

Министерство науки и высшего образования РФ  
Федеральное государственное автономное  
образовательное учреждение высшего образования  
**«СИБИРСКИЙ ФЕДЕРАЛЬНЫЙ УНИВЕРСИТЕТ»**

Институт космических и информационных технологий

Кафедра вычислительной техники

**ОТЧЁТ О ПРАКТИЧЕСКОЙ РАБОТЕ №5**

**по дисциплине**

**«Гибридные вычислительные системы»**

Преподаватель

\_\_\_\_\_  
подпись, дата

С. А. Тарасов

инициалы, фамилия

Студент

КИ22-07Б, 032212677

номер группы, зачетной книжкой

\_\_\_\_\_  
подпись, дата

Л. А. Глушков

инициалы, фамилия

Студент

КИ22-07Б, 032215583

номер группы, зачетной книжкой

\_\_\_\_\_  
подпись, дата

А. М. Коробков

инициалы, фамилия

Студент

КИ22-07Б, 032214653

номер группы, зачетной книжкой

\_\_\_\_\_  
подпись, дата

И. О. Бердин

инициалы, фамилия

Красноярск 2025

## СОДЕРЖАНИЕ

Введение .....	3
Цель работы .....	3
Задание .....	3
Ход работы .....	5
Ключевые фрагменты кода .....	5
Результаты выполнения программы .....	5
Заключение.....	7
Приложение А.....	8

# ВВЕДЕНИЕ

## Цель работы

Закрепить навык работы с разделяемой памятью и примитивами синхронизации нитей. Освоить программирование уровня warp'а и технику ширококестования регистров.

## Задание

1. Разработать кёрнел `kernel_veccred_nobr` (1), который принимает экземпляр `VectorView` по значению и вычисляет сумму элементов вектора. Для хранения частичных сумм блоков использовать разделяемую память.

2. Разработать кёрнел `kernel_veccred_br` (2), который принимает экземпляр `VectorView` по значению и вычисляет сумму элементов вектора, используя функцию `__shfl_down_sync`. Для хранения частичных сумм warp'ов использовать разделяемую память.

3. Используя фреймворк Google Test, разработать модульные тесты для функций (1) и (2) с размерами векторов  $n \in \{1, 2, 3, 127, 129, 512, 541, 1037\}$ . В качестве эталона для сравнения использовать результат аналогичной операции для `Eigen::Matrix (sum)`; для верификации результатов применять макрос `EXPECT_NEAR` с абсолютной точностью 10–4.

4. Используя фреймворк Google Benchmark, разработать бенчмарки для функций (1) и (2) с размерами векторов  $n \in \{8 \cdot 2^0, 8 \cdot 2^1, 8 \cdot 2^2, \dots, 8 \cdot 2^{28}\}$ . Бенчмарки должны игнорировать время, затраченное на выделение, копирование и освобождение памяти. Для корректного измерения времени выполнения CUDA-кода необходимо использовать CUDA Events API.

5. Построить графики реальной вычислительной сложности для (1) и (2).

6. Построить график ускорения (2) относительно (1).

7. Объяснить экспериментальные результаты.

8. Подготовить отчёт, содержащий: • ключевые фрагменты кода; • ссылку на репозиторий с полной реализацией; • графики результатов измерений; • анализ и интерпретацию полученных результатов.

## Ход работы

### Ключевые фрагменты кода

Ключевые моменты кода:

- Реализация VectorOperations (исходный код приведен в приложении А);
- Реализация VectorOperationsKernel (исходный код приведен в приложении А).

### Результаты выполнения программы

На рисунках 1 и 2 представлены графики, основанные на результате выполнения программы.

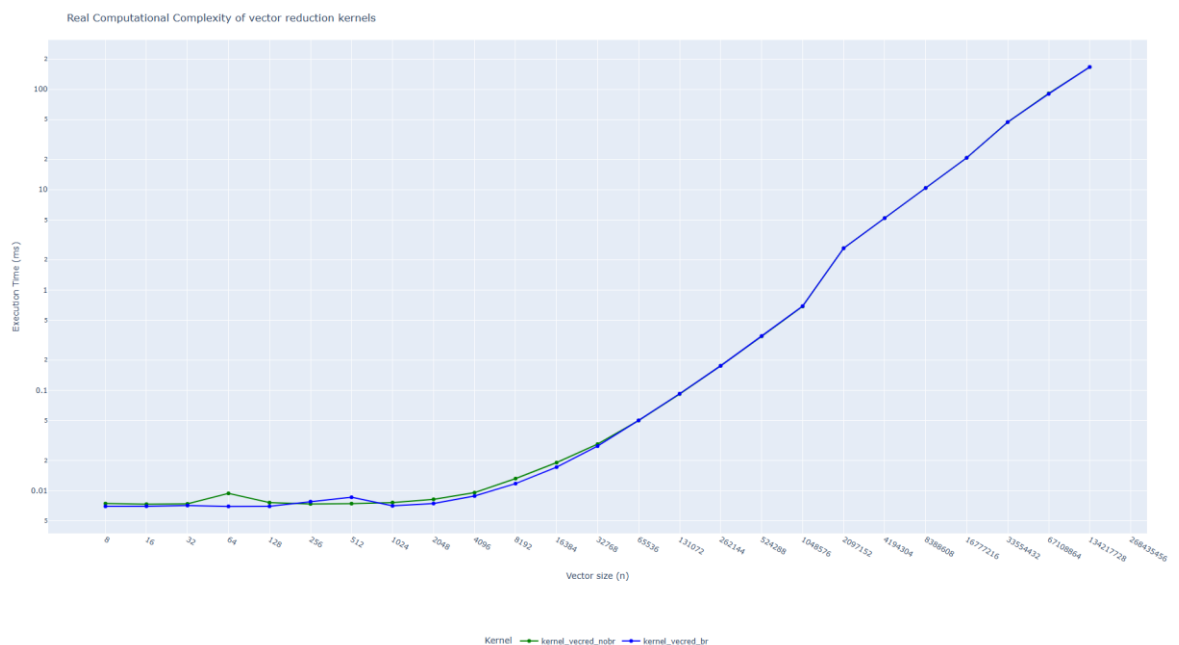


Рисунок 1 – График реальной вычислительной сложности

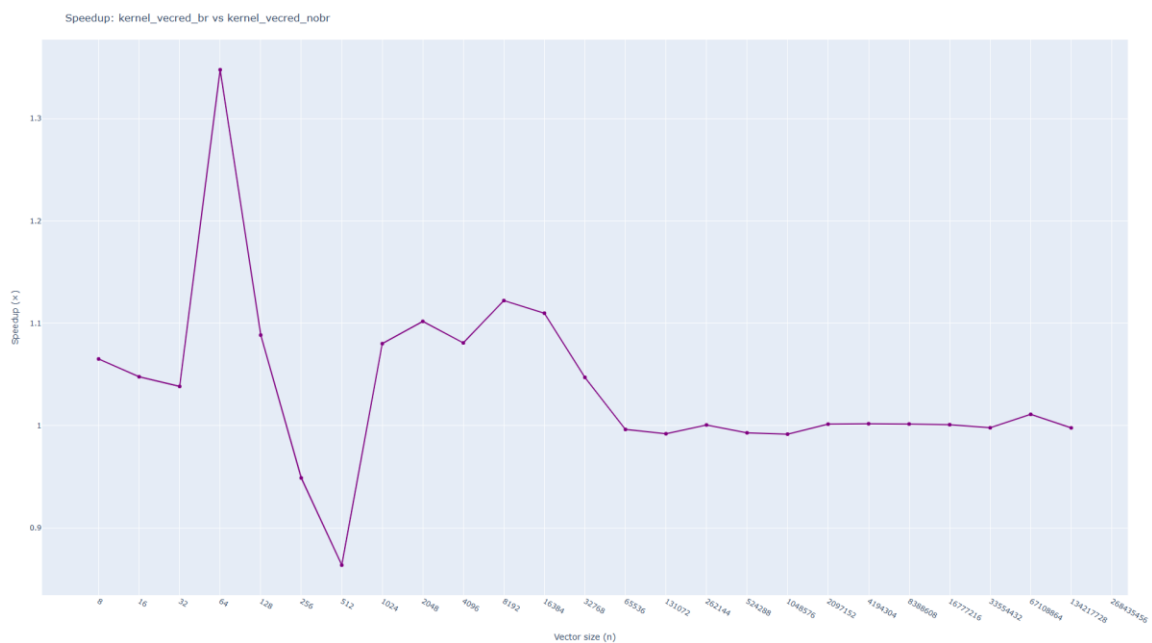


Рисунок 2 – График реальной ускорения

## **ЗАКЛЮЧЕНИЕ**

В результате выполнения данной работы были закреплены навыки работы с разделяемой памятью и примитивами синхронизации нитей. Освоено программирование уровня warp'а и техника широковещания регистров.

## ПРИЛОЖЕНИЕ А

```
1  #pragma once
2  #include "vector.cuh"
3  #include "vectorOperationsKernel.cuh"
4
5  template<typename T, AddAlgorithm Algorithm>
6  T Vector<T, Algorithm>::sum() const {
7      std::size_t n = view().size();
8
9      if (n == 0) return static_cast<T>(0);
10
11     constexpr std::size_t block_size = 256;
12
13     std::size_t nBlocks = ((n + block_size - 1) / block_size < 128) ? ((n + block_size - 1) / block_size) : 128;
14
15     Data<T> blockSum(nBlocks);
16     Data<T> result(1);
17
18     T* blockSumData = blockSum.data();
19     T* resultData = result.data();
20
21     if constexpr (Algorithm == AddAlgorithm::nobl) {
22         vectorAddNoblKernel<T><<<nBlocks, block_size, block_size * sizeof(T)>>>(view(), blockSumData, n);
23     } else if constexpr (Algorithm == AddAlgorithm::br) {
24         vectorAddBrKernel<T><<<nBlocks, block_size, ((block_size + 31) / 32) * sizeof(T)>>>(view(), blockSumData, n);
25     }
26
27     cudaError_t err = cudaGetLastError();
28     if (err != cudaSuccess) {
29         throw std::runtime_error("CUDA error 1: " + std::string(cudaGetErrorString(err)));
30     }
31
32
33     if constexpr (Algorithm == AddAlgorithm::nobl) {
34         finalReductionNoblKernel<T><<<1, block_size, block_size * sizeof(T)>>>(blockSumData, resultData, n);
35     } else if constexpr (Algorithm == AddAlgorithm::br) {
36         finalReductionBrKernel<T><<<1, block_size, ((block_size + 31) / 32) * sizeof(T)>>>(blockSumData, resultData, n);
37     }
38
39     cudaError_t err2 = cudaGetLastError();
40     if (err2 != cudaSuccess) {
41         throw std::runtime_error("CUDA error 2: " + std::string(cudaGetErrorString(err2)));
42     }
43
44     T final;
45     result.copyToHost(&final);
46
47     cudaDeviceSynchronize();
48     return final;
49 }
```

Рисунок А.1 – Реализация VectorOperations



```

5  template<typename T>
6  __global__ void vectorAddNobrKernel(VectorView<T> a, T* blockSum, std::size_t n) {
7      extern __shared__ T sharedMem[];
8
9      std::size_t tid = threadIdx.x;
10     std::size_t blockSize = blockDim.x;
11
12     T threadSum = static_cast<T>(0);
13
14     for (std::size_t i = tid; i < n; i += blockSize) {
15         threadSum += a[i];
16     }
17
18     sharedMem[tid] = threadSum;
19     __syncthreads();
20
21     for (int i = blockSize / 2; i > 0; i >>= 1) {
22         if (tid < i) {
23             sharedMem[tid] += sharedMem[tid + i];
24         }
25         __syncthreads();
26     }
27
28     if (tid == 0) {
29         blockSum[tid] = sharedMem[0];
30     }
31 }

```

Рисунок А.2 – Реализация VectorOperationsKernel

```

33  template<typename T>
34  __global__ void finalReductionNobrKernel(T* blockSum, T* sum, std::size_t n) {
35      extern __shared__ T sharedMem[];
36      std::size_t tid = threadIdx.x;
37
38      if (tid < n) {
39          sharedMem[tid] = blockSum[tid];
40      } else {
41          sharedMem[tid] = static_cast<T>(0);
42      }
43      __syncthreads();
44
45      for (int i = blockDim.x / 2; i > 0; i >>= 1) {
46          if (tid < i) {
47              sharedMem[tid] += sharedMem[tid + i];
48          }
49          __syncthreads();
50      }
51
52      if (tid == 0) {
53          *sum = sharedMem[0];
54      }
55  }

```

Рисунок А.3 – Реализация VectorOperationsKernel

```

57 template<typename T>
58 __global__ void vectorAddBrKernel(VectorView<T> a, T* blockSum, std::size_t n) {
59     extern __shared__ T warpSum[];
60
61     std::size_t tid = threadIdx.x;
62     std::size_t blockSize = blockDim.x;
63     std::size_t lane = tid % 32;
64     std::size_t warp = tid / 32;
65     std::size_t warpPerBlock = (blockDim.x + 31) / 32;
66
67     T threadSum = static_cast<T>(0);
68     for (std::size_t i = tid; i < n; i += blockSize) {
69         threadSum += a[i];
70     }
71
72     T warpS = threadSum;
73     unsigned mask = 0xFFFFFFFF;
74
75     for (int j = 16; j > 0; j /= 2) {
76         warpS += __shfl_down_sync(mask, warpS, j);
77     }
78
79     if (lane == 0) {
80         warpSum[warp] = warpS;
81     }
82     __syncthreads();
83
84     if (warp == 0) {
85         T val = (lane < warpPerBlock) ? warpSum[lane] : static_cast<T>(0);
86
87         for (int i = 16; i > 0; i /= 2) {
88             val += __shfl_down_sync(mask, val, i);
89         }
90
91         if (lane == 0) {
92             blockSum[tid] = val;
93         }
94     }
95 }

```

Рисунок А.3 – Реализация VectorOperationsKernel

```

97  template<typename T>
98  __global__ void finalReductionBrKernel(T* blockSum, T* sum, std::size_t n) {
99      extern __shared__ T warpSum[];
100
101      std::size_t tid = threadIdx.x;
102      std::size_t blockSize = blockDim.x;
103      std::size_t lane = tid % 32;
104      std::size_t warp = tid / 32;
105      std::size_t warpPerBlock = (blockDim.x + 31) / 32;
106
107      T threadVal = static_cast<T>(0);
108      if (tid < n) {
109          threadVal = blockSum[tid];
110      }
111
112      T warpS = threadVal;
113      unsigned mask = 0xFFFFFFFF;
114
115      for (int j = 16; j > 0; j /= 2) {
116          warpS += __shfl_down_sync(mask, warpS, j);
117      }
118
119      if (lane == 0) {
120          warpSum[warp] = warpS;
121      }
122      __syncthreads();
123
124      if (warp == 0) {
125          T val = (lane < warpPerBlock) ? warpSum[lane] : static_cast<T>(0);
126
127          for (int i = 16; i > 0; i /= 2) {
128              val += __shfl_down_sync(mask, val, i);
129          }
130
131          if (lane == 0) {
132              *sum = val;
133          }
134      }
135  }

```

Рисунок А.4 – Реализация VectorOperationsKernel