# Computational Science on Many-Core Architectures: Exercise 9

Viktor Beck, 11713110

December 2022

## Task 1: Libraries

Below we can see the comparison of the performances of the different libraries and the custom CUDA and OpenCL implementations. The time was measured such that only the computation of the dot product and the addition / subtraction of the vectors x and y is taken into account.

It seems that the custom CUDA implementation (from exercise 2) is the fastest but the amount of code written for such a small task is quite a lot in comparison to the few lines which were necessary for ViennaCL. The custom OpenCL implementation has the most lines of code and still performs second worst. In general one could say that the ViennaCL library has the best "effort to performance" ratio - at least in this case.
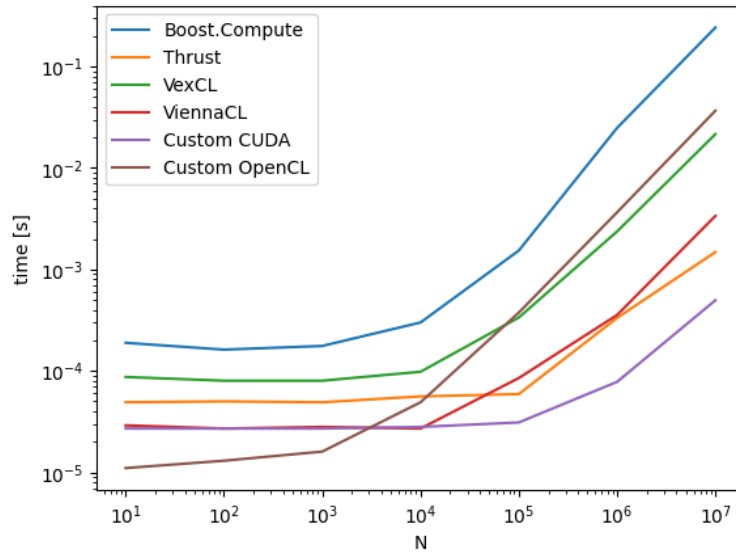


Figure 1: Comparing execution times for different libraries and custom implementations of a dot product for different vector sizes N on the RTX GPU

## Task 2: HIP

Somehow, I was not able to install and use HIP therefore, and converted it manually as described in lecture 9 (with good old "ctrl+F"). Also, I could not test it... You can find the respective code snippet in the appendix.

## Appendix

### Task 1:

**Boost.Compute**

```cpp
// specify use of OpenCL 1.2:
#define CL_TARGET_OPENCL_VERSION  120
#define CL_MINIMUM_OPENCL_VERSION 120

#include <vector>
#include <algorithm>
#include <iostream>

#include "timer.hpp"

#include <boost/compute/algorithm/transform.hpp>
#include <boost/compute/algorithm/inner_product.hpp>
#include <boost/compute/container/vector.hpp>
#include <boost/compute/functional/math.hpp>

namespace compute = boost::compute;

int main()
{

    // get default device and setup context
    compute::device device = compute::system::default_device();
    compute::context context(device);
    compute::command_queue queue(context, device);

    Timer timer;

    for (int k = 0; k <=7; k++){

        int N = pow(10, k);



        // generate data on the host
        std::vector<float> h_x(N);
        std::fill(h_x.begin(), h_x.end(), 1);
        std::vector<float> h_y(N);
        std::fill(h_y.begin(), h_y.end(), 2);

        // create a vector on the device
        compute::vector<float> d_x (N, context);
        compute::vector<float> d_y (N, context);

        // transfer data from the host to the device
        compute::copy(h_x.begin(), h_x.end(), d_x.begin(), queue);
        compute::copy(h_y.begin(), h_y.end(), d_y.begin(), queue);
```

```
48
49              timer.reset();
50
51              // x+y
52              compute::vector<float> x_plus_y(N, context);
53
54              compute::transform(
55                  d_x.begin(),
56                  d_x.end(),
57                  d_y.begin(),
58                  x_plus_y.begin(),
59                  compute::plus<float>(),
60                  queue
61              );
62
63              // x-y
64              compute::vector<float> x_minus_y(N, context);
65              compute::transform(
66                  d_x.begin(),
67                  d_x.end(),
68                  d_y.begin(),
69                  x_minus_y.begin(),
70                  compute::minus<float>(),
71                  queue
72              );
73
74              float res = compute::inner_product(x_plus_y.begin(), x_plus_y.end(),
                        x_minus_y.begin(), 0.f, queue);
75
76              //std::cout << res << std::endl;
77
78              std::cout << timer.get() << std::endl;
79          }
80
81          return 0;
82      }
```

## Thrust

```
1   #include <thrust/host_vector.h>
2   #include <thrust/device_vector.h>
3   #include <thrust/sort.h>
4   #include <thrust/transform.h>
5   #include <thrust/inner_product.h>
6   #include <cstdlib>
7
8   #include <vector>
9   #include "timer.hpp"
10
11  int main(void)
12  {
13    Timer timer;
14
15    for (int k = 0; k <= 7; k++)
16    {
17      int N = pow(10, k);
18
19      // x and y on host
20      thrust::host_vector<int> h_x(N);
21      thrust::fill(h_x.begin(), h_x.end(), 1);
```

```
22      thrust::host_vector<int> h_y(N);
23      thrust::fill(h_y.begin(), h_y.end(), 2);
24
25      // transfer data to the device
26      thrust::device_vector<int> d_x = h_x;
27      thrust::device_vector<int> d_y = h_y;
28
29      timer.reset();
30
31      // x+y
32      thrust::device_vector<float> x_plus_y(N);
33
34      thrust::transform(
35          d_x.begin(),
36          d_x.end(),
37          d_y.begin(),
38          x_plus_y.begin(),
39          thrust::plus<float>());
40
41      // x-y
42      thrust::device_vector<float> x_minus_y(N);
43
44      thrust::transform(
45          d_x.begin(),
46          d_x.end(),
47          d_y.begin(),
48          x_minus_y.begin(),
49          thrust::minus<float>());
50
51      float res = thrust::inner_product(x_plus_y.begin(), x_plus_y.end(), x_minus_y
            .begin(), 0.f);
52
53      // std::cout << res << std::endl;
54
55      std::cout << timer.get() << std::endl;
56  }
57
58  return 0;
59 }
```

### VexCL

```
1
2  // The following three defines are necessary to pick the correct OpenCL version
       on the machine:
3  #define VEXCL_HAVE_OPENCL_HPP
4  #define CL_HPP_TARGET_OPENCL_VERSION 120
5  #define CL_HPP_MINIMUM_OPENCL_VERSION 120
6
7  #include <iostream>
8  #include <stdexcept>
9  #include <vexcl/vexcl.hpp>
10
11 #include "timer.hpp"
12 #include <vector>
13
14 int main()
15 {
16     vex::Context ctx(vex::Filter::GPU && vex::Filter::DoublePrecision);
17
```

```
18        std::cout << ctx << std::endl; // print list of selected devices
19
20        Timer timer;
21
22        for (int k = 0; k <= 7; k++)
23        {
24
25            int N = pow(10, k);
26
27            std::vector<double> h_x(N, 1.0), h_y(N, 2.0);
28
29            vex::vector<double> d_x(ctx, h_x);
30            vex::vector<double> d_y(ctx, h_y);
31
32
33            timer.reset();
34
35            vex::vector<double> x_plus_y = d_x + d_y;
36            vex::vector<double> x_minus_y = d_x - d_y;
37
38            vex::Reductor<double, vex::SUM> sum(ctx);
39            double res = sum(x_plus_y * x_minus_y);
40
41            // std::cout << res << std::endl;
42
43            std::cout << timer.get() << std::endl;
44        }
45        return 0;
46 }
```

## ViennaCL

```
1
2  #include <vector>
3  #include <iostream>
4  #include "timer.hpp"
5
6  #define VIENNACL_WITH_CUDA
7
8  #include "viennacl/vector.hpp"
9  #include "viennacl/linalg/inner_prod.hpp"
10
11 int main()
12 {
13   Timer timer;
14
15   for (int k = 0; k <= 7; k++)
16   {
17
18     int N = pow(10, k);
19
20     viennacl::vector<double> x = viennacl::scalar_vector<double>(N, 1.0);
21     viennacl::vector<double> y = viennacl::scalar_vector<double>(N, 2.0);
22
23     timer.reset();
24
25     viennacl::vector<double> x_plus_y = x + y;
26     viennacl::vector<double> x_minus_y = x - y;
27
28     double res = viennacl::linalg::inner_prod(x_plus_y, x_minus_y);
```

```
29
30      // std::cout << res << std::endl;
31
32      std::cout << timer.get() << std::endl;
33    }
34
35    return EXIT_SUCCESS;
36  }
```

## Task 2:

```cpp
1   #include "poisson2d.hpp"
2   #include "timer.hpp"
3   #include <algorithm>
4   #include <iostream>
5   #include <stdio.h>
6
7   // y = A * x
8   __global__ void hip_csr_matvec_product(int N, int *csr_rowoffsets,
9                                          int *csr_colindices, double *csr_values,
10                                         double *x, double *y)
11  {
12    for (int i = hipBlockIdx_x * hipBlockDim_x + hipThreadIdx_x; i < N; i +=
13        hipBlockDim_x * hipGridDim_x)
14    {
15      double sum = 0;
16      for (int k = csr_rowoffsets[i]; k < csr_rowoffsets[i + 1]; k++)
17      {
18        sum += csr_values[k] * x[csr_colindices[k]];
19      }
20      y[i] = sum;
21    }
22  }
23
24  // x <- x + alpha * y
25  __global__ void hip_vecadd(int N, double *x, double *y, double alpha)
26  {
27    for (int i = hipBlockIdx_x * hipBlockDim_x + hipThreadIdx_x; i < N; i +=
28        hipBlockDim_x * hipGridDim_x)
29      x[i] += alpha * y[i];
30  }
31
32  // x <- y + alpha * x
33  __global__ void hip_vecadd2(int N, double *x, double *y, double alpha)
34  {
35    for (int i = hipBlockIdx_x * hipBlockDim_x + hipThreadIdx_x; i < N; i +=
36        hipBlockDim_x * hipGridDim_x)
37      x[i] = y[i] + alpha * x[i];
38  }
39
40  // result = (x, y)
41  __global__ void hip_dot_product(int N, double *x, double *y, double *result)
42  {
43    __shared__ double shared_mem[512];
44
45    double dot = 0;
46    for (int i = hipBlockIdx_x * hipBlockDim_x + hipThreadIdx_x; i < N; i +=
47        hipBlockDim_x * hipGridDim_x)
48    {
49      dot += x[i] * y[i];
```

```
46     }
47
48     shared_mem[hipThreadIdx_x] = dot;
49     for (int k = hipBlockDim_x / 2; k > 0; k /= 2)
50     {
51       __syncthreads();
52       if (hipThreadIdx_x < k)
53       {
54         shared_mem[hipThreadIdx_x] += shared_mem[hipThreadIdx_x + k];
55       }
56     }
57
58     if (hipThreadIdx_x == 0)
59       atomicAdd(result, shared_mem[0]);
60   }
61
62   ////////////// CG KERNEL 1: //////////////
63
64   __global__ void hip_cg_1(int N,
65                            double alpha,
66                            double beta,
67                            double *x,
68                            double *r,
69                            double *p,
70                            double *Ap,
71                            double *result_rr)
72   {
73
74     __shared__ double shared_mem[512];
75     double dot = 0;
76
77     for (int i = hipBlockIdx_x * hipBlockDim_x + hipThreadIdx_x; i < N; i +=
          hipBlockDim_x * hipGridDim_x)
78     {
79       // line 2,3,4: get x, r, p
80       x[i] = x[i] + alpha * p[i];
81       r[i] = r[i] - alpha * Ap[i];
82       p[i] = r[i] + beta * p[i];
83
84       // line 6: get dot(r,r)
85       dot += r[i] * r[i];
86     }
87
88     __syncthreads();
89     shared_mem[hipThreadIdx_x] = dot;
90
91     for (int k = hipBlockDim_x / 2; k > 0; k /= 2)
92     {
93       __syncthreads();
94       if (hipThreadIdx_x < k)
95       {
96         shared_mem[hipThreadIdx_x] += shared_mem[hipThreadIdx_x + k];
97       }
98     }
99
100    if (hipThreadIdx_x == 0)
101    {
102      atomicAdd(result_rr, shared_mem[0]);
103      // printf("%g", r[0]);
104    }
105  }
106
```

```
107  /////////////// CG KERNEL 2: ///////////////
108  __global__ void hip_cg_2(int N,
109                           int *csr_rowoffsets,
110                           int *csr_colindices,
111                           double *csr_values,
112                           double *p,
113                           double *Ap,
114                           double *result1,
115                           double *result2)
116  {
117
118    __shared__ double shared_mem1[512];
119    __shared__ double shared_mem2[512];
120
121    for (int i = hipBlockIdx_x * hipBlockDim_x + hipThreadIdx_x; i < N; i +=
122        hipBlockDim_x * hipGridDim_x)
123    {
124      // line 5: get Ap
125      double sum = 0;
126      for (int k = csr_rowoffsets[i]; k < csr_rowoffsets[i + 1]; k++)
127      {
128        sum += csr_values[k] * p[csr_colindices[k]];
129      }
130      Ap[i] = sum;
131    }
132
133    // line 6: get dot(Ap,Ap) and dot(p,Ap)
134    double dot1 = 0;
135    double dot2 = 0;
136    for (int i = hipBlockIdx_x * hipBlockDim_x + hipThreadIdx_x; i < N; i +=
137        hipBlockDim_x * hipGridDim_x)
138    {
139      dot1 += Ap[i] * Ap[i];
140      dot2 += p[i] * Ap[i];
141    }
142
143    shared_mem1[hipThreadIdx_x] = dot1;
144    shared_mem2[hipThreadIdx_x] = dot2;
145    for (int k = hipBlockDim_x / 2; k > 0; k /= 2)
146    {
147      __syncthreads();
148      if (hipThreadIdx_x < k)
149      {
150        shared_mem1[hipThreadIdx_x] += shared_mem1[hipThreadIdx_x + k];
151        shared_mem2[hipThreadIdx_x] += shared_mem2[hipThreadIdx_x + k];
152      }
153    }
154
155    if (hipThreadIdx_x == 0)
156    {
157      atomicAdd(result1, shared_mem1[0]);
158      atomicAdd(result2, shared_mem2[0]);
159    }
160  }
161
162  /** Implementation of the conjugate gradient algorithm.
163   *
164   *  The control flow is handled by the CPU.
165   *  Only the individual operations (vector updates, dot products, sparse
166   * matrix-vector product) are transferred to hip kernels.
167   *
168   *  The temporary arrays p, r, and Ap need to be allocated on the GPU for use
```

```cpp
 * with hip. Modify as you see fit.
 */
void conjugate_gradient(int N, // number of unknows
                        int *csr_rowoffsets, int *csr_colindices,
                        double *csr_values, double *rhs, double *solution)
//, double *init_guess)   // feel free to add a nonzero initial guess as needed
{
  // initialize timer
  Timer timer;

  // clear solution vector (it may contain garbage values):
  std::fill(solution, solution + N, 0);

  // initialize work vectors:
  double alpha, beta, residual_norm_squared, dot_pAp, dot_ApAp;
  double *hip_solution, *hip_p, *hip_r, *hip_Ap, *hip_scalar, *hip_dot_pAp, *
      hip_dot_ApAp;
  hipMalloc(&hip_p, sizeof(double) * N);
  hipMalloc(&hip_r, sizeof(double) * N);
  hipMalloc(&hip_Ap, sizeof(double) * N);
  hipMalloc(&hip_solution, sizeof(double) * N);
  hipMalloc(&hip_scalar, sizeof(double));

  hipMalloc(&hip_dot_ApAp, sizeof(double));
  hipMalloc(&hip_dot_pAp, sizeof(double));

  hipMemcpy(hip_p, rhs, sizeof(double) * N, hipMemcpyHostToDevice);
  hipMemcpy(hip_r, rhs, sizeof(double) * N, hipMemcpyHostToDevice);
  hipMemcpy(hip_solution, solution, sizeof(double) * N, hipMemcpyHostToDevice);

  // get residual_norm_squared
  const double zero = 0;
  hipMemcpy(hip_scalar, &zero, sizeof(double), hipMemcpyHostToDevice);
  hipLaunchKernelGGL(hip_dot_product, 512, 512, 0,0, N, hip_r, hip_r, hip_scalar)
      ;
  hipMemcpy(&residual_norm_squared, hip_scalar, sizeof(double),
      hipMemcpyDeviceToHost);

  double initial_residual_squared = residual_norm_squared;

  // line 1: get alpha0, beta0, Ap0
  hipLaunchKernelGGL(hip_csr_matvec_product, 512, 512, 0,0, N, csr_rowoffsets,
      csr_colindices, csr_values, hip_p, hip_Ap);

  hipMemcpy(hip_scalar, &zero, sizeof(double), hipMemcpyHostToDevice);
  hipLaunchKernelGGL(hip_dot_product, 512, 512, 0,0, N, hip_p, hip_Ap, hip_scalar
      );
  hipMemcpy(&dot_pAp, hip_scalar, sizeof(double), hipMemcpyDeviceToHost);
  alpha = residual_norm_squared / dot_pAp;

  hipMemcpy(hip_scalar, &zero, sizeof(double), hipMemcpyHostToDevice);
  hipLaunchKernelGGL(hip_dot_product, 512, 512, 0,0, N, hip_Ap, hip_Ap,
      hip_scalar);
  hipMemcpy(&dot_ApAp, hip_scalar, sizeof(double), hipMemcpyDeviceToHost);

  beta = (alpha * alpha * dot_ApAp - residual_norm_squared) /
      residual_norm_squared;


  int iters = 0;
  hipDeviceSynchronize();
```

```
222   timer.reset();
223   while (1)
224   {
225
226     hipMemcpy(hip_scalar, &zero, sizeof(double), hipMemcpyHostToDevice);
227
228     // std::cout << alpha << ", " << residual_norm_squared << std::endl;
229     hipLaunchKernelGGL(hip_cg_1, 512, 512, 0,0, N, alpha, beta, hip_solution,
                hip_r, hip_p, hip_Ap, hip_scalar);
230
231     hipMemcpy(&residual_norm_squared, hip_scalar, sizeof(double),
                hipMemcpyDeviceToHost);
232
233     // std::cout << residual_norm_squared << std::endl;
234
235     hipMemcpy(hip_dot_ApAp, &zero, sizeof(double), hipMemcpyHostToDevice);
236     hipMemcpy(hip_dot_pAp, &zero, sizeof(double), hipMemcpyHostToDevice);
237
238     // std::cout << dot_ApAp << ", " << dot_pAp << ", " << residual_norm_squared
                << std::endl;
239
240     hipLaunchKernelGGL(hip_cg_2, 512, 512, 0,0, N, csr_rowoffsets, csr_colindices
                , csr_values,
241                               hip_p, hip_Ap, hip_dot_ApAp, hip_dot_pAp);
242
243     hipMemcpy(&dot_ApAp, hip_dot_ApAp, sizeof(double), hipMemcpyDeviceToHost);
244     hipMemcpy(&dot_pAp, hip_dot_pAp, sizeof(double), hipMemcpyDeviceToHost);
245     // std::cout << alpha << ", " << residual_norm_squared << std::endl;
246
247     // line 7:
248     alpha = residual_norm_squared / dot_pAp;
249
250     // line 8:
251     beta = (alpha * alpha * dot_ApAp - residual_norm_squared) /
                residual_norm_squared;
252
253     // std::cout << dot_ApAp << ", " << dot_pAp << std::endl;
254
255     // check for convergence
256     if (std::sqrt(residual_norm_squared / initial_residual_squared) < 1e-6)
257     {
258       break;
259     }
260
261     if (iters > 10000)
262       break; // solver didn't converge
263     ++iters;
264   }
265   hipMemcpy(solution, hip_solution, sizeof(double) * N, hipMemcpyDeviceToHost);
266
267   hipDeviceSynchronize();
268   /*
269   std::cout << "Time elapsed: " << timer.get() << " (" << timer.get() / iters <<
          " per iteration)" << std::endl;
270
271   if (iters > 10000)
272     std::cout << "Conjugate Gradient did NOT converge within 10000 iterations"
273               << std::endl;
274   else
275     std::cout << "Conjugate Gradient converged in " << iters << " iterations."
276               << std::endl;
277   */
```

```
278
279    std::cout << timer.get() / iters << "," << std::endl;
280
281    hipFree(hip_p);
282    hipFree(hip_r);
283    hipFree(hip_Ap);
284    hipFree(hip_solution);
285    hipFree(hip_scalar);
286  }
287
288  /** Solve a system with 'points_per_direction * points_per_direction' unknowns
289   */
290  void solve_system(int points_per_direction)
291  {
292
293    int N = points_per_direction *
294            points_per_direction; // number of unknows to solve for
295
296    // std::cout << "Solving Ax=b with " << N << " unknowns." << std::endl;
297
298    //
299    // Allocate CSR arrays.
300    //
301    // Note: Usually one does not know the number of nonzeros in the system matrix
302    // a-priori.
303    //        For this exercise, however, we know that there are at most 5 nonzeros
304    //        per row in the system matrix, so we can allocate accordingly.
305    //
306    int *csr_rowoffsets = (int *)malloc(sizeof(double) * (N + 1));
307    int *csr_colindices = (int *)malloc(sizeof(double) * 5 * N);
308    double *csr_values = (double *)malloc(sizeof(double) * 5 * N);
309
310    int *hip_csr_rowoffsets, *hip_csr_colindices;
311    double *hip_csr_values;
312    //
313    // fill CSR matrix with values
314    //
315    generate_fdm_laplace(points_per_direction, csr_rowoffsets, csr_colindices,
316                         csr_values);
317
318    //
319    // Allocate solution vector and right hand side:
320    //
321    double *solution = (double *)malloc(sizeof(double) * N);
322    double *rhs = (double *)malloc(sizeof(double) * N);
323    std::fill(rhs, rhs + N, 1);
324
325    //
326    // Allocate hip-arrays //
327    //
328    hipMalloc(&hip_csr_rowoffsets, sizeof(double) * (N + 1));
329    hipMalloc(&hip_csr_colindices, sizeof(double) * 5 * N);
330    hipMalloc(&hip_csr_values, sizeof(double) * 5 * N);
331    hipMemcpy(hip_csr_rowoffsets, csr_rowoffsets, sizeof(double) * (N + 1),
332        hipMemcpyHostToDevice);
        hipMemcpy(hip_csr_colindices, csr_colindices, sizeof(double) * 5 * N,
333        hipMemcpyHostToDevice);
        hipMemcpy(hip_csr_values, csr_values, sizeof(double) * 5 * N,
            hipMemcpyHostToDevice);
334
335    //
336    // Call Conjugate Gradient implementation with GPU arrays
```

```cpp
337    //
338    conjugate_gradient(N, hip_csr_rowoffsets, hip_csr_colindices, hip_csr_values,
           rhs, solution);
339
340    //
341    // Check for convergence:
342    //
343    /*
344    double residual_norm = relative_residual(N, csr_rowoffsets, csr_colindices,
           csr_values, rhs, solution);
345    std::cout << "Relative residual norm: " << residual_norm
346               << " (should be smaller than 1e-6)" << std::endl;
347    */
348    hipFree(hip_csr_rowoffsets);
349    hipFree(hip_csr_colindices);
350    hipFree(hip_csr_values);
351    free(solution);
352    free(rhs);
353    free(csr_rowoffsets);
354    free(csr_colindices);
355    free(csr_values);
356 }
357
358 int main()
359 {
360    std::vector<int> N_vec = {10, 25, 50, 75, 100, 250, 500, 750, 1000};
361    for (const auto &N : N_vec)
362    {
363      solve_system(N); // solves a system with 100*100 unknowns
364    }
365
366    return EXIT_SUCCESS;
367 }
```