

# Теория параллелизма

## Отчёт

### Уравнение теплопроводности на CUDA

Выполнил: Лейсле Александр Геннадьевич, группа 21932

05.04.2023

## 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 CUDA
3. Операцию редукции (вычисление максимального значения ошибки) на графическом процессоре реализовать через вызовы функций из библиотеки CUB
4. Сравнить скорость работы для разных размеров сеток на GPU текущей реализации и прошлой реализации через OpenACC
5. Произвести профилирование программы с использованием NsightSystems.

Используемый компилятор: nvcc

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 через OpenAcc 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

# Реализация через OpenAcc 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

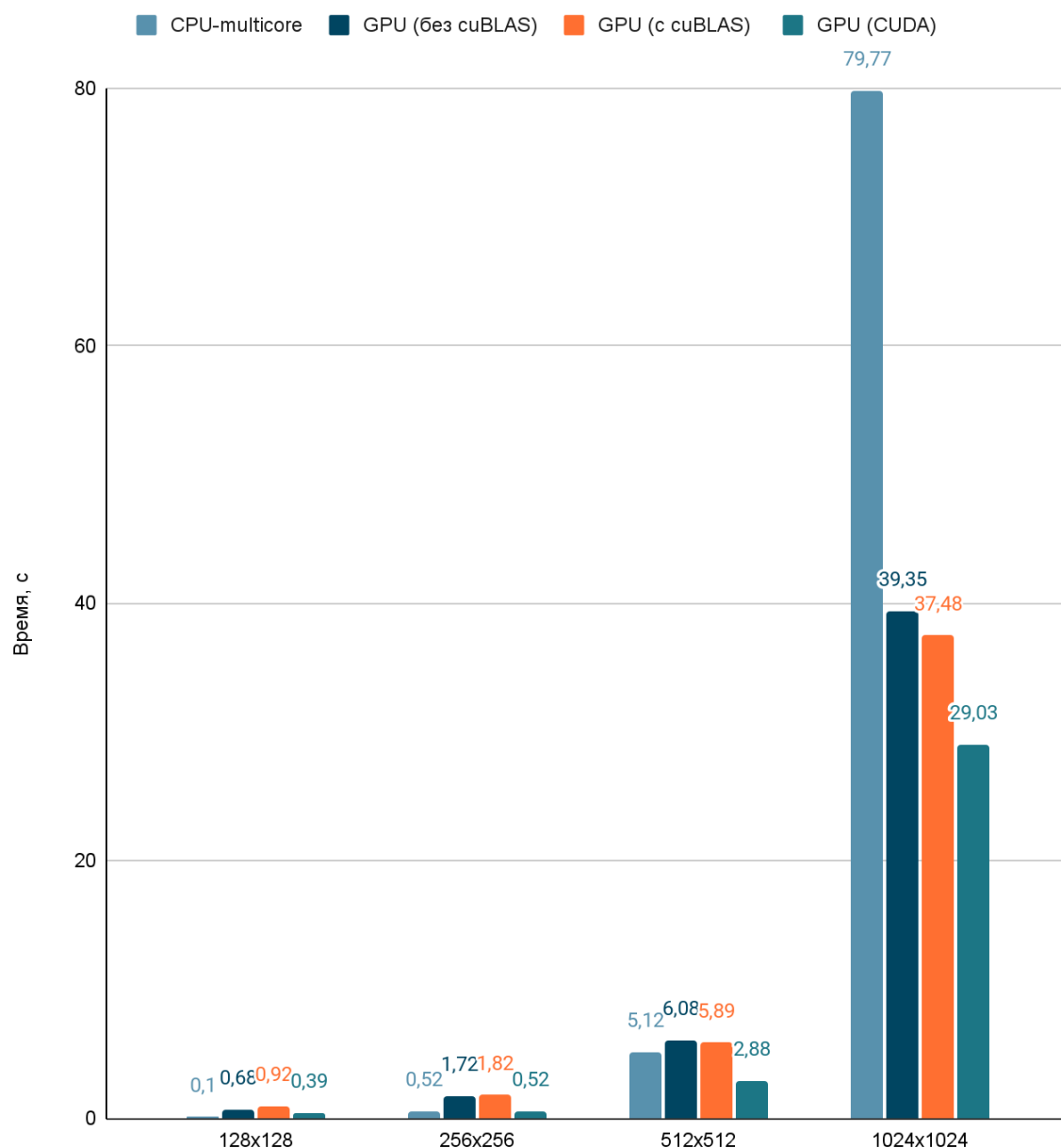
# Реализация через CUDA

## GPU

Размер сетки	Время выполнения	Точность	Количество итераций
128x128	0.39с	$9.9 * 10^{-7}$	30080
256x256	0.52с	$9.9 * 10^{-7}$	102912
512x512	2.88с	$9.9 * 10^{-7}$	339712
1024x1024	29.03с	$1.37 * 10^{-6}$	1000448

# Диаграмма сравнения времени работы CPU-multicore и GPU

## Общее время работы



# Выполнение на GPU

## Этапы оптимизации на сетке 512x512

№	Время выполнения	Точность	Кол-во итераций	Комментарии (что было сделано)
1	0.21с	0.04	256	<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 acc data copyout. По выходе из секции данные массивов копируются обратно на CPU;</p> <p>В CPU происходит выделение памяти под массив из 4 элементов, представляющих собой размеры шагов для линейной интерполяции между углами сетки по границам. Под данный массив происходит выделение памяти на GPU с помощью acc data create;</p> <p>В память GPU копируются значения: общего количества элементов сетки, размера стороны сетки и максимальная ошибка, которая</p>



				<p>высчитывается между итерациями (в начальный момент превышает заданную ошибку на 1);</p> <p>Обращение к индексам массива происходит с помощью макроса ID, 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 async. При этом циклы, в которых происходит 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 на CPU, а затем с помощью</p>
--	--	--	--	---

				<p>функции <code>acc_attach</code> происходит пересвязывание адресов этих сеток в памяти GPU;</p> <p>Происходит обновление значения максимальной ошибки в памяти CPU, если прошедшее количество итераций кратно половине размера стороны сетки или пройдено заданное количество итераций, с помощью <code>acc_update host</code></p>
2	0.56с	0.04	256	<p>Реализация с использованием директив OpenACC 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 <code>cuBLAS</code>;</p> <p>Текущее состояние сетки сохраняется во временный массив путем копирования сетки с помощью функции <code>cublasDcopy</code>;</p> <p>Из нового состояния сетки вычитается скопированный массив с помощью <code>cublasDaxpy</code>, результат записывается на место скопированного массива;</p>

				<p>В конце с помощью функции <code>cublasIdamax</code> в переменную <code>id</code> записывается адрес самого большого по модулю значения, которое представляет собой значение наибольшей ошибки между двумя итерациями.</p> <p>В переменную <code>max_error</code> записывается значение наибольшей ошибки</p> <p>После обновления текущего состояния сетки происходит обновление максимальной ошибки в памяти CPU.</p>
3	0.32с	0.04	256	<p>Реализация с использованием функций CUDA;</p> <p>В памяти CPU инициализируются сетка на текущем и новом шагах, а также значение максимальной ошибки;</p> <p>Помимо этого объявляются указатели для соответствующих переменных в памяти GPU, а также для временного массива, который будет хранить разницу между сетками на новом и текущим шагами;</p> <p>Границы сеток на текущем и новом шагах заполняются на CPU;</p> <p>Для сеток, максимальной ошибки и временного массива в памяти GPU выделяется память с помощью <code>cudaMalloc</code>;</p>

			<p>Данные сеток и максимальной ошибки переносятся в память GPU с помощью функции <code>cudaMemcpy</code>;</p> <p>Находится необходимое количество байт для нахождения максимума в массиве разниц двух сеток с помощью <code>cub::DeviceReduce::Max</code>, при этом массив для временных вычислений равен <code>NULL</code>;</p> <p>Выделяется память в GPU для массива временных вычислений с помощью <code>cudaMalloc</code>;</p> <p>Происходит объявление графа, а также происходит инициализация размера сетки и размера блока сетки;</p> <p>Количество потоков в блоке устанавливается равным <code>8x8</code>;</p> <p>Происходит инициализация графа: начинается режим записи в граф с помощью <code>cudaStreamBeginCapture</code>, затем создается цикл, в котором помещается дважды вычисление нового шага сетки. Запись прекращается с помощью <code>cudaStreamEndCapture</code>, затем создается экземпляр графа с помощью <code>cudaGraphInstantiate</code>;</p> <p>В цикле запускается проход по графу, после чего находится абсолютная разница между текущим и предыдущими шагами, находится максимальная ошибка с помощью <code>cub::DeviceReduce::Max</code> и в конце</p>
--	--	--	---

				<p>происходит обновление максимальной ошибки в памяти CPU;</p> <p>По выходу из цикла данные сетки на текущем шаге копируются в память CPU.</p>
--	--	--	--	--

## Вывод

В ходе работы было реализовано решение уравнения теплопроводности на сетках размера 128x128, 256x256, 512x512, 1024x1024. Реализация была перенесена на GPU с помощью функций CUDA. Нахождение максимальной ошибки было реализовано с помощью функций библиотеки cub. Было произведено сравнение скорости выполнения реализаций через OpenAcc, cuBLAS и CUDA на разных размерах сеток. Для полученных реализаций была выполнена профилировка с помощью NsightSystem.

Использование функций CUDA, а в частности использование графов, позволяет значительно ускорить выполнение программы, за счет того, что kernels записываются в граф и в дальнейшем не происходит расходов на их вызов, следовательно, между выполнениями отдельных kernels находятся очень маленькие промежутки. Помимо этого, использование CUDA позволяет использовать вычислительные способности GPU по максимуму за счет применения блоков с потоками, что позволяет получить выигрыш во времени как на маленьких, так и на больших сетках.

# Приложение

[https://github.com/HerrPhoton/Heat\\_equation\\_on\\_CUDA](https://github.com/HerrPhoton/Heat_equation_on_CUDA)

```
● 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 через OpenAcc

```
● 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 22.1428 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 через cuBLAS

```

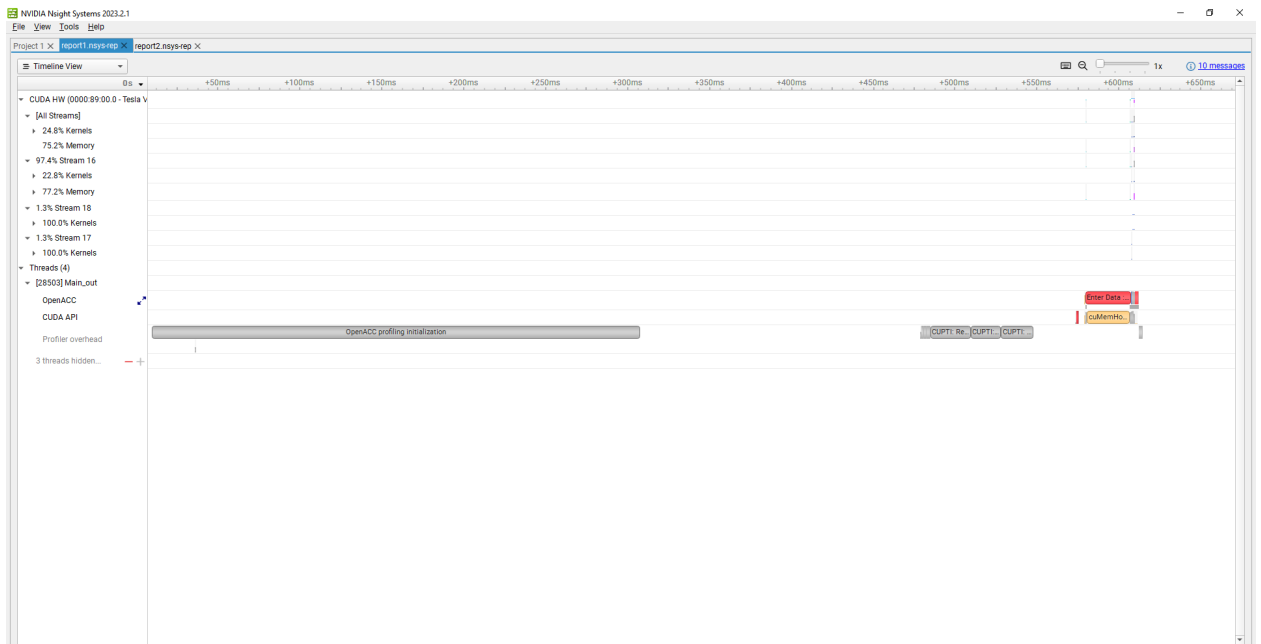
a.leisle@b14e941a1b5c:~/Lab4$ ./out -side 15 -error 0.000001 -iters 1000000 -show
Passed iterations: 476/1000000
Maximum error: 8.18288e-07/1e-06
Total execution time: 279.783 ms
10 10.7143 11.4286 12.1429 12.8571 13.5714 14.2857 15 15.7143 16.4286 17.1429 17.8571 18.5714 19.2857 20
10.7143 11.4286 12.1429 12.8571 13.5714 14.2857 15 15.7143 16.4286 17.1429 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 16.4286 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 16.4286 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 17.1429 17.8571 18.5714 19.2857 20 20.7143 21.4285 22.1428 22.8571 23.5714 24.2857 25 25.7143 26.4286
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
a.leisle@b14e941a1b5c:~/Lab4$

```

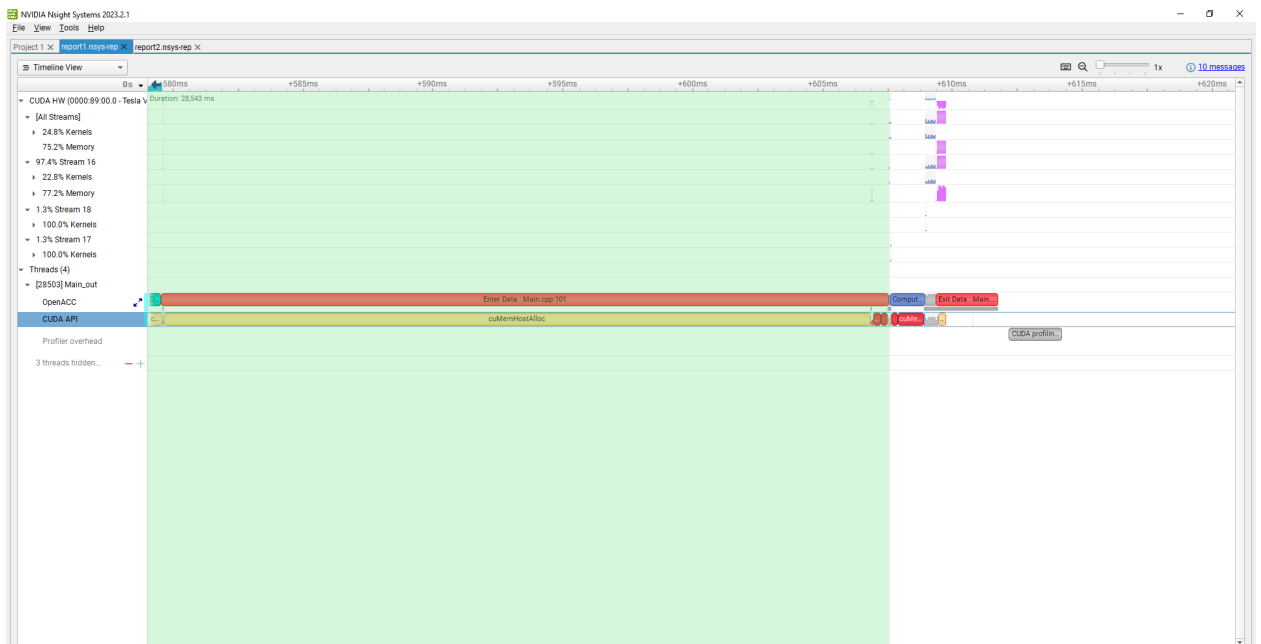
Вывод сетки 15x15 на реализации через CUDA



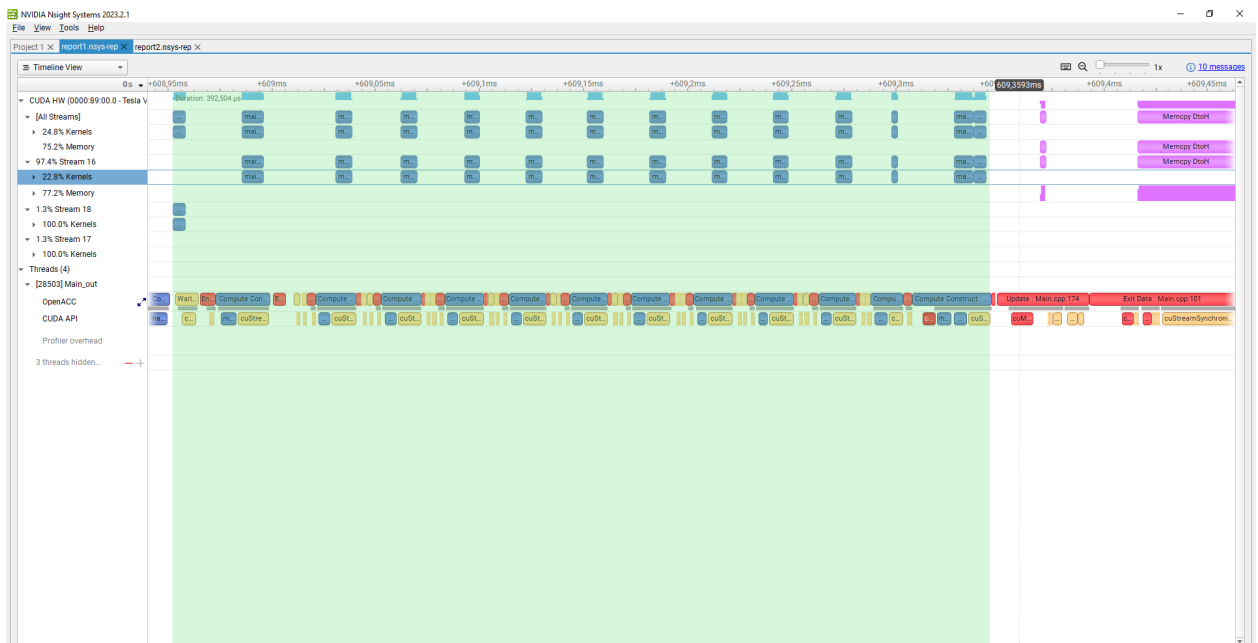
# Профилирование реализации через OpenAcc



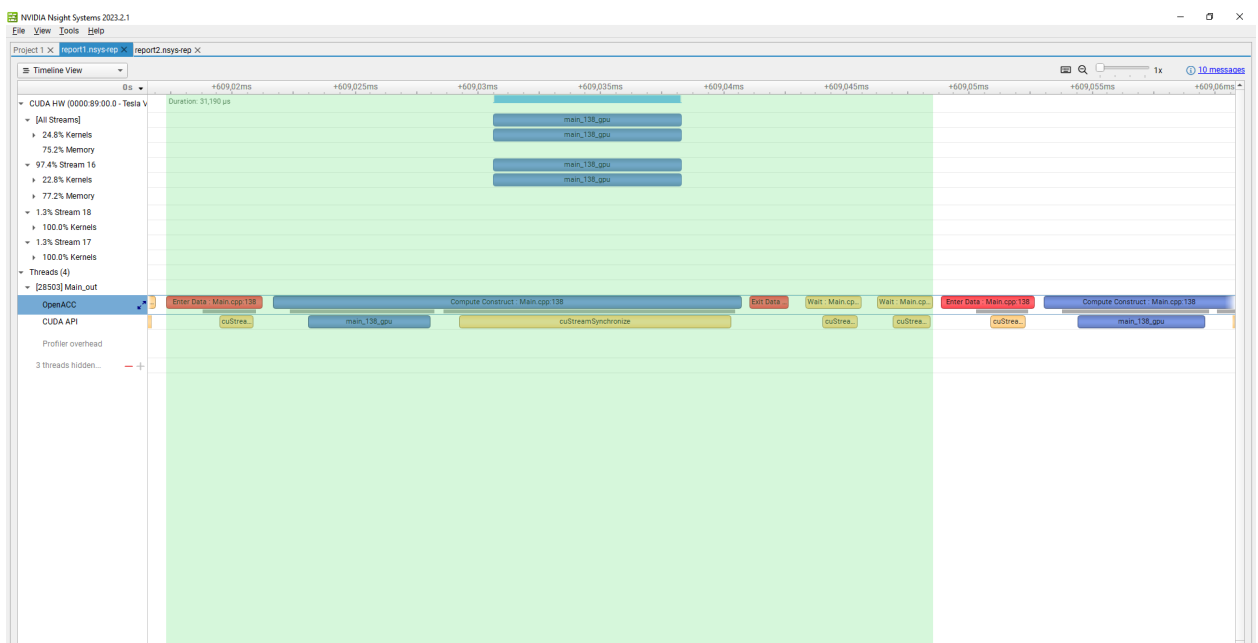
## Общий вид профилирования



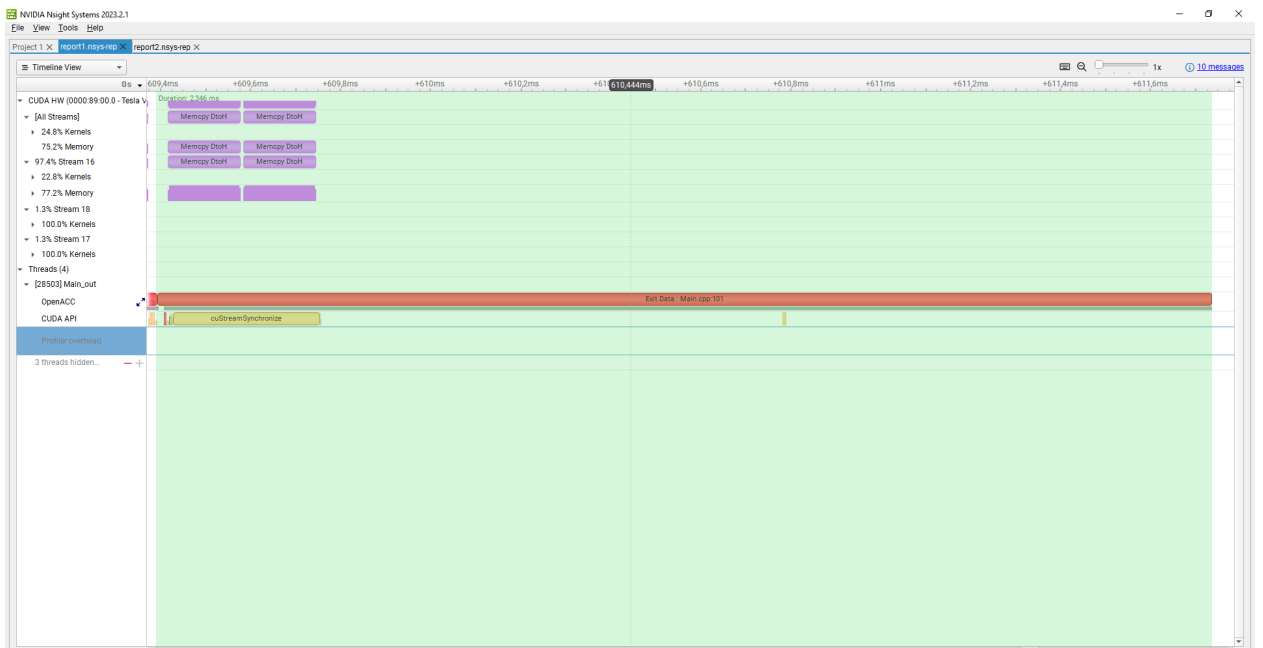
## Перенос данных на GPU + выделение памяти



### Общий вид процесса расчета сетки 512x512 с 10ю итерациями



## Одна итерация расчета сетки



Копирование данных с памяти GPU

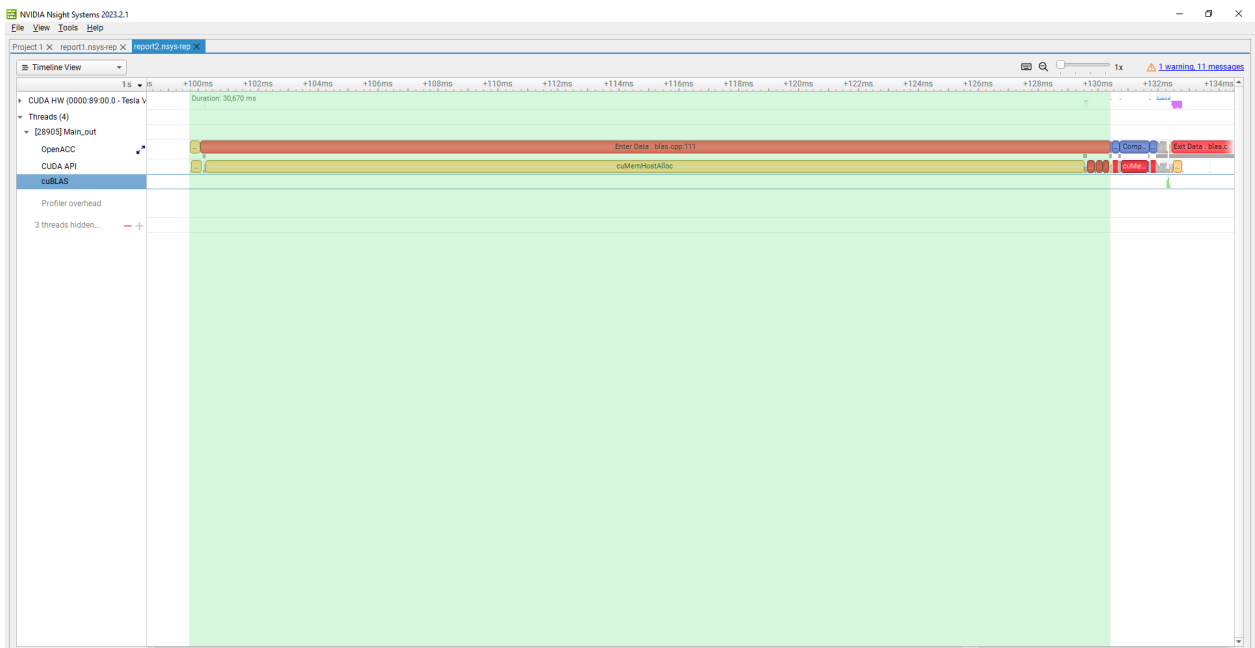
# Профилирование реализации через cuBLAS



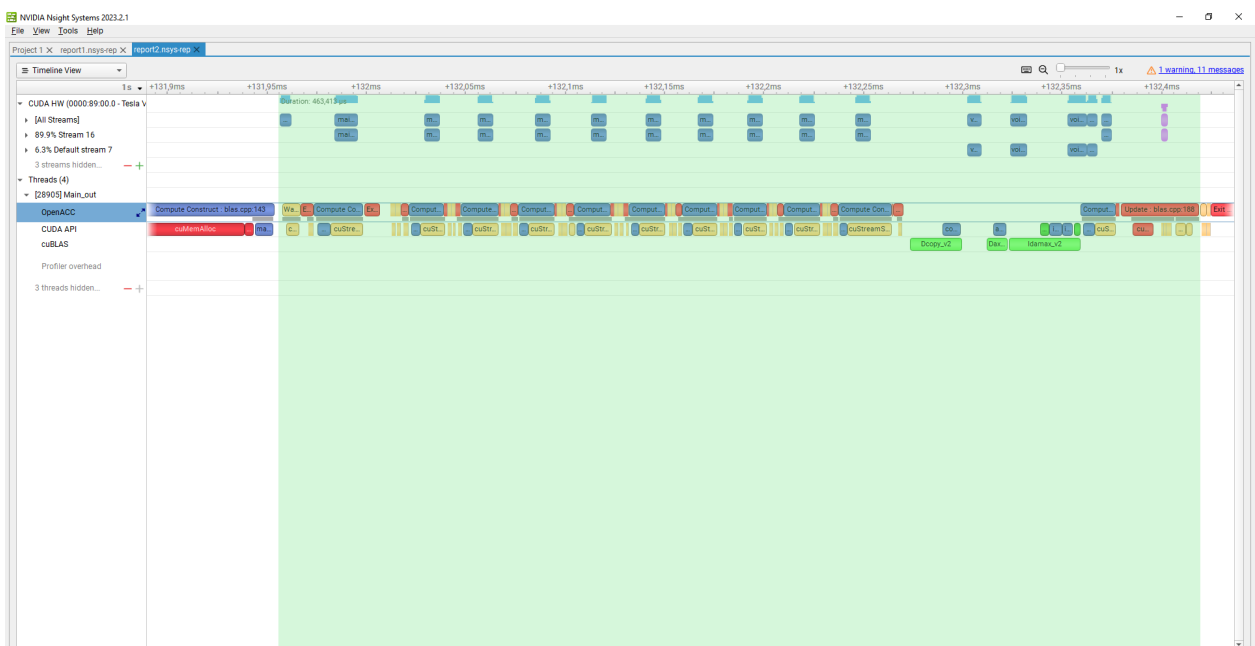
## Общий вид профилирования



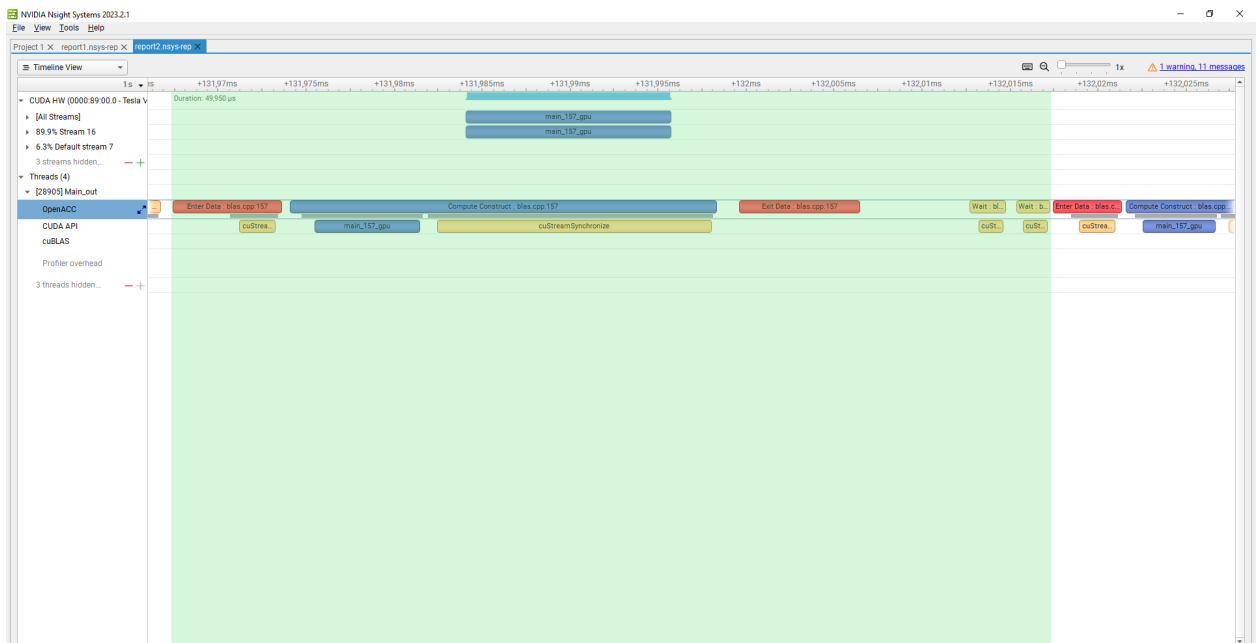
## Инициализация контекста cuBLAS



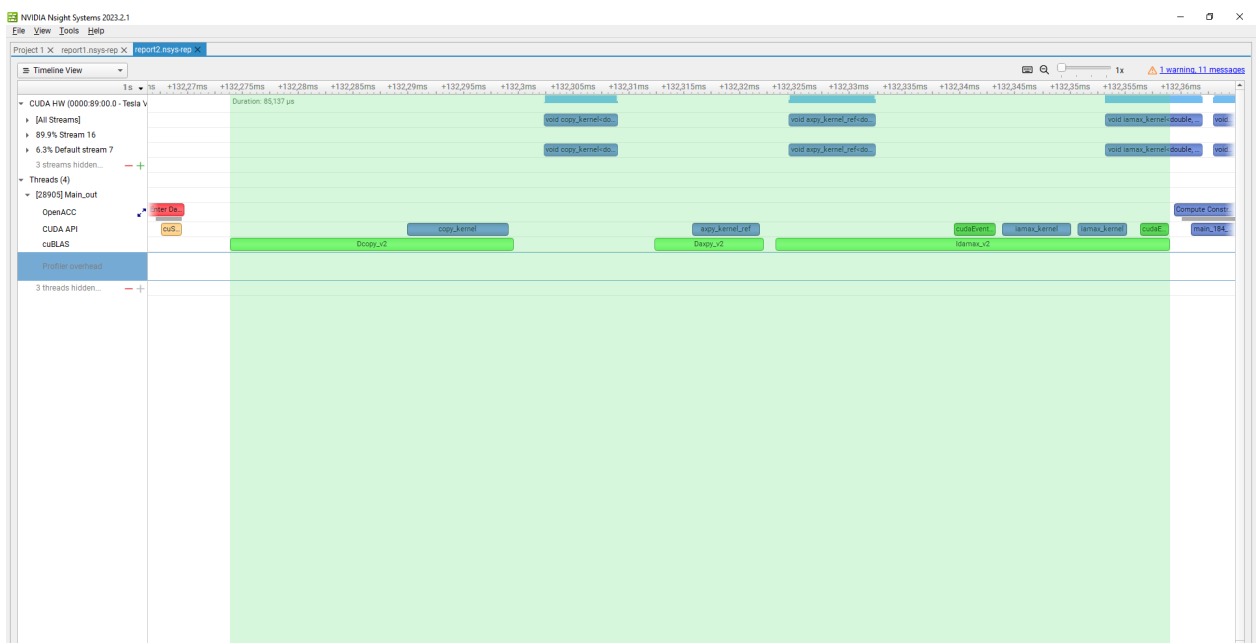
Перенос данных на GPU + выделение памяти

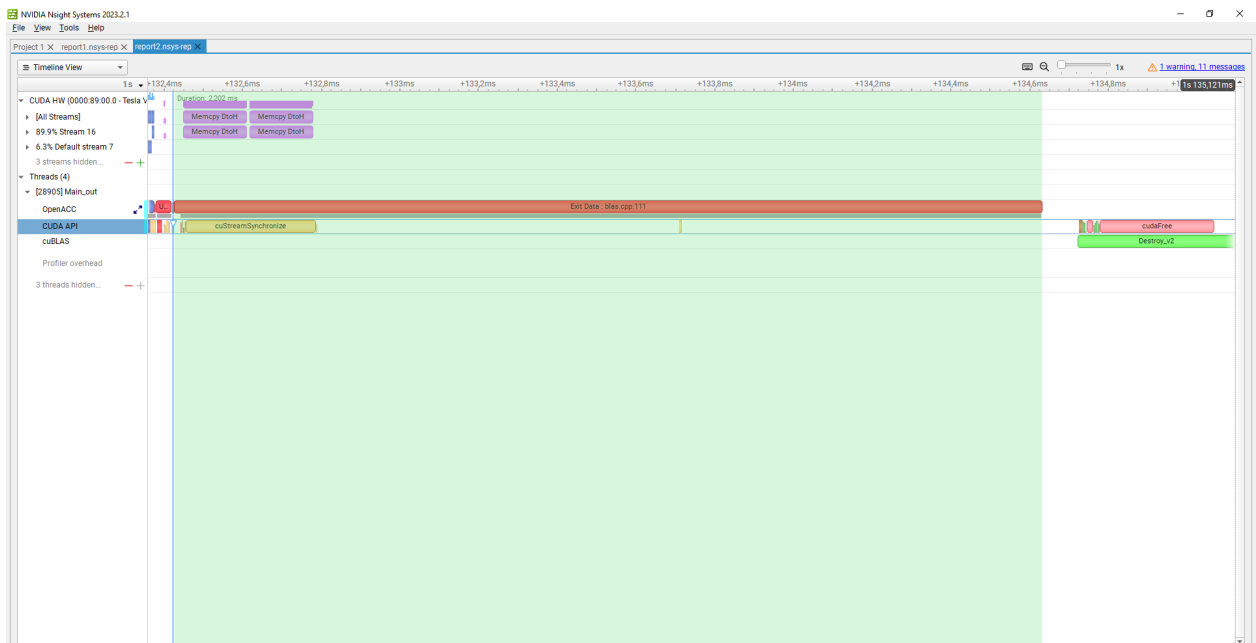


Общий вид процесса расчета сетки 512x512 с 10ю итерациями

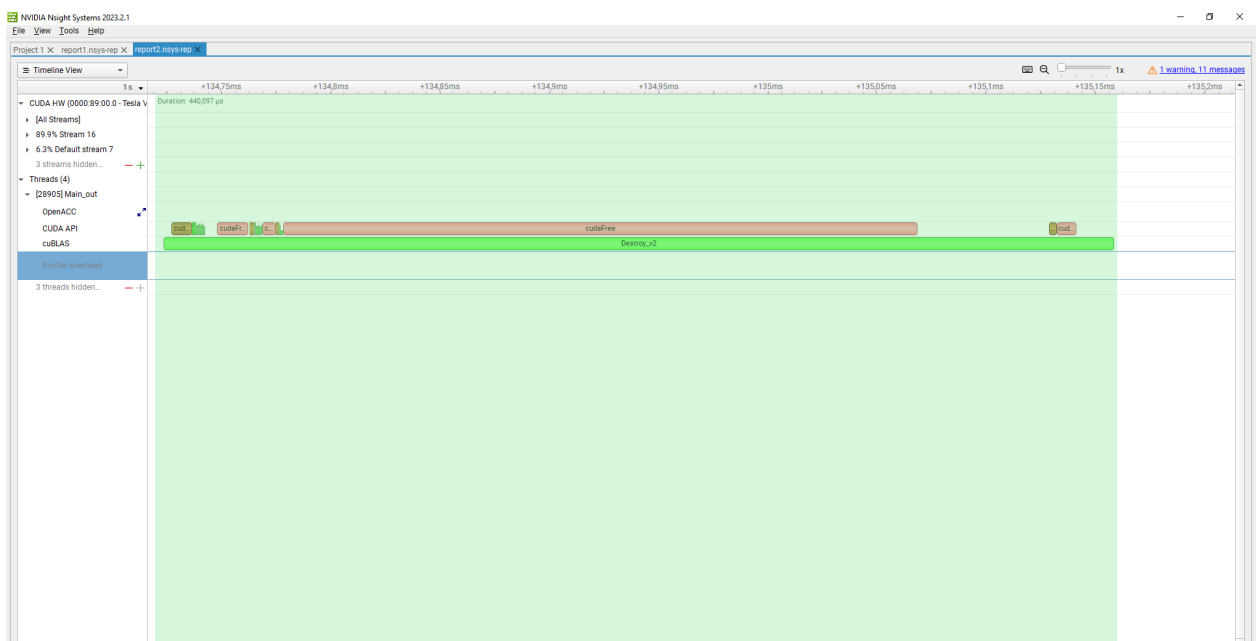


## Одна итерация расчета сетки



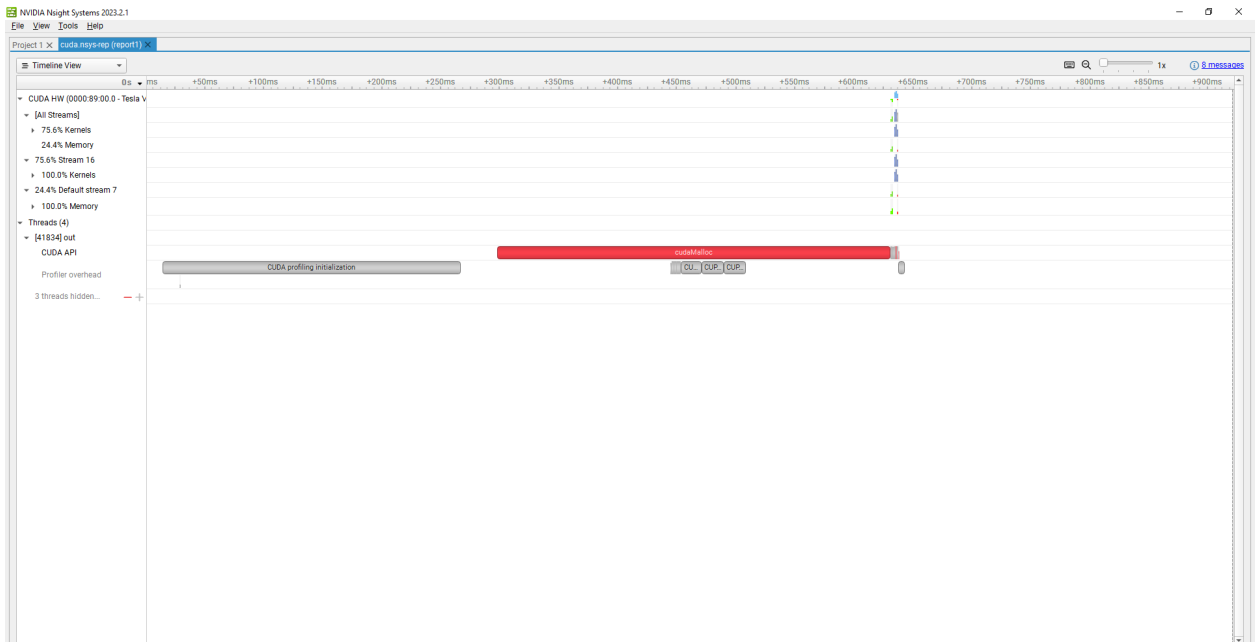


## Копирование данных с памяти GPU

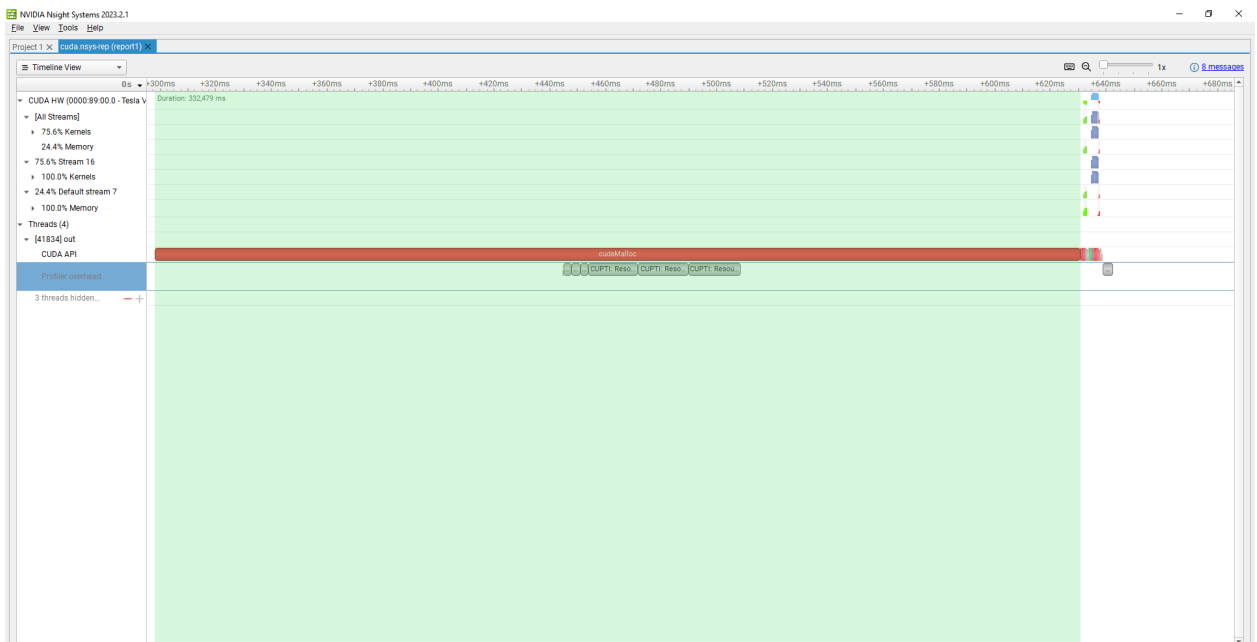


## Освобождение памяти контекста cuBLAS

# Профилирование реализации через CUDA



## Общий вид профилирования

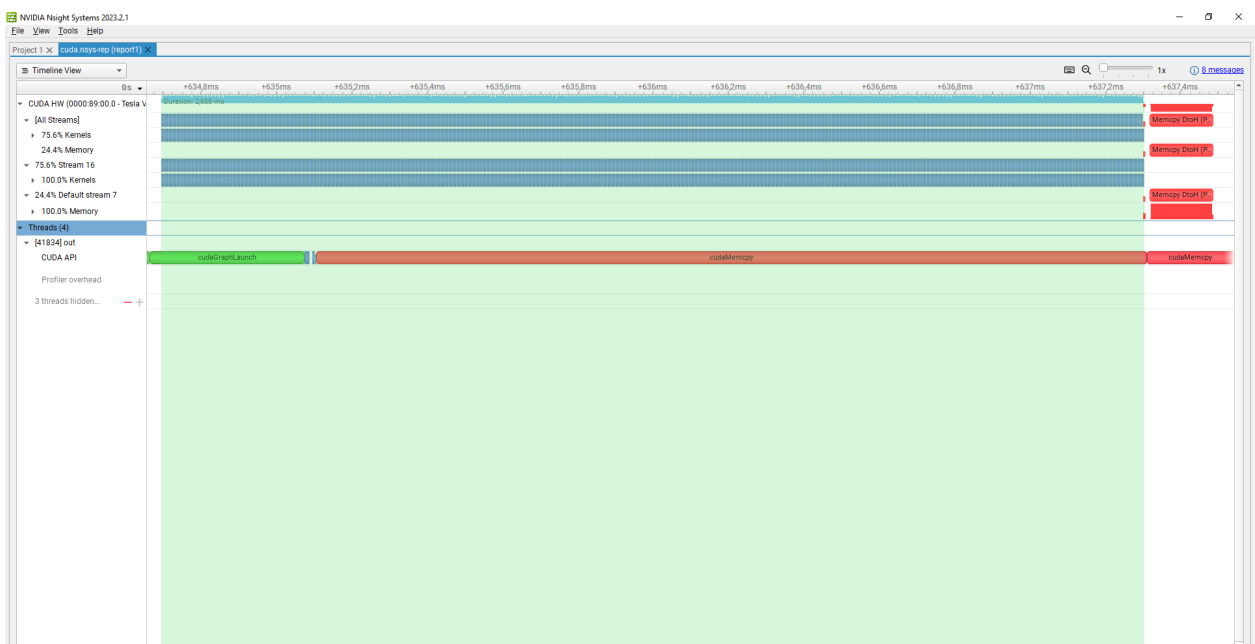


## Перенос данных на GPU + выделение памяти

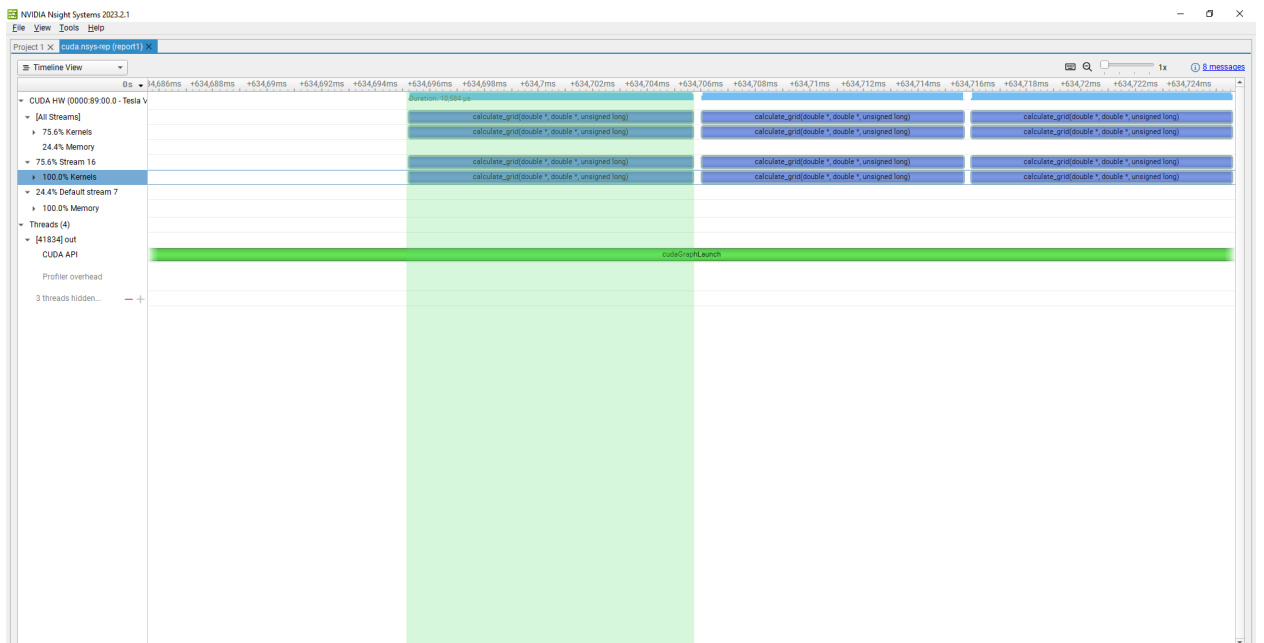




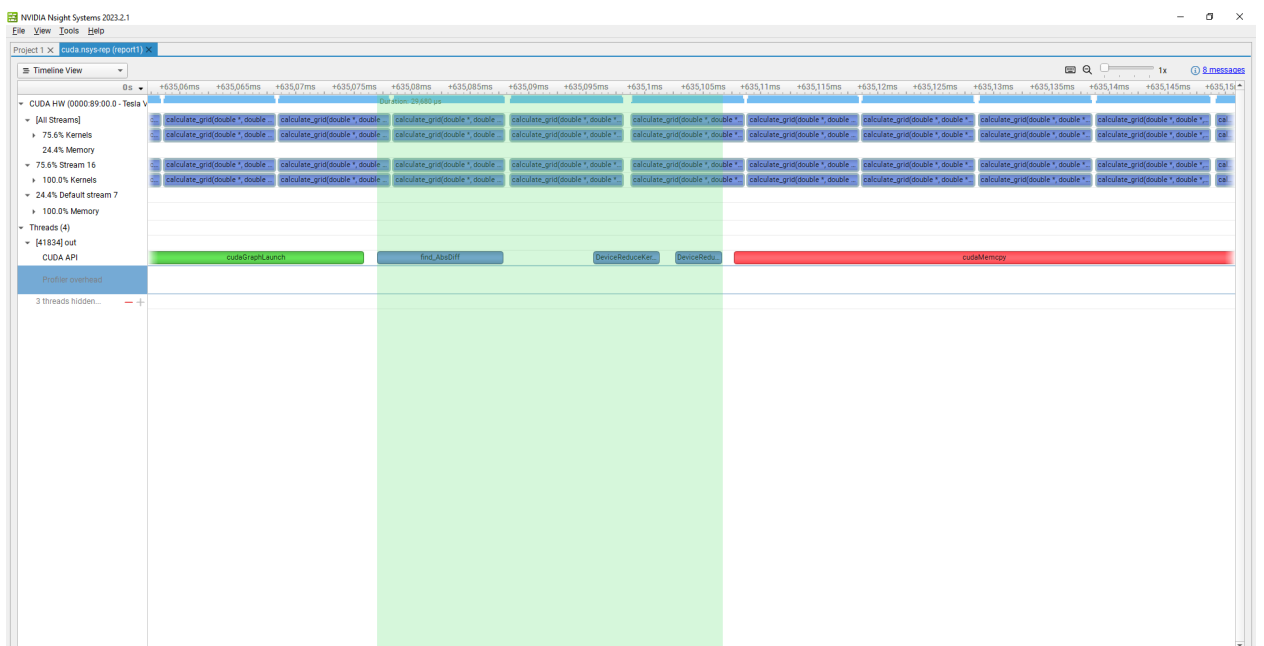
Инициализация графа



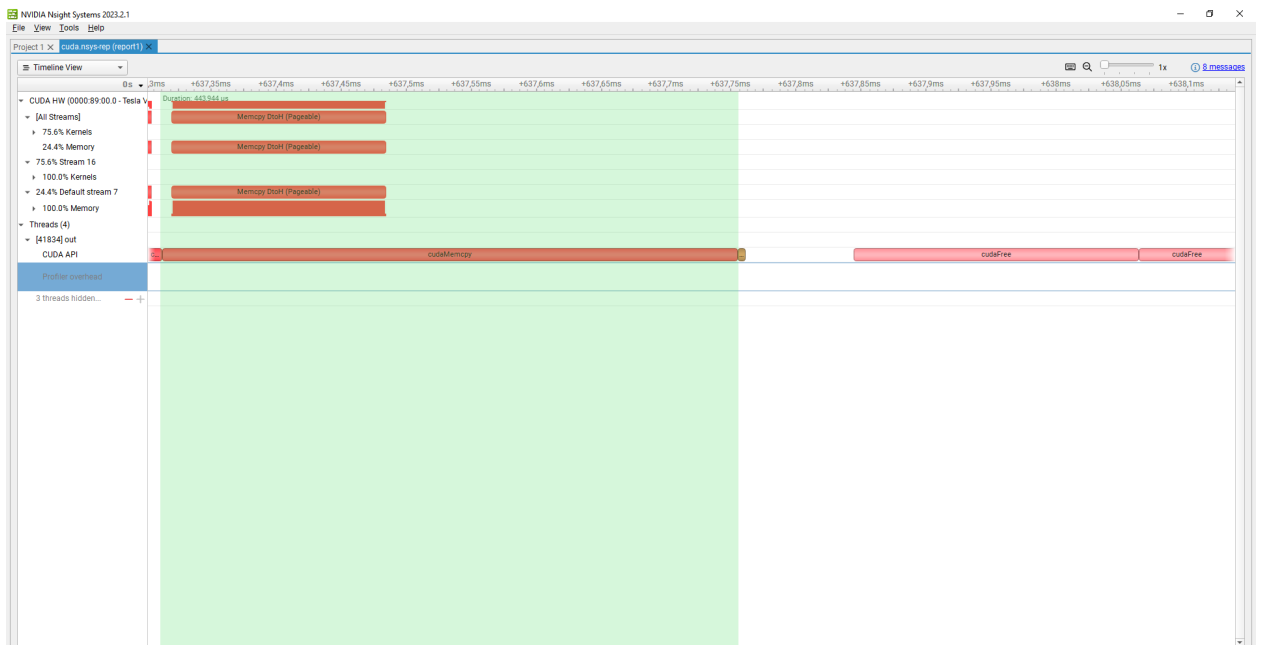
Общий вид процесса расчета сетки 512x512 с 256ю итерациями



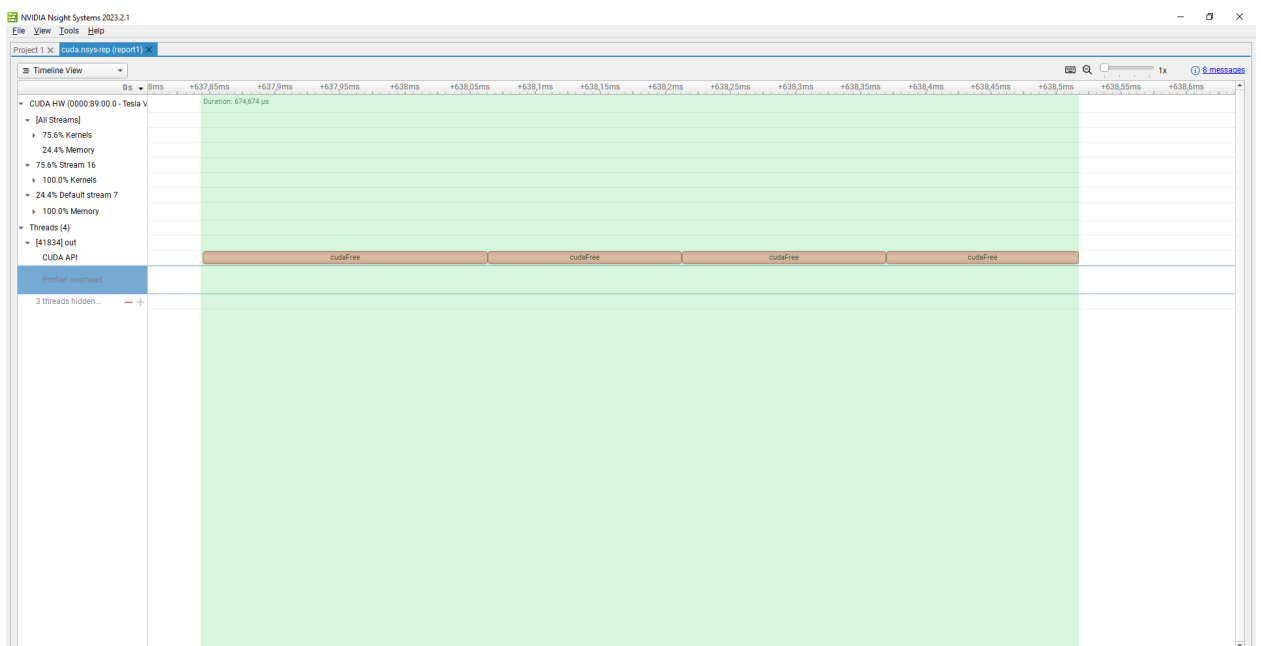
Одна итерация расчета сетки



Нахождение максимальной ошибки



## Копирование данных с памяти GPU



## Освобождение памяти на GPU