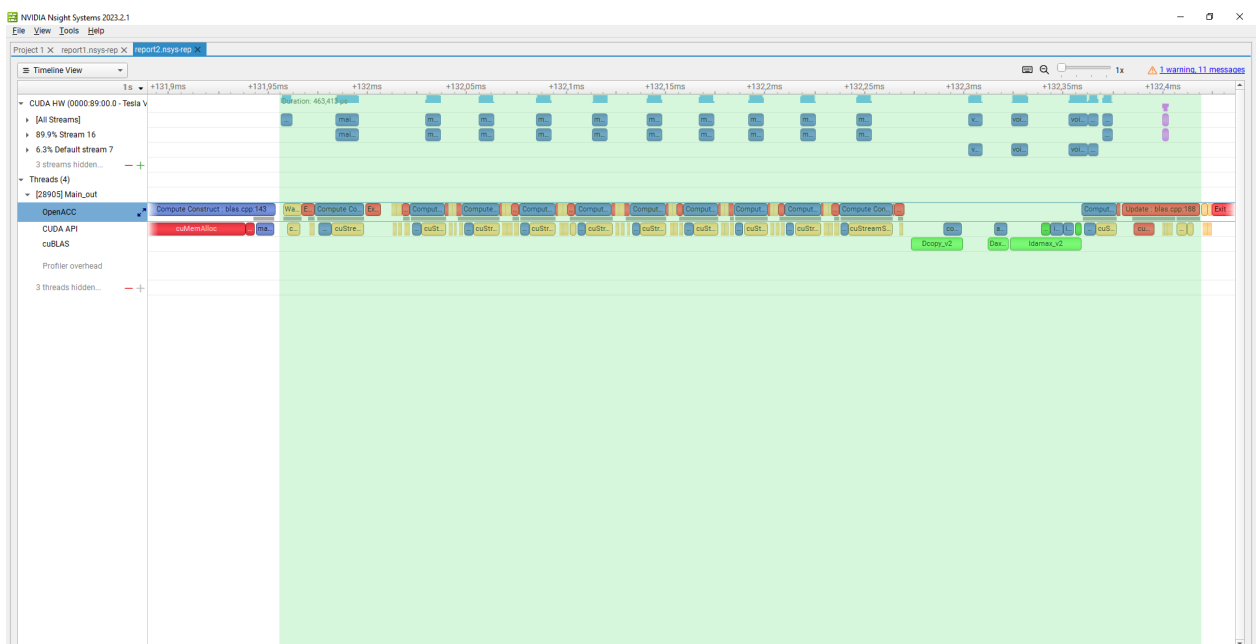
Общий вид профилирования



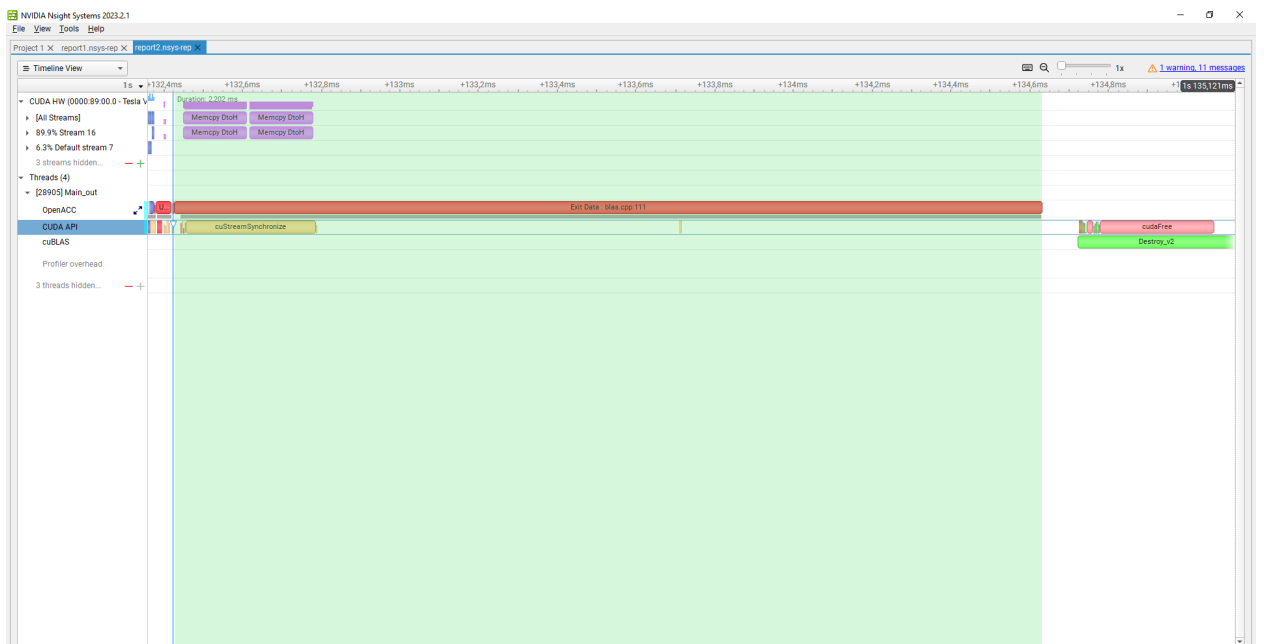
Инициализация контекста cuBLAS



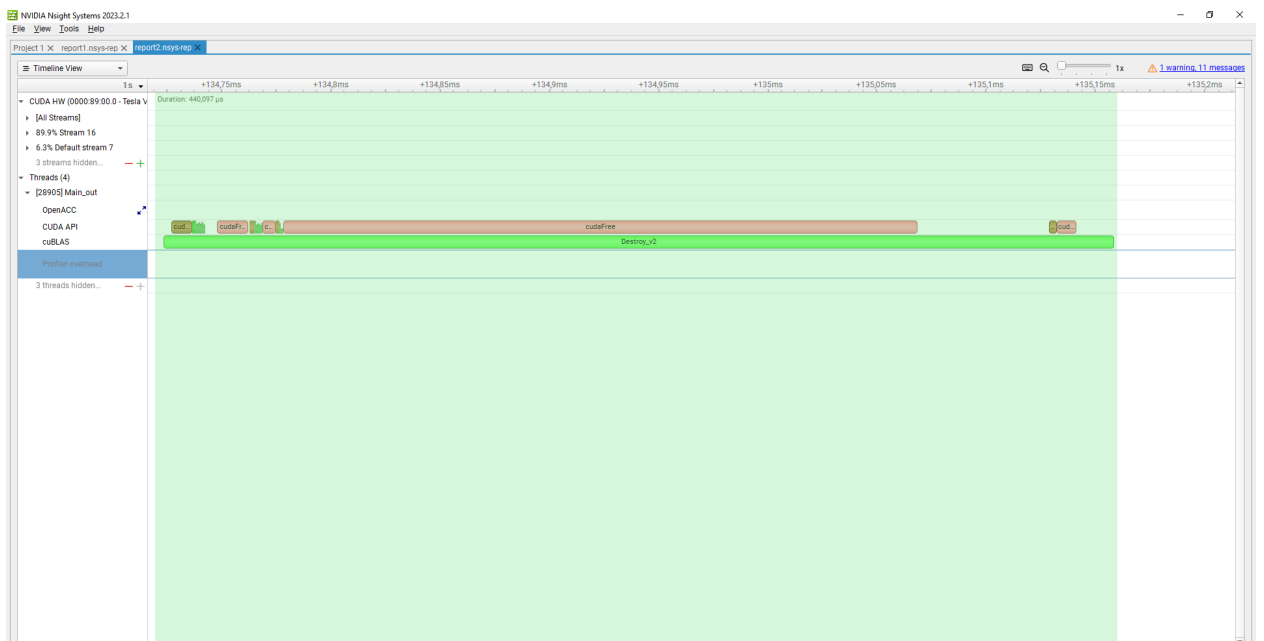
Перенос данных на GPU + выделение памяти



Общий вид процесса расчета сетки 512x512 с 10ю итерациями



Копирование данных с памяти GPU



Освобождение памяти контекста cuBLAS